



Processing Raw Images Efficiently with the MAX78000 AI Neural Network Accelerator

Mehmet Gorkem Ulkar

Principal Engineer, Machine Learning

Analog Devices



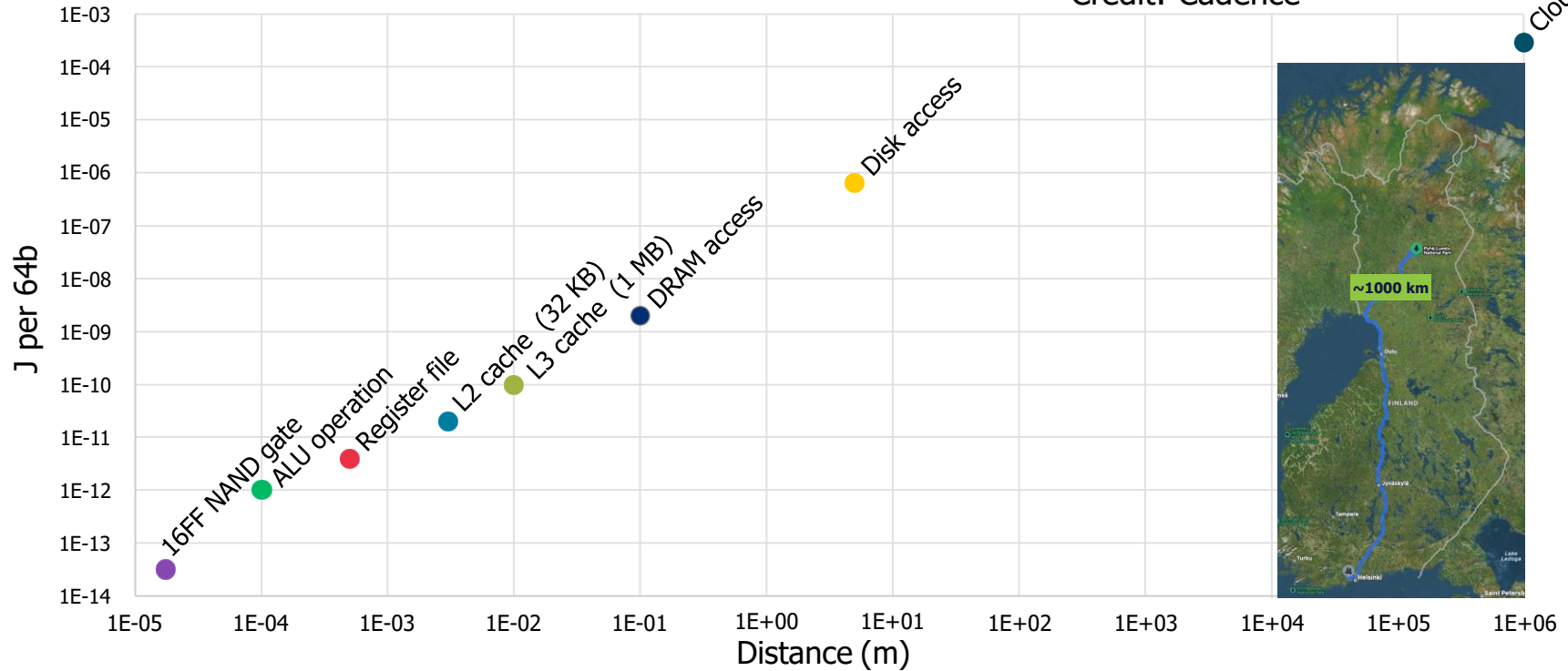


Mehmet Gorkem Ulkar,
PhD
Dallas, TX
Principal ML Engineer

1. Challenges of AI at the edge
2. MAX78000 overview
3. MAX78000 sample applications
4. Energy requirements for data manipulation
5. Proposal: CNN based de-bayerization
6. Results
7. Q&A

Keep Your Data Close: The Physics of Data

Credit: Cadence



Sources:

- Rick Zarr, TI, 2008, The True Cost of an Internet "Click" - estimate of transfer cost for 30KB page from server <http://energyzarr.typepad.com/energyzarrnationalcom/2008/08/the-true-cost-o.html>
- J Kunkel et al, University of Hamburg 2010, Collecting Energy Consumption of Scientific Data
- Horowitz ISSCC 2014, 1300-2600 pJ per 64b access
- Chris Rowen, Cadence Design Systems, January 2016, Get Real! - Neural Network Technology for Embedded Systems

