



Joana Margarida Dias **Caracterização de utilizadores em sistemas P2P**
de Bragança Gonçalves **Peer-level characterization of P2P systems**



Universidade de Aveiro Departamento de Electrónica, Telecomunicações e
2008 Informática

Joana Margarida Dias de Bragança Gonçalves **Caracterização de utilizadores em sistemas P2P**
de Bragança Gonçalves **Peer-level characterization of P2P systems**

Dissertação apresentada à Universidade de Aveiro para cumprimento dos requisitos necessários à obtenção do grau de Mestre em Engenharia Electrónica e Telecomunicações, realizada sob a orientação científica do Doutor António Nogueira, Professor Auxiliar do Departamento de Electrónica, Telecomunicações e Informática da Universidade de Aveiro e do Doutor Rui Valadas, Professor Associado com Agregação do Departamento de Electrónica, Telecomunicações e Informática da Universidade de Aveiro.

aos meus pais
ao meu irmão
à minha família
aos meus amigos

o júri

Doutor João Nuno Pimentel da Silva Matos

Professor Associado da Universidade de Aveiro

Doutor Joel José Puga Coelho Rodrigues (Arguente)

Professor Auxiliar do Departamento de Informática da Faculdade de Ciências da Engenharia da Universidade da Beira Interior

Doutor António Manuel Duarte Nogueira (Orientador)

Professor Auxiliar da Universidade de Aveiro

Doutor Rui Jorge Morais Tomaz Valadas (Co-Orientador)

Professor Associado com Agregação da Universidade de Aveiro

Acknowledgments

First of all, I want to thank everyone that made this thesis possible, emphasising:

- my supervisor, co-supervisor and collaborator, Doctor António Nogueira, Doctor Rui Valadas and Doctor Paulo Salvador respectively, for the opportunity they gave me to make this thesis and for the motivation, encouragement and support they have always deposited on me;
- my real friends for their simple presence in my life, being very important for my personal evolution and essential during this Master graduation period;
- finally, to the most important people in my life, my great family, with a special place for my parents, my brother and my cousin that always believed in my capabilities and inner strength, supporting all my choices, being responsible for my well-fare, never let me going down on the hardest periods of my life.

Palavras-chave

Redes *Peer-to-Peer*, partilha de ficheiros, BitTorrent, Vuze, disponibilidade dos *peers*, caracterização dos *peers*, *Round Trip Time*.

Resumo

As aplicações para a partilha de ficheiros através de redes *Peer-to-Peer* (P2P) são cada vez mais populares, sendo responsáveis por uma grande parte de todo o tráfego gerado na Internet. As facilidades que estas aplicações proporcionam no que se refere à sua fácil utilização e ao acesso a um grande número de ficheiros são o motivo de atracção dos seus utilizadores, cujo número tem vindo a aumentar cada vez mais. No entanto, porque consomem uma quantidade excessiva de recursos de rede, revelam-se prejudiciais para outras aplicações concorrentes que necessitam de aceder à rede. Deste modo, é fundamental proceder-se a uma caracterização completa destes sistemas no sentido de otimizar as políticas de encaminhamento da rede, os mecanismos e as estratégias de gestão da Qualidade de Serviço (QoS), melhorar e otimizar os mecanismos de segurança e planejar as infra-estruturas de rede, fazendo *upgrade* ou *downgrade* de certos recursos segundo os parâmetros medidos e as exigências pretendidas.

Esta tese foca-se num dos clientes do BitTorrent, o Vuze, devido à sua popularidade neste tipo de aplicações P2P. Com esta aplicação foi possível obter uma grande e variada amostra de *peers*, através da partilha de diferentes categorias de ficheiros, com o objectivo de os caracterizar segundo diversas perspectivas, em particular a sua distribuição geográfica, a evolução no número de *peers* activos ao longo do período de análise e em diferentes horas do dia e as variações do *Round Trip Time* (RTT) de acordo com as distâncias entre o *peer* de origem e o *peer* de destino. Foram também alvo de estudo as dependências do RTT com diferentes tipos de ligações de acesso à Internet e para diferentes horas do dia.

No final da realização desta tese é possível concluir que: (i) existe uma grande utilização deste tipo de serviço em todo o mundo, especialmente na Europa, América e Ásia, sendo detectado um número elevado de *peers* que varia em função da categoria de ficheiros que está a ser partilhada; (ii) existe uma grande quantidade de *peers* que se encontram protegidos através de Firewalls ou atrás de mecanismos de NAT/PAT; (iii) existe uma utilização contínua deste tipo de aplicações de partilha de ficheiros; (iv) há uma forte dependência do RTT com a distância entre *peers* e com o tipo de ligação de acesso à Internet.

Keywords

Peer-to-Peer networks, file sharing, BitTorrent, Vuze, peers' availability, peers-level characterization, Round Trip Time.

Abstract

Peer-to-Peer (P2P) file-sharing applications are becoming very popular, being responsible for the major fraction of the current Internet traffic. The facilities given by these applications in terms of their easy utilization and ability to access and share a large number of files are a strong attraction for users, whose number has been increasing all the time. However, these applications consume a vast amount of network resources, which can be critical for other applications that need to access the network. Therefore, a complete characterization of these systems is needed in order to optimize network routing policies, Quality of Service (QoS) mechanisms and management strategies, to improve and optimize the network security mechanisms and to plan network infrastructures, upgrading or downgrading certain resources according to measured parameters and predicted requirements.

This thesis focuses on Vuze, one of the BitTorrent clients, since BitTorrent is one of the most popular P2P file sharing applications. By sharing different types of files using this client application, it was possible to obtain a huge and varied sample of peers aiming to obtain a complete peer-level characterization of this system based on different parameters, such as the geographical distribution of involved peers, the availability of peers during the whole period of analysis and for different daily periods and the variability of Round Trip Time (RTT) according to distances between origin and end-hosts. The study of the RTT dependences with the type of Internet connectivity used and the period of the day was also one of the objectives of this thesis.

At the end of the work, it was possible to verify: (i) the high number of peers existing around the world, mainly in Europe, America and Asia, which depends on the file category under analysis; (ii) the high proportion of peers that are protected by Firewalls or located behind Network/Port Address Translation (NAT/PAT) mechanisms; (iii) the continuous usage of this application by peers and (iv) the strong dependence of RTT with the distance between hosts and the type of Internet connection that is used.

Contents

1. Introduction	1
1.1 Thesis proposal	2
1.2 Thesis contributions	3
1.3 Related work	3
1.4 Thesis outline	5
1.5 Notation used	6
2. State of Art	7
2.1 Introduction.....	7
2.2 P2P networks	9
2.2.1 P2P architecture	9
2.2.2 P2P generations.....	13
2.3 BitTorrent protocol	14
2.4 Vuze	18
2.5 Summary	19
3. Methodology.....	21
3.1 Introduction.....	21
3.2 Files selection.....	21
3.3 Peers localization	24
3.4 Availability and Round Trip Time.....	25
3.5 Shell scripts	26
3.6 Summary	27
4. Peers localization and availability	29
4.1 Geographical distribution.....	29
4.2 Peers' availability.....	43

4.3	Summary	58
5.	Round Trip Time	59
5.1	Summary	85
6.	Conclusions	87
6.1	Future work	88
	Acronyms	91
	References	93
	Appendix A - Shell scripts	97

List of figures

Figure 2.1 - Relative P2P traffic volume [32].	7
Figure 2.2 – Traffic volume distribution for most popular P2P protocols in five different world regions [32].	8
Figure 2.3 – Design of a centralized architecture.	10
Figure 2.4 – Design of a decentralized architecture.	11
Figure 2.5 - Design of a hybrid architecture.	13
Figure 2.6 – Vuze logging option field.	19
Figure 4.1 - Geographical distribution of peers per continent – <i>2008 movies</i> category.	30
Figure 4.2 - Geographical distribution of peers per country – <i>2008 movies</i> category.	30
Figure 4.3 - Geographical distribution of peers per continent – <i>Music</i> category.	31
Figure 4.4 - Geographical distribution of peers per country – <i>Music</i> category.	31
Figure 4.5 - Geographical distribution of peers per continent - <i>Animated movies</i> category.	32
Figure 4.6 - Geographical distribution of peers per country – <i>Animated movies</i> category.	32
Figure 4.7 - Geographical distribution of peers per continent – <i>French movies</i> category.	33
Figure 4.8 - Geographical distribution of peers per country – <i>French movies</i> category.	33
Figure 4.9 - Geographical distribution of peers per continent – <i>Indian movies</i> category.	34
Figure 4.10 - Geographical distribution of peers per country – <i>Indian movies</i> category.	34
Figure 4.11 - Geographical distribution of peers per continent – <i>Linux distribution</i> category.	35
Figure 4.12 - Geographical distribution of peers per country – <i>Linux distribution</i> category.	35
Figure 4.13 – Peers’ distribution normalized by country population – <i>2008 movies</i> category.	37
Figure 4.14 – Peers’ distribution normalized by country population – <i>Music</i> category.	38
Figure 4.15 – Peers’ distribution normalized by country population – <i>Animated movies</i> category.	39
Figure 4.16 – Peers’ distribution normalized by country population – <i>French movies</i> category.	40
Figure 4.17 – Peers’ distribution normalized by country population – <i>Indian movies</i> category.	41

Figure 4.18 – Peers’ distribution normalized by country population – <i>Linux distribution</i> category.	42
Figure 4.19 – Evaluation of available peers – <i>2008 movies</i> category.	44
Figure 4.20 – Evaluation of available peers – <i>Music</i> category.....	44
Figure 4.21 – Evaluation of available peers – <i>Animated movies</i> category.....	45
Figure 4.22 – Evaluation of available peers – <i>French movies</i> category.....	45
Figure 4.23 – Evaluation of available peers – <i>Indian movies</i> category.....	46
Figure 4.24 – Evaluation of available peers – <i>Linux distribution</i> category.....	46
Figure 4.25 – Evaluation of available peers – Israel and eight European countries.....	49
Figure 4.26 – Evaluation of available peers – United States, Canada and Brazil.	50
Figure 4.27 – Evaluation of available peers – Philippines, Australia, Malaysia and Singapore.	51
Figure 4.28 – Evaluation of available peers – France.	52
Figure 4.29 – Evaluation of available peers – India.	53
Figure 4.30 – Evaluation of peer availability in a short period of time - <i>Music</i> category.....	55
Figure 4.31 – Evaluation of peer availability in a short period of time - <i>Animated movies</i> category.	55
Figure 4.32 – Evaluation of peer availability in a short period of time - <i>French movies</i> category.	56
Figure 4.33 – Evaluation of peer availability in short periods of time - <i>Indian movies</i> category.	56
Figure 4.34 – Evaluation of peer availability in short periods of time - <i>Linux distribution</i> category.	57
Figure 5.1 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – <i>2008 movies</i> category.	60
Figure 5.2 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – <i>2008 movies</i> category.	60
Figure 5.3 – Round Trip Time distribution of the CATV 12 Mbps Internet connection –	61
Figure 5.4 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection –.....	61

Figure 5.5 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – <i>Animated movies</i> category.....	62
Figure 5.6 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – <i>Animated movies</i> category.....	62
Figure 5.7 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – <i>French movies</i> category.....	63
Figure 5.8 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – <i>French movies</i> category.	63
Figure 5.9 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – <i>Indian movies</i> category.....	64
Figure 5.10 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – <i>Indian movies</i> category.....	64
Figure 5.11 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – <i>Linux distribution</i> category.....	65
Figure 5.12 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – <i>Linux distribution</i> category.	65
Figure 5.13 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connection types - <i>Movies</i> category.....	67
Figure 5.14 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - <i>Music</i> category.	68
Figure 5.15 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - <i>Animated movies</i> category.....	68
Figure 5.16 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - <i>French movies</i> category.	69
Figure 5.17 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - <i>Indian movies</i> category.	69
Figure 5.18 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - <i>Linux distribution</i> category.	70
Figure 5.19 – RTT distribution in one day of analysis - <i>2008 movies</i> category with CATV 12 Mbps Internet connection.....	72

Figure 5.20 – RTT distribution in one day of analysis - <i>2008 movies</i> category with ADSL 4 Mbps Internet connection.	72
Figure 5.21 – RTT distribution in one day of analysis - <i>Music</i> category with CATV 12 Mbps Internet connection.	73
Figure 5.22 – RTT distribution in one day of analysis - <i>Music</i> category with ADSL 4 Mbps Internet connection.	73
Figure 5.23 – RTT distribution in one day of analysis - <i>Animated movies</i> category with CATV 12 Mbps Internet connection.	74
Figure 5.24 – RTT distribution in one day of analysis - <i>Animated movies</i> category with ADSL 4 Mbps Internet connection.	74
Figure 5.25 – RTT distribution in one day of analysis - <i>French movies</i> category with CATV 12 Mbps Internet connection.	75
Figure 5.26 – RTT distribution in one day of analysis - <i>French movies</i> category with ADSL 4 Mbps Internet connection.	75
Figure 5.27 – RTT distribution in one day of analysis - <i>Indian movies</i> category with CATV 12 Mbps Internet connection.	76
Figure 5.28 – RTT distribution in one day of analysis - <i>Indian movies</i> category with ADSL 4 Mbps Internet connection.	76
Figure 5.29 – RTT distribution in one day of analysis - <i>Linux distribution</i> category with CATV 12 Mbps Internet connection.	77
Figure 5.30 – RTT distribution in one day of analysis - <i>Linux distribution</i> category with ADSL 4 Mbps Internet connection.	77
Figure 5.31 – RTT cumulative distribution – <i>2008 movies</i> category with CATV 12 Mbps Internet connection.	78
Figure 5.32 – RTT cumulative distribution – <i>2008 movies</i> category with ADSL 4 Mbps Internet connection.	79
Figure 5.33 – RTT cumulative distribution – <i>Music</i> category with CATV 12 Mbps Internet connection.	79
Figure 5.34 – RTT cumulative distribution – <i>Music</i> category with ADSL 4 Mbps Internet connection.	80

Figure 5.35 – RTT cumulative distribution – <i>Animated movies</i> category with CATV 12 Mbps Internet connection.	80
Figure 5.36 – RTT cumulative distribution – <i>Animated movies</i> category with ADSL 4 Mbps Internet connection.	81
Figure 5.37 – RTT cumulative distribution – <i>French movies</i> category with CATV 12 Mbps Internet connection.	81
Figure 5.38 – RTT cumulative distribution – <i>French movies</i> category with ADSL 4 Mbps Internet connection.	82
Figure 5.39 – RTT cumulative distribution – <i>Indian movies</i> category with CATV 12 Mbps Internet connection.	82
Figure 5.40 – RTT cumulative distribution – <i>Indian movies</i> category with ADSL 4 Mbps Internet connection.	83
Figure 5.41 – RTT cumulative distribution – <i>Linux distribution</i> category with CATV 12 Mbps Internet connection.	83
Figure 5.42 – RTT cumulative distribution – <i>Linux distribution</i> category with ADSL 4 Mbps Internet connection.	84
Figure A. 1 - <i>If condition</i> structure.	101
Figure A. 2 – <i>Case condition</i> structure.	101
Figure A. 3 – <i>For loop</i> structure (a).	102
Figure A. 4 – <i>For loop</i> structure (b).	102
Figure A. 5 – <i>While loop</i> structure.	102
Figure A. 6 – Shell script to extract information from Log files.	105
Figure A. 7 – Shell script for the geographical localization of peers and calculation of the RTT.	106
Figure A. 8 – Shell script for the geographical localization of peers and calculation of the RTT.	106
Figure A. 9 – Shell script to measure RTT values.	107
Figure A. 10 – Gnuplot program to create normalized maps.	108
Figure A. 11 – Shell script for treating the results in order to create polar maps.	109

Figure A. 12 – Gnuplot program to create polar maps	110
Figure A. 13 – Gnuplot program to create 3-Dimensional plots	112

List of tables

Table A. 1 – Mathematical comparison.	103
Table A. 2 – String comparison.....	103
Table A. 3 – Shell test for files or directories.....	104
Table A. 4 – Logical operators.	104

1. Introduction

In recent years, traffic corresponding to Peer-to-Peer (P2P) networks became a significant percentage of the total Internet traffic, turning its detailed analysis into a very important issue. Unlike client/server networks, P2P networks don't have a Central Server, becoming easy for a common user to participate on it. The facility of using these networks to share a great variety of file contents, eventually having large sizes, using applications such as Napster, Gnutella and BitTorrent, attracts a huge number of users resulting in a very fast growing utilization. Nowadays, these networks are evolving towards a real time multimedia content distribution infrastructure that is able to provide reliable Internet Protocol Television (IPTV) and Video on Demand (VoD) services.

The characterization of P2P networks can be helpful for network operators in order to be able to manage their infrastructures: using complete and reliable information, it is possible to prevent bad usage of network resources identifying wrong behaviours and provide a better network capacity arrangement; it is also possible to improve network security, develop new load models and efficient Quality of Service (QoS) mechanisms. Peer-level characterization is also important for the development of new technical and marketing design plans and for the deployment and management of new services.

However, this complete characterization is not easy to fulfil. There are some factors that make a fair characterization of P2P networks impossible or at least difficult, such as the difficulty to identify and control hosts wrapped in the network due to the fact that they are protected by Firewalls or NAT/PAT (Network Address Translation/Port Address Translation) mechanisms and because some P2P protocols use specific ports.

A good characterization of P2P networks must rely on a solid know-how of their applications and protocols, which is a hard task due to the constant evolution of these protocols and applications.

The current and fast growth of P2P networks has definitely attracted the attention of Internet network operators that are improving traffic characterization methods and networks management procedures. It is really important to anticipate eventual P2P networks problems

and solve them from the beginning in order to avoid that this great evolution in our lives becomes a true nightmare.

In this thesis, our analysis of P2P networks will be focused on BitTorrent [33] since it is one of the most popular P2P systems, involving a huge number of peers all around the world that usually share multiple kinds of files. Thereby, using this P2P application it is possible to obtain a big and varied sample of peers for a more complete and conclusive peer-level characterization. In this way, this work will evaluate the BitTorrent network making an analysis of peers involved in the download of different kind of files from different perspectives: localizing their geographic areas, studying their availability, studying the Round Trip Time (RTT) between origin and destinations hosts and its variability and finally trying to identify and analyse main causes for obtained results, whether they are network-related, socio-economic or user behaviour-related causes.

1.1 Thesis proposal

The growing popularity of Peer-to-Peer networks, that are responsible for a major fraction of the current Internet traffic, called the attention of service providers for the absolute need of assuring its good functioning. Everyone responsible for this task has to know real needs of P2P consumers, to predict their future needs and eventual failures that can occur on these systems.

Since some protocols that make use of this kind of systems for sharing files, such as BitTorrent, are relatively recent and undergo repeated modifications and since new protocols are continuously appearing, an analysis of their functioning behaviour is necessarily incomplete.

Having all these issues in mind, this study intends to present a detailed analysis of BitTorrent systems, focusing on different aspects of their functioning behaviour:

1. The temporal evolution of the geographical localization of peers and their availability;

2. Evaluation of Round Trip Time values, observing its dependencies particularly on distances between origin and end-hosts, the hourly period of the day and the Internet connectivity.

1.2 Thesis contributions

Due to actual popularity of BitTorrent applications, Vuze [37] - one of its clients - was used to obtain the necessary data for the study that will be carried out in this thesis. Using Vuze, it was possible to obtain log files with IP address information, as well as the Port number used by each identified peer on this application. These log files were a very important raw material for the development of this thesis, allowing the geographical localization of involved peers and the calculation of the average Round Trip Time values between hosts, among other important tasks.

We expect that results achieved and conclusions that were taken in this thesis regarding the peer level characterization of P2P networks can be useful to improve the performance and the QoS mechanisms of these systems in order to provide a better service and avoid future congestion problems. Furthermore, these conclusions can also play an important role on the technical design, deployment and management of new real-time multimedia content distribution services based on P2P architectures, like Internet Protocol Television (IPTV) and Video on Demand (VoD), as well as on the definition of their commercial and marketing plans.

1.3 Related work

Due to the significant growth of P2P networks over the last few years, some works [1-7] have been done in order to better understand how they work, trying to optimize their performance and avoiding future overload problems on Internet traffic.

In [1] a survey and comparison of structured and unstructured P2P networks are presented. In [2,3] methods to identify P2P traffic through transport layer behaviour are suggested. Furthermore, in [4,5,6,7,21] several P2P measurement studies focused on the analysis of the topological characteristics of these networks were conducted.

Since BitTorrent is one of most popular P2P applications, it is also a big target of evaluation. Therefore, several studies on modelling and analysis of BitTorrent applications have been conducted, aiming to improve the performance and to characterize these systems at a peer level. In this way, some new mathematical models for BitTorrent were proposed in [8-12].

In [8] a simple fluid model for the BitTorrent application is presented and the steady state network performance was studied. In [9] a new strategy of the peers' selection is proposed in order to make the download even faster. With the aim of studying the performance of piece scheduling decisions made at the seed, a stochastic approximation was proposed in [10].

Reference [11] presents BitProbes, a system performing measurements to Internet end-hosts, analysing their geographical distribution and upload capacity. Some changes in BitTorrent systems are suggested in [12] in order to facilitate efficient streaming between peers as well as providing soft Quality of Service guarantee.

As it was already mentioned, some BitTorrent simulation analysis were also made in [13-17], aiming essentially the BitTorrent traffic characterization.

In [16] high level characteristics and users behaviour of BitTorrent were tested analysing the activity, availability, integrity, flash-crowd and download performance of this application.

In [17], using the Multiprobe framework, measurements of BitTorrent and Internet were conducted and statistically correlated with location, route, connectivity and traffic.

In [13], the peers' performance corresponding to the share of a unique torrent file of 1.77 GB of content, the Linux Redhat 9 distribution, was evaluated during 5 months in terms of throughput per client during the download and the ability to sustain high flash-crowd. In this work, a geographical distribution of involved peers was also done.

A geographical, organizational, temporal, network robustness and peer activity analysis of BitTorrent application for a significant number of files, 120 files shared from December 2003 to March 2004 was presented in [14]. In [15] several studies on the popularity, availability, content lifetime and injection time, download speed and the pollution level for this application were also presented.

In [22], a study with some similarities to the one that was made in this thesis was conducted at Instituto de Telecomunicações, University of Aveiro, focusing on the

geographical distribution and availability of BitTorrent, but restricting its utilization only to video files sharing and comprising a shorter period of time - 8 days only.

Traffic characterization of other P2P applications was also made on several works, such as [4, 18] for Gnutella and [19, 20] for Kazaa.

Since BitTorrent characterization works are still incomplete and inconclusive, this thesis tries to bring a more complete analysis of peers involved in this application, having in mind main goals that we want to achieve: a peer-level characterization of the geographical distribution, availability and RTT characteristics.

1.4 Thesis outline

This thesis is divided in 6 chapters and 1 appendix disposed as follows:

➤ Chapter 1 – Introduction

This is the current chapter and it is subdivided in 5 sub-chapters. It starts with a brief introduction about the thesis, followed by the proposal and the main contributions of this study. Finally, some related works are mentioned.

➤ Chapter 2 – State of Art

This chapter is divided in 4 sub-chapters. After a brief introduction, a detailed description of Peer-to-Peer networks is presented, specifying file sharing as one of its application; then the BitTorrent Protocol functioning behaviour is discussed in detail and Vuze is presented as the BitTorrent client that was chosen for this work.

➤ Chapter 3 – Methodology

This chapter explains how measurements were conducted and what was the methodology used to obtain the data that will constitute the basis for the proposed work.

➤ Chapter 4 – Peers localization and availability

This chapter and the next one are both dedicated to results. In this chapter, the geographical localization of the peers involved on the downloaded files is presented, together with the availability of those peers during the period of analysis.

➤ Chapter 5 – Round Trip Time

This chapter presents the Round Trip Time values measured between BitTorrent peers, as well as some relationships that exist between variations on obtained values and possible causing factors, such as the Internet connection type and the time of the day when those measurements were made.

➤ Chapter 6 –Conclusions

This chapter focuses on main conclusions that can be extracted from this work, as well as some suggestions about future work that still needs to be developed.

➤ Appendix A

This appendix contains most important shell script commands that were developed as well as an explanation of their functionality.

1.5 Notation used

The acronyms used in this thesis are explained after the *Conclusions* chapter, on the *Acronyms* section, as well as on their first occurrence in the text.

Bibliographic references used in this work are invoked in straight parenthesis and are presented in the *References* section at the end of the document.

2. State of Art

2.1 Introduction

Nowadays, P2P networks are very popular and should contribute with a significant fraction of the Internet traffic due to their favorable architecture characteristics like scalability, efficiency and performance. In fact, previously conducted studies have proved that P2P systems are responsible for a major portion of all Internet traffic. According to Sandvine's research, Peer-to-Peer (P2P) traffic remains dominant in the upstream direction with 61% of the network traffic and more than 22% per cent of the downstream bandwidth consumption around the world [38]. Ipoque's Internet Study 2007 [32] presents an analysis of the Internet traffic realized between August and September 2007 in Australia, Eastern Europe, Germany, Middle East and Southern Europe. As can be seen from Figure 2.1, taken from this study, P2P network contents were the most popular among all the Internet traffic.

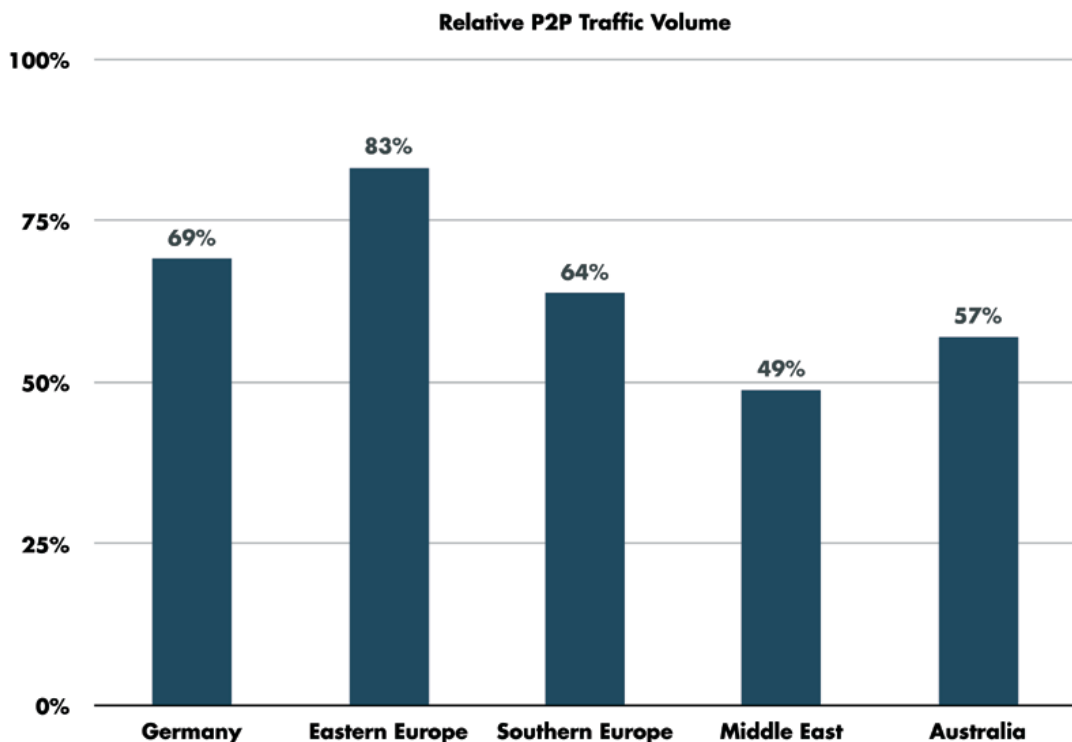


Figure 2.1 - Relative P2P traffic volume [32].

This study has also analysed the most popular P2P protocols: as can be seen from Figure 2.2, BitTorrent is the most used P2P application in almost every world region, with the exception of Southern Europe where eDonkey was identified as being even more popular than BitTorrent.

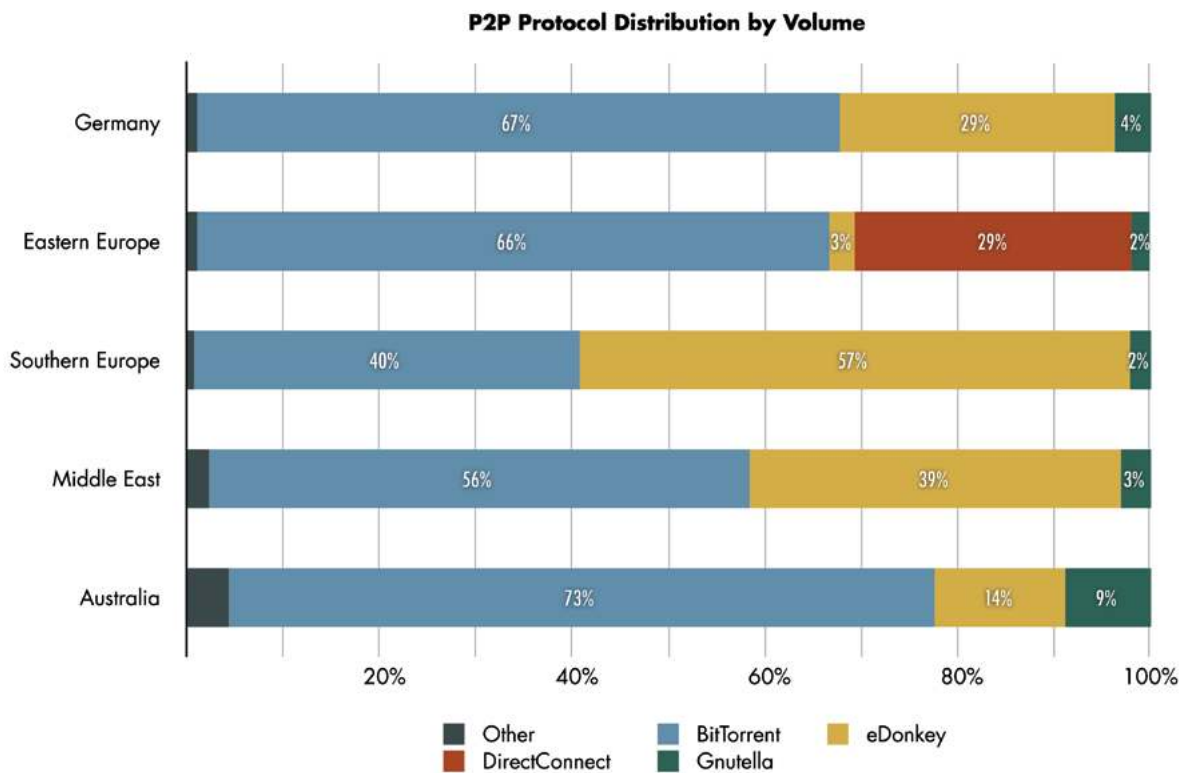


Figure 2.2 – Traffic volume distribution for most popular P2P protocols in five different world regions [32].

P2P applications became so common due to the facility they provide to obtain for free all kind of files, such as music, video, games and books. However, they also bring some legal and management problems:

- Systems and networks security become more vulnerable, since hackers can use this application to make attacks in a wide scale.
- These applications can prevent other kind of traffic and critical applications accessing the network, because the traffic they generate can easily congest the network.
- Users can easily download copyrighted files causing big legal troubles.

In a way to avoid first two problems just mentioned above, Intrusion Prevention Systems (IPS) are used to control those applications and the network traffic generated by them. These systems provide network protection, identifying and blocking threats in real time.

In the next sections of this chapter, P2P networks will be characterized, particularly the BitTorrent P2P file sharing system that was chosen as the target application to be evaluated in this thesis.

2.2 P2P networks

Peer-to-Peer networks are computer connections through the Internet that allow users to share their resources, such as computation power, data and bandwidth, acting as clients and servers simultaneously. In opposition to the Client/Server model, there is no need for a central control.

As a client, each peer can query and download what needed from others peers and as a server, it can supply objects to other peers. In a P2P network, each peer follows basically four phases:

- query for objects using the P2P Routing Protocol and the P2P Location Protocol, in order to find peers that contain such objects;
- join to the P2P system to obtain some information about their neighbours and to inform what objects it has;
- start downloading objects;
- leave the system.

2.2.1 P2P architecture

Depending on the type of connection between peers, P2P networks can be characterized on three different connection architectures that will be described in the next paragraphs: centralized, decentralized and hybrid.

➤ Centralized architecture

In this architecture each network has a Central Server where all peers have to log in if they want to access the network. Figure 2.3 presents a typical design of a centralized P2P network architecture. The Central Server contains information about all files in the system. Therefore, when a peer sends a request, the server answers with the list of available files and the host contact for information. After selecting the claimed file, the peer will directly contact intended peers to download it.

However, when a pair of peers for some reason, can not establish connection with each other by themselves the Central Server will help, in some specific centralized architectures, to establish the connection between them, for example when the host peer resides behind a personal Firewall. In this situation, the server contacts the host peer in order to establish a connection with the downloading peer.

In some advanced P2P applications, it is also possible to connect with more than one peer simultaneously and therefore, download from multiple hosts.

This architecture presents a good performance for search requests and it is commonly used in small networks. However, it is not scalable enough and the server is a single point of failure and a bottleneck for big networks. Hackers can easily attack these networks disabling the Central Server.

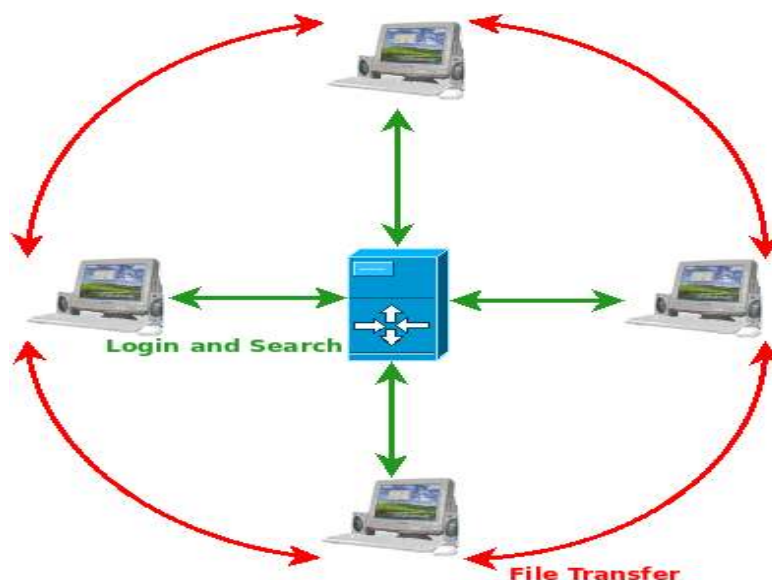


Figure 2.3 – Design of a centralized architecture.

➤ Decentralized architecture

Decentralized architectures do not have a Central Server to log in. In this case, peers send a request to all peers in the network, which will reply with detailed contact information. When the user selects a file to download, directly contacts the host peer and it starts downloading. As it happens in centralized architecture modality, peers can also communicate with more than one peer.

An important point in this structure is that it scales for large networks and it is more difficult to be attacked by hackers because of its distributed control. The figure below depicts the main design aspects of such architecture.

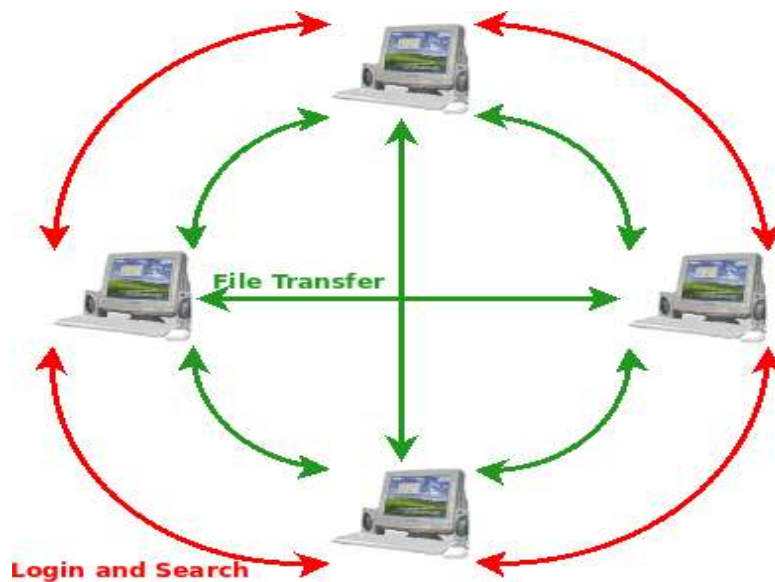


Figure 2.4 – Design of a decentralized architecture.

