# Deploying CNN based vision solutions on a $3 microcontroller

Greg Lytle

May 22, 2018

# Au-Zone Technologies Inc.

**Au-Zone Technologies** is a leading provider of development tools and enabling IP used for the design of intelligent embedded vision products and solutions.

Our architecture agnostic development tools (eCV SDK and DeepView) enable our customers quickly to develop and securely deploy machine learning solutions on a range of embedded hardware.
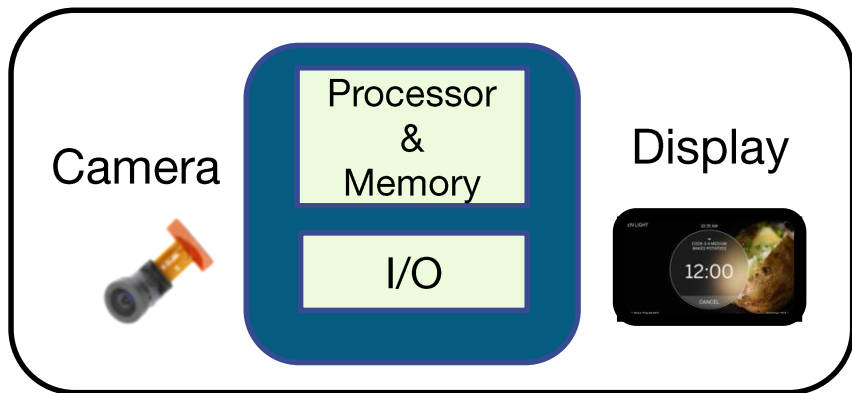
# Design Challenge

Requirements:

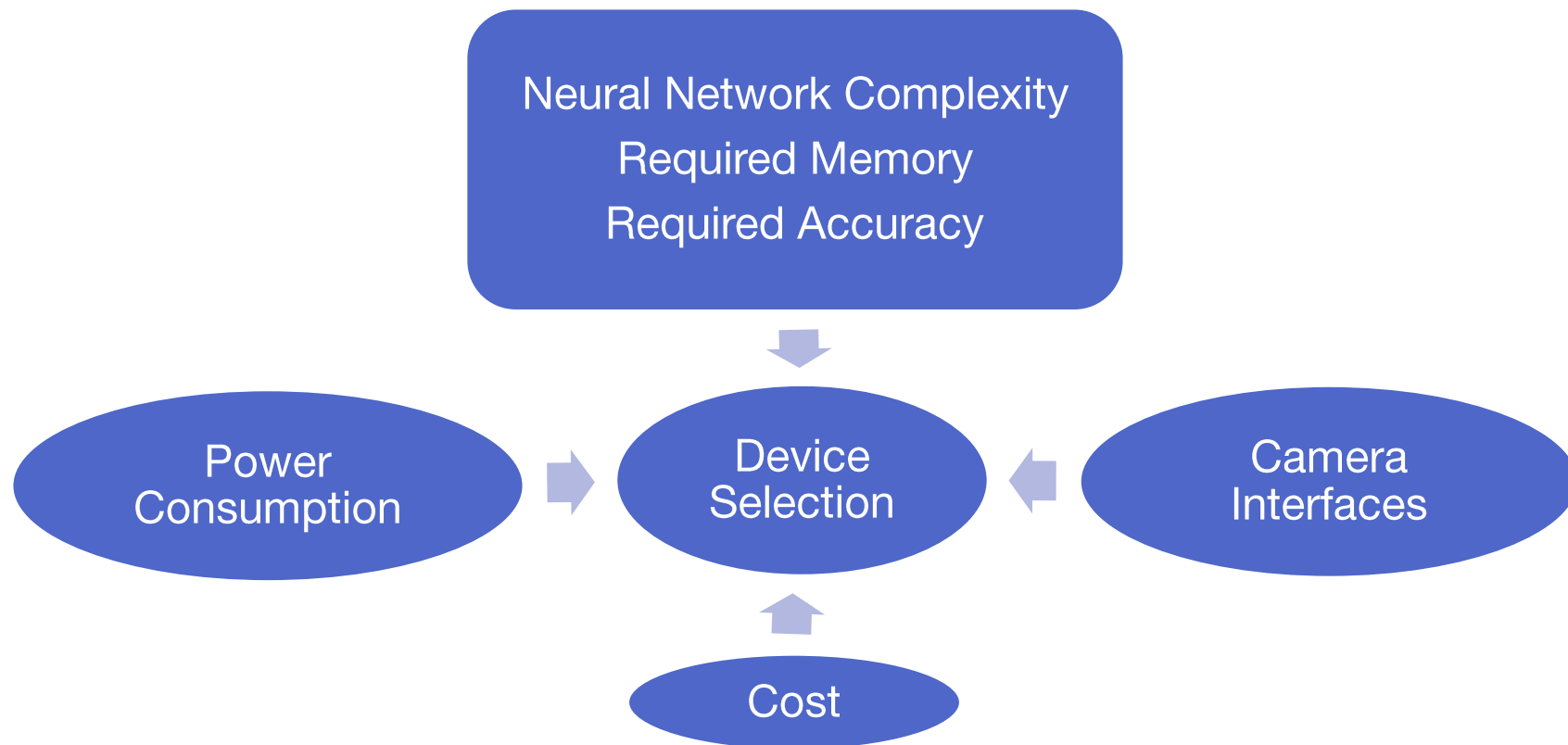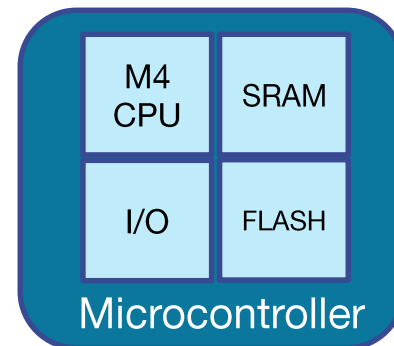- Embedded image capture and classification
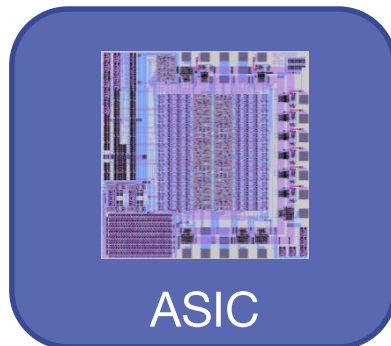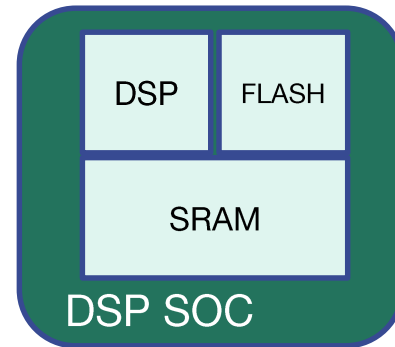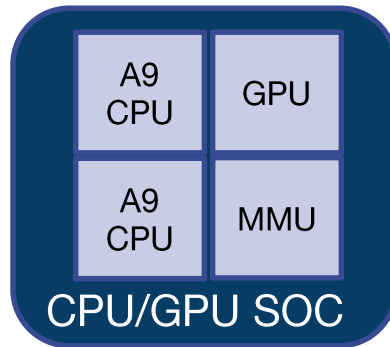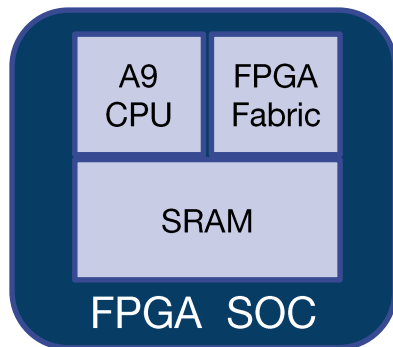
- Image Sensor

