# *Infrastructure for Predictive Materials Science Simulations*

François Gygi

*Department of Applied Science*
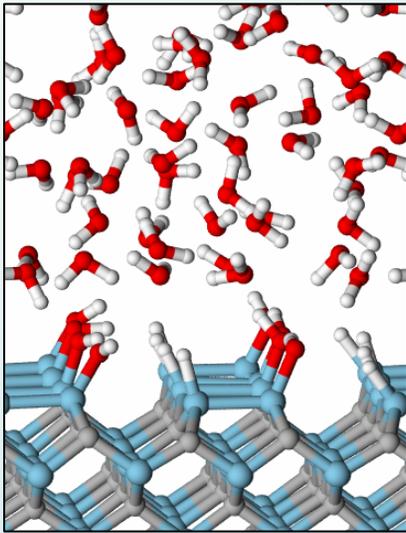*University of California, Davis*
*Davis, CA 95616*
fgygi@ucdavis.edu

NSF Cyberinfrastructure Workshop, Aug 4, 2006

# *First-Principles Molecular Dynamics*

- Molecular dynamics: an atomic-scale simulation method
  - Compute the trajectories of all atoms
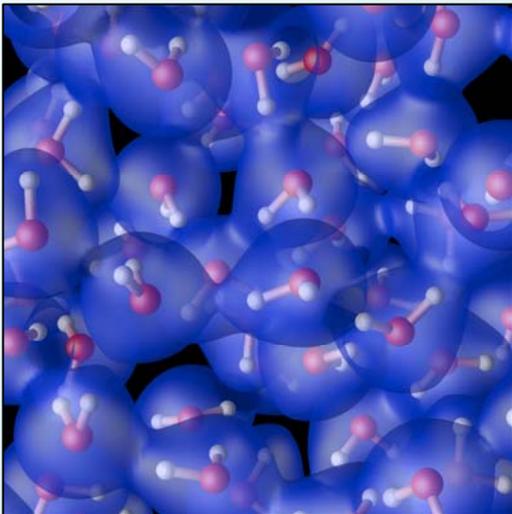  - extract statistical information from the trajectories

Water on a silicon carbide surface

Atoms move according to Newton's law:

$$m_i \ddot{\mathbf{R}}_i = \mathbf{F}_i$$

# *First-Principles Molecular Dynamics*

- Why "First-Principles"?
  - Avoid empirical models of interatomic forces
  - Use fundamental principles instead: Quantum Mechanics
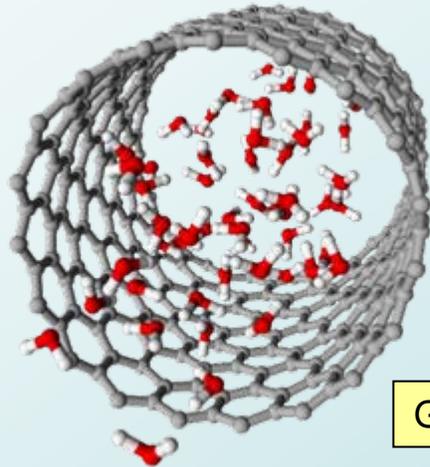  - Must describe ions and electrons consistently and simultaneously



