

# Performance Evaluation for Petascale Quantum Simulation Tools

---

*Stan Tomov*

Innovative Computing Laboratory ( ICL )  
The University of Tennessee

joint work with *Wenchang Lu*<sup>1,2</sup>, *Jerzy Bernholc*<sup>1,2</sup>, *Shirley Moore*<sup>3</sup>, and *Jack Dongarra*<sup>2,3</sup>  
( NCSU<sup>1</sup>, ORNL<sup>2</sup>, UTK<sup>3</sup> )

**CUG 09: Compute the Future, Atlanta GA**  
May 4<sup>th</sup> – May 7<sup>th</sup>, 2009

# Outline

- Background
  - Simulation of nano materials and devices
  - Challenges of future architectures
- Electronic structure calculations
- Performance evaluation
- Performance analysis
- Bottlenecks and ideas for their removal
- Conclusions

# Electronic properties of nano-structures

- Semiconductor Quantum dots (QDs)
  - **Tiny** crystals ranging from a few hundred to few thousand atoms in size; made by humans

At these small sizes electronic properties **critically depend** on **shape** and **size**

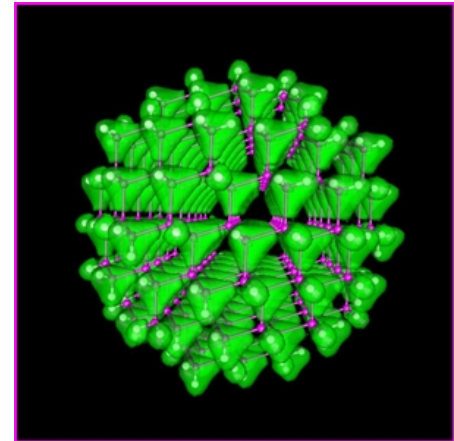
⇒ electronic properties can be tuned  
⇒ enables remarkable applications

The dependence is **quantum mechanical** in nature and can be modelled

- can not be done on macroscopic scales
- has to be at **atomic and subatomic level (nanoscale)**

- Quantum wires (QWs) and devices

- their conducting properties are affected by build-in nano-materials

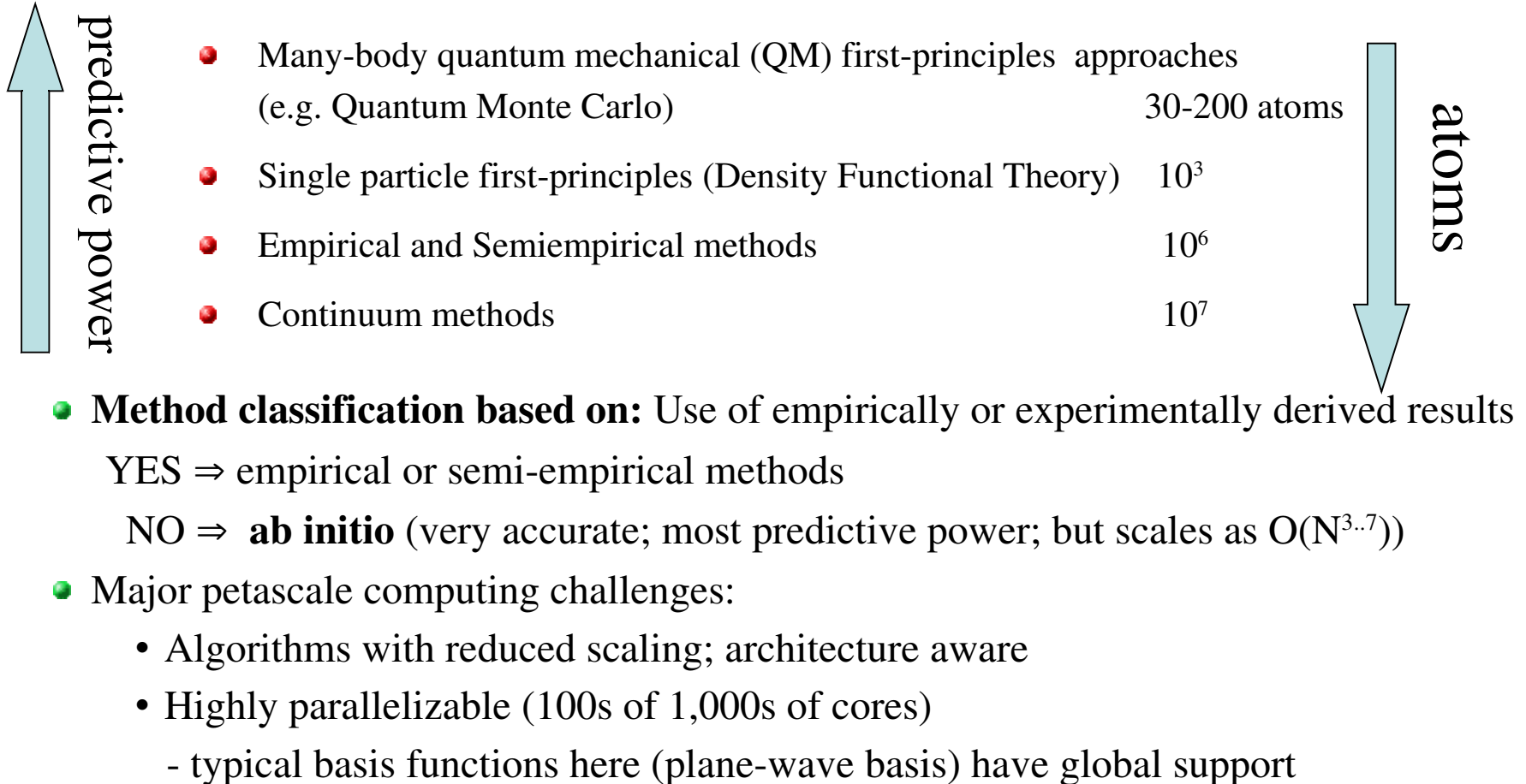


Total electron charge density of a quantum dot of gallium arsenide, containing just 465 atoms.



