

INNOVATIONS IN THE INTERNET'S ARCHITECTURE THAT CHALLENGE THE STATUS QUO

CHRISTOPHER S. YOO*

INTRODUCTION	79
I. THE ARCHITECTURE OF THE EARLY INTERNET	81
A. <i>The Topology of the Early Internet</i>	81
B. <i>Business Relationships on the Early Internet: Peering and Transit</i>	84
II. THE EVOLUTION OF THE INTERNET'S TOPOLOGY	85
A. <i>Private Peering, Multihoming, and Secondary Peering</i>	86
B. <i>Server Farms and Content Delivery Networks</i>	88
III. THE EVOLUTION OF BUSINESS RELATIONSHIPS	90
A. <i>The Growing Importance of Peer-to-Peer Architectures</i>	90
B. <i>The Emergence of Partial Transit and Paid Peering</i>	95
CONCLUSION	99

INTRODUCTION

Despite having received sustained attention from both policymakers and academic commentators for the past several years, network neutrality shows no signs of retreating from the forefront of the policy debate. It has remained a central focus for Congress,¹ the Federal Communications Commission (FCC),² and both presidential candidates during the last election.³ As President, Barack Obama has effectively ensured that network neutrality will remain at the top of the policy agenda by including provisions in the stimulus package requiring that the FCC

* Professor of Law and Communication and Founding Director, Center for Technology, Innovation, and Competition, University of Pennsylvania. The author thanks the Milton and Miriam Handler Foundation for its financial support.

1. See *The Internet Freedom Preservation Act of 2008: Hearing on H.R. 5353 Before the Subcomm. on Telecomm. and the Internet of the H. Comm. on Energy & Commerce*, 110th Cong. (2008); *Net Neutrality and Free Speech on the Internet: Hearing Before the Task Force on Competition Policy and Antitrust Laws of the H. Comm. on the Judiciary*, 110th Cong. (2008).

2. See Formal Complaint of Free Press and Public Knowledge Against Comcast Corp. for Secretly Degrading Peer-to-Peer Applications, *Memorandum Opinion & Order*, 23 FCC Rcd. 13,028 (2008); *En Banc Hearing on Broadband Network Management Practices Before the FCC* (Feb. 25, 2008), <http://www.fcc.gov/realaudio/mt022508v.ram>.

3. See Lee Gomes, *Debugging Obama-McCain*, FORBES, Oct. 13, 2008, at 72.

formulate a national broadband plan and through requiring that grants made by the National Telecommunications and Information Administration comply with the network neutrality principles articulated by the FCC in 2005.⁴

Although pinning down a precise definition of network neutrality has proven elusive,⁵ the most common position appears to be that network providers should route traffic without regard to the source or content of the packets, the application with which the packets are associated, or the sender's willingness to pay. In the words of leading network neutrality proponent Lawrence Lessig, "Net neutrality means simply that all like Internet content must be treated alike and move at the same speed over the network."⁶

Some commentators have questioned whether this description of network neutrality represents an accurate description of the Internet's past.⁷ Indeed, it would be surprising if any two similar packets would be treated exactly alike when traveling through a network consisting of more than thirty thousand autonomous systems that each determine their terms of interconnection through arms-length negotiations. There are, however, some systematic changes in the architecture of the Internet that have largely been overlooked by both commentators and policymakers. These changes are largely the result of network providers' attempts to reduce cost, manage congestion, and maintain quality of service.

4. American Recovery and Reinvestment Act of 2009, Pub. L. No. 111-5, § 6001(j)–(k), 123 Stat. 115, 515–16.

5. See Rachelle B. Chong, *The 31 Flavors of Net Neutrality: A Policymaker's View*, 12 INTELL. PROP. L. BULL. 147, 151–55 (2008) (identifying five distinct versions of network neutrality); Eli Noam, *A Third Way for Net Neutrality*, FIN. TIMES, Aug. 29, 2006, <http://www.ft.com/cms/s/2/acf14410-3776-11db-bc01-0000779e2340.html> (identifying seven distinct versions of network neutrality).

6. See, e.g., Lawrence Lessig & Robert W. McChesney, *No Tolls on the Internet*, WASH. POST, June 8, 2006, at A23.

7. See, e.g., Robert W. Hahn & Robert E. Litan, *Portioning Bit by Bit: The Myth of Network Neutrality and the Threat to Internet Innovation*, MILKEN INST. REV., 1st Qtr. 2007, at 28, 31–33; Jonathan E. Nuechterlein, *Antitrust Oversight of an Antitrust Dispute: An Institutional Perspective on the Net Neutrality Debate*, 7 J. ON TELECOMM. & HIGH TECH. L. 19, 36–37 (2009); Douglas A. Hass, Comment, *The Never-Was-Neutral Net and Why Informed End Users Can End the Net Neutrality Debates*, 22 BERKELEY TECH. L.J. 1565, 1576–77 (2007); Kai Zhu, Note, *Bringing Neutrality to Network Neutrality*, 22 BERKELEY TECH. L.J. 615, 634–36 (2007); Michael Grebb, *Neutral Net? Who Are You Kidding?*, WIRED, May 31, 2006, <http://www.wired.com/news/technology/internet/0,71012-0.html>; ANDREA RENDA, I OWN THE PIPE, YOU CALL THE TUNE: THE NET NEUTRALITY DEBATE AND ITS (IR)RELEVANCE FOR EUROPE 9–11 (2008), available at http://shop.ceps.eu/download.php?item_id=1755; Craig McTaggart, *Was the Internet Ever Neutral?*, 34 RES. CONF. ON COMM'N, INFO. & INTERNET POL'Y 1, 4–14 (2006), available at <http://web.si.umich.edu/tprc/papers/2006/593/mctaggart-tprc06rev.pdf>; David Clark, Written Statement to the En Banc Public Hearing on Broadband Network Management Practices (Feb. 25, 2008), available at http://www.fcc.gov/broadband_network_management/022508/clark.pdf ("The Internet is not neutral and has not been for a long time.").