Software Inference: Slow and Power Hungry

- ▶ In inference, computational effort is in **forward propagation**
 - On classic hardware, almost all spent in a triple nested matrix multiplication loop
 - $O(n^3)$ to $O(n^{2.8})$ *
- ▶ Very energy intensive even with fast matrix multiply using integer math on DSP or GPU
 - large number of memory accesses

```

// Main loop
for (l = 1; l < NLAYERS-1; l++) {
    // Compute z = w * a
    matrix_mul(&w[l-1], &a[l-1], &z[l]);

    // Add the bias values : z = w * a + b
    matrix_add(&b[l], &z[l]);

    // Compute a = g(z)
    nn_activate(z[l].elements, a[l].elements, lv[l]);
}

matrix_mul(matrix_f32_t *a, matrix_f32_t *b, matrix_f32_t *c)
{
    uint32_t m = a->nrows;
    uint32_t n = a->ncols;
    uint32_t p = b->ncols;

    c->nrows = m;
    c->ncols = p;

    int i, j, k;

    for (i = 0; i < m; i++) {
        for (j = 0; j < p; j++) {
            f_t sum = 0;
            for (k = 0; k < n; k++) {
                sum += a->elements[i * n + k] * b->elements[k * p + j];
            }

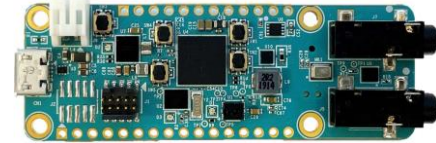
            c->elements[i * p + j] = sum;
        }
    }
}

```

*Strassen's algorithm

CNN Accelerator: MAX78000/MAX78002

- The conv operation is parallelizable in the channel dimension.
 - 64 processors in total, more channels are processed in a multi-pass fashion
- Proper architecture that minimizes data movement provides energy efficiency
 - Each input channel is processed in parallel using different processors to minimize data movement
 - Each processor uses dedicated memory



