

# Neuromorphic computing using phase-change memory devices

[Introduction to SNU class]

2021. 04. 27.

Uicheol Shin

Neuromorphic Device, Science & Technology, IBM Research - Tokyo

Department of Materials Science and Engineering, Seoul National University

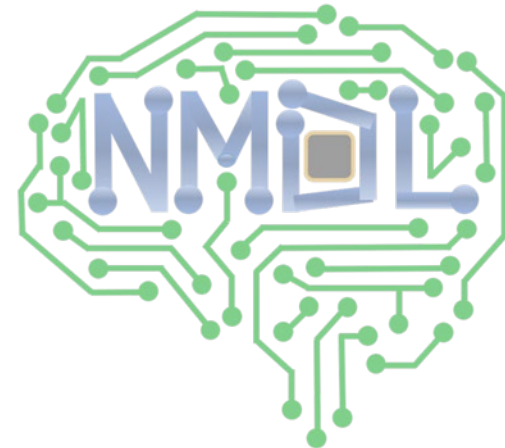
1. Introduction

2. Computational memory

3. Deep learning co-processors

4. Spiking neural networks (SNN)

5. Summary



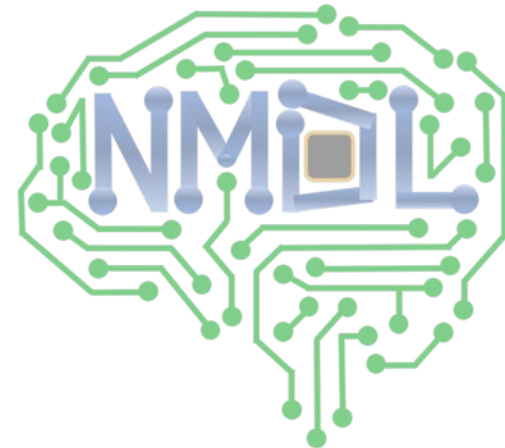
## 1. Introduction

## 2. Computational memory

## 3. Deep learning co-processors

## 4. Spiking neural networks (SNN)

## 5. Summary



# Introduction

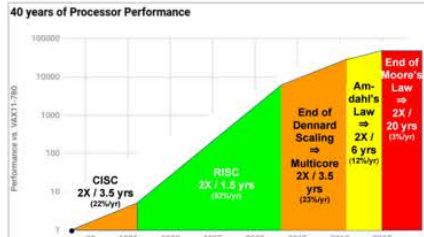
-4-

nature  
electronics

Editorial | Published: 17 April 2018

## Does AI have a hardware problem?

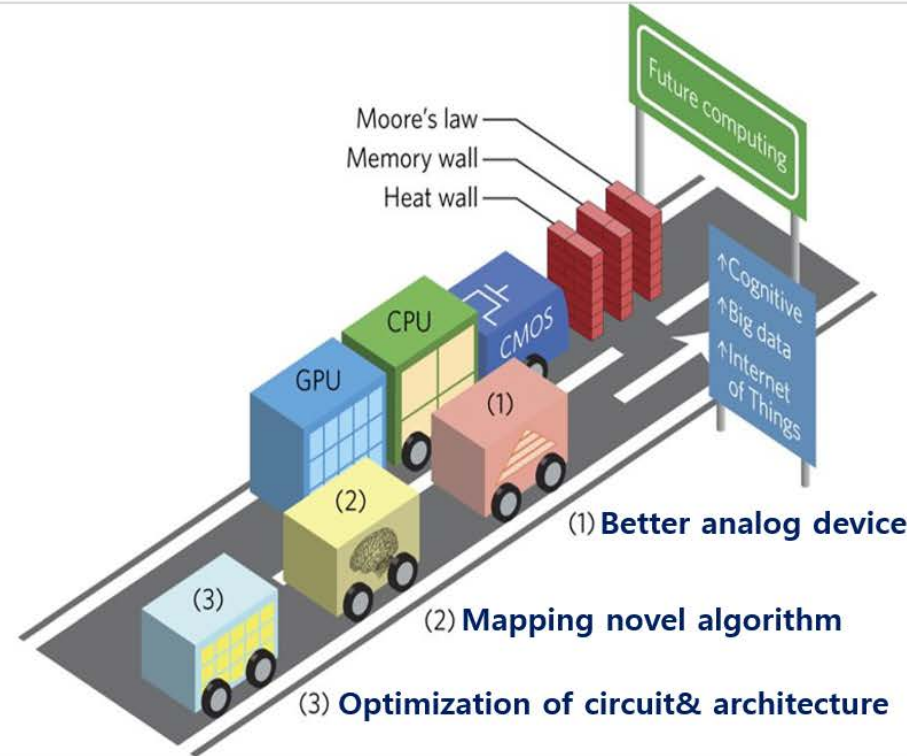
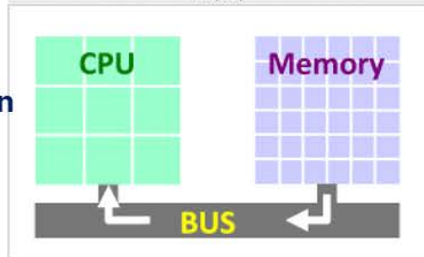
End of  
Moore's Law



Explosion of  
data usage

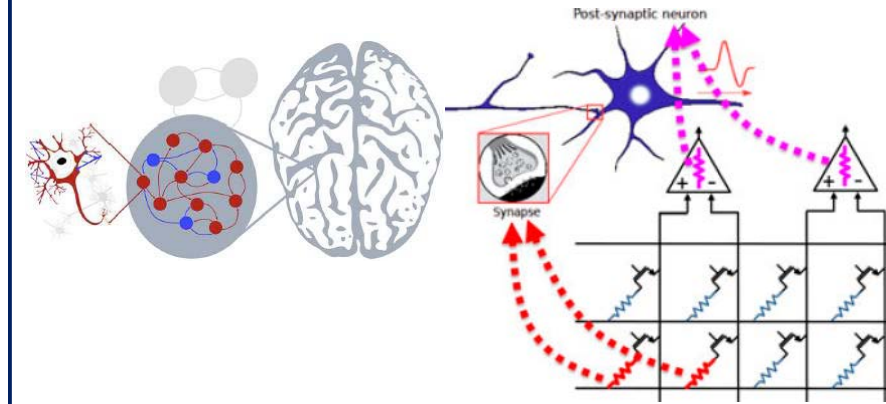


Von Neumann  
bottleneck



## Neuromorphic computing system

→ One of the non-von Neumann approach,



G. Burr, Adv. in Phys.:X, 2017, Vol. 2 No.1, 89-124

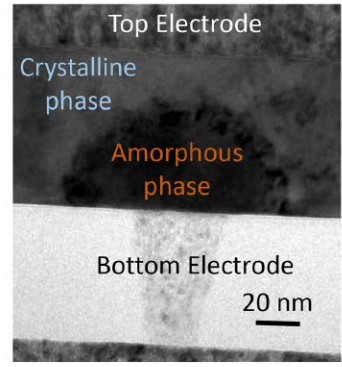
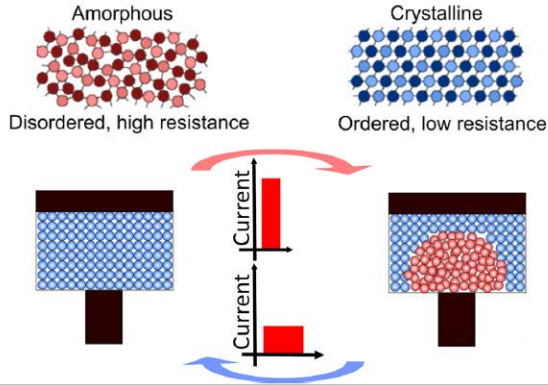
Novel architectures where memory and processing are collocated can perform efficiently in area/energy

**Phase-change memory devices** are used to build **computational memory**.

This novel architecture can accelerate the training of **deep neural networks(DNN)**, and also it can play a key role in **spiking neural networks(SNN)**.

# Introduction – PCM device

-5-

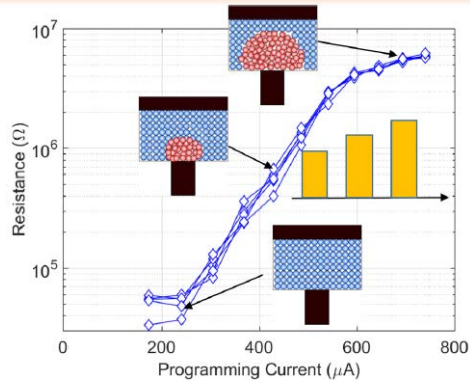


Phase change memory is one of the 'mature' non-volatile memory devices.

Typically, using Ge, Sb, Te chalcogenide compounds to program as memory.