Part I frames the subsequent developments by describing the architecture and business relationships that defined the early Internet. Part II analyzes the architectural changes that have made the Internet's topology increasingly heterogeneous, including the emergence of multihoming, secondary peering, private networks, and content delivery networks. Part III describes the changes in ways that networks interconnect and price their services, focusing on the emergence of peer-to-peer applications and pricing innovations that go beyond the traditional bipartite distinction between peering and transit. Far from representing some network provider's efforts to promote its self interest at the expense of the public, as some network neutrality proponents have suggested, these changes have the potential to yield substantial benefits both to individual consumers and to society as a whole.

I. THE ARCHITECTURE OF THE EARLY INTERNET

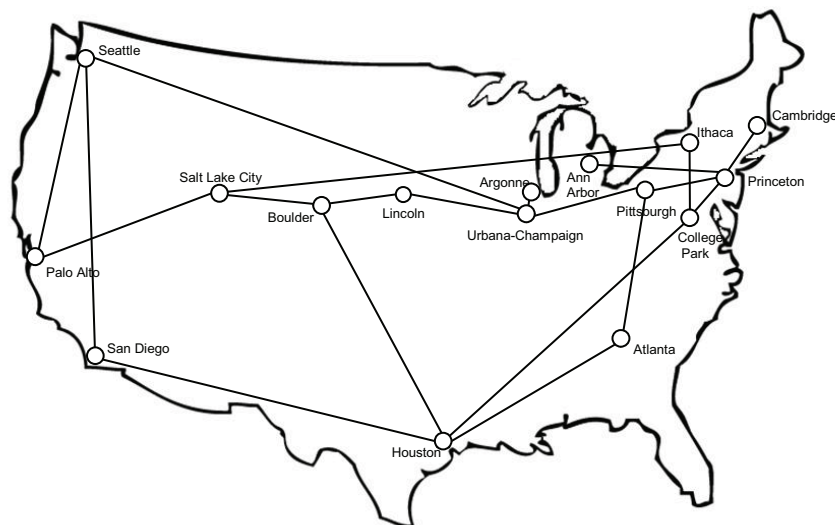
This Part reviews the architecture of the early Internet. Section A reviews the tripartite hierarchical structure that characterized its topology. Section B describes the peering and transit relationships that governed the way individual networks interconnected with one another.

A. The Topology of the Early Internet

When the Internet first emerged, its topology and the business relationships comprising it were relatively simple. As is widely known, the Internet evolved out of the NSFNET backbone, which was created in 1986 and eventually decommissioned in 1997 to provide universities all over the country access to federally funded supercomputing centers located in five universities. The primary architects of the NSFNET decided to give it a tripartite structure. At the top was the NSFNET backbone, which at its peak connected sixteen research facilities across the country. At the bottom were the campus networks run by individual universities. In the middle were regional networks (typically operated by university consortia or state-university partnerships) that linked the campus networks to the major computing centers.⁸

8. MERIT NETWORK, INC., NSFNET: A PARTNERSHIP FOR HIGH-SPEED NETWORKING, FINAL REPORT 1987-1995, at 11-12 (1996), available at http://www.merit.edu/documents/pdf/nsfnet/nsfnet_report.pdf; Juan D. Rogers, *Internetworking and the Politics of Science: NSFNET in Internet History*, 14 INFO. SOC'Y 213, 219 (1998).

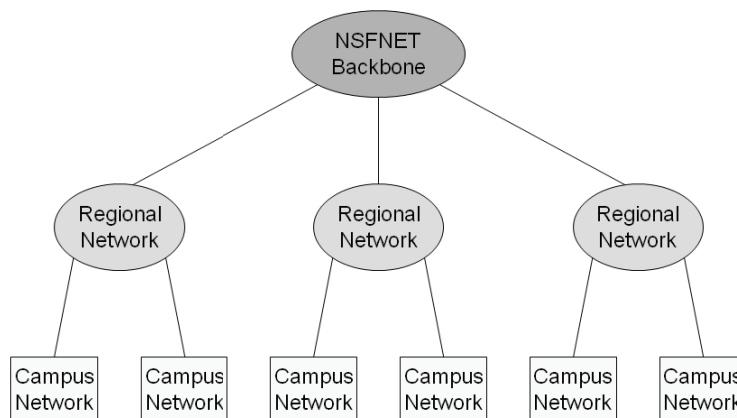
Figure 1: The NSFNET Backbone circa 1992-1993



Every packet had to travel through a parallel path traversing each level of the hierarchy. For example, traffic originating on one campus network would have to connect to the regional network with which it was associated, which handed off the traffic to the NSFNET backbone, which in turn handed it off to the regional network that served the destination campus network. The result was to create a series of parallel hierarchies through which all traffic had to traverse.

The network retained this same basic architecture when it was privatized during the mid-1990s. The NSFNET backbone at the top of the hierarchy was replaced by a series of private backbone providers that interconnected with one another at four public network access points (NAPs) established by the NSF. The campus networks at the bottom of the hierarchy were replaced by last-mile providers that transported traffic from local distribution facilities maintained in individual cities (which in the case of digital subscriber lines (DSL) is usually called a central office and in the case of cable modem systems is usually called a headend) to end users' residences and places of business. The regional networks evolved into regional Internet service providers (ISPs) that transported traffic between the NAPs served by backbone providers and the central offices and headends maintained by last-mile providers.

