

Analog circuits for mixed-signal neuromorphic computing architectures in 28 nm FD-SOI technology

Ning Qiao

Institute of Neuroinformatics

University of Zurich and ETH Zurich

Zurich, Switzerland

Email: qiaoning@ini.uzh.ch Giacomo Indiveri Institute of Neuroinformatics

University of Zurich and ETH Zurich

Zurich, Switzerland

Email: giacomo@ini.uzh.ch

Abstract—Developing mixed-signal analog-digital neuromorphic circuits in advanced scaled processes poses significant design challenges. We present compact and energy efficient sub-threshold analog synapse and neuron circuits, optimized for a 28 nm FD-SOI process, to implement massively parallel large-scale neuromorphic computing systems. We describe the techniques used for maximizing density with mixed-mode analog/digital synaptic weight configurations, and the methods adopted for minimizing the effect of channel leakage current, in order to implement efficient analog computation based on pA-nA small currents. We present circuit simulation results, based on a new chip that has been recently taped out, to demonstrate how the circuits can be useful for both low-frequency operation in systems that need to interact with the environment in real-time, and for high-frequency operation for fast data processing in different types of spiking neural network architectures.

Keywords—Sub-threshold analog, neuromorphic computing, low leakage, spiking neural networks, low power, IoT, ReLU.

I. INTRODUCTION

As computing systems implemented in advanced Very Large Scale Integration (VLSI) processes are facing more and more stringent requirements, mainly related power consumption, circuit designers and system engineers are starting to explore solutions that are alternative to the standard approach of the von Neumann computing paradigm. Neuromorphic computing represents one of these approaches, that proposes to use brain inspired neural network architectures for signal and data processing [1], [2]. One of the main features of neuromorphic computing architectures is their co-localization of memory and computing elements [3]: the synapse elements in these neuromorphic architectures represent at the same time the site of memory (that store the synaptic weight value), and the site of computation (which in the simplest case is the multiplication of the input signal with the stored weight). From the computing architecture point of view, this has the large advantage of avoiding the von Neumann bottleneck problem [3], [4]. Since memory transfer is typically the largest power consuming operation, this approach represents already a large step toward the development of ultra-low power computing systems. Another brain-inspired approach that is extremely

helpful in reducing power consumption and increasing circuit density is that of representing signals using pulse-frequency modulation: if input/output signals are represented as pulses (spikes), the multiplication operation between input signals and synaptic weights reduces to a gating operation at the synapse level, that typically produces a weighted current at the arrival of the pre-synaptic spike that is integrated by the post-synaptic neuron. The higher the frequency of the input spikes, the larger the integrated value that the neuron sees. Furthermore, if many synapses receive input spikes in parallel, the weighted sum operation is implemented directly at the input node of the post-synaptic neuron by Kirchhoff's current law. Power consumption can be further reduced by implementing this spike- or event-based signal representation using asynchronous logic. In this case, the representation is denoted as Address-Event Representation (AER). Given these features, and given that this representation is also optimal for transmitting signals across long distances or chip boundaries, most of the recent state-of-the-art neuromorphic computing approaches use AER [5]–[9]. The last step that can be taken to further minimize power consumption is that of adopting a mixed-signal design approach, and using analog circuits that directly exploit the physics of the devices to implement the desired neural network computational primitives [1]. As these primitives are mainly composed of exponential and logarithm functions, the best approach to follow is that of using sub-threshold analog circuits [10].

In this paper we present sub-threshold analog synapse and neuron circuits that have been designed to implement large-scale multi-neuron multi-core neuromorphic computing architectures using a 28 nm Fully-Depleted Silicon on Insulator (FD-SOI) process [11]. We show how it is possible to implement complex bio-physically realistic synaptic and neural dynamics using ultra-low power compact analog circuits in advanced scaled processes, by using mismatch-reducing and leakage-canceling techniques. Although the circuits proposed can be configured to run at speeds that are much higher than biologically plausible ones, they have been optimized for reproducing time-constants that can be as long as tens or hundreds of milli-seconds. In this way these circuits can be embedded in massively parallel neural architectures optimally suited for processing live streaming data

coming from natural sensory signals, such as auditory signals representing speech, visual signals representing gestures, or bio-signals measured from real neurons or muscles. Although the time constants of the individual computing elements are long, the asynchronous nature of the AER protocol used to process incoming data ensures that latency and response time at the system level are extremely fast, i.e., ranging from hundreds of nano-seconds to milliseconds, depending on the complexity of the networks implemented. The circuits proposed represent a natural extension of similar circuits already fabricated, tested, and validated in less advanced processes [7], currently being used for implementing brain-machine interfaces [12], [13], deep network prototypes [14], and autonomous driving robotic applications [15].

