# Fast Lanes on the Internet

Traffic jams on the networks are slowing scientific collaboration. Possible solutions range from reservations-only service on the existing Internet to high-speed links just for scientists

When the Internet was being promoted just a few years ago as the tool of the future for scientific collaborations, Paul Woodward was just the kind of researcher the system's architects had in mind. An astronomer who is director of the University of Minnesota's Laboratory for Computer Science and Engineering, Woodward is part of a team at Minnesota and the University of Colorado collaborating on studies of convection in the sun. The Internet was supposed to provide a way for the two groups to exchange data and work closely together without ever leaving their labs. But it hasn't worked out that way. "When we want to sit down and look at data sets critically and brainstorm, either we go there or they come here," says Woodward. The reason: traffic congestion on the information superhighway.

"We'd like to be able to have visualization of this data and point things out to each other as if we were in the same room," says Woodward. But such graphical, real-time simulation is now impractical over the network. In principle, the Internet is just about capable of providing the 2 megabits per second (Mbs) of bandwidth (capacity) that Woodward would need to send his pictures back and forth. But in practice, he can count on only 0.2 to 0.7 Mbs, far too little for the electronic collaboration he envisions.

It's a problem familiar to every World Wide Web user who has waited in frustration as Netscape endlessly displays the message: "Host contacted. Waiting for reply." And it's only going to get worse. Most universities and research institutes are connected to the Internet via cables that carry 1.5 or 45 Mbs, but the Internet is so clogged that everything slows down. It's like a highway with a speed limit of 55 mph where everyone ends up going 35 because there's too much traffic—an example, says computer scientist David Farber of the University of Pennsylvania, of how "success tends to get you into trouble." The networks have far more users now than even

real time to run models or analyze large data sets—are on hold. Says Tom DeFanti, director of the Electronic Visualization Laboratory at the University of Illinois, Chicago, "The Internet has become a mass-migration highway. Anyone trying to get work done is looking at alternatives."

The alternatives DeFanti refers to can be as simple as changing work habits. At AT&T Research, says Steve Crandall, a staff scientist, "people modify their behavior. They come in in the early morning, six or seven o'clock, to get some of their high-bandwidth work done." Universities and research laboratories can also address the problem by increasing the bandwidth of their connections to the Internet. But high-bandwidth connections are expensive, and they provide no solution to congestion elsewhere on the Internet.

That is why more and more researchers are concluding that what's needed are strategies that make distinctions among users, offering high bandwidth and prompt service to those who need it while leaving e-mail and other services that can tolerate delays to fend for themselves on congested lines. Such solutions would benefit commercial users as well as scientists, because the

2 years ago, and a larger fraction of them are running video and audio programs requiring lots of bandwidth. The proliferation of glitzy Web sites with sound and color graphics has only added to the traffic jams as people try to download elaborate pictures from popular sites. "These things simply take more capacity than old e-mail," says Mark Luker, director of the National Science Foundation's (NSF's) networking program.

For most scientists, like everyone else, the delays in e-mail and Web browsing are an inconvenience. But for researchers like Woodward, the holdups are intolerable. Indeed, some of scientists' most ambitious visions for Internet use—operating telescopes or other instruments from a distance or collaborating in

delays that are now mostly a nuisance could become a real impediment to commercial expansion of the Internet.

One set of strategies would create fast lanes on the existing Internet by prioritizing data traffic, allowing scientific and commercial traffic to bypass the congestion; another would create "private roads": experimental high-speed networks reserved exclusively for scientists who do intensive computing. Neither solution will be easy to implement. The idea of giving some traffic preferential treatment conflicts with the egalitarian culture of the Internet and would probably require changes in pricing for Internet services. Specialized high-speed networks, meanwhile, are themselves vulnerable to

# Will Pricing Be the Price of a Faster Internet?

Clever technology may succeed in opening some fast lanes on the Internet for scientific users who need high capacity (see main text). But many Internet researchers say that keeping those fast lanes from clogging like the rest of the Internet will take something more than technology: some form of economic incentive—pricing, in other words—so that when the network is congested, bandwidth will go to the users who pay for it. At the moment, a surgical team doing real-time surgery over a video link, say, doesn't have any more claim on Internet resources than does a teenager using up at least as much bandwidth by watching recreational videos. "Pricing is a time-honored, tried and true method of dealing with this," says economist Jeffrey MacKie-Mason of the University of Michigan.

