# Architecture of Participation:
## The realization of the Semantic Web, and Internet OS

by

## Shelley Lau

Submitted to the System Design and Management Program
in Partial Fulfillment of the Requirements for the Degree of

## Master of Science in Engineering and Management

at the
Massachusetts Institute of Technology
June 2007
⌈Feb 2008⌉
© 2007 Shelley Lau
All rights reserved

The author hereby grants to MIT permission to reproduce and to distribute publicly paper
and electronic copies of this thesis document in whole or in part.


Signature of Author_____

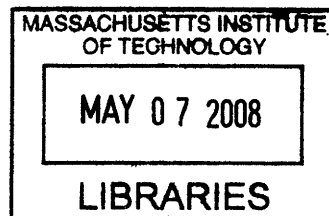Shelley Lau
System Design and Management Program
June 2007


Certified by_____

Prof. Michael Cusumano
Thesis Supervisor
Sloan Management Review Distinguished Professor, MIT Sloan School of Management


Certified by_____

Patrick Hale
Director
System Design and Management Program

TABLE OF CONTENTS

## List of Figures

## List of Tables

## Abstract

The Internet and World Wide Web (WWW) is becoming an integral part of our daily life and touching every part of the society around the world including both well-developed and developing countries. The simple technology and genuine intention of the original WWW, which is to help researchers share and exchange information and data across incompatible platforms and systems, have evolved into something larger and beyond what one could conceive. While WWW has reached the critical mass, many limitations are uncovered. To address the limitations, the development of its extension, the Semantic Web, has been underway for more than five years by the inventor of WWW, Tim Berners-Lee, and the technical community. Yet, no significant impact has been made. Its awareness by the public is surprisingly and unfortunately low. This thesis will review the development effort of the Semantic Web, examine its progress which appears lagging compared to WWW, and propose a promising business model to accelerate its adoption path.

Thesis Supervisor: Michael Cusumano
Title: Professor, MIT Sloan School of Management

## Acknowledgements

First of all, I would like to thank my thesis advisor, Prof. Michael Cusumano, for his patience, guidance and the freedom he provided to bring this thesis to fruition. I greatly appreciate his interest in my research topic and his encouragement throughout the process.

I would also like to thank Pat Hale, the director of the SDM program for his advice, considerations, and flexibility throughout my studies at MIT. The enjoyable and fruitful time I have had at MIT, Sloan, and the SDM program would not have been possible without the continuous support of the SDM staff the SDM '06 cohort. I am very thankful for the friendship and to be part of the program.

The successful completion of this thesis is contributed greatly by all the individuals whom I interviewed with during my research. Many thanks for their time, insights, advice and knowledge.

My sincere appreciation goes to my family and friends. Those whom I met at MIT over the past years have provided tremendous support and have stimulated my intellectual curiosity. They keep me going forward.

Finally, and most importantly, I would like to dedicate this thesis and my achievements to my mom. She was and will continue to be my inspiration and strength for many years to come. Thanks her for everything.

# Chapter 1: Introduction

Although the World Wide Web (WWW) has revolutionized our daily communications and become a critical channel and infrastructure for business, its underlying key technology is rather simple, a crucial characteristic for its quick adoption. Its continuous development and innovation are achieved collaboratively and collectively among corporate and public communities.

As the current Web reaches its limitations, its extension, the Semantic Web, is emerging. The Semantic Web, the ultimate dream and vision of the WWW inventor, Sir Tim Berners-Lee, is to enable the machine-comprehensible web, and true machine-to-machine and machine-to-people "collaborations" beyond the capabilities of the current web.

With the prior success of and experience in WWW, people who follow the status of the Semantic Web have high expectations in its adoption and its commercial performance. Even though the Web took a decade long to mature, it finally became established in the market. Unfortunately, for the Semantic Web, it seems to be in a stalling state or had achieved little progress after more than five years of development. The existing Web thrives on the openness and simplicity of its underlying technology; and the richness and diversity of data and applications. The adoption and development of HTML, one of the key technologies of WWW, has helped the Web prosper and expand in scope. The participation of many decentralized, independent parties was the main driver for the widespread adoption which could not have been achieved if it was driven by one or a limited number of centralized organizations or communities.

To help the Semantic Web move forward and reach the critical mass, this thesis takes a holistic view from diverse standpoints and strives to develop a balanced strategy for standardization, innovation and commercialization. It begins with an overview of *the architecture of participation* which has been materialized in software development and in latest media trend and movement. It examines how the social and cultural phenomenon would help accelerate the development of the Semantic Web. It then presents an evolutionary process starting with the aggregator applications and technology, an intermediate step, before reaching a wide adoption of the Internet OS, which would be enabled by the Semantic Web.

Finally, a proposed solution and architecture is presented along with its technical and commercial potentials. It builds on the premise that the proliferation of WWW and the Semantic Web is sustained by value-added data and links, rather than the standards and technology itself. The solution would facilitate collection and publication of useful metadata using an interdisciplinary framework which leverages various leading-edge technologies and open participation. The benefits of an open participation have been seen in the success of Open Source development and the new media of "wiki" phenomenon. The proposed solution will be illustrated through a thorough business case analysis for the prospective solution providers.

## Chapter 2: Semantic Web

2.1: Background

Semantic Web, coined as Tomorrow's Web, is an extension of the current World Wide Web which has changed the livelihood of communications and information exchange in the past decade. Not only would this new technology revolutionize the future of the online world, it is also the vision of the Web inventor, Sir Tim-Berners-Lee. Currently the Web only provides people-to-people communications and an information sharing medium. His vision is to enable complete machine-to-machine and machine-to-people "collaborations" by giving meaning (semantics) to the content of online documents. Ultimately, machines would become capable of analyzing all the data on the Web including the content, links and transactions between people and computers[1].

Meanwhile, Service Oriented Architecture (SOA) is a new technology architectural framework to provide visibility, access, and interoperability amongst disparate enterprise systems. SOA is built on a stack of technologies including Process Composition, Messaging Infrastructure and Metadata Registry and so on. Semantic Web aims at addressing the Metadata layer.

2.2: Problem with Today's Web

The beauty of the Web lies in its simple elements of protocols which enable communications among computers. The elements include the Universal Resource Identifiers (URIs), the Hypertext Transfer Protocol (HTTP) and the Hypertext Markup Language (HTML). Its simplicity puzzled people initially. There is no central computer "controlling" the Web, no single network on which these protocols work, not even an organization anywhere that "runs" the Web. The Web is not a physical "thing" that existed in a certain "place." It is a "space" in which information could exist. The proliferation of the current Web framework has created a universal information space and channel for people to communicate with each other through the use of computers on the networks.

However, it provides little help in analyzing the context and meaning of the content. It serves very well if you know where to get or send the information and how to get there. Such limitation has created opportunities and high demand for search engines which is to help users find the relevant information online in the right place and at the right time. However, regardless of how good a search engine is, it is primarily based on word occurrence and does not have any understanding of the content. Without knowledge of the context, search engines would not be able to return fully accurate results.

An intermediate solution has been developed called automated brokerage services (e.g., matching a buyers and sellers in the search) to address part of the problem. It executes a product-related search request in a collection of online catalogues on users' behalf and extracts product data, price and other information. Afterward, it synthesizes the results before presenting them to users. The technique applied for

information extraction is called screen scraping. It locates useful information on different websites based on manually-predefined criteria. However, this method is mainly estimation and is neither accurate nor adaptive to future changes. For example, if the price information was moved to a different location in a page or deleted, the automated broker would be confused. Furthermore, it is labor intensive to handle the pre-processing of website in a non-systematic manner.

## 2.3: Solution

As Semantic Web comes into being, the day-to-day mechanisms of trade, bureaucracy and our daily lives will be handled by machines talking to machines, leaving humans to provide the inspiration and intuition. By that, it should not be mistaken for the greater challenge of Artificial Intelligence (AI). The idea behind Semantic Web is to integrate the disparate, incompatible network of machines more seamlessly and to eliminate manual tasks. It sets out to transform the Internet and automate analysis of the Web by giving meaning in a form that machine can understand.

Ultimately, the Web would be much more powerful and useful when it can understand the common sense in human languages and answer some basic questions such as "Did any baseball teams play yesterday in a place where the temperature was 22 degrees?" With data stored in machine-understandable format and processed with mathematical reasoning logic, the future search engine together with a semantic agent would be able to deduce the correct response and filter out irrelevant ones. The Semantic Web enables machines to process, transform, assemble and even act on the data in useful ways.

## 2.4: How It Works

Technically, Semantic Web refers to how the bits of data are related to each other. Like the relational databases which consist of columns of entries of information, the relationships between columns indeed are the semantics or the meaning of the data (i.e. the concepts that the data represent within a particular context, and the relationships between those concepts.) Unlike relational databases, content on Web pages are freeform and HTML only represents the presentation and formatting of the pages.

In the Semantic Web, an agent does not have artificial intelligence but it relies on structured sets of information and inference rules that allow it to "understand" the relationship between different data resources. The computer does not really understand information the way a human can, but it has enough information to make logical connections and decisions.

In order to provide similar structure, a new descriptive language is needed. The W3C consortium[2] is developing a language called Resource Description Framework (RDF), RDF Schema (RDFS) and Web Ontology Language (OWL) which are based on eXtensible Markup Language (XML)[3]. Like HTML which is used to represent hypertext, RDF provides a framework to represent which bits of Web content are

data and how to find the meaning of the data, i.e. metadata – information about information.

In a sense, RDF is a model that teaches the Web English and the basic grammar. It is a framework for describing the resources existing on the Web and eliminating ambiguity. It builds on the existing XML and URI technologies. Like the English sentence, an RDF statement consists of three parts called triples: subject, predicate, and object which correspond to a resource, a property and a property value respectively and can be written with XML tags. For example, "The name of the secret agent is Niki Devgood."[4]

The graphical representation of the above sentence is as follows:

Figure 1: RDF Triples

Building on this simple structure, N-triples can be created for a more complicated and descriptive sentence to include other attributes such as an email address. An elaboration to the previous example is as follows:

Figure 2: RDF n-Triples

The idea is that by putting the information into a structure this way, the framework enables machines to make logical assertions based on associations between subjects and objects. Since RDF uses URI to identify resources, each resource is tied to unique definition available on the Web. However, RDF only provides structure; the underlying meaning is defined by using RDFS and OWL.

In a nutshell, RDFS is a simple vocabulary language for expressing the relationships between resources. Building upon RFDS is OWL, which is a much richer, more expressive vocabulary for defining Semantic Web ontologies which define the hierarchies and relationships between different resources. Ontologies consist of taxonomy and a set of inference rules from which machines can make logical conclusions. Basically, the taxonomy is a system of classification, for example,

classifying plants and animals that groups resources into classes and sub-classes based on their relationships and shared properties.

All the detailed relationship information defined in OWL ontology allows applications to make logical deductions. For instance, given an ontology that Goose is a type of Dark Meat Fowl, and Dark Meat Fowl is a subset of the class Fowl, which is a subset of the class Edible Thing, a Semantic Web agent could infer that Goose is an Edible Thing.

2.5: Development or Evolution

Like many new technologies especially one that requires cooperation in the standardization process among different industries, corporations and development community, it takes many cycles of development and evolution. The goal of the Web is to maintain interoperability and evolvability and to avoid forcing anything building from the ground up.

In speaking with Prof. David Karger[5], the principal investigator of a research group, Haystack[6], at the Computer Science and Artificial Intelligence Lab at MIT, there is a number of development projects underway including a joint project called Simile[7]. A number of tools such as RDFizers and Semantic Bank are being developed to help with the migration of some existing technologies to the future one. The RDFizer helps transform existing data into an RDF representation whereas Semantic Bank is accumulating a collection of semantic data in the meantime.

While the real impact of Semantic Web is yet to be seen, there are many current commercial projects or initiatives underway that are using RDF including MIT's DSpace (a digital research materials repository), Mozilla browser, RSS 1.0, Stanford's TAP project and many others.

2.6: Commercial Potential

Semantic Web is not only about solving the search problem but to tackle the grand challenge of integrating data in personal computers, devices and the enterprise backend storage in order to simplify the daily life of data management and processing. Just imagine booking a business trip which would require access to information of the person, credit card, calendar, personal travel preferences, corporate reimbursement restriction, flight schedule, itinerary, hotel, etc. If all this scattered information can be retrieved, analyzed and acted on automatically with minimal manual procedure, it would greatly simplify the tedious work.

A commonly used and well-known scenario is the Travel Butler (TB). It does not exist yet but services like it would be appealing to companies like Orbitz and AT&T[8]:

> Let's say it is 4 p.m. TB knows you have a flight scheduled for 6 p.m. because it regularly prowls the Web sites you use for travel and found you booked a ticket on Orbitz. TB can tell by checking your online calendar that you are at a meeting downtown.

The service cross-checks with a map service such as MapQuest to find the route you would have to take to get to the airport. Once it knows that, TB goes out on the network to monitor traffic on your route and finds the streams of data on the Department of Transportation Web site, which monitors road cameras and sensors.

TB might see that accidents have backed up traffic for miles. It sends you a message on your BlackBerry e-mail to urge you to leave now in order to catch the flight on time. In addition, TB shows you an Orbitz listing of later flights.

You decide to go on a later flight, so you click on the one you want. TB rebooks you, sends an e-mail to your spouse and contacts the car service in your destination city to change the time to pick you up.

That is an experience that rises above a particular technology. "People really do not want to buy technology," says Lisa Hook, head of America Online's broadband unit. "They do want to buy experiences."

The list of scenarios that could potentially benefit from Semantic Web technologies as they continue to evolve is limited only by the imagination. Other possibilities include everything from crime investigation, scientific research, and literary analysis to shopping, finding long-lost friends and vacation planning once computers can find, present, and act on data in a meaningful way.

2.7: Technology Comparisons

In light of the recent launch of the SOA initiative at IBM and the plan of leading the market of Business Integration and Transformation, the Semantic Web would be the choice for data modeling and formatting standard. Instead of competing with or superseding the SOA technology, Semantic Web complements the SOA and aligns with the business strategy.

The ability of Semantic Web technologies to access and process enterprise data in relational databases together with data of other sources (Web sites, other databases, XML documents and systems) will help grow the amount of useful data faster than ever. In addition, relational databases already include a great deal of semantic information.

Data integration applications offer the potential for connecting disparate sources, but they require one-to-one mappings between elements in each different data repository. The Semantic Web, however, allows a machine to connect to any other machines, and process data efficiently based on built-in, universally available semantic information that describes each resource. In effect, the Semantic Web will allow us to access all the information as one huge database.

Furthermore, the current implementation of SOA is based on XML; however it leaves document interpretation up to consumers. Because semantics can be formally modeled in Semantic Web ontology, the usage and meaning of data are explicitly captured in a machine-interpretable format. It also allows machines to

automatically discover relevant content sources based on business concepts, and enables the framework to expose and reuse the interpretation rules coded in currently existing systems.

From a purely technological standpoint, the Semantic Web is progressively making some impacts and influences on corporations and technical community since it is being actively pursued by the W3C consortium. W3C has successfully developed such technologies as XML and Web Services. In addition, it would also be cost effective and result in potential savings from reusing and leveraging on this open data modeling framework.

In reality, however, as discussed in the Chapter 4, not only the adoption of Semantic Web has been slow, its awareness is surprisingly low even among technical professionals and a group of corporate IT executives in general. In addition to technical hurdles, the Semantic Web faces many execution difficulties and challenges. As shared by some researchers and technology experts, a general comment is that the Semantic Web might not succeed in its current form but will continue to evolve. It would be realized differently, in a way that is similar to the realization of Artificial Intelligence which has not lived up to its original vision on its own but has been embedded in a bigger system such as the Google's search engine and made impacts indirectly.

## Chapter 3: The Future: Internet OS

Internet started out as a free-form medium for a web of information to link with each other. It is to enable efficient data references, retrieval and connectivity. Initially, many individual Web sites were created and they fell into one of the two types: hobbyist site or marketing tool for commercial companies. Some of the static corporate "homepages" have grown into functional e-commerce sites for online transactions. After a few years, some Internet-based commercial sites such as eBay, Amazon and others were emerging.

In a sense, the Internet is a big open source project with everyone contributing content and information to the Web. The vast, diverse, and easy access of information makes the Internet an interesting communication medium but, at the same time, causes the problem of information overflow. As a result, technologies and services like the Internet search tool from Google are not only useful but highly necessary. For commercial Web sites, being ranked high in Web search results has become a critical part of online marketing plan and strategy at the big corporations. That has created a tremendous opportunities for the small startups such as Google and Yahoo to emerge and become "giants" in the Internet era. Unfortunately, such technologies do not provide highly accurate and relevant results. They might not live up to the expectation of average users. The challenges are two-fold: technical and political.

On the technical front, the progress of Artificial Intelligence or machine learning algorithm has staggered. They are not yet capable of understanding and processing free-form natural languages. In addition, the amount of information on the Web is too scattered and fragmented. Some automated software engines or Web crawlers have been trying to scout all the Web sites and content on the entire Web. Unfortunately, there are many still have not been exposed to the search 'bots' (computer programs). Therefore, useful and product information could be buried in the quiet corner of the Web.

In addition, there are many Web sites which serve only one function or one type of information, for example, weather, stock prices, email, and map. These Web sites are like software gadgets. However, none provides, intentionally or not, integrated set of tools and information, like an Operating System does. To the end-users, data needs to be processed to become useful information. To the corporate enterprises, although there are many data stored in structured formats in their IT systems, the formats are mostly proprietary. Typically, the issues are due to political and business reasons since many valuable data such as customer, product and technologies provide competitive advantage to the companies who own them. Most companies protect and preserve their data in certain formats. However, as mentioned previously in Chapter 2, incompatible data formats prevent cross-sales and other business opportunities.

In light of that, an evolution is happening on the Internet. It is transitioning from a model of stand-alone company-specific Web sites to theme or category-based Web portals or aggregators. To further integrate data and functionality from various

enterprises, the development and growing adoption of Web Services is underway. Finally, as we are moving toward Internet-centric and highly connected environment, the boundaries between your local desktop, the internet browser, cell phones, consumer devices like the iPod are blurring[9]. If we could break the barriers of IT system boundaries behind the firewall of small and large enterprises, all these disparate systems would converge into an Internet Operating System (Internet OS) at large. This is the ideal cross-application gateway and cross-corporation integration in dynamic manner building upon open architecture, API, and data interoperability. The cross-application gateway should be incorporated into the standard part of all software, rather than "hacked" on after the fact. The Semantic Web which strives to standardize data structure and adds meaning to the data will play a fundamental role in the development of the future operating system.