II. SUB-THRESHOLD ANALOG/DIGITAL SYNAPSE & NEURON CIRCUITS

Event-based synaptic circuits typically translate pre-synaptic voltage pulses into post-synaptic currents and source them into the target neuron circuit with a gain that corresponds to the synaptic weight. In Fig. 1 we show the schematic diagram of a circuit that comprises 64 programmable synapse blocks with common leakage compensation and temporal dynamics blocks.

Each synapse has 4 current branches with shared analog bias settings $w_{ht}x!$, that can be programmed with a 10-bit temperature compensated bias generator [16]. In the circuit diagrams, all signals ending with a "!" represent programmable bias settings. Upon the arrival of a pre-synaptic spike, the input AER event is decoded to select the activation of one or more of the synapse branches, therefore allowing to set as synaptic weight one of 16 possible analog currents.

In advanced scaled processes, such as the one used in this work, the off-channel leakage current of each branch is of the order of pico-Amperes or more. Considering that there are 64 synapse blocks with a total of 256 current branches, the total leakage current produced is non-negligible. The "Leakage Canceling" block of Fig. 1 attempts to produce the same leakage current, with a 4-to-1 copy of 16 leak cells with 4 branches each that are biased with the same settings used for the synapse blocks. All dark currents from the leak cells branches are summed and copied by $M_{L5} - M_{L12}$ of Fig. 1 to produce the I_{leak} current. The "Leakage Canceling" block subtracts this current from the total synaptic current I_{sum} , to produce a resulting output current I_{wht} which represents the compensated weighted contribution of all 64 afferent synapses. If one assumes that all synapses share the same dynamics, it is possible to use the superposition principle and model the temporal dynamics of all synapses using one single low pass filter. This is indeed the case for the circuits of Fig. 1, that use the Differential Pair Integrator (DPI) block on the bottom right of the figure to implement a current-mode low-pass filter that exhibits synaptic dynamics [17]. The time constant of this circuit is directly proportional to its capacitance, and inversely proportional to the current through $M_{D5} - M_{D6}$. In order to obtain large time constants, while keeping the size of the capacitors to a minimum, it is necessary to generate bias currents that can be as small as pico-Amperes. To achieve this goal in this process, we had

to resort to the use of the "pseudo-cascode" split-transistor sub-threshold technique [18] (e.g., see also $M_{L1} - M_{L4}$ in Fig. 1). The diode connected transistors $M_{D2} - M_{D3}$ are added to reduce the V_{D5} of M_{D1} so as to reduce its early effect, and to improve the circuit's linear performance. The *NMDA* block of Fig. 1 ($M_{N1} - M_{N3}$) models the voltage-gating mechanisms of N-Methyl-D-Aspartate (NMDA) synapses, which can be useful for coincidence detection and precise spike-timing signal processing strategies in spike-based neuromorphic computing applications. The area of the synapse block layout is $3 \mu m^2$, and the active area of the DPI circuit is $12.5 \mu m^2$. The DPI capacitor, implemented using a MIMCAP structure covering the neighboring circuits, measures 1 pF.

