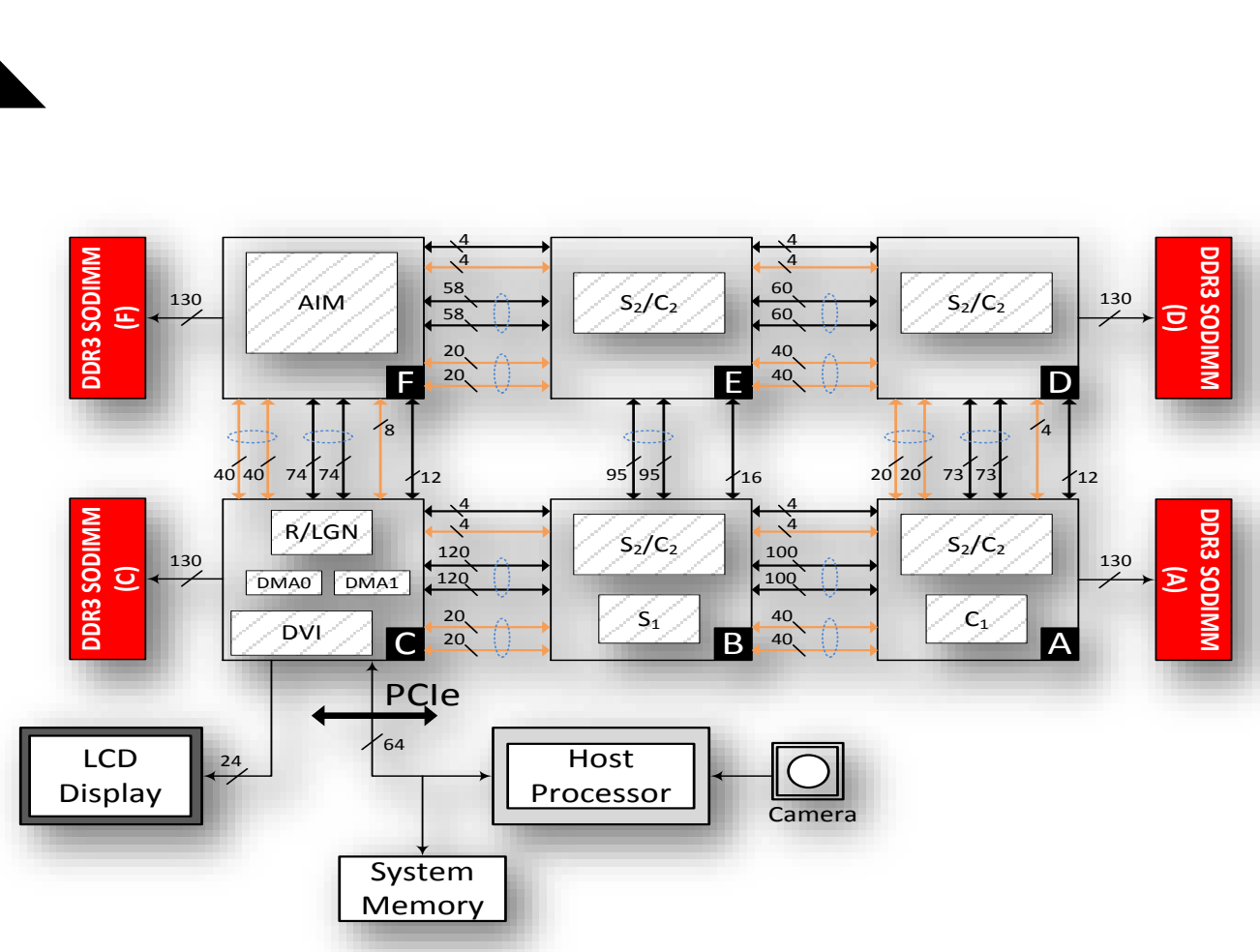# Heterogeneous Architectures for Intelligent Vision Systems
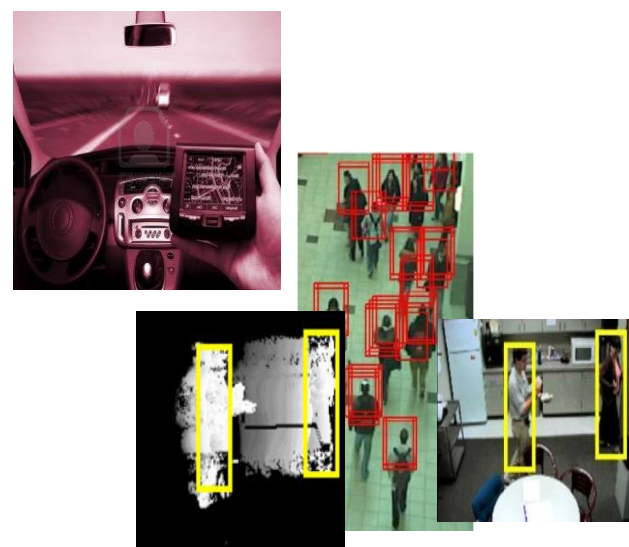
Microsystems Design Lab (MDL), Department