

# Neuromorphic computing : AI needs new hardware



Julie Grollier<sup>1</sup>

Mathieu Riou<sup>1</sup>, Philippe Talatchian<sup>1</sup>, Jacob Torrejon<sup>1</sup>, Miguel Romera<sup>1</sup>, Flavio Abreu Araujo<sup>1</sup>, Paolo Bortolotti<sup>1</sup>, Juan Trastoy<sup>1</sup>, Vincent Cros<sup>1</sup>, Guru Khalsa<sup>2</sup>, Mark Stiles<sup>2</sup>, Sumito Tsunegi<sup>3</sup>, Kay Yakushiji<sup>3</sup>, Akio Fukushima<sup>3</sup>, Hitoshi Kubota<sup>3</sup>, Shinji Yuasa<sup>3</sup>, Damir Vodenicarevic<sup>4</sup>, Tifenn Hirtzlin<sup>4</sup>, Maxence Ernoult<sup>4</sup>, Nicolas Locatelli<sup>4</sup>, Damien Querlioz<sup>4</sup>

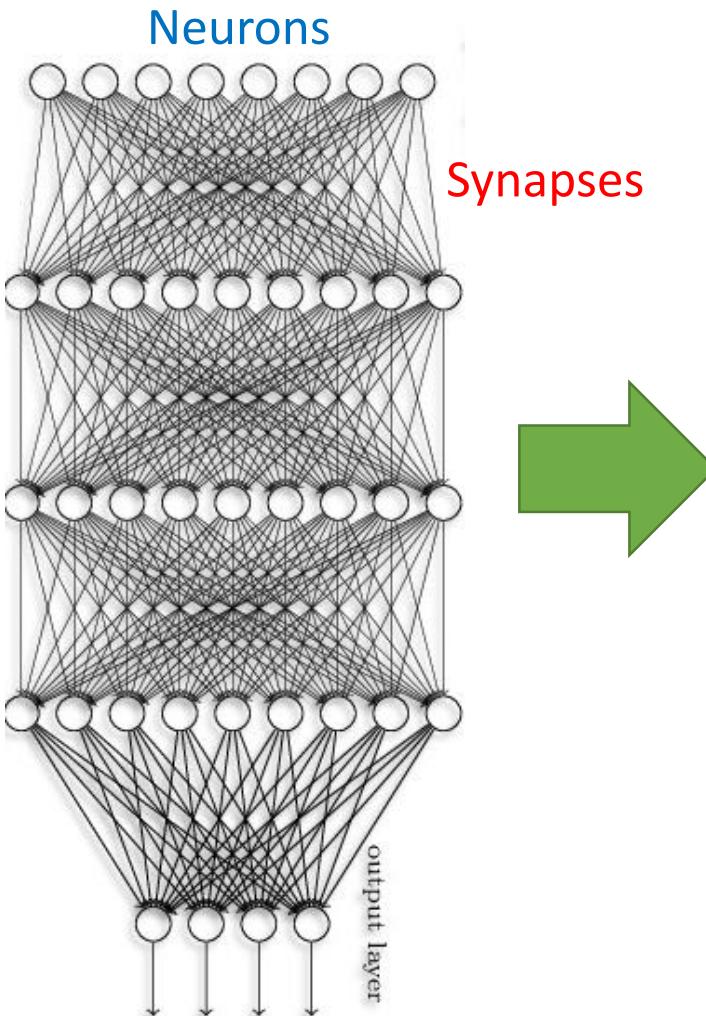
<sup>1</sup>CNRS/Thales, France <sup>2</sup>NIST, USA <sup>3</sup>AIST, Japan <sup>4</sup>C2N, France



THALES



# Today, neural networks are translated to code, binarized, then sent to a digital processor



```
internal XML or external
* external is needed when running in static mode
*
* @var boolean
*/
define('PSI_INTERNAL_XML', false);

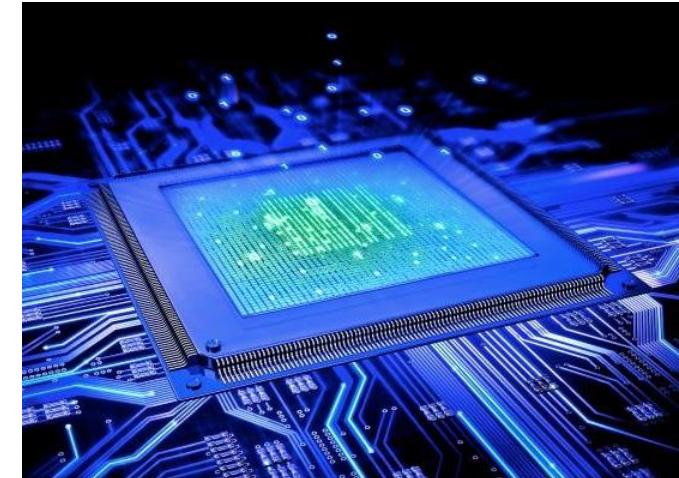
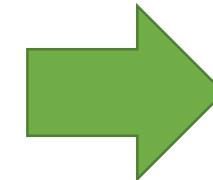
if (version_compare("5.2", PHP_VERSION, ">")) {
    die("PHP 5.2 or greater is required!!!");
}
if (!extension_loaded("pcre")) {
    die("phpSysInfo requires the pcre extension to php in order to work
        properly.");
}

require_once APP_ROOT .'/includes/autoload.inc.php';

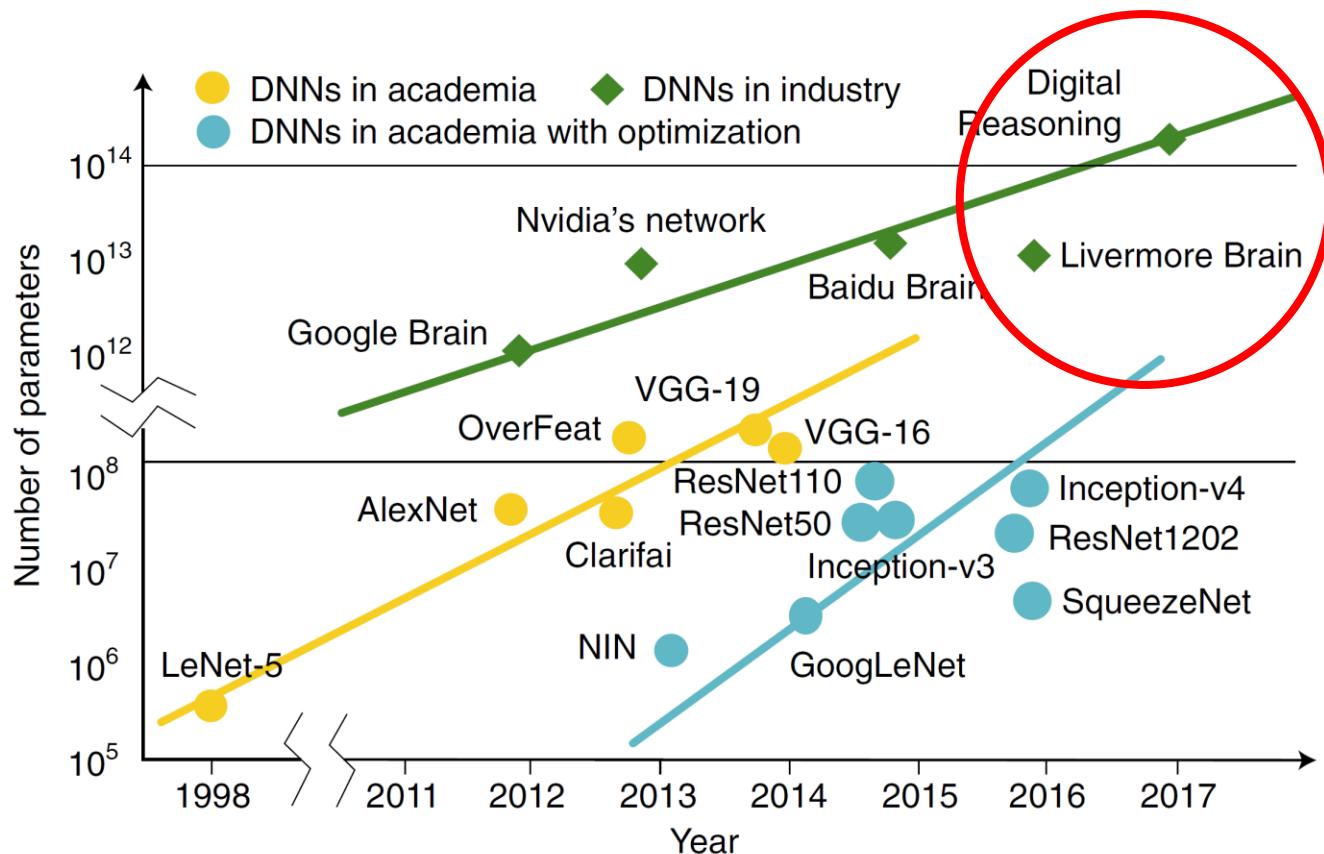
// Load configuration
require_once APP_ROOT .'/config.php';

if (!defined('PSI_CONFIG_FILE') || !defined('PSI_DEBUG')) {
    $tpl = new Template("/templates/html/error_config.html");
    echo $tpl->fetch();
    die();
}
```

0011011



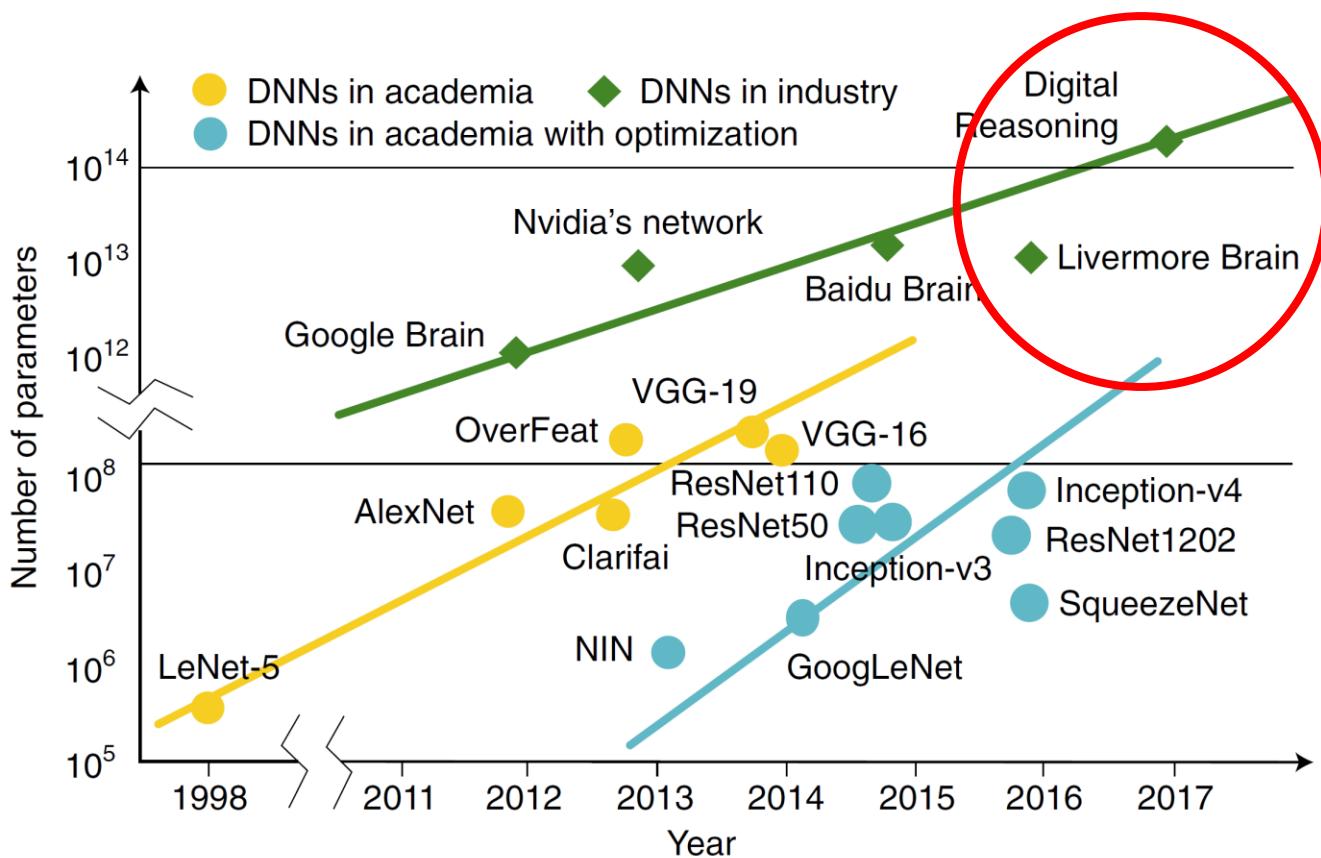
# Current processors cannot sustain AI demands



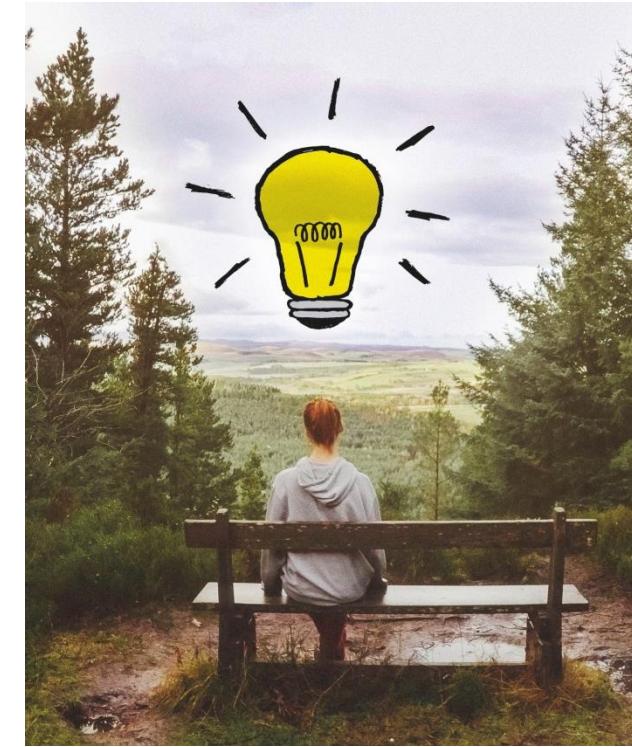
MWatts, slow, no real time learning



# The brain is larger than all these neural networks and consumes only 20W



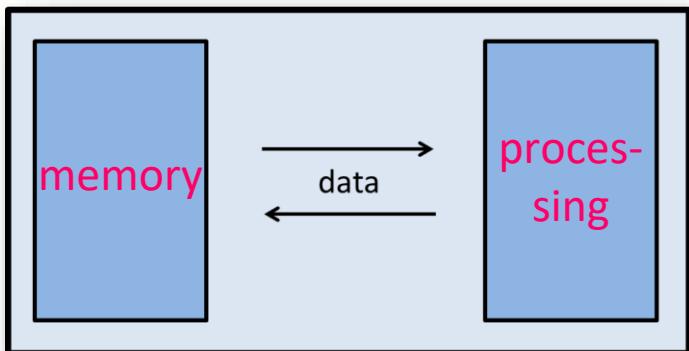
10<sup>11</sup> neurons, 10<sup>15</sup> synapses: 20 Watts



# Entangling memory and processing allows for fast and energy efficient computing

Digital computer

*CPUs, GPUs, TPUs*



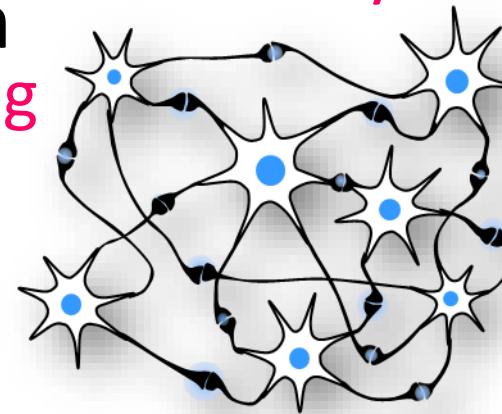
100 W/cm<sup>2</sup>

Brain

*synapse  
memory*

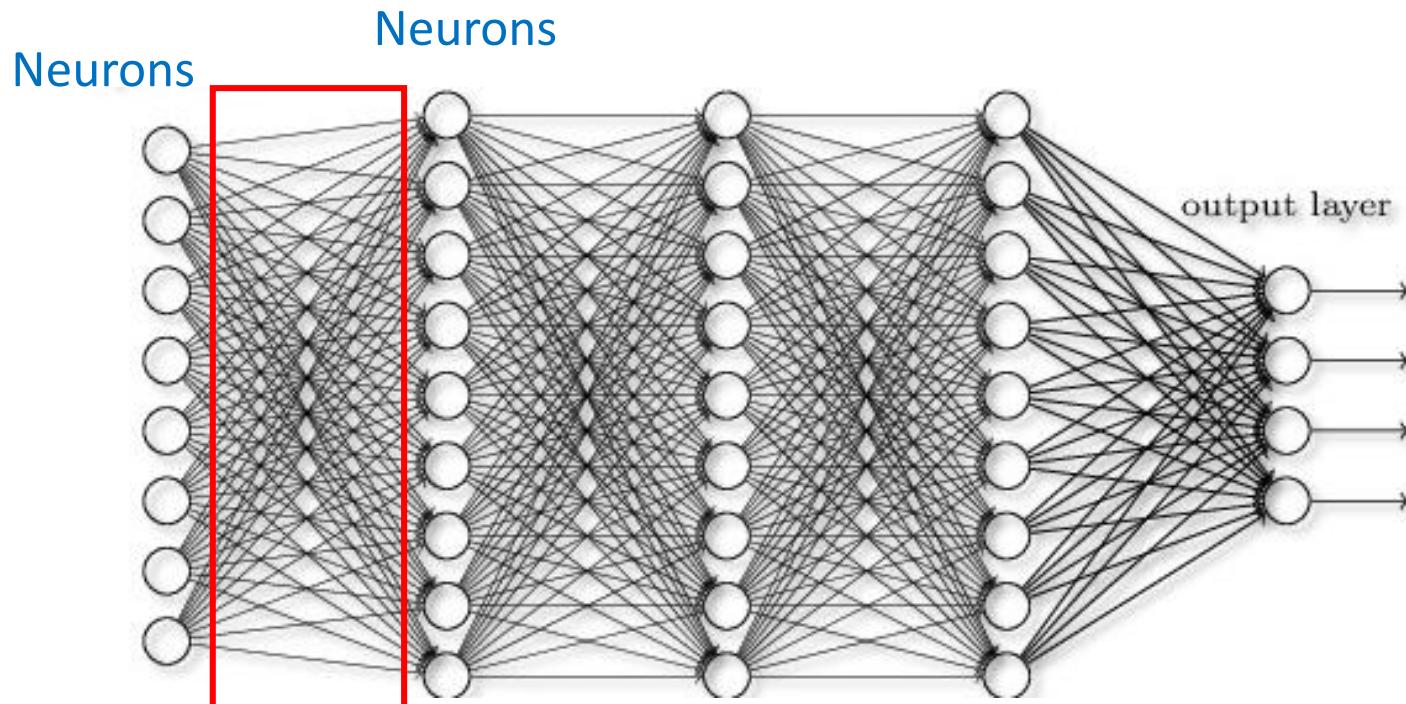
*neuron  
processing*

*neuron  
processing*



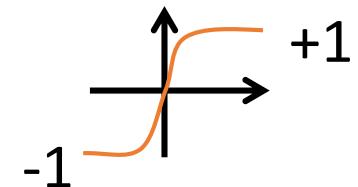
20 W in total !

# Chips implementing the layered structure of the network can enable fast, low energy computing with real-time learning

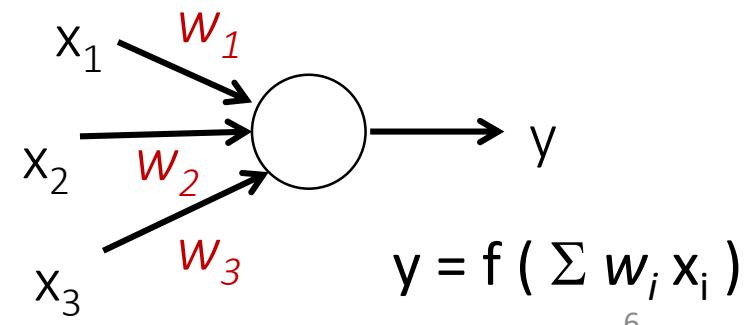


Synapses = in-situ memory

**Hardware neurons:**  
non-linear



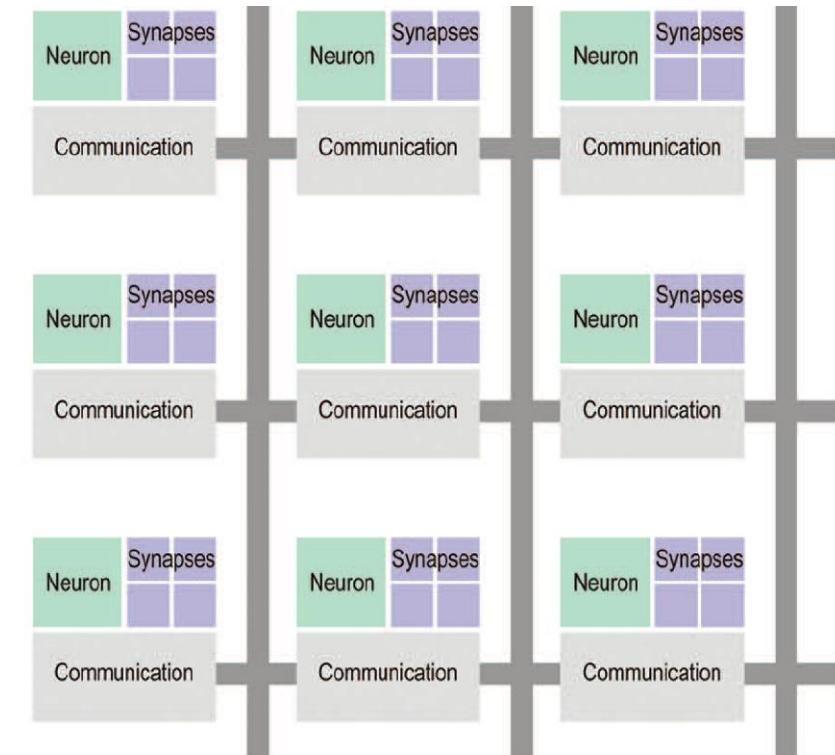
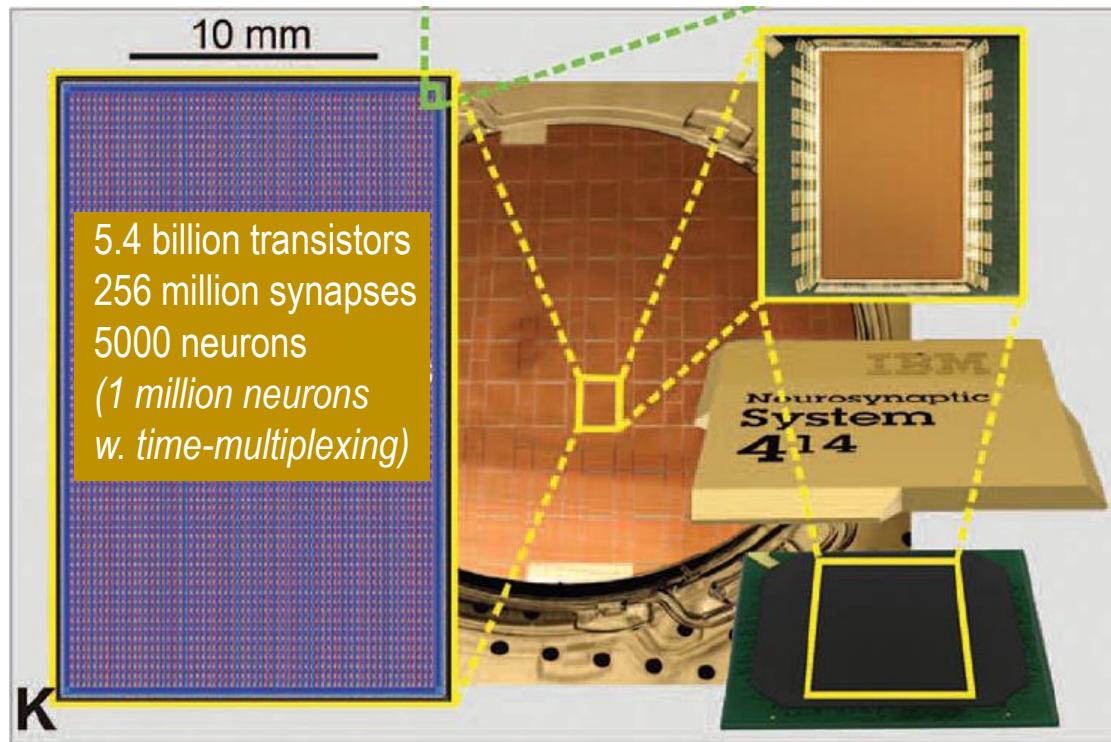
**Hardware synapses:** analog  
memory values (weights w)



- Limitations of CMOS chips
- Current efforts to implement low energy deep learning with emerging nanodevices
- Futuristic ideas

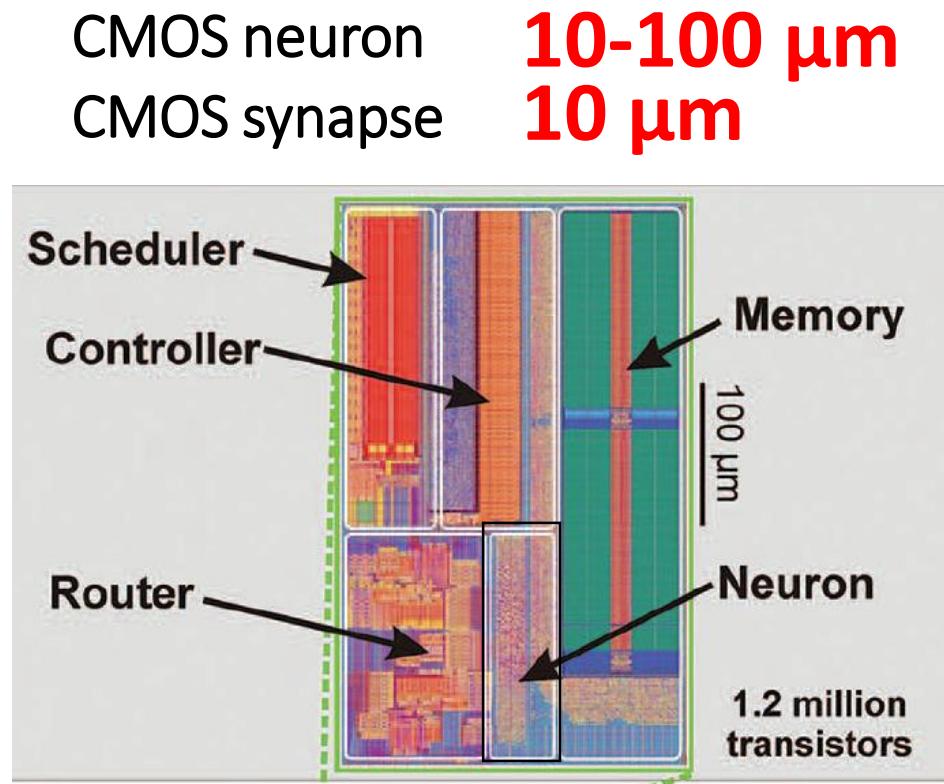
# IBM's TrueNorth neuromorphic CMOS chip

- Highly parallel, colocalized memory and processing
- Low power consumption 20 mW/cm<sup>2</sup> (processor 100 W/cm<sup>2</sup>)
- Cannot learn



# CMOS neurons and synapses are complex circuits

- A transistor is nanoscale but it is just a switch
- CMOS does not provide memory (volatile)



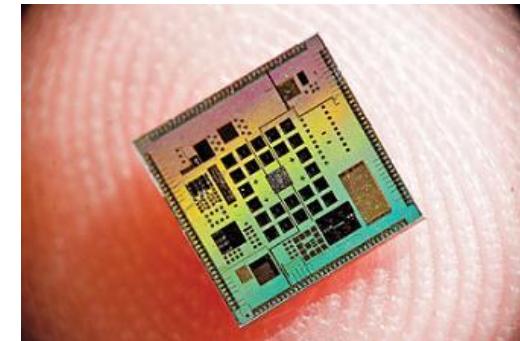
Merolla et al, *Science* 345, 668 (2014)



Brainscales 20 wafer machine. 4M neurons, 1B synapses

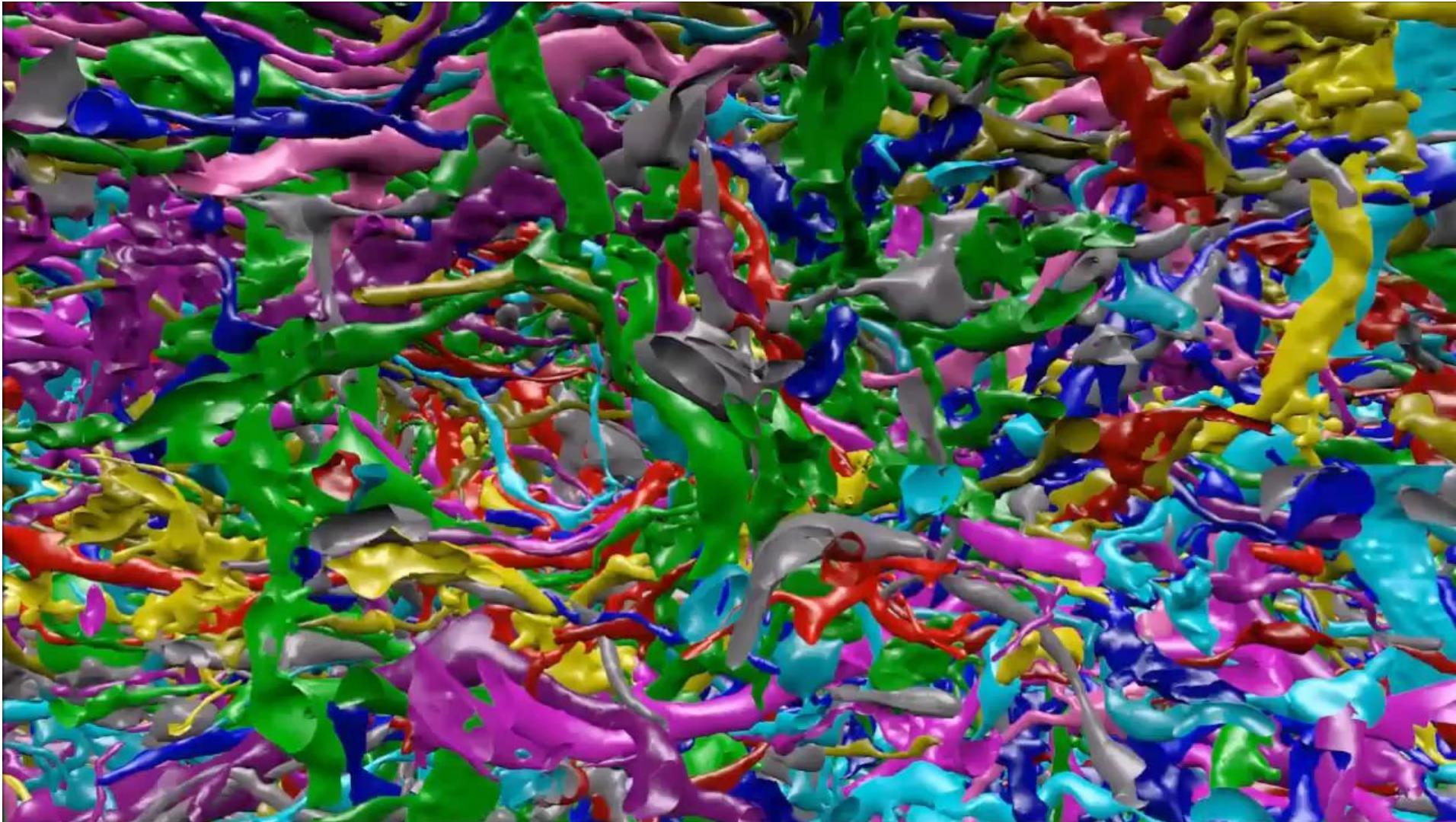
# **Building small neuromorphic chips requires to complement CMOS with nanoscale devices emulating synapses and neurons**

