

**A HOLISTIC REDUNDANCY- AND INCENTIVE-BASED  
FRAMEWORK TO IMPROVE CONTENT AVAILABILITY IN  
PEER-TO-PEER NETWORKS**

by

**Octavio Herrera-Ruiz**

BSc Eng Electronics, National Autonomous University of México, 1997  
(Ingeniero Mecánico Electricista-Electrónico, UNAM, México, 1997)

M.S. Telecommunications, University of Pittsburgh, 2001

Submitted to the Graduate Faculty of  
School of Information Sciences in partial fulfillment  
of the requirements for the degree of  
**Doctor of Philosophy**

University of Pittsburgh

2012

UNIVERSITY OF PITTSBURGH  
SCHOOL OF INFORMATION SCIENCES

This dissertation was presented

by

Octavio Herrera-Ruiz

It was defended on

December 9th, 2011

and approved by

Dr. Prashant Krishnamurthy, Associate Professor, Telecommunications, SIS

Dr. Daniel Mosse, Professor, Computer Science, School of Arts and Sciences

Dr. David Tipper, Associate Professor, Telecommunications, SIS

Dr. Martin Weiss, Professor, Telecommunications, SIS

Dissertation Advisor: Professor Taieb Znati, Telecommunications, SIS

Copyright © by Octavio Herrera-Ruiz

2012

# **A HOLISTIC REDUNDANCY– AND INCENTIVE–BASED FRAMEWORK TO IMPROVE CONTENT AVAILABILITY IN PEER-TO-PEER NETWORKS**

Octavio Herrera-Ruiz, PhD

University of Pittsburgh, 2012

Peer-to-Peer (P2P) technology has emerged as an important alternative to the traditional client-server communication paradigm to build large-scale distributed systems. P2P enables the creation, dissemination and access to information at low cost and without the need of dedicated coordinating entities. However, existing P2P systems fail to provide high-levels of content availability, which limit their applicability and adoption. This dissertation takes a holistic approach to devise mechanisms to improve content availability in large-scale P2P systems.

Content availability in P2P can be impacted by hardware failures and churn. Hardware failures, in the form of disk or node failures, render information inaccessible. Churn, an inherent property of P2P, is the collective effect of the users' uncoordinated behavior, which occurs when a large percentage of nodes join and leave frequently. Such a behavior reduces content availability significantly. Mitigating the combined effect of hardware failures and churn on content availability in P2P requires new and innovative solutions that go beyond those applied in existing distributed systems. To address this challenge, the thesis proposes two complementary, low cost mechanisms, whereby nodes self-organize to overcome failures and improve content availability. The first mechanism is a low complexity and highly flexible hybrid redundancy scheme, referred to as Proactive Repair (PR). The second mechanism is an incentive-based scheme that promotes cooperation and enforces fair exchange of resources among peers. These mechanisms provide the basis for the development of distributed self-



organizing algorithms to automate PR and, through incentives, maximize their effectiveness in realistic P2P environments.

Our proposed solution is evaluated using a combination of analytical and experimental methods. The analytical models are developed to determine the availability and repair cost properties of PR. The results indicate that PR's repair cost outperforms other redundancy schemes. The experimental analysis was carried out using simulation and the development of a testbed. The simulation results confirm that PR improves content availability in P2P. The proposed mechanisms are implemented and tested using a DHT-based P2P application environment. The experimental results indicate that the incentive-based mechanism can promote fair exchange of resources and limits the impact of uncooperative behaviors such as "*free-riding*".

## TABLE OF CONTENTS

ACKNOWLEDGEMENTS .....	XVI
1.0 INTRODUCTION .....	1
1.1 CONTENT AVAILABILITY .....	2
1.1.1 Physical errors.....	6
1.1.2 Logical errors .....	7
1.2 THESIS STATEMENT .....	8
1.3 RESEARCH OVERVIEW AND CONTRIBUTIONS .....	8
1.4 DISSERTATION ORGANIZATION .....	12
2.0 BACKGROUND AND LITERATURE REVIEW .....	14
2.1 P2P TECHNOLOGY .....	14
2.1.1 P2P Taxonomies and Features .....	15
2.1.1.1 Unstructured P2P Overlays.....	17
2.1.1.2 Structured P2P Overlays .....	18
2.1.2 P2P Architectures Comparison .....	21
2.1.3 Content Availability.....	23
2.1.3.1 Effect of Churn .....	23
2.1.3.2 Effect of User Behavior .....	25
2.1.4 Content Ownership in P2P Networks .....	26
2.2 BAMBOO .....	28
2.3 LITERATURE REVIEW .....	33
2.3.1 Churn Models.....	33
2.3.2 Redundancy Methods.....	35
2.3.3 Incentives and Content Availability.....	38

3.0 CONTENT AVAILABILITY FRAMEWORK .....	41
3.1 FRAMEWORK STRUCTURE .....	41
3.2 FRAMEWORK MODELS .....	43
3.2.1 User Behavior Component.....	43
3.2.2 Node Availability Component.....	44
3.2.3 Network Structure Component .....	46
3.2.4 Content Component .....	46
3.2.5 Incentives and Redundancy Component.....	48
3.2.6 Response Variables.....	49
3.3 FRAMEWORK PARAMETERS .....	51
3.3.1 User Behavior Component.....	52
3.3.2 Node Availability Component .....	52
3.3.2.1 Session Lengths .....	53
3.3.2.2 Off-line Interval .....	55
3.3.2.3 Oscillating Network Size .....	56
3.3.3 Network Component.....	57
3.3.3.1 Modelnet .....	58
3.3.4 Content Component .....	60
3.3.5 Incentives and Redundancy Component.....	62
3.3.5.1 Redundancy.....	62
3.3.5.2 Incentives-Based Mechanism .....	63
4.0 REDUNDANCY.....	64
4.1 REDUNDANCY SCHEMES .....	66
4.2 FILE AVAILABILITY .....	67
4.2.1 Heterogeneous node availabilities .....	68
4.2.2 Homogeneous node availabilities .....	73
4.2.2.1 Storage Overhead.....	74
4.2.2.2 Fragment Availability Model.....	76
4.3 REDUNDANCY REPAIR .....	79
4.3.1 Failed Repairs .....	83
4.3.2 Repair Cost.....	84

4.3.2.1 Proactive Replication with Different ( $S, r$ ) Parameters .....	89
4.3.3 Maintenance Epochs .....	91
4.4 REDUNDANCY MAINTENANCE .....	95
4.4.1 System Architecture .....	98
4.4.2 Agents .....	100
4.4.3 System Processes .....	101
5.0 INCENTIVES .....	110
5.1 ECONOMIC MODELS .....	110
5.2 INCENTIVE-BASED MECHANISMS FOR CONTENT AVAILABILITY .....	112
5.2.1 Objective .....	112
5.2.2 Design Guidelines .....	113
5.2.3 Agents .....	114
5.2.4 Control Metric .....	115
5.2.5 Incentive and Penalty Functions .....	116
5.2.5.1 Replication Probability Function .....	119
5.2.5.2 Transmission Bandwidth Function .....	120
5.3 INCENTIVE-BASED MECHANISM CONSTRUCTION .....	122
5.3.1 Redundancy-Maintenance Finite State Machine .....	122
5.3.1.1 Requestor State .....	124
5.3.1.2 Holder-Publisher State .....	124
5.3.1.3 Holder-Provider State .....	125
5.3.1.4 Index State .....	127
5.3.1.5 Candidate State .....	128
5.4 ENHANCED CONTRIBUTION METRIC .....	129
5.5 ADDITIONAL CONSIDERATIONS .....	133
6.0 EVALUATION .....	136
6.1 SIMULATION-EMULATION PLATFORM .....	136
6.2 BASE CONFIGURATION .....	137
6.3 PERFORMANCE METRICS .....	139
6.4 REDUNDANCY EFFECT .....	142
6.5 PARAMETERS EXPLORATION .....	144

6.5.1 Proactive Replication Redundancy Scheme .....	145
6.5.1.1 PR Redundancy with Different Coding Gains ( $S$ ) .....	148
6.5.1.2 PR Redundancy with Different Reception Efficiencies ( $k$ ) .....	149
6.5.2 Redundancy Maintenance .....	149
6.5.2.1 Adaptive Maintenance Epochs .....	150
6.5.2.2 TARGET Fragment Availability .....	151
6.5.2.3 Minimum Number of Segments Available (MIN_SEG) .....	153
6.5.3 Incentive-Based Mechanism .....	154
6.5.3.1 Sigmoid Functions Parameters .....	155
6.5.3.2 Contribution Gain function .....	160
6.6 REDUNDANCY-SYSTEM FLEXIBILITY .....	164
6.6.1 Churn Rate .....	164
6.6.2 Nodes Participation .....	165
6.6.3 Benefactor Nodes .....	170
6.6.4 Content Space Size .....	171
7.0 CONCLUSIONS AND FUTURE WORK .....	174
APPENDIX A .....	178
BIBLIOGRAPHY .....	181

## LIST OF TABLES

Table 1. Peer-to-Peer Systems Properties .....	16
Table 2. Comparison of P2P Architectures.....	22
Table 3. Response Variables .....	50
Table 4. Content Availability Framework Parameters .....	51
Table 5. User Behavior Parameters .....	52
Table 6. Node Availability Component Parameters .....	53
Table 7. Churn Measurement Studies .....	54
Table 8. Peers User Class Sample Session Distributions .....	55
Table 9. Bamboo Parameters .....	58
Table 10. File Type/Size Distribution.....	61
Table 11. Content Parameters .....	62
Table 12. Redundancy Mechanism Parameters .....	63
Table 13. Incentive-Based Mechanism Parameters .....	63
Table 14. Redundancy Notation: Redundancy Scheme(s) Parameters.....	65
Table 15. Redundancy Notation: File Availability .....	65
Table 16. Redundancy Notation: Repair Cost .....	66
Table 17. Repair Cost Calculation Parameters .....	85
Table 18. Redundancy Methods with Flexible vs Fix Coding Gain.....	88
Table 19. IndexEntry Data Structure .....	104
Table 20. Content Contribution Metric.....	116
Table 21. Sigmoid Function Parameters .....	117
Table 22. Replication Probability Function Parameters .....	119
Table 23. Transmission Bandwidth Function Parameters .....	121
Table 24. IndexTable Data Structure .....	128
Table 25. Base Content/Networking Configuration .....	138

Table 26. Response Variables for Alternative Redundancy Schemes.....	147
Table 27. Size and Popularity Components in Contribution Gain Functions.....	160

## LIST OF FIGURES

Figure 1. Content Availability System States and Transitions Diagram .....	5
Figure 2. Content Availability System Proposed in this Dissertation .....	10
Figure 3. Dissertation Organization.....	12
Figure 4. Unstructured P2P Overlays Taxonomy .....	17
Figure 5. Structured P2P Overlays Taxonomy .....	19
Figure 6. Overlays Network Design Decisions.....	20
Figure 7. Churn's Effect on Content Availability.....	24
Figure 8. Neighbors in Bamboo.....	29
Figure 9. Distributed Hash Table.....	30
Figure 10. Bamboo's Performance Under Churn .....	32
Figure 11. Content Availability Framework.....	42
Figure 12. User Behavior Model: User Profiles .....	43
Figure 13. Node Churn Model.....	45
Figure 14. Content Component Factors .....	47
Figure 15. Overlay Network Size vs Simulation Time.....	57
Figure 16. Modelnet Network Topology .....	59
Figure 17. Redundancy Data Structure.....	68
Figure 18. Algorithm 1 Sample Output .....	73
Figure 19. Minimum Storage Overhead .....	75
Figure 20. Redundancy Repair Cost.....	85
Figure 21. Repair Cost for Flexible and Fixed Coding Gain.....	87
Figure 22. Repair Cost for PR with different <i>Coding Gain</i> ( $S$ ) values.....	90
Figure 23. Repair Cost for PR with different <i>Replication Gain</i> ( $r$ ) values.....	90
Figure 24. PR Cost vs MDS and Ideal MDS Cost.....	94



Figure 25. System Architecture for Redundancy Maintenance .....	98
Figure 26. System Processes Overview .....	102
Figure 27. Effect of Incentive-based Mechanism: a) Normal System, b) System with Incentives .....	113
Figure 28. Sigmoid Function .....	118
Figure 29. Replication Probability Function.....	120
Figure 30. Transmission Bandwidth Function.....	122
Figure 31. Redundancy-maintenance process FSM.....	123
Figure 32. Time Contribution Gain .....	131
Figure 33. Sample Simulation's Cost Results Output .....	140
Figure 34. Sample Simulation's Efficiency Results Output .....	141
Figure 35. Effect of Redundancy on Content Availability .....	144
Figure 36. Cost and Efficiency for Alternative PR Redundancy Settings.....	146
Figure 37. Adaptive Maintenance Epochs and Publication Rate: a) $\text{Adapt} > 0$ , b) $\text{Adapt} \geq 0$ , c) Smooth with $\alpha = 0.4$ and d) Smooth with $\alpha = 0.8$ .....	151
Figure 38. Cost and Efficiency for Different TARGET Values .....	152
Figure 39. Cost and Efficiency for Different MIN_SEG Values.....	154
Figure 40. Cost and Efficiency for Different Incentive-based Mechanism Parameters .....	156
Figure 41. Transmission Bandwidth CDF for Various TB-Util Shape Parameter Values .....	157
Figure 42. Transmission Bandwidth for different TB-Util Shape Parameter Values (Scatter Diagram) .....	158
Figure 43. Fair Allocation of Resources in Incentive-based Mechanism.....	159
Figure 44. Cost and Efficiency for Contribution Gain Functions.....	161
Figure 45. Download's Transmission Bandwidth CDF for Contribution Gain Functions .....	162
Figure 46. RP-Cost Function for Contribution Gain Functions.....	163
Figure 47. Cost and Efficiency for Various Churn Rates .....	165
Figure 48. Cost and Efficiency for Different Percentage of Non-Participants .....	167
Figure 49. Distribution of Transmission Bandwidth versus Content Contribution for: a) compliant nodes b) non-compliant nodes. ....	169
Figure 50. Cost and Efficiency for Different Percentage of <i>Benefactors</i> .....	170
Figure 51. Cost and Efficiency for Different Content Space Sizes .....	172

Figure 52. Content Availability for Different Content Space Sizes .....	173
Figure 53. Alternative Node Initializations for Transient Removal .....	180

## LIST OF ALGORITHMS

Algorithm 1. Optimization of Redundancy Resources .....	72
Algorithm 2. Items Registration Process at Holder Agent .....	103
Algorithm 3. Items Registration Process at Index Agent .....	105
Algorithm 4. Redundancy Evaluation by Index Agent.....	106
Algorithm 5. Smooth Maintenance Epoch Mechanism at Index Agent .....	107
Algorithm 6. Adaptive Maintenance Epoch Mechanism at Index Agent.....	108
Algorithm 7. RedundancyFix at Index Agent.....	109
Algorithm 8. RedundancyFix at Holder Agent.....	109
Algorithm 9. Function getSpeed() .....	126
Algorithm 10. Function Utility().....	127
Algorithm 11. Procedure AcceptReplica() .....	128
Algorithm 12. Function Cost().....	129
Algorithm 13. Procedure PerfectSimulation() .....	179

## ACKNOWLEDGEMENTS

“Utopia lies at the horizon.  
When I draw nearer by two steps,  
it retreats two steps.  
If I proceed ten steps forward, it  
swiftly slips ten steps ahead.  
No matter how far I go, I can never reach it.  
What, then, is the purpose of utopia?  
It is to cause us to advance.”

— Eduardo Hughes Galeano  
*Author of Open Veins of Latin America*

It has been a very long road to reach this milestone. Too many times it seemed even utopic that I would make it, but here I am. At this point, I believe it is important to stop for a few minutes and reflect on why and how I got here. The why has always been easy, I want to use my field of expertise (i.e., telecommunications) to improve the livelihood conditions in my home country, México. The how is more complicated, but I believe that at least I have to acknowledge and thank all the people who has been fundamental in the completion of this milestone. They have provide me with guidance, support, inspiration and love throughout (and before) this process. I will dedicate the rest of this section solely to this.

I want to thank my advisor Dr. Taieb Znati for his guidance in the development of this research and his continuous encouragement to improve it. In addition, I want to thank my

dissertation committee for their input on my research work and their teachings during my PhD and Master program.

I am thankful for “my team” of editors that kindly donated their time to help me improve my written work. Sandy, Bedda, Emilio, Ruth, Julie and Daniel have certainly played a fundamental role in improving the quality of this document.

I am grateful for all the organizations and individuals that provided me with financial support during my graduate studies. First of all, my wife, who kept me afloat during the last two terms of this endeavor. My alma mater, the National Autonomous University of México (UNAM) and the former Central Academic Computing Division (DGSCA) now DGTIC where I grew professionally, academically and personally. I will always cherish the years I spent as an intern, staff, faculty and manager. The CONACYT and the Fulbright-Garcia Robles scholarship program for their financial support. The graduate program in telecommunications for granting me a full tuition waiver during my master studies. The European Union Center of Excellence and European Studies Center (EUCE/ESC) at the University of Pittsburgh, for the Graduate Student Assistant appointment (GSA) that supported me during most of my PhD studies. Special thanks to Dr. Alberta Sbragia, whose leadership defined the GSA policies of the EUCE/ESC and who always trusted in my technical expertise. Jorge Reyes, who trusted me to work with him despite my limited handyman experience. From him, I learned not only to be proud of what you do, but also that even the messiest project can be done orderly and cleanly. Bosch’s Research and Technology Center, where I was an intern. The projects I developed there were professionally fulfilling. The Vira I. Heinz (VIH) Program for Women in Global Leadership for granting me the opportunity to work freelance, editing their website.

I also want to thank the extraordinary group of people that have provided me with emotional support one way or another throughout this endeavor: my wife Sandy (obviously), my family at large (incluyendo a todos los enanos), my closest friends Bedda and Emilio (always generous and supportive) and my friends of the squash federation (lead by the brave David Brumble).

I definitively want to thank those people who inspire me and provided me with role models. My parents first of all. Next, my brother Gilberto, who bought for me my first computer and was the first in the family to get a PhD; then he became the Dean of Engineering and now the future Chancellor at his university. He certainly has been a great career model (even though matching his career path appears utopic). My sister (Dra. nena) who has managed to sprung a successful academic career together with two bright and beautiful *chilpayates*. Also, professor Alberta Sbragia, who managed to balance her teaching, research and leadership position at the EUCE/ESC. Finally, Professor Daniel Mosse, for his exemplary academic attitude, always willing to help students, thought provoking comments and his upbeat personality.

## 1.0 INTRODUCTION

The Peer-to-Peer (P2P) networking paradigm comprises several alternative distributed architectures to build large-scale virtual networks (e.g., unstructured and DHT-based P2P networks<sup>1</sup>) on top of an existing network infrastructure, generally the Internet. The primary objective of these networks is to enable the sharing of services and content among participants. The resources and tasks needed for end-to-end communication are contributed and performed (with several possible levels of decentralization) by the participating nodes. Generally, P2P networks are open access self-organizing systems formed by volunteer nodes; nonetheless, individual nodes are autonomous and self-interested in nature.

P2P technology has emerged as a viable solution for the deployment of large-scale distributed applications such as content distribution and file sharing. Currently, deployed P2P systems reach several million users, with several hundred thousand peers<sup>2</sup> connected simultaneously [1]. Recent measurement studies indicate that P2P applications still make up the majority of the Internet traffic worldwide [2] and more than half of all upstream traffic in the USA is still attributed to P2P applications [3]<sup>3</sup>.

---

<sup>1</sup> We present a brief description of these two architectures in Chapter 2.

<sup>2</sup> The terms user, node and peer are used interchangeably throughout this work in reference to a uniquely identifiable autonomous entity (e.g., a single computer) participating in the P2P network.

<sup>3</sup> However, real-time entertainment traffic constitutes the majority of the total traffic nowadays.

The sheer size and open access nature of the P2P application environment generate important challenges for the design of reliable and efficient P2P networks. Intermittent node connectivity, also known as churn, is one of them. Not only because nodes' joins and leaves increase the network's overhead, but also because reliable access to resources (e.g., files, in a file-sharing P2P network) can be limited by this intermittent availability pattern. Another key challenge is how to minimize the incidence of non-cooperative behaviors, which can lead to unfairness, performance degradation and limited scalability. A well known example for this class of behavior is *free-riding*, which has been documented as a major problem in many deployed P2P networks [4] [5].

In this dissertation, we undertake the specific problem of content availability in large-scale P2P networks. The following section describes the challenges mentioned above in the context of content availability in P2P networks and outlines our proposed solution. We also present our thesis statement and describe how this dissertation proves it. For an overview of P2P technology and its operation we refer the reader to section 2.1.

## 1.1 CONTENT AVAILABILITY

Content availability refers to the ability of the system to make content (e.g., files) readily accessible to users. Presently, deployed P2P systems exhibit important content availability performance issues. For example, in BitTorrent 40% of the swarms<sup>4</sup> lack publishers (i.e., seeds<sup>5</sup>)

---

<sup>4</sup> In BitTorrent, a swarm is a set of nodes that cooperatively download a single file by exchanging portions of it among them.

<sup>5</sup> In BitTorrent, a seed is a member of the swarm (i.e., node) that has a complete copy of the file being downloaded.



during more than 50% of their lifetime [6], and in 86% of the cases the participants (i.e., leechers<sup>6</sup>) are unable to reconstruct the original data [7].

P2P has characteristic features that demand different content availability solutions to those applied in other distributed systems. Large-scale, open access, heterogeneous resources and churn are just a few of these features. The case of redundancy is a perfect example. Traditionally, erasure coding is used in distributed systems to prevent data loss due to failures. However, it is generally assumed that performing repairs after such failures is not of major scalability or cost concern. In P2P systems, redundancy still is a plausible solution to protect the system against failures, but repairs have major cost and scalability implications. First, due to the rate of content errors generated in the system. When nodes leave the network or when hard-disks fail, the network needs to be reorganized to compensate for the content lost. Second, because access bandwidth is limited. Nodes commit their access bandwidth to repairs without impacting their performance. Consequently, redundancy for P2P networks must address not only availability as key design principle, but repair bandwidth cost as well. In addition, distributed self-organizing processes and data structures to manage the redundancy are key architectural requirements to device reliable and scalable access to content in P2P networks. In particular, we argue that a hybrid redundancy scheme, combining traditional MDS (Maximum Distance Separable) erasure coding and replication redundancy, can provide content availability reliably and efficiently. Our proposed redundancy scheme, named Proactive Repair (and described in detail in Chapter 4) uses erasure coding to achieve a good storage-availability tradeoff. In addition, our scheme generates replicas for each fragment proactively, so that when a single fragment is lost, it can be repaired using minimum bandwidth.

---

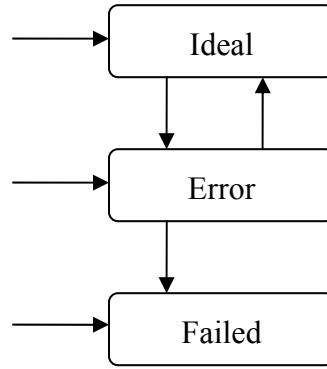
<sup>6</sup> In BitTorrent, leechers are members of the swarm that have not completed the transfer of the file.

Our proposed solution also includes self-organizing distributed algorithms to monitor the state information of our proposed redundancy scheme. We name redundancy maintenance process to this component. Building the algorithms for the redundancy maintenance process is easy using a Distributed Hash Table (DHT) P2P architecture (described with more detail in Chapter 2), which we assume in our solution. The basic lookup functionality in a DHT can be used to define a single point of contact (i.e., the index node) where nodes can gather information for each file. Nodes storing fragments would announce their availability at this point and nodes requesting a file could retrieve a list of targets from the same location. In addition, the node responsible for this location could evaluate the file availability and determine if a repair is needed. In which case, it would instruct the proper nodes to start a repair.

The correct operation of the components above is conditioned upon the full participation of nodes storing and repairing redundancy. However, to assume full cooperation limits the applicability of any solution. Thus, we embed an incentives-based mechanism into the operation of the system to foster cooperation and regulate fair exchange of resources among nodes, with the purpose of maximizing the feasibility of our solution. Our proposed incentive is a bartering mechanism for content availability (i.e., storage) versus performance (i.e., transmission bandwidth). Nodes receive a level of service proportional to their contribution towards the content availability of the system.

To characterize content availability in P2P networks, we present a multi-level failure model, similar to the one introduced in [8]. The purpose of this model is to describe the possible states of the system, on a per file basis (i.e. file availability), and the type of events that can cause a transition between these states.

The top level represents the ideal state. Whenever a piece of this information is lost, the system enters into the error state. We say that a content error (or simply an error) has occurred. If content errors accumulate without taking any corrective measures, the system will reach the failed state; at which point, the original data can not be retrieved from the system. Thus, in simple terms, the purpose of our research is to develop data structures and algorithms that would allow us to avoid the failed state efficiently.



**Figure 1. Content Availability System States and Transitions Diagram**

In a P2P system, multiple events generate content errors. Some of these errors originate at the network and application layers while others occur at the lower layers. To the best of our knowledge, there is no initiative consolidating content failure models for different layers into a comprehensive content availability model for P2P. In that regard, our research presents a multifaceted content availability framework integrating the effects of content errors at different layers of the application stack and presents a robust and efficient set of mechanisms to alleviate such failures.

There is extensive research on how to recover from content errors at lower layers (e.g., disk failure and nonrecoverable read errors), but these initiatives ignore other sources of content errors prevalent in P2P networks. For instance, in the work presented by Weatherspoon and Kubiatowicz [9] content availability is assumed to be dominated solely by disk failure rates. The research on

how to handle content failures at the network and application layers is also extensive, but weak content availability requirements are usually analyzed [10] and content errors at the lower layers are completely ignored. In contrast, our research goal is to provide strong content availability guarantees while considering a diverse set of content errors. We categorize all lower layer errors as physical errors, and all upper layer errors as logical errors. Next, we briefly describe this taxonomy.

### **1.1.1 Physical errors**

One of the key features of P2P systems is their support for large scale deployments. The sheer number of nodes embedded into the system implies that storage resources are typically constituted by a very large pool of independent storage devices. As a result, despite the low probability associated with nonrecoverable disk read errors and disk failures, the large number of components embedded into the system translate this probabilities into an error rate that can impact the reliability of content [11]. Furthermore, the distributed nature of P2P system prevents the use of data protection mechanisms typically used in centralized storage facilities (such as RAID). Thus, the effect of disk failures (and other hardware related errors) is likely to be higher than in other computing architectures. In general, addressing this type of failures can be done easily with traditional reliability measures. For example, given a set of hardware related failure rates –which are usually modeled as i.i.d. events– and a desired availability level, we can achieve a required reliability figure by defining a minimum level of data redundancy to be hardcoded into the system.

### 1.1.2 Logical errors

In addition to the physical errors inherent in any computing system, P2P networks have many more sources of errors. Thus, reliability in P2P networks is more complex and dynamic than in some other distributed systems. For instance, PlanetLab's nodes [12] have server-like availabilities; thus, the redundancy requirements to achieve a desired level of content availability are much lower than in a traditional P2P system [10]. Furthermore, reliability in a P2P system can not be engineered using the same static reliability metrics used in some other distributed systems. We categorize as logical errors a wide variety of content error events such as routing inconsistencies (when these prevent access to content), nodes departures (i.e., churn) and even content-related user behaviors (e.g., users canceling data transfers or free riding<sup>7</sup>).

Addressing the content availability effects of logical errors in P2P networks is the most challenging task of our research due to the number, complexity and dynamism of the factors involved. Furthermore, the temporal, quantitative and qualitative contributions of individual nodes towards the overall content availability of the network are typically highly heterogeneous. Some users contribute large amounts of resources consistently during long periods of time, while others may contribute only a few or no resources at all to improve content availability. As a result, any content availability mechanism to be used should use a holistic approach to manage this diversity in behaviors. Mechanisms to manage the intermittent connectivity nature of peers should be complemented with mechanisms to compensate the heterogeneity in user-behaviors prevalent in P2P networks.

---

<sup>7</sup> In the context of this dissertation, free riders are nodes consuming resources without sharing a fair portion of their resources with other nodes in the overlay.

## 1.2 THESIS STATEMENT

The thesis of this dissertation is that *the combination of erasure coding and replication redundancy constitutes a flexible and cost effective way to provide reliable access to content in P2P overlay networks. Moreover, this hybrid redundancy scheme can be integrated seamlessly with incentive-based mechanisms to support fair allocation of resources and to improve the scalability of the system while recognizing the participants' autonomy and diversity.*

To validate the thesis statement above, we take an approach that combines analytical and experimental methods. The analytical portion shows that our proposed hybrid redundancy mechanism outperforms other redundancy schemes in terms of its repair bandwidth requirements. We also develop distributed self-organizing algorithms to automate the redundancy repair for our proposed scheme. In addition, we augment these algorithms with an incentive-based mechanism to promote cooperation and achieve fair exchange of resources among nodes. Together, these mechanisms constitute a holistic framework to improve content availability in P2P networks. The experimental work presented demonstrates the feasibility of our framework. In addition, it realistically incorporates the complex nature of a dynamic P2P routing architecture into our validation process. Our experimental work is based on a deployed DHT-based P2P application stack, named Bamboo [13].

## 1.3 RESEARCH OVERVIEW AND CONTRIBUTIONS

Redundancy and economic models for P2P networks are addressed extensively in the research literature [14], [15], [16], [17], [18]. However, to the best of our knowledge, the analysis and

integration of these components has not been proposed earlier in the literature to address the diverse set of content failures present in P2P networks. To do so, this work pursues two research thrusts:

- i) Creation of a robust content availability model capturing the heterogeneous and dynamic nature of P2P networks
- ii) Development of a multifaceted redundancy-maintenance mechanism to improve content availability in the presence of diverse sources of errors

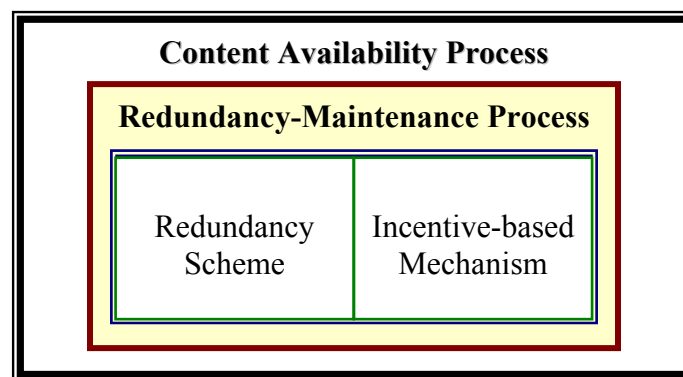
The aim of the first thrust is to capture the effects of two types of content errors into our analysis: physical errors and logical errors. In addition, we want to capture the diverse habits, capabilities, needs, and interests of individual P2P participants. From this content availability model, it will be possible to portray how these heterogeneous capabilities and attitudes are reflected by the content error events in the network.

The goal of the second research thrust is to create and integrate scalable redundancy and incentive-based mechanisms capable of improving content availability in P2P networks.

To study content availability in P2P networks, we develop analytical models for file availability and redundancy maintenance. We use these models to demonstrate the superior performance of our proposed redundancy scheme against erasure coding and network coding redundancy. In addition, we present a content availability evaluation framework that defines the factors and levels to be analyzed experimentally.

This dissertation proposes the use of redundancy to compensate the negative effects of dynamic node membership and other content error events, together with a utility/cost economic model to achieve fair allocation of resources and self-organization. The redundancy-maintenance method proposed has low maintenance bandwidth, with minimal complexity and superior flexibility properties. In the incentive-based mechanism we propose, we assume that peers transfer

the items they possess upon request. Thus, peers' content availability contribution can be expressed simply as a function of the number of items hosted by a node. Nonetheless, we explore the use of additional qualifiers to improve the fairness of the mechanism, such as size, popularity and the amount of time each item has been shared. Nodes receive an incentive in the form of performance (i.e., download bandwidth, as in BitTorrent) in exchange for hosting entire or fragments of files to improve the content availability of the network. In addition, the incentives-based mechanism presented provides a mechanism to achieve fair allocation of resources (i.e., each node receives a level of service proportional to its content availability contribution) while preserving node autonomy. The functional components for our system are illustrated in Figure 2.



**Figure 2. Content Availability System Proposed in this Dissertation**

The relevant features of the proposed solution are simplicity, scalability, self-organization and flexibility. Simplicity is achieved by engineering a proactive redundancy-maintenance process that does not require complex data structures or elaborate scheduling algorithms. Scalability is achieved in two ways: by using less maintenance bandwidth than other redundancy schemes, and by implementing the redundancy-maintenance process in a completely distributed fashion. In addition, the redundancy-maintenance mechanism adapts its operation automatically in order to maintain its performance when the network conditions change through time.



The proposed solution is implemented as a DHT-based content availability system. The underlying routing architecture employed is named Bamboo [19] and it has been used earlier to test other distributed applications. Our system is tested using an emulated<sup>8</sup> wide area networking environment, named Modelnet [20], which captures the delay and bandwidth restrictions of P2P networks deployed over the Internet. This system architecture (Modelnet + Bamboo + Redundancy Maintenance Process) is used to conduct extensive simulations to demonstrate the content availability features of the proposed solution.

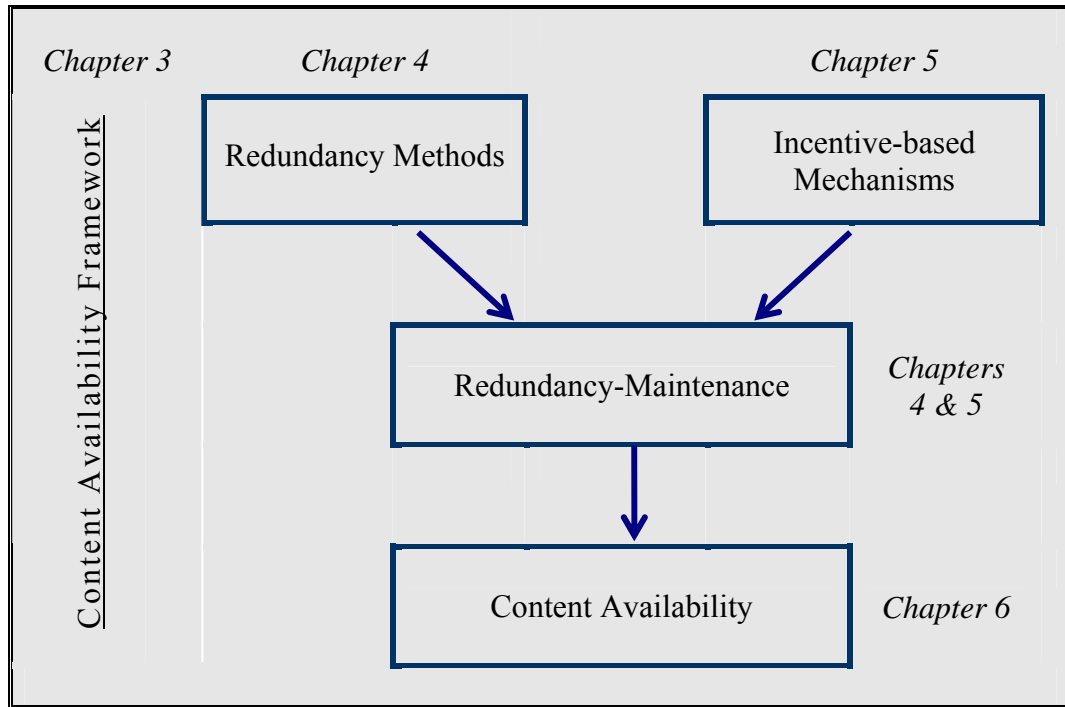
In summary, the contributions of this dissertation are:

- The creation of a broad content availability framework (Chapter 3.0)
- The definition of a new low-complexity and highly-flexible redundancy scheme that requires small amounts of repair bandwidth (Sections 4.1 and 4.3.2)
- The formulation of an analytical model to determine fragment availability, for code-based redundancy schemes, in P2P networks (Section 4.2.2.2)
- The formulation of analytical models to assess the redundancy repair cost for different code-based redundancy schemes (Section 4.3)
- The formulation of a redundancy-maintenance process (Section 4.4) with incentive-based mechanisms (Section 5.3) that:
  - Uses a proactive redundancy maintenance methodology that consumes the same or less maintenance bandwidth than other methods (Section 4.3.2)
  - Improves the content availability of items independently of their popularity
  - Is completely distributed (Section 4.4.1) and self-organizing (Section 5.2.2)
  - Adapts automatically to the dynamic conditions in the network (Section 4.4.3)
- A prototype implementation of the redundancy-maintenance process with incentive-based mechanisms as proof of concept (Chapter 6.0)

---

<sup>8</sup> However, we use the term simulation in our experimental work in reference to the entire simulated P2P networking environment..

Figure 3 illustrates how the topics presented in each chapter are integrated to create a content availability system for P2P networks under churn.



**Figure 3. Dissertation Organization**

## 1.4 DISSERTATION ORGANIZATION

This dissertation is organized in seven chapters.

Chapter 2 presents an overview of P2P and the research literature related with our research. As part of the background section on P2P, we include a description of the operation of the DHT routing substrate employed in the implementation of the proposed content availability system. In the literature review section, we describe research initiatives organized in four topics: churn models, redundancy methods, incentives mechanism and content availability.

Chapter 3 presents a framework describing the major factors affecting content availability in a P2P system. In addition, this section presents the levels used for each of these factors in the experimental portion of this work.

Chapter 4 introduces different redundancy schemes that can be used to improve content availability in P2P networks. In this chapter, we assess the effectiveness of our proposed redundancy scheme and compare its performance with other code-based redundancy methods.

Chapter 5 introduces economic models for P2P networks and describes the implementation of an incentive-based mechanism to promote the participation of nodes in a redundancy-maintenance process.

Chapter 6 presents an evaluation of our prototype implementation of a content-availability system that uses our proposed redundancy scheme in conjunction with incentive-based mechanisms to improve content availability.

Finally, Chapter 7 presents the conclusion of this study and possible directions for future work.

## **2.0 BACKGROUND AND LITERATURE REVIEW**

The purpose of this chapter is to present an overview of the features of P2P technology as well as a review of a variety of subjects related with content availability in P2P networks. The rest of this chapter is organized as follows. Section 2.1 presents a description of P2P technology. Section 2.1.4 describes Bamboo [21], which is the DHT routing architecture employed for the construction of the redundancy maintenance system presented in this dissertation. Section 2.1.4 describes the different content ownership modes in P2P networks. Finally, Section 2.3 presents an overview of the research literature related with our research.

### **2.1 P2P TECHNOLOGY**

Peer-to-Peer (P2P) computing or networking is a distributed application architecture where participants self-organize into virtual networks of nodes and logical links on top of an existing network infrastructure, generally the Internet. The participants are autonomous, self-interested nodes that share a portion of their resources to provide the services and content offered by the network.

P2P technology has emerged as a viable solution for the deployment of large-scale distributed applications such as content distribution (with BitTorrent as the most prominent example) and file sharing (with Napster first, and Gnutella, and many others later on). Currently,

deployed P2P systems reach several million users, with several hundred thousand peers<sup>9</sup> connected simultaneously at any time. The number of individuals, groups and organizations that have adopted the P2P computing paradigm is growing steadily, as does the amount of traffic exchanged. In recent years, content distribution sites exhibited tremendous growth; Mininova's<sup>10</sup> site for example, doubled its number of downloads in 2008 to 7 billion in a year [22]. In addition, recent measurement studies indicate that P2P applications make up the majority of the Internet traffic worldwide [2]. The emergence of real-time video traffic in the entertainment industry (Netflix in particular) has surpassed P2P downstream traffic during peak periods in North America [23], but more than half of all upstream traffic is still attributed to P2P applications [3]<sup>11</sup>.

### 2.1.1 P2P Taxonomies and Features

P2P computing covers a broad range of systems and applications that can be categorized using multiple taxonomies and characterizing features. For example, Milojevic et al., in [24] present taxonomies for the classification of P2P computing from a system, application, and market perspective; Keong et al., in [25] classify P2P networks as *unstructured* and *structured*, and employ a nine dimension taxonomy to compare them. More recently, Buford and Yu in [26] classify P2P networks in *unstructured* and *structured* categories, and also describe other classes of P2P overlays, such as hierarchical, service, semantic and sensor overlays. This section summarizes the properties of P2P systems and describes the fundamental characteristic of *unstructured* and

---

<sup>9</sup> The terms user, node and peer are used indistinctively throughout this work in reference to a uniquely identifiable autonomous entity (e.g., a single computer) participating in the overlay network.

<sup>10</sup> Mininova is an indexer site for the popular BitTorrent network.

<sup>11</sup> However, real-time entertainment traffic constitutes the majority of the total traffic nowadays.

structured overlays. For a comprehensive review of different classes of P2P overlays we refer the reader to [26].

**Table 1. Peer-to-Peer Systems Properties**

<b>Property</b>	<b>Description</b>
<b><i>Decentralization</i></b>	The responsibility for the operation of the system is equally distributed among participants. There is no central point of control. However, in many designs this property is relaxed and some nodes assume special roles.
<b><i>Resource Sharing</i></b>	The resources required to support the services and content offered by the network are provided by its participants. The nature and amount of resources depends on the application and the system's architecture. As a result of the decentralized and resource sharing properties, nodes in P2P overlay networks (i.e., peers) play dual roles as providers and consumer of resources.
<b><i>Scalability</i></b>	The minimum amount of resources that each node needs to commit to the overlay grows less than linear with respect to the network size. In addition, the performance metrics of the system, such as response time, does not degrade significantly with increased system size or load.
<b><i>Autonomy</i></b>	Participation in the overlay is voluntary. Every node decides unilaterally the extent of its participation, including the amount of resources committed and when to join and leave the system.
<b><i>Self-organization</i></b>	Nodes use local knowledge to make decisions that over time result in a better organization of the system. Nodes continuously adapt to the dynamic networking conditions in the overlay, assuming the responsibility for maintaining their own neighbor relationships despite the continuous arrival and departure of other nodes to and from the overlay.
<b><i>Fault resiliency</i></b>	The decentralized overlay architecture precludes the existence of single points of failure. This allows the system to continue its normal operation without significant performance degradation in the presence of node and communication errors. After a failure, the system is capable of reorganizing itself without the intervention of a centralized coordination entity (i.e., self-organization). This property also translates, at some degree, into a capability to avoid censorship and resist other types of attacks.
<b><i>Cost of ownership</i></b>	A premise of the P2P paradigm is that the aggregated value of the participants' resources is greater than the sum of the individual resources. Furthermore, given that these resources are contributed voluntarily, the deployment and/or maintenance cost of the system is amortized by its members. As a result, P2P technology constitutes a cost effective alternative for the deployment of network services, features and functionalities that would require modifications in the underlying network infrastructure.

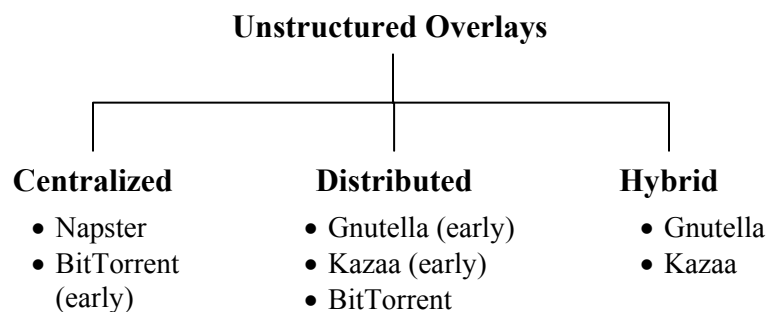
Table 1 describes the fundamental properties applicable to most P2P systems. Depending on specific implementation decisions, different systems favor some properties over others. For

example, Napster used a pure P2P data transfer principle, but its content indexing implementation was completely centralized, thus not P2P.

There are two major P2P overlay architectures, *structured* and *unstructured*. Most P2P systems can be classified into one of these architectures based on the way relationships between nodes are built and maintained.

