# Conghui **He**

Ph.D. Candidate · Tsinghua University

*Room S814, Meng Minwei Science Building, Tsinghua University, Haidian, Beijing, China*

☐ (+86) 153-1177-5057 | ✉ heconghui@gmail.com | 🏠 http://conghui.github.io/ | 🐙 conghui

## **Edu**cation

**Imperial College**  *London, UK*

Visiting PhD student in Custom Computing Group  *Nov. 2016 - PRESENT*

- Supervisor: Wayne Luk (https://www.doc.ic.ac.uk/ wl/)
- Topic: Ultra-low-latency market data generators on reconfigurable platforms

**Stanford University**  *California, USA*

Visiting PhD student in Stanford Exploration Project Group (SEP)  *Jun. 2016 - Sep. 2016*

- Supervisor: Bob Clapp (http://sepwww.stanford.edu/sep/bob/)
- Topic: Approximating Q propagation in 3D elastic modeling on HPC platforms

**Tsinghua University**  *Beijing, China*

Ph.D. Candidate in Department of Computer Science and Technology  *Aug. 2013 - PRESENT*

- I am a member of High Performance Geoscience Computing Group (HPGC) in Tsinghua University (http://thuhpgc.org)
- My research mainly focus on computational geophysics and financial applications on reconfigurable platforms

**Sun Yat-sen University**  *Guangdong, China*

B.S. in Software Engineering  *Sep. 2009 - Jul. 2013*

## **Ski**lls

| | |
|---|---|
| **Programming** | C/C++, Python, Java, Matlab, Bash |
| **HPC** | Profiling, Optimization, Pthread, OpenMP, MPI, GPU, FPGA |
| **Geophysics** | Modeling, Reverse Time Migration (RTM), Full Waveform Inversion (FWI), Beam, Seismology |
| **Finance** | Market Data Generator, Financial Information EXchange (FIX) protocol |

## **Hon**ors & Awards

| | | |
|---|---|---|
| 2014 | **Scholarship**, Schlumberger Scholarship for Computing Earth Science | *Schlumberger* |
| 2013 | **1st place**, IEEE/IBM International Smarter Planet Challenge | *IEEE/IBM* |
| 2012 | **4th place**, International Supercomputing Challenge | *ISC12* |
| 2011 | **Scholarship**, Sun Yat-sen's First Prize Student Scholarship | *Sun Yat-sen* |
| 2010 | **Scholarship**, IBM Outstanding Student Scholarship | *IBM* |
| 2009 | **Scholarship**, National Scholarship | *China* |

## **Wor**k Experience

**Maxeler Ltd**  *London, UK*

Intern in Networking Group  *Nov. 2016 - PRESENT*

- Following the MaxMPT project collaborated with Chicago Mercantile Exchange (CME)

**National Supercomputing Center in Wuxi**  *Wuxi, China*

Software Engineer & HPC Researcher  *Feb. 2016 - Nov. 2017*

- Optimize geophysics applications on the new Sunway manycore supercomputer
- 15-Pflops Nonlinear Earthquake Simulation on Sunway TaihuLight: Enabling Depiction of Realistic 10 Hz Scenarios
- Refactor and optimize the Community Atmosphere Model (CAM) on the new Sunway manycore supercomputer

**Statoil (Beijing) Technology Service Co, Ltd**  *Beijing, China*

Intern in Seismic Imaging R&D Group  *Jul. 2014 - Sep. 2014*

- Design and implement a CPU-GPU hybrid parallel strategy for beam migration

**IEEE Tsinghua Student Branch** *Beijing, China*
Student Chair *Sep. 2013 - Jul. 2014*

## Research/Projects

My interests include computational geophysics and parallel algorithms. I'm experienced in parallel algorithm designs on heterogeneous architectures like GPU, multi-core CPU, and FPGA processors to solve computational challenges raised from geoscience applications. I am also interested in novel designs for financial applications on reconfigurable platforms. Participated projects include:

### An Extremely Low-latency Market Server on Reconfigurable Platforms *Project*

Project leader, cooperated with China Financial Future Exchange (CFFEX) *May. 2015 - Jun. 2016*

- We design an FPGA-based accelerated approach to market data processing, with an FPGA connected directly to the network to parse, split, filter the financial packets, and then push the market data feeds directly to the network after reconstructions of order books. Such a solution offers flexibility, as the FPGA can be reconfigured for different protocols and market processing logic, and high throughput with extremely low latency by eliminating the operating system's interrupts and network stacks. This work proposes a hybrid sorted table design for minimizing electronic trading latency, with three main contributions. First, a hierarchical sorted table with two levels, a fast cache table in reconfigurable hardware storing megabytes of data items and a master table in software storing gigabytes of data items. Second, a full set of operations, including insertion, deletion, selection and sorting, for the hybrid table with latency in a few cycles. Third, an on-demand synchronization scheme between the cache table and the master table. An implementation has been developed that targets an FPGA-based network card in the environment of the China Financial Futures Exchange (CFFEX) which sustains 1-10Gb/s bandwidth with latency of 400 to 700 nanoseconds, providing an 80- to 125-fold latency reduction compared to a fully optimized CPU-based solution, and a 2.2-fold reduction over an existing FPGA-based solution.
- Related work is published in FPGA17', FCCM17'

### 15-Pflops Nonlinear Earthquake Simulation on Sunway TaihuLight: Enabling Depiction of Realistic 10 Hz Scenarios *Research*

Project leader *Nov. 2016 - Apr. 2017*

- This paper reports our work on building a highly efficient earthquake simulation platform on Sunway TaihuLight, with 125 Pflops computing power and over 10 million cores. With the platform originated from AWP-ODC and CG-FDM, a large part of our efforts focuses on redesigning the velocity, stress, and plasticity processing kernels for the completely different microarchitecture and significantly increased parallelism of Sunway TaihuLight. By a combined approach that includes (1) an optimized parallelization scheme, (2) the most suitable blocking configuration, (3) fusion of co-located arrays, (4) register communication with CPE ID remapping for halo exchanges, and (5) customized ROM-less evaluation of elementary function, we manage to achieve an efficient utilization of over 12.5% of the theoretical peak of the entire system. Our simulation program provides a sustained performance of over 15 Pflops, and enables the simulation of the Tangshan earthquake with a spatial resolution of 25 m and a frequency of 10 Hz.
- Related work is submitted to SC17', Gordon Bell Prize

### A Fully-Pipelined Hardware Design for Gaussian Mixture Models *Research*

Topic leader *Apr. 2015 - April. 2017*

- Gaussian Mixture Models (GMMs) are widely used in many applications such as data mining, signal processing and computer vision, for probability density modeling and soft clustering. However, the parameters of a GMM need to be estimated from data by, for example, the Expectation-Maximization algorithm for Gaussian Mixture Models (EM-GMM), which is computationally demanding. This paper presents a novel design for the EM-GMM algorithm targeting reconfigurable platforms, with five main contributions. First, a pipeline-friendly EM-GMM with diagonal covariance matrices that can easily be mapped to hardware architectures. Second, a function evaluation unit for Gaussian probability density based on fixed-point arithmetic. Third, our approach is extended to support a wide range of dimensions or/and components by fitting multiple pieces of smaller dimensions onto an FPGA chip. Fourth, we derive a cost and performance model that estimates logic resources. Fifth, our dataflow design targeting the Maxeler MPC-X2000 with a Stratix-5SGSD8 FPGA can run over 200 times faster than a 6-core Xeon E5645 processor, and over 39 times faster than a Pascal TITAN-X GPU. Our design provides a practical solution to applications for training and explores better parameters for GMMs with hundreds of millions of high dimensional input instances, for low-latency and high-performance applications.
- Related work is accepted by IEEE Transactions on Computers