<u>At each time step:</u>

1) Compute the electronic structure

2) Derive interatomic forces

3) Move atoms
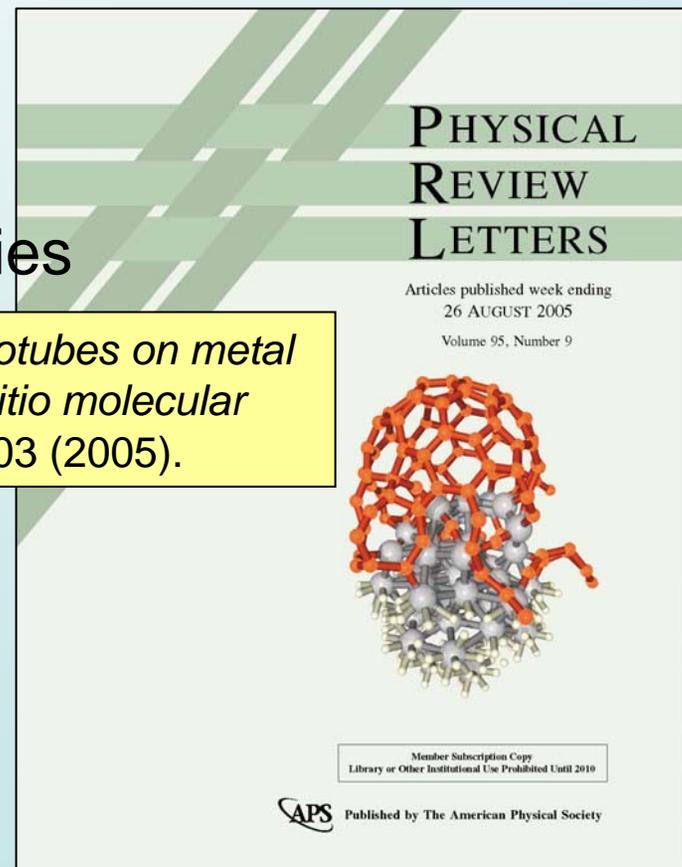
# *First-Principles Molecular Dynamics Applications*

- ## Nanoscience
  - FPMD captures the complex interplay between chemical, electronic and optical properties

J.Y. Raty, F. Gygi, and G. Galli, *Growth of carbon nanotubes on metal nanoparticles: A microscopic mechanism from ab initio molecular dynamics simulations*, Phys. Rev. Lett. **95,** 096103 (2005).



PHYSICAL REVIEW LETTERS

Articles published week ending
26 AUGUST 2005

Volume 95, Number 9

Member Subscription Copy
Library or Other Institutional Use Prohibited Until 2010

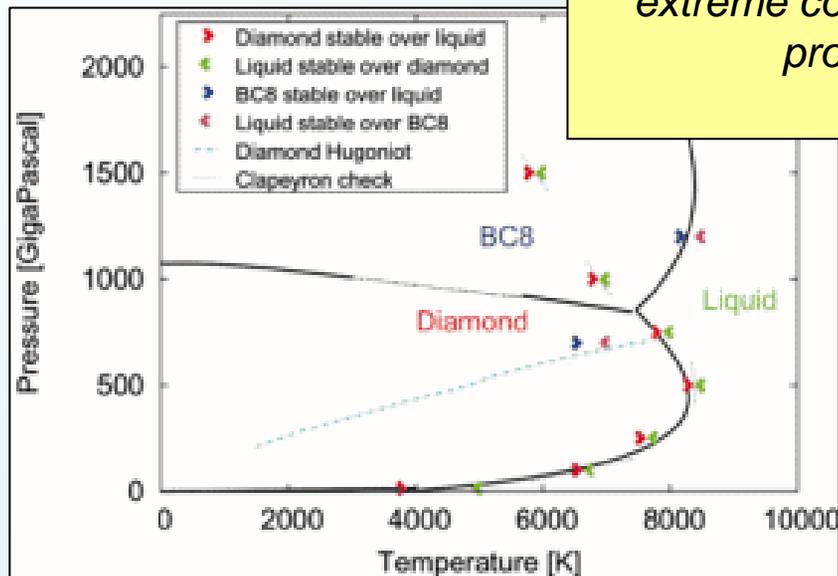APS   Published by The American Physical Society

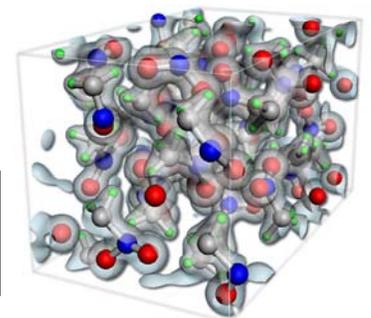G. Cicero, E. Schwegler, J. Grossman, F.Gygi, G. Galli (2006)