### 2.1.1.1 Unstructured P2P Overlays

In *unstructured* P2P overlays, peers join the network without any prior knowledge of the topology and form a random graph in a flat or hierarchical manner that usually exhibits small world phenomena. Nodes rely solely on their adjacent nodes for delivery of messages to other nodes in the overlay. This type of P2P systems usually supports data location using complex queries, but generally there are not guarantees (i.e., data location is a probabilistic process). *Unstructured* overlays can be further classified according to the data location/distribution model used in *centralized*, *distributed* and *hybrid* systems. Figure 4 illustrates this taxonomy.



**Figure 4. Unstructured P2P Overlays Taxonomy**

Napster emerged in 1999 as a file-sharing application based on a centralized indexing mechanism and direct content transfer among participants. The popularity of P2P file-sharing applications rapidly transformed the distribution of Internet's traffic, making up more than half of

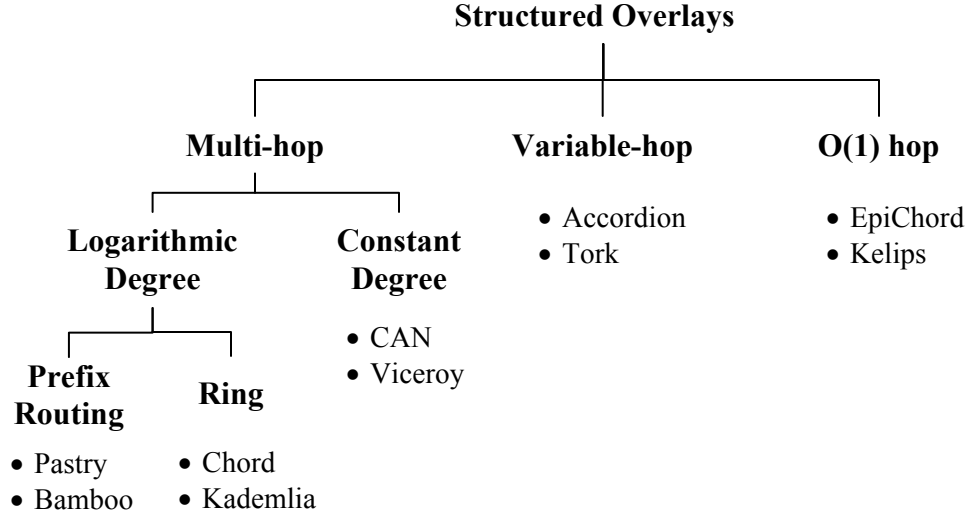
it as early as 2006 [26]. The content distribution application BitTorrent (before version 4.2.0) used a similar approach, but today its architecture includes alternative distributed mechanisms for coordinating data upload and download. The early implementations of file sharing applications Gnutella and Kazaa were completely distributed (i.e., nodes were organized in a random flat topology), but scalability problems with the message propagation technique (i.e., flooding) fostered the adoption of hierarchical architectures (i.e., hybrid) and alternative message forwarding mechanism, such as random walks, that have allowed these systems to scale successfully to hundreds of thousands of nodes.

#### **2.1.1.2 Structured P2P Overlays**

In *structured* P2P overlays, the network topology is formed according to specific criteria and algorithms to achieve robustness and improve performance. *Structured* overlays use a key-based virtual addressing space for node identification and for data placement or location. Nodes cooperate to maintain routing information about how to reach other members of the overlay and support a consistent exact-match data location functionality. Xuemin, et al., in [26] classify *structured* networks according to the number of hops needed to reach other nodes in multi-hop, variable-hop and  $O(1)$  hop.

Multi-hop *structured* overlays are further subdivided in logarithmic and constant degree systems. Logarithmic degree systems are subdivided in prefix-based routing and ring geometry. Figure 5 presents this taxonomy together with the name of two sample systems in each category.





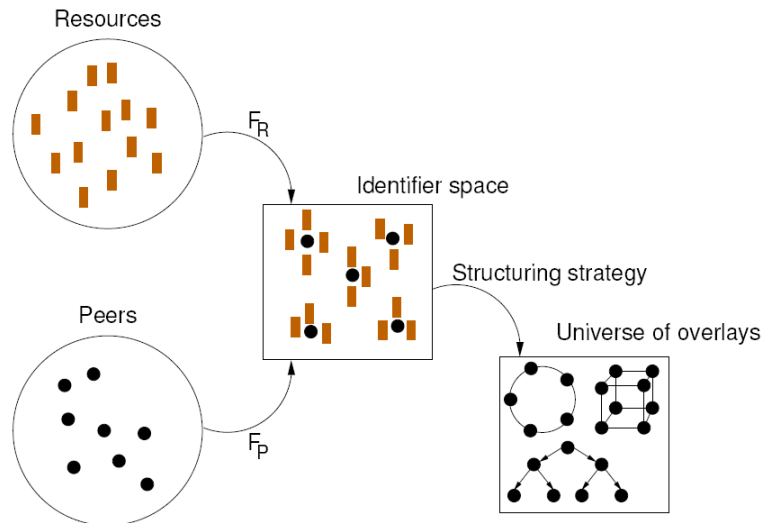
**Figure 5. Structured P2P Overlays Taxonomy**

Aberley et al., in [27] present a characterization framework for P2P overlays based on six key design aspects:

1. Identifier space
2. Mapping of resources and peers
3. Management of the identifier space
4. Graph embedding
5. Routing strategy
6. Maintenance strategy

This framework is presented in Figure 6 to illustrate the overall construction of *structured* P2P overlays. In *structured* systems, both nodes and resources are mapped into a common identifier space (functions  $F_P$  and  $F_R$  in the diagram). In the figure, nodes identifiers are represented by circles and resource identifiers are represented by rectangles. The terms `nodeId` and `key` are commonly used in reference to the identifiers of nodes and data elements respectively. Keys are associated to the node with a `nodeId` numerically closer to their own value. In addition, nodes organize themselves to form a geometry defined by the system architecture (e.g., ring, tree, etc.) and continuously maintain this topology to adapt to the dynamic networking conditions in the

overlay. The overlay's structure facilitates data location and provides performance bounds on the number of hops required to reach any node in the overlay, however, it does not support complex queries.



**Figure 6. Overlays Network Design Decisions<sup>12</sup>**

Most *structured* overlays use a Distributed Hash Table (DHT) to create the system's address space. This is accomplished by employing a globally known hash function to map nodes' IP socket and data objects' identifiers into a key-based address space (generally in the order of 128 bits long), and distributing the responsibility of managing a portion of this address space among the participating nodes.

Examples of structured overlays include mostly academic works, such as Bamboo [21] and Chord [28], but there are also deployed systems like eMule [29]. Most deployed *unstructured* overlays have adopted DHT-based indexing services to improve their data discovery mechanisms. In Section 2.2 we present a more detailed description of Bamboo's DHT implementation, which we use in the experimental portion of this work.

---

<sup>12</sup> Original figure taken from 27. Karl Aberery, et al., *The essence of P2P: A reference architecture for overlay networks*, in *Fifth International Conference on Peer-to-Peer Computing*. 2005: Konstanz, Germany.

### 2.1.2 P2P Architectures Comparison

During the last decade of research and development, the P2P application ecosystem grew enormously. Multitudes of P2P system architectures were introduced and in many cases, these architectures have evolved into multifaceted complex systems that are difficult to compare due to the diversity of features and functionalities involved. For example, presently there are at least twenty software implementations available for the Gnutella network with no less than a dozen optional features [30]. In addition, the Gnutella specification evolved from a flat topology with a broadcast-based content discovery mechanism into a hierarchical architecture with random walks and hash-based content discovery mechanisms. Thus, the reader should keep in mind that when comparing among different overlay architectures, the literature refers to the properties that characterize the fundamentals of each architecture, and not their current, state of the art implementations. In the paragraphs that follow, we compare several properties of structured overlay architectures versus unstructured overlay architectures and present a summary of our observations in Table 2.

*Structured* overlays are fully decentralized and most *unstructured* architectures are not. Decentralization is important in terms of fault tolerance and scalability, but presently, this difference does not seriously jeopardized the viability of *unstructured* overlays.

P2P architectures use diverse resource discovery paradigms. *Structured* systems provide an exact-match consistent data location functionality, which is analogous to their routing functionality. *Unstructured* systems on the other hand, use a probabilistic data location service (such as limited scope broadcast and random walks) that supports the usage of keywords and wildcards. In early *unstructured* overlays, successful data discovery was intimately related to the items' popularity, but presently many systems have incorporated hash-based indexing

functionalities in their architectures to provide consistent data location services (independent of the items' popularity).

**Table 2. Comparison of P2P Architectures**

<b>Features</b>	<b>Architecture</b>	
	<i>Structured</i>	<i>Unstructured</i>
<i>Sample Systems</i>	Bamboo & Chord	Kazaa & Gnutella
<i>Network topology</i>	Deterministic	Random (two layer) hierarchy
<i>Decentralization</i>	Full	Limited
<i>Scalability</i>	Yes	Yes
<i>Resource Discovery</i>	Distributed w/ exact-match (same as routing)	Limited-broadcast w/ keyword & wildcard support
<i>Performance guarantees</i>	Yes	Limited
<i>Fair allocation of resources</i>	n/a <sup>13</sup>	Yes <sup>14</sup>

Fair allocation of resources and cooperation among nodes is assumed by most P2P architectures. That is, nodes are expected to benefit from the system only as much or in proportion to what they contribute [18]. However, most deployed systems suffer from an uneven distribution of resources. To address this problem, systems incorporate additional mechanisms to promote cooperation and fair exchange of resources among peers. Currently, the inclusion of this type of mechanisms is an open issue and is not fundamentally limited by the architecture being used in the overlay.

---

<sup>13</sup> Bamboo and Chord do not enforce fair allocation of resources, but other deployed DHT-based systems do.

<sup>14</sup> Kazaa and Gnutella do not enforce fair allocation of resources, but BitTorrent does.

### **2.1.3 Content Availability**

Deployed P2P networks organize large amounts of resources across the Internet. However, the availability of these resources is hindered (among other factors) by the occurrence of *i*) intermittent node participation, which is an unavoidable feature of the open nature of this environment and *ii*) the autonomous operation of participants, expressed as an heterogeneous, and sometimes disproportionate, contribution and consumption of resources among peers.

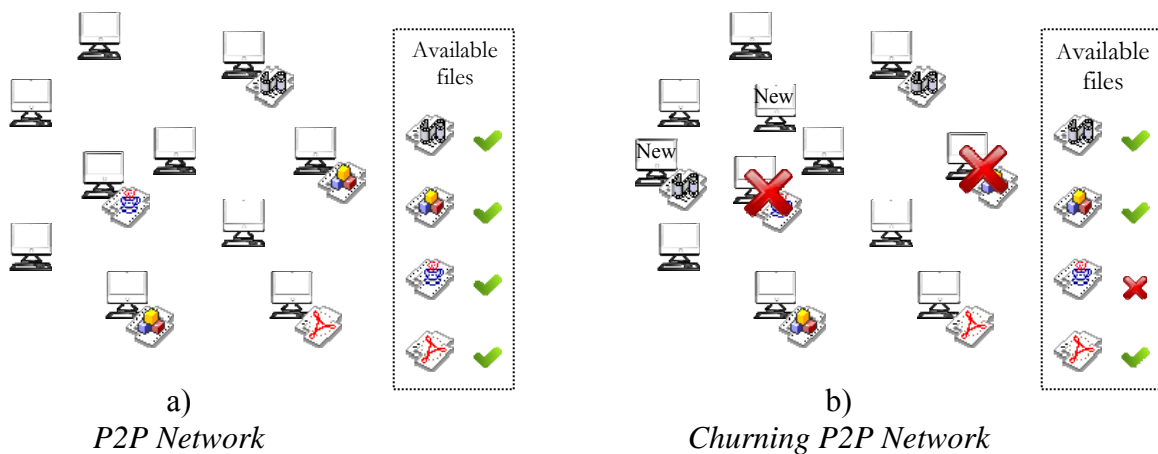
#### **2.1.3.1 Effect of Churn**

In most P2P networks a node's connectivity is transient. This phenomenon, called churn, is the main source of dynamism in the network. Understanding the impact of churn on content availability is the keystone to improve content availability in P2P overlay networks.

Participation in P2P networks is open, meaning that any node across the Internet can join the system. Churn is an unavoidable feature of the open nature of this application environment because users decide autonomously when and for how long to join the network, when to leave and whether to return. Thus, P2P communities are a dynamic conglomerate of heterogeneous hosts and resources where the participation of individual hosts is voluntary and transient.

Previous studies have shown that churn in P2P networks is prevalent across different applications and highly heterogeneous [31]. Nodes exhibit session lengths that can vary from a few minutes [32] to several hours or even days. The impact on content availability occurs when users leave the network, taking their content with them. In addition, users returning to the system get to decide unilaterally whether to share the content they previously obtained from the system. As a result, after a node departure, the network needs to reorganize the remaining peers to maintain reliable access to content.

Figure 7 presents a basic scenario that illustrates the effect of churn on content availability. In Figure 7.a every node can access the information stored in the network. There are five nodes sharing four unique data items. In Figure 7.b two of the nodes storing data have left the network and two new ones have joined. At this point, only three of the original data items are still available and the system has no means to recover the lost data.



**Figure 7. Churn's Effect on Content Availability**

Figure 7.b also illustrates that out of the two items lost; only one of them is now inaccessible. This is due to the extra copy stored at another node. This suggests the trivial solution to the problem: generate copies of all data items in every node, but this is certainly inefficient and not scalable.

Instead, a scalable solution to this problem should:

- i) Generate only enough redundancy in the system to handle the departure of nodes between maintenance epochs and
- ii) Reorganize the remaining nodes to regenerate any lost data before items become inaccessible due to additional churn events

While this solution is expressed in simplistic terms, its realization involves major complexity. For example, how do we determine how much is *enough* redundancy? And how do we reorganize the remaining nodes?

Redundancy is a necessary component to achieve content availability in P2P networks [10]. However, the selection of a redundancy method cannot be performed using the same cost-performance tradeoffs as in other distributed systems. Reliability needs to be continuously repaired. That is, any data loss due to churn (and other logical or hardware related failures) needs to be regenerated before the content availability degradation becomes irreversible. Performing this redundancy-maintenance efficiently has been described as the key limiting factor for the scalability of distributed storage P2P applications [10]. In particular, the amount of bandwidth needed to maintain the reliability of a redundancy method is a fundamental scalability concern in P2P environments.

#### **2.1.3.2 Effect of User Behavior**

Participants in P2P networks are usually expected to voluntarily share resources towards a common goal<sup>15</sup>. Nevertheless, peers are autonomous and decide the extent of their participation in the network unilaterally. In the absence of proper incentives, users acting rationally in their own self-interest do not commit enough resources towards the network's overall objective. Improving content availability in P2P networks requires mechanisms capable of dynamically managing not only transient and heterogeneous node connectivity, but also diverse individual tradeoffs between the peers' goals and the network's content availability objective.

---

<sup>15</sup> In that regard, P2P file-sharing networks can be modeled using economic models as an instance of private provisioning of a public good.

This dissertation uses economic models to address the user behavior aspects of the content availability in P2P networks. In these models, the term participant is equivalent to the terms users or node. Thus, the terms user, peer, node and participant are used indistinctively throughout this work in reference to the computers that form the overlay network.

Economic and networking processes have multiple similarities: complexity, autonomous self-interested participants and dynamic time-varying conditions, are just a few of them. Economic concepts and models have been used effectively in the study and construction of distributed systems, including P2P. In the economic incentives model in particular, the basic principle is that for an efficient and fair allocation of resources, there must be incentives for providers to share their resources as well as encouragement for consumers to maximize the utility of the received resources [33]. In P2P networks, nodes play a dual role as both consumer and providers of resources and consequently, the objective of incentive-based mechanisms is to allow nodes to reach a balance between their provider and consumer roles.

#### **2.1.4 Content Ownership in P2P Networks**

In structured P2P networks, nodes and content are both mapped deterministically to IDs in the same key-naming space. According to their IDs, each participating node becomes responsible for a section of the key-naming space (namely, Distributed Hash Table). The network supports three basic content operations: *query*, *indexing* and *storage*. In a DHT-based system, queries are functionally the same as routing. Therefore, a query operation consists in routing for a given key and finding the node responsible for that portion of the key-naming space. A node managing the section of the key-naming space including the ID of content  $c_i$  is called the *root* or *home* node for content  $c_i$ . *Indexing* is a redirection service that links content-IDs with the nodes storing the



item(s), namely *holder* nodes. The operation of *indexing* can take two forms: register and fetch. The register operation is when a node adds or updates its information to the index and fetch is when nodes request the list of nodes storing a given item. The *indexing* operation consists of several steps. First, nodes *query* for the root-node of the file they are interested in. Second, nodes contact that node directly to either register or fetch information. Finally, a *storage* operation is simply a request to download or upload data. For example, if node A sends a download request to node B, then A is consuming resources from B, and if the request is an upload, then A is publishing content to B.

Nodes can play four alternative roles in the network, namely *publisher*, *index*, *holder* and *requestor*. The *publisher* node is the original creator of a data item. An *index* node is the *root* node for a specific item. A *holder*, is a node storing a partial or complete copy of an item. Finally, *requestor* nodes represent users downloading information from other nodes.

Content ownership can take two forms; nodes keep copies of the files for which they are *root* nodes (assuming temporal ownership of the item), or they simply keep pointers to the actual location of the file, namely the *publisher* node. We will refer to these two variants as *root ownership* and *node ownership* respectively.

For retrieval of content, the *node ownership* model implies that *publisher* nodes periodically contact their *root* nodes to refresh their indexing information, and *requestors* perform a two-stage process to retrieve content from the overlay. First, they ask the *root* node for a list of candidates, and second, they attempt to contact one or more of these candidates directly to retrieve the item. Successful content retrieval is a dynamic process that involves the *publisher* node, the item's *root node* and the *requestor* node. For a successful retrieval of content two conditions must be met *i*) the item's indexing information (at the *root* node) is updated and *ii*) the remaining

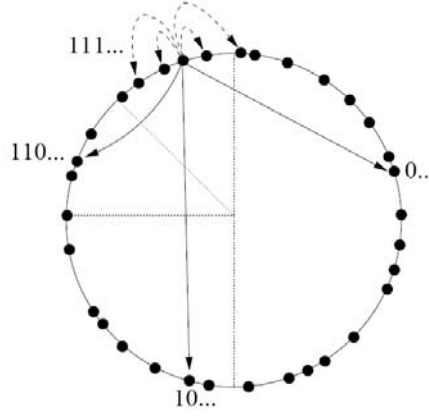
session length of the *publisher* and *requestor* nodes are long enough to complete the transfer. Both of these conditions are influenced by churn. The first condition requires *publisher* nodes to continuously monitor the availability of their *root* nodes and find new ones as needed.

In the *root ownership* case, items are first uploaded to their *root nodes*. Still, *publisher* nodes might be required to monitor their root nodes to guarantee the availability of their contents. However, *holder* nodes do not interact with peers requesting an item. Instead, items are downloaded directly from the *root* nodes. This mechanism decouples, to some extent, the availability of items with the liveliness of their *publisher(s)*, since the *publisher* node can be offline while the *root* node is still uploading an item to another peer. On the other hand, if the *root* node leaves the network, the item will need to be transferred once more into a new *root* node; regardless of any present or future demand for the item. This model is assumed in distributed storage applications such as CFS [34] and OceanStore [35]. The problem with this approach is that in order to minimize the maintenance overhead, the participant nodes must have high availability, which is not the case for open P2P application environments.

## 2.2 BAMBOO

Bamboo [21] is a reengineered version of Pastry [36], a DHT-based P2P routing architecture that uses a circular identifier space, with IDs 160 bits long organized as a sequence of digits base  $2^b$ . In Bamboo, nodes maintain two sets of neighbors, the *left set* and the *routing table*, illustrated in Figure 8 by dashed and solid arrows respectively. The *leaf set* contains the  $k$  nodes preceding and the  $k$  nodes following the current node in the circular identifier space. The routing table is a set of nodes organized in matrix form. All the node identifiers (nodeIDs) in row  $l$  coincide in  $l$  digits with

the current nodeId and the value of the next digit determines their column in the routing table. That is, a node in row  $l$  and column  $i$  shares  $l$  digits with the current nodeId and its  $l+1$  digit has a value equal to  $i$ .



**Figure 8. Neighbors in Bamboo.**

The basic functionality in Bamboo (and other DHT-based routing architectures) is a distributed (*key*, *value*) lookup service. Given a target key  $D$ , nodes follow the following algorithm:

- 1) Check *leaf set*. If  $D$  lies within its *leaf set*, then it forwards the query to the nodeId numerically closest to  $D$ . If that node is the local node, routing terminates. If not, continue with step 2.
- 2) Route message. The node computes the longest matching prefix between  $D$  and its own nodeId, denoted by  $L$ . The value of the first digit in  $D$  different from the local nodeId can be denoted as  $D(L+1)$ . If the node has a non-empty routing entry at row  $L$ , column  $D(L+1)$  it forwards the request to that node; otherwise it goes to step 3.
- 3) Forward to leaf. The message is forwarded to the member in the leaf set numerically closest to  $D$ .

The process above is performed by every node that the request is forwarded to. When the final destination is reached (step 1) a message is sent back to the originating node with the nodeId and the network address of the destination. In DHT-based systems, the processes of routing and

data location are functionally the same. That is, the lookup service we just described is used for both routing and resource discovery.

Figure 9 illustrates the basic procedure employed to map nodes and content into the overlay's identifier space. For nodes, their IP socket (or other form of unique identification) is used as input of the hash function to obtain a key value, namely the nodeId, which is used for routing and to claim responsibility for a portion of the identifier space. For data items, the identifier (e.g., file name) is used as input of the hash function, and the key value obtained determines which node is responsible for managing the lookup service for that item. The node with the numerically closest nodeId to the object's key value is called the root node or home node for the object.

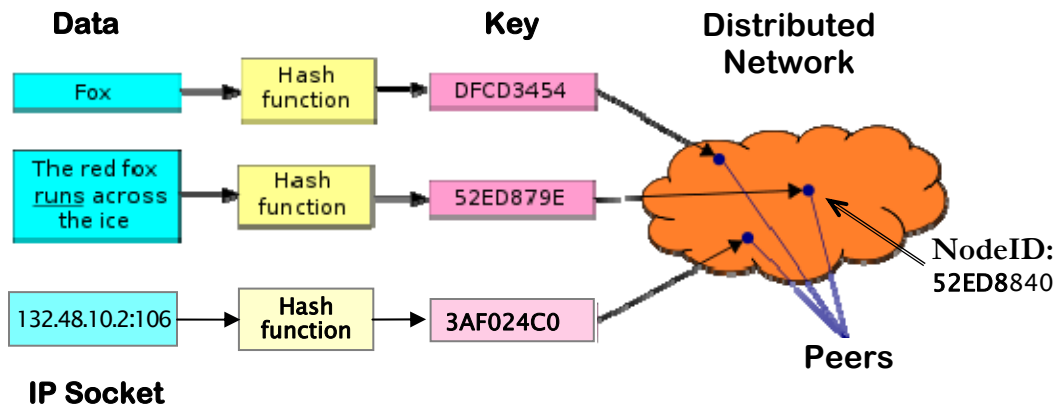


Figure 9. Distributed Hash Table

The scalability of DHT-based P2P system is derived from the structure of their routing tables. The distance between the local nodeId and the target key is reduced logarithmically at every step throughout the lookup process. Thus, DHT-based networks are able to provide a deterministic location service in  $O(\log N)$  steps, where  $N$  is the number of nodes in the network. Other DHT-based systems have similar scalability metrics and also share other functional and structural similarities with Bamboo, but their description is beyond the scope of this work and comparative descriptions are available elsewhere [25, 37].

The resiliency of DHT-based architectures is derived from the routing geometry embedded in their neighbor set. That is, the pattern of neighbor links across the overlay, independently of the routing algorithms or state management algorithms used [38]. The work of Gummadi, et. al. presented in [38] describes that the addition of sequential neighbors significantly increases the *static resiliency*<sup>16</sup> of a system, but increases its latency. In addition, the flexibility in neighbor selection (FNS<sup>17</sup>) is better than using flexible route selection (FRS<sup>18</sup>) to improve the performance of the system (i.e., latency). In that regard, Bamboo’s *leaf set* provides good static resiliency, and *proximity neighbor selection* is used to boost performance when more than one node can be used to fill in a routing table entry.

The creators of Bamboo summarize the factors allowing the system to handle high levels of churn efficiently in three functionalities: periodic recovery, adaptive timeouts and proximity neighbor selection. After a node failure, the remaining nodes need to reorganize their *routing table* and possible their *leaf set*. To avoid overwhelming the system with route repair messages, Bamboo opts for a periodic route maintenance policy that bounds the amount of bandwidth consumed. In addition, Bamboo maintains different timers for each of its neighbors (i.e., timeouts), which allow it to discern judiciously between node failures and network congestion or processor load. Lastly, Bamboo uses a two steps process to fill in its routing table entries. First, it performs a lookup for a random identifier with a required prefix  $p$  corresponding to a hole in its routing table and uses the returned value to fill it in. Secondly, nodes query their neighbors’ routing tables to find alternative entries providing better latencies. The objective of the first process, called *global tuning*, is correctness, while the goal of the second one is performance. The system prioritizes correctness

---

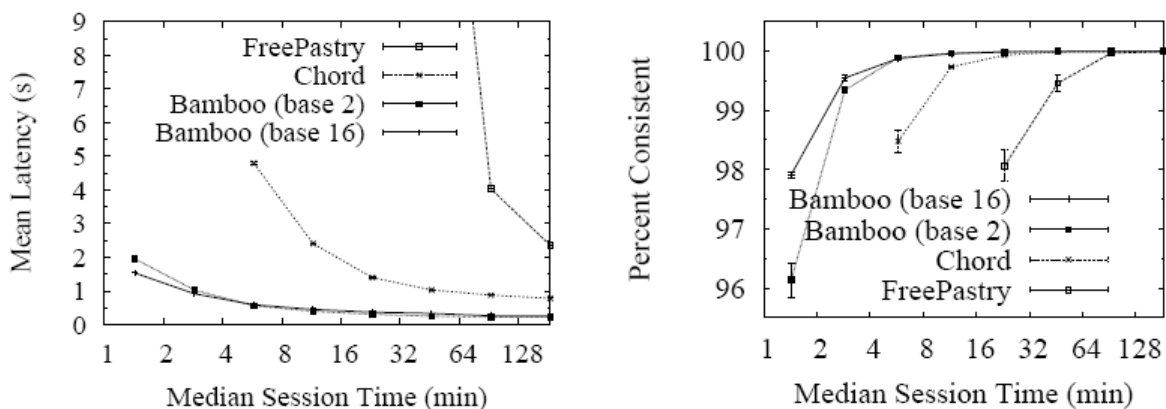
<sup>16</sup> Static resiliency describes the capacity of the network to route messages after failures and before route repairs are performed. This is a fundamental property in dealing with churn.

<sup>17</sup> Flexible neighbor selection is a measure of the level of flexibility nodes have to fill in entries in their routing tables.

<sup>18</sup> Flexible route selection is a measure of the level of flexibility of the routing algorithm to select the next hop.

and improves its performance opportunistically when the amount of traffic being managed by the node allows it.

The results presented in Figure 10 are taken from Bamboo’s technical report [13]. In these graphs, the capacity of Bamboo to handle extremely high levels of churn (i.e., small median session times) can be easily appreciated. Even when the median session length falls below eight minutes, Bamboo manages to perform most requested lookups with a mean latency that outperforms other systems. Nonetheless, we showed in [39] that if the maintenance intervals of Chord (or other DHT-based system) are tune up according to the expected level of churn, the performance level obtained is comparable to Bamboo.



**Figure 10. Bamboo’s Performance Under Churn**

Bamboo was selected over other DHT-based architectures to analyze the effects of churn on content availability because of its software architecture and its probed performance in churning environments. Bamboo is a mature open source DHT implementation written in java using an event-driven single-thread programming style. Thus, we can obtain accurate measurements of a deployed DHT<sup>19</sup> rather than the simplified DHT implementations available in P2P simulators [40]. In addition, there is a reasonable amount of documentation available to guide our

<sup>19</sup> Bamboo was offered as a public DHT service in PlanetLab during 2005-2009.

development efforts [19] without having to deal with detailed implementation issues related with the operation of a DHT application substrate.

## **2.3 LITERATURE REVIEW**

### **2.3.1 Churn Models**

The continuous and unsynchronized arrival and departure of nodes to and from the overlay network is called churn. When nodes leave the network, the content they were contributing also disappears. Thus, understanding the content availability problem space requires analyzing the fundamental properties of churn in P2P networks. Some research initiatives have tackled this subject using analytical models, and some others using empirical studies.

Yao, et al. in [41] highlight that the heterogeneity of lifetimes and off-line intervals is a fundamental property of P2P networks. The authors develop a model for heterogeneous user churn, from which they derive multiple closed form expressions that characterize the properties of nodes in unstructured P2P networks, including their residual lifetime distribution and isolation probability. In our work, we use a variant of this model to recreate the churn behavior of existing P2P networks. Their work is extended in [42] to accommodate non-stationary arrivals of nodes, and create a sampling technique to measure session lengths in unstructured P2P systems. These papers provide insightful findings with regard to the connectivity and availability properties of nodes in P2P networks, but their scope does not contemplate the implications of churn at the content level.

There is a long list of measurement-based research initiatives using probabilistic models to characterize the distribution of session lengths in deployed P2P networks. These characterizations are important to understand the properties of these systems and to synthesize churn events for simulation studies. Pareto [43, 44] and Weibull [45, 46] are the most commonly reported distributions in the research literature, but the results presented in these studies do not include node level characterizations. Yao, et al. in [41] explain that a perfectly shaped Pareto or any other distribution can be the product of a mix of multiple independent exponential distributions. Thus, it is not possible to conclude the real distribution of session lengths of individual nodes by just observing their aggregated behavior.

Luo, et al. in [47] present an alternative view for the characterization and synthesis of churn. Instead of analyzing the network as a single black-box, the authors divide the user population in multiple geographical regions and assign different user behaviors to them –as in our work presented in [48].– This model is used to generate the cyclical node membership exhibited in deployed P2P systems. The authors also provide a set of MatLab tools that automate the generation of churn logs for simulation studies.

Presently, the research community has not yet agreed in a general model for the distribution of session lengths of individual nodes or their aggregate distribution. As an alternative, Fernández-Casado, et al. in [49] developed a tool integrating several competing models to build a general purpose churn generator capable to generate complex churn logs for simulations. The availability of this type of tools facilitates the comparative analysis on the performance of P2P systems using different churn models. Nonetheless, most of the recent research literature in P2P networks uses public network traces [50, 51] to recreate the node dynamics of deployed systems. In our case,



these measurements lack the appropriate granularity and flexibility for analysis at the content level. Furthermore, some of them could present significant measurement bias [52].

In summary, there has been much progress on modeling and characterizing churn in P2P networks. However, the sheer size and heterogeneity in this type of networks make it quite difficult to obtain definite fine grained measurements and models. Our research combines the results presented in diverse measurement studies with analytical models to build a flexible churn framework to analyze the fundamental content availability properties of P2P networks and possible methods to improve it.

### **2.3.2 Redundancy Methods**

Redundancy is a mechanism commonly used to improve the reliability of systems by provisioning excess resources. For distributed data storage, there are multiple methodologies to achieve this reliability. Each of these methods uses different tradeoffs between data reliability and other system properties such as storage space or maintenance bandwidth. The research literature on redundancy in P2P environments can be categorized into three categories: *i*) studies matching a given networking condition (i.e., average node availability) with a redundancy method *ii*) studies defining the optimal parameter settings for a redundancy method given a set of requirements (e.g., required file availability and average node availability) and *iii*) studies defining new redundancy data structures. In addition to these categories, it is important to highlight that a major concern of using redundancy in P2P environments is the amount of bandwidth required to regenerate any data loss due to node departures. This is regularly referred as the *repair problem* [10], [53]. This concern generates a new cost vs performance space for the design of redundancy methods. The

fundamental constraint for this problem is that in P2P networks access bandwidth is scarcer and more expensive than storage space [10].

With respect to the first category, Rodriguez and Liskov in [10] argue that for system with high node availability (such as PlanetLab) replication redundancy should be preferred and furthermore, for other system settings the storage overhead savings of erasure coding might not be worth the associated cost, due to the added system complexity. Lin, et al, in [54] present an analytical expression in terms of storage overhead ( $S$ ) and node availability ( $a$ ) to determine which redundancy method provides better file availability. If  $S*a < 1$ , replication performs better and for  $S*a > 1$  erasure coding performs best.

With respect to the second category, If the system needs to be engineered to achieve a required level of file availability (i.e. 0.99), Rodriguez and Liskov in [10] or Bhagwan et al. in [55] present analytical expression to obtain the optimal storage overhead for both erasure coding and replication. This optimal value is such that the product  $S*a$  is always greater than one; thus, according to the arguments presented by Lin, et al, in [54], erasure coding redundancy is always preferred. In [56], Dimakis, et.al., contrast the redundancy repair cost of erasure coding, with two variants of network coding. Their results indicate that it is possible to construct coding methods with significant repair bandwidth savings over erasure coding over a wide range of target file availabilities. However, their results also indicate that as the network becomes unstable (i.e., lower average node availability) the performance of network coding can become inferior to the performance of erasure coding.

With respect to the third research category, that is, the introduction of new redundancy methods, the research literature is abundant. A recent survey of the field of network coding by Dimakis, et al. in [53] present recent advancements in network coding that reduce the redundancy

repair problem by orders of magnitude compared with standard erasure codes. In addition, this paper describes the different repair modalities of erasure coding, highlighting that minimum bandwidth with exact-repairs is the best suited network coding modality for distributed storage applications. In that regard, Rashmi, et al. in [57] present a minimum bandwidth exact-repair (MBR) explicit construction code for any combination of the parameters  $(m, k, d)$  where  $m$  is the total number of nodes,  $k$  is the number of blocks needed for the file reconstruction and  $d$  is the number of nodes required for the reconstruction of lost fragments. This is the first explicit code construction that allows the selection of the number of nodes ( $m$ ) independently of other system parameters. Dominuco and Biersack in [58] present an alternative construction of erasure codes, called Hierarchical Codes. The authors argue that cost is improved because the average number of transfers needed is lower than in traditional erasure coding (despite needing a higher number of repairs). However, the authors account the number of transfers as cost rather than the product of the number of repairs times the average number of blocks transferred (i.e., bytes transmitted). For the mechanism presented by the authors, the task of choosing which blocks should be downloaded to perform a repair has been reduced compared to the cost of traditional erasure coding, but it is still a non trivial problem. Williams, et al. in [59] evaluate different redundancy techniques for P2P storage. They advise the use of hybrid schemes combining replication and either erasure coding or bucketing (data bundling) as a reasonable compromise between maintenance cost and availability. The authors conclude that replication should be used to simplify data access (most of the time) and erasure codes should be employed to achieve the last nines in the desired availability level. In that regard, Wu, et al. in [60] and Xu, et al. in [61] propose hybrid redundancy mechanisms that combine erasure coding and replication. In the work of Wu, et al. [16, 60] content resides on the user-nodes. These nodes regularly contact a set of  $M$  Indexer nodes to register their contents

(whole-file and fragments) and based on this information, Indexer nodes decide when to send redundancy repair instructions back to the user-nodes. Whole-file replication is assumed as a by-product of user activity. Assuming the existence of at least one whole-file replica available, the system only employs erasure coding to reach a target file availability level for those items with insufficient replicas available. In the work of Xu, et al. [61] on the other hand, the location of content is determined by a hash function. Both nodes and content are mapped into a virtual identification space in which several physical nodes share the same virtual ID and are responsible for maintaining a target number of replicas for each object. One shortcoming of all the initiatives mentioned above is the assumption that all nodes cooperate fully in the redundancy maintenance process when needed. In our research, the structure of our redundancy method is different. We use erasure coding as the foundation to achieve reliability, and replication as a mean to minimize (and simplify) the redundancy repair problem. Furthermore, we integrate the redundancy maintenance problem with economic models to overcome the problems of cooperation and fairness in the context of distribution of content in P2P networks.

### **2.3.3 Incentives and Content Availability**

The next bodies of research related with our work are the use of incentives mechanism in P2P networks and the focus of our research, content availability itself. The research literature on economic-based mechanism for P2P systems is vast. The issues commonly addressed are fair allocation of resources and prevention of free-riding [62, 63], but content availability is rarely mentioned as a desirable property or by-product of these mechanisms. Furthermore, the resource most commonly managed by these mechanisms is bandwidth, which is a short-lived property of the system [64] and does not translate directly into our objective of improving content availability.

Geels and Kubiawicz [14] argue that solutions for large-scale replica management should be based on economic models and they outline the benefits of adopting this approach. They introduce the term Replica Management Economy (RME) to describe this type of systems and highlight that automatic resource management, scalability and guarantees through mechanism design are the key advantages of using economic models to deal with this problem. The authors state that RME allow the level of node autonomy that is necessary in a heterogeneous environment like the Internet; regulating the interactions between nodes while fostering cooperation across domains. The authors explain that one of the directions of future research in this field is the design of utility functions to rate the worthiness of alternative actions by a RME player. Our research advances on this path.

To the best of our knowledge, the only two research initiatives targeting content availability in P2P networks specifically are incentive-based mechanisms. Antoniadis, et al. in [18] present an incentive mechanism to control the minimum amount of time that nodes should participate in the system, as well as the minimum number of files that they should share throughout that time. Bai, et al. in [33] on the other hand, analyze the use of bundling<sup>20</sup> to improve the availability of contents (in BitTorrent in particular) improving even the download time experienced by peers when publishers exhibit high unavailability. Our research differs from these works in three aspects: *i*) we want to investigate the use of incentive-based mechanisms that would not require nodes to change their churn behavior, *ii*) we want to explore the system design tradeoffs in the context of DHT-based P2P systems and *iii*) our system does not differentiate between different types of content; thus, it foster content diversity (regardless of popularity), minimize node accountability and preserves node autonomy. That is, the system allows nodes to

---

<sup>20</sup> Bundling consist on handling more than one object together to be transferred as a unit.

decide unilaterally how much information they want to store; in the understanding that their performance is a function of their contribution.

An additional aspect related with the construction of incentives mechanisms is the definition of redundancy maintenance strategies. In that regard, Data and Aberer in [65] employ a Markov model to analyze the performance of different redundancy maintenance strategies in P2P networks. They determine that a randomized lazy redundancy maintenance mechanism offers significant advantages over existing deterministic and procrastination mechanisms. Yet, they assume full cooperation among nodes in their analysis. Sit, et al. in [66] present a proactive replication system that is capable to maintain high content availability using only several kilobytes per second of bandwidth. The authors argue that “creating redundancy constantly at a limited rate is simple, flexible and effective approach to maintain data durability”. However, their evaluation is based on the availability of server-like distributed systems (PlanetLab) and full cooperation among peers. In contrast, our work uses randomized redundancy maintenance strategies integrated with incentives mechanisms. In addition, the structure of our redundancy method can be understood as a proactive redundancy repair strategy, which in the long-run minimizes the system’s repair bandwidth.

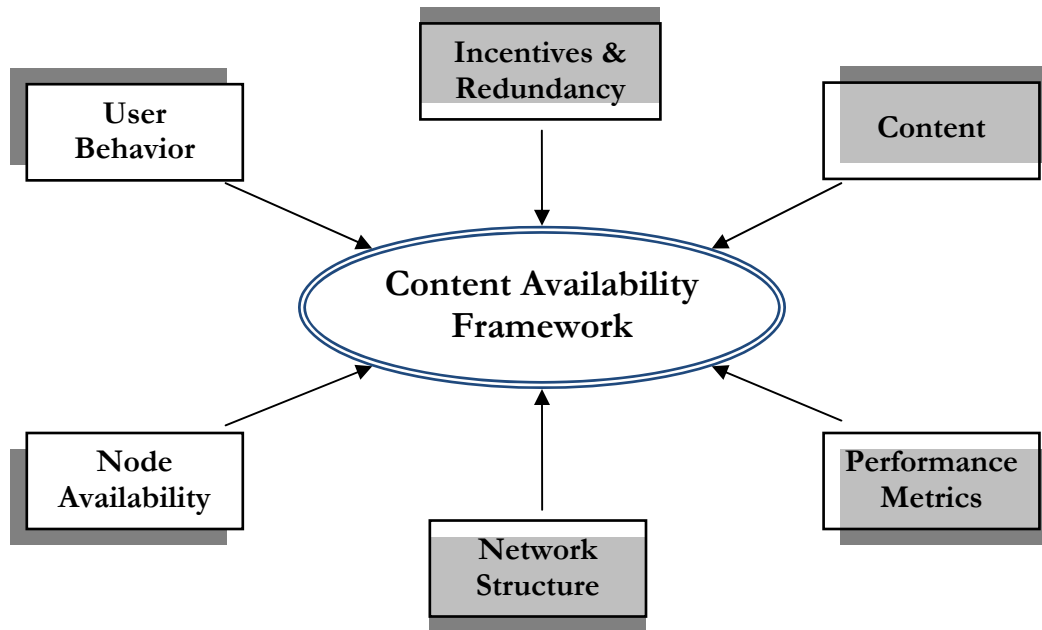
### 3.0 CONTENT AVAILABILITY FRAMEWORK

In this chapter we present a framework for studying content availability in P2P networks. The construction of this framework pursues two fundamental goals: First, to guide our analysis by organizing the diverse factors that can determine the performance of a P2P network and second, to describe the fundamental properties and models that characterize the architecture and operation of a P2P network. The elements portrayed into each component of this framework determine the settings for the various setups in the experimental portion of this work.

The layout of this chapter is as follows. Section 3.1 describes the overall structure of our proposed framework. Section 3.2 describes the models embedded into each of its components, and Section 3.3 presents the parameters to be used in each component during the experimental portion of this study and describes the rationality for the settings selected.

#### 3.1 FRAMEWORK STRUCTURE

Our content availability framework has six components, named *User Behavior*, *Node Availability*, *Network Structure*, *Content*, *Performance Metrics* and *Incentives and Redundancy*. Figure 11 depicts our framework.



**Figure 11. Content Availability Framework**

In conjunction, the six components of our framework constitute a complete characterization of a P2P system, from a content availability perspective. The properties and models included in our framework define not only the fundamental characteristics of the network (e.g., routing architecture) but also the performance metrics to be used as a measure of the success or failure of the system.

The first component in our framework is called *User Behavior*. It models the human attitudes that affect the performance of P2P systems. The second component, *Node Availability*, models the intermittent connectivity of individual nodes. That is, this component models when nodes are online and offline. The third component, *Network Structure*, models the elements related to inter-node communication. Consequently, the network routing architecture, the underlying Internet topology and other network-related metrics fall within this component. The fourth component, *Content*, characterizes the properties of the content being offered in the P2P network. The fifth component, *Incentives and Redundancy*, characterizes and models the two fundamental



mechanisms we propose as solution to improve content availability in P2P networks. Finally, the sixth component, named *Performance Metrics*, describes the response variables used to measure the performance of the system and the effectiveness of the mechanisms we propose. The following sections will further describe the elements incorporated into each component.

## 3.2 FRAMEWORK MODELS

### 3.2.1 User Behavior Component

The *User Behavior* component models the fundamental user attitudes toward P2P networking that could affect the performance of a P2P system. It has been widely documented that P2P networks are highly heterogeneous [67]. For instance, some users share lots of resources while others share little to nothing at all [68]. In this work, we recreate the heterogeneous nature of users behaviors in P2P networks using a two class user profiling technique. We name these two classes *Benefactors* (*Be*) and *Peers* (*Pe*). Figure 12 illustrates our user profiles model and the properties that characterize each of the two profiles we use in our framework.