Quantum dots of the same material but different sizes have different band gaps and emit different colors

# Nano Materials Simulations



# Challenges of Future Architectures

- Increase in parallelism
  - Multicores, GPUs, hybrid architectures, etc
- Increase in communication cost (*vs* computation)
  - Gap between processor and memory speed continue to grow (exponentially)  
[e.g. processor speed improves 59%, memory bandwidth 23%, latency 5.5%]

# Approach

- **Basis selection: Plane-waves, grid functions, or Gaussian orbitals, etc.**
- Plane-waves: 
$$\psi_{nk}(r) = \sum_{g, |g| < E_{cut}} C_g^n(k) e^{i(g+k) \cdot r}$$
  - Good approximation properties
  - Can be preconditioned easily (and efficiently) as the kinetic energy (the laplacian) is diagonal in Fourier space, the potential is diagonal in real space
  - Usually codes are in Fourier space and go back and forth to real with FFTs
  - Concern may be scalability of FFT on 100s of 1,000s of processors as it requires global communication
- Grid functions: e.g. finite elements, grids, or wavelets
  - Domain decomposition techniques can guarantee scalability for large enough problems
  - Interesting as they enable algebraically based preconditioners as well
  - Including multigrid/multiscale
    - e.g. **real-space multigrid methods (RMG) by J. Bernholc et al (NCSU)**

# Goal of this work

- **Performance evaluation** of petascale quantum simulation tools for nanotechnology applications
  - Based on existing real-space multigrid method (RMG)
  - In-depth understanding of their performance on Teraflop leadership platforms
  - With the help of tools such as TAU, PAPI, Jumpshot, KOJAK, etc
- **Identify performance bottlenecks** and ways/ideas for their removal
- **Aid the development of algorithms**, and in particular petascale quantum simulations tools, that effectively use the underlying hardware

# Software/Hardware Environment

- We consider 2 methodologies [implemented so far in our codes]
  - Global grid method
    - Wave functions are represented in the real space uniform grids
    - Most time consuming is orthogonalization and subspace diagonalization
    - Massively parallel, good flops performance, but scales in  $O(N^3)$  with system size
  - Optimally localized orbital method
    - Scales in nearly  $O(N)$  but has computational challenges
- Hardware: we consider **Jaguar**, a Cray XT4 system at ORNL
  - Based on quad-core 2.1 GHz AMD Opteron processors

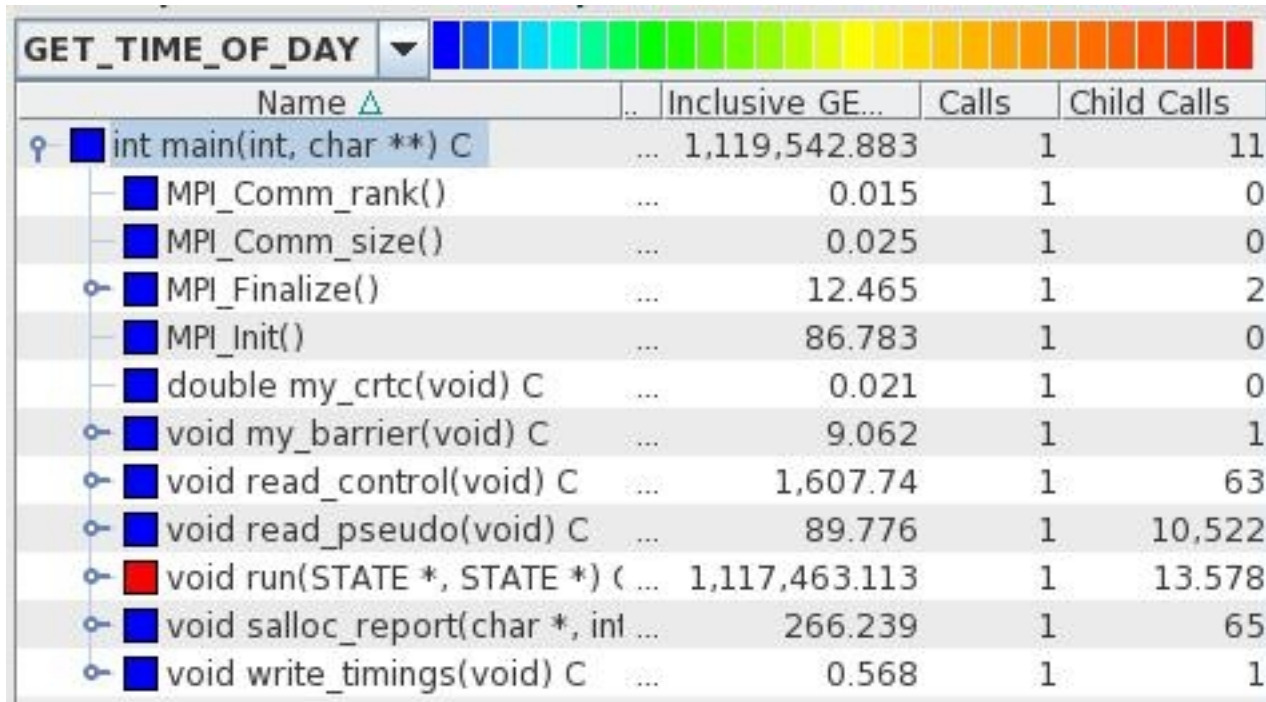


# Performance evaluation

- Techniques that we found most useful
- Profiling [using TAU with PAPI]
  - To get familiar with code structure
  - To get performance profiles
  - To identify possible performance bottlenecks
- Tracing [using TAU]
  - To determine exact locations and cause of bottlenecks

