

Fast plane wave density functional theory molecular dynamics calculations on multi-GPU machines

Weile Jia¹, Jiyun Fu¹, Zongyan Cao¹, Long Wang¹, Xuebin Chi¹, Weiguo Gao², Lin-Wang, Wang³

1. Supercomputing Center, Chinese Computer Network Information Center, Chinese Academy of Science
2. School of Mathematical Sciences, Fudan University
3. Material Science Division, Lawrence Berkeley National Laboratory



中国科学院超级计算中心
Supercomputing Center of Chinese Academy of Sciences



Introduction

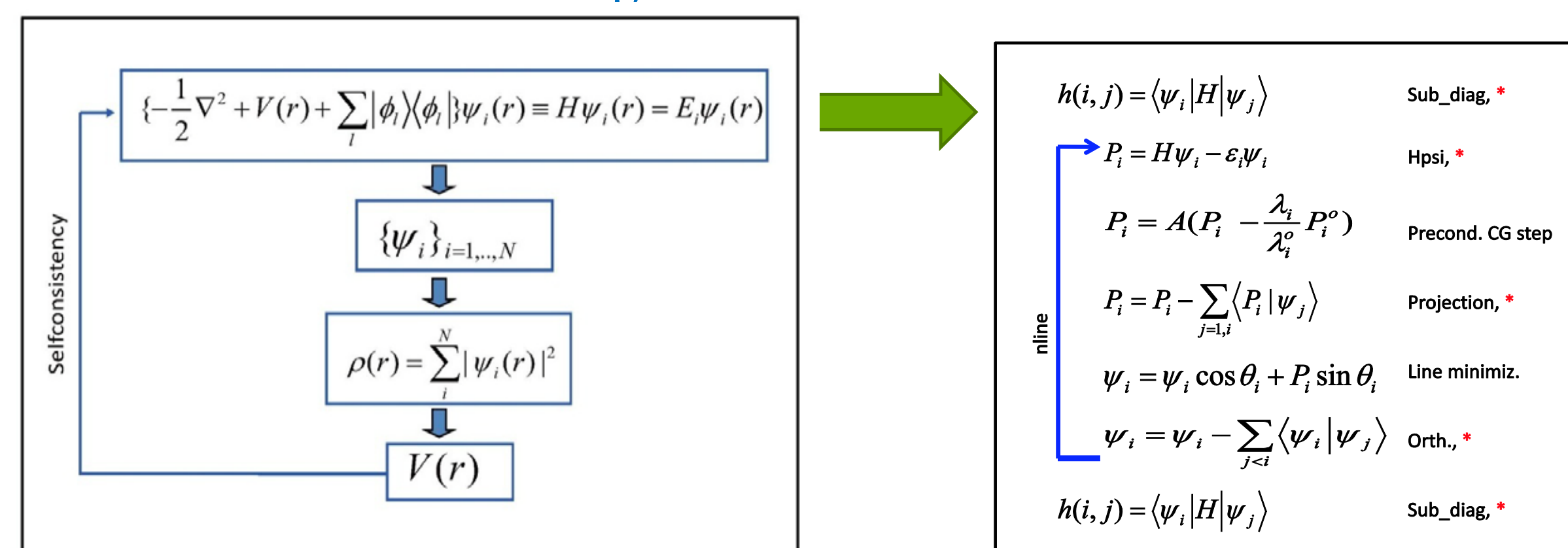
Plane wave pseudopotential density functional theory (PWP-DFT) calculation is one of the most important science simulations in the high performance computing field. In this poster, we present our drastic redesign of the algorithm and moving all the major computation parts into GPU, and reached a speed of 12 seconds per molecular dynamics(MD) step for a 512 atom system using 256 GPU cards. **This is about 7 times faster than the existing CPU runs regardless of the number of CPU cores used.** Our approach represents a shift of the computation paradigm, treating GPU as the main computing resource, instead of as an accelerator to CPU.