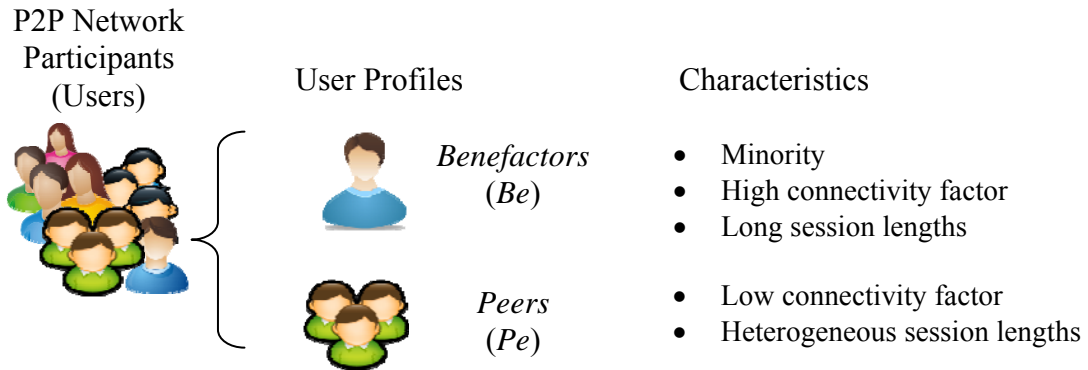


Figure 12. User Behavior Model: User Profiles

Given a set of nodes, a small fraction of them is assigned at random to the *Benefactors* user class and the rest are assigned to the *Peers* user class. Both profiles are characterized using two fundamental properties: *connectivity factor* and a probability density function for the average session length of the nodes in each class. We define the term *connectivity factor* as the probability of finding a node online during its lifetime. The average session length is defined using a set of probability density functions to be described in the next component, node availability.

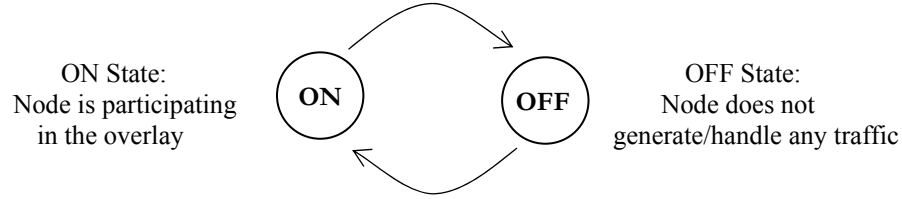
With the *Benefactors* user class, we want to portrait a set of users that contributes resources generously and consistently to the network. Thus, we assume that *Benefactors* (*Be*) are users who stay connected to the overlay network during long periods of time and exhibit a high *connectivity factor*. In our analysis we want to investigate if the stability of the overlay network depends on the contributions of the *Benefactors* user class. That is, we want to analyze if the portion of the network population within this class can play a crucial role in the performance and the type of applications that can be deployed in P2P networks.

We use the *Peers* (*Pr*) user class to complement the *Benefactors* user class recreating the heterogeneous node participation pattern reported in deployed P2P networks. To accomplish this goal, we adapted the heterogeneous user churn model presented by Yao, et. al. in [41]. In their model, both the online and offline intervals are selected independently using two probability density functions (pdf). In our model, we use a single pdf to define the online interval and the offline interval is obtained using the *connectivity factor* parameter.

### 3.2.2 Node Availability Component

The *Node Availability* component captures the churn behavior of individual nodes. That is, the rates at which nodes join and leave the overlay network during their lifetime. An On/OFF state

model is used to characterize this behavior. Figure 13 presents this model. The “ON” state represents the time interval during which nodes are active members of the overlay network. During this time, nodes generate lookup traffic and attempt to download items from other nodes. The “OFF” state represents the time elapsed offline until the node rejoins the network.



**Figure 13. Node Churn Model**

The term node availability ( $a$ ) is used in the research literature to characterize the probability of nodes being active. To avoid confusion with other availability metrics that we use throughout this dissertation, the term *connectivity factor* ( $Cf$ ) is used instead. For the ON/OFF model presented above, connectivity factor is measured as the quotient of the average ON interval over the sum of the average ON and OFF intervals of a node. Throughout this work the average online interval is called average session length. Within the *Node Availability* component  $\overline{ON}$  represents the average session length, and  $\overline{OFF}$  represents the average offline interval. Thus, the *connectivity factor* can be expressed as follows:

$$Cf = \frac{\overline{ON}}{\overline{ON} + \overline{OFF}} \quad (3.1)$$

In our description of the User Behavior component, we mention that  $\overline{ON}$  and  $Cf$  are the characterizing factors for each user class. Consequently, rather than using equation (3.1), as presented above, we employ the following variant in our model:

$$\overline{OFF} = \overline{ON} \cdot \left(1 - \frac{1}{Cf}\right) \quad (3.2)$$

### 3.2.3 Network Structure Component

The *Network Structure* component models the way nodes are organized and how they interact with each other to maintain the overlay. The fundamental elements that can determine the behavior of the system under churn include the underlying network topology, the network size ( $N$ ) and the P2P network architecture (i.e., routing mechanism). For this component, we assume that the properties of the underlying physical network and the overlay network topology can be modeled independently.

For the properties of the physical network, we assume a WAN setting and we employ existing tools to define the characteristic bandwidth and delay properties of the network. In particular, we employ Inet-3.0 [69] to generate our model of the physical network.

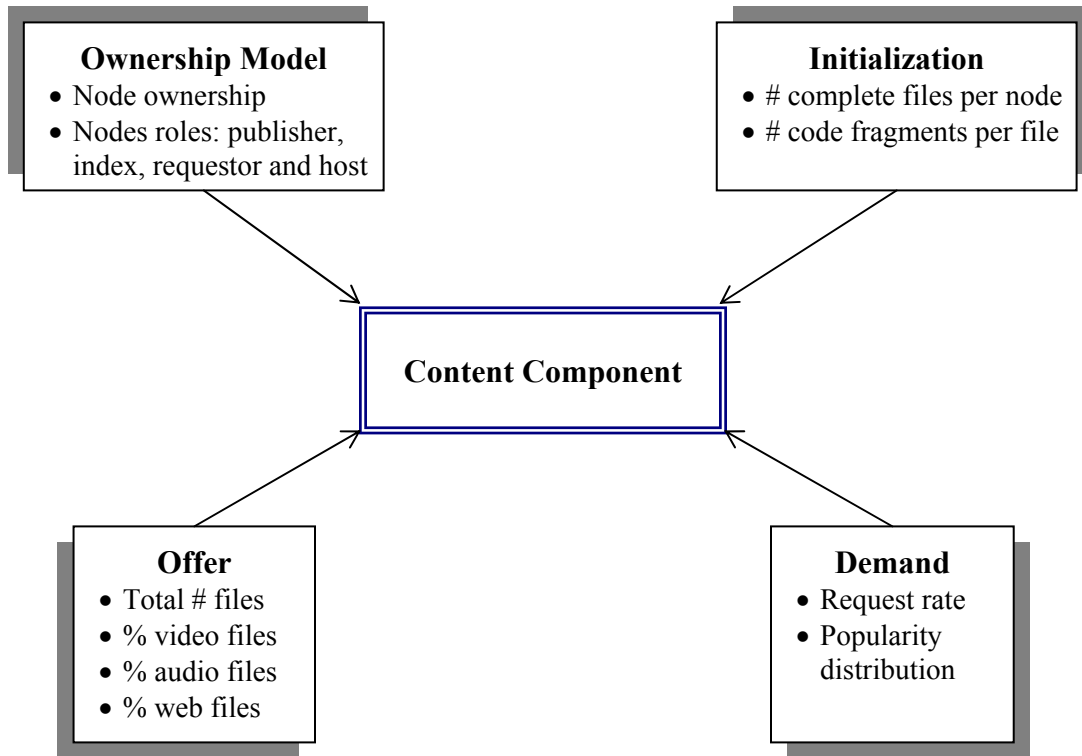
In our model, Network size ( $N$ ) represents the total number of nodes participating in the overlay network and the term overlay network size ( $N_o$ ), is used to denote the average number of simultaneously active nodes in the network. These two elements are related by a constant, called *churn factor* ( $Ch$ ). That is,  $N_o = Ch \cdot N$ .

The structure of our framework does not preclude the analysis of multiple P2P network architectures, but the scope of our current work is limited to structured systems. Thus, the architecture of the network is assumed to be DHT-based.

### 3.2.4 Content Component

The *Content* component models the data items stored in the overlay network and the exchanged patterns of these items among nodes. For our characterization of content we categorize factors into subcomponents named *offer*, *demand*, *ownership model* and *initialization*. In addition, this

component describes the different roles that nodes can play with respect to content. The main factors included in this component are presented in Figure 14.



**Figure 14. Content Component Factors**

The factors included in the *offer* subcomponent are quantitative in nature. We use the term content space in reference to the total number of files ( $M$ ) to be assigned to the nodes participating in the network. These files are categorized by type into video, audio and web-like items and initialized in size according to results presented in the research literature [70]. This type/size breakdown is expressed by the tuple  $\langle A, V, W \rangle$  where  $A$  is the percentage of audio files,  $V$  is the percentage of video files and  $W$  is the percentage of web-like content. The *demand* subcomponent is much simpler. It defines the item selection process (i.e., request rate and items popularity).

The content component defines the initial assignment of data items to individual nodes. That is, before starting the simulation. In addition, the relative frequency at which each individual

file is requested by peers and the average time elapsed between requests are also defined in this component.

In Section 2.1.4 we described two content ownership models to describe the process involved in the retrieval of content (creation, storage, indexing, etc.) and the roles that nodes can play with regard to content management (*publisher*, *index*, *holder* and *requestor*). In our research, the *node ownership* model is assumed. In addition, the following content indexing and retrieval model is assumed for the rest of the dissertation. Nodes play three alternative roles, namely *holder*<sup>21</sup>, *index* and *requestor*. Content *holders* have one or more items in their storage space and they regularly contact their corresponding *index nodes* for indexing purposes. Consequently, *index nodes* have a list of all the nodes with items (i.e. *holder nodes*) that map to the portion of the key-naming space they manage. When a *requestor node* searches the overlay network for an item, it only needs to reach the *index node* to obtain a list of all the available nodes holding a copy of that item.

### 3.2.5 Incentives and Redundancy Component

The purpose of this work is to improve content availability in P2P networks. To do so, we propose the use of redundancy and incentive-based mechanisms. The first of these two mechanisms, redundancy, is a common practice in multiple fields as the mean to build reliable systems out of unreliable components. The second mechanism, incentives, is commonly used to model interactions among autonomous self-interested participants and as a mean to accomplish cooperation among them towards a common goal.

---

<sup>21</sup> The term *publisher* node is replaced by *holder* node to generalize data possession to the case where content availability is provided using redundancy.

In our research, we use redundancy to accomplish reliability under the assumptions that, in P2P networks, idle repair bandwidth is scarcer and more expensive than idle storage space [10]. Thus, we have a redundancy optimization problem, which needs to be formulated as a repair bandwidth minimization problem. In addition, we assume that given the heterogeneity of resources and behaviors among P2P networks' participants, incentives-based mechanisms is a proper tool to model their interactions, foster their cooperation and achieve self-organization (for the redundancy repair process of the redundancy scheme we propose).

Given the importance of these mechanisms in our work, the analysis and description of these components is presented in separate chapters. Redundancy is discussed in Chapter 4.0 and Incentives in Chapter 5.0.

### **3.2.6 Response Variables**

For complex systems such as P2P networks it is not possible to fully characterize their operation using a single metric. Li et. al. in [71] propose a two-metric framework named cost-performance to evaluate the design tradeoffs of DHT-based P2P networks. The cost metric accounts for the amount of traffic exchanged among nodes to maintain the operation of the network and the performance metric reflects lookups delay. In our previous research [39], we shown that in some settings, the lookups success rate is an vital response metric, which we named efficiency. If the system efficiency is poor, the metrics in the cost-performance framework become irrelevant. In addition, for our content availability analysis these metrics (cost, performance and efficiency) need to be qualified at the content level. We need to define the additional traffic that should be considered

in the cost metric and whether performance should be still measured using lookups delays or other time-delay metric, such as the delay between the start of a file lookup (i.e., root lookup) and the initial transfer of data. Furthermore, since our proposed mechanisms to improve content availability in P2P networks are distributed in nature, we need metrics reflecting not only the aggregate behavior of the system (i.e., mean cost), but also the distribution of these metrics among nodes as a function of their level of participation.

In our study, we define two metric categories that we named aggregate and spread. In the aggregate category, we use three metrics: lookup success rate, content retrieval success rate and redundancy repair bandwidth. The first two metrics characterize the efficiency of the system and the third one measures cost. In the spread category, we use only one metric, which is transmission bandwidth between nodes versus content contribution, where content contribution is measured as a function of the number of data items (both fragments and complete files) stored by nodes. Therefore, the spread metric is in fact the combination of two system properties. In addition to the metrics mentioned above Table 3 presents one more response variable that is important to characterize the operation of our system, the level of churn in the network. We measure this property as the mean (or median) session length of nodes throughout a simulation.

**Table 3. Response Variables**

Metric	Description	Units
Cost	Redundancy-maintenance cost	Bytes/file/second Bytes/file/node/second
Efficiency: <ul style="list-style-type: none"> <li>• <i>Routing-level</i></li> <li>• <i>Content-level</i></li> </ul>	lookup success rate content downloads success rate	%
Spread	Distribution of transmission bandwidth versus content contribution	n/a
Churn	Mean/median session length	Minutes



In addition to the metrics mentioned above, which are quantitative in nature, the evaluation of our system requires a set of metrics for which we can not derive simple quantitative results. So is the case of fairness for the mechanisms we propose. In particular, we expect that our spread metric will reflect that the performance (i.e., transmission bandwidth) of nodes remains proportional to their contribution towards the reliability of the system (i.e. content contribution).

### 3.3 FRAMEWORK PARAMETERS

Our content availability framework can be used to study the impact of *i)* the overlay community composition in terms of classes of user behaviors (*Benefactors* and *Peers*), *ii)* churn in terms of the node's average session length and their distribution, *iii)* the content characteristics in terms of number of resources and their query distribution and *iv)* content management strategies to improve content availability. The last topic is the focus of this dissertation.

**Table 4. Content Availability Framework Parameters**

Component	Number of elements	Factor(s)	Description
<i>User Behavior</i>	1	$\Omega_{Be}$	Network make (user type distribution)
<i>Node Availability</i>	2 pairs	$E[g_n(t)], a_i$	Average session length and average availability for each user-class
<i>Network Structure</i>	1	$N$	Network size
<i>Content</i>	1	$M$	Content space size
<i>Incentives &amp; Redundancy</i>	2 sets	<i>Redundancy-Set</i> <i>Incentives-Set</i>	Redundancy method parameters & Incentive-based mechanism parameters

A full factorial design considering all the elements described in this framework has not been completed due to the extensive number of variants possible. Instead, only a reduced number

of elements have been selected to conduct our content availability study. Table 4 summarizes the elements we use as factors for the experimental portion of this dissertation.

The following paragraphs describe the levels selected for each of the factors listed in Table 4; as well as the initialization of other elements not considered as factors. Chapter 6.0 presents results for the different combinations of parameters selected.

### 3.3.1 User Behavior Component

Every node in the overlay network is mapped to one of the two classes used in our framework. Thus, a network make (i.e., portion of nodes in each user class) can be expressed simply with percentages.  $\Omega_{Be}$  and  $\Omega_{Pr}$  denote the percentage of nodes in the *Benefactors* and *Peers* user classes, respectively. Since  $\Omega_{Be} + \Omega_{Pr} = 1$ , only one of these values is needed to fully characterize the user population. In our experimental work we use three levels for  $\Omega_{Be}$ . The selection of levels for the percentage of *Benefactors* assumes that this user class is small in most P2P application environments. The values considered are presented in Table 5.

**Table 5. User Behavior Parameters**

Parameter	Levels
Benefactors ( $\Omega_{Be}$ )	5, 10 and 20

### 3.3.2 Node Availability Component

To define the churn behavior of every node we use two parameters, session length and *connectivity factor*. The *connectivity factor* is a constant parameter value for each user class and for the session length we use a probability distribution function (pdf). For the *Benefactors* user class, we use a

Pareto pdf and for the *Peers* user class we use an exponential pdf. The parameters for the pdf of individual nodes are obtained as follows. For *Benefactors*, all nodes use the same constant parameter value and for the *Peers* user class, the average session length of individual nodes is obtained using a Pareto distribution. These parameters are presented in Table 6.

**Table 6. Node Availability Component Parameters**

Node Class	Parameter	
	Average Session Length: $E[f_i(x)]$	Average Node Availability: $a$
<i>Benefactors</i>	$f_{Be}(x) \sim \text{Pareto with parameters:}$ $\alpha=1.09 \beta=0.85 \text{ hrs}$	0.75
<i>Peers</i>	$f_{Pr}(x) \sim \exp(\lambda_x): \lambda_x=G(t)$ $G(t) \sim \text{Pareto with parameters:}$ $\alpha=1.5/1.09 \beta=900/1350 \text{ sec}$	0.25

The following paragraphs describe in more detail our selection of parameters for the node availability component.

### 3.3.2.1 Session Lengths

Table 7 presents a list of measurement studies of deployed P2P networks and the most relevant churn properties reported in them. The nodes' median session time ranges from one minute to one hour and the distribution of session lengths reported includes Pareto, Weibull and Lognormal distributions. Some of these results show extremely high levels of churn, with median session lengths in the order of minutes [32, 72]. These results contrast with results assigning a value on the order of half an hour to the location parameter of a Pareto distribution [42], mis-estimating the relevance of short sessions. In either case, we adopt a conservative approach in our experimental settings by selecting the session lengths of users in such a way that the experimental distribution of session lengths fits with several of the empirical results of high churn reported in the literature.

Assuming that the behavior of super-peers in the Gnutella network constitute a behavioral match for our *Benefactors* user class, we use the characterization presented by Wang, et al. in [44] to set parameters for this type of nodes. According to this research, the session length of the *Benefactors* class follows a Pareto distribution with a shape parameter  $\alpha=1.09$  and a location parameter  $\beta=0.85$  hrs.

**Table 7. Churn Measurement Studies**

Reference	Network(s)	Churn Metrics/Comments	
Active Measurements Studies			
Bolla [73]	Gnutella	Lognormal distributed Session lengths: $\alpha$ =1.2, $\sigma$ =1.25 Offline intervals: $\alpha$ =3.42, $\sigma$ =2.0	90% < 2 hrs
Wang [42]	Gnutella	Sessions are power-law distributed with $\alpha$ =1.15, $\beta$ =0.69	
Wang [44]	Gnutella	Super-peers’s sessions are Pareto distributed with $\alpha$ =1.09, $\beta$ =0.85	5% < 8 min
Steiner [74]	KAD	Session are Weibull distributed with parameters (169.53, 0.6151)	50% < 155 min
Stutzbach [45]	Gnutella (G) BitTorrent (B) Kad (K)	Sessions are not exponential or heavy tailed	50% < 30 min for G, K 50% < 12 min for B 25% >2hrs for G, K 15% > 2hrs for B
Bustamante [43]	Gnutella	Sessions are Pareto distributed: $4.338x^{-1.0607}$	50% < 1hr
Bhagwan [52]	Overnet	IP measurement underestimate availability by factor of 4	50% < 1 hr
Chu [75]	Gnutella	Significant time of day effect	31% < 10 min 20% > 2 hrs
Saroiu [76]	Gnutella		50% < 1 hr 30 % > 2 hr
Passive Measurement Studies			
Gummadi [32]	Kazaa		50% < 2.4 min 90% < 28.25 min
Sen [72]	Fast-Track Gnutella DirectConnect		50% < 1 min 60% < 10 min

An exponential distributed model is used for the session lengths of the *Peers* user class, but their aggregate behavior follows a Pareto distribution. This is accomplished by initializing the average session length of each node independently at random using a Pareto distribution, called the

initialization function. The parameters of this distribution can take the following values: shape ( $\alpha$ ) can be 1.5 or 1.09, and the location parameter ( $\beta$ ) can be 900 or 1,350 seconds. This initialization is adapted from the churn model presented by Yao, et.al. in [41]. Table 8 presents the resulting quantiles of the distribution of session lengths using different parameters for the initialization function. The selection of these parameters represents a compromise between different measurement studies. For instance: Wang et al. in [42] state that session lengths are Pareto distributed with a location parameter of 0.69 hrs; Stutzbach et al. in [45] report that the median session length is 30 minutes; and other authors report at least 31% of the nodes leave before 10 minutes and approximately 20% staying connected longer than two hours [75, 76].

In Section 2.1 we argue that for a conservative evaluation of the system, higher levels of churn should be preferred, and in that regard, exponential session lengths exhibit higher churn levels than Pareto or Weibull models. In addition, the results presented by Zhonghong et al. in [77] indicate that the lookup performance of a DHT-based system using exponential, Pareto and Weibull distribution models does not exhibit significant performance differences. Thus, our selection of exponential session lengths for the *Peers* user class seems acceptable.

**Table 8. Peers User Class Sample Session Distributions**

Initialization Distribution	Session Lengths Quantiles
Pareto (1.50, 900)	[ 4.55 7.37 10.63 14.45 19.22 25.23 33.48 46.15 72.38 ]
Pareto (1.09, 900)	[ 4.68 7.90 11.60 16.03 21.62 28.70 39.32 56.49 96.58 ]
Pareto (1.50, 1350)	[ 5.63 9.78 14.69 20.35 27.53 36.53 49.02 68.95 110.77 ]

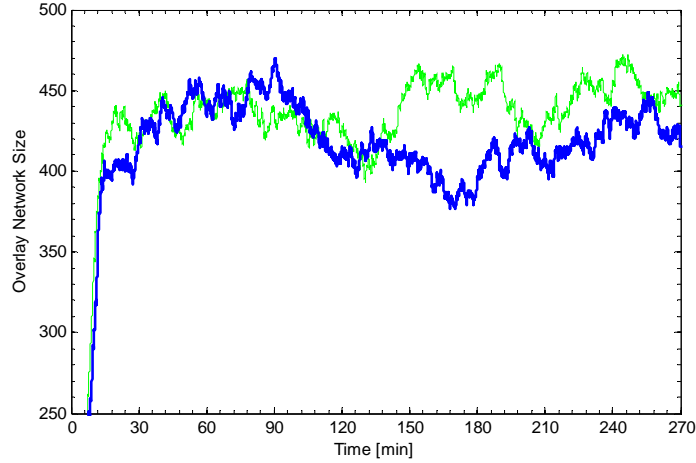
### 3.3.2.2 Off-line Interval

In order to fully characterize the churn behavior of a node we need to define the average interval nodes stay off-line. This is done using the average session length and the *connectivity factor* of each user class.

For the selection of the *connectivity factor*, *Benefactors* are assumed to participate actively in the overlay during 75% of their lifetime, and for the *Peers* user class lower connectivity factors are employed. In particular, a value of 0.25 is considered for the experimental work presented in this dissertation, which is in the same order of magnitude used by other research initiatives [10, 78].

### 3.3.2.3 Oscillating Network Size

One aspect of churn not being reflected into our experimental setting is the possible correlation of arrivals and departures of nodes. The existence of this phenomenon is suggested by the overlay network size fluctuations exhibited in deployed networks. Some researchers have tackled this problem by using non-stationary arrival models [42] and others by breaking down the user population into time zones and assigning them periodic activity patterns [47]. Either way, these models focus on reproducing the oscillating network behavior across a 24 hour cycle. In our case, we are evaluating the negative effects of churn on content availability over shorter periods of time, a few hours in particular. Still, the churn model presented in this dissertation does generate an oscillating network size and represents a more realistic node activity model than those presented earlier in the literature, such as Rhea, et al. in Bamboo [21] that uses a constant network size, or Li, et al. in [71] that uses a single exponential distribution to generate the session length of nodes. Figure 15 presents the evolution of the overlay network size in two experiments selected at random. It is clear that our churn model does produce oscillations in the overlay network size, but by no means are we trying to reproduce the cyclic behavior reported in the literature. The objective is simply to stress the dynamic properties of the routing mechanism beyond what has been done previously in the research literature.



**Figure 15. Overlay Network Size vs Simulation Time**

### 3.3.3 Network Component

The network component parameters use three elements: network size, DHT-routing architecture and Internet topology. Two of these elements are in themselves a complex entity and multiple parameters are required to fully characterize them.

It is unfeasible to conduct a detailed simulation analysis for realistic network sizes. Deployed P2P networks reach several hundred thousand nodes, but the simulation platforms for P2P systems available do not support content level analysis at such scale [40]. Nonetheless, we strive to use the biggest network size possible in all of our simulations because to some degree, route flexibility and content availability depend on the size of the overlay network (in the sense that it implies diversity). Thus, the scale of our experiments depends on the capacity of the simulation platform being used and the hardware resources at our disposal. In particular, the simulation results presented in this work implement a network of 1,840 nodes.

With respect to the DHT routing architecture, Bamboo [21] is employed with its default parameters in all the experiments reported in this dissertation, as shown in Table 9.

**Table 9. Bamboo Parameters**

Parameter	Value	Update interval
<i>Leaf set</i>	size = 3	4 sec
<i>Routing table</i>	base = 2	far alarm 10 sec near alarm 20 sec
<i>Ping interval</i>	-	20 sec

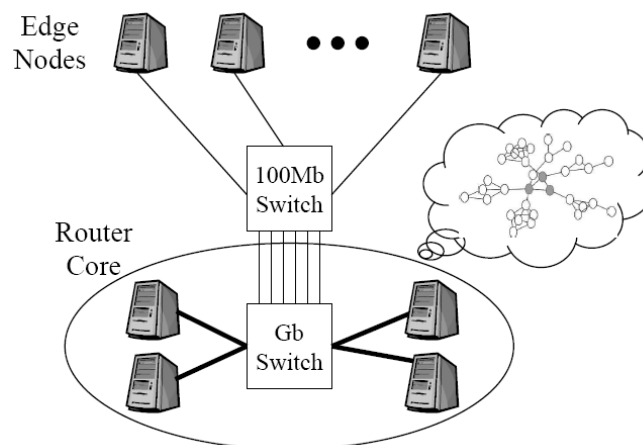
### 3.3.3.1 Modelnet

The Inet topology generator [69] is used in this dissertation to create the underlying network topology for all experiments. The topology consists on 1,344 edge nodes distributed across 836 distinct AS-level stub networks in a 4,000 node wide area network. From this point forward, this network topology model is referred as *net-model*. The bandwidth and delay restrictions for the *net-model* are enforced using Modelnet [20], which emulates a wide area networking environment using a cluster of machines interconnected by a LAN infrastructure.

Modelnet defines two types of machines: emulators and hosts. The emulator is nothing more than a packet forwarding agent that redirects packets according to the topological descriptors provided in a *net-model*. Host machines perform the function of edge-nodes, with their networking characteristics also defined in the *net-model*. For scalability, each host machine is multiplexed logically –according to its capacity– into multiple independent elements named Virtual Nodes (VN). Each VN is equivalent to an edge-node in the *net-model*. An application can then be deployed by executing instances of the application software being tested in each VN. From the application point of view, a VN instance is an end-host with its unique IP address. Each application instance is independent of all other VN instances running in the same physical node. The communication between these application instances is accomplished by forwarding all traffic through the emulator node.



For our evaluation of content availability mechanisms for P2P networks using Bamboo and Modelnet we employ the following resources: a cluster of twelve Red Hat Enterprise Linux 4 machines and one FreeBSD 4.9 system. The machine running FreeBSD 4.9 is a Dell PowerEdge 2550 Server with Dual PIII 1.4GHz, 2GB RAM and Gigabit Ethernet NIC. The Red Hat systems have the following mix of hardware resources: eight have Pentium III 1.4 GHz processors, 512 MB RAM and Fast Ethernet NICs; three have Core i7 2.66 GHz processors with 12 GB RAM and Gigabit Ethernet NICs; and one machine has AMD Athlon 64 X2 Dual Core 2.6 GHz processor with 6 GB RAM and Gigabit Ethernet NIC.



**Figure 16. Modelnet Network Topology**

The diversity in the architecture of the machines used in our experimental cluster creates some system complexities. Besides having to install a set of libraries for each architecture (32 and 64 bits CPUs) we also had to specify Java's TCP/IP stack version (by default Java uses IPv6 and Modelnet's library only works with IPv4). Moreover, we had to fine tune our data structures to maximize the CPU and memory utilization of every machine. The first system setting we used in our experiments supported a total 1,128 nodes; and after additional fine tuning of all the system parameters (in Modelnet, Java and the OS) we scale our experimental platform to 1,840 nodes.

Depending on their resources, each cluster machine hosts 88, or 176 virtual machines (VN) running one or two instances of Bamboo for a total of 1,840 overlay nodes. The FreeBSD machine uses Modelnet to enforce the wide-area delay and bandwidth restrictions of a 4000-node wide-area network topology with 1,344 client nodes connected to 836 distinct stubs by 10 Mbps to 100 Mbps links. Figure 16 illustrates the structure of the *net-model* employed in the experimental portion of the dissertation.

### 3.3.4 Content Component

The content component has three elements: content space size, query distribution, and type/size breakdown. The first two elements are a single value parameter, while the third one is a tuple with three values plus a model characterizing the file sizes for each data type. In addition, this section describes the procedure used to initialize the content space of individual nodes before a simulation starts.

The content space used in the experimental portion of this work is constructed as follows. A maximum number of unique items is defined for each simulation setting. In most cases, a value of 500 items is used, but 400, 750 and 1,500 unique data items are also employed in some experiments. Each item is assigned a unique numerical identifier and is assigned type and size properties according to the distributions indicated in Table 10, as reported in [70].

The selection of items each node holds at the beginning of a simulation is crucial for the performance and the evolution of the system's content. The scenario that we want to evaluate focuses on the long-tail portion of the popularity distribution. That is, those rare items for which the system has just a few whole-copies. For items highly demanded, content availability is a natural by-product of demand. Injecting additional redundancy into the system for these items

would be wasteful. Thus, in this study only a small fraction of nodes is initialized with whole-copy items

**Table 10. File Type/Size Distribution**

Audio	80%
<i>Distribution</i>	Lognormal
<i>Lognormal Parameters</i>	$\sigma^2=0.12$ $\mu = 1.42$ MB
Video	05%
<i>Distribution</i>	Lognormal if $x < 6$ MB, Pareto Otherwise
<i>Lognormal Parameters</i>	$\sigma^2=1.23$ $\mu = 1.55$ MB
<i>Pareto Parameters</i>	$a = 6.0$ $b=0.12$
Web-like	15%
<i>Distribution</i>	Lognormal w/max 843 MB
<i>Lognormal parameters</i>	$\sigma^2=1.978$ $\mu = 7.763$ B

At the beginning of every simulation, content is assigned randomly among nodes as follows: 5% of the *Peers* are assigned 6 unique whole-copy items and 10% of the *Benefactors* receive 1 whole-copy item. The *Peers* user class is assigned a bigger portion of the content space to stress the effect of churn on content availability. In addition, each node is assigned a uniformly random distributed number of fragments between zero and  $F$ ; where  $F$  is equal to  $3*m$  for nodes that were assigned a whole-copy item, and  $5*m$  otherwise. The value of  $m$  is the total number of fragments used by the coding mechanism. As a result, the content space of each experiment varies. The effects of this content initialization process are discussed further in Chapter 6.0.

When a node is in the online state, it generates content request at exponentially distributed intervals. The average time between these requests is 5 minutes for all nodes. The items requested are chosen uniformly at random and individual nodes do not request the same item twice during their lifetime.

Table 11 summarizes the parameters used for the content component of the content availability framework.

**Table 11. Content Parameters**

Parameter	Level(s)/Distributions	Description
<i>Resource space (M)</i>	400, 500, 750 and 1,500	Number of unique items
<i>File types</i>	$T(0.6, 0.1, 0.3)$	Percentage of Audio, Video, and Web-like files
<i>Whole-copy items</i>	<i>Peers</i> : 95% 0, rest 6 <i>Benefactors</i> : 90% 0, rest 1	Initial number of whole-copy items assigned to each node
<i>Fragments</i>	w/whole-copy: Uniform( $3*m$ ) w/o whole-copy: Uniform( $5*m$ )	Initial number of fragments assigned to each node
<i>Query distribution</i>	Uniform	Selection of items for downloads

### 3.3.5 Incentives and Redundancy Component

#### 3.3.5.1 Redundancy

The redundancy subcomponent includes four elements: redundancy method, transmission bandwidth, maintenance epoch policy and target file availability. The first element, redundancy method, defines whether the system employs replication redundancy or the hybrid redundancy method introduced in this dissertation; that is, proactive redundancy. The second element, transmission bandwidth, defines the maximum amount of bandwidth for a single data transfer between two nodes. The third element, maintenance epoch policy, indicates the algorithmic variant employed to adjust the interval between redundancy repair decisions taken by *Index* nodes. Finally, the fourth element, target number of fragments, defines the redundancy metric<sup>22</sup> (i.e., MIN\_SEG and TARGET) that controls how the maintenance epoch is adapted, and whether redundancy repair messages are to be generated.

---

<sup>22</sup> These factors are presented in Section 4.5.3.

The proactive redundancy method utilized in the experimental portion of this work has three parameter values: the number of blocks a file is originally split ( $n$ ), the storage overhead used by the mechanism ( $S$ ) and the number of replicas that the hybrid mechanism will generate for each unique fragment ( $r$ ). In the case of replication redundancy, which is also tested experimentally for comparison purposes, only the storage overhead parameter ( $S$ ) is used.

**Table 12. Redundancy Mechanism Parameters**

Factors	Levels
<i>Proactive Redundancy</i>	$n = \{6, 8\}$ $S = \{2, 3, 4\}$ $r = \{2\}$
<i>Replication Redundancy</i>	$S = 5$
<i>Transmission bandwidth</i>	80, 125 [kBps]
<i>Maintenance Epoch Policy</i>	Static, Adaptive<, Adaptive≤, Smooth $\alpha=0.4$ and Smooth $\alpha=0.8$
<i>Target number of fragments</i> ( <i>MIN_SEG</i> , <i>TARGET</i> )	0.5, 0.6, 0.7, 0.8, 0.9

### 3.3.5.2 Incentives-Based Mechanism

The incentive-based mechanism presented in this dissertation employs two variants of a sigmoid function. The first one is called TB-Utility and the second one is called RP-Cost. Table 13 presents the factors and levels employed for these two equations (described in Section 5.2.5) during the evaluation of the redundancy-system presented here.

**Table 13. Incentive-Based Mechanism Parameters**

Function	Parameters' Levels	
	<i>Shape</i>	<i>Shift/Bias</i>
<i>RP-Cost</i>	3.5, 5.5	1.25, 1.85
<i>TB-Utility</i>	1.25, 2.25, 3.25	0.05

## 4.0 REDUNDANCY

The most popular applications of P2P technology are content distribution and file-sharing. In both cases, content durability is limited by the fading popularity of items and the diverse failure content error patterns present in this type of networks; namely intermittent peer connectivity disk failures and network errors. Content redundancy can improve reliability of these systems, but given the large scale and high levels of churn in P2P networks, redundancy needs to be maintained in a timely manner. In this chapter, we analyze different aspects of redundancy for P2P networks and present an automated redundancy repair mechanism for P2P networks under churn.

Redundancy is the provision of excess resources to improve the reliability of a system. In the context of P2P networks, redundancy can be used to improve both the availability and the durability of content by distributing whole or code-based copies of a file among the network's participants.

In this chapter, we address three questions with regard to redundancy. First, how can we determine the proper level of redundancy to guarantee a desired file availability level? Second, what redundancy schemes perform better in P2P environments in terms of their repair cost? And third, how should we automate the maintenance of such redundancy? The term repair cost describes the amount of information exchanged between peers to restore the redundancy state information lost due to content errors. We tackle the questions above by building an analytical

framework that demonstrates the cost and performance advantages of our proposed redundancy scheme versus other schemes; in particular, erasure coding and network coding.

Our analytical formulation is organized in two stages. First, we describe file availability for a system where data is encoded and stored using a  $k$ -out-of- $N$  redundancy scheme with replication. Second, we analyze the redundancy repair cost for different constructions of the  $k$ -out-of- $N$  redundancy scheme with replication. The file availability analysis portrays the static properties of the redundancy system (i.e., resiliency) and the repair cost formulation focuses on how to maintain these properties in the long term (i.e., maintainability).

**Table 14. Redundancy Notation: Redundancy Scheme(s) Parameters**

Parameter	Description
$k$	<i>Reception efficiency.</i> Minimum number of unique fragments needed to reconstruct a file
$N$	Total number of unique fragments generated.
$S$	<i>Coding Gain.</i> $S=N/k$
$M$	Total number of fragments generated (including replicas)
$d$	<i>Repair degree.</i> Nodes needed to reconstruct a fragment using network coding redundancy. $k \leq d < N$
$r$	Replication gain. Number of per-fragment replicas maintained. $r \geq 1$
$R_i$	Replication gain for fragment type $i$ . $R_i \geq 1$

**Table 15. Redundancy Notation: File Availability**

Parameter	Description
$A_{\mathcal{F}}$	File Availability (for a $k$ -out-of- $N$ redundancy system)
$A_f$	Fragment availability (for a $l$ -out-of- $R_i$ redundancy system)
$\Phi$	Target file availability (i.e., $A_{\mathcal{F}} \geq \Phi$ )
$\phi$	Target fragment type availability (i.e., $A_f(R_i) \geq \phi$ )
$f_{i,j}$	Fragment $i$ , replica $j$ . $i = 1, \dots, N; j = 1, \dots, R_i$
$x_{i,j}^m$	Fragment $i$ , replica $j$ stored at node $m$ . $m = 1, \dots, M$
$q^m$	Availability of node $m$ . $0 < q^m < 1.0$

**Table 16. Redundancy Notation: Repair Cost**

Parameter	Description
$a$	Availability. Can be used in reference to nodes, disks or other system component (e.g., $a_n$ = node availability, $a_f$ = fragment-replica availability)
$\delta$	Maintenance epoch. Time interval at which repairs are performed
$g_n$	Distribution of nodes session lengths, $\mu_n$ =mean
$g_d$	Distribution of disk lifetimes, $\mu_d$ =mean
$\alpha$	Fragment size
$\beta$	<i>Repair unit</i> . Data exchanged between two nodes during a repair
$p_e^r$	Probability of error-free disk read operation
$p_e^t$	Probability of error-free repair unit transmission
$b_e^r$	Non-recoverable disk error rate
$b_e^t$	Transmission bit error rate
$L$	<i>Redundancy loss</i> . Number of nodes that left the system
$\Omega$	<i>Repair cost</i> . Total amount of information transferred during a maintenance epoch
$T$	Long-run redundancy repair interval (10 hrs)
$\tau$	Time needed to transfer a repair unit
$F$	File size
$B$	Repair bandwidth. Average repair bandwidth between a pair of nodes.

#### 4.1 REDUNDANCY SCHEMES

In our analysis we consider three different code-based redundancy schemes. Maximum Distance Separable (MDS) erasure coding, exact minimum bandwidth regenerating (exact-MBR) network coding and our proposed scheme, named Proactive Replication (PR). The basic structure of all these methods is a  $k$ -out-of- $N$  redundancy scheme.

In the MDS scheme, the file content is encoded to generate a set of  $N$  unique fragments, each of which is stored at different nodes. The encoding scheme is such that any  $k$ -out-of- $N$



fragments are sufficient to reconstruct the original file. In such scheme, a file of size  $F$  results in each node storing a fragment of size  $\alpha=F/k$  bytes. The regeneration of a lost fragment typically requires regenerating the original file, resulting in a repair overhead cost of  $\alpha*k$  bytes.

Similarly to MDS, exact-MBR uses a  $k$ -out-of- $N$  code. However, exact-MBR differs from MDS in the amount of information stored at each node and in the number of nodes required to reconstruct a single lost fragment. We base our analysis in the exact-MBR construction presented by Rashmi in [79]. In this scheme, nodes store fragments of size  $\alpha=2*d*F/(k*(2d-k+1))$  bytes, where  $k \leq d < N$  represents the repair degree. To reconstruct a single fragment, information from  $d$  distinct nodes is required, resulting in  $d$  different transfers of size  $\beta=\alpha/d$  bytes from each node.

Our proposed redundancy scheme, PR, also uses  $k$ -out-of- $N$  erasure codes. Our scheme, however, produces  $r$  replicas of each fragment. The replicated fragments are then stored at  $N*r$  different nodes. This strategy obviates the need to recreate the original file to repair a single lost fragment. Thus, the cost for such repairs is only  $F/k$  bytes. To calculate file availability for this new redundancy scheme, we analyze two different scenarios. First, the case when node availability is heterogeneous and second, for homogeneous node availability.

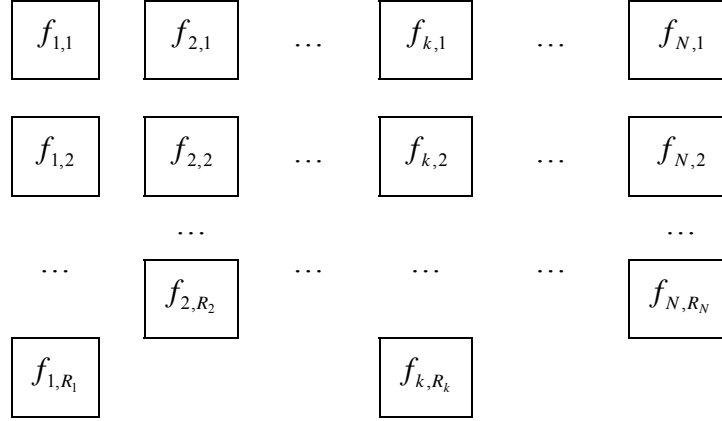
## 4.2 FILE AVAILABILITY

We assume the existence of a community of nodes that cooperatively store information. This community is formed by a fully connected set of nodes  $M=\{n^1, n^2, \dots, n^M\}$  that are available with probabilities  $P=\{q^1, q^2, \dots, q^M\}$  respectively. Data is stored in the system using a  $k$ -out-of- $N$  redundancy scheme with replication. The structure of this redundancy scheme is as follows. For a file of size  $F$ , we define the set  $F = \{f_1, f_2, \dots, f_N\}$  as a collection of  $N$  unique fragments of equal

size (e.g.,  $F/k$ ) such that any subset of  $k$  elements from  $F$  suffices to reconstruct the original file.

Let  $G$  be the set of all fragments contained in the system, which is a superset of replicas of fragments in  $F$ :

$$G = \left\{ \bigcup_{f_i \in F} \bigcup_{j=1}^{R_j} f_{i,j} \right\} \quad (4.1)$$



**Figure 17. Redundancy Data Structure**

Figure 17 shows a graphical representation of set  $G$ . There is a total of  $N$  unique fragment types and each fragment type  $i$  ( $f_i$ ) is replicated a total of  $R_i$  times, where  $R_i$ s are not necessarily the same. When  $R_i=1 \ \forall i$ , and node's availability is homogeneous, file availability can be obtained using the well-known formulation for a  $k$ -out-of- $N$  redundancy scheme, the binomial distribution. For our proposed redundancy scheme, PR, there are not analytical models available (to the best of our knowledge). Consequently, we develop our own to analyze the file availability properties of PR.

#### 4.2.1 Heterogeneous node availabilities

The availability of a file stored using the data structure described above (for  $R_i > 1$ ) is given by:

$$A_{\mathfrak{F}} = \Pr[\text{at least } k \text{ fragment types are available}] \quad (4.2)$$

This is a  $k$ -out-of- $N$  redundancy scheme, where the availability of each fragment type is a 1-out-of- $R_i$  redundancy scheme. For this redundancy structure we want to solve the following optimization problem:

$$\text{minimize:} \quad \text{Cost} = \sum_{i=1}^N C(R_i) \quad (4.3)$$

$$\text{subject to:} \quad A_{\mathcal{F}} \geq \Phi \quad (4.4)$$

The availability of the file ( $A_{\mathcal{F}}$ ) for a set of nodes,  $\mathcal{M}$ , with heterogeneous availabilities has been studied extensively and does not have a closed form solution. The methods reviewed by Kuo and Zuo in [80] for evaluating the reliability of a  $k$ -out-of- $N$  system are enumerative in nature; thus, using them to solve our optimization problem is computationally expensive. Alternatively, we adopt a decomposition strategy to obtain a feasible closed form solution for our optimization problem. By doing so, we obtain an alternative formulation that can be easily implemented in an iterative algorithm.

