

Yonghao Tan CV

Yonghao Tan

- Address: HKUST, Clear Water Bay, Hong Kong, China
- Phone: +852 53946864
- E-mail: ytanaz@connect.ust.hk, yonghaot1017@gmail.com
- Personal Homepage: <https://yonghao-tan.github.io/>

Education & Academic Performance

- **Ph.D. in Electronic and Computer Engineering**, The Hong Kong University of Science and Technology (HKUST), Hong Kong, China (Sept. 2023 - Present)
Supervised by Prof. Tim Kwang-Ting CHENG.
- **B.E. in Microelectronics**, Southern University of Science and Technology (SUSTech), Shenzhen, Guangdong, China (Sept. 2019 - Jun. 2023)
Supervised by Prof. Fengwei An.
- Overall GPA: **3.77 / 4.0**
- Weighted Average Score: **90.38**
- Rank: **11 / 77**

Research Interests

- Software/hardware co-design
- Model compression
- 3D processing

Research Experience

5nm UCle-Enabled Multi-Chiplet Generalizable Rendering Processor

- **Date:** Jan. 2025 - Present
- **Affiliation:** AI Chip Center for Emerging Smart Systems, Hong Kong, China
- Architected a 5nm 4-chiplet GeNeRF processor for generalizable rendering to address the heavy external-memory traffic induced by multi-view feature fetching.
- Introduced a UCle-enabled cross-die unified cache, distributed source-view management, and balance-aware scheduling to maximize source-view reuse while reducing off-chip and cross-die data movement.
- Silicon results reached 91.43 TOPS/W, 55.43 FPS real-time rendering, and 0.29 uJ/pixel through hierarchical sparsity and hybrid NeRF-SR execution in a 45 mm x 45 mm MCM package footprint.

55nm ReRAM-on-Logic Stacked LLM Accelerator for Speculative Decoding

- **Date:** Apr. 2024 - Present
- **Affiliation:** AI Chip Center for Emerging Smart Systems, Hong Kong, China
- Architected a 55nm edge LLM accelerator whose logic die is stacked with four ReRAM dies via face-to-face bump bonding to support in-stack storage of draft-model codebooks.
- Developed block-clustered weight compression, local-rotation-based outlier-free W4A8 quantization, and adaptive parallel speculative decoding to reduce both target-model EMA and rejected-draft overhead.
- The prototype delivered 14.08 to 135.69 token/s and a 4.46x to 7.17x throughput speedup over a BF16 speculative-decoding baseline on a 55.98 mm² logic die.

28nm CNN-Transformer Accelerator for Semantic Segmentation

- **Date:** Nov. 2021 - Sept. 2024
- **Affiliation:** AI Chip Center for Emerging Smart Systems, Hong Kong, China
- Architected a 28nm memory-compute-intensity-aware accelerator for high-resolution ConvFormer and SegFormer semantic-segmentation workloads.
- Combined hybrid-attention processing, data-reuse-oriented layer fusion, and cascaded feature-map pruning

to reduce attention-side EMA and eliminate redundant KV and weight movement across fused blocks.

- Silicon results achieved 0.22 uJ/token on SegFormer-B0 and up to 52.90 TOPS/W peak efficiency in a 13.93 mm² chip.

Publications

- **2026** Tan, Y.*, Ma, S.*, Dong, P., Luo, P., Lei, Z., Lu, W., Ying, G., ... & Cheng, K. T.
A 5nm 91.43 TOPS/W 4-Chiplet Generalizable-Rendering-Processor with UCIe-Enabled Cross-Die-Cache and Balance-Aware Progressive Multi-Level Sparsity.
In 2026 IEEE Custom Integrated Circuits Conference (CICC), IEEE.
- **2026** Dong, P., Tan, Y., Liu, X., Luo, P., Liu, Y., Pang, D., Ma, S., ... & Cheng, K. T.
A 14.08-to-135.69Token/s ReRAM-on-Logic Stacked Outlier-Free Large-Language-Model Accelerator with Block-Clustered Weight-Compression and Adaptive Parallel-Speculative-Decoding.
In 2026 IEEE International Solid-State Circuits Conference (ISSCC), IEEE.
- **2025** Dong, P.*, Tan, Y.*, Liu, X., Luo, P., Liu, Y., Liang, L., ... & Cheng, K. T.
A 28nm 0.22uJ/Token Memory-Compute-Intensity-Aware CNN-Transformer Accelerator with Hybrid-Attention-Based Layer-Fusion and Cascaded Pruning for Semantic-Segmentation.
In 2025 IEEE International Solid-State Circuits Conference (ISSCC), IEEE.
- **2024** Dong, P.*, Tan, Y.*, Zhang, D., Ni, T., Liu, X., Liu, Y., ... & Cheng, K. T.
Genetic Quantization-Aware Approximation for Non-Linear Operations in Transformers.
In 2024 61st ACM/IEEE Design Automation Conference (DAC), IEEE.
- **2022** Tan, Y.*, Deng, H.*, Sun, M., Zhou, M., Chen, Y., Chen, L., ... & An, F.
A Reconfigurable Coprocessor for Simultaneous Localization and Mapping Algorithms in FPGA.
IEEE Transactions on Circuits and Systems II: Express Briefs, 70(1), 286-290.
- **2022** Tan, Y.*, Sun, M.*, Deng, H., Wu, H., Zhou, M., Chen, Y., ... & An, F.
A Reconfigurable Visual-Inertial Odometry Accelerator with High Area and Energy Efficiency for Autonomous Mobile Robots.
Sensors, 22(19), 7669.

* Authors marked with an asterisk contributed equally to the corresponding work.

Honors and Awards

- Aug. 2025: Best Teaching Assistant Award, Department of Electronic and Computer Engineering, HKUST
- May 2023: Outstanding Graduate (School Level), SUSTech
- Sept. 2022: First-Class Outstanding Students Scholarship with the highest score
- Apr. 2022: Undergraduate Innovation and Entrepreneurship Training Program
- Dec. 2021: Shenzhen Longsys Electronics Company Award (Top 2% in the School of Microelectronics)
- Dec. 2021: First Prize, 2021 National College Students FPGA Innovation Design Competition (Top 22 out of 1,341 teams)
- Oct. 2021: First Prize, 2021 International Competition of Autonomous Running Robots (1st place out of 34 finalist teams)

Fundings

- Sept. 2023 - Present: Postgraduate Studentship (PGS) Award in HKUST
- Apr. 2022: Undergraduate Innovation and Entrepreneurship Training Programs (Provincial Level)
- Jul. 2021: Guangdong College Students' Scientific and Technological Innovation (Provincial Level)

Skills

- Programming language: C, C++, Java, Python, System Verilog, Verilog HDL, VHDL
- Professional software: AutoCAD, Cadence, Design Compiler, IC Compiler II, MATLAB, Multisim, Silvaco

Languages

- English (fluent), Mandarin (native), Cantonese (native)