PWP-DFT MD Algorithm

PWP-DFT solves the Schrödinger's equation in electronic structure calculations and applies the Newton's second law in MD calculations. Electronic structure calculation is the most computation intensive part, which involves the wave function to get the force for the atoms.

$$\left\{ \begin{aligned} -\frac{1}{2}\nabla^2 + V(r) + \sum_i |\phi_i\rangle\langle\phi_i| \psi_i(r) &= E_i \psi_i(r) \\ \mathbf{F} &= m\ddot{\mathbf{a}} \end{aligned} \right. \quad \begin{array}{l} \text{Schrödinger's equation} \\ \text{Newton's second law} \end{array}$$

PWP-DFT SCF and CG algorithm



Self-consistent Filed(SCF) iteration CG method to improve wave-function

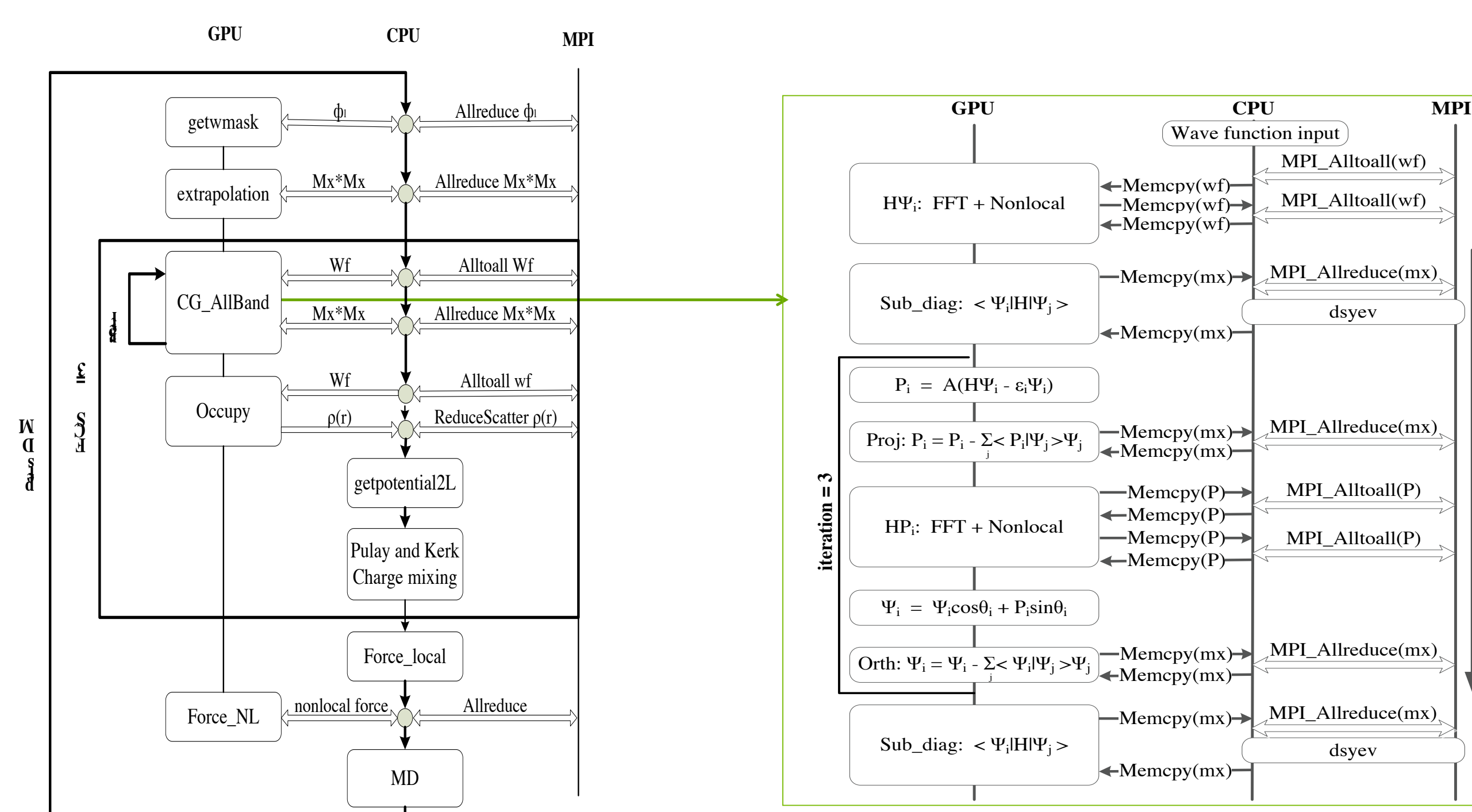
SCF calculation, which solves the Schrödinger's equation iteratively, consumes most of the computational time. Among all kernels, the wave function improvement using CG method costs more than 90% of the total computing time. The algorithms listed above are very typical SCF and CG algorithms.