Let  $\kappa_i$  denote the probability that at least one fragment type  $i$  ( $f_i$ ) is available:

$$\begin{aligned} \kappa_i &= \Pr[\text{at least one fragment type } i \text{ is available}] \\ &= 1 - \Pr[\text{no fragment type } i \text{ is available}] \\ &= 1 - \prod_{j=1}^{R_i} (1 - q^m x_{i,j}^m) \end{aligned} \quad (4.5)$$

where,  $x_{i,j}^m$  is an indicator function such that

$$x_{i,j}^m = \begin{cases} 1 & \text{if } f_{i,j} \text{ is hosted in node } m \\ 0 & \text{otherwise} \end{cases} \quad (4.6)$$

Let  $\phi$  be a target fragment availability that we want to guarantee for all fragment types  $f_i$  such that  $\kappa_i \geq \phi$ . Now, using (4.5) and (4.6) in (4.4) we can obtain a closed form solution for the availability of the file:

$$A_{\mathcal{F}} = \sum_{i=k}^N \binom{N}{i} \cdot \phi^i \cdot (1-\phi)^{N-i} \quad (4.7)$$

Since equation (4.7) has the form of the binomial distribution, we can use the normal approximation to the binomial distribution to derive an optimal value (i.e., minimum) for  $\phi$  that would satisfy the original constraint of our optimization problem (i.e.,  $A_{\mathcal{F}} \geq \Phi$ ):

$$\phi = \left( \frac{\sigma_{\varepsilon} \sqrt{S/k} + \sqrt{\sigma_{\varepsilon}^2 (S/k) + 4(S + \sigma_{\varepsilon} \sqrt{S/k})}}{2(S + \sqrt{S/k})} \right)^2 \quad (4.8)$$

Where  $S$  is the coding gain of the MDS erasure coding scheme (i.e.  $N/k$ ),  $k$  is the reception efficiency and  $\sigma_{\varepsilon}$  is the number of standard deviations for the required level of file availability. For example, for a target file availability of two nines (0.99),  $k=8$  and  $N=12$ ,  $\phi=0.8142$ .

Now, our initial optimization problem can be rewritten as follows:

$$\text{minimize:} \quad \text{Cost} = \sum_{i=1}^N C(R_i) \quad (4.3)$$

$$\text{subject to:} \quad A_f(R_i) \geq \phi \quad \forall i \quad (4.9)$$

What we have accomplished is to reduce the complexity of our file availability problem. Now, we have to guarantee the availability of each fragment type independently. Furthermore, if we consider the cost function to be a non-decreasing concave function of  $R_i$  and the availability of the set of nodes storing the same fragment type, the original optimization problem is reduced to a set of  $N$  independent optimization problems; one for each  $f_i$  plus an additional set of constraints to guarantee that each node stores at most one fragment (4.12) and that the total number of fragment does not surpass the total number of nodes (4.13):

$$\text{minimize:} \quad C(R_i) = d * \sum_{j=1}^{R_i} (1 - q^m \chi_{i,j}^m) \quad \forall i \quad (4.10)$$

$$\text{subject to: } A_f(R_i) = 1 - C(R_i) = d * \sum_{j=1}^{R_i} (1 - q^m \chi_{i,j}^m) \geq \phi \quad \forall i \quad (4.11)$$

$$\sum_{m=1}^M \chi_{i,j}^m \leq 1 \quad \forall i, j \quad (4.12)$$

$$\sum_{i=1}^N R_i \leq M \quad (4.13)$$

In addition, if the availabilities ( $q^m$ ) of all the nodes is assumed to be i.i.d. and equal to  $a$ , equation (4.11) can be used to obtain a closed form solution for  $R_i$ :

$$\begin{aligned} A_f(R_i) &= 1 - \prod_{j=1}^{R_i} (1 - a) \geq \phi \\ &= 1 - (1 - a)^{R_i} \geq \phi \\ \text{thus} \\ R_i &\geq \frac{\log(1 - \phi)}{\log(1 - a)} \end{aligned} \quad (4.14)$$

Algorithm 1 presents the pseudo-code for finding a feasible solution for the optimization problem defined by equations (4.10) through (4.13). Notice that the fragment availability is calculated using a recursive formula ( $A_f(R_i+1) = A_f(R_i) * (1 - q) + q$ ) to minimize the processing cost of the algorithm.

In order to minimize cost, nodes have to be assigned to a fragment type subject to:

$$\begin{aligned} \min \Delta \text{cost} &= C(R_i + 1) - C(R_i) \\ &= d * \sum_{j=1}^{R_i+1} (1 - q^m \chi_{i,j}^m) - d * \sum_{j=1}^{R_i} (1 - q^m \chi_{i,j}^m) \\ &= d * (1 - q^m \chi_{i,R_i+1}^m) \end{aligned} \quad (4.15)$$

and

$$\begin{aligned} \max \Delta A_f &= A_f(R_i+1) - A_f(R_i) \\ &= \left[ 1 - \left( \prod_{j=1}^{R_i} (1 - q^m \chi_{i,j}^m) \right) \cdot (1 - q^m \chi_{i,R_i+1}^m) \right] - \left[ 1 - \prod_{j=1}^{R_i} (1 - q^m \chi_{i,j}^m) \right] \end{aligned}$$

$$\begin{aligned}
&= \left( \prod_{j=1}^{R_i} (1 - q^m \chi_{i,j}^m) \right) \cdot (1 - (1 - q^m \chi_{i,j+1}^m)) \\
&= (\Psi(R_i)) \cdot q^m \chi_{i,j+1}^m
\end{aligned} \tag{4.16}$$

**Procedure optimizeRedundancy()**

*Purpose:* Find a mapping of fragment replicates to nodes

{*N*: total number of unique fragments}

{*M*: set of nodes with their availability probabilities  $q^m$ }

{MIN\_AVA: minimum availability needed for each fragment type }

```

1:  for fi = 1 to N do
2:    node = select next node
3:    assign node to fi
4:    Afi = availability node
5:  end for
6:  for fi = 1 to N do
7:    while true
8:      if Afi < MIN_AVA then
9:        node = select next node
10:       assign node to fi
11:       p=availability of node
12:       Afi = Afi*(1-q)+q
13:     else
14:       break
15:     end if
16:   end while
17: end for

```

**Algorithm 1. Optimization of Redundancy Resources**

Equations (4.15) and (4.16) translate into a simple, and intuitive, node selection strategy: use the nodes with the greatest node availability first. For a given set of node availabilities, Algorithm 1 can determine quickly if a target level of file availability is feasible. We implemented our algorithm in Matlab and by synthetically generating a random set of node availabilities, we can easily determine the minimum node-set required to achieved a required level of file availability. For example, Figure 18 shows some sample results assuming uniformly distributed node availability between 0.15 and 0.45. The target file availability (of 0.99) is not feasible until the node set is larger than 63 nodes.

```

>> p=unifrnd(.15,.45,50,1);
>> [sol p q]=optimizeRedundancy(p);
Using provided node availability vector...
    Target file availability:0.9900
    Target fragment (type) availability = 0.8142
.....
    Unfeasible solution!! ran out of nodes
>> p=unifrnd(.15,.45,63,1);
>> [sol p q]=optimizeRedundancy(p);
Using provided node availability vector...
    Target file availability:0.9900
    Target fragment (type) availability = 0.8142
.....
DONE, Cost=42.5836
Solution:
    0.44    0.44    0.44    0.42    0.41    0.40    0.40 ...    0.39    0.39    0.38
    0.38    0.38    0.37    0.35    0.34    0.31    0.30 ...    0.25    0.23    0.18
    0.38    0.37    0.36    0.35    0.33    0.31    0.30 ...    0.24    0.21    0.18
    0.38    0.37    0.35    0.34    0.32    0.30    0.29 ...    0.24    0.20    0.17
    0        0        0        0        0        0.30    0.29 ...    0.24    0.19    0.17
    0        0        0        0        0        0        0 ...    0.23    0.19    0.17
    0        0        0        0        0        0        0 ...    0        0.19    0.17
    0        0        0        0        0        0        0 ...    0        0        0.16

```

Figure 18. Algorithm 1's Sample Output

#### 4.2.2 Homogeneous node availabilities

File availability has been used to define the fundamental relationships between the parameters of a redundancy scheme and the reliability requirements of a system. Assuming i.i.d. node availabilities, file availability is traditionally presented as:

$$A_{\mathfrak{F}} = \sum_{i=k}^N \binom{N}{i} \cdot a^i \cdot (1-a)^{N-i}, N=S*k \quad (4.17)$$

where:

*File availability,  $A_{\mathfrak{F}}$ .* Measures the reliability of the redundancy scheme. It is frequently expressed as one or more nines of availability (e.g.,  $A_{\mathfrak{F}} \geq 0.99$ ).

*Node availability,  $a$ .* Captures the unreliable nature of the components of the system.

*Reception Efficiency,  $k$ .* Denotes the number of fragments needed to reconstruct the original data

*Coding gain,  $S$ .* Is the fraction of the total number of fragments in the system over the number of fragments needed for reconstruction. Assuming all fragment to be unique,  $S = N/k$ .

For the construction of a redundancy system,  $A_{\mathcal{S}}$  and  $a$  are usually given.  $A_{\mathcal{S}}$  in the form of a system requirement and  $a$  as an environmental condition. The optimal values for  $S$  and  $k$  can be engineered to achieve different design objectives. For example, MDS provide an optimal tradeoff between reliability and storage space [79]. Given the tuple  $a$ ,  $k$  and  $A_{\mathcal{S}}$ , an optimal value (i.e., minimum) for  $S$  can be obtained (either numerically or analytically [10]) and similarly, given a value of  $a$ ,  $S$ , and  $A_{\mathcal{S}}$  an optimal value for  $k$  can be determined [54].

In a system without maintenance, the file availability formulation is used to define both reliability and cost for redundancy scheme, but for a system with maintenance, file availability only determines the minimum parameter settings to achieve a required level of reliability. Cost on the other hand, has to be engineered using additional metrics and guidelines. In particular, we need to define a fragment availability model capable to capture the redundancy repair costs for our redundancy scheme as well as for other redundancy schemes. In the next two sections we introduce the cost metric traditionally used to compare different redundancy schemes and then we describe our fragment availability model.

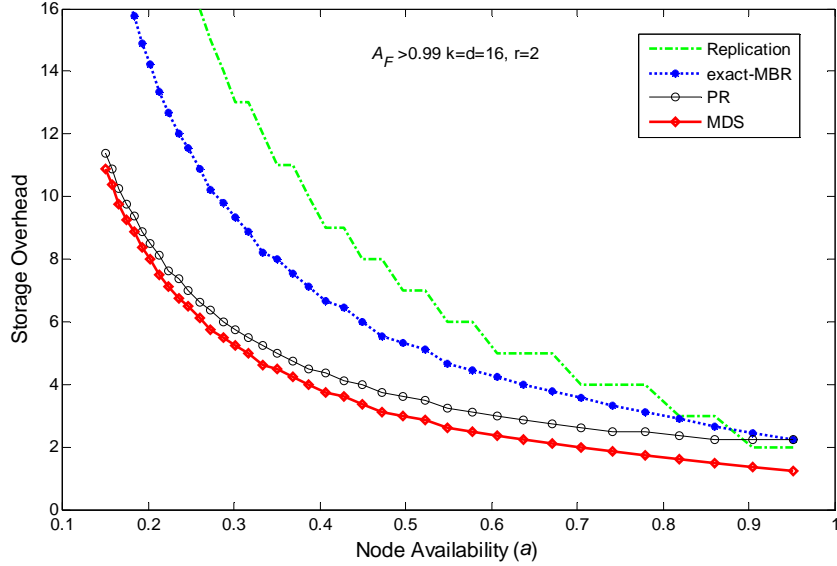
#### 4.2.2.1 Storage Overhead

In a  $k$ -out-of- $N$  redundancy system for data storage, cost is measured by the ratio of the total amount of storage used by the redundancy scheme and the storage space occupied by the original data. In a P2P environment though, repair cost is more important because access bandwidth is



scarcer and more expensive than storage space [10]. Nonetheless, storage overhead remains a fundamental performance measure of the system that needs to be taken into consideration.

The next figure presents the minimum storage overhead for four redundancy schemes: MDS erasure coding, exact-MBR network coding, our proposed hybrid redundancy method, PR, with  $r=2$ , and Replication.



**Figure 19. Minimum Storage Overhead**

These results illustrate that code-based redundancy methods use less storage overhead than replication redundancy, especially at low average node availabilities. For exact-MBR network coding, the coding gain parameter,  $S$ , uses the same value as MDS, but the storage overhead at each node is higher. For a file of size  $F$ , the amount of information stored by each node is  $F/k$  in MDS, while for exact-MBR each node stores  $2*F*d/k(2*d-k+1)$  [81]. For PR, each node stores  $F/k$  bytes, but the total amount of storage used is larger than MDS' storage. The PR scheme stores  $r$  copies of each unique fragment versus the single copy stored by MDS redundancy. However, for low node availabilities the additional storage overhead of PR is minimum compared with MDS. In terms of storage overhead savings, MDS is the most efficient method, but its redundancy

maintenance cost is high. Exact-MBR network coding redundancy and PR redundancy utilize additional storage in order to minimize repair maintenance cost, which we analyze in the following sections.

In addition to storage overhead, the results presented in Figure 19 can be used to determine another important system metric. The total number of nodes needed to store a file. We refer to this property as node-set. For replication redundancy, this parameter's value is the same than storage overhead. For exact-MBR and MDS, the node-set is the product of the coding gain times the reception efficiency,  $N = S*k$ . For PR, the node-set is the product of the coding gain, the reception efficiency and the replication factor,  $m = N*r = S*k*r$ . Consequently, in terms of the number of nodes needed to store a file, replication redundancy requires the least number of nodes, followed by MDS and exact-MBR, and finally, by PR. For MDS and exact-MDS, small node availabilities require larger node-set values which translate into an added coding/decoding complexity. For PR, this increment in coding/decoding complexity also applies, but at a lower degree. In our opinion, large node-set values should be avoided not only for their added coding/decoding complexity, but for their practical implications. As the number of nodes needed to store a single file increases, managing the node-set becomes more complex. For example, for  $\alpha=0.3$ , MDS and exact-MBR require a minimum coding gain of  $S=5.25$ , which implies a redundancy-set with at least 84 nodes.

#### 4.2.2.2 Fragment Availability Model

For the calculation of file availability using a  $k$ -out-of- $N$  redundancy system most of the research literature equates the availability of individual fragments with the availability of the host holding them. In our case, we use a broader model for the availability of these fragments. This model reflects the probability of hardware errors in addition to the traditional host availability component. Furthermore, our formulation derives host availability based on an analytical model of the node's

session length probability distribution, rather than system traces [10, 56, 58]. This approach has two advantages. First, results can be obtained for a wider range of network conditions. That is, we are not limited to the availability of traces. Second, the relationship between the system maintenance epochs and the average node availability is expressed explicitly, which allow us to define a repair frequency for a redundancy maintenance algorithm that would maximize the performance of the system.

In our model,  $a_f$  denotes the availability of a single fragment. This metric is the product of two probabilities. First, the probability that the node storing the fragment is available. Second, the probability that if that node is available, no hardware errors will prevent access to the fragment. Other research initiatives consider host availability [56] or disk failures [9] as the only sources of content errors. Conversely, our model is broader. It can be customized to model environments where host availability is dominant as well as environments where hardware and communication errors are important (e.g., wireless networks). We consider the probability of disk failures and the probability of data corruption as the main sources of hardware errors. In turn, we use disk read errors and transmission errors to construct the data corruption component of our model. For a given system maintenance epoch  $\delta$ , fragment availability is given by:

$$a_f^\delta = \Pr[\text{node is available} \mid \delta] * \Pr[\text{no hardware errors} \mid \delta] \quad (4.18)$$

The first term in the right side of equation (4.18) measures the residual lifetime probability of the node storing the fragment given that the node has survived one maintenance epoch. In other words, this probability measures if the fragment can be used during the next maintenance interval to reconstruct the original file or to perform repairs. We use  $a_n$  to denote this probability. The probability that a fragment is available after a maintenance epoch, can be calculated using the expression derived in [9]:

$$a_n^\delta = \Pr[ x - \delta \mid x > \delta ] = \int_{\delta}^{\infty} \frac{x \cdot f_n(x)}{\mu_n} \cdot \frac{x - \delta}{x} dx = \frac{1}{\mu_n} \int_{\delta}^{\infty} f_n(x) \cdot (x - \delta) dx \quad (4.19)$$

where the term  $(x-\delta)/x$  reflects the probability of storing a fragment early enough in a node's session so that it is still available after the next maintenance epoch,  $g_n$  is the node's session length probability distribution, and  $\mu_n$  is the expected value of this distribution (i.e., mean value).

The second term in the right side of equation (4.18) has three factors: the availability of disks storing the fragment, the probability of non-recoverable read errors, and the probability of corrupted data transmissions. We use  $a_e$  to denote the resulting probability,  $a_d^\delta$  denotes disk availability,  $p_r$  denotes uncorrupted read operations and  $p_t$  denotes successful data transmissions.

$$\Pr[\text{no hardware errors}|\delta] = \Pr[\text{disk available}|\delta] \cdot \Pr[\text{no read errors}] \cdot \Pr[\text{no tx errors}]$$

$$a_e^\delta = a_d^\delta \cdot p_r \cdot p_t \quad (4.20)$$

Substituting (4.19) and (4.20) in (4.18) we have

$$\begin{aligned} a_f^\delta &= a_n^\delta * a_e^\delta \\ &= a_n^\delta * a_d^\delta \cdot p_r \cdot p_t \end{aligned} \quad (4.21)$$

The disk availability component,  $a_d^\delta$ , in equation (4.20) can be obtained using a similar approach to the one used for  $a_n^\delta$ . That is, we can use equation (4.19) by replacing  $g_n$  and  $\mu_n$  with the disk's lifetime probability distribution and its respective mean value, which we denote by  $g_d$  and  $\mu_d$  respectively.

$$a_d^\delta = \int_{\delta}^{\infty} \frac{x \cdot f_d(x)}{\mu_d} \cdot \frac{x - \delta}{x} dx = \frac{1}{\mu_d} \int_{\delta}^{\infty} f_d(x) \cdot (x - \delta) dx \quad (4.22)$$

To derive the probabilities of uncorrupted disk read operations and successful data transmissions we use a simple binomial model and assume that bits fail independently. Let  $\alpha$

denote the size of a fragment,  $b_r$  denote the disk's non-recoverable read error rate, and  $b_t$  denote the transmission bit error rate. Then,  $p_r$  and  $p_t$  can be obtained as follows:

$$p_r = (1 - b_r)^\alpha \quad (4.23)$$

$$p_t = (1 - b_t)^\beta \quad (4.24)$$

We assume that to retrieve the data needed for repairs, nodes must read the whole fragment. The size of a fragment is denoted by  $\alpha$ , and  $\beta$  denotes the amount of information transferred to another node during a repair. In MDS and PR redundancy  $\beta = \alpha$ , and in exact-MBR redundancy  $\beta = \alpha/d$ . Placing (4.23) and (4.24) in (4.21) we get the following expression for the availability of a fragment.

$$a_f^\delta = \Pr[\text{node is available} \mid \delta] * \Pr[\text{no hardware errors} \mid \delta]$$

$$a_f^\delta = a_n^\delta * a_d^\delta * (1 - b_r)^\alpha * (1 - b_t)^\beta \quad (4.25)$$

With this, our fragment availability model is complete. It comprises the system maintenance epoch ( $\delta$ ) the availability of nodes ( $a_n^\delta$ ) and disk ( $a_d^\delta$ ), and the probability of disk-read errors ( $(1 - b_r)^\alpha$ ) and finally the probability of failed transmissions ( $(1 - b_t)^\beta$ ). With this model, we can formulate a repair cost metric for our proposed redundancy scheme and compare it against other redundancy schemes.

### 4.3 REDUNDANCY REPAIR

When nodes leave the system, the reliability of the redundancy method deteriorates. To avoid permanent data loss, the system needs to reconstruct any redundancy lost regularly. The amount of

data transferred between nodes to maintain the reliability of the system is referred as repair cost.

We use  $\Omega_{\mathfrak{F}}$  to denote this metric.

The two fundamental system components that determine  $\Omega_{\mathfrak{F}}$  are redundancy scheme, and maintenance epoch. The structure and parameters of the redundancy scheme determine the amount of information to be transmitted to reconstruct the redundancy lost. The system maintenance epoch determines the frequency and number of repairs to be performed during each maintenance interval. For longer maintenance epochs, a larger number of nodes leave the system.

We assume that at time  $t=0$  we have  $m$  hosts cooperatively storing a file's fragments. We also assume that nodes store single fragments and that their availability,  $a_f^\delta$ , is homogeneous and stationary. Let  $L_{\mathfrak{F}}^\delta$  denote the expected number of nodes lost during a maintenance interval (with length  $\delta$ ). The expected value of  $L_{\mathfrak{F}}^\delta$  can be obtained as the expected value of a binomial distribution with a total of  $m$  items:

$$\begin{aligned} E[L_{\mathfrak{F}}^\delta] &= \sum_{i=1}^m l_i \cdot \Pr[L_{\mathfrak{F}}^\delta = l_i] \\ &= \sum_{i=1}^m i \cdot \binom{m}{i} \cdot (a_f^\delta)^{m-i} \cdot (1 - a_f^\delta)^i = m \cdot (1 - a_f^\delta) \end{aligned} \tag{4.26}$$

The number of hosts storing the file (i.e.,  $m$ ) is determined by the parameters of the redundancy method. Notice that the expected value of the binomial distribution is obtained for  $1 - a_f^\delta$  because we want the number of offline nodes, rather than the number of nodes that remain active.

Once the number of fragments lost is known, we can calculate the cost of repairing them. We assume that all the redundancy lost during a maintenance epoch is fully repaired during the next maintenance interval. Let  $T$  denote a long period of time (e.g. 10 hours). We use the term

long-run repair cost to refer to the average redundancy repair cost of the system during  $T$ . Thus, repair cost is a recurrent process performed  $T/\delta$  times in average. Given  $T$  and  $\delta$ , the long-run maintenance cost of a  $k$ -out-of- $m$  system using MDS, exact-MBR and PR redundancy is given by:

$$\Omega_{\delta}^{\delta} = \begin{cases} (L_{\delta}^{\delta} + k) \cdot \beta_{EC} \cdot \frac{T}{\delta^2} & \text{MDS} \\ (L_{\delta}^{\delta} \cdot d) \cdot \beta_{MBR} \cdot \frac{T}{\delta^2} & \text{exact-MBR} \\ (L_{\delta}^{\delta} + \gamma \cdot k) \cdot \beta_{PR} \cdot \frac{T}{\delta^2} & \text{PR} \end{cases} \quad (4.27)$$

In all the expressions in equation (4.27), the last term correspond to the frequency of repairs (1/seconds) and the remaining terms correspond to their size (bytes). The cost contributions of these components increase in opposite directions. For short maintenance epochs, the frequency component is high and size is small<sup>23</sup>. Conversely, for long maintenance epochs the frequency component gets smaller, but the size component becomes larger. Thus, a fundamental problem is to determine if there is an optimal setting to balance the cost of these components for each redundancy scheme.

The first expression in equation (4.27) corresponds to MDS erasure coding redundancy when no nodes are available with a complete copy of the item. In this scenario, nodes stores  $\alpha = F/k$  bytes; where  $F$  is the size of the file. For repairs, nodes connect to  $k$  available nodes and transfer  $\beta = \alpha = F/k$  bytes from each node to reconstruct the original data and then transfer the  $L$  fragments lost to new nodes. The research literature refers to this scenario as the redundancy repair problem due to the potential high overhead cost [82].

The second expression in equation (4.27) corresponds to exact minimum-bandwidth regenerating (exact-MBR) network coding redundancy. Authors present this type of redundancy as

---

<sup>23</sup> However, for short maintenance epochs, the size efficiency of EC is worst (i.e., highest overhead).

an alternative to traditional MDS erasure coding [79, 82]. Exact-MBR codes can repair individual fragments without reconstructing the original file first, but they incur in an additional storage overhead in every node (compared with MDS). Some of these solutions are not practical because the repair degree is  $m - 1$ , like in [78]. However, the exact-MBR code presented by Rashmi, et.al., in [79] allows the selection of the repair degree,  $d$ , independently of the other parameters. For this exact-MBR code construction, the size of each fragment is  $\alpha = 2*d*F/k(2*d-k+1)$ . For repairs, nodes contact  $d$  available nodes and transfer  $\beta = \alpha/d$  bytes from each.

The third expression in equation (4.27) is another alternative to the redundancy repair problem of MDS erasure coding, a hybrid redundancy method that we name Proactive Replication, PR. Our approach differs from previous analysis of hybrid redundancy mechanisms in several aspects. First, we are presenting a full analytical formulation for the properties of this redundancy method (i.e., the file availability and repair cost equations); as opposed to using it as a simplified redundancy repair scenario where a complete file copy is assumed to be available all the time, like in [10] or [59]. Second, the structure of our hybrid redundancy method is different from other hybrid methods presented in the literature. Again, in [10] a complete copy of the file exists in parallel to the EC redundancy data. In [60] MDS is an auxiliary mechanism to replication; when not enough replicas are available, the system generates MDS blocks to reach a target file availability. In our approach, we use replication on top of MDS. The MDS component is used for reliability and replication is used to minimize repair cost. This mechanism can be conceptualized as either a new hybrid redundancy method, or as a proactive redundancy repair policy. In either case, the main feature of the mechanism is a significant reduction of the repair cost. In PR, the size of each fragment is  $\alpha = F/k$ . For repairs, if there is at least one fragment replica available, nodes contact a single node to transfer  $\beta = \alpha = F/k$  bytes. When there is no fragment replica available,



nodes contact  $k$  available nodes to reconstruct the original file first (i.e., regular MDS). The probability of using the regular MDS repair mechanism is captured by the parameter  $\gamma$ .

Given a redundancy loss  $L_i$  (i.e., number of nodes lost during a maintenance epoch), repairs are done using the traditional MDS mechanism when all  $r$  replicas of a fragment are lost simultaneously. This happens with a probability equal to the portion of  $r$ -nodes arrangements that can be made with  $L_i$  elements out of the total number of  $r$ -nodes arrangements that can be obtained with  $N \cdot (r-1)$  nodes. In other words, the probability of using EC repair goes from zero to one as the number of fragment lost goes from zero to  $N \cdot (r-1)$ . The top value for  $L_i$  comes from the maximum number of fragments that can be lost (i.e., all replicas) by the redundancy structure before only MDS repair is possible.

$$\gamma = \frac{\binom{L_i}{r}}{\binom{N \cdot (r-1)}{r}} \quad (4.28)$$

#### 4.3.1 Failed Repairs

The cost formulation presented in equation (4.27) assumes a perfect repair scenario. To capture more realistic scenarios where repairs can fail due to churn and hardware errors we enhanced our repair cost model with a retransmission factor  $R(\tau)$  which reflects the average number of transmission attempts needed before a successful repair.

$$\Omega_{\mathfrak{F}}^{\delta} = \begin{cases} \left[ L_{\mathfrak{F}}^{\delta} + k \right] \cdot R(\tau) \cdot \beta_{EC} \cdot \frac{T}{\delta^2}, & \tau = \frac{\beta_{MDS}}{B} & \text{MDS} \\ \left[ L_{\mathfrak{F}}^{\delta} \cdot d \right] \cdot R(\tau) \cdot \beta_{MBR} \cdot \frac{T}{\delta^2}, & \tau = \frac{\beta_{MBR}}{B} & \text{exact - MBR} \\ \left[ L_{\mathfrak{F}}^{\delta} + \gamma \cdot k \right] \cdot R(\tau) \cdot \beta_{PR} \cdot \frac{T}{\delta^2}, & \tau = \frac{\beta_{PR}}{B} & \text{PR} \end{cases} \quad (4.29)$$

where

$$\begin{aligned}
R(\tau) &= \Pr[\text{success repair} \mid \tau]^{-1} \\
&= (\Pr[\text{nodes available} \mid \tau] \cdot \Pr[\text{no hardware errors} \mid \tau])^{-1} \\
&= ((a_n^\tau)^2 * a_d^\tau \cdot p_r \cdot p_t)^{-1} \\
&= \rho^{-1}
\end{aligned} \tag{4.30}$$

and

$$B = \text{Repair bandwidth between a pair of nodes} \tag{4.31}$$

Individual repairs are successful with probability  $\rho$ . This probability depends on the probability that each pair of nodes exchanging data remains active long enough to complete the data transfer (i.e.,  $\tau$  seconds) and the probability of no hardware errors; which in turn depends on the amount of information stored and transferred (i.e.,  $\alpha$  and  $\beta$ ). To obtain the value of  $\tau$ , we assume a constant value for the amount of bandwidth committed to repairs. Let  $B$  denote this value, then  $\tau = \beta/B$ . The probability  $\rho$  is obtained using a similar approach to the formulation of  $a_f^\delta$ , with the exception that in this case, two nodes must remain active:

$$R(\tau)^{-1} = (a_n^\delta)^2 * a_d^\delta * (1 - b_r)^\alpha * (1 - b_t)^\beta \tag{4.32}$$

### 4.3.2 Repair Cost

The following figure presents the repair cost across a wide range of average node availabilities for MDS, exact-MBR, PR and the ideal version of MDS redundancy. The ideal version of MDS assumes that fragments can be repaired transferring only  $\beta$  bytes (i.e., without reconstructing the original file first). The parameters used for these calculations are listed in Table 17. For the expected value of the disk lifetime distribution we use the annual failure rate (AFR) reported in [83]. For the probability of bit read errors, we assume the disk specifications of a consumer-class SATA hard-drive [84] and for the transmission bit error rate, we assume the value reported in [85].

The procedure used to obtain these results is as follows. First, we obtain the minimum coding gain,  $S$ , needed to achieve a file availability of 0.99 for each redundancy scheme for different average fragment availabilities. Then, repair cost is calculated using equation (4.29), for each redundancy method using the  $S$  value obtained.

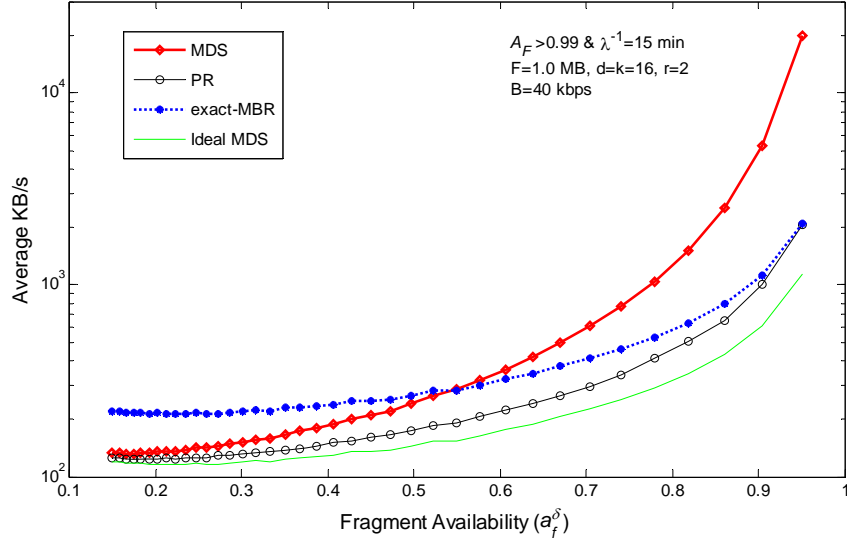


Figure 20. Redundancy Repair Cost

Table 17. Repair Cost Calculation Parameters

	File Availability		Repair Cost
	Free Parameter	Fixed Parameter(s)	$F=1.0MB$ $g_n \sim exp, \mu_n=15 \text{ minutes}$ $T=10\text{hrs}, B=40\text{kbps}$ $g_d \sim exp, \mu_d=100\text{k hours}$ $b_r=1.1 \times 10^{-14}$ and $b_i=1 \times 10^{-13}$
<i>MDS</i>	$S$	$k=16, \Phi=0.99$	
<i>Exact-MBR</i>	$S$	$k=d=16, \Phi=0.99$	
<i>PR</i>	$S$	$k=16, r=2, \Phi=0.99$	

The most important point to be highlighted from the results presented in Figure 20 is that the performance of our proposed redundancy method, PR, outperforms exact-MBR. In addition, all plots indicate that repair cost is an increasing function of average node availability. This suggests that the system should be tune up to work at low node availabilities (i.e., long maintenance epochs). However, this presents two drawbacks. First, the system requires higher coding gains, which implies a larger node-set (i.e., number of nodes) to store files. As the node-set increases, so

does the complexity of selecting and scheduling fragments for repairs. Second, our model simplifies the heterogeneous and dynamic nature of P2P networks. For longer maintenance epochs, (i.e., low node availability) the reliability of the system is more likely to be overwhelmed by churn and other dynamic processes. Consequently, we consider that a reasonable reliability vs cost performance compromise could be achieved at the higher average fragment availability values.

For the results in Figure 20, MDS performs better than exact-MBR redundancy for average node availabilities below 0.5. This result is consistent with the observations presented by Dimakis, et. al., in [78]. The authors indicate that as the network becomes less stable, the performance of network coding can be “very slightly worse” than erasure coding. Our availability model allow us to determine precisely at which point substituting MDS with other redundancy schemes is a viable solution. Furthermore, our model provides a solid and flexible framework to compare the performance of redundancy schemes using different combination of parameters. For example, we can easily determine the performance gain obtained when exact-MBR uses different repair degrees,  $k < d < N$ .

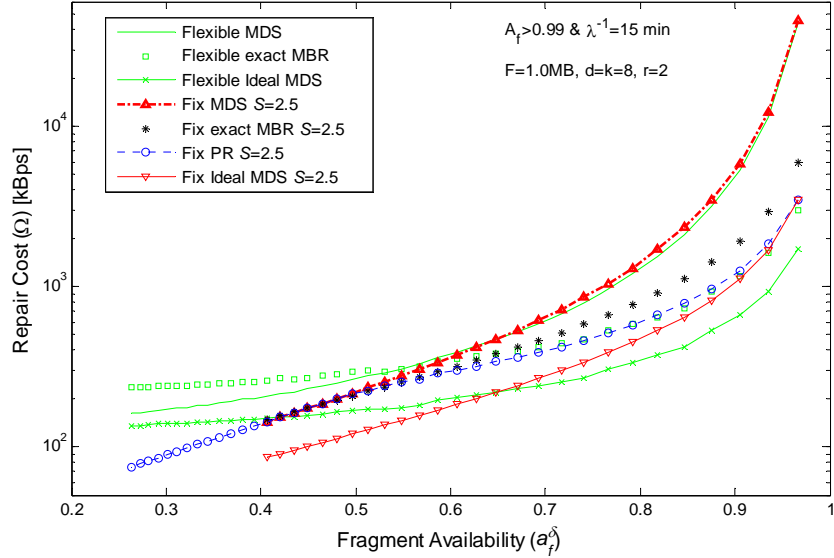
The overall structure of our proposed mechanism does not exclude the use of exact-MBR network coding instead of MDS for the coding component of the scheme. In that case, the repair cost for this alternative redundancy construction is:

$$\Omega_{\mathfrak{F}}^{\delta} = \begin{cases} \left[ L_{\mathfrak{F}}^{\delta} \cdot d \right] \cdot R(\tau_{\beta}) \cdot \beta_{MBR} \cdot \frac{T}{\delta^2}, & \tau_{\beta} = \frac{\beta_{MBR}}{BW_r} \\ \left[ L_{\mathfrak{F}}^{\delta} \right] \cdot R(\tau_{\alpha}) \cdot \alpha_{MBR} \cdot \frac{T}{\delta^2}, & \tau_{\alpha} = \frac{\alpha_{MBR}}{BW_r} = d \cdot \frac{\beta_{MBR}}{BW_r} \end{cases} \quad (4.33)$$

The first expression in equation (4.33) corresponds to regular exact-MBR repair (i.e., contact  $d$  nodes and transfer  $\beta$  bytes from each one) and the second expression correspond to our proposed proactive repair mechanism (i.e., contact a single node and download  $\alpha$  bytes from it).

PR is performed with probability  $1-\gamma$  (which is calculated using equation (4.28) and regular exact-MBR repair is performed with probability  $\gamma$ .

Coding gain is changed to obtain each of the points in the cost versus fragment availability performance space presented in Figure 21. We label these results as “flexible” in reference to the changes in coding gain. In a real deployment scenario, fragment availability can vary, but coding gain can not be adjusted dynamically; instead, it must be engineered to sustain the mean (or worst) fragment availability scenario.



**Figure 21. Repair Cost for Flexible and Fixed Coding Gain**

Figure 21 presents the cost versus fragment availability performance of MDS, ideal MDS, exact-MBR and PR redundancy. Coding gain remains constant to  $S=2.5$ . We label these results “fixed” in reference to their coding gain. Additional plots are included in this graph to illustrate the performance bounds for each scheme. In particular, we have included three flexible schemes: flexible MDS, flexible ideal MDS and flexible exact-MBR. For the fix redundancy schemes, all the performance values to the left of the points where the fix and flexible variants intersect do not satisfy the file availability requirement of 0.99. Thus, the cost performance of the fix schemes is

bounded by their flexible scheme variants. For example, the fix ideal MDS and flexible ideal MDS schemes intersect at an average fragment availability value of 0.65. For average fragment availabilities below this value, the fixed scheme consumes less repair bandwidth than its flexible counterpart, but the resulting file availability is no longer above 0.99.

Some key differences and similarities between the flexible and fixed schemes presented in Figure 21 are summarized in Table 18. In both sets, MDS and Ideal MDS constitute the upper and lower performance bounds of the system. When the repair cost performance of the fix set is better (i.e., lower) than its flexible counterpart, its file availability is below the required level. For the points where both sets achieve the required level of file availability, the repair cost of the fix set is higher than the flexible set. The repair cost performance of exact-MBR and PR redundancy is better than MDS across most of their feasible range, but it is inferior to the performance of Ideal MDS.

**Table 18. Redundancy Methods with Flexible vs Fix Coding Gain**

	Flexible Coding Gain	Fixed Coding Gain
<i>Feasible range</i>	full	limited
<i>File Availability, <math>A_{\mathcal{F}}</math></i>	0.99	Varies
<i>Highest Reliability</i>	high $a_f^\delta$	high $a_f^\delta$
<i>Lowest Cost</i>	low $a_f^\delta$	low $a_f^\delta$
<i>MDS preferred</i>	$a_f^\delta < 0.55$	n/a

Computing the repair cost ratio for exact-MBR and PR redundancy versus ideal MDS we obtain the following. Exact-MBR's performance is a constant factor below ideal MDS (i.e., higher cost), 1.7334 to be specific. This value comes mainly from the ratio of the  $\beta$  values for these two redundancy schemes. That is  $d/k(2d-k+1)$  vs  $1/k$ . For PR redundancy, this ratio varies for different average fragment availabilities. At the highest fragment availability, the repair cost of PR is

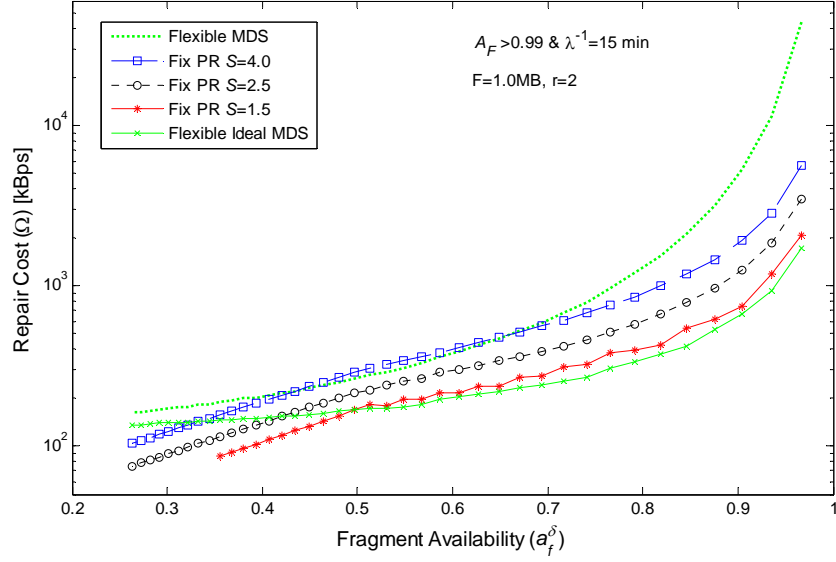
slightly higher (10%) than ideal MDS, and it increases for lower fragment availabilities, reaching the same overhead as exact-MBR at  $a_f^\delta=0.55$ .

The next set of figures illustrates the effects of changing the parameters of exact-MBR and PR redundancy.

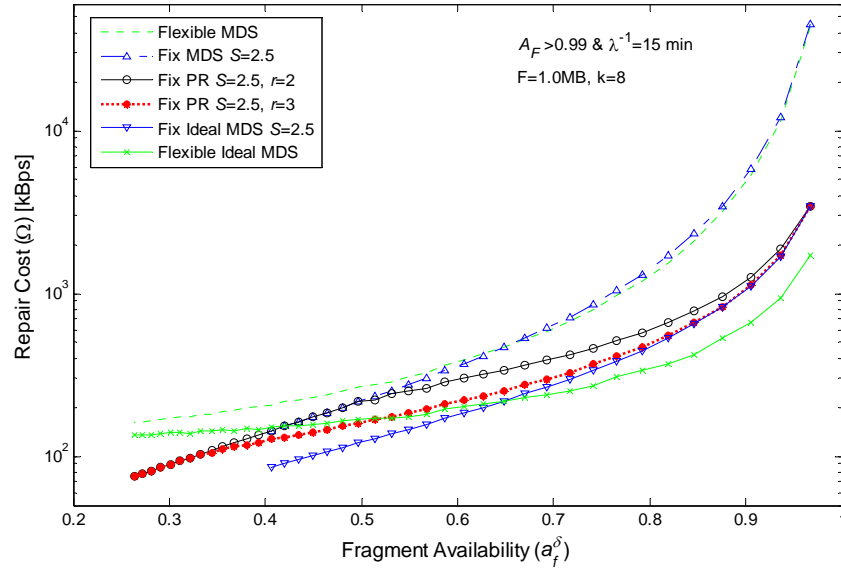
#### 4.3.2.1 Proactive Replication with Different $(S, r)$ Parameters

The next two figures illustrate the repair cost performance trend for PR when coding gain,  $S$ , and redundancy degree,  $r$ , are changed. Figure 22 presents results for different  $S$  and Figure 23 illustrate the performance trend using different  $r$ .

The results in Figure 22 indicate that the best repair cost is achieved when the redundancy scheme uses smaller coding gains. The only drawback is that the range of feasible average fragment availabilities is reduced (for a required level of file availability). For example, given  $A_{\mathfrak{F}} \geq 0.99$ ,  $S=1.5$  and  $k=8$ , the system achieves the required file availability level only for average fragment availabilities above 0.6. Larger coding gain values improve the reliability of the system and increase the feasible range of PR. For  $S=2.5$ ,  $a_f^\delta \geq 0.35$  and for  $S=4.0$ ,  $a_f^\delta \geq 0.21$  to achieve  $A_{\mathfrak{F}} \geq 0.99$ . However, there is an increase in the repair cost associated with increasing the coding gain parameter.



**Figure 22. Repair Cost for PR with different *Coding Gain* ( $S$ ) values**



**Figure 23. Repair Cost for PR with different *Replication Gain* ( $r$ ) values**

The results in Figure 23 indicate that larger replication gains,  $r$ , allow PR to perform closer to ideal MDS at the high end of the availability spectrum. This is a promising result because it suggests that additional storage used in the replication component of the PR mechanism has a net effect of reducing the repair bandwidth of PR over a wider range of fragment availabilities.



