# Mapping of Local and Global Synapses on Spiking Neuromorphic Hardware

Anup Das[‡†], Yuefeng Wu[†], Khanh Huynh[†], Francesco Dell'Anna[§], Francky Catthoor[§], Siebren Schaafsma[†] [†]*Stichting IMEC Nederland*, [§]*IMEC Belgium*, [‡]*Department of ECE, Drexel University, Philadelphia, USA*

Correspondance email: akdas@ieee.org

## Abstract

Spiking Neural Networks (SNNs) are widely deployed to solve complex pattern recognition, function approximation and image classification tasks. With the growing size and complexity of these networks, hardware implementation becomes challenging because scaling up the size of a single array (crossbar) of fully connected neurons is no longer feasible due to strict energy budget. Modern neromorphic hardware integrates small-sized crossbars with time-multiplexed interconnects. Partitioning SNNs becomes essential in order to map them on neuromorphic hardware with the major aim to reduce the global communication latency and energy overhead. To achieve this goal, we propose our instantiation of particle swarm optimization, which partitions SNNs into local synapses (mapped on crossbars) and global synapses (mapped on time-multiplexed interconnects), with the objective of reducing spike communication on the interconnect. This improves latency, power consumption as well as application performance by reducing inter-spike interval distortion and spike disorders. Our framework is implemented in Python, interfacing CARLsim, a GPU-accelerated application-level spiking neural network simulator with an extended version of Noxim, for simulating time-multiplexed interconnects. Experiments are conducted with realistic and synthetic SNN-based applications with different computation models, topologies and spike coding schemes. Using power numbers from in-house neuromorphic chips, we demonstrate significant reductions in energy consumption and spike latency over PACMAN, the widely-used partitioning technique for SNNs on SpiNNaker.

## I. INTRODUCTION

Spiking Neural Networks (SNNs) [1] are powerful and biologically realistic computation models, inspired by the dynamics of human brain. From implementation viewpoint, SNNs are collection of neurons that communicate by sending short pulses (spikes) across connections

(synapses) to other neurons. SNNs are trained to perform variety of tasks, where the training process involves adjusting connection strengths between neurons. SNNs are increasingly being deployed to solve complex tasks such as pattern recognition, function approximation and image classification. Another reason for widespread success of SNNs are their efficient VLSI implementations, such as TrueNorth [2], CxQuad [3] and SpiNNaker [4] among others. Our work is based on CxQuad, which is an analog neuromorphic hardware with 1024 neurons clustered into four crossbars of 256 neurons each. The framework proposed in this work can be extended with reasonable effort reusing the same concepts to other neuromorphic architectures.

With increasing deployment of SNNs in complex scenarios, the field is rapidly maturing in terms of applications, algorithms, computation models and hardware. Although significant research activities in these domains are being conducted separately, the task of a system designer is to build synergy between these domains i.e., to map realistic SNN-based applications on neuromorphic hardware. Although some efforts are made in this direction, the published approaches still remain mostly ad-hoc [5], specific to a given hardware [6] and are often limited to trivial applications [7]. The closest technique to our approach is that of `PACMAN` [8], which is used to map SNNs on SpiNNaker neuromorphic hardware. Following are the limitations of `PACMAN` that motivated this paper. First, `PACMAN` is primarily targeted for architectures supported on SpiNNaker such as Deep Belief Networks [9] and Convolution Neural Networks [10]. `PACMAN` offers limited flexibility to implement evolving computation models such as the Liquid State Machine (LSM) [11] or Hierarchical Temporal Memory (HTM) [12]. Second, `PACMAN` requires significant modification to support clustered neuromorphic architectures with local and global synapses. This is because the `PACMAN` tool is natively developed for SpiNNaker hardware only, where computing elements are ARM cores with conventional Von-Neumann architecture. Third, `PACMAN` determines neuron mapping without considering spike latency related performance distortions and interconnect energy consumption.

We propose a systematic approach to map trained SNNs on a neuromorphic hardware. Fundamental to this is our instantiation of Particle Swarm Optimization (PSO) [13], which partitions a given SNN into local and global synapses. Local synapses are mapped on fully connected crossbars and global synapses on time-multiplexed interconnect between the crossbars. The objective of this optimization is to minimize spike communication on the time-multiplexed interconnect, saving energy and improving performance (such as accuracy) of the overlaying application by reducing spike disorder count and inter-spike interval distortion. Our approach
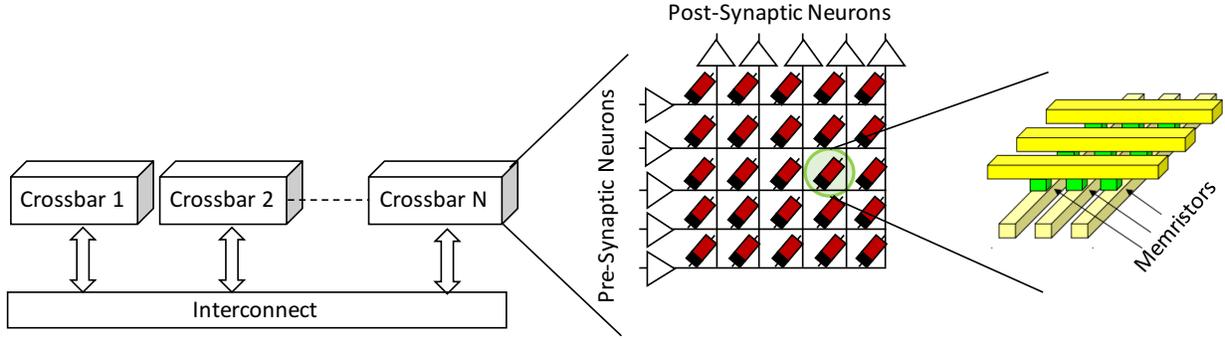
Fig. 1. A reference neuromorphic hardware

can be used with simulators as well as with real hardware. Although energy consumption for CxQuad hardware is used to demonstrate the power/performance improvement, our conceptual approach is however, generic and can be used for a range of devices with memristor-based synaptic elements.