The privatization of the Internet did not change the hierarchical nature of the basic architecture. Each regional ISP still connected to a single backbone, and each last-mile provider still connected to a single regional ISP. Indeed, the early versions of the protocol employed by the backbones (known as border gateway protocol or BGP) would not

Figure 2: The NSFNET Three-Tiered Network Architecture

support more complex topologies.⁹

The one-to-one relationship conferred a number of advantages. This architecture constituted a “spanning tree” that connected all of the nodes with the minimum number of links.¹⁰ Furthermore, the fact that the path between any two nodes was unique greatly simplified determining the path along which traffic should be routed. That said, tree architectures are also subject to a number of drawbacks. The uniqueness of the path connecting any two nodes means that the failure of any link or node in the network will inevitably disconnect part of the network. Even when all network elements are operating properly, if the rate at which traffic arrives exceeds any particular element’s capacity to route the traffic, that network element will become congested and the quality of service provided will deteriorate.¹¹ In addition, the hierarchical structure made each network participant completely dependent on the players operating at the level above them, which in turn provided backbones with a potential source of market power.¹²

9. Christopher S. Yoo, *Network Neutrality and the Economics of Congestion*, 94 GEO. L.J. 1847, 1860–61 (2006) [hereinafter Yoo, *Economics of Congestion*]; Christopher S. Yoo, *Network Neutrality, Consumers, and Innovation*, 2008 U. CHI. LEGAL. F. 179, 195–96 (2008) [hereinafter Yoo, *Consumers and Innovation*].

10. Daniel F. Spulber & Christopher S. Yoo, *On the Regulation of Networks as Complex Systems: A Graph Theory Approach*, 99 NW. U.L. REV. 1687, 1696 (2005).

11. *Id.* at 1699–700.

12. See Stanley Besen et al., *Advances in Routing Technologies and Internet Peering Agreements*, 91 AM. ECON. REV. (PAPERS & PROC.) 292 (2001).

B. Business Relationships on the Early Internet: Peering and Transit

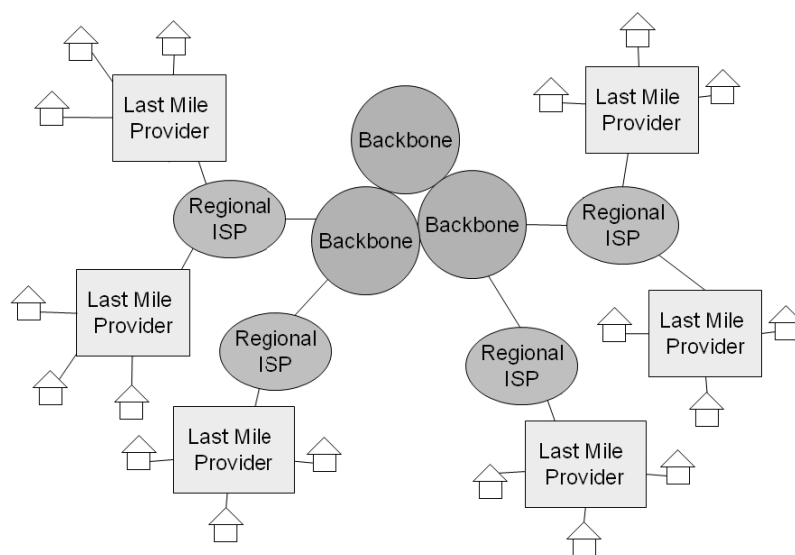
The early Internet was also characterized by relatively simple business relationships. End users typically purchased Internet access through some form of all-you-can-eat pricing, which allowed them to consume as much bandwidth as they would like for a single flat rate. Relationships between network providers typically fell into two categories. Tier-1 ISPs entered into *peering* relationships with one another, in which they exchanged traffic on a settlement-free basis and no money changed hands. The primary justification for foregoing payment is transaction costs. Although the backbones could meter and bill each other for the traffic they exchanged, they could avoid the cost of doing so without suffering any economic harm so long as the traffic they exchanged was roughly symmetrical. Such arrangements would not be economical with when the traffic being exchanged by the two networks was severely imbalanced. Thus tier-1 ISPs will not peer with other networks that are unable to maintain a minimum level of traffic volume. In addition, peering partners typically require that inbound and outbound traffic not exceed a certain ratio. Networks that cannot meet these requirements must enter into *transit* arrangements in which they pay the backbone to provide connectivity to the rest of the Internet.¹³

Most early analyses focused on the financial terms of these arrangements.¹⁴ What is often overlooked is that interconnection agreements performed two distinct functions. Network providers enter into interconnection agreements not only to send and receive traffic. They also enter into interconnection agreements to announce to the rest of the Internet where the IP addresses that they control are located.

Consider this from the perspective of a small network, *A*, which serves a small number of its own customers and purchases access to the rest of the Internet through another ISP. The transit agreement between *A* and the ISP would not only require the ISP to receive traffic sent by *A* and to deliver traffic bound to *A*. It would also require the ISP to announce to the rest of the Internet how to reach the IP prefixes associated with *A*'s customers. In addition, *A* can maintain a very simple routing table. It need only keep track of the prefixes of the customers that it serves. For all other IP addresses, *A* can enter a "default route" into its routing table that directs all other traffic to the other ISP.

13. Yoo, *Economics of Congestion*, *supra* note 9, at 1877; Michael Kende, *The Digital Handshake: Connecting Internet Backbones* 5 (FCC Office of Plans and Policy Working Paper No. 32, 2000), available at http://www.fcc.gov/Bureaus/OPP/working_papers/oppwp32.pdf; Peyman Faratin et al., *The Growing Complexity of Internet Interconnection*, 72 COMM'NS & STRATEGIES 51, 55-56 (2008).

14. *See, e.g.*, Kende, *supra* note 13, at 5.

Figure 3: The Architecture of the Early Internet

The existence of default routes creates a potential problem, however. If none of the routing tables involved in a particular routing session contained the location of the destination, by default the networks would simply hand the packets back and forth, and the packets would never reach their final destination. The only way to avoid these problems is for one or more network providers to maintain routing tables that map the entire Internet without employing any default routes. Thus, tier-1 ISPs are defined not only by the fact that they engage in settlement-free peering with one another, but also by the fact that they maintain routing tables that contain no defaults.¹⁵ Peering contracts also include a number of other requirements to guard against free riding and to ensure the proper functioning of the network.¹⁶

II. THE EVOLUTION OF THE INTERNET'S TOPOLOGY

Over the past decade, ISPs have begun to enter into a more complex set of interconnection arrangements that violate the strict tripartite hierarchy that characterized the early Internet. In addition, content providers have begun to experiment with a variety of ways to locate their content closer to end users. Both types of changes have

15. Paul Milgrom et al., *Competitive Effects of Internet Peering Policies*, in *THE INTERNET UPHEAVAL* 175, 179–80 (Ingo Vogelsang & Benjamin M. Compaine eds., 2000).