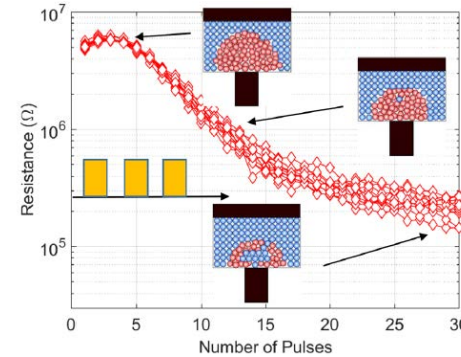
These certain materials exhibit drastically different electrical characteristics depending on their atomic arrangement.

## Key properties for brain-inspired computing



1. Achieve not just two levels, but a continuum of resistance or conductance values

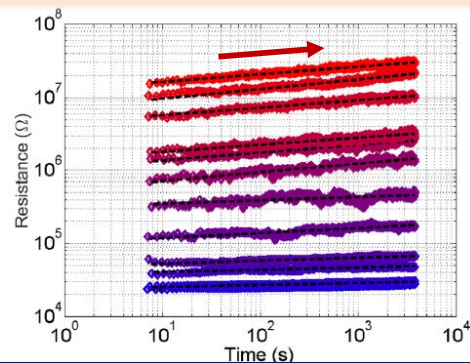
→ Analog synaptic weight



2. Accumulative behavior, to achieve a linear increase in G as a function of the # of SET pulses

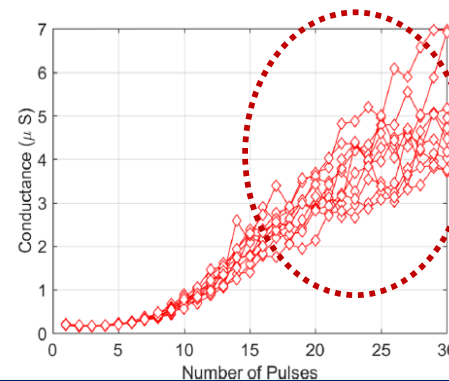
→ Controllable weight change

## Drawback characteristics of PCM devices



- Resistance drift,

Temporal fluctuations of R from a spontaneous 'structural relaxation' of the amorphous phase



- Cycle-to-cycle randomness,

Due to the inherent stochasticity associated with the crystallization process

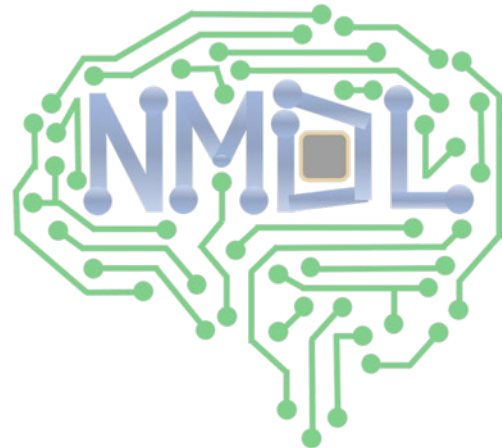
1. Introduction

2. Computational memory

3. Deep learning co-processors

4. Spiking neural networks (SNN)

5. Summary





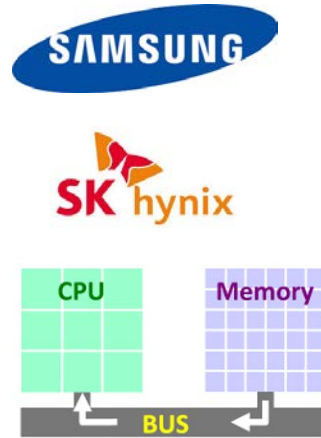
# Novel computer architecture for AI

-7-

H.S. Shin, IEEE TRANSACTIONS ON COMPUTER-AIDED DESIGN OF INTEGRATED CIRCUITS AND SYSTEMS, VOL. 37, NO. 11, NOVEMBER 2018



Samsung AI Forum, 2019



In-memory computing,  
Near-memory computing  
NPU(Neural processing unit)

⋮

Why do we need a new computing architecture?

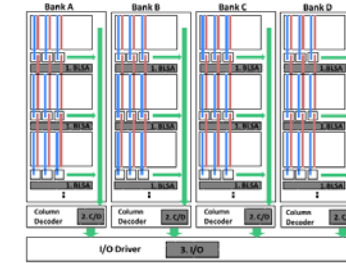


Fig. 1. MAC locations (BLSA, column decoder, and I/O driver).

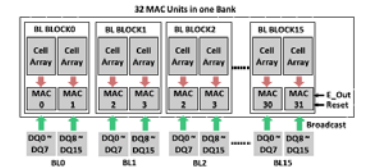
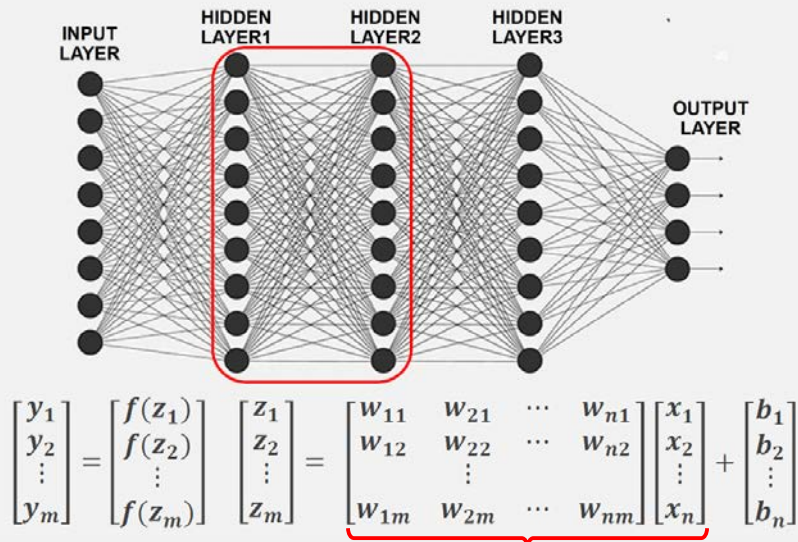


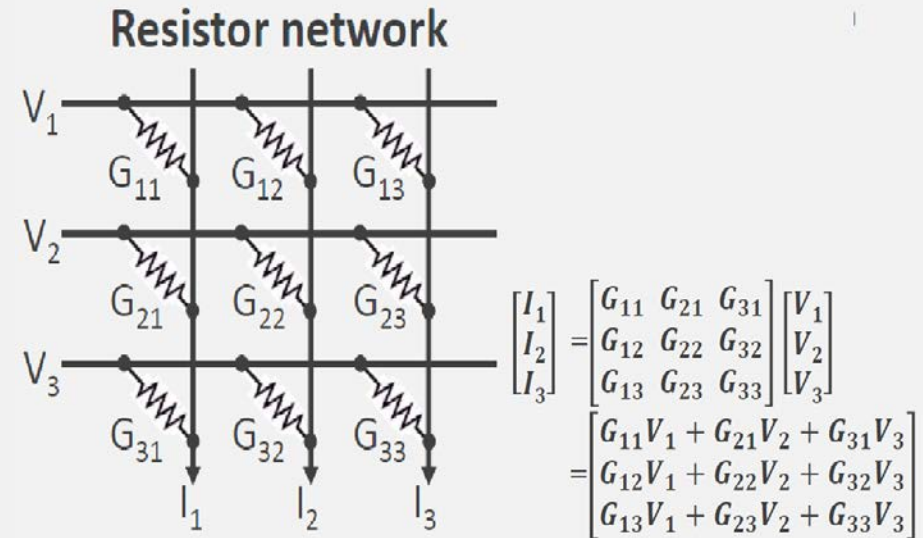
Fig. 4. 32 MAC units in a bank.

## Computation in Artificial Neural Network(ANN)



Key Operation : Multiply-Accumulate, "MAC"

## Key ideas in analog neuromorphic hardware

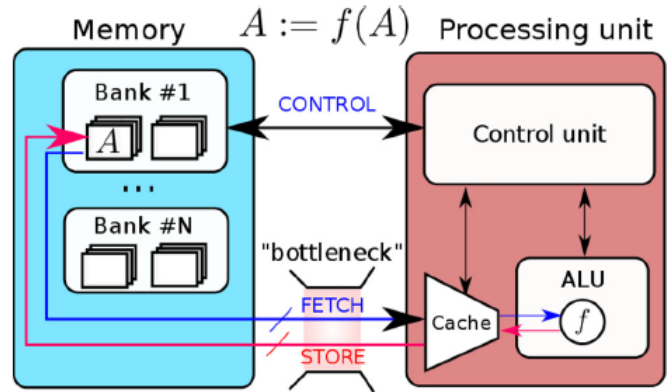


The physics of Ohm's law and Kirchhoff's current law allow the implementation of analog MAC operations in parallel.