PWP-DFT MD implementation and optimization

One principle of the PETot GPU design is to efficiently utilize the computing power of the GPU, so all the computing-intensive parts in MD algorithm are moved into GPU. Among all these kernels, CG_AllBand, which involves the wave function based on the Kohn-Sham equation, consumes most of the computational time. In the next session, we show the PWP-DFT MD on GPU, especially, the CG_AllBand implementation on the GPU cluster in detail.

Description for main kernels:

getmask: Nonlocal projector calculation by the position of the atoms.
CG_AllBand: Using all band CG method to improve the wave function.
occupy: calculate the charge by occupying the wave functions.
getpotential2L: Calculate output potential $V(r)$.
Pulay and Kerker potential mixing: Generate the potential (hence the Hamiltonian H) for next SCF iteration.
Force_local and **Force_NL:** Local/nonlocal atomic forces calculation respectively (by Hellmann-Feynman theory).
MD algorithms: Classic MD algorithm (Velocity Verlet, Langevin, or Nose-Hoover).



PWP-DFT MD implementation

CG method implementation

Four things matters most in optimization:

1. A hybrid parallelization scheme is essential for GPU acceleration of the PWP-DFT MD code[1].
2. Data compression for the MPI_Alltoall is critical for the performance of the multi-GPU code.
3. Keeping the wave function inside GPU, reducing CPU-GPU memory copy operation is important for the performance.
4. Using CUDA libs, such as CUFFT, MAGMA and CUBLAS.