# Profiling

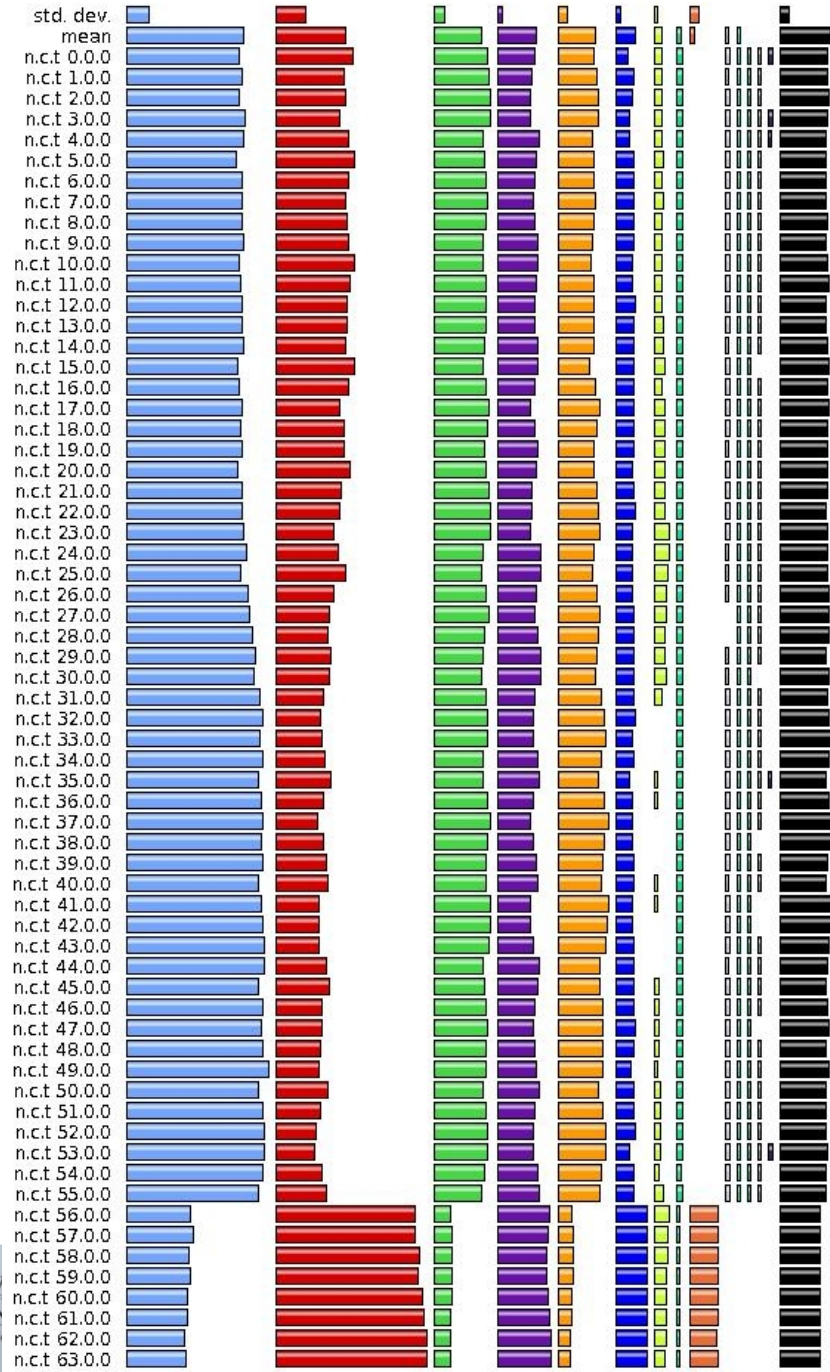
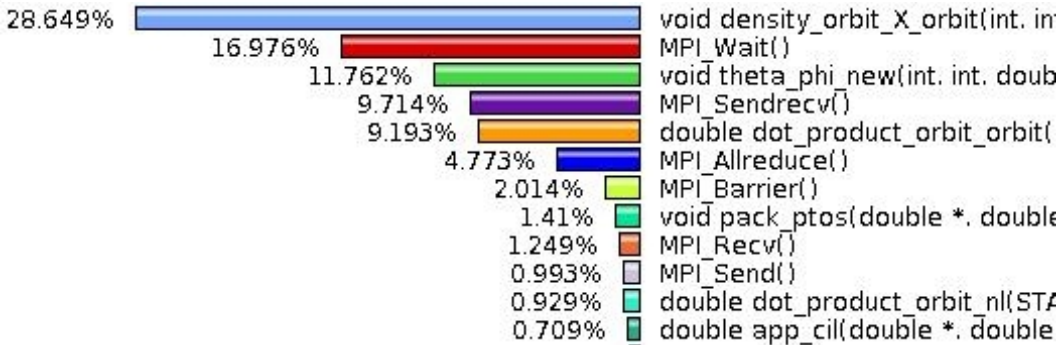
- Getting familiar with code structure [by generating callpath data]