This type of networks can be divided into structured and unstructured networks, as it will be described next.

- *Structured*

In structured networks, the nodes graph and the data location are tightly controlled to increase the efficiency of data discovery. For such control, Distributed Hash Tables (DHT) are used to discover data location providing a key-based route. Therefore, these

topologies are efficient to look up for a file, mainly for a rare one, since for any given key it is possible to locate the peer that is responsible for the correspondent file of such key.

- *Unstructured*

In opposition to structured networks, these networks have no control over nodes connections and data location; thereby nodes are free to choose their neighbours and store data by their own.

This network topology uses a flooding method to find out data stored by peers. Each node sends a search request to intermediate peers that will then resent it to intermediate peers on neighbourhood and to all peers on the network, through replication and forwarding. They will check out the requested data in their lists of stored items data and will reply to the search request.

The advantage of the unstructured P2P architecture is the easy way for localizing a common file and for peers to join and leave the system. However, in this architecture it is harder to localize rare files when compared to structured networks. Furthermore, peers easily become overloaded due to the growth of their loads with the increasing number of requests and the system size, being a non-scalable system.

- Hybrid architecture

This architecture is composed by Supernodes, which have basically the same function as the Central Server of the centralized architecture. They are distributed over the network, providing a larger network. The P2P application establishes communication between Supernodes, transmitting available files on the peer system.

In this case, a peer sends a request to a Supernode that will forward it to other Supernodes of the network. Responses are given to the primary server that compiles and sends them to the peer. In a similar way to others architectures, when a user decides to download a file it contacts the host peer directly and the transfer is ready to start.

As mentioned before, in advanced P2P applications it is possible to establish communication between various peers.

As in the centralized architecture, if for some reason the downloading peer can not connect with the host peer, the Supernode will contact the host peer to do it. Then, the downloading peer is able to connect to the host peer through the Supernode.

Therefore, this architecture is a combination of both centralized and decentralized architectures, acquiring the advantages of both. The hybrid architecture provides a good performance for search requests, scales to large networks and it cannot be easily attacked by hackers due the distributed and dynamic nature of the Supernodes.

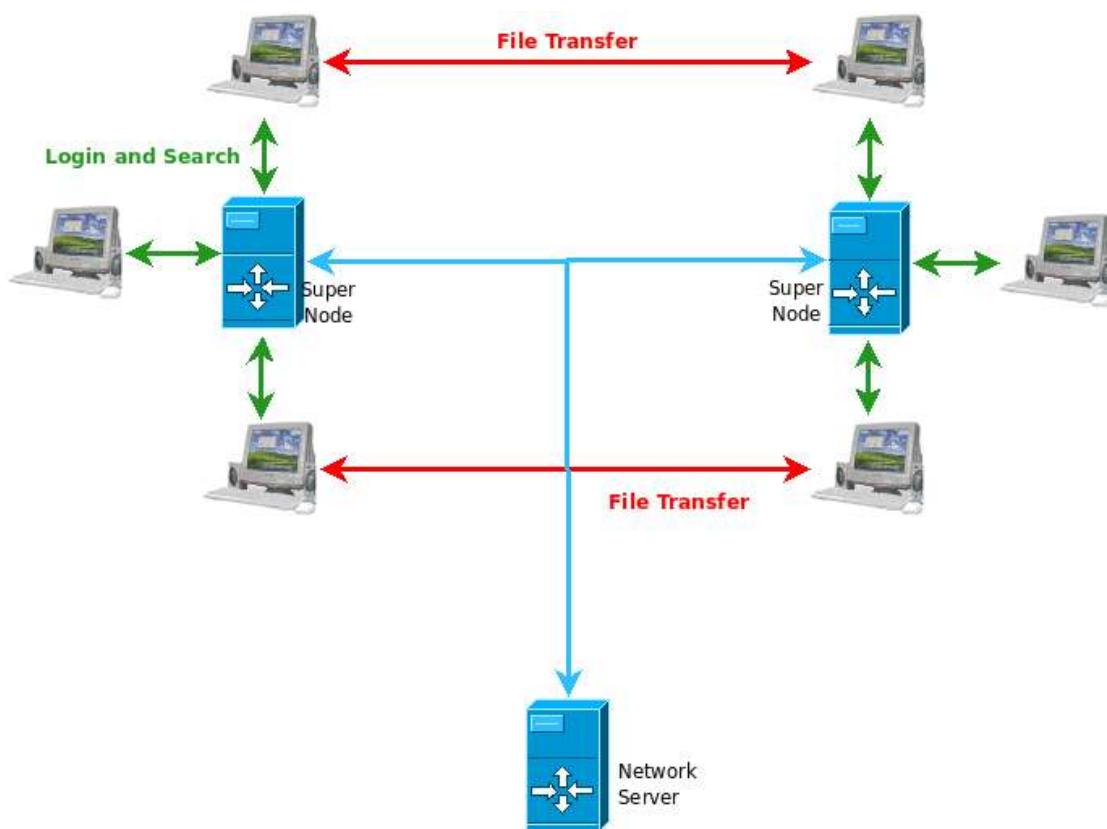


Figure 2.5 - Design of a hybrid architecture.

2.2.2 P2P generations

These networks can also be divided according to their time period of appearance and to their purpose, leading to three generations of P2P systems:

- **1G** is designated as the first P2P network generation, which aim was to be easy and quick in looking up for files, based on a Central Server. They have the big advantage of being simple but they are not scalable and efficient regarding query processing. Napster belongs to this generation, but it was obligated to shutdown by the Recording Industry Association of America (RIAA) based on the argument that some files located at the Central Server could contain illegal contents, such as copyright violations.
- In order to solve the problems of the first generation, **2G** networks use the decentralized architecture and Distributed Hash Table (DHT) technique providing both load balancing and deterministic search guarantees. However, this generation still presents some few problems, including its lack of scalability, implementation difficulties and weak security capabilities, being more disposable to copyrighters' violations. This generation includes Gnutella, Plaxton, Chord and Can.
- **3G** generation aims to provide resilience, through the use of object replication, expanding the connection number between nodes and using a special structured topology. It also aims to provide higher performance to P2P networks in terms of security and integrity. BitTorrent system belongs to this generation.

2.3 BitTorrent protocol

BitTorrent is one of most popular P2P networks, being responsible for more than 50% of the P2P network traffic. Some BitTorrent studies present several advantages for this system, such as its good scalability, good offer of download quality to users, high level of integrity on the contents and meta-data and its capability to handle efficiently flash-crowds effects. This system aims the optimization and fairness on the contents distribution, which is greatly facilitated by the Tit-for-Tat methodology and based on two important algorithms to choose the best pieces and nodes: the Rarest First Algorithm and the Choke Algorithm.

BitTorrent is an unstructured and decentralized P2P overlay network application for file sharing. To download a file, users have firstly to access the web site that contains the *.torrent* files, such as Suprnova.org, Youceff.com, Piratebay.org and download files they want. The

.torrent file contain important meta-data information and it is responsible to point to the respective tracker, which is the Central Server of an intern network whose users are sharing the same file, called swarm. The tracker contains a list of all users of the swarm. In this way, when a new user appears in a swarm the tracker gives the information of some peers, approximately 50, that have the requested file and with whom the user will directly connect to start downloading. The list of known peers from each user is called peer set. Each user sends periodically to the tracker, each 30 minutes, information of his status, as well as each time instant where he had left the system. Besides, when a user has less than 20 peers in its peer set, it will request the tracker a new list of peers.

Each file is divided into pieces of same size, usually of 256 Kbytes and each piece is sub-divided in smaller ones, called sub-blocks, allowing users to start uploading as soon as they complete the first piece and then avoiding free-rider users which are very common in other systems.

Users are divided in two types, seeds and leechers. Seeds are users that already have the entire file and leechers are those that do not have the entire file downloaded yet. The first seed is the one that has the first copy of the file.

➤ **Choose the proper piece**

The choice of proper pieces is an important factor for the system efficiency. It is necessary to avoid the problem of the last piece, which happens when the pieces distribution is asymmetric. In this situation, when users having rarest pieces leave the system, those pieces are inalienable, making the end of the download impossible to other users.

➤ **Local Rarest Piece First policy**

To avoid this situation, BitTorrent uses the Local Rarest Piece First policy that gives all users information about the number of existing copies for each piece. Using this information, users will randomly choose pieces that have less number of copies, rarest pieces. The pieces information is updated every time a piece is added or removed. The BitTorrent system knows three exceptions for this algorithm, called Random First piece, Strict Priority and Endgame Mode policies.

➤ **Random First Piece policy**

This exception is used when a new user does not have four complete pieces. In this case the choice of pieces to download is randomly made, in such a way that the user quickly bootstraps it and starts sharing its pieces with other users. When a node completes its fourth piece, it switches to the Local Rarest First policy.

➤ **Strict Priority policy**

This policy is used when a node requests a block of a specific piece of the file. In such situation, all blocks of that piece have the highest priority, being requested before any other one, in order to complete the piece as soon as possible. This policy is very important because recently completed pieces can be transferred, thus minimizing the number of incomplete downloaded pieces.

➤ **End Game Mode policy**

This policy is used by nodes that are finishing the download of a file, making the end of the download faster. In these situations, the node can request all missing blocks simultaneously to multiple peers without any restriction. Each time the node completes a block, all pendant requests are obviously cancelled.

➤ **Choose the proper node**

It is of node responsibility to maximize its own download rate preferentially uploading to neighbours that provide best download rates, following the Tit-for-Tat policy and avoiding free-riders. In this way, the node sends data information to its favourite nodes and receives the data information back from the nodes that want to do so.

BitTorrent follows a reciprocity method in order to maximize network resources and download rates. This method is obtained by the Choking Algorithm and differs if the user is a leecher or a seed.

➤ **Choking algorithm**

- *For Leechers*

Each node can unchoke up to four connections, uploading to 4 different nodes at the same time. To choose which nodes will be unchoked, a choking algorithm is used. Each 10 seconds the node receives a list of the nodes that want to connect to it, ordered on a decreasing order of the download rates. The first 3 of this list are chosen and a fourth one is added randomly each 30 seconds.

- *For Seeds*

In initial versions of BitTorrent the seed unchoked nodes according to its upload capacity, choose nodes each 10 seconds as it happens in the leechers case. This algorithm works well for maximizing network resources, such as the bandwidth offered by the seed, but it is not fair on favouring the download capacity of independent nodes, thus stimulating free-riders.

New versions follow a similar procedure to the one that is used by leechers. Each node has up to four unchoked nodes and each 10 seconds receives a list ordered according to the time that each node becomes unchoked. Each 20 seconds the node will choose the first 3 on the list and each 10 seconds, randomly it chooses a fourth node. After the 20 seconds period and during next 10 seconds, four nodes on the top of the list are chosen.

➤ **Optimistic unchoked policy**

In addition to nodes that were already unchoked by the choking algorithm policy, each 30 seconds the node will randomly choose one more node to unchoke all neighbours. This policy evaluates the download rate of the nodes, allowing to find out those that can offer higher rates but also giving to new nodes, those that don't have any completed piece, the opportunity to start downloading.

➤ **Anti-snub policy**

Whenever a node unchokes another one and if after a 60 seconds period the unchoked node does not receive any complete piece, it will consider it a snub node and will choke the

connection. If the node is a leecher, the connection just would be unchoked by the optimistic unchoked policy. The seed will never unchoke a snub node.

In case of a leecher, if the snub node sends a sub-block of the file in a certain period of time, usually 45 seconds, the node stops seeing it as a snub node. The seed, on the other hand, never stops to consider it as a snub.

2.4 Vuze

Vuze is the BitTorrent client used in this thesis to download desired files. It is the world's most popular entertainment platform for high resolution digital content [37].

Using Vuze to download a file, the first thing to do is to search on a web site, such as the TorrentZ [36], for a torrent that contains the needed file and download it. With this file, we are able to use this P2P application to start sharing files with other peers involved, being only necessary to open the downloaded *.torrent* file.

With this BitTorrent client it is also possible easily to obtain the necessary information about the involved peers, such as their IP addresses and the Port number used on this application. For this goal, it is only necessary to set some simple configurations, as will be explained next. On the options field, the logging to a file with the maximum allowed size (500MB in this case) was enabled as can be seen in Figure 2.6. Vuze splits this total size in two parts, in order to save a copy of the file and do not loose information when it reaches the maximum size, writing relevant information on the log file. Using appropriate shell script programs to manage this file (that were developed in this thesis), the IP addresses and Port number of each involved peer were easily extracted. With this information and using again shell script programs, it was possible to geographically localize peers and study their availability and the corresponding Round Trip Times (RTTs) during the temporal period under analysis.

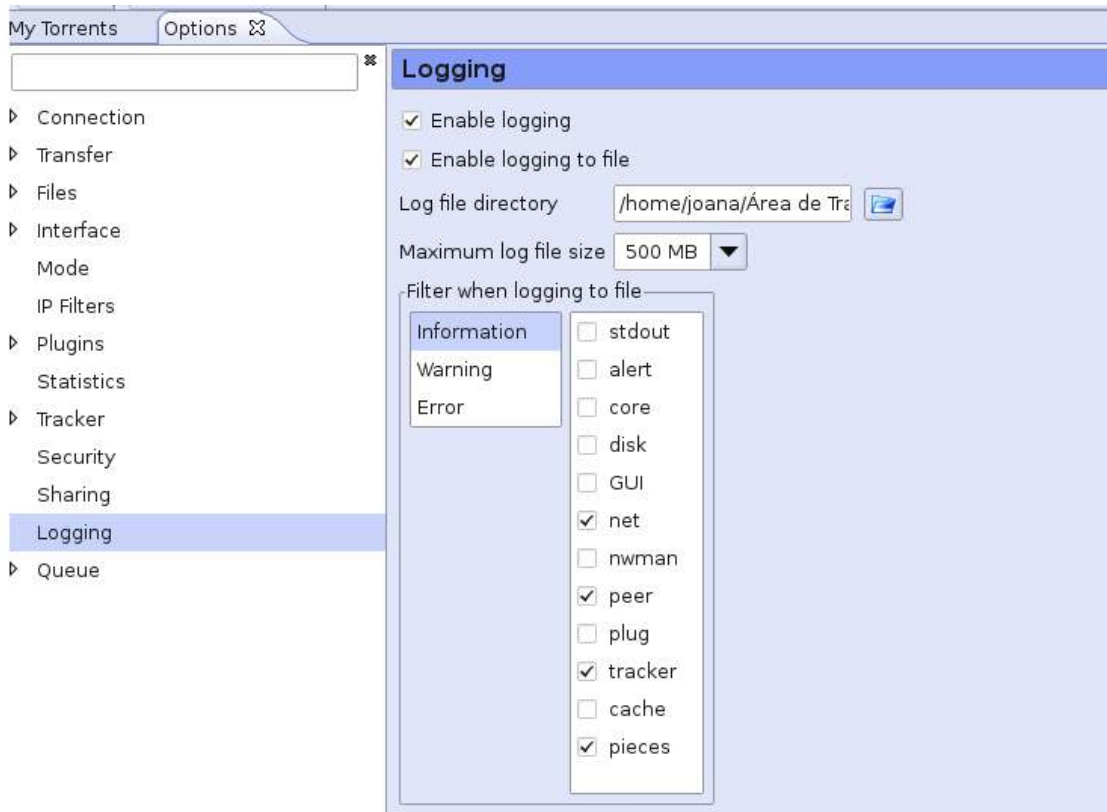


Figure 2.6 – Vuze logging option field.

2.5 Summary

This chapter presented an overview about P2P networks, since they are the basis for the study conducted on this thesis. This part of the thesis aimed to give an explanation about the basic functioning principles of P2P networks and the main involved concepts.

Finally, the P2P application that was intensively used in this thesis, the BitTorrent application, was described in more detail together with its Vuze client - the client that was chosen to share files on our experimental set up.

3. Methodology

3.1 Introduction

The main objective of this thesis is the characterization of P2P traffic and their users behaviours, focusing mainly on the characterization of the peers' availability per country and continent and the analysis of the Round Trip Time (RTT) corresponding to hosts involved on the download of specific chosen files.

Results obtained in this study were extracted from January to August 2008 in two distinct Portuguese cities, Aveiro and Coimbra, using two different Internet connections on each case: CATV 12 Mbps Internet access and ADSL 4 Mbps Internet access, respectively.

Since BitTorrent is currently the most used protocol for P2P file sharing systems, one of its clients, Vuze, was chosen to download a set of pre-specified files. The choice of Vuze was also related to its high popularity as a file sharing software in the major part of the world. Furthermore, with Vuze it is very simple to obtain log files containing the IP addresses and the application TCP Port used by involved peers, which was very helpful for obtaining results needed for this thesis.

In next sub-chapters, the process of selecting files for download will be explained, as well as the procedure that was used to obtain results regarding the localization of peers, peers availability and RTT.

3.2 Files selection

After deciding the best BitTorrent client to use, the next decision was the selection of which files to download. Two criteria were considered for this decision. First, six different categories were chosen and for each one, files having more peers involved at the moment, seeders and lechers, were chosen. With this purpose, a search was done in web sites like Torrenz, which is a search engine that finds out torrents from major torrent sites such as Mininova, Demonoid and The Pirate Bay.

The download of large numbers and different types of files let us to access to an extended sample of the P2P world population, enabling a better evaluation and understanding about preferences and costumes of each country and continent. With this information, it can also be possible to analyse population interests and to conclude how Internet is used and what are their users' needs. Besides, in order to contribute to a better management and planning of such P2P systems and infrastructures, it is essential to understand the relationship between RTT values and distances between hosts and the Internet connection, physical wire, distance, traffic load, link layer technologies and other network characteristics.

Six chosen file categories were then: *2008 movies*, *Music*, *Animated movies*, *French movies*, *Indian movies* and *Linux distribution*. As it was already mentioned, for each category chosen files were ones having more peers involved. Next, we will present selected files as well as the total number of peers involved on each category:

➤ ***2008 movies***

Number of involved peers: 52111

Chosen files:

1. 21
2. Iron Man
3. Loose Change
4. Street Kings
5. The Forbidden Kingdom
6. Untraceable

➤ ***Music***

Number of involved peers: 3622

Chosen files:

1. Dido
2. Florida Mail on Sunday
3. Justin Timberlake recrimination
4. Mariah Carey

5. Miles Davis Love songs
6. Top 100 Trance and Techno Party songs of all time

➤ ***Animated movies***

Number of involved peers: 4303

Chosen files:

1. Enchanted
2. Ratatouille
3. Shrek: The third

➤ ***French movies***

Number of involved peers: 4433

Chosen files:

1. Angles
2. Cloverfield
3. Disco

➤ ***Indian movies***

Number of involved peers: 18645

Chosen files:

1. Aamir
2. Mere Baap Pehle Aap
3. Sarkar Raj

➤ ***Linux distribution***

Number of involved peers: 2179

Chosen files:

1. Fedora-8-dvd-i386
2. KNOPPIX_V5.1.1CD
3. SabayonLinux-x86_64-3.4f

4. Ubuntu-7.10-desktop-amd64

Six different file categories, covering a total of 25 downloaded files and 85293 peers from the whole world were considered in this study. The download of different categories was not made at the same time in order to avoid an overload of files coming into the Vuze client that would cause a decrease on the download velocity due to bandwidth limitations. However, files of the same category were downloaded at the same time in order to obtain appropriate results of the peers' availability, since this study was made mainly per category and not per file. It is also really important to start the evaluation of the peers' availability exactly on the precise moment the files download has ended up (where the number of active peers should be at its maximum value). This will allow taking further conclusions about peers that are protected by Firewalls or NAT/PAT since they will never answer to TCP probes and appear as non-available, even knowing that it is not really true.

The peers' availability analysis ended at the same time for all categories, exactly on August 12th 2008.

In the next sub-chapter, we will present in detail different steps that were taken to discover the localization of the involved peers.

3.3 Peers localization

While the download of desired files was being done, log files created by Vuze were extracted in order to obtain the IP address of each involved peer and the Port number used by this application.

When the download of each category was complete, the IP address and the application TCP Port were extracted to a specific file using shell scripts for an easier management of log files. The next step was to find out the country corresponding to each peer, using the Geoipllookup tool [34]. Such tool is included in Maxmind C Library and it uses the GeoIp Library and a database to find the country corresponding to a specific IP address or host name. After knowing the IP address, using this tool and again a shell script that was specially developed for this purpose, it was possible to obtain the intended peers localization. This

localization allows a better understanding of the P2P file sharing systems' importance on each country analysing which countries have more peers involved and relating this information with other factors such as the quality of network infrastructures, the facility of Internet access, general economic factors and the development level of such countries.

Conclusions about the peers' localization are explained in the next chapter, where the discussion about obtained results is also made.

3.4 Availability and Round Trip Time

The study on peers availability and Round Trip Time characteristics is very important since these factors determine the quality and speed of the file downloading. Obviously, a higher number of peers involved and lower Round Trip Time values between hosts will make files' downloads to become faster and faster.

The Round Trip Time can be defined as the time a packet takes to go from the origin-host to the end-host and comeback. It is easy to understand then the strong dependence between RTT and download speed on these P2P systems. The Nmap tool [35] was used to obtain RTT values, allowing also to get a better understanding of which factors can influence changes on this parameter.

The Nmap tool is a network exploration tool and a security/port scanner that gives very important information about hosts' characteristics, like hosts availability, services offered, operating systems running by hosts and the type of Firewall they are using. By specifying the IP address and the application TCP Port used on this system by a particular host, this tool can perform several TCP port probes, giving the average of the Round Trip Time values.

When, for some reason, the connection between the origin and end-hosts is not possible, the Nmap tool returns the value *-1* on the *final times* field. This value means that the end-host is not available. Peers availability was calculated with this information, considering the total number of end-hosts for whom the origin host could establish a connection, i.e., the final time value was different from *-1*. Having in mind that some of these impossible connections are due to peers that are protected against detection, this does not mean they are out of the system.

These results are not really an indication of all available peers but only an indication of those that respond to TCP requests.

3.5 Shell scripts

A shell can be defined as the main kernel interface that translates commands written by the user, allowing an automatic communication between the user and the kernel. Therefore, shell scripts are files containing a set of executable commands.

The usage of commands available on shell was very helpful on this thesis, because by simply writing some small programs it was quite easy to make measurements, treat available data and obtain main parameters of interest.

First two shell scripts were created to extract the relevant information given by Vuze on its Log files, namely the IP address and application TCP Port of the peers involved on downloading a specific file, and yet to divide such information in different files according to the name of the file that the peer was correlated to.

After obtaining this information, the next step was to obtain geographical localization and RTT measurements. In this way, two more shell scripts were created. It is important to notice that the geographical localization was made only on the first time that we have run such files: after that, the command line that allows localizing peers was blocked. On the other hand, RTT measurements were made a lot of times and for a lot of people during a long period of time, which would be impossible to measure manually. Thus, in this case the use of shell scripts was very important, turning things easier.

Having obtained the intended results, shell scripts were also used to enter commands from other languages, such as Octave [39] and Gnuplot [40], which were used on this work to make further necessary calculations and to create illustrative plots. These plots are shown in next chapters and were impossible to obtain using basic shell script commands. The possibility to call any one of these applications from a simple command line on the shell script file was also an important advantage for our work, because it was possible to reduce greatly the number of necessary files for measuring and processing available data.

Octave is a high-level language, primarily intended for numerical computations, and Gnuplot is a portable command-line driven interactive data and function plotting utility for UNIX, IBM OS/2, MS Windows, DOS, Macintosh, VMS, Atari and many other platforms.

Shell script files used in this thesis are presented on Appendix A.

3.6 Summary

This chapter described procedures that were used to obtain and manage sufficiently enough data in order to fulfill the main goals/requirements of this thesis. In this way, it started with a discussion about the criteria used to choose files that could give more useful information concerning thesis' goals. According to these criteria, files were divided into six categories and for each one chosen files were those with more peers involved, in order to get information on the highest possible number of peers. After the peer selection, it was explained how it was possible to localize involved peers for each file, as well as how it was possible to measure the Round Trip Time between peers and their availability on the system.

At the end of the chapter, a brief description about written shell script commands and their purpose was also presented.

4. Peers localization and availability

Since it is already of general knowledge, P2P networks popularity has been growing in last few years, having a big importance in people's lives and generating the major portion of the current total Internet traffic.

This chapter intends to study and evaluate the importance of P2P networks around the world, making a deep analysis on the usage levels of this way of sharing files between countries.

Evaluating uploads and downloads made with Vuze, it is possible to obtain information on the number of peers involved in the download and upload of specific chosen files using the Geoipllookup tool, the peers' localization can be easily derived.

4.1 Geographical distribution

This section will analyse and evaluate the geographical distribution of peers around the world. It will start with an analysis of the distribution of peers per continent and per country, mainly those countries having the higher number of identified peers. Then, in order to obtain a more significant and meaningful result about the importance of this P2P system, an analysis of the distribution of identified peers normalized by their country population will be done.

Figures depicted below are divided into six categories, representing each one of downloaded files' types: *2008 movies*, *Music*, *Animated movies*, *French movies*, *Indian movies* and *Linux distribution*. For each category, a graph of the peers' distribution per continent is shown, followed by another one illustrating the peers' distribution among the most important involved countries (in terms of the number of peers, obviously).

➤ **2008 movies category**

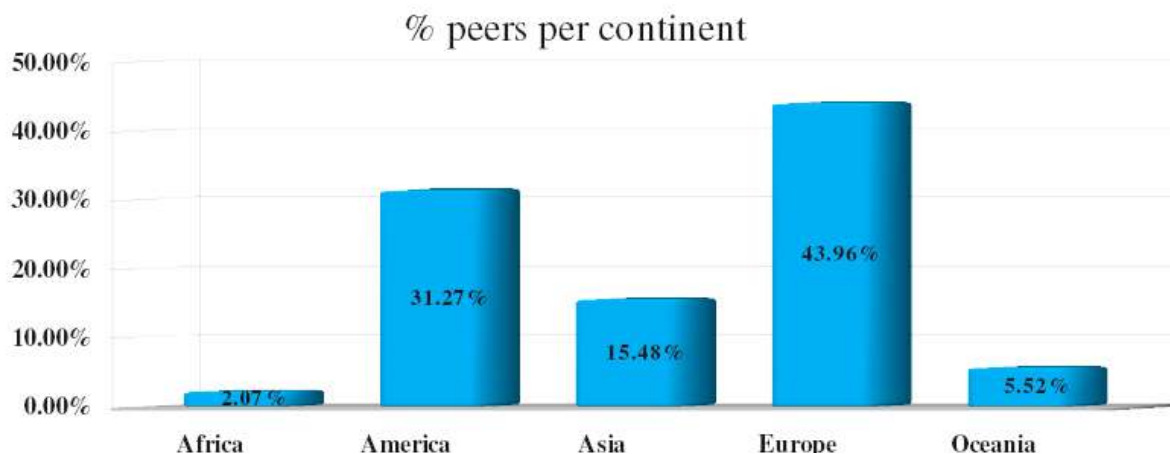


Figure 4.1 - Geographical distribution of peers per continent – *2008 movies* category.

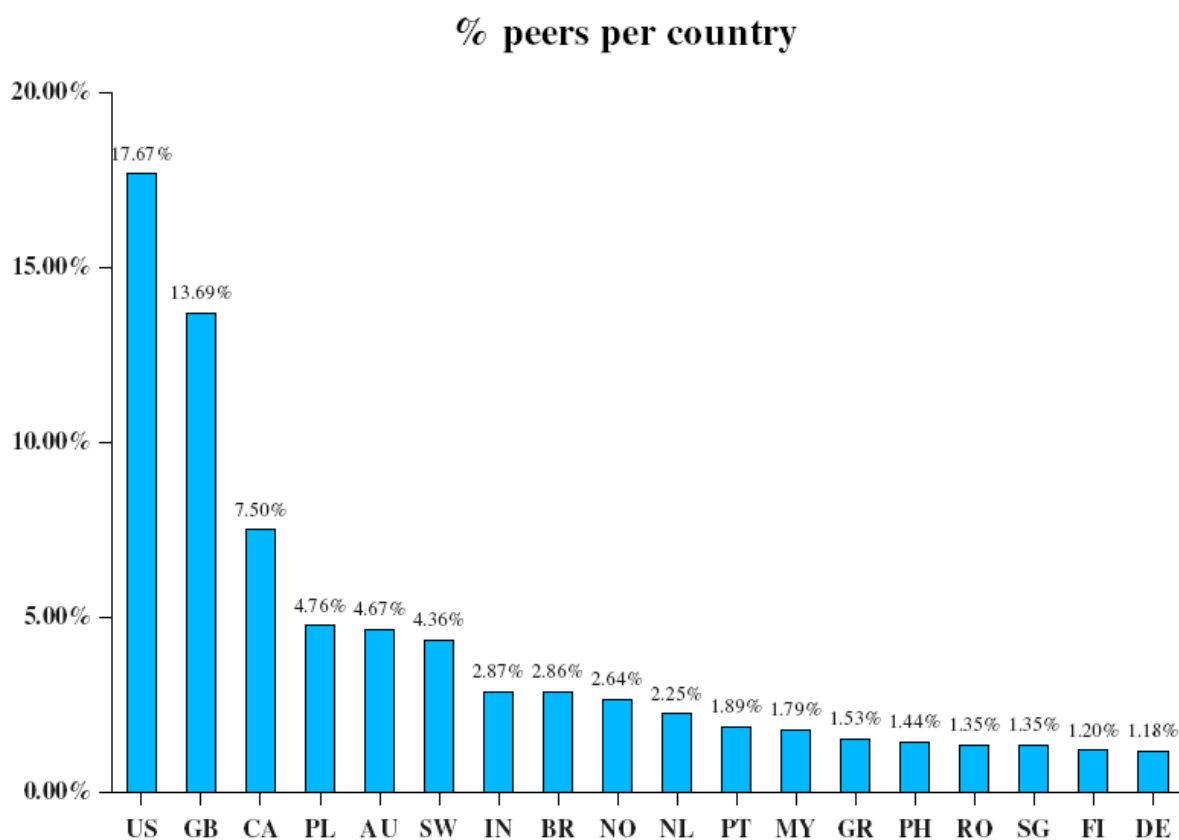


Figure 4.2 - Geographical distribution of peers per country – *2008 movies* category.

➤ *Music category*

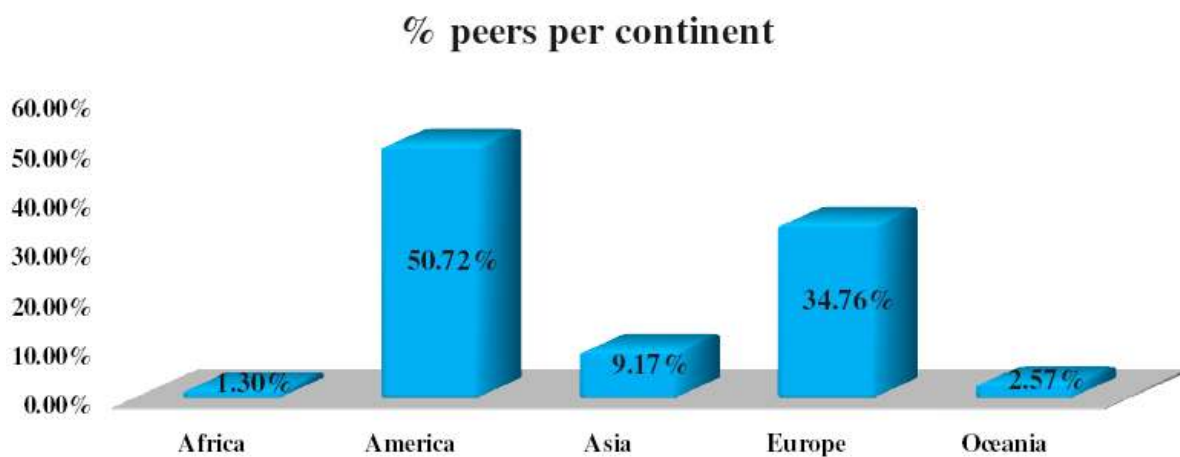


Figure 4.3 - Geographical distribution of peers per continent – *Music category*.

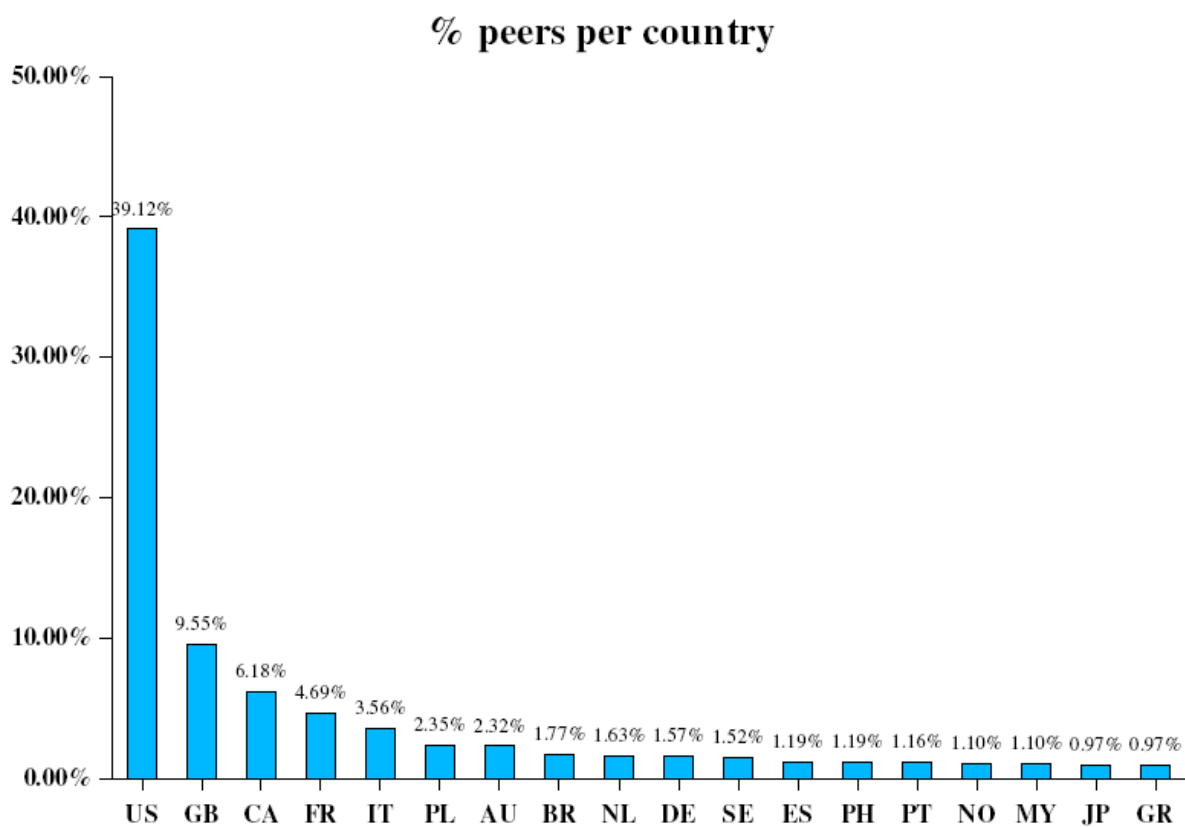


Figure 4.4 - Geographical distribution of peers per country – *Music category*.