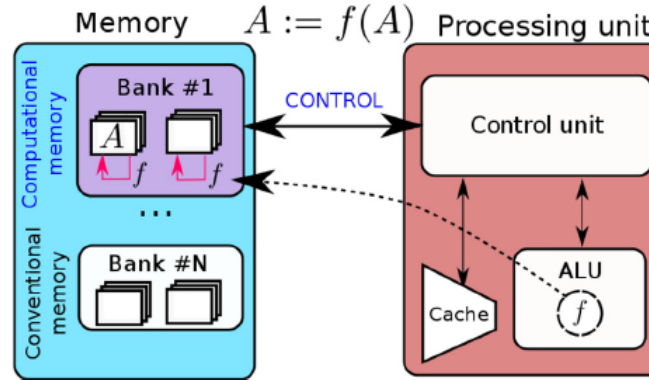
# In-memory computing with PCM devices

-8-

In-memory computing can overcome memory-processing unit bottleneck



Conventional Von Neumann architecture



Non-Von Neumann architecture

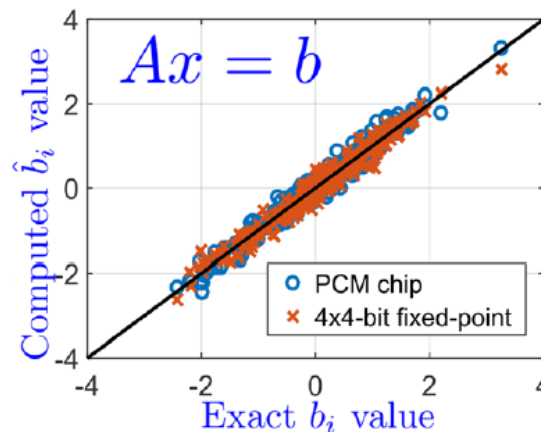
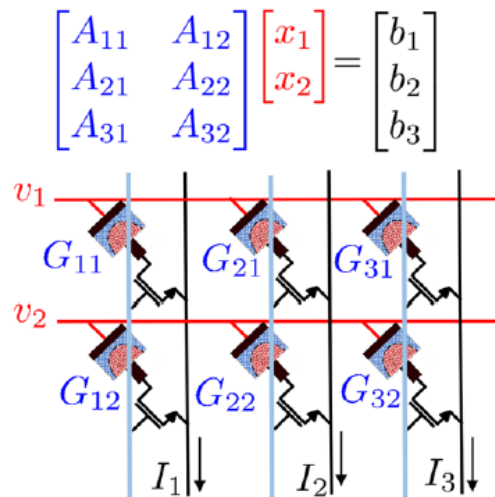
$f(A)$  can be performed in memory which stores data A

**Saves energy and time!**

Especially, in pattern recognition, real-time recognition, and intuitive functions!

Nature Comms., 2017, 8, 1115

## - Matrix-vector multiplication using PCM devices



Accuracy of the computation with PCM is comparable to that of fixed-point digital.



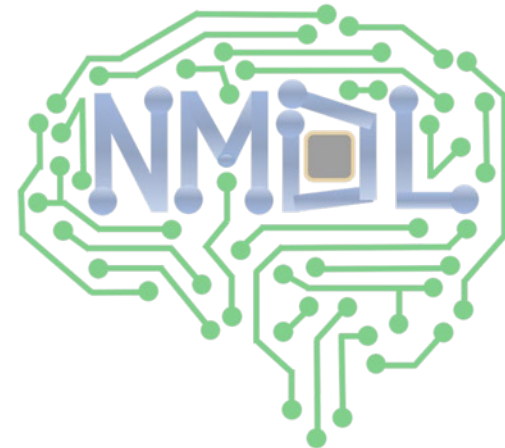
1. Introduction

2. Computational memory

3. Deep learning co-processors

4. Spiking neural networks (SNN)

5. Summary

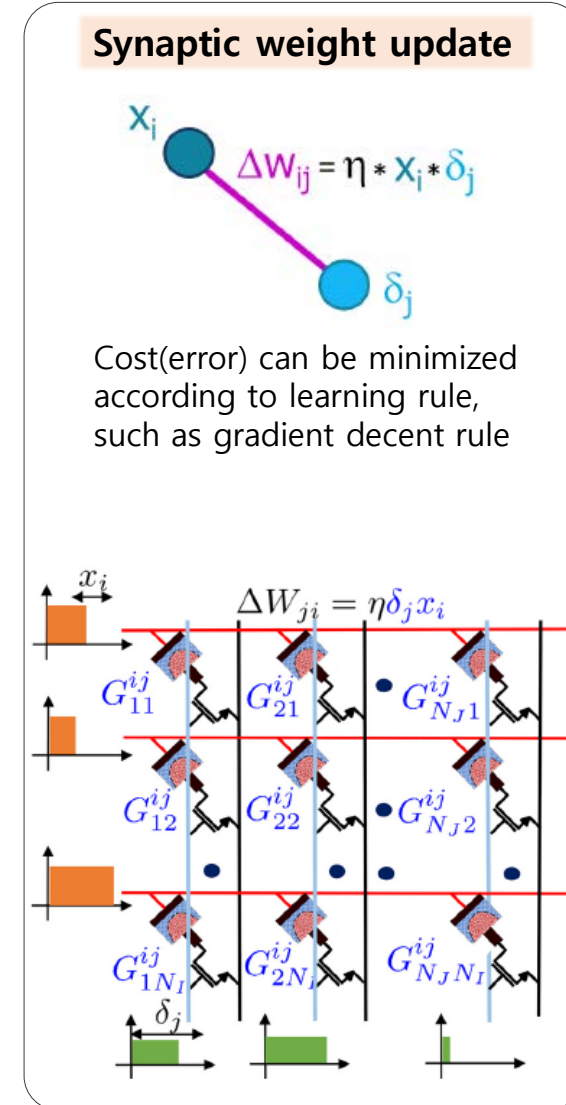
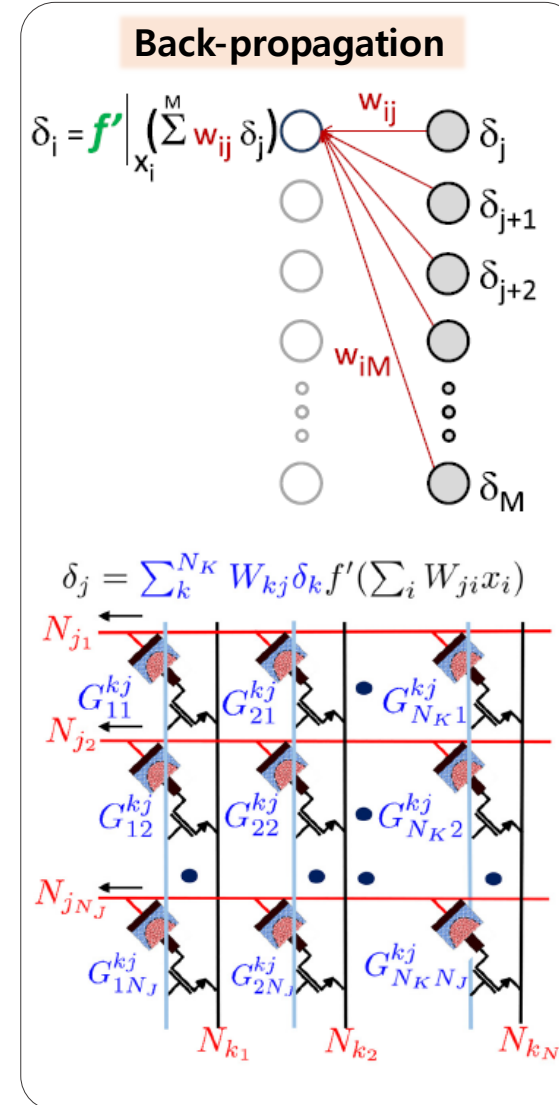
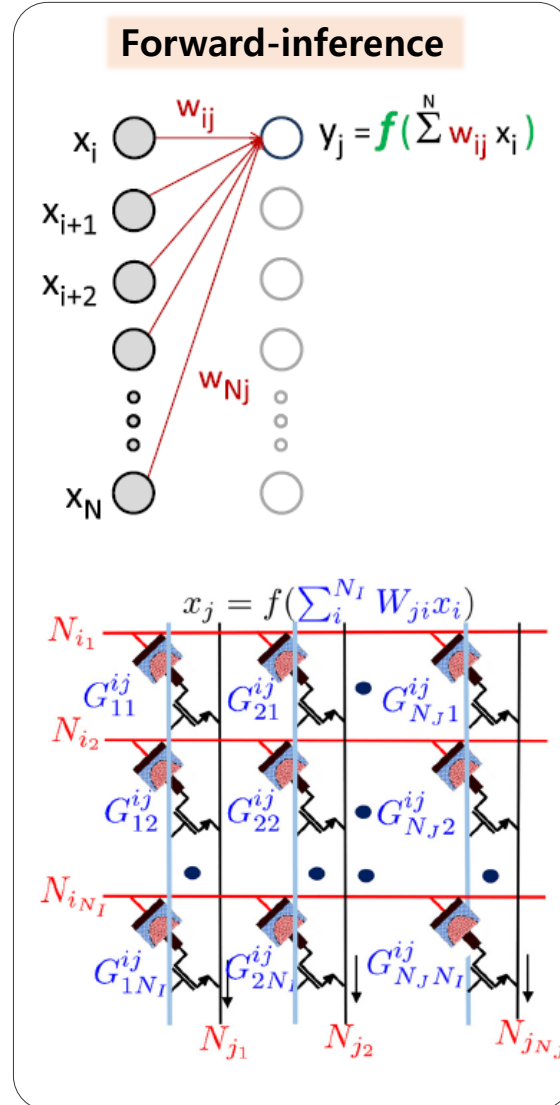
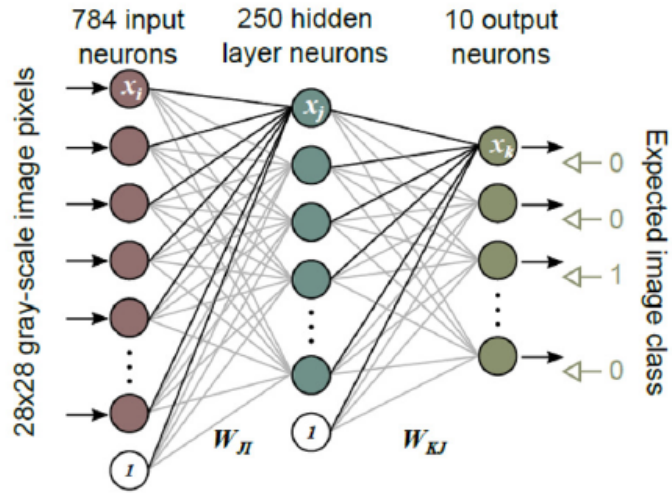