MAX78000 AI Micro - System-on-Chip

Memory

Flash
512 KB

SRAM
128 KB

Ultra Low Power Micro

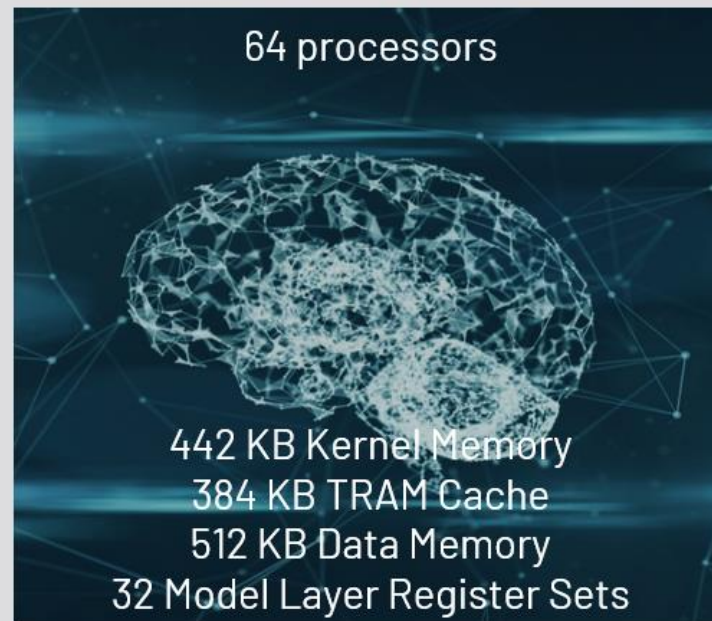
Arm Cortex-M4F
100 MHz

Cache

RISC-V

CNN Accelerator

64 processors



Timers

3 × Timer

Watchdog

External Interfaces

2 × Quad SPI

UART, 2 × I2C, ADC, I2S

Parallel Camera

Security

AES

TRNG

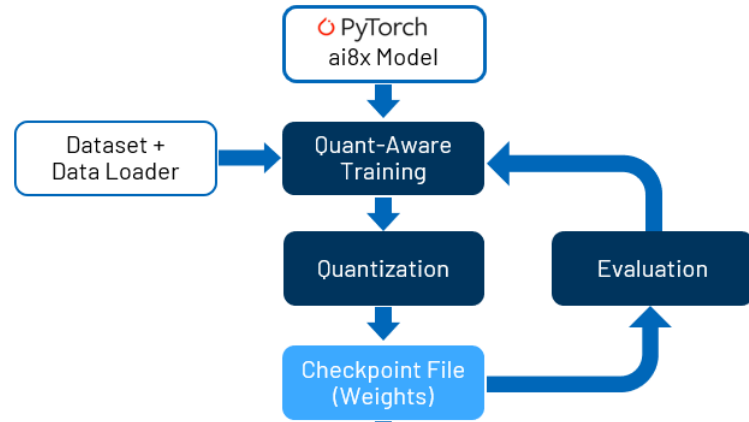
Power

SIMO

Model, Training, Deployment: Development Flow

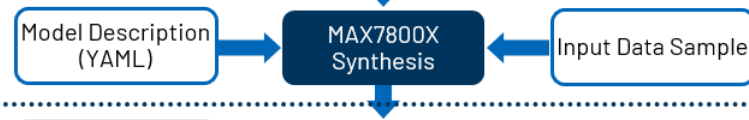
① Training ai8x-training

Machine Learning
Experts

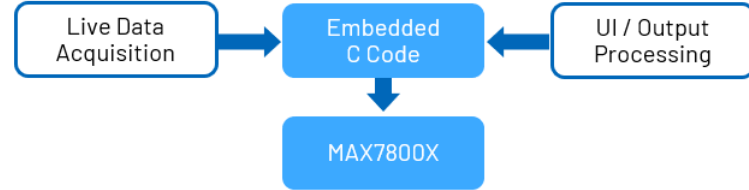


② Synthesis ai8x-synthesis

Embedded
Engineers

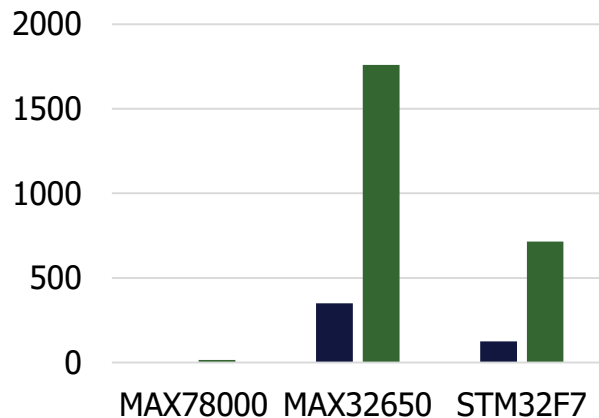


③ Deployment Embedded SDK

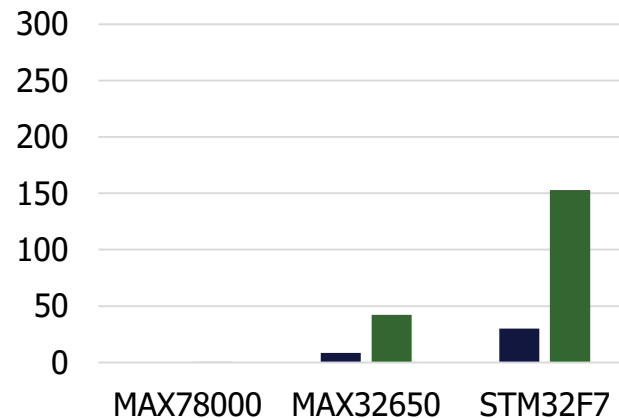


MAX78000 Benchmarks

Inference Time ms



Inference Energy mJ



¹28 billion operations/second,
²ARM DSP with CMSIS-NN, running exact same INT8 network as MAX78000, ³STM722ZE, internal memory,
⁴Includes time to load input,
⁵Does not include time to load input,
⁶STM746NG + external 3.3V SDRAM IS42S32400F-6BL + SDRAM controller