**Contributions:** Following are our novel contributions

- a systematic framework for partitioning and mapping a trained SNN on neuromorphic hardware;
- introduction of performance metrics for mapping SNNs on neuromorphic hardware;
- a PSO-based partitioning of SNN into local and global synapses to reduce spike communication and congestion on time-multiplexed interconnect;
- a Python-based open-source framework for simulating SNN on neuromorphic hardware; and
- thorough experimentation with realistic applications, demonstrating the advantage of our proposed approach.

We build our framework in Python interfacing `CARLsim` [14], a GPU-accelerated application-level spiking neural network simulator with an extended version of `Noxim` [15], for simulating time-multiplexed interconnect. Results demonstrate the energy, latency and performance gains in the global communication network using our approach. The remainder of this paper is organized as follows. Neuromorphic platform description is provided in Section II together with the introduced metrics. PSO-based partitioning is introduced next in Section III. Our proposed systematic framework is described in Section IV. Results and discussions are provided in Section V. Finally, the paper is concluded with future outlook in Section VI.

## II. Neuromorphic Hardware and Related Metrics

Figure 1 shows the general architecture of a modern neuromorphic hardware. The architecture consists of multiple crossbars of fully connected neurons (shown in part b). Implementation wise, a crossbar is a 3D arrangement of nanowires, with pre-synaptic neurons connected to bottom nanowires and post-synaptic neurons to the top nanowires (or vice-verse). Each crosspoint of a top and bottom nanowire, is a two-terminal memristor nanodevice (shown in part c). The synaptic connection strength between a pair of pre- and post-synaptic neurons is encoded as resistance value of the memristor connecting them and is adjusted by regulating the flow of current through it. Crossbars communicate with each other using an interconnect (shown in part a). Traffic on the interconnect is time-multiplexed. Several interconnect alternatives have been explored in literature for neuromorphic computing. The commonly used ones are NoC-tree (CxQuad) and NoC-mesh (TrueNorth, HiCANN).

Analogous to mammalian brain, synapses of a SNN can be classified into local and global synapses based on the distance information (spike) is conveyed. Local synapses are short distance links, where pre- and post-synaptic neurons are located in the vicinity. Global synapses are those where pre- and post-synaptic neurons are farther apart. To reduce power consumption of the hardware implementation of SNNs, two principles are widely adopted in the community:

- the number of point-to-point local synapses is limited to a reasonable dimension (size of a crossbar); and
- instead of point-to-point global synapses (which are of long distance) as found in a mammalian brain, the hardware implementation usually consists of time-multiplexed interconnect shared between global synapses.

CxQuad for example, consists of four crossbars, each with 128 pre- and 128 post-synaptic neurons implementing a full 16K (128x128) local synapses per crossbar. The crossbars are interconnected using a NoC-tree, which time-multiplexes global synaptic connections. The spike communication protocol for the global synapse interconnect is `Address Event Representation` (AER) [16]. An example is shown in Figure 2 to explain the principles behind AER. Here four neurons in the input group (a crossbar) spikes at time 3, 0, 1 and 2 time units, respectively. The encoder encodes these four spikes in order to be communicated on the global synapse interconnect. As can be clearly seen from this figure, a spike is encoded uniquely on the global synapse interconnect in terms of its source and time of spike.
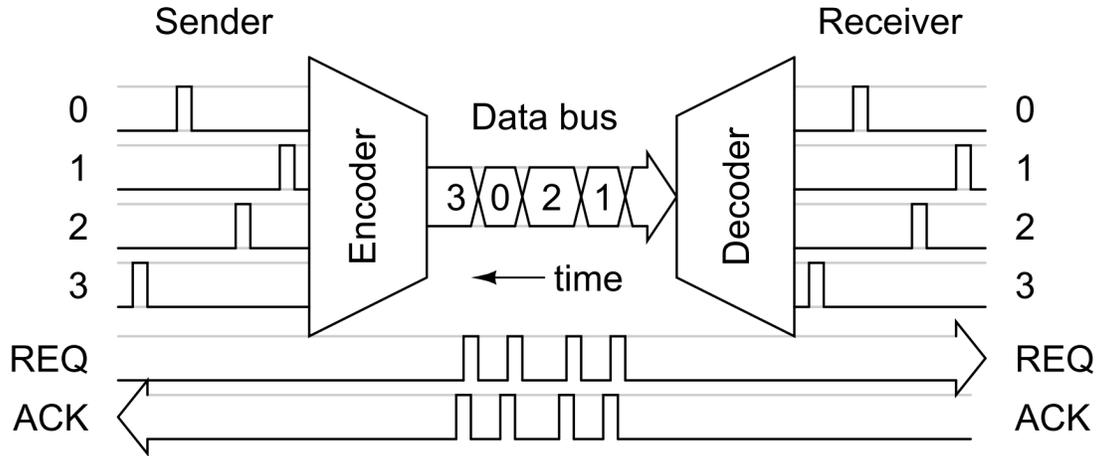
Fig. 2.   An example showing AER protocol (adapted from [16]).

Following are metric for global synapse interconnect.

- Conventional metric for the global synapse interconnect
  - *Latency*: The delay between transmitting a spike packet (AER) by the encoder and receiving the packet by the decoder.
  - *Energy*: The energy consumed by spike communication on the global synapse interconnect.
- Introduced metric for SNN performance on hardware
  - *Spike disorder count*: This is a measure of information loss in a SNN and is calculated based on the difference in spike order between sender and receiver neurons. An example is provided to explain this. Let us assume that neuron A and B need to communicate spikes to neuron C with spikes from A to be received before the spike from neuron B. Also let A, B and C are all mapped to different crossbars. If it happens that the crossbar with B is arbitrated to occupy the interconnect prior to crossbar with A, spikes from B will be received at C before the spike from A causing a disorder of spikes and potential information loss.
  - *Inter-spike distortion*: This is a measure for information distortion in temporally coded SNN and is measured by the difference between the sender neuron's inter-spike interval and receiver neuron's inter-spike intervals. Inter-spike distortion is attributed to spike congestion on the global synapse interconnect, causing some spike packets to be delayed

than others.

Based on the above discussions, the problem we aim to solve in this paper is as follows. Partition a given SNN into local and global synapses, where local synapses are mapped on the crossbar and global synapses on the global synapse interconnect. The objective of this optimization problem is to reduce spike communication (congestion) on the interconnect, which minimizes interconnect energy consumption, spike latency, spike disorder count and inter-spike distortion.

