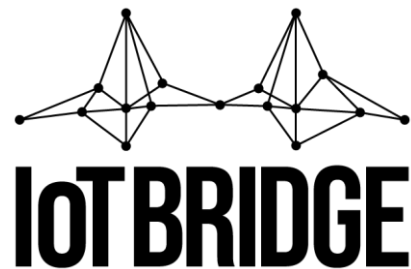


Smart and Secure Gateways for a Secure Internet of Things

The Cybernode Collaboration Conference 2023

Joakim Eriksson, RISE
2023-01-26



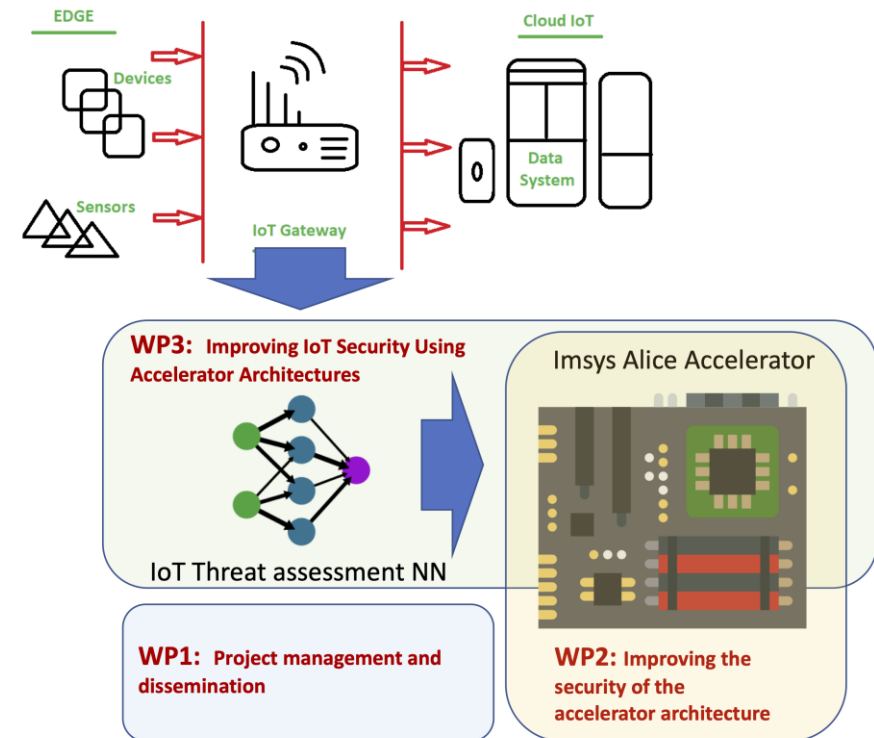
Smart and Secure IoT Gateway

- Partners

- RISE - project lead and use-case provider:
IoT – wireless device fingerprinting
- Uppsala University – research and development of security mechanisms for the AI accelerator
- Imsys – design and implementation of secure AI accelerator
- IoT Bridge – IoT company – use-case provider:
Bridge Safety IoT Application Using AcceleratorArchitecture
- Wittra – IoT company – use-case provides:
Secure IoT for asset tracking and asset lock system