16. Faratin et al., *supra* note 13, at 54.

significant policy implications that have largely been overlooked in the policy debate.

A. Private Peering, Multihoming, and Secondary Peering

One of the first problems to emerge in the early Internet was congestion in the NAPs, which often caused throughput times and network reliability to degrade. Some estimate that congestion in the NAPs caused packet loss at times to run as high as 40%.¹⁷ As the NAPs became increasingly congested, backbones began to find it advantageous to exchange traffic at private interconnection points.¹⁸

In addition, regional ISPs have begun to connect to more than one backbone, a practice known as *multihoming*, in part to protect against service outages and in part to limit their vulnerability to any exertion of market power by a backbone.¹⁹ Regional ISPs that did not have sufficient volume to peer with the tier-1 backbones also began to find that they did have sufficient volume to peer with other regional ISPs, a practice known as *secondary peering*. Enabling regional ISPs to exchange traffic on a settlement-free basis reduced the costs borne by end users. In addition secondary peering would often shorten the number of hops needed for particular packets to reach their final destination and make them subject to bilateral (as opposed to multiparty) negotiations, both of which should increase networks' control over quality of service.²⁰ Secondary peering and multihoming also made the network more robust by creating multiple paths through which network nodes could interconnect. In fact, as much as seventy percent of the nodes in the Internet can now communicate with one another without passing through the public backbone.²¹ This had the additional benefit of weakening the market position of the top-tier backbones, since any breakdown in the business relationship would not necessarily disconnect the ISP from the network and the ability to route along different paths places a natural limit on the backbones' ability to engage in supracompetitive pricing.²²

17. See *InterNAP Wakes Up Transmission Quality*, RED HERRING, Apr. 21, 1999, <http://redherring.com/Home/1744>; see also Kende, *supra* note 13, at 6 (citing reports that packet loss in the NAP located in Washington, D.C., ran as high as 20%).

18. Kende, *supra* note 13, at 6–7; Faratin et al., *supra* note 13, at 62.

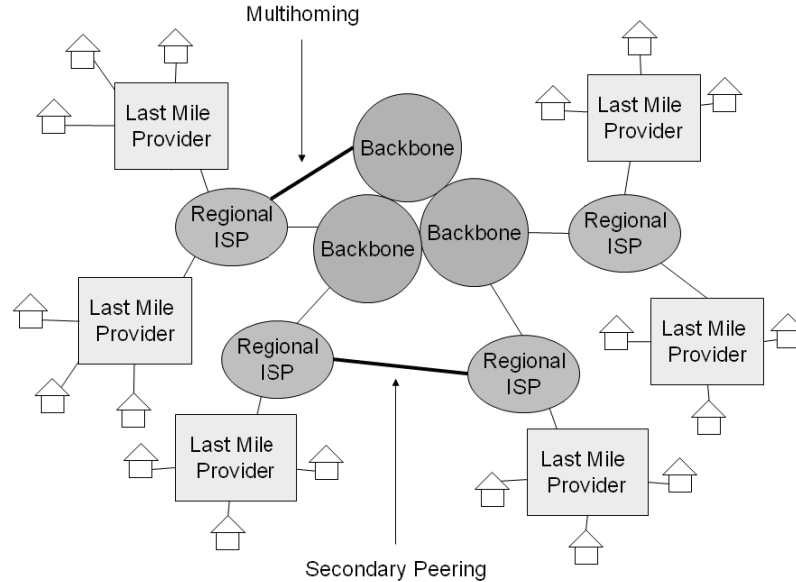
19. See Nicholas Economides, "Net Neutrality," *Non-Discrimination, and Digital Distribution of Content Through the Internet*, 4 I/S: J.L. & POL'Y FOR INFO. SOC'Y 209, 220 (2008).

20. See OECD, WORKING PARTY ON TELECOMMUNICATIONS AND INFORMATION SERVICES POLICIES, INTERNET TRAFFIC EXCHANGE: MARKET DEVELOPMENTS AND MEASUREMENT OF GROWTH 21–22 (2006), <http://icttoolkit.infodev.org/en/Publication.3081.html>; Faratin et al., *supra* note 13, at 55–56.

21. See Shai Carmi et al., *A Model of Internet Topology Using k-Shell Decomposition*, 104 PROC. NAT'L ACAD. SCI. 11,150, 11,151 (2007).

22. See Besen et al., *supra* note 12, at 294–95.

Figure 4: The Emergence of Multihoming and Secondary Peering



The emergence of interconnection relationships that deviate from the strict hierarchy that characterized the early Internet represents a substantial deviation from network neutrality. For example, assume that an end user is downloading content from both CNN.com and MSNBC.com. Assume further that the end user's regional ISP has a secondary peering relationship with the regional ISP serving CNN.com, but does not have a secondary peering relationship with the regional ISP serving MSNBC.com. The absence of a secondary peering relationship means that traffic from MSNBC.com will have to pay transit charges, while traffic from CNN.com will not. The result is that traffic that is functionally identical will end up paying different amounts. The differences in topology may also allow the traffic from CNN.com to maintain greater control over the quality of service.