The circuit that implements the silicon neuron functionality is shown in Fig. 2. This is a current-mode circuit composed of multiple "compartments" or blocks. The *LEAK* block models the neuron's passive leak conductance, producing exponential sub-threshold dynamics in response to constant input currents. The *AHP* block models the generation of the after hyperpolarizing current in real neurons, responsible for their spike-frequency adaptation behavior. The *Na+* and *K+* blocks model the effect of Sodium and Potassium channels, responsible for generating action-potentials (spikes) in real neurons. The *HS* block implements handshaking with following encoder block for encoding spike events following AER protocol. We used an optimized Traff's current comparator (see *CC* box in the *Na+* block) to make an accurate comparison between the neuron I_{mem} current and a programmable I_{ref} threshold current, which sets the neuron's spiking threshold. Current limitation transistors (M_{C7}, M_{C8}) are included to reduce static power consumption. We used the same split-transistor sub-threshold technique used in the synapse circuits for enhanced current-mirror operation and for precise control of small currents. We added also several reset transistors, such as M_{L3} and M_{NA1} to further reduce power consumption during the spike reset phase. The active area of the neuron is $20 \mu m^2$; the neuron capacitance is also implemented using a MIMCAP structure and measures approximately 1.5 pF. Fig. 4 shows the expected response of the neuron circuit to a constant input current. By tuning the biases that control the neuron's integration time constant, firing threshold, refractory period and spike-frequency adaptation dynamics, the proposed circuit can reproduce a wide range of spiking behaviors [7].

III. SIMULATION RESULTS

In Fig. 4 we show circuit simulation results of both synaptic and neuron currents, while they are being stimulated by a 100 Hz input spike train with pulse width of 200 us and pulse amplitude of 10 nA. The spiking threshold reference current was set to 20 nA. As shown, both synaptic and neuron output signals exhibit biologically plausible time temporal dynamics with time constants of the order of milli-seconds.

Fig. 5 shows the combined synapse-neuron transfer function, consisting of the neuron output firing rate as a function of the synapse input firing rate, for three different synaptic efficacy levels (which are inversely proportional to the $dpi_tau!$ bias setting [7]). In these simulations the synaptic pulse width was set to 1 ms, the pulse amplitude to 10 nA, the neuron spiking

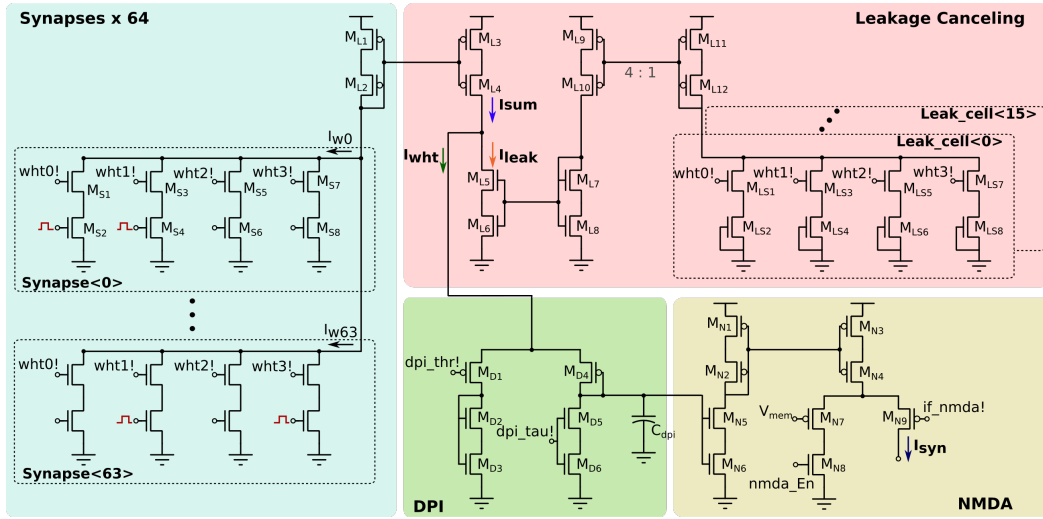


Fig. 1. Schematic diagram of synapse and integrator circuits.

