

## **A Scalable Parallel Implementation of the Plane Wave Time Domain Algorithm on Graphics Processing Unit-Augmented Clusters**

Yang Liu<sup>(1)</sup>, Abdulkadir C. Yucel<sup>(1)</sup>, Vitaliy Lomakin<sup>(2)</sup> and Eric Michielssen\*<sup>(1)</sup>

(1) Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI 48105, USA

(2) Department of Electrical and Computer Engineering, University of California, San Diego, CA 92093-0407, USA

The computational complexity and memory requirements of direct schemes for evaluating transient electromagnetic fields produced by  $N_s$  dipoles active for  $N_t$  time steps scale as  $O(N_t N_s^2)$ . These costs can be reduced to  $O(N_t N_s \log^2 N_s)$  by the multilevel plane wave time domain (PWTD) algorithm [A.A. Ergin et al., *Antennas and Propagation Magazine*, IEEE, vol. 41, pp. 39-52, 1999]. Not surprisingly, the PWTD scheme has been used to accelerate marching on in time (MOT) -based integral equation solvers for analyzing transient scattering from complex structures. The majority of these implementations have been on serial computers though, limiting their applicability to real-life problems. To advance the capabilities of fast time domain integral equation solvers, parallel PWTD schemes and associated MOT solvers are called for.

Recently, CPU clusters with one or more graphics processing units (GPUs) on each compute-node have emerged as promising platforms for many supercomputing tasks. In principle, these hybrid CPU-GPU clusters permit code developers to take advantage of the huge computational power and unique architecture of GPUs, without being constrained by their scarce memory resources. Unfortunately, porting an existing code to a hybrid CPU-GPU cluster is no sinecure and often calls for a complete reorganization of the computational and memory management tasks.

This paper presents a parallel implementation of the multilevel PWTD scheme on a distributed-memory hybrid CPU-GPU cluster. Our work extends that in [M. Lu et al., *Antennas and Propagation Society International Symposium*, IEEE, vol. 4, pp. 4212-4215, 2004] and [Y. Liu et al., the 28<sup>th</sup> International Review of Progress in Applied Computational Electromagnetics, 2012] on pure CPU- and GPU-parallel PWTD implementations. The memory and computation loads are divided among compute-nodes by a hierarchical partitioning strategy that extends the scheme in [J. Fostier et al., *Electronics Letters*, vol.44, no.19, pp.1111-1113, 2008; O. Ergul et al., *Antennas and Propagation*, IEEE Transactions on, vol.57, no.6, pp.1740-1750, 2009] to the time domain. The proposed parallelization strategy is provably scalable and exhibits favorable load balance and computation-to-communication ratios. Most PWTD calculation stages, viz. the construction of outgoing rays, the projection of incoming rays, field translations and near-field calculations are carried out on GPUs, while their CPU hosts control host-to-host and host-to-device communications.

The proposed technique has been used to evaluate transient electromagnetic fields generated by large-scale temporally bandlimited dipole constellations. Numerical results will verify the efficiency and scalability of the hybrid CPU-GPU implementation, and compare its computational speed and memory requirements to those of pure CPU and GPU-parallel implementations.