Hundred millions of neurons and synapses in a  $1 \text{ cm}^2$  chip  
→ Each device smaller than  $1 \mu\text{m}^2$



→ **nanosynapses, nanoneurons**

$10^4$  synapses / neurones =  $10^4$  wires/neurons



Moritz Helmstaedter lab, retina flight 2013

- Limitations of CMOS chips
- **Current efforts to implement low energy deep learning with emerging nanodevices**
- Futuristic ideas

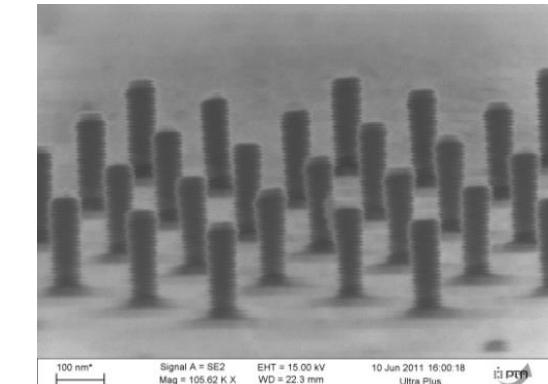
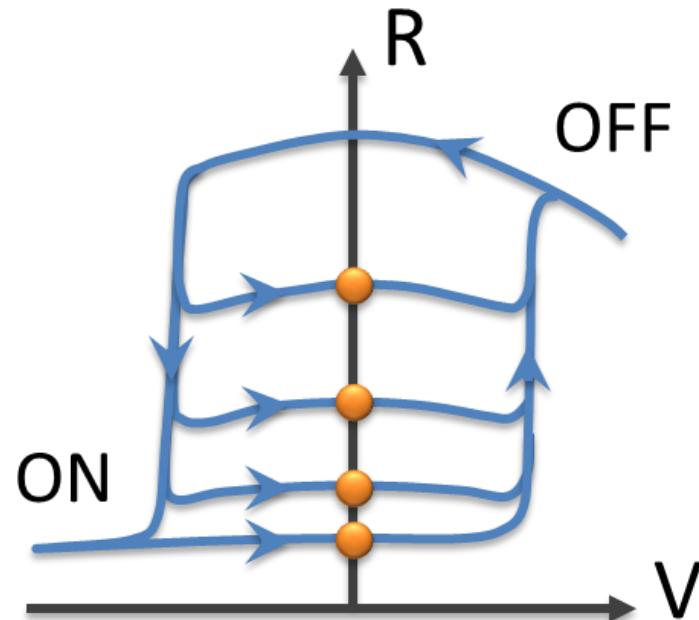
- Limitations of CMOS chips
- **Current efforts to implement low energy deep learning with emerging nanodevices**
- Futuristic ideas

*This talk: electronic nanodevices*

# Memristors implement the fundamental ingredients of synapses

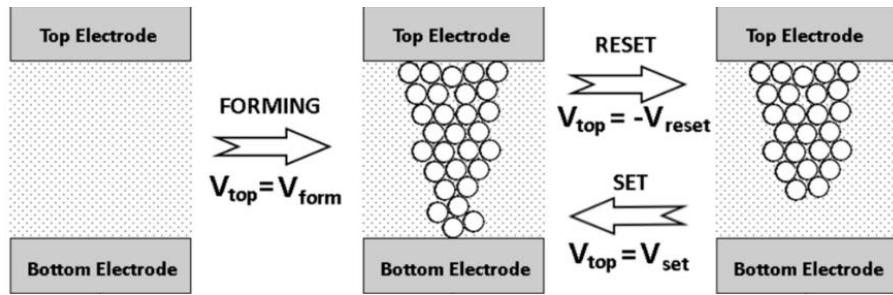
Tunable nano-resistors with memory

Chua, *IEEE Trans.  
Circuit Theory* (1971)



# There are many different flavors of memristors

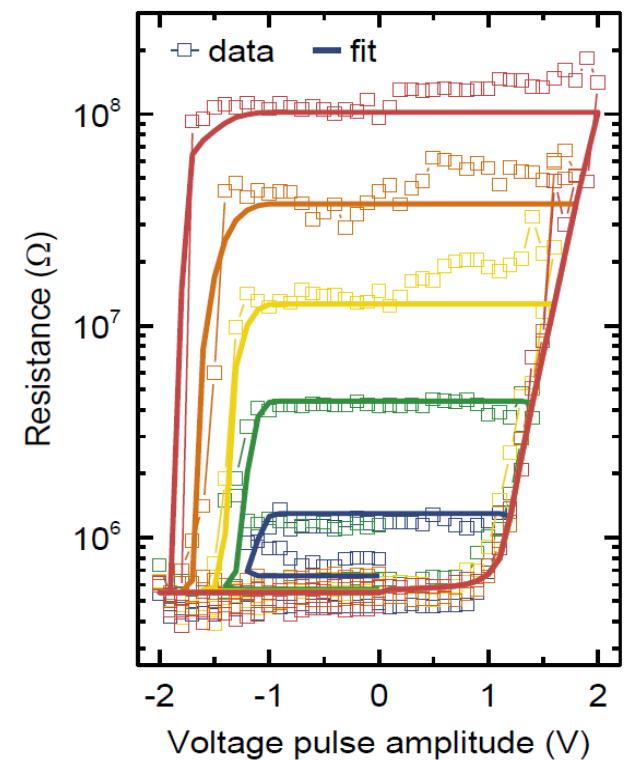
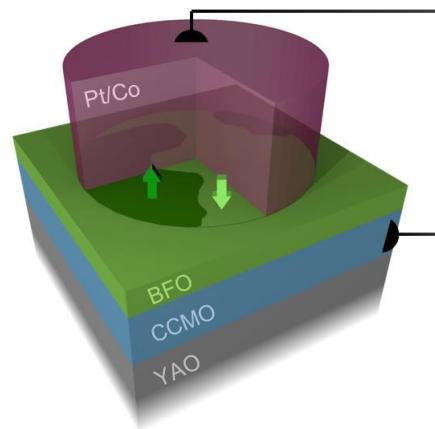
## Standard memristors



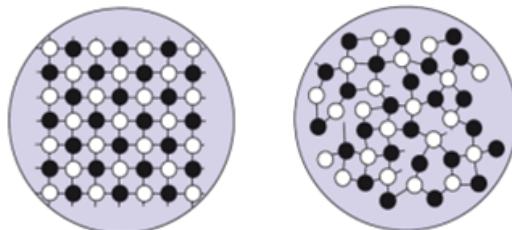
Yang et al., *Nat. Nano.* 8, 13 (2013)

Red-Ox

## Ferroelectric memristors



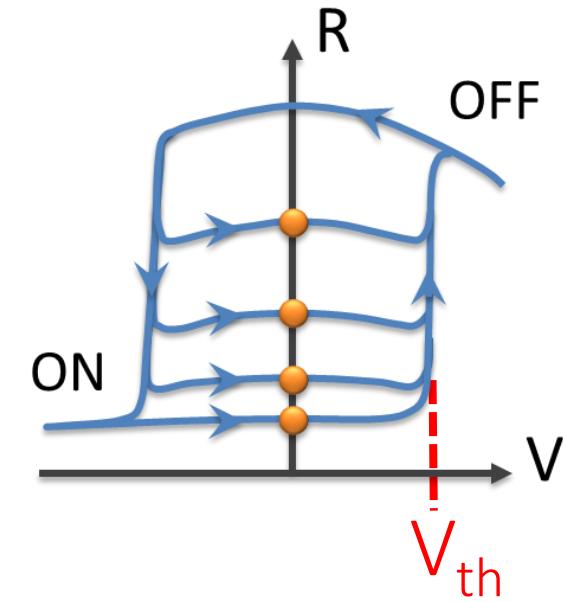
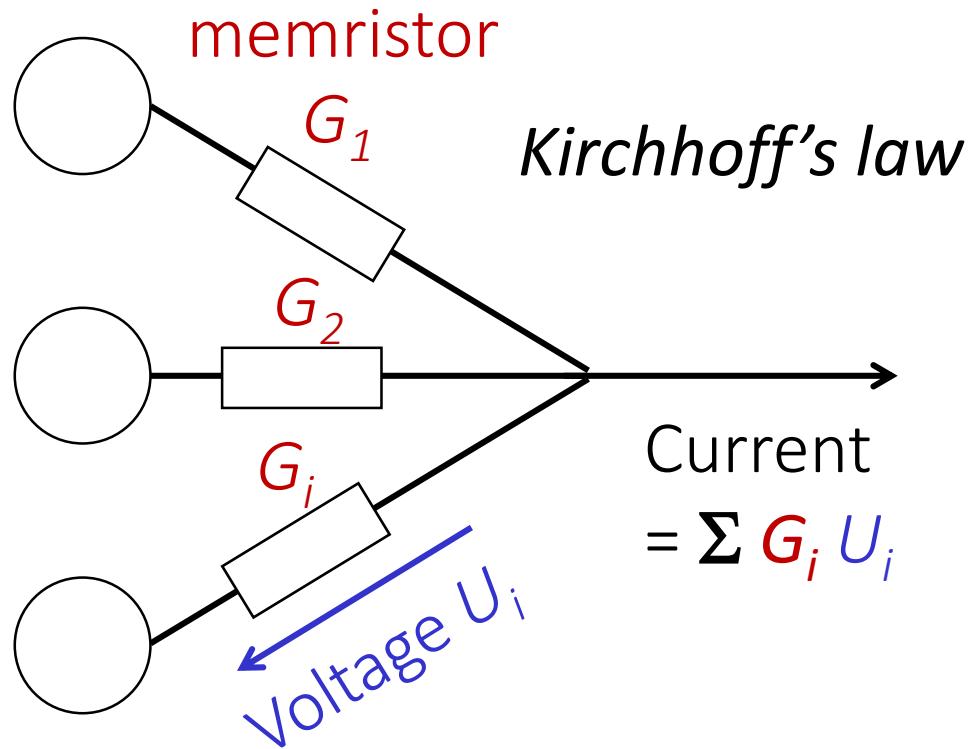
Phase change



Kuzum et al, *Nanotechnology* 24, 382001 (2013)

André Chanthbouala, JG et al, *Nat. Mat.* 11, 860 (2012)

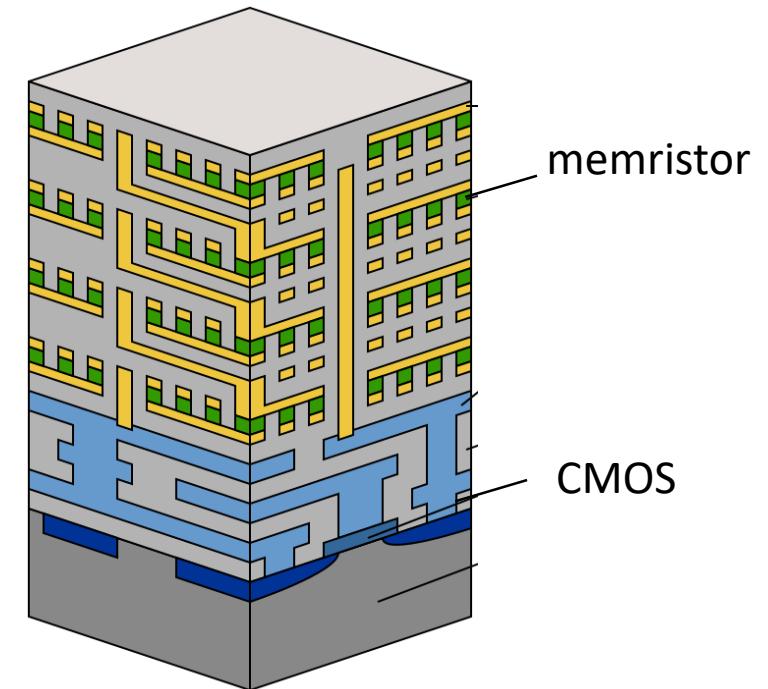
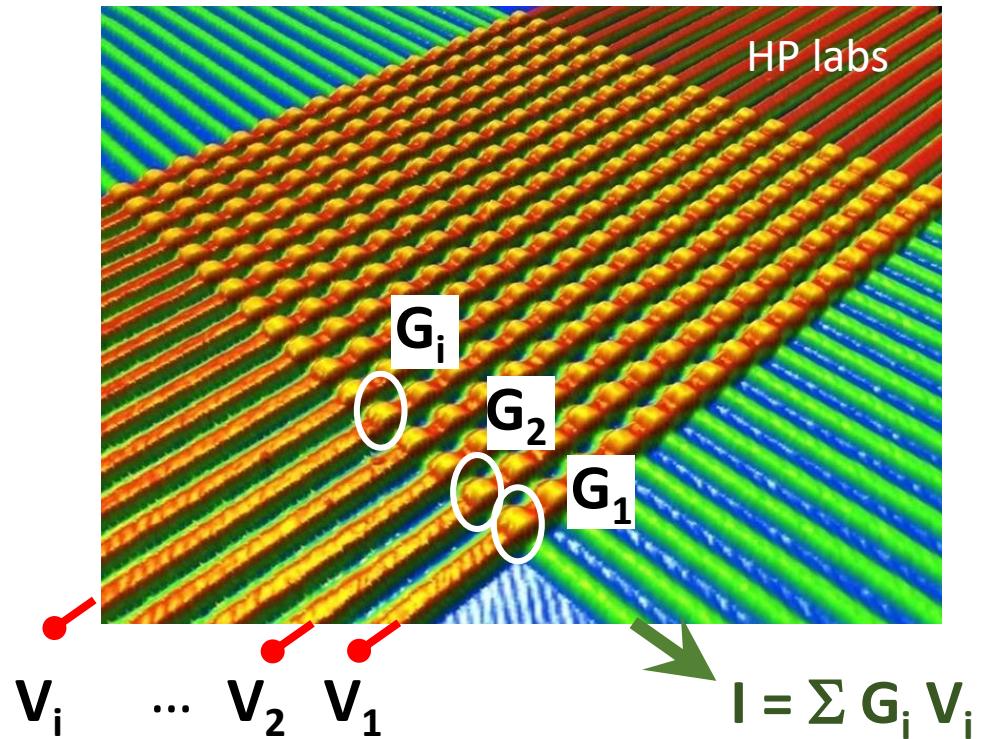
# Memristors emulate electronic synapses: the weight is their tunable conductance $G$



$U_i < V_{th}$  : calculating mode

$U_i > V_{th}$  : learning mode

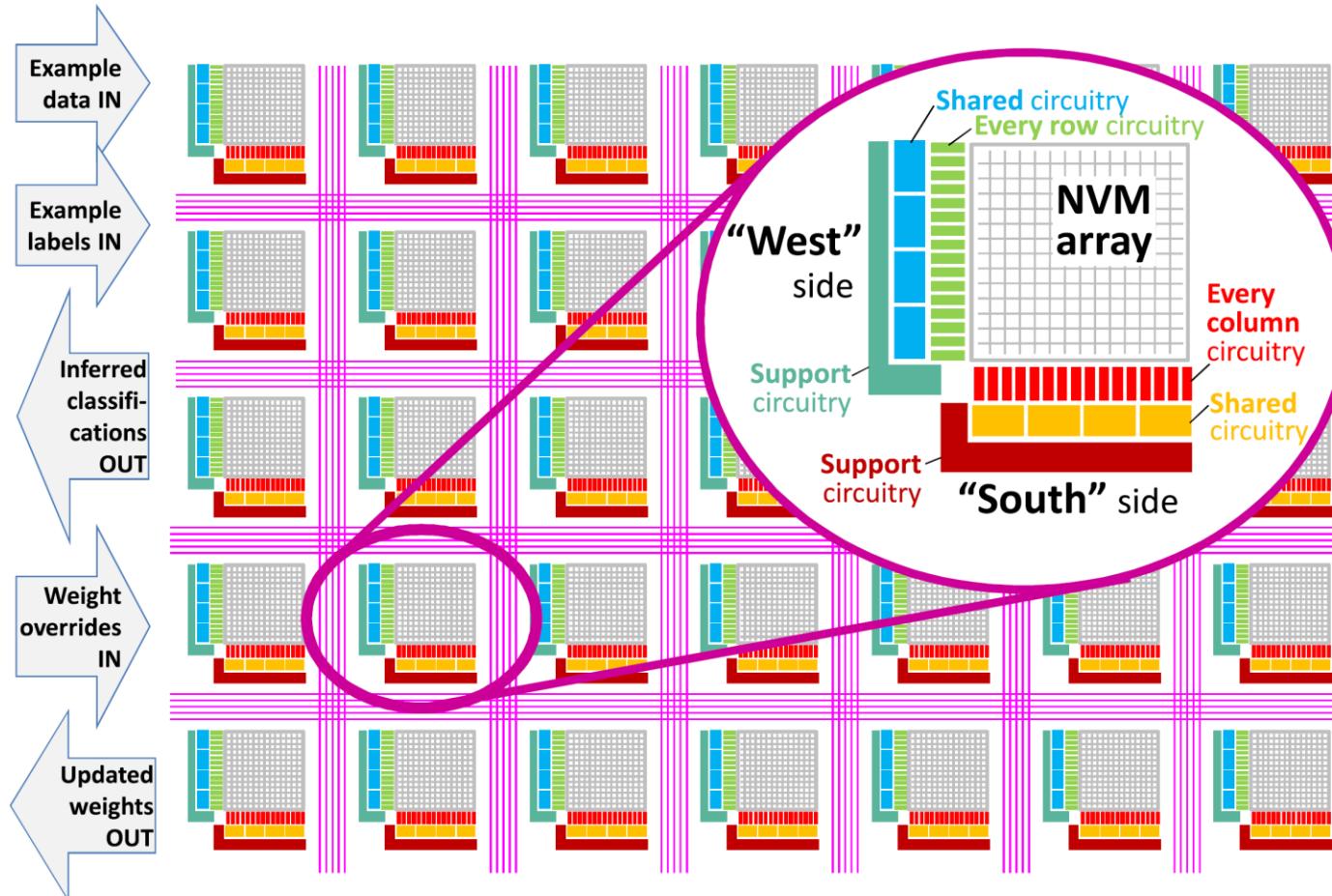
# Memristor crossbar arrays can fully connect layers of > 100 neurons



10 000 synapses per neuron ?

# Memristors open a path to energy-efficient real-time learning, but there is no chip yet

Expected : 100 x faster than GPUs and 100 x less energy consumption



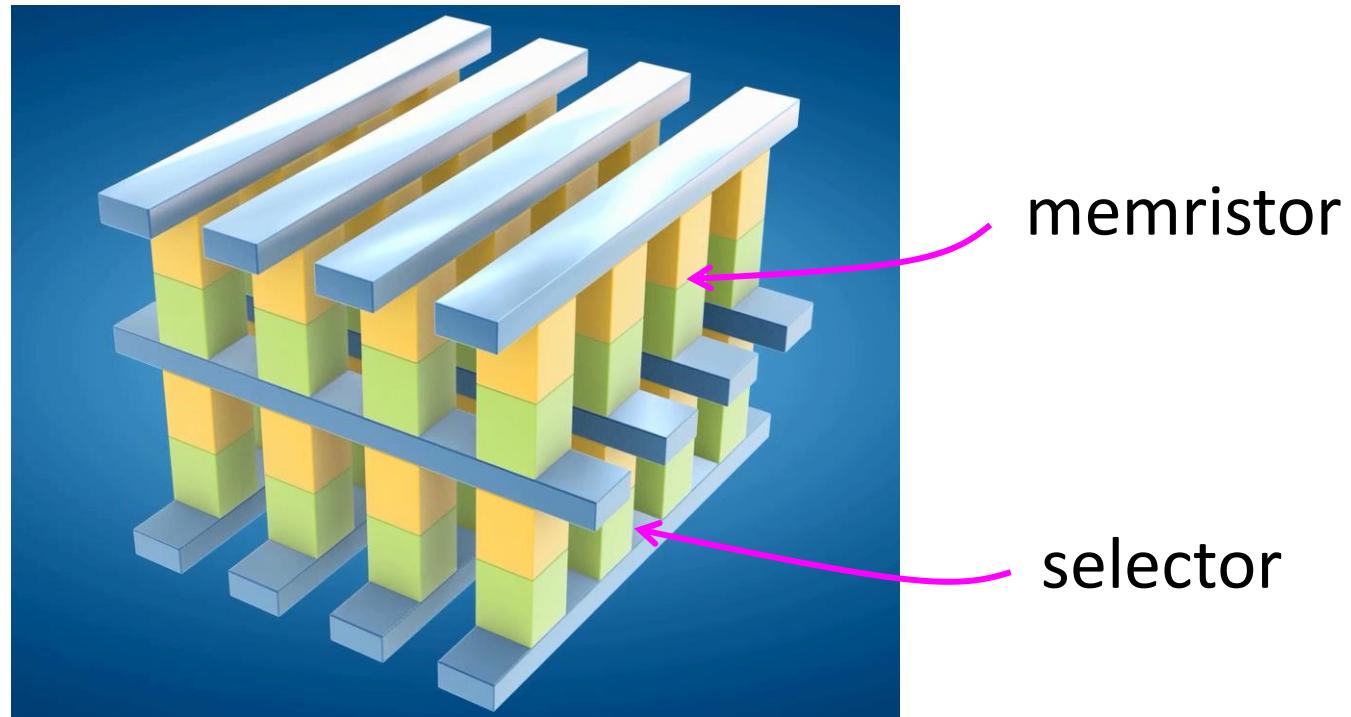
Narayan et al, IBM J.  
RES. & DEV. 61, NO.  
4/5, 11 (2017)

Ambrogio et al,  
Nature 558, 60  
(2018)

# One challenge being currently tackled is building large arrays of memristors

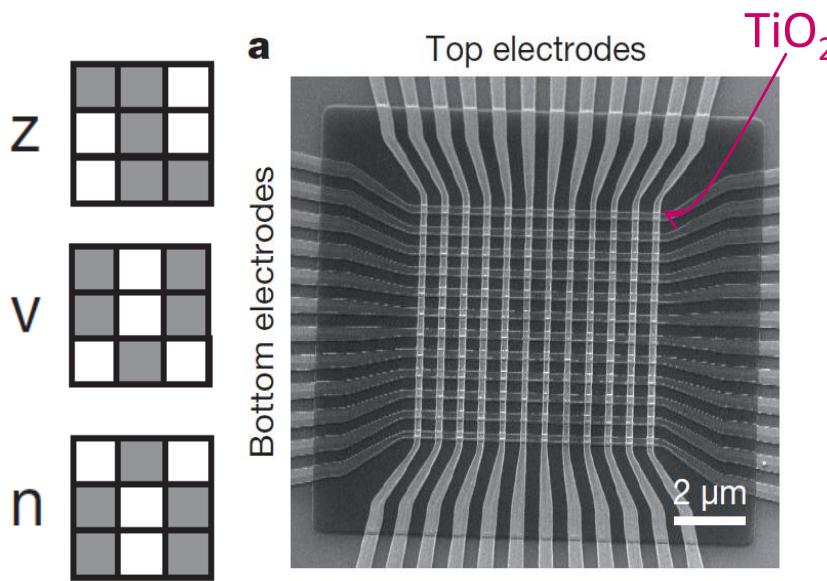
Recent progress has been achieved towards the commercialization of binary memories made of memristor arrays

- 3D Xpoint,  
Intel/Micron  
Optane Lenovo  
32 Gbits

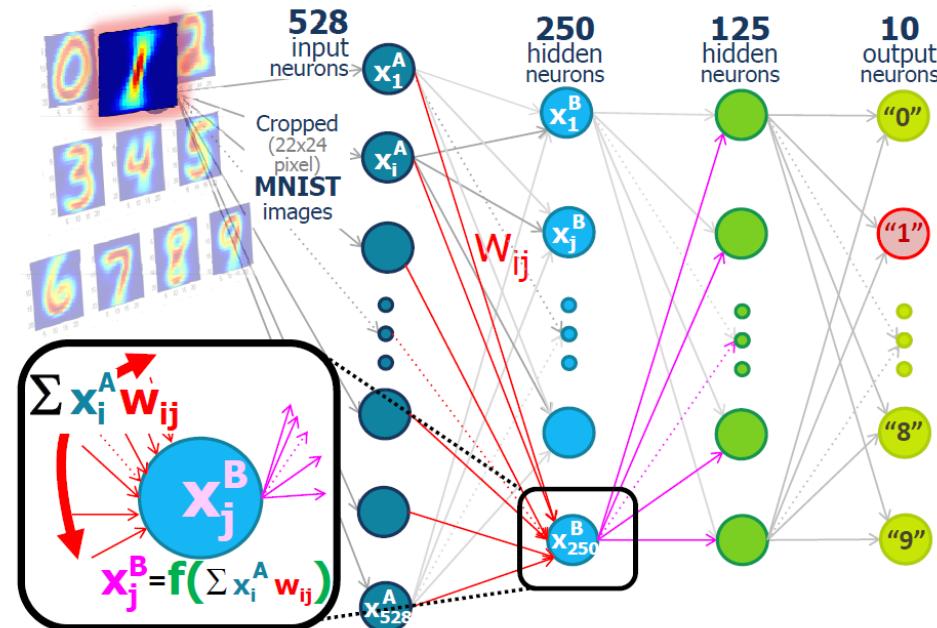


# First results of supervised learning through back-propagation have been obtained but the accuracy is low due to devices imperfections (non-linear, imprecise, noisy)

Prezioso et al, *Nature* 521, 61 (2015)



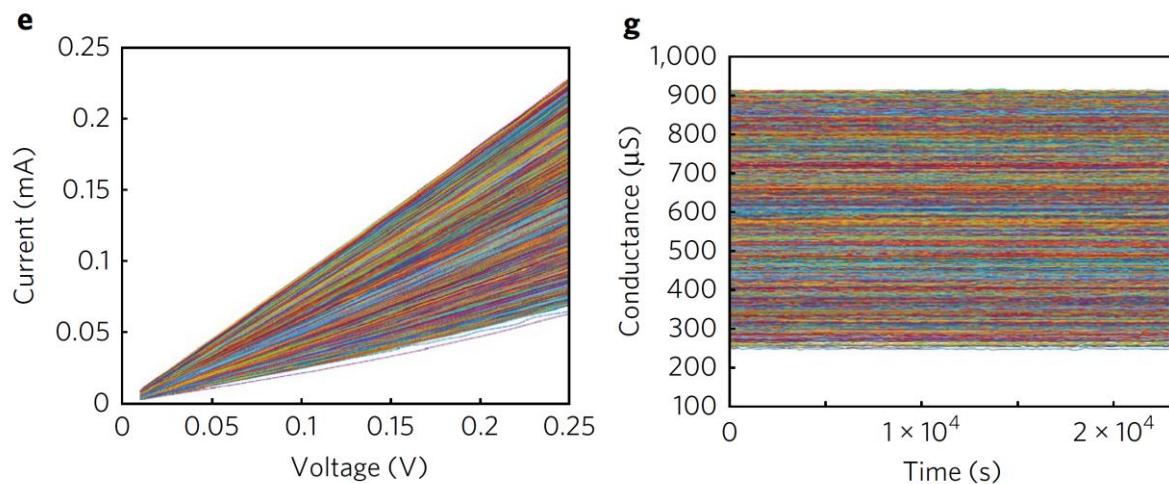
Burr et al, *IEEE IEDM* (2014)



IBM experiments: Handwritten digit recognition with ~ 165 000 synapses (phase change with selector)

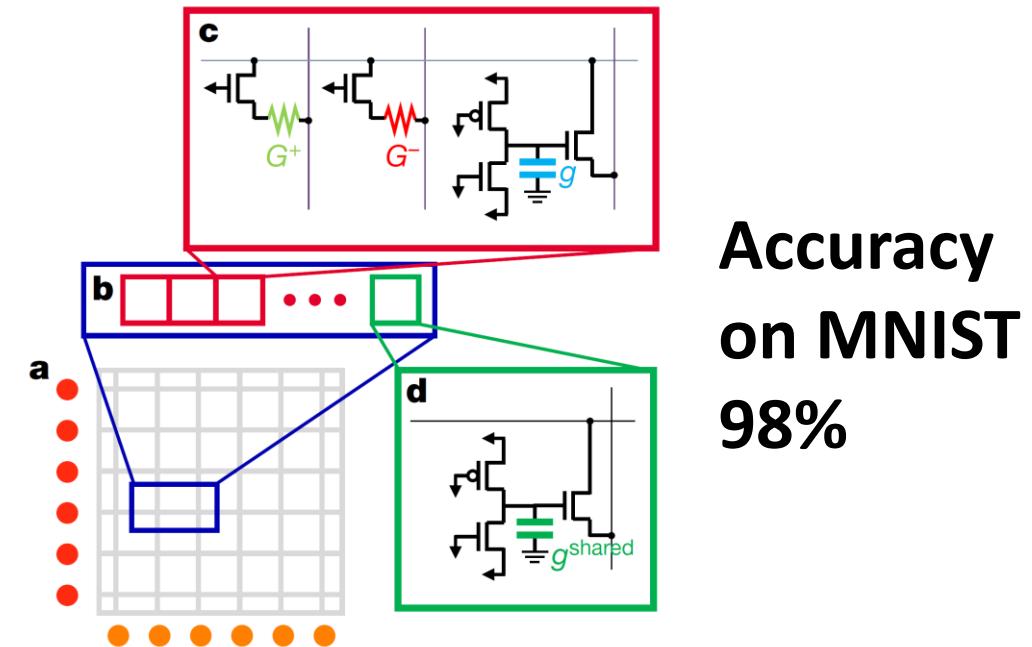
# A solution towards implementing backpropagation is to improve nanodevice properties

- Optimizing memristor properties



Li et al, Nature electronics 1, 52–59 (2018)

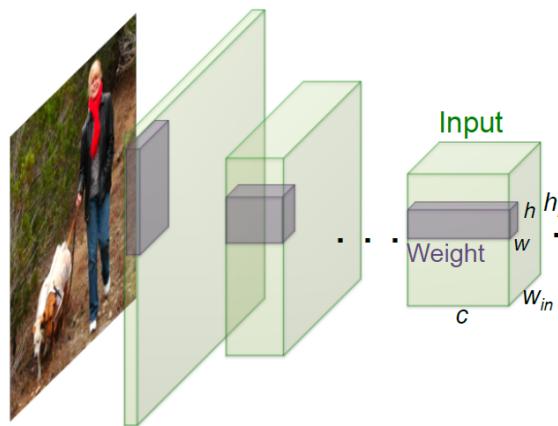
- Complementing memristors with high accuracy weight

