

Deutsche Forschungsgemeinschaft

Commission on IT Infrastructure

Study Tour 2011

USA

I. Summary.....	3
II. Background of the Study Tour	9
III. General Information	11
Participants	11
Agenda (March 20 – April 2, 2011).....	12
IV. Summaries of the Visits.....	14
1. NERSC, Berkeley, CA.....	14
2. NVIDIA, Santa Clara, CA.....	15
3. Apple, Cupertino, CA	16
4. Google, Mountain View, CA.....	17
5. Intel, Santa Clara, CA	18
6. IBM, Almaden, CA.....	19
7. ICBM at UCLA, Los Angeles, CA	21
8. AMD, Austin, TX.....	23
9. ICES, University of Texas, Austin, TX.....	23
10. TACC, Austin, TX.....	24
11. IBM, Austin, TX.....	25
12. NCSA, Urbana, IL.....	27
13. Argonne National Laboratory, Chicago, IL.....	28
14. DOE, Washington, DC.....	29
15. NSF, Arlington, VA.....	30
16. National Cancer Institute, Rockville, MD.....	33
17. Harvard Medical School, Boston, MA	34
V. Conclusions.....	37

I. Summary

In the recent years, the Commission on IT Infrastructure (KfR) of the German Research Foundation (Deutsche Forschungsgemeinschaft or DFG) went on study tours to the USA and Asia to gain insights into the latest developments in IT infrastructures at renowned research laboratories and universities as well as into the IT trends envisioned by the IT industry and into the IT funding schemes of the respective national funding agencies. In this fruitful tradition, the KfR chose the USA as the destination of its most recent study tour for the second time after 2005. A broad spectrum of institutions were visited from March 20 until April 2, 2011, with a topical emphasis on the IT at universities, in university hospitals and medical research, on (super-) computing centers and grids, on IT funding bodies, and on the IT industries.

The first large group of target institutions consisted of university or research lab institutions, respectively, with a focus on the topics mentioned above. The study tour's agenda contained visits to the National Energy Research Scientific Computing Center (**NERSC**) at Lawrence Berkeley National Lab (Berkeley, CA), the Laboratory of Neuro-Imaging (**LONI**) at the University of California at Los Angeles (UCLA, Los Angeles, CA), the Institute for Computational Engineering and Sciences (**ICES**) at the University of Texas (**UT**) at Austin (Austin, TX) as well as the Texas Advanced Computing Center (**TACC**, Austin, TX), the National Center for Supercomputing Applications (**NCSA**) at the University of Illinois at Urbana-Champaign (Urbana, IL), the Argonne National Lab (**ANL**, Chicago, IL), the National Cancer Institute (**NCI**, Rockville, MD), and **Harvard Medical School** (Boston, MA).

NERSC is one of the leading high-performance computing (HPC) facilities funded by the United States Department of Energy (DoE) and currently runs the 1.25 petaflop system Hopper. NERSC's research activities are oriented towards exascale computing, with special emphasis on energy-efficiency, predictive computational science, or computer science methodology for computing applications, e.g. via the Magellan project, cloud computing has also become a main research topic. **LONI** is an independent research center affiliated with UCLA's School of Medicine, and it is one of the main players in the International Consortium for Brain Mapping. With about one third of its researchers having an informatics background, LONI is remarkably interdisciplinary. Over the years, LONI has built up a very large (image) data repository and developed a large amount of processing and analysis tools. **ICES** is an interdisciplinary research and academic institution of UT Austin that brings together applied mathematics, informatics, and science and engineering to foster simulation as a key technology. Having been established in 2003 (with roots going back to 1993), ICES has been

a pioneering institution in terms of establishing structures for computational science and engineering at universities since then. Today, it involves 4 schools and 17 departments with an interesting mixed faculty-financing model (partly by ICES, partly by a single department). The core ICES funds have been provided via private donations. ICES offers an interdisciplinary graduate program, entitled *Computational Science & Engineering and Mathematics*, and it is extremely successful in acquiring funds and renowned faculty. It can be considered as one of the most attractive places for simulation worldwide. Closely related to ICES, **TACC** is the HPC center of UT Austin and the University of Texas system in general. Founded in 2001, TACC managed to position itself as a leading node in the US TeraGrid network, and it runs several large-scale machines. It is noteworthy that, as a part of TeraGrid, TACC is also the recipient of infrastructure and service-oriented projects funded by the National Science Foundation (NSF).

NCSA is a cutting-edge HPC site of the National Science Foundation (NSF). Remarkable investments have been made (by the NSF and the University of Illinois) towards an (academic) exascale facility, which was initially planned to host the IBM-designed Blue Waters computing system. Following IBM's and NCSA's termination of the contract, the path to a greater-than-10 petaflop sustained-performance Blue Waters system has recently been taken over by Cray. NCSA follows a clear "racks & brains" strategy: besides the infrastructure, NCSA is heavily involved in science and engineering projects addressing HPC methodology as well as compute-intense applications. **ANL** is part of the United States Department of Energy (DoE), but is administrated by the University of Chicago. As a so-called National Leadership Computing Facility, it employs about 3,000 people, 50% of whom are scientists. The upgrade of its computing infrastructure to a 10 petaflop system is planned for 2013. A noteworthy issue is the budget of the planned exascale software center, which will be about \$50 M per year, with additional funds for the development of applications. The supercomputing resources are provided as a major service to external scientists for free. The visit of **NCI** or, to be precise, its Center for Biomedical Informatics and Information Technology was centered around caBIG, a very large consortium or collaborative information network aiming at fostering and accelerating cancer research by exploiting complementary information and developing new data mining and image processing tools. Launched in 2004, caBIG has received increasing budgets since then. The **Harvard Medical School** or, to be precise, its Center for Biomedical Informatics, a research center that promotes and facilitates collaborative activities in biomedical informatics, was visited in Boston, MA. It was impressive to see how many resources and how much continuous support are dedicated to research, structures, and infrastructure components all aiming at clinical decision support to increase the patient care as well as to advance clinical and translational research.

Regarding the IT industries as the second large group of target institutions, the agenda comprised visits to the following companies: **NVIDIA** (Santa Clara, CA), **Apple** (Cupertino, CA), **Google** (Mountain View, CA), **Intel** (Santa Clara, CA), **IBM** (Almaden, CA, and Austin, TX), as well as **AMD** (Austin, TX).

NVIDIA considers the handheld market as an advantageous starting point for compute-intensive architectures (i.e., “the handheld is the future supercomputer”). They see GPU together with the ARM processors as an answer to the energy-efficiency issue. Concerning programming paradigms, future algorithm development will see the recalculation of data as a far more attractive alternative to data movement (i.e., “flops for free”). In contrast to that, **Intel** adopts an opposite strategy: “From the accelerators, we learned that hybrid programming works”, while seeing the future in “x86 goes heterogeneous”, with the final goal of having one common programming model. The breadth of Intel’s approach is impressive and appears to be appropriate for playing a leadership role also in the future. The combination of hardware development, software provision, and high-end manufacturing technology is considered as a crucial strategic advantage. While **AMD**’s basic strategy points in the same direction, given the convergence of processor and accelerator technologies and programming models, herein called “fusion” instead of “integration”, AMD has started serious efforts concerning software only recently.

Apple is generally seen as one of the most innovative companies and has undergone a kind of transformation from a computer company to an “expert in computers, applications, and entertainment”. The company’s focus is on usability and ease-of-use. As a consequence, there is not much room for scientific applications. An interesting issue was the iTunesU platform – an e-learning platform that allows non-profit organizations such as schools, museums, or universities to provide learning content. The visit to **Google** was probably the most unconventional. Google considers itself as a software and Internet company, and its corporate philosophy strives to generate a highly creative and efficient working atmosphere. There are various characteristics of a young and deliberately non-traditional IT company that made and make Google a highly attractive and successful place of employment. Finally, **IBM** was visited twice: at their research labs in Almaden, CA, and in Austin, TX. The worldwide network of large research labs certainly constitutes one of the strengths of IBM as a technology company. A major focus of the Almaden lab is on memory technology, ranging from material science via memory devices to parallel file systems. Furthermore, a couple of interesting new technologies are also explored such as new hardware designs emulating the brain’s capabilities (cognitive computing chips). The Austin lab is oriented more towards chip,

processor, and computer system technology, with the project PERCS (productive, easy-to-use, reconfigurable computing system) reflecting this focus, among other activities.