of Computer Science & Engineering,
The Pennsylvania State University, University Park, PA 16802
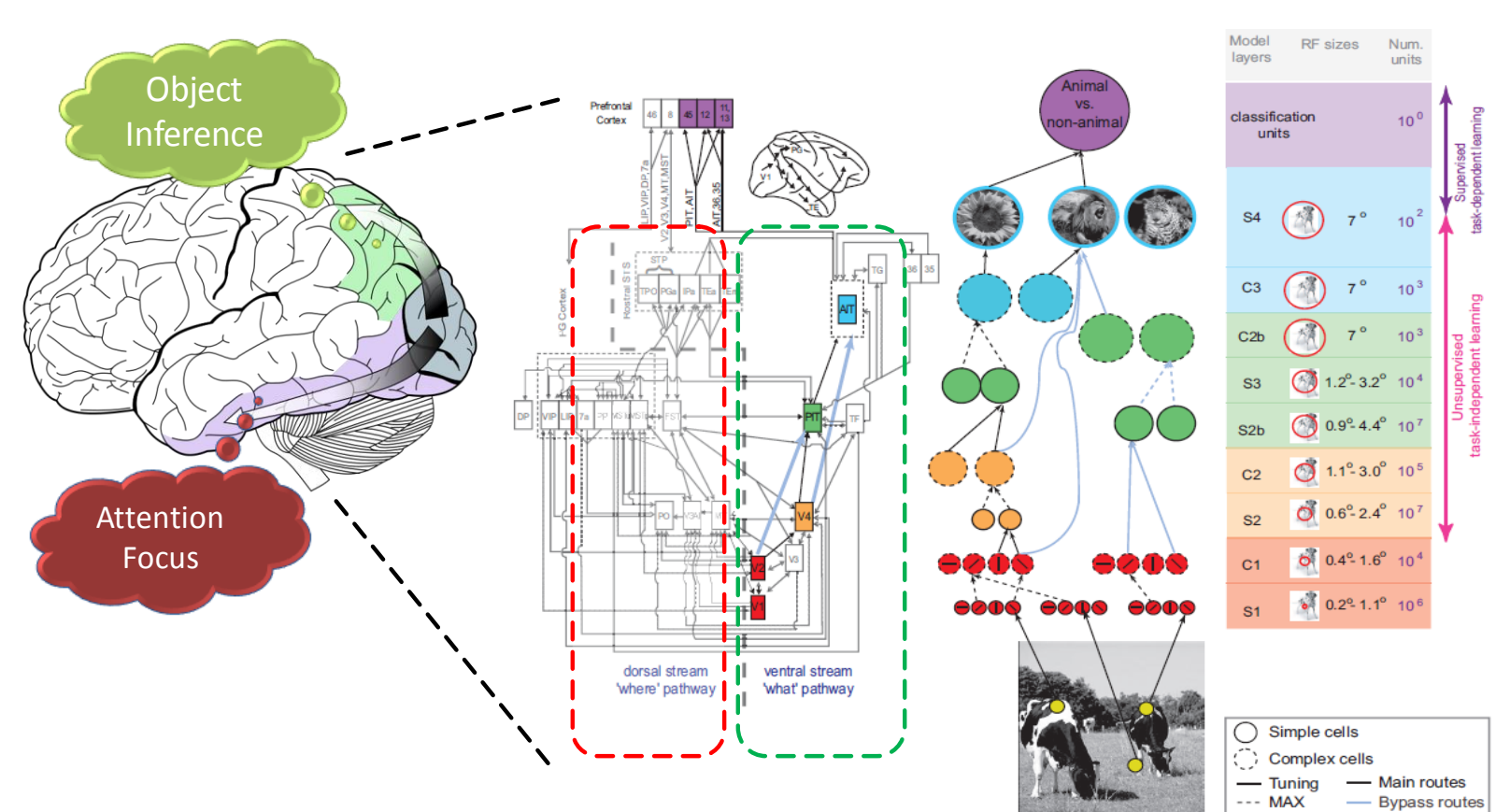
## Biologically-inspired Vision

