

Project Name	Brief Project Description	Primary Advisor
RAG on the EDGE	Implement Retrieval-Augmented Generation (RAG) on AMD NPU, Qualcomm QIDK, and NVIDIA Jetson platforms, followed by a comparative performance evaluation across these edge accelerators.	Suresh Purini
RAG Accelerator on Xilinx Alveo U50	Implement a hardware accelerator for Retrieval-Augmented Generation (RAG) on the Xilinx Alveo U50 FPGA, optimizing compute and memory usage to improve latency, throughput, and energy efficiency. Compare its performance against GPU- and CPU-based implementations using available frameworks, without requiring custom accelerator development on those platforms.	Suresh Purini
RAG Accelerator on Xilinx Versal VCK5000	Develop a Retrieval-Augmented Generation (RAG) accelerator on the Xilinx Versal VCK5000, leveraging AI Engine (AIE) tiles for high-performance matrix operations and programmable logic for data movement. Benchmark its latency, throughput, and energy efficiency against GPU- and CPU-based implementations using existing frameworks, without building them from scratch.	Suresh Purini
BitNet 1.58 on Xilinx Alveo U50	Implement the BitNet 1.58 large language model with ternary weight quantization on the Xilinx Alveo U50 FPGA. Focus on optimizing compute and memory efficiency to enable scalable inference. Performance will be evaluated in terms of latency, throughput, and energy efficiency, with comparisons to CPU and GPU baselines. Refer the latest paper TerEffic: Highly Efficient Ternary LLM Inference on FPGA .	Suresh Purini
BitNet 1.58 on Xilinx Versal VCK5000	Deploy the BitNet 1.58 large language model with ternary weight quantization on the Xilinx Versal VCK5000, leveraging AI Engine (AIE) tiles for high-performance matrix operations and programmable logic for efficient dataflow. The design will be benchmarked for latency, throughput, and energy efficiency, and compared against CPU and GPU implementations. Refer the latest paper TerEffic: Highly Efficient Ternary LLM Inference on FPGA .	Suresh Purini
BitNet 1.58 on AMD NPU	Implement the BitNet 1.58 large language model with ternary weight quantization on the AMD NPU platform, leveraging its specialized matrix engines for efficient inference. Optimize for latency, throughput, and energy efficiency, and compare against FPGA (U50, VCK5000), GPU, and CPU implementations. Refer the latest paper TerEffic: Highly Efficient Ternary LLM Inference on FPGA .	Suresh Purini

Project Name	Brief Project Description	Primary Advisor
Large Language Processor (LLP)- all digital	Develop a streamlined all-digital accelerator hardware, custom instruction and dataflow for high-speed, long-context , mixture-of-experts LLM inference. Platform(s): gem5, firesim, hdl/System C Ref: LPU: A Latency-Optimized and Highly Scalable Processor for Large Language Model Inference	Priyesh Shukla
Large Language Processor (LLP) - IMC	Develop a streamlined in-memory computing -based accelerator hardware, custom instruction set and LLM-specific dataflow for high-speed, long-context LLM inference. Platforms: gem5, hspice/cadence, hdl/system C Ref: LPU: A Latency-Optimized and Highly Scalable Processor for Large Language Model Inference	Priyesh Shukla
Vision AI accelerator for AR/VR	Develop a generic DNN hardware prototype to accelerate low-power computations in 3D vision pipeline for AR/VR glasses. Refer to CICC2022 paper from Meta Reality Labs: System-Level Design and Integration of a Prototype AR/VR Hardware Featuring a Custom Low-Power DNN Accelerator Chip in 7nm Technology for Codec Avatars	Priyesh Shukla
Visual SLAM AI accelerator for Micro-robotics	ASIC/FPGA/EdgeAI/RISC-V implementation of deep learning based SLAM algorithm for low-power, accurate and on-device navigation of tiny drones. Ref: NanoSLAM: Enabling Fully Onboard SLAM for Tiny Robots	Priyesh Shukla
On-device training of DNNs	Design accelerator for backpropagation in Convolutional and fully-connected layers to carry out incremental training locally at the edge Ref: DARKSIDE: A Heterogeneous RISC-V Compute Cluster for Extreme-Edge On-Chip DNN Inference and Training Mandheling: Mixed-Precision On-Device DNN Training with DSP Offloading	Priyesh Shukla
Deep learning operator library and RISC-V custom extensions	Design multi-specification (low power, high speed, varying precision...) library of deep learning operators - 2D/3D convolution (point-wise, depthwise separable...), Self-attention, Multi-head attention for Transformers	Priyesh Shukla

