# 23 The Emergence of Grid Information Infrastructures

## Larry L. Smarr

The information technology revolution is going to affect everything about research and education. What is emerging in many people's minds is a transformation similar in impact to previous infrastructure transformations, such as the birth of the electrical power grid. In that era, we went through a series of science and engineering experiments (real world test beds with competing standards such as Edison and Tesla fighting over DC versus AC), similar to what we have today in the battles over various flavors of operating systems or web browsers. Eventually, over the course of decades, the marketplace settled these disputes and what rapidly developed was a ubiquitous electrical technology infrastructure.

One rarely thinks about the subcomponents of the electrical power system anymore (unless it fails). But the subcomponents of the information technology system are still in development. We are a long way from having ubiquitous information delivery. The goal should be "where there is air, there are bits." But the wireless revolution is only beginning, even though a lot of the fixed fiber and cable are in place.

While such historical analogies are useful, the information technology system is fundamentally different from the other aspects of technology that change society, which are typically linear in nature. This modern infrastructure system is driven by an unrelenting exponential. Most people, even scientists, have a hard time understanding how radical such

exponential change is. When an exponential sweeps through a pre-existing threshold, it changes the impossible into the routine.

As an example, when I founded the National Center for Supercomputing Applications (NCSA), our first supercomputer cost $10 million. We had to dig a trench to the Illinois Power electrical substation to run it. There were only four other such supercomputers available to the National Science Foundation (NSF) funded academic research community. Today, the exponential of Moore's Law has enabled the average PC laptop to have more memory and to run faster than that 1985 supercomputer. This move from the very elite edge of science to mass consumption is happening on time scales of ten to fifteen years.

At public meetings in the early 1980s, I can remember telecom company representatives telling the National Science Foundation that the Federal Government had no business funding the creation of a general Internet for academics like we have today. They said that would be an unacceptable interference by the government in the private sector. So the government said, "Fine. But surely you don't mind if we hook up the five new NSF Supercomputer Centers with a national backbone network of 56 kilobits per second, do you?" And I recall the telecom representatives responding, "Well, let's see. No, there's no revenue there. Sure, you can do that. It won't harm anything." Well that NSFnet supercomputer backbone led directly to the regional and local networks that created the NSFnet that was finally commercialized a few years ago to form the global internet we all live on.

The point is that government research at the edge of the marketplace can make huge changes when you are dealing with an exponential driver. For the ten years prior to the creation of the NSFnet backbone, the Defense Advanced Research Projects Agency (DARPA) supported researchers to create the ARPANET testbed and the TCP/IP protocols, which define today's Internet. This seminal work set off the exponential which the NSFnet backbone picked up, quickly growing from 56 kilobits/sec to 1.5 megabits per second to 45-megabits/sec to today's NSF vBNS backbone of 2.4 gigabits/sec. Applications enabled started with rudimentary email, remote logon, and file transfer, but grew more complex as the greater bandwidth allowed the graphical Web to come into being, now followed by the digital video explosion.

The World Wide Web emerged in a similar pattern to the internet. Research on how to hyperlink physics articles led Tim Berners-Lee, a physicist at CERN, to develop the protocols for the World Wide Web. In 1993 NCSA staff developed the graphical web browser NCSA Mo-

saic and the open source web server NCSA HTTPd, making them available for free to several million users worldwide. As more web servers went up to host content, there were more things to look at which led to more downloadings of the browsers, setting off exponential growth. Once again, this rapidly (in this case less than two years, rather than the ten with NSFnet) led for a need for commercialization.   Many of the NCSA Mosaic developers left to help found the technical core of Netscape, while Microsoft licensed Mosaic to create Internet Explorer. Similarly, Apache built their open source server software based on the NCSA web server code and created the most widely used server today. Again the NSF funded supercomputer centers played an intermediate transformational role between the protocol development and early experiments and the final mass market commercialization.