The presence of multiple routes between these two points also complicates routing decisions. The presence of multiple paths connecting two points naturally means that someone must decide along which path to route the traffic. Although most networks choose routes that minimize the number of hops, networks may sometimes find it beneficial to route traffic in order to satisfy other requirements of their interconnection relationships. For example, a network may seek to enhance efficiency by balancing the loads between the two links. Multihomed entities can also monitor the quality of service provided by each connection and route the most delay-sensitive traffic along the link

with the lowest latency.²³

In addition, transit contracts call for customers to pay a flat fee up to a predetermined peak volume (known as the committed rate) and pay additional charges for any volume that exceeds that level. For the same reason that consumers with two mobile telephones have the incentive to use up all of the prepaid minutes on both lines before incurring any additional per-minute charges, multihomed entities have the incentive to utilize all of their committed rate before paying additional fees. This lowers overall transit cost, but requires diverting some traffic along a path that is longer than the one stored in the routing tables.²⁴ For similar reasons, a network may intentionally route traffic over a more costly path if doing so will help it maintain its traffic within the ratios mandated by its peering contract.²⁵ Again, the effect is to introduce significant variance in the speed with which similarly situated packets will arrive at their destination and the cost that similarly situated packets will have to bear. This variance results not from anticompetitive motives, but rather from networks' attempts to minimize costs and ensure quality of service in the face of a network topology that is increasingly heterogeneous.

B. *Server Farms and Content Delivery Networks*

Large content providers have begun to employ other means to reduce cost and manage latency. One solution is to forego maintaining a single large server and instead to deploy multiple points of presence in carrier hotels across the country. Doing so allows these content providers to avoid paying transit charges to reach the public backbone and instead transmit their traffic through secondary peering arrangements with tier-2 ISPs. Greater reliance on private networks also gives the content providers greater control over network security and performance.²⁶ Indeed, a recent study indicates that Google, Yahoo!, and Microsoft have been able to use server farms to bypass the backbone altogether for roughly a third of their traffic and to keep their number of hops for traffic that had to pass through the backbone to no more than one or

23. Fanglu Guo et al., *Experiences in Building a Multihoming Load Balancing System*, IEEE INFOCOM CONF., 2004, available at http://www.ieee-infocom.org/2004/Papers/26_4.PDF.

24. INTERNAP NETWORK SERVS. CORP., *ECONOMICS OF MULTI-HOMING AND PREMISE-BASED OPTIMIZATION 10* (2008), available at http://internap.com/pdf/white-papers/WP_FCP_Economics_of_MultiHoming_0208.pdf.

25. Faratin et al., *supra* note 13, at 64–65.

26. See Stephanie N. Mehta, *Behold the Server Farm! Glorious Temple of the Information Age!*, FORTUNE, Aug. 1, 2006, available at http://money.cnn.com/magazines/fortune/fortune_archive/2006/08/07/8382587/index.htm; R. Scott Raynovich, *Google's Own Private Internet*, LIGHT READING, Sept. 20, 2005, http://www.lightreading.com/document.asp?doc_id=80968.

content to their caches, they are best regarded as an overlay to the existing network. Increasingly, however, CDNs and server farms are bypassing the public backbone altogether and connecting to their caches through private networks, in the process transforming CDNs into a fundamentally different architecture.³⁰

All of these developments represent innovative solutions to adjust to the realities of the Internet. The differences in topology means that traffic that is otherwise similar may travel through the network at different speeds, with different costs, and with different levels of quality of service.

III. THE EVOLUTION OF BUSINESS RELATIONSHIPS

The evolution of the Internet has not been restricted to topology. Network participants have also been experimenting with an increasingly broad range of business arrangements. As I discuss in Section A, some of these innovations have been driven by the increasing significance of peer-to-peer technologies. Section B discusses the emergence of alternative business arrangements known as partial transit and paid peering.

A. The Growing Importance of Peer-to-Peer Architectures

One of the primary forces causing business relationships to change is the growing importance of applications using peer-to-peer technologies. The traditional Internet employed what is known as a client-server architecture, in which files are stored in large computers at centralized locations (servers) and end users (clients) request files from those computers. The relationship is generally regarded as hierarchical. In addition, the amount of data uploaded by clients is very small relative to the amount of data downloaded by servers. In the classic example of the World Wide Web, client traffic consists solely of uniform resource locators (URLs), the short bits of code identifying a particular website address. Server traffic, which consists of the data comprising the requested website, is much larger. For this reason, the technologies that took the early lead in broadband deployment (cable modem service and DSL) adapted an asymmetric architecture, allocating a larger proportion of the available bandwidth to downloading than to uploading. Newer technologies, such as fiber and wireless broadband, follow the same pattern.³¹

Peer-to-peer technologies follow a very different approach. Edge computers in a peer-to-peer architecture are not divided into those that

30. See Dave Clark et al., *Overlay Networks and the Future of the Internet*, 63 COMM'NS & STRATEGIES 109, 123-25 (2006).

31. Yoo, *Consumers and Innovation*, *supra* note 9, at 191.

host files and those that request files. Instead, computers simultaneously perform both functions. Because this relationship is regarded as less hierarchical than client-server relationships, the computers in this architecture are known as *peers* and communications between them are known as *peer-to-peer*. Peer-to-peer is thus not synonymous with file sharing or user-generated content, as is often mistakenly assumed. On the contrary, many peer-to-peer applications (such as Vuze) support commercial broadcast services, and many platforms for user-generated content (such as YouTube) employ centralized servers. The real significance of the term peer-to-peer lies in the nature of the network architecture.

It is not yet clear what proportion of network traffic will follow each architecture. For example, peer-to-peer traffic had consistently outstripped client-server traffic for several years leading up to 2007. In 2007, however, client-server traffic staged a comeback, thanks primarily to the expansion of streaming video services like YouTube, and exceeded peer-to-peer traffic 45% to 37%.³² Many industry observers now predict that although peer-to-peer will remain important, it will decline as a percentage of total Internet traffic over the next several years.³³ Even so, it is clear that peer-to-peer traffic is likely to remain a more important component of network traffic than it was during the Internet's early years.