# Analog memory-based DNN process

-10-

- Deep Learning (Deep neural network, DNN) on PCM hardware(cross-bar array)

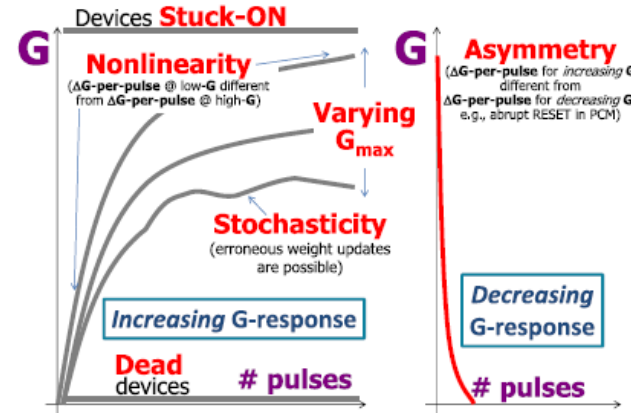
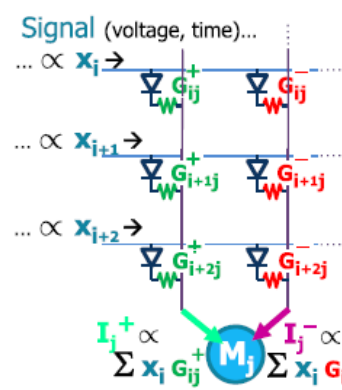
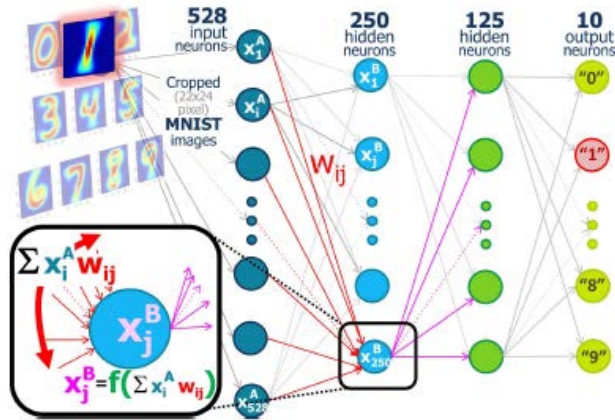
J. Appl. Phys., 2018, 124, 283011



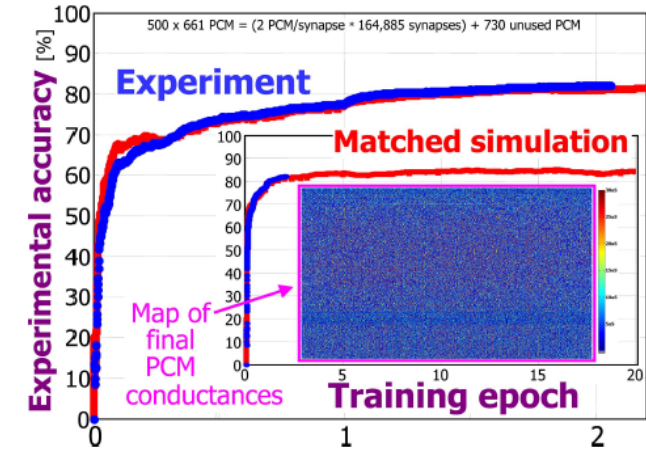
# Demonstration of DNN using PCM devices

-11-

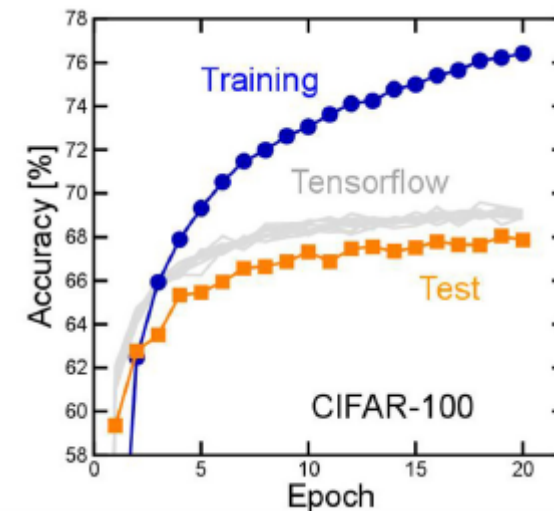
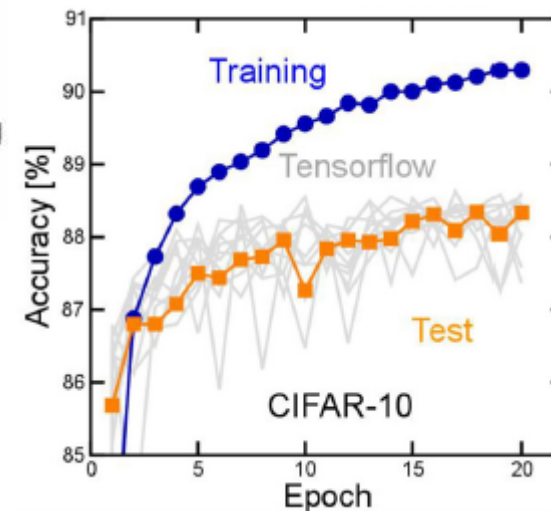
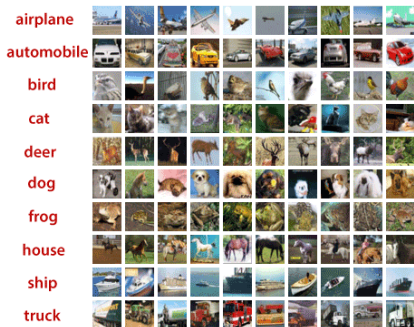
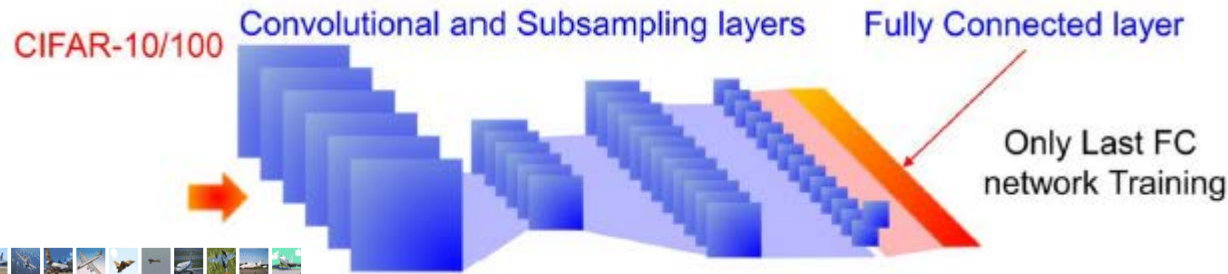
## - MNIST handwritten image recognition