Network	MACs	MAX78000 CNN at 50 MHz ¹ , 1.2V	MAX32650 ² Cortex-M4, 120 MHz, 1.2V	STM32F7 ² Cortex-M7, 216 MHz, 2.1V
▣ KWS20	13,801,088	2.0 ms, 0.14 mJ	350 ms, 8.37 mJ	125 ms, 30.1 mJ ³
▣ FaceID	55,234,560	13.89 ms ⁴ , 0.40 mJ	1760 ms ⁵ , 42.1 mJ	714 ms ⁵ , 153 mJ + 59 mJ ⁶

Battery Life Leader in Independent Benchmarks

ETH zürich

Proof-of-Concept standalone smart camera

Assumptions:

- Trigger time: once per minute.
- Battery capacity: 8.64J.
- Energy per camera image captured: 0.5 mJ.

Platform	Energy per inference (mJ)	Battery Lifetime
SAMD51	5.34	24h30'
Apollo3	1.31	80h00'
Spresense	3.80	46h00'
GAP8	0.52	140h15'
VEGA	0.14	225h00'
xCORE.ai	1.26	81h50'
MAX78000	0.09	244h00'



Perpetual work with only 100Lux with less than 1 second starting time

Best



PROJECT BASED LEARNING



Michele Magno | 29.05.2020 | 18



A Battery-Free Long-Range Wireless Smart Camera for Face Detection: An accurate benchmark of novel Edge AI platforms and milliwatt microcontrollers. Michele MAGNO, Head of the Project-based learning Center, ETH Zurich, D-ITET, EMEA TinyML Talks June 2021

Thinking About Edge AI Use Cases...

If my [application] $\left\{ \begin{array}{c} \text{sees} \\ \text{hears} \\ \text{senses} \end{array} \right\}$ _____
object/sound/event/situation/...

then do _____
action

If my [camera] sees a bear, then take a high-resolution picture and send over cell network

If my [thermostat] hears glass break, then send a text message to the owner

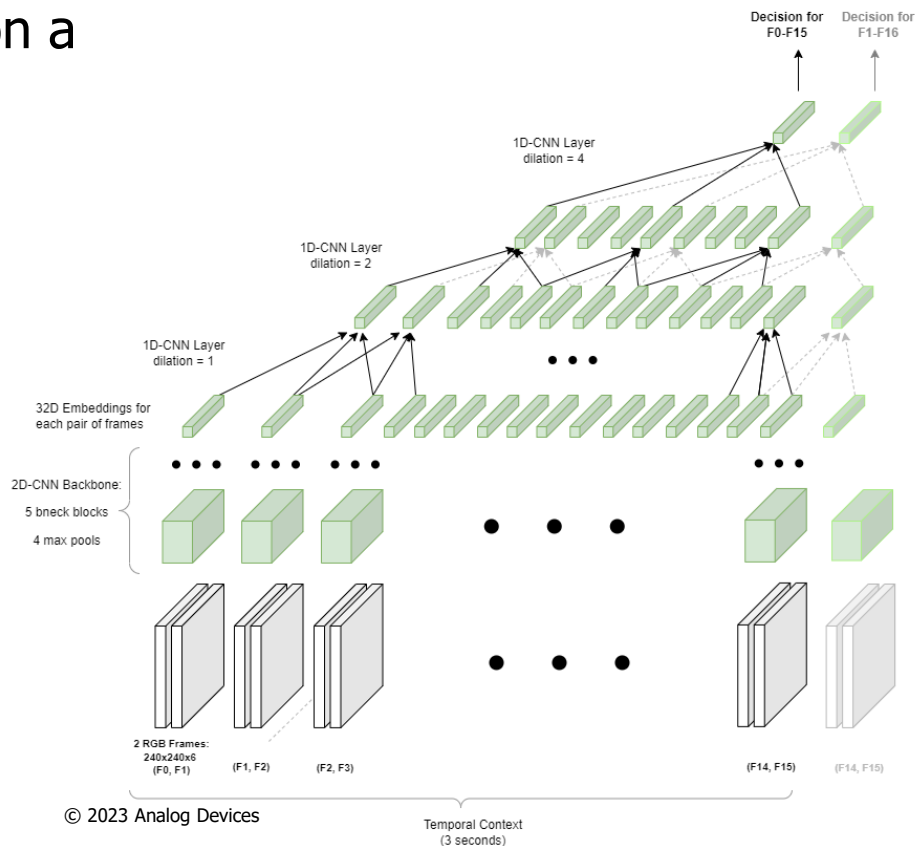
If my [factory robot] sees a person nearby, then shutdown until they leave

