

Xu, RuQing

Deep learning algorithm engineer focusing on fast CUDA kernels, op-fusion, and world foundation models.

<https://github.com/xrq-phys/>
<https://rxu.purineko.me/>

+81 70 3760 4882
rxu.sol@gmail.com

Education

The University of Tokyo Ph.D. in Physics received on 24-SEP-2024 Focus: <i>Statistical Physics (Computer Simulations: Algorithmic and Platform Optimizations)</i>	01-OCT-2021 – 24-SEP-2024 Tokyo, JP
The University of Tokyo M.S. in Physics received on 24-SEP-2021 Focus: <i>Statistical Physics: (Computer Simulations & HPC Basics)</i>	20-SEP-2019 – 30-SEP-2021 Tokyo, JP
University of Science and Technology of China B.S. in Physics received on 01-JUL-2019 Dept: <i>Theoretical Physics</i>	01-SEP-2015 – 01-JUL-2019 Hefei, CN

Employment History

Snr. Deep Learning Algorithm Engineer NVIDIA	29-JUL-2024 – Present Tokyo, JP
• Developing multiple variants of high-throughput attention kernels for NVIDIA's DiT-based research workloads. • Ensuring throughput-critical kernels work as fast as possible at our current and next-generation GPU models. • Taking charge of the optimizations before launch of Cosmos DiT models as a NVIDIA Inference Microservice. • Deliver Tensor-Core-emulated SGEMM & CGEMM kernels to cuBLAS & cuTENSOR for NVIDIA Blackwell architectures.	
Deep Learning Algorithm Engineering Intern (Deep Learning) NVIDIA (<i>Intern</i>)	31-JAN-2023 – 28-JUL-2024 Tokyo, JP

• Helps development & performance verification for workload-specific kernels.
• Synchronize with research teams to make their workloads run as fast as possible.
• Integrate algorithms for Tensor-Core-based fast matrix multiplication from published materials into cuTENSOR. Polish the implementation so as to align with the best that the latest hardware (Hopper GMMA) can provide.

Deep Learning Algorithm Engineering Intern (cuTENSOR)
NVIDIA (*Intern*)

• Special techniques for medium-size performance improvements.
• Tackled multiple unusual underperforming cases.
• In-depth L2 bandwidth analysis & optimizations.

Publications

<i>SCA/HPCAsia 2026</i> , Angelika Schwarz, Anton Anders, Cole Brower, Harun Bayraktar, John Gunnels, Kate Clark, RuQing G. Xu , Samuel Rodriguez, et. al., <i>Guaranteed DGEMM Accuracy While Using Reduced Precision Tensor Cores Through Extensions of the Ozaki Scheme</i>
<i>SIAM SISC (in review)</i> , Ishna Satyarth, Chao Yin, Devin A. Matthews, Maggie Myers, Robert van de Geijn, RuQing G. Xu , <i>Performant Tridiagonal Factorization of Skew-symmetric Matrices</i>
<i>arXiv:2311.10700</i> , Robert van de Geijn, Maggie Myers, RuQing G. Xu , Devin Matthews, <i>Deriving Algorithms for Triangular Tridiagonalization a (Skew-)Symmetric Matrix</i>
<i>ICS '23: Proc. 37th Intl. Conf. Supercomputing: pp. 111–121</i> , RuQing G. Xu , Field G. Van Zee, Robert A. van de Geijn, <i>Towards a Unified Implementation of GEMM in BLIS</i>
<i>Comput. Phys. Commun. 277, 108375</i> , RuQing G. Xu , Tsuyoshi Okubo, Synge Todo, Masatoshi Imada, <i>Optimized Implementation for Calculation and Fast-Update of Pfaffians Installed to the Open-Source Fermionic Variational Solver mVMC</i>
<i>Phys. Rev. Research 3, 023048</i> , Xinliang Lyu, RuQing G. Xu , Naoki Kawashima, <i>Scaling dimensions from linearized tensor renormalization group transformations</i>
<i>J. Chem. Theory Comput. 2019, 15, 3, 1728-1742</i> , James S. Spencer, Nick S. Blunt, ..., William A. Vigor, RuQing Xu , Alex J. W. Thom, <i>The HANDE-QMC project: open-source stochastic quantum chemistry from the ground state up</i>

Skills

Programming Languages: C, C++, CUDA, NVPTX, x86 Assembly, Arm64 Assembly, Julia, Python, OpenAI Triton
Natural Languages: Chinese (native), English (GRE: 330; TOEFL: 108), Japanese (JLPT N1)