## III. PSO-BASED MAPPING OF TRAINED SNN ON NEUROMORPHIC HARDWARE

In this work we propose to use evolutionary techniques in order to solve the optimization problem of reducing spike communication on the interconnect, saving energy and improving SNN performance. We use particle swarm optimization (PSO) [13], an evolutionary computing technique inspired by social behaviors such as bird flocking and fish schooling. Evolutionary computing techniques are efficient in avoiding to stuck at local optima. Additionally, PSO is computationally less expensive with faster convergence compared to its counterparts such as genetic algorithm (GA) or simulated annealing (SA).

In general, PSO finds the optimum solution to a fitness function $F$. Each solution is represented as a particle in the swarm. Each particle has a velocity with which it moves in the search space to find the optimum solution. During the movement, a particle updates its position and velocity according to its own experience (closeness to the optimum) and also experience of its neighbors. We introduce the following notations for PSO.

$$D = \text{dimensions of the search space}$$

$$n_p = \text{number of particles in the swarm}$$

$$\boldsymbol{\Theta} = \{\theta_l \in \mathbb{R}^D\}_{l=0}^{n_p-1} = \text{positions of particles in the swarm}$$

$$\mathbf{V} = \{\mathbf{v}_l \in \mathbb{R}^D\}_{l=0}^{n_p-1} = \text{velocity of particles in the swarm}$$

Position and velocity updates are performed according to

$$\boldsymbol{\Theta}(t+1) = \boldsymbol{\Theta}(t) + \mathbf{V}(t+1) \tag{1}$$

$$\mathbf{V}(t+1) = \mathbf{V}(t) + \varphi_1 \cdot \left( P_{\text{best}} - \boldsymbol{\Theta}(t) \right) + \varphi_2 \cdot \left( G_{\text{best}} - \boldsymbol{\Theta}(t) \right)$$

where $t$ is the iteration number, $\varphi_1, \varphi_2$ are constants and $P_{\text{best}}$ (and $G_{\text{best}}$) is the particles own (and neighbors) experience. Figure 3 shows the iterative solution of PSO, where position and
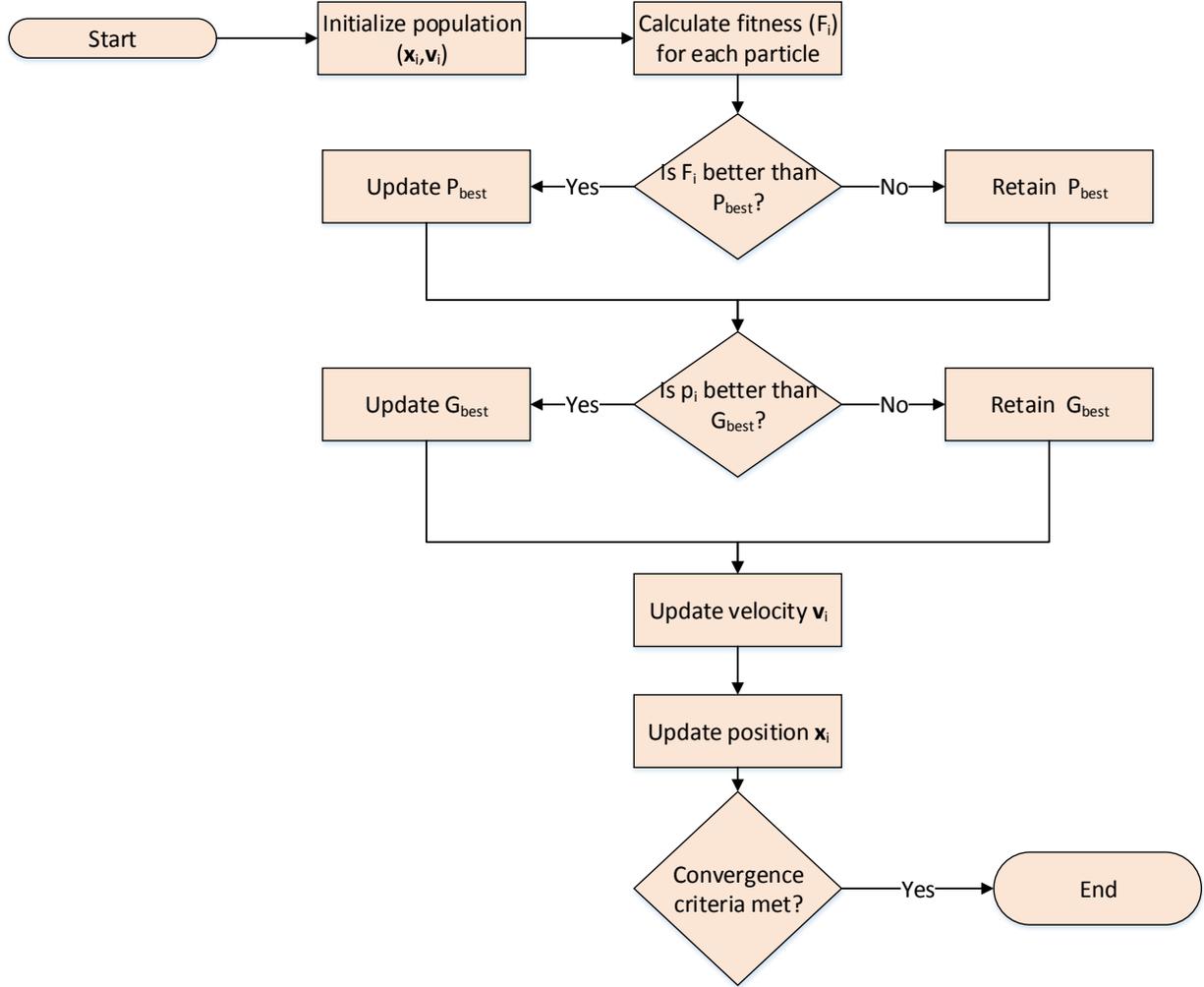
Fig. 3. Particle swarm optimization to find an optimum solution.

velocity are updated until a predefined convergence is achieved. In the following we describe the transformation of our partitioning problem to the PSO domain.