- LCD & Buttons

- Compute & Memory for vision algorithm and application

- Low BOM cost



Camera

Processor & Memory

I/O

Display

FPGA SOC

A9 CPU | FPGA Fabric
SRAM

CPU/GPU SOC

A9 CPU | GPU
A9 CPU | MMU

DSP SOC

DSP | FLASH
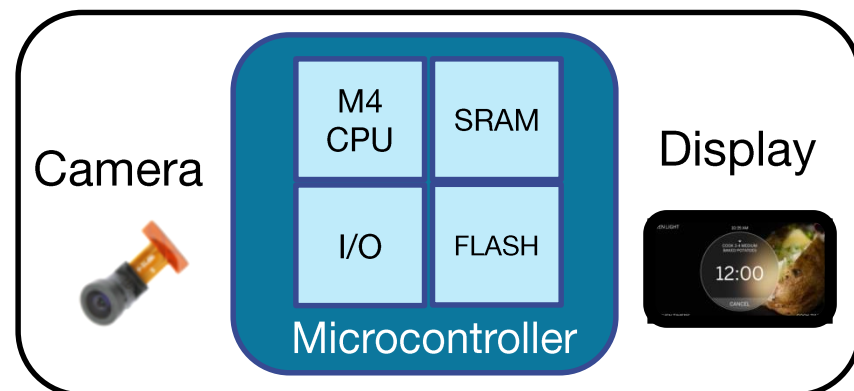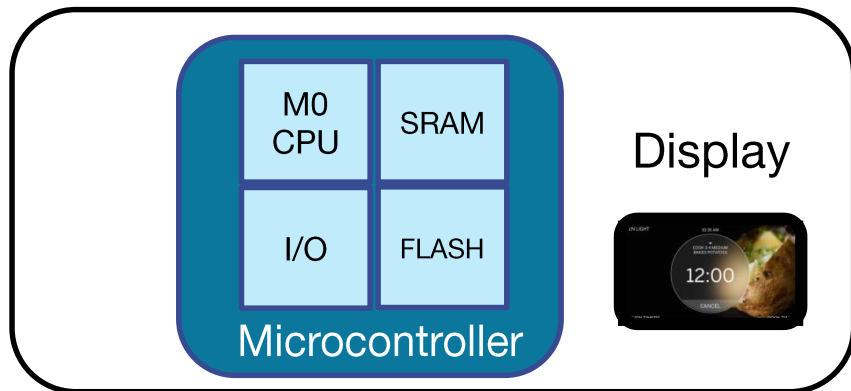SRAM

Many Options?

ASIC

Vision Processor

Microcontroller

M4 CPU | SRAM
I/O | FLASH

# Market Drivers

- Is this practical?
- Why should you care?

Image Classification enabled
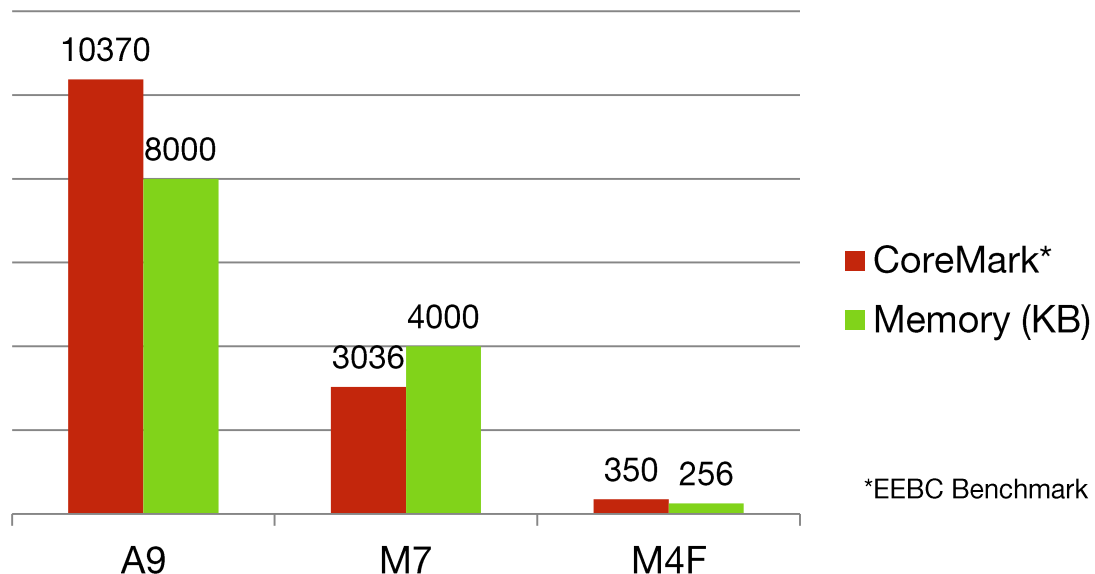


Low Incremental Cost

Much Lower raw
processor performance

Significant Memory
Constraints

## CPU Performance & Memory



- CoreMark*
- Memory (KB)