Finally, the IT-related branches of two of the major US funding bodies were also visited: the Office of Advanced Research of the Department of Energy (**DoE**, Washington, DC) and the Office of Cyber-Infrastructure (OCI) of the National Science Foundation (**NSF**, Arlington, VA).

Through its Office of Science, the **DoE** provides more than 40% of the total federal funding in the physical sciences in the USA. Within the Office of Science, the Advanced Scientific Computing Research program aims at discovering, developing, and deploying computational and networking capabilities to analyze, model, simulate, and predict complex phenomena important to the DoE. A particular challenge of this program is to fulfill the science potential of emerging computing systems, requiring significant modifications of today's tools and techniques to deliver on the promise of exascale science. The SciDAC (Scientific Discovery through Advanced Computing) activities have been a significant lighthouse measure, with structural aspects that can serve as an example almost worldwide. The **NSF** is an independent federal agency. It funds research and education in most fields of science and engineering and accounts for about 25% of federal support to academic institutions for basic research. Its directorate for Computer and Information Science and Engineering (CISE) strives to uphold a world leadership position in computing, communications, and information science and engineering. In addition to that, the Office of Cyber-Infrastructure co-ordinates and supports cyber-infrastructure resources, tools, and services such as supercomputers, high-capacity mass-storage systems, or system software suites and programming environments. The *Cyber-Infrastructure Framework for the 21st century* (CIF) regards cyber-infrastructure as an ecosystem and aims at building a national infrastructure for computational and data-intensive science and engineering, which also can serve as another model of activity for other countries.

In summary, the main findings appear to be:

- There is a very large **diversity of enterprise strategies** and **culture** between the IT companies visited.
- In **processor technology**, a couple of trends became obvious, such as on-chip-parallelism, hybrid architectures, and increasing investments in 3D processor technology.
- The success of **accelerators** such as **GPUs** in particular has led to the fact that the underlying hardware paradigms and the resulting programming models are on the

agenda of all large chip manufacturing companies, with a hybrid strategy pointing towards convergence. Nevertheless, there are different views on the future of supercomputers.

- The development of **high-end HPC installations** happens as a co-development of commercial providers and public research institutions. This involves very large investments in the systems and system software, with an increase also in usability in the widest sense, i.e. in particular applications and application software.
- In addition to being the drivers of the development of petascale and exascale systems (investments in “racks”), it is generally accepted that **algorithms and software** well suited for such systems do not automatically arise as a side-product, they also need large efforts and funding, too (investments in “brains”).
- The **data issue** is more and more considered as a crucial topic for Computational Science and Engineering (CSE) and HPC. The new notion of Computational and Data-Intensive Science and Engineering (CDSE) expands the established notion of CSE.
- **IT in medicine** has seen and sees a significant amount of large research consortia mainly in creating and maintaining data repositories as well as in providing distributed platforms and data analysis tools for patient-centered clinical and translational research and health care issues. The sensitivity for data security issues is increasing especially with respect to collecting medical and socio-economic data from social networks.
- **Medical IT infrastructures** include many people with informatics background, but are developed rather in medical environments with little exchange with dedicated CSE or HPC institutions. Once again, it became obvious that placing long-term tasks such as data repositories on the basis of short-term and project-based funds is highly problematic with respect to their sustainability. However, thorough reviews of long-term project progresses are nevertheless mandatory.
- **“Green IT”** draws more and more attention as power consumption increasingly constraints IT developments, especially but not exclusively within HPC.
- The large **funding agencies** in the US – DoE, NSF, and NIH – have strategic programs and significant funding for both CSE and HPC. In this regard, the situation in Germany is less developed. Although, meanwhile, DFG has started opportunities

for fundamental research in 2011, the recently launched priority program SPPEXA being an excellent starting point in that direction, there is still a strong need to continue the BMBF's more application-oriented HPC software program, just to make the investment in hardware successful. So far, the discrepancy between investments in "racks" and "brains" must be considered as problematic.

II. Background of the Study Tour

In the recent years, the Commission on IT Infrastructure (KfR) went on study tours to the USA and Asia to gain insights into the IT infrastructures at renowned research laboratories and universities as well as into the IT trends envisioned by the IT industry and into the IT funding schemes of the respective national funding agencies. These study tours proved to be effective and efficient for tapping into the scientific, technological, and administrative trends as well as the societal environment in which these developments take place. Crucial for the success of these tours were the on-site personal discussions with scientists and leadership of the research laboratories and the funding programs.

The KfR regards the success of these study tours as a central element for the quality assurance of its work within DFG, namely the continuation of the recommendations on IT infrastructure at universities and university hospitals, the assessment of IT concepts at these institutions, as well as the reviewing of proposals for IT procurements.

A delegation of the Commission revisited for the second time after its study tours in 2005 renowned universities, research institutes, funding agencies, and companies in the US. The intermediate six years easily span two IT innovation cycles. Concepts that were innovative in 2005 have now become standard procedures or have even been overtaken by newer ideas.

On its 2011 study tour, the Commission aimed at obtaining first-hand information about current IT-related scientific and technological developments. Focal points of the visit were:

- *IT at universities*
 - scientific infrastructure and environment
 - computer systems
 - concepts for the improvement of software development in scientific disciplines
 - trends in IT-enabled research fields such as Computational Science and Engineering
 - IT infrastructure and governance
- *IT in hospitals and medical research*
 - scientific and medical applications
 - computer systems in medical applications
 - organization and accessibility of very large medical data repositories
 - platforms for translational and basic research and health care

- *IT in computing centers and networks*
 - High-Performance Computing (HPC)
 - distributed computing
 - challenges of the extreme parallelization of both hardware and software
- *IT funding bodies*
 - strategies and trends in research and education
 - setup and management of distributed resources (grids, clouds)
- *software and hardware industries, start-ups*
 - technology trends and developments
 - corporate strategies of IT companies
 - collaboration with academia
- *(international) co-operation between universities, research institutes and industries, as well as funding agencies*

III. General Information

Participants



For the Commission on IT Infrastructure

Prof. Dr. Jörg Becker, University Münster

Prof. Dr. Dr. Johannes Bernarding, Otto-von-Guericke University Magdeburg

Prof. Dr. Hans-Joachim Bungartz, Technical University Munich (Head of the delegation)

Prof. Dr. Markus Clemens, University Wuppertal

Prof. Dr. Christel Marian, University Düsseldorf

Prof. Dr. Hans Ulrich Prokosch, University Erlangen-Nürnberg

Prof. Dr. Wolfgang E. Nagel, Technical University Dresden

Prof. Dr. Thomas Tolxdorff, Charité Berlin (former member of the commission / guest)

For the DFG Headquarter

Dr. Werner Bröcker, Scientific Instrumentation and Information Technology

Dr. Max Vögler, director and host, DFG Office North America/Washington

Dr. Marcus Wilms, Scientific Instrumentation and Information Technology (organizer)

Agenda (March 20 – April 2, 2011)

#	Date	Location	Institute/company/agency
1	21.03.	Berkeley, CA	NERSC National Energy Research Scientific Computing Center
2	21.03.	Santa Clara, CA	NVIDIA
3	22.03.	Cupertino, CA	Apple
4	22.03.	Mountain View, CA	Google
5	22.03.	Santa Clara, CA	Intel
6	23.03.	Almaden, CA	IBM
7	23.03.	Los Angeles, CA	LONI Laboratory of Neuro-Imaging, UCLA
8	24.03.	Austin, TX	AMD
9	25.03.	Austin, TX	ICES Institute for Computational Engineering and Sciences The University of Texas at Austin
10	25.03.	Austin, TX	TACC Texas Advanced Computing Center
11	25.03.	Austin, TX	IBM
12	28.03.	Urbana, IL	NCSA National Center for Supercomputing Applications
13	28.03.	Chicago, IL	ANL Argonne National Laboratory
14	29.03.	Washington, DC	DOE Department of Energy, Office of Advanced Research
15	30.03.	Arlington, VA	NSF National Science Foundation, Office of Cyberinfrastructure
16	30.03.	Rockville, MD	NCI National Cancer Institute
17	01.04.	Boston, MA	HMS Harvard Medical School

“Disclaimer”

All information and data given in this report were carefully collected and reflect the results of our on-site discussions as seen by the DFG Commission on IT Infrastructure.