In summary, our analytical models for fragment availability and redundancy repair cost indicate that our proposed Proactive Replication (PR) redundancy scheme can reduce significantly the repair cost of a  $k$ -out-of- $N$  redundancy system in a P2P environment. The simple design of the proposed method can use either MDS erasure coding or exact-MBR network coding redundancy for its coding component. In addition, its replication component adds not only to the efficiency of the redundancy method, but to its flexibility. Replication gain can be adjusted dynamically to satisfy different file availability requirements without increasing repair cost. Our analytical results indicate that both repair cost and reliability increase with average fragment availability. Consequently, an automated redundancy maintenance system to improve content availability in P2P networks should avoid both ends of the fragment availability spectrum. Low fragment availability should be avoided due to reliability concerns, and high fragment availabilities should be avoided due to repair cost concerns. As an initial system engineering guideline, we propose that the system should be set up to operate at average node availabilities between 0.65 and 0.85. We believe that this is a reasonable compromise between cost and reliability. In addition, at lower node availabilities, the performance gain of PR versus MDS redundancy is less significant.

### 4.3.3 Maintenance Epochs

The results presented earlier are time-invariant. That is, each point in the performance space is valid for all systems regardless of their particular churn regimes. Tuning a particular system to achieve a particular point in the performance space requires mapping the system dynamics and fragments' availability. In practice, the only system parameter that can be adjusted to tune the system's operation for a target fragment availability metric is the system maintenance

epoch,  $\delta$ . In our fragment availability model, the average fragment availability and the system maintenance epoch can be derived from equation (4.25):

$$a_f^\delta = a_n^\delta * a_d^\delta * (1 - b_r)^\alpha * (1 - b_t)^\beta \quad (4.25)$$

In most P2P settings, the above probability is likely to be dominated by node availability. If we assume that this factor is exponentially distributed (with mean  $\lambda$ ), the relationship between a required average fragment availability value,  $a_i$ , and the system's maintenance epochs can be obtained as:

$$\delta = -\lambda * \log(a_i) \quad (4.34)$$

In fact, given a required level of file availability  $A_g \geq \phi$  and the redundancy parameters  $S$  and  $k$ , we can derive an optimal formulation for the redundancy maintenance of MDS using the normal approximation to the binomial distribution [10, 55]. Other researchers have used this method to obtain the optimal coding gain,  $S$ , but to the best of our knowledge the following formulation has not been presented earlier in the literature:

$$a = \left( \frac{\sigma_\epsilon \sqrt{S/k} + \sqrt{\sigma_\epsilon^2 (S/k) + 4(S + \sigma_\epsilon \sqrt{S/k})}}{2(S + \sqrt{S/k})} \right)^2 \quad (4.35)$$

Where  $\sigma_\epsilon$  is the number of standard deviations for the required level of file availability. For our content availability problem, the formulation above defines the minimum average node availability needed to achieve the required file availability. Using this formulation, we can determine the maximum maintenance epoch for a given parameter tuple  $(\phi, S, k)$ .

From the results presented earlier (see Figure 21~Figure 23 ), we know that repair cost increases with average fragment availability. Consequently, the formulation above also defines the average fragment availability,  $a_f^\delta = a_{\min}$ , that results in a minimum repair cost. Nonetheless, we

have also argued that at  $a_{min}$  the reliability of the system is the lowest. Thus, it should be avoided. Furthermore, for the PR redundancy mechanism we are proposing, the closer the system operates to  $a_{min}$ , the smaller performance gain we get with MDS. As a result,  $a_{min}$  must be considered as an absolute minimum, and the redundancy maintenance for the system should be tune up to operate at larger  $a_f^\delta$  values.

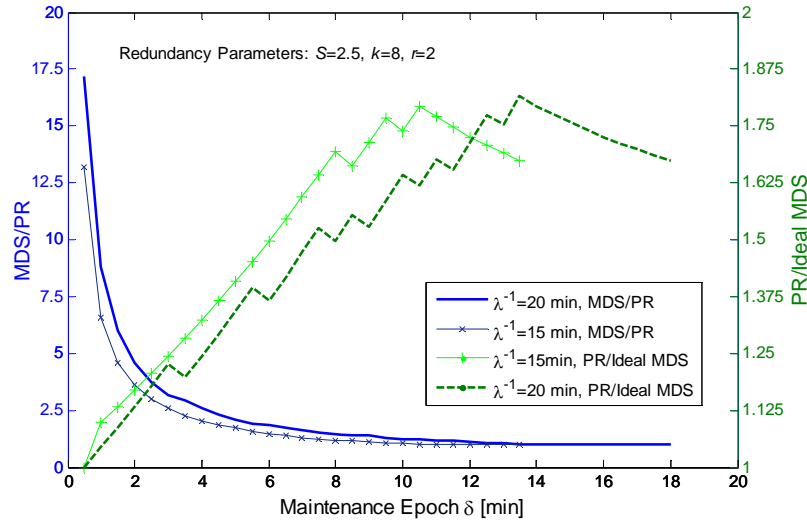
The results in Figure 23 provide additional insights regarding the cost vs performance tradeoffs of PR redundancy. The figure presents the redundancy repair cost ratio of PR redundancy with respect to MDS and Ideal MDS. The y-axis to the right is for the two plots that increase their values for longer maintenance epochs. We name this set the Ideal ratio. The scale to the left correspond to the other two plots, which exhibit an exponential decay. We name this set the MDS ratio.

The Ideal ratio is obtained as the ratio of the PR redundancy repair cost over the cost of ideal MDS. For the shorter maintenance epochs this ratio is close to the ideal value (i.e., 1.0) and for larger maintenance epochs it degrades until it reaches 1.8. At this point, the performance of the PR redundancy method is the same as regular MDS and any further improvement in this ratio are due to the convergence of the repair costs of MDS and Ideal MDS.

The MDS ratio is obtained as the ratio of the repair cost of MDS erasure coding over the repair cost of PR. The large values shown in the figure for short maintenance epochs illustrate the significant cost savings of PR redundancy versus regular MDS. Although shorter maintenance epochs produce the best repair ratios, the absolute repair cost of the system is also the highest (see Figure 22). Consequently, the system needs to be engineered to achieve a compromise between the absolute repair cost and the best repair cost ratios. In particular, we believe that the system should

not operate with an Ideal ratio larger than 1.5. That is, the repair cost of PR should not exceed the repair bandwidth of ideal MDS by more than 50%.

The results in Figure 24 also illustrate an important adaptability requirement for the automated redundancy repair mechanism. For example, if the system is tune up to operate with an Ideal ratio no larger than a given value, its maintenance epoch is different for node sets with different churn behaviors. For an overlay network with exponentially distributed average node session lengths and  $\lambda^{-1} = 15$  minutes, the maintenance epoch is 6 minutes and for  $\lambda^{-1} = 20$  minutes it is 8 minutes (33% longer). A misconfiguration on either case would cause the system to consume more bandwidth (shorter epochs) or lower Ideal ratio performance than desired (larger epochs). The key conclusion from these results is that for networks with dynamic and heterogeneous node behaviors (i.e., most P2P networks), an automated redundancy maintenance mechanism must be capable of self-tuning in order to achieve a predefined cost performance tradeoff.



**Figure 24. PR Cost vs MDS and Ideal MDS Cost**

The next question we need to tackle to define the best parameter settings for a redundancy maintenance mechanism is; what metric should we use to control the adaptation of maintenance

epochs? The obvious answer is average node availability, but this implies tracking down the sessions (or residual lifetimes) of peers. This presents its own set of challenges [44] and might not reflect accurately fragment availability at the file granularity level. Instead, we propose the use of the redundancy loss model presented earlier in equation (4.26):

$$L_{\mathfrak{g}}^{\delta} = \sum_{i=1}^m i \cdot \binom{m}{i} \cdot (a_f^{\delta})^{m-i} \cdot (1 - a_f^{\delta})^i = m \cdot (1 - a_f^{\delta}) \quad (4.26)$$

This model establishes a direct relationship between average fragment availability and the number of fragments lost during a maintenance interval. Therefore, instead of measuring node availability, we can simply count the number of fragments lost during a maintenance epoch. Furthermore, we can account for this metric at different levels of the granularity. For instance, we can track only the total number of missing fragments (i.e., the MDS component) or the total number of fragments with 1, 2 or  $r$  copies available (i.e., MDS plus replication). Using the number of fragment lost (or alternatively, the percentage of fragments still available) is simpler to implement and manage than nodes' sessions lengths. For the rest of this work, it is assumed that the availability properties of the overlay network are obtained using this metric.

Next, we address the architectural and system-level design issues related with the creation of an automated redundancy maintenance mechanism.

#### 4.4 REDUNDANCY MAINTENANCE

The previous section provides us with insights about the selection of a redundancy method for P2P environments under churn. For the remainder of this work, we assume the use of a Proactive Replication, PR, redundancy scheme. The coding component of this method is a

maximum distance separable (MDS) erasure code and replication redundancy is applied on top of it to achieve low repair bandwidth. This section describes the system and architectural considerations to build a scalable and efficient redundancy maintenance system to improve content availability in P2P networks under churn.

The system's goal is to improve the availability of every single file managed by the overlay network regardless of their popularity. We refer to this set of files as the system content space. To accomplish this goal, the system fundamental building block is a distributed hash table (DHT) lookup functionality that supports an efficient indexing and retrieval of individual items. DHTs are characterized by their decentralization, fault tolerance and scalability properties, which can be capitalized to build a completely distributed redundancy maintenance mechanism.

The architecture of the system we are proposing is completely distributed. That is, all nodes have equal functionalities and they share responsibilities for the operation of the system. The design objectives of this system are the following:

- 1) Minimize repair bandwidth
- 2) Scalability and self-organization
- 3) Flexibility and Simplicity

The first objective is crucial for the performance of the system. In P2P networks, access bandwidth is more limited and expensive than storage space [86]. Limiting the amount of information transferred to restore the data lost when nodes leave the system is not a trivial task. Our proposed PR redundancy scheme is just the first component toward achieving this objective. Other important factors needed to minimize the repair bandwidth of the system were presented earlier in this chapter. For example, according to the analytical results presented earlier (see Section 4.4.2), we determined that average fragment availabilities between 0.65 and 0.85 can be a good compromise between reliability and cost.

The second objective, scalability and self-organization, is crucial to accommodate the large scale nature of P2P systems. The system we are presenting pursues these goals by being completely distributed. The responsibility of managing the system's redundancy information is assigned uniformly at random across the overlay. In addition, we decouple the repair and control components of the redundancy system. Nodes performing control functions decide when repairs are needed, but they are not involved in the actual repair process (transfer of data); thus, avoiding the creation of a common bottleneck.

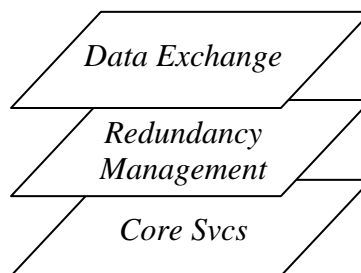
The third objective, flexibility and simplicity, is embedded into several of the components of the system, such as the redundancy method and the availability metric. We avoid the use of complex data structures whenever possible. The hybrid redundancy scheme we are proposing implements a simple and yet highly flexible design. Its structure and repair mechanism are much simpler than other schemes (such as exact-MBR). In addition, its replication component can be adapted dynamically, as opposed to the fixed parameters that must be used with other redundancy schemes. For the availability metric, we simply count for the number of fragments lost between maintenance epochs, avoiding the complexity of measuring the system churn rate. At the same time, the system is capable to adapt to the unique availability conditions of each item (i.e., the set of nodes storing the file).

The following sections describe the components of the system we are proposing. First, from an architectural perspective, to describe the fundamental functionalities required using a bottom-up design and second, using an agent-based model to describe the different roles that nodes can play and the processes governing their interactions.

#### 4.4.1 System Architecture

The system is modeled as a three layers architecture, shown in Figure 25. At the bottom layer, we have a set of core functionalities. The middle layer manages the redundancy state information of the system and the upper layer handles the exchange of data among nodes. The operation of the upper layers depends on the functionality provided by the underlying level(s).

The bottom layer, core services, implements the fundamental functionalities on top of which the system is build. In particular, a key-based lookup service for location/indexing of content. The objective of this component is to define a single (still dynamic) point of contact for content retrieval and redundancy management. All structured P2P systems support this functionality, and some of the most important deployed unstructured P2P systems (such as Gnutella and BitTorrent) have it implemented. This functionality maps content identifiers to a key-based naming space (key space) that provides an easy and consistent access to contents independently of their popularity rank. Nodes cooperatively manage this location/indexing service by assuming responsibility for a portion of the key space, which works as a rendezvous point for the nodes storing content and those requesting it.



**Figure 25. System Architecture for Redundancy Maintenance**

The middle layer, redundancy management, uses the information stored in the key space to monitor the redundancy state information of the system. We assume that all nodes storing items



notify the system of their availability using keepalive messages or similar methods. The system verifies regularly whether repairs are needed or not. If repairs are needed, the proper set of nodes gets notified to start a repair process. In addition, during each evaluation the system adjusts the periodicity of its verification process in order to minimize the repair bandwidth consumed.

The top layer, data exchange, initiates, monitors and terminates the data transfer between nodes. For repairs, the process is purely P2P and for content retrieval, the system supports up to  $k$  concurrent data transfers to improve performance. In the case of failed or staled transmissions, this layer takes care of initiating alternative data transfers.

At each layer, the system performs a set of processes that control the interaction between nodes. For the description of these processes we use an agent-based model, but before discussing this model and the processes regulating the interaction among agents, we need to precise some terms and concepts used in the following section.

The system's content space is a finite set of data objects. We use the term item in reference to a complete copy or a fragment of these objects. Given a set of items  $\{i_1, i_2, \dots, i_m\}$  with unique identifiers  $\{l_1, l_2, \dots, l_m\}$  associated with the complete copy of the original data, there is a unique node in the overlay named the root node of item  $i_i$ . Let  $r(i_i)$  denote this relationship. For a DHT-based overlay network routing infrastructure,  $r(i_i)$  is defined as follows:

$$r(i_i) = n_x \mid \text{hash}(id_{n_x}) \text{ is the closest numerical value to } \text{hash}(l_i) \quad (4.36)$$

where

$n_x$  is an active overlay node

$id_{n_x}$  is a unique identifier of node  $n_x$  (e.g., a node's socket).

$\text{hash}$  is a well known hash function, such as SHA-1<sup>24</sup>.

---

<sup>24</sup> <http://en.wikipedia.org/wiki/SHA-1>

A lookup is a procedure that determines the current root node for an item. That is, it returns the current  $r(i_i)$  value. We want to emphasize the use of the word current, because the transient nature of the participation of nodes in a P2P system causes the value of  $r(i_i)$  to change, eventually. We use the term key in reference to the identifier of nodes and data object in the DHT. For instance, the key of item  $i_i$  is  $hash(l_i)$  and the key of node  $n_x$  is  $hash(id_{n_x})$ . The processes we describe next are performed regularly by the active nodes in the overlay network in order to cope with churn, but in order to make an efficient use of their resources and to minimize cost, their periodicity is adjusted dynamically.

#### 4.4.2 Agents

Nodes performing multiples client/server functionalities, either alternatively or simultaneously is at the core of the P2P application paradigm. The operation of the P2P redundancy maintenance system we are proposing is not the exception. These functionalities can be conceptualized as agents. An agent is an instance of a system functionality being performed by individual nodes. Thus, the terms node and agent are used interchangeably throughout this section. The interactions between agents constitute the system. These interactions are regulated by predefined processes that generate/trigger events to activate specific agents on nodes. Depending on their role in the redundancy repair mechanism, nodes can play three fundamental roles for a specific item  $i_i$ .

- *Holder*. These are nodes with a whole or partial copy of item  $i_i$ .
- *Index*. This is the root node of item  $i_i$ .
- *Candidate*. These are nodes not currently associated with item  $i_i$ .

Nodes alternate between roles according to the event being handled. For example, if the node receives an item registration message, the index agent gets activated to process the request. Agents operate at the content level, which implies that nodes can execute multiple instances of the same agent concurrently, each associated with a different item. For instance, a peer can be an Index for items  $i_i$  and  $i_{i+1}$ , a Holder of item  $i_j$  and a candidate for item  $i_k$ .

Holder nodes maintain a list of index nodes for the items they currently store, contacting them periodically to indicate the availability of these items. Index nodes maintain a list (i.e., index) of all holder nodes with copies of a particular item. This indexing information reflects the current redundancy state information of the system; thus, it can be used to determine when to perform repairs. When a repair for item  $i_i$  is needed, holder nodes for this item are notified to start a repair process, which consist on contacting one or more candidate nodes and transfer data to them. Upon finalizing the data transfer, candidate nodes become holder nodes for item  $i_i$ .

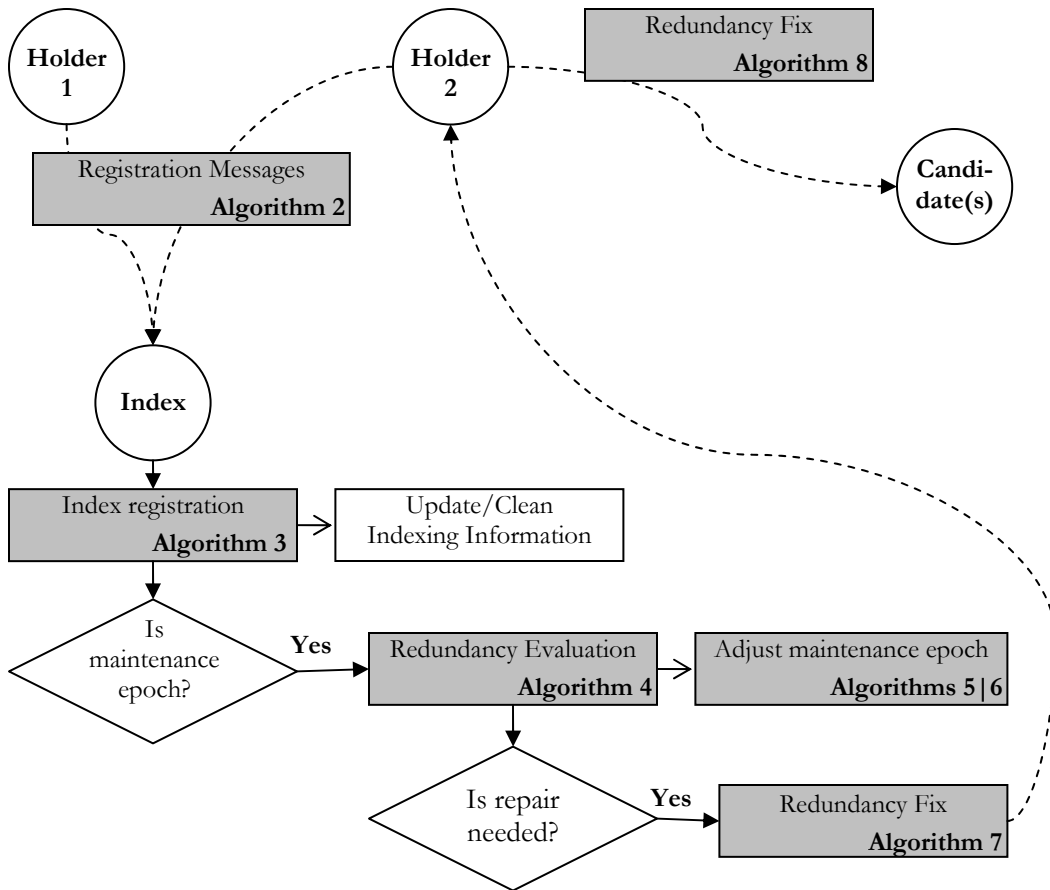
#### 4.4.3 System Processes

The redundancy maintenance system we are proposing can be modeled as a dynamic system of autonomous interacting agents. The interaction of these agents is circumscribed to three fundamental (interconnected) processes:

- Registration
- Redundancy Evaluation (maintenance)
- Redundancy Repair

Figure 26 shows the complete procedural flow of the automated redundancy maintenance system we are proposing. Peers are indicated by circles and labeled according to the agent being executed at each stage. The procedures and functions for the algorithms we are presenting are

indicated by rectangular boxes, together with some sub-processes executed at Algorithm 2 and Algorithm 3. Network messages between nodes are indicated by dashed lines. In the diagram, it is assumed that all messages correspond to a single content item. In addition, it is assumed that the registration message from Holder 1 does not reach the Index node within the maintenance epoch's maintenance window<sup>25</sup>, but the registration message of Holder 2 does.



**Figure 26. System Processes Overview**

The mechanisms presented above assume that all nodes cooperate voluntarily in the redundancy maintenance process. In practice however, we know that P2P networks are formed by autonomous nodes with heterogeneous resources and diverse objective functions. As a result,

<sup>25</sup> This parameter is presented in Algorithm 4. Its purpose is to limit during how many seconds the system performs repairs at each maintenance epoch.

without the proper incentives, nodes will avoid cooperating in the redundancy maintenance process, which could lead to the overall unfairness and performance degradation of the network. In other words, individual nodes will try to obtain a maximum gain from the system committing the least amount of resources possible unless we define cooperation rules to balance the performance and capabilities of individual nodes with the overall performance goal of the system. We address this concern in the next chapter. For the rest of this section we assume that nodes cooperate fully.

#### Procedure Registration()

*Purpose:* Send keep\_alive messages to item's root nodes

{*n*: total number of items this nodes holds}

{*item\_key*: key value of item *i*}

{*item\_seg*: fragments of item *i* that this node holds}

```

1:  for item = 1 to n do
2:      Get item_key
3:      Get item_seg
4:      while ( next_registration_time(item)-now <= 0 ) do
5:          if root(item) == null then
6:              root(item)=lookup(item_key)
7:          end if
8:          send registration_msg(root(item), item_key, item_seg)
9:          success = wait_for_ack()
10:         if success then
11:             set next_registration_time(item) = minimum + random(mean)
12:         else
13:             root(item) = null;
14:             set next_registration_time(item) = now
15:         end if
16:     end while
17: end for

```

#### **Algorithm 2. Items Registration Process at Holder Agent**

The processes we are presenting define the set of rules and algorithms that nodes should follow to cooperatively improve the availability of content. The first process, registration, is the only self-initiated process. The other two processes are triggered in response to events generated by the registration process. Registration starts when *Holder* agents attempt to contact their index nodes to announce the availability of the content they hold. The procedure performed by *Holders*

is presented in Algorithm 2. The objective of this algorithm is simple, to push the state information of each node sharing content into the overlay.

Unless noted otherwise, it is assumed that key values are used to uniquely identify items within the algorithms presented.

To improve the efficiency of the registration process, the system allows the aggregation of items. That is, when several items share the same root node, a single registration message is used. At the receiving end of the registration process, *Index* agents process each registration message and trigger the redundancy evaluation process, if needed. Algorithm 3 shows the pseudo code for the *Index* agents. The interaction between the *Index* agent for item  $i_i$  and all its *Holder* agents provide the fundamental functionality of the system, a dynamic distributed content indexing service; as described earlier for the core services layer of the system architecture.

**Table 19. IndexEntry Data Structure**

Field	Description						
<i>File_ID</i>	Contains a hash key						
<i>Last_evaluation</i>	Contains a time stamp						
<i>Last_unique</i>	Contains a file availability metric						
<i>Redundancy-SetTable</i>	Contains an array of: <table border="1" data-bbox="701 1289 1302 1402"> <tr> <td><i>Peer</i></td><td>Contains an IP socket</td></tr> <tr> <td><i>Fragment(s)</i></td><td>Contains a bit map</td></tr> <tr> <td><i>Last_registration</i></td><td>Contains a time stamp</td></tr> </table>	<i>Peer</i>	Contains an IP socket	<i>Fragment(s)</i>	Contains a bit map	<i>Last_registration</i>	Contains a time stamp
<i>Peer</i>	Contains an IP socket						
<i>Fragment(s)</i>	Contains a bit map						
<i>Last_registration</i>	Contains a time stamp						

*Index* agents maintain an IndexEntry data structure, shown in Table 19, for each of the items they are responsible for. The information, per file, consists of four elements: the item ID in the key-naming space, one time stamp, the result of the last file availability evaluation and a table with all the nodes in the redundancy-set of the item. For this redundancy-set table, three values are recorded for each node: peer ID (i.e., socket), a bit-map describing the portion of the original data stored by each peer and a time stamp of the last registration message received from each node.

When an *Index* agent determines that a node (which recently joined the overlay), is the root node of one or several of the items it is currently indexing, it transfers the respective IndexEntry data structure to the new node. On the other hand, when an *Index* node leaves the overlay, the IndexEntry data structure is lost. In this case, the data structure is regenerated once the *Holder* agents detect the failure of their current root node and contact a new one.

The redundancy evaluation process is triggered in an *Index* agent only at the beginning of each maintenance epoch. The MAINT\_WINDOW constant defines this interval. By default, the system uses a value of 30 seconds. Only the registration messages received within this interval trigger the redundancy evaluation process; all others are just used to update the IndexEntry data structure.

**Procedure Handle\_Registration(registration\_msg)**

***Purpose:*** Update IndexEntry for all items in registration\_msg and evaluate if repairs are needed

{x: number of items contained in registration\_msg}  
 {EPOCH: default redundancy maintenance interval}  
 {shift\_epoch: adapts maintenance epoch dynamically. Updated by Evaluate\_Redundancy process}  
 {MAINT\_WINDOW: interval during which redundancy evaluation can be triggered}

```

1:  peer = getHolder(registration_msg)
2:  now = current time
3:  for item = 1 to x do
4:    id = get_key(item)
5:    segments = get_fragments(item)
6:    IndexEntry = getIndexEntry(id)
7:    update IndexEntry.Redundancy-SetTable using (id, peer, segments, now)
8:    if ( now - IndexEntry.last_evaluation ) > (EPOCH + shift_epoch) then
9:      IndexEntry.last_evaluation = now
10:   end if
11:   elapsed = now - IndexEntry.last_evaluation
12:   if elapsed < MAINT_WINDOW then
13:     initiate redundancy evaluation using peer
14:     IndexEntry.last_evaluation = now;
15:   end if
16: end for

```

**Algorithm 3. Items Registration Process at Index Agent**

The pseudo code for the redundancy evaluation process for *Index* agents is presented in Algorithm 4-6. Algorithm 5 and Algorithm 6 are mutually exclusive; these are two variants of the

algorithm used to adapt the system maintenance epochs dynamically. Algorithm 5 uses a smoothed average of the number of unique fragments available to control the maintenance epochs of the mechanism, while Algorithm 6 uses the number of fragments lost (or gained) during a single maintenance epoch.

The pseudo code for the redundancy repair procedure performed at *Index* nodes is presented in Algorithm 7. The process is quite simple; the *Index* agent gets a random candidate node for each of the bits the holder node *peer* can repair. These candidates are obtained from the IndexEntry data structures of other items *Index* node is root, or directly from its routing table. The only condition for this selection is that the candidate node must not be listed in the IndexEntry data structure of the item being repaired. In addition, a set of backup targets is added to the message in case that any of the preselected candidates is not longer available.

**Procedure Evaluate\_Redundancy(*peer*)**

*Purpose:* Evaluate if repairs are needed, and update the next maintenance epoch.

{*peer*: node that triggered this process}

{MIN\_SEG: minimum number of unique segments required}

{DO\_SMOOTH: controls which algorithm is used to update the maintenance epoch}

```

1:  remove stalled entries
2:  unique = unique segments available
3:  if unique > MIN_SEG then
4:    if DO_SMOOTH then
5:      shift_epoch = Update_MaintenanceEpoch()
6:    else
7:      shift_epoch = Update_MaintenanceEpoch_L()
8:    end if
9:    if cpl < 1.0 then
10:     missing = get IndexEntry missing fragments
11:     peer_seg = get fragments peer has
12:     if (missing & peer_seg > 0) then
13:       initiate redundancy repair using peer and missing & peer_seg
14:     end if
15:   end if
16: end if

```

**Algorithm 4. Redundancy Evaluation by Index Agent**



The redundancyFix message generated by *Index* nodes is received by a holder node and processed using the pseudo code presented in Algorithm 7. Upon receiving this message, the holder node will attempt to push a replica of the segments requested. The initial target for this data transfer is the candidate node previously selected by the *Index* node (Algorithm 8). Holder nodes use an auxiliary data structure named pending replicas to monitor the progress of each repair. If the replication fails or stales for any reason, the holder node will retrieve an alternate candidate node from the pending replicas data structure and to complete the redundancy repair with it. Once the repair request is completed, or all target entries have been exhausted, the pending replicas data structure for item  $i$  is cleared.

**Function Update\_MaintenanceEpoch()**

*Purpose:* Adjust the length of the maintenance epochs of the system.

{ $m$ : total possible number of unique segments}  
 {alpha: smoothing average factor, default = 0.4}  
 {TARGET: number of fragments available [ $k...m$ ]; default  $0.3*m$ }

```

1:  unique = unique segments available
2:  avg = alpha*last_unique+(1-alpha)*unique
3:  if avg < TARGET then
4:      result = -10 seconds
5:  else
6:      if avg == m then
7:          result = 0
8:      else
9:          if avg > 0.8*m then
10:             result = 5 seconds
11:          else
12:             result = 2.5 seconds
13:          end if
14:      end if
15:  end if
16:  last_unique = avg
17:  return shift_epoch + result

```

**Algorithm 5. Smooth Maintenance Epoch Mechanism at Index Agent**

In summary, the automated redundancy maintenance mechanism we are proposing is composed of three fundamental processes: registration (with instances at *Holder* and *Index* nodes), evaluation (performed exclusively by *Index* nodes) and repair (initiated by *Index* agents, but

executed entirely by *Holder* agents). The effectiveness of these mechanisms depends greatly on nodes participating on these processes. To promote this participation while preserving node autonomy, we propose the use of economic incentives. These incentives, define a set of simple, yet effective, rules for fair exchange of resources in the overlay that can be easily integrated into the automated redundancy maintenance process describe above.

**Function Update\_MaintenanceEpoch\_L()**

*Purpose:* Adjust the length of the maintenance epochs of the system.

{*m*: total possible number of unique segments}

{TARGET: number of fragments available [*k...m*]; default  $0.3 * m$ }

```

1:  unique = unique segments available
2:  loss = last_unique-unique
3:  if loss > 0 then
4:      if loss > TARGET then          result = -10 seconds
5:      else
6:          if loss > TARGET/2 then    result = -2.5 seconds
7:          end
8:      end
9:  else
10:     if loss == 0 then
11:         result = 0                // in Adapt>=0, result=2.5
12:     else
13:         if loss < -TARGET/2 then    result = 5 seconds
14:         else result = 2.5 seconds
15:         end if
16:     end if
17: end if
18: last_unique = unique
19: return shift_epoch + result

```

**Algorithm 6. Adaptive Maintenance Epoch Mechanism at Index Agent**

**Procedure RedundancyFix(*peer, segs*)**

*Purpose:* Send a redundancy repair request to *peer*.

{*peer*: holder node that will perform the repair}

{*segs*: bit-map with the list of fragments to be repaired}

{BACKUP: it controls whether to include backup targets or not}

```
1:  for  $i=1$  to num_bits_in(segs) do
2:      target = select random candidate
3:      seg = select random bit in segs
4:      add (target, seg) to redundancyFix_msg
5:  end for
6:  if BACKUP then
7:      for  $i=1$  to 6 do
8:          target = select random candidate
9:          add(target, nil) to redundancyFix_msg
10:      end for
11:  end if
12:  send redundancyFix_msg to peer
```

**Algorithm 7. RedundancyFix at Index Agent**

**Procedure handle\_RedundancyFix(*msg*)**

*Purpose:* Perform redundancy repairs.

{*msg*: redundancyFix message received}

```
1:  id = get item key from msg
2:  for  $i=1$  to size(msg) do
3:      target = get target i
4:      seg = get fragment i to be repaired
5:      if seg <> nil then
6:          push seg to target
7:          store (id, target, seg) in pending replicas
8:      else
9:          store (id, target, nil) in pending replicas
10:      end
11:  end for
```

**Algorithm 8. RedundancyFix at Holder Agent**

## **5.0 INCENTIVES**

This chapter presents an incentive-based mechanism to promote the participation of nodes in a redundancy-maintenance process to improve content availability in P2P networks and describes its implementation. The structure of this chapter is as follows. Section 5.1 introduces the use of economic models in P2P networks. Section 5.2 outlines the components and design guidelines of the incentive-based mechanism. Section 5.3 describes its construction and in Section 5.4 we describe the formulation of an enhanced contribution metric for our mechanism. Finally, in Section 5.5 we conclude with a discussion of additional considerations regarding the implementation of our incentive-based mechanism.

### **5.1 ECONOMIC MODELS**

Economic models are simplified abstract representation of complex economic processes in the form of logical, qualitative and/or quantitative relationships between a set of variables. Problems in the fields of economics and distributed systems have multiple structural similarities, which allow both disciplines to borrow tools, terminology and models from each other to conduct sophisticated analysis. In that regard, economic inspired models have been used in multiple research initiatives pursuing fair allocation of resources in P2P networks [14, 15, 17, 62]. The key motivation behind these initiatives is that the benefits of P2P networking depend heavily in cooperation among peers,

but in reality most systems suffer from the *tragedy of the commons* problem [87]. This dilemma arises from the situation in which multiple users, acting independently and rationally according to their own self-interest, contribute as few resources as possible to the network [88]; even when these actions deteriorates the long-term properties of the system from which they benefit. For example, *free-riders*<sup>26</sup> is an instance of this type of problem in file-sharing networks. Measurement studies of deployed systems reported that in 2000 about 66% of the participants in P2P networks were *free-riders* [4] and recent studies confirm that this phenomenon still persist, with more than 70% of the nodes being *free-riders* [89].

One branch of economics that has been used extensively to alleviate the problem of fair allocation of resources in P2P networks is economic incentives. This type of mechanism assesses, coordinates and controls the outcome of a system by influencing the behavior of its participants<sup>27</sup> in a certain way. We use this approach in our research.

There are three classes of incentive-based mechanisms: monetary-payment, fixed contribution and reciprocity-based [17]. A monetary-payment scheme relies on the existence of a central banking authority, which could result in a central point of failure and scalability concerns. In fixed contribution schemes, nodes are forced to make minimum contributions to the system, which does not take node heterogeneity into account. Reciprocity-based schemes use rules to control the allocation of resources as a function of nodes' contribution to the network's overall objective. An example of this technique is BitTorrent's bandwidth bartering mechanism. This dissertation adopts the reciprocity-based approach to build an incentive-based mechanism because

---

<sup>26</sup> *Free-riders* are nodes consuming resources in the overlay while sharing few to no resources with the rest of the network.

<sup>27</sup> The term participant in economic models is equivalent to node in P2P terminology.

it supports a completely distributed implementation to achieve scalability, and it is also capable of accommodating node heterogeneity.

## 5.2 INCENTIVE-BASED MECHANISMS FOR CONTENT AVAILABILITY

The description of the incentive-based mechanism we propose is organized in five sections entitled: *objective*, *design guidelines*, *agents*, *control metric*, and finally, *incentive functions*. The following sections describe in detail each of these stages.

### 5.2.1 Objective

The goal of the incentive-based mechanism is to promote the participation of nodes in a redundancy-maintenance process by regulating the allocation of resources between nodes in proportion to their contribution. In other words, our proposed incentive-based mechanism is a bartering economy of resources (e.g., bandwidth) vs contribution, in the form of redundancy state information (e.g., fragments). Chapter 4.0 describes the components and operation of the redundancy-maintenance process. It is assumed that by participating in the redundancy-maintenance process, nodes can attain better levels of service, statistically speaking, than those that do not collaborate.

Figure 27 illustrates the desired effect of the incentive-based mechanism on the behavior of nodes. These figures portrait the levels of contribution of three random nodes, through time, toward the global objective of the system, which in our case, is content availability. The right

figure shows the results we expect to achieve. Notice that the behavior of nodes still is random and heterogeneous in nature, but their content availability contribution exhibits an upward trend.

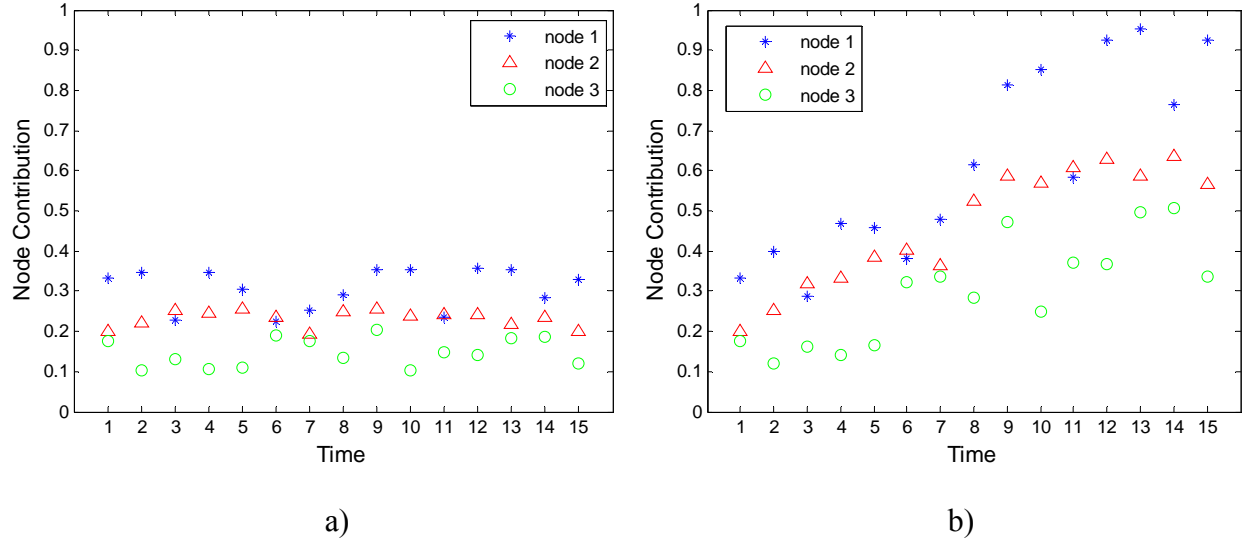


Figure 27. Effect of Incentive-based Mechanism: a) Normal System, b) System with Incentives

### 5.2.2 Design Guidelines

Four design properties are sought for the incentive-based mechanism presented in this dissertation: *incentive*, *penalty*, *decentralization*, and *adaptability and lightweight* [17]. The *incentive* property refers to the creation of a positive encouragement component that can affect the behavior of nodes so that the utility of the system increases. In contrast, the *penalty* property defines the negative effect that non-compliant nodes may face by contributing few to no resources towards the common goal. The *decentralization* property is adopted as a fundamental architectural consideration to improve the scalability and robustness of the system. Last, the *adaptability and light weight* property refers to the capability of the incentive-based mechanism to adapt efficiently to different networking and node behavior scenarios while minimizing the produced overhead.

To minimize overhead, the incentive-based mechanism needs to be embedded into existing system functionalities, which includes the redundancy maintenance process. The redundancy maintenance process presented in Chapter 4.0 can be easily augmented with incentives. From a communications perspective, the implementation of the incentive-based mechanism we propose only requires a contribution metric to be piggybacked into existing systems messages.

### 5.2.3 Agents

In economic models, it is common practice to conceptualize different stakeholders in the economic process as agents. Agents represent autonomous entities participating in the economic process and can affect its outcome. In a P2P economic model each physical node performs multiple agent roles. The incentive-based mechanism we propose for content availability employs the following agents: *Holder*, *Index*, *Requestor* and *Candidate*. *Holder* nodes store information with the purpose of sharing it with the rest of the network, improving the content availability of that specific piece of information. *Requestors* are nodes that consume resources from the network. In other words, *Requestors* are nodes downloading information. *Index* nodes manage the list of all available *Holder* nodes storing data for a specific data item and provide a relaying service between *Requestors* and *Holders*. When a *Requestor* node wants to obtain an item, it contacts an *Index* node to get a list of *Holder* nodes to download it from. During the redundancy-repair process, nodes available to store data for a specific item are called *Candidates*.

*Holder* nodes have two modalities: provider and publisher. The provider modality answers download requests from *Requestor* nodes. The publisher modality is in charge of communicating to *Index* nodes its presence as well as restoring lost redundancy information.



### 5.2.4 Control Metric

The incentive-based mechanism we propose measures the level of participation of nodes in the content availability of the network and uses it as a control metric. Assuming that redundancy is being used to provide reliable access to content, the level of participation is measured in terms of the primary elements defined in the redundancy scheme. For replication redundancy, these elements are copies of files, and for code-based redundancy schemes the primary elements are file's fragments. The term *content contribution* or simply *contribution* is used in the rest of this work to denote the control metric of our incentive-based mechanism.

With respect to the redundancy scheme to be used in the system, we believe that code-based mechanisms should be preferred because in addition to a fine-grained content contribution metric, the system gains additional properties. For instance, nodes are not accountable for the fragments they hold (since these are randomly selected), censorship becomes an extremely hard task to accomplish, and code-based mechanisms represent a good tradeoff between storage space and reliability. In the following sections we assume that the redundancy scheme used in the network is our proposed redundancy scheme, Proactive Replication (PR). That is, redundancy is code-based.

Let  $r_m$  denote the contribution of node  $m$ , which is defined as follows:

$$r_m = F(\chi^m) \cdot G(\chi^m, t, s, p) \quad (5.1)$$

where:

$$F(\chi^m) = \sum_{\forall i} \sum_{\forall j} \chi_{i,j}^m = n_m \quad (5.2)$$

We name this function *fragment count*.  $\chi^m$  is an indicator function (defined as before in Section 4.2.1 for our PR redundancy scheme):

$$x_{i,j}^m = \begin{cases} 1 & \text{if } f_{i,j} \text{ is hosted in node } m \\ 0 & \text{otherwise} \end{cases} \quad (5.3)$$

and  $G(\chi^m, t, s, p)$  is named *contribution gain*, which is a non-decreasing function of fragments' age (i.e., time), size and popularity. For this function, we define two variants named *basic contribution gain* and *enhanced contribution gain* as follows:

$$G(\chi^m, t, s, p) = \begin{cases} 1 & \text{Basic} \\ H(t) \cdot K(s) \cdot L(p) & \text{Enhanced} \end{cases} \quad (5.4)$$

For simplicity, we use the first formulation (i.e., basic contribution gain) to describe the structure and implementation of our incentive-based mechanism in the following sections. In Section 5.5 we describe the *enhanced contribution gain* case.

**Table 20. Content Contribution Metric**

<i>Parameter</i>	<i>Name</i>	<i>Description</i>
$r_m$	<i>Content contribution</i>	measures the content availability contribution of node $m$
$F(\chi^m)$	<i>Fragment Count</i>	counts the number of fragments node $m$ stores
$G(\chi^m, t, s, p)$	<i>Contribution Gain</i>	augments <i>fragment count</i> as a function of items' age ( $t$ ), size ( $s$ ) and popularity ( $p$ )

In the redundancy system presented in this dissertation, nodes are expected to include their content contribution when communicating with other nodes. Depending on the nature of the communication, the content contribution metric could be used to assign a level of service to the requesting node in proportion to its contribution.

### 5.2.5 Incentive and Penalty Functions

The incentive-based mechanism translates the node contribution metric in two system properties that promote cooperation in the redundancy-maintenance process by providing positive

and negative feedback to the end-user. These properties are realized using two functions, namely the *utility* and the *cost* functions. In the mechanism we propose, the utility function is used as incentive and the cost function as penalty. In turn, these two functions are based on a sigmoid function, which has the form:

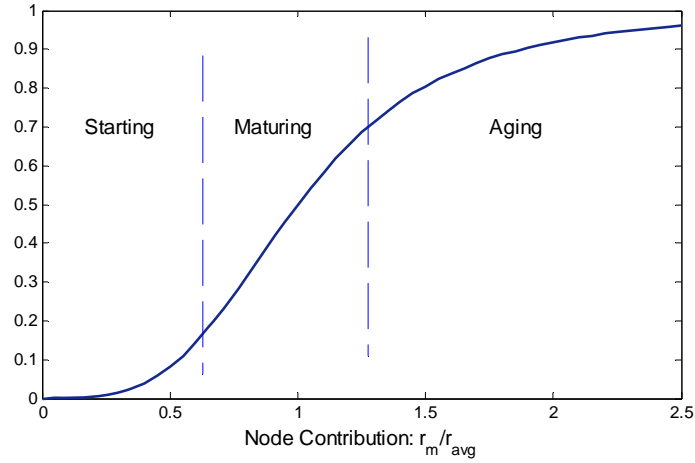
$$S(r_m) = \frac{(r_m/r_{tgt})^\sigma}{1 + (r_m/r_{tgt})^\sigma} \quad (5.5)$$

where:

**Table 21. Sigmoid Function Parameters**

<i>Parameter</i>	<i>Description</i>
$r_m$	is the content contribution of node $m$
$r_{tgt}$	is the target (average) content contribution
$\sigma$	is the shape parameter, $\sigma > 1$

Sigmoid functions have a tilted S-shaped curve that return values between zero and one. The return value can be used to assign a level or probability of service to requesting nodes. The use of a normalized content contribution value in this function allows us to model the system behavior independently of the scale of the content space. Sigmoid functions exhibit three characteristic phases that portray the desired evolution in the behavior of peers for the incentive-based mechanism: *starting*, *maturing* and *aging* [33]. These phases are illustrated in Figure 28. During the starting phase a node's contribution is small and so is its payoff. In the maturing phase, as nodes increase their contribution, their payoffs increase rapidly. Finally, when nodes reach the aging phase, additional increments to their contribution do not increase their payoff significantly.