Let us consider a SNN with $N$ neurons organized in a topology, specific to a given application (more details in Section V). The SNN can be represented as a graph $\mathcal{G} = (A, S)$, where $A = \{\mathbf{a_i} \mid 0 \leq i \leq N - 1\}$ is the set of nodes of the graph, representing neurons and $S = \{s_{i,j}\}$ is the set of synapses representing connections between the neurons. Each synapse $\mathbf{s_{i,j}}$ is a tuple $\langle \mathbf{a_i}, \mathbf{a_j}, \mathbf{T_{i,j}} \rangle$, where $\mathbf{a_i}$ is the pre-synaptic neuron, $\mathbf{a_j}$ is the post-synaptic neuron and $\mathbf{T_{i,j}} = \{t_i^l \mid 0 \leq l \leq L_i\}$ are the spike times of the pre-synaptic neuron $\mathbf{a_i}$. This graph represents initial specification of a trained SNN in terms of synaptic weights and spike times. This graph is generated from `CARLsim` [14].

Let $C$ be the number of crossbars in the architecture, with $N_c$ being the maximum number of neurons per crossbar. The specification $C$ is usually provided by a designer for a given architecture. In Section V-C we discuss how to obtain this specification for a given set of applications. Let $x_{i,k} = \{0, 1\}$, be the variable indicating if neuron $\mathbf{a_i}$ is allocated to crossbar $\mathbf{c_k}$, where $0 \leq k \leq C$. These variables $x_{i,k}$'s are dimensions of our PSO with $D = N \cdot C$. In order to transform real-valued $x_{i,k}$ to binary values (for 0-1 assignments), the velocity and position updates need to be binarized. This is achieved using the following set of equations

$$\hat{\mathbf{v}_{i,k}} = \texttt{sigmoid}(\mathbf{v}_{i,k}) = \frac{1}{1 + e^{-\mathbf{v}_{i,k}}} = \begin{cases} 0 \text{ if } \mathbf{v}_{i,k} < 0 \\ 1 \text{ otherwise} \end{cases} \tag{2}$$

$$\hat{x_{i,k}} = \begin{cases} 0 \quad \text{if } \texttt{rand()} < \hat{\mathbf{v}_{i,k}} \\ 1 \quad \text{otherwise} \end{cases} \tag{3}$$

Constraints for our PSO are as follows

- every neuron is allocated to one crossbar only i.e.,

$$\sum_k \hat{x_{i,k}} = 1 \quad \forall i \tag{4}$$

- assignment must satisfy the dimensions of crossbar i.e.,

$$\sum_i \hat{x_{i,k}} \leq N_c \quad \forall k \tag{5}$$

To determine the number of spikes communicated between crossbars $k_1$ and $k_2$, we define two sets $K_1$ and $K_2$, where $K_1$ ($K_2$) is the set of neurons allocated to crossbar $k_1$ ($k_2$). Elements of these sets are determined as follows

$$K_1 = \{(i \cdot \hat{x_{i,k_1}})\}_{i=0}^{N-1} \text{ and } K_2 = \{(j \cdot \hat{x_{j,k_2}})\}_{j=0}^{N-1} \tag{6}$$

The total spikes between crossbars $k_1$ and $k_2$ is

$$\text{spikes}(k_1, k_2) = \begin{cases} 0 & \text{if } k_1 = k_2 \\ \sum_{\substack{i,j \\ i \in K_1 \\ j \in K_2}} \mathbf{T_{i,j}} & \text{otherwise} \end{cases} \tag{7}$$

Total spikes on the global synapse interconnect is

$$F = \sum_{k_1, k_2} \text{spikes}(k_1, k_2) \tag{8}$$

The objective of the PSO is to minimize $F$ satisfying constraints Equation 4–5. The assignment information (i.e., outcome of PSO) is used by the global synapse simulator (discussed in Section IV) to compute global statistics of spike communication. This statistics are then used to compute
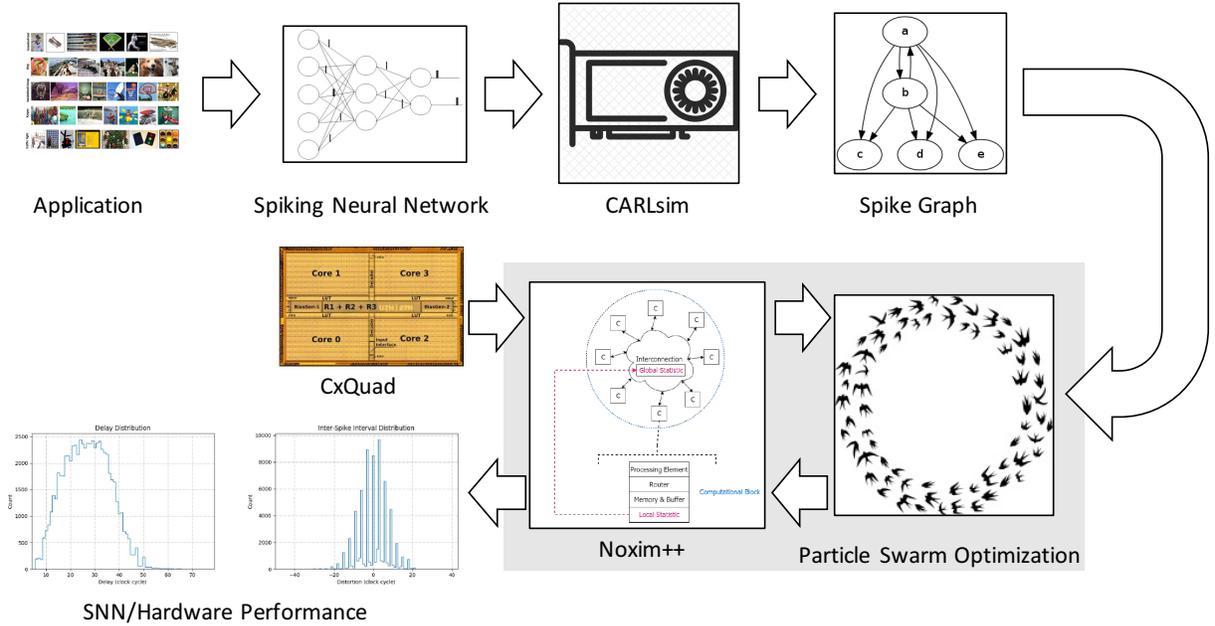
Fig. 4. Systematic partitioning of SNN on neuromorphic hardware.