IV. Summaries of the Visits

1. NERSC, Berkeley, CA

The **National Energy Research Scientific Computing Center** (NERSC) is a division of the Berkeley Lab of the Office of Science of the U.S. Department of Energy (DoE). NERSC employs 65 people and has an annual budget of \$65 M, of which \$3 M is spent for its power consumption. NERSC offers computing power to research institutions and other non-governmental organizations. Every three years, NERSC invests in a new computing system. The current system NERSC 6 (“Hopper”) was installed in October 2010. It is a two 12-core Magny-Cours per node, 68 racks Cray system with 150,000 cores, 1.25 petaflops, and 3.3 MW power consumption. NERSC supports about 400 projects, 4,000 users, and 500 code instances. A total of 65% of its available computing power is granted to universities, while 25% of the compute power is shared among organizations funded by the DoE.

The main research projects lay in the fields of HPC, especially, energy-efficient computational science: “more science per Watt”), end-to-end computational methods (i.e. from experimentation to analysis: data-enabled science), predictive computational science (i.e. validating and quantifying uncertainty), Applied Mathematics (i.e. mathematical modeling), Computer Science (i.e. algorithms and data structures, as well as languages, compilers, and tools), and Computational Science (i.e. biology, climate, astrophysics, environment, cosmology).

The **Magellan project** focuses on the possibilities and limits of cloud computing; the needs and the features of a science cloud are explored therein, as well as the appropriate role for commercial and/or private cloud computing services for the DoE. It tries to answer various questions, such as: What is the appropriate type of application? What are the best fitting programming models? And: What are useful cost calculation models? NERSC is one of the very few institutions visited where cloud computing is a main research topic.

NERSC faces the difficulty to transform “research codes” into “production codes”. Researchers are interested in quick research results, not in widespread dissemination. This seems to be a worldwide problem.

NERSC combines its offer of computing power with one of the most powerful computing systems. Being a remarkable research institution, the publication record of the investigators

in the 400 projects is excellent (1600 publications based on the usage of the NERSC system in 2009).

2. NVIDIA, Santa Clara, CA

The commission's visit to **NVIDIA** included presentations on the Tesla hardware roadmap, on the GPU-oriented C-dialect CUDA roadmap, and on further software development tools. A detailed presentation on quantum chemistry served as an example of how HPC applications can be adapted to NVIDIA GPUs. Furthermore, a description of NVIDIA's exascale and ARM strategy as well as its CUDA architecture roadmap were given and concluded with an open discussion with a particular emphasis on the human resource strategy of NVIDIA.

In the remarkably open presentations the emphasis was on the envisaged eminent role of the **handheld devices** based on NVIDIA's CUDA architecture (represented by the NVIDIA Tegra processor family), whose future significance ("The handheld is the new supercomputer") was seen in a strong competitive situation with the current Intel desktop/laptop orientation.

This stance was justified by the higher sales figures on the consumer electronic market, the specific requirements on the energy efficiency of processors in this application segment, to be achieved through the use of GPU processor technology, particularly in combination with energy-efficient ARM processors and the special dynamics of development based on open systems such as Android and the NVIDIA proprietary CUDA parallel architecture. CUDA is not seen as a proprietary programming language alone, but also as a heterogeneous co-processor architecture with applications scaling from "super phones" to supercomputers.

In the future, an opening of the CUDA standard is eventually being considered, while the OpenCL standard is hampered by its slow development process and the number of different interests involved in its definition and standardization process. The recent open GPU directive standard OpenACC seems to provide an easy to use way to accelerate applications.

Particularly in the handheld device market NVIDIA assumes that Intel and AMD lack a suitable strategy, while in contrast, NVIDIA's strategy is clearly manifested technically and economically, and on this basis further HPC developments can be derived and financed. According to NVIDIA the HPC sector benefits from a further increase of the GPUs' efficiency, especially regarding the double-precision performance and the GPUs' enhancements towards a highly energy-optimized general-purpose processor.

With regards to the HPC programming paradigms with an emphasis on power consumption, mathematical algorithm development for the recalculation of data is to be favored in the future over data movement (“flops for free”).

Besides short-term additions to the annually improved NVIDIA GPU processor lines involving Fermi (current architecture), FERMI+ (2011), Kepler (>1TF double-precision 250+ Gb/s, announced for 2012), and Maxwell (15 GF/W double precision, announced for 2014), further developments were presented from the DARPA-financed “Ubiquitous High-Performance Computing” project Echelon. This project involves NVIDIA GPUs as part of future GPU-based exascale-HPC systems, with double-precision performance of 16 TF and a data throughput of up to 1000 GB/s applying the 3D-stacking technology. Due to their higher flops/W ratio, the GPU-based systems are seen as the best candidates for setting up an exascale system.

Approximately every two years larger redesigns of the processors are implemented, and the close contacts of NVIDIA with game developers are considered as a strategic advantage also for other applications.

In the area of human resource development NVIDIA stated that with a workforce of currently 6,000 employees there are 715 vacancies in the development division. Concerning this matter, the NVIDIA recruitment policy prerequisites high-level professional excellence of applicants, while still facing competition from rival companies and a declining percentage of students in information technology on the U.S. labor market.

3. Apple, Cupertino, CA

By employing 46,600 people and reaching sales of \$65 B in 2010, **Apple** is one of the big players. It is seen as one of the most innovative companies (in most rankings among the top three) and is regarded as number one in customer support (according to the journal Consumer Report). In the recent years, Apple has changed from being a computer company to being an “Expert in computer, applications, and entertainment”. Nonetheless, it has drastically reduced the product portfolio down to four product groups: the iPhone with a market share of 39% (growing), Music with a market share of 20%, the iPad with a market share of 8% (growing), and the Mac with a market share of 30% (shrinking, yet still very profitable, having a faster growth than the market). The iPhone, the iPad, and the iPod run the same iOS Operating System, and Apple has sold 160 million devices of these products (up to the time of our visit). This platform strategy is one of the success factors of Apple.

Apple defines itself as a product company with a strong focus on usability, user friendliness, and ease-of-use. The main functionality for a “normal user” should be as easy as possible. There is hardly any room for scientific applications. The success in the application area is based on the so-called apps, software applications mainly provided by third-party developers.

Although Apple does not develop scientific or clinical applications, their handhelds and smartphones seem to be suitable devices for accessing such applications.

Apple’s research strategy focuses on topics closely related to their products and their philosophy (user interface). It co-operates only with very few universities.

Concerning the Open Source policy, Apple supports the BSD concept.

Apple offers **an e-Learning-platform** (iTunesU), which allows non-profit organizations such as schools or museums to provide learning content including video and audio (from free access to limited access, depending on the providers’ choice). A number of 800 providers in 26 countries offer e-Learning content to 90 countries as users. A number of German universities are already users of this system.

Apple’s success is based on its innovation in user friendliness and the combination of information (computer) and entertainment (music, games), while having a high innovation rate in this area. This strong focus results in the fact that hardly any scientific, clinical, or business applications are being offered by Apple.

4. Google, Mountain View, CA

Google considers itself as a software and Internet company with a main focus on Internet search services. Recently, additional internet-based services such as Google Maps, Google Earth, Street View, Google Mail, and others are also provided. Google has about 25,000 employees worldwide (about 7,000 working on the campus in Mountain View in more than 60 buildings). The internal organization is characterized by numerous employees services (free food, shuttle service to San Francisco, and sports facilities) and a high degree of flexible working conditions including a home office. Google’s corporate philosophy strives to create a highly creative working atmosphere with meetings usually lasting not longer than 30 min. Energy-efficiency and optimized computer infrastructures are important topics for Google, and a large amount of effort is spent on researching the transfer and the efficient use of electrical power. Google states to require about 50% less energy compared to other standard

computer centers. Each of the 10 large computer centers each require about 5 MW and are installed near population dense areas reaching a power efficiency of 1.25.

Google is interested in co-operating with scientists ([Link](#)) and grants awards and stipends as well as smaller joint research projects (requires a partner within Google) towards this interest. Cloud Computing is seen in the realm of other large companies such as Amazon. Search algorithms are still the main focus of Google research. Other research themes include operating systems (Android) and browser technology (Chrome). There is no specific scientific software, while a free email service is offered to universities.