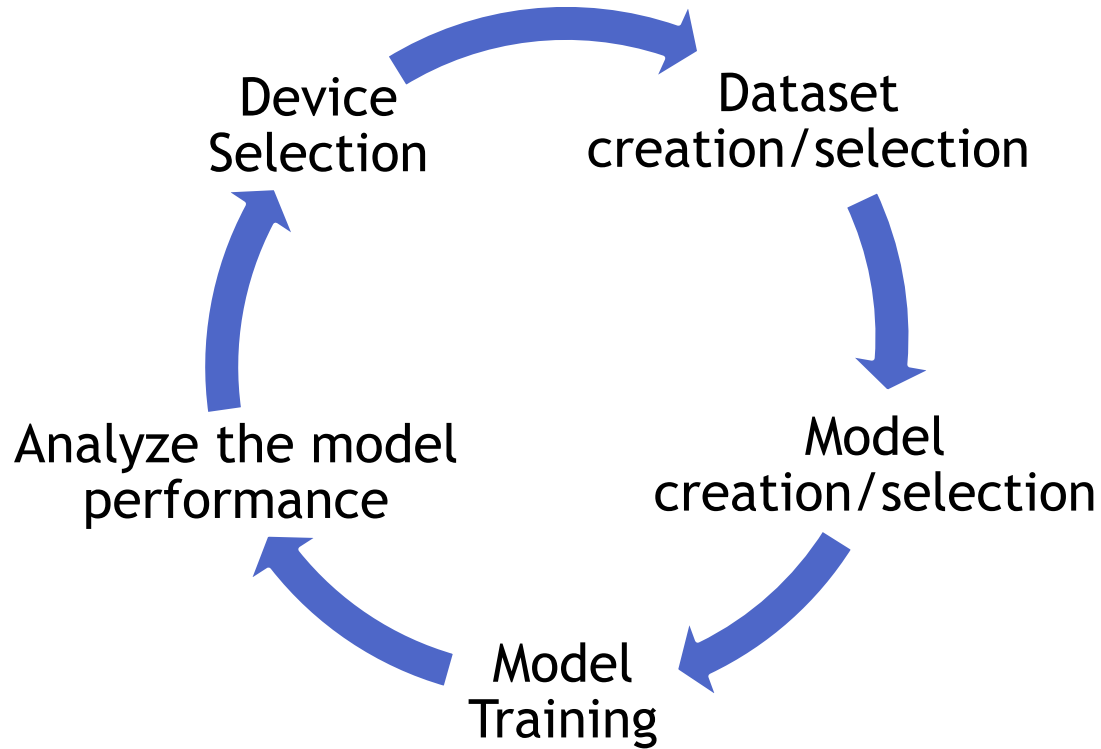*EEBC Benchmark

Chart data: A9 — CoreMark 10370, Memory 8000; M7 — CoreMark 3036, Memory 4000; M4F — CoreMark 350, Memory 256
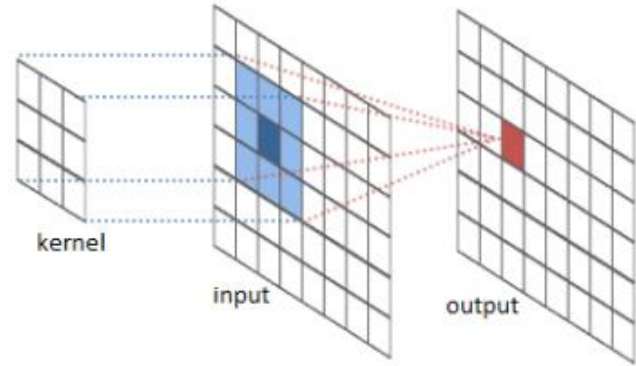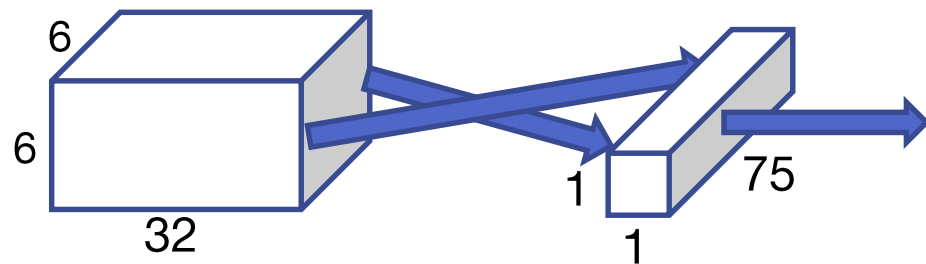
# Network Design

Fully Connected Layer

Convolutional  Layer

# Fully Connected Operation



One Weight required for every connection

8-bit quantization saves on FLASH size with low impact on error

| Parameter | Resource | Size | 8bit | 32 bit float | Typical |
|---|---|---|---|---|---|
| Input Buffer | RAM | 6 x 6 x 32 | 1.1 KB | 4.6 KB | 4.6 KB |
| Weights | FLASH | 6 x 6 x 32 x 1 x 1 x 75 | 86.4 KB | 345.6 KB | 86.4 KB |
| Output | RAM | 1 x 1 x 75 | .075 KB | 3.00 KB | 3 KB |
| Total Memory | | | 87.6 KB | 351 KB | 7.6 KB RAM 86.4 KB FLASH |
| MAC Operations | CPU | (6 x 6 x 32) x (1 x 1 x 75) | | 86.4 K | 86.4K |