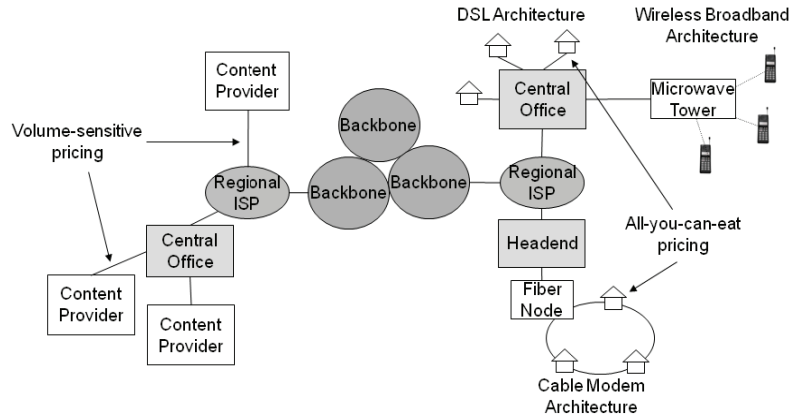
The growing importance of peer-to-peer technologies is causing significant congestion in certain areas of the network and is putting pressure on the traditional approach to pricing network services. The emergence of end users as important sources of data is putting severe pressure on the limited bandwidth allocated to upload traffic. In addition, unlike in a client-server architecture, where end users usually only generate traffic when a person is seated at the keyboard, edge computers in a peer-to-peer architecture can generate traffic for as long as the computer is left running. The result is that the lion's share of upload traffic is generated by a small number of superheavy peer-to-peer users. As few as five percent of end users may be responsible for generating more than 50 percent of all Internet traffic.³⁴

32. See Press Release, Ellacoya Networks, Inc, Ellacoya Data Shows Web Traffic Overtakes Peer-to-Peer (P2P) as Largest Percentage of Bandwidth on the Network (June 18, 2007), (on file with the author), available at <http://www.ellacoya.com/news/pdf/2007/NXTcommEllacoyamediaalert.pdf>.

33. CISCO SYS., INC., CISCO VISUAL NETWORKING INDEX: FORECAST AND METHODOLOGY 2008-2013, at 1-2, 5-6 (June 9, 2009), http://www.cisco.com/en/US/solutions/collateral/ns341/ns525/ns537/ns705/ns827/white_paper_c11-481360.pdf.

34. See Steven Levy, *Pay per Gig*, WASH. POST, Jan. 30, 2008, at D1; DAVID VORHAUS, YANKEE GROUP, CONFRONTING THE ALBATROSS OF P2P 1 (May 31, 2007); Comments of CTIA - The Wireless Association, in the *Petition to Establish Rules Governing Network Management Practices by Broadband Network Operators*, WC Docket No 07-52, 12 (Feb. 13,

Figure 6: The Traditional Approach to Internet Pricing



The most recent generation of peer-to-peer technologies can exacerbate congestion still further. In the first generation of peer-to-peer technologies, each end user stored the entirety of the files they hosted. As a result, anyone requesting those files was limited by the total bandwidth and the level of congestion associated with the network connection attached to that end user's computer. Technologies such as BitTorrent follow a different approach. Instead of storing entire files in one location, BitTorrent divides each file into pieces and distributes them at multiple locations around the Internet. When a BitTorrent user requests a file, the software then retrieves the various pieces from multiple computers at the same time. Reducing the amount of bandwidth required from any one peer improves download performance. BitTorrent also dynamically reallocates requests for pieces away from the slowest connections and toward the fastest connections, thereby placing the heaviest burden on those peers with the fastest connections.

The congestion caused by peer-to-peer technologies weighs heaviest on last-mile technologies that share bandwidth locally, such as cable-modem and wireless broadband systems. For example, cable modem technology requires that subscribers share bandwidth with the other households operating through the same neighborhood node. As a result, cable modem customers are significantly more vulnerable to the downloading habits of their immediate neighbors than are telephone-based broadband systems, which offer dedicated local connections.

Service can slow to a crawl if as few as fifteen of the five hundred or so users sharing the same node are using peer-to-peer applications to download files.³⁵

The classic economic solution to congestion is to set the price of incremental network usage equal to the congestion costs imposed on the network by that usage. However, determining the congestion cost imposed by any particular user at any particular time can be quite complex. Subscribers that use large amounts of bandwidth can contribute very little to network congestion if they confine their usage to hours when network usage is low. Conversely, a subscriber that only uses small amounts of bandwidth may nonetheless impose significant congestion costs on the network if they generate traffic at peak times. The contribution of any particular usage cannot be determined simply by counting the number of bits being transmitted. The overall impact of any particular increase in network usage can only be determined in light of other subscribers' Internet usage.³⁶ Thus it may make sense to charge different amounts to users who are using the Internet to access the same content or application if a sufficient number of other users sharing the same bandwidth are using the network at the same time.

The growth of peer-to-peer technologies has also heightened the pressure on the models that network providers have used to price their services. As noted earlier, the traditional approach charges content and application providers prices that increase with the peak bandwidth consumed, while end users are charged on an unmetered, all-you-can-eat basis. The fact that every download had to pass through one link that charged on a volume-sensitive basis allowed this pricing approach to serve as a reasonable approximation of efficient congestion pricing. For example, one hundred downloads of a 700 megabyte movie would generate 70 gigabytes of traffic from the server, which in turn would be reflected in the price paid by the content provider to its ISP.

The situation is quite different under peer-to-peer architecture. In that case, the movie could be downloaded once from the server, and the remaining ninety-nine downloads could be served by other end users running the same peer-to-peer software. Because end users are provided with service on an all-you-can-eat basis, the additional ninety-nine downloads served by the peer-to-peer network do not generate any additional revenue. The only revenue received by the network is for the

35. See James J. Martin & James M. Westall, *Assessing the Impact of BitTorrent on DOCSIS Networks*, IEEE BROADNETS, Sept. 2007, available at <http://people.clemson.edu/~jmarty/papers/bittorrentBroadnets.pdf>; see also Leslie Ellis, *BitTorrent's Swarms Have a Deadly Bite on Broadband Nets*, MULTICHANNEL NEWS, May 8, 2006, <http://www.multichannel.com/article/CA6332098.html>.