➤ *Animated movies category*

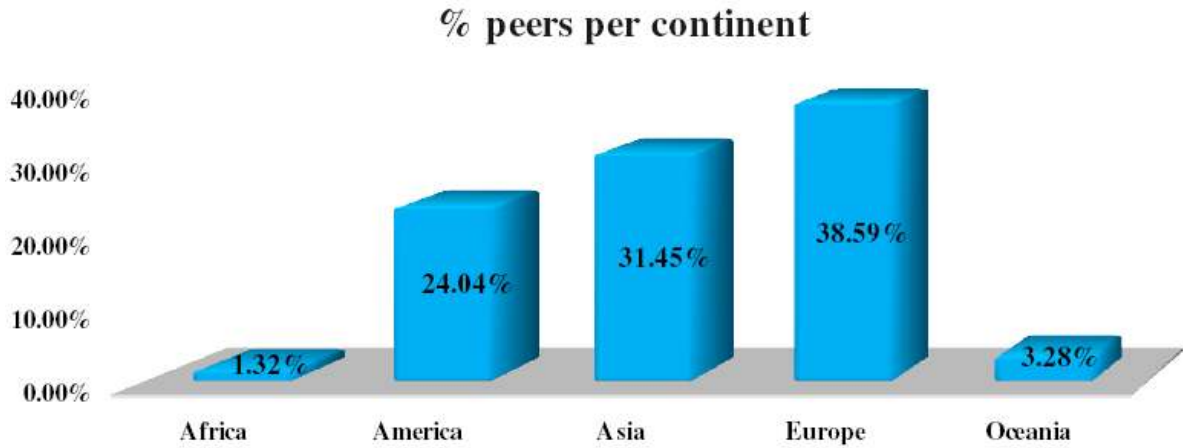


Figure 4.5 - Geographical distribution of peers per continent - *Animated movies* category.

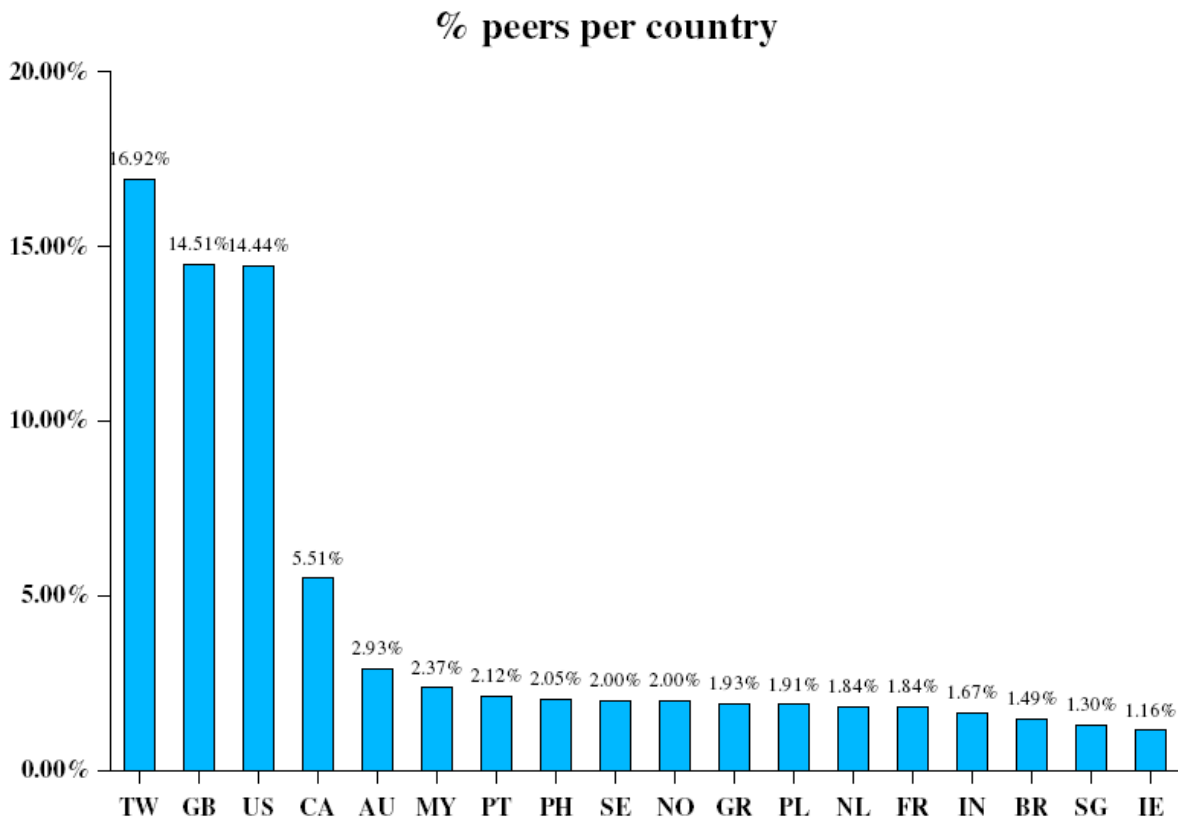


Figure 4.6 - Geographical distribution of peers per country – *Animated movies* category.

➤ *French movies category*

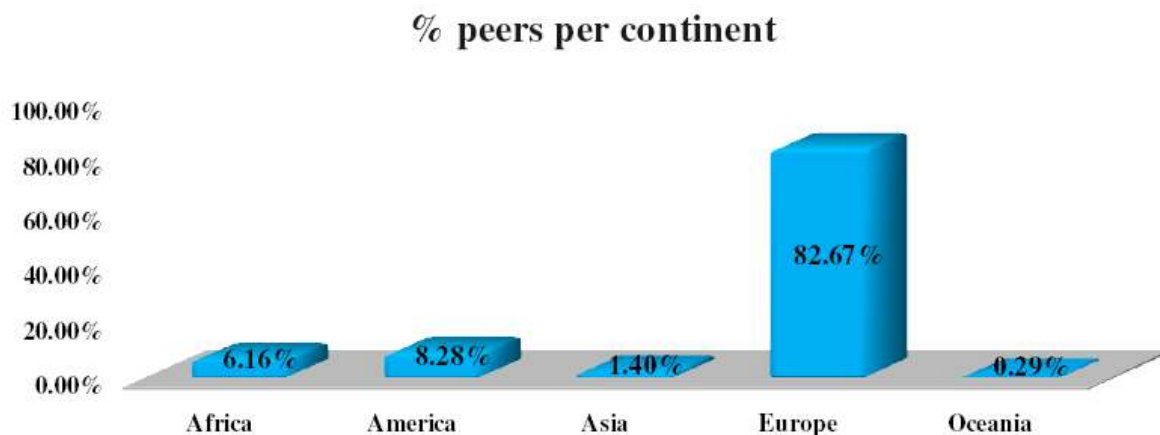


Figure 4.7 - Geographical distribution of peers per continent – *French movies category*.

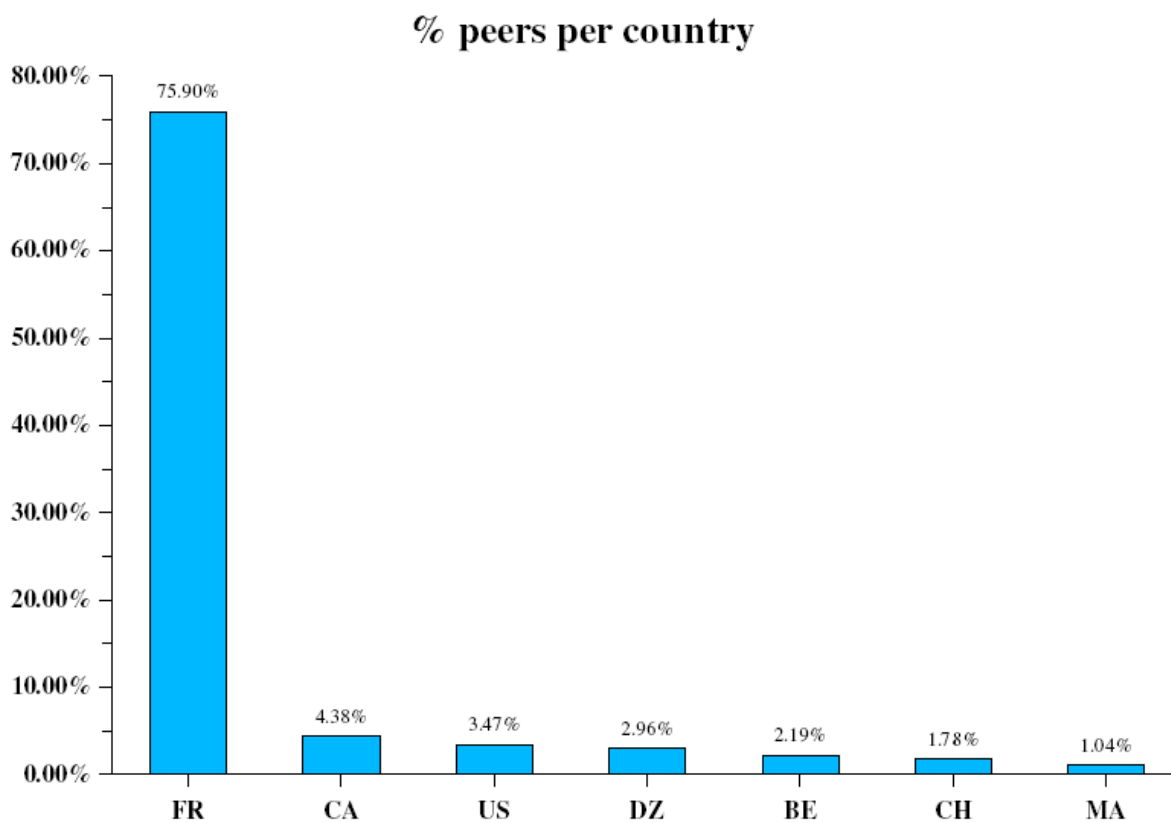


Figure 4.8 - Geographical distribution of peers per country – *French movies category*.

➤ **Indian movies category**

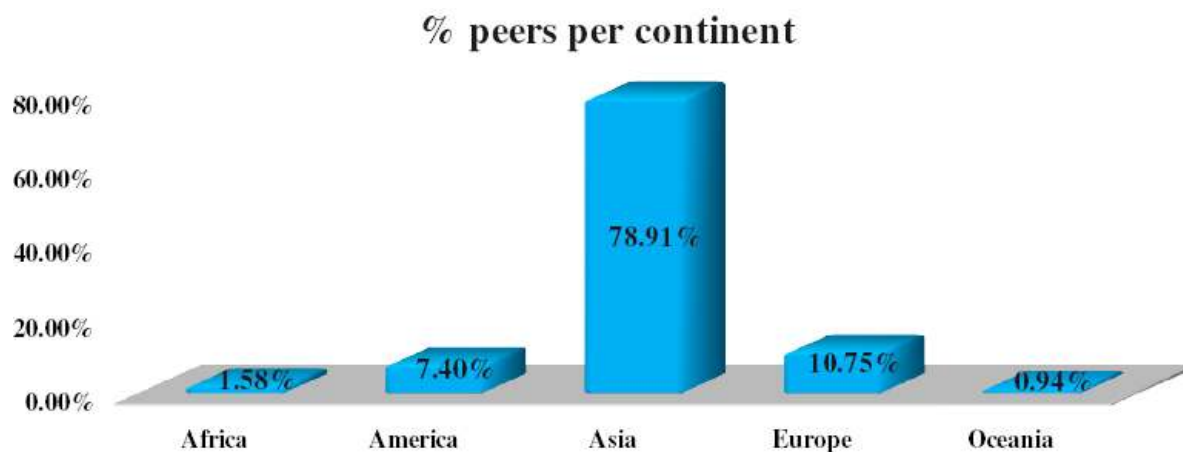


Figure 4.9 - Geographical distribution of peers per continent – *Indian movies* category.

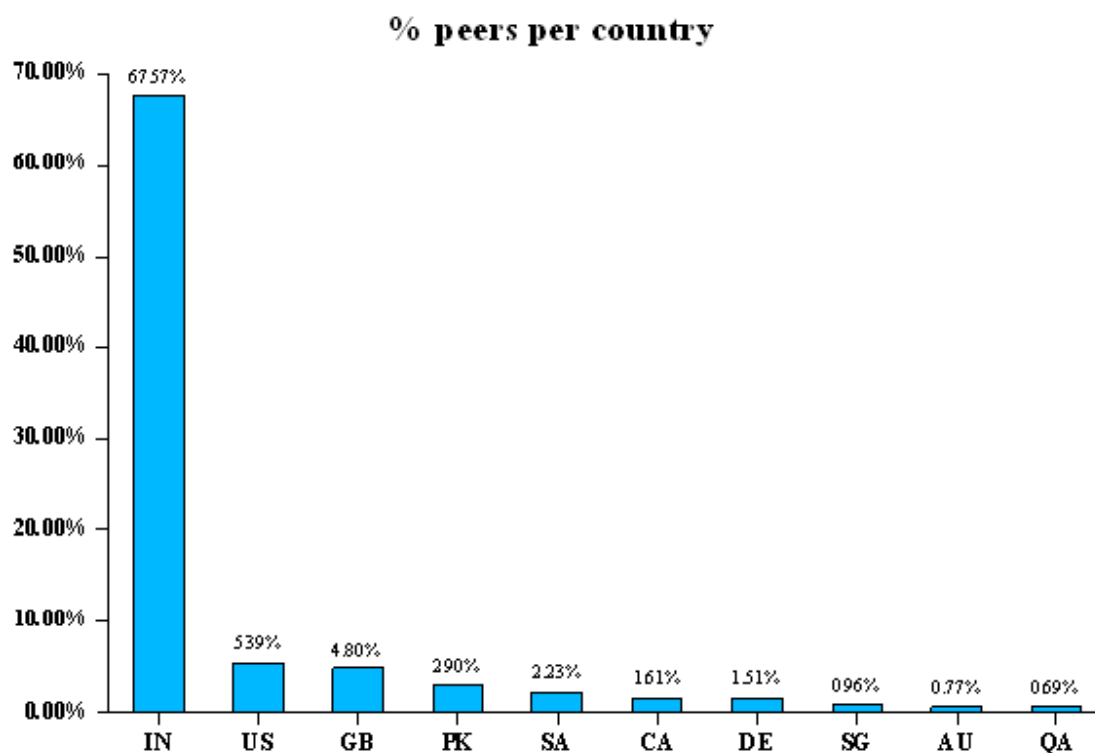


Figure 4.10 - Geographical distribution of peers per country – *Indian movies* category.

➤ *Linux distribution category*

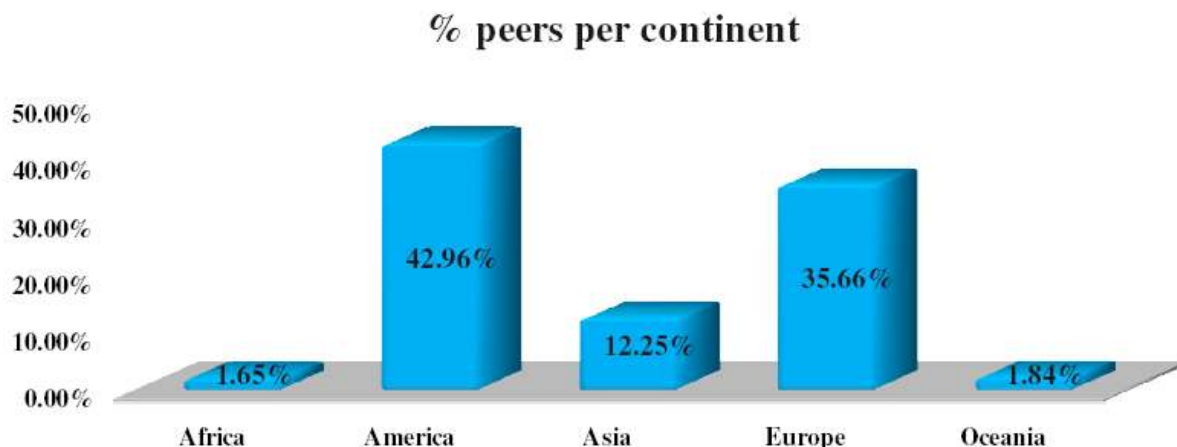


Figure 4.11 - Geographical distribution of peers per continent – *Linux distribution category*.

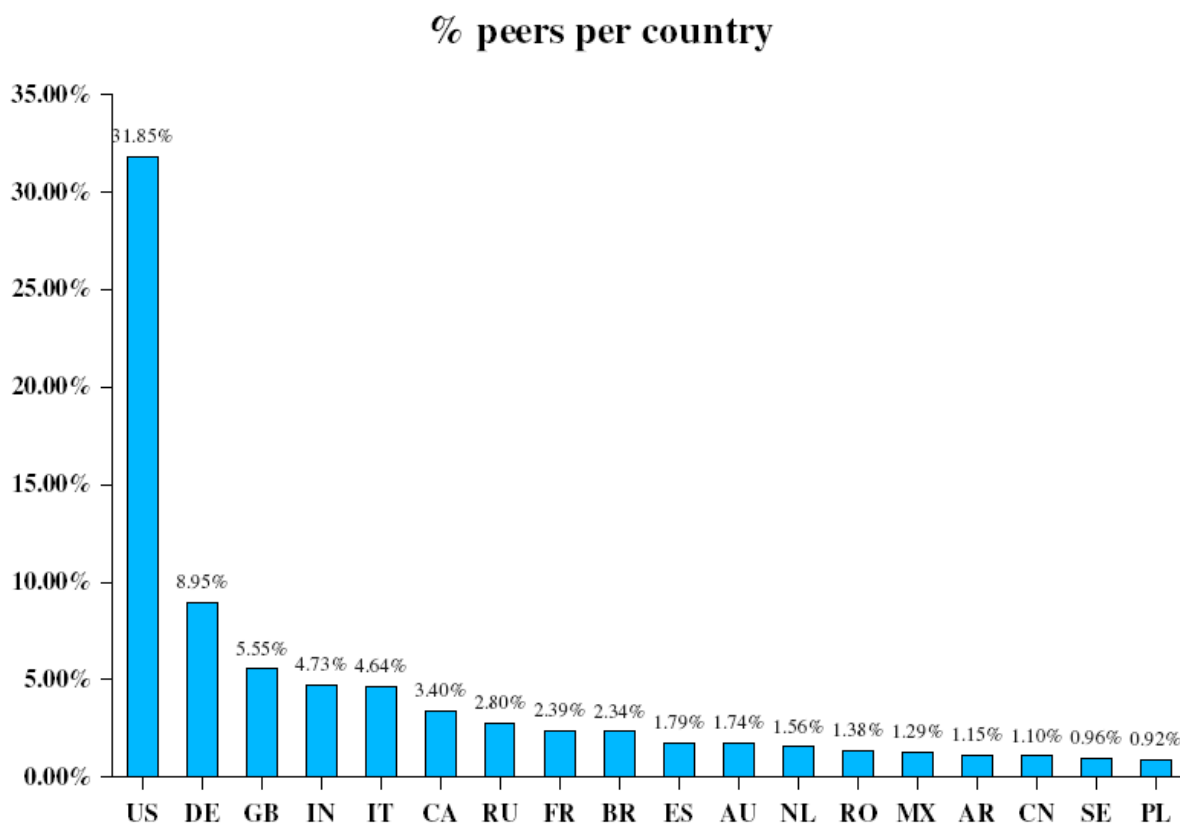


Figure 4.12 - Geographical distribution of peers per country – *Linux distribution category*.

It is not so strange to see United States of America in a place corresponding to the highest percentage of peers, having in mind their characteristics, interests and cultural variability, including people from different origins and races. This cultural diversity, conjugated with the technological development of the country, including excellent network infrastructures and easy Internet access for a major portion of population when compared to other countries explain the strong usage of this type of file sharing. (For example, India has a much higher population but Internet access is not cheap and so fully deployed for the whole population).

It is also important to mention that categories where USA is not in the first position are those that correspond to files with regional content/interest: the *French* and *Indian movies* categories where, with no surprise, France and India respectively, have the domain on the number of involved peers.

China is the country with the highest percentage of population, approximately 22% of the total world population and it is one of the best countries in terms of technological development, with great network infrastructures and facilities in the Internet access. However, this country does not seem to be very relevant in any type of downloaded file. The small use of Vuze in China can be explained by the fact that downloaded files are not of big interest for Chinese people. Note that China has a lot of national production that certainly is more attractive for Chinese people when compared to the international production. Another justification for this occurrence is that China has its own file sharing protocols, which probably makes their torrents to be indexed in websites that are completely distinct from ones considered in this study. Therefore, it is normal to observe a lower usage of BitTorrent from the Chinese users [23].

With the purpose of better localizing countries where P2P file sharing systems have a relevant importance, a normalization of the number of peers by the country population was done. It is important to notice that the percentage of population on each country represents a trivial factor influencing the number of localized peers: the larger is the country population, the bigger will be the probability of having more peers involved.

Thereby, following six normalized maps (corresponding to the six studied categories) depict the relationship between the total number of involved peers on each country and the proportion of its population regarding the whole world population.

Normalized maps confirm the existence of different results when compared with non-normalized ones. The decline of USA after making this normalization is clearly visible, allowing us to conclude that the high number of peers that this country involves it is not so relevant when compared to its population that reaches almost 5% of the world population.

On the other hand, it is interesting to refer that some countries have very low number of inhabitants that places them on top of the ranking even when normalization is done: however, this situation does not mean that they are most relevant countries in terms of file sharing. These are the cases of Dominica, Iceland, Faroe Islands and Slovenia, as can be observed in the above maps.

Next, an individual description of obtained results from such normalization will be done for each category.

➤ **2008 movies category**

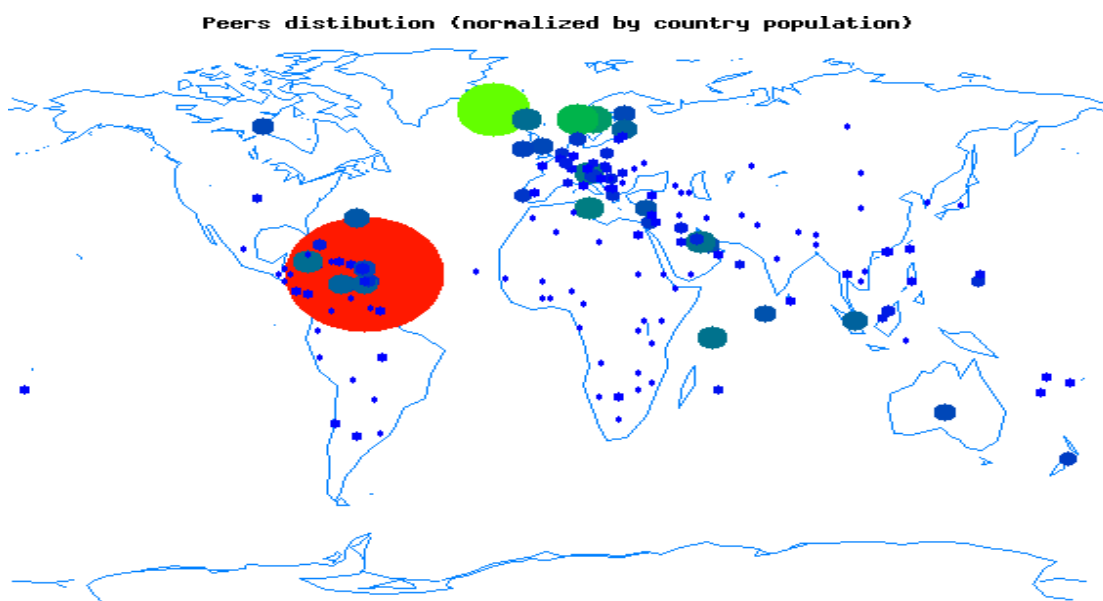


Figure 4.13 – Peers' distribution normalized by country population – *2008 movies* category.

The normalized map of the *2008 movies* category shows that Dominica is the country with the highest relationship between the percentages of identified peers and the percentage of population, followed by Iceland, Norway and Sweden. As it was previously mentioned, for

example Dominica and Iceland have a low population weight in the total world population, 72000 (0.001%) and 296750 (0.004%) inhabitants respectively, so even a small number of involved peers appear highly relevant when normalized. Notice that Dominica has only 96 peers involved and Iceland 171 peers, for a total number of 52111 peers involved.

Following this idea, we can say that Norway and Sweden are countries with more meaningful values of such relationship between percentage of involved peers and percentage of the whole world population. Canada, Australia and United Kingdom also deserve to be mentioned, because these are countries with more involved peers and also with important population proportions, showing us that these countries make a quite intensive use of file sharing applications.

Finally, as it was already mentioned, United States loose its position in the overall ranking in spite of their high quantity of peers, because that number of peers is not so relevant when compared to the percentage of population this country has in the whole world population.

➤ **Music category**

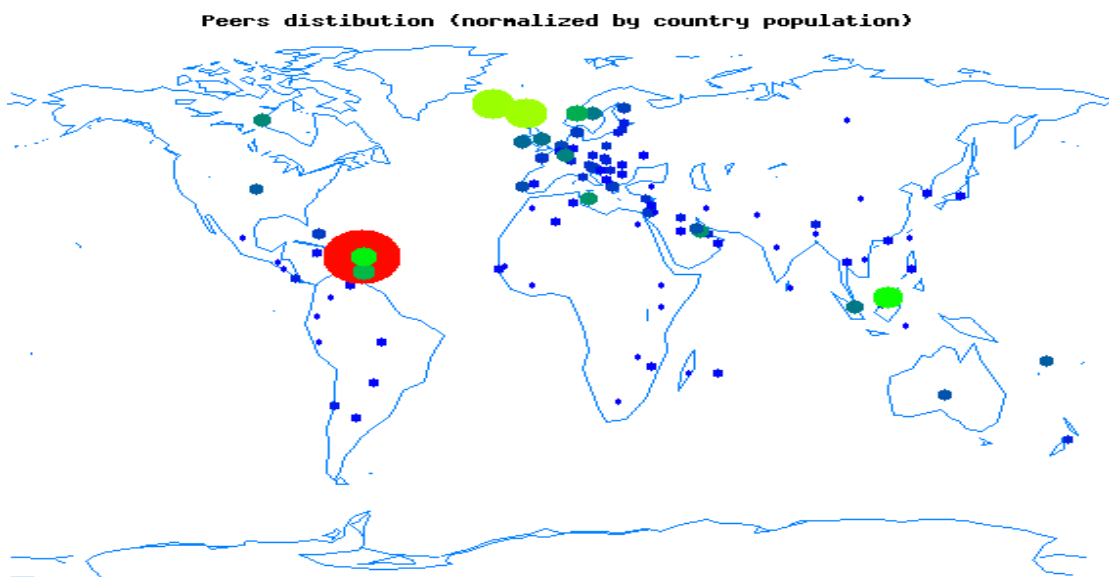


Figure 4.14 – Peers' distribution normalized by country population – *Music* category.

As can be observed on the map above, in the case of *Music* files it is also possible to conclude that top countries are those with very low population number values, such as Saint

Kits and Nevis, Faroe Islands, Iceland, Brunei Darussalam, Antigua and Bermuda. Then, we can see countries having a better balance between the total number of involved peers and percentage of population on the world, like Canada, Sweden and United Kingdom. Even in this case, United States of America and Australia occupy the 17th and 19th positions respectively.

➤ *Animated movies category*

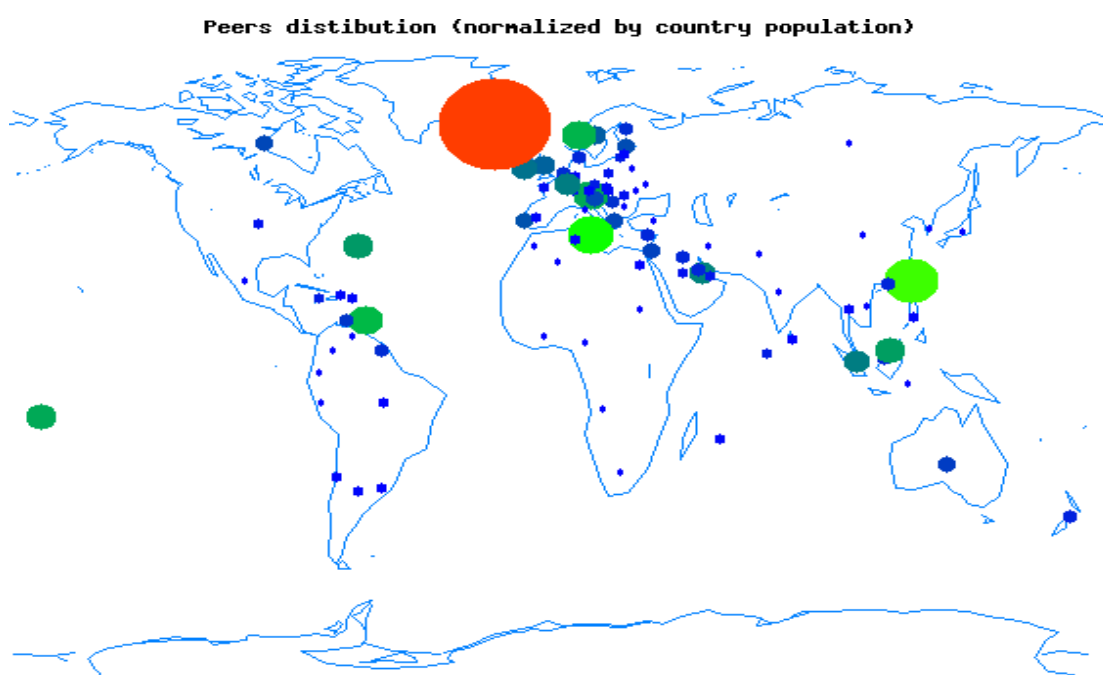


Figure 4.15 – Peers' distribution normalized by country population – *Animated movies* category.

Regarding the *Animated movies* category, Taiwan keeps a relevant place, occupying the second position, just after Iceland. United Kingdom also has a good position in the number of downloaded files from this category, as well as Canada and Australia.

➤ *French movies category*

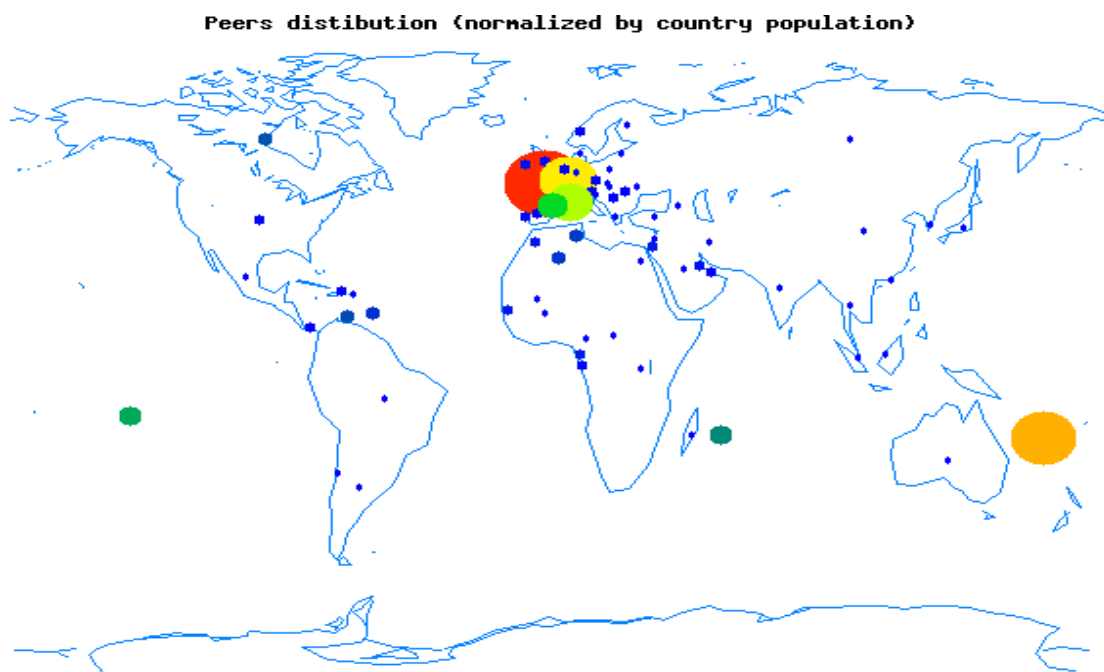


Figure 4.16 – Peers’ distribution normalized by country population – *French movies* category.

As it happened with the *Animated movies* category, in the case of *French movies*, France maintains its top position in the number of involved peers when compared with the quantity of population of the country. It is important to notice that results for this category are tightly related to the French language, which makes some francophone countries appear in the top of the list: New Caledonia, Luxembourg, Monaco, Andorra, French Polynesia, Switzerland, Mauritius, Belgium and Canada.

➤ *Indian movies category*

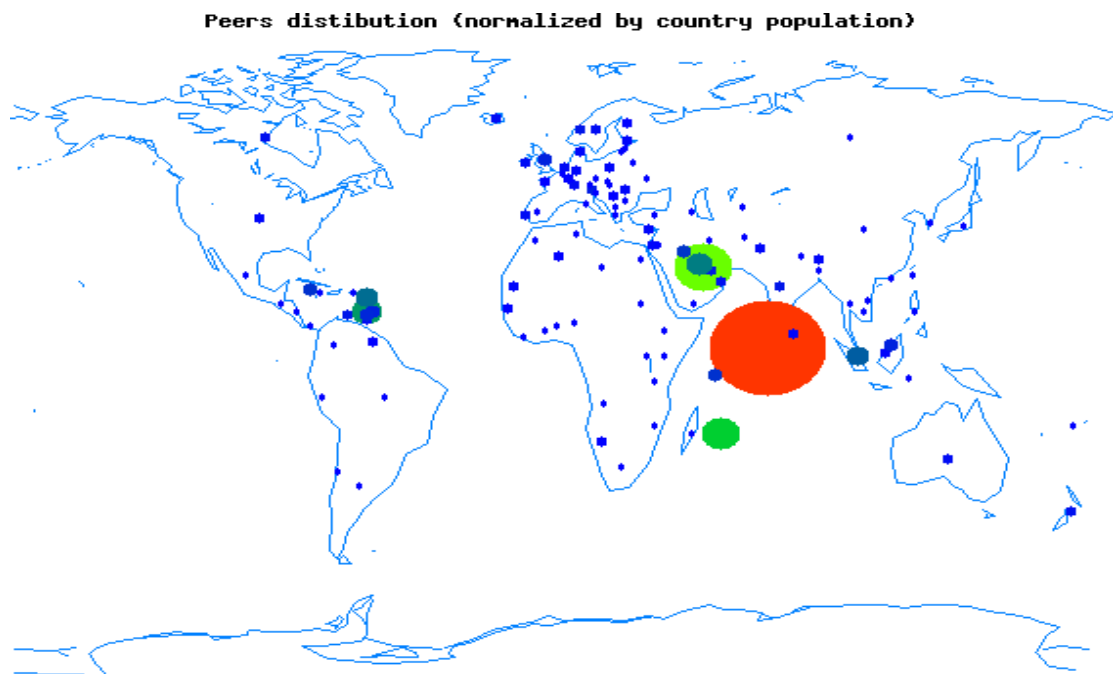


Figure 4.17 – Peers’ distribution normalized by country population – *Indian movies* category.

For *Indian movies*, India also lost its top position, as it was expected, since it is one of the densest countries in the world, contributing with more than 18% of the global population. It is also possible to observe on the normalized map shown above the strong concentration of peers in Asia.

