

Fig. 2. MLA enables various machine-learning kernel functions.

Data Computation Unit (DCU)

Fig. 3 shows the DCU. Multiplications/additions/subtractions are structured in stages with selectable operands, enabling specialized data paths to be configured for the various kernels. The two-stage pipeline mitigates active energy due to propagating glitches. In order to maximize hardware usage, the operands from memory are provided through an interface (MMU) capable of $2\times$ clock-boosting. From application profiling, a key challenge identified is the ability to optimize the fixed-point precision across the framework configurations. Two approaches overcome this. First, like the MAXMIN, the DPU can be configured as one 32b module or two parallel 16b modules wherein each input (in1 and in2) can be configured into two operands; configurability is achieved with low-overhead via the 33b adder shown, where the PAD bit allows separation/merging of two 16b paths. Second, a configurable barrel shifter enables in-line truncation to support any Q-format. For 16b mode, a reduce stage is used to combine the parallel outputs into a final result (e.g., for vector dot products, the outputs from sub-divided vectors are added together, as shown).

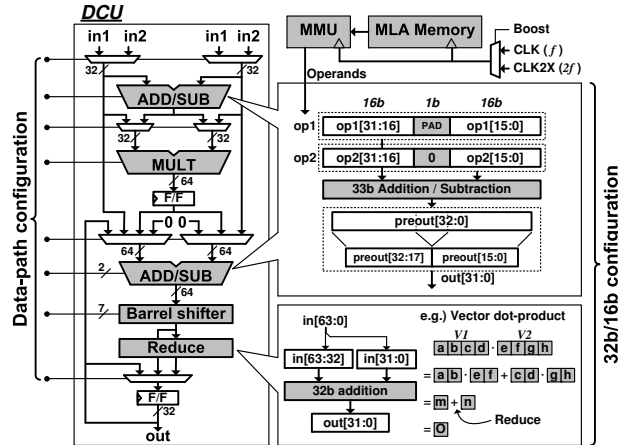


Fig. 3. DCU can be configured into data-paths specialized for the MLA kernels, and enables precision scaling via (1) configurability into 32/16b modules, and (2) in-line barrel shifting and truncation.

Compression Decompression Accelerator (CDA)

Fig. 4 shows the CDA. Various compression algorithms were investigated. ADPCM was selected due to its low computational complexity, enabling substantial compression ($4\times$) with low energy overhead ($<8\%$). Though ADPCM is lossy, various applications were profiled, validating negligible impact in performance; the results from a seizure detector are shown. The CDA consists of two parallel encoders and decoders to support the 32/16b data formats. The DCU can be configured for upto four operands, but the kernel computations involve up to three operands from memory; local buffers thus enable data from three vectors to be stored locally so that the iterations required for estimation in the ADPCM algorithm can proceed simultaneously. Compression then involves encoding deltas (with respect to the running estimates) using values/indices from steptable/indexable LUTs; decompression involves using the encoded bits in parallel to retrieve decodings from the LUTs.

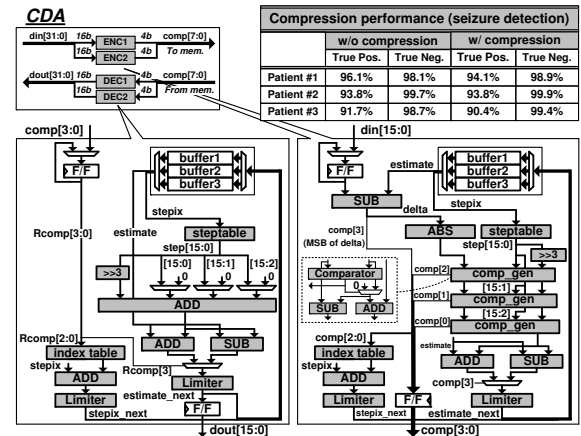


Fig. 4. CDA enables low-overhead on-line data compression (encoding) and decompression (decoding) using ADPCM; seizure-detector performance illustrates negligible impact of lossy ADPCM encoding.

Measurements and Demonstrations

The processor is implemented in 130nm CMOS from IBM. Fig. 5 shows a die photo, hardware measurements, and a comparison with the state of the art, illustrating the advances in hardware support achieved for machine-learning algorithms. The processor voltage scales from 1.2-0.7V for real-time operation of the applications tested. The applications, tested with real medical data, include the following: (1) ECG-morphology-based arrhythmia detection using SVM classification [4]; (2) ECG-wavelet-based arrhythmia detection using GMMs [5]; (3) patient-adaptive ECG-morphology-based arrhythmia detection using AL; (4) EEG-based seizure detection using SVM classification [6]; (5) EEG-based sleep-stage prediction using hybrid HMMs and GMMs [7]; and (6) spike/LFP-based motor-intention decoding for brain-machine interface using SVM regression [8]. From application profiling, energy savings of $3.1\text{-}497\times$ are achieved compared to CPU-only implementations (the total power of the applications on the processor is denoted in the last column of the table).

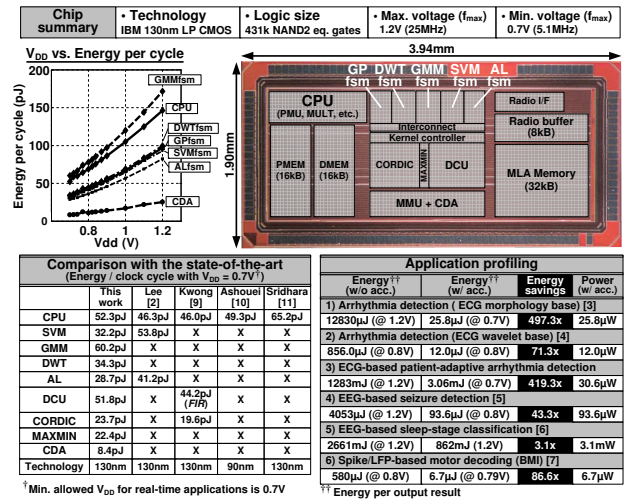


Fig. 5. Die photo, prototype measurements, and comparison.

References

- [1] A. Csavoy, et al., *VLSI Symp.*, Jun. 2009, pp. 4-7.
- [2] K. H. Lee and N. Verma, *ESSCIRC*, Sep. 2012, pp. 285-288.
- [3] J. Park, et al., *JSSC*, Nov. 2012, pp. 2711-2723.
- [4] P. de Chazal, et al., *IEEE Tran. Biomed. Eng.*, Jul. 2004.
- [5] R. J. Martix, et al., *J. Med. Syst.*, 2012, pp. 677-688.
- [6] A. Shoeb and J. Gutttag, *ICML*, Jun. 2010.
- [7] A. B. Rossow, et al., *BRC*, Jan. 2011, pp. 1-6.
- [8] L. Shpigelman, et al., *NIPS*, 2008, pp. 1489-1496.
- [9] V. K. Chippa, et al., *DAC*, Jun. 2010, pp. 555-560.
- [10] J. Kwong and A. P. Chandrakasan, *ESSCIRC*, Sept. 2010, pp. 526-529.
- [11] M. Ashouei, et al., *ISSCC*, Feb. 2011, pp. 332-334.
- [12] S. R. Sridhara, et al., *VLSI Symp.*, Jun. 2010, pp. 15-16.