# Audio and Image Cross-Modal Intelligence via a 10TOPS/W 22nm SoC with Back-Propagation and Dynamic Power Gating

Zichen Fan, Hyochan An, Qirui Zhang, Boxun Xu, Li Xu, Chien-wei Tseng, Yimai Peng, Ang Cao, Bowen Liu, Changwoo Lee, Zhehong Wang, Fanghao Liu, Guanru Wang, Shenghao Jiang, Hun-Seok Kim, David Blaauw, Dennis Sylvester

University of Michigan, Ann Arbor, USA.    zcfan@umich.edu

**Abstract** We present an ultra-low-power multimedia signal processor (MMSP) SoC that integrates a versatile deep neural network (DNN) engine with audio and image signal processing accelerators for cross-modal IoT intelligence. The proposed MMSP features 2MB MRAM to store all DNN weights on-chip with an energy-efficient dataflow using an MRAM-cache and dynamic power gating. The SoC achieves up to 3-10 TOPS/W peak energy efficiency and consumes only 0.25-3.84 mW. Being the first to demonstrate CNN, GAN, and back-propagation (BP) on a single accelerator SoC for cross-modal fusion, it outperforms state-of-the-art DNN processors by 1.4 - 4.5× in energy efficiency.

**Introduction:** DNN-based image and/or audio processing has been widely adopted in intelligent IoT systems. However, the traditional processing flow (Fig. 1 top) that offloads DNN processing to the cloud suffers from bandwidth, energy and privacy problems for resource-constrained IoT applications. To tackle these challenges, the proposed SoC adopts a fully-at-edge processing flow for *audio and image multimodal intelligence* with integrated pre-/post-processing accelerators (JPEG, H.264, FFT, and Mel-spectrum), a GAN/BP-capable neural engine, and 2MB on-chip MRAM for non-volatile weight/data/code storage.

Fig. 2 depicts target application scenarios of the proposed SoC. It can simultaneously receive 12b per pixel VGA images and 8kHz 8b per sample audio signals using dedicated interfaces. The SoC can perform image- or audio-only applications such as face detection or keyword recognition, and also (conditionally) execute image-audio fusion applications such as cross-modal verification and active speaker detection to further increase the detection and recognition credibility. To accomplish this, the neural engine supports different types of DNN operations such as (depth-wise) convolution and fully connected layers (FCLs) for CNN, deconvolution for GAN, back-propagation for BPGAN [3], and on-chip weight retraining (transfer learning of the last layer).

**Architecture:** Fig. 3 shows the overall SoC architecture. All sub-blocks communicate with each other using an AHB bus. The audio interface performs audio feature extraction using frame buffers, a 256-point FFT unit, a 32/64 channel Mel-filter unit, and a power-to-dB $\log_2$ unit. The image interface features a change detection and on-the-fly JPEG compressed-memory [4] to temporarily store VGA frames in a compressed format. Only the change detected macro-blocks are stored and processed as regions-of-interest (RoIs). The H.264 engine performs image compression [4] on non-rectangular RoIs for compact storage in on-/off-chip memory. The neural engine (NE) is a reconfigurable DNN accelerator that supports various operations including (depth-wise) convolution, deconvolution, FCL, and back-propagation. NE supports DNN weights both in uncompressed (8b/weight) and Huffman-compressed (~2b/weight) formats. The SoC integrates a 2MB MRAM macro to store all weights on-chip, and a 1.5MB multi-bank SRAM activation memory to store all feature maps for backpropagation. For simple applications that only require partial activation memory, unused SRAM macros are dynamically power gated to save leakage power. The main computation unit is an 8×8×8 processing element (PE, each with an 8-bit MAC) array, which enables activation and weight reuse via inter-PE connections. The top 1×8×8 PEs are multi-functional PEs (MPEs) that supports both MAC and max/average pooling operations. Furthermore, ping-pong memory structures for the local memory and row/col buffer enable non-blocking pipelining between data movement and computation to increase PE utilization. The MRAM macro is accompanied with an MRAM cache to enable a dynamic power gating scheme (details in later section).

**Power Domain:** Fig. 4 shows the power domain design of the SoC. We implement power gating for each SRAM block and MRAM macro to decrease standby leakage power up to 83%. Because the MRAM can retain all data, SRAMs can be power gated when NE is inactive. Most of the computation logic can also be power gated while the always-on block (pads, headers, state registers, and power sequence control logic) consumes only 460nW.