$G^+$ ,  $G^-$  approach used to overcome PCM's inherent shortcomings



## - CIFAR-10/100 image recognition (Using Convolutional layer)

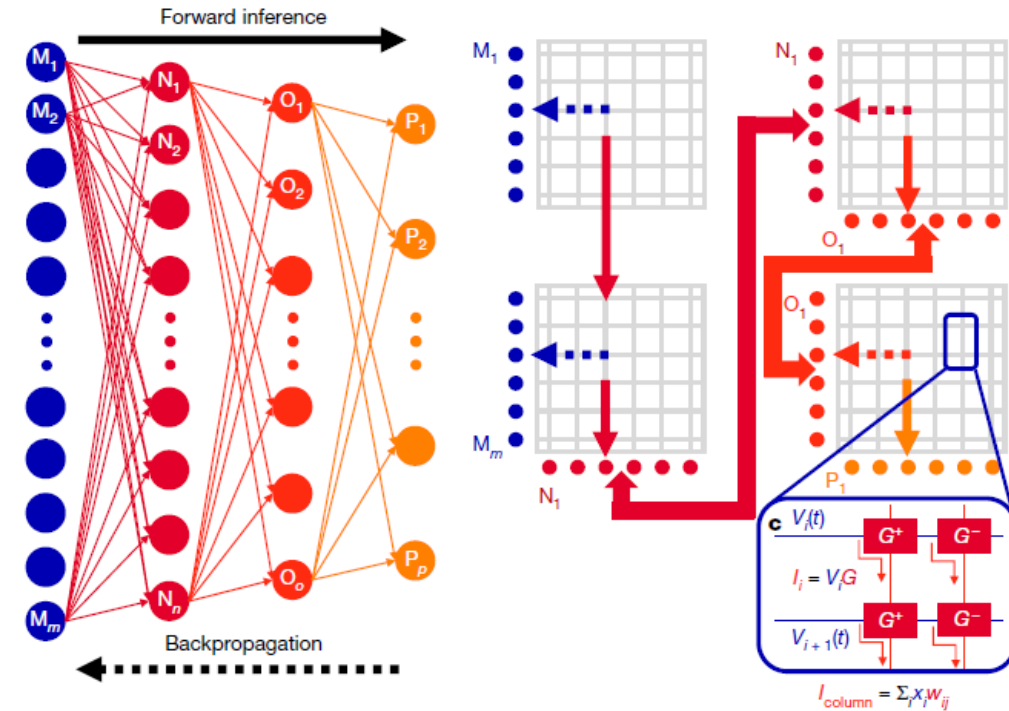
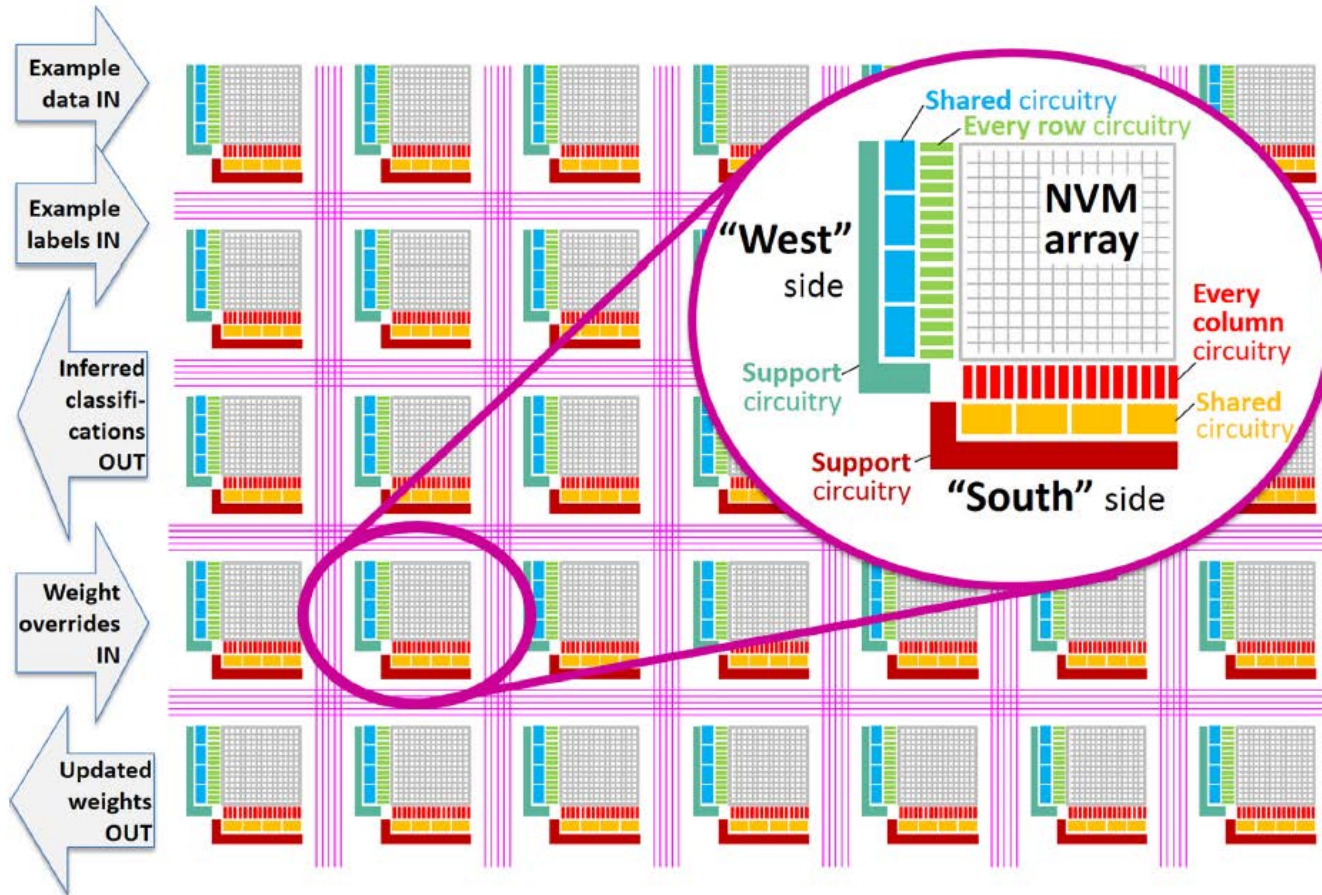




# A proposed chip architecture

-12-

- Chip architecture for a co-processor for deep learning based on PCM arrays



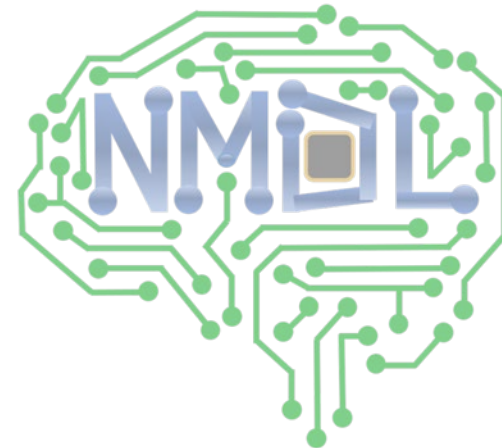
1. Introduction

2. Computational memory

3. Deep learning co-processors

4. Spiking neural networks (SNN)

5. Summary





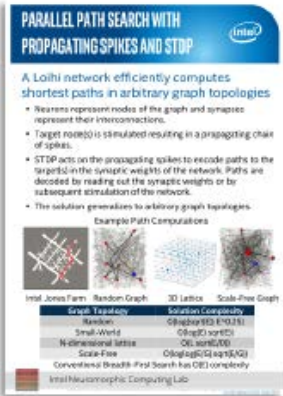
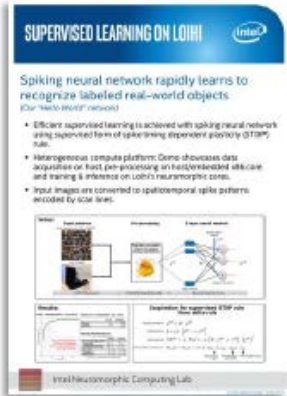
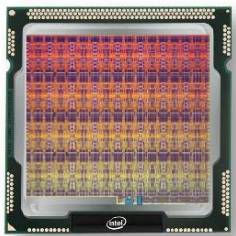
# Spiking Neural Networks (SNN)



## To Spike or Not to Spike: That Is the Question

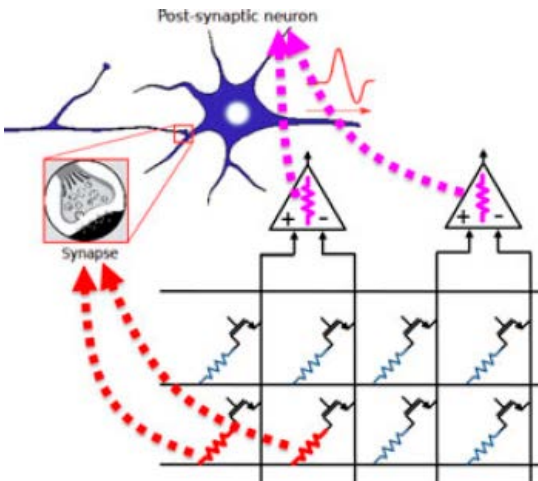
By **WOLFGANG MAASS**  
*Institute for Theoretical Computer Science,  
Graz University of Technology, Graz 8010, Austria*