Perhaps we can get a glimpse of the future information infrastructure by observing what scientists are doing now, namely creating a ubiquitous digital infrastructure in which to carry out science and education. Because of the volume of data that science deals with this infrastructure must have very high speed networks, appropriate levels of computer power distributed throughout the Net, and large data caches strategically positioned. Effectively, academic scientists are saying that we need to do this because trying to live in the filled commercial pipes of the Internet is not the way to get science done very quickly.

Of course, with this vast increase in data, we have to be able to analyze and visualize the data, using virtual reality technologies. Furthermore, every modern digital scientific instrument has a host computer with it. With a high speed networking card, such local instruments can become shared global superinstruments.  Finally, given that scientific research and education is a collaborative activity, researchers are creating digital video mediated electronic meeting rooms and teleclassrooms.

In order to systematically explore this new digital workspace, the National Science Foundation has undertaken a bold experiment called the PACI Program (Partnerships in Advanced Computational Infrastructure). About 280 universities have faculty and students using one of the remote NSF supercomputers over the last couple of years, with almost 1,000 projects. As a result of the national PACI competition, NSF established two leading edge centers: the San Diego Supercomputer Center at UCSD and NCSA at UIUC.  Roughly 50 partnering universities are connected to these leading-edge sites by the NSF vBNS and Internet2 high speed networks. The PACI programs involve computer scientists, computational scientists, deployers and supporters of infrastructure,

and people in education. They have formed (there is not even an English word for it) an interlocking web of virtual partnerships. As we would say in the Midwest, these folks are working together to raise a "cyberbarn."

The "cyberbarn" that the PACI program is building will be a prototype of the commercial 21st century infrastructure, termed a Grid. Just as with the NSFnet, these Grid experiments are not the commercial broadband Internet, but may directly lead to new services on it. This high speed research network is much faster than the commercial broadband network (enabled by cable modems and DSL) that you keep hearing about which will roll out commercially to your home and business over the next five or ten years. It is also being used intensively right now. NSF funds it to see what kind of applications will emerge in this futuristic information infrastructure, which has a thousand times the bandwidth of the ubiquitous modem based network of today. Why do we care about this experiment? Because it all comes back to the exponential growth-in ten years, by Moore's law, you will get somewhere between a 100- and 1,000-fold increase. What PACI is experimenting with now, by 2010 it will be in your living room. The fact that our country and the world are getting an early view of what this is like will drive the marketplace in terms of preparing for it.

NSF has also funded a digital port-of-entry to the United States, the STARTAP, in Chicago at the Ameritech Network Access Point. At that interconnection point, the NSF funded grids link with the grids that the other agencies are funding, as well as research networks from many different foreign countries. A lot of experiments emerging in this next phase will be inherently global in extent.

It is instructive before we explore the new world of the Grid, to recall why the supercomputers from which they sprang were so important to some researchers. Supercomputers like radio-telescopes are national facilities with unique capabilities. They enable projects that cannot be done anywhere else. Through peer review, anyone with a good idea can get access to them. That's why NSF has traditionally built such national facilities.

We can see an example of why you might want to use a supercomputer from a breakthrough calculation done earlier this year. Atmospheric phenomena are excited on many different space and time scales. Getting only an overall view (like on the TV news) is not enough. The Center for Analysis and Prediction of Storms (CAPS) program, which is a NSF Science and Technology Center located at the University of Oklahoma, has been using supercomputers to predict on a much finer

level grid, for example a few kilometers (compared with current National Weather Service computational grids of 20 or 30 kilometers). You will not see this today on the evening news, but if these PACI experiments work out perhaps you will benefit from this routinely in five to ten years.

A year ago we started a pilot of this approach by dedicating one of NCSA's largest supercomputers to the CAPS group during the week-long American Meteorological Society meeting in January 1999. Every morning the CAPS team ran simulations of the weather faster than real time. Where severe weather shows up they used nested grids to resolve features as a variety of spatial and temporal scales. Because of the parallelism of the supercomputer they were able to carry out multiple runs in which they made variations in the uncertainties in the initial data. This "ensemble computing" created much more accurate forecasts. They also improved regional forecast accuracy by the use of data on velocities of the atmospheric water obtained from the new NEXRAD Doppler radars in Oklahoma.