Name $\Delta$	Inclusive GE...	Calls	Child Calls
int main(int, char **) C	1,119,542.883	1	11
MPI_Comm_rank()	0.015	1	0
MPI_Comm_size()	0.025	1	0
MPI_Finalize()	12.465	1	2
MPI_Init()	86.783	1	0
double my_crtc(void) C	0.021	1	0
void my_barrier(void) C	9.062	1	1
void read_control(void) C	1,607.74	1	63
void read_pseudo(void) C	89.776	1	10,522
void run(STATE *, STATE *) C	1,117,463.113	1	13,578
void salloc_report(char *, int ...)	266.239	1	65
void write_timings(void) C	0.568	1	1

# Profiling

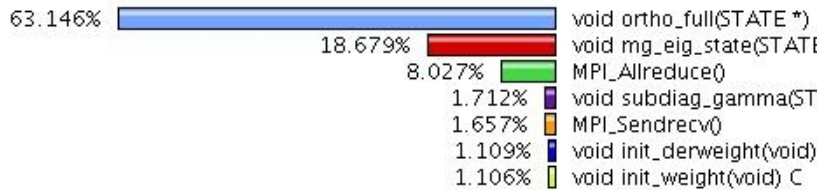
- Performance profiles  
[load balance, what to optimize, etc]



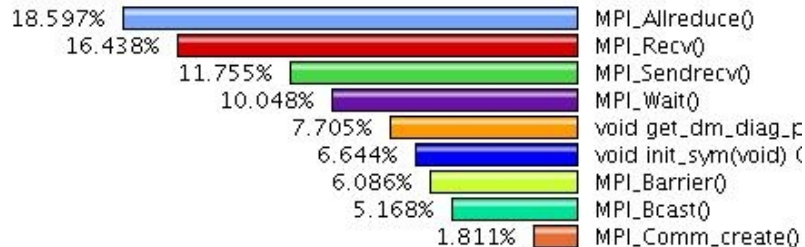
# Profiling

- Performance evaluation results  
 [PAPI counters; example with PAPI\_FP\_INS; right]  
 [Profiles for the 2 codes on large problems;  
 1024 cores; below ]  
 [ 1<sup>st</sup> code about 6 x faster; 2<sup>nd</sup> more sparse op.]

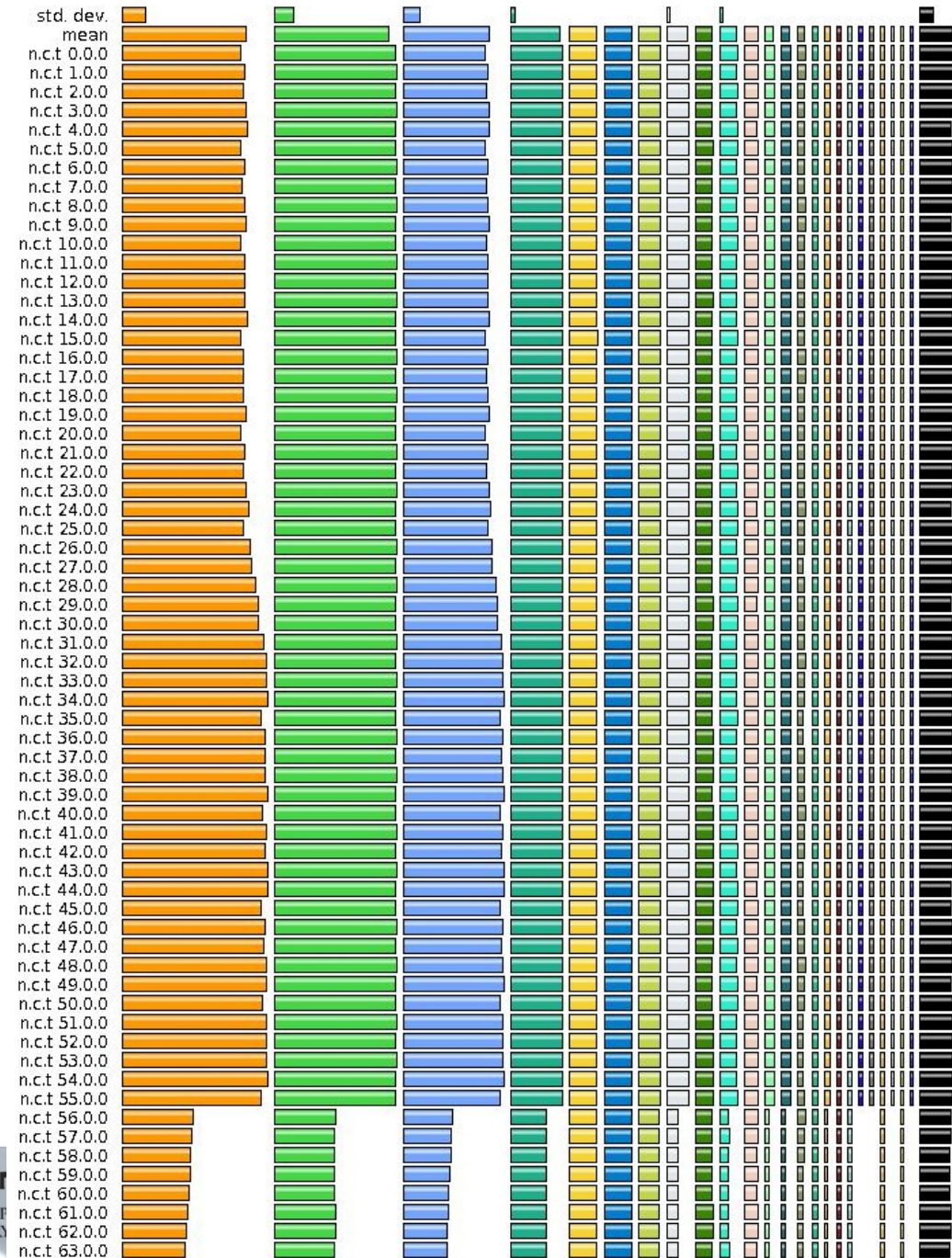
Metric: GET\_TIME\_OF\_DAY  
 Value: Exclusive percent