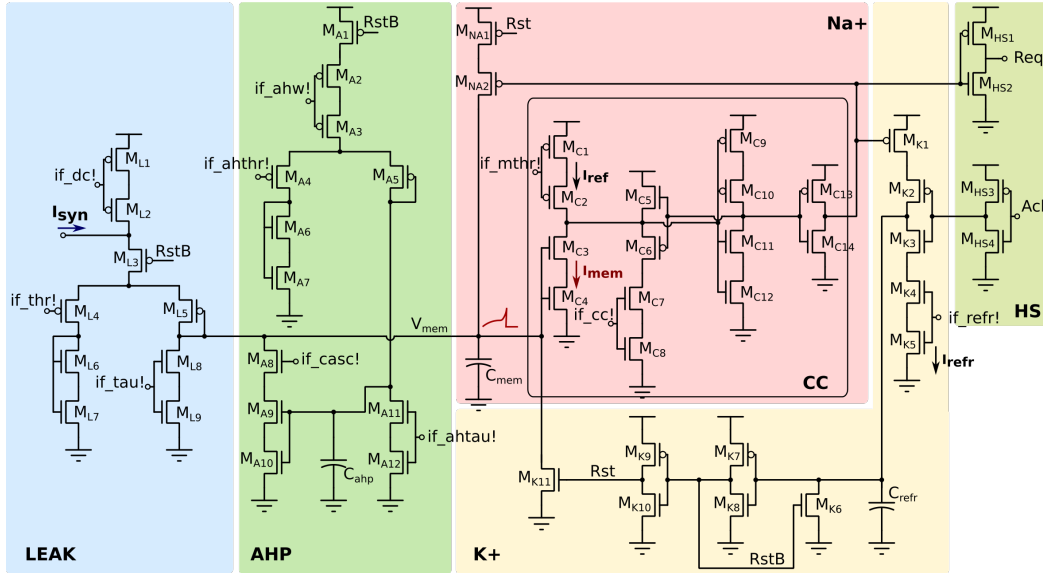


Fig. 2. Schematic diagram of an analog Integrate-and-Fire (I&F) neuron.

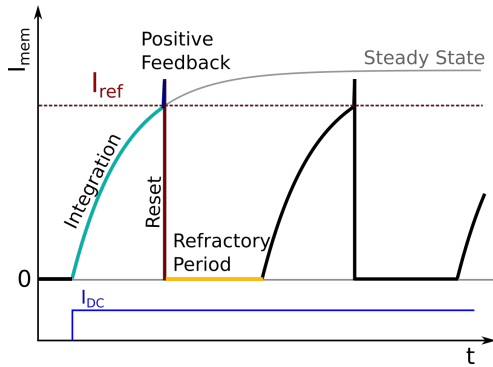


Fig. 3. Membrane current trace over time.

threshold reference to 20 nA, and its refractory period to 5 ms. Setting a refractory period to long intervals forces the neuron circuits to saturate at low frequencies, therefore reproducing the behavior of real neurons and limiting the bandwidth requirement for spiking neural networks. By changing the bias settings that affect the neuron refractory period, it is possible to configure the circuit to operate in a linear manner over a much larger range of output frequencies. Fig. 6 shows the same transfer function with bias settings tuned to reproduce the function of Rectified Linear Unit (ReLU) units, typically used in deep-networks, over a wide range of fast input/output frequencies, thus making these circuits also suitable for high-speed spiking deep network models. The data of Fig. 6 was obtained by setting the width of the input synaptic pulses to 50 us, their amplitude to 10 nA,

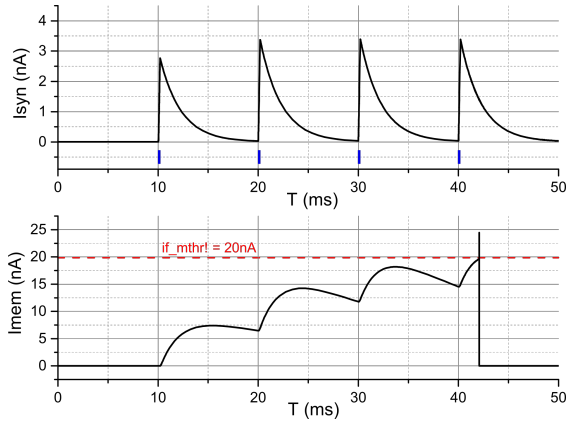


Fig. 4. Synapse and neuron response to a 100Hz spike train.