## 3.1: First, Aggregator

In the publishing business, publishers have long been in the business of aggregator: researching, collecting, and synthesizing the information and putting them together in the context of the topics of the publications such as magazine and book. Such a middleman has been critical in bringing the producers and consumers together. In the rise of the Internet, such role has become even more and more important. Google, Yahoo, and other Web-based consumer services are serving the consumers with the desired content. Although automated search is an integral part of the business at most of the big search engine companies such as Yahoo, they reposition themselves to lead the online media publishing and strive to be the large digital media and entertainment company. As aforementioned, the current search technologies do not entirely satisfy human needs. As a result, Yahoo employs an army of administrators and programmers continually rebuilding the product. Dynamic content is not just automatically generated; it is often hand-tailored, typically using an array of quick and dirty scripting tools. According to Jeffrey Friedl, the author of the book Mastering Regular Expressions and a full-time Perl programmer at Yahoo, Yahoo does not create content. They aggregate it[10]. They have feeds from thousands of sources, each with its own format. So, they do massive amounts of "feed processing" to clean up the content and find out where to put it on Yahoo. For example, to link appropriate news stories to tickers at quotes.yahoo.com, Friedl needed to write a "name recognition" program to search for more than 15,000 company names. It was made possible by a scripting language, Perl which has he ability to analyze free-form text with powerful regular expressions. Essentially, this aligns with the magazine publishing without editorial staff or like a "retailer of content" with services free of charge.

On the other hand, the globalization of information on the Internet presents significant opportunities and challenges at the same time. The prospect of combining information from diverse sources for superior decision making is plagued by the challenge of semantic heterogeneity, as data sources often adopt different conventions and interpretations when there is no coordination. For example, a European site lists airfare in Euros, while the US-based one lists them in Dollars; the airfare in one contains all the taxes and fees, while in another it does not. Furthermore, users of these Web sites have their own assumptions about what the data means, which sometimes does not correspond to the intention of the actual Web

sites. This is essentially what the Semantic Web is trying to address by developing an ontology as a standard data model for a domain of interest, and then defining the correspondences between the data sources and this common model to eliminate their semantic heterogeneity. However, for the Semantic Web to work, it requires ontology development to standardize the exact meaning and representation of ontological terms. Such requirement turns ontology development and adoption into a standardization process which is a notoriously arduous and a lengthy process which requires cooperation among different parties including big corporations. The adoption of the Semantic Web is facing a number of challenges which will be discussed more in detail in the following chapters.

To address that, there is another school of research effort at the MIT Sloan School of Management under the supervision of Dr. Stuart Madnick. In his group, they primarily take the automatic and pragmatic approach by the technique called "screen scraping"[11]. By their definition, Web aggregators are entities that collect information from a wide range of Web sources, with or without prior arrangements. They add value by post-aggregation services using the new Web-based extraction tools which allow them to easily and *transparently* gather information from multiple sources with or without permission or knowledge of the underlying data sources. In the process, the aggregator resolves the semantic or contextual differences in the information, such as different currencies or price information which include or exclude shipping charges. Under this definition, search engines such as Google and Lycos, and personalized Web portals, such as MyYahoo are not aggregators since they provide little contextual transparency or analysis. Their research is called Context Interchange and is taking the problem one step further by automatically determining semantics or the meaning of data. They have developed the mediation technologies to perform semantics extraction. They define two aggregator types and sources: comparison (consumer education shopbot services which compare different products) and relationship (one that helps users manage relationship more effectively). For comparison aggregator, the intent is to enable information comparisons across the Internet (such as book prices, bank balances, shipping rates, and intelligence information). Some examples are Maxmiles, Priceline, Orbitz, etc. On the other hand, relationship aggregator consolidates all one's frequent flyer or financial accounts information, e.g., Yodlee, and CashEdge. From inter-operational and organizational point of view, the relationship aggregator consolidates all information about each customer from the company's separately maintained Web sites across function (accounting, and service) and geography (domestic and international).

In the case of Yahoo or COIN research group, their approach is still very much ad-hoc or non-systematic. They do not conform to any standard format for interoperability or extensibility for future reuse. In light of that, Web Services has emerged to address the problem.

3.2: Then, Web Services

In the process of making Web sites more useful to people, businesses find that more and more applications and data need to be connected behind the scenes. Web sites can operate as stand-alone platforms as long as they offer limited information, but

they would quickly run into problems when they try to bring together a broader range of information from different sources for users to buy products. Web Services is an emerging technology architecture based on a set of common standards and protocols which enable connecting applications and data directly with each other, automating connections without human intervention. The architecture is designed to ensure that applications and data can be accessed by authorized entities, regardless of location or underlying technology platforms. The key innovation is to enable loosely coupled resources to have dynamical access to each other across multiple entities and diverse technology platforms. Loosely coupling means that the connections can be established without being tailored to the specific functionality embedded in the applications, the issues of which will be discussed more in detail in the following. Web Services fundamentally alters the definition of the edge of the firm[12]. Dell integrated its entire supply chain into the company's IT infrastructure and automated the coordination with its supply-chain partners using the Web Services which resulted in reduced component inventories from 26 to 30 hours of production to three to five hours, a reduction of more than 80 percent[13]. Furthermore, Amazon provides a Web Services Developer Kit at no cost and enables smaller online retailers with the power to access Amazon's transactional technology and catalog platform in only a few hours. Not only does Amazon serve as a large e-retailer, but also as an e-retailing platform on which enables thousands of new e-retailers to build complementary e-retailing services.

The Web Services architectures are built on a foundation of standards and protocols. Some of these standards and protocols are well defined but still evolving. It is important to note that although all are being defined and maintained by broadly supported public standard bodies, in reality, a few big companies such as IBM and Microsoft have come up with variations which is tailored to and optimized for their own technologies platform. Figure 3 below shows the Web Services Architecture and its components.
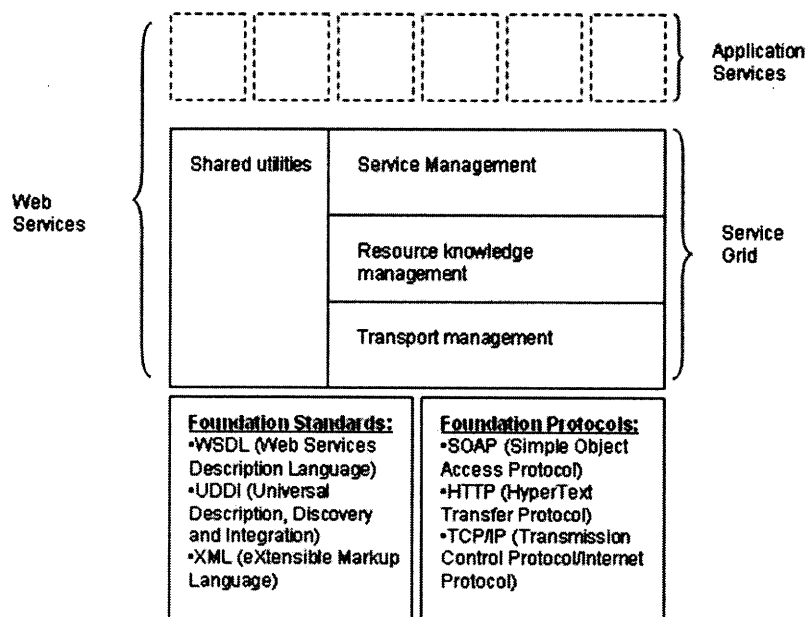


Figure 3: Web Services Architecture: Standards and Protocols[14]

In addition, many of the product distributors setting Web sites receive information on product specs and pricing from hundreds, if not thousands, of product vendors. Organizing this information so that it could be easily searched and compared by customers represents a serious technical challenge. Product vendors supply the information in their own unique formats. Somehow, these diverse formats have to be converted or translated into a single, uniform format. Conventional technology solutions like Electronic Data Interchange (EDI) networks or early generations of Enterprise Application Integration (EAI) technology focus on establishing point-to-point connections. These technologies are poorly suited to the complex and shifting web of relationships emerging around Internet Web sites. In practice, as aforementioned, many opt for a no-tech solution by employing an army of people manually generating information out of one set of applications or databases and manually converting it into another format, and then inputting the reformatted information into other applications or databases. Not only is this labor-intensive process very expensive, it is also time-consuming and prone to errors. It provides the experience of integration to end-users, but it is neither sustainable nor scalable. The point-to-point approach causes the number of connections growing exponentially for 3 or more applications. Furthermore, in the ever changing business world, managers have little foresight into the long-term network of relationships. It becomes unbearable if the nature and number of connections are being continuously redefined.

The original intent of the World Wide Web (WWW) is to help researchers share and exchange knowledge across incompatible information technologies. While Hypertext Markup Language (HTML), the foundation of WWW is enormously useful when a person is accessing material on a Web site, but it becomes less relevant when computers are accessing information from other computers. Essentially, the basis of the Internet is to help applications communicate directly with each other. The problems are two-fold. First, there is no easy way to enable the enterprise's procurement application to search a number of supplier electronic catalogs to find out if a product is available in a certain size. Second, the heterogeneous data formats have posed numerous technical challenges. All that lead to the development of a new standard called eXtensible Markup Language (XML). It helps connect diverse applications and data of all kinds, targeting the technology foundations of all businesses. It provides the standards which are the common formats for representing entities like documents, products, and customers, and actions like "ship" or "confirm", similar to how HTML describes the graphics and text of a document should be displayed when they are delivered to a user. However, HTML is static in a sense, while XML is as extensible as its name implies. XML provides a way for "tags" to be created in documents or messages so that other applications can quickly locate the information they need. For example, the managers of a supplier's catalogs could implement an XML tag designating where product size information is located. By providing a common format for representing this information, XML makes it easier to automate connections across applications. Where a human had to search for information, the information can now be automatically accessed and delivered. XML provides much of the same functionality that the formats defined by various EDI vendors in earlier years provided. The key difference is that XML

formats are broadly adopted, whereas EDI formats are proprietary to specific vendors. XML also serves as the foundation for other, more specialized standards.

Every application ought to expose some version of its data as an XML Web Services feed via some well-defined and standard access mechanism. It should be as simple as the naming scheme that fuels the success of the early Web, e.g., www.amazon.com, which made the web addresses eminently guessable. If we could do it for Web sites, the potential of doing for local applications is even bigger. It is becoming more important to have all applications, whether local or remote, be set up for two-way interactions. That is, they can be either a source or receiver of online data. Since data is coming from multiple sources, we need to understand who owns what, and come up with mechanisms that protect the legitimate rights of individuals and businesses to their own data, while creating the liquidity and free movement of data that will fuel the next great revolution in computer functionality.

3.3: Finally, Internet OS

Although Web Services is striving to eliminate some of the scalability, ad-hoc development and labor- and cost-intensive work with the approach of aggregators, aggregators are not an absolute replacement of Web Services. The two approaches co-exist. The idea of aggregators and Web Services are converging into the development of the future internet-centric Operating System (OS) for which the existing OS will be part of this larger environment[15].

Hackers are building unauthorized interfaces via Web spidering and screen-scraping. The Web-facing databases are treated as software components to be recombined in new ways, and with new interfaces. At O'Reilly, they have built a whole custom set of market research interfaces to Amazon by spidering every computer book on Amazon every three hours, collecting vital statistics like author, publisher, page count, price, sales rank, and number and nature of customer reviews. This early stage Web Services tend to be inefficient, brute-force spiders. Recently, Amazon, Google, eBay and other sites have started to offer XML-based APIs realizing that it is in their own interest to enable users to become their developers by providing the data they request and letting them to re-use it in creative new ways.

Some of which are the emerging "mashup" services of which Websites are integrating content from various sources to provide more comprehensive services or experience. It is analogous to systems of modular design with individual components integrated through public interfaces in the new medium of Internet. For example, Trulia[16] is a real estate market trend or house searching service using Google map to display location information.

One of the challenges is not only to integrate the different applications into one that is accessible through the Internet browser, but to also enable machine "agents" to work with each other dynamically through open APIs. Data-driven web sites need to be seen as software components. Google's API is a good step in the right direction, but it's hoped that all data-driven Web sites would serve as program-callable software components accessible by developers. One good example is Apple's

integration of Yahoo! Finance and MapQuest right into Sherlock (a file and Web searching tool made by Apple Inc.), rather than just through browser. What is particularly interesting is the non-controlling integration into the Apple's OS, instead of replacing or competing with the work of existing players like what Microsoft did with Netscape, Real Networks, and AOL. If Microsoft's Windows is "just a bag of drivers" as referred by Marc Andreesen[17], the challenge of the Internet OS is to have a set of right drivers and seamless integration, much like the plug-n-play in Windows. As the Internet is primarily data-driven, the Web databases are the part of what we need standard "drivers" for.

The challenge is that the current development of the Web Services focuses primarily on modularization of service layers as needed by the industry but ignores the expressiveness of the service descriptions and data addressed by Semantic Web. Although the Web Services provide many benefits, key concerns of the initiatives are in short-term applicability and scalability[18]. Typically, the procedure of setting up and consuming Web Services requires that the human requestors (users) searching for the services in the registry to which the providers must register and expose their service and manually "bind" with users' applications. However, the service description is based upon standard taxonomy and the data exchange is syntactic-based which are essentially the issues being faced in the current Web. In the world of the Semantic Web, services are published along with semantic descriptions to enable crawlers to find them. The services should be available without registration and searchable automatically by a personalized machine agent who will take over the role of a service requestor.

Traditionally, Operating System (OS) is a special software program that manages application software and is build for specific variety of hardware in the computer systems including the Central Processor Unit, and other input and output peripherals. On the other hand, the Internet Operating System (Internet OS) is a defined high level set of functions and Internet-based applications or services which is operated in uniform manner independent of the underlying hardware systems and provides a distributed computing environment. As discussed in this chapter, the Semantic Web plays a central role in data and functions connectivity. Essentially, the Semantic Web is the key enabling technology for the data-driven architecture and the implementation of the Internet OS. A more detailed technology diagram in Chapter 6 illustrates that the Internet OS will be built on the foundation of the Semantic Web which serves as the data pipeline for the information flow among applications.

## Chapter 4: Problem with the Progress of the Semantic Web

Development of a new technology and its success often follow or resemble the path of entrepreneurship in a new venture. Many factors affect the results. The approach used for the development and market penetration is usually influenced by the background and discipline of the principal investigators or founders of the venture.

One might define success differently from an engineering perspective than from a business perspective. However, it is useful to look at it holistically and find a balance instead of creating a great divide between the two. In order to analyze the progress and judge the success of the Semantic Web, it is important to define what success means in this context. Although this thesis maintains a balanced and objective view and evaluation based on engineering excellence, commercialization and practicality, success is considered in terms of adoption rate, improvements and its impact (economic, social and commercial). For Google, it is clear that they have succeeded in the marketplace and are able to monetize their research and technology. More importantly, they achieved what they set out to do: make information more accessible. This is becoming more critical than ever before due to information overflow and exploding growth of the Internet. Not only could technologies commercialization bring in financial reward, it could also make a great impact like what the Google founders, Larry Page and Sergey Brin have achieved – a huge financial gain along with benefits to the general public. Otherwise, their research projects would have been put on the shelf once they graduated.

While the Semantic Web is the process of leaving research labs and entering commercial development and production environments, it is facing numerous hurdles to adoption. In particular, a number of areas will be addressed:
- The hurdles to adoption of the current Semantic Web (Chapter 4)
- The potential of new trends to adoption of the Semantic Web (Chapter 5)
- The way in which the Semantic Web enhances the existing Web infrastructure and what may serve as a model and roadmap for adoption. (Chapter 6 & 7)

### 4.1: Why the current form of the Semantic Web might fail

1) Big Bang mentality.

   "Hard core" scientists and researchers tend to think invention and innovation are a matter of the Big Bang phenomenon while many of the innovations are, in fact, incremental but "disruptive" in nature.

   In general, there are two types of inventions: launch and growth as coined by Don Clausing[19]. Although launch inventions, which are like Big Bang phenomena, do exist and have successes in the marketplace like the Chester Carlson's invention of Xerography which became a multi-million dollar business, they rarely exist and the success rate is unpredictable and low. Often, the success factor lies in the effective architectural innovations and incremental improvements rather than innovation at the level of basic components[20]. For

example, the newest Ford car is vastly superior to the Model T, an automobile first produced by Henry Ford in 1908. However, since the invention, the new models or generations of cars are the result of myriad smaller innovations. For example, Ford and other manufacturers in the industry have produced big improvements in disk brakes, fuel injection, radio, and air conditioning. These are growth innovations.

As conceived in the Disruptive Innovation Theory by Clay Christensen, the idea of disruption is a matter of finding simpler, more convenient and relevant products that sell for less money and appeal to a new or unattractive customer set (low margin profit). [21] On the other hand, talented and creative technologists and scientists tend to pursue the next big, radical thing which might be technologically challenging and a breakthrough; it is characterized as sustaining innovations and causes a head-to-head competition with the incumbent.[22]

<u>Analysis</u>

Although the Semantic Web has been advocating as the extension of the existing Web and infrastructure, the mandate appears to be pushing the Big Bang phenomenon. Data format compliance and data uniformity are the critical aspects for the Semantic Web to become successful. To a great extent, it is similar to requiring everyone to speak the same language in order to do business in the global market. Even having a universal measurement metric system or universal currency is hard to achieve, the challenge to have universal spoken language across the globe would be insurmountable. If globalization demands everyone to change some fundamental behaviors and underlying culture, it would be nearly impossible, at least in the near term.

In the case of traffic congestion, many would complain that the existing system is neither efficient nor effective in addressing the issues. We might dream of solving the problem by tearing down the whole highway and road infrastructure and rebuilding it. Unfortunately, that is not possible. Similarly, in WWW, it's an established system which has become an integral part of many people's life. The "information highway" has been built, therefore transforming it to achieve universal data format would be very difficult in the foreseeable future. In order to move forward, the Semantic Web would need a more flexible framework.

2) One-size-fits-all syndrome

Misapplying technique or theory can be the result of the one-size-fits-all mindset. One's prior success might lead to failure in another instance because they often consciously or subconsciously apply previous theory or practice to completely different and irrelevant circumstances. The result would be analogous to giving standard drug prescription to patients with different illnesses. The consequences could be fatal[23].