# Convolution Operation



3x3 Convolution x 16 channels

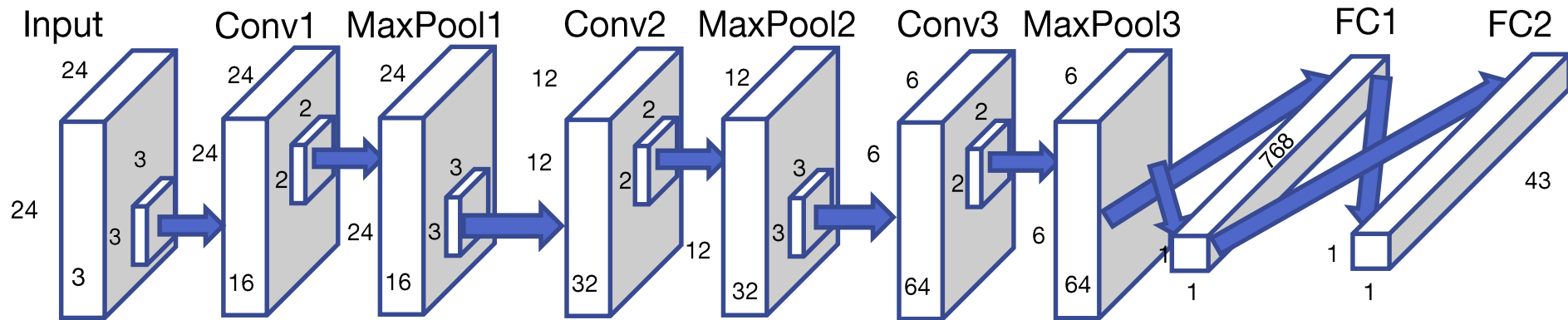| Parameter | Resource | Size | 8bit | 32 bit float | Typical |
|---|---|---|---|---|---|
| Input Buffer | RAM | 26 x 26 x 3.<br>(24x24x3 plus padding) | 2.1 KB | 8.1 KB | 8.1 KB |
| 3x3 Kernel | FLASH | 3 x 3 x 3 x 16 | 0.42 KB | 1.73 KB | 0.42 KB |
| Output | RAM | 24 x 24 x 16 | 9.2 KB | 36.8 KB | 36.8 KB |
| Total Mem | | | 11.6 KB | 46.7 KB | 0.42 KB FLASH<br>44.9 KB RAM |
| MAC Ops | CPU | (24 x 24 x 3) x (3 x 3 x 16) | | 248.8 K | |

# GTSRB –

- German Traffic Signs in the wild
- 39K       Training Samples
- 12.6K   Testing  Samples
- 43 classes of signs
- Cropped images 24 x 24 x 3

# CNN with Fully Connected Layers



| | | | | |
|---|---|---|---|---|
| Model FLASH | **2 MB** | MAC Operations | | ~2M |
| Minimum RAM | 74 KB | Weights | | 488K |
| Minimum Cache | 2 KB | Layers | | 9 |
| Optimum Cache | 74 KB | Accuracy | | 94.9 |

(example for 32-bit float implementation)

# Reduced complexity

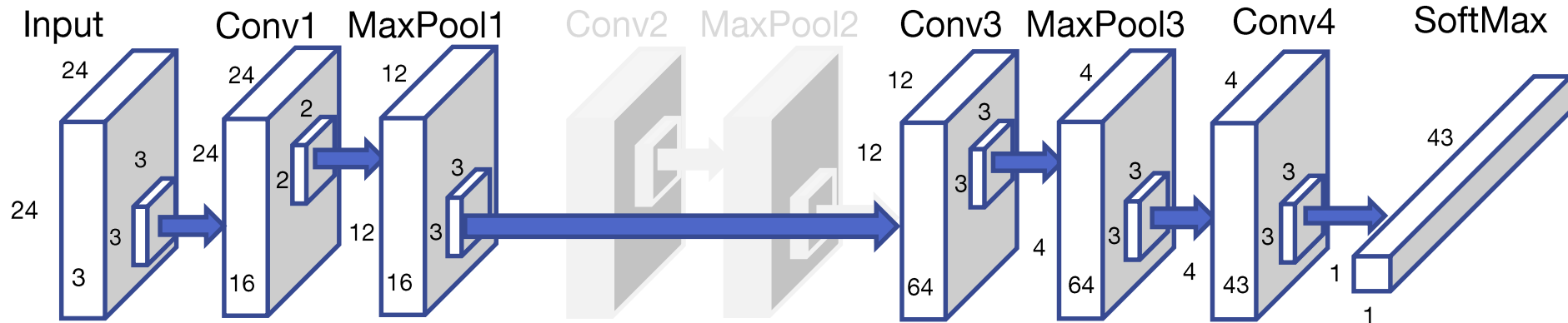| | | |
|---|---|---|
| Model  FLASH | 2 M → 195 KB | MAC Operations ~2M |
| Minimum  RAM | 74 KB | Weights 488K→ 48K |
| Minimum Cache | 2 KB | Layers 9 → 8 |
| Optimum  Cache | 74 KB | Accuracy 94.9 → 92.5% |