**Perceptual systems as efficiently as humans**
Recognition Accuracy – 1000s of classes
Performance –milliseconds
Power – < 20 Watts

**Leverage knowledge of Visual Cortex**
Retinal Enhancement
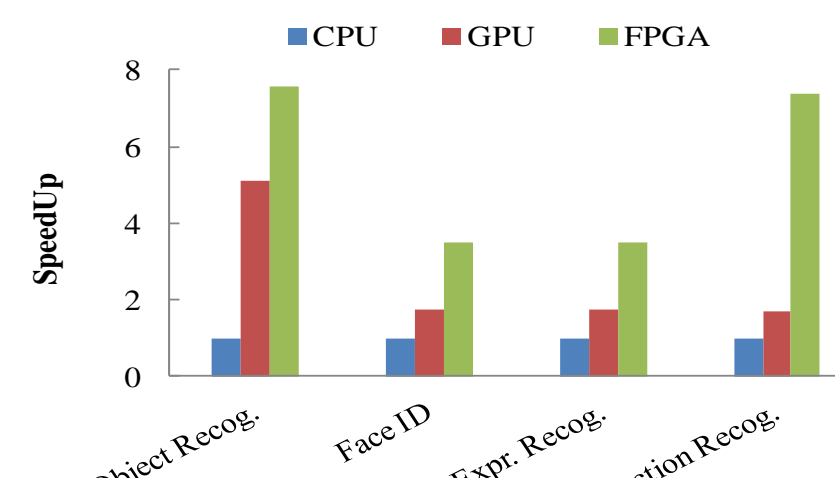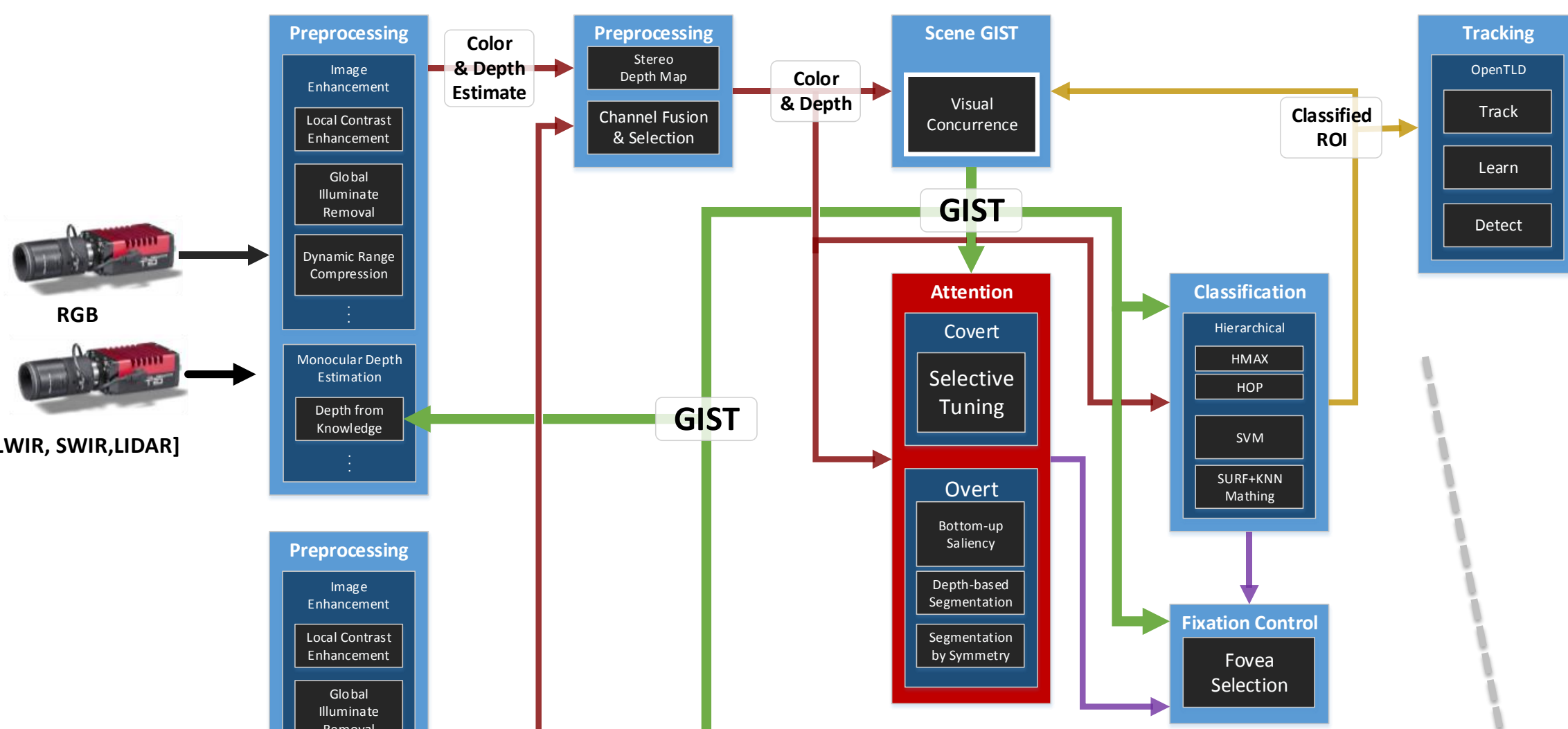Visual Attention
Inference