Such a syndrome can be witnessed in traditional marketing paradigm in which companies often segment market by product attributes and customer demography without proper consideration of circumstance-based

categorization[24]. The problem with attribute-based categorization is that it ignores a critical question, "what circumstance would fail?" It mixes up the correlation with causality between attributes and outcomes. For example, a study might show that males in the age group of 30-40 drink beer daily. However, the study does not conclude that the gender and age cause the desire and the need for beer. Predictable strategy for bringing new technology and products to the marketplace requires understanding of the circumstances in which the technologies are adopted. The circumstances lead to the definition of features, functions, and positioning. The circumstance-based analysis and job-to-be-done view has helped Research in Motion (RIM) uncover a non-trivial insight into their positioning and competitors who turn out to be Wall Street Journal, CNN, and Airport News instead of ones in the trivial and parallel category including notebook computer vendors, and Personal Digital Assistance (PDA) market (e.g., Palm Pilot).[25]

<u>Analysis</u>

Although the Semantic Web (SW) possesses similar characteristics as WWW, SW lives in a different time and environment. While the "competitors" of WWW are fax, telephone, library and other communication or knowledge management systems, SW competes with many of the proprietary and established standards and data formats used in existing enterprise systems. Unfortunately or not, the revolutionary success of WWW has set the expectation and precedent path for the Semantic Web. As technology and business environments are changing rapidly, it is hardly possible that the same path would lead to the same success. Although prior experience in innovation and technology commercialization would be tremendously valuable, the same path that led to the success of WWW could not be applied to achieve the same result. Before WWW existed, the challenges were to gain buy-in of the new potentials and possibilities that it could bring. As the Information and Web Technologies become more mature and get to the "harvesting" stage, much functionality is available. Enterprises providing and relying on technologies would have fewer incentives and would be more risk averse to change. Since WWW has already enhanced communications and knowledge sharing, pursuing a similar path and making subtle and marginal enhancement would not provide enough differentiation or compelling reason for adoption.

3) "If you build it and they will come" mentality.

"If you build it and they will come" mentality is often found in projects dominated by engineers, especially in the software world. There is never a lack of good and creative ideas, but thorough customer needs and market analysis are often overlooked. It is a common thinking that good ideas and technology that solve problems will attract end-users and customers. Therefore, by default, there would be no barrier to capture market share since a small percentage of a large user-base would be a good start. Such mentality does not provide a strategic go-to-market plan. It is due partially to inexperience, lack of knowledge and information, deliberate focus on technical side, ignorance or occasionally arrogance or naiveté.

<u>Analysis</u>

Although a lot of Semantic Web-related application development is underway, Tim Berners-Lee tends to give it a lower priority and suggests his colleagues not to think about killer-application ("killer app") for the Semantic Web. The technology will become successful and be adopted when new links among information begin to emerge, he said[26]. Interestingly, however, Mosaic and VisiCalc are the two great killer apps that enable WWW and Apple II respectively to leap into the mainstream markets (early and late majority adoption in the Adoption Life Cycle[27]).

Regardless of how good a product or idea is, having such a mentality is risky since there is no strategic plan for reaching a certain target. The future of a product would become very unpredictable and vulnerable to market uncertainty.

4) Over-design

Despite high quality, overly-designed and engineered products might still fail to gain the adoption or achieve the critical mass since engineering is squarely concerned with design excellence and engineering practice. For example, we can compare UNIX and Windows. Although UNIX is perceived as a greatly designed, implemented and engineered software operating system with better quality; Windows is more widely used. Undoubtedly, UNIX has superior design and implementation and receives great reviews from sophisticated users, but it does not satisfy the need of general end-users and consumers, i.e., simple configuration and ease of use. Even though Windows might not have lived up to some people expectations and engineering standards, it provides a few simple but useful applications which are critical to its success. Simplicity is often overlooked by some exceptional engineers.

<u>Analysis</u>

For the Semantic Web, the burning issue is that the syntax is very complicated and not intuitive. The RDF statement, in the triple structure, is difficult to develop by hand. As RDF is an integral component of the Semantic Web, the key is to keep it simple and readable so that more application developers would adopt it and help the technology spread out more quickly.

As Tim Bray, the founder of OpenText and pioneer in SGML which has been evolved into XML, puts in his blog, the problem with RDF is the syntax[28].

> "Conceptually, nothing could be simpler than RDF. You have Resources, which by definition are identified by URIs. The resources have Properties, which by convention are identified by URIs. The properties have Values, which can be strings or numbers or Resources.
>
> ...Speaking only for myself, I have never actually managed to write down a chunk of RDF/XML correctly, even when I had the triples laid out quite

clearly in my head. Furthermore—once again speaking for myself—I find most existing RDF/XML entirely unreadable. And I think I understand the theory reasonably well."

The quotation above illustrates that if the complexity of RDF syntax challenges an expert, it would impose even more obstacles for the general users.

5) "Pioneer" mentality

The problem is not with being the pioneer. Instead the competitive characteristic of brilliant scientists or engineers might hinder them from developing and executing the most effective new product strategy.

By definition, a pioneer is someone who goes into unexplored territory. Unlike in the Olympic Games, the first company in the marketplace would likely end up losing while the one who comes in second would win. That idea is probably counter-intuitive at first glance because many think that the first one who enters the market would have a chance to win big. Unfortunately, as the record shows, that assumption is actually not true, since those who are successful now are mistakenly perceived as the first ones in the market when their predecessors failed and vanished without notice.

To understand the point, one needs to distinguish the subtle difference between innovation and invention. By convention, both words are used interchangeably and mean new and original creation. But in the business world, innovation is a crucial differentiation factor for a company to gain leadership in the marketplace, while invention alone usually does not produce financial benefits. In an internal communication, Steve Mills, the Senior Vice President and Group Executive of IBM Software Group, stated that innovation is creative application and integration of new and existing technology. It is also a process of generating new ideas to better serve the customer needs but not building things from ground zero. Technological perfection and pioneer usually do not necessarily translate into a profit, a key measurement of success. In essence, innovation is science while invention is art. In the marketplace, innovation would be a strategic component that drives profit while invention is a means to be the first and lead frontier which is important on a different level. But a business blindly striving to be the first one in the marketplace is doomed to failure because the quest for novelty takes precedence.

The search engine business provides a prominent example of a runner-up taking the prize. AltaVista and OpenText probably pioneered the field, but they let Google steal the thunder. The same holds true for the invention of the light bulb; few people know that the first light bulb was invented by Joseph Swan[29], an English physicist and chemist in the 1820s. He was not able to create a viable commercial product[30] as Thomas Edison did.

Analysis

Although Tim Berners-Lee invented WWW, he did not intend to bring it to the market initially, nor did he achieve its current success by himself. It was a compound effect of lead users adoption and the simplified user interface in the creation of Mosaic, the first user-friendly Internet browser. On the other hand, with the Semantic Web, Tim Berners-Lee and his research group at W3C consortium set out a grander mission to address the issues of WWW and create a more seamless and effective collaborative Web environment. His deep knowledge of undesirable behaviors of the Web and his great desire to fix them have driven him to focus on inventing the first and best new standards of the next Web. However, he focuses primarily on the system side of the technologies. As it turns out, the rise of the Web 2.0 has provided the collaborative environment such as the Wikipedia or wiki-sites. Although it does not fully deliver the same capabilities as the Semantic Web envisions, it enables users to publish information and interact with people in a two-way online medium. It focuses more on simple user interface and integration of technologies but less on the technologies itself.

6) Too-grand a vision

A vision is intended to be a high level big picture of a project or someone's ideas on products and future technologies. However, if the picture is too big, vague and futuristic, it would be difficult to devise a step-by-step well-defined and actionable plan. In other words, it would be very difficult to realize a vision which intends to solve a very big problem (e.g. fully intelligent programs and robots) or imposes many interdependencies (e.g. electric cars.)

There is also the "lure of the Horizontal" as termed by Dr. Michael Cusumano[31]. It could be the result of the lack of focus and unclear direction. In addition, it could be due to the appeal of vast market and overestimated market potential. People often think that they could easily capture many or most of the customers. Sometimes, engineers tend to develop one-for-all solution with a good intention to solve everyone's problems. However, most often than not, over-generic products do not provide any specific functions and require too many supplemental tools, or customization.

Analysis

Tim Berners-Lee's success in WWW lies in its initial, manageable focus on creating solutions for easy information exchange among scientists in the community. The subsequent expansion and departure from the initial purpose has been the work and effort of many people. But the initial focus and subsequent relevant applications have helped launch WWW, and set the stage for growth and revolution. However, the Semantic Web seems to be a much greater vision which imposes many interdependencies: existing and future data are required to conform to one standard. Not only does it demand extra costs and technology upgrade or migration, it also demands a tremendous effort in overcoming barriers in business differences and politics before companies would work together toward a common platform.

7) Lack of complementary products and scalable infrastructure.

In addition to economic climate and competitive landscape, the driving force behind many successful technologies requires the right sustainable and scalable infrastructure. Secondly, a good mix of complementary technologies and products are required to create and drive the demand. For example, the advances in Internet bandwidth have recently enabled the consumption of high-quality rich media on the Internet, while it failed and was impractical a few years ago. Many Internet Protocol TV (IPTV) and online video distribution technologies pioneers failed to penetrate the market a few years ago when broadband Internet was not capable of delivering acceptable-quality images in a reasonable time. However, according to Parks Associates, the number of US households that downloaded videos online at least once a month was already 11 million in December 2005[32].

Analysis

Some people purchase an iPod to listen to songs in MP3 format while others switch from CDs to Mp3 because of the coolness of owning and using an iPod. Therefore, iPod has been the driver for the growing adoption of MP3. Unfortunately, there is not an obvious driver for adoption of the Semantic Web. As of now, the Semantic Web is perceived as good to have but not a mandatory technology required by critical operations. It is often a difficult hurdle for any new technologies to penetrate the market if there are not clear complementary applications or real demand drivers. In addition to better performance and economic value, some would adopt new technologies to enjoy the benefits of the associated applications but rarely for the new technologies alone.

8) Bad-timing

A good idea today could have been a threat or bad idea in the past. Napster would have survived today since its model is similar to YouTube and avoided the legal battle over piracy if they came out in the market today instead of a few years ago. While the content distribution model and mechanism created by Napster were perceived as threats by big conglomerates, they have been transformed into great opportunities for YouTube. The key difference is that the competitive landscape has changed tremendously over the past years. The big media content owners have realized that content ownership and distribution capabilities through traditional channels do not provide them greater competitive advantage nowadays. The potential of online distribution channel has convinced the media companies to open up some of their content and distribute them in this new channel.

Analysis

The idea of the Semantic Web would be more appealing to the public if Google did not alleviate the search and information overflow problem so well. They have successfully integrated their products into everyone's life seamlessly. Although their technologies might not be perfect – 100% accuracy, the functionality is

"good enough" that people feel satisfied and do not have the urgency for a better solution. Certainly, the Semantic Web is not only about search even though it is one of the problem areas that it is capable to address. Unfortunately, information overflow is one of the biggest pain points that users suffer. Furthermore, problems in data integration, application interoperability, and communications among autonomous machine agents have been tackled by other means such as screen scraping, Web aggregators, Web Services and other specialized techniques, as described in Chapter 3. Those intermediate solutions require less upfront capital and resources investment. They are considered as better strategic solutions even though they are indeed tactical work-around. Therefore, while everyone is enjoying the "harvest" of technological accomplishments, they would have less motivation to tackle the problem immediately.

9) "Bad Money" – bad source of funding

The source of fundings and customers often dictate or dominate the direction of a technology development.

In the high-tech sector, a product and technology might get rave reviews by early adopters. However, the sales and revenue might stall once the early adopter market is exhausted if the technologies are not designed and positioned to leap to the mainstream customers who are fundamentally a different group with different needs. The situation is a symptom or marketing illusion phenomenon described as the Chasm.[33] In an attempt, a company might become more sales-driven to achieve a certain financial target and to keep up the momentum. Unfortunately, such tactical defense would lead the company heading in a wrong direction and captured irrelevant customer base which might provide short term gain but the right strategy and path to cross the chasm and leap to the intended mainstream market.

The source of money often influences the marketing strategy. For example, in the case of WebVan which provides online grocery shopping and delivery service, a large amount of initial venture capital funding had led the company to expand into a bigger market in suburban regions with an enormous spending on infrastructure in order to leverage the economy of scale and scope[34]. Unfortunately, spreading the service coverage has done more harm to the business than good due to the wrong targeted customers. Its service value such as convenience is not worth the extra cost and change of behavior to households in suburbs. It is more relevant to young working professionals in Manhattan or other metropolitan areas. As a result, expanding the service networks caused the company a huge lost and eventually ran out of business.

Analysis

Unlike a startup, the research and development of the Semantic Web had been funded by various research grants and industry partnerships. Therefore, in theory, it is not constrained by any particular commercial factors. However, the initial corporate sponsorship and affiliated organizations or industry do have

influences on its direction to a certain extent. It might have associated with a wrong "public image" and pursue the wrong type of sponsors initially. First, Travel Butler has become the most or overly used example so far by the technical community and Tim Berners-Lee himself. Secondly, the current Semantic Web development focuses on the pharmaceutical and bioinformatics industry. Such an industry is very specialized and would not be able to serve as a strong business and reference case for other industries. For example, financial services companies would not adopt new technologies if they could not identify direct benefits to their operations and to enhance the bottom line[35].

## 4.2: Lessons Learned

Since its introduction by Tim Berners-Lee, the Semantic Web initiatives have been undertaken and evangelized by the W3C consortium for about six years. However, it has neither achieved significant progress and work result nor gained sufficient awareness in the technical community and industry. To do an acid test, an information technology expert and a MIT Sloan professor, Stuart Madnick[36], did a quick survey while conducting a lecture to a group of 50 mid-level executives from the IT industry. He asked the audience whether or not they knew what the Semantic Web was. Only 5% said yes but could vaguely define what it was. One of the managers described it as a new Web-related project initiated by Tim Berners-Lee. Certainly, such candid survey could not be an accurate means to draw an objective conclusion. It is a sign of warning and problem. We can analyze the situation through lenses of various theories and compare with examples of past successes and failures to gain some new insights. Madnick has been studying and building applications to solve the problems of data integration and integrity. One of his immediate feedbacks on the progression of the Semantic Web was not about technical issues but rather business or "PR" issues as he put it. Tim Berners-Lee has been the Semantic Web evangelist himself since the idea was conceived. His book, Weaving the Web, published in 1999 discussed in great details what his vision was. However, he often associates and demonstrates the idea of the Semantic Web with Travel Butler example (see chapter 2) regardless of the background of the audience. Although the example is a good illustration of its functionality and potential, for IT executives from big corporations, they are not able to see the immediate relevancy and benefits to their organization and business.

## 4.3: Is Commercial Attention and Stakeholders a Make or Break

As discussed in the previous section, the Semantic Web had fallen into traps, that some of the failed technologies had in the past. However, as the research and development continue, Tim Berners-Lee believes that the research is ready to leave the lab and enter the marketplace. In a recent interview with Tim, he stated that early commercial attention is both a blessing and a curse. On one hand, it has provided the momentum and a head-start which make gaining funding and other support much easier. The press coverage and visibility have helped attract exceptional talents and engineers to work on their projects. The recruiting process would otherwise have been difficult given its early stage. Apparently, Tim's affiliation with the project has increased credibility. On the other hand, Tim also mentioned having high visibility for an initially research-oriented project causes

unexpected pressure and sets different expectations and agenda. The initiative was able to gain sponsorship initially (through membership) from various sources and corporations. However, commercial sponsorship and partnership influenced the project direction and they demand more near-term financial returns and full-fledged end-products rather than experimental prototypes.

It is important to make a strategic decision on initial position and adoption path for the Semantic Web. From a technology standpoint, the Semantic Web is able to solve problems of data integration such as Enterprise Application Integration (EAI). In particular, EAI is critical in e-commerce applications, Merger and Acquisition and cross-corporation data sharing. It might seem logical to "follow the money" and target domains with high financial returns. Unfortunately, pushing the Semantic Web to handle large IT problems pre-maturely would be a setup for failure since commercial IT systems have high expectations in performance and financial returns which would not be delivered initially. Incapable of delivering near-term value reduces endorsement by big corporations, discredits its true value and increases roadblocks for future developments. Therefore, a suitable domain should be able to tolerate the risk and evolutionary nature of the technology. The domain will be discussed in Chapter 6 along with a pragmatic adoption approach.

## Chapter 5: Trend and Leverage of Participatory Architecture

In Tim Berners-Lee's own words, he measures the success of the Semantic Web (SW) based on "serendipity or unexpected reuse". In a sense, SW should facilitate the discovery of information and knowledge which would have been unfound or inaccessible. Based on a scale of 1 to 10, Tim commented that the serendipity reuse enabled by SW was about 2 to 3 and remained minimal. He believes that it is important to expand the scope of adoption from Bioinformatics to the mass consumption in order to benefit the general public.

This chapter will discuss how the unique characteristic of the Internet and Semantic Web would define the adoption strategy. Then, it will provide data on how the new trend in New Media consumption and distribution and Open Source development movement will provide opportunity for the Semantic Web to thrive on openness and accelerate the adoption.

### 5.1: Next Step: Strategy to Drive and Achieve the Critical Mass

The Internet offers a very unique characteristic in that it resides in a networked environment that can increase the value of individual components or participants far beyond what they could be on their own. If only one member enhances its online presence, it is of limited value to other members. However, if many members do so, the enhancement enacted individually will be compounded in a networked environment to produce an effect greater than the cumulative sum of the individual actions, hence, the "network effect." A network effect is a characteristic that causes goods or services to have a value to a potential customer which depends on the number of other customers who own the goods or are users of the service. One consequence of a network effect is that the purchase of a good by one individual indirectly benefits others who own it. For example, by purchasing a telephone, a person makes other telephones more useful. [37]

There are very strong network effects in the business model for computer software. One prominent example is Microsoft Windows and Office due to its compatibility with a variety of hardware and software, interoperability with other users and users' familiarity with the user-interfaces. Other examples include the recent Web Services such as the online auction service, eBay, the recent growing social networking sites or Web 2.0 applications such as MySpace, YouTube and Wikipedia, and communication tools such as Instant Messaging software and services MSN Instant Messenger. All these rely on effective strategy and leverage on the network effect business models. Achieving this enhanced network state requires reaching the tipping point or critical mass in adoption. Once reached, the incentive for all members to follow suit increases, since the benefits of doing so are likely to accrue to each new member as well as any existing members. For a standard data format, a crucial aspect of the Semantic Web, it's clear that not only does the network effect play a critical role in its successful adoption and future growth, it also determines its usefulness and potential value.