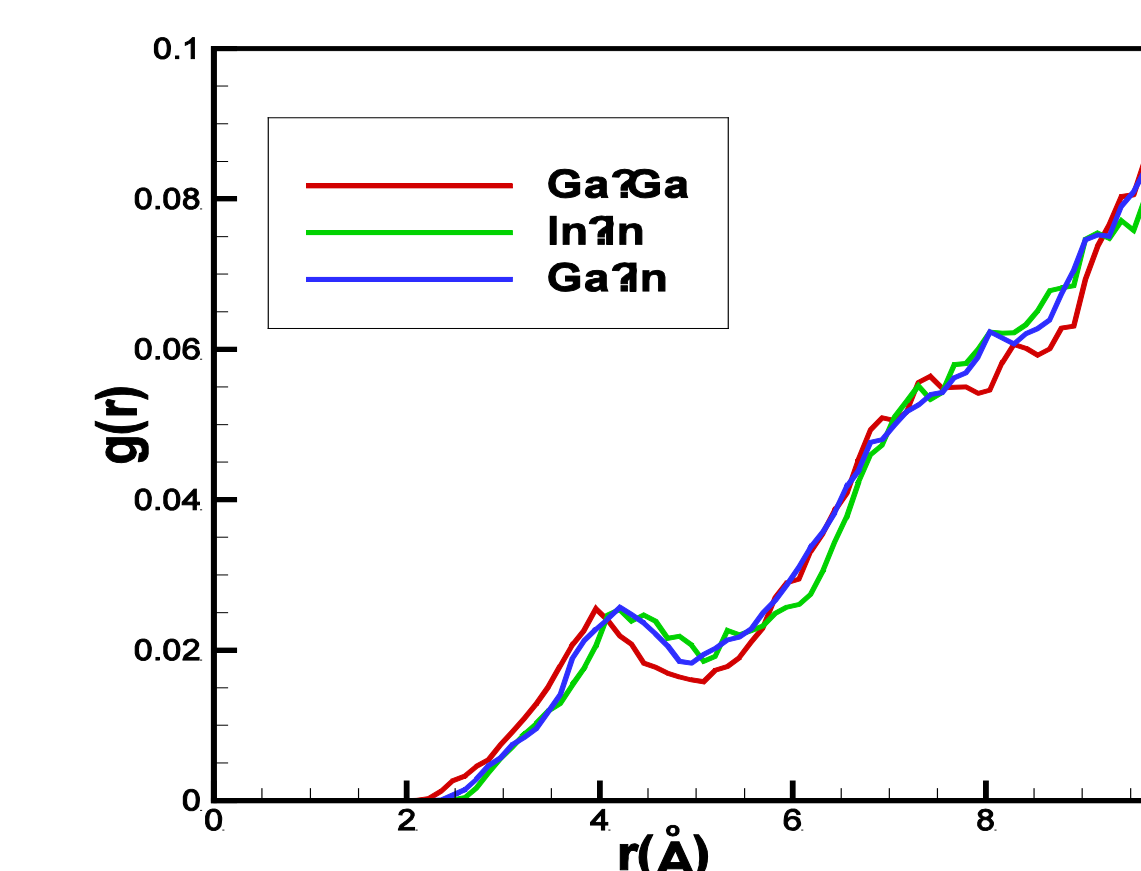
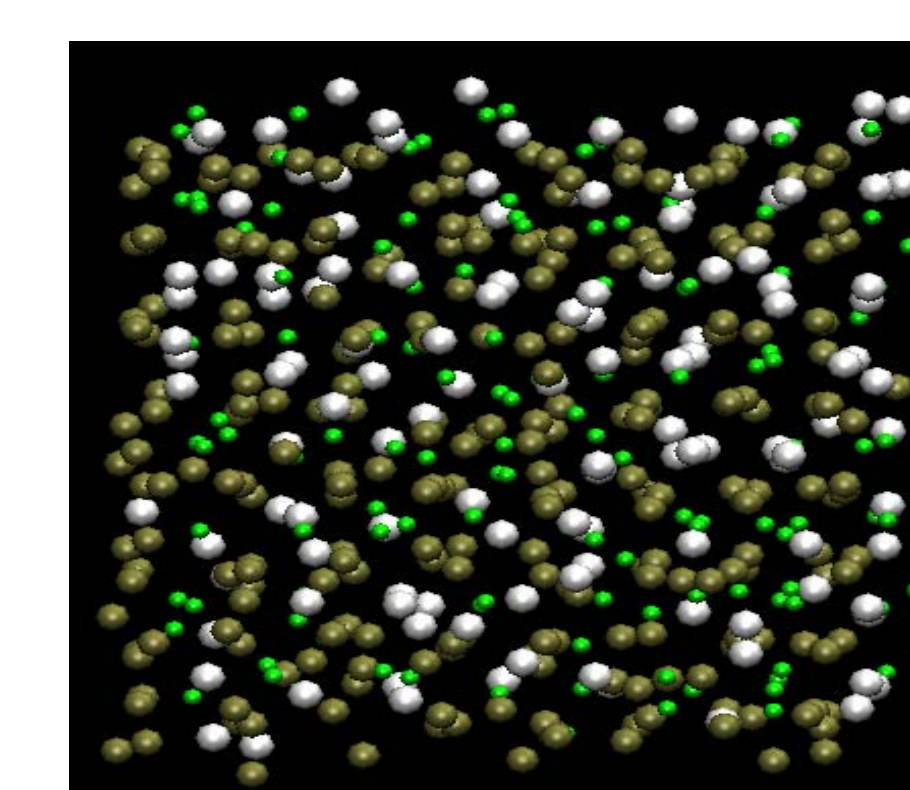
Testing results

The first testing system: 512 atom GaAs bulk system

No. of CPU core	32 × 16	64 × 16	128 × 16	256 × 16
*PETot_CPU(NP)	277s(8)	223s(8)	203s(8)	216s(8)
No. of GPU	32	64	128	256
PETot_GPU(Titan)	31.6s	20.8s	13.2s	11.4s