(example for 32-bit float implementation)

| | | |
|---|---|---|
| Model Size FLASH | 109 KB | |
| Minimum  RAM | 72 KB | |
| Minimum Cache | 2 KB | |
| Optimum  Cache | 74 KB | |

| | | |
|---|---|---|
| MAC Operations | ~2 M | |
| Weights | 26 K | |
| Layers | 9 | |
| Accuracy | 93.3% | |

(example for 32-bit float implementation)

| Model FLASH | 109 →54 KB | MAC Operations | ~2M |
| Minimum RAM | 72 KB | Weights | 26K→13K |
| Minimum Cache | 2 K→.57 KB | Layers | 9→7 |
| Optimum Cache | 74 K→37 KB | Accuracy | 93.3 → 92.7% |

(example for 32-bit float implementation)

# Convolutional Network  FCN_20



| Input | Conv1 | MaxPool1 | Conv2 | MaxPool2 | Conv3 | MaxPool3 | Conv4 | SoftMax |

16→6    32→8    64→24

| Model FLASH | 109 → 20 KB | MAC Operations | ~2M→ 222K |
| Minimum  RAM | 72 → 28 KB | Weights | 26 → 4K |
| Minimum Cache | 2 → 0.2 KB | Layers | 9 |
| Optimum  Cache | 74 →   8 KB | Accuracy | 93.3 → 91.5% |

(example for 32-bit float Implementation)

# Target Implementation

How to test and deploy on target hardware?

Requires optimized inference implementation

Vendor framework with generic network support

- Design flow to benchmark and iterate quicklywith different network designs
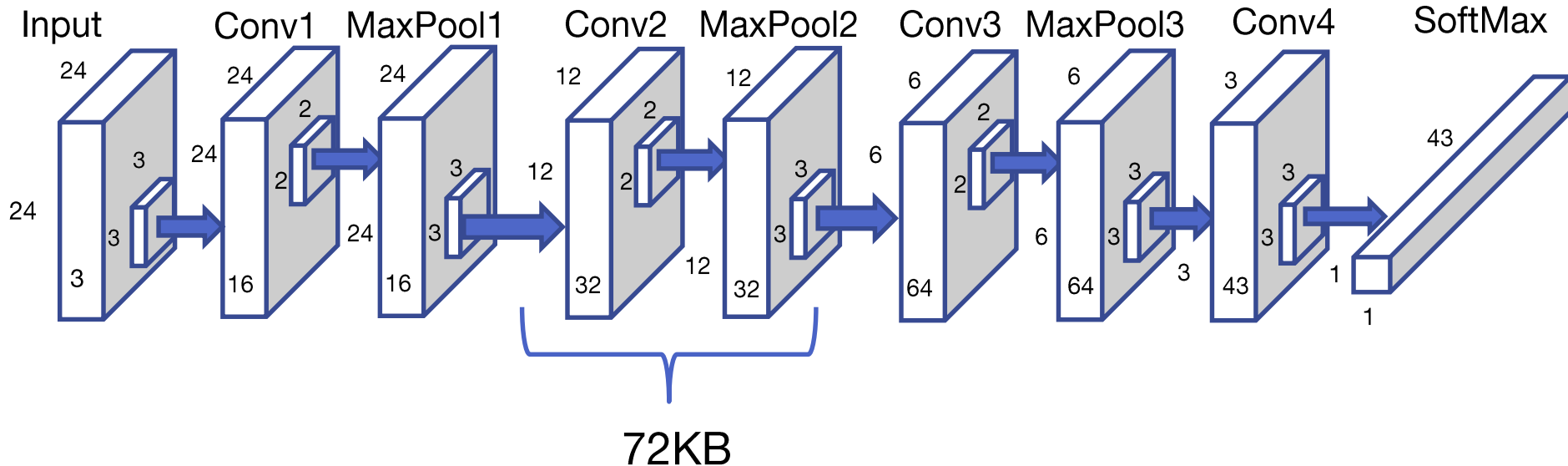- Profile tools for fitting designs
- Support for network update