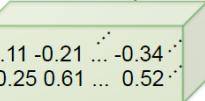
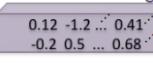
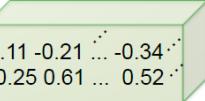
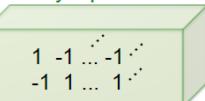
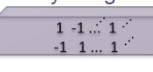


Ambrogio et al, Nature 558, 60 (2018)

# In parallel, there is a lot of work on finding algorithms that can be accelerated on hardware

Ex : XNOR nets, Rastegari et al, arXiv:1603.05279



	Network Variations	Operations used in Convolution	Memory Saving (Inference)	Computation Saving (Inference)	Accuracy on ImageNet (AlexNet)
Standard Convolution	Real-Value Inputs  Real-Value Weights 	+ , - , ×	1x	1x	%56.7
Binary Weight	Real-Value Inputs  Binary Weights 	+ , -	~32x	~2x	%56.8
BinaryWeight Binary Input (XNOR-Net)	Binary Inputs  Binary Weights 	XNOR , bitcount	~32x	~58x	%44.2

Dates n/a. @ None

***Training Deep Neural Networks with 8-bit Floating Point Numbers***

Naigang Wang · Jungwook Choi · Daniel Brand · Chia-Yu Chen · Kailash Gopalakrishnan

NIPS 2018

Poster

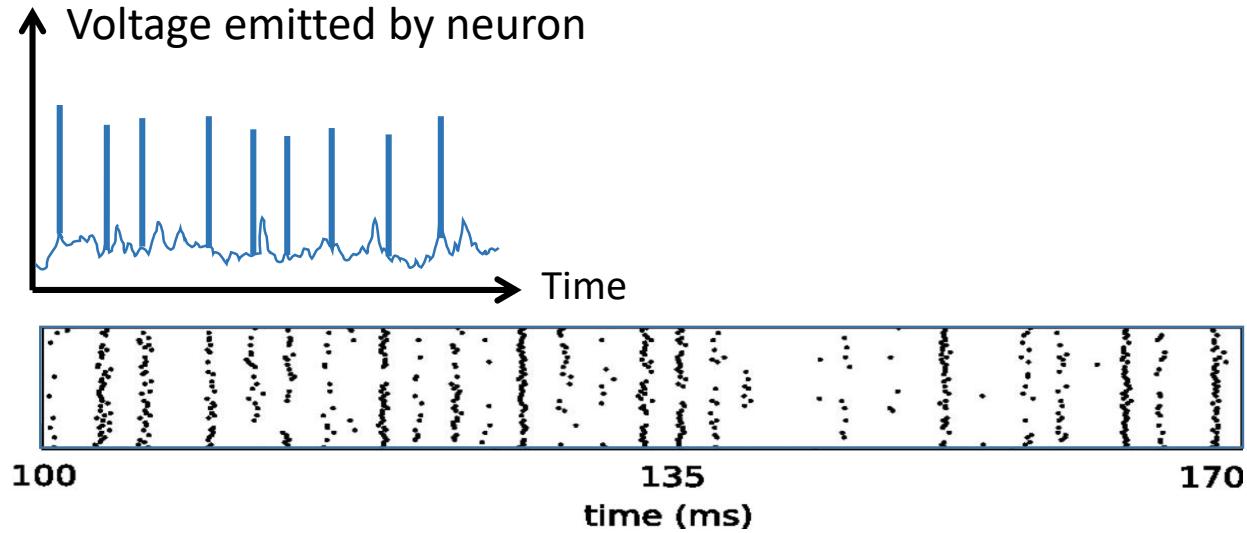
But nanodevice noise and non-linearity remain a problem

- Limitations of CMOS chips
- Current efforts to implement low energy deep learning with emerging nanodevices
- **Futuristic ideas**

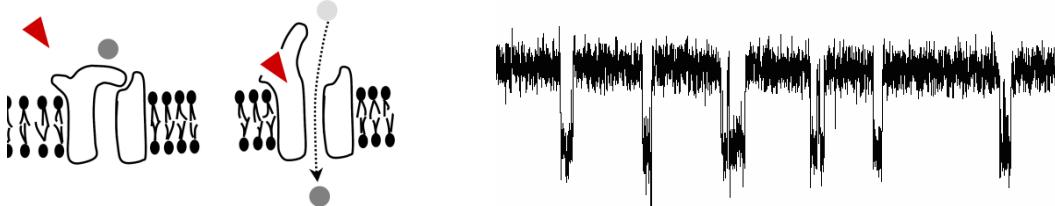
**Can noise be used instead of mitigated ?**

# Biological synapses and neurons are noisy: the brain seems to operate at the thermal limit to minimize its power consumption

Neural spikes in response to the same input recorded 50 times

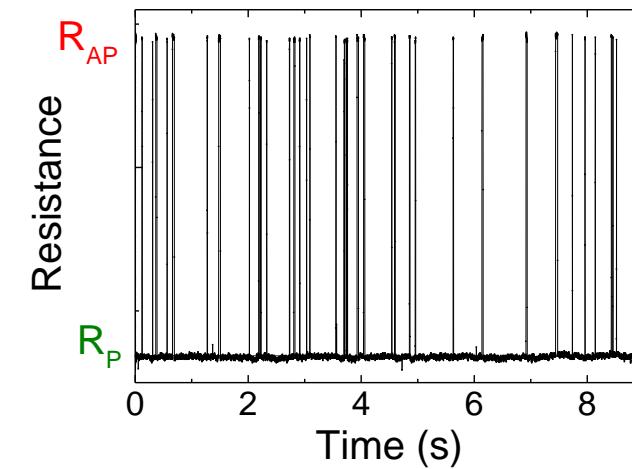
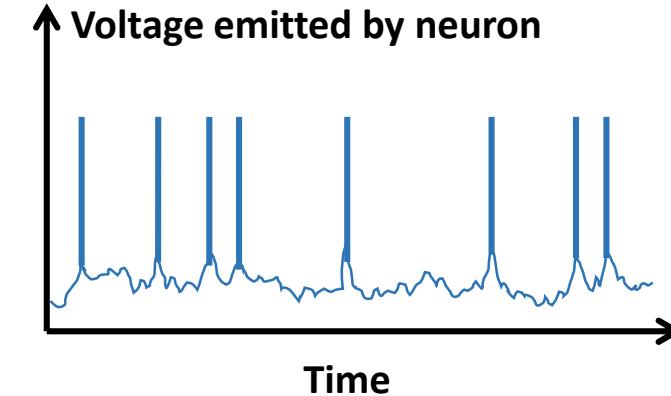
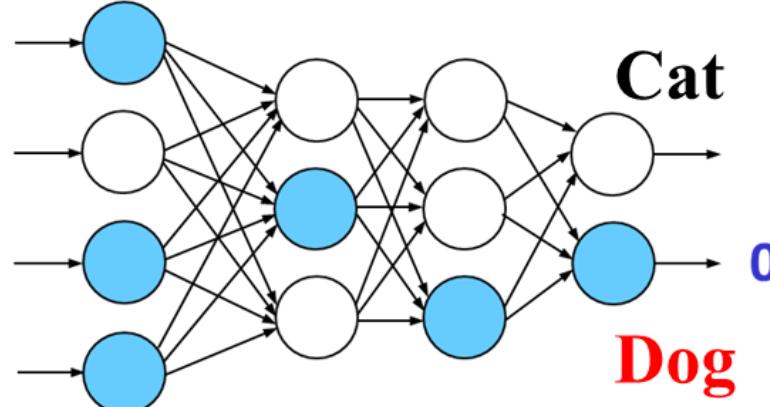
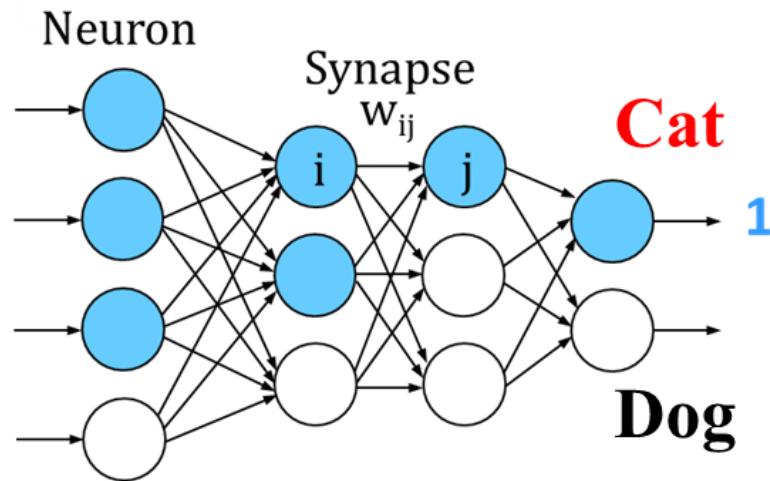


Synaptic ionic channels



Computing at low power with stochastic components is possible

# Can we achieve reliable computations with noisy devices ?

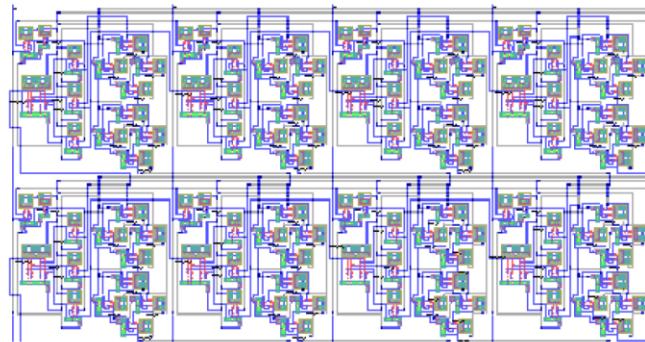


Low energy  
consumption

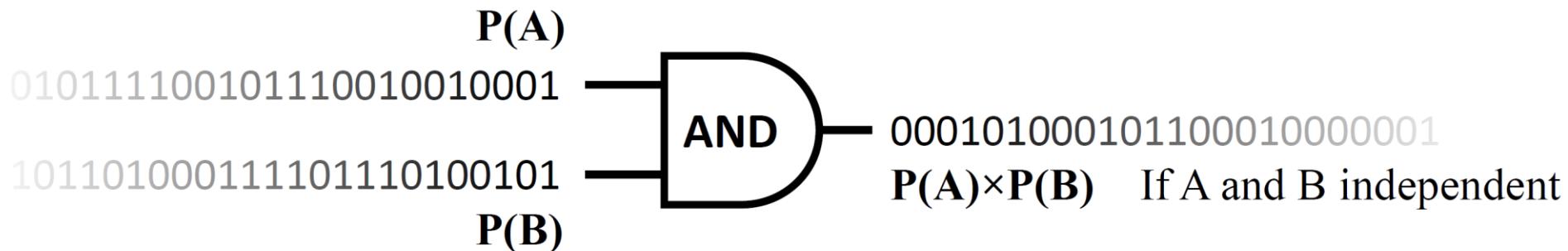
# First strategy for computing with noisy devices: accumulating evidence in time by sampling the same device several times

Example: stochastic computing

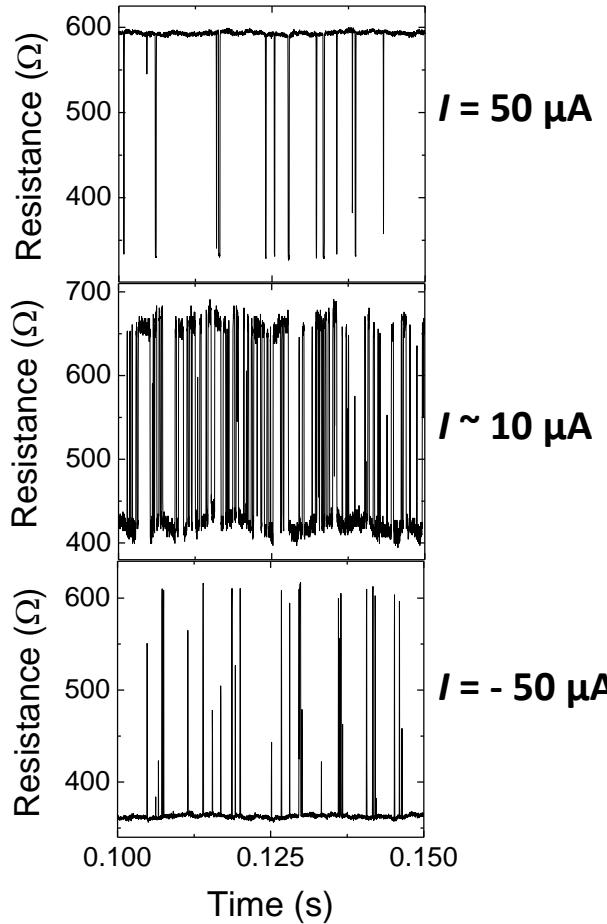
Conventional multiplier:



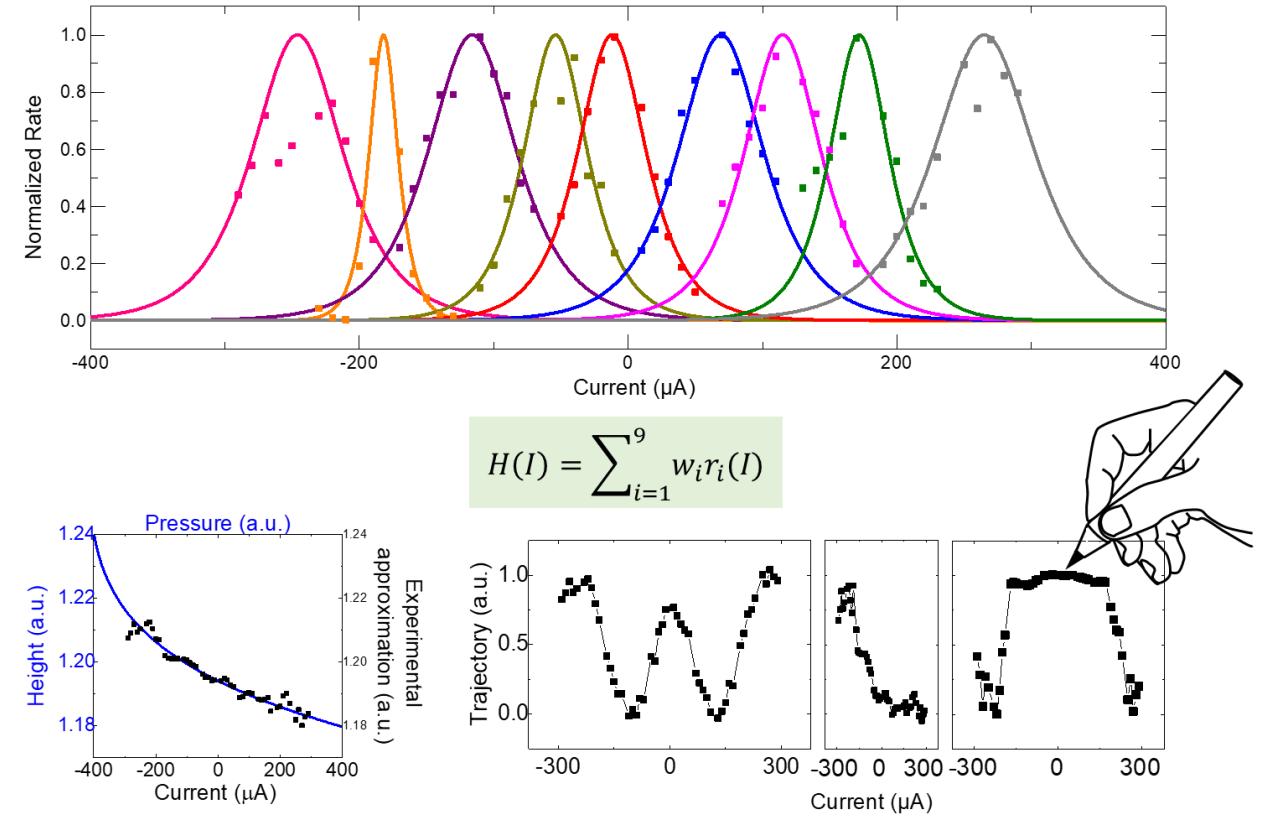
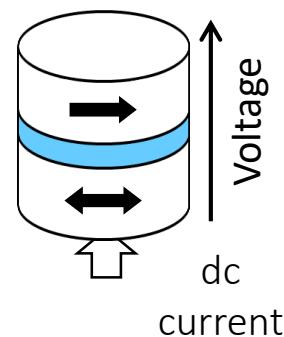
Stochastic computing multiplier:



# 2nd strategy for computing with noisy devices: sampling from populations of nanodevices (redundance)



Superparamagnetic  
tunnel junctions



# Could we leverage noise for computing ?

Randomness can be useful - Example: stochastic resonance



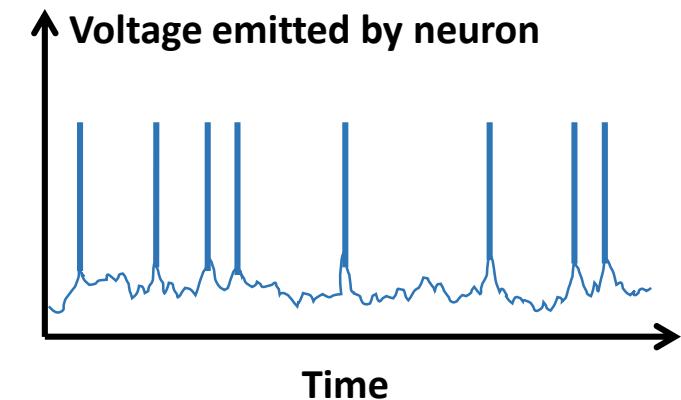
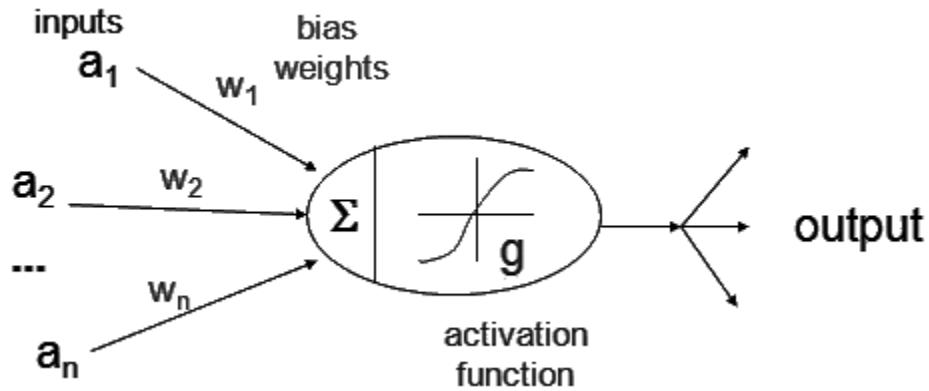
Low noise

Optimal noise

High noise

Gammaitoni et al, *Reviews Of Modern Physics*, 70(1), 223–287 (1998)

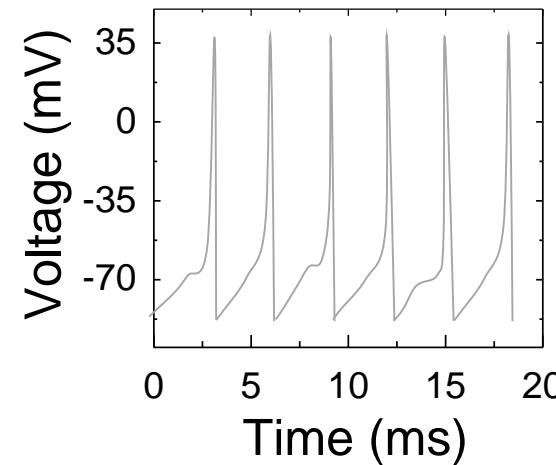
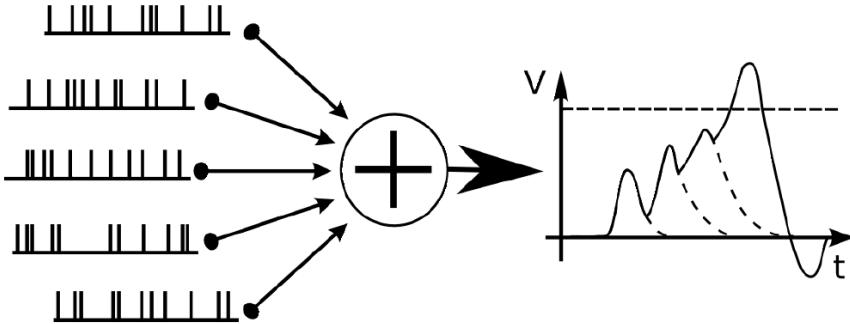
Locatelli, Mizrahi, JG et al, *Phys. Rev. Appl.* 2014



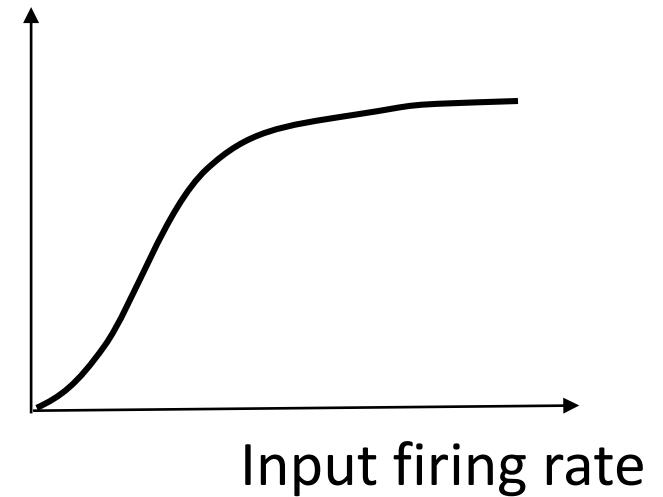
## Can spikes help computing with nanodevices?

# Biological neurons are more than non-linear functions : they spike !

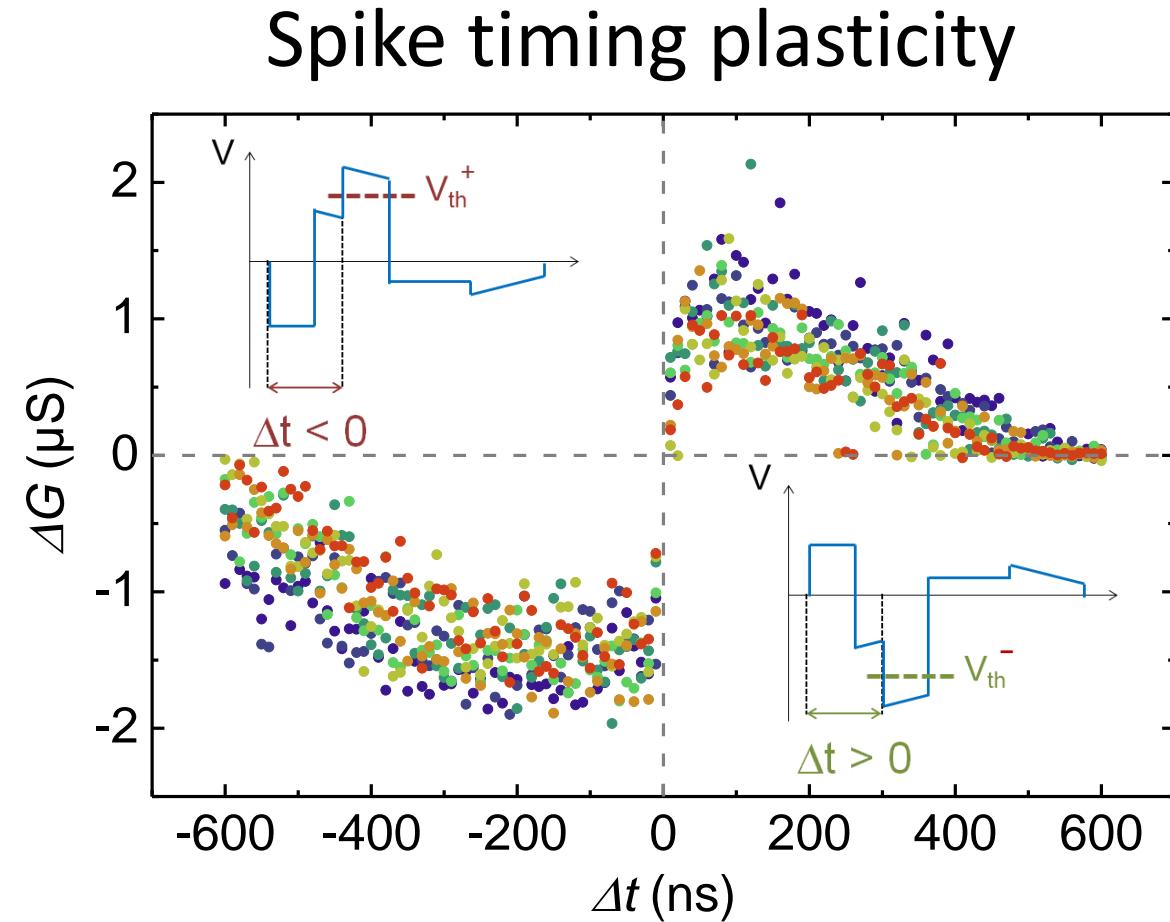
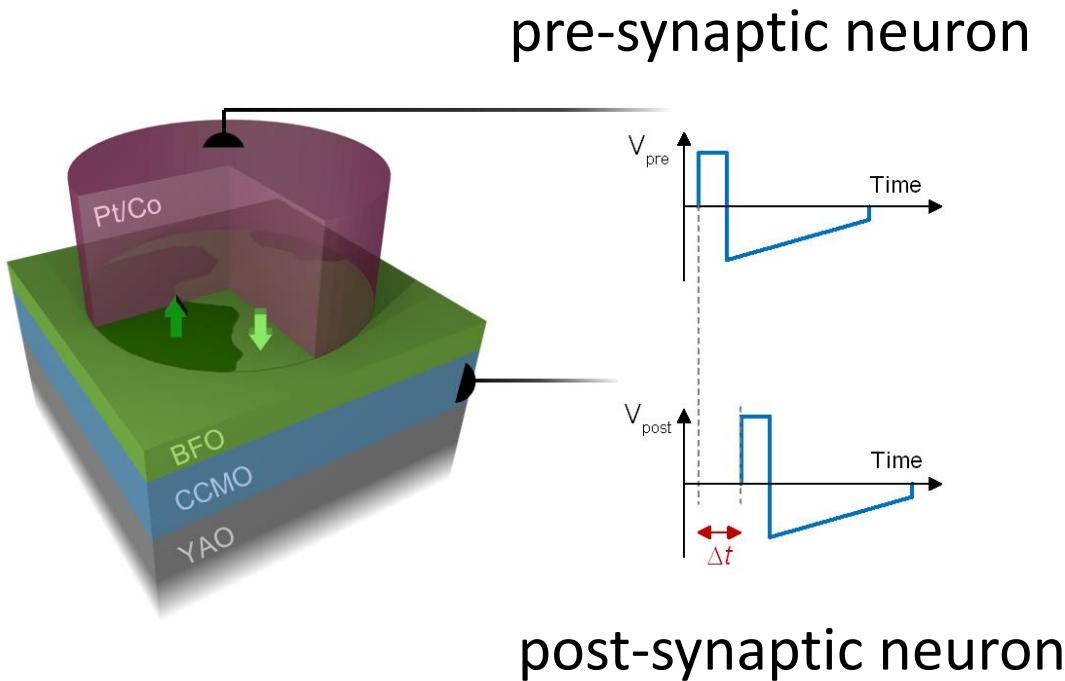
*Integration, spikes*



Output firing rate

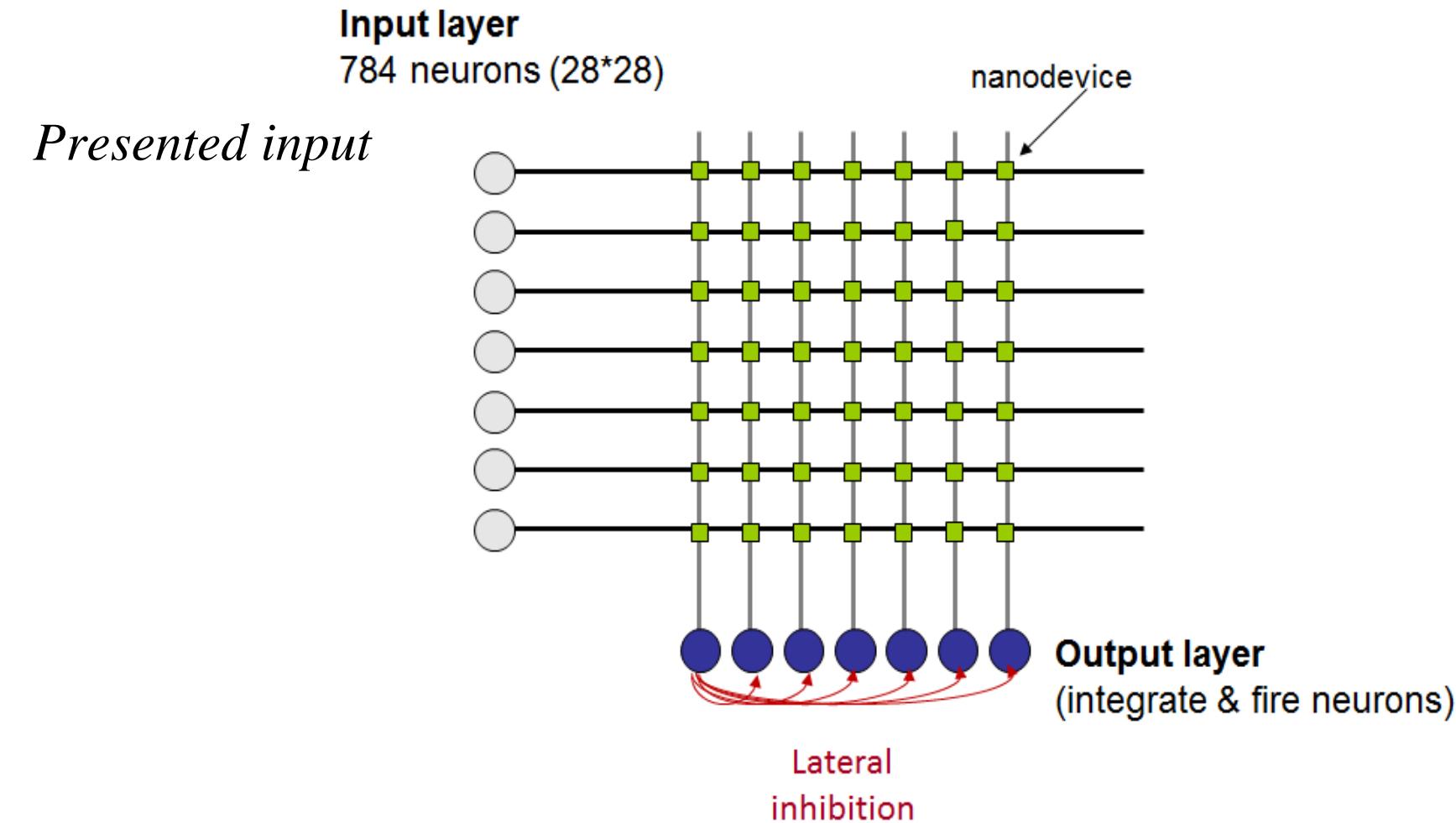


# Spikes allow 1) unsupervised learning 2) autonomous evolution of memristor conductance



Jo et al, *Nanoletters* 10, 1297 (2010) Zamarrenos-Ramos et al, *Frontiers in Neuroscience* 5, 26 (2011)  
Sören Boyn, JG et al, *Nature Com.* 8, 14736 (2017)