**Internet economists.** Hal Varian (above) and Jeffrey MacKie-Mason (right).

The basic idea is simple, MacKie-Mason explains. Varying prices by time of day or type of service (e-mail, video, Web traffic, etc.) will help control demand because only those people who really need the highest bandwidth service will pay for it. "For congestion purposes, [the function of] pricing isn't to raise money but to allow people to express a preference of how much they value [the services]," he says. Already, researchers are developing the accounting software and payment schemes that would be needed. They are also debating the administrative and sociological aspects of pricing—who should administer it, and how it will affect the current Internet free-for-all.

Internet experts who favor differential pricing say that it is fairer than the current system, in which an institution pays a fixed annual fee to its Internet service provider for unlimited use of its connection to the Internet, no matter how many users the institution has or what Internet services they favor. While Internet providers such as Compuserve and America Online bill individuals by duration of connection and sometimes by type of service, most universities and laboratories don't impose similar charges on their users. At the University of California, Berkeley, says computer scientist Pravin Varaiya, "20% of the users account for 90% of the traffic." As a result, he says, "light users subsidize heavy users."

What's more, the revenues that universities or Internet providers would derive from congestion pricing could be used to add capacity to the Internet. "You have to have the money to pay for expanding the capacity, and it's better for that to come from the users who actually use it," says Hal Varian, dean of the School of Information Management and Systems at Berkeley. That's already happening in the commercial world: MCI recently announced a venture to establish new high-speed backbone connec-

tions that will be available to customers—mainly businesses—willing to pay for premium service to insure that their traffic gets through when the Internet is too crowded.

The everyday Internet doesn't yet have a system for metering usage, but Varian says "I feel the problem is more tractable than people are willing to admit." Some accounting software is already available. In New Zealand, the University of Waikato operates the single Internet gateway to the United States, via Hawaii, on behalf of the New Zealand universities. It meters international Internet traffic from the universities according to type of service (such as ftp or e-mail), time of day, and number of bytes, and bills them accordingly. Varaiya and his colleagues have developed software that also authenticates the user, a measure to help prevent fraud. The accounting and authentication only add about 150 milliseconds of overhead time to the operation. "Accounting doesn't add too much in terms of time and money," Varaiya says.

The system includes a purchasing agent that can ask the user to decide if a requested service is too expensive. A user might set up the agent to accept charges of less than 50 cents, but to pop up a window asking for the user's okay on charges greater than that. Varaiya is now planning to try out the system on a group of about 200 users on campus. "We want to give them variable rates of service and a pricing structure and see how they react," he explains.

Along with software for doing the accounting, a pricing scheme requires a means of payment. The traditional method is centralized billing, as is done by telephone companies. That's how New Zealand's system works. But other proposals would eliminate billing and require users to pay as they go, perhaps by attaching "digital stamps," purchased in advance from their Internet service provider, to each message. Varian suggests that Internet tolls could also be collected by the micropayment technologies currently being developed for commercial uses of the Internet—systems capable of coping with price increments of thousandths of a cent. A user might have $50 in a micropayment "card," and a few thousandths of a cent would be deducted automatically each time the user sent out e-mail or browsed the Web.

The big unknown, the pricing enthusiasts agree, is how Internet users would respond to congestion pricing. "Would they turn off the images in Netscape?" Varian wonders. "No one knows." Nor does anyone know whether Internet users and administrators could even be persuaded to adopt a pricing scheme—something to which Internet culture is historically hostile. Says MacKie-Mason, "There's a cultural resistance to allowing people to buy their way to the head of the line." The question is how much congestion users will tolerate before that cultural resistance breaks down.     –E.G.

---

being overwhelmed by traffic growth. But somehow, says Hans-Werner Braun, a staff scientist who works on networking at the San Diego Supercomputer Center, "we need to make the current Internet more predictable and more verifiable. If you're a telephone or airline customer, you have certain expectations [about performance]. We need to have similar expectations for the Internet."