# *First-Principles Molecular Dynamics Applications*
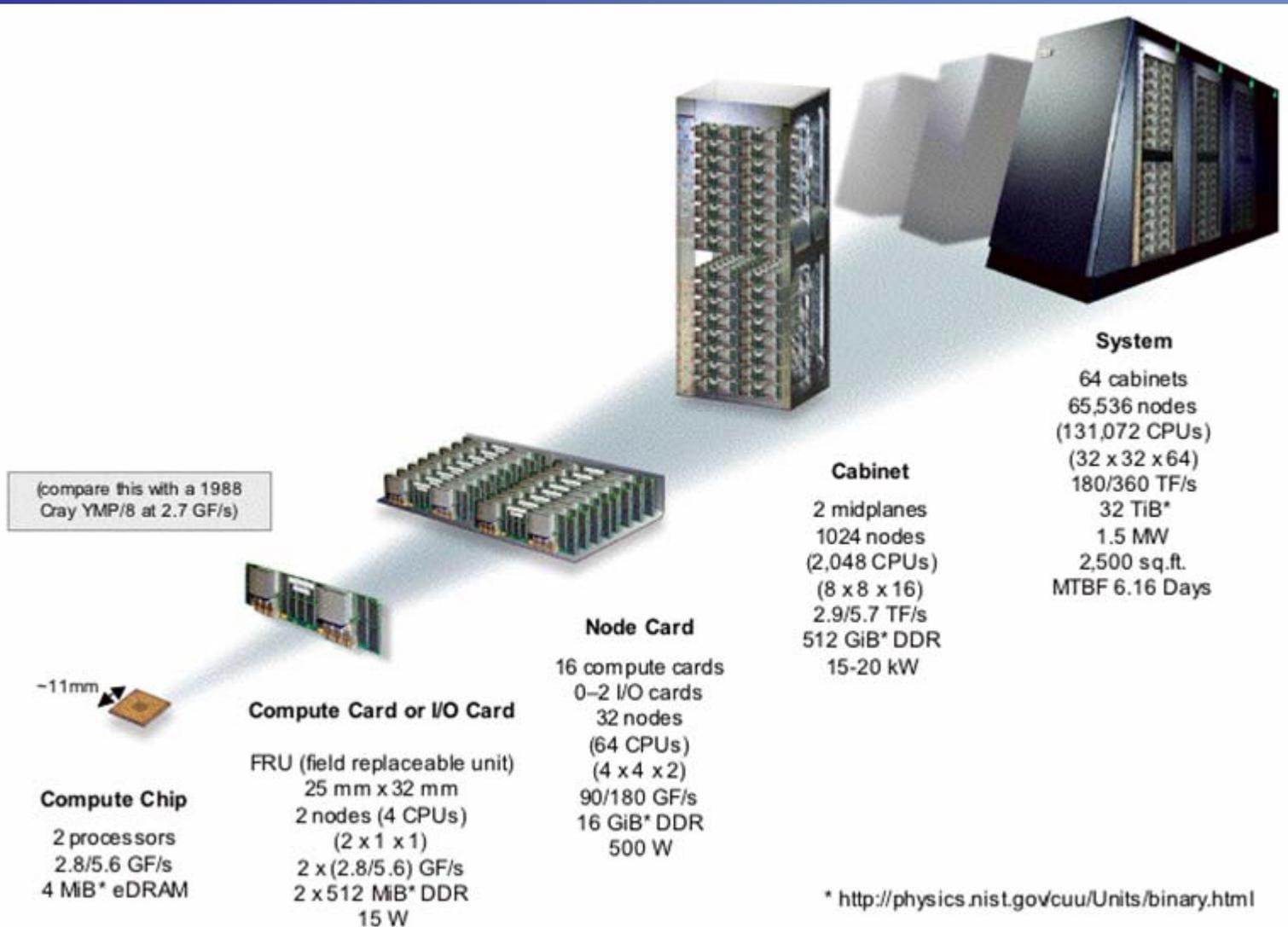
- ## Materials in extreme conditions

A. A. Correa, S. A. Bonev and G. Galli, *Carbon under extreme conditions: Phase boundaries and electronic properties from first-principles theory,* PNAS **103**, 1204-1208 (2006).