# How memristor learn without supervision through STDP

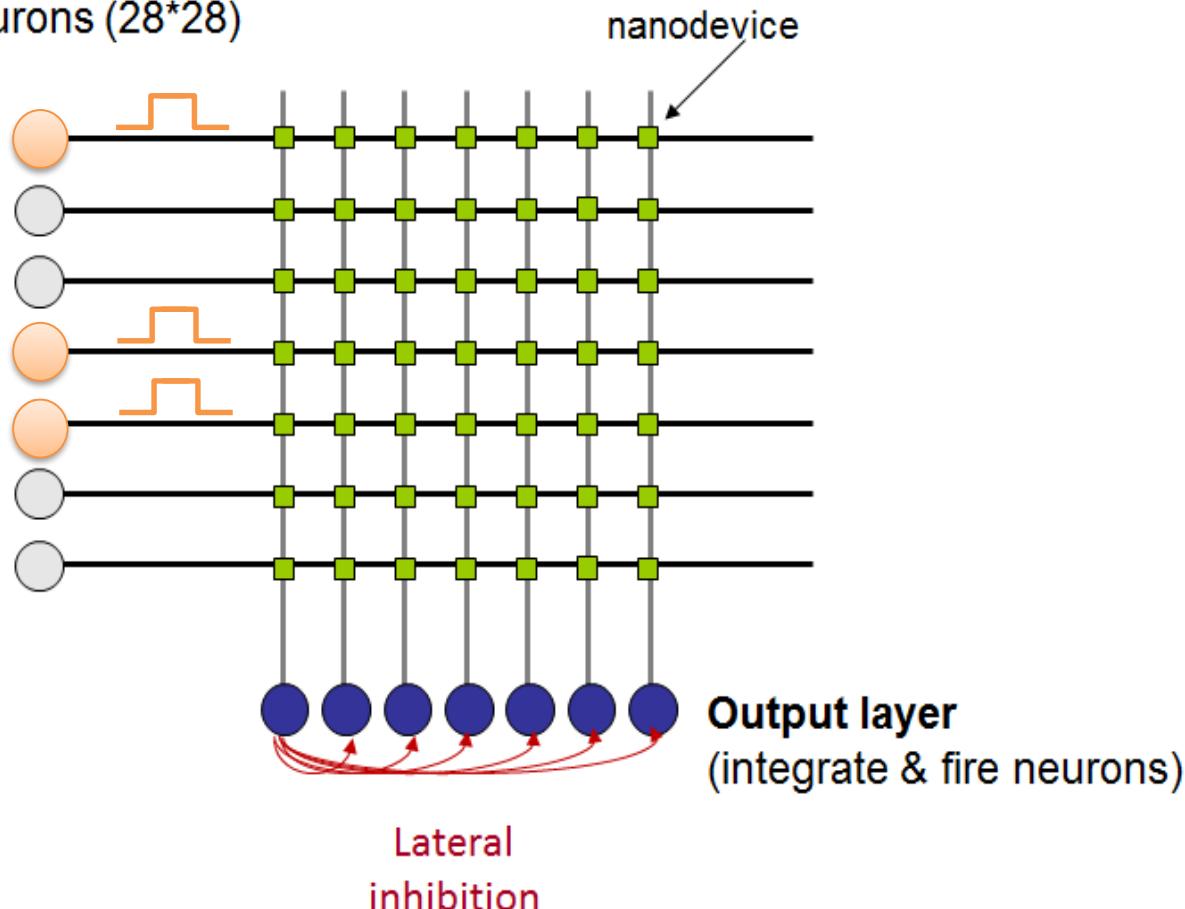


# How memristor learn without supervision through STDP

*Presented input*



**Input layer**  
784 neurons ( $28 \times 28$ )

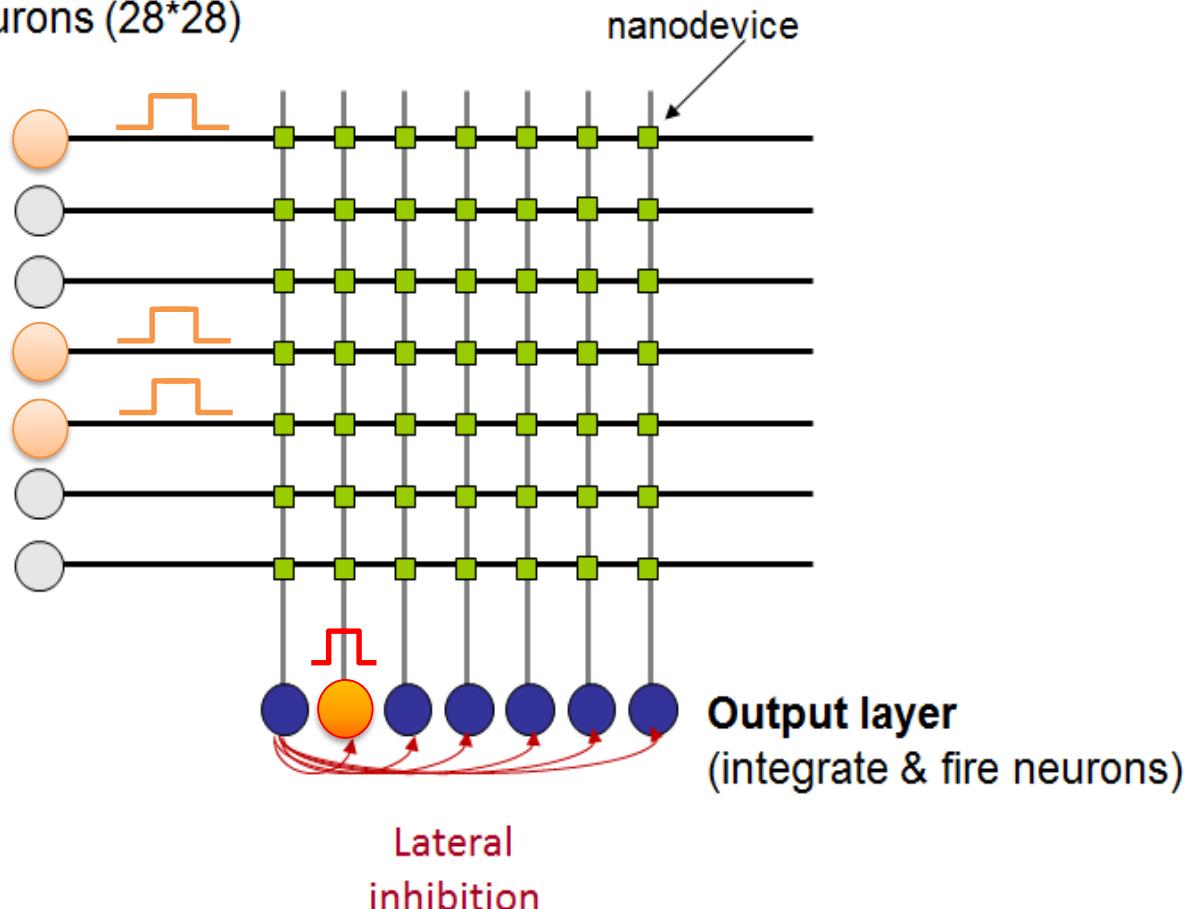


# How memristor learn without supervision through STDP

*Presented input*



**Input layer**  
784 neurons ( $28 \times 28$ )

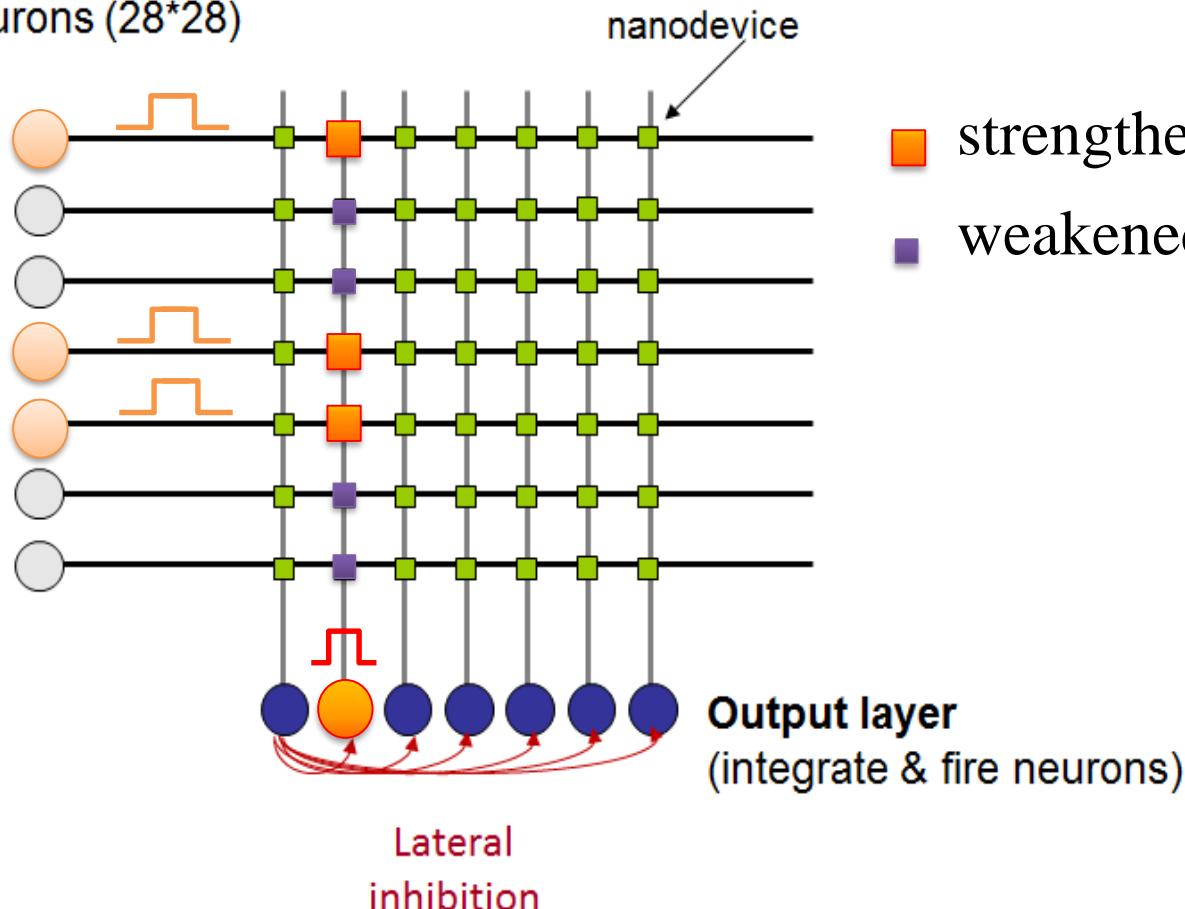


# How memristor learn without supervision through STDP

*Presented input*

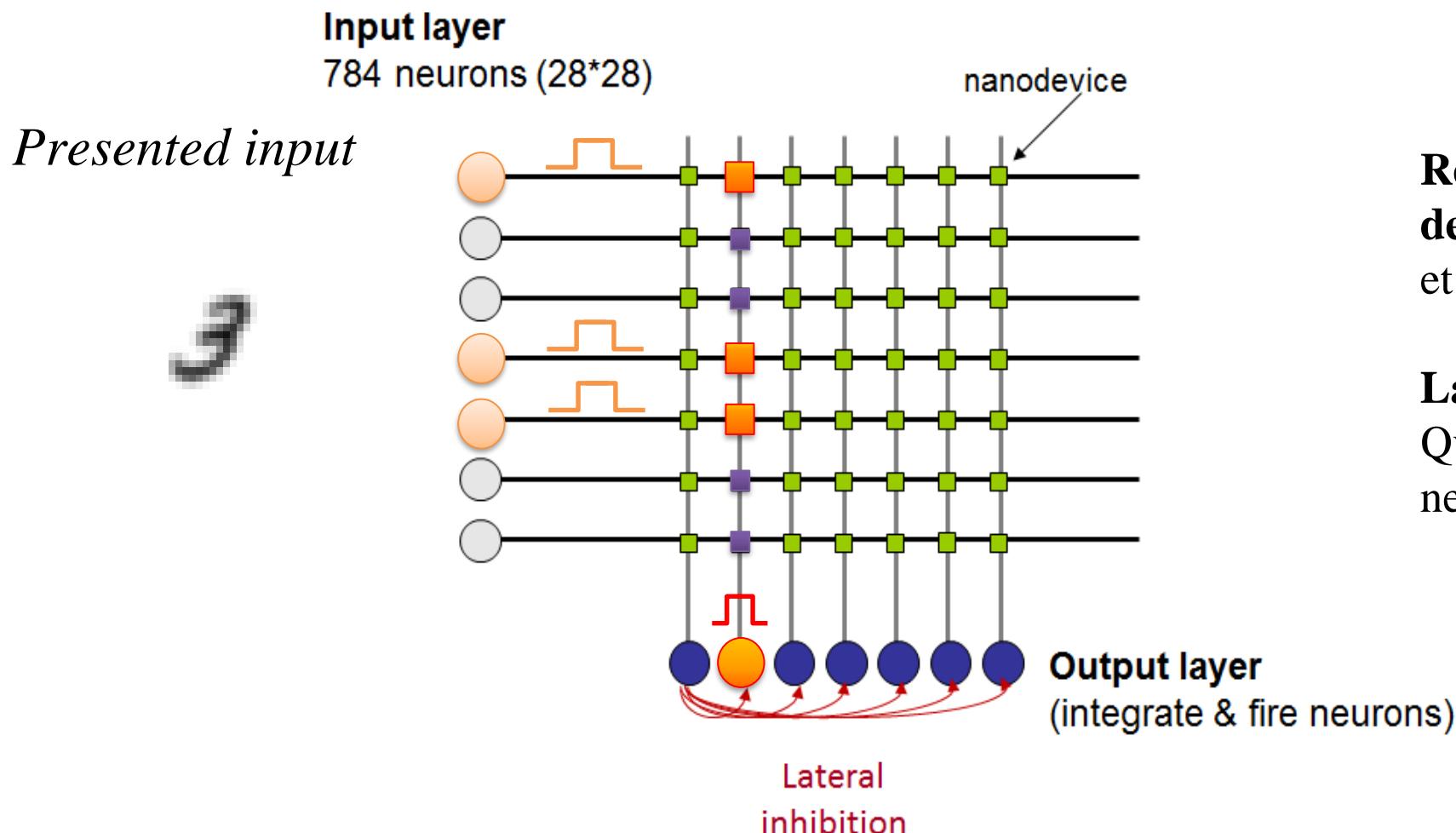


**Input layer**  
784 neurons ( $28 \times 28$ )



- strengthened synapse
- weakened synapse

# How memristor learn without supervision through STDP



**Recent small scale demonstration:** Pedretti, Ielmini et al, Sci Rep 2017

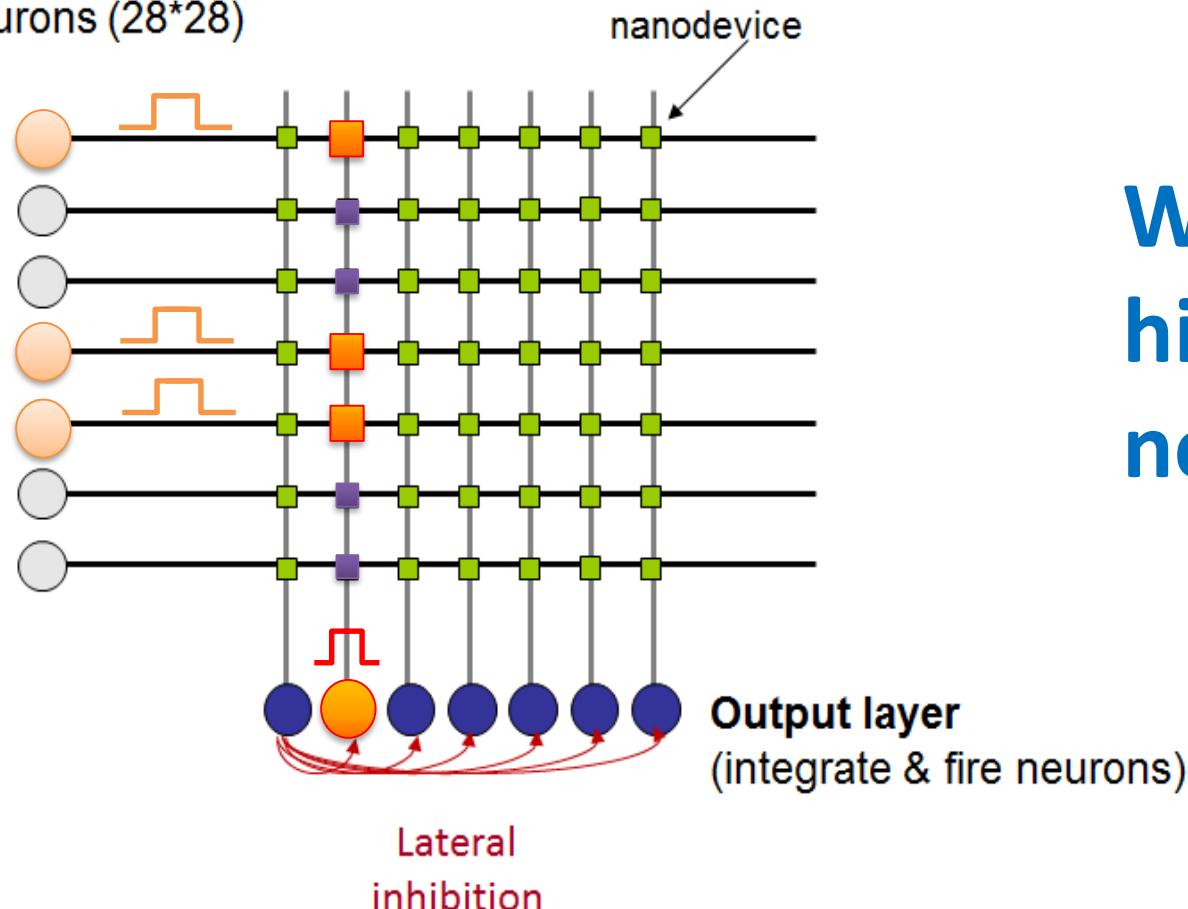
**Large scale simulations:**  
Querlioz, Bichler et al, Neural networks 2012

# How memristor learn without supervision through STDP

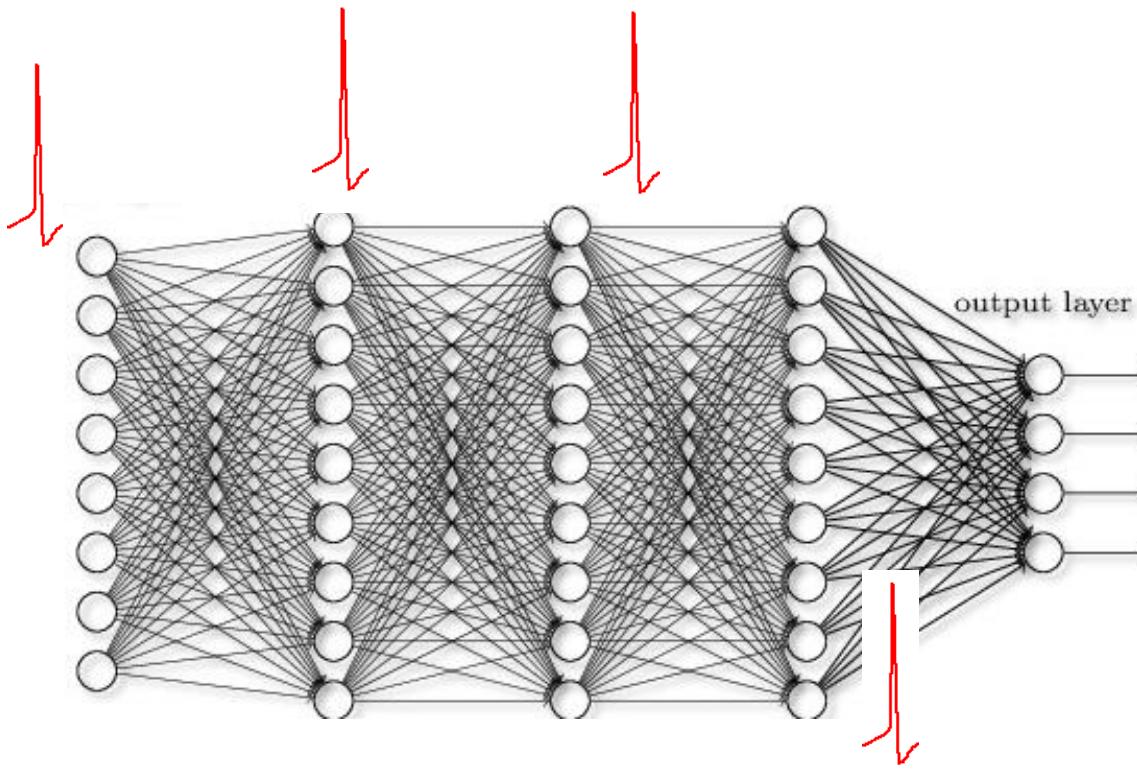
*Presented input*



**Input layer**  
784 neurons ( $28 \times 28$ )



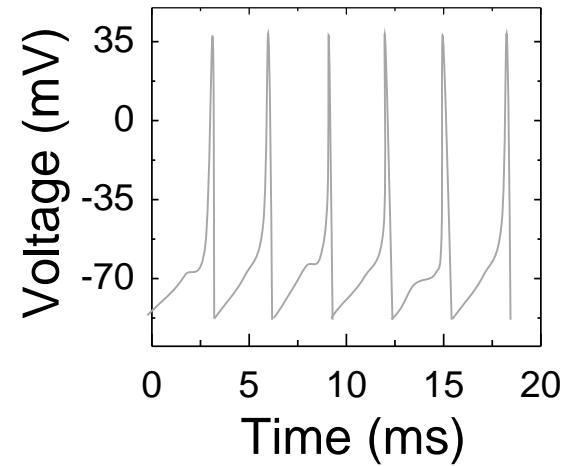
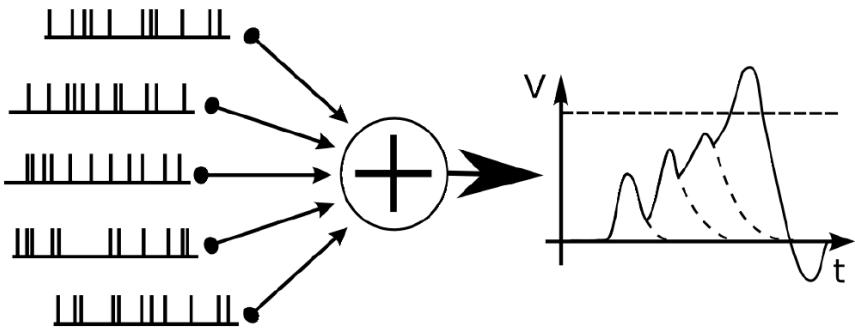
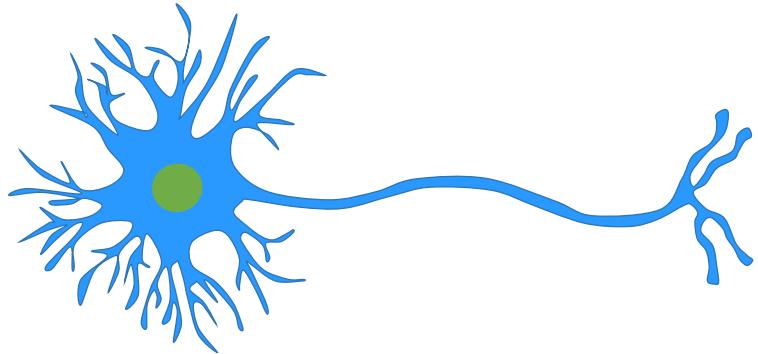
**Works even with  
highly non-linear,  
noisy devices !**



**Can non-linear dynamics help ?**

# Biological neurons are auto-oscillators: can we enhance performances of neural networks by giving them dynamical properties ?

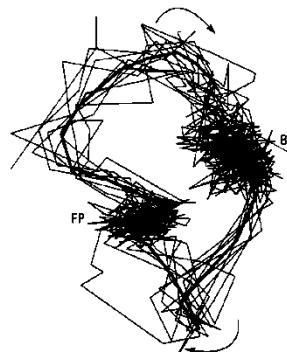
*Spikes, Non-linear oscillations, Synchronization*



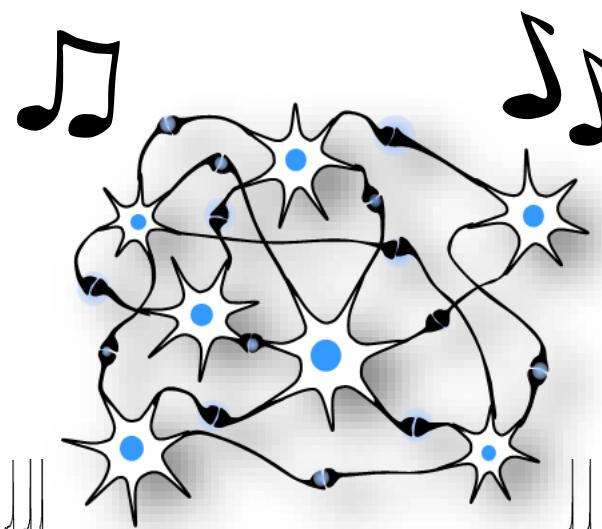
Can we reduce the number of nanodevices by giving them dynamical features ?

# Non-linear dynamics in the brain is probably useful for computing

## Complex transients

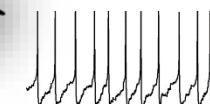
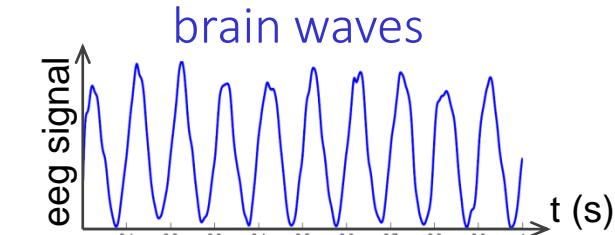


neuron: oscillator



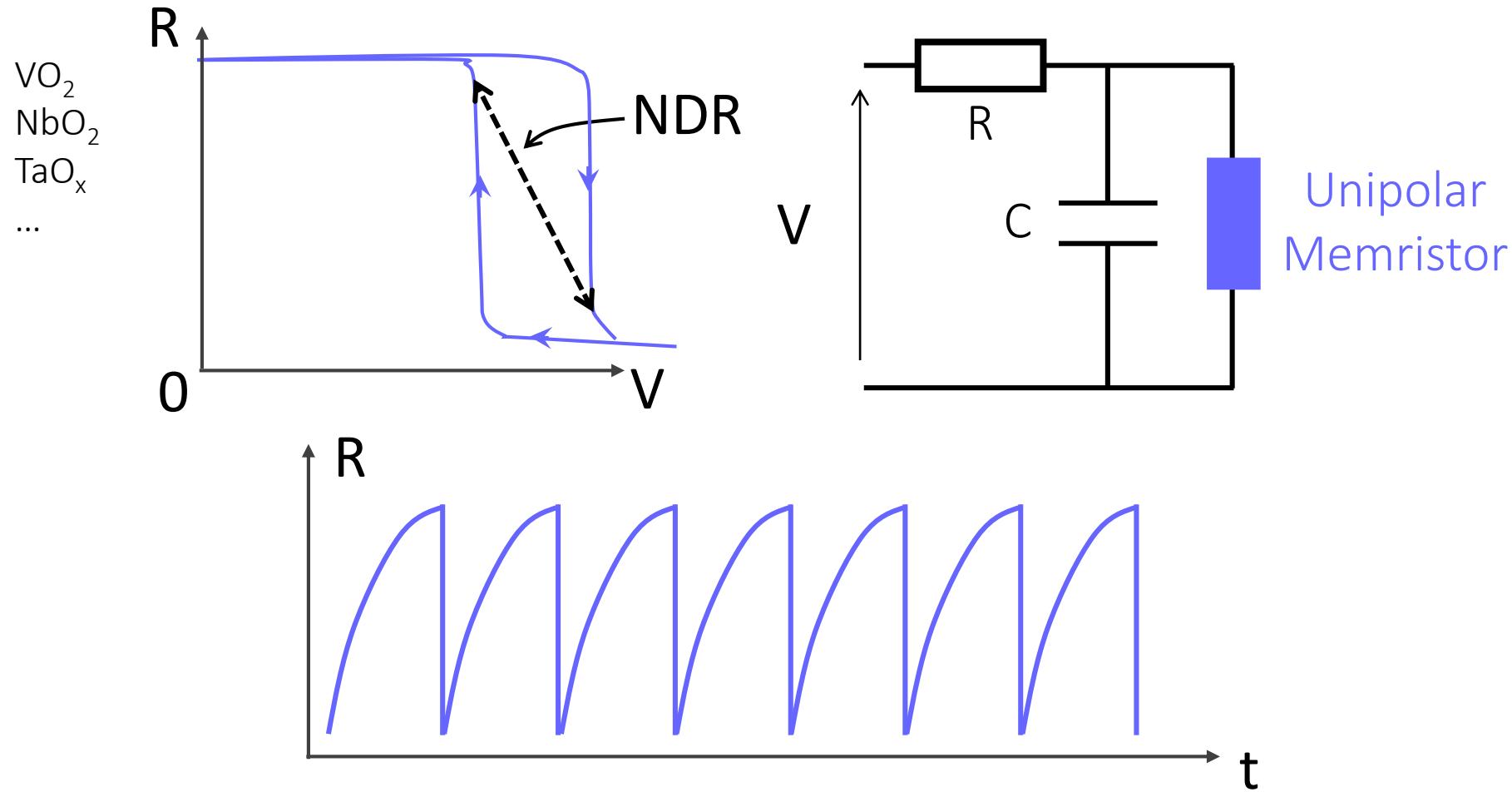
synapse:  
coupling

## Synchronization



neuron: oscillator

# There is currently a lot of work towards creating nanoscale dynamical neurons



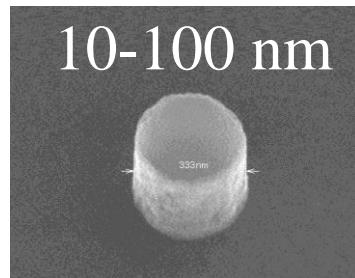
Pickett et al, Nat. Mater. (2013), N. Shukla et al, Sci. Rep. (2014)