Vol 103, No. 12, December 2015 | PROCEEDINGS OF THE IEEE



### - Third generation of neural network?

	ANN : Non-spiking NN	Brain : Spiking NN
Inputs & Outputs	Real-valued numbers	Spikes (encoded)
Neuron Operation		

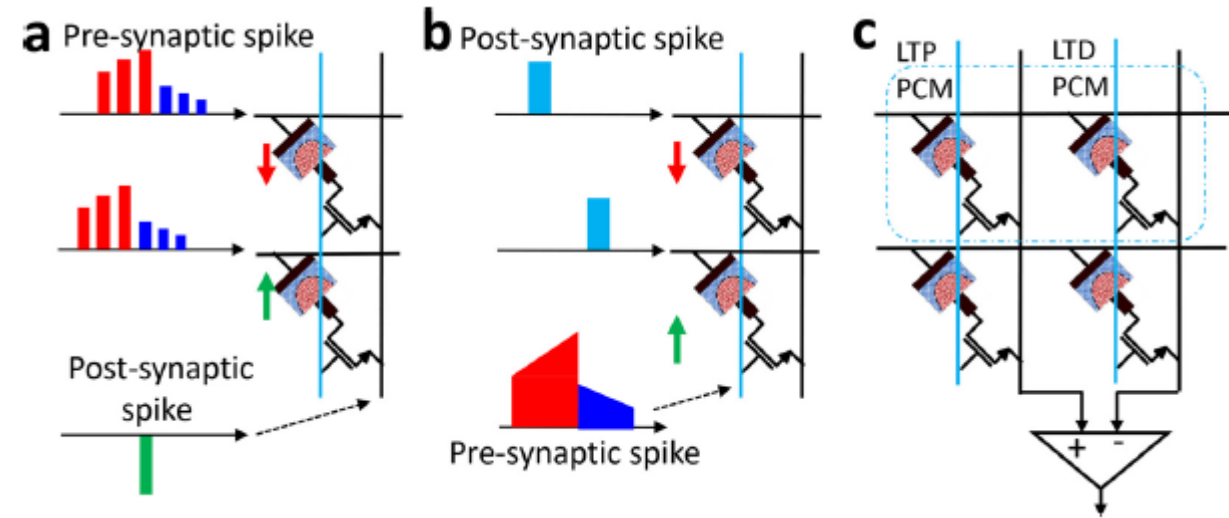
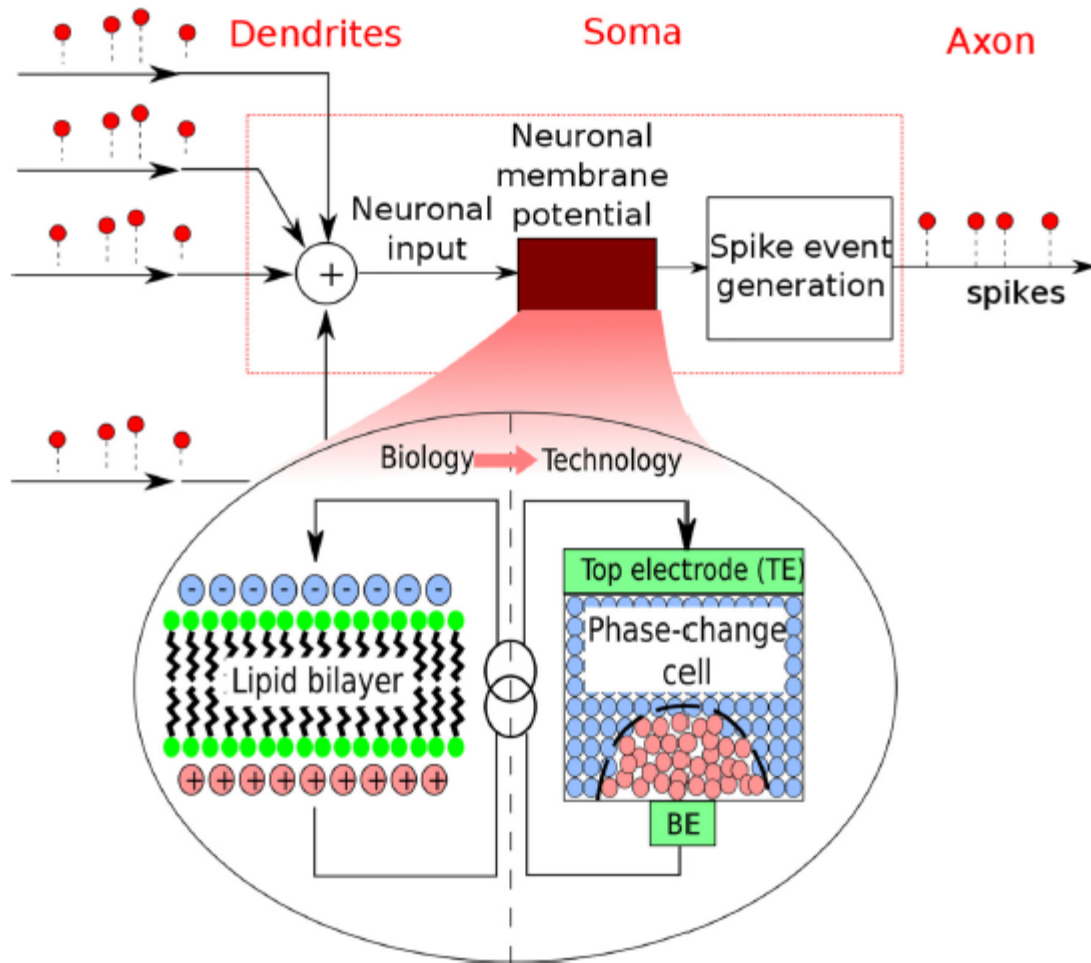


Neuron?  
Synapse?

M. Davies, NICE, 2018, Loihi intro talk

# Neuronal dynamics implementation using PCMs

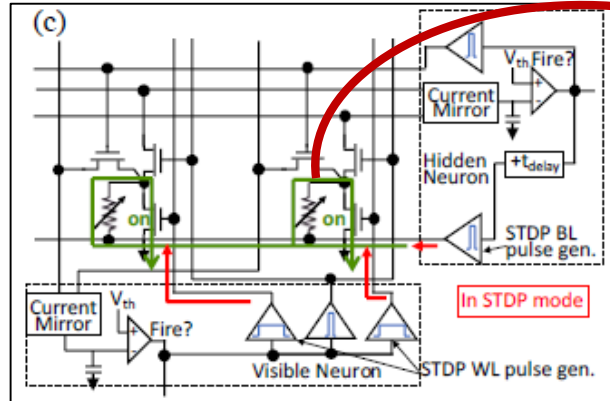
-15-



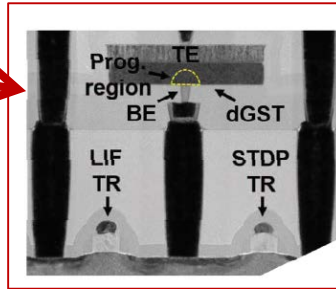
# Hardware concept for PCM programming & LIF

-16-

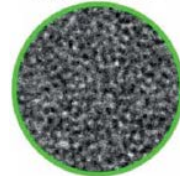
## [ Phase change memory ]



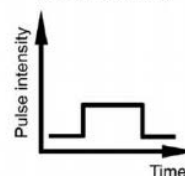
M. Ishii *et al.*, in *IEDM*, 2019, pp. 14.2.1-4.