F. Gygi and G. Galli, *Ab Initio Simulation in Extreme Conditions*, Materials Today **8**, 26-32 (2005).
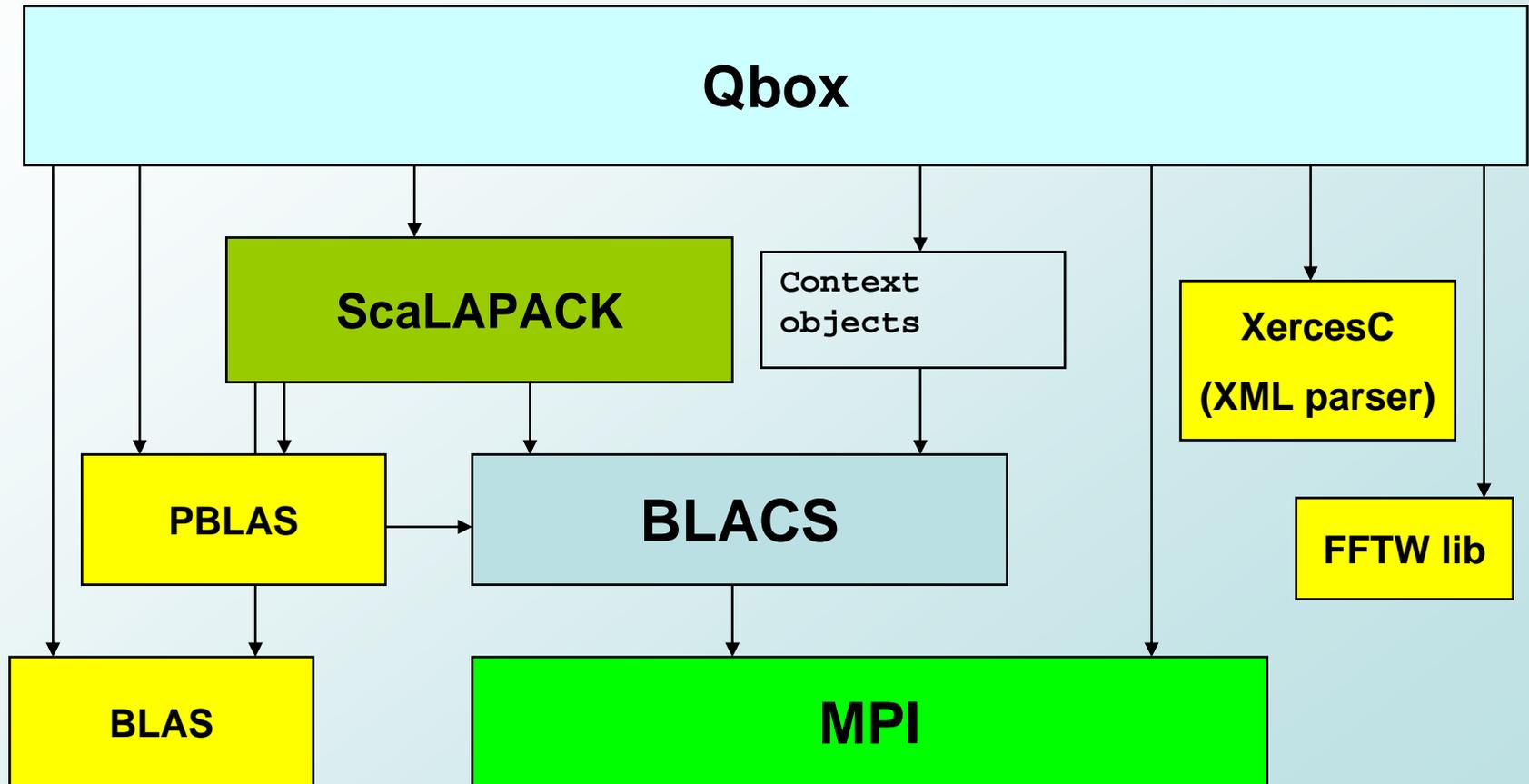
# *Large parallel platforms: BlueGene/L*



(compare this with a 1988 Cray YMP/8 at 2.7 GF/s)

~11mm

**Compute Chip**

2 processors
2.8/5.6 GF/s
4 MiB* eDRAM

**Compute Card or I/O Card**

FRU (field replaceable unit)
25 mm x 32 mm
2 nodes (4 CPUs)
(2 x 1 x 1)
2 x (2.8/5.6) GF/s
2 x 512 MiB* DDR
15 W

**Node Card**

16 compute cards
0–2 I/O cards
32 nodes
(64 CPUs)
(4 x 4 x 2)
90/180 GF/s
16 GiB* DDR
500 W