The D-Day analogy used by Geoffrey A. Moore[38] provides a good insight into how to reach the mainstream and wide adoption for the Semantic Web. Similar to a startup venture, the Semantic Web research group is trying to take technology to the marketplace. A common problem new startups face is that they tend to pursue a market segment based on the initial interest or hype which might in turn leads to a dead-end rather than a long-term end-goal. Essentially, the D-Day strategy is to help ensure our effort and resources to focus on one niche segment or application which will lead us to appropriate reference user-base and build up the credibility and viral effects in the field. Most technology failed to cross the chasm because companies are confronted with the immensity of opportunity that presents itself, and find themselves unable to deliver salable proposition to any true pragmatic adopter[39].

Typically, the initial challenge with any new technology or product is getting the first customer and then to cross the chasm and reach the mainstream for sustainable growth. Initial user base and implementation often serve as the flagship demonstration of the capability and value of the technology. It is important to identify the right domain to pursue. The decision on associating the technology and applying it to the appropriate domain is a strategic choice because it determines the future direction and whether or not it will be successful.

Although the Semantic Web did not face a big hurdle in getting initial adopters, applying the Semantic Web to EAI (Enterprise Application Integration) domain and the Travel Butler was either pre-mature or inappropriate. Consequently, the W3C and Tim are switching gear and targeting the Bioinformatics field which is a more appropriate starting point and provides a better platform for experimentation because the industry can tolerate a longer research and development cycle. Although drug development is a big expenditure business and requires high ROI (Return on Investment), it has a longer development period than the IT industry – years vs. months. Therefore, it would allow the Semantic Web to evolve until it is ready for prime-time. To a certain extent, it follows the footsteps of WWW technology which was first introduced to the science research community initially, except that it is not the end-goal for the Semantic Web.

To Tim Berners-Lee, the Web is about people, and the interactivity and interconnectivity among them. HTML pages enable people to link information and documents with each other. The process to grow the universe of WWW is a social process rather than a pure technological evolution and revolution. The idea is proven with the recent growth trend in the Web 2.0 phenomenon and CGM (Consumer Generated Media). Technically speaking, the Web 2.0[40] is not a new technology but a collection of technologies combined together and applied to the social interaction in a networked environment. The grassroots way of building a community-driven encyclopedia has demonstrated an application of the collective power and wisdom through a social process. Although Tim's current focus is to build a web of semantic data for the Bioinformatics field, he admitted that the rise of the Web 2.0 has raised the awareness of the value and significance of metadata and tagging technology in the general public. Unfortunately, the metadata in the Web 2.0 seems to be loosely defined and could be abused or misused. Nevertheless, there

is a good opportunity and synergy to leverage the participatory architecture and processes as seen in the Wikipedia and other social networking community.

In order to help the Semantic Web expand into the mainstream from the Bioinformatics domain, it is important to identify an appropriate new application domain. To benefit the general family households, applying the Semantic Web to multimedia is the most rewarding technologically and commercially. In the next few sections, a market analysis as well as a business case will be presented. In particular, New Participatory Media and Open Source development are the two new important trends and movements that will drive and inspire new opportunity for the Semantic Web.

5.2: Trend in New Participatory Media

In general, the Semantic Web community envisions the adoption process occurring along two paths with rough simultaneity. One path will follow closely defined, highly targeted projects that have the effect of embedding the technology within a corporate environment. The second path will be broader and Web oriented since more efficient data exchange, interoperation and resources sharing in inter-corporate settings offer high potentials.[41] In the rise of Web 2.0 applications and the new participatory Media phenomenon, it is believed that the second path resembles the adoption process of the World Wide Web. Since the adoption cycle in corporate setting is typically longer, open environment and consumer applications would offer the opportunity for the Semantic Web reach out wider audience and accelerate growth.

There are several ways to approach this issue. For instance, it is conceivable that a Web author could embed appropriate tags (XML, in the case of the Semantic Web) on a site that would enhance the likelihood of discovery using a Semantic Web-based search. The reason for taking an extra step to define XML tags is to enhance "discovery," or enable information easier to be found. The music industry provides a reasonable example.

Suppose a fan of mid-19[th] century American folk music were to use the major search engines to find new titles to add to his collection. Presently, this kind of search yields hundreds of thousands of entries, which is far too many to sift through. Now, suppose a Semantic Web search engine was available in the near future, a musician specializing in this genre might hire a thoughtful Web author to embed tags that define his music with greater precision than currently possible. As a result, the music seeker is much more likely to find the music maker and in turn, a mutually satisfactory exchange or relationship may ensue. Of course, as one musician experiences success, others will follow suit. The importance and success as a result of superior search results will subsequently gain attention and adoptions. This kind of grassroots pattern would closely resemble early personal Web pages, published by individuals who wanted to express themselves in a very public and discoverable way. The existence and growth of small and successful corporate projects and grassroots adoption will lead to a critical mass for Semantic Web sites, content, and tools. Reaching this point is critical in the growth and use of this technology. In many respects, tagging of objects will be optional and based on the desire to be discovered.

Years ago, pioneering individuals published their personal pages in a barren Web landscape hoping they would be discovered. With the advent of search engines, it was possible to register a site for the purpose of increasing the likelihood of being discovered. Just as the Web developed, the desire to be discovered becomes more significant, appropriate tagging will follow and grow.

All this is in fact happening to a certain extent. One project that goes beyond the typical annotation of music is called Pandora[42], a music discovery project grown out of MGP (Music Genome Project).[43] Over the past six years, people on the project carefully listened to the songs of over 10,000 different artists and captured the essence of music at the most fundamental level. They assembled hundreds of musical attributes or "genes" into a very large Music Genome. Taken together, these genes capture the unique musical identity of a song – everything from melody, harmony and rhythm, to instrumentation, orchestration, arrangement, lyrics and the singing and vocal harmony. It is not only about what a band looks like, or what genre they belong to.

Furthermore, studies and analysis shows the significant commercial needs and interests to utilize the Semantic Web in the Video Content Management space.

The following two sub-sections will first analyze the potentials of two New Media phenomena and how they will drive the adoption of the Semantic Web as a rich-media Web platform. Then, Section 5.3 will discuss a probable service-oriented model to capitalize these New Media trends and markets.

1) CGM (Consumer Generated Media movement)

> Personalized media development or CGM have shown a tremendous amount of growth recently: community-based Web sites such as MySpace, Facebook, Flickr and YouTube have proven effective in amassing large user-bases and maintaining high user-retention rates. The graph below shows the audience growth for CGM sites in the six-month period from July to December of 2005[44]. Essentially, the advantages of IPTV (Internet Protocol Television) include two-way capability lacked by traditional TV distribution technologies, as well as point-to-point distribution allows viewer to view individual broadcasts and possibly with other viewers with the full internet connectivity and via the browsers on the TV set, personal computer or other internet-enabled devices. A surge in use of YouTube and other hybrid social networking and content distribution service is an evidence of how the social interactivity and such behavioral change is crucial and set the stage to the growing adoption of Internet TV.

**Audience Growth for CGM Sites, July-Dec. 2005 (Note: 6 month growth period)**
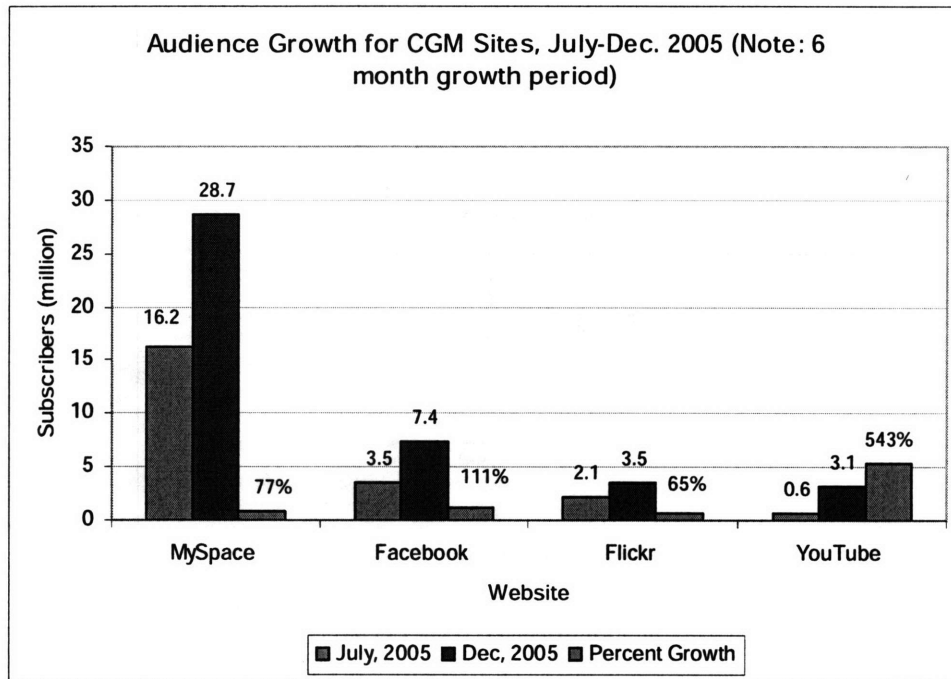
Figure 4: Audience Growth for CGM Sites, July-Dec 2005

There are currently thousands of Web 2.0 sites and hundreds more everyday. Some of which are the emerging "mashup" services which are Websites integrating content from various sources to provide more comprehensive services or experience. It is analogous to the system of modular design with individual components integrated through public interfaces in the new medium of internet. For example, Trulia[45] is a Real Estate market trend or house searching service with the location's information displayed using Google map. All of which require end-user participations and contributions of content and information in order to be useful, meaningful, and successful. Otherwise, they would be a site with empty pages. The kind of phenomenon of leveraging collective intelligence, participatory go along with the open source software development or the construction of Wikipedia, an online encyclopedia, for which almost everyone contributes free-of-charge, yet high quality is greatly maintained. Such paradoxical phenomenon has been puzzling the business world and a call for change in corporate hierarchy[46].

The figure below shows a selection of Web 2.0 site logos as of January, 2006. However, more are launching each day while some go out of business or services as well.

Figure 5: Web 2.0 Logos[47]

There are also services like Fourio[48] tracking the geographical location and the type of Web 2.0 services that are launched. Based on their tracking status, it appears the US is the most active region geographically. However, the influence and growing trend is certainly spreading around the world.
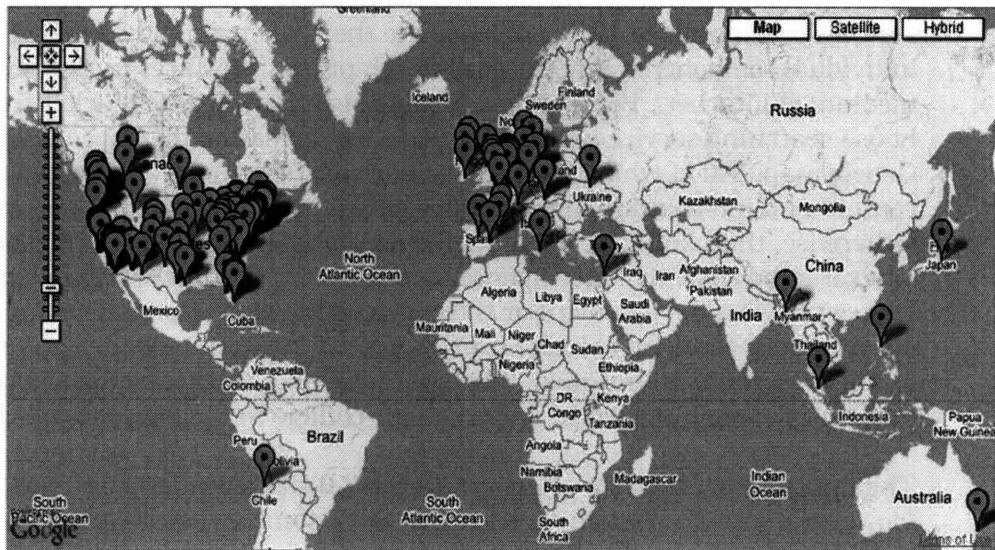


Figure 6: Web 2.0 Innovation Map

2) Trends in online television broadcasting and video distribution

Furthermore, the online video distribution is displaying growth and tremendous traction.

In addition to factors such as convenience, accessibility, viewing-on-demand, and affordability, the online video distribution market is largely driven by broadband penetration in the home and the improving quality of Web-based movies and television programs. With U.S. broadband penetration surpassing 50% mark for all web users, rich media continues to play an increasingly important role in determining future market strategies and managing digital assets for all Fortune 5000 companies[49]. The number of residential broadband users has grown from 21 million in 2003 to 40 million in 2006 and maybe as high as 100 million households worldwide. While currently many broadcasters are adding an Internet presence as a part their content-delivery channels and marketing tool to drive growth of viewers of this channel; the technology will potentially create a new market and business model which could eventually displace the traditional cable TV model.
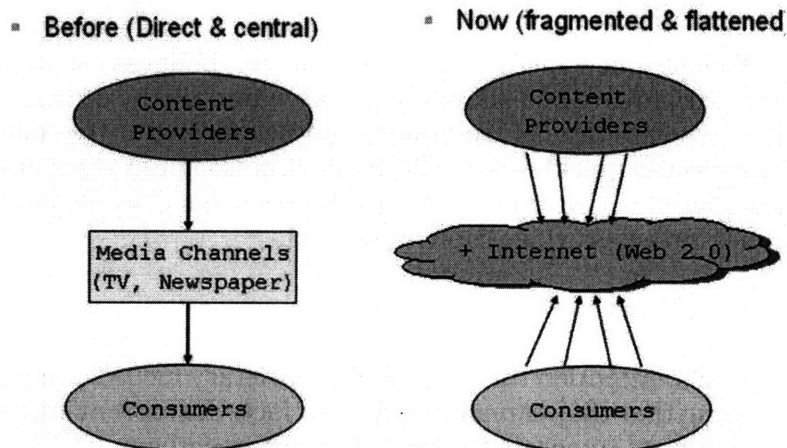


Figure 7: New Media Phenomenon

New video network start-ups such as Brightcove, DaveTV and their European counterparts, Akimbo and ITV, are showing traction and proving themselves in the marketplace of Internet Television and as a technology platform. These companies are joining the big video aggregators such as Google and Yahoo to redefine the television and video media channels. For example, Brightcove recently signed a deal with The New York Times to enable the distribution of streaming video content across all of The Times Company's online properties. According to Parks Associates, the number of US households that downloaded videos online at least once a month was 11 million in December 2005[50].
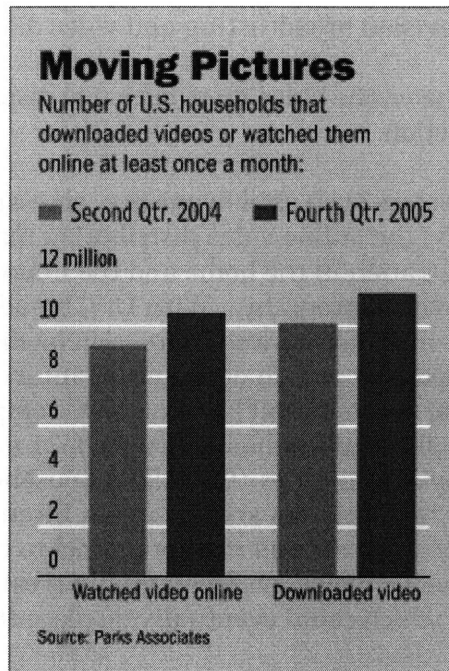
Figure 8: U.S. Online Video Download

Allen Weiner, a new-media analyst at Gartner, believes companies like DaveTV, Brightcove and Akimbo can thrive initially by drawing early adopters who are tired of traditional network fare. But the competitive landscape could quickly change if the likes of Google, Yahoo! and Microsoft can leverage their Internet search capabilities and provide portals for consumers to access video content, he said.

5.3: Open vs. Proprietary

There has been misconception that open and proprietary models are mutually exclusive, meaning that they cannot co-exist. In fact, the current status and success criteria focus on how to integrate the two effectively together.

The growing trend in Open Source development demonstrates the shift from product-oriented business model to service or utilization model in the new Web 2.0 environment which appears to offer everything for free. For the Semantic Web, it means that it is important to leverage the new participatory culture and to open up or "outsource" the development of semantic data to the community in order to accelerate the growth and adoption of the technology. Along the way, commercialization would be realized by a service model.

1) Open is free but not free

First, we will take a look at the value proposition to maintain openness. The evolution and recent success of online distribution have provided a few lessons and insights. Some people, at least those in the music and film industries, used to think that the peer-to-peer file sharing network would

destroy their industries.  In fact, creative force and publishers should not spend time to investigate whether or not they would be undermined by such technologies; instead they should leverage these new technologies to increase the visibility of their work.

Technology itself does not change the fundamental dynamic by which millions of products reach millions of potential consumers since the role of publishing would not change and would always be the middleman and more so as aggregator to help bring the producers and consumers together.  In the Internet time, Google and Yahoo are emerging as big players serving the traditional role of aggregator although the means of doing it might be different.  Google's use of implicit peer recommendation in its page ranking algorithm plays much the same role as the large retailers' use of detailed sell-through data to help them select their offerings.

New media does not replace, but rather augments or expands existing media marketplace.  Opportunities exist to arbitrage between new distribution medium and the old.  In the case of file sharing networks, it fuels the trading of records and CDs on eBay which are unavailable through traditional channels.

As seen in the rich media distribution model, the paradigm shift from copyright economy to moral economy, as described by Henry Jenkins[51], is changing the competitive landscape.  According to the definition in Wikipedia, a moral economy is an economy that is based on goodness, fairness, and justice[52].  Certainly, there have been many arguments what works and what does not.  The theory has been proven with the success of iTune that music fans are willing to pay for their favorite music at a fair price of $0.99 through a convenient channel rather than going through the hassle to download songs from illegitimate sources.  The shutdown of Napster has in fact inspired a lot of succeeding applications or similar models with tremendous successes, one of which is YouTube, a free online video sharing service which was subsequently acquired by Google for $1.65B[53].