Eventually we would like to run high resolution local forecasts on regional supercomputers and tie them all together into a national or global low-resolution forecast with this high-speed Grid linked into the Doppler radar system. Already the results of these weather forecasting experiments are being distributed over the World Wide Web. In fact, the FAA and some airlines are now using the results of these experiments to reroute planes away from severe storms.

So what we see here is a glimpse into a more general transition from using stand alone supercomputers to using the Grid to enable scientific research. The Grid is not just a supercomputer, but rather a network of computers, sensors for input, with results output to the Web. It is a grid-based, not just a supercomputing-based, approach to computational science.

The Grid is also transforming the way we hold meetings. It may be a while before large groups like the AAAS can meet on the Grid, but smaller groups are doing it now as experiments on the Access Grid. The Alliance just opened an ACCESS facility next to NSF building in Arlington, Virginia. We recently had a remote lecture from a scientist at Argonne National Laboratory to an audience at the ACCESS center. On the right-hand side of the floor to ceiling screen were the power points and on the left-hand side was a video window in which the speaker looked larger than life. The speaker could see the audience and they could see him. And the audience could ask questions and he could respond.

Furthermore, he was only one of a half dozen sites simultaneously linked, all mutually seeing and interacting with each other.

Within a few years you will see similar services coming over broadband to in-home high-definition television (HDTV), with life-like human figures will be projected using digital video. Alliance researchers are already adding interactive real time virtual reality capabilities to the Access Grid. And, of course, the entire World Wide Web of information is digitally available to bring into these virtual discussions, lectures, and group meetings. The PACI program is setting up a series of prototype ACCESS centers across the country as cyberports into the Grid. With our long range government funding we are focusing on driving applications from science and education.

As a young theoretical astrophysicist, I had always assumed that almost all scientists were theorists. It was pointed out to me that in fact, maybe only five or ten percent of academic scientists are theorists and all the rest were experimentalists or observers. So until Grid technology impacts experimentalists it will not really impacted science as a whole. This critical phase is still in the very early stages. For instance for years the San Diego Supercomputer Center has worked with one of the largest electron microscopes in the country, linking it through the network to high-power visualization software and to supercomputers to turn batch observations into interactive ones. This effectively creates a "super instrument" that is available from anywhere on the Net.

The Alliance has a partner experiment, the Berkeley-Illinois-Maryland Array, which has as its "lens" our supercomputer at NCSA. This millimeter array synthesis telescope, which is located in a high desert site in California, sends data across the network to the NCSA SGI Origin supercomputer. Then the result, which is 2000 images of an object in the sky, each taken at a different millimeter wavelength spectral line, is stored as a data cube in a digital library on the Internet. You can interactively rotate the image and zoom in and out to see details inside the cube. Here we see the practice of science actually beginning to be impacted by the Web.

What I have found most peculiar about this whole emergent global cyber-enterprise is that scientists, the people who invented the Web, do not use it nearly as much as nonscientists. In the average household (with only 1/000 of the bandwidth of a university researcher), citizens are using portals which have pages individually customized just the way the user wants them. People spend hours a day on the web trading stocks, researching purchases, chatting with friends, and buying goods and services. They are coming to *live* on the web. As I visit universities around

the country I find that your typical scientist is way behind the Web power curve compared with the typical home cyberpioneer. We have to figure out how to get the scientists to actually use the results of this revolution they created.

Another big change that is coming has to do with the human computer interface and in particular, how we visualize massive datasets. Two-dimensional megapixel screens may be acceptable for personal PCs, but if your dataset was generated by a high performance computer with 1,000 times as much computing power as a PC, then the PC screen is clearly not going to be adequate for analysis of the data. For instance, if you have computed the details of a complex physical/biological environment, how are you going to represent that other than the way the real environment is represented? It must be three-dimensional and time dependent. That is why people have been creating virtual reality theatres and large screen displays. Furthermore, you can create virtual worlds and then link them over the Internet. That is, you see your remote colleague as an avatar, a software representation of that person, in the appropriate place in the shared virtual space where they are standing. The computer also knows where hands and eyes are pointing, so the avatar has a moving head and hand. You jointly go into this virtual space with other scientists, essentially like living with talking ghosts, and then you do your science.