➤ *Linux distribution category*

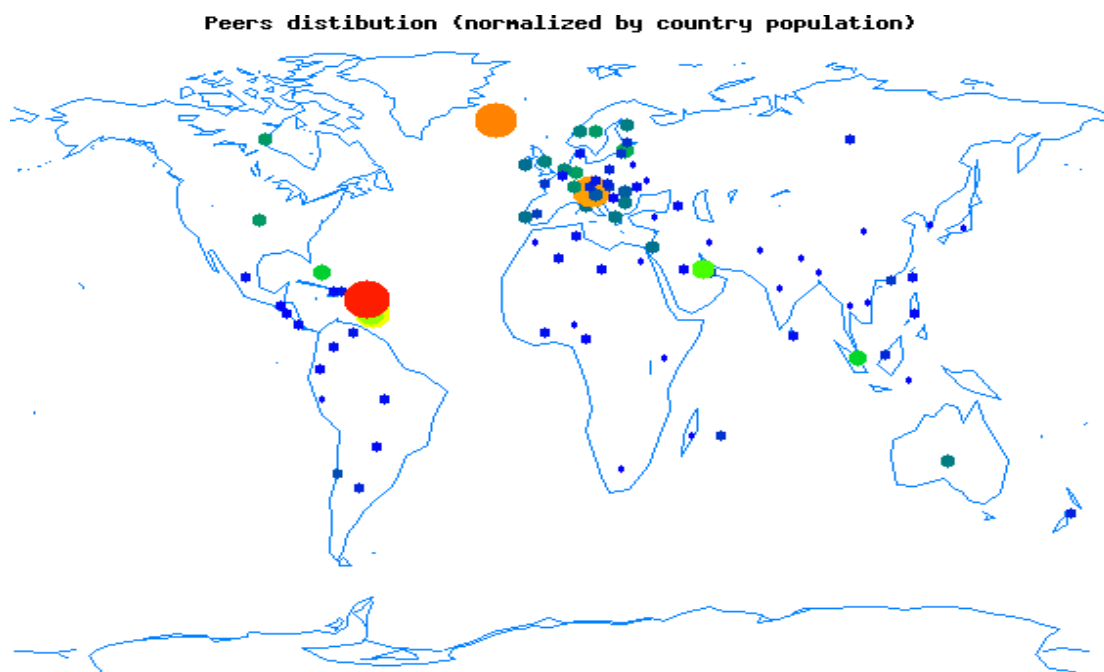


Figure 4.18 – Peers' distribution normalized by country population – *Linux distribution* category.

For the last category, the *Linux distribution*, it was verified that most relevant countries are Germany, United States of America, Sweden and Canada, maintaining results similar to those obtained before normalization, which can be explained by the small number of involved peers.

From this analysis, it is possible to conclude that the localization of most relevant peers is strongly correlated with the kind of file that is being shared: for example, for the *French movies* category countries with more interest are those that have some relationship with the French language, whereas for *Indian movies* peers are more concentrated in the Asian continent, mainly in India.

Besides, the importance of Canada, United Kingdom and Australia in this P2P file sharing system can be clearly seen, since these are developed countries with good networking infrastructures, easy Internet access and a good life quality standard.

Furthermore, the normalized number of peers shows that for the *2008 movies*, *Animated movies* and *Music* categories, Nordic and Central Europe countries have an higher predominance of peers, which reveals a cultural aptness of these populations to consume multimedia contents at home, maybe due to climate constrains as suggested in [24].

It is also important to remind that, except for *2008 movies* and *Indian movies*, categories have a small sample of involved peers, between 2000 and 5000, and in the most part of the cases they are strongly concentrated in one country. This fact imposes an additional difficulty on our study. As it was already concluded, for these categories results obtained before and after normalization were quite similar, in opposition to other categories, that have much more peers involved.

4.2 Peers' availability

After localizing different countries, the next step is to analyse the peers' availability, which is defined as the percentage of peers that answer to a resource sharing request (the answer is evaluated from obtained RTT values). When RTT is different from -1 , it means that the peer is sharing a resource.

In this section the evolution of the number of available peers during the period of analysis will be evaluated. Firstly, the analysis will be made for each category and then for each country, in order to obtain a better understanding of what are daily periods involving more and less peers.

Next plots show obtained results for the peers' availability during different periods of the day and for different days.

➤ *2008 movies category*

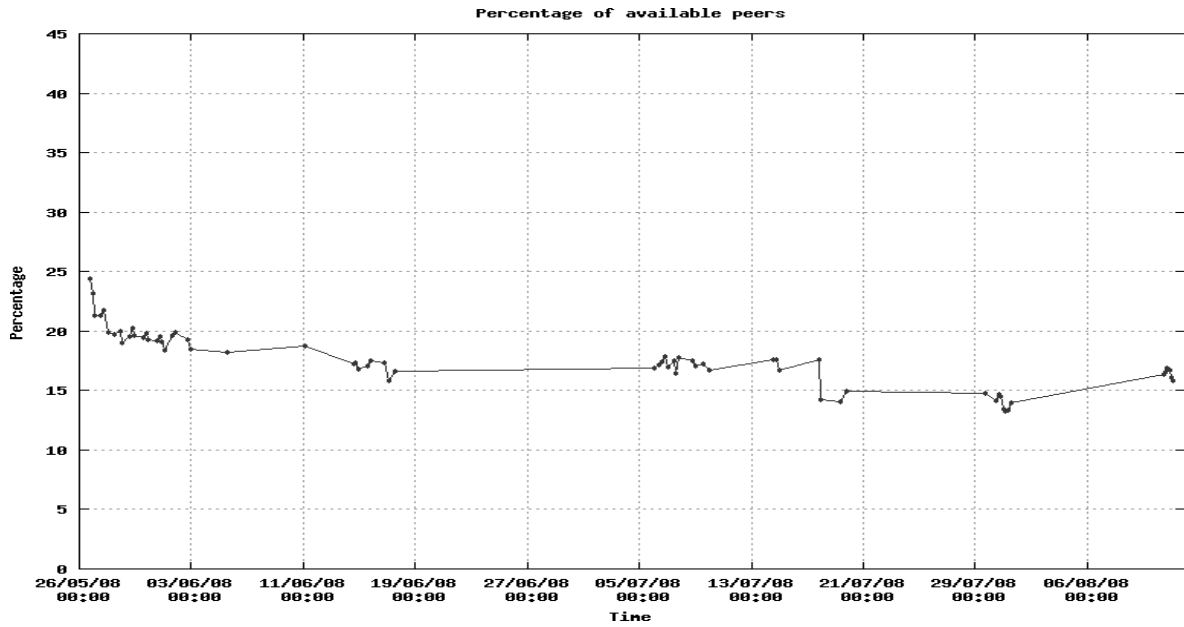


Figure 4.19 – Evaluation of available peers – *2008 movies category*.

➤ *Music category*

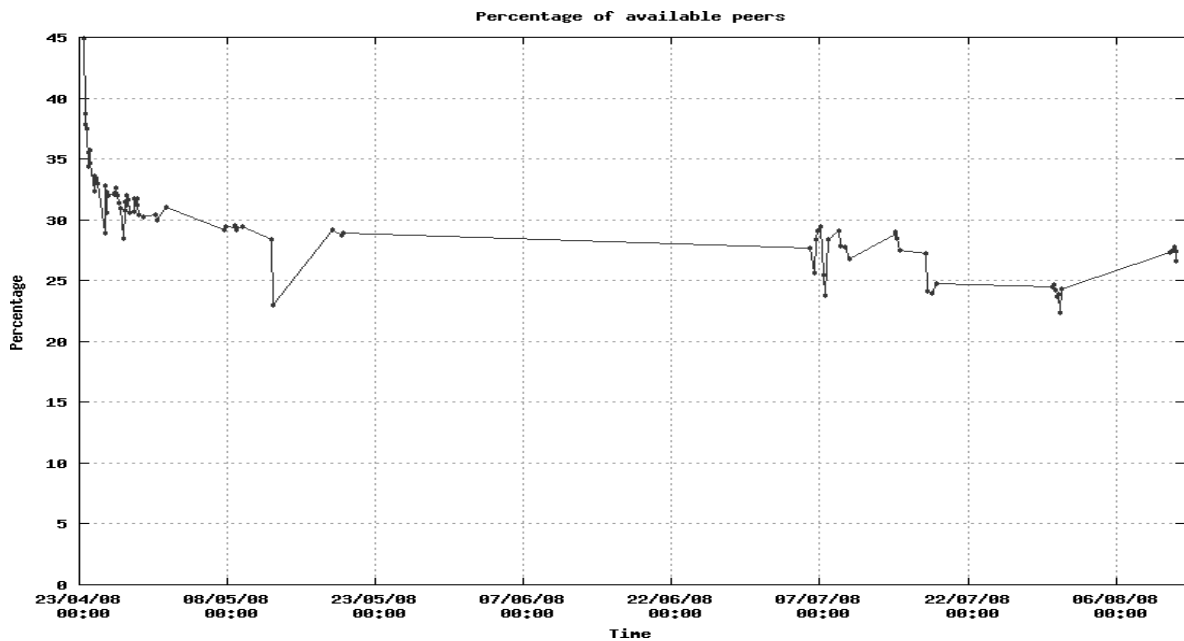


Figure 4.20 – Evaluation of available peers – *Music category*

➤ *Animated movies category*

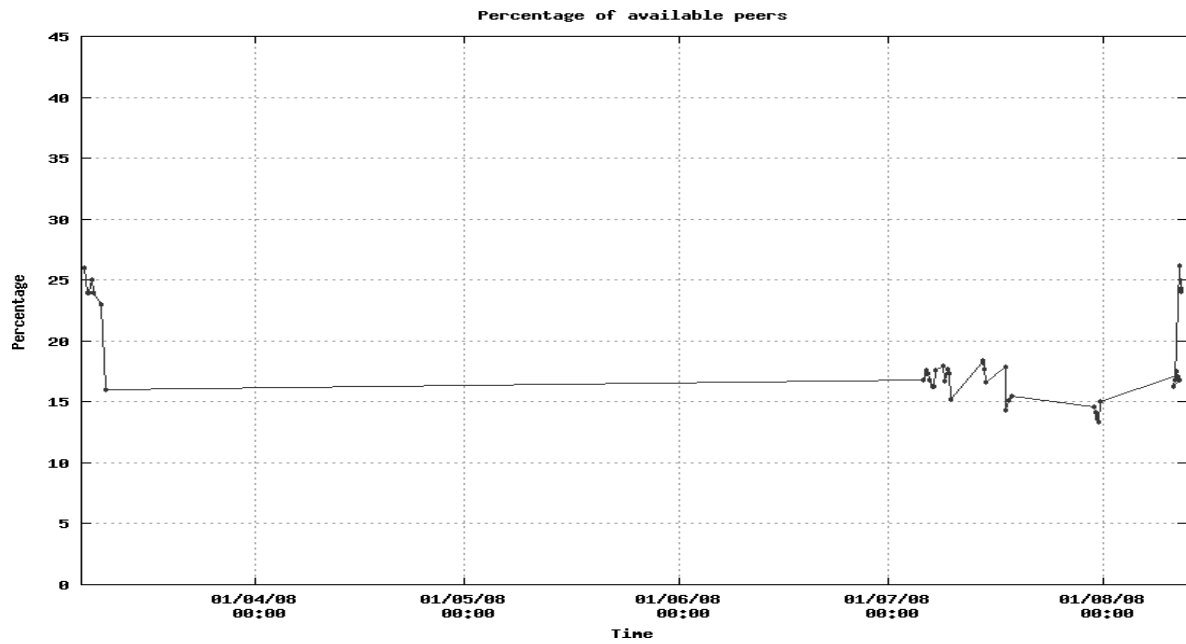


Figure 4.21 – Evaluation of available peers – *Animated movies category*.

➤ *French movies category*

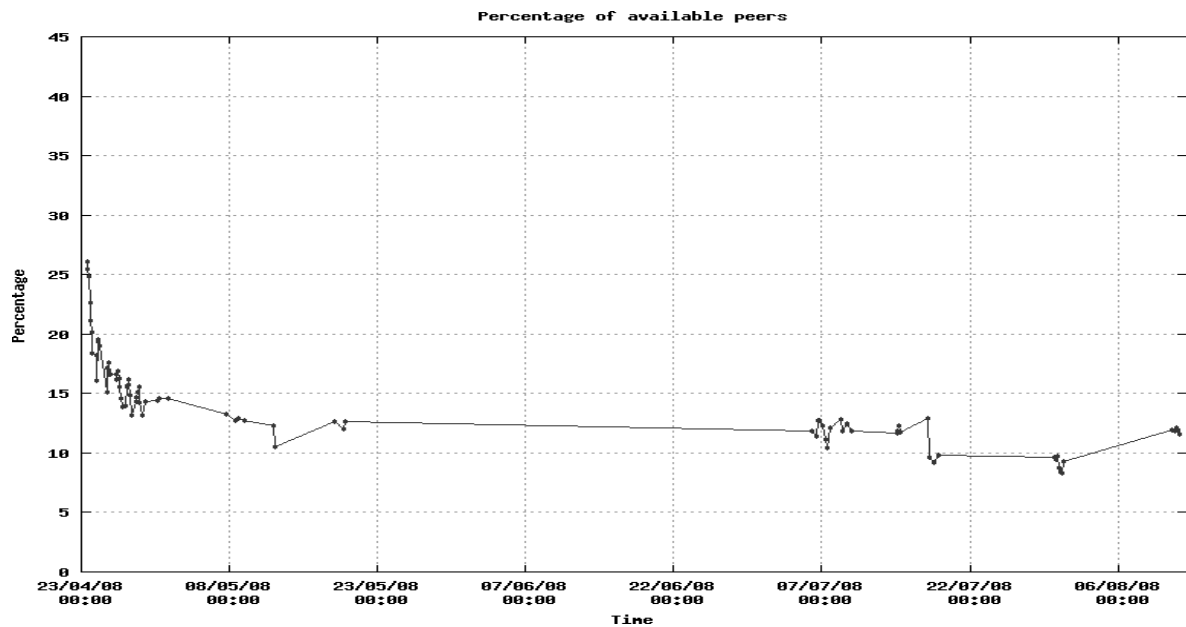


Figure 4.22 – Evaluation of available peers – *French movies category*.

➤ *Indian movies category*

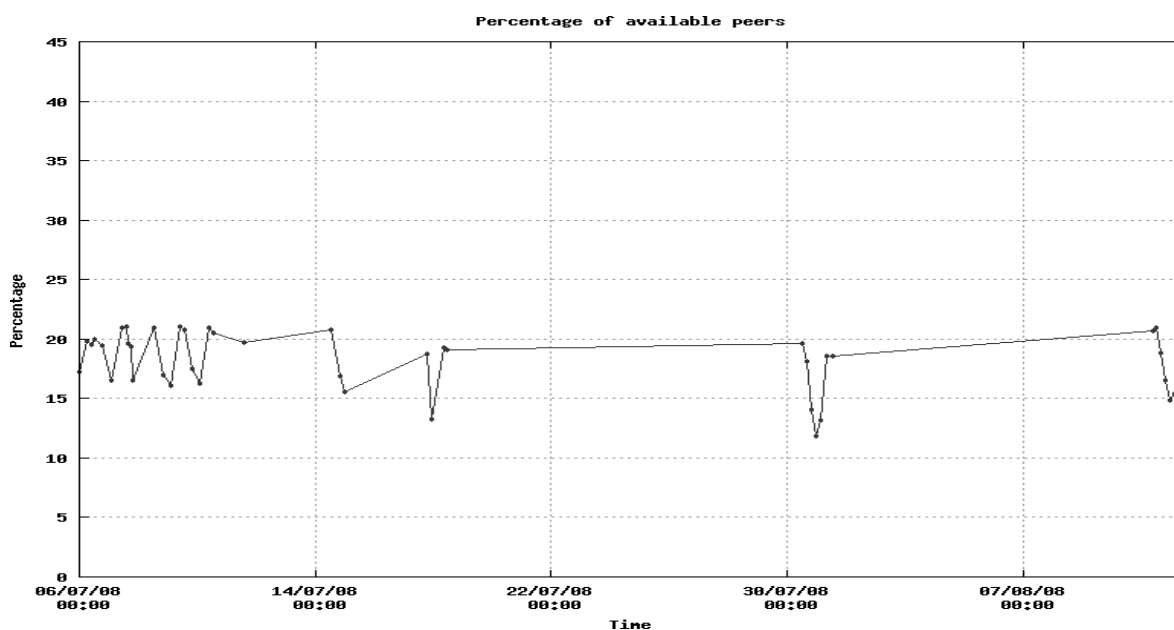


Figure 4.23 – Evaluation of available peers – *Indian movies category*.

➤ *Linux distribution category*

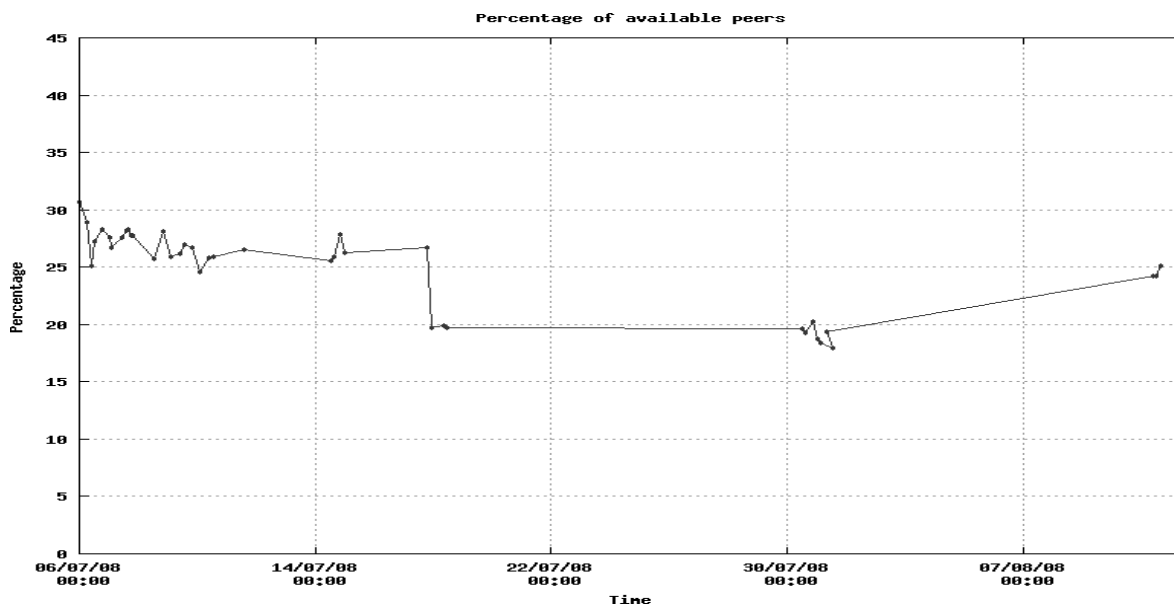


Figure 4.24 – Evaluation of available peers – *Linux distribution category*.

As it can be observed in previous graphs, the percentage of available peers in almost all categories changed very slowly during the period of study. This behaviour shows that in general, people keep using this way of sharing files; it is not just a specific file that they are searching for. If this was not true, the detected TCP Port number would not be available any more after finishing the file download. Therefore, we can conclude that P2P file sharing is becoming more and more common in people lives and it is not a sporadic attitude, which is comprehensible due to its simple utilization, fast speed on file downloading and good variety of available contents. Furthermore new files are firstly available on these systems than in commercial circuits in every country. Another attraction of these systems is that people don't have to pay/buy to have access to the desired files contents: they do not need to pay to watch them on the cinema or to buy a CD/DVD containing such subjects.

From figures above, it is also possible to observe that the peers' availability percentage never reaches very high levels, starting in almost all cases with 25% of the whole peers sample except for the *Music* category that starts with 45% of the total number of peers. This can be justified by the fact that an increasing number of hosts are firewalled or located behind NAT boxes and proxies in order to avoid detection and consequently, do not answer to any probing test.

Such conclusions can be taken into consideration since probing tests were done just after the knowledge of peers' ports: obviously at least at that moment a great portion of those peers should be available and if not, that behaviour can be explained by the fact that those ports are protected against intrusions.

Observing the peers' availability during the daily period, we can see that there are not huge changes, which can be explained by the fact that our analysis involved peers from all over the world, having different time zones according to their localization. Nonetheless, for the *Indian movies* category differences on the number of available peers for different hours are visible. The daily period with lowest values is, as it can be observed, between 8.30 pm and 6.30 am, corresponding to the GMT between 4.00 pm and 2.00 am. Thus, it is possible to conclude that the period of time with lower number of users coincide with the break day, which can be justified by:

- people that does not have unlimited traffic nor have any period of the day with free Internet access (or if they eventually have free Internet access, it should be during the working period) are encouraged to turn off their connections and computers when they are not needing it, for example when they are sleeping;
- a great fraction of companies are closed during the night, computers are turned off and consequently, Internet is disconnected;
- Internet access in India is not cheap, being not accessible at home for a major slice of the population.

With the aim of evaluating the relationship between the peers' availability and the time of the day, an analysis of such availability was done by distributing peers by their countries, then allowing us to pay attention to the time zone of each case. Countries chosen were those with more peers involved. Such plots can be seen in figures below, where the general decline on the number of peers available during the break day can be confirmed.

Figures 4.25, 4.26 and 4.27 correspond to the peers' availability for the *2008 movies* category: since this category has more peers involved, it will allow us to make a more fair evaluation of this parameter. In the same way, Figure 4.28 corresponds to peers involved in the *French movies* category and the last one corresponds to peers discovered on the *Indian movies* category. Figures 4.25, 4.26 and 4.27 include data corresponding to more than one country, due to their closely related time zones. Figure 4.25 presents results for Israel and eight European countries: Poland, Sweden, Finland, Norway, Netherlands, Portugal, Romania and Greece. Figure 4.26 corresponds to United States, Canada and Brazil and Figure 4.27 corresponds to Philippines, Australia, Malaysia and Singapore. Note that the time shown at all plots is the GMT.

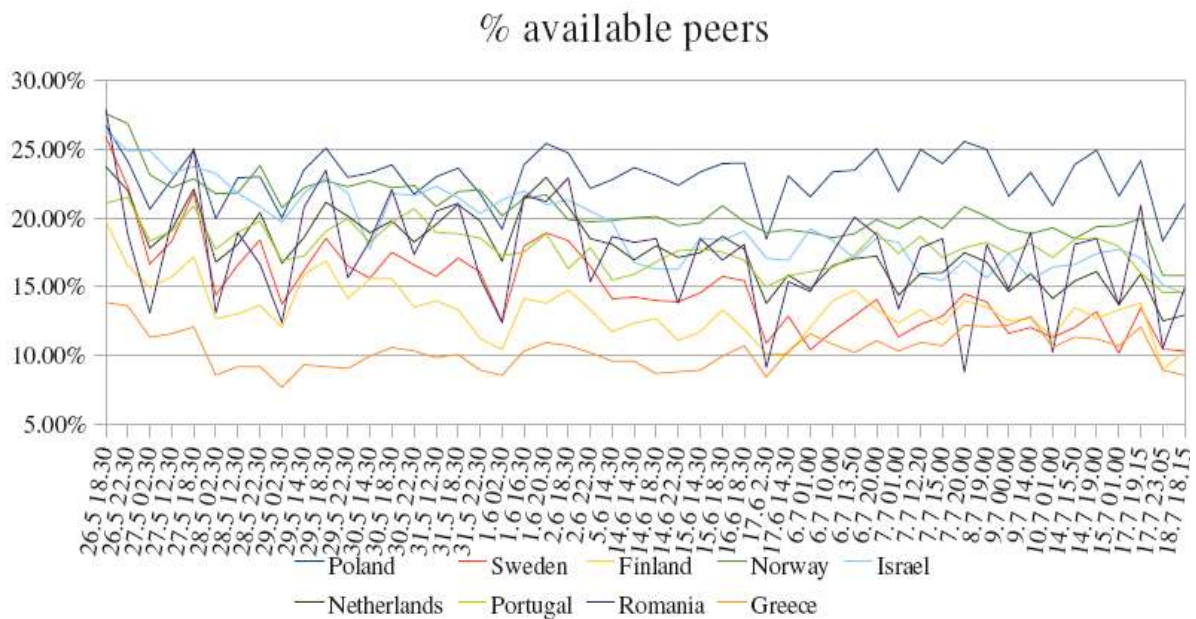


Figure 4.25 – Evaluation of available peers – Israel and eight European countries.

When observing the plot above it is obvious that, in the major part of the time, the peers' availability decrease on breaking day periods and lowest values are observed between 10.30 pm and 2.30 am GMT (corresponding to one or two hours more for other countries). It can also be observed that highest values of the availability percentage occur mainly at the afternoon, more or less at 6.30 pm. These high values that occur at this period of the day can be justified since this is the end of the daily working period for the majority of the population, when people arrive at home and they can finally execute their desired downloads. On the other hand, lowest values observed generally at breaking day were expected, since this is the time of the day when people are usually sleeping or not working.

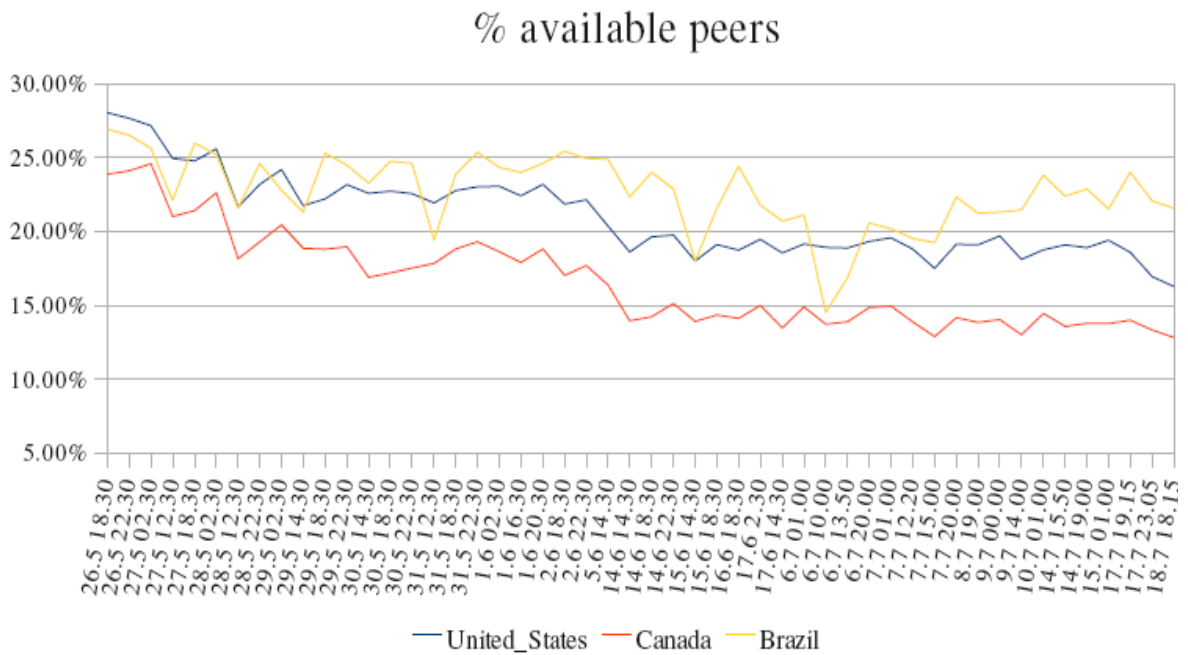


Figure 4.26 – Evaluation of available peers – United States, Canada and Brazil.

The plot above presents the evaluation of the peers' availability for three most important American countries, concerning the number of peers: United States, Canada and Brazil. As it was already said, the presented time zone is GMT, which for those countries - depending on the state/province - corresponds to a local time of four up to nine hours less than GMT. Having these time zones into consideration, it is visible that minor values appear on the early afternoon, between 00.30 pm and 2.30 pm GMT, which corresponds to the breaking day on these countries. On the other hand, highest values were verified at the GMT night period, corresponding to afternoon on these American countries.

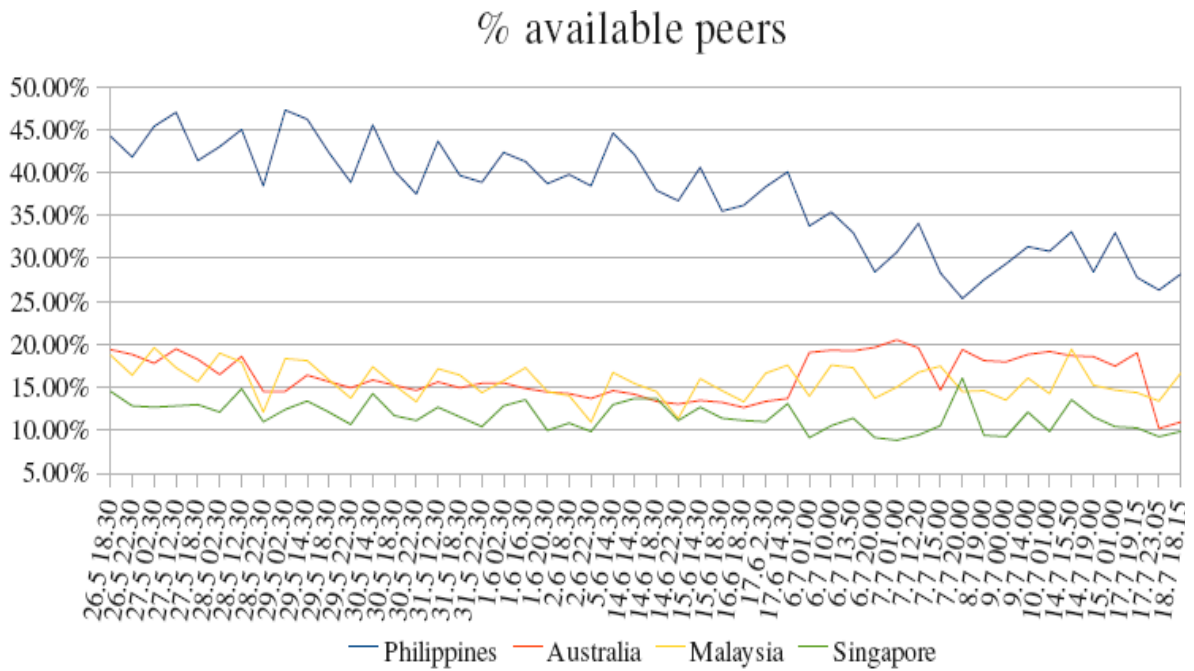


Figure 4.27 – Evaluation of available peers – Philippines, Australia, Malaysia and Singapore.

On Figure 4.27, Australia and three Asian countries, Philippines, Malaysia and Singapore are analyzed. On these countries, GMT corresponds to seven hours less than local time. Having this in mind and observing the plot, it is possible to check that highest values of available peers occur at GMT lunch time, which once again corresponds to afternoon on these countries. The lowest values appear mostly at GMT night period, corresponding to the breaking day period on these countries, as it was expected.

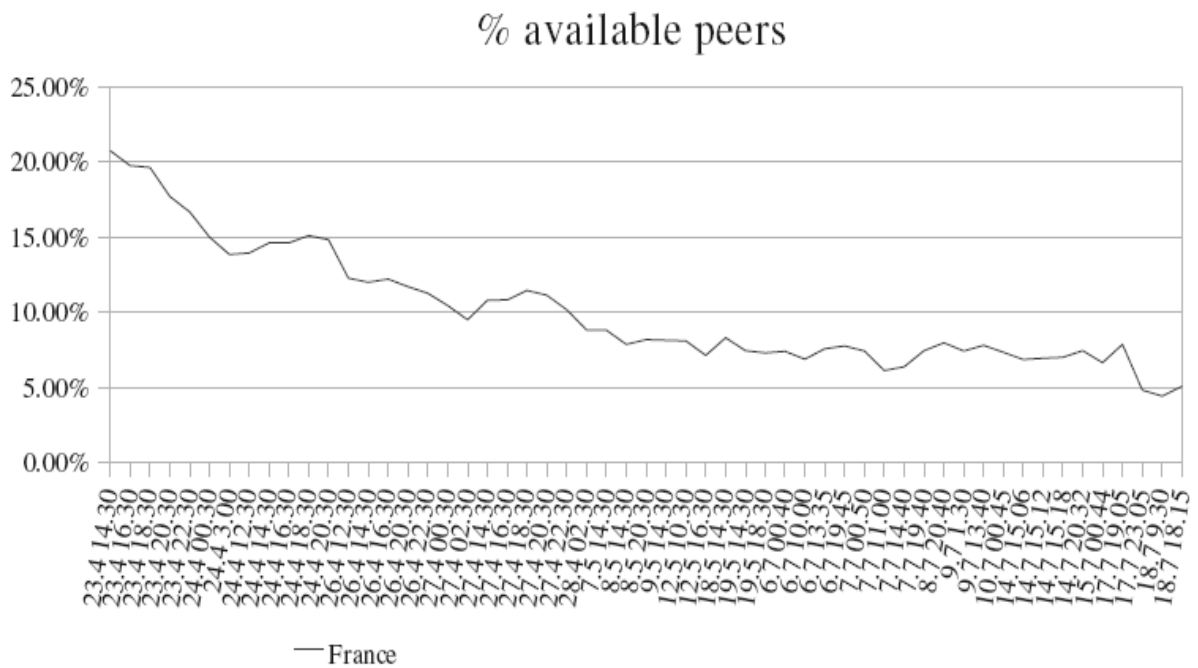


Figure 4.28 – Evaluation of available peers – France.

On the figure above, availability of French peers is shown for the *French movies* category. As known, France has one hour more when compared to GMT zone. In opposition to other analyzed countries, this one presents a high and progressive decrease on the percentage of peers available during the studied period, turning the analysis of their availability during the day period even harder. Even so and taking the risk of not being so rigorous as before, it is also possible to verify a decrease on the number of available peers during the night period and an increase during the afternoon.

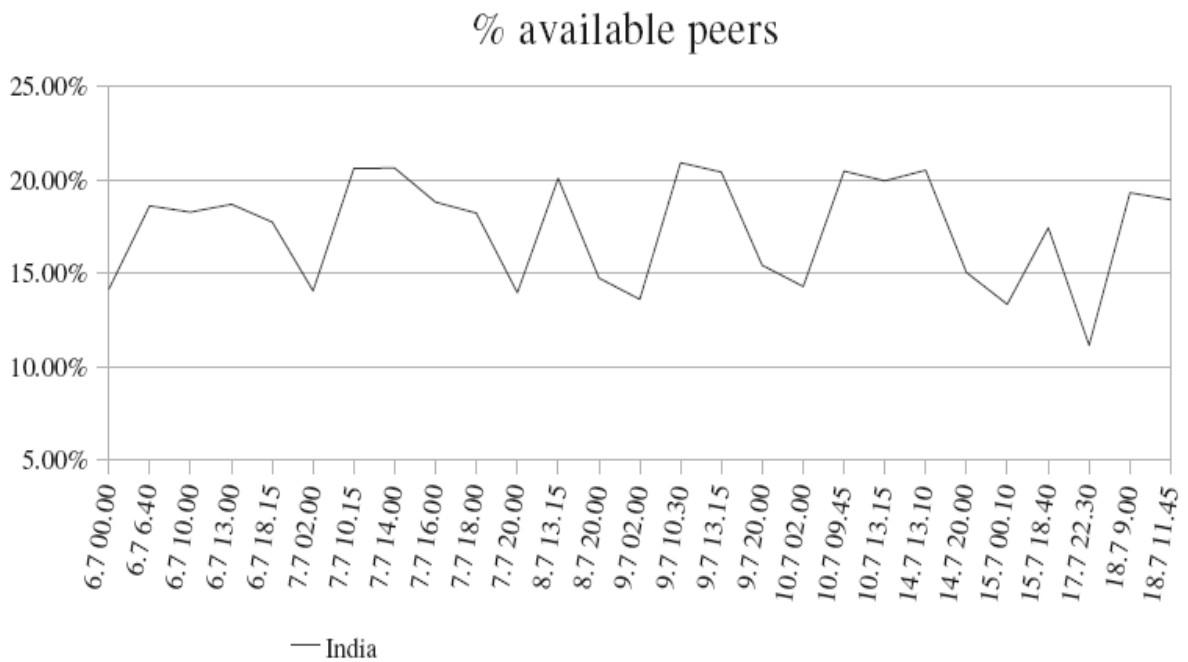


Figure 4.29 – Evaluation of available peers – India.

The last country evaluated was India, which distribution can be seen on Figure 4.29. This country was separated from other Asian countries due to the fact that information about its peers' availability has been extracted from another category of files. Besides, this country has a different time zone when compared to others (four hours and half more than GMT zone at that season).

Analysing this plot it can be seen that this country is the one that presents more differences on the peers' availability during the day. In fact, it is easy to observe the peers' availability break down during the GMT afternoon and night periods, corresponding once again to the night period in India. On the other hand, the peers' availability starts to increase from the GMT early morning on, corresponding to their lunch time and maintains high values during afternoon until dinner time.