In the past, copyright was the primary mechanism for protecting the ownership of intellectual property (IP) and creating barriers to competitions.  However, as we have seen in the past decade, piracy of IP has flourished ever more.  On the other hand, building upon the theory of moral economy, Tim O'Reilly does not see piracy as too much of a threat but rather progressive taxation[54].  Essentially, obscurity is a bigger threat to creative artists than piracy.  Not only does the online distribution provide established companies an outlet to wider audience quicker but it also empowers the smaller, less well-known producers the chance to get the exposure which otherwise would have been impossible as explained in the Long Tail theory.  The phrase, Long Tail was first used in the business context by Chris Anderson who observed the phenomenon of the Long Tail in today and future business environment[55].  In his argument, products that are in low demand or have low sales volume, hence the Long Tail, can collectively make up a market that rivals or exceeds a handful of bestsellers and blockbusters, if there is enough variety and the

aggregated volume is high. Examples of such mega-stores include the Amazon.com, an online retailer and Netflix, a web-based video rental service. The Long Tail has a market potential and, as the examples illustrate, the Internet has enabled businesses to tap into that market successfully.

Many media companies resist adapting to the online channels which they think would hurt their businesses. In fact, it is only when artists or authors are well-known enough before would their work be pirated. At that point, piracy is kind of a progressive taxation, which may shave a few percentage points off the sales of records of well-known artists in exchange for massive benefits of greater exposure which may lead to increased revenues. At the same time, customers want to do the right thing if they can. If they conceive the artwork or product in general as good value and high quality, they become evangelists for the brand. For example, enthusiastic readers reported the infringing copies to support O'Reilly Media Business and recognized that free copies were illegitimate. Ultimately, consumers are willing to pay for online services for a subscription fee of $9.99 a month.

## 2) New model: Services

Although the Internet is open and free, leveraging it to empower service-oriented businesses are proven to be the new model for growth and sustaining revenue. Free is eventually replaced by a higher-quality paid service[56]. The transition from product orientation to services model is becoming ever more significant, especially in the "Age of Commodization" in high-tech businesses[57] when the improved performance does not generate higher financial returns. It is due to the fact that improved functionality of a system is performing more than well enough to address the needs of customers in mainstream applications. When that situation occurs, companies overshoot what customers, in a given market, can utilize. Companies trying to differentiate themselves from the competitors on the margin of performance improvement would fail miserably and end up with performance surplus[58]. In high-tech businesses, price has dropped tremendously for PCs, memory capacity, mainframes and many other hardware. Furthermore, China, India, Southeast Asia and Eastern Europe have the competitive advantage of low-cost labor sourcing. In software business, the cheap overseas competition, the freeware and open source development are putting tremendous pressure on price and driving down the profit.

The resonating theme is to encourage transitioning from product orientation to services-focused, a successful model which has been advocated by Prof. Michael Cusumano in the context of software business. In his latest book, The Business of Software, he provided guidelines and insights into deriving optimal strategy with the right mix of products and services for software companies[59]. It consists of extensive studies of different models and practices with company examples of both great success and failure. Two of the success stories are Infosys, and IBM. In particular, IBM is well known for its

dramatic corporate turn-around with a transformation from product business to a service company led by Louis Gerstner in the early 1990s[60]. Not only did it help the company resurrect from the downturn in the business of proprietary hardware and software products but it also helped the company tap into the high margin, recurring revenue of service businesses including consulting and maintenance. In addition, the services have helped increase sales of its own software and hardware using the product-bundling strategy.

In the consumer side, the Web 2.0[61] has created a number of business model based on free Web services. First, one popular model is online advertising, known as the "Google model" in which Internet-based services are offered to end-users free of charge by selling advertising space on their Web sites. However, the advertising revenue model would not provide sustainable and sufficient revenue, particularly for early entrants who have not reached the critical mass to achieve a high volume of page-views and user-clicks[62]. In addition, the issues of online-advertising frauds and click-through effectiveness have not been resolved. A brief qualitative article, dated back in 2002, touched on this issue. The skepticism on the current online advertising model had led to the discussions of rich-media advertising, one of the hot topics today. As discussed later in Chapter 7, rich-media advertising is the next face of online advertising and can offer a great strategy to transform the technology-centric Semantic Web to a consumer-driven solution.

Second, another promising model is the Software as a Service (SaaS) model which is gaining traction and is well-known for its association with salesforce.com. SaaS is an Internet-based software delivery model, in which maintenance, daily technical operation and software support are performed entirely through the Internet.

Third, a new subscription-based model is also growing although it is still in the Early Adopter Stage in the Technology Adoption Life Cycle[63]. It has been successful in a number of niche markets including the online music download service with monthly subscription. Services like Kazaa flourished in the absence of competitive alternatives. If there is a service that provides access to all the same songs, freedom from onerous copy-restriction, more accurate metadata and other added value, there will be hundreds of millions of paying subscribers. Many Web portals like MSN, Yahoo! and others have already started offering such services with success. Essentially, their business is built on the free Web. The current online file sharing services are mediocre and their quality is low and uneven. The potential of the subscription model is yet to be unleashed. Other "early adopters" exist in the marketplace, for example, in the computer and technical book publication industry (O'Reilly's Safari Books Online) and online sport games broadcasting (MLB.com). In the case of MLB.com, they provide readers and sports fans exclusive access of content which might not have been economical and viable using traditional channels. As the record in the television industry shows, people prefer subscriptions to pay-per-view. Ultimately, the Web should be used as a great and effective marketing tool and channel. Take Safari subscription service as

an example. The company publishes substantial number of free publications online on their advertising-supported Website under the "open publication licenses" where free distribution is allowed. The rationale behind it is to build product awareness, loyalty and a community of audience. Doing so would help drive sales of hard copy of the books and other items[64].

Finally, the professional services business in technology consulting (Infosys or McKinsey Business Technology Practice), maintenance and customization is well established. For open source operating system such as Linux, vendors like Red Hat follow the service model in which they sell customized version of the "raw" Linux with enhanced help support and user-friendly interfaces. Although their businesses are still in the early stage and yet to be proven in the marketplace, Bill Gates and Microsoft are watching them closely since they pose potential threats to their business in the foreseeable future[65].

In sum, focusing on the broader, Web oriented, and rich-media tagging for the adoption of the Semantic Web are driven by two forces. First, there is tremendous growth in New Participatory Media including online video distribution and the movement towards grassroots media development and social communities. Second, the rise of Open Source development demonstrated the shift from product-oriented business model to services or utilization model in the new Web 2.0 environment which appears to offer everything for free. For the Semantic Web, it means that it is important to leverage new participatory culture and to open up or "outsource" the development of semantic data to the community in order to accelerate the growth and adoption of the technology. Along the way, commercialization would be realized by a service model.

The previous two sections first presented the potential of the new phenomena that will drive the adoption of the Semantic Web as a rich-media tagging Web platform. Then, business strategy and case will be discussed further in Chapter 6, and 7 on how to monetize and maximize its potentials.

## Chapter 6: Pragmatic Adoption for the Semantic Web

In the previous chapters, we have discussed the benefits and challenges in the development of the Semantic Web. As discussed in Chapter 3, the Semantic Web plays a central role in data and functions connectivity. Essentially, the Semantic Web is the key enabling technology for the data-driven architecture and implementation of the Internet OS. A more detailed technology diagram presented later in this chapter will illustrate that the Internet OS will be built on the foundation of the Semantic Web which serves as the data pipeline for the information flow among applications. In order to overcome the challenges and to move forward in realizing the future Web, a pragmatic approach for the adoption of the Semantic Web will be crucial. The approach to be discussed will help accelerate adoption and achieve the critical mass for the Semantic Web.

### 6.1: Strategy

Business managers who have lived through the implementation of enterprise applications experienced failure of gaining return on new technologies which usually requires a large amount of upfront investment. There is usually significant lead times required before meaningful business impact can be realized. A typical cycle takes five to ten years to implement. That issue significantly prevents and discourages management from investing new technologies. In an interview with Clay Christensen[66], he advised to analyze three aspects when preparing for commercialization of innovative technologies:

- Determine whether there is a new application.

  For any new technologies, it is not the novelty or better performance which secures wide adoption in the marketplace. In certain circumstances, it is important to derive a new application to serve the unconsumption market[67]. For the Semantic Web, the issue is that it is targeting a market segment which has already adopted existing technologies. From the users' point of view, technologies that have achieved a certain satisfactory level are considered as good enough. Therefore, users have fewer incentives to switch to new technologies or systems. In the other words, the switching cost is considered higher than the potential benefits gained from adopting the new technologies.

- Understand the level of interdependency. How many steps in the value chain before it can provide an end-to-end solution?

  To understand the level of interdependency, one should determine the number of steps in the value chain before it can provide an end-to-end solution. Before commercializing a new technology, one should always thoroughly evaluate the readiness and strategy to profit from the investments in the research. An effective strategy should take into account the interdependencies and modularity of the product's architecture. Interdependency means the way one component is designed and made

depends on the way other components are being designed and made. One of major mistakes is that research organizations or companies tend to fail to decouple an integrated value chain and begin selling components into the open marketplace[68]. The figure below shows how different stages of products and technology development changes the nature of competition and the product strategy[69].
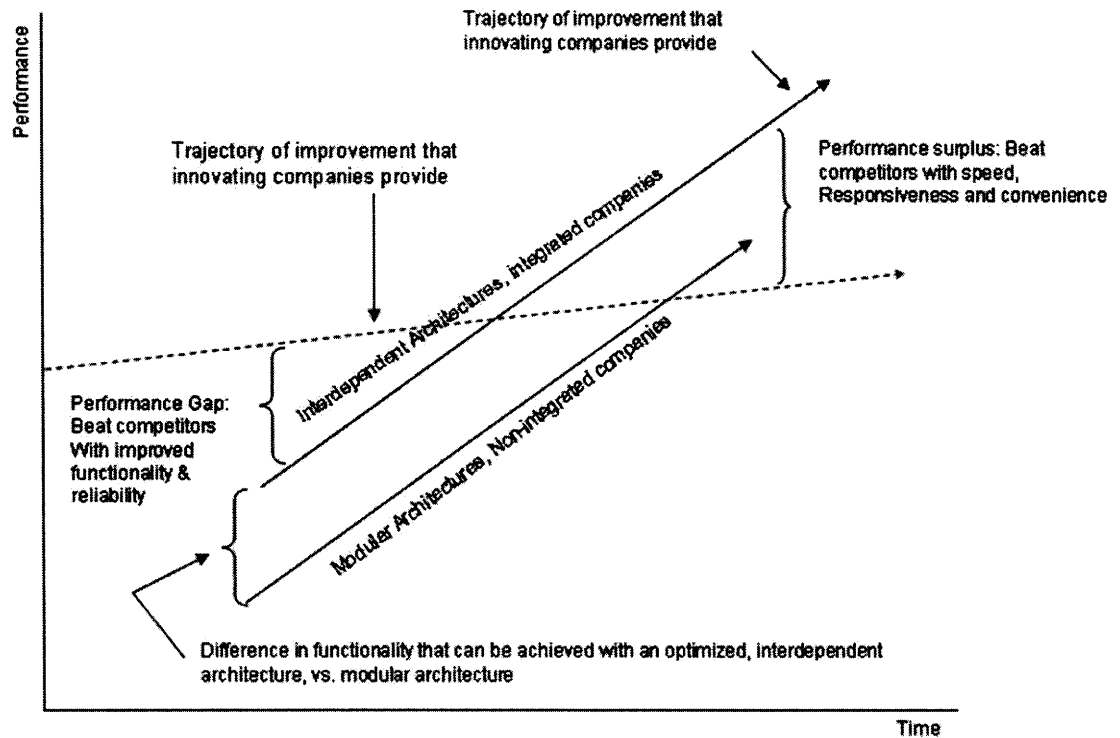


Figure 9: How Technological Progress Changes the Basis of Competition and the Nature of Product Architectures

When the market is in the early stage where most products are just not-good-enough and there is consumer demand in higher performance, companies with the proprietary, highly interdependent technologies enjoy the control and competitive advantage against the competitors if they make the best possible products. In this case, the companies must control the design and manufacture of every critical component of the system. However, as the products become more mature that the performance improvement in products outpaced the desire or ability of customers' usage, proprietary architecture will lose stance to modular architectures since speed-to-market, flexibility and compatibility with the other subsystems in the overall system are the critical factors. If a new player with radical innovation or breakthrough which is not compatible with the existing infrastructure and system in this mature stage, they would fail to get the mass volume of adoption. There are typically many interdependencies that mandate change in other elements of the system before a viable product that incorporates a breakthrough technology can be used. In the case of the Internet, although the existing WWW is not perfect, it has certainly fulfilled many of the needs of the general

consumer and has reached the mature stage and beyond the original intent of the invention. On one hand, the Semantic Web is an extension to the existing Web and does not provide a new obvious functionality or application from the user perspective. On the other hand, the effort required to migrate the existing data and construct new ones to conform to the new standard is tortuously long and expensive. Therefore, the Semantic Web requires taking the modular design strategy and ensuring plug-in compatibility with existing systems without major overhaul.

- Find the right job.

  In the traditional marketing model, companies typically define the market segments by product attributes, category or price. Unfortunately, those analysis and models provide little information about the consumption behaviors which change far more often than their demographics, psychographics and attitudes[70]. Demographic data cannot explain why a man takes a date to a movie on one night, but orders in pizza to watch a DVD from Netflix the next night. In the conventional thinking, one would consider Palm Treo and other PDA devices as competitors of Blackberry. However, from job-based market analysis, consumers, particularly working professionals, use Blackberry to receive news, kill time when waiting for the next flight, in addition to accessing emails. A key insight is that some of their competitors turn out to be news publishers, e.g., New York Times and TV programs in airport waiting areas. The problem with the Semantic Web is that it has been perceived as a technology framework without clear association what the job it is supposed to fill. From a technical point of view, the idea of connecting and automating the communications between machines does not clearly explain what purpose it serves. While the Semantic Web is a general-purpose standard and framework, focusing on certain domains or applications would be critical in achieving wide adoption. The initial focus or "job" should lay out a pragmatic adoption path.

6.2: Implementation

The experiences from the early adopters of Web Services such as Dell[71] have lent some valuable lessons which applicable to the Semantic Web. A key insight from these lessons is that the strategy should follow a list of criteria as follows:

- Leverage existing technology investments
- Implement incrementally
- Focus on tangible early wins
- Plug-in elements over time

To follow the above criteria, an application domain will be first identified based on the analysis and service strategy discussed in Chapter 5. A pragmatic interdisciplinary framework will be derived to transform the technology into a successful product in the marketplace.

With the above criteria as a guideline, rich media metadata tagging solution is poised to be a promising application domain. Metadata is text-based information that describes the essence of a media asset. As the "Web 2.0" moves into rich media, every content owner is trying to make their media searchable. Media assets such as sound, images and motion pictures are not natively searchable by computer. At the moment the world is full of proprietary, manual-labor intensive, balkanized systems for attaching metadata. Without a common lingua franca, valuable assets remain frozen in their silos unable to be found by interested parties. In parallel, the traditional channel structure of television and radio is being replaced by searchable VOD (Video On Demand), IPTV (Internet Protocol TV), podcasts, and videos on the Web. In this brave new world, all content must be wrapped with descriptions in order to be marketed.

Among all the rich media assets, the challenges in organization, retrieval and accessibility of video content are the most rewarding opportunities technologically and commercially. Once the problems with video content are resolved, solution to other rich media types could be easily obtained. In a number of correspondences and interviews with management executives at MLB.com, a sports and entertainment company, and Broadway Video, a digital video archiving studio house, the media industry is in desperate need for assembling quality metadata from multiple sources – some gleaned from past searches (the Google method), some generated from automated image analysis tools, some penned by professionals and some generated by aggregating public "folksonomy" [72] tagging such as Flickr and del.icio.us. Most promising marketing technologies such as targeted ads, recommendation engines, and preference-based television have their adoption rate regulated by the quality and readability of their metadata, so there is huge pressure for metadata to evolve and standardize. There have been some successful but simple and open metadata schemes such as MP3 tags, Flickr tags, and del.icio.us tags. Conversely, most complex metadata systems are proprietary and weak extensions of library schema. Various open architecture attempts have been made for cataloging, but nothing addresses the searching media globally.

Along with the extensive and established standardization and development effort, the Semantic Web is well positioned to solve the problem. Among media industry companies, a private special interest group initiative called OMNI (Open Metadata Nomenclature Initiative) was created[73]. It is modeled on the open source movement, to promote extensible standards for describing images, motion images, games and non-text objects. At the same time, under the Semantic Web research group, there has been a number of special interest group established for different industries. It is in everyone's best interest to consolidate and create synergy among these groups.

Furthermore, many media companies are diversifying their business model by establishing a distribution network with online media companies including search engine giants like Google, online entertainment portals like Yahoo and/or the new growing players of online social networking sites like MySpace. Instead of concentrating investments in company's own Web site and online channel, the traditional media companies are distributing content with online media companies. "The pricing around video advertising is similar to that of cable TV which makes it attractive to advertisers. That's clearly a growth area, and one that we're investing

in heavily", said Jon Miller, CEO of AOL in <u>Business 2.0 Magazine</u> on July 10, 2006[74]. Just one month later on Aug 7, 2006, a headline in Wall Street Journal reads "Google to distribute MTV Clips. Deal for Ad-backed Videos Could Bolster Revenue, Broaden Viacom's Reach."[75] A paradigm shift is undergoing in the TV and video distribution space.

The figure below illustrates a new online content delivery stack which creates tremendous opportunity for solution providers of video metadata tagging technologies.
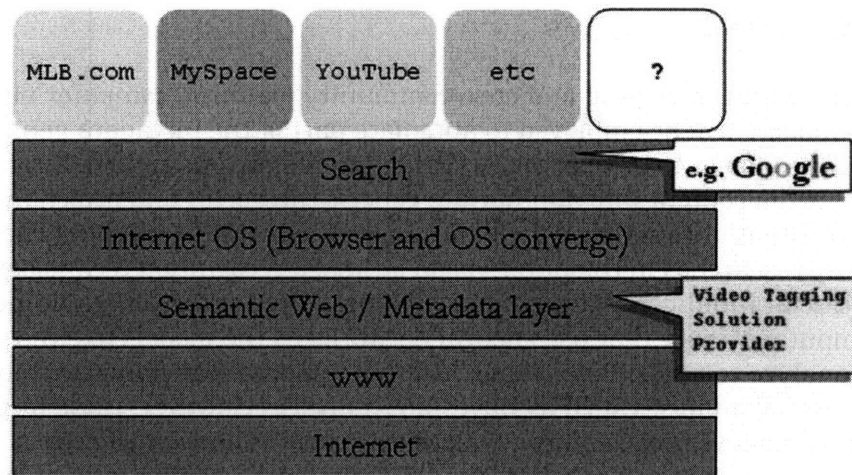


Figure 10: Online Content Delivery Stack

With abundance of metadata already generated and stored in the video content, it provides a great opportunity in leveraging the existing technology and investment. The following process will be applied:

- First, bootstrap a repository of semantic data by converting the existing taxonomy and metadata in the video content
- Then, use the clustering technology to build the metadata context and ontology.
- Finally, leverage open, collaborative decentralized community-based platform (Wiki or user participatory model) to refine and transform the taxonomy into "folksonomy" – methodology consisting of collaboratively generated, open-ended labels that categorize content such as Web pages, online photographs, and Web links[76]. It is intended to make a body of information increasingly easier to search, discover, and navigate over time. A well-developed folksonomy is ideally accessible as a shared vocabulary that is both originated by and familiar to its primary users.