Project Name	Brief Project Description	Primary Advisor
KV Cache Optimization	Investigate various algorithmic/deep learning architecture (CNNs/LLMs) optimization schemes for efficient processing at the edge Refs: DEEP COMPRESSION: COMPRESSING DEEP NEURAL NETWORKS WITH PRUNING, TRAINED QUANTIZATION AND HUFFMAN CODING	Girish Varma
Neural Architecture Search	An Efficient Hybrid Deep Learning Accelerator for Compact and Heterogeneous CNNs DeepSeek-V3 Technical Report	
Quantization aware training and inference	EIE: Efficient Inference Engine on Compressed Deep Neural Network Integer quantization for deep learning inference: Principles and empirical evaluation	
Pruning and Compression Schemes		

Design Space Exploration of GEMM on Alveo U50 (FP32 & INT8)	Implement and optimize GEMM kernels on the Xilinx/AMD Alveo U50, exploring FP32 (via floating-point IP) and INT8 (INT32 accumulation) variants.	Suresh Purini	Chipyard Develop RISC-V Co-processor (GEMM accelerator) in Chipyard to accelerate DNN training and inference Ref: https://chipyard.readthedocs.io/en/stable/index.html	Priyesh Shukla
Design Space Exploration of GEMM on Versal VCK5000 (FP32 & INT8)	Implement and optimize GEMM kernels on the Xilinx Versal VCK5000, exploring FP32 and INT8 variants. The project leverages AI Engine (AIE) tiles for high-throughput matrix operations and the programmable logic (PL) for memory orchestration and data movement.	Suresh Purini	Chipyard Develop In-memory computing macro as RISC-V co-processor in Chipyard for DNN training and inference Ref: https://chipyard.readthedocs.io/en/stable/index.html	Priyesh Shukla
Sparse Matrix-Dense Matrix Multiplication (SpMM) Accelerator on Alveo U50	Design and implement a hardware accelerator for Sparse Matrix-Dense Matrix Multiplication (SpMM) on the Xilinx Alveo U50 FPGA. Explore different sparse storage formats (CSR, CSC, Blocked-CSR) and scheduling strategies to optimize compute efficiency and memory bandwidth utilization. Evaluate performance in terms of latency, throughput, and energy efficiency, and compare against optimized CPU (MKL) and GPU (cuSPARSE) implementations. Refer FPGA-Based Sparse Matrix Multiplication Accelerators: From State-of-the-Art to Future Opportunities	Suresh Purini	In-memory computing ADC-free IMC macro for DNN training and inference capable of performing analog computations and digitization by exploiting memory bank characteristics Ref: HCIIM: ADC-Less Hybrid Analog-Digital Compute in Memory Accelerator for Deep Learning Workloads A Brain-Inspired ADC-Free SRAM-Based In-Memory Computing Macro With High-Precision MAC for AI Application	Priyesh Shukla
Sparse Matrix-Dense Matrix Multiplication (SpMM) Accelerator on Versal VCK5000	Develop a SpMM accelerator on the Xilinx Versal VCK5000, leveraging AI Engine (AIE) tiles for dense matrix computation and programmable logic for sparse indexing and dataflow control. Investigate design choices such as tiling, load balancing across sparse rows, and hybrid AIE-PL partitioning. Benchmark latency, throughput, and energy efficiency, and compare against CPU and GPU baselines to highlight the benefits of Versal's heterogeneous architecture. Refer FPGA-Based Sparse Matrix Multiplication Accelerators: From State-of-the-Art to Future Opportunities	Suresh Purini	High-bandwidth processing-in-memory (HBM-PIM) simulation and characterization Gem5 simulation framework for HBM-PIM simulation and characterization for deep learning inference Ref: https://www.gem5.org/assets/files/hpca2023-tutorial/gem5-tutorial-hpca-2023.pdf Gem5-X: A Many-core Heterogeneous Simulation Platform for Architectural Exploration and Optimization https://www.gem5.org/documentation/general_docs/stdlib_api/gem5.components.memory.hbm.html	Priyesh Shukla
Graph Neural Network (GNN) Accelerator on Alveo U50	Design and implement a GNN inference accelerator on the Xilinx Alveo U50 FPGA. The focus is on optimizing the sparse-dense matrix multiplication and neighbor aggregation kernels that dominate GNN workloads. Explore different sparse formats (CSR, COO) and tiling strategies to balance memory bandwidth and compute utilization. Evaluate the accelerator on standard GNN models (e.g., GCN, GraphSAGE) and benchmark performance against CPU and GPU baselines. Refer A Review of FPGA-based Graph Neural Network Accelerator Architectures	Suresh Purini	Bitnet (1-bit LLM) orchestration on multi-core RISC-V accelerator-vector cores Enhance PULP's ARA multicore RISC-V architecture with custom accelerator extension and implement 1-bit LLM on the RISC-V architecture Ref: Ara2: Exploring Single- and Multi-Core Vector Processing With an Efficient RVV 1.0 Compliant Open-Source Processor	Priyesh Shukla