Similar conclusions were obtained when evaluating available peers for the *Indian movies* category as expected, since India represents 67.5% of the whole peers involved on this category.

After analysing peers' availability changes during the day, obviously taking time zones of involved countries into account, it is possible to conclude that there is a clear break down on peers' availability values during night periods. As it was already mentioned, these changes were not surprising if we keep in mind that this is the time of the day when generally people are sleeping and for situations where users don't have unlimited Internet traffic, users are encouraged to shut down their computers because they are not using them. Besides, on this period the majority of companies are closed, with their computers turned off. Therefore, this is the time period that has the largest number of inactive people.

In order to perform a more deep analysis on the peers' availability, several consecutive availability tests were done, corresponding to small periods of time. The analysis was done for the *Music*, *Animated movies*, *French movies*, *Indian movies* and *Linux distribution* categories. The *2008 movies* category was excluded from this part of the study because it involves a huge number of peers, demanding a long period of time in order to obtain RTT results from all peers (approximately one hour). In this way, figures shown below are area plots presenting the percentage of available peers for each category analysed in short periods of time.

Before taking any conclusions regarding on obtained plots, it is important to refer that *Indian movies* and first two periods of *Linux distribution* analysis were made when the download of these files were still in progress and then, not involving the total number of peers that were localized at the end of the files' download for each category. An advantage of this analysis during the download period is that some of these evaluated peers were found on the system on that moment, so it is more probable that such peers remain on the system on the following minutes.

In opposition to these two categories, others were intentionally analysed just once, having in consideration that this would be enough to take intended conclusions. The main reason for this was that the download of files from these categories were done too much time before the analysis. Therefore, involved peers were also contacted a long time ago and they are not so useful for the current research.

➤ **Music category**

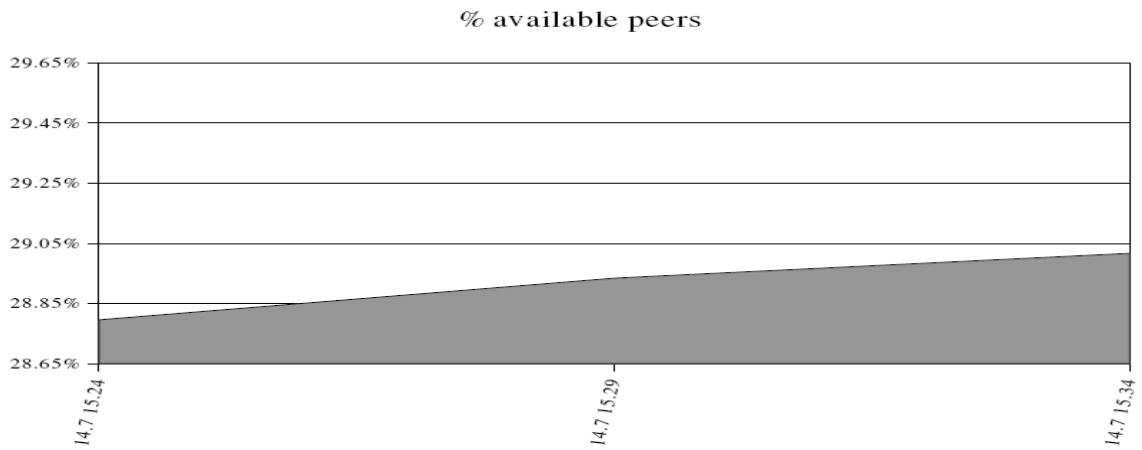


Figure 4.30 – Evaluation of peer availability in a short period of time - *Music* category.

➤ **Animated movies category**

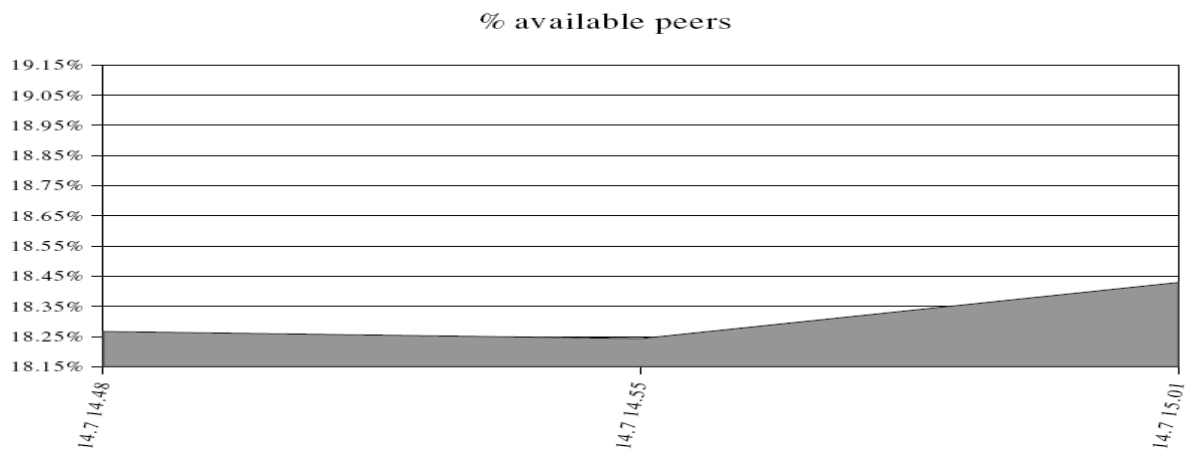


Figure 4.31 – Evaluation of peer availability in a short period of time - *Animated movies* category.

➤ ***French movies category***

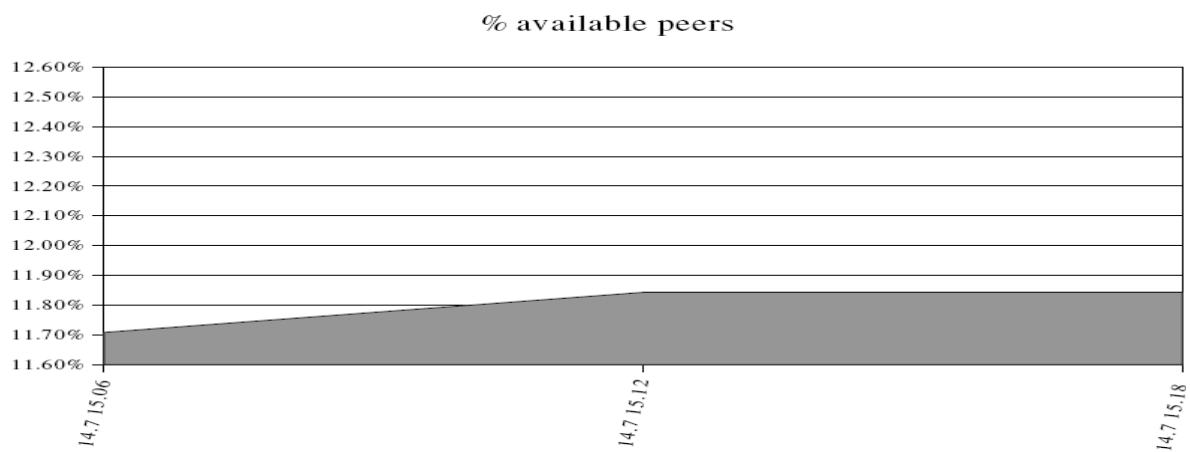


Figure 4.32 – Evaluation of peer availability in a short period of time - *French movies* category.

➤ ***Indian movies category***

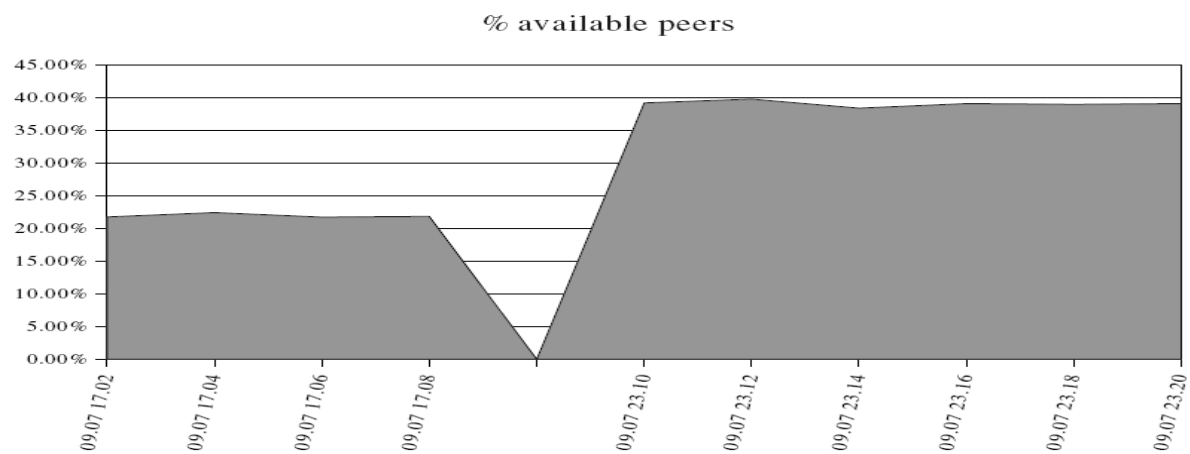


Figure 4.33 – Evaluation of peer availability in short periods of time - *Indian movies* category.

➤ **Linux distribution category**

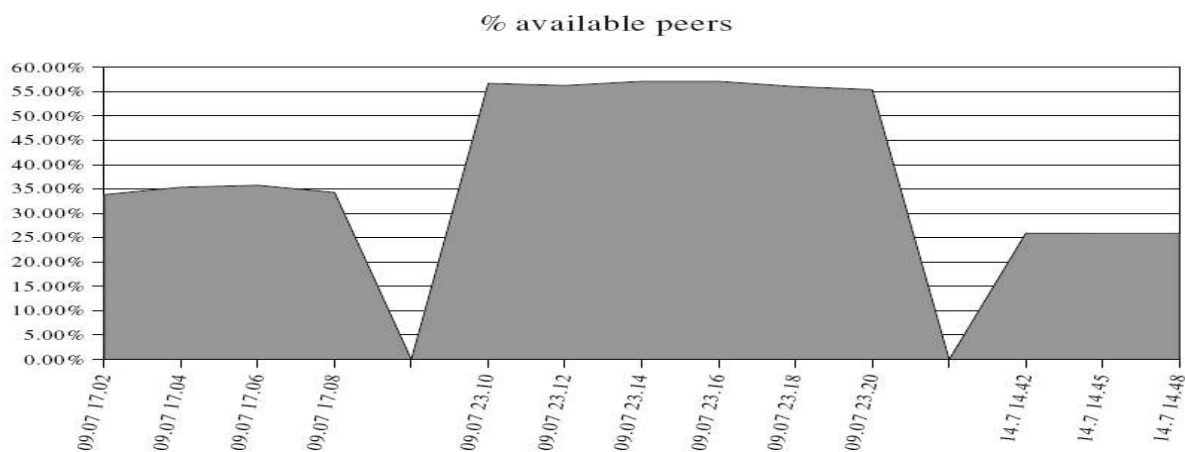


Figure 4.34 – Evaluation of peer availability in short periods of time - *Linux distribution* category.

After observing figures above, it can be seen that for short periods of time there are no significant changes on any of analysed categories.

Meanwhile, it is important to see that peers' availability changes for *Music*, *Animated movies* and *French movies* categories were too much smaller when compared to changes that occurred for *Indian movies* and *Linux distribution* categories. The first three categories did never achieve a variation of 0.5%, being the maximum variation equal to 0.22%. On the other hand, although with small variations, *Indian movies* reached a variation of 1.39% and *Linux distribution* a variation equal to 1.94%.

These differences on obtained results can be justified by the fact that peers evaluated for *Indian movies* and *Linux distribution* were found on the moment that this analysis was done, which did not occur for other three categories, where peers were localized more than one month before this evaluation. Therefore, it is more probable that peers found at that moment have more activity, with more frequent and successive entrances and exits from the system.

4.3 Summary

This chapter started making a geographical localization of all peers that were discovered as being part of this P2P file sharing system. After the peers' localization by country and by continent and in order to get more detailed information about the distribution of peers around the world, another analysis was made regarding the distribution of the number of peers normalized by the corresponding country population. From this analysis, a study on the peers' availability was also conducted and based on it, several analysis were made and their corresponding results were discussed: the variability on the peers' availability during the time of the day and during the period of study.

5. Round Trip Time

Round Trip Time can be defined as the time a packet takes to go from one host to another and return to the original host. This parameter has a strong dependence on the distance between hosts, traffic load, time of the day, Internet connectivity and several other factors.

This chapter will present a study on the variability of RTT in order to find out its dependences, especially from the distance between hosts, Internet connection quality and the time of the day.

For such an analysis, firstly we had to use the information about IP addresses and application TCP ports used by peers, obtained through log files generated by Vuze. Thereby, performing TCP probes of port numbers using the Nmap tool, RTT average values were registered and evaluated.

As it was already mentioned, RTT is the time a packet takes to go and comeback from a host to another, so variations are expected according to the distance that separates the origin from the end-host. A study of such relationship between RTT and host distances was also made in this chapter.

Distribution of RTT values by the continent each available peer belongs to is shown below through polar maps. These maps represent the RTT distribution for each one of the different studied file type categories: *2008 movies*, *Music*, *Animated movies*, *French movies*, *Indian movies* and *Linux distribution*. For each type of file, four polar maps are shown: a pair of maps for each Internet connection, CATV 12 Mbps and ADSL 4 Mbps, where the second map is simply a zoom view of the first one in order to get a better visualization of most significant RTT values and differences.

Note that, since the origin-host is located in Portugal, countries located in Europe should present lower values of RTT, immediately followed by countries from America, Africa, Asia and Oceania, in decreasing order of magnitude.

➤ *2008 movies category*

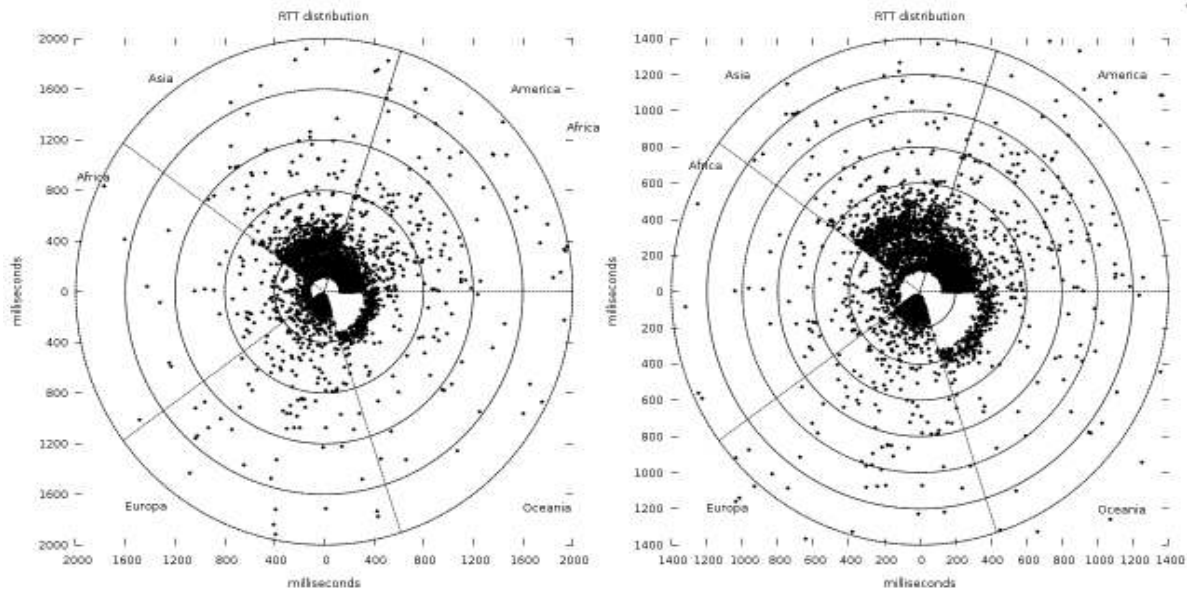


Figure 5.1 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – *2008 movies category*.

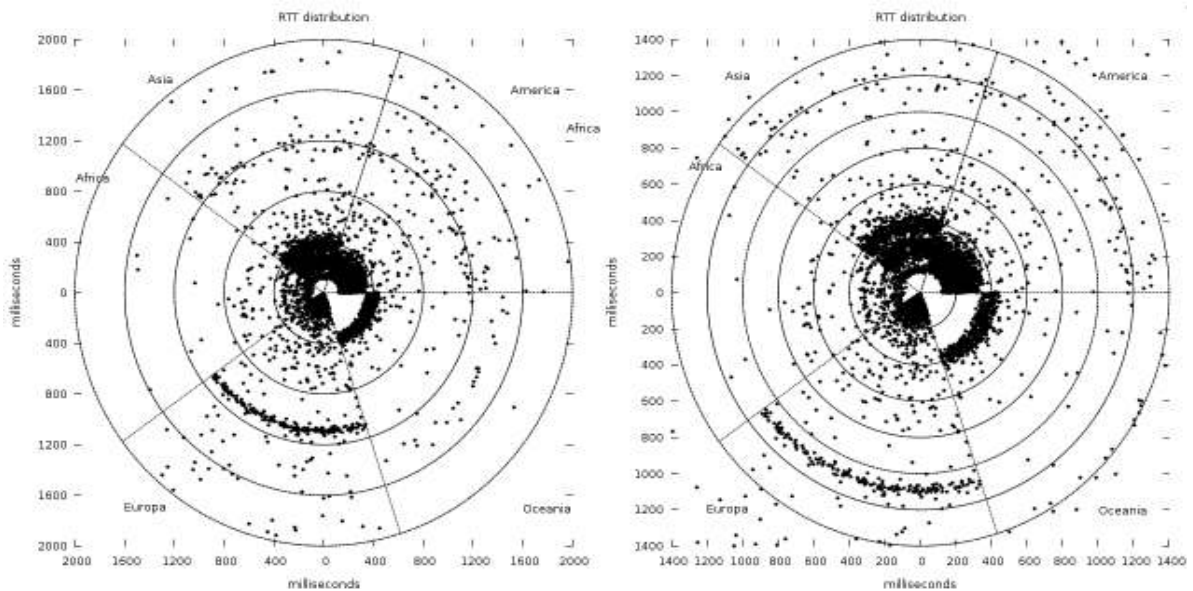


Figure 5.2 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – *2008 movies category*.

➤ *Music category*

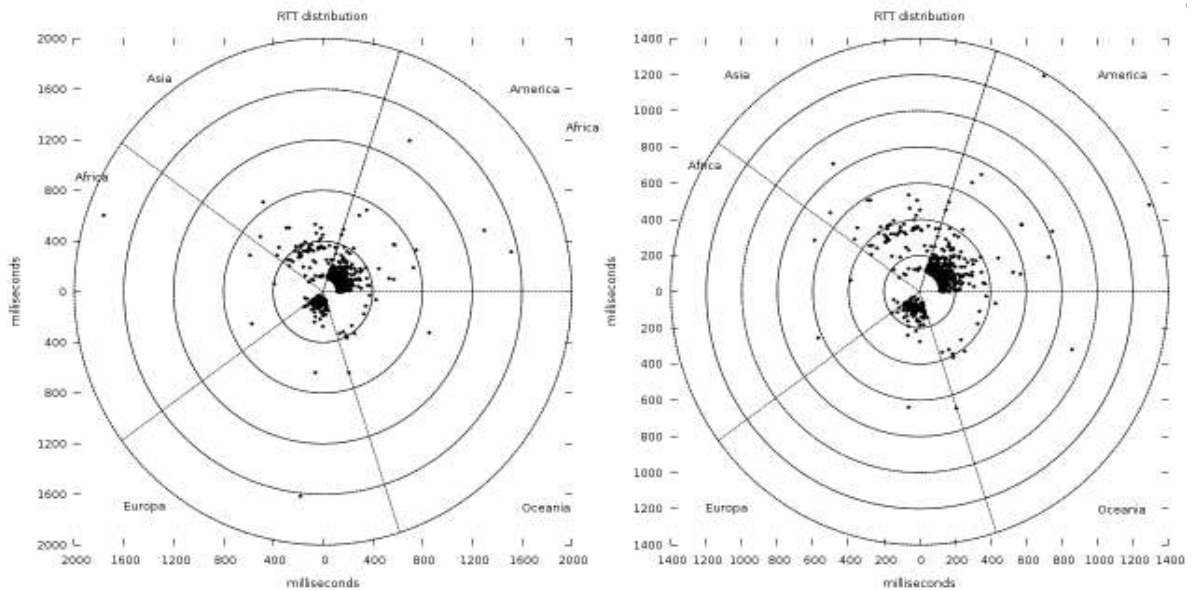


Figure 5.3 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – *Music category.*

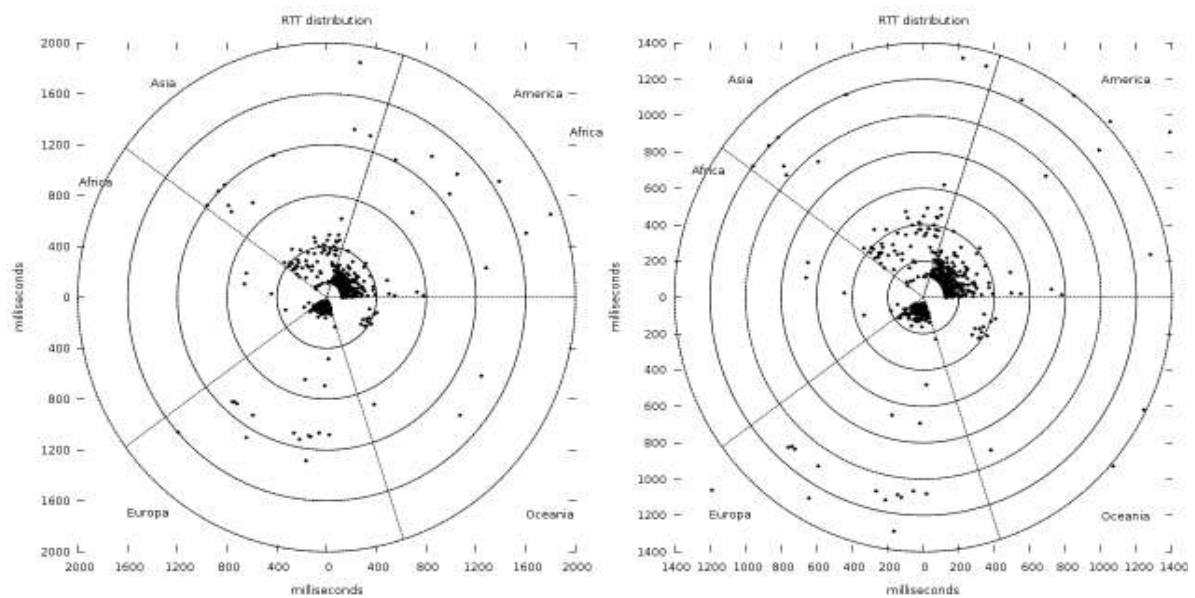


Figure 5.4 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – *Music category.*

➤ *Animated movies* category

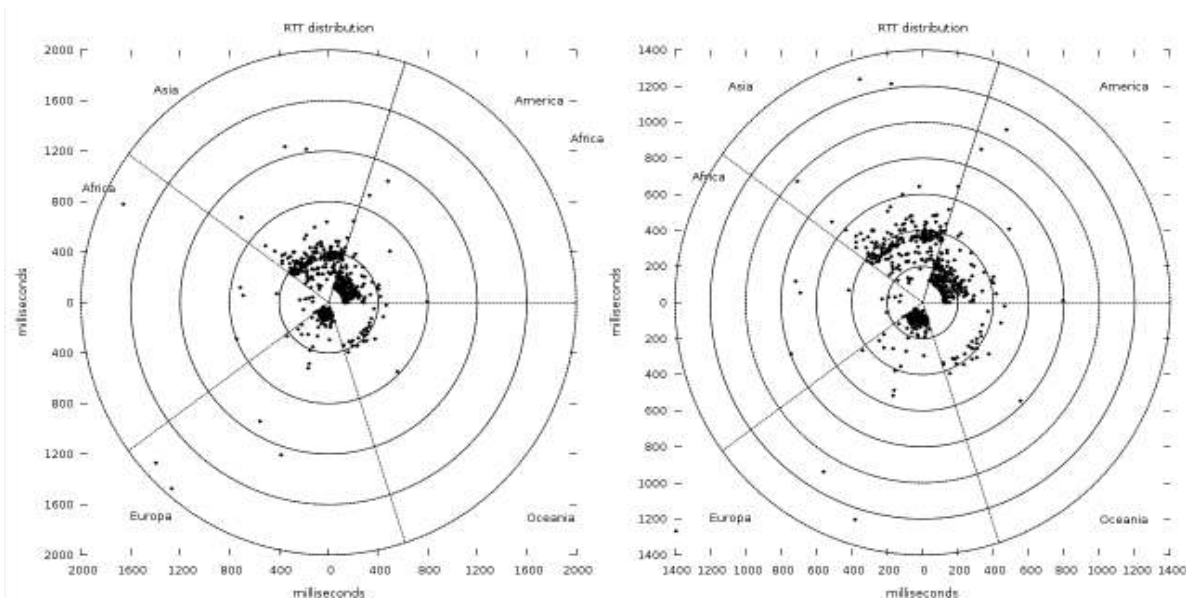


Figure 5.5 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – *Animated movies* category.

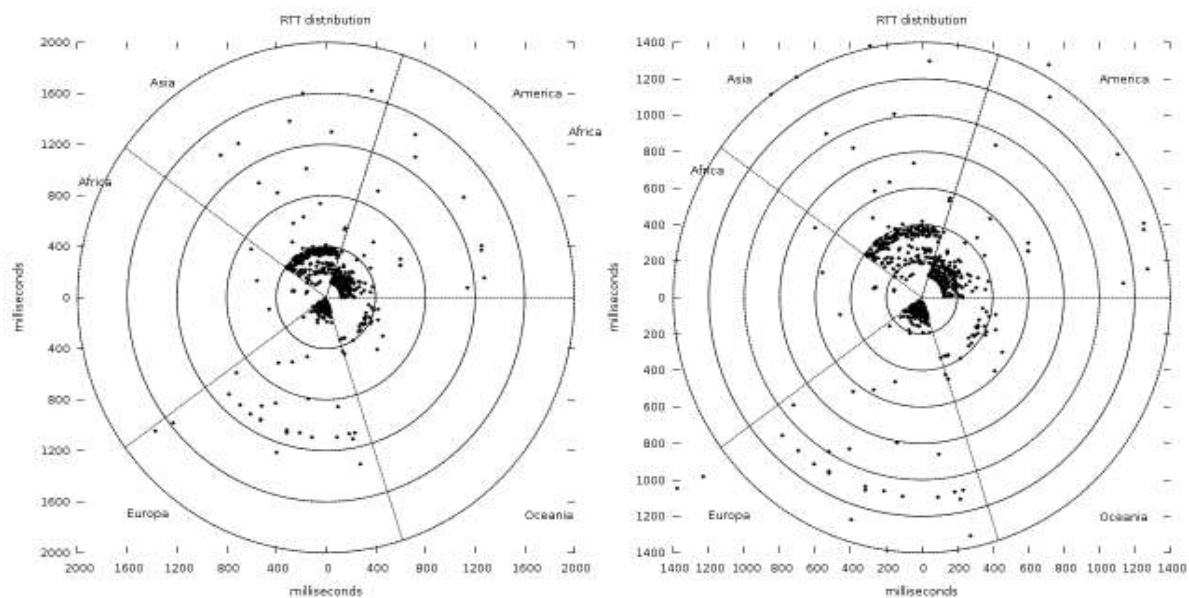


Figure 5.6 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – *Animated movies* category.

➤ *French movies category*

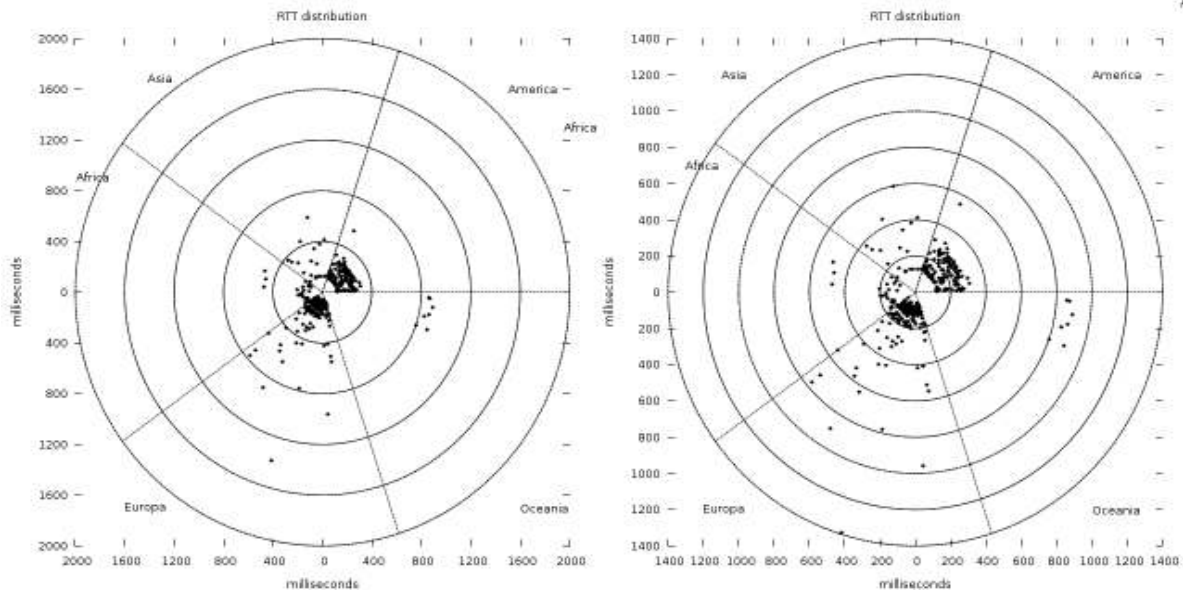


Figure 5.7 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – *French movies category*.

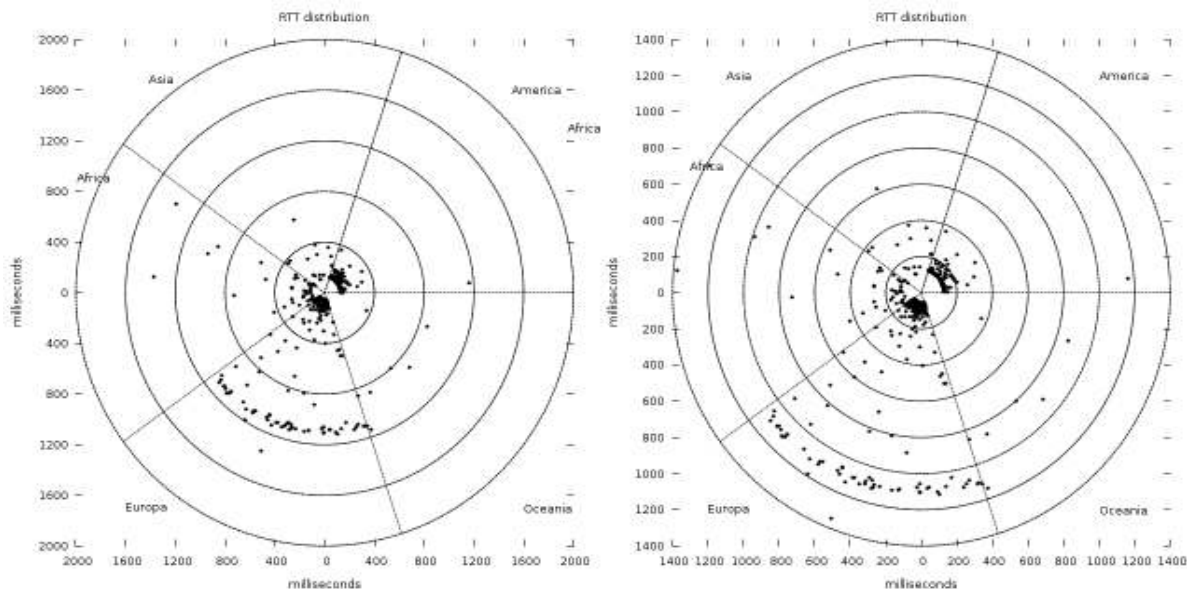


Figure 5.8 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – *French movies category*.

➤ **Indian movies category**

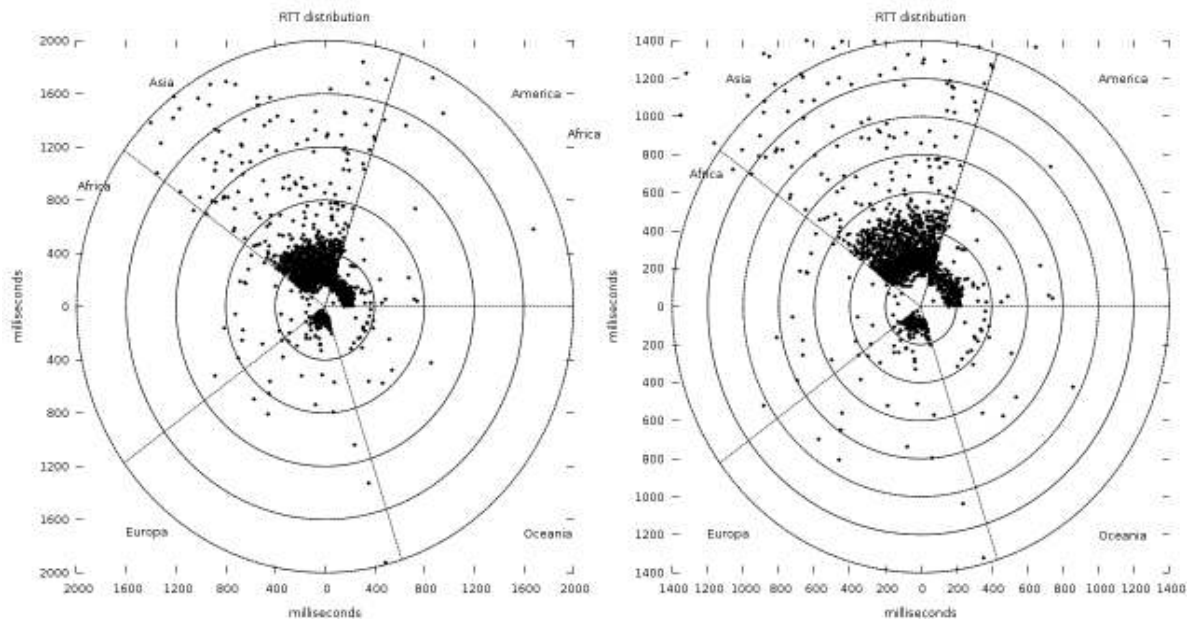


Figure 5.9 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – *Indian movies category*.