The idea is to start by compiling the current popular XML schema, taxonomy, constrained vocabularies used for motion images today. Then a core group will offer the first stab at other pressing needs—multiple element collections, partial encryption of sections of the metadata, version history, and temporal relationship problems. From there, the community will refine and proselytize the standard. The

standards will be made open and extensible. This will create an abundance of service and software opportunities.

6.3: Interdisciplinary Framework

To realize the aforementioned process and solution above, a multi-stage process and interdisciplinary framework will be discussed below. The framework will make use of two technologies: Computer Vision and Clustering. It combines the two to implement a rich media tagging platform.

1) Computer Vision Technology

Along with the tagging and open community platform, computer vision technology will be employed to automate part of the metadata generations. The idea is not to have precise tagging completely automated, but rather put the media content into appropriate generic categories and attach metadata to the content. From there, the data is enhanced with user input. There are many researches in computer vision or image recognition technology based on Artificial Intelligence algorithms or Neural Network theory, a branch of computer science that uses neural networks as the models to either simulate or analyze complex phenomena. All these approaches primarily rely on statistical analyses such as Bayesian Analysis. However, those methods rarely achieve great accuracy. Recently, there is a new research breakthrough in computer vision. The challenge is being tackled by researchers at MIT's CBCL (Center for Biological & Computational Learning Lab) led by Tomaso Poggio. One of his former students, now post-doctoral researcher, Stanley Bileschi has made some promising progress. His invention and technologies has successfully surpassed the performance of traditional approaches. For example, one of their systems can detect people and cars in a street scene about 95 to 98 percent of the time. The research could be used in surveillance camera in an office building or military base, eliminating the need for a human to watch monitors or review videotapes. Other applications might automate computer editing of home movies, or sort and retrieve photos from a vast database of images[77].

Processing visual data is computationally complex, noting that people use about 40 percent of their brains just on that task. There are many variables to take into account when identifying an object: color, lighting, spatial orientation, distance and texture.

As an example, the diagram below shows the major component of a face detection system.
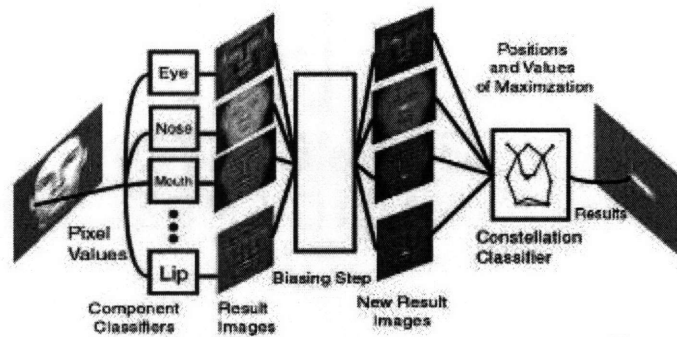
Figure 11: Major Components of a Face Detection System[78]

Instead of using statistical learning systems to teach computers to recognize objects, the CBCL researchers used another approach: they looked at how neurons are acting. The programmers make a mathematical model of those patterns, tracking which neurons fire (and how strongly) and which do not. They tell the computer to reproduce the right pattern when it sees a particular pixel and train the system with positive and negative examples of objects. For example, it is able to identify which one is tree and which one is not.

But instead of learning about the objects themselves, the computer learns the neuron stimulation pattern for each type of object. Essentially, it is learning patterns of patterns: the patterns of neural reactions not just to pixels but to groupings of pixels which is very powerful. Later, when the system sees a new image of a tree, it will see how closely the resulting neuron pattern matches the ones produced by the other tree images. According to Dr. Poggio, the process is similar to the way a baby's brain gets imprinted with visual information and learns about the world around it.

There are numerous research labs around the world trying to push the state of the art of the technology. The latest technical research report based on "Caltech101 data set" has published[79]. The CBCL lab has been leading the league in performance. Figure 11 shows the benchmarking test result among different groups.
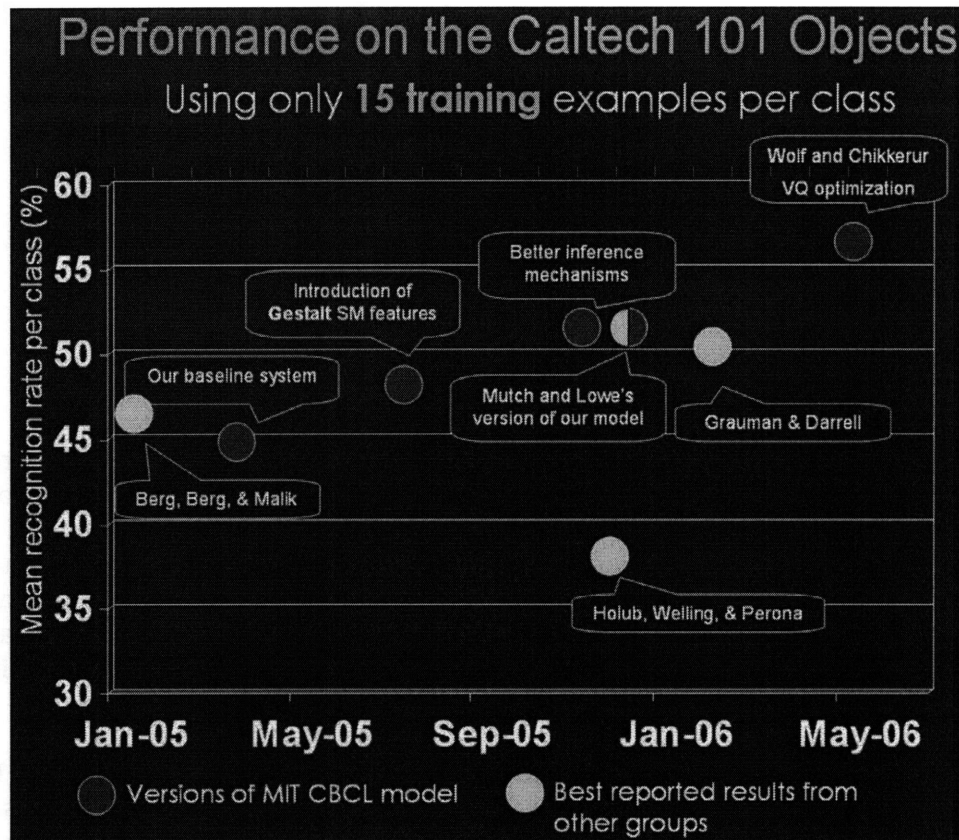
Figure 12: Performance on the "Caltech 101 Objects" test dataset[80]

The following two figures demonstrate the application of their technology by integrating it with Picasa, a photo organizer software from Google. First, an end user need to find some photos to train the system and go through three simple steps to annotate or add tagging data to the training sample photos. After inserting the meta-tags, the end user can use it in a production environment for the rest of the photo in the database. The demo is integrated into the Picasa software from Google to display the search results in a user-friendly atmosphere.

With this significant progress in image recognition technologies, applying it to video frame-by-frame would be quite straightforward since video editing techniques already exist to extract images frame-by-frame and link metadata to each frame.

## Identity Based Search

**Step 1. Upload Images**

Images are uploaded

Faces are automatically and rapidly detected during the uploading process.

**Step 2. Train Detector**

By assigning names to a few of the detections, the user teaches the machine to associate facial appearances to identities.

The machine is learning what these people look like.

**Step 3. Recognition**

Now the machine takes over and assigns names to all the remaining faces.

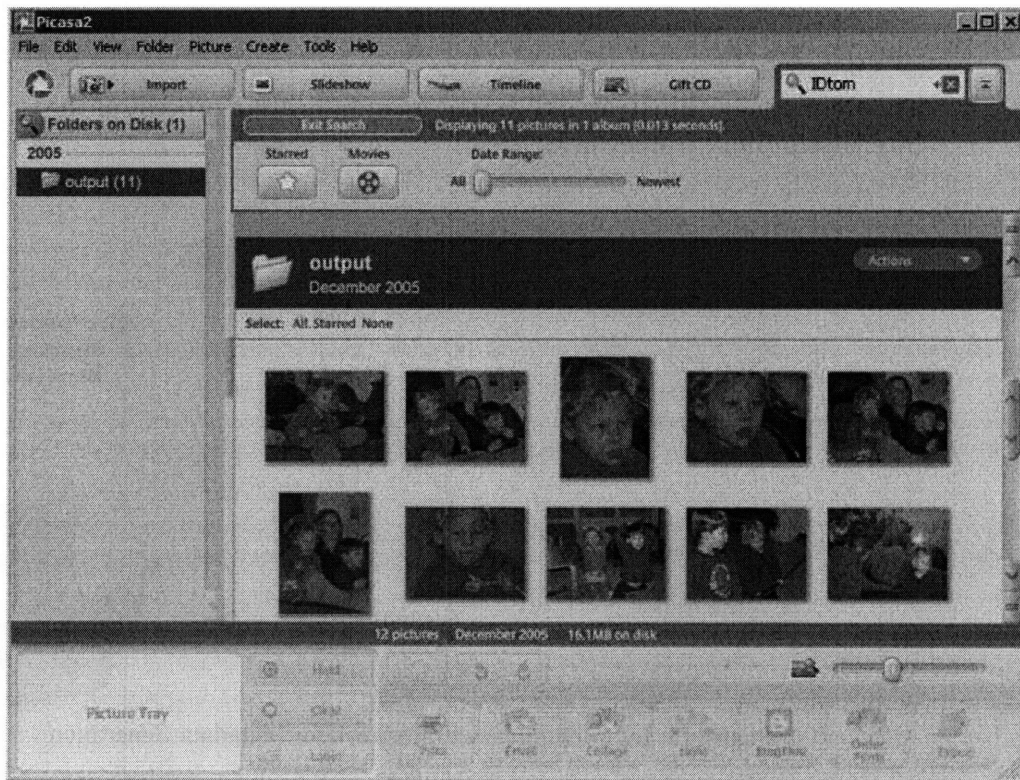Figure 13: Simple Training Sample Photos and Steps[81]



Figure 14: Picasa Integration[82]

## 2) Clustering technology

Clustering technology can help automate the generation of video metadata tags. One promising and superb technology is EigenCluster[83] developed by a research group at MIT Computer Science Research Lab under the supervision of Professor Santosh Vempala. EigenCluster is an experimental search-and-cluster engine based on a spectral algorithm. The divide-and-merge methodology used in the engine can be applied to the generation of a hierarchical taxonomy and for image segmentation. It is effective in generating ontology for video content and development of inference engine. Therefore, applying EigenCluster to the preliminary video metadata generated by an image recognition algorithm would be very effective.

A few companies such as Applied Semantics (acquired by Google and now part of AdSense) have developed an auto-categorizer software. In addition, a technology called "visual crawler" is used in video search engines which is capable of finding an abundance of contextual information, or metadata that relates to each video. All these examples provide a great reference for applying EigenCluster to rich-media content categorization since it provides superior performance over the peer technologies in the field. Although currently, it is solving text-based information, it could be tailored and optimized for media content.

## 3) Rich-media tagging framework

The figure below illustrates the staging and iterative process of the interdisciplinary framework.
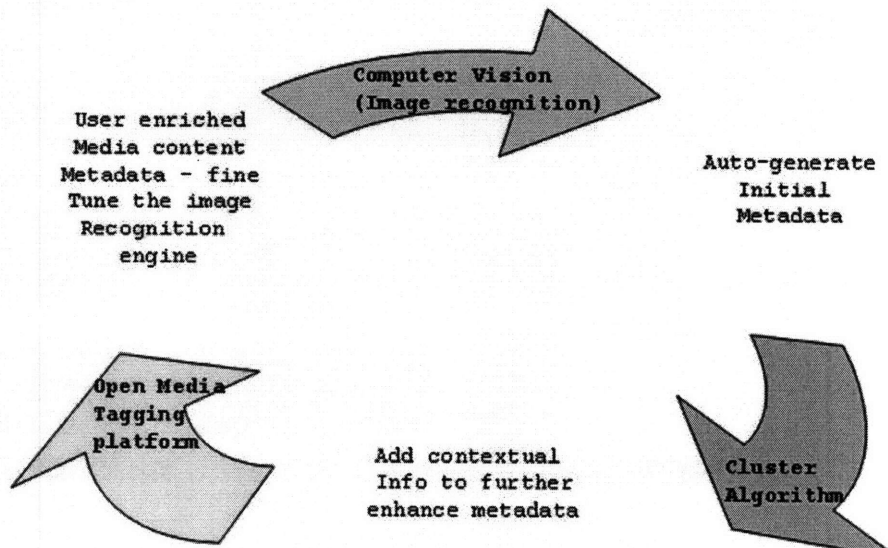


Figure 15: Interdisciplinary Framework for Metadata Generation

The proposed solution is diverging from the extremes of the traditional closed-environment manual-labor model and the fully artificial-intelligent approach to metadata content generation. Instead, the proposed solution is

focusing on an automated solution to metadata generation through open, collaborative solutions. This technology is disruptive to conventional and existing approaches because a knowledge-sharing model for rich-media content management reduces maintenance costs and achieves higher production efficiency for metadata generation. The proposed solution simultaneously captures customers' activities and preferences on the content provider's own Website.

Over the past decade, mass-media has been transformed by the digitization of content and by the introduction of new channels to access, assemble, and distribute digital media throughout the Internet. Rather than being passive consumers of content, waiting for companies to push out new information, consumers increasingly consume what they want, when they want it, and how they want it. On the production side, a new kind of "remix culture" has emerged. The result is a culture of participation where every individual has the potential to be a publisher, photographer, performer, or programmer.

For more than a hundred years, organizations across all sectors have been driving towards efficiency by creating routine or automated processes. This single-minded focus has allowed many organizations to create an array of impressive core competencies, but it has also blinded many to see where true value lies. Moving forward, it is not effective to separate the leaders from the followers; companies should organize and enable their talent. The proposed solution will facilitate this shift towards amplifying the practices of consumer, especially as they seek to collaborate.

## 6.4: Product Architecture

Through the use of tagging technology, the proposed solution can create searchable, descriptive text for any digital video asset. Tagging information is generated and edited by the end-user of the proposed solution software, thus automatically creating relevant digital video search indexes. To ensure credible and consistent information, a community voting system will keep the casual user from incorrectly or nefariously editing the descriptive digital video text. Additionally, some metadata can be restricted from user-editing by the content holders themselves (e.g., the audio transcript of a video file, or statistics for a particular sporting event). By giving client corporations control over their digital video metadata, corporations can maintain a productive community atmosphere without sacrificing the quality of their searchable text. Since tagging information will not only be file-specific, but *time*-specific within a particular file, the proposed solution is to enable search *within* a video file.

The figure below illustrates the system architecture of the proposed solution in the context of MLB.com as an example.
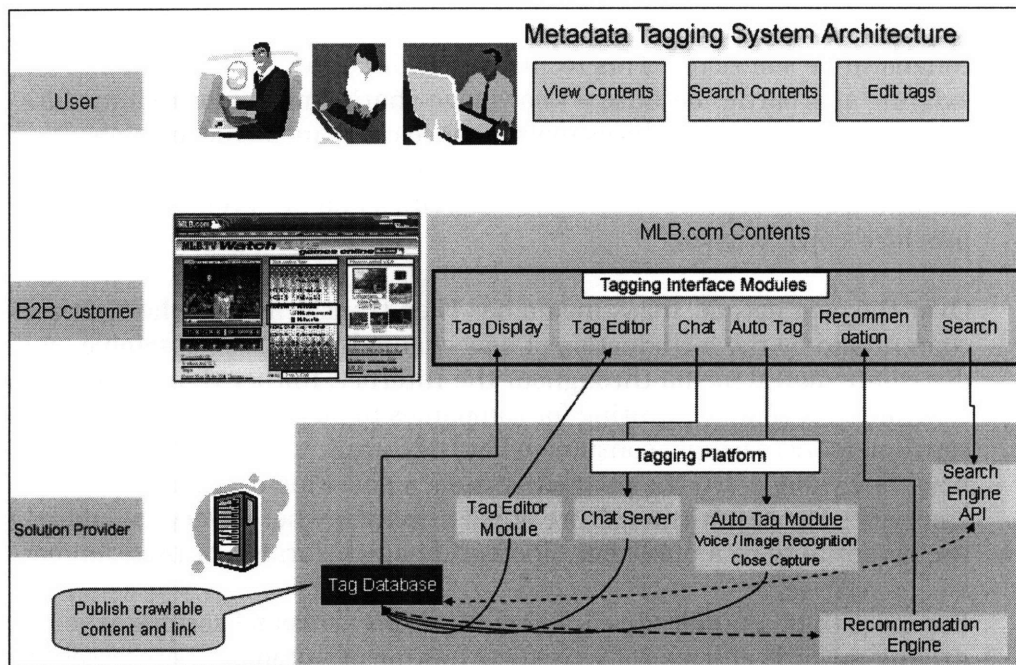
Figure 16: System Architecture of Proposed Solution

The software architecture functions according to the above diagram. Key to development and implementation will be a distributed software architecture model, allowing for a fault-tolerant system that balances load while minimizing latency. Since the size of the infrastructure will be proportional to the traffic generated, it enables solution providers to implement revenue-sharing and pay-per-click pricing which will be critical for maintaining infrastructure while creating company revenue. The business case will be discussed more in details in the next chapter.

In this chapter, we discussed the importance of focusing on a pragmatic adoption approach by taking the strategy that determined relevant and new application. It is important to ensure that the strategy does not require high level of interdependency in technology deployment and the number of stakeholders involved. The kind of application to be focused on should effectively and clearly associate with a job that the technology is supposed to fill. To many, the Semantic Web appears to be abstract and does not fulfill a "job". As we identified the significance of searchable rich-media, solutions developed upon the Semantic Web should fulfill the job of enabling searchable content. However, as we have noticed, metadata tags generation is both tedious and laborious process, using complementary technologies to enhance the process is crucial. In the past, many have attempted to fully automate the tagging process using image recognition technologies but have not gained great success. Many overlooked the power of multi staging process as discussed above in the Interdisciplinary framework section. It is more efficient and effective combining the human wisdom and the machine intelligence.