Hardware Architecture



**Diverse Application Domains**
Vehicle Navigation Assist
First Person Analytics
Biometrics

Object Inference

Attention Focus

Neuromorphic Vision Models

## Visual Perception Pipeline on a SoC



Preprocessing
- Image Enhancement
- Local Contrast Enhancement
- Global Illuminate Removal
- Dynamic Range Compression
- Monocular Depth Estimation
- Depth from Knowledge

Color & Depth Estimate

Preprocessing
- Stereo Depth Map
- Channel Fusion & Selection

Color & Depth

Scene GIST
- Visual Concurrence

GIST

Classified ROI

Tracking
- OpenTLD
- Track
- Learn
- Detect

Attention
- Covert: Selective Tuning
- Overt: Bottom-up Saliency
- Depth-based Segmentation
- Segmentation by Symmetry

Classification
- Hierarchical
- HMAX
- HOP
- SVM
- SURF+KNN Matching

Fixation Control
- Fovea Selection

RGB

[LWIR, SWIR,LIDAR]

Memory

LW Core — Inst Mem (×4), CPU, I/D Cache

uAccel (type 1), Macro Accelerator, uAccel (type 1)

uAccel (type 1), uAccel (type 2), uAccel (type 2)

uAccel (type 3), uAccel (type 3)

Macro Accelerator

Speedup: up to 7.6X (4.3) compared to CPU (GPU)

Legend: CPU, GPU, FPGA — Speed'p — Object Recog., Face ID, Facial Expr. Recog., Action Recog.

## Visual Perception Platform – Vortex Network Infrasturcture

- ✓ Accelerators as shared services
- ✓ Native support for dataflow processing
- ✓ Support for low latency control-flow
- ✓ Accelerator composition and configurability

S — C-NIF, S-NIF — On-chip Router — S-NIF — C, S-NIF, S, M-NIF, M

**C – Switch Attached Processor (Core)**
**S – Stream Operator (accelerator)**
**M – Memory**
**NIF – Network Interface**