- Main contribution

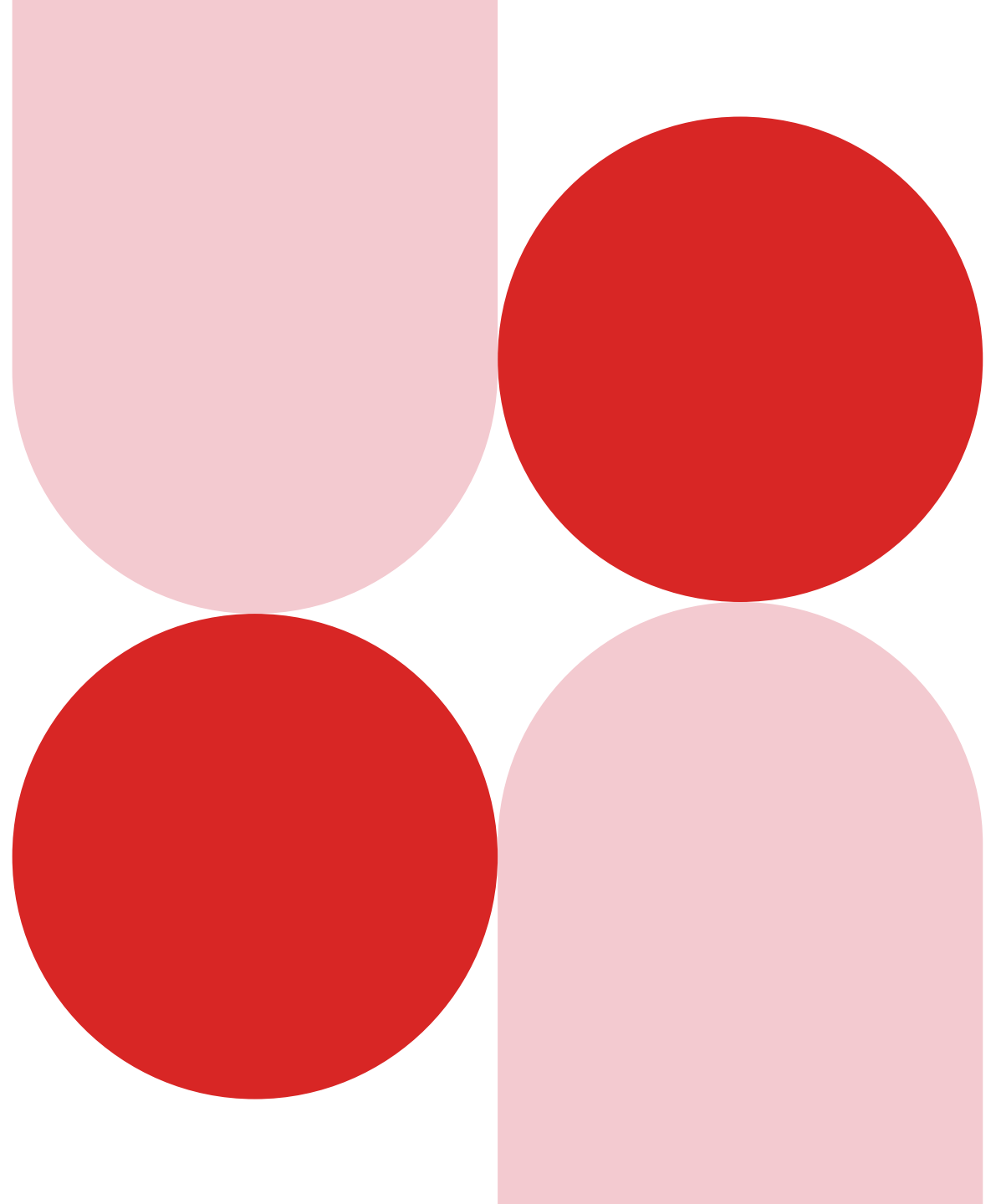
- Smart and Secure IoT Gateway based on Imsys Alice AI-accelerator
- Use-cases evaluating Alice