## Chapter 7: Business Case

With U.S. broadband penetration surpassing the 50% mark for all Web users, rich-media continues to play an increasingly important role in determining future market strategies and managing digital assets for all Fortune 5000 companies[84]. These companies currently spend up to $40 million to manage their digital assets; the proposed solution on the Semantic Web architecture provides solutions for the associated 15% maintenance costs for keeping these digital assets up-to-date[85]. While helping with these maintenance costs, the solution also adds value to existing rich-media assets. By allowing end-users to access rich-media and edit customer-specific tagging information, companies can expect an increase in sales of at least 0.75%[86] of several billion dollars in (CPG) Consumer Packaged Goods. This same activity by end-users also renders a company's digital assets searchable, and enables marketers to target specific advertising to their customers.

The digital video files contain an abundance of information that cannot be extracted easily, but is desired by both video content-providers and consumers. As broadband penetration continues to increase, greater amounts of rich media are being created and converted from traditional analog and digital handheld formats. Unfortunately, digital video content providers have been unable to provide searchable video content to consumers without investing heavily in cost-prohibitive, tedious metadata management. Often, this tagging must be done by hand, with individuals having to watch each video and record all relevant tagging information. Maintenance activity routinely amounts to 15% of all digital asset management costs, which can soar to $40 million. Consequently, there is an identifiable need for an alternative solution that is autonomous, self-maintaining, and relevant to the end-consumer.

A community-based, video information-editing Web platform exposes and associates descriptive and searchable text with any digital video file. This is accomplished using digital "tagging" technology. Leveraging online content creation and (CGM) consumer generated media, our solution overcomes the textual barriers that have rendered rich-media assets ineffectual in terms of their market potential.

The signs of vast opportunity are readily apparent: adding rich-media to advertisements doubles user activity rate per impression[87]; 44% of Internet users have helped create content for the online world[88]; 60% of consumers trust other consumers' online postings[89].

### 7.1: Target Market

Initially, the proposed solution is to target online sports video content providers because of their large fan-base and content paying audience. The sports industry is one of the largest and fastest growing industries in the United States: last year the size of the industry was estimated to be $213 billion[90]. Sports entertainment is more than twice the size of the U.S. auto industry and is seven times the size of the movie industry[91]. Broadcast advertising is the largest source of the revenue for the industry, consisting of $27 billion[92]. Internet revenue accounts for $239.1 million, which includes $230 million in advertisement revenue and $9 million in subscription

service[93]. Spectator spending is valued at $26.17 billion, sports multimedia (including DVD's, magazines, etc) is valued at $2.12 billion, and internet gambling and legal sports booking is valued at $9.7 billion[94]. Although the primary focus is in multimedia subscription, our technology enables greater marketing and sales growth in all of these revenue generating activities. By enabling faster and more relevant access to desired digital sports content, content providers can lead users more precisely to purchases of DVD's, tickets and memorabilia. Advertisers can reach their target audience more efficiently based on profiling data the proposed solution provides. The potential market for solution providers only grows with the increase in broadband penetration and rich-media use.

Besides the spectacular revenue growth of the sports industry, the initial focus is primarily the sports-enthusiast market segment because of their high involvement with digital content. Sports content has a very strong appeal to the online male demographic: more than 38% of male online users visit sports-related Web sites monthly or more frequently, making it the second most popular Internet activity[95]. While visiting online sporting Web sites, users' main activities include finding scores and statistical information, reading recaps for sporting events, participating in fantasy-sports leagues, and watching video highlights or live videos[96]. In summary, the online sports fan is the ideal community member to help propagate the tagging information.

The community-based tagging technology will leverage the sports fan, online multiplayer game-players and other niche aficionados to help populate rich metadata such as scores, rankings and game statistics. Focusing initially on sports-media and its devoted following will help to ensure a successful market entry into online video distribution.

According to this market-penetration strategy, the adoption of the proposed solutions will first be by online sports information providers such as MLB and ESPN, followed by mass-market content distributors and aggregators such as CNN. Depending on the evolution of the online-video and television markets, late adopters would include feature film distributors and other television content providers. The graph below shows the adoption curve for the solution providers.

-Online News Broadcast (e.g. CNN)
-Content aggregator (Google, Yahoo)

-Feature films
-Episodic TV shows

- MLB.com /
ESPN.com Online
Sport broadcast
- New Internet TV
(e.g. Bridghtcove)

Adoption / Innovation Curve

Early
Majority

Early
Adopters

Late
Majority

2.5%    13,5%    34%    34%    16%

Innovators              Laggards
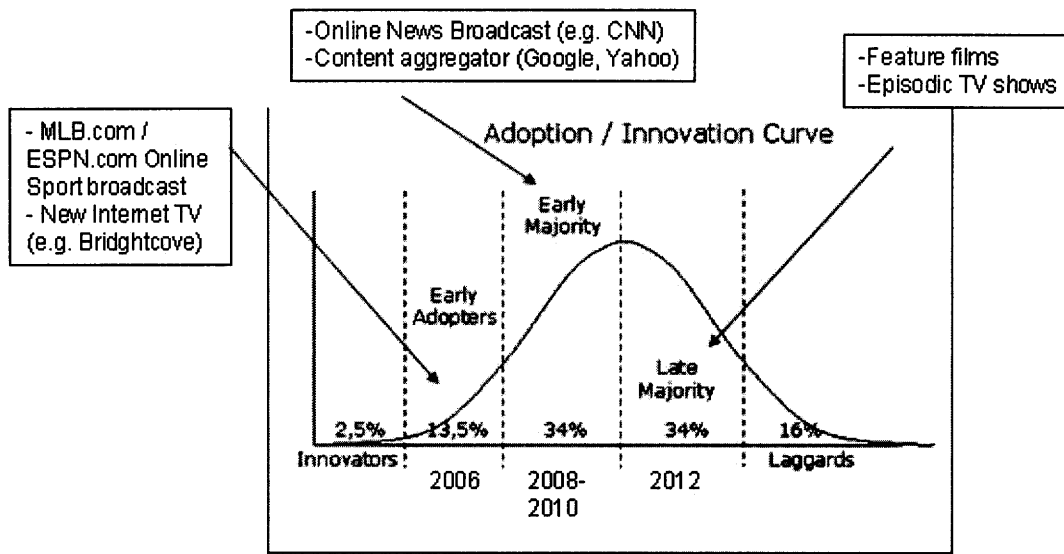
2006    2008-    2012
        2010

Figure 17: Video Tagging Technology Adoption Curve

## 7.2: Business Model

A strategic two-phased approach is devised for generating short-term and long-term growth and sales. Initially, the digital video asset community is charged through a licensing and pay-per-performance pricing model. Included with the licensed software will be customization services specific to each customer. In the second phase of the growth model, the solution providers will position itself as a service company by providing specific search capabilities to both business clients and consumers based on metadata (information about the video content) generation.

The solution providers will roll out their products and services by focusing on one regional market segment at a time. This phased approach is necessary to build up a customer base and will help increase the scalability of solution providers' operating and the marketing efforts. The first phase of the marketing strategy began in 2006 and will focus on professional sports organizations and their online operations. This includes Web sites such as MLB.com, NBA.com and NFL.com. These companies lend themselves specifically to proposed solution due to their avid statistical fan bases, which can help propagate video metadata quickly and accurately.

Once a reference customer base has been established, the second phase of the marketing strategy will commence: solution providers will focus on entering the emerging Internet TV market segment, a space currently occupied by companies such as Brightcove and Revver. Internet TV companies will need to incorporate effective search capabilities into their digital video content to stay competitive with major search and digital video providers such as Google, Yahoo, and MSN. The second phase of the marketing plan will also focus on new geographical markets, covering customers located on the West Coast. Target segments will include entertainment networks and news broadcasting companies who are clustered throughout Los Angeles and the greater San Francisco Bay area.

As solution providers continue to increase their client base, the company will shift focus to marketing product services. This will allow the solution providers to penetrate new market opportunities with corporations who organize and distribute digital rich media assets. Because the solution providers will initially rely on one universal technology platform, initial software development costs can be distributed across all clients, facilitating a service transition. Operational costs will be optimized by leveraging economies of scale and scope.

For initial customers, the solution providers will offer discounted licensing fees but also provide pay-for-performance pricing strategies to help mitigate risk for the client. Such a pricing strategy reduces risk for early adopters and gives the solution providers an important portfolio of blue-chip clients.

Original Equipment Manufacturer

> The solution providers' marketing strategy relies heavily on an OEM (Original Equipment Manufacturer) model targeting major digital media aggregators across multiple sectors, including Google, iTunes, Major League Baseball, YouTube, and Broadway Video. Allowing these companies to keep their existing webpage layouts, the solution providers will offer video-tagging and metadata management solutions using a licensing fee and cost-per-action pricing model. A customer action could be a purchase, download, registration, or click-through use of our solution. Pay-for-performance pricing allows profit and risk sharing between clients and the solution providers, and ensures high client ROI. By allowing customers to generate time-consuming, cost-prohibitive tagging information, businesses can better deliver specific advertising and market relevant content to consumers. For those rich media assets that do not generate a sufficient amount of customer activity or interest, a rewards program will be established to ensure adequate metadata content creation.

Software as a Service

> The SaaS (Software as a Service) model is a long-term growth strategy designed to establish an online presence for service provision through a Web-based application platform, similar to the outsourcing service model used by Salesforce.com. Under the SaaS model, digital video content would be hosted and stored on the clients' server while the solution providers' Web services provide the functionality for customers to associate and host metadata within the video content. Since video content is the essential commercial asset for content owners and providers, the fact that the client continues to host the content is of high importance. While having the client host the video content is ideal, the solution providers will offer customers external hosting of metadata information using the solution provider's own storage and servers.

> The SaaS model also includes an online community site for users to share, view and tag videos. Video content will be provided by community members themselves. Eventually, this video file-sharing will be commercialized to an online marketplace, where users can sell their own consumer generated

media. A recommendation system will be incorporated based on user activity and the popularity of certain content. The recommendation system will return relevant user results and dynamically propagate targeted content on specific user pages. For those media companies who choose to make metadata visible to the community, their content can be included in the same community recommendation system. Corporate customers will enjoy substantial ROI with exposure to the large membership community.

Finally, the proposed solution will incorporate non-intrusive rich media as well as text-based advertising based on an efficient recommendation system and search results. Therefore, revenue streams will be based on corporate sales, advertising revenue from the Web site and commission fees from community video marketplaces.

7.3: Value Chain

The proposed solution's strategic goal: to establish a leading a community of service providers of digital metadata solutions and improved search technology. To accomplish these goals, the initial focus is a direct OEM sales approach, targeting businesses that have large digital video content libraries with active and enthusiastic customers.

Utilizing initial advisor contacts, the founders will help form the original sales force. The group will specifically target major sports organizations with a strong online presence, including MLB.com, NBA.com and NFL.com. Leveraging advisors who are familiar with the sporting and digital asset communities, solution providers will focus on establishing a relationship with these firms through accredited introductions. The sales force will focus specifically on establishing a rapport with the Senior Vice President of Technology in these firms, since these officials represent the critical decision making entity for these types of firms. By initially signing sporting companies with a major Internet presence and large digital asset libraries, the solution providers will gain substantial credibility to help enter markets where customers may be more skeptical of new technology.

The value proposition to potential customers is to address the client's pain points as follows:

1. **Increased sales:** latent digital video assets are now marketable
2. **Increased site traffic:** digital video assets can be indexed into search engines and generate site traffic at a fraction of traditional advertising click-through costs
3. **Additional advertising revenue:** due to increased search traffic
4. **Improved search:** more relevant search results provided by metadata information
5. **Increased customer loyalty:** interactive nature of metadata use and generation creates more involved customers
6. **Reduced costs:** labor-intensive metadata management has been outsourced to the consumer

With the direct sales approach, the solution providers will not need to form relationships with retailers, distributors or wholesalers. Rather, the digital video industry's value-chain includes the following key stakeholders: content creators, content aggregators, search providers and end-consumers. The proposed solution is poised as a search-provider alternative that is using technology and community involvement to lower the cost of click-through traffic to a company's Website.

Additionally, the proposed solution increases the total profit of entire value chain by enabling additional revenue and sales opportunities through improved search capability of latent digital assets. The metadata information hosting will allow the solution providers to secure competitive advantages and establish barriers to entry in the market. Since content creators and content aggregators have many players in each sector, competition among companies will result in high demand for the metadata technology and strong leverage in negotiation situations.
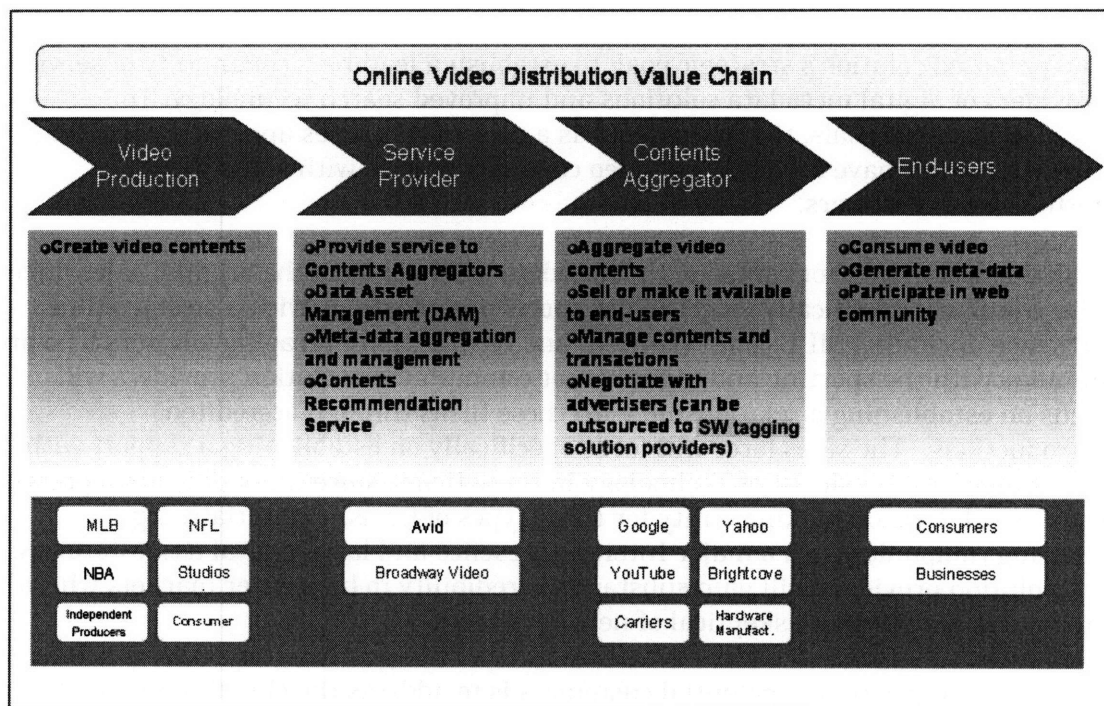


Figure 18: Online Video Distribution Value Chain

## 7.4: Financial Analysis

Pricing is modeled after the consumer benefit strategy. Specifically, a licensing model combined with a revenue sharing approach is used to ensure adequate returns to both the customer and the solution providers. A licensing model helps ensure cash flow during the initial start-up phase of the company, whereas a revenue sharing approach helps to dilute risk for both parties. Concerning revenue sharing, the consumer does not have to risk heavy, up-front costs for technology that might not create any benefit for the company. The solution providers, on the other hand, will be able to take advantage of rapid growth and consumer use that will lead to significant revenue.

By comparing the number of successful hits using the proposed solution and metadata versus the number of successful hits using the company's old search technology, the solution providers can determine exactly the added search benefit of the proposed solution. Existing cost rates for click-through revenues are set as an industry standard at approximately five cents per click-through[97]. The solution providers' initial clients will feature Web-pages with heavy daily traffic (characterized here as more than three million page views per day). Using the proposed solution, it is estimated an additional 50% of users will benefit from better search results, leading to more significant click-through on Web sites. With all of these variables, a price-by-benefit model can be created.

The solution providers will charge 50% of the five cents revenue sharing fee for its improved search capability. This means the company will earn an additional one cent for every click-through search, and also provides an easy means of creating a value proposition for clients. Given an estimate of three million daily page views generated on the customer's Website, and a click-through improvement of 50%, this means the proposed solution will generate 5,475,000 click-through clients per year[98]. At 2.5 cents per click-through, the solution providers will thus earn $136,875 per year of use for each large client (see Figure 15). While the revenue is heavily volume related, these revenue streams will only become more consistent as more clients are obtained.

Because the proposed revenue sharing model is the same model used by search and advertising competitors, the solution providers can leverage the above established pricing precedents when marketing its own product. Providing services at only a fraction of the benefit ensures an easily identifiable value proposition for customers. The high-volume nature of the metadata technology enables the proposed model sustainable, and the fact that this metadata information is generated free of charge by the consumer.



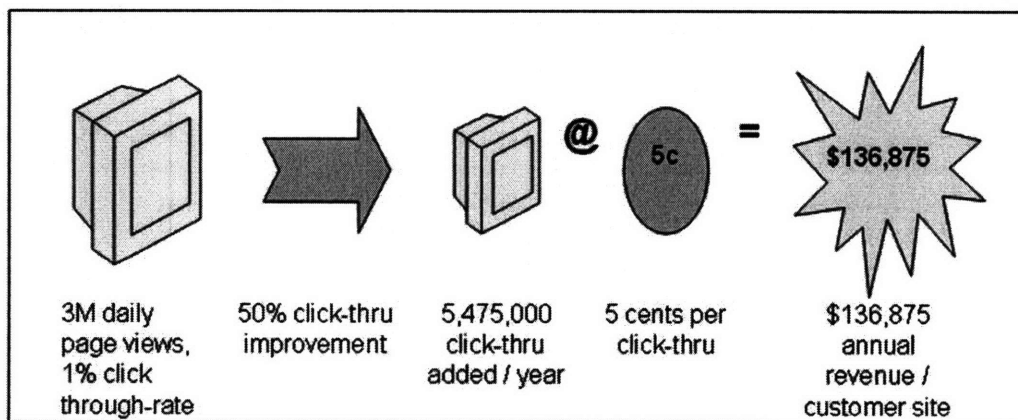| 3M daily page views, 1% click through-rate | 50% click-thru improvement | 5,475,000 click-thru added / year | 5 cents per click-thru | $136,875 annual revenue / customer site |

Figure 19: Revenue Sharing Model

Another analysis is conducted based on the large online social networking site MySpace who reported an amazing figure of 1.5 billions daily page views in

February 2006[99]. In addition, comScore Media Metrix, the leader in digital media measurement, also reported that MySpace led in the number of stream among U.S. Internet users. MySpace ranked first among all sites in individual video streams initiated by U.S. users with nearly 1.5 billion streams followed by Yahoo and YouTube with 812 million and 649 million respectively[100].