(1) the `spike disorder count` as the fraction of total spikes arriving out of order at the neurons and (2) the `Inter-spike distortion` as the maximum difference between the inter-spike interval of source and destination neurons.

## IV. SYSTEMATIC PARTITIONING FRAMEWORK FOR SNNS

Figure 4 shows our systematic framework. This framework takes an application implemented using SNN as input (shown at the top left). Table I provides a set of applications used to evaluate our partitioning methodology. It is to be noted that the first three applications are based on rate-coding, while the last one is with temporal coding. This provides a range of applications with conventional and evolving topologies with different spike coding schemes. Apart from these realistic applications, we also considered synthetic applications by varying the depth and width of SNN layers. All applications are first simulated using `CARLsim` [14], a GPU-accelerated library for simulating spiking neural network models with a high degree of biological detail.

The output of `CARLsim` is a trained SNN with spike times of each neuron in the SNN. This is converted into a dataflow graph as shown at the top right hand of Figure 4. The SNN graph (format is explained earlier) is presented to the PSO algorithm to partition into local and global synapses and map them on a neuromorphic hardware. The neuromorphic hardware architecture

TABLE I

REALISTIC APPLICATIONS USED FOR EVALUATING OUR APPROACH.

| Approach | Application | Topology |
|---|---|---|
| CARLsim native [14] | hello world (HW) | Feedforward (117, 9) |
| CARLsim native [14] | image smoothing (IS) | Feedforward (1024, 1024) |
| Diehl et al. [17] | handwritten digit (HD) | Unsupervised, recurrent (250, 250) |
| Das et al. [18] | heartbeat estimation (HE) | Unsupervised, LSM (64, 16) |

is input to `Noxim++`, which can incorporate details of a real chip (CxQuad in Figure 4). It is also possible to replace `Noxim++` with actual CxQuad or other neuromorphic chip directly, making the framework similar to `PACMAN`. The PSO formulation is discussed in the previous section.

The `Noxim++` simulator is an extended version of the originally proposed `Noxim` simulator [15], which is highly configurable network-on-chip simulation based on mesh architecture. The configurable parameters include buffer size, network size, packet size, packet injection rate, routing algorithm, selection strategy, among others. For the power consumption simulation, users can modify the power values in external loaded YAML file to benefit from the flexibility. During a simulation, `Noxim` calculates latency, throughput and power consumption automatically based on the statistics collected during runtime. The original `Noxim` simulator is extended with the following features for our framework

- addition of different interconnect models for representative neuromorphic hardware
- incorporation of SNN-related metrics (spike disorder count and inter-spike distortion)
- multicast feature, where spike packets can be communicated to a selected subset of crossbars.

Overall, the framework generates different results, a snapshot of which is presented to the bottom left corner of Figure 4. The full framework will be released upon acceptance of this work for the benefit of research community.

## V. RESULTS AND DISCUSSIONS

All experiments are conducted on Google Cloud Platform configured with 4 CPUs, 26 GB RAM and NVIDIA Tesla K80 GPU. The platform runs Ubuntu 14.04. Apart from the realistic applications of Table I, we considered synthetic applications with different number of neural
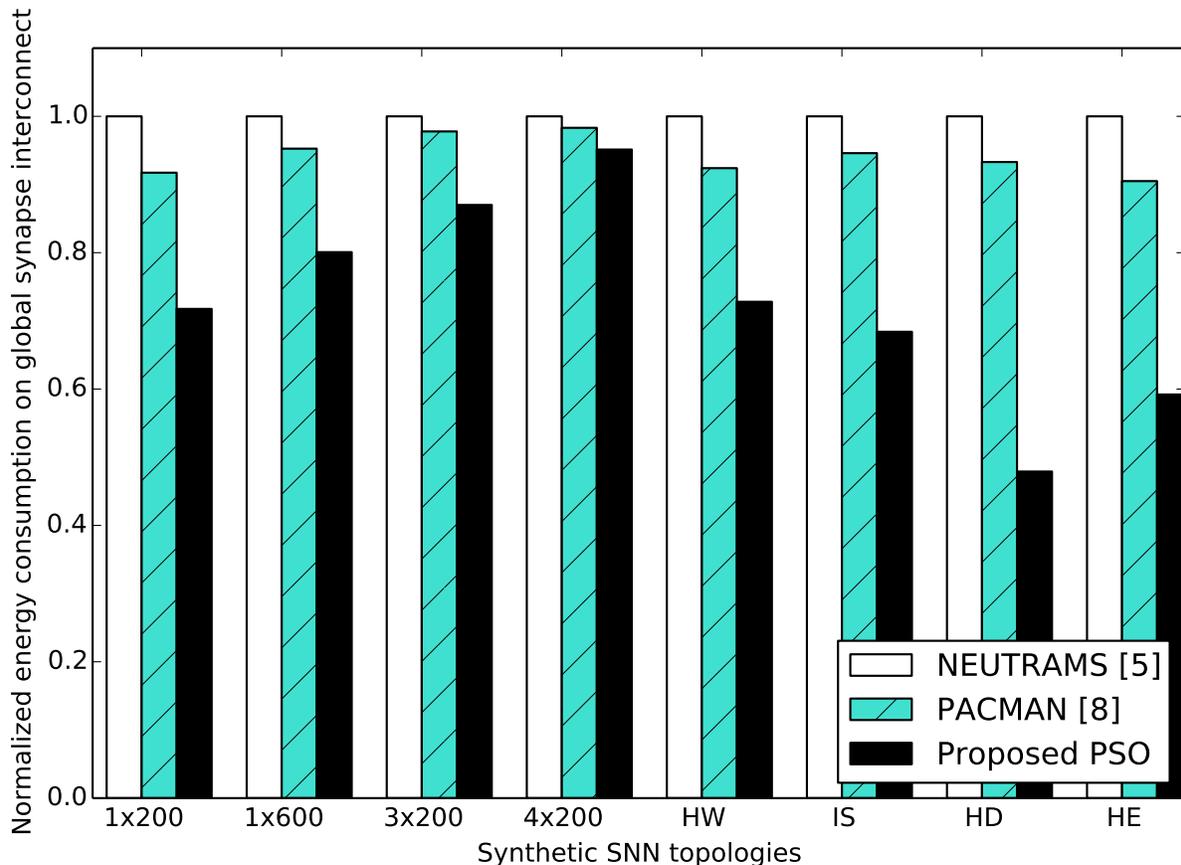
Fig. 5. Exploration with synthetic and realistic SNN-based applications.

network layers and number of neurons per layer. Our systematic framework is compared with PACMAN [8], adapted for CxQuad architecture. Additionally, we also compare our approach with the ad-hoc mapping technique NEUTRAMS [5], which uses a Network-on-Chip simulator to determine energy consumption on a neuromorphic architecture, without solving the local and global synapse partitioning problem and incorporating SNN performance.