Metric: LINUX\_TIMERS  
 Value: Exclusive percent



Metric: PAPI\_FP\_INS  
 Value: Exclusive



# Performance analysis

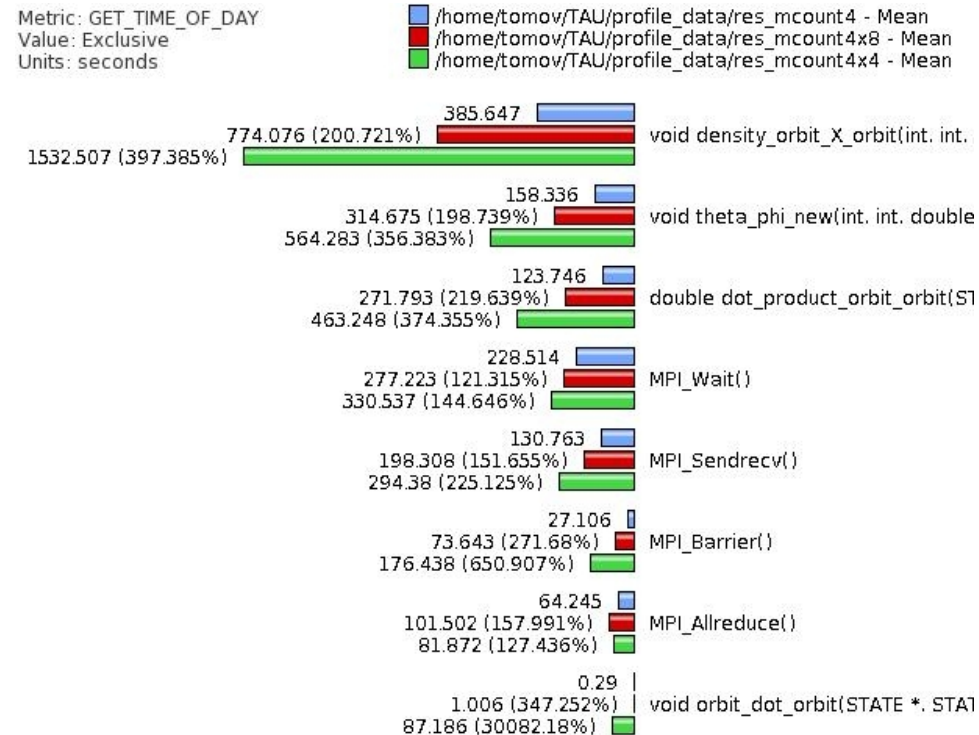
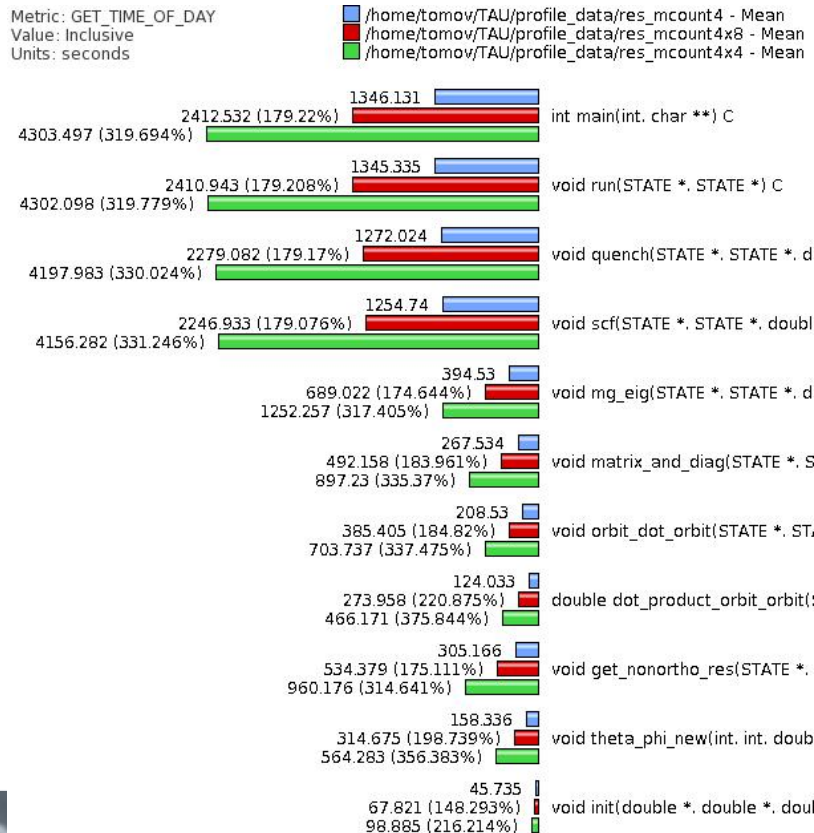
- Tracing