A. Sharma, et al, IEEE J. Explor. Solid-State Comput. Devices Circuits (2015)

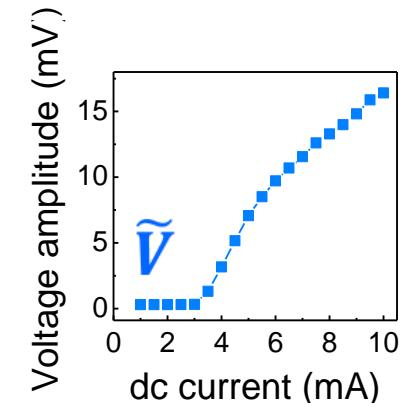
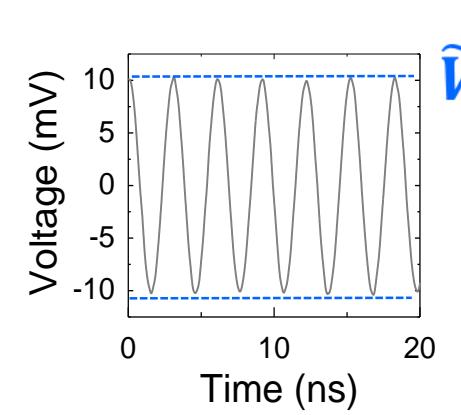
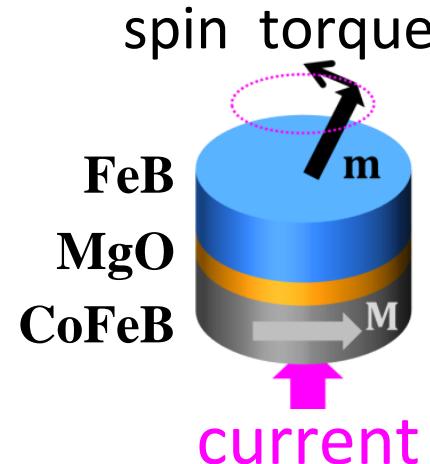
# Magnetic nano-oscillators are non-linear nano-radios

Nanoscale, fast (GHz), non-linear and easily measurable

magnetic tunnel junction



compatible with CMOS

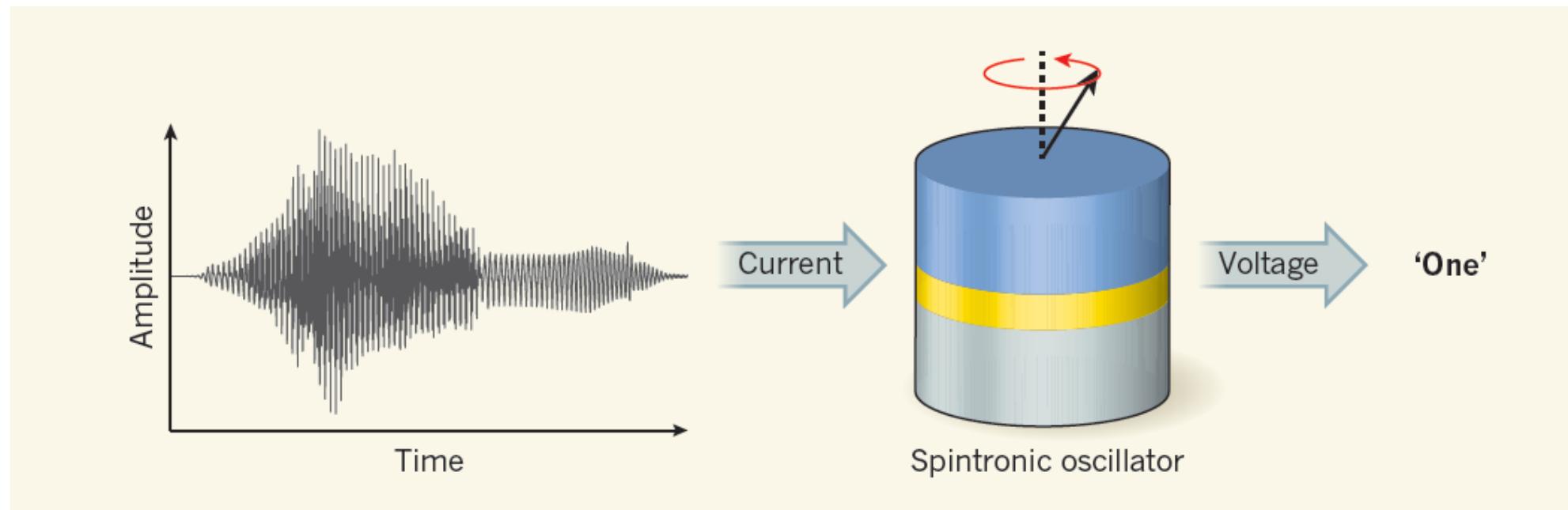


Same structure as magnetic memories

N. Locatelli, V. Cros and J. Grollier, Spin-torque building blocks, Nature Mat. 13, 11 (2014)

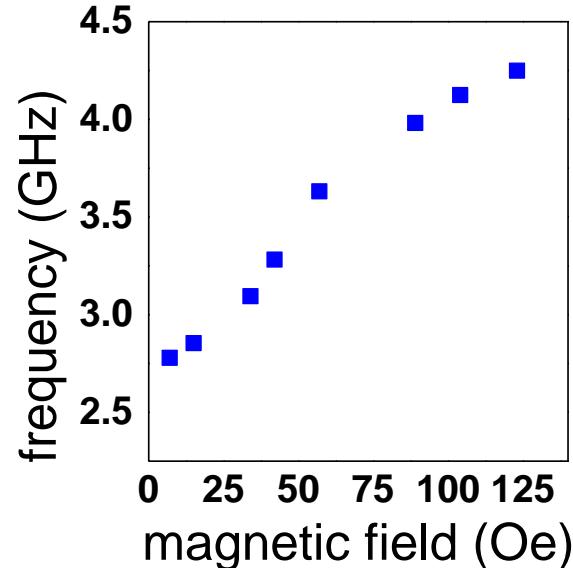
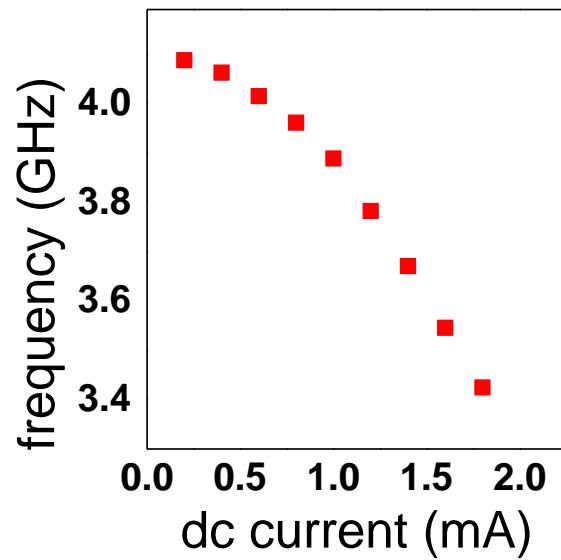
**Due to its stability and non-linearity, a single magnetic oscillator can emulate an assembly of neurons and perform neuromorphic computing**

Spoken digit recognition through reservoir computing



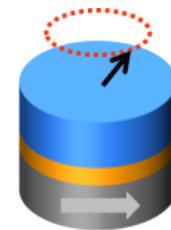
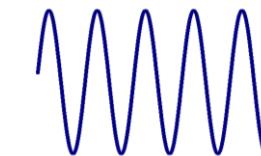
J. Torrejon, M. Riou, F. Abreu Araujo et al, Nature 547, 428 (2017)

# Spin-torque nano-oscillators have a high tunability : they are radio-receivers



Enhanced sync ranges

AC signals

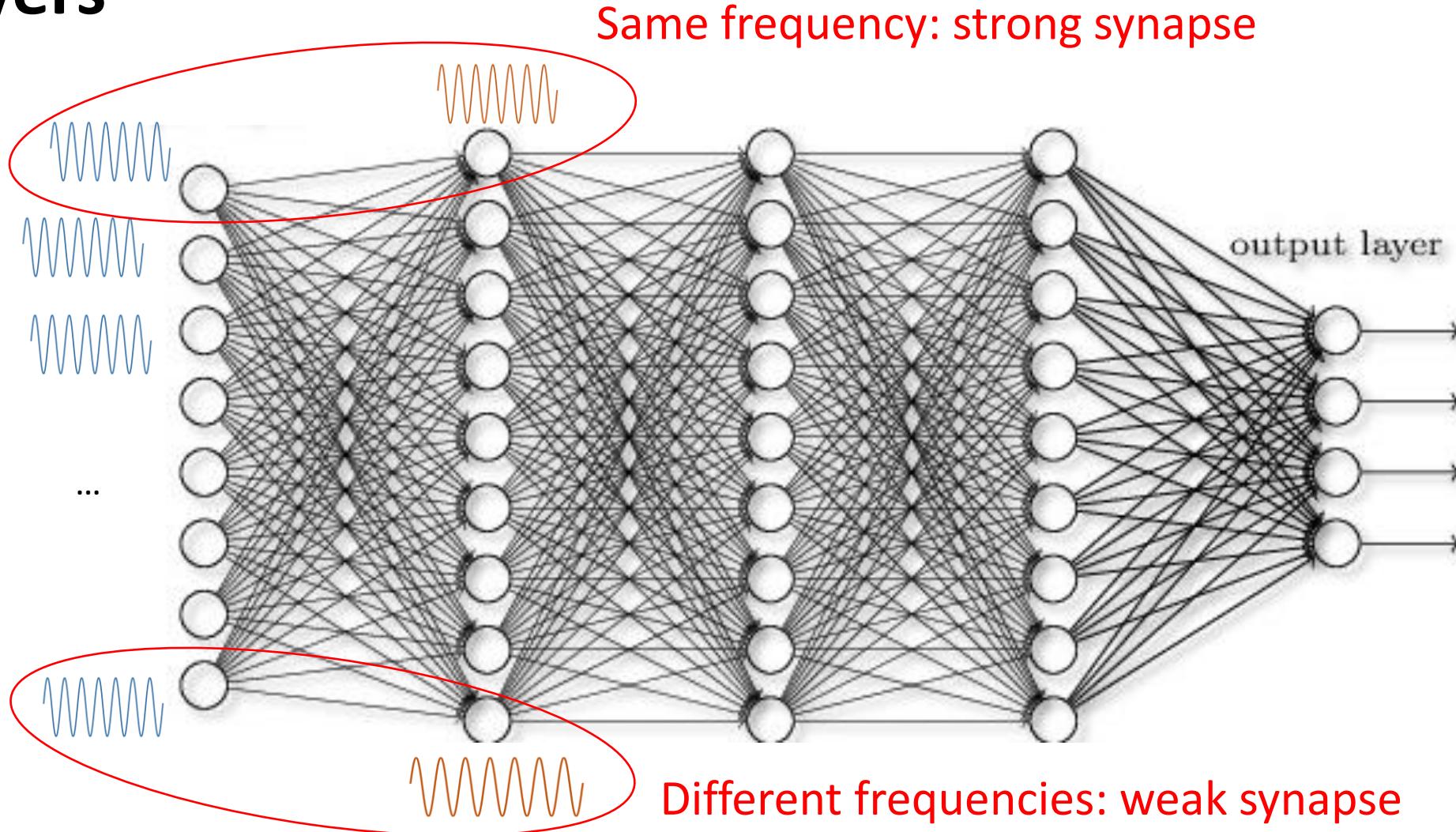


A. Slavin and V. Tiberkevich, IEEE TM 45, 1875 (2009)

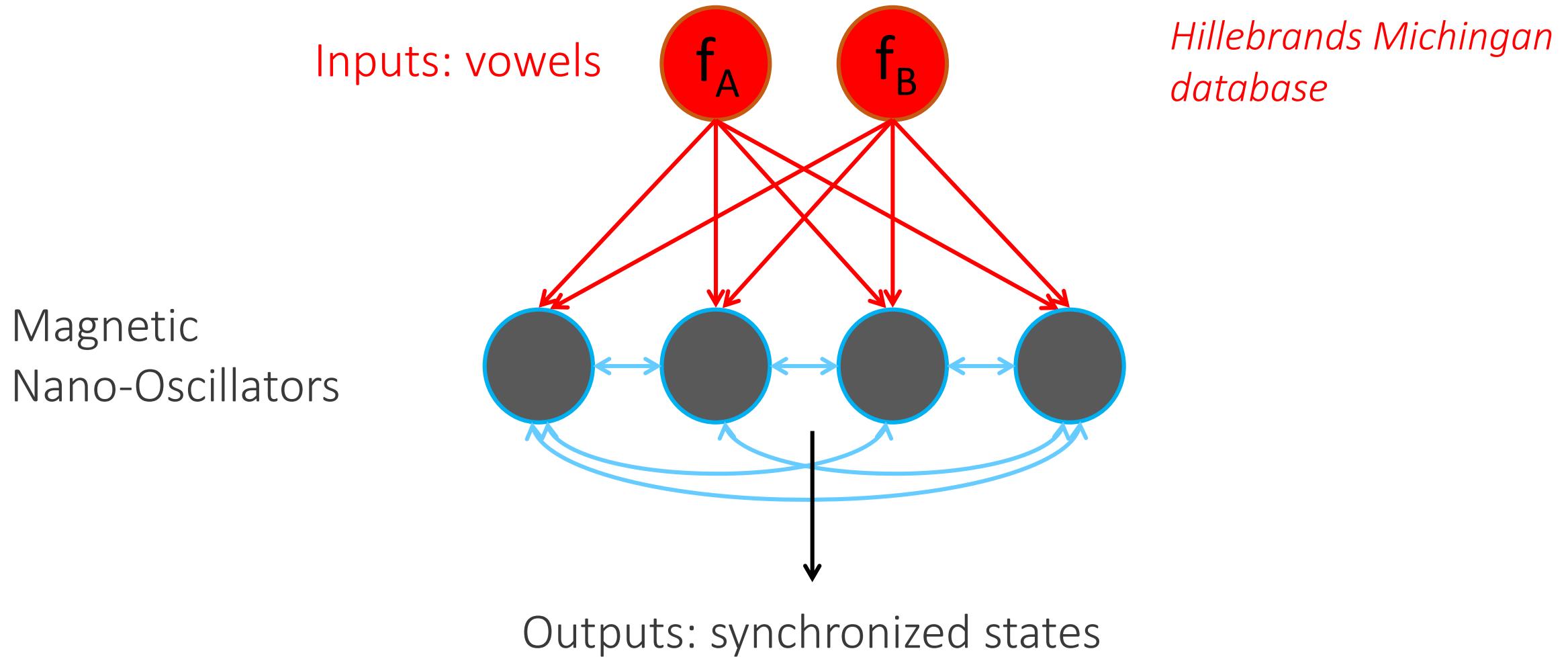
W. H. Rippard et al., PRL. 95, 067203 (2005), R. Lebrun, JG et al, PRL 115, 017201 (2015)

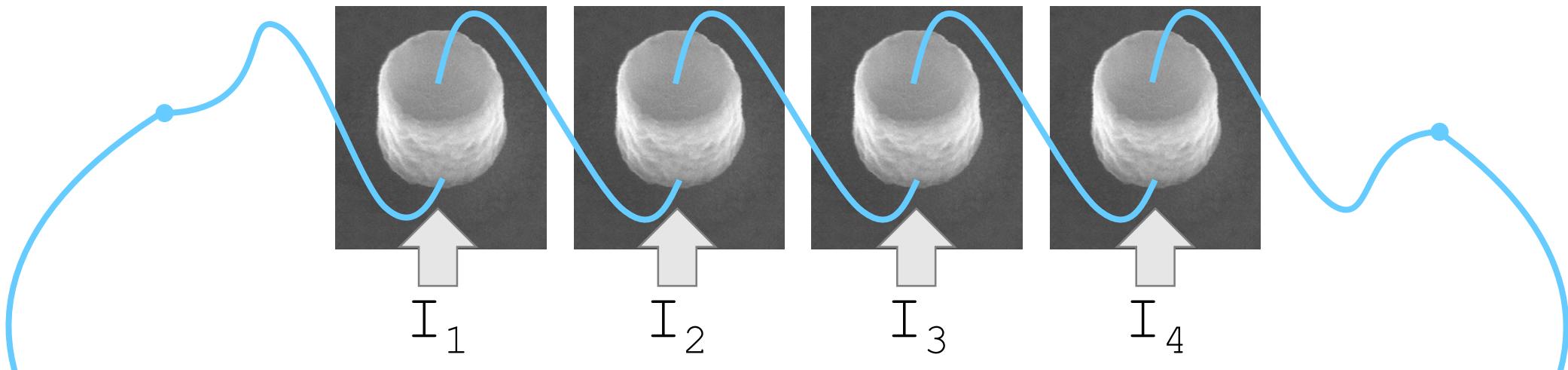
A. A. Awad et al, Nature Phys. (2016)

# The oscillators ability to mutually interact opens the path to RF on-chip communication between neuron layers

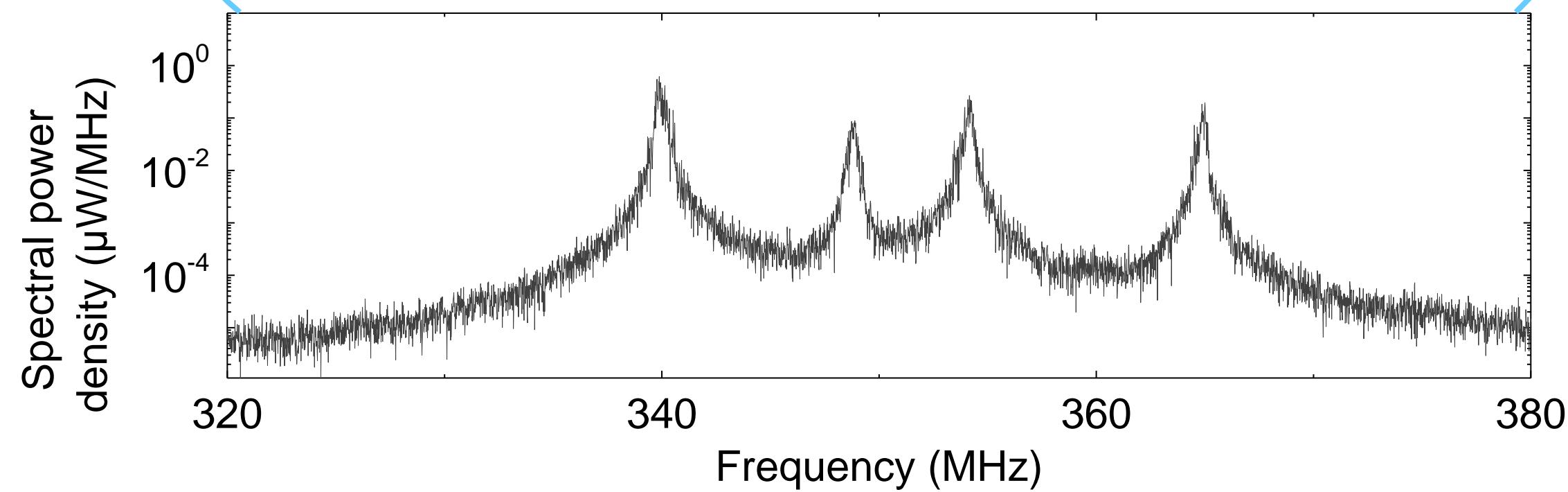


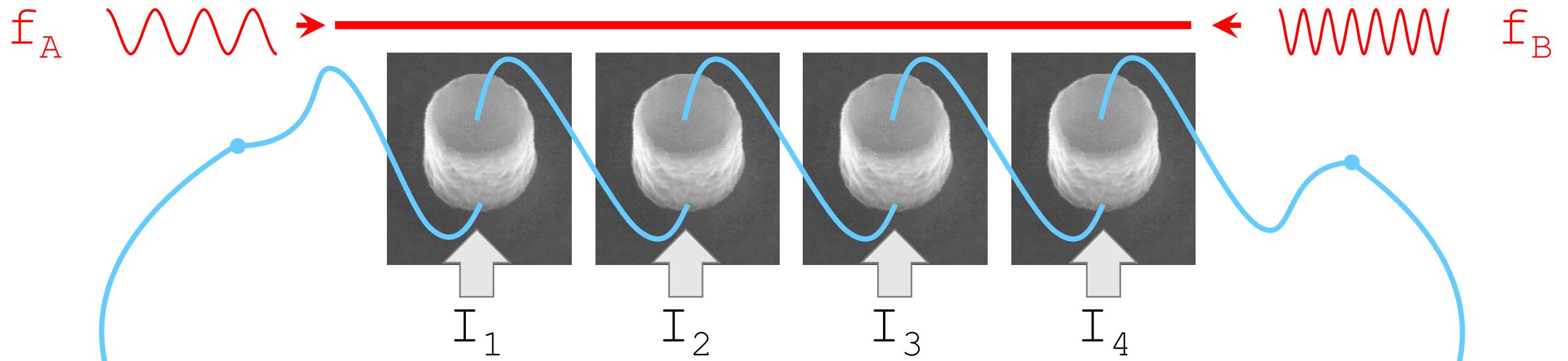
# Vowels classification with spin-torque oscillator neural network



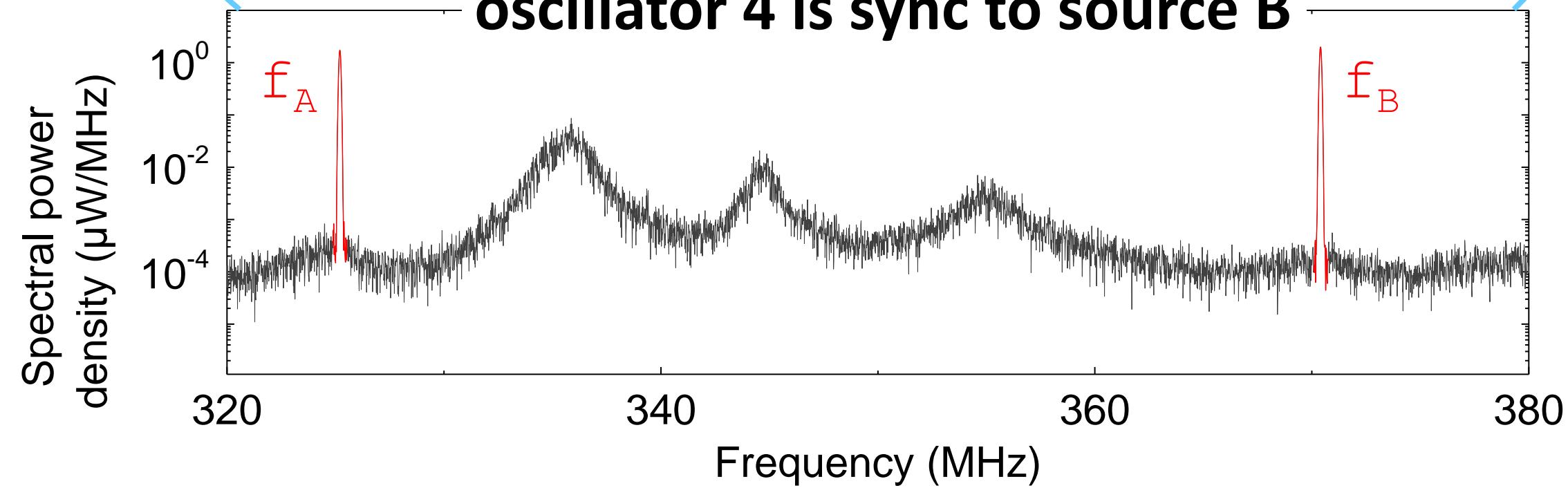


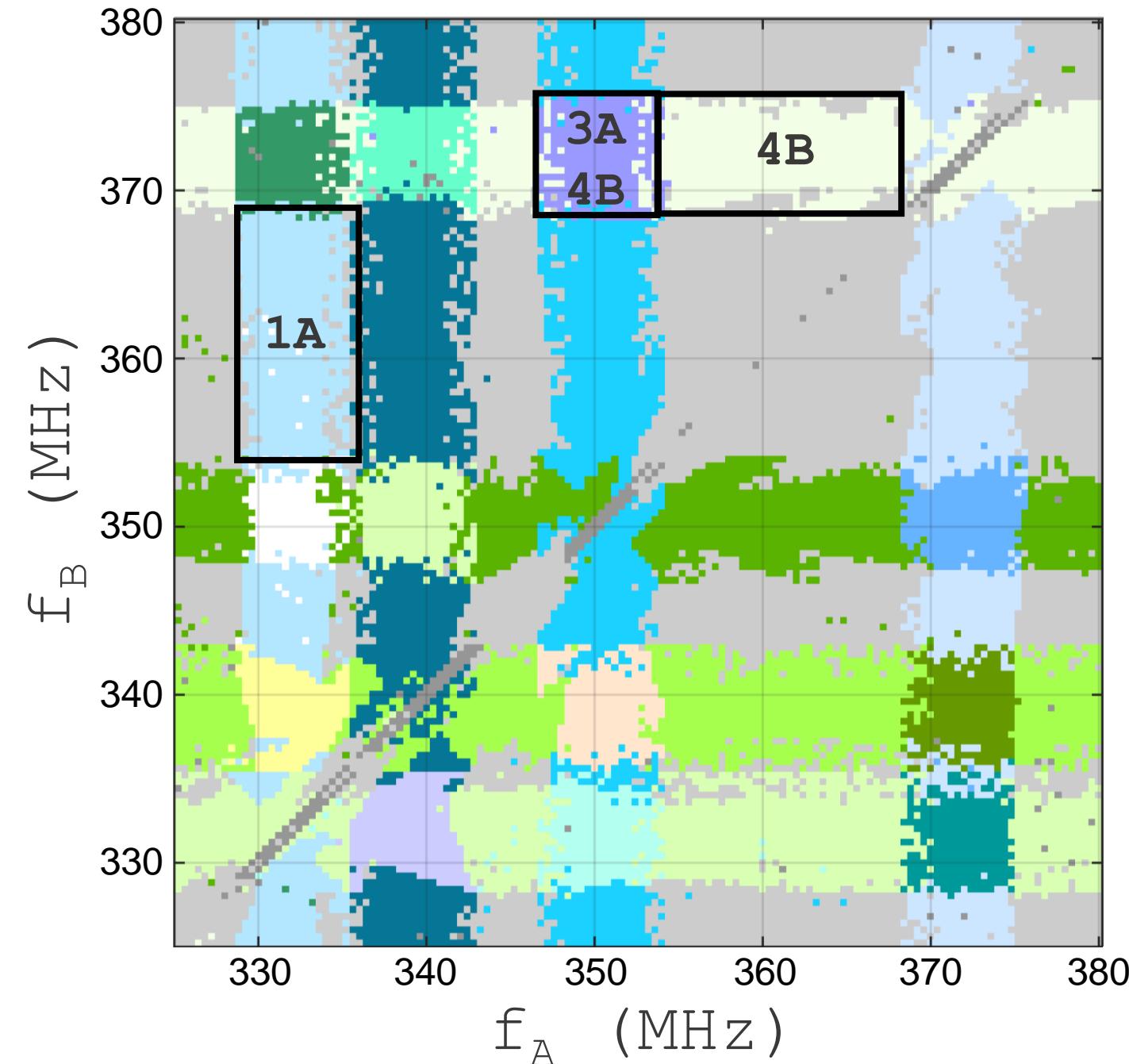
## Response of the neural network without inputs



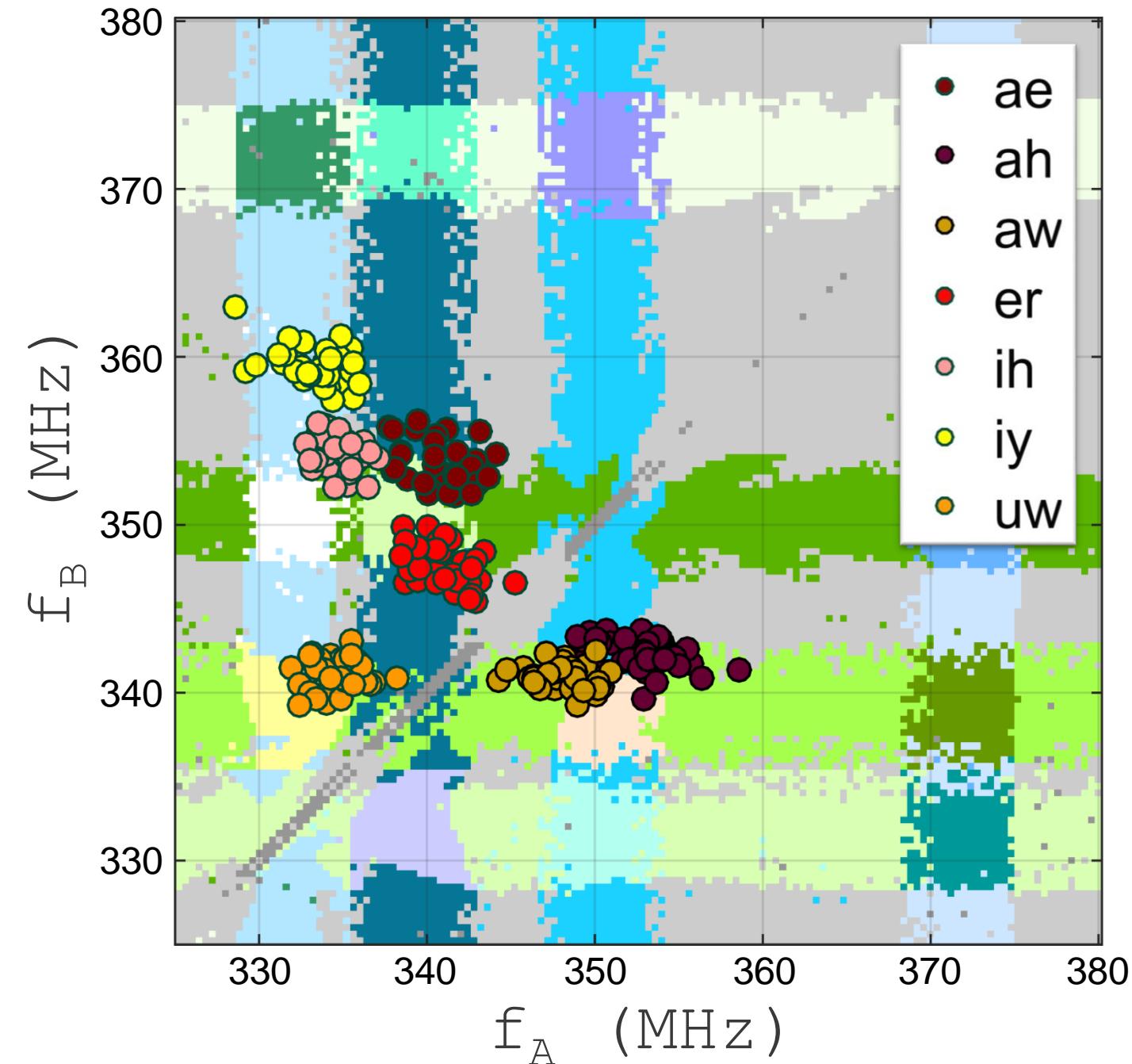


**The inputs modify the oscillator responses:  
oscillator 4 is sync to source B**



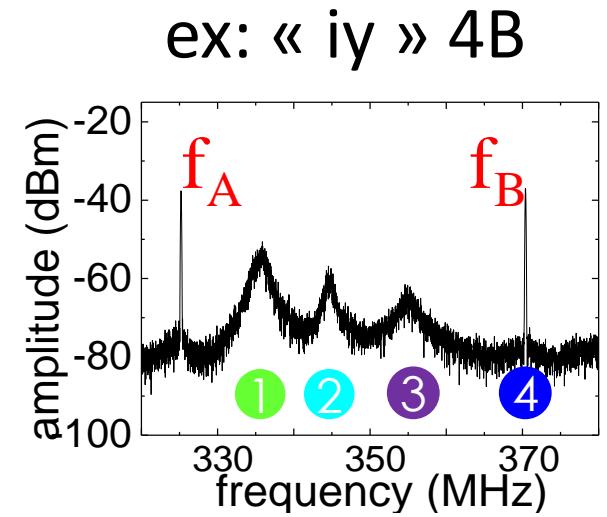
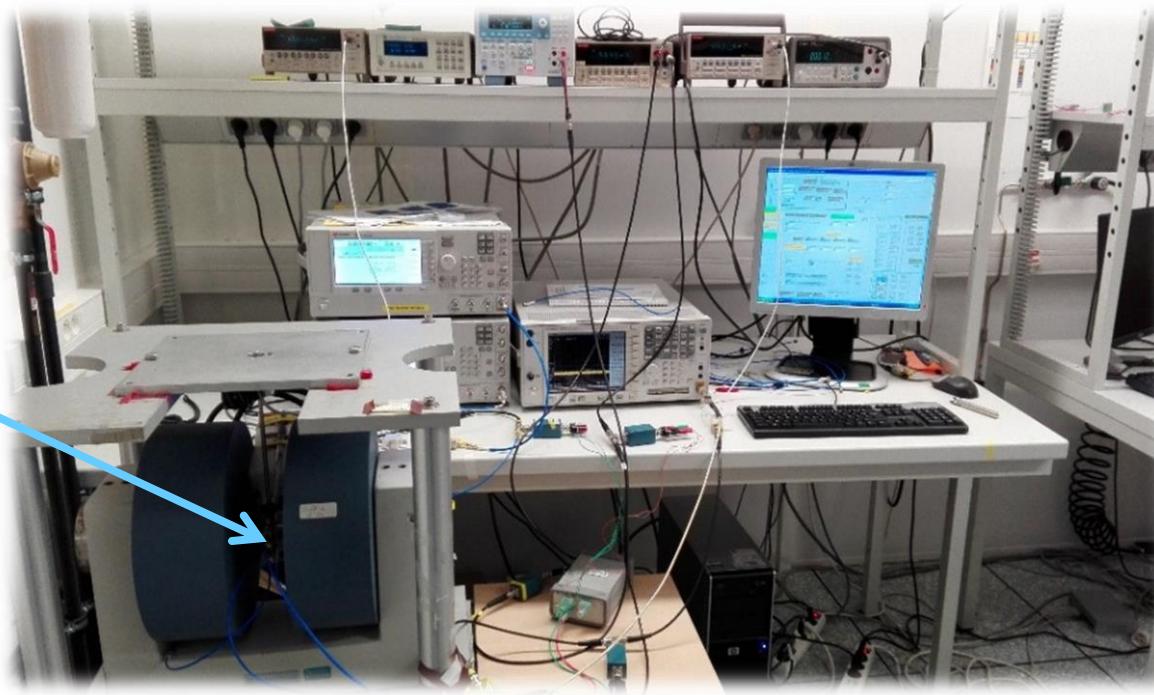
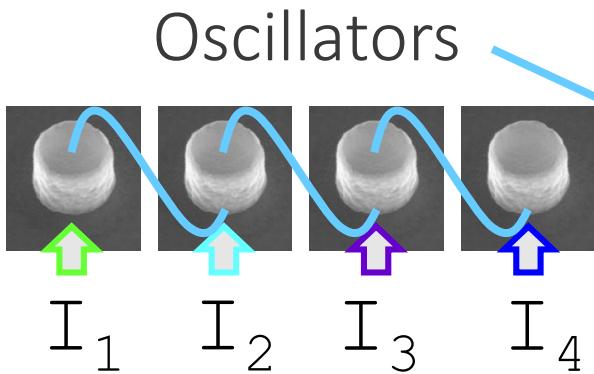