### Approximating Q propagations for elastic modeling on GPUs

*Research*

PROJECT LEADER, COLLABORATED WITH STANFORD UNIVERSITY

*July. 2016 - Jan. 2017*

- Propagating wavefields using the explicit finite difference method is the kernel of reverse time migration (RTM) and high-end velocity algorithms in seismic applications. In recent decades there has been a significant increase of interest in the seismic exploration community to invert the image of the subsurface in larger regions and higher resolutions in the elastic media, which brings tremendous computing challenges. As a result, the optimizing methods for improving the performance of the wavefield propagation are in great demands. This work manages to boost the performance of the wavefield propagation in 3D elastic scenarios by approximating the Q propagation and using the multi-GPU platform. We first extend the constant-Q formulation from the 2D viscoelastic case to the 3D viscoelastic case. Different optimization techniques on GPUs are then described for an efficient modeling kernel. We also propose a set of schemes to reduce the computation to further improve the performance. The experimental results show that we can achieve significant performance speedups of 60 to 200 times with 4 GPUs over the CPU-based solution as a function of Q.
- Related work is published in EAGE17'

### Refactoring and Optimizing the Community Atmosphere Model (CAM) on the New Sunway Manycore Supercomputer

*Project*

CORE MEMBER

*Jul. 2015 - Apr. 2016*

- We refactor and optimize the Community Atmosphere Model (CAM) on the new Sunway many-core supercomputer of China, which is the rank 1 supercomputer in the latest Top 500 announcement. It uses a many-core processor that consists of management processing elements (MPEs) and clusters of computing processing elements (CPEs). To tackle the major challenges of mapping the large code base of CAM to the millions of cores on the Sunway system, we take OpenACC-based refactorization as the major tool, and apply source-to-source translator tools to generate the most suitable parallelism for the CPE cluster, and to fit the intermediate variable into the limited on-chip fast buffer. For single kernels, when comparing the originally ported version using only MPEs and the refactorized version using both the MPE and CPE clusters, we achieve up to 22x speedup for the compute-intensive kernels. For the 25km resolution CAM global model, we manage to scale to 24,000 MPEs, and 1,536,000 CPEs and achieve a simulation speed of 2.81 model years per day.
- Related work is published in SC16'

### Ensemble Full Waveform Inversion with Source Encoding

*Research*

TOPIC LEADER

*Sep. 2014 - Jun. 2015*

- Full waveform inversion (FWI) suffers from convergence toward local minima because of the inaccuracy of the initial model and the lack of low frequency data. Noises in seismograms further deteriorate the imaging quality. To relax the dependency on high-quality low-frequency data, we present an ensemble full waveform inversion method with source encoding (EnFWI), which is an ensemble approximation of the total inversion proposed by Tarantola. The method refines the velocity model iteratively by incorporating the observation, while the nonlinear evolution of the covariance is approximated by ensemble covariance. Encoded simultaneous-source FWI (ESSFWI) is applied to improve the representation for the low rank ensemble approximation, and to increase the rate of convergence. Experiments show that EnFWI achieves larger convergence range and better tolerance to data noise with less computational costs than traditional FWI methods.
- Related work is published in SEG16'

### A CPU-GPU Hybrid Parallel Design for Beam Migration

*Project*

PROJECT LEADER

*Sep. 2013 - Dec. 2014*

- The Kirchhoff beam-stack migration is quite popular in production with both better image quality and faster speed compared to Kirchhoff migration. However, the beam forming step and beam mapping step are still expensive. Meanwhile, continuous High Performance Computing (HPC) developments offer new opportunities for the industry to further enhance the efficiency of beam migration methods. We present a design of a highly efficient CPU-GPU hybrid beam migration. By parallelizing both the beam forming and the beam mapping routines with millions of GPU threads and using an asynchronous IO scheme, we derive a parallel beam migration design that fits current CPU-GPU hybrid clusters. Then, we test our GPU-based beam migration on the SEG/EAGE salt model and the SEAM salt model for different generations of GPU architectures, presenting accurate imaging results with 4-12 times speedup compared to a parallel 16-core CPU design. The significant performance improvement would further close the gap to an interactive migration engine.
- Related work is published in EAGE15'