## Sharing the pain

The Internet now is a best effort system in which all data packets are treated equally. Current Internet routers—the way stations along the network, which read the addresses of incoming packets and send them along to their destinations—work on a first-in, first-out scheduling algorithm. "As packets arrive, they're stuck in a queue, and as bandwidth becomes available they're shipped out,"

explains Sally Floyd, a network researcher at Lawrence Berkeley National Laboratory (LBNL). When the networks become congested, service degrades equally for all packets, whether they are e-mail, which can tolerate delay, or real-time video, which cannot.

One barrier to tackling congestion on the Internet is the lack of good statistics on how data is flowing, where the crunch points are, and how much different services contribute

to the congestion. "There are no widely accepted metrics to assess service," Braun says. The fundamental dilemma confronting network planners is clear, however: "If you have more people using the network than you have bandwidth for, should you push people off, or should you let everyone suffer?" asks Scott Shenker, a computer scientist at Xerox PARC.

On the U.K.–U.S. trans-Atlantic link, one of the most heavily congested segments of the Internet, network researchers are testing a scheme for pushing off some users—with a minimum of pain. "At the moment the international link is unusable," says Jon Crowcroft of University College, London, and a member of the Internet Architecture Board. Then he corrects himself. "No, there was a recent upgrade. Before that it was unusable. Now it's workable about a third of the day."

To improve that figure, the network community in Britain is trying to reduce the amount of Web traffic. When U.K. users taking part in the experiment want to visit a Web site in the United States, they point their browsers at one of six Web cache sites within Britain. If another U.K. user has recently visited the U.S. Web site, it will have been stored on the cache and won't have to be called up from across the Atlantic again. If the site requested isn't yet in the cache, the system will go out and fetch it, then store it for later users. It also checks to see if the page has been modified since a given date.

"Caches aren't great, but they will do as a stopgap," says Crowcroft. No one in the United Kingdom will be forced to go through caches, but as an incentive, part of the international link will be set aside for traffic routed through them. "If people go [to U.S. sites] via the U.K. national Web cache sites, they will get very good performance, but otherwise they will get bad performance. With wide publicity we hope this will fix a large proportion of the problem," Crowcroft explains. At least 30 British universities are already using the caches. A similar system has been in place in Israel since November 1995, where all university Web traffic is required to go through cache servers. By relieving congestion due to Web use, these measures should open the bandwidth needed for critical applications.

Still under development are congestion-cutting schemes that would carry such prioritizing much further. One would rely on the routing computers to assign priorities automatically to different kinds of Internet traffic. At LBNL, for example, Floyd is working on a scheduling algorithm called class-based queuing. "The idea is that the router can classify arriving packets into classes," Floyd says, then give them different grades of service. For example, traffic for the Web might be put into a class of its own and restricted when the network is congested. However,

given the huge commercial push to expand Web use and make it more user-friendly, any scheme to improve scientific use at the expense of Web use may well conflict with growing commercial interests.

## Reservations required

What's more, automatic schedulers can't evaluate the traffic's urgency by distinguishing, say, a recreational video from a scientific visualization effort. But another set of approaches to reducing congestion—so-called reservation schemes—lets users divide their own Internet traffic into priority classes. The leading reservation scheme is a system called Resource Reservation Protocol (RSVP), which Shenker, Steve Deering, and others at



**Pooling computer power.** To simulate eddies in a high-speed fluid flow, researchers would like to link distant supercomputers.

Xerox PARC and elsewhere helped develop.

Computers use RSVP to request a specific "quality of service" from the network on behalf of an application such as video teleconferencing or real-time simulation. That request might be something like "I need 2 Mbs of bandwidth for the next 2 hours between New York and Phoenix." The system then sends a quality-of-service request from router to router along the data's path, and the rout-
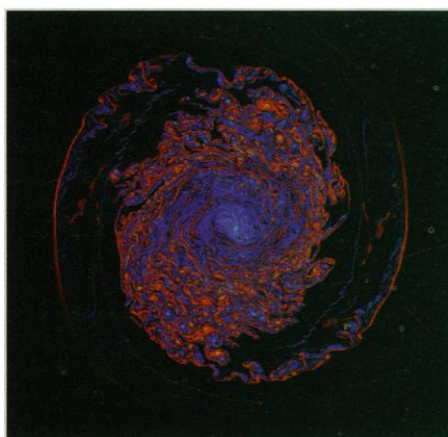
ers try to reserve the necessary bandwidth. If some routers along the path have not yet implemented RSVP, the request is ignored for those segments of the trip. Thus RSVP can't guarantee a particular quality of service for an entire trip across the Internet. But even if only some routers take part, they will improve the average performance of the network considerably. "The feeling in the Internet community is that most users would rather have better service than an absolute guarantee," says Robert Braden, head of the RSVP project at the University of Southern California's Information Sciences Institute.

