# A Mixed-signal Time-Domain Generative Adversarial Network Accelerator with Efficient Subthreshold Time Multiplier and Mixed-signal On-chip Training for Low Power Edge Devices

Zhengyu Chen[1], Sihua Fu[2], Qiankai Cao[1], Jie Gu[1]

[1]Northwestern University, Evanston, IL, USA; [2]Google Inc., Mountain View, CA, USA

**Abstract** - This work presents a low-cost mixed-signal time-domain accelerator for generative adversarial network (GAN). A significant reduction in hardware cost was achieved through delicate architecture optimization for 8-bit GAN training on edge devices. An area efficient subthreshold time-domain multiplier was designed to eliminate excessive data conversion for mixed-signal computing enabling high throughput mixed-signal online training demonstrated in a 65nm CMOS test chip.

## Introduction

GAN is rendered as one of the most interesting and challenging applications in deep learning space. As shown in Fig. 1, GAN contains two deep neural networks (DNN), i.e. a generator and a discriminator, contesting and evolving with each other [1]. Despite its broad real-time applications in gaming, authentication, VR, there is a lack of dedicated low power GAN accelerator due to the tremendous challenges on resource-limited edge devices. From the algorithm aspect, GAN is extremely difficult to train due to model collapses from unbalanced models and high sensitivity to hyper-parameters. From the hardware aspect, GAN involves two DNNs with complex training sequences, e.g. 41 different training stages as in this work. Moreover, the typical floating-point training and complex calculation, e.g. batch normalization and optimizers, are very expensive for a resource-limited edge device [1]. This work, through significant architecture improvement and hardware adaptation, presents a mixed-signal GAN accelerator with 8-bit resolution for cost-effective implementation on edge device. The contributions include: (1) for the first time, a complete GAN training core was implemented on an 8-bit low-power ASIC chip consuming only 39mW; (2) An efficient subthreshold time-domain (TD) multiplier was designed with significant area saving compared to digital design; (3) On-chip training was performed in mixed-signal TD for the first time. The design eliminated 94% overhead from domain conversion, leading to the state-of-art throughput for a mixed-signal based accelerator which normally suffers from slow operation speed.

## GAN Accelerator Design with Time-domain Circuits

Fig. 1 shows the implemented GAN architecture with model compression and hardware adaptation techniques used in this work. For fitting with a small chip budget on edge device, we targeted a low-budget architecture implementation of DCGAN [1] using greyscale image with a size of 28 x 28 pixels. The following techniques were specially developed: (1) model balancing and adaptive training were utilized to enable 8-bit training versus conventional floating-point training, leading to a 5X reduction in hardware cost; (2) The challenging and memory consuming operations of batch normalization were simplified by disabling low-impact runtime operations, rendering a 77% removal of the associated operations; (3) The expensive ADAM optimizer was replaced by a succinct momentum stochastic gradient descent optimizer suitable for integer implementation with an 11X reduction of the optimizer's computation; (4) The number of layers and channels were further minimized to reduce the computation load by 6X to 9X. Overall, a 6X reduction of training complexity, a 6.5X hardware cost reduction, and an 11X reduction of on-chip memory were achieved through the algorithm simplification with about a 3% loss of accuracy.

Fig. 2 shows the training sequence. Each training iteration consists of 7 unique phases (e.g. forward prop., loss cal.) with 5 phases for the generator and 4 phases for the discriminator. Each phase also contains 4 to 6 sub-tasks (e.g. Conv, FC, pooling, etc.). To avoid model collapsing, an adaptive training and model strength control scheme was implemented which ceases the training of discriminator if its strength is too high and adaptively increases the magnitude of the gradients during backpropagation. The training sequence is managed by an ASIC training management unit (TMU) shown in Fig. 2. A total of 41 training stages were implemented in the TMU as a finite state machine. Special operations such as pooling, sigmoid, data transpose etc. were handled by the dedicated hardware modules inside the TMU. Register files were used to store temporary weights and feature map outputs, bridging the throughput mismatch between SRAM and MAC arrays.