Google was the most extreme example of a young and deliberately non-traditional IT company. It seems that it has started to realize that such a company model might have problems in scaling with respect to size.

5. Intel, Santa Clara, CA

Intel is the world-wide leading CPU/processor manufacturer with a turnover of \$43.6 B, an annual R&D budget of \$5 B, and 82,500 employees. Intel's activities are structured into four main areas: (1) IA/x86-based products (comprising all hardware activities, from "ultra-mobile" to "data centers" and HPC), (2) software and solutions (a sector of increasing technological and business relevance), (3) technology and manufacturing ("in-house production", also a characteristic feature), and (4) the Intel labs (i.e. corporate R&D).

The **breadth of Intel's approach** is impressive and seems to also be appropriate for playing a leadership role also in the future. The combination of hardware development, software provision (including substantial investments), and front-end manufacturing technology is considered as a crucial strategic advantage compared to its competitors.

There have been several statements and questions that nicely illustrate Intel's strategy: *"Our processor research always has to be four generations ahead"*, and *"our strategy is to anticipate trends, not to react"*, *"power is the key challenge, data movement is the key contributor"*, *"where is the added value of reconfigurable architectures?"* The guiding paradigm seems to be *"integration"*, i.e. driving both classical CPU (Sandy Bridge, Haskell, and others) and accelerator technologies (Knights X), while integrating them with the goal of enabling future processors to support the full range of serial and parallel programming models to minimize developer re-training and to allow for more code re-use – *"x86 goes heterogeneous"*, but with the goal of *"one programming model"*.

Concerning **HPC** strategies, Intel has identified the interplay of hardware and software as well as energy consumption as core challenges (energy per operation, extreme concurrency and locality, resiliency, memory (capacity, bandwidth, and power, and hardware-software co-design). HPC is considered as a field of very high strategic relevance for Intel.

As a leading global player, Intel is also engaged in numerous university collaborations, mainly organized through the Intel labs. Such collaborations occur at a corporate level (mainly with leading US universities, such as Berkeley, Stanford, Carnegie Mellon, or Illinois, but also with selected institutions in Europe and Germany), through individual grants, or through internships.

Intel has been and is exposed to several potential threats (such as AMD as a classical competitor, NVIDIA as a recent competitor questioning the “classical CPU approach”, Chinese companies having entered the business and others), while certainly having a strong response strategy.

6. IBM, Almaden, CA

IBM has 10 research labs worldwide with approximately 3,000 scientific researchers, including several Nobel laureates. The IBM Research Labs at Almaden, CA employ about 400 scientific researchers (more than 50% of whom hold a Ph.D. degree) and are considered to be among the most innovative IBM research centers. Each IBM research lab concentrates on specific areas of research.

The Almaden IBM Research Labs develop important new results for progress in information technology and aims at an in-depth understanding of the underlying basic natural sciences. IBM has stated its desire to maintain its leading role in the field of information technology. This involves the application-driven research of chemical processes and materials (materials for high-resolution photograph-based lithography, e.g.), new materials and concepts for future technologies (e.g. memory cells based on magnetic tunnel resistances or organic polymers) as well as fundamental research in areas with high prospective potential for applications, such as for instance nano-scale technology or nano-scale medical research.

IBM has shown thus far a certain preference for prestigious projects that draw public attention. For example, in November 2010, the supercomputer "Watson" won against human competitors in "Jeopardy", an American TV quiz show, which requires considerable language cognition abilities and a reliable assessment of the answer quality.

The United States National Academy of Engineering (NAE) identified in 2010 the grand challenges of the 21st century in 14 engineering sciences and asked companies, scientific institutions, and individuals to enter into a competition for solutions ([Link](#)). The IBM Research Lab Almaden addresses several of these "Grand Challenges for Engineering". This includes projects aiming at the understanding and reverse engineering of the human brain (with the purpose of building a biologically inspired chip) as well as enabling access to clean water and to efficient energy storage.

A central research effort of the Almaden IBM Research Labs concentrates on the **optimization of memory systems**. This involves research projects concerning memory devices, controllers and networks, memory and system management software and parallel file systems. One of the latest IBM flagship projects is the "General Parallel File System" (GPFS), a scalable, parallel file system that was originally planned by IBM as a video service system. The GPFS is well known for its excellent I/O-throughput capability and its high-level system availability. Its most recent installation at the Lawrence Livermore National Laboratory, on a 1536-node cluster with 100 teraflops compute performance and a total of 2 PB GPFS disc storage, achieves an I/O throughput rate of up to 126 GB/s. GPFS conforms with the POSIX standard and, hence, is also suitable for several different platforms, such as AIX, Linux, and Windows. GPFS market segments are up to 50% in research labs and in companies with small compute clusters.

The current IBM roadmaps includes the Panache Global File System (scalable, high-performance, clustered file system cache for parallel data intensive applications), the Perseus Software Raid (activities for recovery after data loss distributed over several discs and thus considerably faster than established standard RAID-systems; planned to be available starting May 2011), improved solid-state storage (short latency time and thus an increased throughput and higher reliability), and support for Apache Hadoop (open-source software for reliable and scalable distributed computing). However, IBM indicates that this roadmap may always be subject to change.

Acknowledging the economic importance within the IT sector, Almaden Services Research gave the impulse for a new interdisciplinary study course curriculum named "Services Sciences, Management and Engineering (SSME)". IBM establishes this course together with universities and other companies: In Germany, the Karlsruhe Service Research Institute at the Karlsruhe Institute for Technology (KIT) in co-operation with IBM Germany offers an interdisciplinary Ph.D. course in SSME. Its objective is to build up capabilities for IT-based services (to provide compute performance and memory capacity via introduction of cloud computing, e.g., but also via software solutions).

IBM has an established tradition of offering integrated solutions both in hardware and in software solutions. In the IBM "Smart Planet" project persons, computer systems and objects equipped with measurement instruments are supposed to communicate with each other. Computer simulations based on realistic models result in situation-optimized, intelligent solutions. An exemplary scientific study shows that cars waste approximately 10% of their fuel while their drivers are in search of a parking place. An intelligent communication system helps in avoiding this waste and in preserving time and fuel resources.

IBM research labs have world-wide co-operation with about 5,000 universities, approximately one third of which are located in the USA, another third are located in other industrial countries, and the last third are located in developing countries. Within the framework of the IBM University Award program IBM annually invests about \$100 M in co-operation projects (\$10 M in cash, \$60 M in hardware, and \$30 M in software and grants). As a result of this program, IBM observes a 30-fold return of investment after five years. This success could also inspire German companies to increase their co-operation with universities. Strategic analyses of IBM also show that the amount of HPC investments of (American) universities have a direct relation to the university's ranking which is determined based on factors such as the total amount of third party funding, the number of publications in high-ranking journals, and the number of spin-off companies, to name a few.

7. ICBM at UCLA, Los Angeles, CA

The **Laboratory of Neuro-Imaging** (LONI) is an independent research center affiliated with the Department of Neurology (School of Medicine at University of California, at Los Angeles, UCLA) having separate financing. The initiator of the project and several of his co-workers presented the project and research topics in a roundtable talk. The main goal of the 17-year-old project is the continuous build-up of a database containing the anatomy, the histology, and the functions of normal and pathologic brains (**International Consortium for Brain Mapping**, ICBM). The main medical focus is on the Alzheimer and Parkinson diseases.

Only two of the 115 positions are permanent (financed by UCLA), while the remaining positions are being financed via third party grants (85% from NIH, various other donators and foundations such as Microsoft and MJ Fox Foundation, while few from NSF funding). The overall annual budget was formerly in the range of \$8-10 M, which recently dropped to \$6.5-8 M. An amount of 45% of the grants is transferred to UCLA. A percentage of 30-40% of the researchers has a background in informatics. Similar to the situation in other countries including Germany, the publications of IT researchers rather appear as conference abstracts

than as journal articles, which is seen as problematic since grant reviewers with a medical background require publications in high-impact journals. International co-operation includes Forschungszentrum Jülich, Germany, and the McConnell Brain Imaging Centre (BIC) of the Montreal Neurological Institute, McGill University, Canada, and others. There are no co-operation projects between LONI and larger IT companies. Due to the high quota of external funding, the project requires permanent large amounts of successful third party support.