RSVP has been running in experimental test beds for the past couple of years, and a number of vendors are starting to produce routers that implement the scheme. RSVP and other reservation systems require more than hardware, however; they also need some mechanism for deciding who should get reservations for high-quality resources. "As soon as you provide any kind of preferred service, you need a mechanism to prevent abuse," says Braden.

The obvious mechanism is pricing: Make it more expensive to obtain a higher priority (see box). But the prospect of charging more for first-class service makes some network researchers reluctant to embrace resource reservations. Even Deering, who worked on RSVP, says he now has his doubts about it: "Reservations are expensive and complex and haul in a kind of charging system, and what's the criterion for saying who gets reservations and who doesn't?" He acknowledges that "special critical applications like remote surgery need reservations. It's [reservations for] every phone call that I don't approve of." One major concern is that reservations might price out those with less money, such as schools.

The same fears apply to another technology for opening up fast lanes on the Internet: a system known as ATM (asynchronous transfer mode). Unlike RSVP, which offers a first-class version of existing Internet 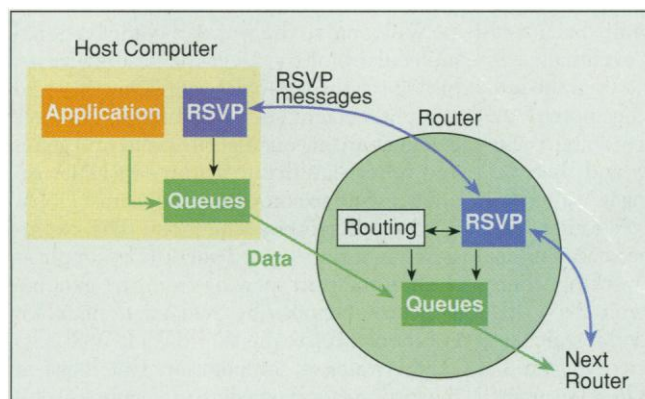service, ATM is a new way of packaging and sending data. It replaces the various-sized data packets that travel over the existing Internet with data cells of a fixed size. And instead of allowing the packets in a single message to wend their way to their destination via many different routes, it opens up a single "virtual circuit" from source to destination. Both features make the system more predictable than the standard Internet. In the existing network, unpredictable delays can result when,



**Asking for priority.** RSVP can speed the passage of messages through the Internet by sending requests for bandwidth from the host computer to the routers along the way.

for example, a small packet gets stuck behind a large one while the large packet waits for enough bandwidth to become available. But an ATM-based network can predict, based on traffic, how long a transmission will take.

Moreover, because ATM sets up a virtual circuit for each transmission, it allows a user to request a specific quality of service in advance. A user running a multimedia application, for example, would request—and presumably pay for—service with very little cell delay because it's a real-time application, while e-mail users would be content with cheaper, slower service.

With software changes, ATM can run over existing Internet cables and routers, says Mark Laubach, who chairs an Internet Engineering Task Force working group on ATM, "but it doesn't work well unless done in hardware," meaning expensive new cables and other equipment. Still, ATM networks are up and running now: for example, the Bay Area Gigabit Testbed, an experimental high-speed ATM network connecting 15 sites in Northern California. It is used for a variety of collaborative scientific experiments including remote studies with optical and electron microscopes. MCI's Internet backbone (a major Internet pathway linking smaller users, the way an interstate highway connects feeder roads) also uses ATM.

Some lucky scientists stymied by the congestion on the Internet don't have to bother with caching, RSVP, or ATM. They can move off the existing network altogether. One "private roadway" already available to scientists is the NSF-sponsored vBNS (very high-speed Backbone Network Service), which connects five NSF supercomputing centers at 155 Mbs on an ATM network and provides bandwidth for cutting-edge network applications and research. It is not meant to be used for day-to-day operations such as e-mail and ftp, but that restriction may be difficult to maintain because the NSF is tying more universities and other sites into the vBNS.