36. Yoo, *Economics of Congestion*, *supra* note 9, at 1868–69.

initial 700 megabyte download. Thus, in a peer-to-peer architecture, the amounts that content providers pay under the traditional pricing regime no longer serve as a workable approximation of the total traffic they impose on the network. Moreover, the failure to charge network participants prices that reflect their incremental contribution to congestion causes excessive consumption of network resources that ultimately harms consumers.

It thus comes as no surprise that the network providers that are most subject to local congestion are experimenting with other means for managing the congestion caused by peer-to-peer applications. For example, Time Warner has recently experimented with bandwidth caps and other forms of metered pricing. Although many network neutrality proponents have no objection to metered pricing,³⁷ recent attempts to impose metered pricing and bandwidth caps have met such a hostile reaction from the network neutrality community that the network providers had to back down.³⁸ That said, metered pricing is far from a panacea. As I have discussed in greater detail, true congestion-based pricing would vary from moment to moment based on the volume of traffic introduced into the network by other users. Not only would such a pricing regime challenge consumers' ability to process the relevant information; the distributed nature of the Internet means that no one entity has the information needed to formulate such policies. As a result, other network providers have turned to proxies that are strongly associated with high-volume activity, which most importantly includes a ban on operating a server as required by peer-to-peer technologies.³⁹

37. *Net Neutrality: Hearing Before the Senate Committee on Commerce, Science & Transportation*, 109th Cong 55, 58, 74 (2006) (statement of Prof. Lawrence Lessig), available at <http://www.gpo.gov/fdsys/pkg/CHRG-109shrg605/pdf/CHRG-109shrg605.pdf>; Tim Wu, *Network Neutrality, Broadband Discrimination*, 2 J. ON TELECOMM. & HIGH TECH. L. 141, 154 (2003).

38. For criticism of Time Warner's January 2008 attempt to impose metered pricing, see Catherine Holahan, *Time Warner's Pricing Paradox: Proposed Changes in the Cable Provider's Fees for Web Could Crimp Demand for Download Services and Hurt Net Innovation*, BUS. WK., Jan. 18, 2008, http://www.businessweek.com/technology/content/jan2008/tc20080118_598544.htm; Posting of Marvin Ammori to Save the Internet, Time Warner Goes Back to the Future, <http://www.savetheinternet.com/archive/2008/01/25/back-to-the-future-time-warner-broadband-plan-recalls-aols-walled-garden/> (Jan. 25, 2008); Posting of Lynn Erskine to Save the Internet, Time Warner Metered Pricing: Not the Solution, <http://www.savetheinternet.com/blog/2008/01/17/time-warner%e2%80%99s-metered-pricing-not-the-solution/> (Jan. 17, 2008); Posting of Fred von Lohmann to DeepLinks, Time Warner Puts a Meter on the Internet, <http://www.eff.org/deeplinks/2008/01/time-warners-puts-meter-internet> (Jan. 22, 2008). For criticism of Time Warner's January 2009 attempt to impose bandwidth caps, see Press Release, Free Press, Free Press Wary of Internet Caps (Feb. 4, 2009), <http://www.freepress.net/node/47855>; Press Release, Public Knowledge, Public Knowledge Statement on Time Warner Halt to Broadband Caps (Apr. 16, 2009), <http://www.publicknowledge.org/node/2100>.

39. Yoo, *Economics of Congestion*, *supra* note 9, at 1871.

Although this would constitute a violation of network neutrality by discriminating against a particular type of application, even network neutrality proponents acknowledge that such a restriction represents a good proxy for bandwidth-intensive activity.⁴⁰

B. *The Emergence of Partial Transit and Paid Peering*

Network providers have also begun to enter into business relationships that go beyond peering and transit relationships that dominated the early Internet. Some are driven by the emergence of secondary peering relationships discussed above.⁴¹ Before such relationships existed, a tier-2 or tier-3 ISP would have to buy transit from a tier-1 ISP that had obtained access to all of the IP addresses that it did not serve. In other words, a tier-2 or tier-3 ISP's transit relationships would cover the entire Internet (except for its own customers).

The advent of secondary peering reduces the scope of transit services that the ISP needs to purchase. In short, the ISP no longer needs to buy transit to the entire Internet. The secondary peering relationships already provide it with the ability to reach those customers served by its secondary peering partners. As a result, these ISPs have begun to purchase *partial transit* that covers less than the entire Internet (i.e., those portions of the Internet not already covered by its secondary peering relationships). In addition, an ISP with inbound traffic that far exceeds its outbound traffic may run the risk of having traffic ratios that put it in violation of its peering contract. Under these circumstances, it may attempt to cover its deficit in outbound traffic by selling partial transit contract that covers only outbound traffic, but not inbound traffic. Alternatively, it may reduce its inbound traffic by buying partial transit for inbound traffic.⁴²

Another interesting development is the emergence of *paid peering*.⁴³ Paid peering involves all of the same aspects as conventional peering relationships. Peers announce to the rest of the Internet the addresses that their peering partners control, maintain a sufficient number of interconnection points across the country, and maintain the requisite total volume and traffic ratios. The key difference is that one peering

40. Brett M. Frischmann & Barbara van Schewick, *Network Neutrality and the Economics of the Information Superhighway: A Reply to Professor Yoo*, 47 JURIMETRICS J. 383, 409 (2007).

41. See *supra* Part II.A.

42. Faratin et al., *supra* note 13, at 60–61.

43. For earlier discussions, see Christopher S. Yoo, *Network Neutrality after Comcast: Toward a Case-by-Case Approach to Reasonable Network Management*, in NEW DIRECTIONS IN COMMUNICATIONS LAW AND POLICY: THE WAY FORWARD 55, 71–76 (Randolph J. May ed., 2009) [hereinafter Yoo, *Toward a Case-by-Case Approach*]; Yoo, *Consumers and Innovation*, *supra* note 9, at 222–27.