Acceleration in SecureGW

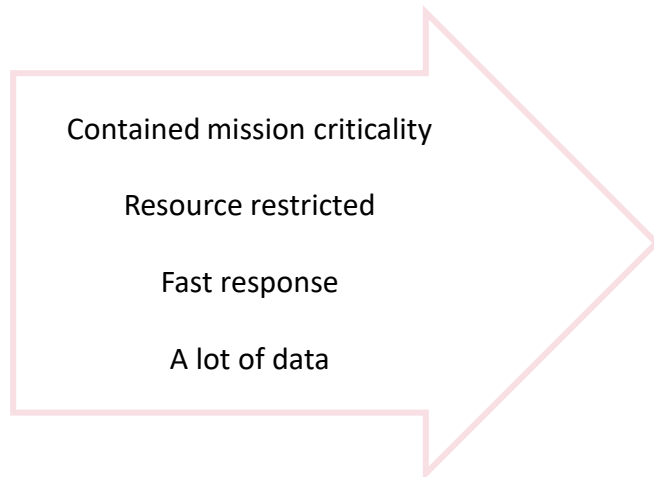
Dag Helmfrid, CTO

Imsys take on AI acceleration with power efficiency



Inference close to data

- High capacity important for fast response



Alice Accelerator Architecture

- **Accelerator controller, IM4000 (GPP)**

Full software stack
High level language support
Open and open source based

- **Network on Chip (NoC)**

High speed data and control
Application controlled peer2peer

- **Processing Element Cluster (PEC)**

Multiple Processing Elements with shared memory
Processing near memory in PE

- **I/O**

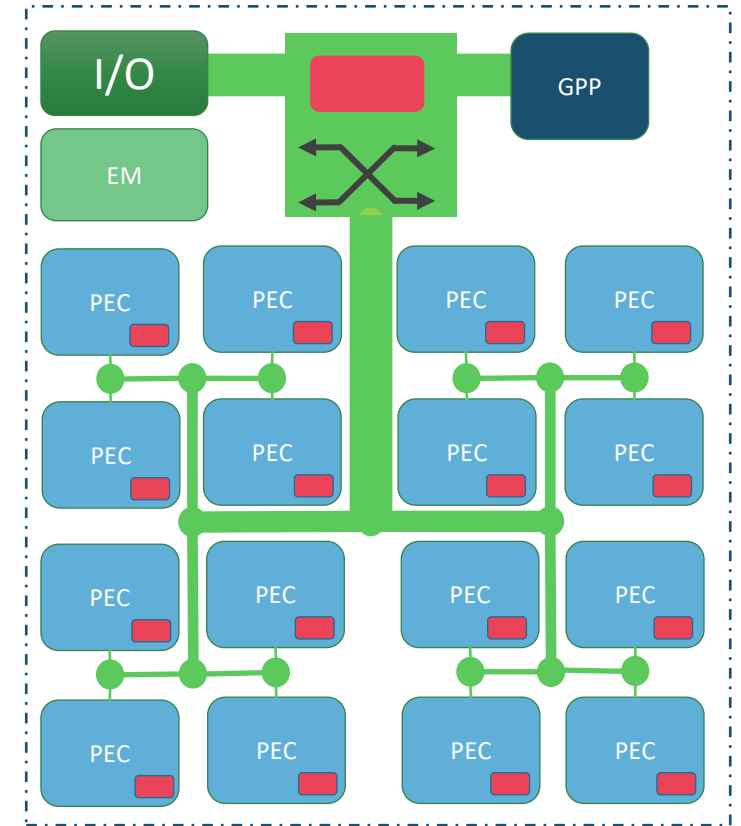
External memory, Highspeed i/f, Ethernet ...

- **Energy manager (EM)**

Sleep modes, Performance, External power source



Flexible
Deep
Neural
Network
Instruction
Set

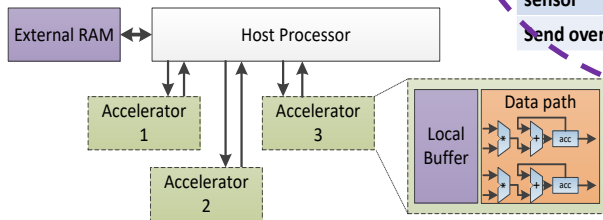


The Imsys Accelerator Design for Low energy

Accelerator Challenges

- Data Movement: Get parameters + activations from RAM
- Data movement is expensive
 - Energy, latency, bandwidth
 - You need data to compute
- Focus on data locality

Action	Energy	Relative
ALU op	1 pJ – 4 pJ	1x
SRAM Read	5 pJ – 20 pJ	5x
Move 10mm across chip	26 pJ – 44 pJ	25x
Send to DRAM	200 pJ – 800 pJ	200x
Read from image sensor	3.2 nJ – 4 nJ	4,000x
Send over LTE	50 uJ – 600 uJ	50,000,000x