The **data repository** contains image data of brain scans including Diffusion Tensor Imaging as well as brain histology and genomic information. Data were acquired at different sites. LONI solved the difficult problem of consistent data acquisition by defining the acquisition protocols. Data are merged at LONI (presently 5000 scans of 50 sites including genomic information). Data security is realized by eliminating all personal information and by data agglomeration (for genomic data). For the histologic specimens about 30 post-mortem brains were sent to the co-operation partners in Jülich to prepare the slices using special cell and receptor staining techniques before analyzing the specimens. New probabilistic algorithms and descriptions were developed at LONI and in Jülich. The project is interested in general patterns of a population rather than in individual representations. ICBM also collaborates also with the *Mouse Brain Project (Connectathon)*, in which nerve connections between different brain parts are analyzed and visualized, among others, using diffusion tensor imaging.

LONI has a 1000 nodes cluster (administrated independently of UCLA). Data and software tools to analyze and visualize the data are provided by LONI, while third parties are now developing software tools. Formerly, LONI developed an image processing tool based on a pipeline concept. No special HPC software is developed, but Grid concepts are used for data analysis.

Altogether, there are 30 projects that need large data storage resources. For example, the ADNI project alone requires 3 TB of storage. However, to query and retrieve the data, meta-data have to be generated and attached to the image, while data are derived prior to storage. There are no solutions for *image-* or *content-driven* data retrieval. Experience with external software development also shows that co-ordination among the external developers is required to guarantee harmonized and cost-efficient software frameworks that can consistently be used by the entire community. Similar to the situation in other countries, this project shows that data repositories require reliable and long-term financial support.

8. AMD, Austin, TX

AMD is the world-wide second largest CPU/processor manufacturer. After some organizational restructuring, namely the outsourcing of production sites, AMD now concentrates on the development of chip designs as well as on software aspects. Promising design roadmaps include architectures down to 28 nm in 2013 based on the “bulldozer” monolithic dual-core building block that supports two threads of execution. The bulldozer family will have 16 or more cores on one chip (Interlagos) in future.

Together with other important players, AMD stresses the increasing importance of HPC applications. In order to provide the necessary compute power, AMD introduces its AMD Fusion Architecture **Accelerated Processing Unit** (APU) that integrates CPU, GPU, and other accelerators onto a single die. This seamless integration might be a nice feature for programming and could be an architectural advantage, compared to other vendors. Regarding the software, AMD strongly supports the OpenCL standard.

The design optimization includes AMD’s **Application Power Management** (APM), which controls the power consumption of the CPU or core, depending on the actual needs of active applications.

AMD co-operates with companies from the graphics, games, and entertainment sector.

Maintaining its own education and developer outreach program, AMD provides free software tool kits (e.g. OpenCL courses, libraries, training, support) in order to stimulate and propagate new uses of GPU computing and graphics, especially at universities. There was clear interest to intensify co-operation between AMD and German universities.

9. ICES, University of Texas, Austin, TX

The **Institute for Computational Engineering and Sciences** (ICES) at The University of Texas at Austin is an interdisciplinary research and academic institution that brings together applied mathematics, informatics, and science and engineering to foster simulation as a key technology. Established in 2003, the ICES thus extends the focus of its predecessor, the Texas Institute for Computational and Applied Mathematics (TICAM) that goes back to 1993. The ICES mission is stated as: *“To provide the infrastructure and intellectual leadership for developing outstanding interdisciplinary programs in research and graduate study in areas of computational engineering and sciences and in information technology”*.

The ICES has been a pioneer, both in terms of its interdisciplinary mission and its interdisciplinary organization. It involves 4 schools and 17 departments, which also partially contribute to financing the positions. In fact, the ICES faculty has typically at least one additional affiliation with a school or a department. The relations between ICES and the schools have not always been easy, but are now considered as good. Significant funds for ICES (more than \$113 M) have been provided via private donations (the building, to fund 15 chaired positions, to support an extensive visiting fellow program, to offer scholarship programs, and others), which also marked the beginning of TICAM as well as the re-design towards ICES. Project money primarily comes from public funding agencies (such as NSF, DoE, and NIH), while the number of industry collaborations is smaller (for industries such as aerospace, oil, biomedical, communication, and IT). Overall, ICES is extremely successful in acquiring funds, which is also due to the fact that it runs a strictly quality-based but conservative development policy (“small, but excellent”).

The ICES offers an interdisciplinary graduate program (M.Sc. and Ph.D.), namely the *Computational Science & Engineering and Mathematics* (CSEM), for which an intense recruiting program has been established. The design of the program also reflects ICES’ interdisciplinary mission. The relatively small number of 68 CSEM students must be seen in an overall relation to the 89 faculty members, of which 36 are called ICES core faculty. The ICES operates a number of nine local compute clusters, and it has direct access to the equipment of TACC, the *Texas Advanced Computing Center*, such as terascale computing systems, as well as the ACES 3D visualization lab.

Examples of research groups are Computational Materials, Computational Hydraulics, Distributed and Grid Computing, Computational Geosciences, Predictive Engineering and Computational Sciences, and Computational Life Sciences and Biology.

The ICES is an extremely visible and successful early model for organizing **cross-school and cross-department research and education** in a modern field. Indeed, ICES managed to synergistically combine high reputation in modeling, algorithmics, and IT infrastructure. When thinking of establishing a certain type of a CSE research center, most institutions worldwide take ICES as an example of best practice.

10. TACC, Austin, TX

The **Texas Advanced Computing Center** (TACC) is the HPC center of the University of Texas (UT) at Austin (and of the UT system in general). Founded in 2001, it is operated by UT Austin and currently has more than 80 employees. TACC hosts several large-scale

computing machines (such as the HPC systems *Ranger* (600 teraflops, at time of installation being one of the largest machines for free science access) and *Lonestar* (300 teraflops) or the visualization system *Longhorn*.

TACC managed to position itself as a leading node in the US TeraGrid network, which besides technology leadership and visibility also brings significant NSF funds for the cutting-edge systems. It is noteworthy that NSF also funds mere infrastructure or service-oriented projects. TACC's rise to its current position occurred rather quickly, which also reflects the strong political will of UT Austin (certainly also driven by the ICES success) to become a major player not only in HPC applications and algorithmic, but also in terms of HPC technology and facilities. Consequently, TACC's mission reads rather ambitious: to be a place of most advanced HPC and visualization technology; to provide leadership in that technology; and to enable transformational science. In fact, there is the impression of a very much strategically driven and hierarchical structure.

TACC's role in research is primarily a supportive one. In addition, there are also a couple of TACC-hosted research activities, such as applications (in computational biology), algorithms and tools (GotoBLAS2, MyCluster, and SciVis, ...), and on cooling technology, which leave an impression of heterogeneous research, to some extent.

TACC's involvement in education is rather small, which is typical for a computing center. There are four courses given on a regular basis that are part of existing UT Austin programs, also supporting the idea of dissemination of computing technology.

Problems may arise from TACC's previous close relations with SUN, due to SUN (now ORACLE) having left the HPC field.

11. IBM, Austin, TX

The subjects addressed during the visit to the **IBM** Research Labs in Austin included an overview of the PERCS software and hardware. The PERCS ("Productive, Easy-to-use, Reconfigurable Computing System") project is related to a DARPA project for a 10 petaflops HPC system based on the IBM Power7 chip architecture. Following these presentations, a report on IBM's internal pre-development for 3D chip technology was given, succeeded by a visit to the assembly production facilities of the Power-7 PERCS systems.

The **IBM Software PERCS** focuses on providing complete solutions for all IBM systems in the area of resource and system management, application development (Eclipse-based HPC

Workbench), support software for the operation (load leveler with special functionality for very large clusters), and solutions for data and file management.

IBM presented in detail the existing tools, including in particular the “Eclipse Integrated Developer Environment” (IDE) that aims to help compensate the increasing divergence of "sustained performance" and "peak performance" in HPC systems. The HPC cloud-management software suite, announced for 2011 will also include virtual-machine techniques and tools.

The **IBM PERCS Hardware** is based on the Power7 technology and was presented as a “green data center in a rack” with respect to its performance density and its high energy efficiency. The performance data of one rack add up to approximately 100 teraflops out of 3072 cores within 3 x 128 Power7-processors running at up to 4.0 GHz clock frequency. These processors are collected into three so-called “super nodes” per rack, each consisting of four blocks featuring 32 processors (a total of 256 cores). Each rack may be equipped with a main memory of up to 24.6 TByte/rack. Such a rack has similar performance figures as the total “IBM ASCI Purple” HPC-system of 2005 (with 91 teraflops out of 258 racks equipped with Power5 processors).