- To determine exact locations and causes of bottlenecks
- TAU to generate trace files and analyze them with Jumpshot and tools like KOJAK
- Codes well written
  - Blocked communications, asynchronous, intermixed with computation
- Domain decomposition guarantees weak scalability
  - We have to concentrate on efficient use of multicores within a node
- We found early posting of MPI\_Irecv will benefit our codes
- We found useful to compare traces of different runs – to study effects of code changes
- Generate profile-type statistics for various parts of the codes

# Performance analysis

- Scalability

- Studied both strong and weak [example on strong scalability]



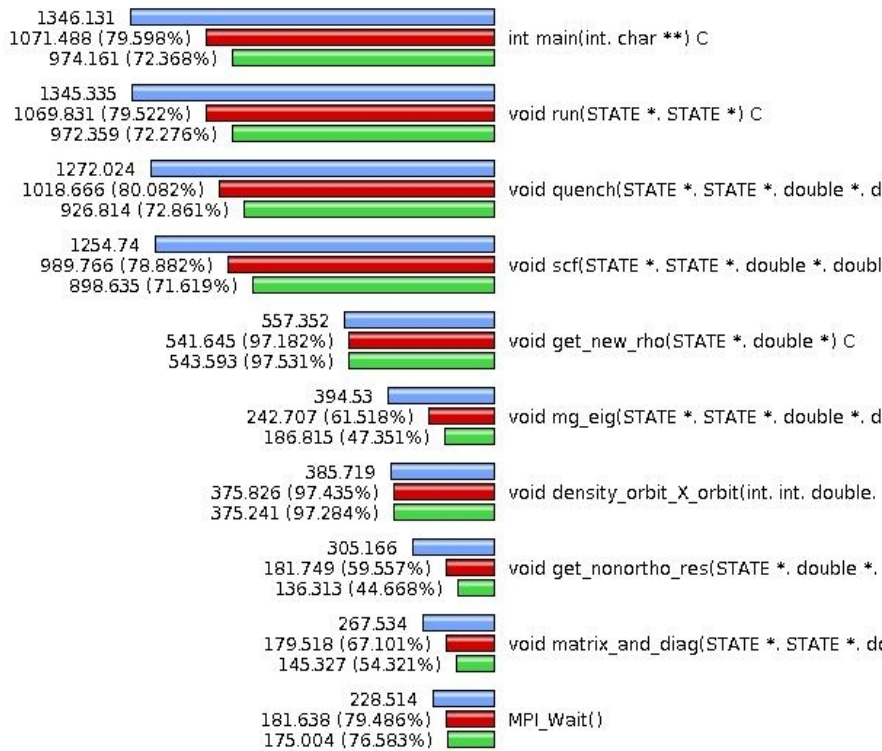
# Performance analysis

- Multicore use

- Measurements in different hardware configurations: runs using 4, 2, and single core of the quad-core nodes [number of cores to be the same]

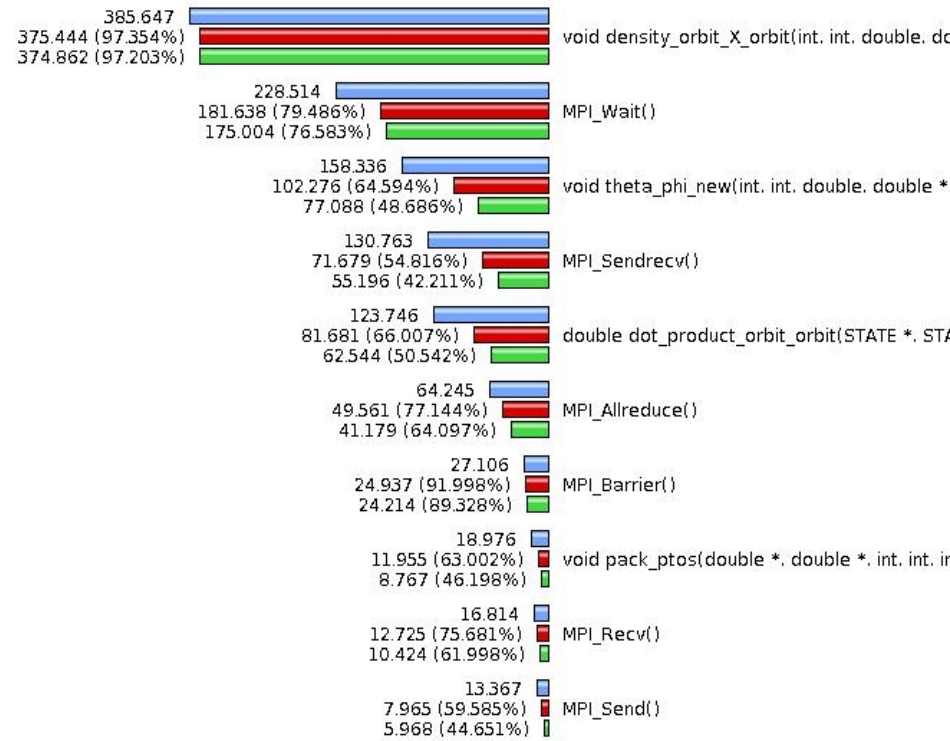
Metric: GET\_TIME\_OF\_DAY  
Value: Inclusive  
Units: seconds