**We summarize all these measurements in a map where the different synchronization states have different colors**

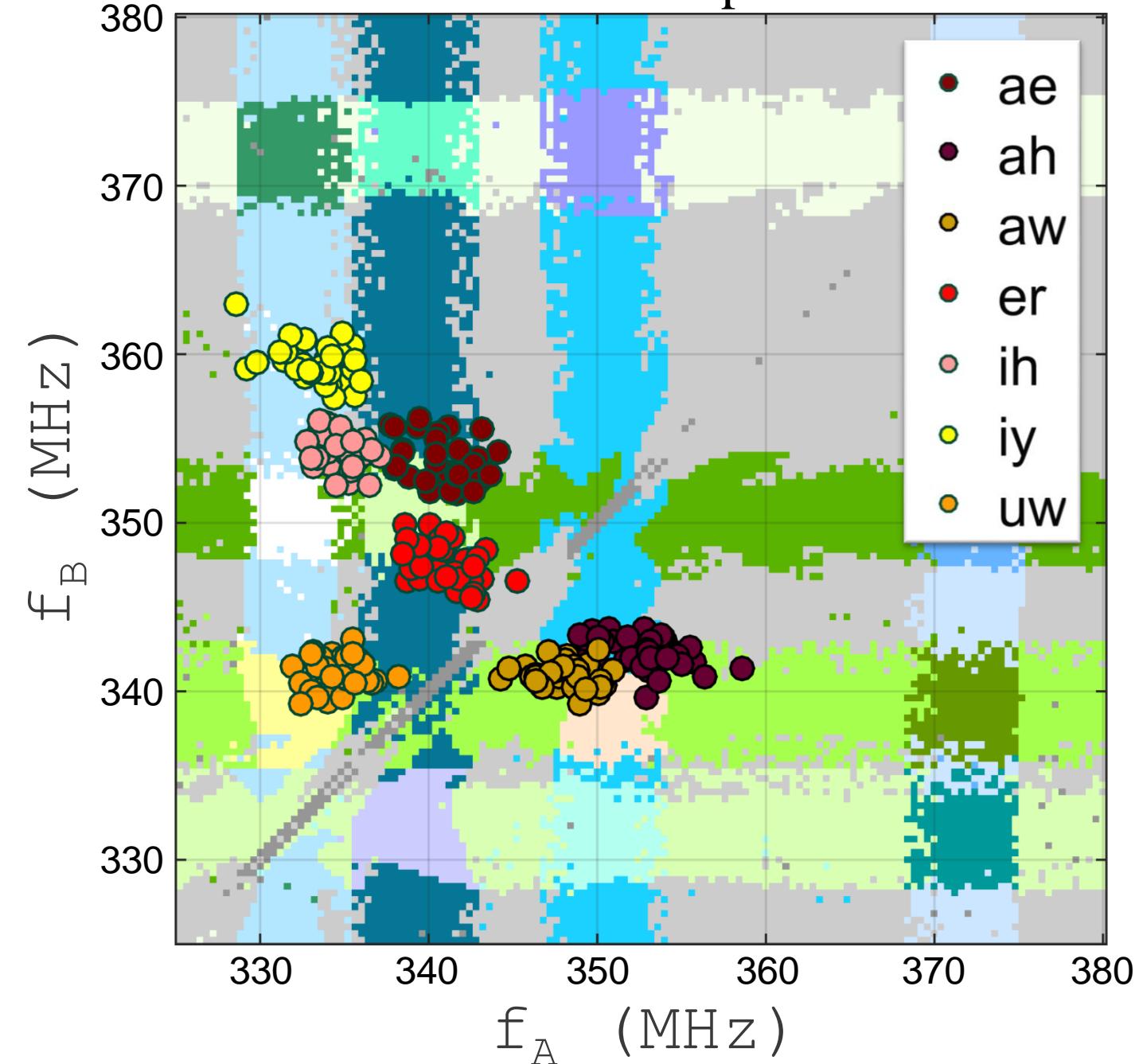


**For classification,  
all the points  
corresponding to one  
vowel should fall in a  
single synchronization  
region**

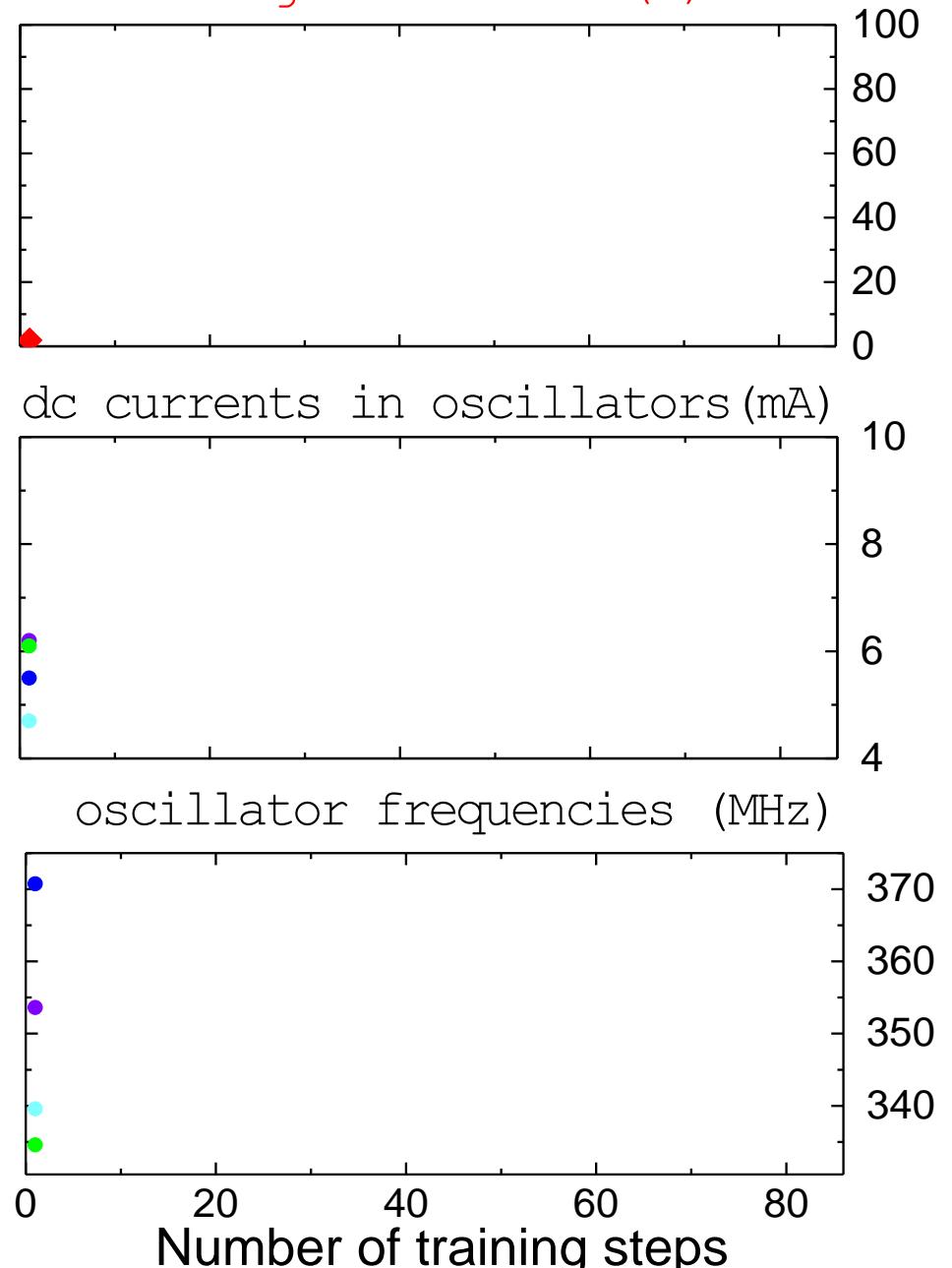
# We train the network by tuning the currents through the oscillators according to an online learning rule



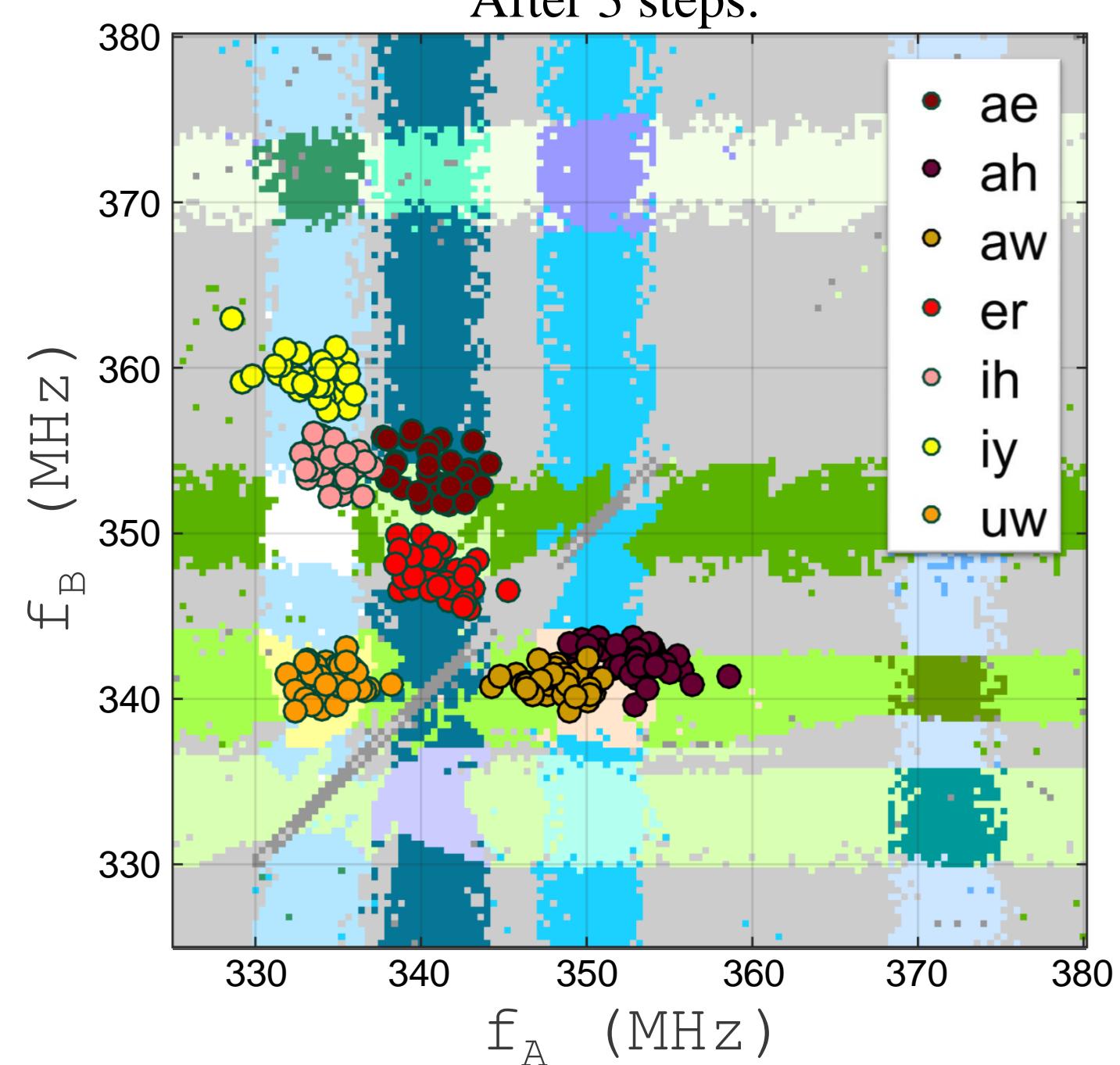
After 1 step:



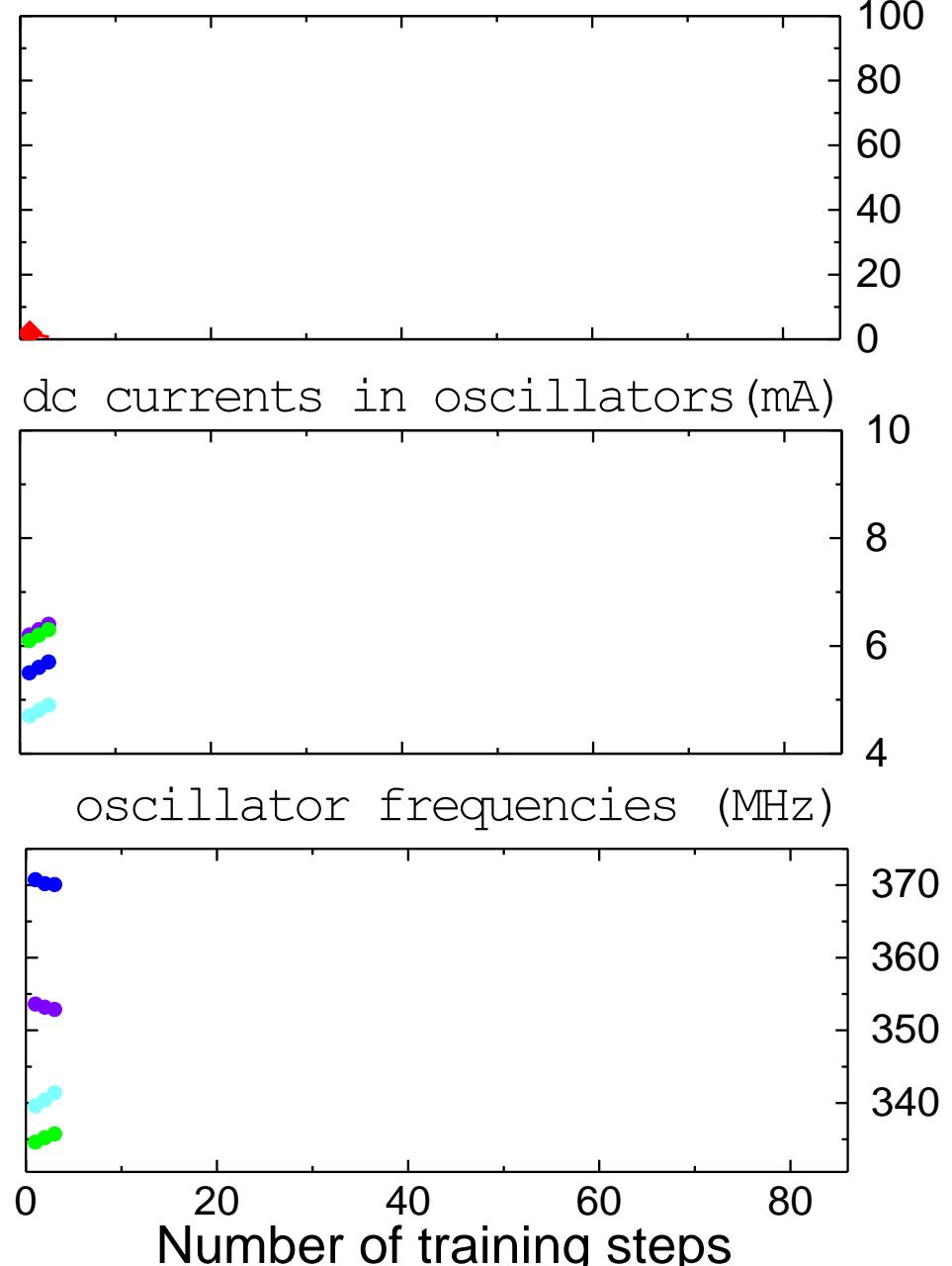
recognition rate (%)



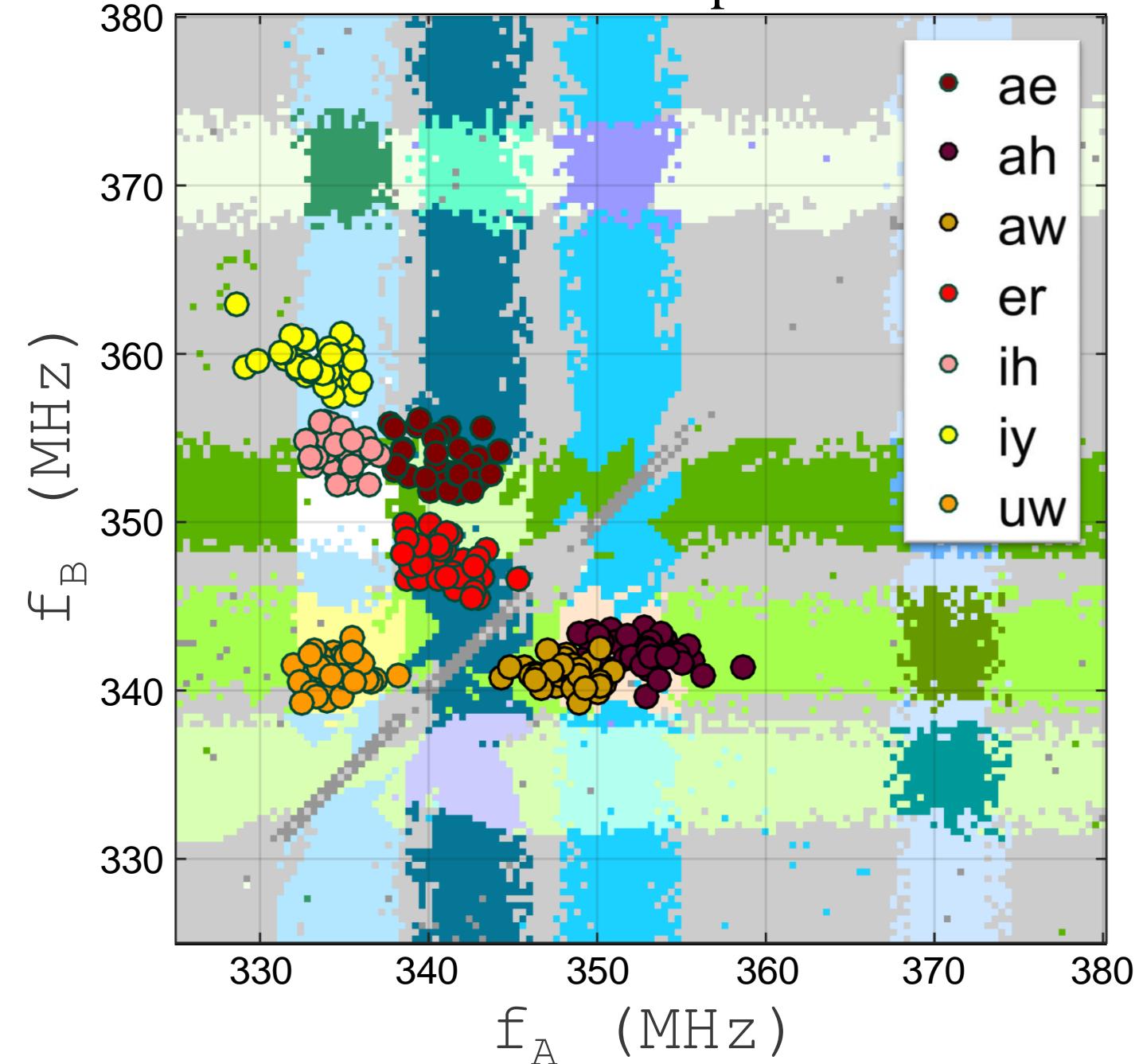
After 3 steps:



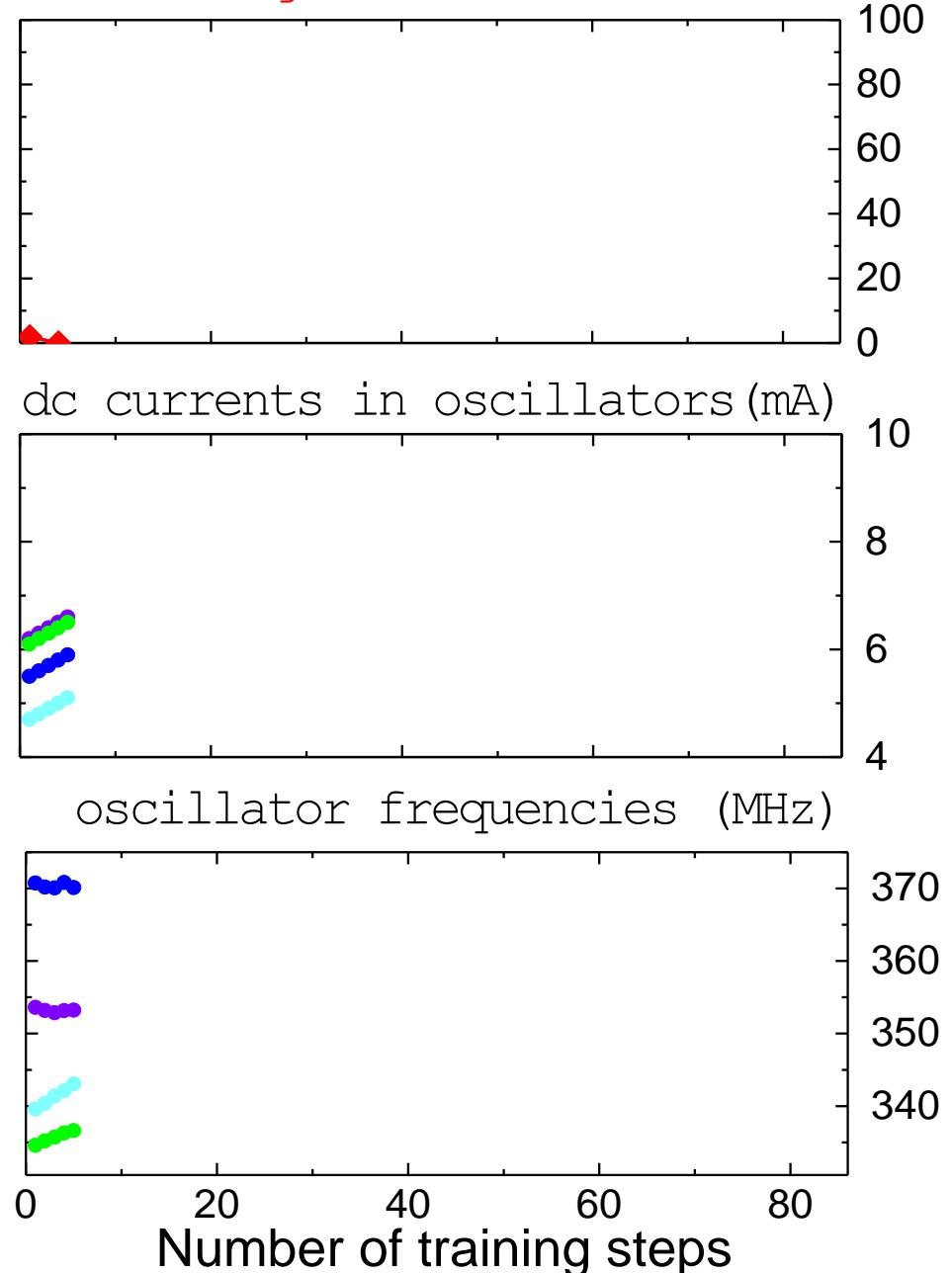
recognition rate (%)



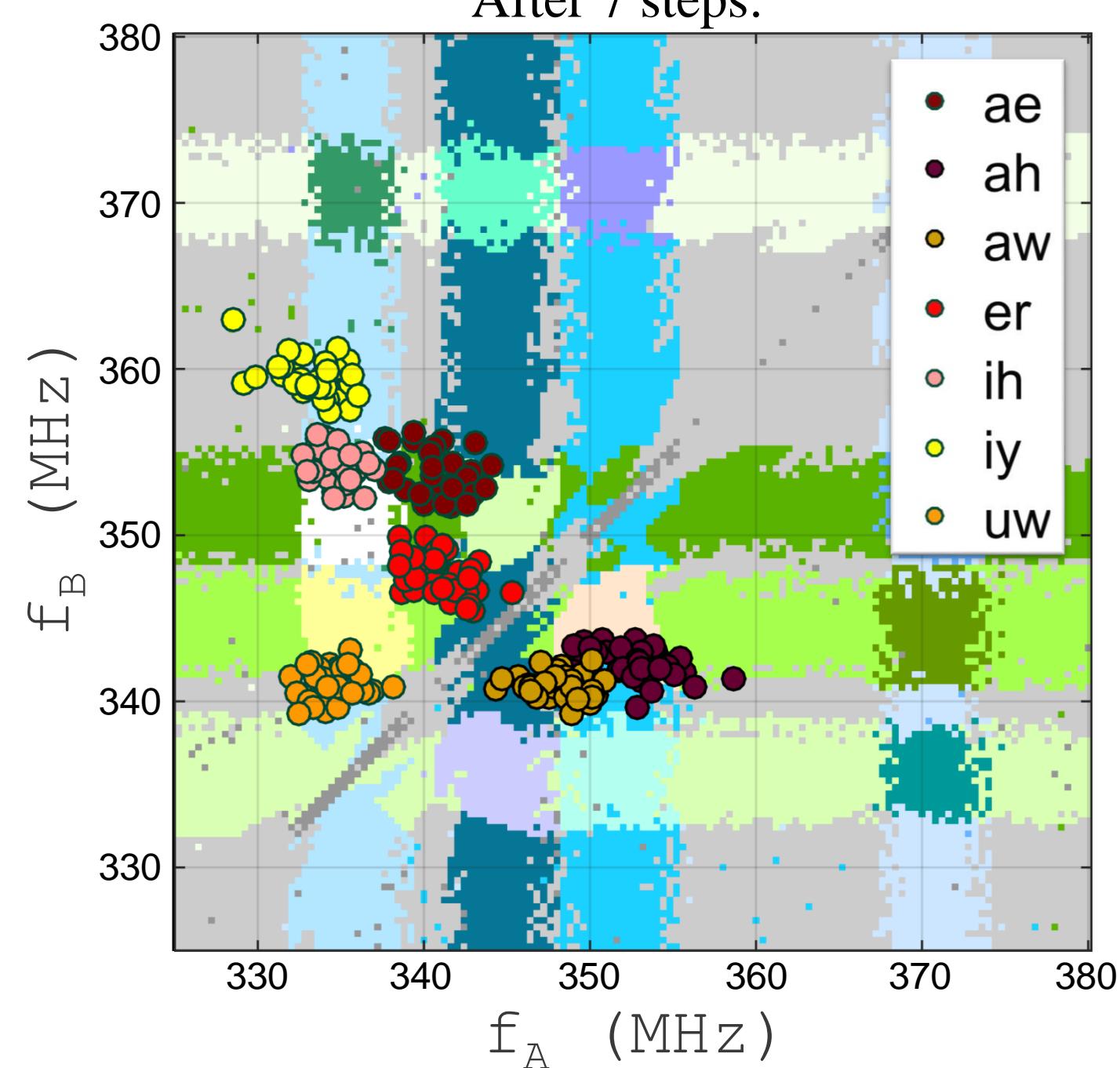
After 5 steps:



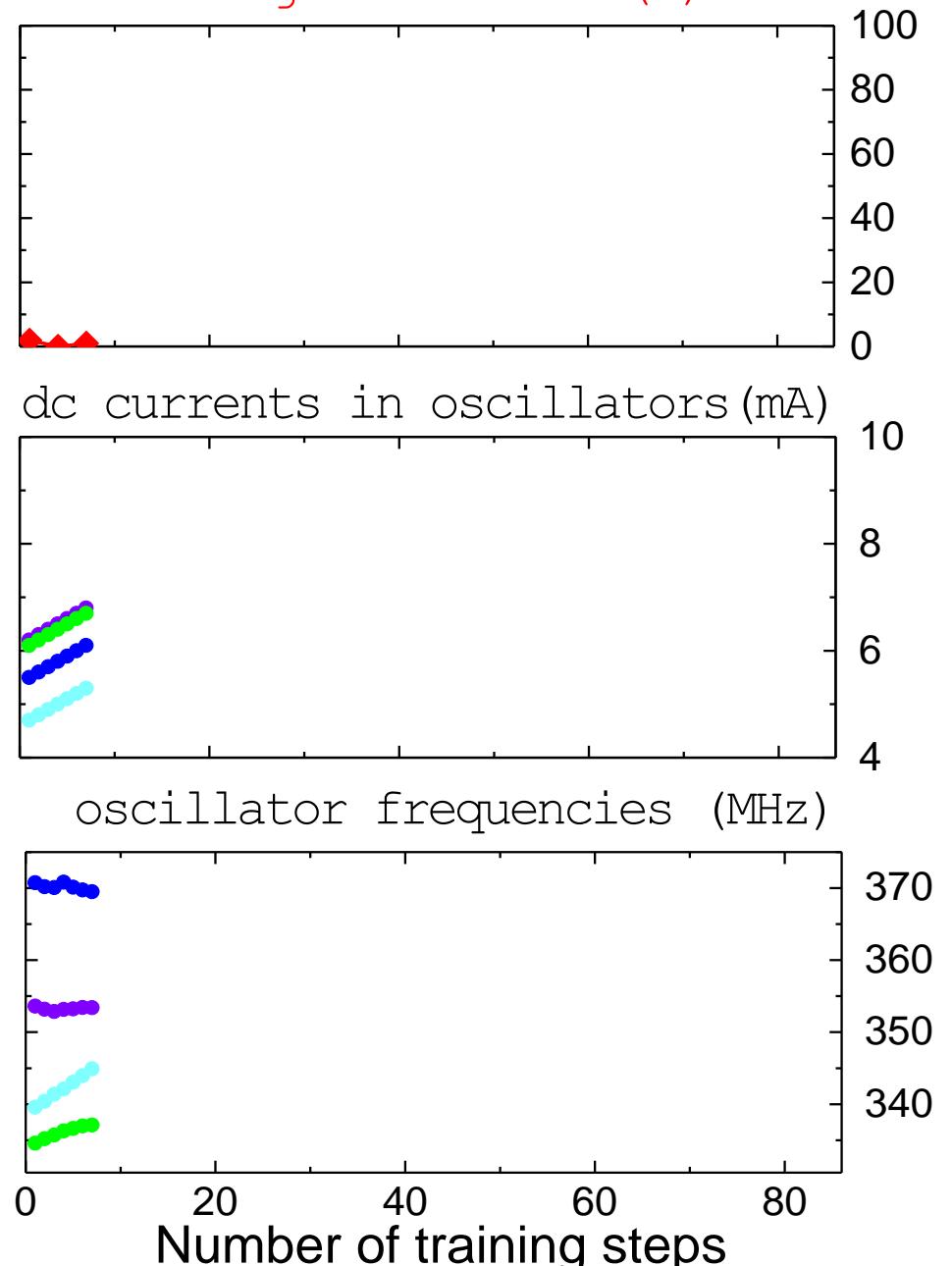
recognition rate (%)



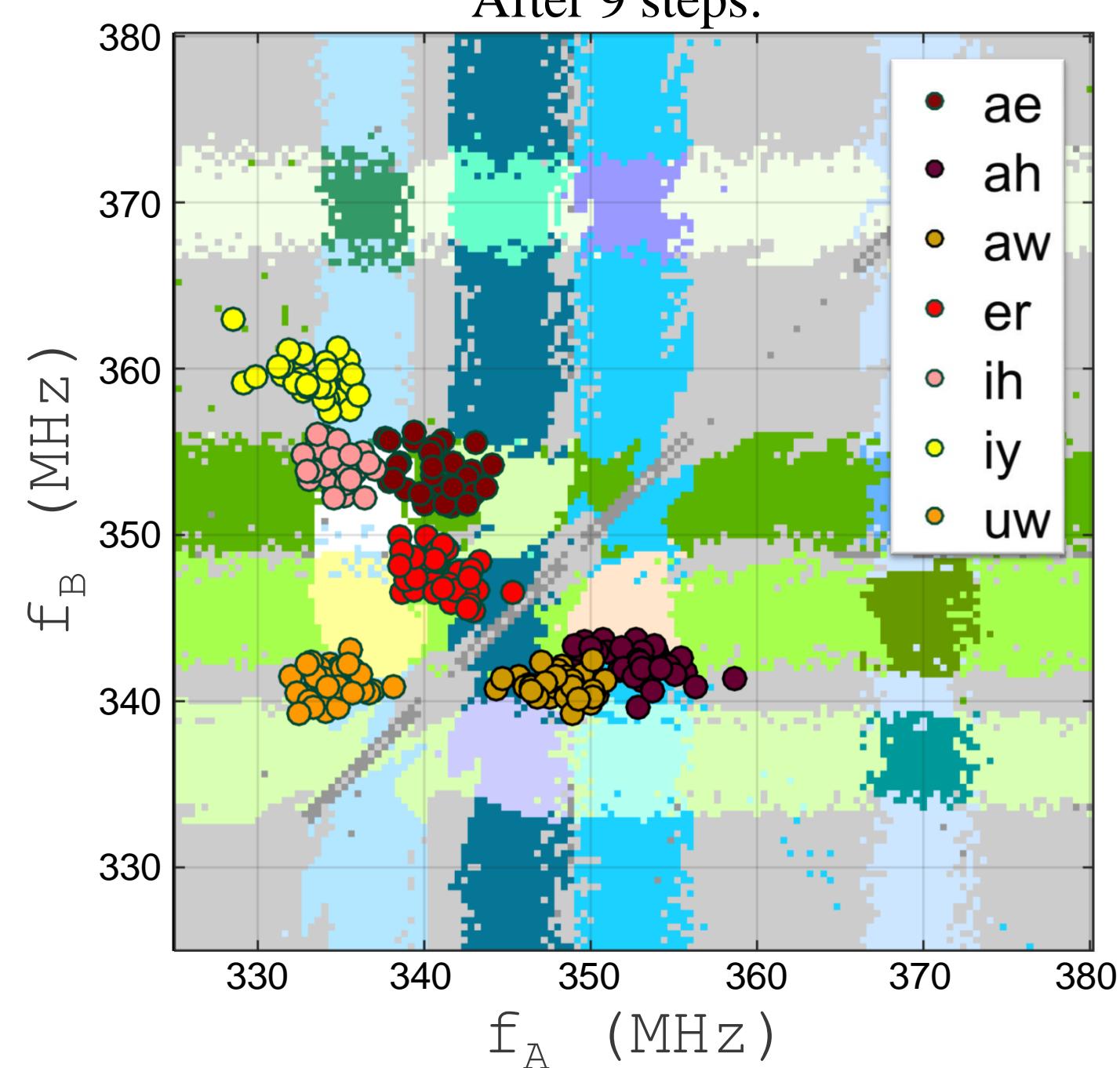
After 7 steps:



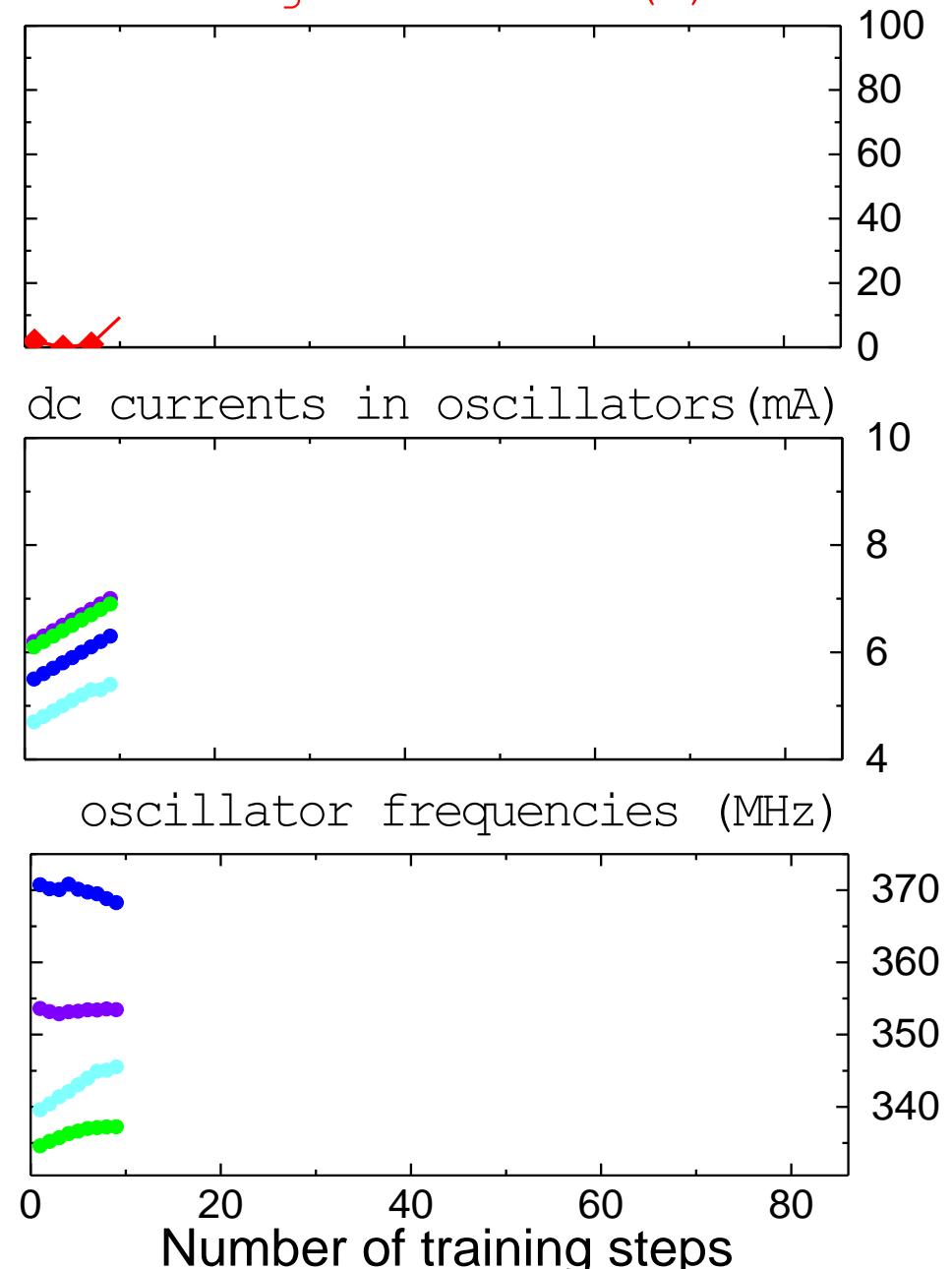
recognition rate (%)



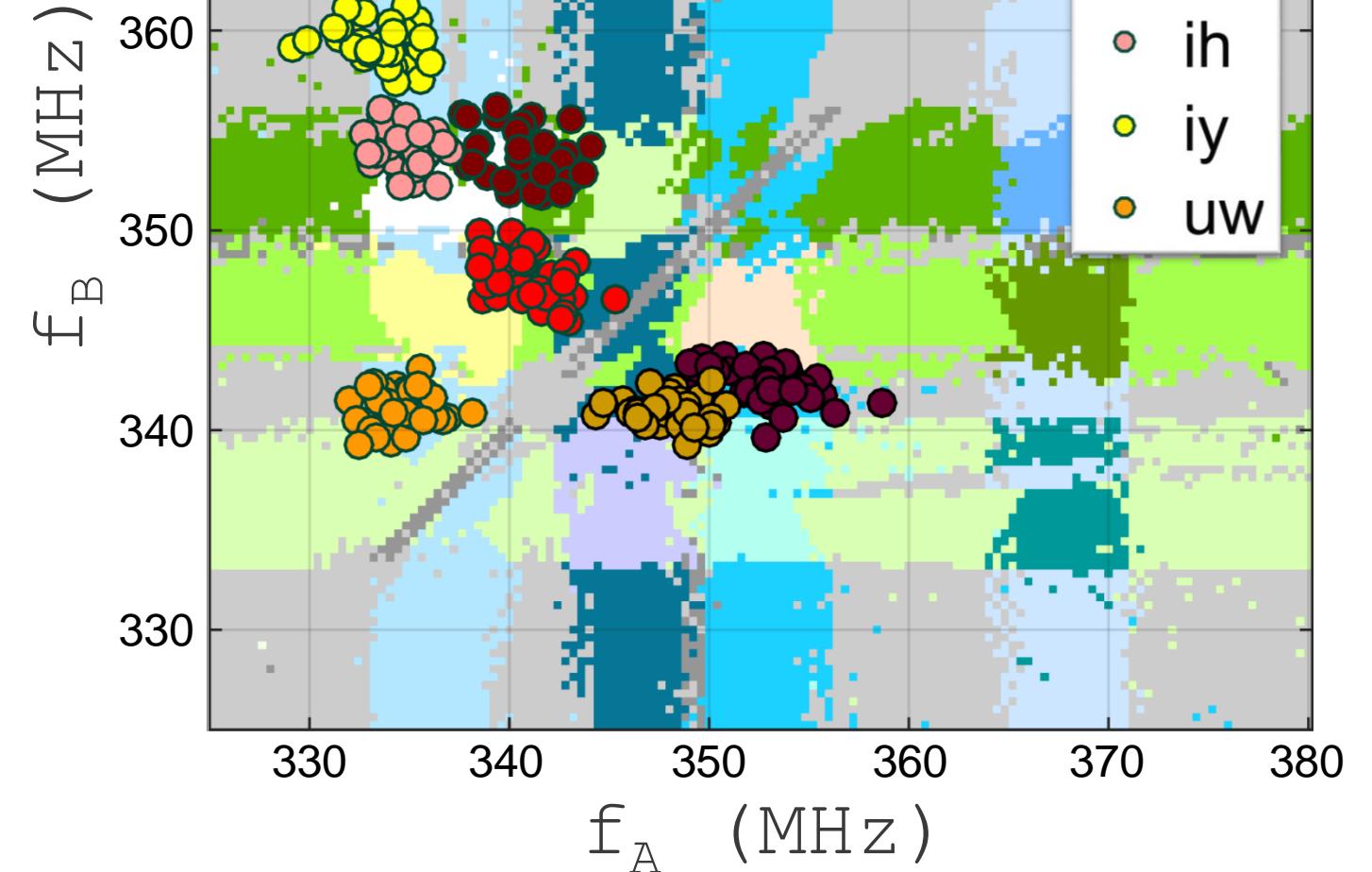
After 9 steps:



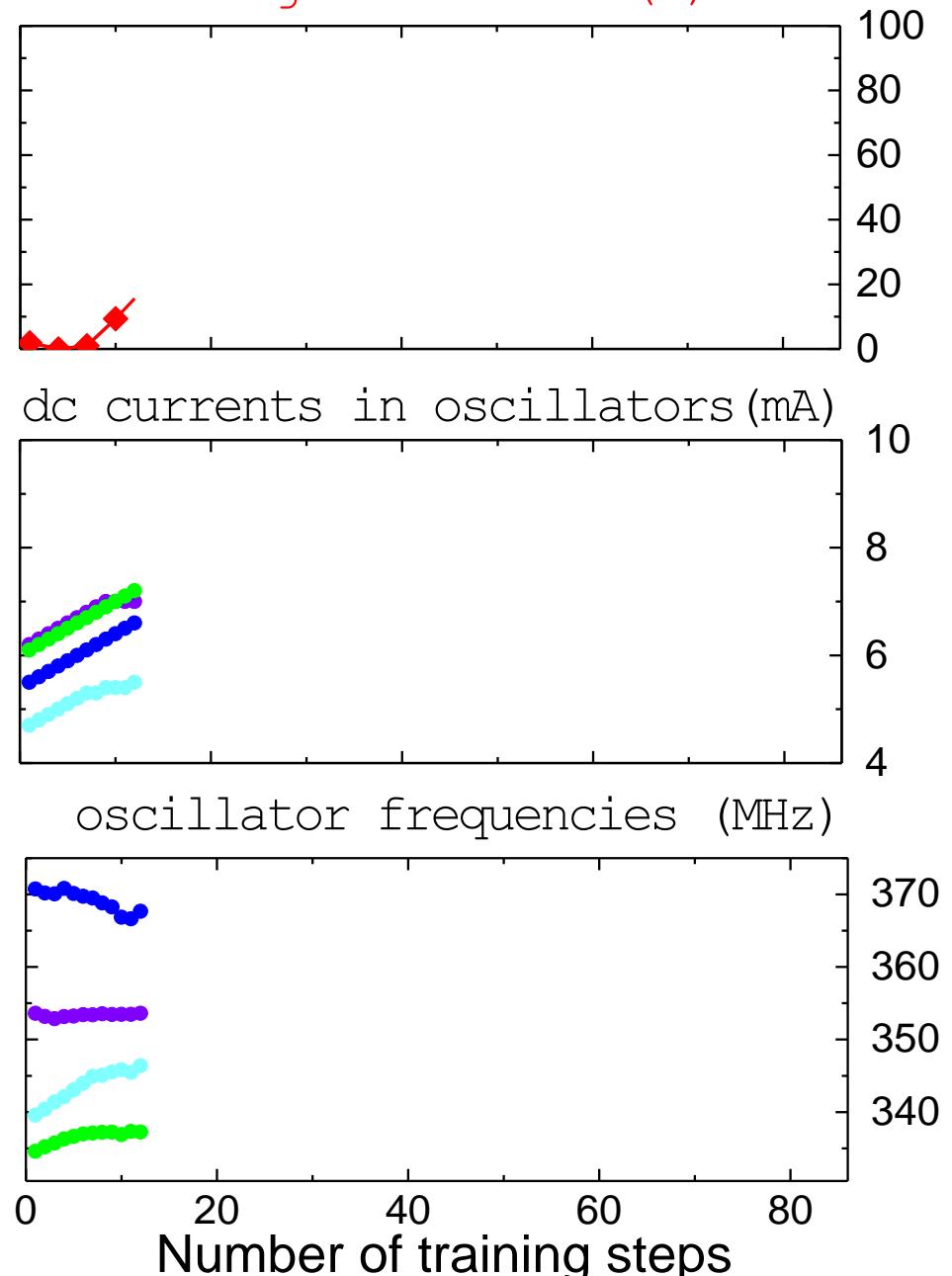
recognition rate (%)



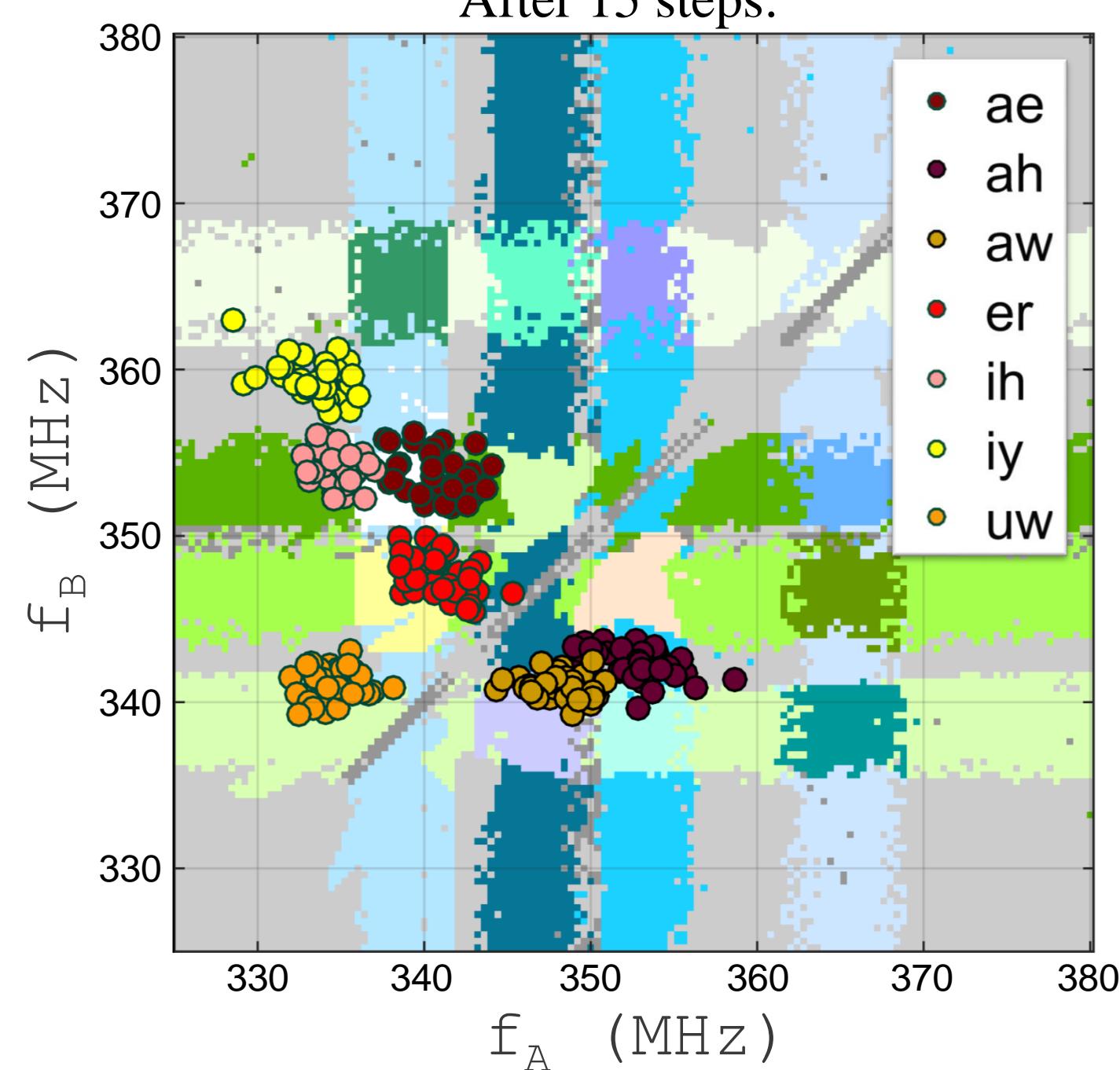
After 12 steps:



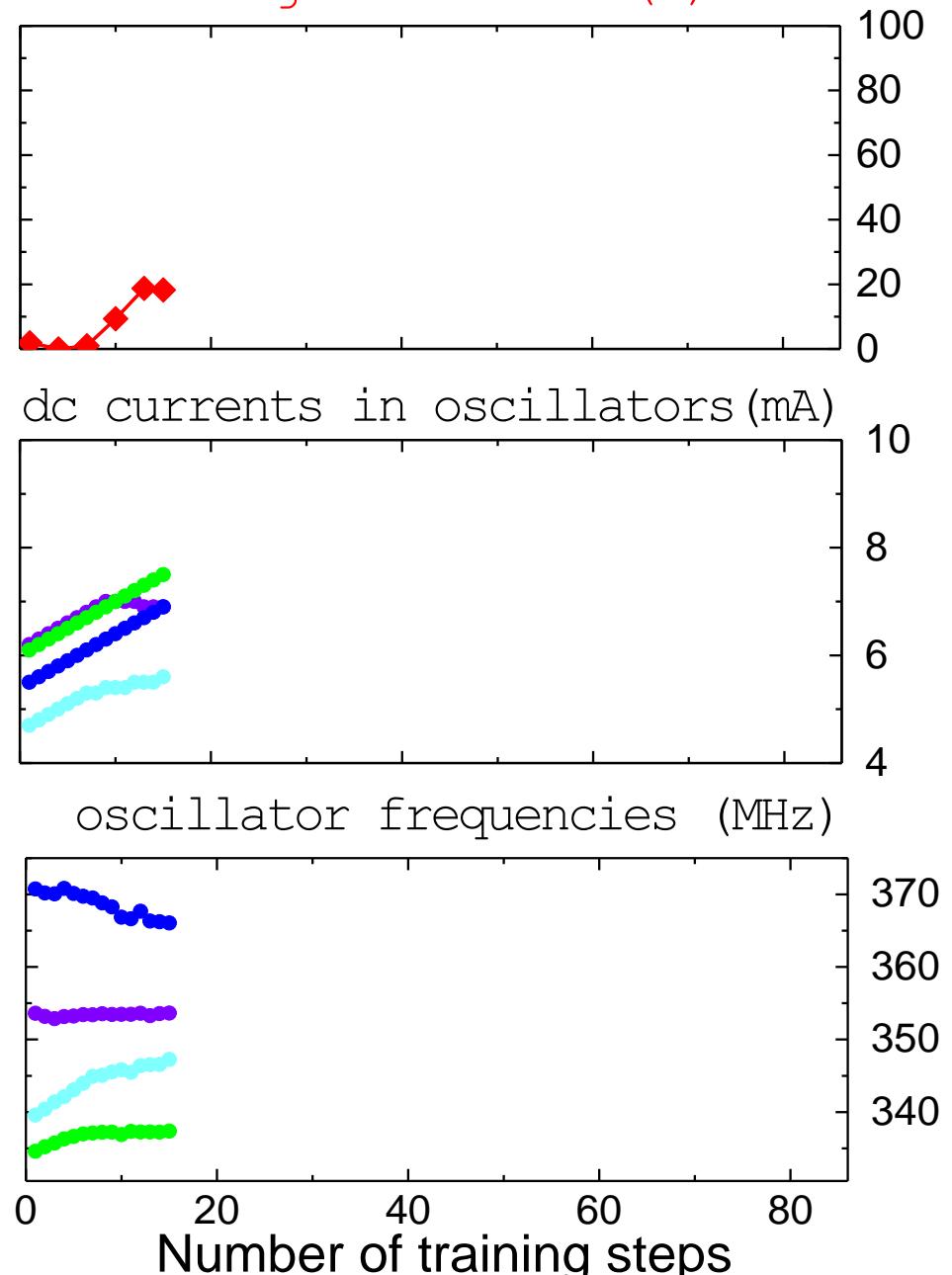
recognition rate (%)



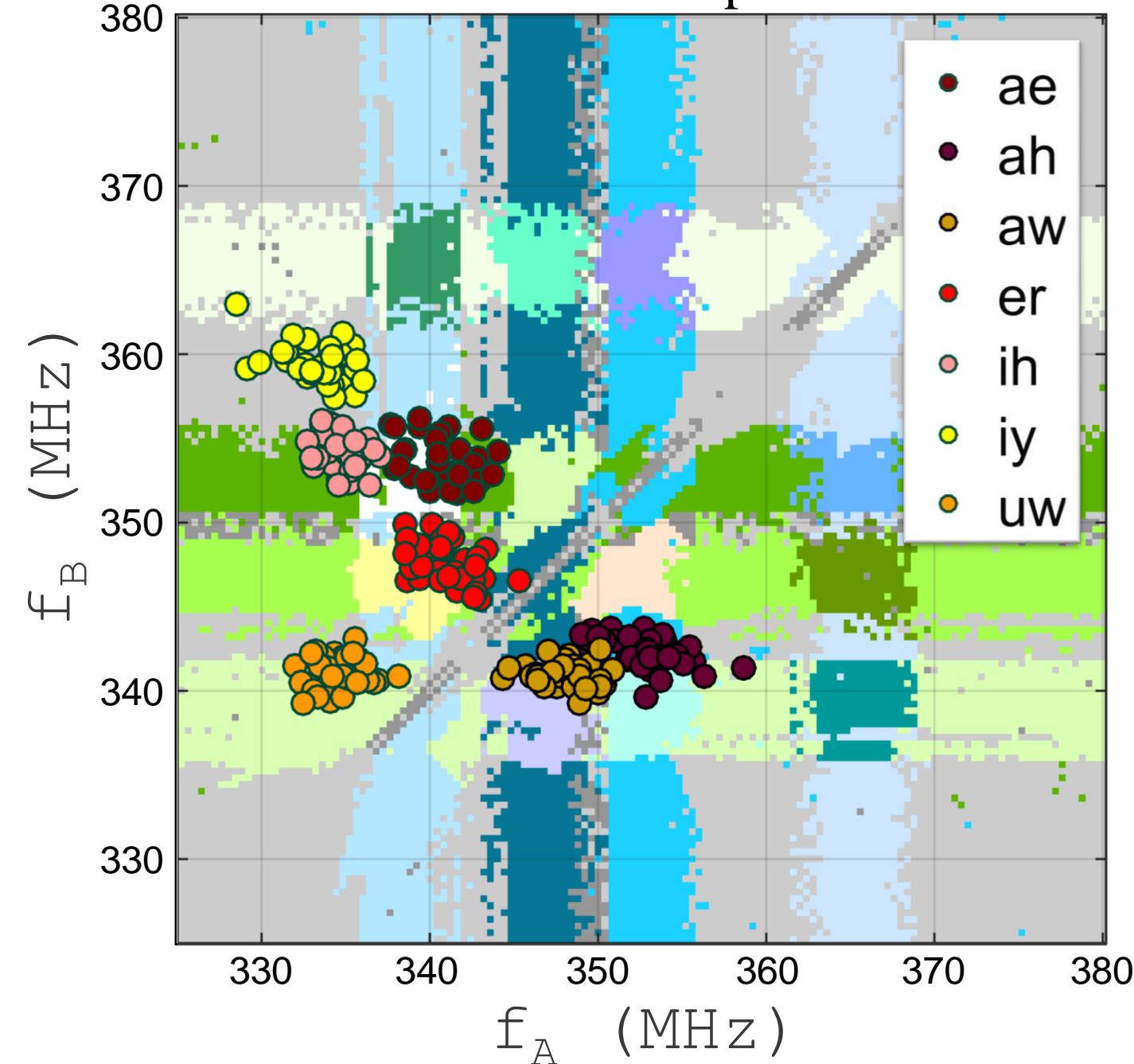
After 15 steps:



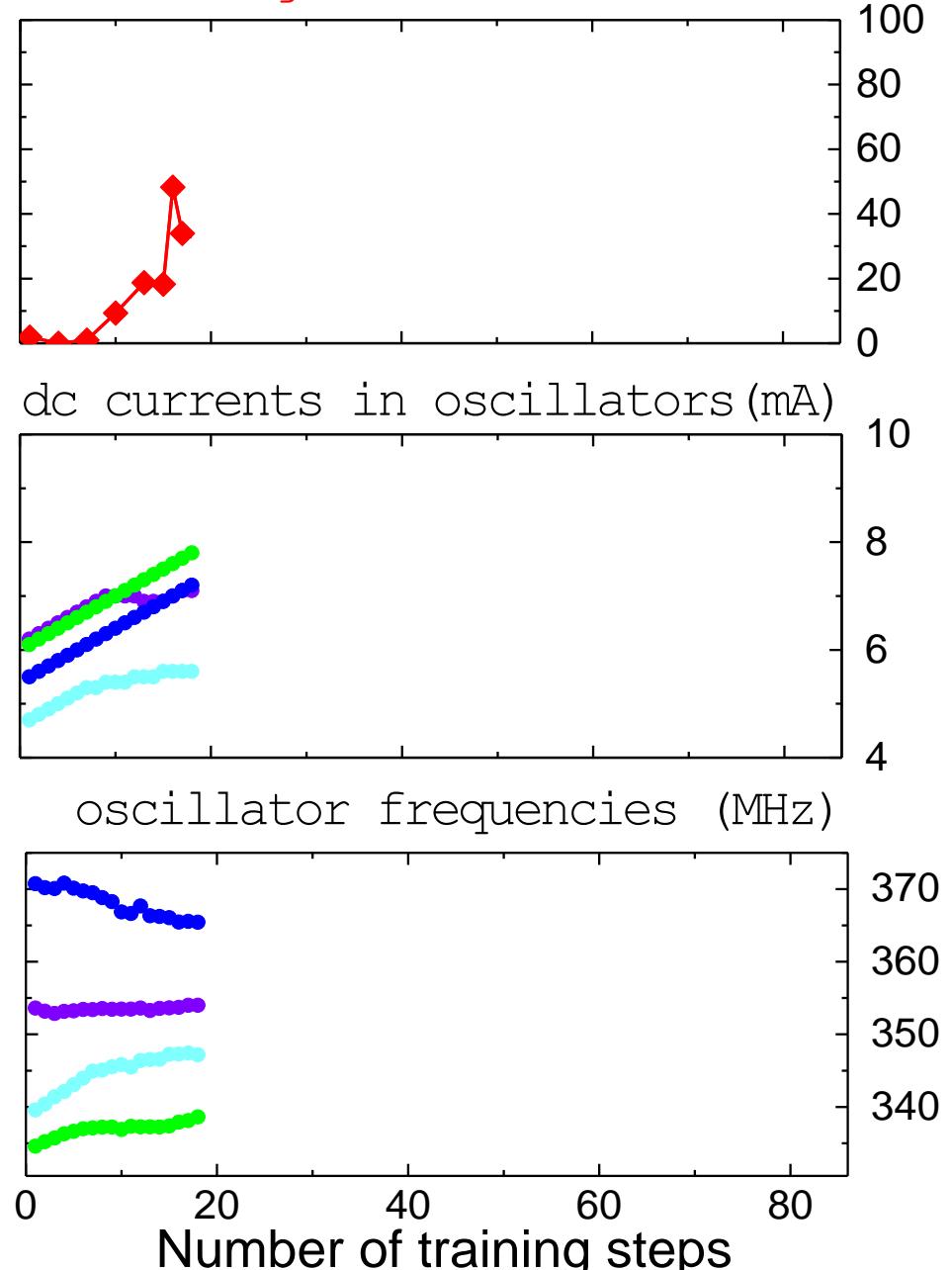
recognition rate (%)