These Power7 processor rack systems specifically feature water-cooling of processors, memory banks, and others, and, thus, have a weight of approximately 2.7 tons per system. At the time of our visit, these racks were the essential system components of the planned “Blue Waters” of above 10 petaflops peak / 1 petaflops sustained performance HPC–system design in the National Petascale Computing Facility of the National Center for Supercomputing Applications (NCSA) at the University of Illinois at Urbana-Champaign, IL. This HPC system was planned to consist of 114 Power7 racks (as a number of 38 3-rack blocks) with a total number of 38,912 Power7 processors and 43,776 hard disc drives with costs of up to \$208 M. However, in August 2011, IBM and NCSA terminated the Blue Waters contract due to difficulties in financing the rather complex and expensive system. IBM pointed out, that this contract termination was an isolated customer-individual decision. It would neither impact the general availability of IBM Power7 HPC solutions nor the development of the entire IBM PERCS software stack.

The presentation on “**3D chip-technology**” that followed highlighted IBM’s research and development efforts concerning an important feature of future hardware, where IBM is currently in a leading position in this field of research. The 3D technology is considered to be able to solve problems related to on-chip data transport. At the moment, however, the emphasis at IBM is still placed on the reduction of energy consumption with a reduced I/O-

load and reduced cooling power requirements. The improved energy-efficiency is yet preferred over a possible increase in memory bandwidth enabled by short interconnect distances and reduced latencies due to a higher memory density. The 3D technology is expected to allow increased cache memories. Open problems, such as those concerning the cooling concepts of these chips, however, require this technology to mature even further over several years before becoming part of an industrial production process. According to IBM the production necessitates access to company-owned fabrication facilities.

12. NCSA, Urbana, IL

The **National Center for Supercomputing Applications** (NCSA) at the University of Illinois at Urbana-Champaign (UIUC) is a leading edge HPC site of the National Science Foundation (NSF). Remarkable investments have been made towards an (academic) petascale computing facility that was initially planned to host the IBM-designed Blue Waters system (cf. previous visit to IBM, Austin). The path to a >10PF peak / 1PF sustained performance Blue Waters system has recently been renewed by Cray Inc. stepping in as an industrial partner ([Link](#)). Besides providing the infrastructure, NCSA is heavily involved in science and engineering projects in close co-operation with the application development groups outside of NCSA, supporting more than 30 applications, while addressing the full spectrum of scalability issues, ranging from I/O, data intensive workflows via programmability and debugging to basic algorithm development. A significant number of NCSA project investigators (PIs) are members of the International Exascale Software Project (IESP) group and work on addressing the software-related challenges in exascale computing. NCSA has hired several PIs to support the whole Blue Waters project on all levels. Overall, about 175 staff members and roughly 100 students are involved in NCSA's project activities.

NCSA will be an important yardstick on the way to **exascale**. It has roughly the same structure and the same objectives as the German Gauss Centre for Supercomputing (GCS). While the investments in infrastructure and computing are on the same level, the focus at NCSA is much more on people and support that is funded by projects, to address the HPC challenges of the future. Overall, NCSA is one of the most important HPC centers in the US, in addition to the Oak Ridge National Laboratory (ORNL), the Lawrence Berkeley National Laboratory (LBNL) with NERSC, and Argonne National Lab (ANL), which are all DoE-funded institutions. Being embedded into a leading university (UIUC) is a definite advantage, even though the IT organization of UIUC does not operate smoothly and is not leading edge. The connection between UIUC and NCSA is organized via UIUC's *Institute for Advanced Computing Applications and Technologies* and the *Center for Extreme-Scale Computation*,

which address an impressive range of relevant topics (software tools and applications, computing technologies, multiscale simulation, information systems, computing and creativity, etc.).

Despite the high reputation of NCSA and of UIUC's Computer Science (CS) department, the overview of the IT infrastructure of UIUC (and of the UI system) revealed several problematic issues, including a change of the CIO's embedding into UI's leadership structure involving a change between reporting to the Chancellor and reporting to the (non-academic) Chief Financial Officer, which in fact caused the campus CIO to resign in April 2011; rather arbitrary decisions on what to operate centrally and what to run locally; a rather high number of UIUC staff (about 700 people working in IT, among which about only one third on a central/university payroll).

There is quite some uncertainty about the time period following the Blue Waters system installation. While the future of NCSA as an institution seems to be ensured, it is unclear whether NSF will continue in supporting leading-edge HPC systems, in parallel to DOE's activities. Discussions with NSF officials allow concluding that NSF funding will target in an even more algorithmic- and application-oriented direction.

13. Argonne National Laboratory, Chicago, IL

The **Argonne National Laboratory** (ANL) as part of the Office of Science of the Department of Energy (DoE) is administrated by the University of Chicago. The ANL is a National Leadership computing facility of the DoE, and consists of several divisions. The Commission on IT Infrastructure met with the Mathematics and Computer Science (MCS) Division. About 3,000 people work at ANL, about 50% of them are scientists. The annual budget is about \$50 M and is funded by the DoE. The ANL does not provide academic degrees while accepting summer students from all over the country. The institution is also open to all scientific co-operation partners while the project results must be accessible to the public. A main focus is on environmental science.

About 350 active users use the IT infrastructure, consisting of the Intrepid system with BlueGene/P CPUs and the Eureka system with 200 NVIDIA GPUs. An upgrade to the BlueGene/Q system (10 petaflops peak performance, 49,152 nodes, 786,432 cores, 786 TB memory) is planned for 2013. Other projects include energy storage systems for electrical cars, short pulse X-rays, and development for new catalysts for propylen oxidation. The Laboratory's Computing Resource Center was founded in 2002 to develop and support the HPC-projects and the HPC facilities. About 450 users and 45 projects, each project having

an ANL-affiliated PI, use the actual hardware (320 nodes: 304 regular nodes with 36 GB memory each, 16 big nodes with 96 GB memory each).

The applications cover modeling and simulation at petascale including a roadmap to exascale (exascale program with hardware platform for about \$200 M in 2015 and 2018). The annual budget for the planned exascale software center is about \$50 M, with additional funds for development of applications (in nuclear engineering, fusion, computational chemistry, climate modeling, combustion, high-energy density physics, and materials, and in collaboration with National Nuclear Security Administration). A percentage of 50% of the main users are physicists, while the other users come from material sciences, engineering sciences, life sciences, and chemistry. There is a small number of biomedical HPC projects. Another important project targets encapsulated test beds for non-expert users to develop HPC applications (immediate account of ~8,000 hours) including hands-on training, expert consulting, and user support. Since 2002, there is also a long expertise in Grid- and Cloud computing ([Link](#)). Additional projects investigate mathematical problems, data analysis and visualization methods, including stochastic and probabilistic methods. Several libraries for HPC applications have also been developed.

The projects are divided into three categories depending on the project duration (1-3 years, 1 year, less than 1 year) and are evaluated by an internal panel ([Link](#)). External users receive support from internal co-operation partners. Despite the excellent hardware and software environment, certain applications, especially in physics, impose new challenges that require new solutions for real-time data analysis and for fast data storage capabilities.

14. DOE, Washington, DC

The United States **Department of Energy** (DoE) with a budget of \$29.5 B for FY 2012 ([Link](#)) is a cabinet-level department of the US government concerned with the US' policies on energy and nuclear safety. Through its Office of Science, which has a budget of \$5.4 B for FY 2012 ([Link](#)), DoE provides more than 40 percent of total Federal funding in the physical sciences in the US. It oversees, and is the principal Federal funding agency of, the Nation's research programs in high-energy physics, nuclear physics, and fusion energy sciences ([Link](#)). As a mission agency, DoE is lead by strategic priorities to advance discovery science and invest in science for the national needs in energy, climate, and the environment. It does so by calling for proposals in areas that fit the current and the prospected national needs. The DoE funding commitments are often longer-term compared to those made by the NSF. This allows for implementing and maintaining large scale national research laboratories.