■ /home/tomov/TAU/profile\_data/res\_mcount4 - Mean  
■ /home/tomov/TAU/profile\_data/res\_mcount2 - Mean  
■ /home/tomov/TAU/profile\_data/res\_mcount1 - Mean






Metric: GET\_TIME\_OF\_DAY  
Value: Exclusive  
Units: seconds

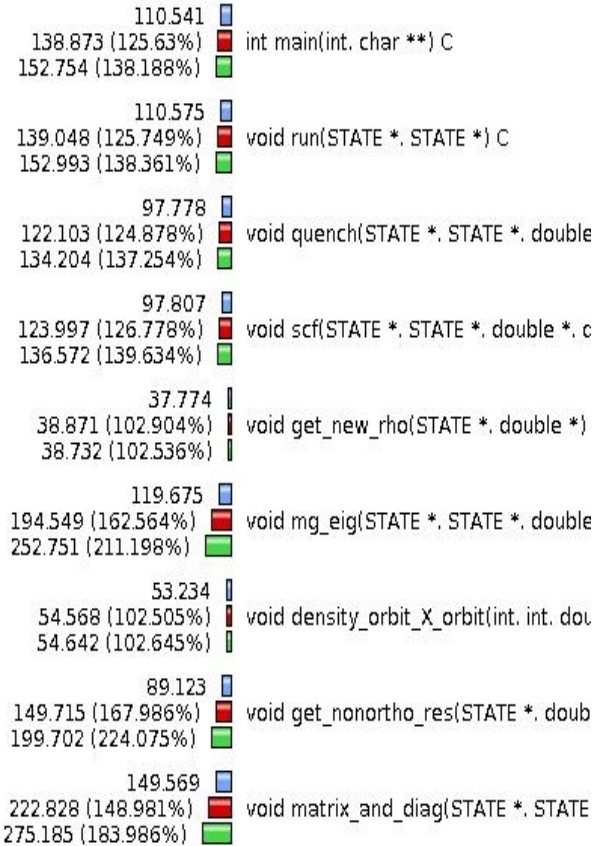
■ /home/tomov/TAU/profile\_data/res\_mcount4 - Mean  
■ /home/tomov/TAU/profile\_data/res\_mcount2 - Mean  
■ /home/tomov/TAU/profile\_data/res\_mcount1 - Mean






# Performance using single, 2, and 4 cores

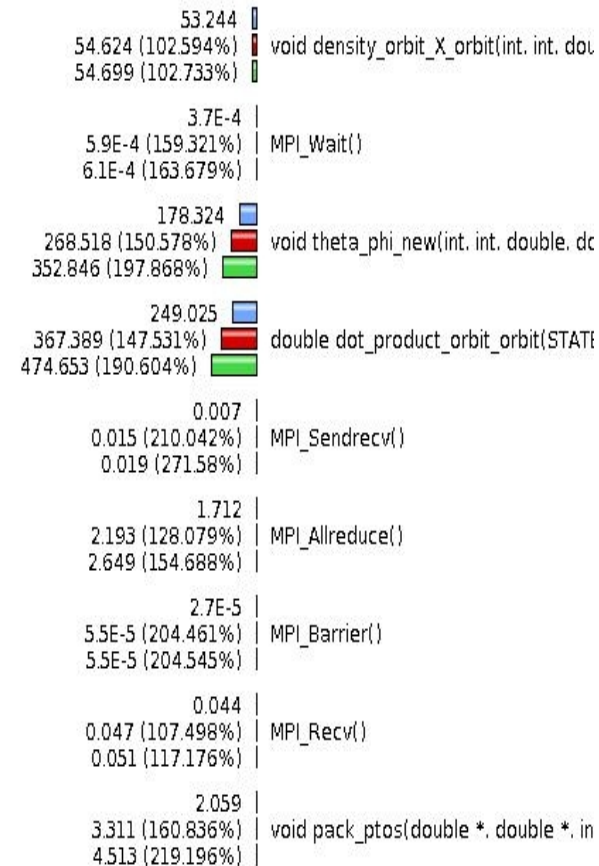
Metric: PAPI\_FP\_INS / GET\_TIME\_OF\_DAY  
 Value: Inclusive  
 Units: Derived metric shown in  
 microseconds format

 /home/tomov/TAU/profile\_data/res\_mcount4 - Mean  
 /home/tomov/TAU/profile\_data/res\_mcount2 - Mean  
 /home/tomov/TAU/profile\_data/res\_mcount1 - Mean



Metric: PAPI\_FP\_INS / GET\_TIME\_OF\_DAY  
 Value: Exclusive  
 Units: Derived metric shown in  
 microseconds format

 /home/tomov/TAU/profile\_data/res\_mcount4 - Mean  
 /home/tomov/TAU/profile\_data/res\_mcount2 - Mean  
 /home/tomov/TAU/profile\_data/res\_mcount1 - Mean