Custom C code for selected network

- High level of optimization is possible
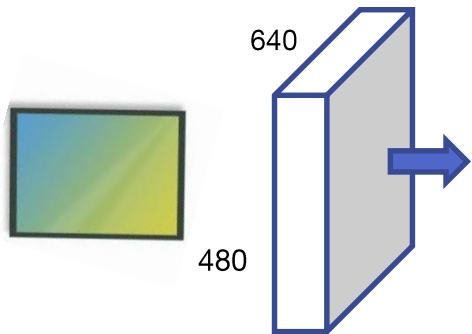- Additional effort to prototype and implement

# Challenges to address

- Fixed memory allocation at compile time

- Available contiguous memory on target microcontroller?

- Buffer management for best RAM utilization

- Efficient image sensor pipeline

- Validation of accuracy on the target

# Fully Convolutional  (FCN) – RAM Optimization



72KB

Minimum RAM is determined by largest layer for linear network
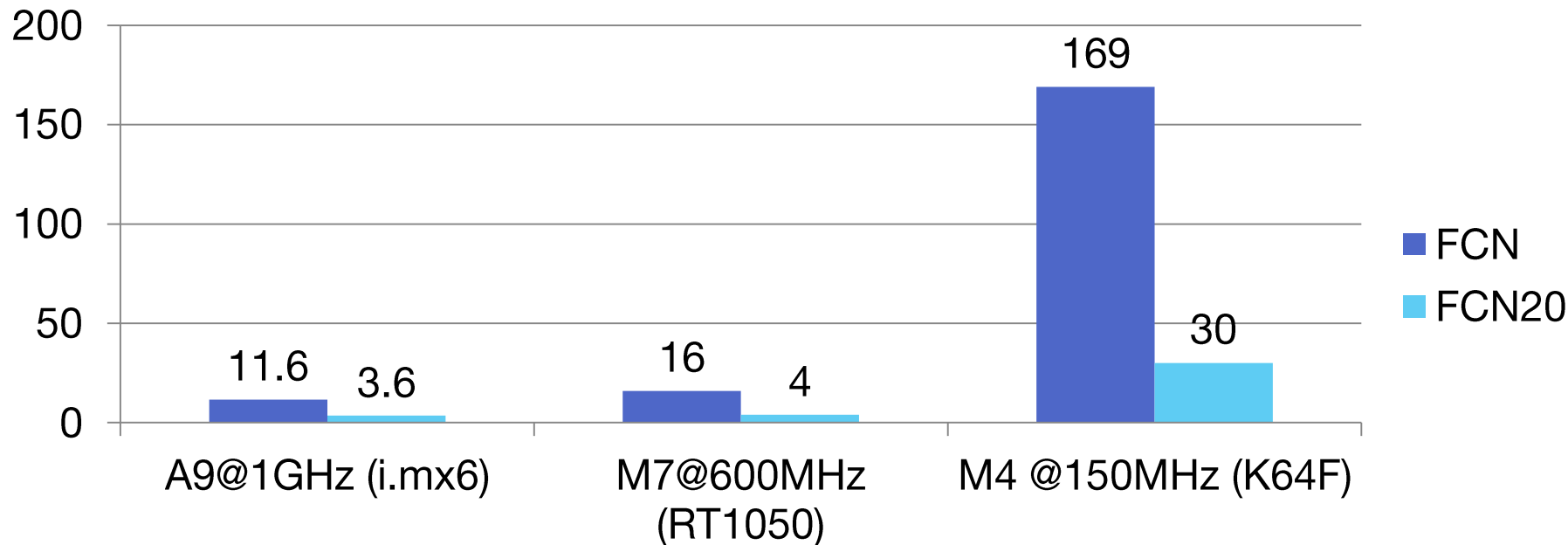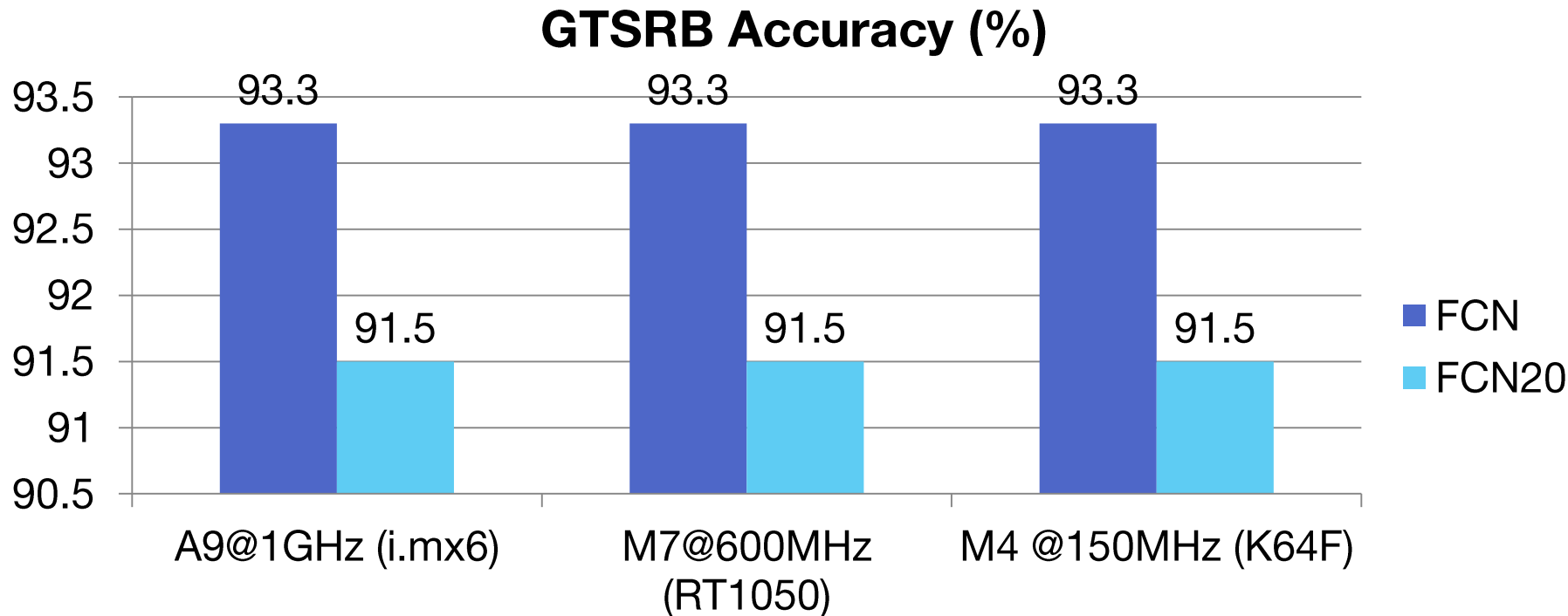Non-linear networks require RAM buffer management