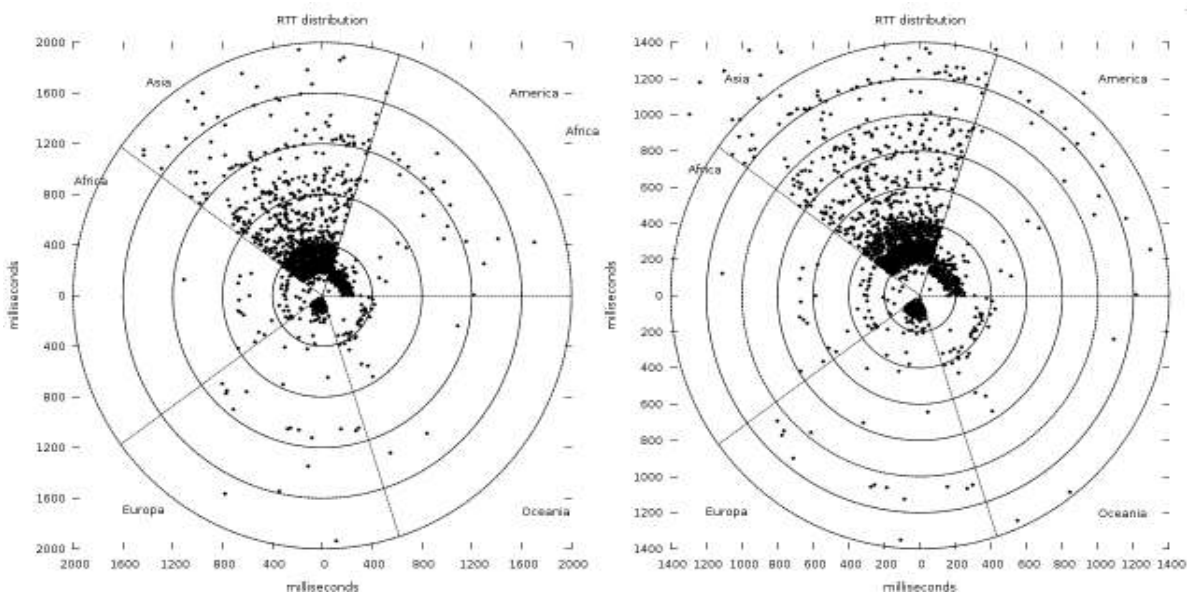


Figure 5.10 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – *Indian movies category*.

➤ *Linux distribution category*

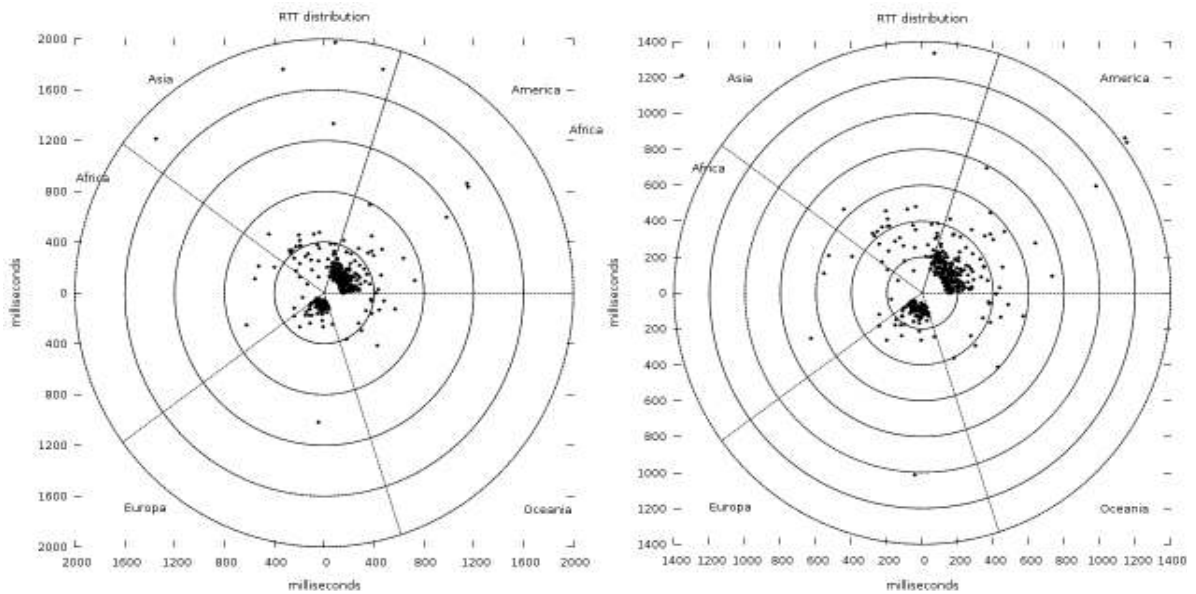


Figure 5.11 – Round Trip Time distribution of the CATV 12 Mbps Internet connection – *Linux distribution category.*

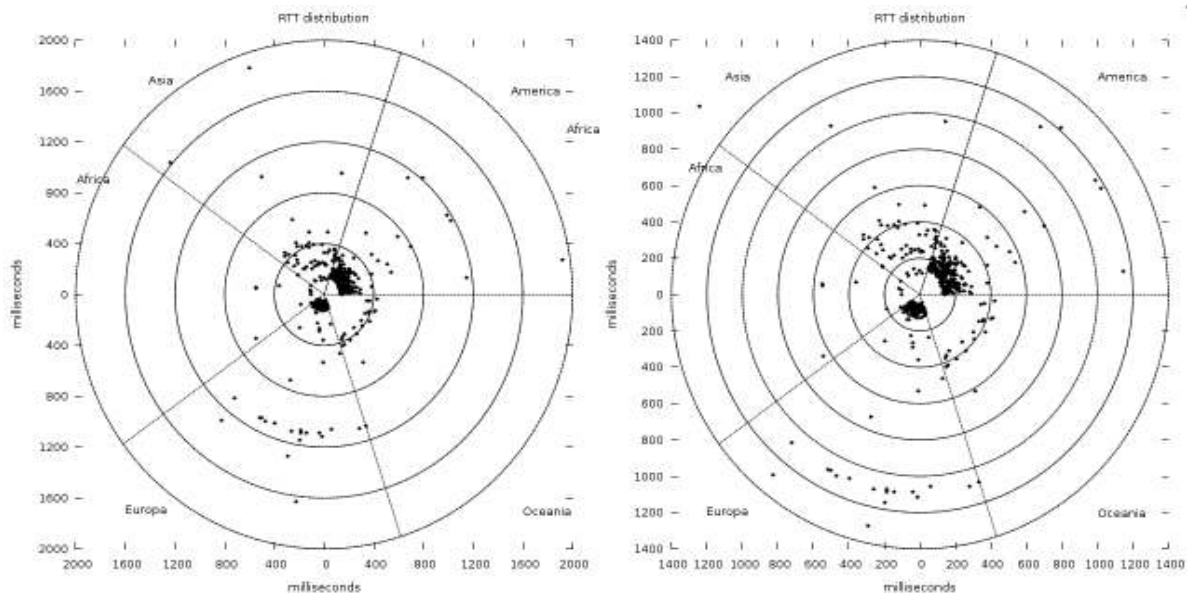


Figure 5.12 – Round Trip Time distribution of the ADSL 4 Mbps Internet connection – *Linux distribution category.*

As it can be observed, it is true that as anticipated European peers present lower values of RTT, which are concentrated at the centre of plots until the 200 ms circle, since the origin host is localized in Portugal. American and Asian peers have mostly of their RTT values located between 100 and 400 ms. Finally, Oceania RTTs are around 400 ms and even more. Unfortunately, due to the low number of African peers it is not so easy to conclude about RTT values to end-hosts localized at this continent but it can be observed that they have similar values for peers located in America and Asia, which is acceptable since distances are similar.

From this analysis it is possible to conclude that there is a strong dependence between RTT and the physical distance between hosts.

In a second analysis, we measured RTT in order to conclude about its dependency on Internet connection types that are used.

Comparing polar maps shown above, a high concentration of peers around 1000 and 1400 ms is visible for the ADSL 4 Mbps Internet connection, especially for peers localized in Europe, but also in America and Asia. Oceania also has more peers located around 1400 ms for the ADSL 4 Mbps Internet connection when compared to the CATV 12 Mbps connection.

In order to better understand obtained results, a new analysis will be made on the number of peers per country that are located between 1000 and 1400 ms, for each Internet connection type and for each type of file that was selected for this study. Obtained results are presented in following plots.

➤ 2008 movies category

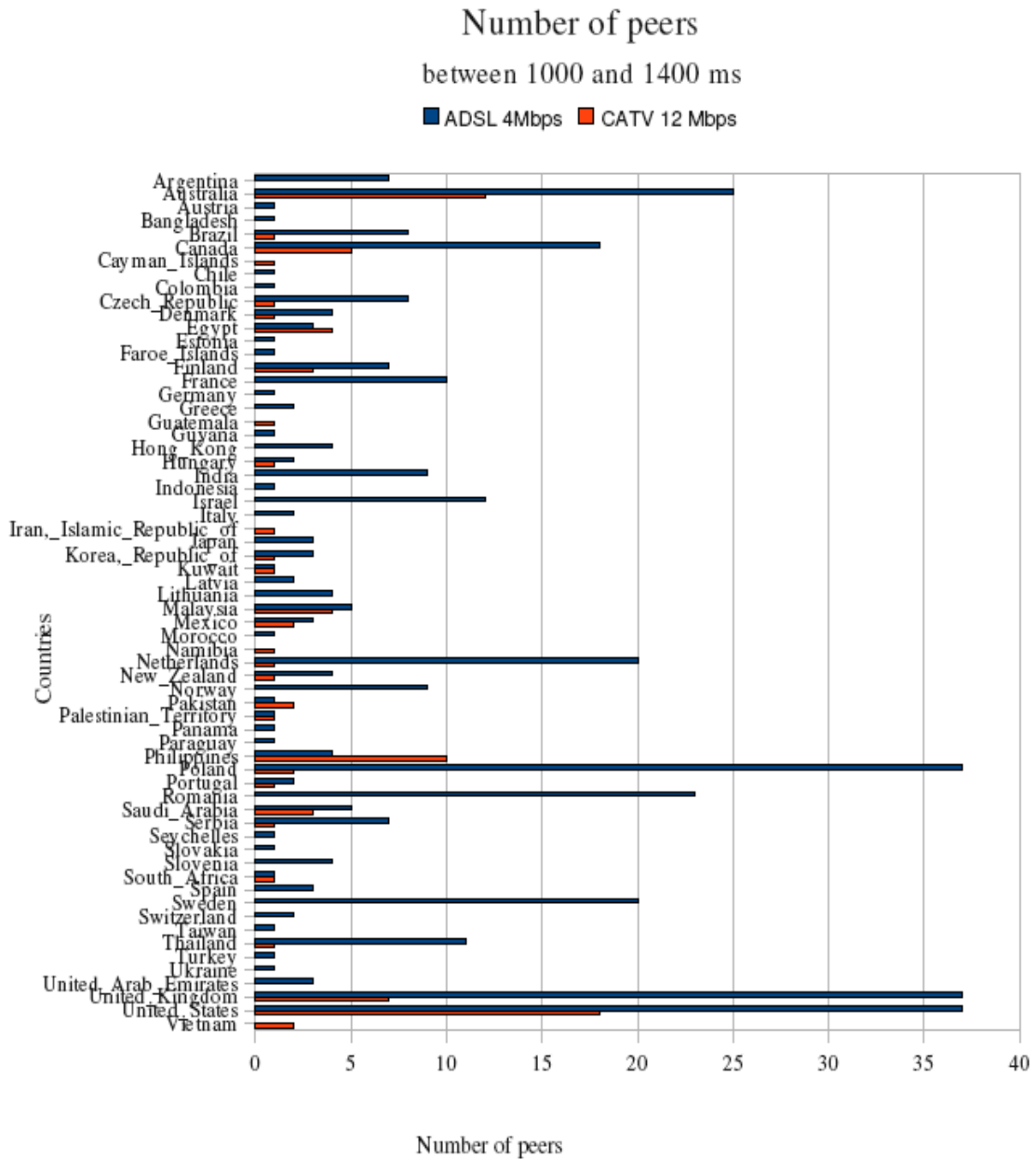


Figure 5.13 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connection types - 2008 *Movies* category.

➤ **Music category**

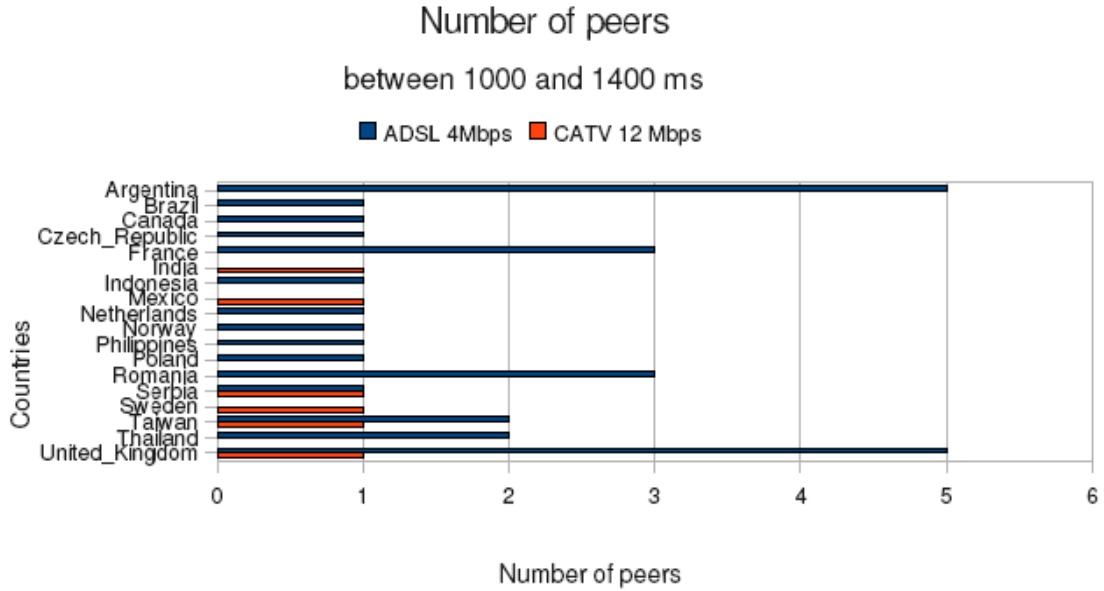


Figure 5.14 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - *Music* category.

➤ **Animated movies category**

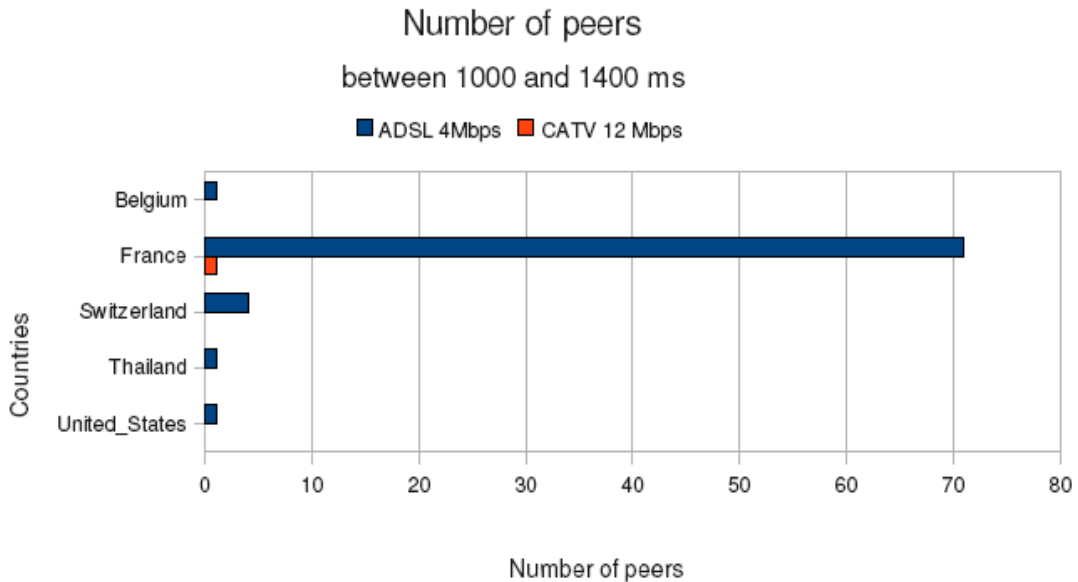


Figure 5.15 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - *Animated movies* category.

➤ **French movies category**

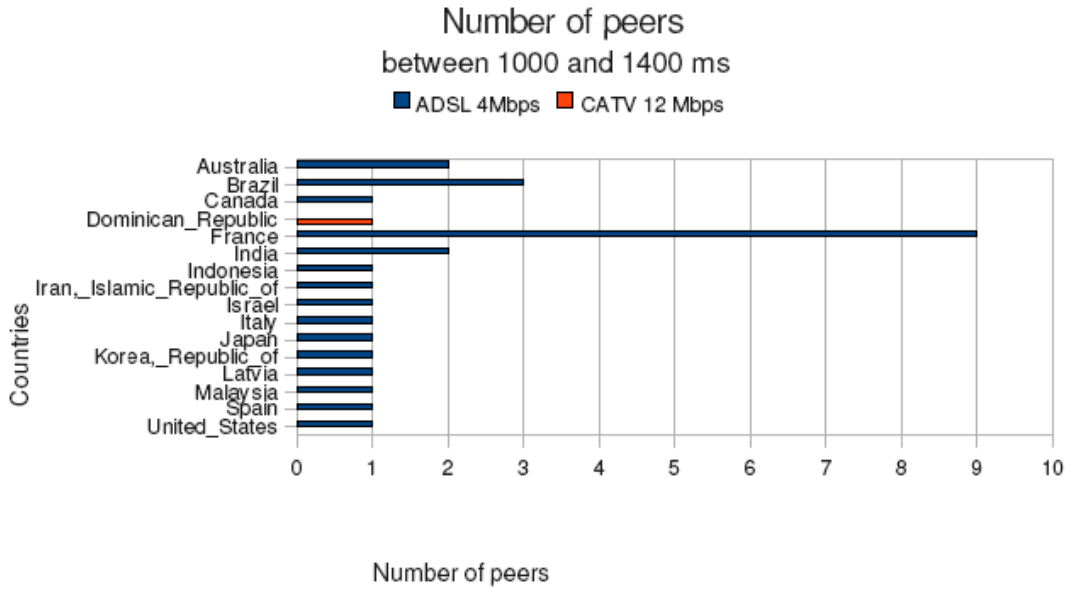


Figure 5.16 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - *French movies* category.

➤ **Indian movies category**

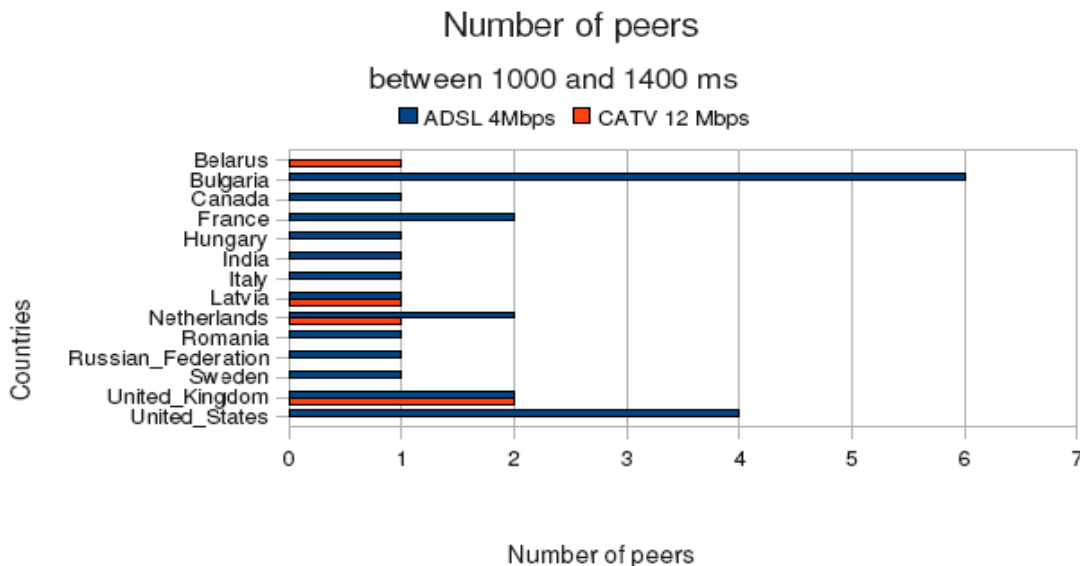


Figure 5.17 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - *Indian movies* category.

➤ *Linux distribution category*

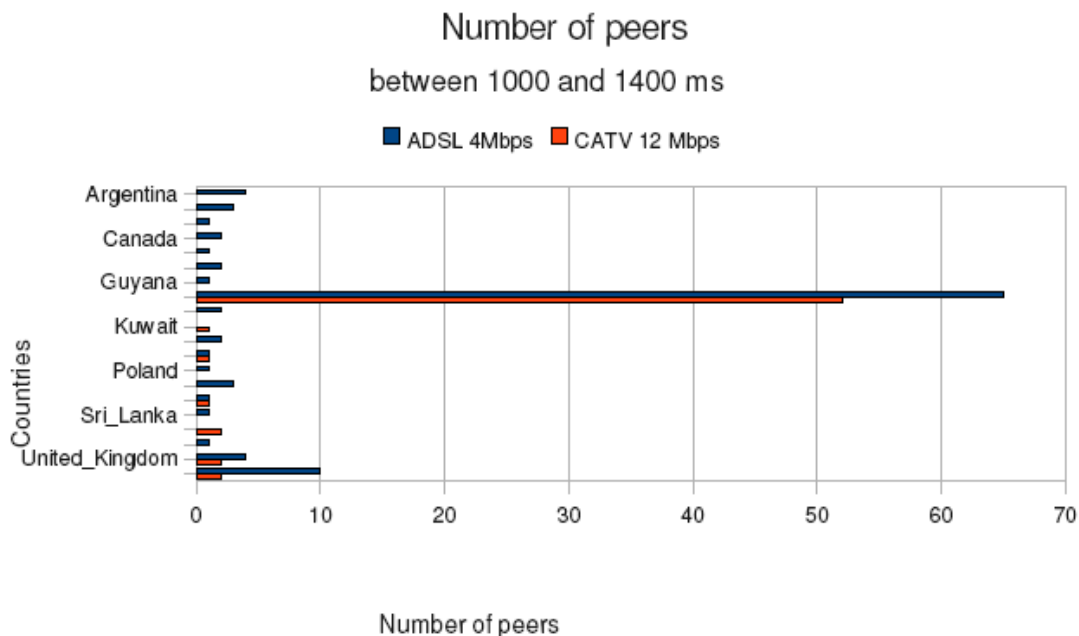


Figure 5.18 – Comparison of peers per country localized between 1000 and 1400 ms for both studied Internet connections - *Linux distribution category*.

Having discriminated peers per country and in order to better understand the peers' RTT distribution between 1000 and 1400 ms, the possible existence of different paths from the origin to the end-host for each connection type, was analysed. The traceroute tool was used to obtain those paths, since it traces packets' routes from one host to another. Since different connections can have different routes, this can probably be the reason for such differences on the obtained RTT values. When a route can not be traced and the corresponding RTT can not be measured, the traceroute tool also informs us about where and why the problem occurred, if it was a network that was shutdown or a router that was not working.

Bar plots clearly show the increase on RTT values in some countries between 1000 and 1400 ms, especially for the *2008 movies* category that has more peers involved and more countries crossed.

Using the traceroute tool, it was possible to verify that different routes really occur depending on the type of Internet connection. It was possible to confirm that, in general, the

CATV 12 Mbps Internet connection gives more direct routes between the origin and end-host than the ADSL 4 Mbps Internet connection.

In ADSL 4 Mbps Internet connections, it was observed that paths to other European countries usually pass through United Kingdom and generally, these connections have longer routes than the corresponding CATV 12 Mbps Internet connection. Tests made for the CATV 12 Mbps Internet connection revealed that they usually use a more direct route from the origin to the end-host, so they usually do not need to pass through any country more than those corresponding to end-hosts. This behaviour did never happen on tests that were made for the ADSL 4 Mbps Internet connection, which obviously explains worse RTT results that were obtained in this case. Tests made for hosts localized in countries like France, Germany, Italy, Netherlands and United Kingdom proved the existence of this direct connection between hosts.

For other continents except Europe, obtained traces were more mixed, since best paths were not always given by CATV 12 Mbps Internet connection, although the passage through United Kingdom still keeps appearing in the ADSL 4 Mbps Internet connection. Countries like Australia, Japan, Malaya, Morocco, Philippines, Taiwan, Turkey and United States have shorter routes in the ADSL 4 Mbps Internet connection.

Furthermore, it is certainly important to notice that in order to reach certain destinations, packets have to pass through several networks, sub-networks, routers and so on, so they can face blocking problems along this path. Thereby, a route change will be needed in order to reach the intended destination. Such problems and eventual route changes cause an increase on the time it takes to arrive to the destination and sometimes these changes do not solve connectivity problems at all. So, it is easy to understand that the longer the route is, the higher is the Round Trip Time and the probability of occurring an error that can make the communication between end points impossible.

Finally, an analysis of the RTT variation along the time of the day was also made. Figures presented below are 3-Dimensional plots, corresponding to one day of analysis, illustrating the probability of different RTT values as a function of the time of the day. For each kind of file, two plots are shown, each one corresponding to one of Internet connection types.

➤ 2008 movies category

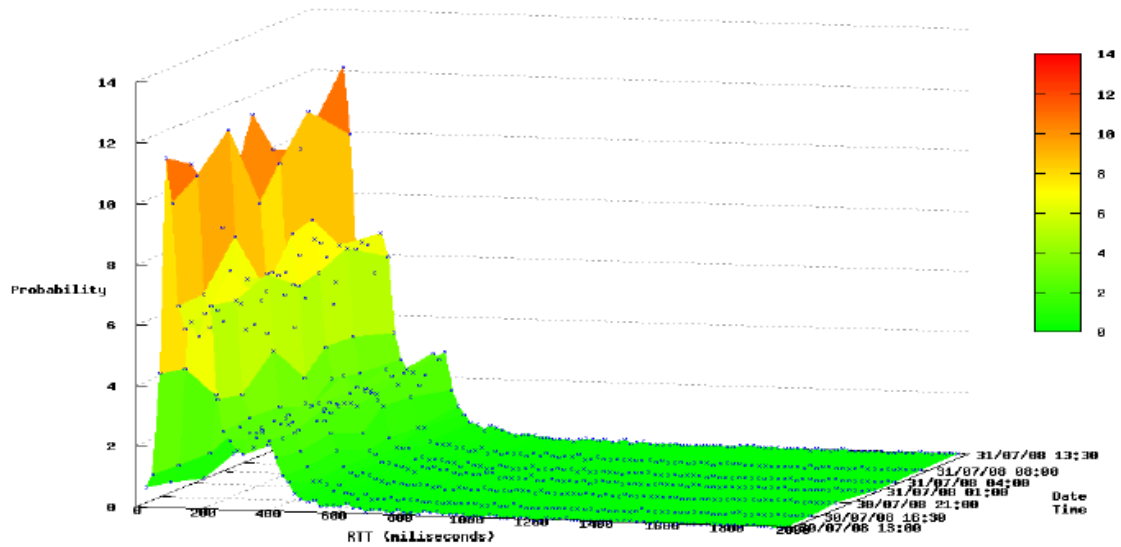


Figure 5.19 – RTT distribution in one day of analysis - 2008 movies category with CATV 12 Mbps Internet connection.

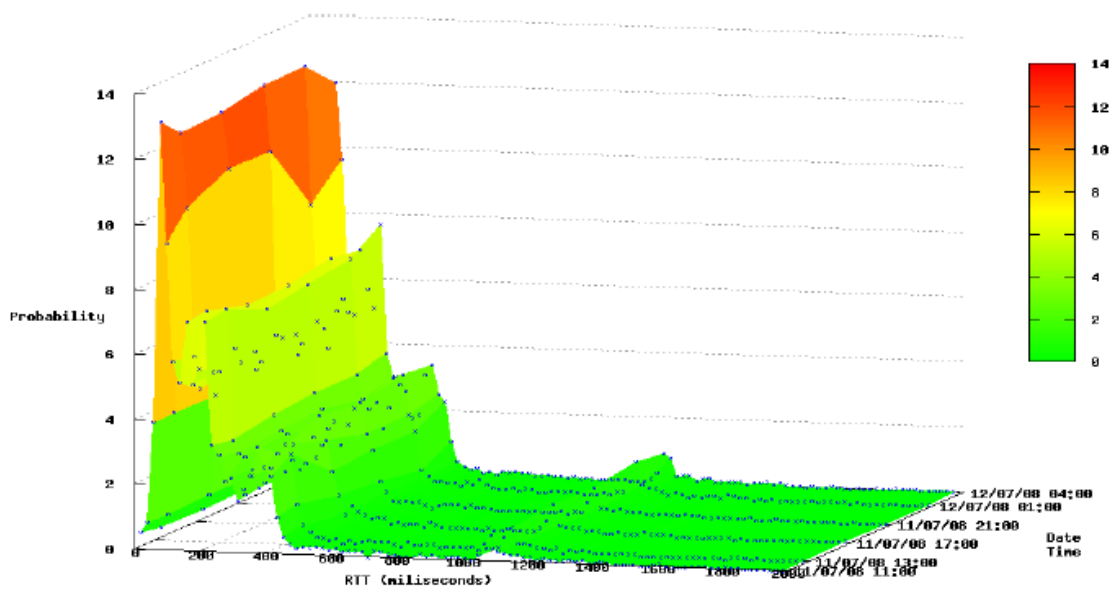


Figure 5.20 – RTT distribution in one day of analysis - 2008 movies category with ADSL 4 Mbps Internet connection.

➤ *Music category*

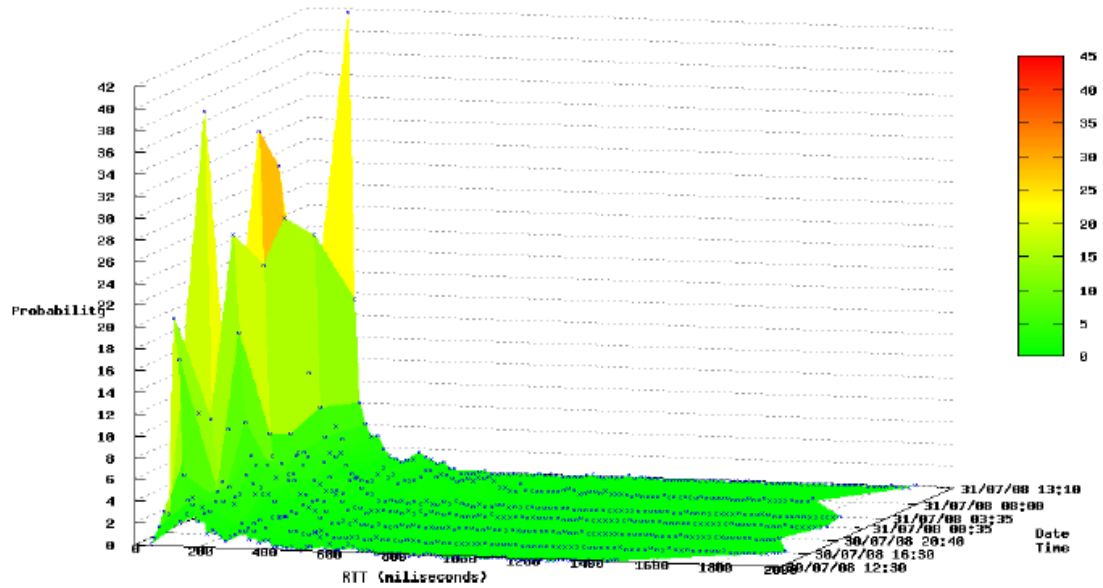


Figure 5.21 – RTT distribution in one day of analysis - *Music* category with CATV 12 Mbps Internet connection.

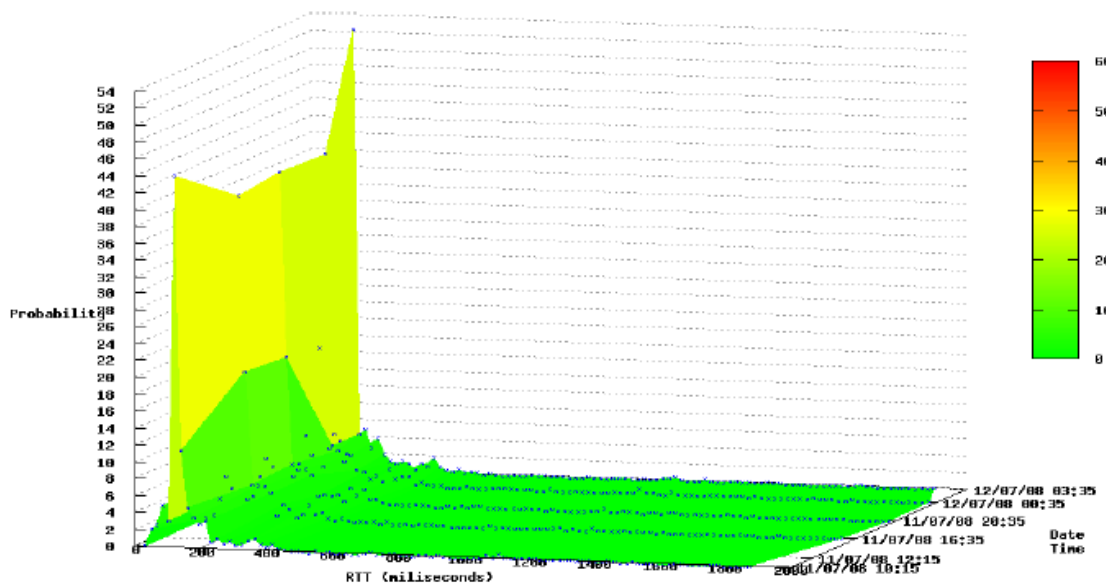


Figure 5.22 – RTT distribution in one day of analysis - *Music* category with ADSL 4 Mbps Internet connection.

➤ *Animated movies category*

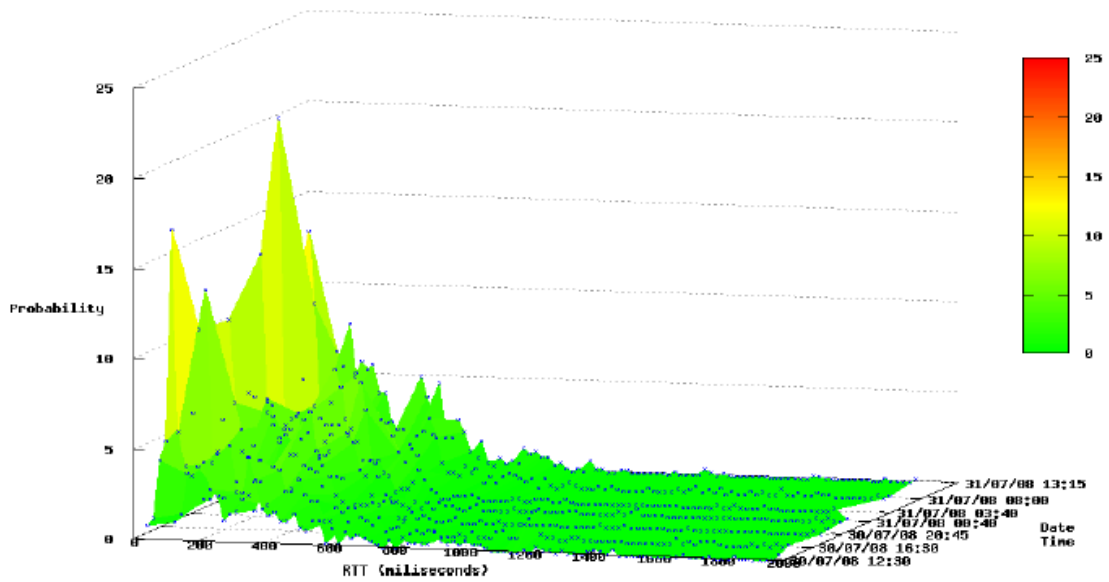


Figure 5.23 – RTT distribution in one day of analysis - *Animated movies* category with CATV 12 Mbps Internet connection.