That's the paradox of the Internet—and the reason that congestion is likely to plague scientists for the foreseeable future. Scientists move to high-speed networks, eventually everyone else jumps on board, and then the scientists have to move up another notch. "A few of us are out on the edge doing these things on very fast machines, and then 10 years later everyone else is doing it," says Paul Bash, a research scientist at Argonne National Laboratory. The Internet began as an experiment in computer networking, then became a popular phenomenon. Now it's groaning under the demand, and researchers are trying to make it safe for science again.

–Ellen Germain

*Ellen Germain is a science writer in Arlington, Virginia.*

# Software Matchmakers Help Make Sense of Sequences

Gene sequencers are spinning out data at a mind-boggling rate. They have already sequenced the complete genomes of several bacteria and brewer's yeast, they will have completed the genome of the roundworm *Caenorhabditis elegans* in a couple of years, and they intend to wrap up the human genome by 2005. A string of the four letters A, G, T, and C, designating the four nucleotides that make up DNA, is unreeling from sequencing labs at an ever-increasing pace, now nearly a million a day. For the human genome alone, the sequence will total 3 billion nucleotides.

All this would be little more than so much genetic ticker tape without some way to decipher its real meaning, which is largely hidden in the genes—the stretches of DNA, amounting to barely 3% of the human genome, coding for the proteins that are the workhorse molecules of life. The first step is to recognize the genes from their distinctive sequences of nucleotides. The next is to infer the function of the proteins they code for—and the key to doing that is to find related genes and proteins whose functions are already known.

As molecular evolutionist Russell Doolittle of the University of California, San Diego, explains, "The structures of all these proteins and the genes that code for them are all related through a big evolutionary expansion—some small number run through biochemical Xeroxes and used over and over in different settings." The challenge of learning the function of a newly generated sequence is the kind of challenge that computer scientists in other fields have been wrestling with for decades: spotting obvious, or less than obvious, similarities in different strings of data.

Welcome to the world of computational molecular biology. Over the past few years, biologists-turned–computer scientists and computer scientists–turned-biologists have begun churning out algorithms to find genes and other significant features in DNA sequences and to compare and contrast DNA, RNA, and protein sequences. The explosion has been triggered not only by supply—the information spewing from the genome projects—but also by demand from biologists hooked up to the World Wide Web, says David States, a computational biologist at the Institute for Biomedical Computing (IBC) at Washington University in St. Louis. "Most biologists in academic settings now have access to the Internet and Web browsers," says

States, and that allows them to send their sequences to on-line analytical tools—or even borrow the tools and wield them on their own workstations (see p. 591).

This past June, States and his colleagues at the IBC hosted the Fourth International Conference on Intelligent Systems in Molecular Biology (ISMB) in St. Louis to survey the explosion. The computational tools under discussion ranged from simple programs that search for similarities between known and unknown sequences to ambitious efforts to find complete genes in DNA sequences and relate the proteins they produce to known protein structures. Many of the new tools rely on techniques developed by researchers in machine learning and artificial intelligence, and the hottest subject of the conference, known as hidden Markov models, springs directly from statistics and linguistics.

Sustaining all these efforts is a sense of mission, says Doolittle. The ISMB researchers "are missionaries and proselytizers, and they have this great esprit de corps." With the tools now under development, biologists "should be able to relate proteins whose relationships weren't detectable and do faster searching of genomes and comparisons of genomes," says Doolittle—"all sorts of things that weren't possible before."

**Make me a match.** One reason sequence comparisons are so powerful is evolution's conservative style. While the 20 amino acid alphabet of proteins could in theory spell out a nearly infinite number of proteins, actual proteins are variations on a limited number of themes. Human beings alone have perhaps 100,000 proteins, but we and other organisms "are dipping into a pool of relatively slowly evolving proteins we all share," says David Lipman, head of the National Center for Biotechnology Information (NCBI). All together, the number of different protein families is "maybe less than 1000." The result is that comparing an unknown gene to known ones has a reasonable chance of coming up with a match—providing the computer algorithm can recognize subtle similarities.

The first problem is to find the genes, which in higher organisms, known as eukaryotes, come interspersed with pieces of noncoding DNA called introns. One approach is to look for the telltale patterns of DNA that mark the boundaries between the coding and noncoding regions. Researchers have come up with various pattern-recognition