### Accelerating the Global Vegetation-Precipitation Correlation Algorithm

*Research*

IMPROVE THE PERFORMANCE OF THE CODE

*Sep. 2013 - Nov. 2013*

- Startup Project for Ph.D. candidate cooperated with a Professor in Remote Sensing field, aiming to accelerate the algorithm taking months to finish. Optimization strategies for it include modifying the algorithm to reduce I/O accessing by utilizing local buffer, adding a memory pool to reduce frequent memory allocation/destruction, overlapping I/O transferring and computing. It gained 20x speedup in the end.
- Related work is published in the journal of Remote Sensing

# International Conference & Visit

| | | |
|---|---|---|
| April, 2017 **Speaker,** IEEE International Symposium on Field-Programmable Custom Computing Machines | | *Napa, USA* |
| Nov, 2016 **Student,** Custom Computing Group in Imperial College | | *London, UK* |
| Jun, 2016 **Student,** Stanford Exploration Project Group (SEP) in Stanford | | *Stanford, USA* |
| Jun, 2016 **Exibitor,** International Supercomputing (ISC) High Performance 2016 | | *Frankfurt, Germany* |
| Oct, 2015 **Speaker,** Society of Exploration Geophysicists (SEG) Annual Meeting 2015 | | *New Orleans, USA* |
| Jul, 2015 **Student,** OpenSPL Summer School Symposium | | *London, UK* |
| Jun, 2015 **Speaker,** European Association of Geoscientists and Engineers (EAGE) 2015 | | *Madrid, Spain* |

# **Pub**lications

[1] Conghui He, Haohuan Fu, Ce Guo, Wayne Luk and Guangwen Yang. "A Fully-Pipelined Hardware Design for Gaussian Mixture Models." IEEE Transactions on Computers (Accepted, not published yet).

[2] Fu, Haohuan, Conghui He, Wayne Luk, Weijia Li, and Guangwen Yang. "A Nanosecond-level Hybrid Table Design for Financial Market Data Generators." The 25th IEEE International Symposium on Field-Programmable Custom Computing Machines.

[3] Fu, Haohuan, Conghui He, Huabin Ruan, Itay Greenspon, Wayne Luk, Yongkang Zheng, Junfeng Liao, Qing Zhang, and Guangwen Yang. "Accelerating Financial Market Server through Hybrid List Design." In Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, pp. 289-290.

[4] Fu, Haohuan, et al. "Refactoring and optimizing the community atmosphere model (CAM) on the sunway taihulight supercomputer." Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE Press, 2016.

[5] Chen, Yushu, Guangwen Yang, Xiao Ma, Conghui He, and Guojie Song. "A time-space domain stereo finite difference method for 3D scalar wave propagation." Computers & Geosciences 96 (2016): 218-235.

[6] Chen, Bingwei, Conghui He, Yushu Chen, Haohuan Fu. "Full Wave Inversion Based on EnKF and Source Encoding" In 2016 SEG Annual Meeting. Society of Exploration Geophysicists.

[7] He, C., Chen, Y., Fu, H., & Yang, G. Ensemble Full Wave Inversion with Source Encoding. In 77th EAGE Conference and Exhibition 2015.

[8] He, Conghui, Haohuan Fu, Bangtian Liu, Huabin Ruan, Guangwen Yang, Hui Yang, and Are Osen. "A GPU-based Parallel Beam Migration Design." In 2015 SEG Annual Meeting. Society of Exploration Geophysicists, 2015.

[9] Clinton, Nicholas, Le Yu, Haohuan Fu, Conghui He, and Peng Gong. "Global-Scale Associations of Vegetation Phenology with Rainfall and Temperature at a High Spatio-Temporal Resolution." Remote Sensing 6, no. 8 (2014): 7320-7338.