DSD 2018 AMDL Keynote, Prof. Dr. Henk Corporaal

Sources of energy consumption challenging the system solution



Don't move data around

- Automated tools for data flow analysis
- Cache less memory access
- Data "reuse"
- Processing near memory

Efficient processing

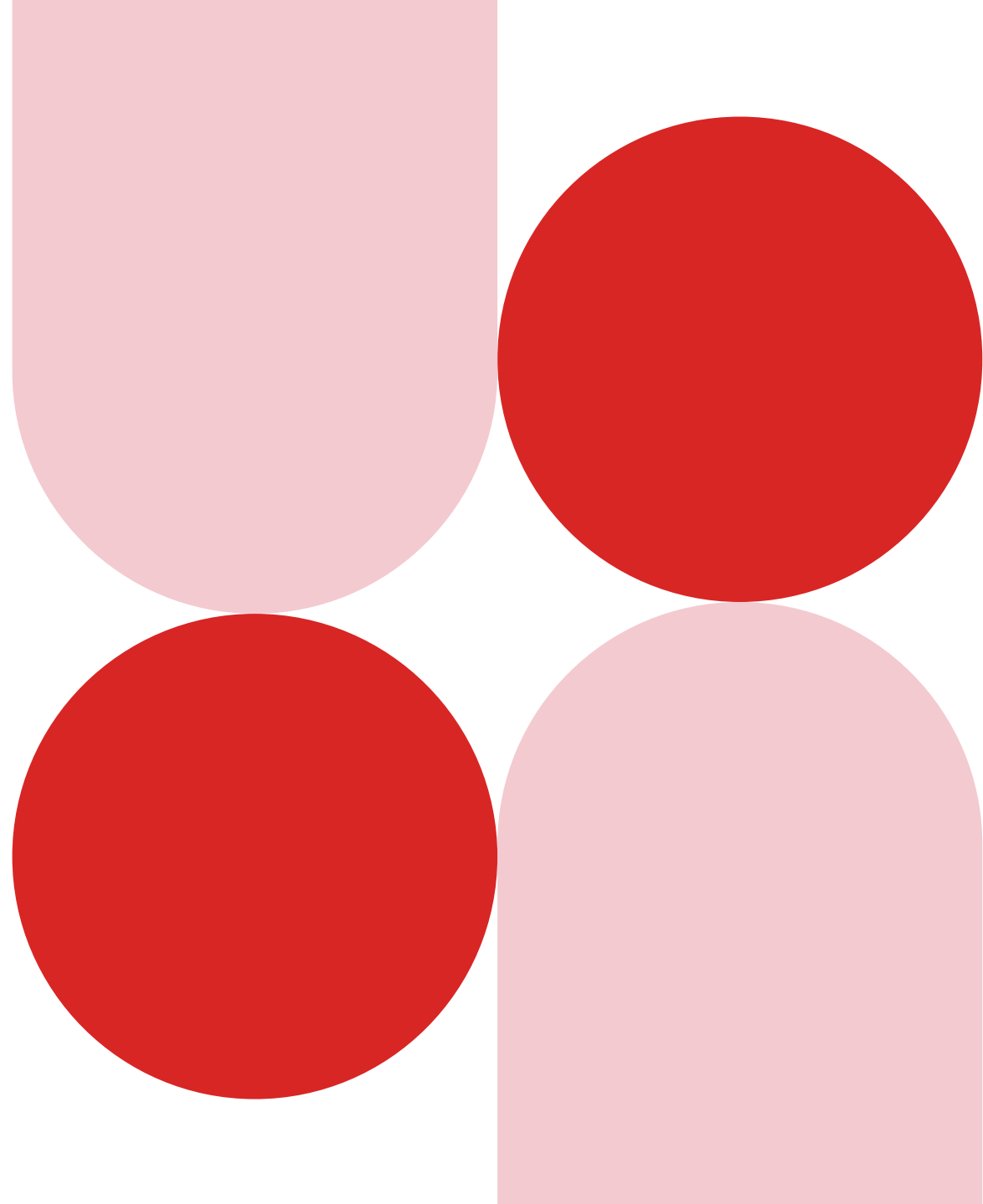
- Lean data types (int8, BFloat16)
- Low power circuit design matching architecture and advanced technology nodes for SoC
- Only local data transport

Important project results: Models have been quantized for energy efficient execution. Same model precision to less than a quarter of the energy Imsys accelerator's vector engines and tools have been updated.