- ✓ Integrated network awareness using data flows
- ✓ Flow = sequence of operations which can be mapped to a node (core or SOP).
- ✓ Flow id represents a unique sequence of operations; Next hop of data flow at each node determined using flow id embedded in each packet by referencing a flow table
- ✓ Efficient direct streaming data transfer between nodes of a data flow thereby "cascading" operations
- ✓ Resource sharing between multiple data flows
- ✓ High aggregate bandwidth to ensure multiple data flows are processed concurrently
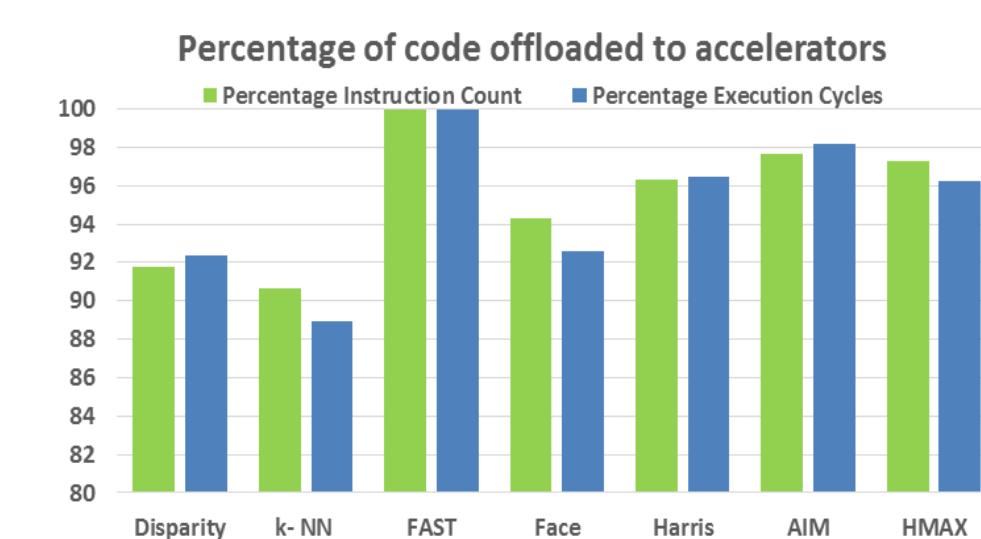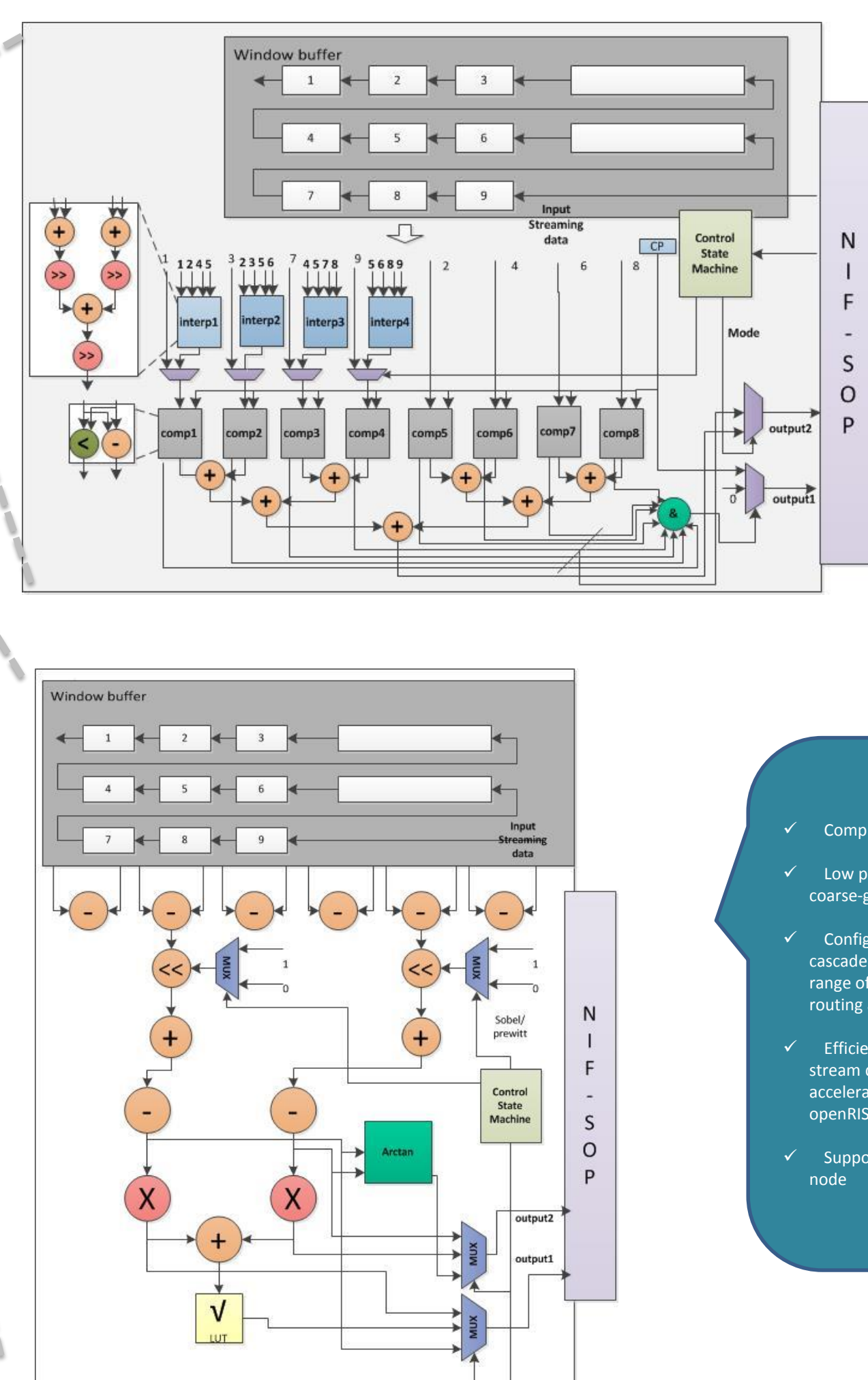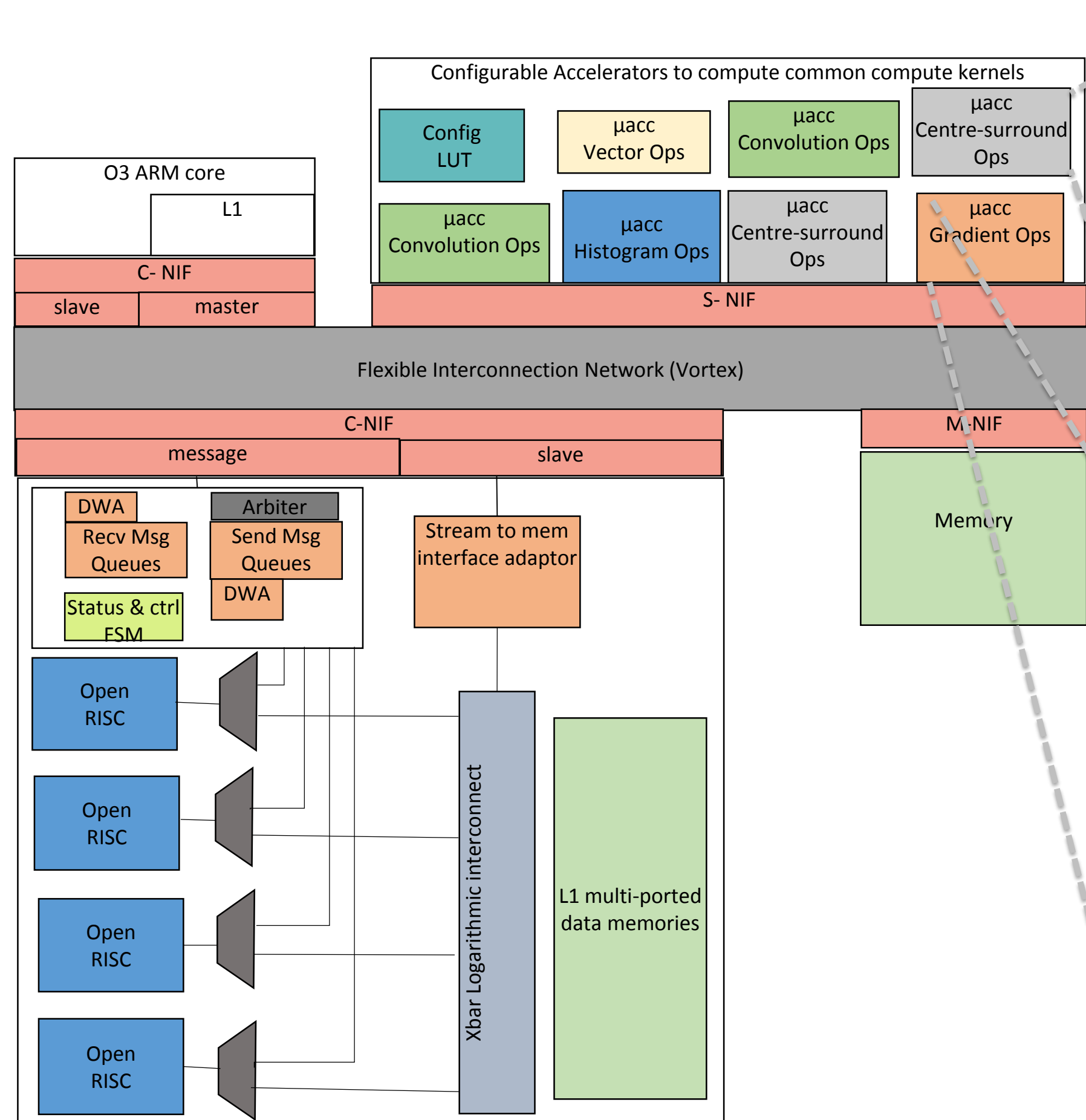