# Bottlenecks

- Maximum performance
  - Jaguar has quad-core Opterons 2.1 GHz
    - Theoretical maximum : 8.4 GFlop/s per core (~ 32 GFlop/s per quad-core)
    - Memory bandwidth : 10.6 GB/s (shared between the 4 cores)
  - Close to peak – only for operations of high enough ratio of **Flops vs data needed**
    - e.g. Level 3 BLAS for large enough N (~200)
  - Otherwise, in most cases, memory bandwidth and latencies are limits for the maximum performance
    - e.g. stream (copy) is ~ 10 GB/s (1 core enough to saturate the bus)
    - dot product ~ 1 GFlop/s (16 Bytes for 2 operations; 1 core saturates bus)
    - FFT ~ 0.7 GFlop/s (2 cores), 1.3 GFlop/s (4 cores)
    - Random sparse ~ 0.035 GFlop/s (2 cores), 0.052 GFlop/s (4 cores)

# Bottlenecks

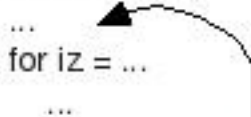
- A list of suggestions for performance improvements
  - Try some standard optimization techniques on the most compute intensive functions
  - Change the current all-MPI implementation to a multicore-aware implementation where communications are performed only between nodes
  - Try different strategies/patterns of intermixing communication and computation (e.g. early MPI\_Irecv)
  - Consider changing the algorithms if performance is still not satisfactory

# Bottlenecks removal

- Example of standard performance optimization techniques  
[e.g. DoxO, runs 29% of time, accelerated **2.6 x**; overall brings 28% acceleration]

## Original (pseudo-)code

```
for ix = ...  
  ...  
  for iy = ...  
    ...  
    for iz = ...  
      ...  
      rho_global[ f(iz) ] += ...
```



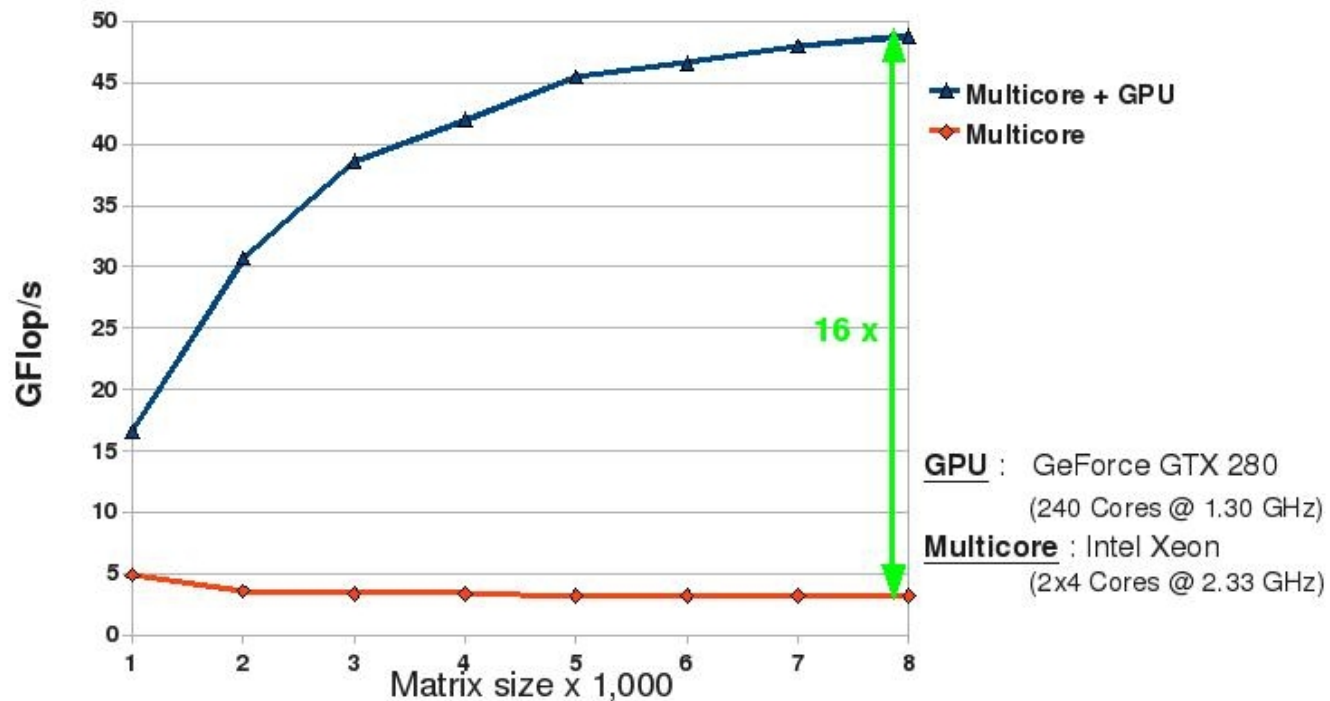
```
where int f(int x){  
  if (x<0)  
    return ...  
  else if ...  
    return ...  
}
```

## Optimized code

- The innermost loops were reorganized
  - The **ifs** of function " **f** " were taken outside of the **iz** loop
  - Unrolling (of 4 iterations)
- New variables were introduced for some recurring index calculations

# Bottlenecks removal

- **New algorithms** – advances from linear algebra for multicore and emerging hybrid architectures [e.g. hybrid Hessenberg reduction in double precision; accelerated 16x; related to the subspace diagonalization bottleneck]



# Conclusions

- We profiled and analyzed 2 petascale quantum simulation tools for nanotechnology applications
- We used different tools to help in understanding the performance on Teraflop leadership platforms
- We identified bottlenecks and gave suggestions for their removal
- The results so far indicate that the main steps that we have followed (and described) can be viewed/used as a methodology to not only easily produce and analyze performance data, but also to aid the development of algorithms, and in particular petascale quantum simulation tools, that effectively use the underlying hardware.