Automated network optimization and programmed flexibility

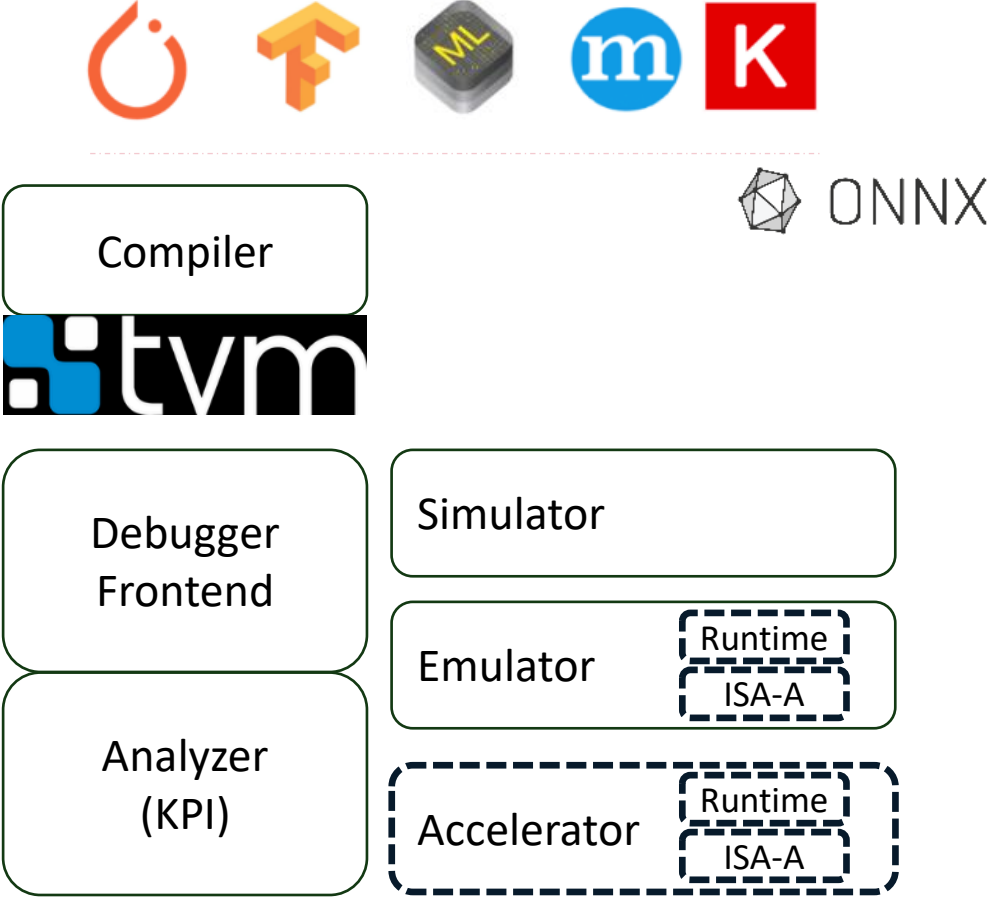
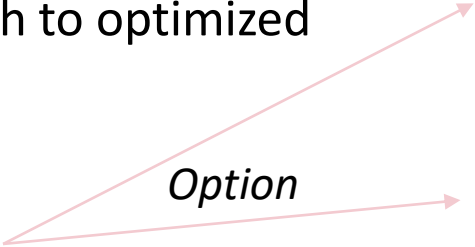
- Minimize memory usage & maximize utilization
- Layer fusion, zero pruning, operator fusion ...