If my [pet door] sees a cat with a mouse in its mouth, then lock the pet door and send me a text message

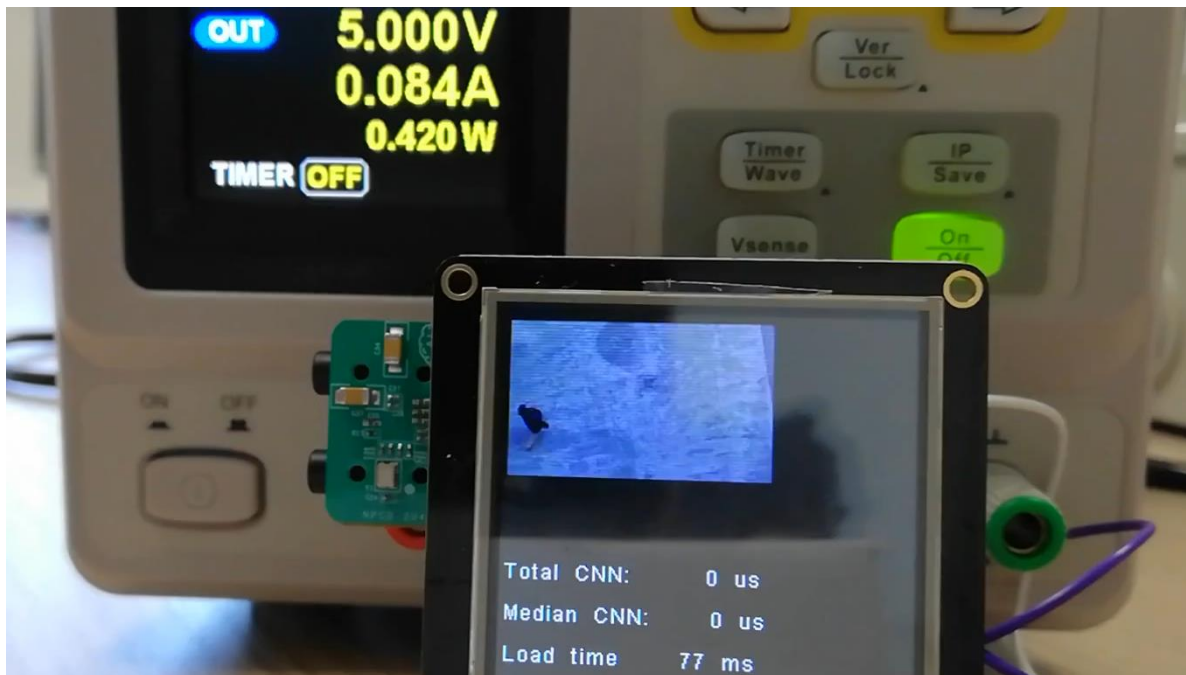
Action Recognition

- Embeddings saved in memory on a rolling basis
 - No redundant calculations

Dataset	Validation Acc.	Parameters
Kinetics-400 (4 classes + other)	79.8%	379k



People Tracking



No Url

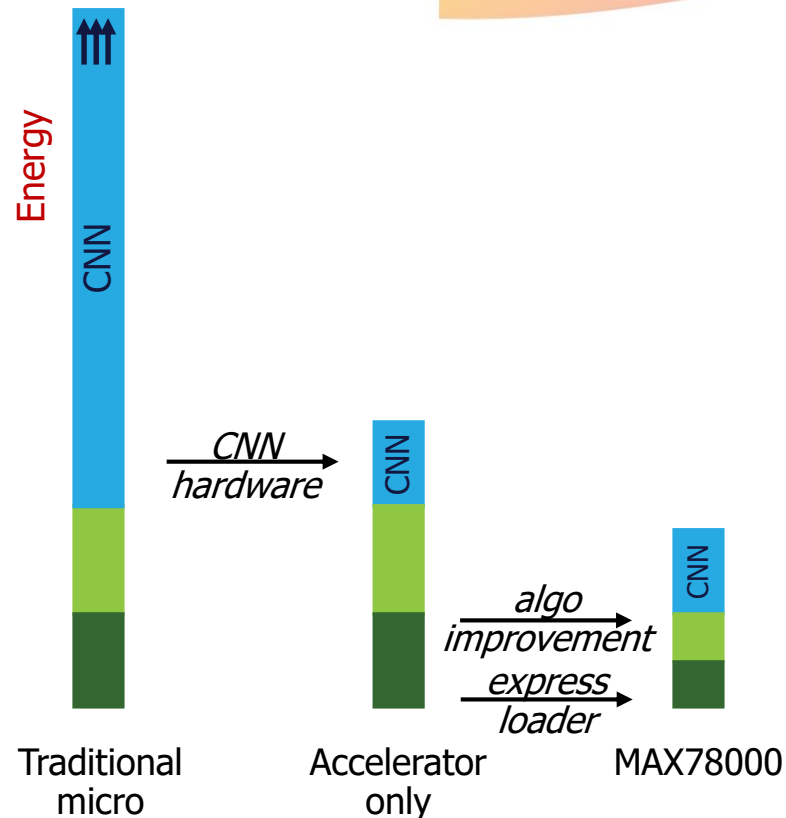
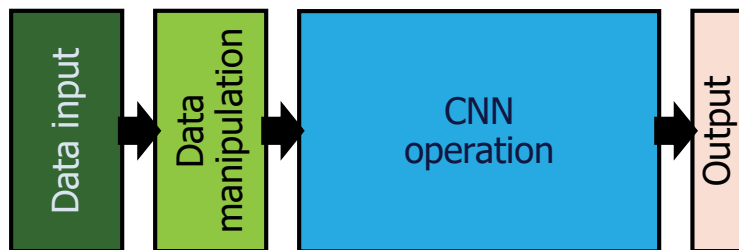
Trail Camera



<https://github.com/MaximIntegratedAI/refdes>

System Energy: From Traditional Systems to MAX78000

- Accelerator drastically lowers CNN energy
- Input and **data manipulation** become much larger *relative* contributors to energy