C- Network Interface

- ✓ Streaming Operators accelerate operations such as 2-d Convolution, Euclidean distance calculation
- ✓ On data arrival from network, Flow id is resolved to Opcode, which determines the mode of operation of accelerator and is forwarded to accelerator
- ✓ Accelerator output data is packetized and Next Hop is resolved based on Flow id
- ✓ Header is updated and data injected into network

SOP- Network Interface

- ✓ Core initiates R/W transactions through master interface's FIFO-like handshaking or DMA-like transactions between local memory space and remote address space accessible through slave interface
- ✓ Message interface allows to send light weight messages to communicate with other nodes.
- ✓ Multiple outstanding non-blocking operations managed by NIF, each with a different flow-id

## Heterogeneous Multi-core Accelerator for Visual Perception (HMAP)

Configurable Accelerators to compute common compute kernels
- Config LUT
- µacc Vector Ops
- µacc Convolution Ops
- µacc Centre-surround Ops
- µacc Convolution Ops
- µacc Histogram Ops
- µacc Centre-surround Ops
- µacc Gradient Ops

O3 ARM core — L1 — C- NIF — slave, master — S- NIF

Flexible Interconnection Network (Vortex)

C-NIF — message, slave — M- NIF

DWA Recv Msg Queues, Arbiter Send Msg Queues, Status & ctrl FSM, DWA

Memory

Open RISC (×4)

Xbar Logarithmic interconnect

L1 multi-ported data memories

Window buffer

Percentage of code offloaded to accelerators
Legend: Percentage Instruction Count, Percentage Execution Cycles
Categories: Disparity Map, k-NN, FAST, Face Recognition, Harris Corner, AIM saliency, HMAX

- ✓ Complex cores to initiate transactions
- ✓ Low power OpenRISC core cluster to extract coarse-grain parallelism
- ✓ Configurable micro-accelerators that can be cascaded to compute common kernels from a range of algorithms with custom storage, routing and DMA memory access
- ✓ Efficient data and control flow that facilitates stream direct data transfers between accelerators that offload computation from openRISC core cluster or complex core
- ✓ Support for multiple transactions at each node

- ✓ Can we achieve performance gain and energy savings comparable to that of an ASIC?
- ✓ What are sources of performance gain and energy?
- ✓ How does Configurability compare with customization?

Estimated performance gain - Comparison of execution times normalized to that of HMAP

640*480, 1024*768, 1280*960

Legend:
- 2-core ARM C-A9 #threads=1 @ 1.5GHz
- 2-core ARM C-A9 #threadv=16 @ 1.5GHz
- Micro-accelerators @100 MHZ
- HMAP @ 500MHZ
- 8-core Xeon #threads=16 @ 2.67GHz