### A. Energy Comparison on Global Synapse Interconnect

Figure 5 reports the energy consumption on the global synapse interconnect for three approaches – NEUTRAMS [5], PACMAN [8] and our proposed PSO-based partitioning. We evaluated 8 synthetic topologies (4 are plotted in Figure 5) and are marked on the X-axis with $m \times n$, where $m$ is the number of layers and $n$ is the number of neurons per layer. Neurons of the first layer in each of these topologies receive their input from 10 neurons creating spike trains,

TABLE II

METRIC EVALUATION FOR REALISTIC APPLICATIONS.

| Metric | hello_world | | image smoothing | |
|---|---|---|---|---|
| | PACMAN [8] | Proposed | PACMAN [8] | Proposed |
| ISI Distortion (cycles) | 6.6 | 2.9 | 8.1 | 4.7 |
| Disorder count (%) | 0.08 | 0.05 | 0.32 | 0.06 |
| Throughput (AER/ms) | 0.2 | 0.16 | 0.4 | 0.3 |
| Latency (cycles) | 108 | 70 | 99 | 69 |
| Metric | digit recog. | | heartbeat estimation | |
| | PACMAN [8] | Proposed | PACMAN [8] | Proposed |
| ISI Distortion (cycles) | 6.2 | 4.3 | 4.9 | 3.9 |
| Disorder count (%) | 0.02 | 0.01 | 0.12 | 0.02 |
| Throughput (AER/ms) | 0.5 | 0.4 | 0.3 | 0.33 |
| Latency (cycles) | 150 | 146 | 216 | 171 |

whose inter-spike interval follows a Poisson process with mean firing rates between 10 Hz and 100 Hz. Additionally, these synthetic SNNs implement fully connected feedforward topologies. Energy numbers for each topology are normalized with respect to NEUTRAMS. It can be seen clearly from Figure 5, that our proposed PSO-based partitioning achieves the minimum energy out of the three techniques. The improvement with respect to NEUTRAMS is between 2.4% and 48.7% (average 20.2%), while that with respect to PACMAN is between 1.5% and 45.4% (average 17.2%). It is to be noted that the energy improvement decreases with increase in the number of synapses. This is observed from results for topology 4x200 (with dense 122000 synapses) where energy consumption for the three approaches is comparable (energy gains are less than 2%). For topology 1x200 (with 2000 synapses), improvements using our approach is more than 40%. These results indicate that our approach is able to find the best partition for sparse and dense synapses, with higher improvements for sparse connections. Energy gains observed for the four realistic applications are in the range (27.0% – 52.1%, average 38%) with respect to NEUTRAMS and (21.2% – 48.7%, average 33%) with respect to PACMAN.

### B. SNN Metric Evaluation on Global Synapse Interconnect

Table II reports results for other metrics relevant to SNNs for realistic applications. Specifically, the table reports average inter-spike interval (ISI) distortion in terms of interconnect clock cycles

in rows 3 & 9; the spike disorder count as a fraction of the total spikes in rows 4 & 10; the average throughput in terms of number of AER packets per ms on the global synapse interconnect in rows 5 & 11; and finally the maximum latency of spike communication on the global synapse interconnect in terms of interconnect clock cycles in rows 6 & 12. Furthermore, the table compare results from our approach with that from `PACMAN` [8]. As can be seen clearly from this table, our approach outperforms `PACMAN` in terms of ISI distortion achieving an average 37% fewer global synapse interconnect cycles. It is to be noted that ISI distortion significantly impacts application performance with temporal information coding such as the `heartbeat estimation`, where we observe that 20% reduction of ISI distortion improves estimation accuracy by over 5%. For the other 3 applications which are based on rate coding, the accuracy improvement due to ISI distortion is insignificant. In terms of spike disorder count, we observe that our approach achieves an average 63% lower spike arrival disorder compared to `PACMAN`. It is to be observed that the throughput in `PACMAN` is usually higher than that of our approach. This is because the number of spikes communicated on the global synapse interconnect is usually higher in `PACMAN`. Finally, the spike propagation latency on the global synapse interconnect is also lower using our approach by 22% (2% − 35%). Improvements are consistent for the 8 evaluated synthetic topologies. These improvements are due to our PSO-based partitioning, which reduces spike congestion on the global synapse interconnect, improving communication energy and latency.