- MAX78000 improves data loading, better algorithms can help with data manipulation: e.g. **better ways of handling raw images**



Data Manipulation: Debayerization

In order to obtain an RGB format, the raw image must be debayerized. There are several debayerization methods*:

- Bilinear Interpolation
 - Sequential Demosaicing
 - Iterative Demosaicing
 - **Machine Learning Methods**
 - Adaptive Color Plane Interpolation
- ▶ Outside the CNN accelerator
- ▶ Increased system energy consumption

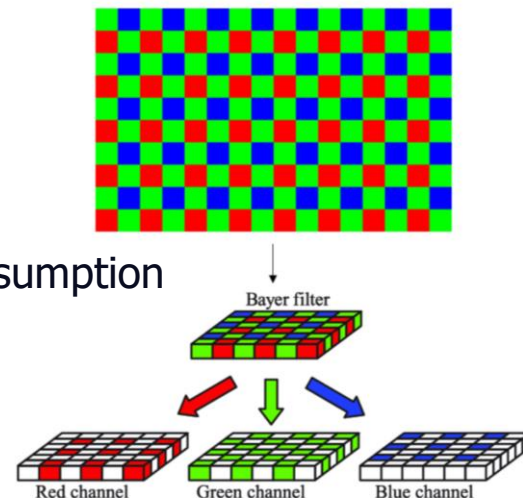


Figure 1. Bayer Filter (Nkansah et. al., 2022)

* Dammer, K., Grosz R., (2017). Demosaicing using a Convolutional Neural Network approach. Lund University, Lund, Sweden.