Amorphous phase  
Low reflectivity  
High resistance



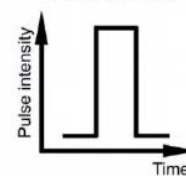
SET  
(long low laser  
or current pulse)



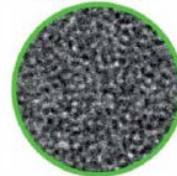
Crystalline phase  
High reflectivity  
Low resistance



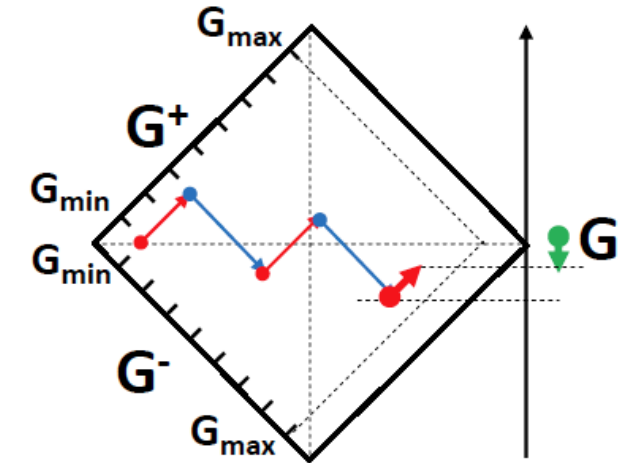
RESET  
(short high laser  
or current pulse)



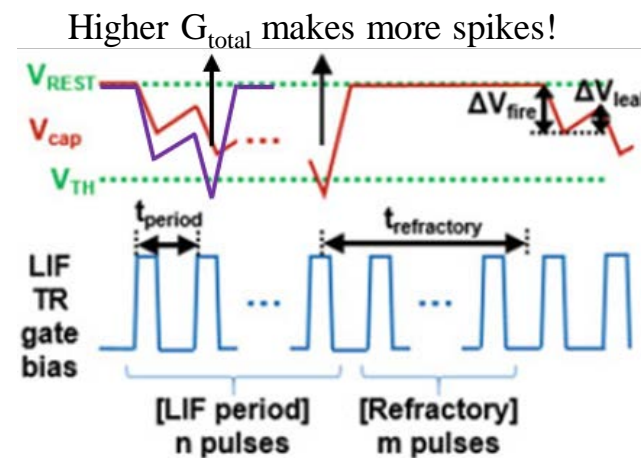
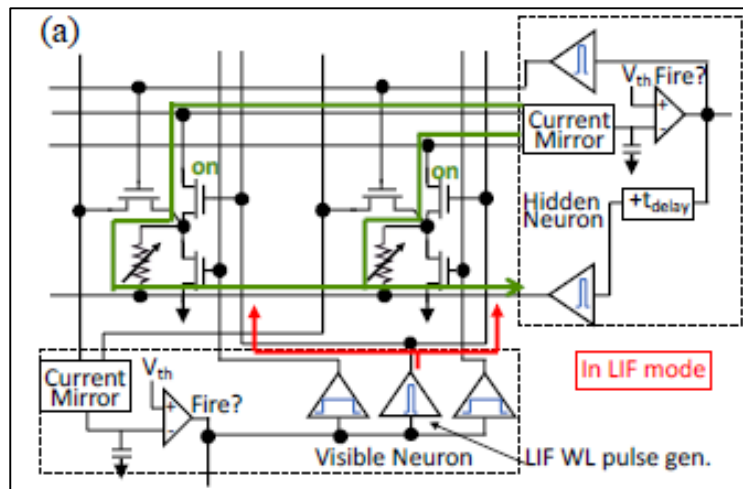
Amorphous phase  
Low reflectivity  
High resistance



## [ Bipolar synaptic weights ]

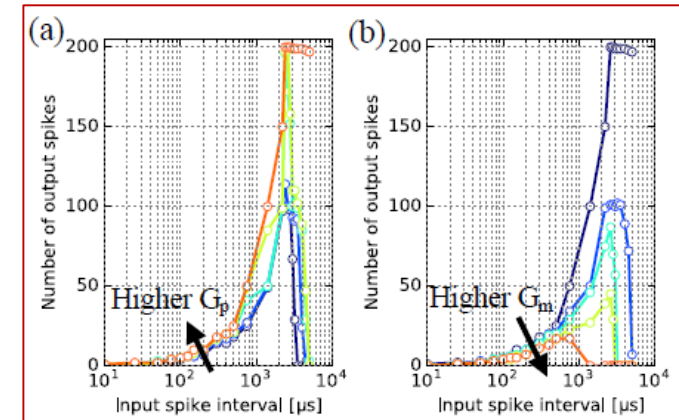


## [ LIF (Leaky-integrate & fire) process ]



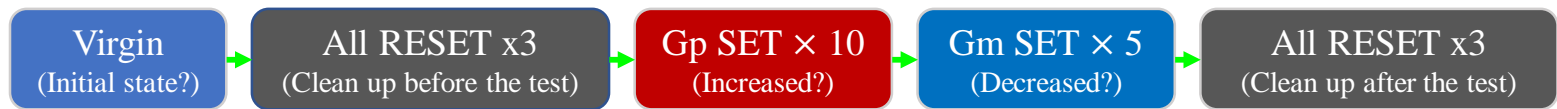
$$\left\{ \begin{array}{l} \text{SET } G_p \rightarrow G_{total} \text{ increase} \rightarrow \text{LIF outputs } \uparrow \\ G_m \rightarrow G_{total} \text{ decrease} \rightarrow \text{LIF outputs } \downarrow \end{array} \right\}$$

$(G_{total} = G_p - G_m)$





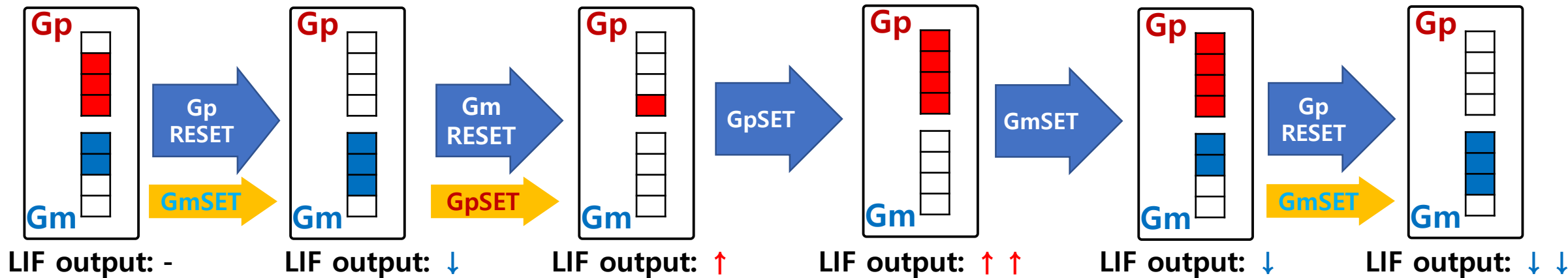
## [ PGM & LIF test flow ]



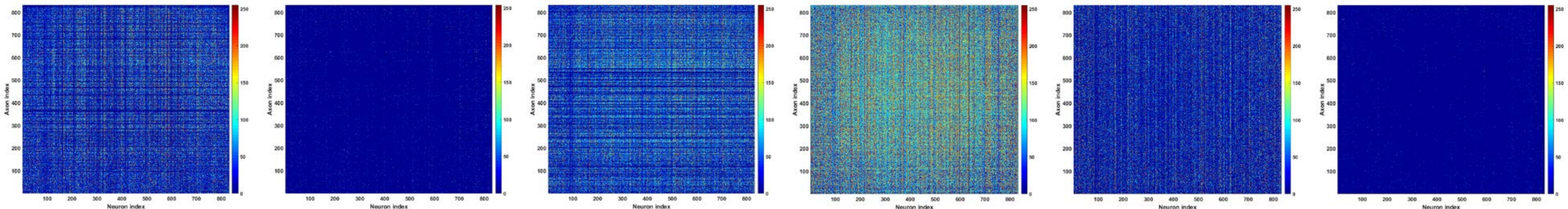
## [Schematics of the programming]



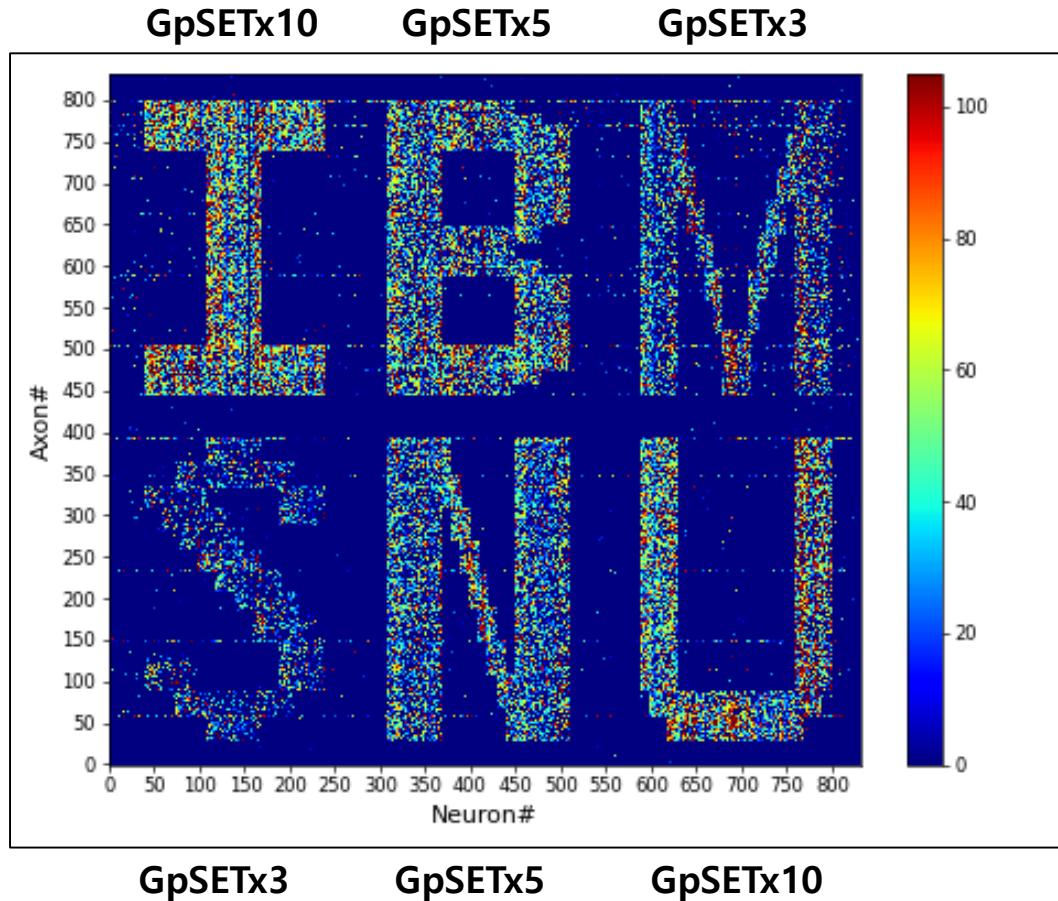
**SET** : SET processes occur inevitably during the RESET process.