## C. Neuromorphic Architecture Exploration

To further demonstrate the broad usage of our approach and framework, we take an application (`digit recognition` [17]) and explore architectural alternatives. We aim to answer the following question: given an application, which is preferred between an architecture with fewer number of large crossbars or an architecture with large number of small crossbars? Figure 6 plots exploration results in terms of local/global synapse energies and latency for spike communication on the global synapse interconnect. The number of neurons per crossbar is increased from 90 to 1440. The local synapse energy is the total energy for spike communication inside all crossbars in the architecture. The global synapse energy is the total energy for spike communication on the global synapse interconnect. The energy numbers are averaged for processing/interpreting a hand-written digit image of 28x28 pixels. Latencies in this figure are worst-case values for all spikes communicated on the global synapse interconnect.
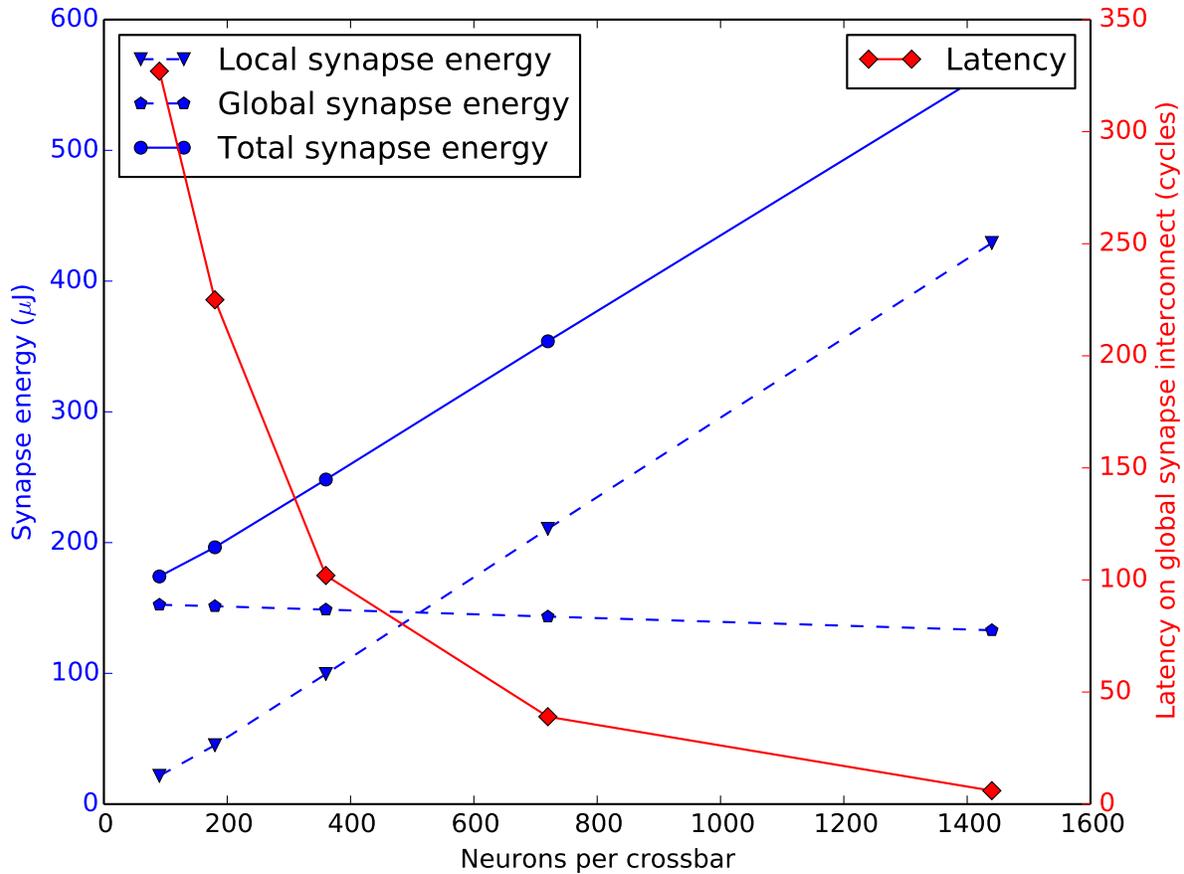
Fig. 6. Architecture exploration with hand-written digit recognition [17].

Following observations can be made from Figure 6. First, as the size of crossbars increases, more synapses are mapped by our approach inside a crossbar reducing the number of spikes and hence energy on the global synapse interconnect. Conversely, the local synapse energy increases due to more spike communication inside a crossbar. Second, the worst-case spike latency on the global synapse interconnect decreases with increase in the size of crossbars. This is because, as more synapses are mapped locally, congestion on the interconnect reduces, reducing the latency. Clearly, the best choice is an intermediate point in between the extremes. So we clearly need our framework to explore this complex search space. Once the best solution is identified, the size (and hence the number) of crossbars for a given application can be determined, together with partitioning the underlying SNN into local and global synapses. The mapping can then be enforced on the CxQuad board at design-/configuration time.
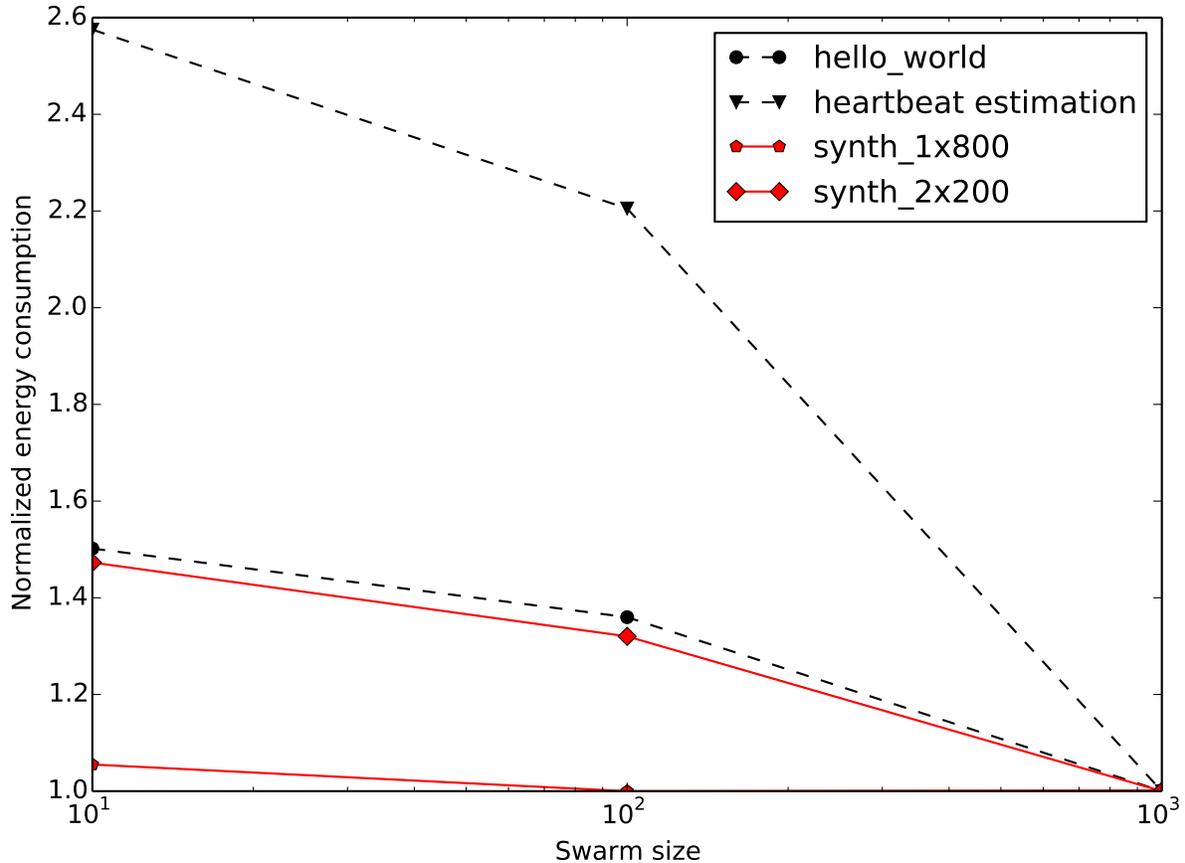
Fig. 7. Exploration with swarm size.