The mission of the Office of Science Advanced Scientific Computing Research (ASCR) program is to discover, develop, and deploy computational and networking capabilities to analyze, model, simulate, and predict complex phenomena important to the DoE. A particular challenge of this program is to fulfill the science potential of emerging computing systems and other novel computing architectures, requiring significant modifications to today's tools and techniques to deliver on the promise of exascale science ([Link](#)).

In order to accomplish this goal, a major initiative of the Office of Science was to initiate the **Scientific Discovery through Advanced Computing** (SciDAC) program in 2001 as a partnership involving all of the Office of Science program offices (\$53 M budget in FY 2010). The SciDAC-funded institutes place a greater emphasis on integrated multi-disciplinary teams from universities and laboratories that work in close partnership. They are meant to facilitate collaborations of applied mathematicians, computer scientists and domain scientists to solve complex problems in physics, chemistry, and biology. The program also includes research on improved mathematical and computing systems software that will allow these codes to effectively and efficiently use modern massively parallel computers. Additionally, the program will develop "collaborative" software to enable geographically separated scientists to effectively work together as a team, to control scientific instruments remotely, and to share data in a more readily immediate manner.

The proposals in exascale-related areas in FY 2010 have requested \$145 M of funding of which \$18 M have been granted. Three SciDAC Co-Design centers for computational partnerships were funded in FY 2011, at a total of \$12 M/year ([Link](#)).

The SciDAC goals closely resemble those put forward by the DFG IT Commission's recent position paper "Software for Exascale Computing" which eventually led to the establishment of a new DFG Priority Program with the same title. The program is scheduled to span six years and will be funded by up to €4 M per year. The funding requirements as well as the appropriate funding strategies developed by the SciDAC program should be monitored and compared to the needs of the German scientific community.

15. NSF, Arlington, VA

The **National Science Foundation** (NSF) is an independent US federal agency. It funds research and education in most fields of science and engineering and accounts for about one-fourth of federal support to academic institutions for basic research. In some fields, such as mathematics, computer science, economics and the social sciences, the NSF is the major source of federal support. The agency operates no laboratories itself but does support

National Research Centers, user facilities, certain oceanographic vessels and Antarctic research stations. The Foundation also supports co-operative research between universities and industry, US participation in international scientific and engineering efforts, and educational activities at every academic level. Compared to the DFG, the NSF has an almost three times higher budget without funding the medical sciences, which in turn are supported by the U.S. National Institute of Health (NIH).

The NSF organizes its research and education support through seven directorates, each encompassing several disciplines. The Directorate for Computer and Information Science and Engineering (CISE) strives to uphold a position of world leadership in computing, communications, and information science and engineering. To achieve this, CISE supports investigator initiated research in all areas of computer and information science and engineering, helps develop and maintain cutting-edge national computing and information infrastructure for research and education generally, and contributes to the education and training of the next generation of computer scientists and engineers.

CISE is organized in three divisions: the Division of Computing & Communication Foundations (CCF); the Division of Computer and Network Systems (CNS); and the Division of Information and Intelligent Systems (IIS). Each division is organized into a small number of programs that are responsible for managing a portfolio of grants and proposal competitions within a broad area of research and education. While individual program directors may be designated as the point of contact for specific sub-disciplines, collaboration takes place within each program, across each division, and between divisions and directorates. ([Link](#))

The NSF also supports research through several offices within the Office of the Director. The **Office of Cyber-Infrastructure** (OCI) co-ordinates and supports cyber-infrastructure resources, tools and related services such as supercomputers, high-capacity mass-storage systems, system software suites and programming environments, scalable interactive visualization tools, productivity software libraries and tools, large-scale data repositories and digitized scientific data management systems, networks of various reach and granularity and an array of software tools and services that hide the complexities and heterogeneity of contemporary cyber-infrastructure while seeking to provide ubiquitous access and enhanced usability. OCI supports the preparation and training of current and future generations of researchers and educators to use cyber-infrastructure to further their research and education goals, while also supporting the scientific and engineering professionals who create and maintain these IT-based resources and systems and who provide essential customer services to the national science and engineering user community. ([Link](#))

Cyber-infrastructure (CI) is a growing field of interest for NSF and is seen as the enabling infrastructure for science and engineering. The Advisory Committee for Cyber-infrastructure (ACCI) established Task Forces, which reported to the NSF in December 2010. It was recommended that *permanent programmatic activities in Computational and Data-enabled Science & Engineering (CDS&E) should be established within NSF*. Furthermore, the NSF should establish processes to collect community requirements and plan long-term software roadmaps. Among others, these recommendations were embraced in the **Cyber Infrastructure Framework for 21st Century (CIF21) Science & Engineering**, which regards CI as an *ecosystem*. This ecosystem is meant to enable and nourish increasingly data- and compute-intensive, integrative, and multiscale as well as multi-disciplinary modern science. CIF21 thus aims at building a national infrastructure for CDS&E, leveraging common methods, approaches, and applications with a focus on interoperability. This will also catalyze and focus other CI investments across NSF. On the administration side, CIF21 is embedded into every directorate and office and managed as a NSF-wide coherent program.

Although DFG is different from NSF it seems worthwhile to assess whether similar measures should be undertaken within DFG. The recently established DFG Priority Program on Software for Exascale Computing points in such a direction. This priority program could well attract as well as co-ordinate the German HPC and CDS&E communities towards collaborative research. In addition, the DFG should think of installing the means to streamline scientific software development, i.e. to make it more efficient, sustainable, and accessible (to avoid re-inventing code over and over again), attractive (code development should be honored in suitable ways if it leads to scientific discoveries), and thereby career-safe (scientists should not be penalized for investing time in developing quality code for breakthrough discoveries).

NSF stresses the new role **data** will play. In this picture, science will increasingly be data-enabled and data-intensive, directing attention away from computing and more to data management (compute, storage, use, and access). Data infrastructures need to be fit to digest the huge amounts of data that the future experimental set-ups will generate. The data center consolidation activities currently pursued by the US Federal government point into this direction.

16. National Cancer Institute, Rockville, MD

In Rockville, the Center for Biomedical Informatics and Information technology (CBIIT) of the **National Cancer Institute** (NCI) was visited in Rockville, MD. The main topic of the visit was caBIG[®], which is intended to be a collaborative information network that accelerates the discovery of new approaches for the detection, diagnosis, treatment, and prevention of cancer. The caBIG[®] initiative was started by the NCI in 2004 and operates through an open development community. Its goals are to connect scientists and practitioners through a shareable and interoperable infrastructure, develop standard rules and a common language to share information more easily, build or adapt tools for collecting, analyzing, integrating, and disseminating information associated with cancer research and care. Since its start, caBIG[®] has been committed to the following principles: *Federated, Open development, Open access and Open source*.

A number of presentations and discussions provided a general overview of caBIG[®] as well as insights into its enterprise architecture, its semantic infrastructure, its imaging infrastructure, the issue of data integration and its clinical and translational infrastructure. Compared to the enormous investment made within caBIG[®] in the last seven years (its budget has grown annually, from approximately \$15 M in 2004 to more than \$47 M in 2010; an additional \$87-\$100 M from the American Reinvestment and Recovery Act (ARRA) in 2009/2010 bring the total cost of the caBIG[®] program to at least \$350 M for seven years) the real impact of all those tools developed and their concrete application in American cancer research institutions remained unclear to us. Even though the portfolio of caBIG[®] components seemed to be immense and the aspect of an open service-oriented application platform was mentioned several times, the impression remained that all those services would work together very well in one grid environment mainly composed by caBIG[®] components, but that the integration of caBIG[®] components as services within other research frameworks or environments has not been accomplished, yet. It was mentioned that one of the major problems of the caBIG[®] initiative was, that it was not so much based on real requirements and input from cancer researchers, but seemed more as a prominent top down approach, too often not matching the real requirements of the research community (*"We have built a large set of tools for cancer research, but we lost the clinical users and researchers on our way. ... We've built a Ferrari for them, but when we delivered it, they told us that they needed a Landrover."*).

Even though the main caBIG[®] idea to centrally develop and provide a large cancer research toolbox as an open-source system to be used nationwide as standardized IT-infrastructure components in clinical and translational research environments is exemplary and deserves to be praised, its realization seems to be less successful. Our impressions of caBIG[®] match

those of the NIH Board of Scientific Advisors Ad Hoc Working Group on the NCI published in its management summary report ([Link](#)).