**Cabinet**

2 midplanes
1024 nodes
(2,048 CPUs)
(8 x 8 x 16)
2.9/5.7 TF/s
512 GiB* DDR
15-20 kW

**System**

64 cabinets
65,536 nodes
(131,072 CPUs)
(32 x 32 x 64)
180/360 TF/s
32 TiB*
1.5 MW
2,500 sq.ft.
MTBF 6.16 Days

* http://physics.nist.gov/cuu/Units/binary.html

# *Qbox: a Large-Scale Parallel Implementation of FPMD*

- Qbox is a C++ implementation of FPMD
- Designed for scalability up to $>10^4$ CPUs
- Development started in 2003
- Built on several optimized libraries developed elsewhere:
  - PBLAS
  - ScaLAPACK
  - FFTW
  - MPI
  - Apache Xerces-C

# Qbox code structure

The Leading Source for Global News, Information and Events on High Performance Computing / June 22, 2006

**Home Page**

**Breaking News:**

# NNSA Announces Record Performance on IBM Blue Gene/L

The Department of Energy's National Nuclear Security Administration (NNSA) and IBM teamed up to announce that a new mark was achieved on the world's fastest supercomputer named Blue Gene/L (BG/L). This world record for a scientific application was set by achieving a sustained performance of 207.3 teraflops on the "Qbox" computer code for conducting materials science simulations critical to national security.

BG/L is an IBM supercomputer housed at NNSA's Lawrence Livermore National Laboratory, and is ranked as the world's fastest supercomputer. It is used to conduct materials science simulations for NNSA's Advanced Simulation and Computing (ASC) program, which unites the scientific computing know-how of NNSA's Los Alamos, Sandia and Lawrence Livermore national laboratories. The computer simulation capabilities developed by the ASC program provide the nuclear weapons analysis that NNSA needs to keep the nuclear weapons stockpile safe, secure and reliable without underground nuclear testing.

"This is an important step on t
ensuring the safety and reliabi
can scale with real-world appl
collaborators, who helped to o
NNSA's ASC Program.
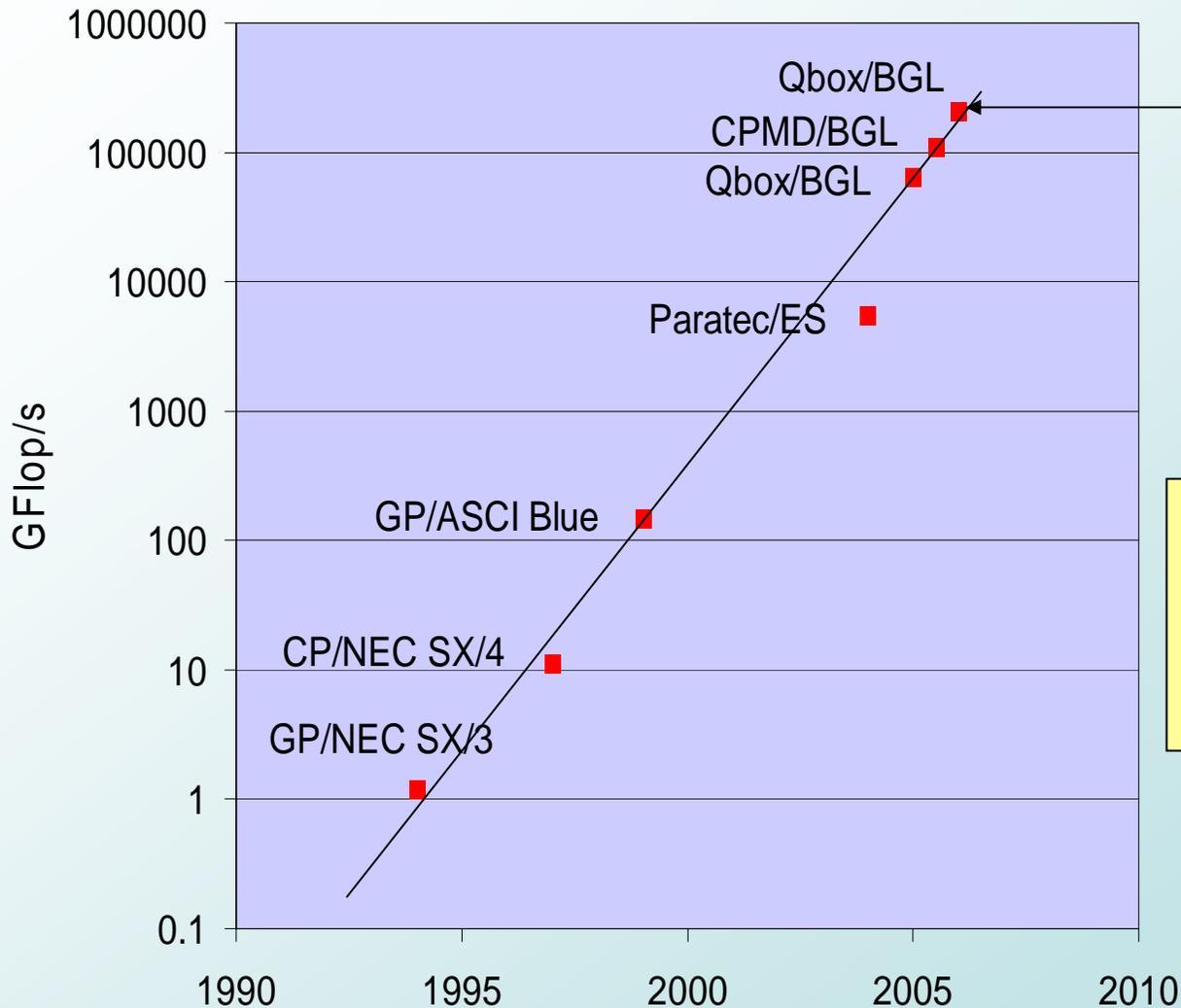
The performance improvemen
software researchers at IBM th