CNN based Debayerization

- Approach 1: Learning the manipulation & interpolation by a CNN model and embedding this network into an accelerator → Efficient way of debayerization

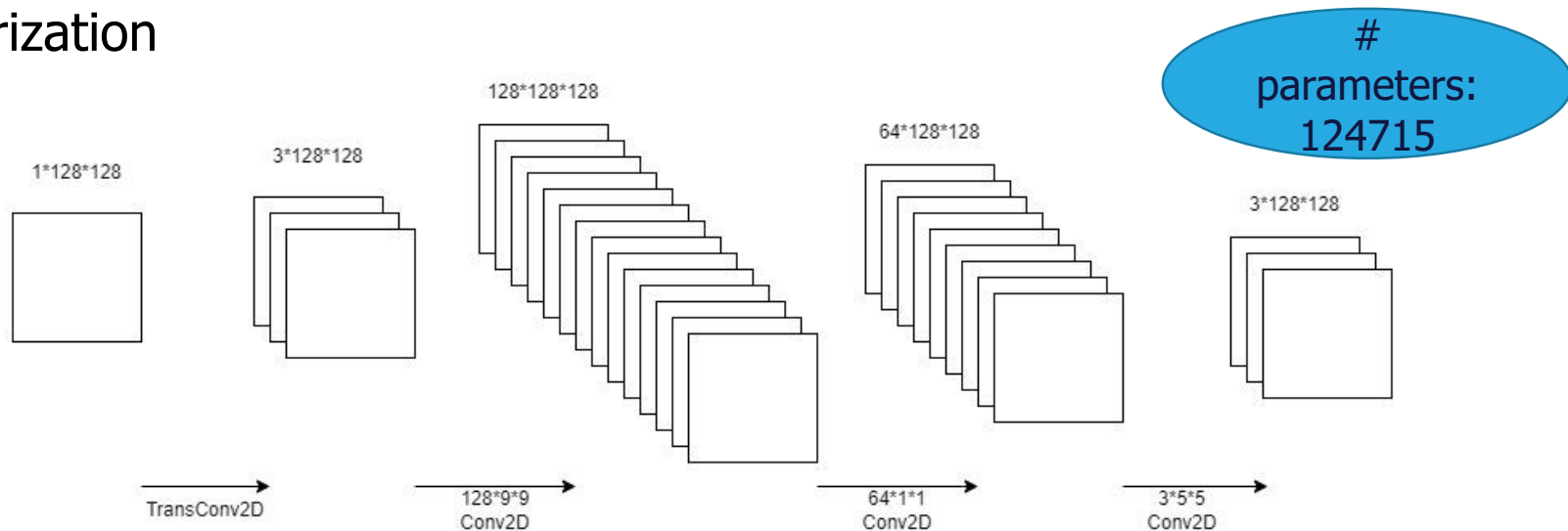
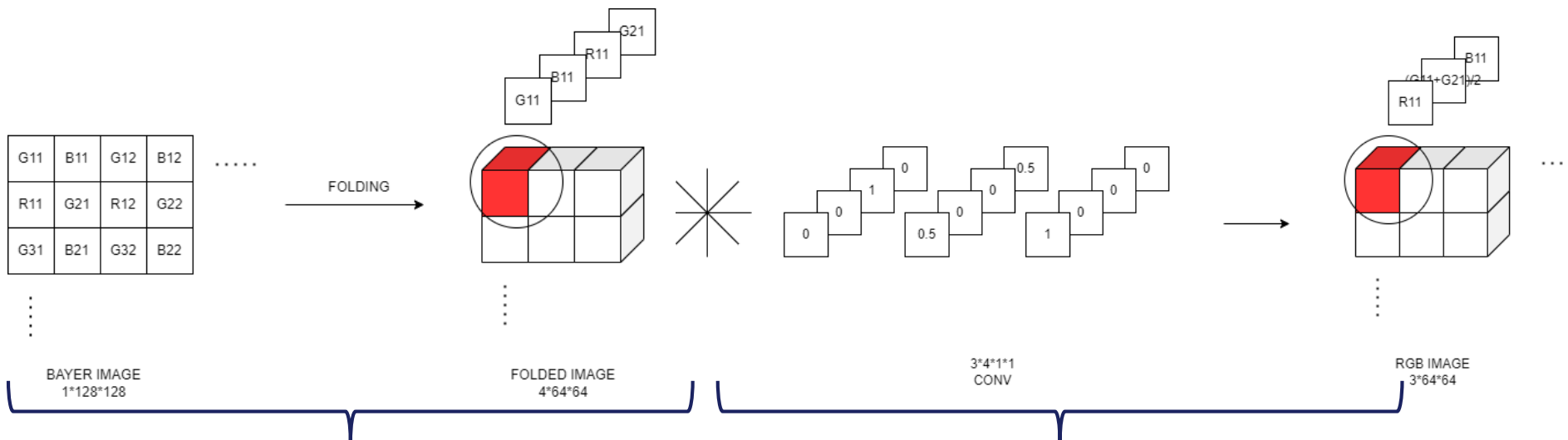