SecureGW demonstrator platform



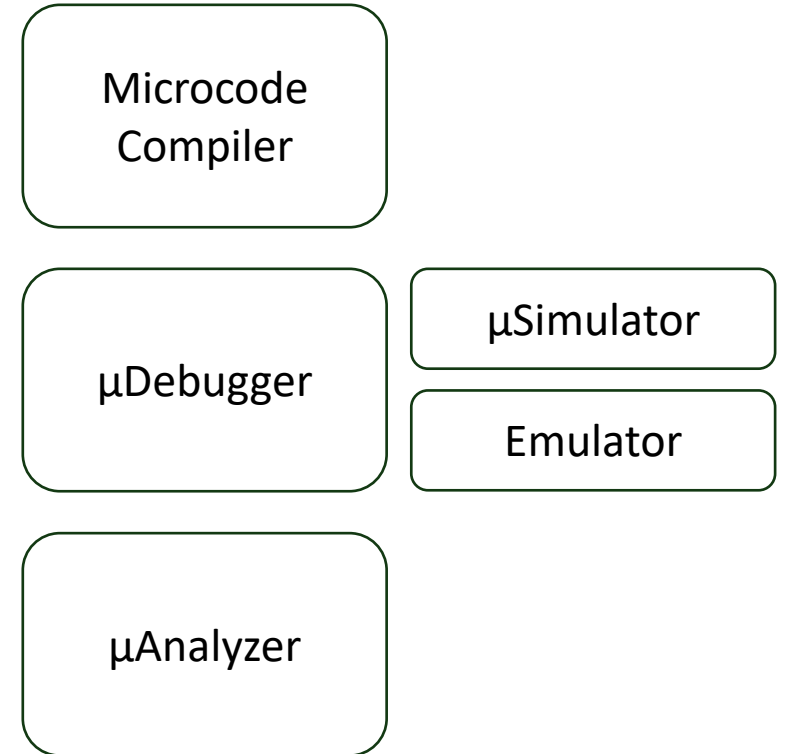
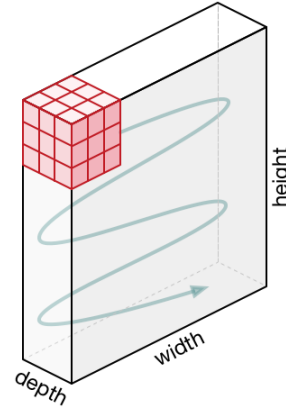
Optimizer, Compiler & Runtime

- Supports development flow from inference application graph to optimized target object code
- Quantization support
- Customizable optimization: pipelining, layer fusing ...
- Seamlessly integrates with existing AI development frameworks