"Results of this magnitude an

Nuclear Security Administration and Lawrence Livermore National Laboratory can the bound far as they've been today. We will continue to work together, pushing the boundaries of insigh shared mission in ways never before possible."

Qbox is a first-principles molecular dynamics (FPMD) code, designed to predict the propertie conditions of temperature and pressure -- a longstanding goal for researchers in materials scie physics. FPMD codes are used for complex simulations at the atomic level in a number of sci solid-state physics, chemistry, biology and nanotechnology.

The "Q" in Qbox is for "quantum," a reference to the quantum mechanical descriptions of ele this type of simulation code. The ability to accurately model changes to the electronic structur codes from classical molecular dynamics codes.
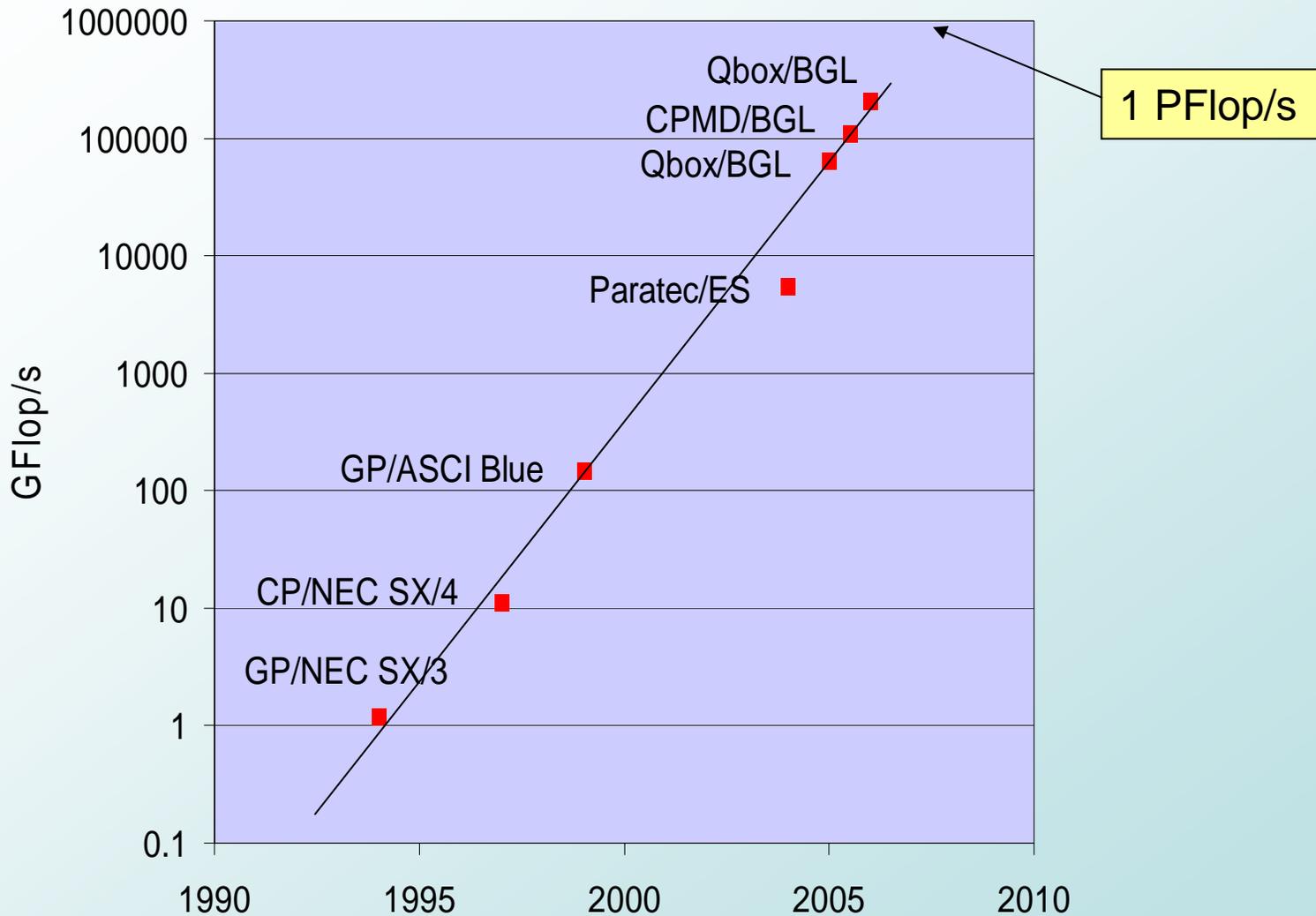