Graph Neural Network (GNN) Accelerator on Versal VCK5000	Develop a GNN inference accelerator on the Xilinx Versal VCK5000, leveraging AI Engine (AIE) tiles for dense feature transformations and programmable logic for sparse neighbor aggregation. Investigate partitioning strategies across AIE and PL, dataflow optimizations, and load balancing for irregular graph structures. Benchmark throughput, latency, and energy efficiency using real graph datasets, and compare against state-of-the-art CPU and GPU implementations. Refer A Review of FPGA-based Graph Neural Network Accelerator Architectures	Suresh Purini
---	--	---------------

Brain-inspired (spiking) neural network on FPGA/IMC	Implement neuromorphic computing accelerator on Xilinx Alveo exploiting spike time-dependent elasticity with higher energy efficiency Refs: DeepFire: A Convolutional Spiking Neural Network Accelerator on FPGAs IzhIRISC-V - a RISC-V-based Processor with Custom ISA Extension for Spiking Neuron Networks Processing with Izhikevich Neurons	Priyesh Shukla
Photonic neural accelerator	Implement a simulation framework for photonic computations to accelerate CNNs Ref: LightML: A Photonic Accelerator for Efficient General Purpose Machine Learning	Priyesh Shukla
NVM (non-volatile memory) crossbars	Explore NVM technology (PCM, ReRAM, MRAM...) to for analog in-memory computing to accelerate CNNs and transformers Ref: Memory devices and applications for in-memory computing	Priyesh Shukla
Chiplets - Wafer-scale computing	Explore chiplets - heterogenous integration, partitioning, scheduling, network routing strategies for efficient LLM/LLM inference Ref: WSC-LLM: Efficient LLM Service and Architecture Co-exploration for Wafer-scale Chips	Priyesh Shukla
KV Cache management for LLMs	Explore strategies to reduce redundant computations and improve memory utilization during long-context real-time LLM inferencing Ref: A Survey on Large Language Model Acceleration based on KV Cache Management	Priyesh Shukla
Toolchain optimization passes for CNNs and LLMs	Compilers (MLIR, IREE, LLVM, GCC)/OS - optimizations for new accelerator co-processor Refs: An MLIR-based Compiler Flow for System-Level Design and Hardware Acceleration https://github.com/iree-org/iree	Priyesh Shukla