Fig. 3 shows the test chip architecture diagram including the TMU, a 10x10 time-domain (TD) MAC matrix, SRAM modules and supporting blocks. All the MAC operations of CNN and Transpose-CNN are performed by a TD MAC matrix to improve area and energy efficiency. The time pulses generated from digital-to-time converters (DTC) are processed by the subsequent multiplication, accumulation and activation all in time domain and are finally converted back into digital domain using time-to-digital converters (TDC). A special 16b time-pulse based time-domain accumulator (TD-ACC) is designed using four 4-b ring-based time accumulators [6] with carry propagation to realize accumulation efficiently. With the special TD-ACC, the TDC is only activated once every 25 MAC operations, removing 94% of the time and power overhead from the expensive TDC operations. Pushing all operations in time domain significantly reduces the cross-domain data conversion, rendering a 160X speed-up in MAC operation compared with previous counter-based TD designs [3]. The 8-b TD multiplication is partitioned into four 4-bit multiplications to improve the computation accuracy and speed.

Fig. 4 shows the detailed circuit design featuring a subthreshold (sub-vth) TD multiplier (TD-MUL) and a DTC-based linearization technique. The TD-MUL takes input time pulses and generates output pulses of the multiplication results. As in Fig. 4, the current starving PMOS transistor is pre-biased at subthreshold region and generates a delay equals to the multiplication results through charge accumulation at the gate with logarithmic addition, i.e. a multiplication is addition in log domain. Compared to the digital implementation, the implemented sub-vth multiplier renders a 4.3X reduction of area. However, as shown in simulation, significant nonlinearity is observed in sub-vth multiplication. The nonlinearity is compensated by a logarithmic encoding of DTC. As shown in both equation and the simulated waveforms in Fig. 4, the linearization technique elegantly removes nonlinearity with negligible overhead. After the multiplication, the resulting time pulses are sent into TD-ACC for accumulation of 25 cycles avoiding time-consuming digitalization as [3, 4, 5].

Simple TD ReLU function is also implemented at each CNN layer except the final layer which uses digital sigmoid function.

**Measurement**

Fig.5 shows the measured linearity from both the TD-MUL and TD-ACC. For the multiplier, although up to 4% error is seen in the result, most of the error is just a small scaling factor shift. Less than 1b error is observed in the TD-ACC design. We trained the GAN with 3 databases, i.e. a digit-MNIST, a fashion, and an emoji database [8-9]. The accuracy of the generated images with conditional GAN from 3 databases shows less than 1% error compared to the ideal integer 8-bit training on CPU and 3% compared with ideal floating-point training (1.6% comes from quantization loss and the rest from process variation of TD circuit). The chip is verified with supply voltages down to 0.7V with up to 5% degradation of accuracy compared with ideal GAN operation. Interestingly, a "self-healing" feature of GAN is observed, recovering most of the error loss from on-chip variations compared with no on-chip training. This intrinsic resiliency presents a merit for training empowered design using mixed-signal circuits. The chip consumes 39mW power with TD-MAC at 90MHz. The total training time of MNIST database takes 4.5 minutes which is 82X less than a high-performance CPU (2.6GHz Intel i7 Quad-core with a power of 197W). The die photo and comparison table with prior analog mixed-signal (AMS) designs are shown in Fig. 6. As most of existing AMS designs suffer from low throughput, this work achieves the highest throughput of 18~5400X [2-5, 7] with similar efficiency. In addition, a low-cost 8-bit on-chip training was realized for AMS design on the very challenging GAN operation.



Fig. 1 GAN algorithm and hardware adaption in this work.



Fig. 3 Top-level architecture diagram with MAC array and TD-Accumulator, TD ReLU circuit, and TD MAC unit.



Fig. 2 GAN training sequence and ASIC TMU core design.



Fig. 4 TD sub-vth multiplier and linearization technique.



Fig. 5 Measurement results.



Fig. 6 Die photo and comparison table.

**References** [1] A. Radford, et al., *arXiv*, 2015. [2] M. Liu, et. al. *CICC*'17. [3] N. Cao, et al., *ISSCC*'19. [4] A. Sayal, et al., *ISSCC*'19. [5] E. H. Lee, et al., *ISSCC*'16. [6] Z. Chen, et al., *ISSCC*'19. [7] K. Yoshioka, et al., *VLSI*'18. [8] FASHION database, https://www.kaggle.com/zalando-research/fashionmnist. [9] EMOJI database, https://getemoji.com/.