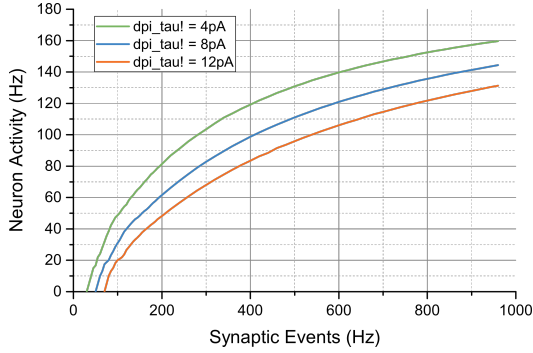


Fig. 5. Combined synapse-neuron transfer function with refractory period of 5ms for synaptic input firing rate from 0-1k Hz.

and the refractory period to micro-seconds. The slope of the transfer function can be modulated by changing the gain of the synapse/neuron DPI circuits, via the corresponding `_thr!` and `_tau!` bias settings of Fig. 1 and of Fig. 2.

IV. CONCLUSION

We proposed novel compact analog/digital neuromorphic circuits optimized for minimizing leakage currents and producing

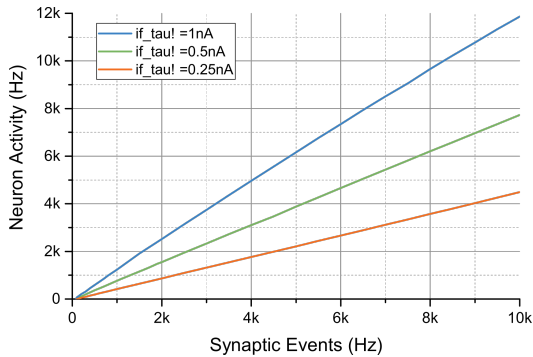


Fig. 6. Combined synapse-neuron transfer function with short refractory period.

long time constants in a 28 nm FD-SOI process, that can enable neuromorphic architectures to process natural sensory signals in real-time. This allows the design of massively parallel ultra-low power mixed-signal neuromorphic computing architectures that would not be affected by the von Neumann bottleneck, as they would not need to time-multiplex shared neuron/synapse circuits and transfer their state information to separate memory blocks. We showed how the circuits proposed reproduce the synapse and neural dynamics expected from theory and can be used to reproduce both biologically realistic dynamics, or fast ReLU transfer functions. They are fully compatible with spike-based learning algorithm/circuits [7] and can be readily integrated in the next generation of large multi-neuron, multi-core neuromorphic architectures.

ACKNOWLEDGMENT

This work is supported by the EU ICT grant “NeuRAM³” (687299).

REFERENCES

- [1] C. Mead, “Neuromorphic electronic systems,” *Proceedings of the IEEE*, vol. 78, no. 10, pp. 1629–36, 1990.
- [2] G. Indiveri and T. Horiuchi, “Frontiers in neuromorphic engineering,” *Frontiers in Neuroscience*, vol. 5, no. 118, pp. 1–2, 2011. [Online]. Available: http://www.frontiersin.org/neuromorphic_engineering/10.3389/fnins.2011.00118/full
- [3] G. Indiveri and S.-C. Liu, “Memory and information processing in neuromorphic systems,” *Proceedings of the IEEE*, vol. 103, no. 8, pp. 1379–1397, 2015. [Online]. Available: http://ncs.ethz.ch/pubs/pdf/Indiveri_Liu15.pdf
- [4] J. Backus, “Can programming be liberated from the von neumann style?: a functional style and its algebra of programs,” *Communications of the ACM*, vol. 21, no. 8, pp. 613–641, 1978. [Online]. Available: <http://doi.acm.org/10.1145/359576.359579>
- [5] B. V. Benjamin, P. Gao, E. McQuinn, S. Choudhary, A. R. Chandrasekaran, J. Bussat, R. Alvarez-Icaza, J. Arthur, P. Merolla, and K. Boahen, “Neurogrid: A mixed-analog-digital multichip system for large-scale neural simulations,” *Proceedings of the IEEE*, vol. 102, no. 5, pp. 699–716, 2014.
- [6] S. Furber, F. Galluppi, S. Temple, and L. Plana, “The SpiNNaker project,” *Proceedings of the IEEE*, vol. 102, no. 5, pp. 652–665, May 2014.
- [7] N. Qiao, H. Mostafa, F. Corradi, M. Osswald, F. Stefanini, D. Sumislawska, and G. Indiveri, “A re-configurable on-line learning spiking neuromorphic processor comprising 256 neurons and 128k synapses,” *Frontiers in Neuroscience*, vol. 9, no. 141, 2015. [Online]. Available: http://www.frontiersin.org/neuromorphic_engineering/10.3389/fnins.2015.00141/abstract
- [8] J. Park, T. Yu, S. Joshi, C. Maier, and G. Cauwenberghs, “Hierarchical address event routing for reconfigurable large-scale neuromorphic systems,” *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2016.
- [9] P. A. Merolla, J. V. Arthur, R. Alvarez-Icaza, A. S. Cassidy, J. Sawada, F. Akopyan, B. L. Jackson, N. Imam, C. Guo, Y. Nakamura, B. Brezzo, I. Vo, S. K. Esser, R. Appuswamy, B. Taba, A. Amir, M. D. Flickner, W. P. Risk, R. Manohar, and D. S. Modha, “A million spiking-neuron integrated circuit with a scalable communication network and interface,” *Science*, vol. 345, no. 6197, pp. 668–673, Aug 2014. [Online]. Available: <http://www.sciencemag.org/content/345/6197/668>
- [10] S.-C. Liu, J. Kramer, G. Indiveri, T. Delbruck, and R. Douglas, *Analog VLSI: Circuits and Principles*. MIT Press, 2002. [Online]. Available: http://ncs.ethz.ch/pubs/pdf/Liu_et al02b.pdf