Figure 3. The Network of B2RGBNet (Syu et. al., 2018)

CNN based Debayerization

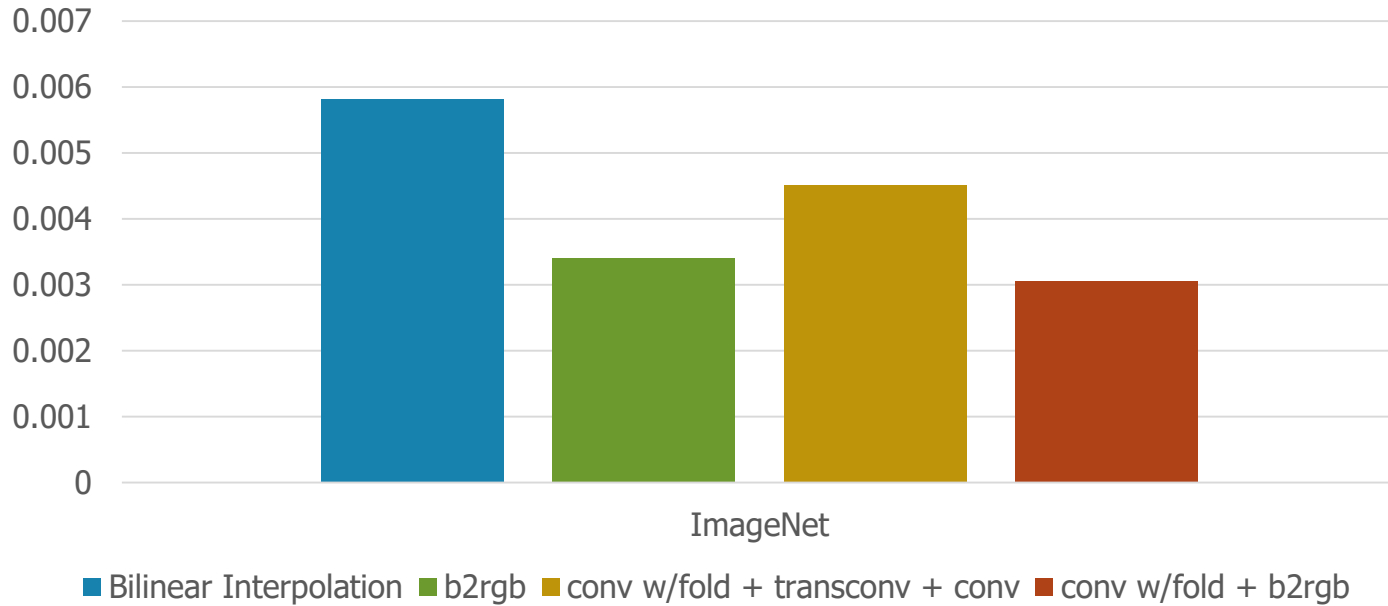
- Approach 2: Using folding and fixed 1x1 kernels



Step 1: Folding the pixels into channels

Step 2: Convolution with the fixed kernel to obtain RGB

Mean Squared Reconstruction Error



Conclusion

- MAX78000 enables battery-powered smart applications at the edge
- Effective data manipulation and preprocessing are much more important when using highly-efficient NN inference engines
- Two methods proposed to perform interpolation inside CNN accelerator, MAX78000
- Results show better accuracies compared to simple conventional interpolation; the work is ongoing

- We are waiting for you at the ADI booth!
- Upper-level AI repo: <https://github.com/MaximIntegratedAI>
- Open-source training repo: <https://github.com/MaximIntegratedAI/ai8x-training/>
- Open-source synthesis repo: <https://github.com/MaximIntegratedAI/ai8x-synthesis>
- Data-folding paper: L3U-net: Low-Latency Lightweight U-net Based Image Segmentation Model for Parallel CNN Processors
<https://arxiv.org/pdf/2203.16528.pdf>
- B2RGBNet paper: Learning Deep Convolutional Networks for Demosaicing
<https://arxiv.org/pdf/1802.03769.pdf>