**NE Dataflow:** Fig. 5 shows the energy-efficient, computation-skipping dataflow scheme implemented in NE. The base dataflow is output stationary, which is used in convolution, strided convolution, and depth-wise convolution. Weights are shared along PE array rows, and input activations are shared across the output channels (OCs). For deconvolution, a zero-skipping dataflow exploits the deterministic and regular pattern of zero padding to increase throughput by 4× by only computing non-zero values in each PE (Fig. 5 left). For back-propagation through a ReLU layer, the zero activation positions are pre-recorded during the forward path to data-gate all computations in the backward path if they (back)propagate to a pre-recorded zero activation position (Fig 5 right).

**MRAM-cache Dynamic Power Gating:** Although MRAM has the advantage of high density and non-volatility, its active leakage and readout power can be significant for ULP applications. To mitigate this issue, we propose an MRAM-cache architecture and dynamic power gating scheme (Fig. 6). In our scheme, MRAM is power gated until NE executes the load weight (LD_WEIGHT) instruction. During LD_WEIGHT operation, MRAM powers up and weights are read and loaded into the SRAM-based weight cache (WC). MRAM then goes back into either 1) sleep (SLP) mode where the MRAM array is powered off but peripherals remain on, or 2) power-down (PD) mode where the peripherals are also power gated. The optimal selection between SLP and PD depends on the reuse factor of the cached weight, as shown in Fig. 6 (bottom, right), which can be identified during NE programing. The measured MRAM VDIO current waveform shows the whole MRAM-cache sequence. During NN processing, weights are read from the weight cache (while MRAM is in either SLP or PD) for reduced memory readout power compared to MRAM accesses. Based on measured results, the combined weight caching and power gating reduces weight readout power by 95.3% with only slightly increased (4.3%) operation time due to MRAM wake-up latency and cache loading time overhead.

The table in Fig. 6 summarizes the three different power modes: PD, SLP, and stand-by (STB). The analysis concludes that when each cached weight is reused for 353.4 MACs, the PD mode is preferable. Otherwise, SLP mode has an advantage because of the lower overhead in power-up energy (offsetting its higher leakage).

**Measurement Results:** The SoC was fabricated in 22nm ULL technology and the die photo is shown in Fig. 9. The peak energy efficiency for various NN instructions is shown in Fig. 7 (left). The efficiency for convolution / deconvolution / stride-convolution-backpropagation (CONV / DECONV / S_CONV_BP) is 3.1 / 10 / 10 TOPS/W. The efficiency for DECONV and S_CONV_BP is significantly higher because of the zero-skipping dataflow. The convolution backpropagation can achieve 3.7 TOPS/W with zero-gating dataflow. Depth-wise convolution (DWCONV) and FCL have lower efficiencies due to only 1/8 of the PE array being utilized. Fig. 7 (right) shows the voltage-frequency-efficiency tradeoff for CONV. The SoC achieves the highest energy efficiency at 0.46V (VDD_MAIN) and 1.2MHz, while the total system power is 387uW. Fig. 7 (right) shows that MRAM dynamic power gating enhances the energy efficiency especially at low voltage as it significantly reduces the leakage.

Fig. 8 demonstrates a person-of-interest (PoI) tracking scenario and chip performance for that cross-modal intelligence scenario. This task first performs face detection (FD) on change-detected regions and then face recognition (FR) if faces are detected. The H.264 engine compresses the changed blocks if a PoI is recognized. In the meantime, the audio interface extracts the audio features (AFE), and cross-modal verification (CMV) is performed using both the audio features and face image. If the CMV model confirms the audio is matched with the target face, the audio signal is compressed (AC) by BPGAN. Fig. 8 (bottom) shows the power consumption and latency of each step in this process. Table 1 compares this SoC and other state-of-the-art prior designs in similar application spaces.

**References:** [1] K.D. Choo *et al.*, ISSCC'19 [2] M. Cho *et al.*, ISSCC'19 [3] B. Liu *et al.*, ICASSP'20 [4] H. An *et al.*, VLSI'20 [5] P. Jokic *et al.*, VLSI'21 [6] J. Giraldo *et al.*, VLSI'19 [7] J. Lee *et al.*, VLSI'19 [8] M. Giordano *et al.*, VLSI'21 [9] S. Kang *et al.*, JSSC'21
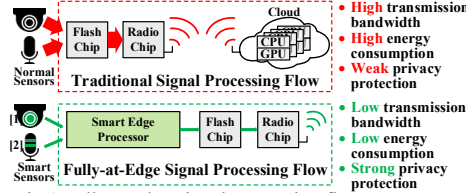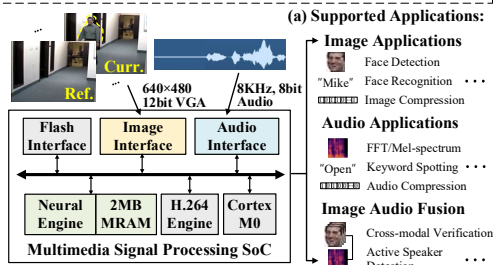
Fig.1 Fully at edge signal processing flow.