| Sensor | H | V | Pixels | 16bit  (RGB565) |
|---|---|---|---|---|
| VGA | 640 | 480 | 307 K | 614 KB |
| QVGA | 320 | 240 | 77 K | 154 KB |
| QQVGA | 160 | 120 | 19K | 38 KB |

- Use Binning or Scaling to reduce Image sensor buffer
- Use DMA  to minimize slow memory buffer copies
- Use a monochrome sensor if color is not required for application

GTSRB Inference Time (ms)

GTSRB Accuracy (%)

# Areas for further Optimization

- Implement networks with 8-bit /16-bit arithmetic for devices that do not have hardware FPU support

- Cloud integration for distributed edge and network-based solutions where connectivity is practical

# Conclusion

1.  IoT-embedded vision using CNNs is practical on very low-cost microcontrollers
2.  An appropriate network design is required for design and target fit
3.  Memory optimization techniques should be carefully implemented
4.  Specific part selection is driven by performance/accuracy  and  other system requirements
5.  Optimized embedded inference code and tools can accelerate implementation

- Drop by Au-Zone  booth 802 for further demos and details

- [Software Frameworks and Toolsets for Deep Learning-based Vision Processing](#)

- [Demonstration of the DeepView ML Toolkit for Embedded Platforms](#)

- [www.embeddedml.com/deepview](http://www.embeddedml.com/deepview)

- [ARM CMSIS NN](#)

- [NXP i.MX RT Series: Crossover Processor](#)

- [NXP Kinetis® K Series:M4 Core MCU](#)

# Backup Slides

# Data Set Curator and Augmentation

# Network Design with profile details

# Network Training

# Network Validation on Host