- [11] N. Qiao and G. Indiveri, "Scaling mixed-signal neuromorphic processors to 28nm fd-soi technologies," in *Biomedical Circuits and Systems Conference, (BioCAS), 2016*. IEEE, 2016, pp. 552–555. [Online]. Available: <http://ncs.ethz.ch/pubs/pdf/QiaoIndiveri16.pdf>
- [12] F. Corradi and G. Indiveri, "A neuromorphic event-based neural recording system for smart brain-machine-interfaces," *Biomedical Circuits and Systems, IEEE Transactions on*, vol. 9, no. 5, pp. 699–709, 2015. [Online]. Available: http://ncs.ethz.ch/pubs/pdf/Corradi_Indiveri15.pdf
- [13] F. Boi, T. Moraitis, V. De Feo, F. Diotalevi, C. Bartolozzi, G. Indiveri, and A. Vato, "A bidirectional brain-machine interface featuring a neuromorphic hardware decoder," *Frontiers in Neuroscience*, vol. 10, p. 563, 2016. [Online]. Available: <http://journal.frontiersin.org/article/10.3389/fnins.2016.00563>
- [14] G. Indiveri, F. Corradi, and N. Qiao, "Neuromorphic architectures for spiking deep neural networks," in *Electron Devices Meeting (IEDM), 2015 IEEE International*. IEEE, Dec. 2015, pp. 4.2.1–4.2.14. [Online]. Available: http://ncs.ethz.ch/pubs/pdf/Indiveri_et al15.pdf
- [15] M. B. Milde, H. Blum, A. Dietmüller, H. Blum, D. Sumislawaska, J. Conradt, G. Indiveri, and Y. Sandamirskaya, "Obstacle avoidance and target acquisition for robot navigation using a mixed signal analog/digital neuromorphic processing system," *Frontiers in Neuroscience*, 2017.
- [16] M. Yang, S.-C. Liu, C. Li, and T. Delbruck, "Addressable current reference array with 170db dynamic range," in *Circuits and Systems (ISCAS), 2012 IEEE International Symposium on*. IEEE, 2012, pp. 3110–3113.
- [17] C. Bartolozzi and G. Indiveri, "Synaptic dynamics in analog VLSI," *Neural Computation*, vol. 19, no. 10, pp. 2581–2603, Oct 2007. [Online]. Available: http://ncs.ethz.ch/pubs/pdf/Bartolozzi_Indiveri07.pdf
- [18] V. Saxena and R. J. Baker, "Compensation of CMOS op-amps using split-length transistors," in *Circuits and Systems (MWSCAS), 2008 IEEE 51st International Midwest Symposium on*. IEEE, 2008, pp. 109–112.