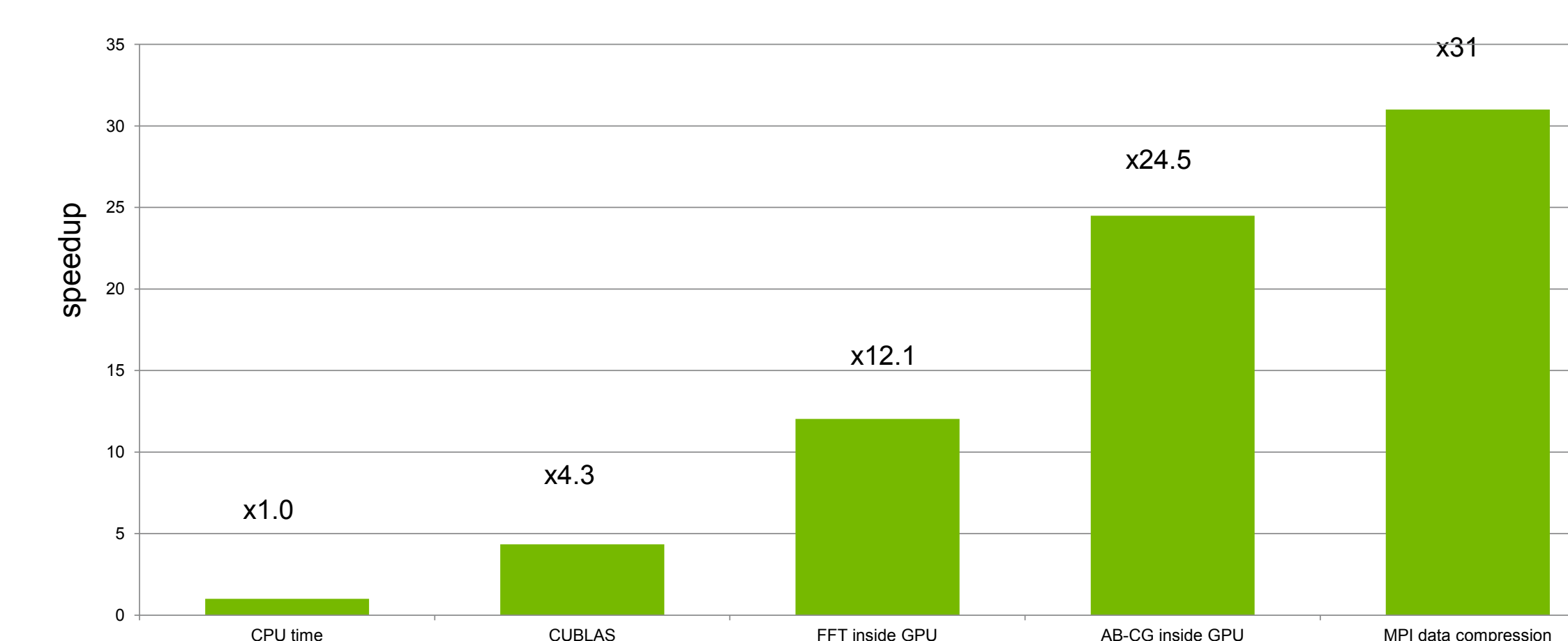
Testing systems: 512 atoms GaAs bulk system with one As replaced by N. The time is for one MD step, which contains 3 SCF. The test is performed on Titan. It has one 16-core AMD CPU and one NVIDIA Tesla X2090 GPU. Please note the results are compared between one GPU and a 16-core CPU. And we reached 12 seconds per MD step on Titan. This is much faster than any reported PWP-DFT CPU code.

The second testing system: 512 atom system of GaInP



Testing systems: The atomic correlation functions of the GaInP system with 512 atoms. The correlation functions are taken from the last 1200 steps out of 1800 total steps of MD simulations. It takes more than 4 days to calculate using CPU code, and it takes 10 hours using PETot GPU code.

Conclusion



The speedup of a multi-GPU PWP-DFT code

Using 256 GPUs, the PETot-GPU code can perform one MD step for a 512 atom system within 12 seconds, and it is about 20 times faster than the corresponding CPU regardless of the number of CPU cores used.

References

1. Long Wang, Weile Jia, Xuebin Chi, Yue Wu, WeiguoGao, Lin-Wang Wang. Large Scale Plane Wave Pseudopotential Density Functional Theory Calculations on GPU Clusters, *The International Conference for High Performance Computing, Networking, Storage, and Analysis*, 2011.
2. Weile Jia, Zongyan Cao, Long Wang, Jiyun Fu, Xuebin Chi, WeiguoGao, Lin-Wang Wang. The analysis of a plane wave pseudopotential density functional theory code on a GPU machine, *Computer Physics Communications*, Available online 14 August 2012, ISSN 0010-4655, 10.1016/j.cpc.2012.08.002.