## D. PSO-Related Exploration Results

In this section we present results related to our particle swarm optimization. Figure 7 plots the energy consumption on the global synapse interconnect with different number of swarm particles (in log scale) for four applications – two real and two synthetic, with the number of iterations fixed to 100. Energy results are normalized with respect to the minimum energy obtained for these applications. Furthermore, there are no improvements in the energy consumption with more than 1000 particles. Hence the x-axis is limited to 1000 particles. As can be seen, with more number of particles, the algorithm is able to find better results for a fixed number of iteration of the algorithm. For applications `synth_2x200`, the minimum is reached for a swarm size of 105. For other three applications, the point of minimum energy is close to 1000. Trends are consistent for other applications considered. Based on these results, we used a swarm size

of 1000 for our experiments. The average wall clock time to find an optimum solution using these settings (swarm size = 1000, iteration = 100) is average 35 minutes on the Google Cloud platform. Detailed timing results are omitted for space limitation.

## VI. CONCLUSIONS AND OUTLOOK

This paper presents an approach to map local and global synapses of a SNN-based application on a crossbar-based neuromorphic architecture. Fundamental to our approach is the use of particle swarm optimization, which partitions SNN-based applications, with local synapses mapped to crossbars and global synapses mapped on a time-multiplexed interconnect. Furthermore, we proposed two metrics – inter-spike interval distortion and disorder count, specific to SNN performance for spatial and temporal information coding. Using realistic and synthetic applications we show that our approach reduces spike communication on the global synapse interconnect, reducing communication energy by an average 33% and spike propagation latency by an average 22%, with respect to PACMAN [8], which is the standard SNN mapping technique for SpiNNaker. Run-time SNN mapping will be addressed in future.

## REFERENCES

[1] W. Maass, "Networks of spiking neurons: the third generation of neural network models," *Neural networks*, vol. 10, no. 9, pp. 1659–1671, 1997.

[2] F. Akopyan, J. Sawada *et al.*, "TrueNorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip," *IEEE transactions on computer-aided design of integrated circuits and systems*, vol. 34, no. 10, pp. 1537–1557, 2015.

[3] G. Indiveri, F. Corradi, and N. Qiao, "Neuromorphic architectures for spiking deep neural networks," in *International Electron Devices Meeting (IEDM)*. IEEE, 2015.

[4] M. M. Khan, D. R. Lester *et al.*, "SpiNNaker: mapping neural networks onto a massively-parallel chip multiprocessor," in *International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2008.

[5] Y. Ji, Y. Zhang, S. Li, P. Chi, C. Jiang, P. Qu, Y. Xie, and W. Chen, "NEUTRAMS: Neural network transformation and co-design under neuromorphic hardware constraints," in *International Symposium on Microarchitecture (MICRO)*. IEEE, 2016.

[6] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE journal of solid-state circuits*, vol. 52, pp. 127–138, 2017.

[7] R. Serrano-Gotarredona *et al.*, "CAVIAR: A 45k neuron, 5m synapse, 12g connects/s aer hardware sensory–processing–learning–actuating system for high-speed visual object recognition and tracking," *IEEE transactions on neural networks*, vol. 20, no. 9, pp. 1417–1438, 2009.

[8] F. Galluppi, S. Davies, A. Rast, T. Sharp, L. A. Plana, and S. Furber, "A hierachical configuration system for a massively parallel neural hardware platform," in *International Conference on Computing Frontiers*, 2012.

[9] E. Stromatias, D. Neil, F. Galluppi, M. Pfeiffer, S.-C. Liu, and S. Furber, "Scalable energy-efficient, low-latency implementations of trained spiking deep belief networks on spinnaker," in *International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2015.

[10] T. Serrano-Gotarredona, B. Linares-Barranco, F. Galluppi, L. Plana, and S. Furber, "ConvNets experiments on SpiNNaker," in *Circuits and Systems (ISCAS), 2015 IEEE International Symposium on*. IEEE, 2015.

[11] W. Maass, T. Natschläger, and H. Markram, "Real-time computing without stable states: A new framework for neural computation based on perturbations," *Neural computation*, vol. 14, pp. 2531–2560, 2002.

[12] J. Hawkins, S. Ahmad, and D. Dubinsky, "Hierarchical temporal memory including HTM cortical learning algorithms," *Techical report, Numenta, Inc, Palto Alto*, 2010.

[13] R. Eberhart and J. Kennedy, "A new optimizer using particle swarm theory," in *International Symposium on Micro Machine and Human Science (MHS)*. IEEE, 1995, pp. 39–43.

[14] M. Beyeler, K. D. Carlson, T.-S. Chou, N. Dutt, and J. L. Krichmar, "CARLsim 3: A user-friendly and highly optimized library for the creation of neurobiologically detailed spiking neural networks," in *International Joint Conference on Neural Networks (IJCNN)*, 2015.

[15] V. Catania, A. Mineo, S. Monteleone, M. Palesi, and D. Patti, "Noxim: An open, extensible and cycle-accurate network on chip simulator," in *International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, 2015.

[16] K. A. Boahen, "Communicating neuronal ensembles between neuromorphic chips," in *Neuromorphic systems engineering*. Springer, 1998.

[17] P. U. Diehl and M. Cook, "Unsupervised learning of digit recognition using spike-timing-dependent plasticity," *Frontiers in computational neuroscience*, vol. 9, 2015.

[18] A. Das, P. Pradhapan, W. Groenendaal, P. Adiraju, R. T. Rajan, F. Catthoor, S. Schaafsma, J. L. Krichmar, N. Dutt, and C. Van Hoof, "Unsupervised heart-rate estimation in wearables with liquid states and a probabilistic readout," *arXiv preprint arXiv:1708.05356*, 2017.