**Figure 28. Sigmoid Function**

The cost function is named *Replication Probability (RP-cost)* and the utility function is named *Transmission Bandwidth (TB-utility)*. Given the dual nature of nodes as providers and consumer of resources, the cost and utility functions can be interpreted as a positive (i.e., incentive) and negative (i.e., penalty) feedback to user behavior. In *RP-cost*, the function is used as a probability by the incentive-based mechanism. It determines the likelihood of a node accepting to store additional elements with the purpose of improving the overall content availability of the network, and increasing its future performance metric (i.e., *utility*). In other words, the *RP-cost* function represents a cost for the downloading node, but also an opportunity to improve its performance. The *TB-utility* function is used to assign a level of service to requesting nodes. It determines how much bandwidth nodes can expect from other peers during a data transfer. This transmission bandwidth can be interpreted as a performance metric (i.e., *utility*) or as the amount of time nodes have to stay connected to finish the download (i.e., *cost*).

To be able to effectively influence the behavior of nodes towards participating in the content availability process, the *RP-utility* and *TB-cost* functions use a modified version of the sigmoid function described above. The modifications presented ahead adjust the shape of the curve

and modify the mapping of the nodes' contributions to the starting, maturing and aging phases of the sigmoid function.

### 5.2.5.1 Replication Probability Function

For the *RP-cost* function, two modifications are done on (5.5). First, the complement of the sigmoid function (i.e.,  $1-S(r_m)$ ) is used because when nodes have a small content contribution the probability of hosting new replicas must be high. The function should allow nodes to improve their *utility*. When nodes have a high content contribution, the replication probability should be close to zero (the aging phase). Second, the incentive-based mechanism uses shifted node contribution values. This modification defines at what contribution level nodes should stop accepting new replicas. In the original sigmoid function, when a node's contribution reaches the target content contribution value ( $r_{tgt}$ ) its replication probability is equal to 50% and it does not reach zero until the contribution value of the node is much larger than the target value (e.g.,  $r_m > 3 \cdot r_{tgt}$ ). Instead, the *RP-cost* function should reach zero when the target content contribution value is reached.

The general form for the *RP-cost* function is:

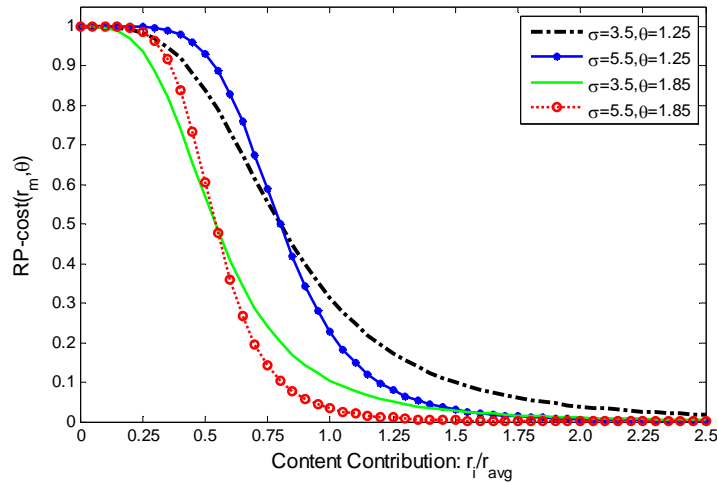
$$RP-Cost(r_m, \theta) = 1 - \frac{(\theta \cdot r_m / r_{tgt})^{\sigma_c}}{1 + (\theta \cdot r_m / r_{tgt})^{\sigma_c}} \quad (5.6)$$

where:

**Table 22. Replication Probability Function Parameters**

<i>Parameter</i>	<i>Description</i>
$r_m$	is the content contribution of node $m$
$r_{tgt}$	is the target (average) content contribution in the network
$\theta$	is the contribution shift parameter, $\theta > 1$
$\sigma_c$	is the cost shape parameter, $\sigma_c > 1$

Figure 29 presents four alternative *RP-cost* functions that combine two shape and two shift parameter values. The value of the shape parameter determines the length of the maturing phase. For larger values of the shape parameter, the upper and lower transition points (i.e., the knees of the curve) from starting to maturing phases and from maturing to aging phases are more pronounced. The shift parameter on the other hand, determines the contribution value at which the function reaches zero. The results presented in Figure 29 indicate that for a smooth transition from accepting new replicas to rejecting all of them, the values of the shape parameter should be kept low.



**Figure 29. Replication Probability Function**

### 5.2.5.2 Transmission Bandwidth Function

For the *TB-utility* function, a small offset value is added to the sigmoid function to avoid starving new nodes. When nodes first join the network, their content contribution is zero, and if these nodes are not allowed to receive at least a minimum transmission bandwidth, they would not be able to participate in the redundancy-maintenance process at all.

The general form for the *TB-utility* function is:

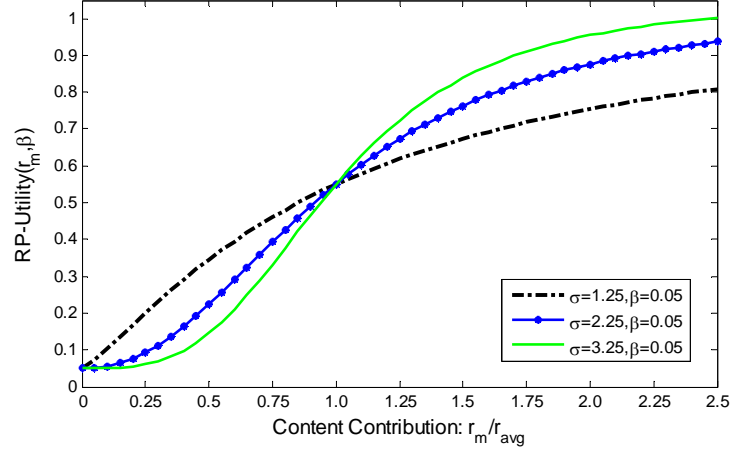
$$TB\text{-}Utility(r_m, \beta) = \beta + \frac{(r_m/r_{tgt})^{\sigma_u}}{1 + (r_m/r_{tgt})^{\sigma_u}} \quad (5.7)$$

where:

**Table 23. Transmission Bandwidth Function Parameters**

<i>Parameter</i>	<i>Description</i>
$r_m$	is the content contribution of node $m$
$r_{tgt}$	is the target (average) content contribution in the network
$\beta$	is the bandwidth bias parameter, $\beta > 0$
$\sigma_u$	is the utility shape parameter, $\sigma_u > 1$

Using a bandwidth bias too small could constitute a barrier to entry for new nodes and using a value too high would diminish the effectiveness of the positive feedback component of the system (i.e., incentive). For the incentive-based mechanism evaluated in Chapter 6.0, the bandwidth bias parameter selected is equal to 0.05. Figure 30 presents the *TB-utility* function employed in the experimental portion of this work for three alternative shape parameter values ( $\sigma$ ) of 3.25, 2.25 and 1.25. When a node's contribution reaches the target value ( $r_{tgt}$ ) the function takes a value of  $0.5 + \beta$  regardless of the shape parameter used. Higher values of the shape parameter produce a more pronounced knee effect before and after the target contribution value ( $r_{tgt}$ ). The effect of the shape parameter can be interpreted as a mean to modify the “length” of the maturing phase of the sigmoid function. For larger values of the shape parameter the function has a more pronounced “S” shape; with the starting and aging phases occupying a larger range of the node's contribution values.



**Figure 30. Transmission Bandwidth Function**

### 5.3 INCENTIVE-BASED MECHANISM CONSTRUCTION

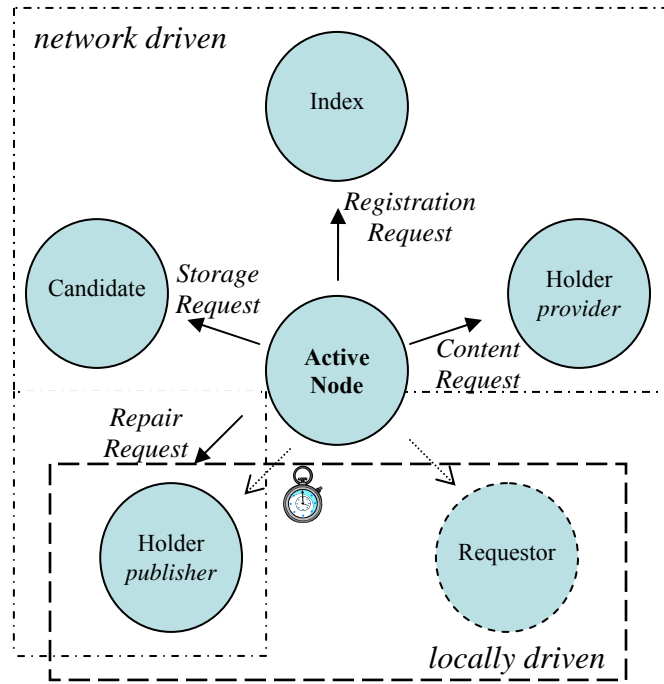
The use of the *RP-cost* and *TB-utility* functions in the redundancy-maintenance process is described using a finite state machine (FSM) representation of this process. The naming and operation of the states of this FSM correspond to the economic agent model presented earlier in Section 5.2.3.

#### 5.3.1 Redundancy-Maintenance Finite State Machine

The FSM in Figure 31 presents the state transitions of a single physical node using the incentive-based mechanism. The physical node could be performing multiple roles (i.e., agents) simultaneously, but for the sake of clarity throughout this description it is assumed that the node can only be in one state at a time. Throughout this description the term node is used in reference to the local physical node and the term peer is used in reference to other physical nodes.



From the redundancy-maintenance process point of view, the FSM can be in six different states: *Active Node*, *Index*, *Holder-provider*, *Requestor*, *Holder-publisher* and *Candidate*. In the first state, *active node*, the system is not performing any redundancy-related activity, but is participating actively in the overlay network. The second state, *holder-provider*, represents the node serving information to other peers (i.e., uploading on request). The third state, *requestor*, is when the node downloads information from the network (i.e., downloading from other peers, which must be in the *holder-provider* state). The fourth state, *holder-publisher*, represents the node pushing information into the network (i.e., generating new replicas). Finally, the sixth state, *candidate*, represents a node that is been considered to receive new replicas (i.e., this node is the target of another peer in the *holder-publisher* state).



**Figure 31. Redundancy-maintenance process FSM**

The transition of the physical node from the *active node* state to a different one is triggered by specific *events*. These *events* are classified in two categories depending on whether the event was generated by the local node or by a peer. In that regard, the states of the FSM are said to be

network driven or locally driven, as indicated in Figure 31. Network driven events are requests generated by peers (i.e., other physical nodes in the overlay network). In this figure, the transitions (i.e., arrows in the figure) indicate the name of the event that triggers the transition from the *active node* state to the respective “network driven” state. The transitions from the *active node* to the locally driven states (*holder-publisher* and *requestor*) are generated by timers or by the end-user.

The following sections describe the operation of each state/agent.

#### 5.3.1.1 Requestor State

The *Requestor* state represents nodes consuming resources (i.e., content) from the overlay. The node transitions to this state when it submits a request for content to the overlay. If the node remains in the system after completing a download, and if it decides to share the item<sup>28</sup>, this state contributes indirectly to the redundancy-maintenance process.

#### 5.3.1.2 Holder-Publisher State

The *Holder-publisher* state performs two activities depending on whether the transition was generated locally or by the network (i.e., a repair request from a peer). The local transition occurs periodically when nodes contact their respective *Index* peers. The occurrence of network transition on the other hand, is random in nature and is triggered when nodes need to perform a redundancy repair operation. The frequency of these repairs depends on the level of churn exhibited in the network.

For local transition, the number of *Index* peers to be contacted depends on the set of items stored by the node. To avoid traffic spikes, registration requests of individual items should be

---

<sup>28</sup> This is the default behavior for most file-sharing applications, but measurement studies have reported that peers do not keep all the items they download in the system’s share-directory.

unsynchronized, unless two or more items share the same *Index* node, in which case their registration requests are aggregated to minimize overhead.

For network driven transition, the node receives a *repair request* from one of its *Index* peers. This *repair request* instructs the node to perform a redundancy repair procedure (see Algorithm 7 and Algorithm 8 in Section 4.4.3). Repair request messages are created by *Index* peers and contain a list of elements for the local node (i.e., *Holder*) to replicate. The selection of the recipients for these replicas is performed by the *Index* peer according to their content contribution metric. If any of the recipients (*Candidate* peers) fail to respond, *Holders* can select a new recipient out of the list of alternative candidates provided by the *Index* peer in the *repair request* message.

#### **5.3.1.3 Holder-Provider State**

The *Holder-provider* state represents the server side of a download. When nodes receive *content request* messages (from *Requestor* peers) they look for a content contribution value in the message. If a contribution value for the peer is found, the node uses it in the *TB-utility* function to determine the amount of bandwidth to allocate for this transmission. A pseudo code for the bandwidth calculation done by *Holder-provider* nodes is presented in Algorithm 9.

### **Function getSpeed(*cont*)**

**Purpose:** calculate maximum bandwidth that can be assigned to this request

**Returns:** maximum bandwidth, speed [kBps]

{ *cont*: is the content contribution of the requesting node }

{ MIN\_TXBW: minimum transmission bandwidth }

```
1: upBW = freeUploadBW()
2: if upBW > MIN_TXBW then
3:   if cont>0 then
4:     speed = MIN_TXBW + upBW*utility(cont)
5:   else
6:     speed = MIN_TXBW + ( upBW*random() )%MIN_TXBW
7:   end if
8: else
9:   speed = 0
10 end if
```

### **Algorithm 9. Function getSpeed()**

There are procedures used in the function getSpeed() for which an algorithm has not been presented. These procedures are freeUploadBW() and random(). The first is a function that returns the total amount of upload bandwidth not yet committed by the node and random() is a uniformly distributed random variable between zero and one. The random() function is used for non-compliant peers and for compliant peers with zero contribution. It is possible to use a different function for non-compliant peers, but this does not represent any additional gain based on our *penalty* design guideline (Section 5.2.2).

The pseudo-code of the function utility() used in the getSpeed() function is presented in Algorithm 10. This function uses a procedure named IndexTable.GetAverageContribution(). The purpose of this procedure is to calculate the average contribution of peers that register their items at the current local node (*Index* agent). If this value is larger than a predefined average node contribution (TARGET\_CONT), the node will use this value instead for normalizing the content contribution of the requesting peer. It is assumed that a larger than expected average content contribution is due to the existence of a larger than expected content space. In such a case, the system would be able to adapt its behavior accordingly.

#### Function Utility(*cont*)

**Purpose:** calculate utility function for requesting node

**Returns:** utility value,  $u$  [UTIL\_BIAS,1]

{ *cont*: is the content contribution of the requesting node }  
{ TARGET\_CONT: default normalization value for content contribution }  
{ SIGMA\_UTIL: shape parameter }  
{ UTIL\_BIAS: bandwidth bias parameter }

```
1: avg = IndexTable.GetAverageContribution()
2: if avg > 1.5*TARGET_CONT then
3:   norm = avg
4: else
5:   norm = TARGET_CONT
6: end if
7: a = (cont/norm)^SIGMA_UTIL
8: u = UTIL_BIAS + a/(1+a)
```

#### **Algorithm 10. Function Utility()**

##### **5.3.1.4 Index State**

When a node receives a *registration request*, it transitions into the *Index* state and evaluates whether the originating peer should receive a *repair request*. Algorithm 4 presents the pseudo code employed by the *Index* node for this evaluation. As described in Section 4.4.3, Index nodes are in charge of decisions in the redundancy repair process. They are the brain of the process and the actual cost associated with the transfer of information is responsibility of other peers. In this way, *Index* nodes do not have a direct incentive to misbehave.

The number of repairs to be made for a single file can vary greatly depending on the set of nodes participating in the storage process. One approach to minimize this maintenance cost is selecting nodes with the highest availability, but this represents a negative incentive for highly available nodes and a scalability concern as well. In our implementation, we use an alternative approach to balance the load of the redundancy mechanism instead. The IndexTable data structure, presented in Table 23, keeps a list of the nodes within its array of IndexEntry elements (see Table

19 in section 4.5.3 for a description of this data structure). When the list of targets is built, nodes with the lowest contributions are added first.

**Table 24. IndexTable Data Structure**

Field	Description
<i>Content</i>	Contains an array of < IndexEntry>
<i>Nodes</i>	Contains an array of <peers>

### 5.3.1.5 Candidate State

Nodes transition to the *Candidate* state when they receive a Storage request from a *Holder-publisher* node. In order to accept/reject the Storage request, *Candidate* nodes evaluate their cost function. If the value is above a randomly generated number, the request is accepted and the *Holder* node is notified that it can initiate its upload (with a maximum bandwidth proportional to the content contribution of the publishing *Holder* node). The pseudo code employed by *Candidate* nodes is listed in Algorithm 11 and Algorithm 12.

#### Procedure AcceptReplica(*cont*)

**Purpose:** accept/reject Replica request from Holder peer

{ *cont*: is the content contribution of the Holder peer }

{ *mycont*: is the content contribution of this node }

```

1: if cost(mycont) > random() then
2:   bw = utility(cont)
3:   send accept message(bw)
4: else
5:   send reject message()
6: end if

```

**Algorithm 11. Procedure AcceptReplica()**

### Function Cost(*cont*)

**Purpose:** calculate cost function for *cont* contribution

**Returns:** cost value,  $c \in [0,1]$

{ *cont*: is the content contribution of this (or requesting) peer }  
{ TARGET\_CONT: default normalization value for content contribution }  
{ SIGMA\_COST: shape parameter }

```
1: avg = IndexTable.GetAverageContribution()
2: if avg > 1.5*TARGET_CONT then
3:   norm = avg
4: else
5:   norm = TARGET_CONT
6: end if
7: a = (cont/norm)^SIGMA_COST
8: c = 1 - a/(1+a)
```

**Algorithm 12. Function Cost()**

## 5.4 ENHANCED CONTRIBUTION METRIC

In Section 5.2.4 we defined node contribution,  $r_m$ , as the product of two functions, *fragment count*,  $F(\chi^m)$  and *contribution gain*,  $G(\chi^m, t, s, p)$ , where  $\chi^m$  is an indicator function (taking the value of one) for the fragments  $f_{i,j}$  stored at node  $m$ . In addition, we defined two alternative formulations for  $G(\chi^m, t, s, p)$ . First, a unitary *contribution gain* (i.e.,  $G(\chi^m, t, s, p)=1$ ) and second, a *contribution gain* with time, size and popularity components. The purpose of our second formulation is to augment our incentives mechanism with the means to induce additional changes in user behavior.

The reasoning behind enhancing our contribution metric with time is that by rewarding nodes for the time they share content, the incentive mechanism can also promote lower churn rates, which would reduce the redundancy maintenance traffic. With the addition of size and popularity, we expect that by rewarding the nodes that store rare, or large, items, we could further reduce the system's redundancy maintenance load. We assume that repairing large files' fragments is more

costly than repairing fragments for smaller files. Furthermore, we assume that highly popular items are replicated naturally in the system. That is, nodes are more likely to perform download requests by themselves (i.e., replicate the complete file, without incentives) for highly popular items than for rare items. Thus, the availability of popular items depends less on our redundancy maintenance process.

For our augmented *contribution gain* function we defined three auxiliary functions,  $H(t)$  for the time component,  $K(s)$  for size and  $L(p)$  for popularity. All of them are defined as staircase functions with the following form:

$$Z(w) = \frac{1}{n_m} \sum_{\forall i} \sum_{\forall j} \chi_{i,j}^m \cdot \sum_{x=1}^N \alpha_x \Psi_{Ax}(w) \quad (5.8)$$

where

$n_m$  is the result of the *fragment count* function,  $F(\chi^m)$ .

$\chi_{i,j}^m$  is an indicator function for fragment  $f_{i,j}$  being stored at node  $m$

$N > 0$  is the number of steps (i.e., levels) in the function

$\alpha_x$  is the gain at level  $x$ .  $\alpha_x < 1.0$  constitute a penalty and  $\alpha_x > 1.0$  is an reward.

$A_x$  are disjoint intervals in  $w$  and  $\Psi_{Ax}$  is the indicator function of  $A$ :

$$\Psi_{Ax} = \begin{cases} 1 & \text{if } x \in A, \\ 0 & \text{if } x \notin A \end{cases} \quad (5.9)$$

To define the time, size and popularity components of the *contribution gain* function we need to define a set of  $\alpha$  values and a set of  $A$  intervals reflecting the desired penalty or reward for each interval.

Let  $T=2$  minutes be a base time interval for the time component of our *contribution gain* function so that

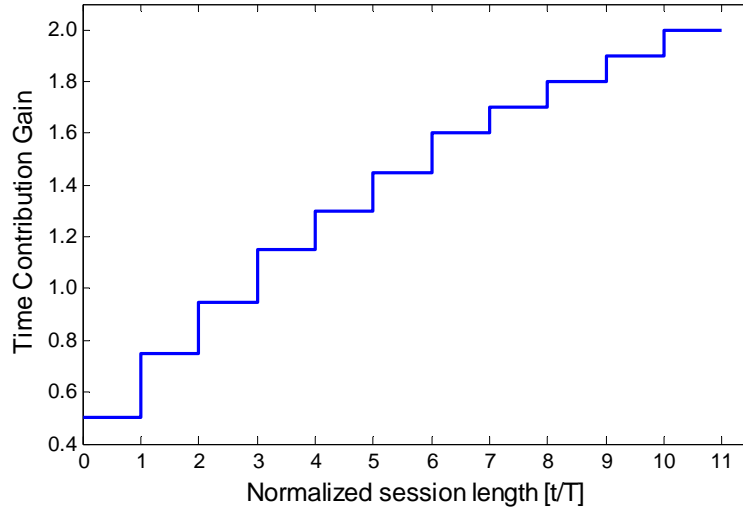


$$A_x^H = \{[0, T), [T, 2T), \dots, [10T, \infty)\} \quad (5.10)$$

for which we define the following penalty/reward value set

$$\alpha_x^H = \{0.5, 0.75, 0.95, 1.15, 1.3, 1.45, 1.6, 1.7, 1.8, 1.9, 2\} \quad (5.11)$$

The equations above define an incentive mechanism where nodes with session lengths shorter than 6 minutes receive a penalty (i.e.,  $\alpha_x < 1.0$ ) and those with session lengths of 6 minutes and larger receive an incentive. In addition, the reward for sessions larger than 10T (20 minutes) remains constant. Figure 32 presents the resulting function  $H(t)$ .



**Figure 32. Time Contribution Gain**

For the size component of our *content gain* function,  $K(s)$ , we use a single step staircase function with  $S_{avg}$  being the fragment size at which nodes start receiving a reward. Thus,  $K(s)$  is defined by

$$A_x^K = \{[0, S_{avg}), [S_{avg}, \infty)\} \quad (5.12)$$

and

$$\alpha_x^H = \{1.0, 1.5\} \quad (5.13)$$

For the popularity component of our *content gain* function,  $L(p)$ , we assume that files' popularity is Zip-f distributed. Let  $p_i$  be the CDF popularity of the item ranked  $i$  and let  $p_0$  be the popularity threshold at which nodes start receiving an incentive. The popularity gain function,  $L(p)$  is defined by

$$A_x^L = \{[0, p_0), [p_0, 1.0)\} \quad (5.14)$$

and

$$\alpha_x^L = \{1.0, 1.25\} \quad (5.15)$$

In addition to the individual formulations presented above for time, size and popularity, we also explore a combined metric for the size and popularity components:

$$KL(s, p) = \frac{1}{n_m} \sum_{\forall i} \sum_{\forall j} \chi_{i,j}^m \cdot \sum_{x=1}^N \Gamma_{Ax}(s, p) \quad (5.16)$$

where

$$\Gamma_{Ax}(s, p) = \begin{cases} g\_min & s^{p+p_0} / s_{avg}^{p+p_0} < g\_min \\ s^{p+p_0} / s_{avg}^{p+p_0} & g\_min \leq s^{p+p_0} / s_{avg}^{p+p_0} \\ g\_max & s^{p+p_0} / s_{avg}^{p+p_0} > g\_max \end{cases} \quad (5.17)$$

The factors in equation (5.17) allow us to define the minimum ( $g\_min$ ) and maximum ( $g\_max$ ) size-popularity gains. The value of  $p_0$  is used to determine at what point of the popularity's CDF the system switches its incentive from a penalty to reward (i.e.,  $p+p_0=1$ ).

We include this additional metric assuming that the time and storage contributions of nodes can be directly influenced by the incentives mechanism. Popularity on the other hand, is independent. In other words, each node could adopt a (selfish) strategy to maximize their benefit by defining maximum time and storage contributions. Furthermore, the amount of traffic associated with a particular item is determined by its size and its popularity. In other words, how often a file is requested and how much information is exchanged to serve these requests is

independent of the time component in our *contribution gain* function. Nonetheless, we want to explore whether integrating size and popularity into a single metric could allow us to steer the contributions of nodes to compensate for any potential bias caused by a skewed popularity distribution.

## 5.5 ADDITIONAL CONSIDERATIONS

Flexibility is a key feature of the PR redundancy method we present in this dissertation. The parameters of code-based mechanisms need to be known by all members of the overlay to enable file reconstruction by any node. This implies that these parameters are fixed and changing them to achieve different storage or reliability tradeoff would be costly if not unmanageable. In contrast, the hybrid redundancy mechanism that we propose in this dissertation allows us to change two of its parameters dynamically without affecting the underlying operation of the underlying code-based redundancy scheme. These parameters are the number of fragments replicated and the replication gain to be used for each of these fragments. The evaluation of the redundancy system presented in the Chapter 6.0 of this dissertation does not explore all the possibilities of adjusting these parameters. We leave this endeavor for future research.

For a redundancy mechanism to be compatible with a market economy, the system should be able to certify the contributions of individual peers. The incentive-based mechanism presents nodes with a tradeoff between their content contribution and the level of service they would get from other peers. Therefore, nodes will be tempted to exaggerate their contributions in order to gain additional performance. To avoid this form of cheating, the system should provide a light weight verification mechanism.

Providing proof of data possession would not be a difficult task in most redundancy methods<sup>29</sup>. An audit mechanism could be used to verify the validity of a node's contribution, counteracting the incentive to cheat. Proof of data possession could be accomplished using some form of self-certifying data for each fragment. Furthermore, this mechanism can also be used by content publishers to secure their contribution against tampering.

The use of redundancy in the system also contributes to the redistribution of requests, and their associated traffic, among the active members of the overlay, with possible performance gains. In that regard, a node performing a redundancy repair has an incentive to distribute any future demand for the object it is currently replicating. To minimize the negative effects of sharing (i.e., performing a repair) additional mechanism could be implemented. For example, prioritizing acknowledge messages over traffic, as suggested in [15]. (But this is beyond the scope of our current work.)

For nodes joining the overlay network with zero contribution, the system should provide an alternative mechanism for increasing their contribution without waiting for an *Index* node to select them as *Candidates*. In that regard, the system implementation that we evaluate in Chapter 6.0 allows nodes to request *Index* nodes for a list of fragments to be replicated. Once nodes obtain this list, they can initiate download requests for as many of items as provided by the *Index* node.

We believe that the structure of the mechanism presented in this dissertation does not preclude the use of additional mechanism to improve the performance of the system. Most incentives mechanism use short-term performance metrics, such as the bandwidth bartering mechanism of BitTorrent, as opposed to our mechanism, which uses a long-term metric. In that regard, the only system modification needed to allow two (or more) mechanisms to allocate

---

<sup>29</sup> Except in the case of the functional repair modalities of network coding, which are not considered in this work.

resources concurrently is the definition of a function to reserve a certain amount of resources for each mechanism. For example, one could assign 25% of the transmission bandwidth to content availability and the rest to another incentive mechanism. In addition, there are other mechanisms that can be added to our proposed system to improve its effectiveness. Gamification is one of them [90]. The contribution metric of the incentive-based mechanism can be presented graphically to the end-user as a rating or score. The purpose of such feedback would be to provide a positive psychological feedback to the end user when s/he is contributing to the content availability of the system. Here again, the effectiveness of such mechanism would depend on the end-user perception of the incentive, but given that social networks and internet gaming are part of our daily lives, we think that there is an opportunity for making this approach work.

## **6.0 EVALUATION**

The structure and operation of the PR redundancy scheme and the incentive-based mechanism we propose in this dissertation share a simple design philosophy. Nonetheless, their interplay in a heterogeneous and dynamic P2P environment is too complex to be evaluated analytically. As an alternative, this chapter presents an experimental evaluation. The layout of this chapter is as follows. Section 6.1 presents the software platform used in this evaluation. Section 6.2 describes a base system configuration, against which the alternative system configurations can be compared. Section 6.3 discusses some of the metrics obtained in the experiments and describes the presentation of results. Section 6.4 examines the impact of different system variants on the content availability of the network. Section 6.5 presents the system's response using different settings for its redundancy and incentives components and Section 6.5.3.2 describes the capacity of the proposed mechanisms to adapt to different networking conditions, such as varying churn rate and content space size.

### **6.1 SIMULATION-EMULATION PLATFORM**

Our redundancy-system has been implemented using Bamboo [21] and Modelnet [20]. Bamboo is our DHT application substrate and Modelnet is used to emulate wide area networking conditions using a cluster of computers. A description of Bamboo operation is presented in Section 2.2 and a

description of the *net-model* employed in Modelnet is presented in Section 3.3.3.1. At the routing layer, the experiments using Modelnet-Bamboo are actually emulations rather than simulations. Each packet generated by Bamboo's application stack is injected into the physical network and forwarded to another node through Modelnet's emulator node. At the application layer, it is unfeasible to transfer the actual content without overwhelming the capacity of the emulator or the physical network. Thus, the actual transfer of content is simulated rather than emulated. For the transfer of files, nodes only exchange control information. Throughout the rest of this section, no distinction will be made between the emulation and simulation component of the experiments and the term simulation is used in reference to both. All experiments are executed for 4.5 hours and replicated to obtain proper confidence intervals of the evaluation metrics being reported. The response variables of the system (described in Section 3.2) are recorded throughout the simulation (excluding the initial transient interval) and unless noted otherwise, the average values reported are for one single hour of operation of the system, between 2:45 and 3:45 hours after the start of the simulation.

## 6.2 BASE CONFIGURATION

The goal of the system we are proposing is to enable the creation of P2P systems that would work as open access content repositories. In a certain way, the overlay network would work as a public library whose content collection has been donated (and maintained) by its members. For example, our redundancy maintenance mechanisms could be used to build a P2P-based wikileaks [91] with better performance and resistance to censorship than the current web-based approach or to build a P2P-based Internet Archive [92] with improved scalability and lower cost of ownership.

The design of the system components (i.e., redundancy and incentives) we are presenting is quite simple. Nonetheless, their interactions in a dynamic and heterogeneous P2P setting is quite complex. The performance of the system can vary dramatically given different content and environmental conditions (e.g., churn). To manage this complexity and to guide our analysis of the proposed redundancy system, we define a reference system configuration, which we name the base configuration. Unless noted otherwise, the content and networking conditions of the system are assumed to be the ones presented in Table 24.

**Table 25. Base Content/Networking Configuration**

Parameter	Value	Description
<i>Unique Items</i>	500	Content space size
<i>Nodes</i>	1,840	Total network population
<i>Network Make</i>	10% <i>Benefactors</i> 90% <i>Peers</i>	Network population breakdown by class
<i>Content Contribution</i> <sup>30</sup>	<i>Benefactors</i> : 0.1 probability of having 1 whole-copy item	Initial content contribution for each node class
	<i>Peers</i> : 0.05 probability of having 6 whole-copy items	
<i>Churn</i>	<i>Benefactors</i> : Pareto(1.09,3060) <i>Cf</i> =0.75, median 100 min	Churn parameters for each node class: session length distribution and average node connectivity factor. Resulting median session length <sup>31</sup>
	<i>Peers</i> : initialized using Pareto(1.50,900) <i>Cf</i> =0.25, median 19.22 min	
<i>Redundancy</i>	Erasur Coding with $(k, S) = (6, 3)$ and proactive replication $r = 2$	Base redundancy method and maintenance strategy
<i>Incentives</i>	TB-Util with $\sigma_u=2.25$ , $\beta=0.05$ RP-Cost with $\sigma_c=3.5$ , $\chi=1.85$	Default parameter settings for incentive-based mechanism

When the overlay network operates without any of the content availability mechanisms we propose, the content contribution defined for the base configuration is such that the performance of the system exemplifies the content availability problems of P2P networks under churn.

<sup>30</sup> A detailed description of the initialization process, including the addition of fragments to the content storage of nodes is presented in Section 3.3.4.

<sup>31</sup> The average mean and median session lengths across 100 of the simulations performed are 34.45 and 18.69 minutes respectively.



In most graphs, a legend is added to specify the redundancy and content parameters used in each setup. For example, “EC(6, 3) 500 Items” indicates that the MDS erasure coding component of the redundancy system uses the parameters  $(k, S)=(6, 3)$  and that the content space size includes 500 unique items. If the redundancy and/or content parameters are not specified, the values presented in Table 24 must be assumed.

### 6.3 PERFORMANCE METRICS

We use two fundamental metrics to characterize the performance of the system, cost and efficiency. The terms cost and efficiency are used, respectively, to describe the amount of traffic generated to maintain the redundancy state information and the effectiveness of the content retrieval process. Cost is measured in bytes per second per file and efficiency is measured as a percentage.

Unless noted otherwise, the cost metric used in this section is the per file redundancy repair cost rate (7<sup>th</sup> column of the line *PUB-Bp* in the *Per\_f\_rates* section in Figure 33). Since the level of involvement of nodes in the redundancy maintenance process can be highly heterogeneous, we believe that the per file cost metric is more adequate to compare the behavior of different system settings.

The simulation’s processing script produces multiple cost results. The results labeled with “#” are a count of the number of repairs, and the results labeled with “B” report bytes transferred. The results without the “p” suffix indicate that measurements are taken throughout the whole simulation. Conversely, the results with the “p” suffix indicate that measurements are taken only during a measurement window of one hour in the second half of the simulation (more specifically,

between 2:45 and 3:45 hours). The number of repairs made during the simulation provides a count for the total number of repairs (T), the successful repairs (G) and the number of failed repairs (F). Besides the absolute count for each of category (columns 1-3), average values are obtained using the mean overlay network size during the simulation (columns 4-6) and the total number of unique files in the simulation (columns 7-9).

Method: HYBRID n=6 m=18 r=2										
Avg File Size [Tot,G,F,Gp]=3931.21 3835.11 5245.31 3841.51 kbytes										
Tot Tx_bytes=23189 Mbytes, Avg Tx_time[T,G,F]=89,72,331 sec										
	Total	Good	Fail	Net_rates [T,G,F]<#,kbytes>			Per_f_rates [T,G,F] <#,kbytes>			
PUB#	36242	33772	2470	148.02	137.93	10.09	0.30	0.28	0.02	#/sec
PUB-B	23745848	21586530	2159318	96983.44	88164.29	8819.15	193.97	176.33	17.64	kB/sec
PUB#p	7217	6665	552	120.28	111.08	9.20	0.24	0.22	0.02	#/sec
PUB-Bp	4760589	4267279	493309	79343.16	71121.33	8221.83	158.69	142.24	16.44	kB/sec

**Figure 33. Sample Simulation's Cost Results Output**

The efficiency metric, namely download rate, can be calculated using different measurements. This is illustrated in Figure 31. The simulation's processing script produces results for five metrics:

1. Total content lookups (*Search* column)
2. Lookups returning a root entry (*Finish* column)
3. Lookups that fail (*Fail* column)
4. Non-empty content responses from Index nodes (*QryR[x]* column)
5. Number of successful content retrievals (*ContSuc* column)

The *QN* line presents a total count for each metric throughout the simulation and the following lines indicate the effectiveness metric for each category. That is, the ratio of each metric over a reference metric. Line *Q1* uses the *Search* count as reference for all calculations; line *Q3* uses the *Finish* count, and line *Q4* uses the *QryR[x]* count to obtain the effectiveness metric. The purpose of the last two metrics is to analyze the efficiency of different sub-processes. For example, line *Q4* calculates the efficiency metric using the non-empty index requests, (which for the sample

values presented in Figure 34 is 73% of the total) resulting in a efficiency metric of 80% for the content retrieval process (instead of 59%, which is obtained using the total count).

Simul wide values					
Q	Search	Finish	Fail	QryR[x]	ContSuc
QN	28486	28208	446	20811	16701
Q1%	1.00	0.99	0.02	0.73	0.59
Q3%		1.00	0.00	0.74	0.59
Q4%				1.00	0.80
Analyzing Queries during [8820...12420] sec simul time					
Qp	Search	Finish	Fail	QryR[x]	ContSuc
QNp	6646	6610	101	5104	4309
Q1p%	1.00	0.99	0.02	0.77	0.65
Q4p%				1.00	0.84

**Figure 34. Sample Simulation's Efficiency Results Output**

The alternative efficiency metrics presented above allow us to focus our development efforts on those sub-processes with the highest potential to improve the performance of the system. In that regard, it is clear that the lookup process of Bamboo is highly efficient (*Finish* column in Figure 34). Consequently, our development efforts were focused on content indexing, replication and retrieval algorithms.

All the efficiency metrics presented in the rest of this section are obtained using the total number of *Search* requests generated (first metric shown in line *QN*, Figure 34). We selected this metric for two reasons. First, to avoid efficiency figures (e.g., line *Q4p*) that would reflect an artificially high system performance. Second, to be able to compare different system settings consistently. The network is able to complete at least 98% of the lookup operations across all simulations, but in the best case scenario, only 80% of these requests culminate with a non-empty targets-list message back from *Index* nodes. Thus, to visualize the maximum indexing efficiency achieved by the system (column *QryR[x]* in Figure 34), the top scale in all the efficiency plots is set to 0.8. As a result, the scale used in the plots of the efficiency metric is absolute, but its

graphical representation also conveys its relative performance with respect to the maximum efficiency of the content indexing mechanism.

Section 3.2 describes that to characterize the behavior of a system under churn we require three evaluation metrics: cost, performance and efficiency. Most of the results to be presented in the following sections are content related, and in that context, we do not require the lookup's performance metric. Nonetheless, this metric has been recorded and verified not to exceed a 500 milliseconds threshold for all simulations. In fact, the median lookup delay does not exceed 270 milliseconds in any of the simulations performed.

## **6.4 REDUNDANCY EFFECT**

The first research question that we need to address to determine the relevance of our proposed redundancy maintenance system is what are the positive and negative effects of using redundancy in a P2P system under churn? Figure 35 illustrates both of these effects. On the positive side, the use of redundancy substantially improves the efficiency of the system (i.e., content retrieval). On the negative side, the use of redundancy requires additional transmission bandwidth and storage resources. We consider that the negative effects are unavoidable. Then, the next question is whether these negative effects can be controlled and to what extent. The analytical results presented in Section 4.3.2 indicate that our proposed PR redundancy scheme consume less resources than MDS erasure coding (and even exact-MBR network coding). The system evaluation presented in this chapter is a proof of concept of the effectiveness and flexibility of the PR redundancy scheme.

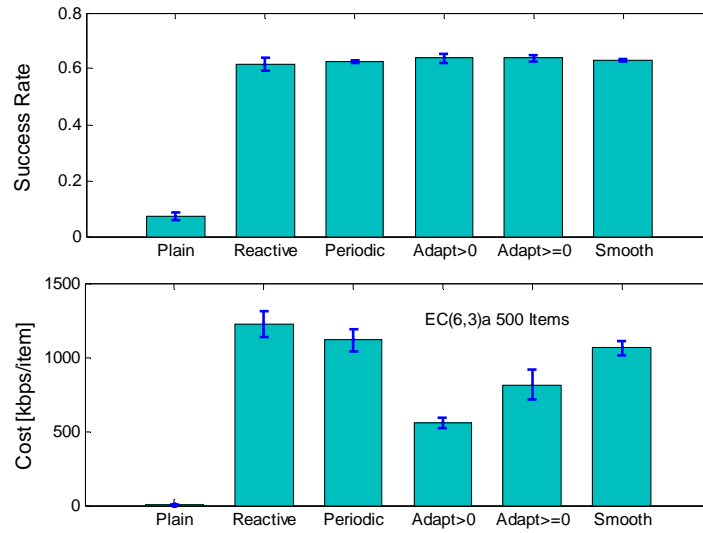
Figure 35 contains two sets of results. The first set, top graph, presents the download rate metric (i.e., efficiency), and the second set, lower graph, illustrates the average cost metric of the system. These cost and efficiency metrics correspond to six different system configurations labeled as follows: Plain, Reactive, Periodic, Adapt $>0$ , Adapt $\geq 0$  and Smooth. The first system configuration, Plain, does not use any form of redundancy. This represents a system with content availability problems. The second system variant, Reactive, corresponds to a maintenance policy that repairs any fragment lost as soon as it is detected. The third system variant, Periodic, uses a fix-epoch redundancy-maintenance policy. That is, repairs are spaced in time by a constant interval. The fourth system variant, Adapt $>0$ , uses an adaptive redundancy-maintenance mechanism that increases its maintenance epochs when the redundancy completeness<sup>32</sup> metric increases. The fifth system configuration, Adapt $\geq 0$ , is an adaptive redundancy-maintenance mechanism that increases its maintenance epochs when the redundancy completeness metric remains stable or increases. Finally, the sixth system configuration, Smooth, uses a smoothed average algorithm to adapt the system's maintenance epochs. The pseudo code for the alternative adaptation algorithms used by the redundancy-maintenance process are presented in Algorithm 5 and Algorithm 6, Section 4.4.3.

The results presented in Figure 35. Effect of Redundancy on Content Availability clearly show that redundancy considerably improves content availability in P2P networks. Without redundancy, the system only completes 7.2% of the content requests; with any of the redundancy-maintenance modalities the system achieves a download success rate above 62%. The difference between the alternative redundancy modalities is the cost to be paid for this improved content availability. The reactive repair policy consumes the most bandwidth, followed by the periodic

---

<sup>32</sup> This metric is a value between 0 and 1 that measures the availability of unique fragments.

repair policy (i.e., fixed maintenance epochs). The three adaptive maintenance mechanisms use a proactive repair policy that achieves different reductions on the system cost metric with respect to the periodic system modality: Adapt>0 scheme reduces cost by 50%, Adapt>=0 by 27% and Smooth does it by only 5%. The adaptive maintenance epoch algorithms presented in this dissertation prove the advantages of this feature but should not be taken as optimal. We believe that additional system considerations should be taken into account to construct an ideal adaptation algorithm, but this goes beyond the scope of our current evaluation.



**Figure 35. Effect of Redundancy on Content Availability**

## 6.5 PARAMETERS EXPLORATION

The redundancy-system presented here possesses two components that improve content availability in churning P2P networks. The first is redundancy, and the second, is an incentive-based mechanism. Both of these mechanisms have a set of parameters that can determine its

effectiveness in improving content availability. The following sections explore an array of parameters settings for these mechanisms and analyze their impact on the system behavior.

