Nikhil Rout

 $\mbox{$\lozenge$}$ Hyderabad, India $\mbox{$\boxtimes$}$ nikhil
rout97@gmail.com $\mbox{$\nwarrow$}$ +91 6303583185 $\mbox{ in }$ nikhil-rout

NikhilRout
NikhilRout.github.io

Education

Vellore Institute of Technology, Chennai Campus

2022 - 2026

Bachelor of Technology, Electronics and Communication Engineering

CGPA: 9.11/10

Research Experience

Research Intern 😯

Remote

Vortex GPGPU, ORCAS Lab, UCLA

May 2025 - Present

- Developed a Tensor Core Unit (TCU) extension for the open-source RISC-V based Vortex GPGPU, enabling efficient WMMA operations for deep learning workloads
- Architected and implemented a high-performance configurable mixed-precision Fused Dot Product (FEDP) unit as part of the TCU in SystemVerilog
- \circ Extended custom float conversions, SGEMM test, and cycle-level simulator implementation for FP8/BF8 support in C++
- Introduced a novel fused pipeline stratergy for supporting Integer arithmetic within the floating-point datapath, maximizing component reuse while incurring minimal overhead
- Performed microarchitecture design space explorations, evaluating 3:2 vs 4:2 compressors in Carry-Save Adders, and format-wise Dedicated vs Shared multiplier schemes
- Reduced dynamic power consumption through clock/enable gating across FEDP pipeline stages during operation in 2:4 structured sparsity mode
- Contributed to TCU specific CI/CD and Verilator-based unit testing verification framework
- \circ Achieved 362.2 MHz operation with 4-cycle latency on the Xilinx U55C FPGA, delivering 1.448 GFLOPS single-cycle throughput—4.3× the open-source baseline Berkeley HardFloat

Summer Research Intern 🖸

Remote

Centre for Nanoelectronics and VLSI Design, VIT Chennai 🗹

May 2024 - July 2024

- Implemented pipelined radix-2 FFT modules using both DIT and DIF approaches
- Developed Number Theoretic Transform (NTT) modules for efficient polynomial multiplication in lattice-based post-quantum cryptography schemes on FPGAs

Research Publications

• A Configurable Mixed-Precision Fused Dot Product Unit for GPGPU Tensor Computation

Presentation at Vortex Workshop and Tutorials, IEEE/ACM MICRO 2025 Preprint 🗹 🗘

Optimizing RGB to Grayscale, Gaussian Blur and Sobel-Filter operations on FP-GAs for reduced dynamic power consumption

Presented at $3^{\rm rd}$ IEEE conference on AIIoT 2024 IEEE Xplore \square

Projects

Carry-Save FAN modification (SIGMA DNN Accelerator) 🗘

- Optimized Forward-Adder-Network (FAN) critical path by extending a Carry-Save-Adder tree structure for multi-vector reduction in sparse and irregular GEMM workloads
- Systolic Array Processing Element integration and synthesis PPA analysis in progress

Flash-Attention CUDA MegaKernel Inference Optimization

- Implemented CUDA kernels for Flash-Attention inference with baseline and Tensor Core variants, leveraging shared memory tiling and online softmax to reduce HBM bottleneck
- Profiled kernel performance on NVIDIA T4 GPU, demonstrating speedup from operator fusion and memory hierarchy optimization

Skills and Competencies

Languages: Verilog/SystemVerilog, C/C++, CUDA, Python, TCL, MATLAB

Tools: Quartus Prime, Xilinx Vivado/Vitis, Synopsys VCS/DC, CMake, Git, LaTeX

Relevant Coursework: Digital Systems Design, VLSI System Design, Computer Architecture, FPGA-based System Design, Heterogeneous Computing Systems, Deep Learning

References

Dr. Blaise-Pascal Tine

UCLA

ORCAS Lab, Assistant Professor

☑ blaisetine@cs.ucla.edu

Dr. Jean Jenifer Nesam J

VIT Chennai

Assistant Professor

☑ jeanjenifernesam.j@vit.ac.in

Dr. Umadevi S

VIT Chennai

CNVD, Associate Professor

☑ umadevi.s@vit.ac.in