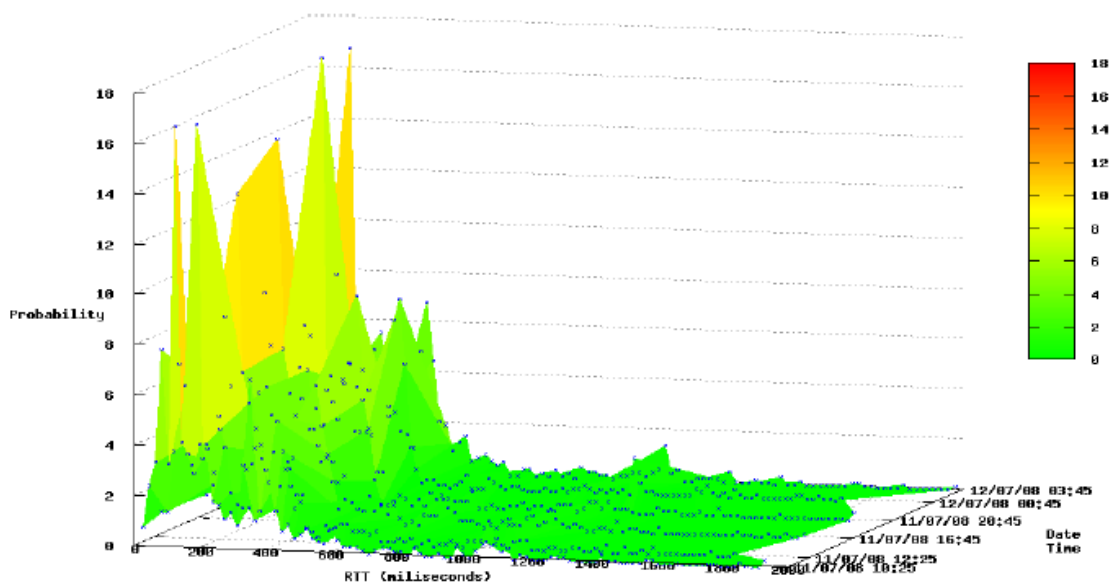


Figure 5.24 – RTT distribution in one day of analysis - *Animated movies* category with ADSL 4 Mbps Internet connection.

➤ *French movies category*

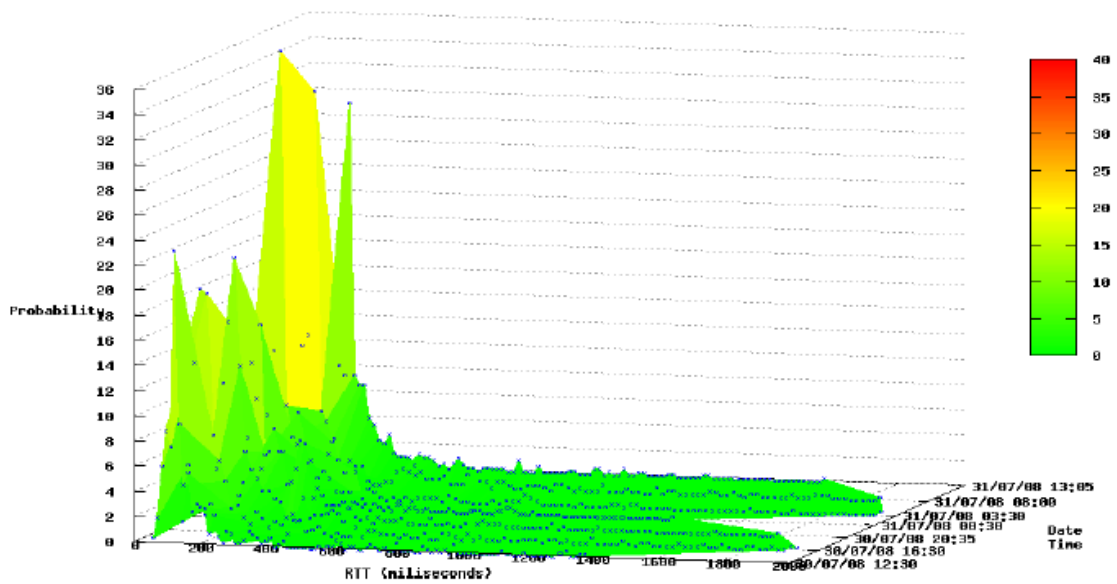


Figure 5.25 – RTT distribution in one day of analysis - *French movies category* with CATV 12 Mbps Internet connection.

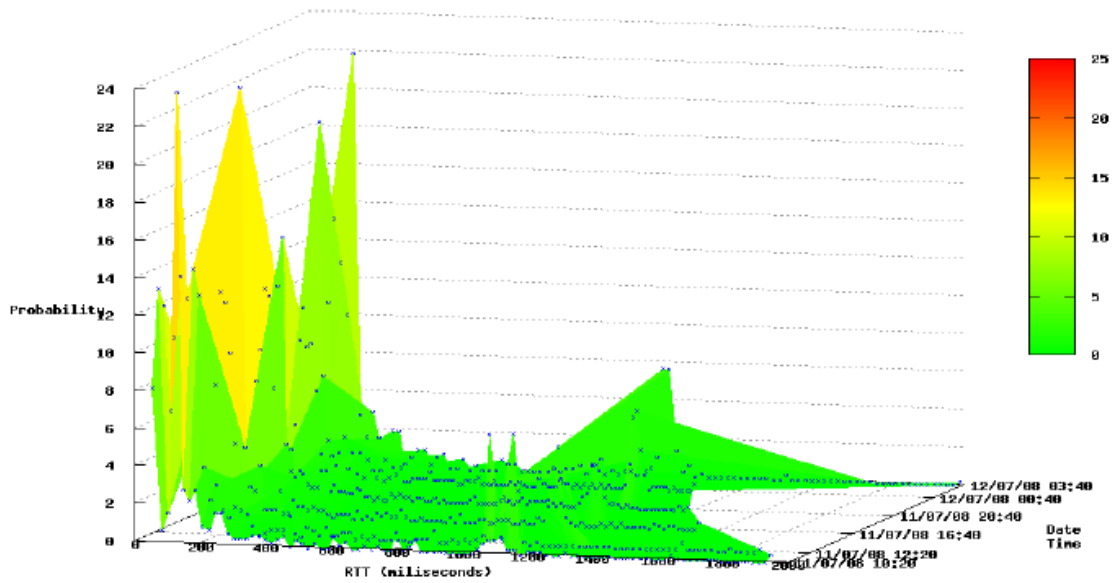


Figure 5.26 – RTT distribution in one day of analysis - *French movies category* with ADSL 4 Mbps Internet connection.

➤ *Indian movies category*

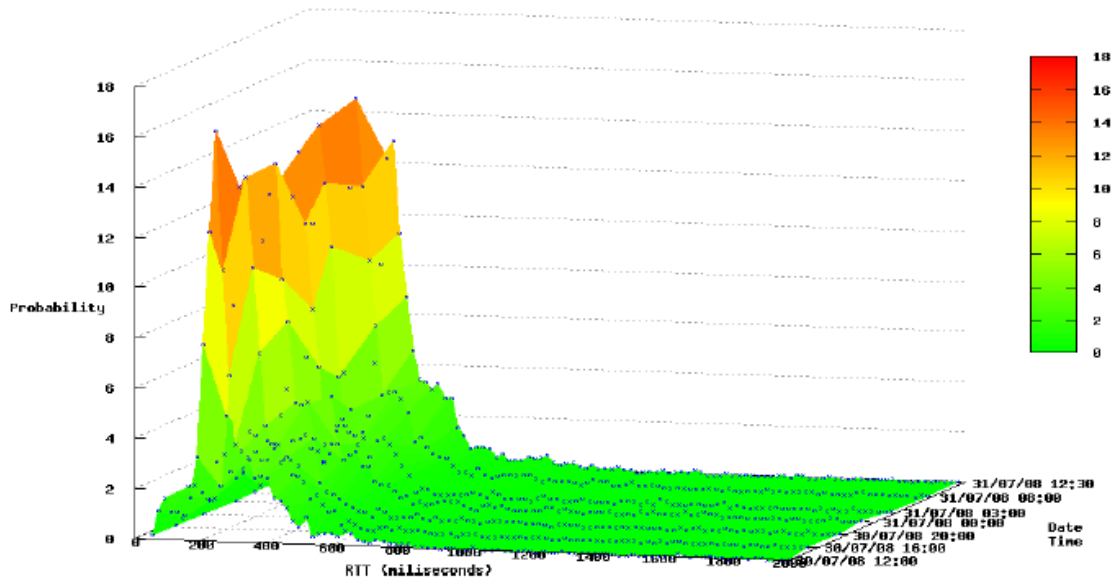


Figure 5.27 – RTT distribution in one day of analysis - *Indian movies* category with CATV 12 Mbps Internet connection.

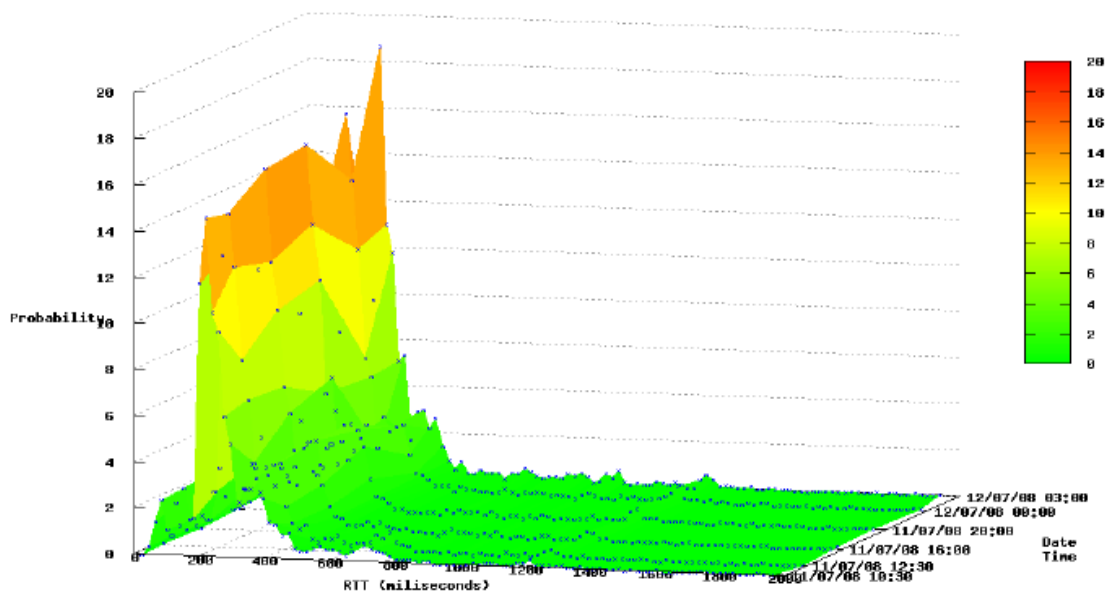


Figure 5.28 – RTT distribution in one day of analysis - *Indian movies* category with ADSL 4 Mbps Internet connection.

➤ *Linux distribution category*

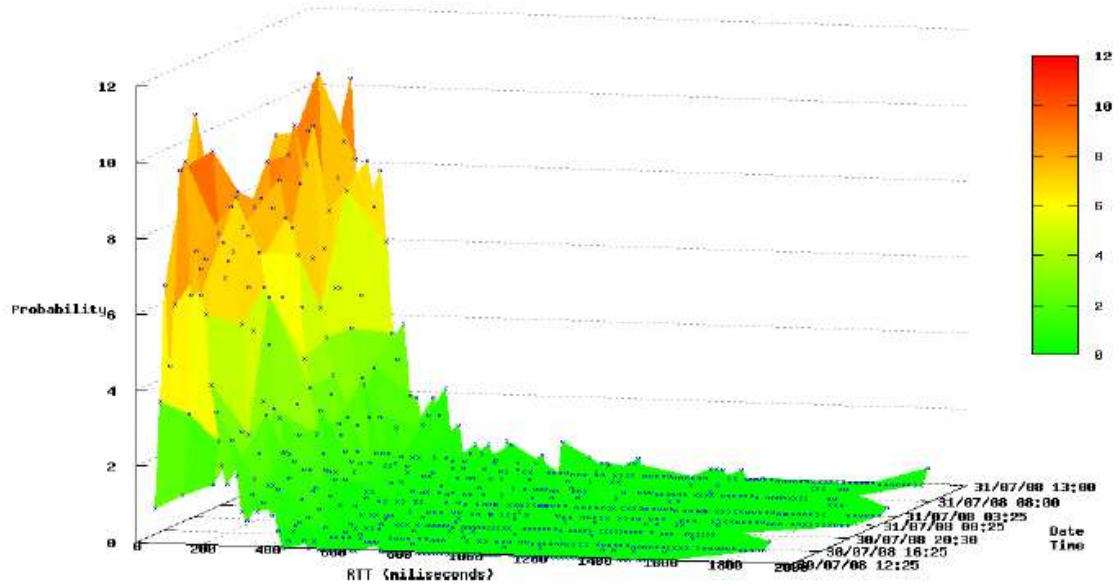


Figure 5.29 – RTT distribution in one day of analysis - *Linux distribution category* with CATV 12 Mbps Internet connection.

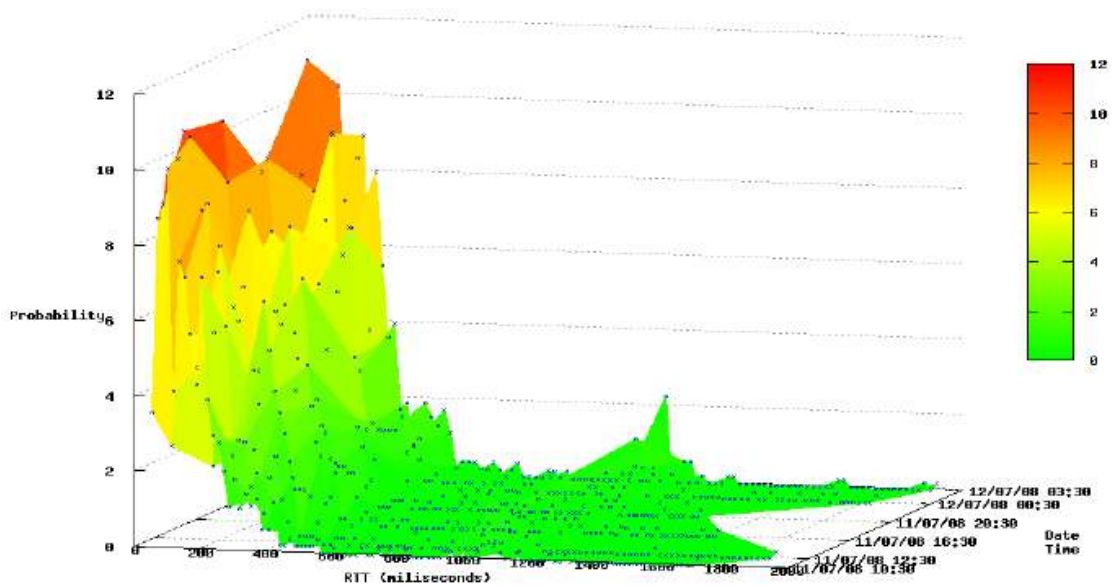


Figure 5.30 – RTT distribution in one day of analysis - *Linux distribution category* with ADSL 4 Mbps Internet connection.

From plots shown above, it is not possible to see significant changes in different daily periods. This fact can be explained by the huge increase on the Internet usage over the last few years, due mainly to easier and faster Internet access, an increase in the amount and quality of the available information and certainly also an increase on the number of available and used P2P networks. These facts allow everybody to perform currently quick and efficient downloads of any kind of files, leading to a significant increase on the Internet traffic and to a flattening behaviour of the P2P traffic profile for both shorter and longer time periods.

Again, it was possible to analyse and study differences for RTT values probability around 1000 ms to 1400 ms for both ADSL 4 Mbps and CATV 12 Mbps Internet connections. Following plots present the cumulative distribution function (CDF) of RTT values of available peers for each type of selected file and connection.

➤ *2008 movies category*

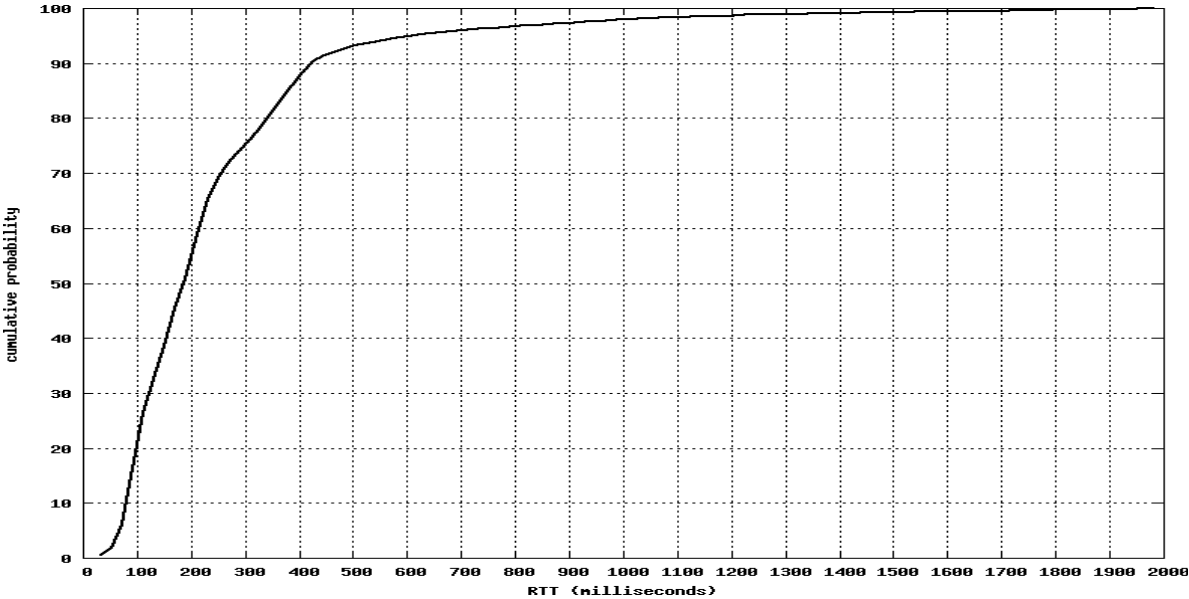


Figure 5.31 – RTT cumulative distribution – *2008 movies* category with CATV 12 Mbps Internet connection.

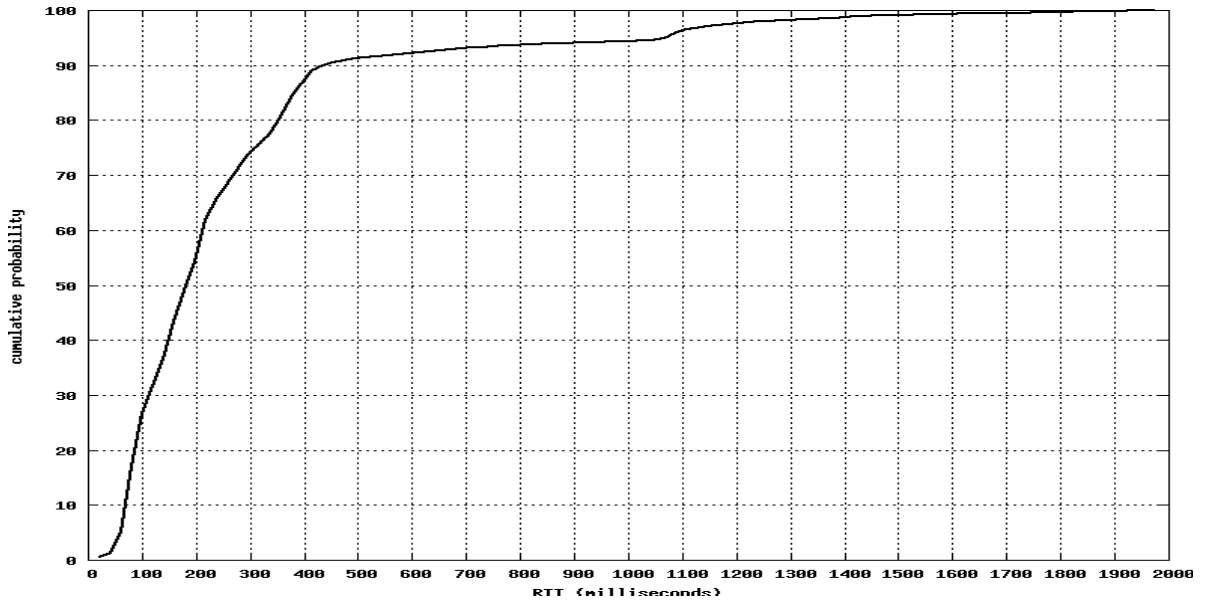


Figure 5.32 – RTT cumulative distribution – *2008 movies* category with ADSL 4 Mbps Internet connection.

➤ *Music* category

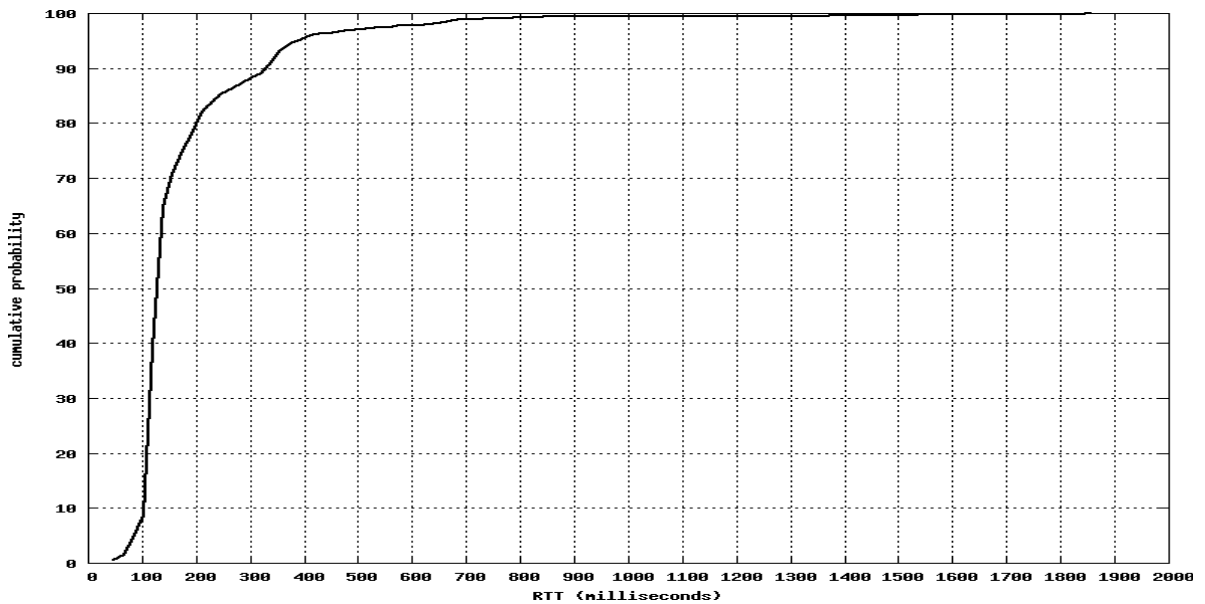


Figure 5.33 – RTT cumulative distribution – *Music* category with CATV 12 Mbps Internet connection.

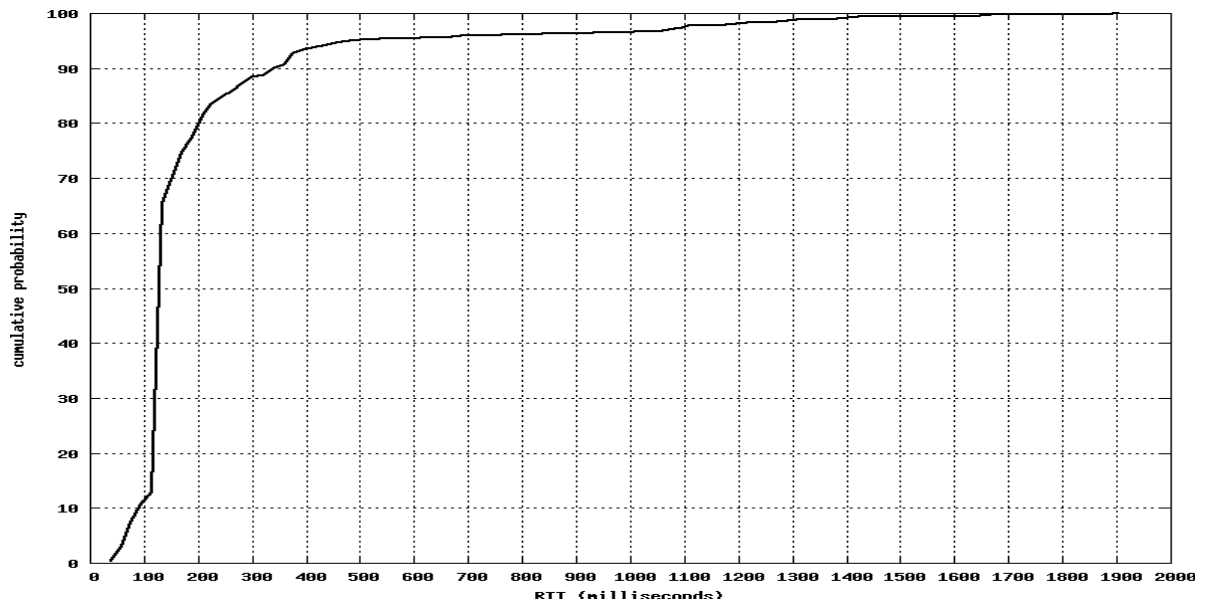


Figure 5.34 – RTT cumulative distribution – *Music* category with ADSL 4 Mbps Internet connection.

➤ *Animated movies* category

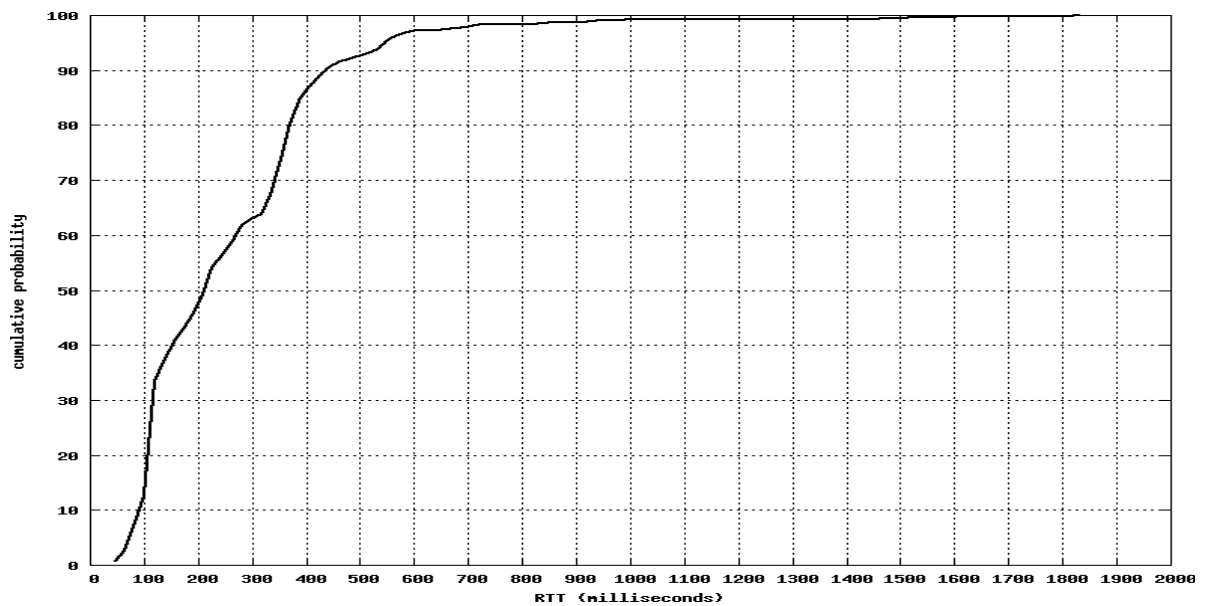


Figure 5.35 – RTT cumulative distribution – *Animated movies* category with CATV 12 Mbps Internet connection.

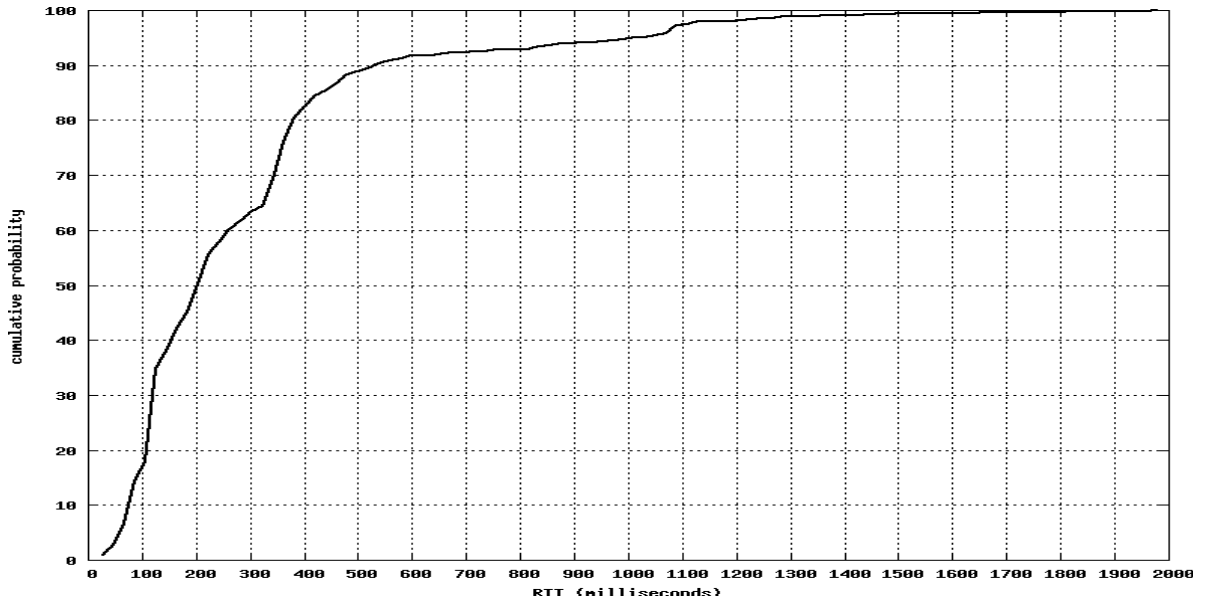


Figure 5.36 – RTT cumulative distribution – *Animated movies* category with ADSL 4 Mbps Internet connection.

➤ *French movies* category

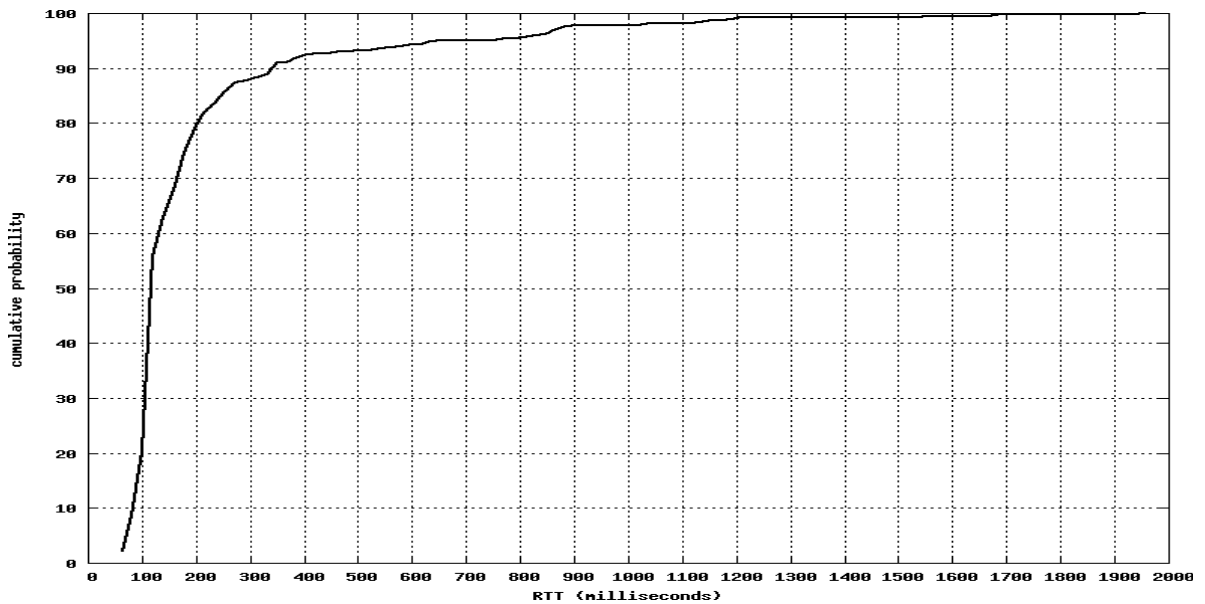


Figure 5.37 – RTT cumulative distribution – *French movies* category with CATV 12 Mbps Internet connection.

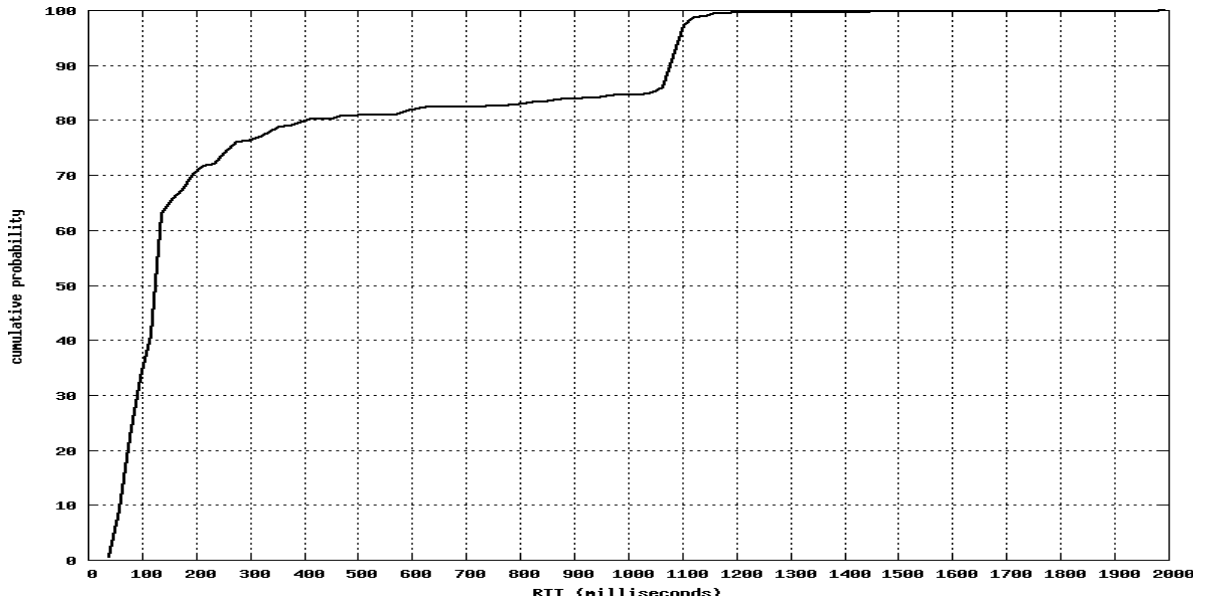


Figure 5.38 – RTT cumulative distribution – *French movies* category with ADSL 4 Mbps Internet connection.