# History of FPMD code performance

# *Future Challenges and CI*

- Moving FPMD to petascale platforms
- Developing better algorithms
- Data management issues / Standards

# *The road to petascale simulation*

# *The road to petascale simulation*

- "At 200 TF we are 1/5 of the way there", 🙂 but…
  - The architecture of petascale computers is still changing
  - Scaling of the current FPMD design to petascale is unlikely
- A complete rewrite of FPMD codes will likely be necessary

# *The road to petascale simulation*

- Petascale code optimization is a new "science"
  - Consider debugging your code on 32,768 CPUs
  - Need new tools for SW development
    - e.g. visualization of message traffic
  - New optimization issues
    - Multicore chips
    - Heterogeneous architectures
    - Data flow becomes more critical than Flops
- Optimization is essential
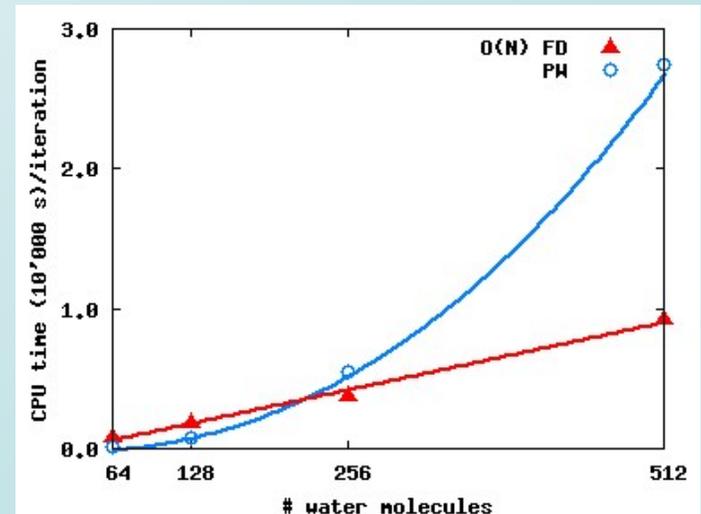  - 1% of a petascale platform is a terrible thing to waste