A pioneering example of this cyber-modality is the Alliance Environmental Hydrology Application Team. NCSA and University of Illinois researchers are investigating the Chesapeake Bay with their colleagues at Old Dominion University. ODU has an ImmersaDesk (a one-wall CAVE) that is linked over the highspeed NSF vBNS to NCSA's Power Wall (which is a 4-screen high-resolution device), a walk-in CAVE and to a desktop. We are able to bring the various scientific and technical specialists that are needed, including computer scientists who help develop and drive the software, into this single collaborative environment.

This may sound like science fiction, but it has already been used for several years by industry. Caterpillar, an NCSA Strategic Industrial Partner, used their computer-aided design databases for earth moving equipment as the data source for the CAVE. The equipment, which may not have physically been built yet, appears as a full three-dimensional object which can be driven through the Peoria Proving Grounds or a Siberian rock quarry, all created by realtime texture mapping of photographs of those environments. By linking up CAVES and/or televideo from Germany to Houston to Peoria to Urbana, Caterpillar can bring together

experts in design, manufacturing, maintenance, and customer support to make joint decisions about design choices.

-The level of research innovation in universities is possible because of the strong support of the Federal Government, such as the National Science Foundation for long-term basic research. With the PACI program, the NSF has also moved to fund "reduction to prototype" of information technology research results. The Federal Government is the only part of our society that is able to fund the 5-, 10-, 15- year-out research. Industry is not able to do that because of the quarterly focus on share holder value. The recent report by the President's Advisory Committee on Information Technology made this point strongly. It also pointed out that the Federal Government's support for long-term high-risk IT research has actually been pulled back over the last few years. For instance, we are seeing an under-investment in IT research compared to what we had with DARPA in its heyday.

Given that 30 percent of the real GDP growth in the last ten years has come from information technology related activities and yet today only one out of 75 of the federal R&D dollars goes into information technology research, the Federal Government has done some soul-searching since the PITAC report came out. Rarely has the Federal Government reacted as quickly with a budget request in response to a study as it did with the PITAC report. This is particularly important because of how quickly this information revolution is moving.

The PITAC report looked in an integrated way at the Internet, supercomputers, virtual reality, database, software, and so forth. Based on this study, the report chose four areas as priorities for research: software; scalable information infrastructure; high-end computing; and, for the first time, socio- and economic-impact studies. With the FY00 actual budget and the FY01 budget request, the Federal Government proposes to roughly double its investment in long-term information technology research, adding another $1 billion a year to what is currently funded. We can expect to see this proposal treated in the traditional bipartisan manner of a national initiative that will be supported by all parties.

This is a good response by the Administration and I expect a strong support in Congress. But there are a few potential roadblocks we need to be worried about. The AAAS analysis of the R&D funding environment in Congress points out that our country may have prematurely celebrated the end of the deficit era. There is some questionable bookkeeping involving how the federal government is treating trust funds as income rather than long-term debt. We still have the discretionary caps

on the books and heretofore the political will has not been there to re-
move them. We must get the universities involved in a new style of R&D
which is appropriate to the information industry rather than the manu-
facturing industry. Finally, we must train a large number of young peo-
ple who are going to carry the innovations of the universities the next
set of companies. If we do not do this we will never meet the kind of
growth rate that it will take for our country to avoid stalling out in this
information revolution. These are all serious barriers, but there are
encouraging signs that most of them will be dealt with. If so, then we
are on the verge of an even greater speed of transformation of the infor-
mation infrastructure than we have seen in the past. This will allow our
country to continue its global leadership in creating the information infra-
structure needed for 21st Century commerce, research, and education.