➤ *Indian movies* category

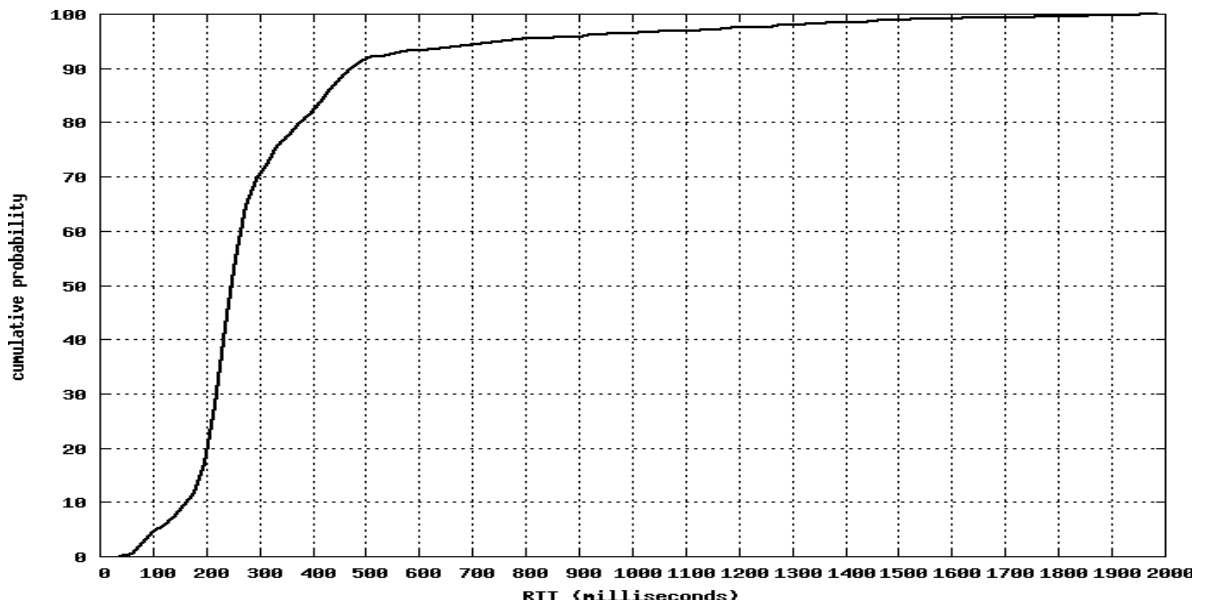


Figure 5.39 – RTT cumulative distribution – *Indian movies* category with CATV 12 Mbps Internet connection

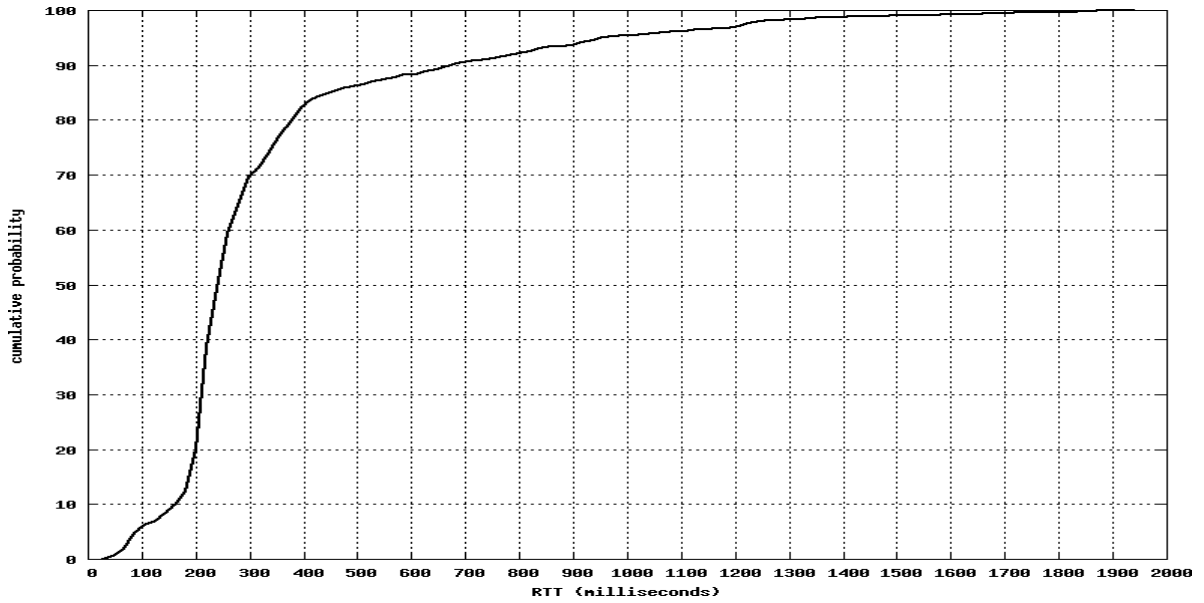


Figure 5.40 – RTT cumulative distribution – *Indian movies* category with ADSL 4 Mbps Internet connection.

➤ *Linux distribution category*

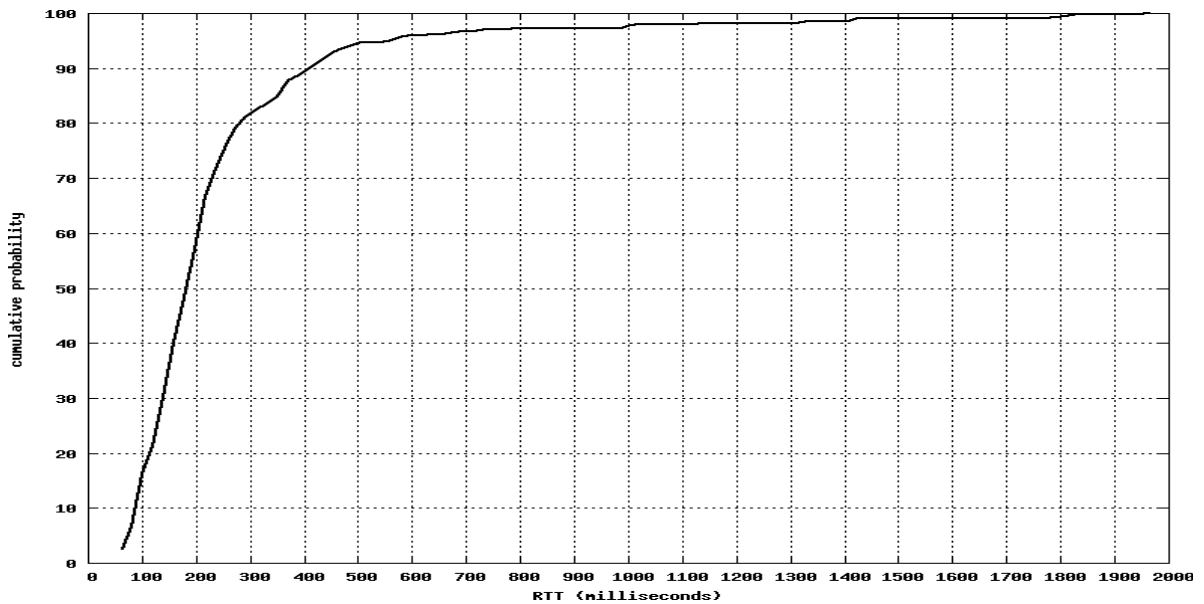


Figure 5.41 – RTT cumulative distribution – *Linux distribution* category with CATV 12 Mbps Internet connection.

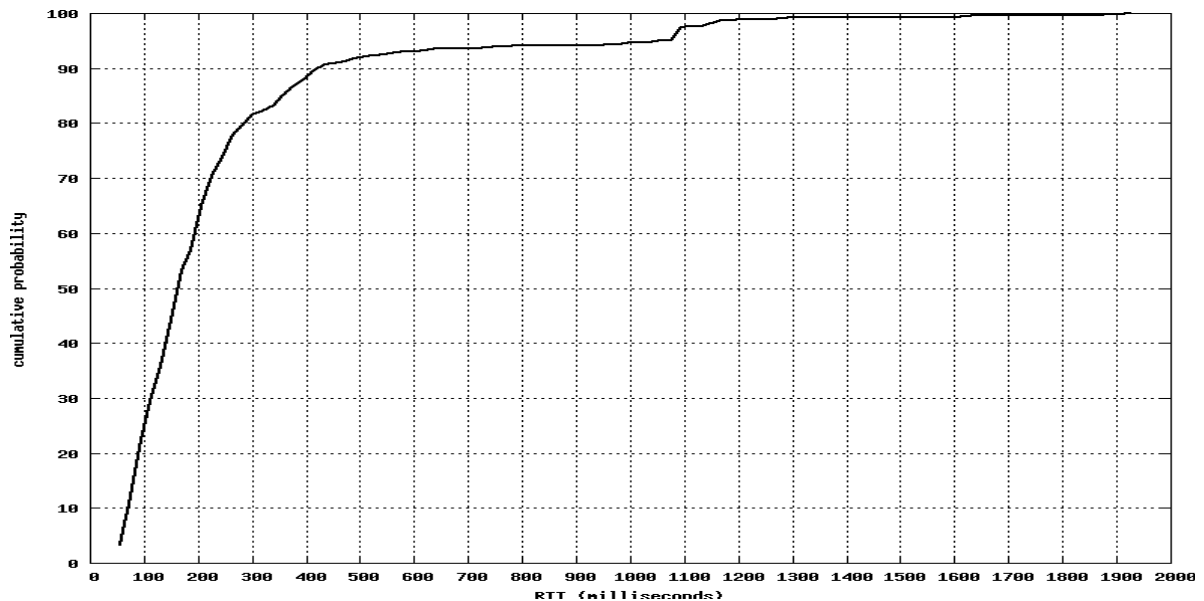


Figure 5.42 – RTT cumulative distribution – *Linux distribution* category with ADSL 4 Mbps Internet connection.

CDF plots help us to visualize how RTT evolves, showing time instants that correspond to higher probability values. It is also possible to verify that RTT between origin and end-hosts are, in general, concentrated from 0 up to 400 ms.

For the *Music category*, 70% of measured RTTs are localized between 0 and 150 ms, which can be justified since almost 40% of peers were localized on United States of America and nearly 35% on Europe. As we had already concluded, RTT values calculated from Portugal to end-hosts placed in Europe present the lowest values, followed by RTT values to hosts located in USA. Also for the *French movies* category, 70% of obtained RTT results are localized between 0 and 200 ms, which was already expected since more than 75% of peers are localized in France.

On the other hand, *2008 movies* and *Linux distribution* categories, due to the major diversity on the localization of involved peers, have a slow evolution on CDF plots.

For the *Indian movies* category, the strongest concentration of peers occurs in the Asian continent especially in India, causing a faster RTT probability growth for RTT values located between 200 and 300 ms, which represent almost 50% of the total number of measured values.

In the *Animated movies* case, two different time periods were identified as corresponding to faster evolutions of the CDF: the range between 50 and 125 ms, corresponding to 35% of the whole measured RTT values and the range between 325 and 375 ms, corresponding to 20% of the whole measured RTT values. These results were expected since almost 17% of peers are from Taiwan, justifying the existence of RTTs between 325 and 375 ms and almost 40% of peers are localized in the European continent, which can explain the first zone of the fast RTT probability growth.

Furthermore, it is also possible to confirm previous conclusions about different Internet connection types. Observing CDF plots, we can see that there is generally a higher concentration of RTT values between 1050 and 2000 ms for the ADSL 4 Mbps Internet connection, whereas the CATV 12 Mbps Internet connection presents an almost constant curve because we could not find a significant number of measured RTTs between such values.

5.1 Summary

This chapter aimed to make a deep analysis about RTT values and factors that could influence their variations. In this way, the chapter started presenting a study on the distance between the origin-host and the localization of the end-hosts. Then, using two different Internet connections, it was possible to identify different characteristics and behaviours between these connections that could be the cause for observed differences on measured RTT values. Finally, the influence of the time of the day on RTT variations was also studied and analysed in some detail.

6. Conclusions

This thesis focused on the characterization of P2P networks, specifically BitTorrent. In order to reach this goal, a BitTorrent client - Vuze - was used to share different kind of files and some important aspects of the BitTorrent network functionality were measured and analysed: dependences of the Round Trip Time with the type of Internet connectivity and the period of the day, as well as several peer-level characteristics such as the geographical distribution of involved peers, the availability of peers during the whole period of analysis and for different daily periods and the variability of RTT according to distances between origin and end-hosts.

For this purpose, the study was developed during an 8 month period, from January to August 2008 and the Vuze client was used to download six different categories of files, involving a total of 25 downloaded files and 85293 peers.

From the analysis of the geographical distribution of involved peers, some conclusions can be drawn:

- P2P file sharing systems comprehend a large number of countries; continents with more peers involved are Europe, America and Asia;
- United States of America, due to its population quantity, cultural diversity and development level, is the country with more peers involved, in absolute numbers, but when normalized by their percentage of the world population, its position in the overall ranking decreases on some places;
- China, being today a huge and technologically developed country does not appear with a significant position on this application, allowing us to conclude that this country has its own file sharing protocols that are more attractive for Chinese people and therefore, more used than BitTorrent;
- Peers distribution depends on the category under analysis, proving that there are different interests for different countries.

After the geographical distribution, the peers' availability was evaluated and some remarks can be made:

- a strong number of peers is protected by Firewall or are localized behind NAT boxes and proxies, avoiding their detection. In fact, the percentage of available peers is far from 100%, being in almost all categories near 25%, with the exception of the *Music* category that starts with an availability of 45%;
- during all the analysed period, the peers' availability remains almost constant, proving high usage levels of this file sharing application.

Finally, concerning obtained Round Trip Time results, it can be concluded that:

- Round Trip Time is strongly dependent on the distance between the origin and destination hosts;
- different types of Internet connections lead to different route paths, resulting in visible differences on RTT results;
- CDF plots showed that the majority of RTT values are located between 0 and 300 milliseconds.

6.1 Future work

Obtained results can be used to improve the performance of P2P networks and develop new and more efficient protocols for these systems. In fact, with a complete and concise information regarding the peers' distribution around the world and RTT dependences between hosts, new protocols can be engineered in order to make files' download even faster, including for example new algorithms that establish as a priority the choice of the closest peers from the origin-host, mainly from the same ISP thus preventing unnecessary capacity consumptions.

This kind of studies can also be helpful in order to improve Quality of Service assurances in these P2P networks as well as to prevent and to avoid overload problems on the Internet itself.

Results obtained on this study can also be helpful on the deployment and management of new real-time multimedia content distribution services based on P2P architectures.

It is also important to mention that the fast growth of P2P file sharing systems, as it was verified in this thesis, attracts a large percentage of the whole population (including some that

are protected by Firewalls or are localized behind NAT/PAT) and it has some associated problems that need to be solved in the future, such as:

- the existence of bad nodes, which can provide polluted contents due to the lack of control (basically, everyone can participate on these P2P systems);
- attacks from one node to another can easily happen due to the facility of obtain information about them (IP address and geographical localization), which constitutes a potential danger to these systems and to the whole network security itself;
- the existence of huge amounts of illegal contents, making the development of mechanisms that are able to control or avoid their distribution very urgent and necessary;
- the traffic load that these systems represent on the whole Internet traffic is also a problem that ISPs have to solve.

Acronyms

DHT	Distributed Hash Table
IPS	Intrusion Prevention Systems
IPTV	Internet Protocol Television
NAT	Network Address Translation
PAT	Port Address Translation
CDF	Cumulative Distribution Function
P2P	Peer-to-Peer
QoS	Quality of Service
RTT	Round Trip Time
VoD	Video on Demand

References

- [1] K. Lua, J. Crowcroft, M. Pias, R. Shama and S. Lima. A survey and comparison of peer-to-peer overlay network schemes. *IEEE communications survey and tutorial*, vol. 7, no. 2, pp. 72–93, 2005.
- [2] T. Karagiannis, A. Broido, M. Faloutsos and K. Claffy. Transport layer identification of P2P traffic. *Proceeding ACM Sigcomm Internet Measurement Conference*, October 2004.
- [3] S. Sen, O. Spatscheck and D. Wang. Accurate, scalable in-network identification of P2P traffic using application signatures. *Proceeding of the 13th international Conference on World Wide Web*, pp. 512–521, New York, USA, May 2004.
- [4] S. Saroiu, P. K. Gummadi and S. D. Gribble. A measurement study of peer-to-peer file sharing systems. *Proceeding of Multimedia Computing and Networking*, 2002.
- [5] N. Leibowitz, A. Bergman, R. Ben-Shaul and A. Shavit. A measurement study of peer-to-peer file sharing systems. *Proceeding of the 7th Int. WWW Caching Workshop*, 2002.
- [6] R. Bhagwan, S. Savage and G. Voelker. Understanding availability. *IPTPS 2003: international workshop on peer-to-peer systems, Lecture notes in computer science, vol. 2735*, pp. 256–267, 2003.
- [7] C. Gkantsidis, M. Mihail and A. Saberi. Random walks in peer-to-peer networks. *Proceeding of INFOCOM 2004, Twenty-third Annual Joint Conference of the IEEE Computer and Communications Societies*, 2004.
- [8] D. Qiu and R. Srikant. Modeling and performance analysis of BitTorrent-like peer-to-peer networks. *Proceeding of ACM SIGCOMM '04, pp. 367–378.*, Portland, Oregon, USA, August - September 2004.
- [9] Y. Tian, D. Wu and K. W. Ng. Modeling, analysis and improvement for BitTorrent-like file sharing networks. *Proceeding of INFOCOM 2006, 25th IEEE International Conference on Computer Communications*, Barcelona, Spain, April 2006
- [10] P. Michiardi, K. Ramachandran and B. Sikdar. “Modeling seed scheduling strategies in BitTorrent. *Proceeding of NETWORKING 2007, Ad Hoc and Sensor Networks, Wireless*

Networks, Next Generation Internet, pp 606-616, *6th International IFIP-TC6 Networking Conference*, Atlanta, USA, May 2007.

[11] T. Isdal. Using BitTorrent for measuring end-to-end Internet path characteristics. Msc. Thesis, The Royal Institute of Technology (KTH), Stockholm, Sweden, October 2006.

[12] D. Erman. Extending BitTorrent for streaming applications. *Proceeding 4th Euro-FGI Workshop on New Trends in Modelling, Quantitative Methods and Measurements*, Ghent, Belgium, 2007.

[13] M. Izal, G. Urvoy-Keller, E. Biersack, P. Felber, A. Al Hamra and L. Garcés-Erice. Dissecting BitTorrent: Five months in a torrent's lifetime. *Proceeding Passive and Active Network Measurement*, pp. 1–11, France, April 2004.

[14] A. Iosup, P. Garbacki, J. A. Pouwelse and D.H.J. Epema. Analyzing BitTorrent: Three lessons from one peer-level view. *Proc. 11th ASCI Conference*, 2005.

[15] J. A. Pouwelse, P. Garbacki, D.H. Epema and H.J. Sips. An introduction to the BitTorrent peer-to-peer file sharing system. *Hand-out at the 19th IEEE Annual Computer Communications Workshop*, October 2004.

[16] J. Pouwelse, P. Garbacki, D. Epema and H. Sips. The BitTorrent P2P file-sharing system: Measurements and analysis. *Proceeding of the 4th International Workshop on Peer-To-Peer Systems (IPTPS'05)*, Ithaca, New York, USA, February 2005.

[17] A. Iosup, P. Garbacki, J. A. Pouwelse and D. H. Epema. Correlating topology and path characteristics of overlay networks and the Internet. *Sixth IEEE International Symposium on Cluster Computing and the Grid Workshops*, May 2006.

[18] D. Zeinalipour-Yazti and T. Foliás. A quantitative analysis of the Gnutella network traffic. Course project for *Advanced topics in networks* with M. Faloutsos, University of California - Riverside, Dpt. Of CS, April 2002.

[19] K. P. Gummadi, R. J. Dunn, S. Saroiu, S. D. Gribble, H. M. Levy, and J. Zahorjan. Measurement, modelling, and analysis of a peer-to-peer workload. *SOSP 2003: Proceedings of the 9th ACM symposium on Operating systems principles*, pp 314–329, ACM Press, 2003.

[20] N. Leibowitz, M. Ripeanu and A. Wierzbicki. Deconstructing the kaza network. *Proceedings of the Third IEEE Workshop on Internet Applications*, page 112, San Jose, California, June 2003.

- [21] S. Sen and J. Wang. Analyzing peer-to-peer traffic across large networks. *IEEE/ACM Transactions on Networking*, pp 219-232, 2002.
- [22] P. Salvador and A. Nogueira. Study on geographical distribution and availability of BitTorrent peers sharing video files. *12th Annual IEEE International Symposium on Consumer Electronics (ISCE 2008)*, Vilamoura, Portugal, April 2008.
- [23] China taking P2P to the next level. <http://www.zeropaid.com/news/>.
- [24] I. Dedinski, H. De Meer, L. Han, L. Mathy, D. P. Pezaros, J. S. Sventek and X.Y. Zhan. Cross-layer peer-to-peer traffic identification and optimization based on active networking. *Proceeding of the Seventh Annual International Working Conference on Active and Programmable Networks (IWAN'05)*, CICA, Sophia Antipolis, France, November, 2005.
- [25] L. Hammerle. P2P population tracking and traffic characterization of current P2P file-sharing systems. Msc. Thesis, Swiss Federal Institute of Technology of Zurich, April 2004.
- [26] R. Kaspar. P2P file-sharing traffic identification method validation and verification. Semester Thesis, Computer Engineering and Networks Laboratory (TIK), Zurich, Swiss, November 2005.
- [26] B. Cohen. Incentives build robustness in BitTorrent. *BitConjurer*, <http://www.bitconjurer.org>, May 2003.
- [27] C. A. Parker. The true picture of peer-to-peer file sharing. *Tech. Rep., Cachelagic Research*, May 2005.
- [28] C. Dana, D. Li, D. Harrison and C.-N. Chuah. BASS: BitTorrent assisted streaming system for video-on-demand. *Proceeding of IEEE 7th Workshop on Multimedia Signal Processing*, October 2005.
- [29] R. Susitaival. Traffic engineering in the Internet: from traffic characterization to load balancing and peer-to-peer file sharing. PhD thesis, Helsinki University of Technology, Espoo, Finland, September 2007.
- [30] R. Susitaival and S. Aalto. Modelling the population dynamics and the file availability in a BitTorrent-like P2P system with decreasing peer arrival rate. *Proceeding of the International Workshop on Self-Organizing Systems (IWSOS 2006)*, pp 34-48, September 2006.

[31] C. Wang Bo Li. Peer-to-peer overlay networks: a survey. The Hong Kong University of Science and Technology, April 2003.

[32] H. Schulze and K. Mochalski. Internet study 2007 - The impact of p2p file sharing, voice over IP, skype, joost, instant messaging, one-click hosting and media streaming such as youtube on the Internet. Ipoque GmbH, Germany, 2007.

[33] BitTorrent. <http://www.bittorrent.com>.

[34] MaxMind - GeoLite Country - Open Source IP Address to Country Database. <http://www.maxmind.com/app/geolitecountry>.

[35] NMAP free security scanner for network exploration & hacking. <http://nmap.org/>.

[36] Torrent Search Engine. <http://www.torrentz.com/>.

[37] Vuze. <http://www.vuze.com/>

[38] Sandvine. http://www.sandvine.com/news/pr_detail.asp?ID=203

[39] Octave. <http://www.gnu.org/software/octave/>

[40] Gnuplot. <http://www.gnuplot.info/>

Appendix A - Shell scripts

Shell scripts are sequences of commands written on simple text files. In this thesis, the bash shell (Bourne-Again shell) was used to write shell scripts.

The first line of the shell scrip must present the script shell name. As all our shells were written in bash shell, all shell scripts created and shown bellow will present on the first line the command:

```
#!/bin/bash.
```

Furthermore, in order to turn the simple text file on an executable shell script, it is necessary to execute one of the following commands:

- ***\$chmod +x script_filename***
- ***\$chmod 755 your-script-name***

Now the file is ready to be executable whenever it is needed. In this way, we only need to call the script with one of the following commands:

- ***\$path./scriptfilename***
- ***bash scriptfilename***
- ***sh scriptfilename***

where path is the route that allows to reach the file, i.e., where it is saved.

Next, we will present a brief description of some commands that were used in shell scripts of this thesis to measure results.

➤ ***Wildcards***

Give the ability to refer more than one file by their name using special characters.

- * Matches any string or group of characters.
- ? Matches any single character.
- [...] Matches any one of the enclosed characters.

➤ **Quotes**

There are three types of quotes “, ’, `

- " "Double Quotes" Anything enclosed in.
- ' 'Single quotes' Enclosed in single.
- ` `Back quote` To execute command.

➤ **Redirection**

There are three main redirection symbols >, >>, <

- > *command > filename* output commands result to the filename. If the file already exists, it will be overwritten, otherwise a new file will be created.
- >> *command >> filename* output commands result to the end of the filename. If the file already exists, it will be opened and new information will be written at the end of such file, without losing previous information.
- < *command < filename* to take input to command from the filename instead of key-board.

➤ **mkdir** creates a directory
mkdir <dirname>

➤ **rmdir** removes a directory
rmdir <dirname>

➤ **rm** removes a folder with files
rm -r -f <dirname>

- ***cd*** change directory
cd<*newpath*>

 - ***vi*** insert data into file
vi <*newfile.ext*>

 - ***cp*** copy file with same name
cp <*sourcedir*>/<*sourcefilename.ext*> <*destinationdir*>

 - ***mv*** move file with same name
mv <*sourcedir*>/<*sourcefilename.ext*> <*destinationdir*>

 - ***Echo*** display a text or a value of variable.
echo [*options*] [*string, variables...*]
- Options***
- n* Do not output the trailing new line.
 - e* Enable interpretation of the following backslash escaped characters in the strings:
 - a* alert (bell)
 - b* backspace
 - c* suppress trailing new line
 - n* new line
 - r* carriage return
 - t* horizontal tab
 - * backslash
-
- ***cat*** concatenate files and print on the standard output
cat <*filename*>

 - ***wc*** count the number of lines, words, characters from the file

- wc <filename>*
- **cut** selects a portion of a file.
cut -f{field number} {file-name}
 - **paste** put the desired lines together.
paste {file1} {file2}
 - **join** joins, lines from separate files.
join {file1} {file2}
 - **tr** translate range (pattern-1) of characters into other ranges (pattern-2)
tr {pattern-1} {pattern-2}
 - **awk** manipulate data.
awk 'pattern action' {file-name}
 - **sed** stream editor for filtering and transforming text
sed {expression} {file}
 - **grep** finds and prints the matching pattern.
grep "pattern" {file-name}
 - **head** show the contents of the firsts n lines of the file
head -n x \$path
 - **tail** show the contents of the last n lines of the file
tail -n x \$path
 - **date** gives the actual date and hour

➤ ***If condition***

```
if condition
then
    condition is zero (true - 0)
    execute all commands up to else statement
else
    if condition is not true then
    execute all commands up to fi
fi
```

Figure A. 1 - *If condition* structure.

➤ ***Case***

```
case $variable-name in
    pattern1) command
        ...
        command;;
    pattern2) command
        ...
        command;;
    patternN) command
        ...
        command;;
    *)    command
        ...
        command;;
esac
```

Figure A. 2 – *Case condition* structure.

➤ **For loop**

```
for { variable name } in { list }  
  
do  
execute one for each item in the list until the list is not finished  
(And repeat all statement between do and done)  
done
```

Figure A. 3 – *For loop* structure (a).

OR

```
for (( expr1; expr2; expr3 ))  
do  
repeat all statements between do and done until expr2 is TRUE  
done
```

Figure A. 4 – *For loop* structure (b).

➤ **While loop**

```
while [ condition ]  
do  
    command1  
    command2  
    ...  
done
```

Figure A. 5 – *While loop* structure.

➤ *Mathematical comparison*

Operator	Meaning
-eq	is equal to
-ne	is not equal to
-lt	is less than
-le	is less than or equal to
-gt	is greater than
-ge	is greater than or equal to

Table A. 1 – Mathematical comparison.

➤ *String comparison*

Operator	Meaning
string1 = string2	string1 is equal to string2
string1 != string2	string1 is NOT equal to string2
string1	string1 is NOT NULL or not defined
-n string1	string1 is NOT NULL and does exist
-z string1	string1 is NULL and does exist

Table A. 2 – String comparison.

➤ *Shell test for files or directories*

Test	Meaning
-s file	Non empty file
-f file	If file exists or normal file and not a directory
-d dir	If the dir exist and is a directory
-w file	If it is a writeable file

-r file	If it is read-only file
x file	If it is an executable file

Table A. 3 – Shell test for files or directories.

➤ *Logical Operators*

Operator	Meaning
! expression	Logical NOT
expression1 -a expression2	Logical AND
expression1 -o expression2	Logical OR

Table A. 4 – Logical operators.

➤ *Pipes*

A pipe is a way to connect the output of one program to the input of another program without any temporary file.

➤ *Extract the needed information from the Vuze Log file*

```
#!/bin/bash

TORS=`cat $1 | grep -e "COMPACT PEER" | cut -d ":" -f 5 | sort -u`
NT=`echo "$TORS" | wc -l`
echo "Torrents:"
echo "$TORS"
echo $NT

for (( t=1; t<=$NT; t++))
do
    TNAME=`(echo "$TORS" | head -n $t | tail -n 1 | cut -d "" -f2)`
    TNAMEF=`echo $TNAME | sed '\[/s/\\[/g' | sed '\[/s/\\[/g`
    FNAME=`echo "ipsports.$TNAME.txt" | sed '/s//.g' | sed '\[/s//.g' | sed '\[/s//.g`

    cat $1 | grep -e "COMPACT PEER" | grep -e "$TNAMEF" | grep -e "tcp_port" | awk
    '{print $6}' | cut -d "," -f1-2 | cut -d "=" -f2-3 | sed ",tcp_port=/s///" | sort -u > $FNAME

    NPEERS=`cat $FNAME | wc -l`
    echo "$TNAME has $NPEERS peers"
done
```

Figure A. 6 – Shell script to extract information from Log files.

- **Shell script created to obtain the geographical localization of peer and calculate RTT**

```
#!/bin/bash
for n in $*
do
IN=(`cat $n`)
DATE=`date +%F`
echo "Processing $n at $DATE"
i=0
while [ $i -lt ${#IN[*]} ]
do
    eval IP=${IN[$i]}
    let i++
    eval PORT=${IN[$i]}
    let i++
    ./BTstats2.sh $IP $PORT >> res.$n.$DATE.ods &
done
done
```

Figure A. 7 – Shell script for the geographical localization of peers and calculation of the RTT.

- BTstats2.sh

```
#!/bin/bash
COUNTRY_S=`(geoipllookup $1 | grep Country | cut -d " " -f4 | sed '/,/s///g')`
COUNTRY_N=`(geoipllookup $1 | grep Country | cut -d " " -f5- | sed '/,/s//_g')`
RTT=`(nmap -sT $1 -p $2 -P0 --initial-rtt-timeout 2000 --max-rtt-timeout 2000 --scan-delay 300 --max-retries 1 -d | grep -e "Final times" | awk '{print $6}')`

echo "$COUNTRY_N $COUNTRY_S $1 $2 $RTT"
```

Figure A. 8 – Shell script for the geographical localization of peers and calculation of the RTT.

➤ **Results measurement**

- Measure RTT values

```
#!/bin/bash

C=`cat $1 | awk '{print $2}' | sed '/(s///g' | sed '/')s///g' | sort -u`
nTP=`cat $1 | wc -l`
echo "$nTP peers"
echo "" > s.tmp

for c in $C
do
    Nc=`cat $1 | grep -e "$c" | wc -l`
    echo "$Nc $c" >> s.tmp
done

cat s.tmp | sort -n -r > s2.tmp

cat s2.tmp
for (( i = 1 ; i <= $N ; i++ ))
do
    nP=`cat s2.tmp | head -n $i | tail -n 1 | awk '{print $1}'`
    pP=`octave -q --eval "$nP/$nTP" | awk '{print $3}'`
    C=`cat s2.tmp | head -n $i | tail -n 1 | awk '{print $2}'`
    nA=`cat $1 | grep -e "$C" | grep -v -e "-1" | wc -l`
    pA=`octave -q --eval "$nA/$nP" | awk '{print $3}'`
    rtt=`cat $1 | grep -e "$C" | grep -v -e "-1" | awk '{print $5}'`
    mrtt=`octave -q --eval "mean([$rtt])/1000" | awk '{print $3}'`
    echo "$C $pP $pA $mrtt" >> a.tmp
done
```

Figure A. 9 – Shell script to measure RTT values.

➤ *Create normalized maps*

To create normalized maps shown on this thesis, the following gnuplot program was used:

```
unset key
unset border
unset yzeroaxis
unset xtics
unset ytics
unset ztics
set term X11

set view map
unset hidden3d

plot 'world.dat' using 1:2:(0) with lines lt 8, \
      'data.cor' using 2:3:($1/$5):(($1/$5)/7+0.8) with points pt 7 ps var lt palette

set terminal png nottransparent crop size 800,600
set output "map.png"
reset
```

Figure A. 10 – Gnuplot program to create normalized maps.

The file world.dat contains the correct coordinates to trace the world map.

The data.cor file contains, on the first column, the percentage of identified peers in each country followed by two more columns with corresponding country coordinates and the name of the country. The fifth column includes the percentage of country population.

➤ *Create polar maps*

```
#!/bin/bash
A=`cat America | awk '{print $1}'`
for a in $A
do
    Bmrtt=`cat $1 | grep -e "$a" | grep -v -e "-1" | awk '{print $5}'`
    echo "$Bmrtt" >> america.tmp
done

C=`cat Asia | awk '{print $1}'`
for c in $C
do
    Dmrtt=`cat $1 | grep -e "$c" | grep -v -e "-1" | awk '{print $5}'`
    echo "$Dmrtt" >> asia.tmp
done

E=`cat Africa | awk '{print $1}'`
for e in $E
do
    Fmrtt=`cat $1 | grep -e "$e" | grep -v -e "-1" | awk '{print $5}'`
    echo "$Fmrtt" >> africa.tmp
done

G=`cat Europa | awk '{print $1}'`
for g in $G
do
    Gmrtt=`cat $1 | grep -e "$g" | grep -v -e "-1" | awk '{print $5}'`
    echo "$Gmrtt" >> europa.tmp
done

H=`cat Oceania | awk '{print $1}'`
for h in $H
do
    Imrtt=`cat $1 | grep -e "$h" | grep -v -e "-1" | awk '{print $5}'`
    echo "$Imrtt" >> oceania.tmp
done
```

Figure A. 11 – Shell script for treating results in order to create polar maps.

```
reset
unset key
set multiplot
set multiplot layout 1,2
set title "RTT distribution"
set polar
set clip points
unset border
set label "America" at 1500, 1600
set label "Asia" at -1417, 1690
set label "Africa" at -1990, 913
set label "Europa" at -1600, -1700
set label "Oceania" at 1600, -1700
set xlabel "milliseconds"
set ylabel "milliseconds"
set yrange [-2000:2000]
set xrange [-2000:2000]
set grid polar 2*pi/5
set xtics 400
set ytics 400
plot 'am.tmp' using 1:2 with points ps 0.5 pt 7 linecolor rgb "black"
set label "America" at 1000, 1200
set label "Asia" at -1100, 1200
set label "Africa" at -1300, 700
set label "Europa" at -1200, -1200
set label "Oceania" at 1000, -1200
set yrange [-1400:1400]
set xrange [-1400:1400]
set xtics 200
set ytics 200
plot 'am.tmp' using 1:2 with points ps 0.5 pt 7 linecolor rgb "black"
```

Figure A. 12 – Gnuplot program to create polar maps.

The RTT.tmp file was created to divide peers by continent in order to obtain the wanted polar maps shown on Figures 5.1 to 5.12. Therefore, the first column contains random values between 0 and 2π , divided into 5 sections corresponding to each continent. The second column of this RTT.tmp file contains the RTT value of each identified peer. With peers distributed by continent, the RTT of each peer will be allocated to the correspondent position on the polar map that is divided into 5 sections corresponding to 5 continents.

➤ For 3-Dimensional plots

```
reset
unset key
set terminal X11
unset hidden3d
set pm3d
set palette defined (0 "green", 6 "yellow", 12 "red")
set ticslevel 0
set grid
set grid z
set view 80, 15
unset mytics
set ytics
set xlabel "RTT (milliseconds)"
set ylabel "Date\nTime"
set ylabel 3
set xlabel "Probability"
set xlabel 10
set timefmt "%d/%m/%y %H:%M"
set ydata time
set ytics ("11/07/08 11:00", "11/07/08 13:00", "11/07/08 17:00", "11/07/08 21:00", "12/07/08 01:00",
"12/07/08 04:00")
set format y "%d/%m/%y %H:%M"
set timefmt "%d/%m/%y %H:%M"
```

```
set ztics 2
splot '3D.dat' using 3:1:4 with points ps 0.3 pt 2 linecolor rgb "blue"
set terminal png nottransparent crop size 1000,800
set output "3D.png"
reset
```

Figure A. 13 – Gnuplot program to create 3-Dimensional plots.