Based on the above revenue scheme, projection is performed based on data as of Aug 2006 with the following parameters: 1.5B pageviews/day, 34.1% sales rate (comScore Networks), $0.05 per click through rate at 1%. The result shows that solution providers could enjoy as much as $70 million in revenue per year.
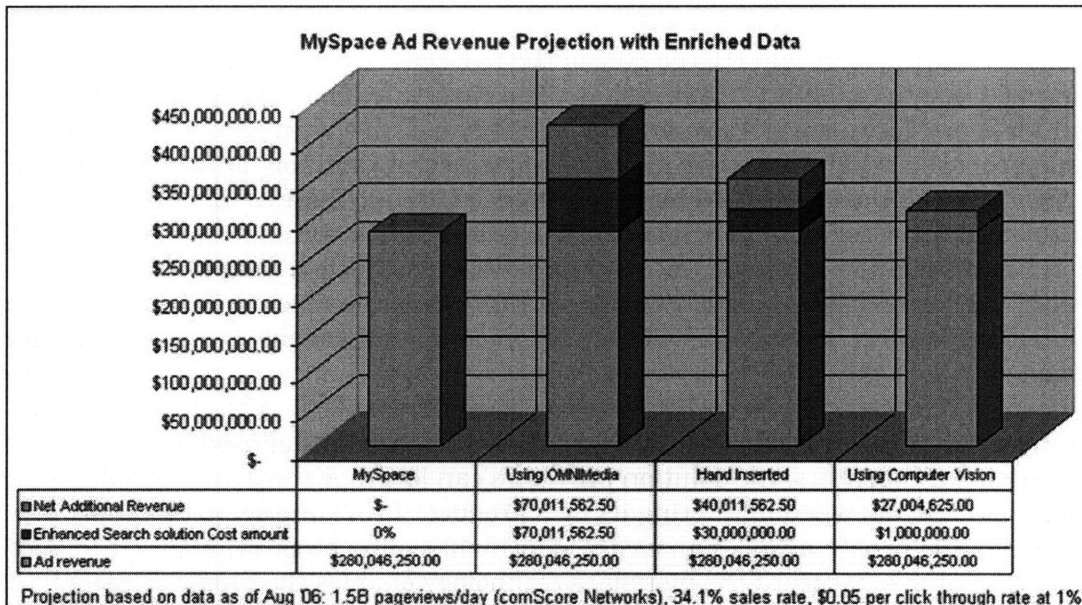
**MySpace Ad Revenue Projection with Enriched Data**

| | MySpace | Using OMNIMedia | Hand Inserted | Using Computer Vision |
|---|---|---|---|---|
| ▣ Net Additional Revenue | $- | $70,011,562.50 | $40,011,562.50 | $27,004,625.00 |
| ▬ Enhanced Search solution Cost amount | 0% | $70,011,562.50 | $30,000,000.00 | $1,000,000.00 |
| ▣ Ad revenue | $280,046,250.00 | $280,046,250.00 | $280,046,250.00 | $280,046,250.00 |

Projection based on data as of Aug '06: 1.5B pageviews/day (comScore Networks), 34.1% sales rate, $0.05 per click through rate at 1%

Figure 20: Advertising Revenue Comparison between Proposed Solution and Other Technologies

| Potential Ad Impression for Revenue | MySpace Revenue Comparison (Aug '06) | | | |
|---|---|---|---|---|
| | MySpace | using Semantic Web solution | Hand Inserted | Using Computer Vision |
| Page viewed / day | 1,500,000,000 | 1,500,000,000 | 1,500,000,000 | 1,500,000,000 |
| Page viewed / year | 547,500,000,000 | 547,500,000,000 | 547,500,000,000 | 547,500,000,000 |
| View Impression per Ad served | 547,500,000,000 | 547,500,000,000 | 547,500,000,000 | 547,500,000,000 |
| Ad space # | 3 | 3 | 3 | 3 |
| Sales rate: | 34% | 34% | 34% | 34% |
| # of view impression sold | 560,092,500,000 | 560,092,500,000 | 560,092,500,000 | 560,092,500,000 |
| cost rates for click-through | $0.12 - $0.16 | $0.12 - $0.16 | $0.12 - $0.16 | $0.12 - $0.16 |
| avg rate | $0.05 | $0.05 | $0.05 | $0.05 |
| click thru rate | 1% | 1% | 1% | 1% |
| ttl clicks | 5,600,925,000 | 5,600,925,000 | 5,600,925,000 | 5,600,925,000 |
| Ad revenue | $ 280,046,250.00 | $ 280,046,250.00 | $ 280,046,250.00 | $ 280,046,250.00 |
| Click-through improvement with more targeted Ad from metadata | 0% | 50% | 25% | 10% |
| Additional ttl clicks | 0 | 2,800,462,500 | 1,400,231,250 | 560,092,500 |
| Additional Ad Revenue | | $ 140,023,125.00 | $ 70,011,562.50 | $ 28,004,625.00 |
| Enhanced Search solution Cost (%) | 0% | 50% | 43% | 4% |
| Enhanced Search solution Cost amount | | $ 70,011,562.50 | $ 30,000,000.00 | $ 1,000,000.00 |
| Net Additional Revenue | $ - | $ 70,011,562.50 | $ 40,011,562.50 | $ 27,004,625.00 |
| Rev / year | $ - | $ 70,011,562.50 | | |
| Rev / day | $0.00 | $191,812.50 | | |

Table 1: Advertising Revenue Comparison between Proposed Solution and Other Technologies

## Chapter 8: Conclusion, Future Work and Research

Like many other new technologies and innovation, the Semantic Web is poised to provide a very promising future. If succeeded, the impact that it makes on the future of communications and information technologies would be profoundly significant. The financial analysis described in previous chapter demonstrates lucrative financial returns based on the market analysis. A key insight gained from the research for this thesis is that focusing and pursuing an appropriate strategy and market would be crucial for the Semantic Web to reach the critical mass, build viable businesses and profit from the extensive research effort. Research that generates the right technologies at the right time is critical to competitive success in many industries. Over the long term, research into the breakthrough technologies can make products perform better, cost less, and generate attractive profits. Unfortunately, the promise and potential cannot always be realized for many companies who invest in the R&D. Among many lessons, the main one stresses the importance of finding and focusing on the right job for the technology.

While this thesis has analyzed and identified a promising strategy for commercialization and a pragmatic adoption path for the Semantic Web, the challenge ahead is to gain sponsorship and to form strategic collaborations across research, academia and media industry. Although the media industry is recognizing the significance of online channel and is very serious and open to investing more in the Web technologies, their development cycle is short and their focus tends to be on the near-term result which does not align with the longer lifecycle in research projects. As briefly mentioned in Chapter 6.2, media companies, who has initiated the OMNI project, hesitate to collaborate formally with pure research groups like the Semantic Web group at W3C. In parallel, research team leaders of the Semantic Web project including Tim Berners-Lee are less interested in commercial projects which have close tie with specific industry. The research group's general- purpose focus and neutral orientation enables them to continue to drive the standardization effort as a central authority. Yet the pure standardization effort without pragmatic adoption in the commercial capacity could slow down the progress. For the next step, it is important to synergize and create a healthy collaboration environment across different commercial and research entities involved in the Semantic Web.

Another recommendation for future work would be to analyze the potential of starting the groundwork in regions which have less established Web infrastructure and development. For example, although India and China are technologically very advanced, their Web infrastructure and e-commerce environment is still less sophisticated and is behind the US. Furthermore, from an empirical observation, the Internet penetration and e-commerce transaction volume are relatively low in many parts of European countries and cities including Italy. Together, they provide a nice platform for experiments, research and development for the Semantic Web. Those regions provide the environment for the Semantic Web to build from ground up, establish the reputation and solid reference cases, and leapfrog the adoption outside of the US.

# Bibliography

1 Tim Berners-Lee, Weaving the Web. San Francisco, CA, HarperSanFrancisco, 1999, chapter 1.

2 http://www.w3.org/2000/01/sw/

3 http://www.xml.com

4 http://www.altova.com/semantic_web.html

5 Interview with Prof. David Karger, Principal Investigator in the Haystack Research at CSAIL, MIT

6 http://haystack.csail.mit.edu/

7 http://simile.mit.edu/

8 http://www.usatoday.com/tech/webguide/internetlife/2004-10-01-cover-web_x.htm

9 Tim O'Reilly, Tim O'Reilly in a Nutshell. Sebastopol, CA, O'Reilly Media, Inc., 2004, p41.

10 Tim O'Reilly, Tim O'Reilly in a Nutshell. Sebastopol, CA, O'Reilly Media, Inc., 2004, p61.

11 http://en.wikipedia.org/wiki/Screen_scraping

12 William T. Shelton, Web Services: A Strategic Analysis. Master's Thesis, MIT Sloan School of Management. Cambridge, MA. 2003

13 John Hagel III, Out of the box, strategies for achieving profits today and growth tomorrow through web services. Boston, MA, HBS Publishing, 2002, p.70.

14 John Hagel III, Out of the box, strategies for achieving profits today and growth tomorrow through web services. Boston, MA, HBS Publishing, 2002, p.28

15 Tim O'Reilly, Tim O'Reilly in a Nutshell. Sebastopol, CA, O'Reilly Media, Inc., 2004, p68

16 http://www.trulia.com/

17 http://www.timeinc.net/b2/subscribers/articles/print/0,17925,1019738,00.html

18 Tanja Sollazzo, Seigfried Handschuh, Steffen Staab, and Martin Frank, Semantic Web Service Architecture, University of Koblenz-Landau. 2002.
http://www.aifb.uni-karlsruhe.de/~sst/Research/Publications/sub-flairs2002.pdf

19 Don Clausing and Victor Frey, Effective Innovation. ASME Press (US), 2004.

20 Rebecca M. Henderson and Kim B. Clark. 1990. Architectural innovation: The reconfiguration of existing product technologies and the failure of established firms. Administrative Science Quarterly.

21 Clayton M. Christensen and Michael E. Raynor, The Innovator's Solution. Boston, MA, Harvard Business School Press, 2003, p.32.

22 Clayton M. Christensen and Michael E. Raynor, The Innovator's Solution. Boston, MA, Harvard Business School Press, 2003, p.34

23 Clayton M. Christensen and Michael E. Raynor, The Innovator's Solution. Boston, MA, Harvard Business School Press, 2003, p.13 and p. 17.

24 Script of Interview with Clay M. Christensen, by Gartner Fellow, Howard Dresner, April 26, 2004

http://www.gartner.com/research/fellows/asset_93329_1176.jsp

25 Clayton M. Christensen and Michael E. Raynor, The Innovator's Solution. Boston, MA, Harvard Business School Press, 2003, p.83

26 Chen, Anne, "Semantic Web is 2 Steps Closer," DevSource (http://www.devsource.com/article2/0,1895,1621521,00.asp) July 6, 2006.

27 Geoffrey A. Moore, Crossing the Chasm. New York, NY, HarperCollins, 1999, p41.

28 Tim Bray's Blog, Ongoing, (http://www.tbray.org/ongoing/When/200x/2003/05/21/RDFNet)

29 http://en.wikipedia.org/wiki/Joseph_Swan

30 http://www.maxmon.com/1878ad.htm

31 Michael A. Cusumano, The Business of Software. New York, NY, Free Press, 2004, p.55

32 Grant, Peter, "Online Video Goes Mainstream, Sparking an Industry Land Grab," The Wall Street Journal Online 21 February 2006.

33 Geoffrey A. Moore, Crossing the Chasm. New York, NY, HarperCollins, 1999, p63.

34 HBS class, Building Successful and Sustaining Enterprise by Clay Christensen, September 18, 2006

35 Interview with Prof. Stuart Madnick by the author, August, 2006

36 Interview with Prof. Stuart Madnick by the author, August, 2006

37 http://en.wikipedia.org/wiki/Network_effect

38 Geoffrey A. Moore, Crossing the Chasm. New York, NY, HarperCollins, 1999, Chapter 3

39 Geoffrey A. Moore, Crossing the Chasm. New York, NY, HarperCollins, 1999, p67.

40 What's Web 2.0? by Tim O'Reilly (http://www.oreillynet.com/pub/a/oreilly/tim/news/2005/09/30/what-is-web-20.html)

41 David Provost, "Hurdles in the Business Case for the Semantic Web" MIT Sloan School of Management, p13

42 http://www.pandora.com/

43 http://www.pandora.com/mgp.shtml

44 "CGM sites tempt brands to share control of image: advertisers weigh risk, reward with consumer-generated media," Adweek 13 February 2006: 47.

45 http://www.trulia.com/

46 Thomas W. Malone, The Future of Work. Boston, MA, HBS Press, 2004, Chapter 4.

47 http://www.flickr.com/photos/stabilo-boss/sets/72057594060779001/

48 http://www.fourio.com/web20map/

49 Nielsen//NetRatings, "U.S. Broadband Connections Reach Critical Mass," NetRatings, Inc. 18 August 2004.

50 Grant, Peter, "Online Video Goes Mainstream, Sparking an Industry Land Grab," The Wall Street Journal Online 21 February 2006.

51 Author's interview with Prof. Henry Jenkins, November 28, 2006

52 http://en.wikipedia.org/wiki/Moral_economy

53 http://www.google.com/press/pressrel/google_youtube.html

54 Tim O'Reilly, Tim O'Reilly in a Nutshell. Sebastopol, CA, O'Reilly Media, Inc., 2004, p68

55 http://en.wikipedia.org/wiki/Long_tail

56 Tim O'Reilly, Tim O'Reilly in a Nutshell. Sebastopol, CA, O'Reilly Media, Inc., 2004, p76

57 MIT Sloan class, 15.358 Software Business, lecture 1's note.

58 Clayton Christensen, Christopher Musso and Scott Anthony, Capturing the Returns From Research, HBS

59 Michael A. Cusumano, The Business of Software. New York, NY, Free Press, 2004, p.25.

60 Michael A. Cusumano, The Business of Software. New York, NY, Free Press, 2004, p.102.

61 http://en.wikipedia.org/wiki/Web_2

62 http://www.clickz.com/showPage.html?page=1402241

63 Geoffrey A. Moore, Crossing the Chasm. New York, NY, HarperCollins, 1999, p12.

64 Tim O'Reilly, Tim O'Reilly in a Nutshell. Sebastopol, CA, O'Reilly Media, Inc., 2004, p68

65 Michael A. Cusumano, The Business of Software. New York, NY, Free Press, 2004, p.122

66 Interview with Clay Christensen by the author, Nov 30, 2006.

67 67 Clayton M. Christensen and Michael E. Raynor, The Innovator's Solution. Boston, MA, Harvard Business School Press, 2003...

68 Clayton Christensen, Christopher Musso, and Scott Anthony, "Capturing the Returns from Research." Harvard Business School

69 Clayton Christensen, Christopher Musso, and Scott Anthony, "Capturing the Returns from Research." Harvard Business School, p3

70 Clayton Christensen, Scott Anthony, Gerald Berstell and Denise Nitterhouse, "Finding the Right Job for your Product." Harvard Business School

71 John Hagel III, Out of the box, strategies for achieving profits today and growth tomorrow through web services. Boston, MA, HBS Publishing, 2002, p.70

72 http://en.wikipedia.org/wiki/Folksonomy

73 Information obtained in personal interview with the GM of digital media services at Broadway Video (a Lorne Michaels company) and the SVP of multimedia and technology at Major League Baseball (MLB.com)

74 http://money.cnn.com/magazines/business2/business2_archive/2006/07/01/8380226/index.htm

75 http://online.wsj.com/article/SB115490043035628128.html?mod=home_whats_news_us

76 http://en.wikipedia.org/wiki/Folksonomy

77 Stanley Bileschi, interview by author, Cambridge, MA, July 25, 2006

78 Image provided by Stanley Bileschi

79 Kristen Grauman and Trevor Darrell. Pyramid Match Kernels: Discriminative Classification with Sets of Image Features. MIT-CSAIL-TR-2006-020. March 18, 2006.

80 Data provided by Stanley Bileschi

81 Image provided by Stanley Bileschi, Lior Wolf

82 Image provided by Stanley Bileschi, Lior Wolf

83 http://eigencluster.csail.mit.edu//cgi-bin/main.cgi?query=apple

84 Nielsen//NetRatings, "U.S. Broadband Connections Reach Critical Mass," NetRatings, Inc. 18 August 2004.

85 "Digital Media Asset Management & Workflow Management in the Broadcast Industry," Multimedia Research Group, Inc April 2004.

86 Hallerman, David, "Why CPGs Need Rich Media," eMarketer 2 September 2003.

87 Hallerman, David, "Why CPGs Need Rich Media," eMarketer 2 September 2003.

88 Lenhart, Amanda, et al., "Content Creation Online," Pew Internet & American Life Project 29 February 2004.

89 Blackshaw, Pete and Mike Nazzaro, "Consumer-Generated Media (CGM) 101," Intelliseek Spring 2004.

90 "Sports Business Marketplace",Smith's SportsBusiness Journal 2005.

91 "Sports Business Marketplace",Smith's SportsBusiness Journal 2005.

92 "Sports Business Marketplace",Smith's SportsBusiness Journal 2005.

93 "Sports Business Marketplace",Smith's SportsBusiness Journal 2005.

94 "Sports Business Marketplace",Smith's SportsBusiness Journal 2005.

95 Matiesanu, Bayriamova, Card, Laszlo and Peach " US Online User Consumer Survey, 2005" Jupitar Research Group, Inc January 2006.

96  Card, Widger, Matiesanu and McLeary "US Fantasy Sports - Programming for and Marketing to Hardcore Fans" Jupitar Research Group, Inc December 2005.

97 Click Affiliate.com. Online. Internet. 9 April 2006. <http://www.clickaffiliate.com/affiliate/advertising/advertising_index.shtml/>.

98 Click Affiliate.com. Online. Internet. 9 April 2006. <http://www.clickaffiliate.com/affiliate/advertising/advertising_index.shtml/>.

99 http://weblogs.asp.net/scottgu/archive/2006/03/25/441074.aspx

100 http://www.comscore.com/press/release.asp?press=1015