# *The road to petascale simulation*

- Petascale software will be more complex
  - complexity is driven by machine architecture, number of CPUs
- Petascale code development must start early
  - Qbox code development took 3 years, started well before BlueGene/L hardware was built
  - Hardware lifetime is becoming comparable to code development time (~3-5 yrs)

# *Developing better algorithms*

- Current FPMD simulation cost is O($N^3$)
  - ($N$ = # of electrons)
- This is currently the main obstacle to applying FPMD to more complex problems
- Need research in linear-scaling algorithms (O($N$))
  - Need for *controlled accuracy*



J.-L. Fattebert and F. Gygi, "*Linear-scaling first-principles molecular dynamics with plane-waves accuracy*", Phys. Rev. B **73**, 115124 (2006).
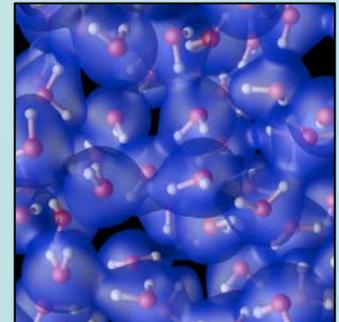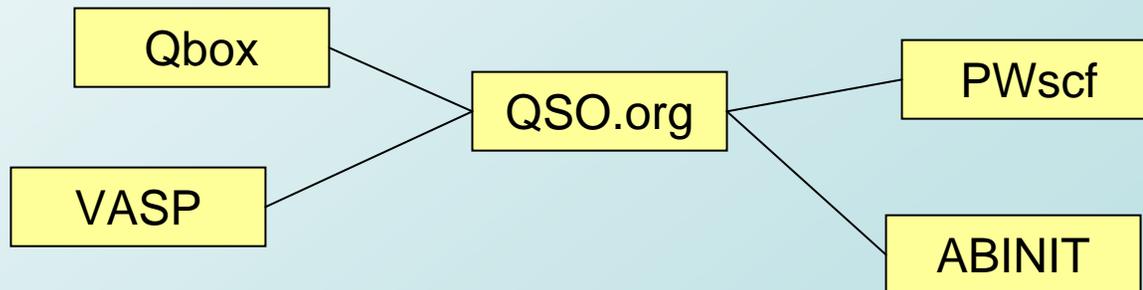
# *Data Management*

- FPMD simulations generate large datasets (many TB)
- Need for large, www-accessible repositories
  - Accessibility is critical to ensure validation of data by different groups (need high bandwidth networks)
  - Reproducibility of simulation data is an important goal
    - a "Quality of Science" issue
  - Need to develop new ways to limit the data flow
    - e.g. application-specific compression algorithms

# *Data Management*

- Need for data standards
  - Today: 4 community codes = 4 different data formats
  - New standards should be XML-based
  - Developing XML standards (schemas definitions, etc.) is a costly, iterative process
  - Use automatic data validation
  - New standards should be code-neutral
  - Our effort: http://www.quantum-simulation.org

# *Data Management*

- Need to encourage the development of XML-aware codes
  - Example: all Qbox I/O is valid, well-formed XML
  - Integrating XML parsers in application codes incurs a maintenance cost
  - May require new CS developments (e.g. parallel XML parsers)
- Need for web-aware codes
  - Example: definition of a carbon pseudopotential in Qbox:

```
[qbox] species carbon http://www.quantum-
simulation.org/examples/species/carbon_pbe.xml
```

# *Is the code doing the right thing? Code verification*

- Encourage the development of well documented test suites
  - time-consuming process
  - need dedicated computing platforms *just for testing!*
- Need for *a few* competing community codes
  - Cross-verification is an invaluable tool
  - We must resist the urge to completely eliminate duplication of effort

# *Open source licenses*

- Variants of open-source licences abound
  - GNU-GPL, BSD, MIT, etc..
- Should NSF help software developers make the "right" choice?
  - NSF endorses one license, or
  - NSF creates the NSF license

# *Summary*

*The successful solution of a realistic problem in applied mathematics requires the fusion of four distinct ingredients:*

*1) knowledge of the subject area of the problem*

*2) knowledge of the relevant mathematics*

*3) knowledge of the relevant computer science*

*4) a talent for selecting just what part of all this knowledge will actually solve the problem, and ignoring the rest.*

George E. Forsythe