“network economic effects,” which cause a network to increase in value as the number of users connected to it increases. To use a classic example, the value of a telephone network to a particular consumer depends on more than just the services provided and the price charged, as is the case with most goods. It also depends on the number of other subscribers connected to the network. The more people you can reach through the network, the more valuable it becomes.

The benefits created by the network economic effect for telephone networks arise with respect to a single class of customers. When a market is two sided, instead of bringing together a single class of similarly situated users, networks bring together two completely different classes of users. In those cases, the value is determined not by the number of users of the same class, but rather the number of users of the other class. A classic example is broadcast television, which brings together two groups: viewers and advertisers. Advertisers gain no benefit (and if anything suffer a detriment) from belonging to a network with a large number of other advertisers. The value of the network for advertisers is instead determined solely by the number of viewers, i.e., the size of the other class of users.

The literature suggests that social welfare would be maximized if the network provider were permitted to price discriminate on both sides of the two-sided market. It also suggests that the prices paid by those on each side of the market can differ widely and that in many cases, it is economically beneficial for one side to subsidize the other side of the market. The fact that the Internet has become increasingly dominated by advertising revenue paid to content and application providers suggest that it may be socially beneficial for content and application providers to subsidize the prices paid by end users. An advertiser’s willingness to pay for an ad on any particular website depends on the number of end users viewing that website. Under these circumstances, the optimal solution may be for the website owner to subsidize the total number of end users by making payments to the network provider to help defray their costs of connection. The costs of subsidizing more users would be more than offset by the additional revenue generated by the fact that advertisers can now reach more potential customers. In the case of broadband, this would be both economically efficient and would be a boon to consumers both in terms of providing service in more geographic areas and in reducing the prices that consumers pay.⁴⁷

These dynamics are again well illustrated by broadcast television.⁴⁸ In many ways, broadcast television and the Internet are analogous. The

222–27.

47. *See id.* at 225–26.

48. *See Yoo, Toward a Case-by-Case Approach, supra* note 43, at 73–75.

movie studios that create television programs play a similar role to content and application providers. Television networks aggregate programs and deliver them nationally in much the same manner as content networks and backbone providers. Local broadcast stations provide last-mile connectivity that is quite similar to the role played by eyeball networks. In addition, the revenue structure is quite comparable, in that television networks receive advertising revenue in much the same manner as content and application providers. Furthermore, the cost structure is somewhat similar in that connecting individual homes is much more costly than distributing programming nationally.

For decades, the standard business arrangement has been for television networks to subsidize the operations of local broadcast stations by paying them to be members of their television networks. The industry's revenue and cost structure make such arrangements quite logical. The cost of paying these broadcast stations to affiliate with a network is more than offset by the increase in advertising revenue made possible by the fact that the network is now able to reach a larger audience. Broadcast television thus represents a prime example of when firms operating on one side of the market find it economically beneficial to subsidize end users on the other side of the market.

Furthermore, the magnitude of the affiliation fees that the networks pay to broadcast stations is anything but uniform. The precise amount varies with the relative strength of the network and the relative strength of the broadcast station. Stronger broadcast stations receive more, while weaker ones receive less. Equally interesting is the fact that in recent years, the cash flow has begun to vary in its direction as well as magnitude, with weaker stations having to pay rather than be paid to be part of the television network. The dynamic nature of this pricing regime benefits consumers by providing incentives for networks to invest in better quality programming and by providing an incentive for stations to provide better carriage.

The two-sided market analysis reveals the potential drawbacks of preventing network providers from charging differential prices. As a general matter, pricing flexibility makes it easier for network providers to recover the costs of building additional bandwidth. Granting network providers pricing flexibility with respect to content and application providers should reduce the percentage of the network costs borne by consumers. Conversely, preventing network providers from exercising pricing flexibility with respect to content and application providers would simply increase the proportion of the network costs that providers must recover directly from end users. This simultaneously raises the prices paid by consumers and decreases the likelihood that the capital improvements

will ever be built.⁴⁹ Charging content and application providers differential prices thus has the potential to increase social welfare and can reduce, not increase, the burden borne by consumers.

CONCLUSION

It is all too easy to forget that the Internet is not a monolith with a brooding omnipresence overseeing the entire system. Instead, it is a collection of autonomous systems that determines the terms of interconnection through a series of arms-length negotiations between individual networks. Given the Internet's essence as a network of networks, it should come as no surprise that no two packets will pay the same amount for the same service.

The developments that I have outlined in this article have made such differences even more likely. The network no longer adheres to the rigid and uniform hierarchy that characterized the early Internet and its predecessor, the NSFNET. Packets can now travel along radically different paths based on the topology of the portion of the network through which they travel. This is the inevitable result of reducing costs and experimenting with new structures. At the same time that network providers are experimenting with new topologies, they are also experimenting with new business relationships. Gone are the days when networks interconnected through peering and transit and imposed all-you-can eat pricing on all end users. That fairly simple and uniform set of contractual arrangements has been replaced by a much more complex set of business relationships that reflect creative solutions to an increasingly complex set of economic problems. Again, these differences mean that the service that any particular packet receives and the amount that it pays will vary with the business relationships between the networks through which it travels. Although many observers reflexively view such deviations from the status quo with suspicion, in many (if not most) cases, they represent nothing more than the natural evolution of a network trying to respond to an ever-growing diversity of customer demands. Imposing regulation that would thwart such developments threaten to increase costs and discourage investment in ways that ultimately work to the detriment of the consumers that such regulation is ostensibly designed to protect.

49. See *Wall Street's Perspective on Telecommunications: Hearing Before the S. Comm. on Commerce, Science, and Transportation*, 109th Cong. 13-16 (2006) (testimony of Craig E. Moffett, Vice President and Senior Analyst, Sanford C. Bernstein & Co.), available at <http://www.gpo.gov/fdsys/pkg/CHRG-109shrg589/pdf/CHRG-109shrg589.pdf>.