### **6.5.1 Proactive Replication Redundancy Scheme**

Chapter 4.0 presented an analytical discussion of the maintenance cost properties of alternative redundancy schemes, particularly MDS erasure coding and exact-MBR network coding (Section 4.3.2.) These analytical results indicate that MDS erasure coding is not adequate for P2P environments due to its high repair cost requirements.

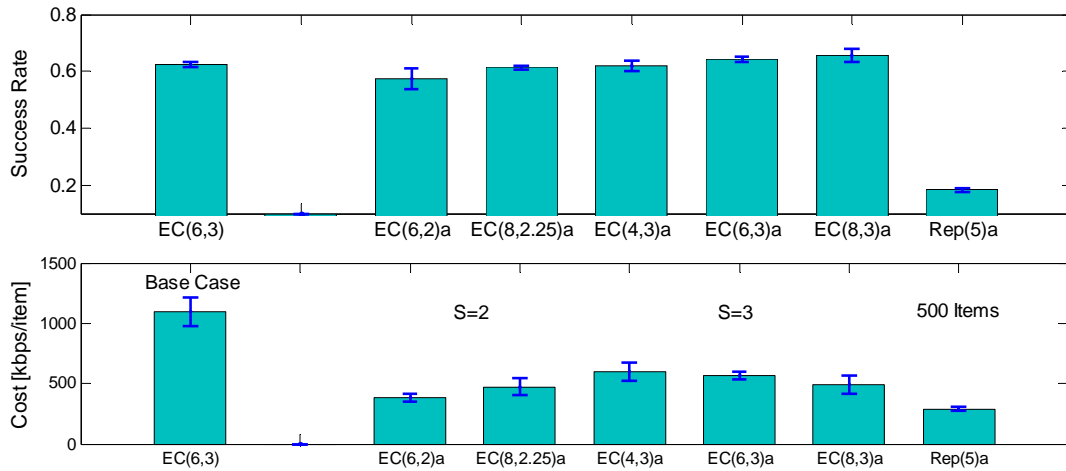
Section 4.1 introduced our proposed hybrid redundancy scheme, Proactive Replication (PR), which combines the reliability properties of MDS erasure coding with the simplicity of replication to achieve low redundancy repair cost. The analytical results represented in Chapter 4.0 are a simplified version of the complex redundancy maintenance process that occurs in a P2P network. Thus, we rely on simulations to make a more realistic assessment of the capacity of a redundancy scheme to improve the content availability of a P2P network. For example, the analytical results assume an ideal content registration process. That is, the content indexing information is assumed to be accurate all the time. During the experimental evaluation, this simplification is removed.

Figure 36 presents the cost and efficiency results for several settings of the PR redundancy as well as for replication redundancy. For PR, the graphs show different parameters settings of its erasure coding component. All settings use the Adapt>0 maintenance policy, except for EC(6,3), which uses a periodic maintenance policy and is included as reference only.

The results clearly indicate that PR redundancy outperforms replication redundancy. The download success rate of replication redundancy (last bar to the right) is quite low, below 20%. In

addition, the maintenance cost of replication redundancy is highly inefficient. At least 25% of the maintenance bandwidth consumed by replication redundancy corresponds to failed repair attempts. In contrast, failed repairs in the system using PR redundancy are always below 7%. This is a clear indication that the amount of data transferred during each repair operation has a significant effect on the efficiency of a churning system. In addition, the download rate (i.e., efficiency) of the proactive redundancy method is superior to replication redundancy because nodes take advantage of concurrent downloads. This is a well known performance improvement mechanism used in deployed P2P file-sharing networks.

All the variants of the PR redundancy scheme presented achieve similar efficiency metrics, but they differ on the amount of maintenance bandwidth consumed.



**Figure 36. Cost and Efficiency for Alternative PR Redundancy Settings**

The experimental settings presented in Figure 37 for PR redundancy illustrate the behavior of the system for various combinations of the parameters  $k$  and  $S$ . Parameter  $k$  is the reception efficiency of the MDS erasure coding component and parameter  $S$  is its coding gain; which determines the total number of unique fragments generated by the coding scheme (i.e.,  $N=S*k$ ). The replication component of PR, parameter  $r$ , uses a constant value of 2 in all the experimental



settings. Table 25 list the success rate and cost results obtained in each experimental setup, together with two additional system attributes, the system's mean maintenance epoch and mean repair rate.

Table 25 also includes results for one system variant not included in Figure 37. This variant,  $EC(6,3)a-7$ , replicates only seven out of the 18 total fragments defined by the erasure coding component of PR redundancy. These results suggest that there are significant cost savings when not all fragments are replicated and furthermore, the effectiveness of the system does not deteriorate with this measure. This is certainly an interesting line of experimentation for future research. For instance, different fragment sets could be replicated using various priorities in order to find different cost/efficiency tradeoffs.

**Table 26. Response Variables for Alternative Redundancy Schemes**

System Setting	Download Rate [%]	Cost [bps/file]	Average Maintenance Epoch [seconds]	Average Repair Rate [repairs/file/sec]
$EC(6, 2)a$	0.572	386.14	148.68	0.67
$EC(8, 2.25)a$	0.613	473.51	148.97	1.01
$EC(4, 3)a$	0.618	596.18	166.54	0.65
$EC(6, 3)a$	0.638	559.15	156.93	1.03
$EC(8, 3)a$	0.656	487.07	155.38	1.17
$EC(6, 3)a - 7$	0.635	349.28	161.75	0.51

In Section 4.3.2.1 we argued that PR redundancy with small parameter values, namely for  $k$  and  $S$ , produce the smallest redundancy state information (i.e., set of fragments). However, this redundancy state information needs more frequent maintenance than other settings with larger overheads (i.e.,  $S$  parameter). Moreover, when longer maintenance intervals are used, the methods with larger overhead values consume the same maintenance bandwidth as methods with smaller parameters at shorter maintenance intervals. We found that the analytical results presented in

Section 4.3.2.1 apply only partially to our experimental results. The following sections describe our observations with respect to the effects of parameter  $S$  first, and with respect to the effects of parameter  $k$  second. The parameter being discussed in each section appears in boldface in the text.

#### **6.5.1.1 PR Redundancy with Different Coding Gains ( $S$ )**

The net effect of augmenting the coding gain parameter,  $S$ , is to increase the total number of nodes storing an item. That is, the redundancy-set. The experimental results indicate that larger values of  $S$  improve download rate, but they also increase the system's cost metric. For example, EC(6,3)a's efficiency is higher than EC(6,2)a by 7%, but EC(6,2)a consumes 31% less bandwidth.

According to our analytical results, (Section 4.4.3) by increasing  $S$ , we expect the system's maintenance epoch to increase<sup>33</sup> while others metrics would remain constant. In the experimental evaluation we observed that the maintenance epochs increase for larger  $S$  values, but so does the maintenance cost. For example, the average redundancy maintenance epoch increases from 149 to 157 seconds when  $S$  increases from 2 (in EC(6,2)a) to 3 (in EC(6,3)a).

Table 26 includes an alternative metric for the efficiency of the redundancy maintenance process, average repair rate per file. In the experiments we performed, larger  $S$  values increase this metric. For example, EC(6,2)a has a repair rate of 0.67 and in EC(6,3)a this metric increases to 1.03. Thus, maintaining a single file becomes more costly for larger  $S$  values. Nonetheless, it is also clear that this additional cost increases the efficiency of the system; from 0.57 for EC(6,2)a to 0.64 for EC(6,3)a.

---

<sup>33</sup> The analytical results suggest that larger overhead values can be used to maximize reliability without significantly increasing the system's maintenance cost. The only requirement is the use of larger maintenance epochs.

### 6.5.1.2 PR Redundancy with Different Reception Efficiencies ( $k$ )

The value of the reception efficiency parameter,  $k$ , determines the number of nodes needed to reconstruct a file as well as the size of each fragment. The experimental results presented in Figure 36 illustrate that increasing  $k$  reduces cost and increases efficiency. For example, EC(8,3)a has a lower cost and better efficiency than EC(6,3)a.

With respect to the number of repairs, increasing the value of  $k$  increases the average system repair rate. The coding gain in EC(8,2.25)a was selected to have exactly the same total number of nodes as EC(6,3)a; that is, 18 nodes. Both configurations use the same total number of nodes to store a single file, but their fragments are of different sizes. The repair rate of EC(8,2.25)a is slightly lower than the repair rate of EC(6,3)a. If this metric is adjusted according to the value of  $k$  (i.e., divide it by  $k$ ), the resulting cost values are 0.17 for EC(6,3)a and 0.13 for EC(8,2.25)a. Consequently, these results corroborate that for a fixed value for the total number of nodes to store a file, higher values of  $k$  reduce the average amount of information exchanged during repairs and should be preferred. However, increasing  $k$ , while  $N$  is kept constant, causes the availability of the redundancy scheme to degrade.

## 6.5.2 Redundancy Maintenance

Adaptability is without a doubt a desirable system property; especially for heterogeneous environments such as P2P networks. While a small set of highly available nodes could maintain the availability of a specific file, there is also the probability that correlated node departures could cause a specific content to suffer sudden availability transitions. In that situation, applying exactly the same redundancy-maintenance mechanism for both scenarios is inefficient. For the mechanisms proposed in this dissertation, the redundancy-maintenance interval is the fundamental

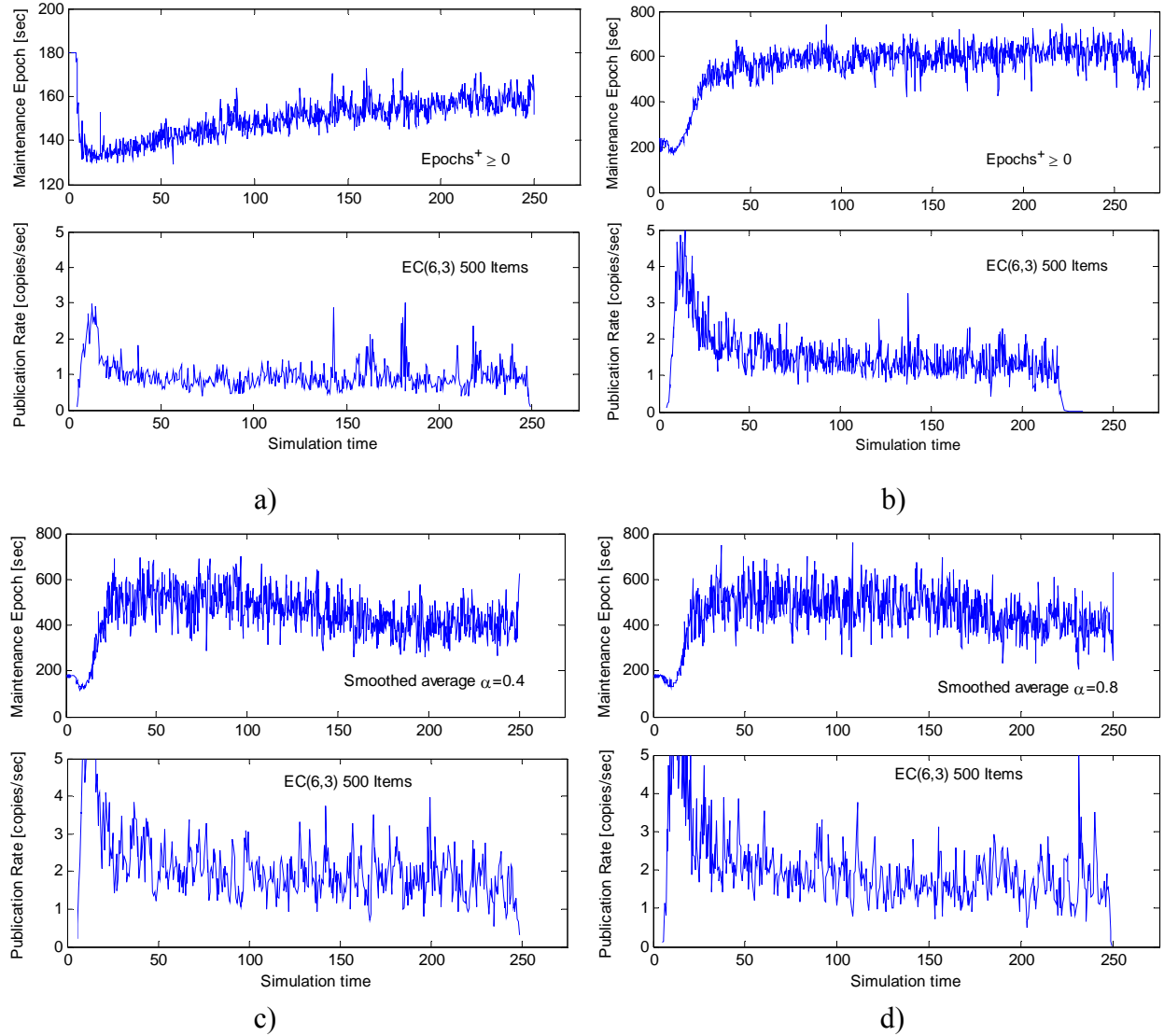
dynamic parameter that adapts the operation of the system to different node membership configurations. The following subsections describe two alternative adaptation mechanisms and the role of the parameters TARGET and MIN\_SEG in the operation of the system.

### 6.5.2.1 Adaptive Maintenance Epochs

Figure 37 includes four sets of results. The top graphs present the evolution of the average maintenance epochs, and the graphs at the bottom illustrate the average redundancy repair rate in the system (i.e., publication rate). The first row of graphs corresponds to the two variants of Algorithm 6 and the second row of graphs corresponds to two variants of Algorithm 5. The graphs in Figure 37.a (left column) correspond to the adaptation mechanism labeled earlier as  $\text{Adapt}>0$ , and the graphs in Figure 37.b (right column) to  $\text{Adapt}\geq 0$ . The graphs in Figure 37.c (second row, left column) correspond to the adaptation mechanism labeled earlier as Smooth, with  $\alpha=0.4$ , and the last set of graphs correspond to the same Smooth algorithm, but with  $\alpha=0.8$ .

The only algorithmic difference between  $\text{Adapt}>0$  and  $\text{Adapt}\geq 0$  is that if two consecutive evaluations of the redundancy completeness metric are equal,  $\text{Adapt}\geq 0$  increase the maintenance epoch by 2.5 and  $\text{Adapt}>0$  does not. The resulting effect is almost a factor of four increase in the average maintenance epoch for  $\text{Adapt}\geq 0$  compared with  $\text{Adapt}>0$ . As a result, the  $\text{Adapt}\geq 0$  system variant also produces a higher average number of repairs. Thus, this system configuration consumes more resources than the alternative  $\text{Adapt}>0$  repair policy. For the Smooth adaptation algorithm, the two  $\alpha$  values used to obtain the smoothed average produce a similar effect on the average redundancy maintenance of the system and its respective repair rate. In terms of cost and efficiency, the Smooth algorithm performs worst than  $\text{Adapt}>0$ . Consequently,  $\text{Adapt}>0$  is the default setting for the rest of our experimental work. After all, previous research has highlighted

that one added benefit of proactive redundancy maintenance mechanisms is that the maintenance bandwidth usage of each node becomes predictable [66].



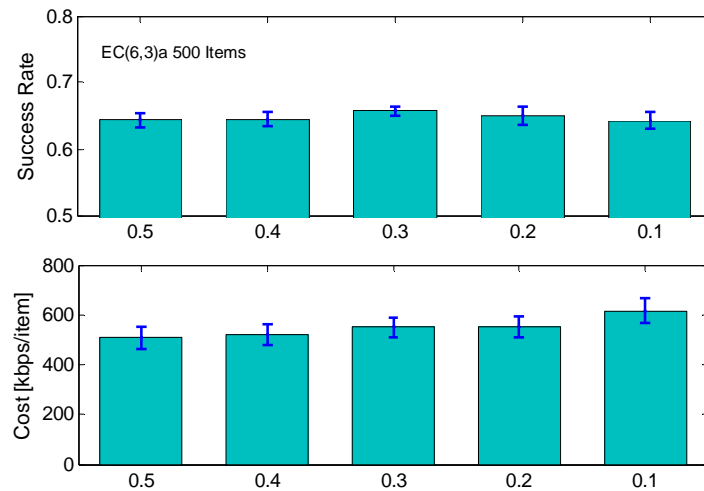
**Figure 37. Adaptive Maintenance Epochs and Publication Rate: a)  $\text{Adapt} > 0$ , b)  $\text{Adapt} \geq 0$ , c) Smooth with  $\alpha = 0.4$  and d) Smooth with  $\alpha = 0.8$**

### 6.5.2.2 TARGET Fragment Availability

The analytical evaluation of the redundancy maintenance cost presented in Section 4.3.3 indicates that for a redundancy method with fixed parameters, the redundancy maintenance cost decreases for longer maintenance intervals. However, increasing the maintenance epochs also deteriorates

the reliability of the redundancy mechanism. To achieve maximum efficiency without jeopardizing the reliability of the redundancy mechanism, the  $\text{Adapt} > 0$  maintenance epoch adaptation mechanism continuously adjusts the system maintenance interval. The parameter TARGET in Algorithm 6 (Section 4.4.3) determines the amount of redundancy gain/loss that triggers a certain maintenance epoch adjustment. For example, when the system increases its redundancy, Algorithm 5 increases the maintenance epoch by 5 seconds if the gain is greater than  $\text{TARGET}/2$  or by 2.5 seconds otherwise. In other words, the TARGET parameter controls the rate and size of adjustments in the maintenance epoch adaptation mechanism.

Figure 38 presents the efficiency and cost metrics for the system using several TARGET values. The upper graph presents the efficiency metric and the bottom one the cost metric. The efficient values do not change significantly among the different parameter settings listed (there is only a 2% difference), and for the cost, longer values of the TARGET parameter result in a lower cost metric.



**Figure 38. Cost and Efficiency for Different TARGET Values**

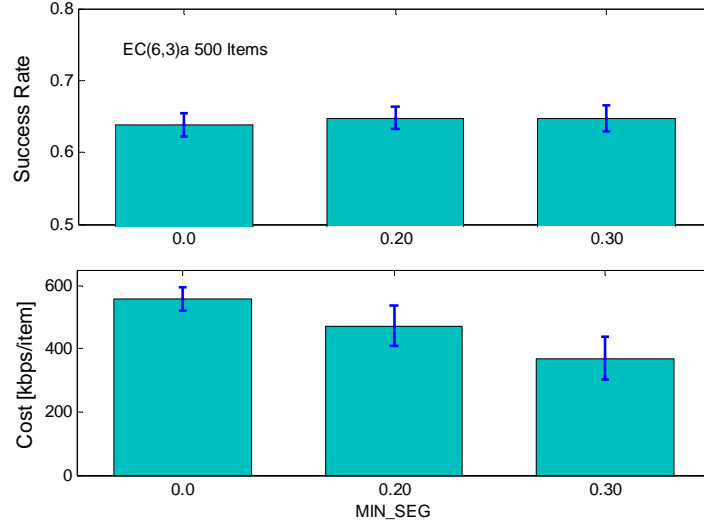
The smallest value of TARGET exhibits the highest maintenance cost metric and the lowest efficiency. Therefore, this extreme parameter value should be avoided. Increasing the value

of TARGET improves both metrics, but efficiency starts to decay after 0.3. For the rest of the evaluation of the redundancy system a value of 0.3 for the TARGET is assumed. This is a conservative setting to preserve the reliability of the system.

### 6.5.2.3 Minimum Number of Segments Available (MIN\_SEG)

Line 3 in Algorithm 4 (Section 4.4.3) controls the adaptation of maintenance epochs based on whether the number of unique segments available is above the MIN\_SEG threshold or not.

Figure 39 presents the cost and efficiency metrics for the system using different values for this parameter. The first column represents an aggressive content availability approach. Replicas are created whenever possible regardless of the current availability of the file (i.e., total number of unique fragments available). This is the default for all the experimental settings presented, unless noted otherwise. In terms of efficiency, there is not a significant advantage of using either variant, but with respect to cost, preventing the generation of new replicas when no enough fragments are available for a single file (i.e. less than  $k$  unique fragments exist) represents a significant cost savings measure. For example, if instead of replicating fragments all the time (i.e. MIN\_SEG=0.0) the system uses a MIN\_SEG threshold of 0.3, cost is reduced by 34%. More importantly, this modification does not degrade the efficiency metric. Notice that the presence of a complete copy of a file implies the availability of all its segments. Thus, the redundancy maintenance mechanism would trigger a maintenance process when there is at least one complete copy of a file or when the number of nodes storing unique fragments is above  $\text{MIN\_SEG} * m$ .



**Figure 39. Cost and Efficiency for Different MIN\_SEG Values**

The cost savings of using a MIN\_SEG value other than 0.0 can also be measured in terms of the average repair rate in the system. For instance, for MIN\_SEG=0.3, the average repair rate is 0.51 and for MIN\_SEG=0.0 the average repair rate is 0.92. With respect to the average maintenance epoch, the value of the parameter MIN\_SEG does not produce significant differences among the settings presented. For example, for MIN\_SEG=0.3 the average maintenance epoch is 155.3 seconds, compared with 154.9 for MIN\_SEG=0.0.

### 6.5.3 Incentive-Based Mechanism

The incentive-based mechanism serves various purposes in the redundancy-maintenance process. First, to promote nodes participation, second, to facilitate self-organization and third, to regulate fair exchange of resources. The last property is also enforced for exchange of resources in processes other than redundancy-maintenance, such as files downloads.

Using different variants of the incentive-based mechanism should not have a significant impact on the cost metric of the system since repairs are mainly determined by churn rather than



cooperation. For the system's efficiency on the other hand, the positive and negative feedback components of these variants (presented in Chapter 5.0) can produce both a positive and negative effects.

We understand that the success or failure of any explicit or implicit incentive-based mechanism depends heavily on the actual perception of the incentives by the system's participants [93]. In that regard, our evaluation does not reflect the perception of end-users towards the proposed incentive-based mechanism. Still, the results illustrate the potential capabilities of the proposed mechanism.

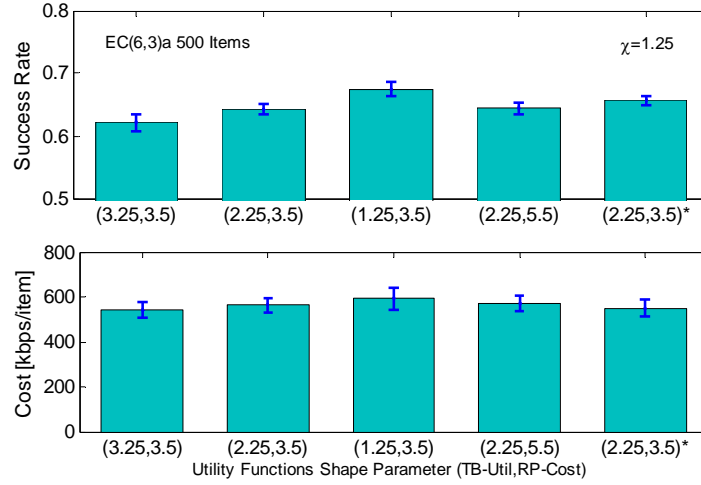
This section, examines different parameter settings for our incentive-based mechanism for content availability. First, we examine the use of different parameters settings for the sigmoid functions used in our mechanism and their effects on the behavior of the system. Second, we analyze the effects of different variants of the *contribution gain* function (presented in Section 5.2.4).

### 6.5.3.1 Sigmoid Functions Parameters

In this section, the content contribution metric is assumed to be the number of fragments stored by nodes. That is, we assume that the system uses an unitary *contribution gain* function (i.e., equation 5.4,  $G(\chi^m, t, s, p)=1$ ).

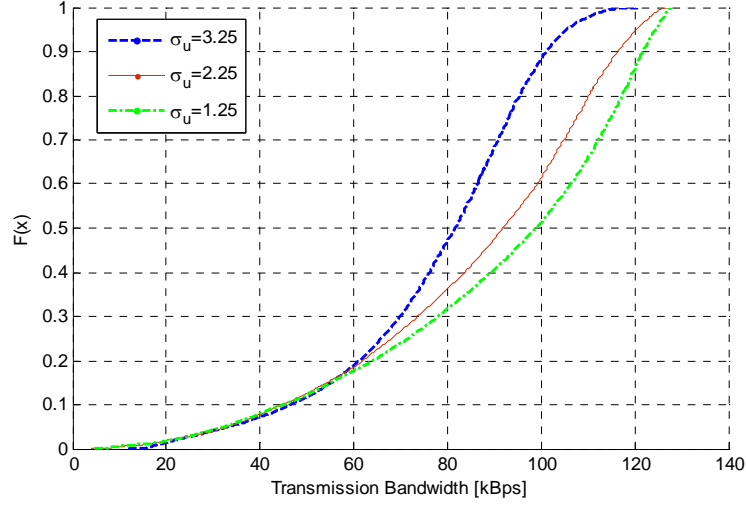
Figure 40 presents the cost and efficiency metrics for different parameter settings of the sigmoid functions of the incentive-based mechanism. Each bar's label (e.g., (3.25,2.5)) indicates the value of the shape parameter for the TB-Util and RP-Cost functions, respectively. The RP-Cost shift parameter,  $\vartheta$ , is equal to 1.85 for all the results presented, except for the last one, which uses 1.25 instead. These results corroborate that the incentive-based mechanism does not determine the

cost of the redundancy maintenance process. The efficiency of the system on the other hand, presents slight variations, but we do not consider this to be a negative effect. That is, this effect on the system's efficiency is one of the intended effects of the incentive-based mechanism. If nodes decide not to participate in the redundancy process, there must be negative consequences.



**Figure 40. Cost and Efficiency for Different Incentive-based Mechanism Parameters**

For the TB-Util function, the experimental results indicate that its shape parameter can determine the efficiency metric of the system. It produces the lowest efficiency metric for  $\sigma_u=3.25$  (first bar) and the highest for  $\sigma_u=1.25$  (third bar). Higher efficiency for lower values of the shape parameter is not surprising because  $\sigma_u$  determines the allocation of bandwidth among peers for all data transfers (i.e., repairs and file downloads).

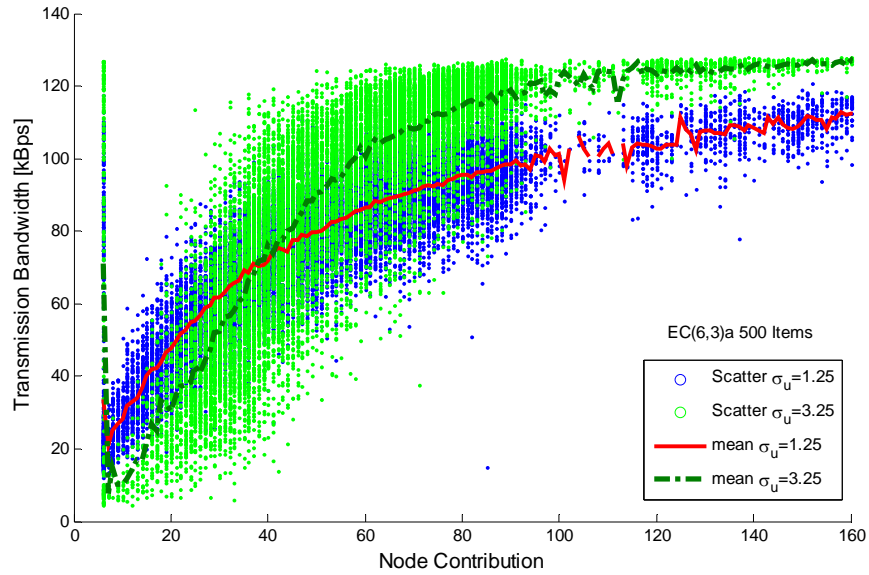


**Figure 41. Transmission Bandwidth CDF for Various TB-Util Shape Parameter Values**

Figure 41 presents sample CDFs of the transmission bandwidths of three simulations settings with different  $\sigma_u$  values. The median transmission bandwidth for  $\sigma_u=1.25$  is 100 kBps and only 80 kBps for  $\sigma_u=3.25$ . Before nodes' content contribution reaches the target average value ( $r_{tgt}$ ), the TB-Util function assigns larger transmission bandwidths sooner (i.e., for smaller contributions) for  $\sigma_u=1.25$ , than for  $\sigma_u=3.25$ . In addition, when nodes' content contribution surpasses the target content contribution, the TB-Util function increases the nodes' payoff slower for  $\sigma_u=1.25$ , than for  $\sigma_u=3.25$ . We believe that the system should sacrifice efficiency (i.e. for nodes with smaller contributions) to avoid nodes contributing only small amount of resources to the system. Thus, small  $\sigma_u$  values should be avoided. The experimental results presented in Figure 42 illustrate the reasoning for this design guideline. This figure presents a scatter diagram of the transmission bandwidths assigned to nodes for each of their download requests versus their contribution. In addition, this figure presents the resulting average transmission bandwidth for each node contribution value. Nodes with smaller contributions (e.g., less than 30) receive better transmission bandwidths when the system uses  $\sigma_u=1.25$  than when  $\sigma_u=3.25$ . Conversely, nodes with larger contributions (e.g., above 80) reach smaller average transmission bandwidths for

$\sigma_u=1.25$  than for  $\sigma_u=3.25$ . Consequently, in order to promote nodes committing more than just the target contribution, the system should use larger  $\sigma_u$  values.

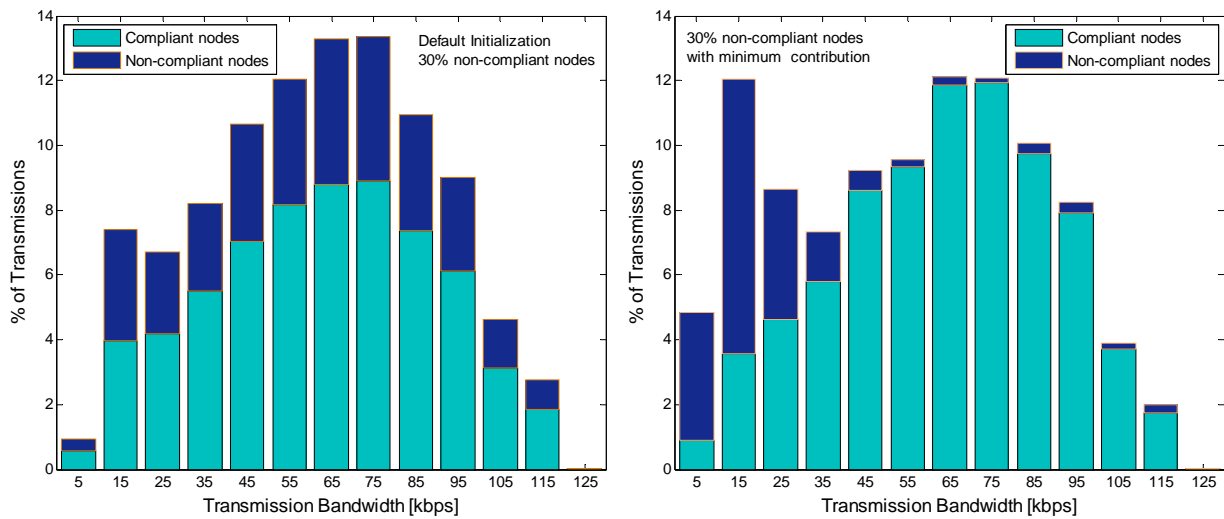
The incentive-based mechanism uses a probabilistic approach for fair allocation of resources. The results in Figure 43 clearly indicate that nodes with higher content contributions receive better service from their peers. Nodes with smaller contributions can obtain good transmission bandwidths from other peers, but the probability of such events is lower than for nodes with higher contributions.



**Figure 42. Transmission Bandwidth for different TB-Util Shape Parameter Values (Scatter Diagram)**

The incentive-based mechanism promotes the participation of nodes in the redundancy repair process in the form of an improved future performance. That is, nodes participate in good faith in the redundancy maintenance process expecting to improve their performance during subsequent participations in the system. Using an economic analogy, the incentive mechanism promotes investments (in the form of storage) given the promise that nodes will receive an economic return in the future (in the form of bandwidth).

The incentives-based mechanism defines a mechanism for fair exchange of resources among members, which is independent of the redundancy maintenance process. That is, nodes are assigned resources in proportion to their current content participation in the system. Even if nodes decide not to participate anymore in the redundancy maintenance process, they get a fair share of resources. This feature is illustrated in Figure 43, which presents two experimental scenarios where 30% of the nodes do not participate actively in the redundancy maintenance process. In the first scenario (left graph), the non-compliant nodes (i.e., nodes that are not performing any redundancy maintenance activity) have a wide array of contribution values. Consequently, the transmission bandwidth they receive is also spread across the full spectrum of transmission bandwidths. In the second scenario (right graph), the non-compliant nodes have only minor content contributions. As a result, the transmission bandwidth they receive from other nodes is at the lower end of the spectrum. Thus, if we assume that the problem of free-riders in P2P networks is associated with nodes that do not contribute content, our incentive-based mechanism does a good job at limiting their impact on the performance of the system. We explore this feature further in the following section.



**Figure 43. Fair Allocation of Resources in Incentive-based Mechanism**

### 6.5.3.2 Contribution Gain Function

The enhanced *contribution gain* function (described in Section 5.4) gives us greater flexibility to promote the cooperation of nodes (to improve content availability) as well as self-organization.

This portion of our experimental work is performed with non-default parameters. In particular, *Peers'* churn is initialized using Pareto(1.5, 1,350) and for MIN\_SEG we use a value of 0.3.

For the different components of the *contribution gain* function we use the following settings. For the time component, all nodes use the same staircase function (equations 5.10 and 5.11, presented in Section 5.4). For the size and popularity components, we use four parameters settings, which are summarized in Table 27.

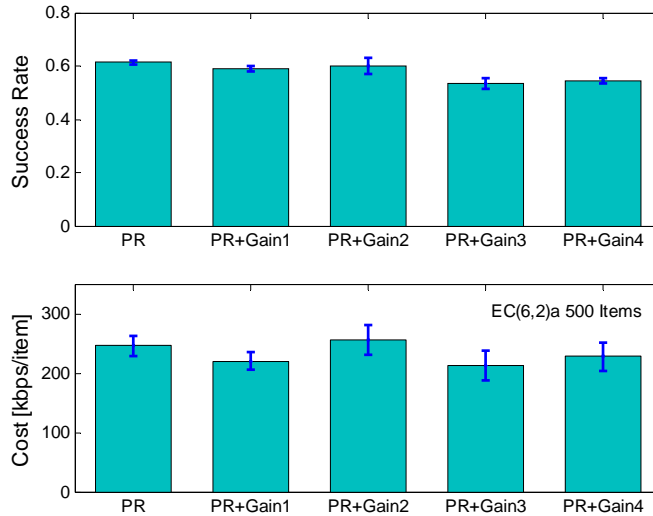
**Table 27. Size and Popularity Components in Contribution Gain Functions**

Set Name	Parameters Values
Gain1	$K(s): A_x^K = \{[0, 4\text{MB}/k), [4\text{MB}/k, \infty)\}$ and $\alpha_x^H = \{1.0, 1.5\}$ $L(p): A_x^L = \{[0, 0.8), [0.8, 1.0)\}$ and $\alpha_x^L = \{1.0, 1.25\}$
Gain2	$K(s): A_x^K = \{[0, 4\text{MB}/k), <s_i>, [12\text{MB}/k, \infty)\}$ and $\alpha_x^H = \{0.5, <s_i/4\text{MB}/k>, 3.0\}$
Gain3	$KL(s,p): \Gamma_{Ax}(s, p)$ with $s_{avg}=4\text{MB}/k, p_0=0.2, g_{min}=1.0, g_{max}=2.0$
Gain4	$KL(s,p): \Gamma_{Ax}(s, p)$ with $s_{avg}=4\text{MB}/k, p_0=0.8, g_{min}=0.5, g_{max}=3.0$

The time component used in all *contribution gain* function variants provides both penalties and rewards. If nodes' sessions are shorter than 2 minutes, their time *contribution gain* is 0.5 (see Figure 32 in Section 5.4). In other words, they are penalized when their sessions are short. Only nodes' sessions lasting longer than four minutes receive a reward.

For the size and popularity components of the *contribution gain* function, Gain1 and Gain3 are reward-only schemes while Gain2 and Gain4 are penalty-reward schemes. Gain1 and Gain3 provide and incentive when nodes' storage size or popularity are above a predefined threshold.

Gain2 and Gain4 schemes on the other hand, include a penalty component. Gain2 uses a unitary popularity gain though. For the size contribution gain, nodes storing less than 4MB/k receive a contribution gain penalty of 0.5, and as nodes increase their average fragment size contribution, they receive larger size contribution gains, up to 3.0 when their average fragment size reaches 12MB/k. In the Gain4 scheme, the popularity and size content gains are merged. Nodes receive a penalty when  $s^{p+p_0}$  is less than  $s_{avg}^{1.6}$  (where  $s_{avg}=4\text{MB}/k$ ) and a reward otherwise, up to a maximum size-penalty *contribution gain* of 3.0. We included a variant without popularity, Gain2, assuming that in practice measuring the popularity of individual items could be a complex task; especially in a dynamic distributed system. Thus, we include this variant to analyze the behavior of a *contribution gain* formulation that does not require any global knowledge to operate.

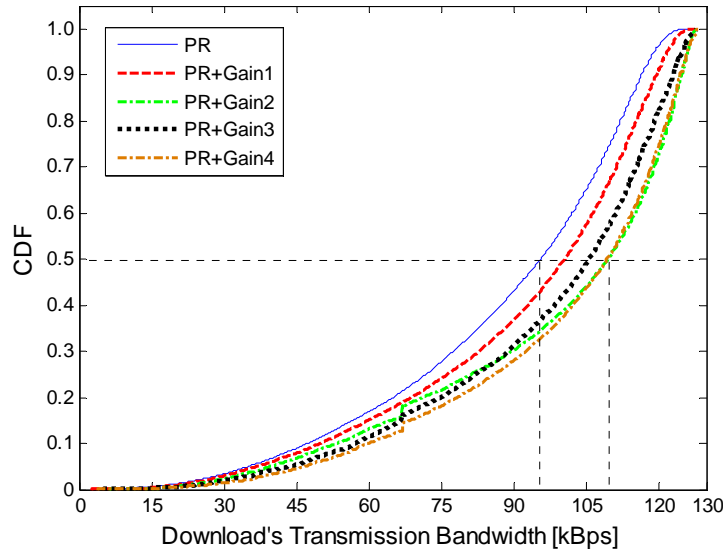


**Figure 44. Cost and Efficiency for Contribution Gain Functions**

Figure 44 presents the cost and efficiency metrics for the four simulations sets in Table 27 and for PR with a unitary *contribution gain*. Considering a confidence interval of 95%, the cost metrics obtained can not be differentiated. That is, the results obtained are statistically equivalent. For the efficiency metric, the differences between sets are minimal. Nonetheless, to determine

conclusively whether these *contribution gain* variants are equivalent or not, we analyze additional system properties.

The objective of the *contribution gain* variants is to modify a node's contribution metric. Thus, we present the effects of the different *contribution gain* variants on the RP-Cost and the TB-Utility functions.

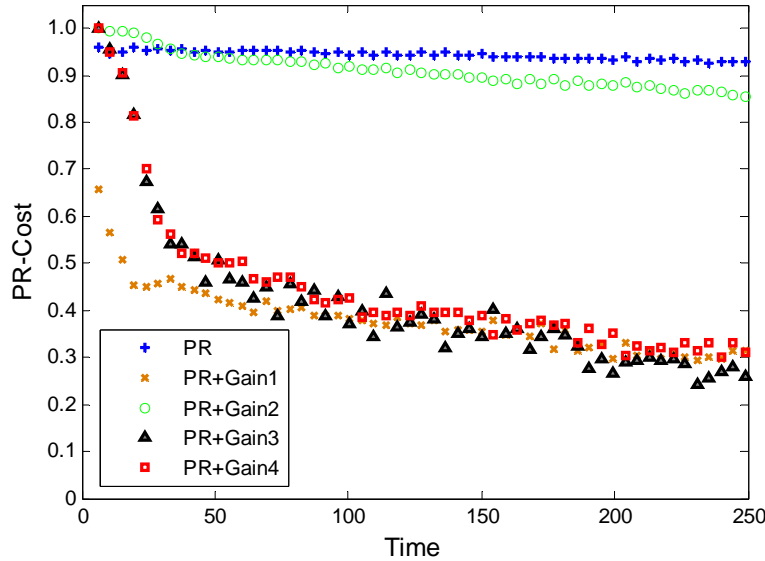


**Figure 45. Download's Transmission Bandwidth CDF for Contribution Gain Functions**

Figure 45 presents the CDF of the transmission bandwidths (i.e., TB-Utility function) that nodes receive for their download requests. PR, (with a unitary *contribution gain*) has the lowest median transmission bandwidth while Gain2 and Gain4 present the highest. Given that these *contribution gain* variants produce similar efficiency metrics, we consider that the best variant is the one resulting in better transmission bandwidths. Consequently, based on the experimental results obtained, we conclude that using an enhanced *contribution gain* function improves the system. All the *contribution gain* variants presented in Figure 45 improve the system's median transmission bandwidth with respect to PR.



In Figure 46, we present the average RP-Cost function in the system during the 4.5 hours that each experiment last. As expected, at the beginning of the simulation this function returns values close to one. That is, nodes accept new replicas with high probability. As the system evolves, nodes increase their individual content contribution, so that new storage requests (targeting the same node) are rejected with increasing probability. We think that this mechanism helps the system to achieve self-organization. Nodes with smaller contributions are “offered” a chance to increase their content contribution first. When nodes reject (or fail) a repair request, other nodes with larger contributions are selected next as targets.



**Figure 46. RP-Cost Function for Contribution Gain Functions**

The Rep-Cost function of PR and Gain2 decreases much slower than in the other schemes. Whether this trend is better or worst than the sharper change in Rep-Cost exhibited in Gain1, Gain3 and Gain4 remains an open question. The experimental results obtained are inconclusive. Furthermore, the redundancy maintenance process (Algorithm 4, Section 4.4.3) selects target nodes (i.e., candidates to receive new replicas) using a simple sorted list of node’s contributions. Consequently, whether these contributions are higher or lower is not important as long as Index

nodes maintain their list of target nodes sorted. Lastly, the fact that all experiments exhibit a decaying RP-Cost trend indicates that nodes behave as expected. That is, they increase their storage contribution as needed.

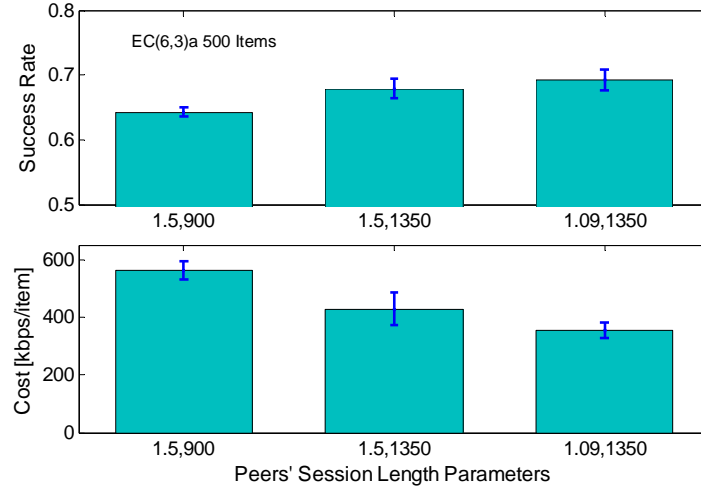
## 6.6 REDUNDANCY-SYSTEM FLEXIBILITY

In addition to exploring the behavior of the system using different parameter settings for the redundancy and incentive-based mechanisms, we study the behavior of the system with different environmental configurations. In particular, results are presented for different churn regimes, for different distribution of nodes classes, and for the presence of nodes that, for some undetermined reason, decided not to participate in the redundancy-maintenance process.