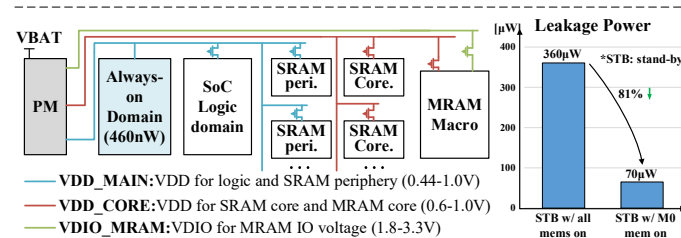Fig.2 Proposed multimedia signal processing SoC and its supported applications.



Fig.3 System overall architecture.



Fig.4 Overall power domain with 460nW always-on domain.



Fig.5 Efficient dataflow for both forward and backward propagation.



Fig.6 MRAM-cache architecture and dynamic power gating scheme for MRAM.



Fig.7 Measured peak efficiency and voltage-freq-efficiency scaling



Fig.8 Person-of-interest (PoI) tracking system scenario and its performance.



Fig.9 Die photo and chip summary.

| | | Specifications |
|---|---|---|
| | Technology | 22nm |
| | Die Area | 4.35mm*2.75mm (12mm²) |
| | Voltage | 0.44V-1.0V(Main)/0.6V-1.0V(Core)/1.8V-3.3V(IO) |
| | On-chip Memory | 1.98MB SRAM + 2MB MRAM |
| | Max Frequency | 70MHz (On-chip CLK gen) |
| | Power Consumption | 460nW @ Always-on block |
| | | 70µW @ Stand-by Leakage w/ M0 mem on |
| | | 0.25mW@0.8MHz, 3.84mW@10MHz |
| | Performance | Sys. Peak Efficiency: 3.1-10TOPS/W |

| | | This Work | [4] | [5] | [6] | [7] | [8] | [9] |
|---|---|---|---|---|---|---|---|---|
| Applications | | Face detection, face recognition, audio compression, cross-modal verification | Person detection, face detection, face recognition | Face detection/ face recognition | Keyword spotting, speaker verification | Super resolution | NN inference & training | Face manipulation |
| Algorithms | CNN/FC | ✓ | ✓ | ✓ | ✓(LSTM) | ✓ | ✓ | ✓ |
| | GAN | ✓ | - | - | - | - | - | ✓ |
| | BP | ✓ | - | - | - | - | ✓(FC only) | ✓ |
| Image processing | | Change detection/ JPEG/H264 | Change detection /JPEG/H264 | - | - | - | - | - |
| Audio processing | | FFT/Mel spectrum | - | - | FFT/MFCC | - | - | - |
| Technology [nm] | | 22 | 40 | 22 | 65 | 65 | 40 | 65 |
| Die Area [mm2] | | 12 | 27 | 3.4 | 2.6 | 16 | 29.2 | 32.4 |
| Non-volatile Memory | | 2MB MRAM | - | - | - | - | 2MB RRAM | - |
| On-chip SRAM [MB] | | 1.98 | 1.13 | 1.2 | 0.1 | 0.56 | 0.50 | 0.66 |
| Off-chip Memory access | | No | No | No | No | Yes | No | Yes |
| Precision | | INT8 | INT8 | INT1/16 | INT8 | INT8 | INT8 | FP8/FP16 |
| Voltage [V] | | 0.44-1.0 | 0.58-0.7 | 0.65 | 0.6-1.2 | 0.75-1.1 | 0.7-1.1 | 0.7-1.1 |
| Max Frequency [MHz] | | 70 | 0.15 | 180 | 12.5 | 200 | 200 | 200 |
| Power [mW] | | 0.25-3.84* | 0.17 | 16.7 | 0.02 | 31-211 | 126 | 58-64† |
| Sys. Peak Eff. [Tops/W] | | 3.1-10** | 1.5 | 1.1 | 0.7 | 1.9 | 2.2 | 1.8 |

*0.25mW is measured @0.44V, 0.8MHz, 3.84mW is measured @0.6V, 10MHz. **10Tops/W is from deconvolution operation.

Table.1 Comparison table versus state-of-the-art.