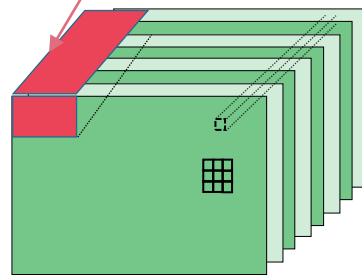
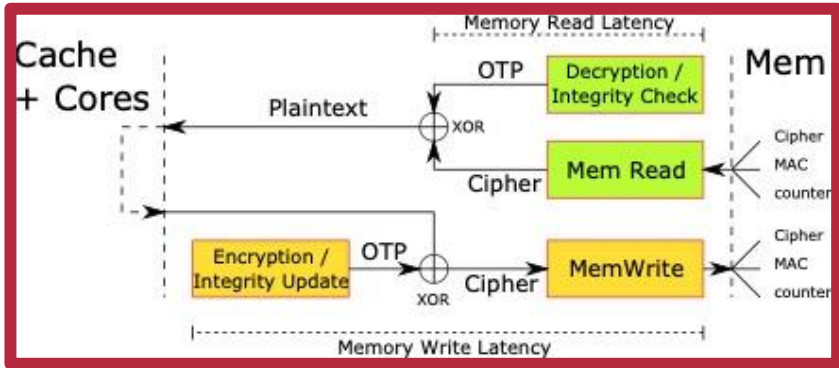
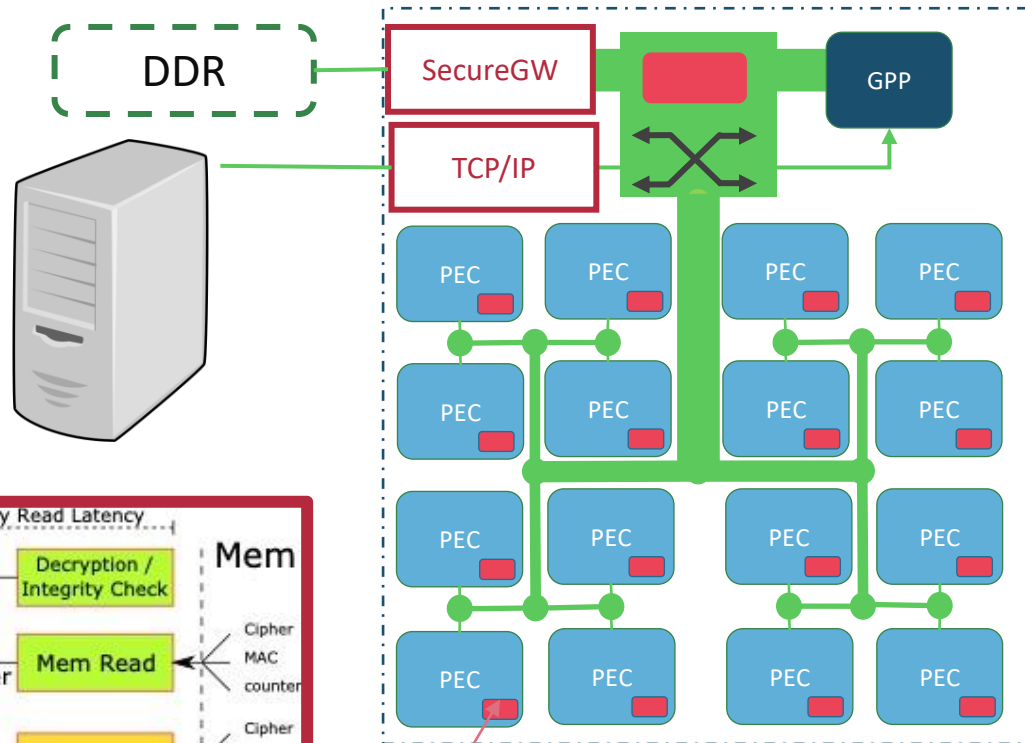
Kernel Library (DNN Instruction set, ISA-A) & SDK

- Kernel libraries support
 - Extensive instructions for quantized neural network operations and other kernel-based operations
- Programmable user customization
- Tools for custom kernel development



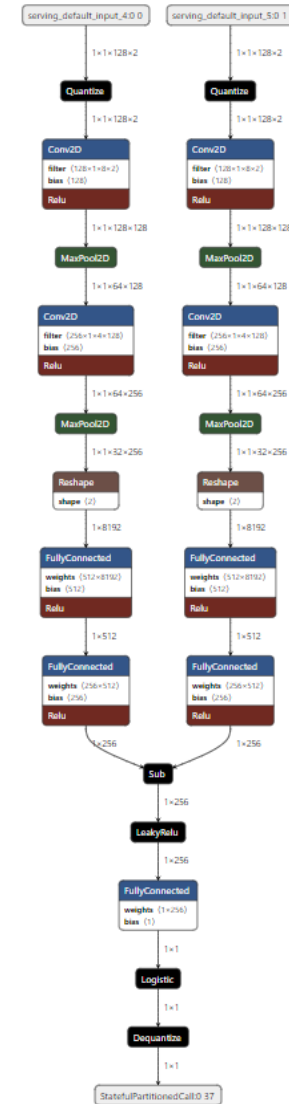
“SecureGW project has triggered new optimizations and kernel library extensions”

SecureGW application demo



MAC:s dominate computational load

- Point-wise convolution
- Depth-wise convolution



Thank you.

imsys