### 6.6.1 Churn Rate

The first condition to evaluate is different churn rates. All of the results presented earlier used the churn parameters listed in Table 24, which result in a median session length of 18.69 minutes. The analytical results presented in Section 4.3.3 suggest the existence of significant redundancy-maintenance savings when the system is able to adapt to even small variations in the average session length of the network. The results in Figure 47 illustrate this effect. The difference among each of the experimental settings presented is the parameters used for the initialization function of the *Peers* user class (see Section 3.3.2.1). The first bar uses  $\text{Pareto}(1.5, 900)$ , the second one uses

Pareto(1.5,1350) and the third one uses Pareto(1.09,1350). For the experiments presented, the resulting median session lengths are 18, 28, and 30 minutes, respectively<sup>34</sup>.



**Figure 47. Cost and Efficiency for Various Churn Rates**

Figure 47 shows that the system is capable to adapt to different churn conditions, improving its response metrics accordingly. For lower churn rates, the system performs fewer repairs and thus consumes fewer resources. In addition, the efficiency of the system increases for lower churn rates, reaching 69% for the slowest churn rate scenario.

### 6.6.2 Nodes Participation

The objective of the incentive-based mechanism is to promote participation in the redundancy process rather than to exclude nodes from benefiting from the network. The incentives mechanism recognizes that participants can have diverse objective functions and that their behaviors do not have to be homogeneous in order to achieve the system goal (i.e., content availability). In other words, each participant can decide autonomously its own level of participation in the system. In

---

<sup>34</sup> The average session lengths for these experiments are 34, 48 and 54 minutes respectively.

that regard, the next set of results, presented in Figure 48, illustrates the behavior of the system when different percentages of nodes do not participate in the redundancy-maintenance process. As before, content is distributed among nodes at the beginning of the simulation, but for these experiments, a fraction of the node population does not accept or publish any new fragments. Nonetheless, these non-participant nodes do perform the regular overlay tasks such as indexing, lookups, and respond to download requests.

The results obtained indicate that larger percentages of non-participating nodes degrade the efficiency of the system by a small degree and reduce the maintenance cost per file. This loss of efficiency is not surprising, as it is caused by the scarcity of data items that are held only by non-cooperating nodes<sup>35</sup>. As a result, the availability of these items can be highly dependant on the availability of the non-cooperating nodes, rather than on the effectiveness of the redundancy mechanism. Nonetheless, the distribution of fragments among nodes at the beginning of the simulation compensates, to some degree, for this dependency; thus the loss of efficiency is limited. The reduction of the cost metric on the other hand, is the result of a more complex process.

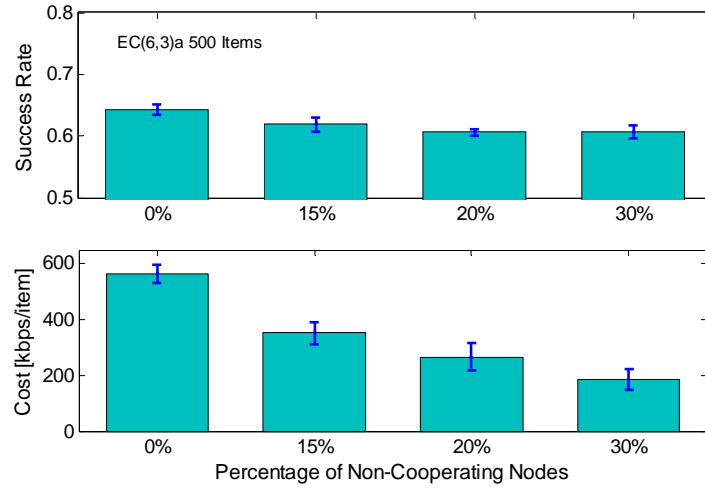
For the content space used in these experiments a total of 9,000 unique data fragments are being exchanged among nodes<sup>36</sup>. Given that the proactive redundancy method uses a replication factor of 2, the redundancy maintenance process will attempt to maintain a total of 18,000 fragments available at all times. In addition, since the average overlay network size is about 415 nodes, each node has to store (on average) 43.4 fragments to maintain complete redundancy state information. In fact, for the experimental setting with 100% of the nodes cooperating in the redundancy maintenance process, the average number of segments in the network is 49.16. That is,

---

<sup>35</sup> The percentage of complete items held exclusively by non-cooperating nodes at the beginning of the simulation varies widely across experiments, from 4% to 23%.

<sup>36</sup> The content space is 500 items and the erasure coding method generates 18 fragments for each file ( $k, S = 6, 2$ ), which results in a total of 9,000 unique fragments.

the system does generate more redundancy of what is ideally needed. Furthermore, in Section 6.5.2.3 we described the significant cost effects of the MIN\_SEG parameter (e.g. 34% savings when MIN\_SEG=0.3 instead of 0.0). In the experimental setting presented in Figure 48, cooperating nodes use MIN\_SEG=0.0. This means that nodes replicate fragments “aggressively.” On the other hand, non-cooperating nodes do not contribute to this “excessive” maintenance traffic, which is reflected in the maintenance bandwidth savings shown in the figure.



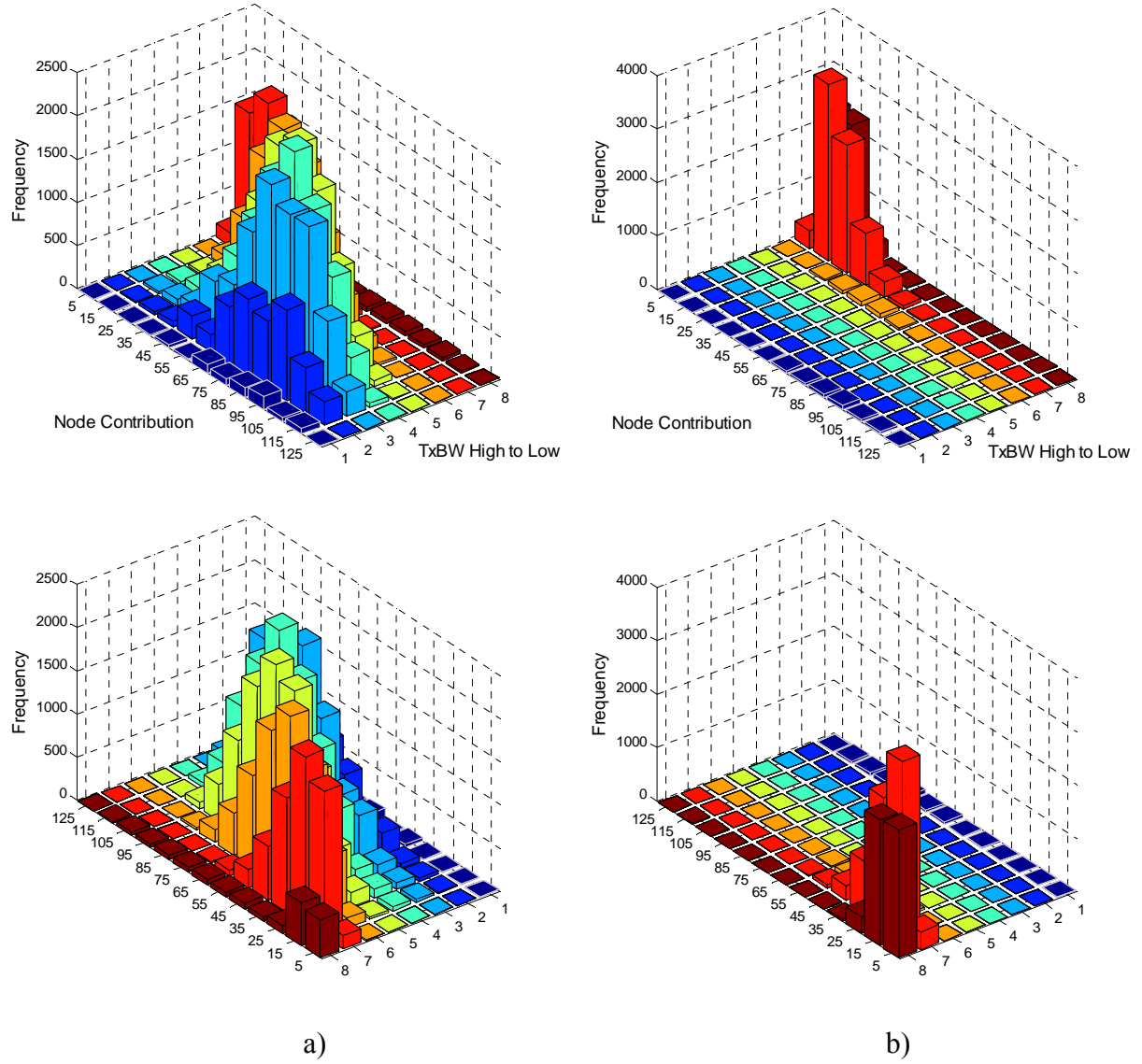
**Figure 48. Cost and Efficiency for Different Percentage of Non-Participants**

For the experimental settings with non-cooperating nodes, we would expect an increase in the average number of fragments of cooperating nodes, but that is not necessarily the case. There are two reasons for this. First, the average number of fragments in non-cooperating nodes at the start of the simulation is close to optimal. That is, 45.13, 44.19 and 42.2 for the experimental setting with 15%, 20% and 30% non-cooperating nodes respectively. In turn, the resulting average number of fragments of the cooperating nodes is 47.6, 47.3 and 45.6 for the experimental setting with 15%, 20% and 30% non-cooperating nodes respectively. This is consistent with the performance trend of the cost metric. Second, the redundancy maintenance process is nothing more than the permutation of fragments across the available nodes, that is, the distribution of up to 9,000

unique items in  $N_0$  bins without repetition; where  $N_0$  is the average number of nodes cooperating in the redundancy process. The system adapts continuously to the dynamic set of  $N_0$  nodes and the fragments they hold. For smaller values of  $N_0$ , the redundancy system has a smaller set of arrangements possible for the 9,000 unique fragments. Thus, it generates repairs less often. Indeed, for the experimental settings with 100%, 85%, 80% and 70% cooperating nodes the average repair rate is 0.91, 0.53, 0.43 and 0.28 repairs per second, respectively. This explains the improved cost metric when the number of non-cooperating nodes increases.

In order to corroborate the ability of the incentive-based mechanism to allocate resources fairly, we repeated the experiment with 30% non-cooperating nodes. This time however, the initialization of fragments among the non-compliant nodes was kept to a minimum. Figure 49 presents the results for this experimental setting. The graphs to the left present the distribution of transmission bandwidth among the participating nodes and the graphs to the right present the results for the non-compliant node population. The y-axis is for the transmission bandwidth. The value of 1 in the y-scale corresponds to the highest bandwidth (120 kbps) and the value of 8 corresponds to 10 kbps. The results show clearly that the non-compliant nodes are assigned small transmission bandwidths (because of their poor content contribution). Nonetheless, it is important to note that the incentive-based mechanism allows them to participate in the network. The set of participating nodes on the other hand, receive transmission bandwidths proportional to their contribution. Thus, they receive a better service than their non-compliant counterparts. Figure 49 presents the results in two different views. The difference between the graphs in the first and second row is that the scales in the x and y axes are reversed. This is done solely for illustrative purposes, since we are not interested in analyzing the response of the system for specific contributions or transmission bandwidth scenarios. Instead, we are interested in showing that the

incentive-based mechanism we have developed has the property of differentiating the performance of nodes according to their contribution.

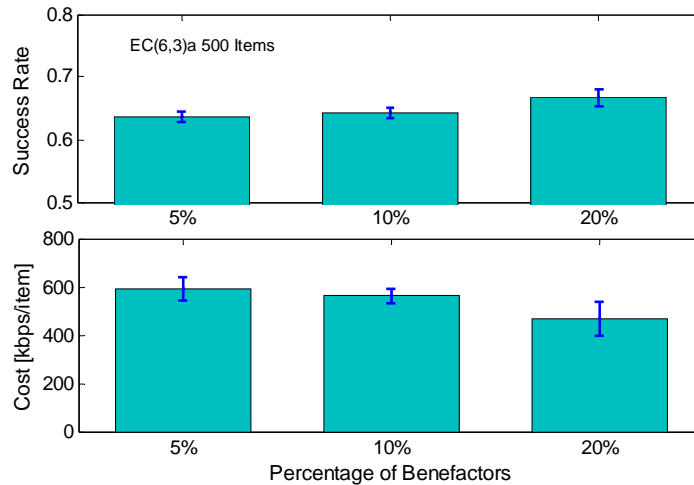


**Figure 49. Distribution of Transmission Bandwidth versus Content Contribution for: a) compliant nodes b) non-compliant nodes.**

### 6.6.3 Benefactor Nodes

Using a class-based user model in the construction of the content availability framework enables us to analyze to what extent the behavior of the *Benefactors* user class determines the efficiency of the system. While controlling the problem of free-riding has been the focus of multiple research initiatives, the presence of altruist behaviors represents another opportunity for improving the performance of the system [93]. The next set of simulations explores this possibility.

Figure 50 illustrates that the system exhibits good metrics even with a small percentage of *Benefactor* nodes present. On the opposite side of the spectrum, a larger number of *Benefactors* in the network only produces only a small improvement in system efficiency (i.e., download rate) and a small reduction in redundancy-maintenance cost. Thus, the efficiency of the system is not dependent on this user class, and using larger numbers of *Benefactor* nodes does not constitute a strategic alternative for improving content availability in P2P networks.



**Figure 50. Cost and Efficiency for Different Percentage of *Benefactors***



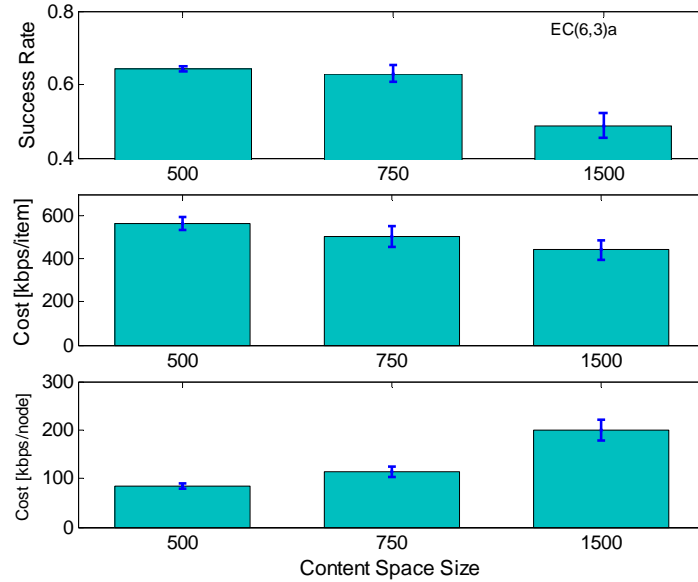
#### 6.6.4 Content Space Size

We mentioned earlier that the main scalability concern for content-sharing networks is the amount of bandwidth required to preserve content availability [10]. In that regard, the proposed mechanisms circumvents the redundancy maintenance problems of erasure coding by proactively replicating fragments. This way, the amount of information needed for a repair is significantly reduced. The next set of results presents the system's efficiency and cost metrics for various content space sizes.

With respect to efficiency, the system can sustain a good efficiency metric for up to 750 unique items, but its efficiency drops to 0.5 for 1500 items. With respect to cost, these results indicate that the system consumes less bandwidth per file in the larger content size experiments, which is important for scalability. However, in the 1,500 unique items experiment, this improvement in cost is not relevant since there is a serious efficiency loss. In addition, to properly establish the scalability of the redundancy system with respect to the content size, the units of the cost metric need to be kbps per node, rather than kbps per file. In that regard, we have included a third set of results in Figure 51 presenting this alternative cost metric. In this case, it is clear that for bigger content space sizes each node needs to commit additional bandwidth, but fortunately, this increment in bandwidth is sublinear with respect to the content space size. For instance, the cost for 750 items is 113.84 kbps per node, which represents an increase of 34.2% with respect to the cost of 84.83 kbps per node for the content space of 500 unique items. For the system configuration of 1,500 items, additional elements at play explain the abrupt drop in efficiency.

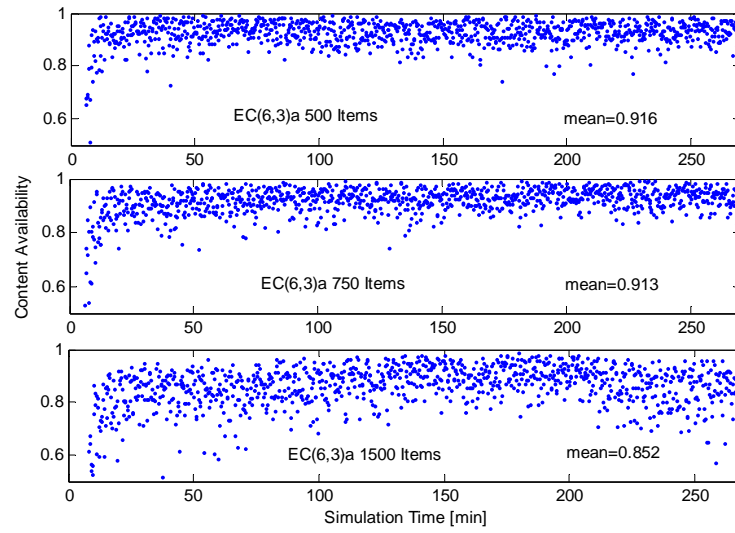
The decay in system efficiency for larger content sizes is not necessarily caused by a degradation of the redundancy repair mechanism. The initial content contribution of each node (described in Section 3.3.4) is performed in two stages. In the first stage, nodes receive complete

files, and in the second, nodes are assigned fragments. The number of items for which at least one node has a complete copy assigned is random, and the distribution of fragments among nodes is random as well. As a result, the redundancy-system might end up with insufficient resources for some items. In the different experiments conducted using a content space size of 500 unique items, the percentage of items with at least one complete copy assigned to nodes ranges between 0.5 and 0.67. We believe that this initial content distribution is a crucial factor for the efficiency of the system. In fact, this explains why none of the experimental setups surpasses the 69% efficiency mark.



**Figure 51. Cost and Efficiency for Different Content Space Sizes**

Figure 52 presents a scatter diagram of the content availability of items throughout the simulation for the three experimental content space sizes mentioned earlier. For the 500 and 750 items configuration, there is little difference in the availability of contents. The availability of each item rarely drops below 0.8. However, for the 1,500 configuration, the graph illustrates the content availability problems in the system.



**Figure 52. Content Availability for Different Content Space Sizes**

## 7.0 CONCLUSIONS AND FUTURE WORK

Three of the most important features of P2P technology versus the traditional client-server paradigm are scalability, fault tolerance and ease of deployment. Nonetheless, P2P has inherent challenges that must be addressed to enable the adoption of this technology in new applications environments as well as to improve the performance of already deployed P2P system. Intermittent node connectivity and non-cooperative user behavior are two of them. These challenges generate a dynamic set of problems not present in existing client server architectures, and we recognize content availability as one of them. This dissertation explores the problem of content availability in P2P networks and proposes a holistic redundancy maintenance framework to alleviate it.

The redundancy maintenance process presented in this dissertation can be conceptualized as a hybrid redundancy method or as a proactive redundancy maintenance policy. This mechanism, named Proactive Replication redundancy, combines the reliability gains of coding mechanisms (i.e., erasure coding or network coding) with the simplicity and flexibility of a proactive replication policy. A distinctive feature of the proposed mechanism compared with MDS erasure coding is a reduction in the amount of bandwidth needed for redundancy repairs. Furthermore, the proposed mechanism is not a replacement of other code-based mechanisms designed for distributed storage applications, such as exact-MBR network coding, it is a mean to augment their effectiveness by providing an alternative view of the redundancy maintenance problem caused by churn. That is, proactive repair.

The redundancy maintenance process presented in this dissertation is augmented with incentive-based mechanism to promote cooperation among participants and achieve fair allocation of resources. The use of incentive-based mechanisms is complementary to the proactive repair methodology proposed and obeys to the open nature of the P2P application environment. In particular, the diversity of resources and objectives among users, which need to be mediated and reconciled into a coherent cooperative framework. The proposed incentive-based extension for the redundancy maintenance process manages this heterogeneity in the network and provides the means to achieve fair allocation of resources.

The evaluation of the proposed redundancy maintenance process explores the relevance of different parameter settings in the behavior of the system; from which we derived the following conclusions:

**Content Availability.** Code-based redundancy with proactive replication of fragments is a redundancy maintenance methodology capable to improve the content availability of P2P networks significantly. For the experimental settings presented, the proposed mechanism improves the content retrieval of the system by a factor of eight. Some of the lines for future research on this topic include analyzing the durability properties of the system and examining how the characteristics of the content space could influence the behavior of the system. With respect to durability, we would like to analyze the capacity of the system to extend the average lifetime of items in the absence of their publishers. That is, for how long the fragment replication mechanism can guarantee the availability of an item in the absence of nodes with a complete copy of the item. With respect to the characteristics of the content space, we would like to analyze the behavior of the system in application scenarios with different content distribution characteristics to the ones evaluated in this dissertation, like BitTorrent or a distributed gaming application. Furthermore, we

would like to explore the use of popularity to control the redundancy maintenance process more efficiently, giving higher priority to highly demanded items.

**Adaptability of Maintenance Epochs.** The heterogeneity in nodes' behaviors and the dynamic conditions of P2P networking demand the implementation of adaptive mechanisms. The algorithms presented in this dissertation adjust the rate of the redundancy maintenance process according to the storage conditions of each individual file; achieving redundancy maintenance savings of up to 56% versus reactive repair alternatives and 50% with respect to a periodic repair methodology (Figure 35). Other research initiatives like [10, 56] use reactive repairs policies and explore different timeout values to reduce the redundancy-maintenance traffic cost. However, they operate at a completely different time scale than our system and they use loose content availability guarantees. For instance, Rodriguez in [10] analyzes the bandwidth requirements for replication and erasure coding using timeouts in the order of hours, which is adequate to solve the problem of nodes departing the overlay definitely, rather than strict content availability. The  $\text{Adapt}_{>0}$  and  $\text{Adapt}_{\geq 0}$  adaptation algorithms presented in this dissertation exhibit almost a factor of four difference on their mean maintenance epochs (156 vs 566 respectively), with  $\text{Adapt}_{\geq 0}$  consuming 46% more bandwidth than  $\text{Adapt}_{>0}$ . We also presented an adaptation algorithm based on smoothed averages, but its performance does not surpass the one of  $\text{Adapt}_{>0}$ . Naturally, a subject of future research on this topic is the exploration of alternative adaptation algorithms to achieve better bandwidth savings.

**Incentives for Content Availability .** The incentive-based mechanism presented in this dissertation barter content contribution for performance as a viable solution to promote content availability in P2P networks. The experimental evaluation presented in this dissertation indicates that the proposed mechanism effectively assign resources among nodes in proportion to their

contribution to the content availability of the overlay. In addition, other factors can be used to further characterize the content contribution of nodes to improve fairness and boost performance. Such is the case of time (i.e., session lengths) and storage (i.e., average fragment size). An interesting subject for future research is the combination of the proposed incentive, for content availability, with other incentive mechanisms for performance, like BitTorrent's tic-for-tac.

**Cost.** The proposed redundancy maintenance process presented in this dissertation reduces the redundancy-maintenance cost problem by orders of magnitude using a low complexity and flexible mechanism. However, for the level of node dynamics evaluated, the system consumes an amount of bandwidth that could be excessive for most end-users' connectivity. For instance, the system variants presented in Figure 37 consume no less than 368 kbps per file. Given that the average file size in the system is 3.5 MB, the redundancy maintenance traffic is almost 10% of the file size every second. Consequently, a natural extension to our current work is the use of the redundancy system presented as a general purpose redundancy maintenance framework to determine feasible P2P application scenarios.

## **APPENDIX A**

### **TRANSIENT REMOVAL**

The experimental portion of this work emulates a wide area P2P network using a cluster of computers interconnected by a switched LAN. Validation of the P2P routing substrate implementation is not required because the system uses exact copies of a deployed DHT application substrate. However, the price to be paid for using this architecture is time. Nodes execute a fully functional P2P routing stack in real time; thus, each experiment takes several hours to complete. In that regard, removing the transient behavior is critical not only for the accuracy of the measurements, but also for the cost of each experiment.

For the average session length of nodes, this work employs two mechanisms to reproduce the heterogeneous and unsynchronized behavior of nodes in a controlled environment. The first mechanism is the initialization process in our two-class node behavior model to recreate the heterogeneous distribution of user sessions seen in deployed systems. The second mechanism is pre-computing the initial session (or departure) for each node before starting the simulation. The second mechanism highly reduces the cost of each experiment, since other techniques to remove the transient component of the simulation require the system to be running for considerably longer intervals.



Overlay network size is the system property that better reflects transient behavior in an experiment. This metric depends on the number of simultaneously active nodes. That is, the superposition of the ON/OFF model of multiple independent nodes. Fast churning nodes become unsynchronized quickly, but nodes with longer session lengths (i.e., *Benefactors*) can take several hours to have their arrival and departures times unsynchronized. To minimize the length of this transient behavior, the initial session or offline interval of nodes is pre-computed according to the pseudo-code shown in Algorithm 13.

**Procedure PerfectSimulation( $N$ )**

**Purpose:** Pre-compute initial node state for simulations

**Returns:** Tuple (*state*, *length*) for each node

```

{  $N$ : Total number of nodes in simulation }
{  $T$ : Time interval used to reach steady state }
{ lifetime[]: Array containing the cumulative ON/OFF intervals of each node }
{ state[]: Array containing ON/OFF state of each node }

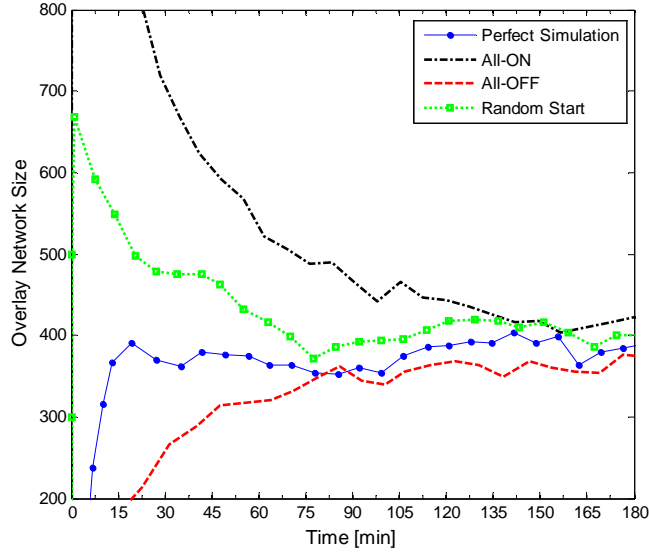
1: for  $i = 1$  to  $N$  do
2:   while lifetime[ $i$ ] <  $T$  do
3:     if state[ $i$ ] = ON then
4:       lifetime[ $i$ ] += nextSession()
5:       state[ $i$ ] = OFF
6:     else
7:       lifetime += nextOffline()
8:       state[ $i$ ] = ON
9:     end if
10:  done while
11:  length[ $i$ ] = lifetime[ $i$ ] -  $T$ 
12: done for

```

**Algorithm 13. Procedure PerfectSimulation()**

The cost savings of this mechanism versus other initialization options are presented in Figure 53. The curve for the *perfect simulation* [94] mechanism reaches steady state in approximately fifteen minutes, while the other mechanisms reach steady state in no less than one hour. In the All-ON variant, every node starts online and the average overlay network size does not reach the steady state for three hours. This is due to the effect of long on-line sessions of the

*benefactors* user class. Conversely, when all nodes start in the offline state (All-OFF) or each node is selected at random to start either in the ON or OFF state (Random Start), the system reaches a steady state faster. Nonetheless, it is still about ninety minutes for both mechanisms, which is at least six times longer than using the *perfect simulation* technique.



**Figure 53. Alternative Node Initializations for Transient Removal**

In all the measurements reported in the experimental section of this work, the values within the first fifteen minutes of each experiment are always excluded.

## BIBLIOGRAPHY

1. Daniel Stutzbach, Reza Rejaie, and S. Sen, *Characterizing Unstructured Overlay Topologies in Modern P2P File-Sharing Systems*. IEEE/ACM Transactions on Networking, 2008. **16**(2).
2. Hendrik Schulze and K. Mochalski, *Internet Study 2008/2009*. 2009, ipoque.
3. *BitTorrent Still Dominates Global Internet Traffic*. 2010 October 10 [cited; Available from: <http://torrentfreak.com/bittorrent-still-dominates-global-internet-traffic-101026/>].
4. E. Adar and B. Huberman, *Free riding on Gnutella*, in *First Monday*. 2000.
5. Jaeok Park and M.v.d. Schar. *Pricing and Incentives in Peer-to-Peer Networks*. in *INFOCOM*. 2010.
6. Daniel S. Menasche, et al., *Content Availability and Bundling in Swarming Systems*, in *CoNeXT'09*. 2009: Rome.
7. Ruben Cuevas, et al., *Is Content Publishing in BitTorrent Altruistic or Profit-Driven?*, in *Proceedings of the 2010 ACM Conference on Emerging Networking Experiments and Technology, CoNEXT 2010*. 2010, ACM: Philadelphia, PA.
8. Parhami, B., *Defect, Fault, Error, ..., or Failure?* IEEE Trans. on Reliability, 1997. **46**(4).
9. Hakim Weatherspoon and John D. Kubiatowicz, *Erasure Coding vs Replication: A Quantitative Comparison*, in *IPTPS*. 2002.
10. Rodrigo Rodriguez and Barbara Liskov, *High Availability in DHTs: Erasure Coding vs Replication*, in *IPTPS*. 2005.
11. Xin, Q., et al. *Reliability mechanisms for very large storage systems in 20th IEEE/11th NASA Goddard Conference on Mass Storage Systems and Technologies (MSST)*. 2003. Santa Cruz, CA.
12. *Planet-Lab*. [cited; Available from: <http://www.planet-lab.org/>].
13. Sean Rhea, et al., *Handling Churn in a DHT*, in *Technical Report*. 2003, University of California, Berkeley.

14. Dennis Geels and J. Kubiawicz, *Replica management should be a game*, in *Proceedings of the 10th workshop on ACM SIGOPS European workshop*. 2002: Saint-Emilion, France.
15. M. Feldman, et al., *Quantifying Disincentives in Peer-to-Peer Networks*, in *1st Workshop on Economics of Peer-to-Peer Systems*. 2003: Berkeley, CA.
16. Hongxing Li and G. Chen, *Data Persistence in P2P Networks with Redundancy Schemes*, in *Sixth International Conference on Grid and Cooperative Computing*. 2007, IEEE.
17. Kan Zhang, Nick Antonopoulos, and Z. Mahmood, *A Review of Incentive Mechanisms in Peer-to-Peer Systems*, in *First International Conference on Advances in P2P Systems*. 2009, IEEE Computer Society.
18. Panayotis Antoniadis, Costas Courcoubetis, and Ben Strulo, *Incentives for Content Availability in Memory-less Peer-to-Peer File Sharing Systems*, in *ACM SIG on Ecommerce*. 2005.
19. Rhea, S. *The Bamboo Distributed Hash Table* 2009 [cited; Available from: <http://bamboo-dht.org/>].
20. *Modelnet*. [cited; Available from: <http://modelnet.ucsd.edu>].
21. Sean Rhea, et al. *Handling Churn in a DHT*. in *USENIX '04*. 2004.
22. *Mininova's Torrent Downloads Double to 7 Billion in a Year*. 2009 January 5 [cited; Available from: <http://torrentfreak.com/mininovas-torrent-downloads-doubled-in-a-year-090105/>].
23. *Netflix is Killing BitTorrent in the US*. 2011 April 27 [cited; Available from: <http://torrentfreak.com/netflix-is-killing-bittorrent-in-the-us-110427/>].
24. Dejan S. Milojevic, Vana Kalogeraki, and et al., *Peer-to-Peer Computing*, in *HP Technical Report*. 2002, HP.
25. Eng Keong Lua, et al., *A Survey and Comparison of Peer-to-Peer Overlay Network Schemes*. Communications Surveys & Tutorials, IEEE 2005. 7(2).
26. Xuemin Shen, et al., *Handbook of Peer-to-Peer Networking*. 2010: Springer.
27. Karl Aberery, et al., *The essence of P2P: A reference architecture for overlay networks*, in *Fifth International Conference on Peer-to-Peer Computing*. 2005: Konstanz, Germany.
28. Ion Stoica, et al. *Chord: A scalable peer-to-peer lookup service for internet applications*. in *Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications*. 2001.
29. *The eMule Project*. 2011 [cited; Available from: <http://www.emule-project.net/>].

30. *Gnutella*. Wikipedia 20 July 2011 [cited; Available from: <http://en.wikipedia.org/wiki/Gnutella#Software>].
31. D. Stutzbach and R. Rejaie, *Understanding Churn in Peer-to-Peer Networks*, in *Technical Report*. 2005, University of Oregon.
32. Gummadi, K.P., et al., *Measurement, Modeling, and Analysis of Peer-to-Peer File-Sharing Workload*, in *SOSP '03 Proceedings of the nineteenth ACM symposium on Operating systems principles* 2003, ACM.
33. Xin Bai A , Dan C. Marinescu A , and et. al., *A Macroeconomic Model for Resource Allocation in Large-Scale Distributed Systems*. Parallel and Distributed Computing, 2008.
34. Dabek, F., et al. *Wide-Area Cooperative Storage with CFS*. in *ACM SOSP* 2001.
35. S. Rhea, et al. *Pond: The OceanStore Prototype*. in *2nd USENIX conference on File and Storage Technologies (FAST)*. 2003.
36. Antony Rowstron and Peter Druschel. *Pastry: Scalable, Distributed Object Location And Routing For Large-Scale Peer-To-Peer Systems*. in *IFIP/ACM Middleware*. 2001. Heidelberg, Germany: ACM
37. Androutsellis-Theotokis, S. and D. Spinellis, *A survey of peer-to-peer content distribution technologies*. ACM Computer Surveys, 2004. **36**(4): p. 335-371.
38. Gummadi and et al. *The Impact of DHT Routing Geometry on Resilience and Proximity*. in *ACM SIGCOMM*. 2003.
39. Octavio Herrera-Ruiz and T. Znati. *Static Resiliency vs Churn-Resistance Capability of DHT-Protocols*. in *ISCA PDCS* 2005. Las Vegas, Nevada.
40. Naicken, S., et al., *The State of Peer-to-Peer Simulators and Simulations* ACM Computer Communications Review, 2007. **37**(2): p. 95-98.
41. Z. Yao, et al. *Modeling Heterogeneous User Churn and Local Resilience of Unstructured P2P Networks*. in *IEEE ICNP*. 2006.
42. X. Wang, et al., *Robust Lifetime Measurement in Large-Scale P2P Systems with Non-Stationary Arrivals*, in *IEEE P2P*. 2009.
43. F. E. Bustamante and Y. Qiao, *Friendships that last: Peer lifespan and its role in P2P protocols*, in *International Workshop on Web Content Caching and Distribution*. 2003.
44. X. Wang, Z. Yao, and D. Loguinov, *Residual-Based Estimation of Peer and Link Lifetimes in P2P Networks*. IEEE/ACM Transactions on Networking, 2009. **17**(no. 3).
45. D. Stutzbach and R. Rejaie. *Understanding Churn in peer-to-peer networks*. in *ACM Internet Measurement Conf. (IMC)*. 2006.

46. M. Steiner, T. En-Najjary, and E. W. Biersack. *A Global View of KAD*. in *7th ACM Internet Measurement, IMC'07*. 2007. San Diego, California, USA.
47. Qiuming Luo, et al., *A Novel Model and a Simulation Tool for Churn of P2P Network*, in *Parallel and Distributed Computing, Applications and Technologies*. 2010.
48. Octavio Herrera and T. Znati, *Modeling Churn in P2P Networks*, in *40th Annual Simulation Symposium*. 2007.
49. Enrique Fernández-Casado, Marc Sánchez-Artigas, and P. García-López, *Affluenza: Towards Universal Churn Generation*, in *P2P*. 2010.
50. Boxun Zhang, Alexandru Iosup, and D. Epema, *The Peer-to-Peer Trace Archive*, in *Parallel and Distributed Systems Report Series*. 2010, Delft University of Technology.
51. Godfrey, B. *Repository of Availability Traces*. 2010 [cited; Available from: <http://www.cs.illinois.edu/~pbg/availability/>].
52. Ranjita Bhagwan, Stefan Savage, and Geoffrey M. Voelker, *Understanding Availability*, in *IPTPS*. 2003.
53. AG Dimakis, K Ramchandran, and Y Wu and C Suh. *A Survey on Network Codes for Distributed Storage*. in *IEEE Surveys*. 2011.
54. W. K. Lin, D.M. Chiu, and Y. B. Lee, *Erasur Code Replication Revisited*, in *IEEE P2P*. 2004.
55. Ranjita Bhagwan, Stefan Savage, and Geoffrey M. Voelker, *Replication Strategies for Highly Available Peer-to-Peer Storage Systems*, in *Tech Report*. 2002, UC San Diego.
56. Alexandros G. Dimakis, et al., *Network coding for Distributed storage Systems*, in *INFOCOM*. 2007: Anchorage, Alaska.
57. K. Rashmi, N. B. Shah, and P.V. Kumar, *Optimal Exact-Regeneration Codes for Distributed Storage at teh MSR and MBR Points via Product-Matrix Construction*. IEEE Trans. Information Theory, 2010.
58. Alessandro Dominuco and E. Biersack, *Hierarchical Codes: How to Make Erasure Codes Attractive for Peer-to-Peer Storage Systems*. in *Proceedings of Peer-to-Peer Computing*, 2008: p. 89-98.
59. Chris Williams, et al. *Redundancy Management for P2P Storage*. in *Proceedings of the Seventh IEEE International Symposium on Cluster Computing and the Grid* 2007.
60. Fan Wu, Tongqing Qiu Yuequan Chen, and Guihai Chen, *Redundancy Schemes for High Availability in DHTs*, in *IPSA 2005*. 2005.

61. Guangping Xu, Gang Wang, and J. Liu, *A hybrid redundancy approach for data availability in structured P2P network systems*, in *International Symposium on Pacific Dependable Computing*. 2007, IEEE.
62. Nishida, H. and T. Nguyen, *A Global Contribution Approach to Maintain Fairness in P2P Networks*. IEEE Transactions on Parallel and Distributed Systems, 2010. **21**(6): p. 812-826.
63. Rameez Rahman, et al. *Improving efficiency and fairness in p2p systems with effort-based incentives* in ICC. 2010.
64. Tsuen-Wan Johnny Ngan , et al., *On Designing Incentives-Compatible Peer-to-Peer Systems* in *2nd Bertinoro Workshop on Future Directions in Distributed Computing (FuDiCo II: S.O.S.)*. 2004: Bertinoro, Italy.
65. Anwitaman Datta and Karl Aberer. *Internet-Scale Storage Systems under Churn – A Study of the Steady-State using Markov Models*. in *Sixth IEEE International Conference on Peer-to-Peer Computing*. 2006.
66. Emil Sit, et al., *Proactive Replication for Data Durability*, in *5th Int'l Workshop on Peer-to-Peer Systems (IPTPS)*. 2006.
67. Stefan Saroiu, P. Krishna Gummadi, and Steven D. Gribble. *A Measurement Study of Peer-to-Peer File Sharing Systems*. in *Multimedia Computing and Networking (MMCN)*. 2002.
68. Shanyu Zhao, Daniel Stutzbach, and Reza Rejaie. *Characterizing Files in the Modern Gnutella Network: A Measurement Study*. in *SPIE/ACM Multimedia Computing and Networking*. 2006. San Jose, CA.
69. *Inet topology generator v 3.0*. [cited; Available from: <http://topology.eecs.umich.edu/inet/>].
70. Mauro Andreolini and Riccardo Lancellotti. *Analysis of peer-to-peer systems: workload characterization and effects on traffic cacheability*. in *12th Annual Meeting of the IEEE / ACM International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS 2004)* 2004. Volendam (NL).
71. J. Li, et al., *A Performance vs. Cost Framework for Evaluating DHT Design Tradeoffs under Churn*, in *INFOCOM*. 2005: Miami, FL.
72. Subhabrata Sen and Jia Wong, *Analyzing peer-to-peer traffic across large networks*. IEEE/ACM Transactions on Networking (TON), 2004.
73. R. Bolla, et al., *A measurement study supporting P2P file-sharing community models*. Computer Networks, 2009(53): p. 485-500.
74. Moritz Steiner, Taoufik En-Najjary, and E.W. Biersack, *Analyzing Peer Behavior in KAD*, in *Research Report*. 2007, Institut Eurecom: Sophia-Antipolis, France.

75. Jacky C. Chu, Kevin S. Labonte, and Brian N. Levine, *Availability and locality measurements of peer-to-peer file systems*, in *SPIE* 2002.
76. S. Saroiu, P. K. Gummadi, and S. D. Gribble. *A Measurement Study of Peer-to-Peer File Sharing Systems*. in *Multimedia Computing and Networking (MMCN)*. 2002. San Jose, CA.
77. Zhonghong Ou, E. Harjula, and M. Ylianttila, *Effects of Different Churn Models on the Performance of Structured Peer-to-Peer Networks*, in *Personal, Indoor and Mobile Radio Communications*. 2009.
78. Dimakis, A.G., P. B. Godfrey, and Y. Wu. *Network Coding for Distributed Storage Systems*. in *Information Theory*. 2010.
79. K. V. Rashmi, et al., *Optimal Exact-Regeneration Codes for Distributed Storage at the MSR and MBR Points via Product-Matrix Construction*. *IEEE Trans. Information Theory*, 2010.
80. Kuo, W. and M. J. Zuo, *Thek-out-of-n System Model*, in *Optimal Reliability Modeling: Principles and Applications*. 2003, Wiley. p. 231-280.
81. K. V. Rashmi, et al. *Explicit Construction of Optimal Exact Regenerating Codes for Distributed Storage*. in *Proc. Allerton Conf., Urbana-Champaign*. 2009. San Diego.
82. Alexandros G. Dimakis, et al. *A Survey on Network Codes for Distributed Storage*. in *IEEE Surveys*. 2011.
83. Pinheiro, E., W.-D. Weber, and L.A.e. Barroso, *Failure Trends in a Large Disk Drive Population*. *USENIX Conference on File and Storage Technologies*, 2007.
84. *Hard drive manufacturer specifications*. 2011 [cited; Available from: <http://www.wdc.com/wdproducts/library/SpecSheet/ENG/2879-701220.pdf>].
85. Lehpamer, H., *Transmission Networks Fundamentals*, in *Microwave Transmission Networks, Second Edition*. 2010, McGrawHill. p. 15.
86. Charles Blake and Rodrigo Rodriguez, *High Availability, Scalable Storage, Dynamic Peer Networks: Pick Two*, in *HotOS IX*. 2003.
87. Hardin, G., *The Tragedy of the Commons*, in *Science*. 1968. p. 1243-1248.
88. *Tragedy of the commons*. 2011 [cited; Available from: [http://en.wikipedia.org/wiki/Tragedy\\_of\\_the\\_commons](http://en.wikipedia.org/wiki/Tragedy_of_the_commons)].
89. S. B. Handurukande, et al., *Peer Sharing Behavior in the eDonkey network, and implications for the design of server-less file sharing systems*. *SIGOPS Oper. Sys. Rev.*, 2006. **40**: p. 359-371.



90. Deterding, S. *Meaningful Play: Getting Gamification Right* Google Tech Talk [Webcast] Jan 24, 2011 [cited 2011; Available from: <http://www.youtube.com/watch?v=7ZGCPap7GkY&feature=relmfu>].
91. Wikileaks. [cited; Available from: [wikileaks.org](http://wikileaks.org)].
92. Schwarz, T., et al., *Disk Failure Investigations at the Internet Archive*, in *NASA/IEEE Conference on Mass Storage Systems and Technologies*. 2006.
93. Zghaibeh, M. and K.G. Anagnostakis, *On the Impact of P2P Incentive Mechanisms On User Behavior*, in *NetEcon+IBC*. 2007: San Diego.
94. Mecke, K.R. and D. Stoyan, *A Primer in Perfect Simulation*. Springer Lecture Notes in Physics, 2000: p. 349 - 378.