## [LIF results of the programming]



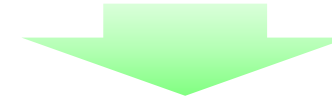
**From the results, we can program all PCM cells gradually on both Gp & Gm**



We can program the synaptic array with accurate position so that able to do 'pixel art'.

And we are able to confirm that how the conductance of each cell affects hardware with basic operation, LIF.

Even we don't know accurate conductance of each synaptic cell, we can expect the 'actual weight'!



This programming technology can be applied to implementation of 'weight transfer', which is necessary for 'off-chip learning'.



## Neuromorphic computing using phase-change memory devices

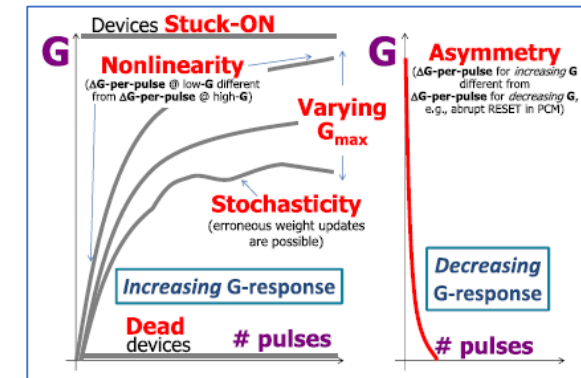
Brain-inspired computing schemes is promising technology to overcome 'von Neumann architecture' in AI computing.

Phase-change memory device, one of the 'mature' non-volatile memory technologies  
can achieve significantly higher performance compared to conventional computation with  
i ) Multi-conductance level, ii ) Accumulative behavior.

Computational memory, In-memory computing, Near-memory computing  
are novel approaches to mitigate 'von-Neumann bottleneck'.

These massively parallel computing units with PCMs can perform in application of

- i ) Deep-learning, Deep Neural Network(DNN)
- ii) Spiking Neural Network(SNN)



But there are challenges towards the adoption of PCM-based computing systems in future AI hardware,  
i ) Phase-change memory device – conductance fluctuation, non-linearity & non-symmetric...  
ii) Spiking Neural Network(SNN) – novel learning algorithms...



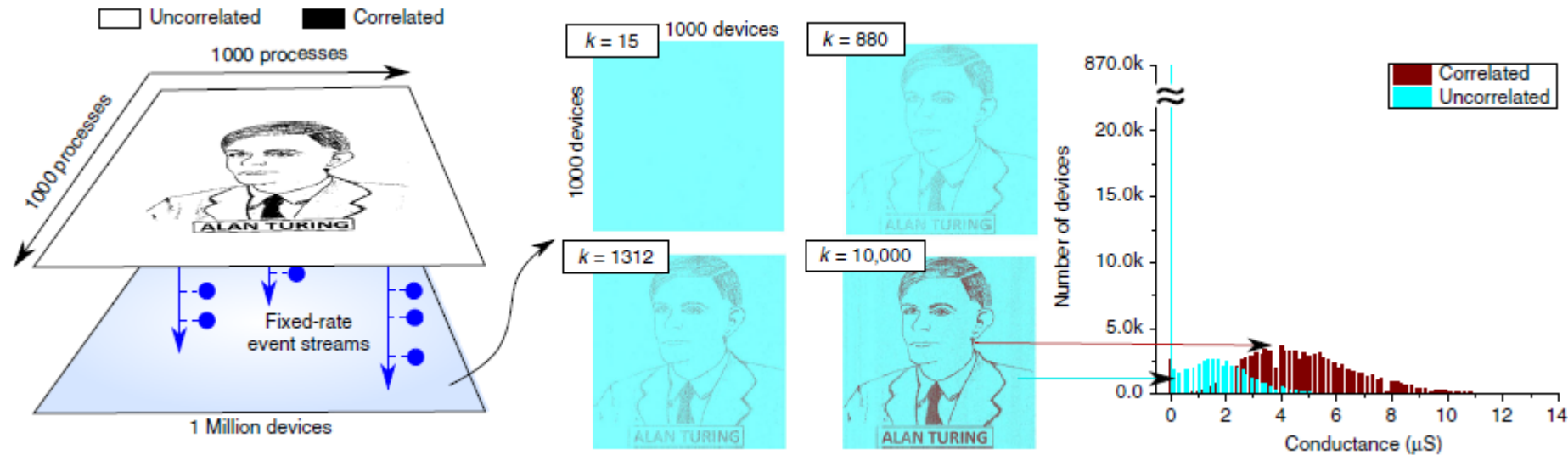
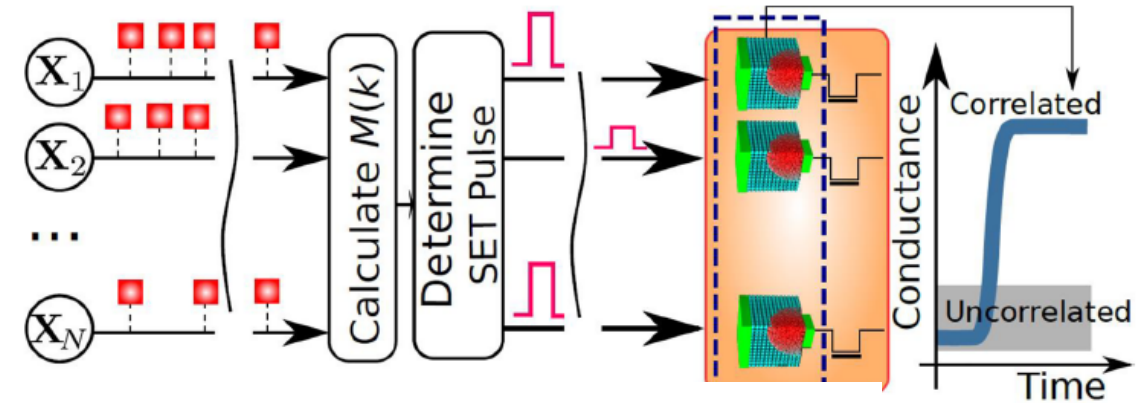
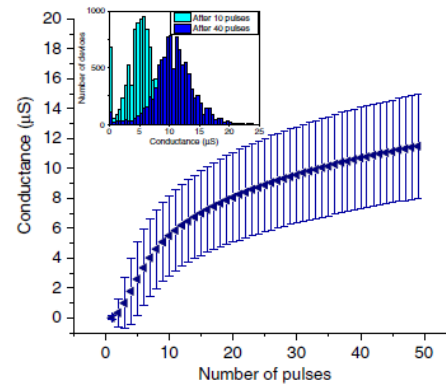
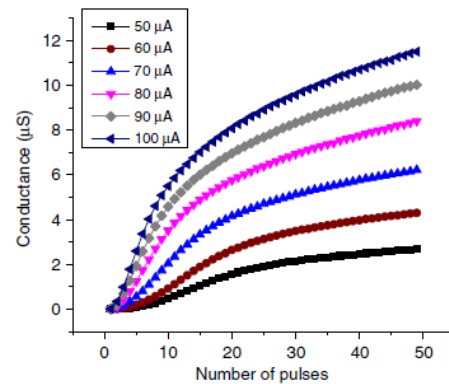
# Thank you for listening!

# Demonstration of in-memory computing with PCM -21-

## - Unsupervised learning of temporal correlations

Nature Comms., 2017, 8, 1115

➡ Using the accumulative behavior of the PCM devices,



➡ Computation can be accelerated by a factor of 200 relative to using 4 GPU devices, And also energy saving is over two orders of magnitude. (Assumed the PCM write latency of 100ns and SET programming energy of 1.5pJ)