17. Harvard Medical School, Boston, MA

The Center for Biomedical Informatics (CBMI) at **Harvard Medical School** (HMS) was visited in Boston, MA. CBMI is a research center within the Harvard Medical School that promotes and facilitates collaborative activities in biomedical informatics among researchers at Harvard Medical School and its affiliated institutions. In an era of biomedical knowledge overload, high-throughput biomedical data generation, increasing consumer access to biomedical information, and demands for real-time, information-based public health, the CBMI was established in 2005 to lead research and educational activities in those areas. It is co-directed by two HMS faculty members. The background and motivation to establish a “center of excellence” for biomedical computing and translational research at HMS were to integrate the various research activities that were already prominently pursued at HMS in the last decades, but in separated organizational units.

A number of presentations have been given by the CBMI staff, the Partners Healthcare (which is an integrated healthcare system founded in 1994 by Brigham and Women's Hospital and Massachusetts General Hospital), the Beth Israel Deaconess Medical Center, and the Children's Hospital Boston. These presentations focused on

- a) **the i2b2 and shrine projects** working on the development of a (federated) data integration platform and scalable informatics framework that will enable clinical researchers to use existing clinical data for discovery research and combine it with genomic data in order to facilitate the design of targeted therapies for individual patients with diseases having genetic origins.
- b) **the indivo project**, which has a long history at Children's Hospital Boston and is a free and open-source personally controlled health record (PCHR) system, enabling individuals to own and manage a complete, secure, digital copy of their health and wellness information. Indivo has interfaces to electronic medical records within hospitals, thus being able to integrate health information across sites of care and over time. It is actively deployed in diverse settings and comprises an easy to use Application Programming Interface (API), making the Indivo platform extensible in many diverse use cases.

- c) **the healthmap project**, which brings together disparate data sources to achieve a unified and comprehensive view of the current global state of infectious diseases and their effect on human and animal health. This freely available website integrates outbreak data of varying reliability, ranging from news sources (such as Google News) to curated personal accounts (such as ProMED) to validated official alerts (such as World Health Organization). It provides a take-off point for real-time information on emerging infectious diseases and has particular interest for public health officials and international travelers.
- d) **the SMART Initiative** which aims at the development of Substitutable Medical Apps, based on the Reusable Technology (SMART) initiative, which is developing an open source API to access any electronic health record (EHR) system. Examples of such applications are currently medication management tools, health risk detectors, and e-prescribing applications.
- e) the **work scope and agenda of the Clinical Informatics Research and Development group** at Partners Healthcare System and
- f) the **US National Healthcare IT Program**

Even though all projects focused on - at first appearance – completely different topics, they all followed two major strategic lines:

1. Using **open-source software** as much as possible and providing own developments as open-source software to the community as well. This led to a widespread usage of the developments not only at HMS but also in many other organizations and research institutions throughout the U.S. (and for i2b2 even worldwide) and fostered input, extensions and enhancements to those developments from a large user community. The latter was even more supported by **clearly documented APIs** of the core development components that supported **extensions** for many different purposes and scenarios. Core components and underlying concepts were typically implemented as **reusable software modules**, so that different system developments could build on those core components.
2. **Trying to reuse data that have already been documented during routine clinical care** and stored within electronic medical records or personal electronic health records **for other purposes**, such as clinical/translational research and clinical decision support.

The presentations provided an excellent comparison of the American healthcare system with several European healthcare systems and insight into the motivation and goals of two immense funding initiatives initiated by the Obama legislation (HITECH and ARRA) to support the introduction and meaningful use of health information technologies within American hospitals and general practitioner practices. With presentations and examples from his “own development of an electronic medical record system at Partners Healthcare” the aims of those funding initiatives have been illustrated.

It was impressive to see how many resources and how much continuous support were dedicated at HMS and Partners Healthcare to research, structures and infrastructure components aiming at clinical decision support in order to increase the quality of patient care as well as clinical and translational research within the Harvard medical faculty. The presentations were given by different people belonging to different medical informatics / biomedical computing organizations at HMS and illustrated a high level of interactivity between those groups and interdisciplinary collaborative research approaches combining competencies including clinicians, biologists, geneticists, public health scientists, and computer scientists.

From some of the projects presented (i2b2, shrine), it was clear that collaborations between German medical informatics groups (Göttingen, Erlangen) and CBMI have already been established. Following these examples to further use of open-source systems developed at HMS within German research projects and as core IT infrastructure components for translational research, as well as establishing even broader collaborations can be recommended.

V. Conclusions

Given the large number and wide range of institutions visited during the study tour, there is also a broad spectrum of impressions, and it is probably impossible to reduce this spectrum to a few conclusions. Nevertheless, a few important issues that were observed, addressed, or discussed on several visits shall be mentioned in the following.

There is a very broad **diversity of enterprise strategies** and **culture** between the IT companies visited. While Intel or IBM constitute examples of long-term-oriented IT technology industry, Apple successfully bridges the gap between IT technology and lifestyle. Despite all differences, both Intel and NVIDIA appear as classical companies, whereas Google, for instance, goes different ways.

In **processor technology**, a couple of trends became obvious: an ongoing focus on on-chip parallelism, even though the number of cores might not increase that far as predicted recently; despite of the overall “convergence” strategies, hybrid architectures will be predominant for the next generations; there are increasing investments in 3D processor technology.

The success of **accelerators** in general, and **GPU** in particular, has led to the fact that the underlying hardware paradigms and the resulting programming models are on the agenda of all chip manufacturers, with a hybrid strategy pointing towards convergence (Intel with MIC, AMD with ATM, and NVIDIA with the ARM processor). Nevertheless, there are different views on the future of supercomputers (cf. NVIDIA’s statement “*handhelds are the supercomputers of the future*”). Programmability, i.e. the early availability of a complete software stack, has been addressed with different intensity thus far, but is now considered as crucial everywhere.

The development of **high-end HPC installations** (driven and funded by DoE at ANL or by NSF at NCSA) takes place as a co-development of commercial providers and public research institutions. This involves huge investments in the systems and system software, but increasingly also in usability in the widest sense, namely in particular applications and application software. While the DoE envisions that task on a long-term agenda, the NSF, in the future, will rather focus on the application, methodology, and algorithmics issues only.

In addition to driving the development of petascale and exascale systems (investments in “racks”), it is generally accepted that **algorithms and software** well-suited for such systems do not come automatically as a side-product, but need large efforts and funding, too

(investments in “brains”). DoE’s SciDAC centers devoted to specific aspects on this behalf (algorithms, data, and others) can serve as an example. Both DoE and NSF have respective programs specially addressing this issue.

The **data issue** is more and more considered as a crucial topic for Computational Science and Engineering (CSE) and HPC. This is reflected in the new notion of CDSE (Computational and Data-intensive Science and Engineering), which expands the established CSE. Aspects of such a data-centered perspective are data management, storage, data transport, and data exploration.

IT in medicine has seen and sees numerous large research consortia. There are both difficulties and success stories. One central aspect is the IT support of a patient-centered clinical and translational research and healthcare. The sensitivity for data security issues is increasing especially with respect to collecting medical and socio-economic data from social networks.

Medical IT infrastructures include many people with informatics background, but are developed rather in medical environments with little exchange with dedicated CSE or HPC institutions. Once again, it became obvious that placing long-term tasks such as data repositories on the basis of short-term and project-based funds is very problematic with respect to sustainability. However, thorough reviews of long-term project progresses are nevertheless mandatory. **IT frameworks** and projects in a medical environment always reflect the specific socio-economic boundary conditions of the respective national laws and medical organizations. Medical IT solutions that are tailored to the US may thus be difficult to map onto a German context.

“**Green IT**” draws more and more attention as power consumption increasingly constraints IT developments, especially but not exclusively within HPC.

The large **funding agencies** in the US, namely DoE, NSF, and NIH, all have strategic programs and significant funding for both CSE and HPC. In this respect, the situation in Germany is less developed, and the discrepancy between investments in “racks” and “brains” is problematic. As a consequence for Germany, BMBF’s HPC software program should be continued, while DFG should focus on and offer opportunities for fundamental research, the recently launched priority program SPPEXA being an excellent starting point in that direction. DFG should launch a fundamental-research-oriented initiative, while the CSE and HPC communities should contemplate on the organizational structures to represent their interests. Furthermore, a respective DFG review board for CSE and HPC would certainly be helpful.