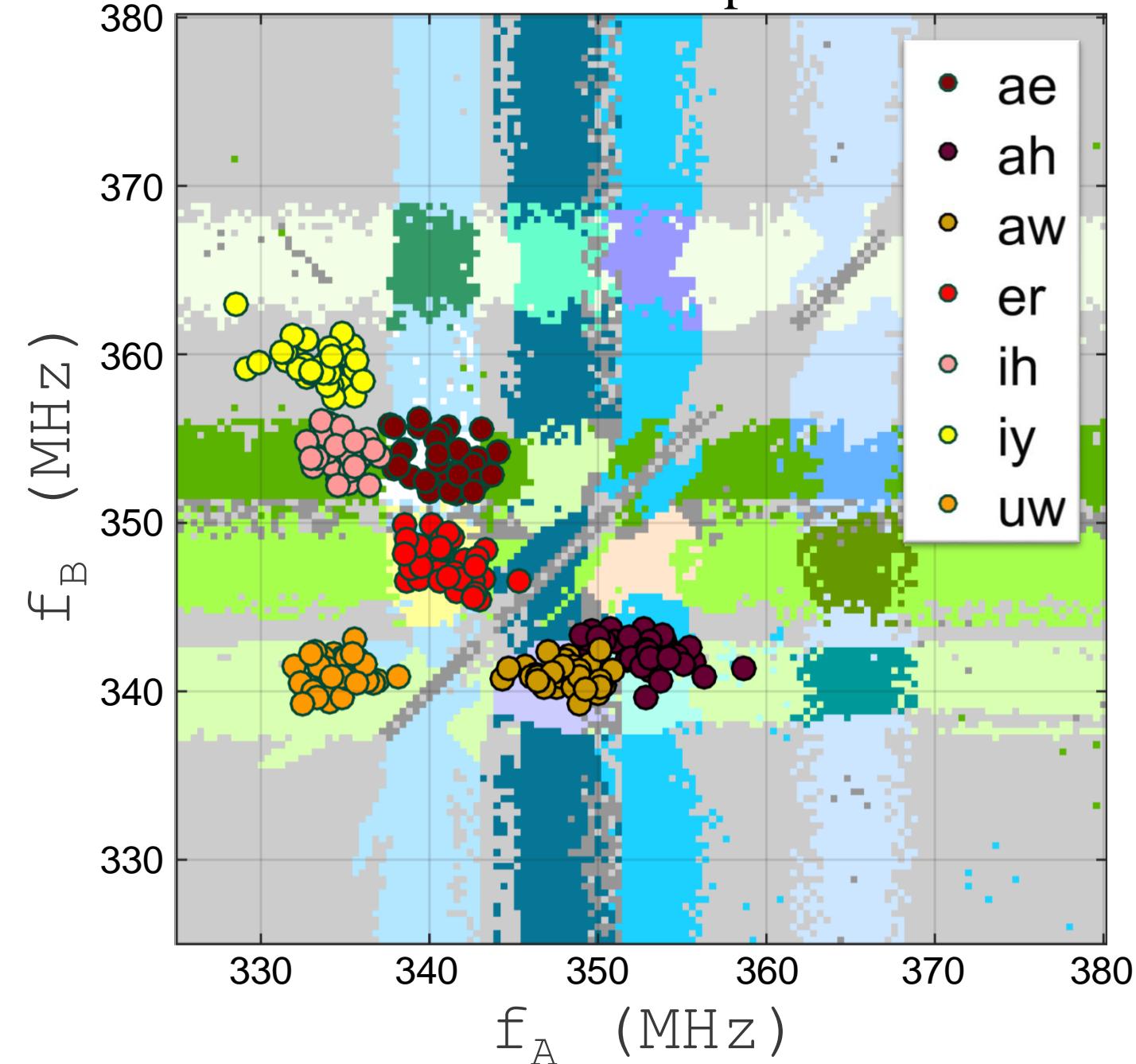
After 18 steps:



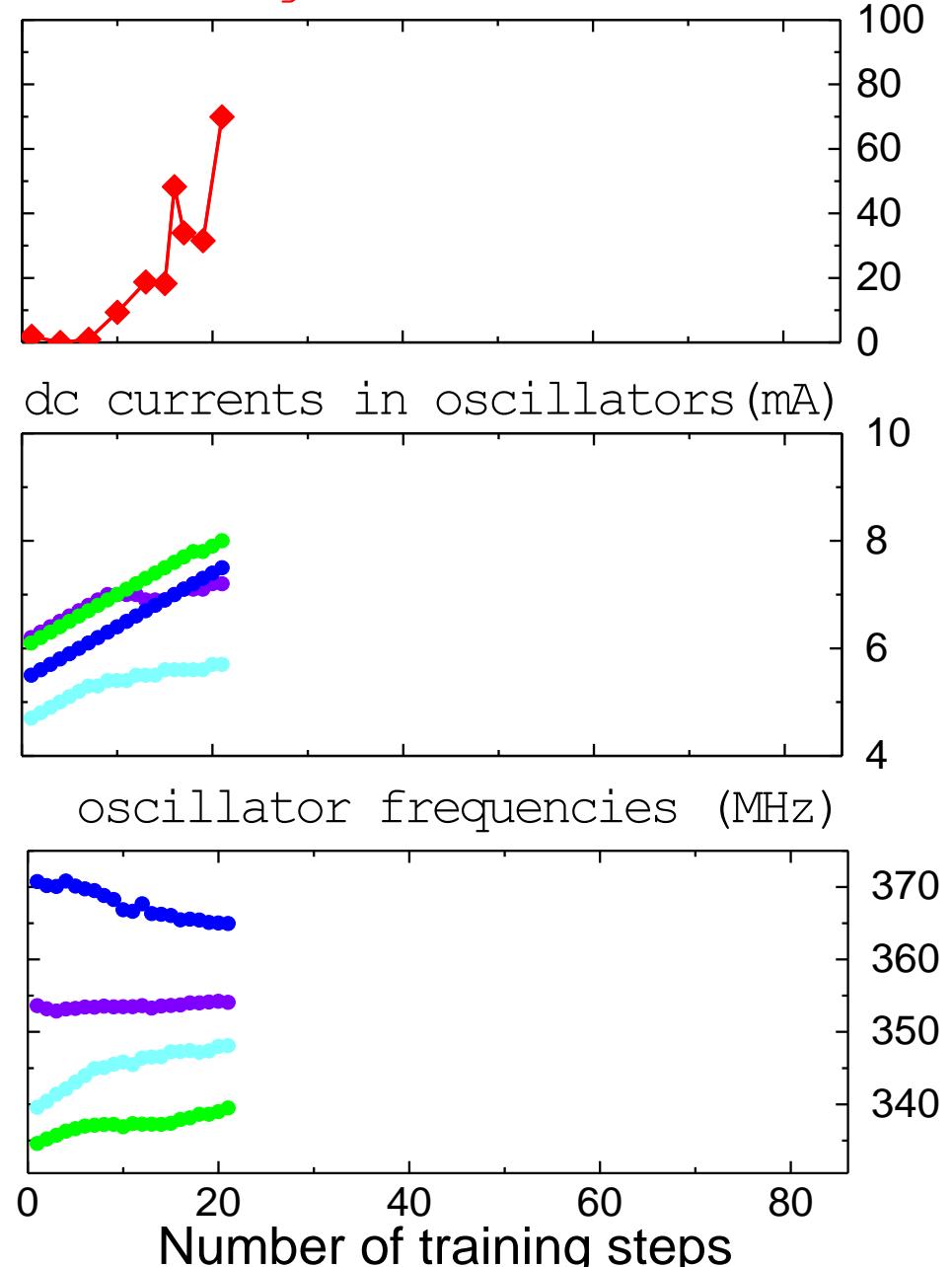
recognition rate (%)



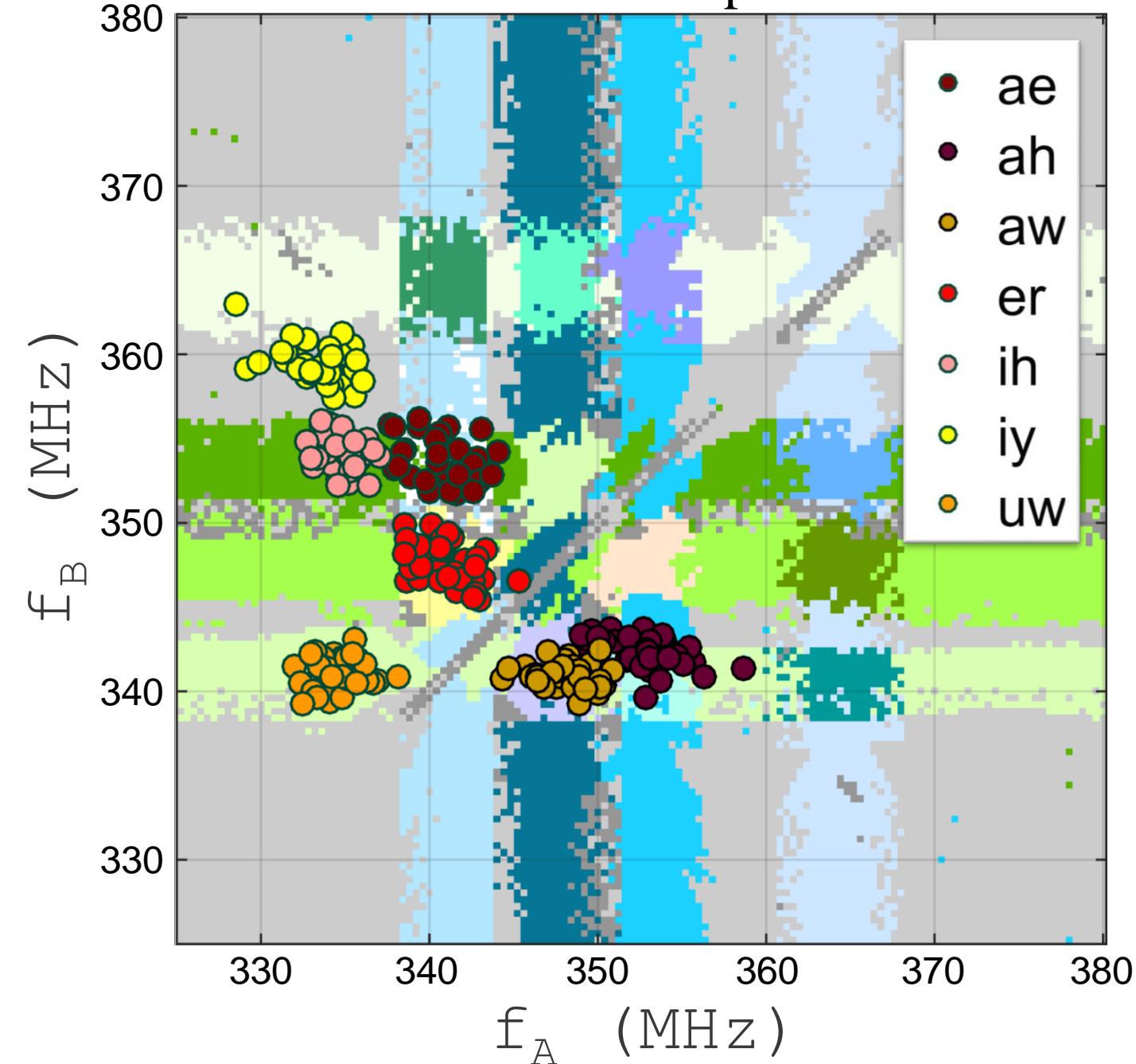
After 21 steps:



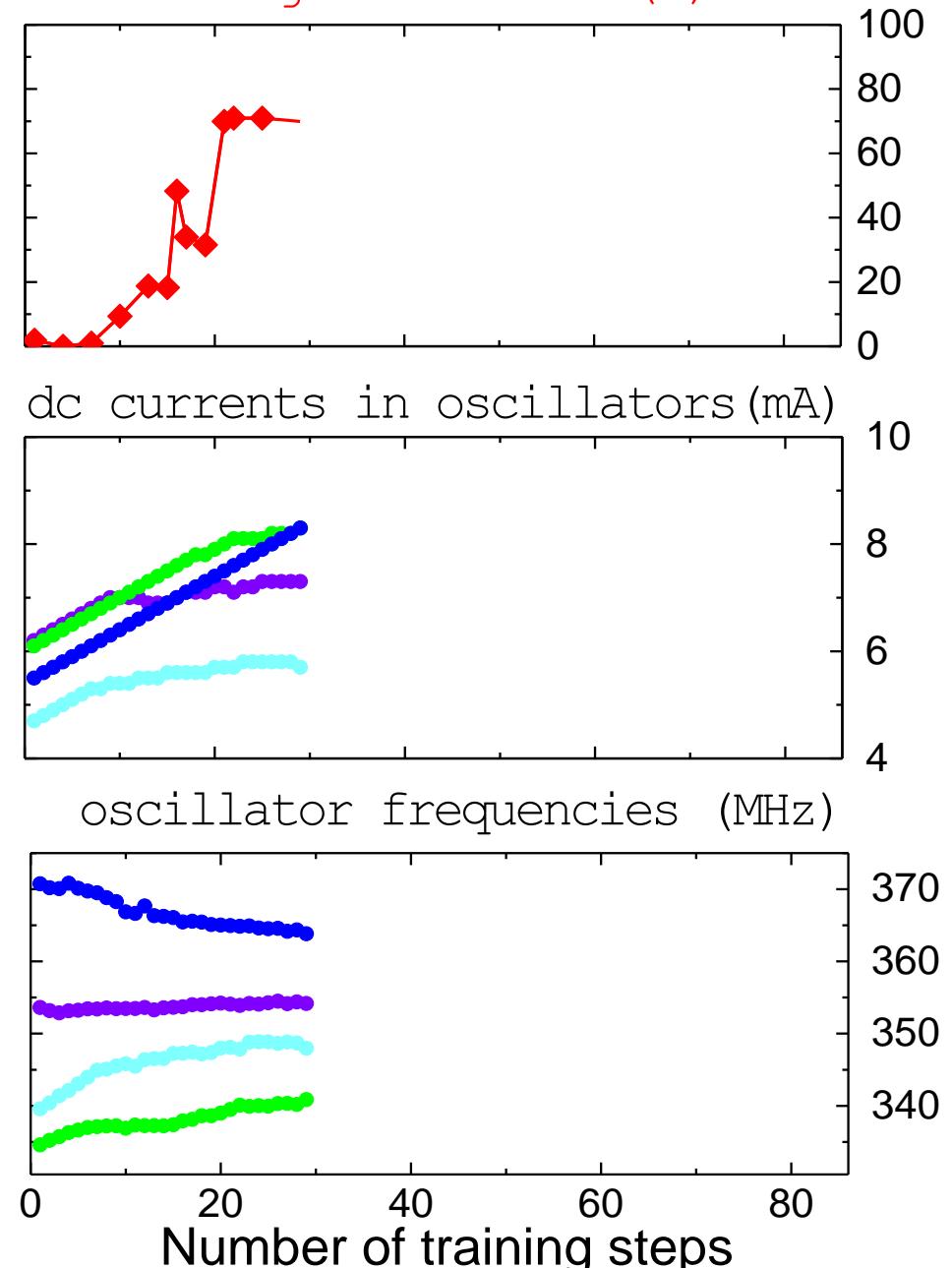
recognition rate (%)



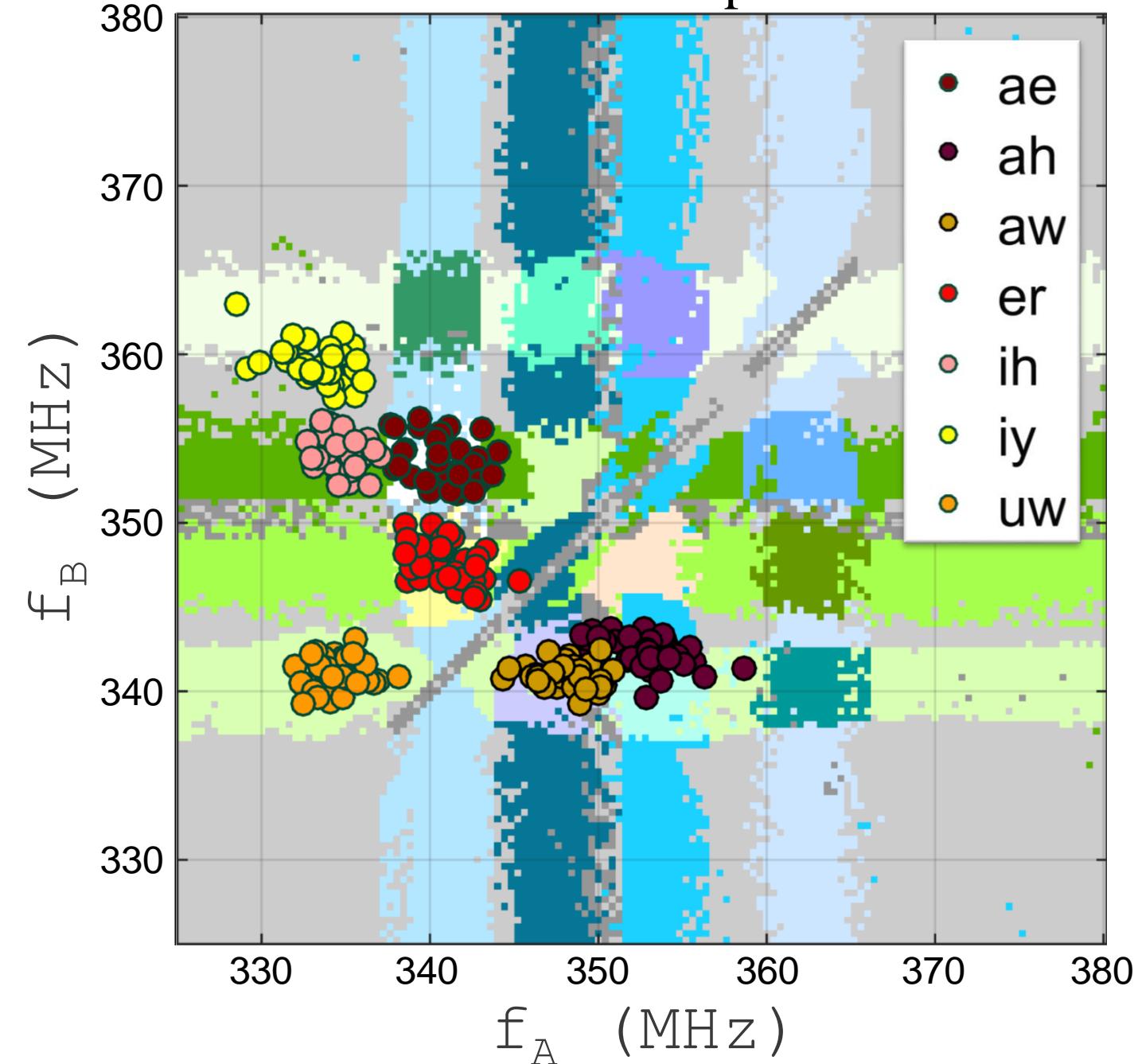
After 29 steps:



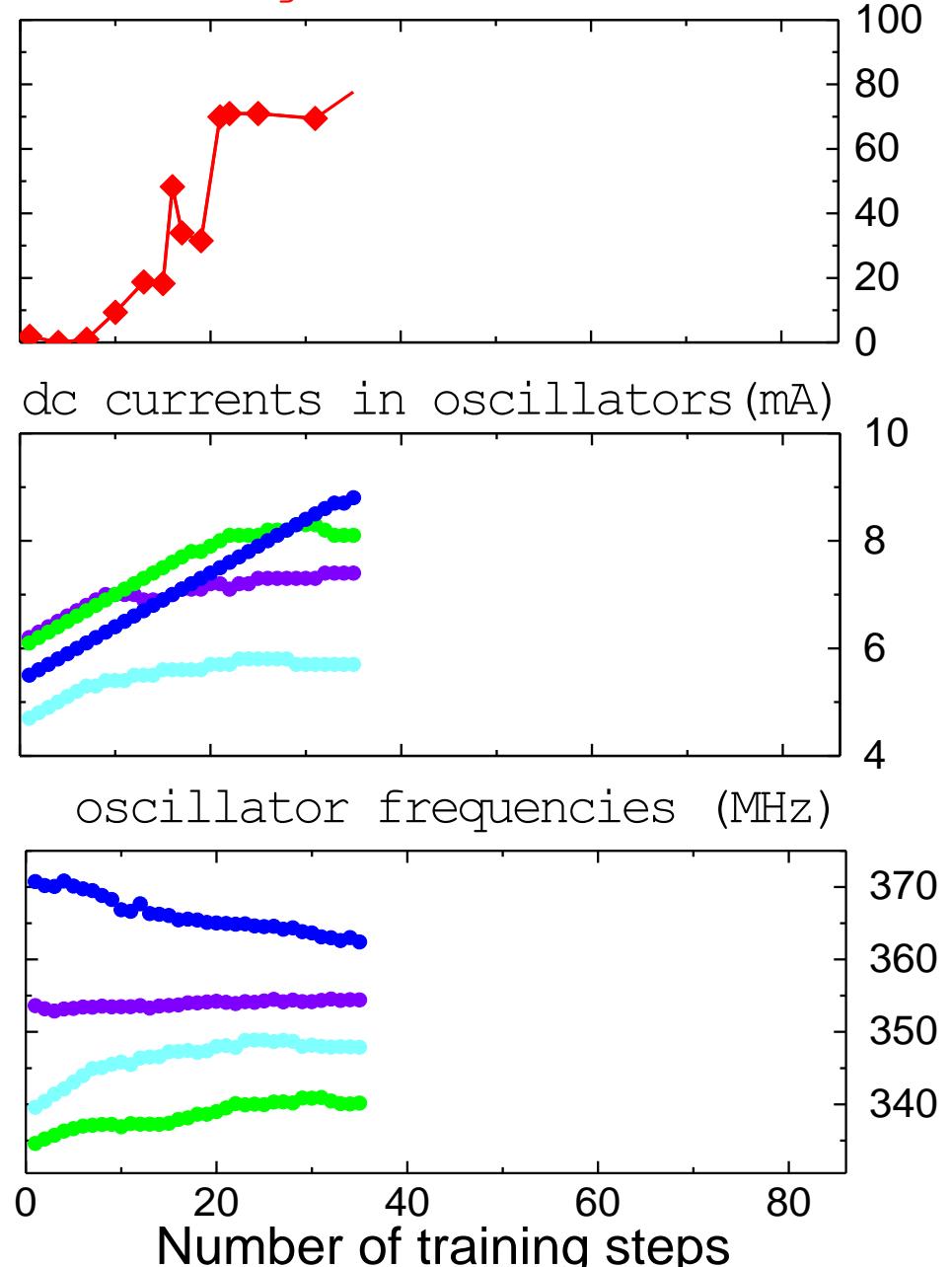
recognition rate (%)



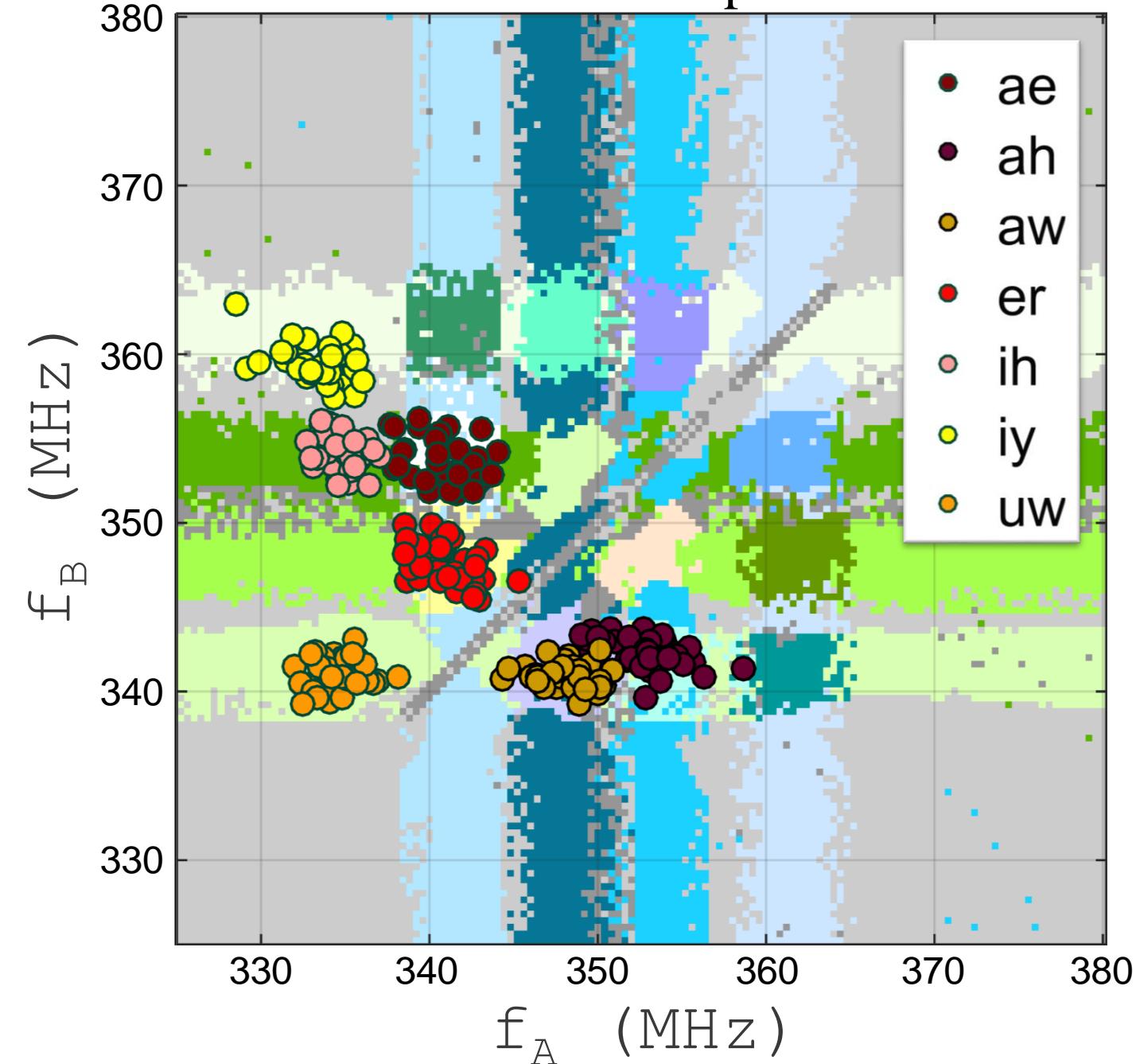
After 35 steps:



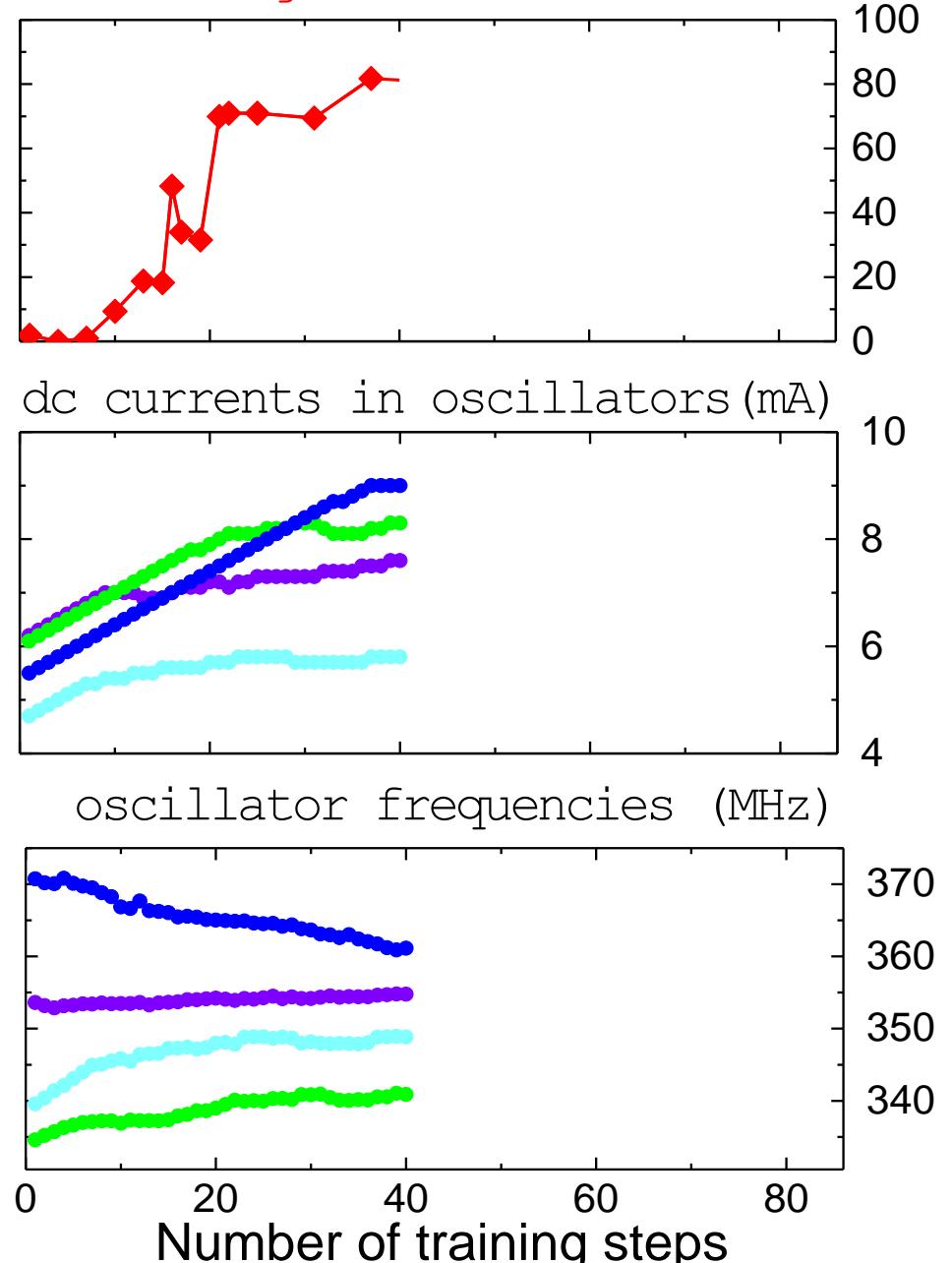
recognition rate (%)



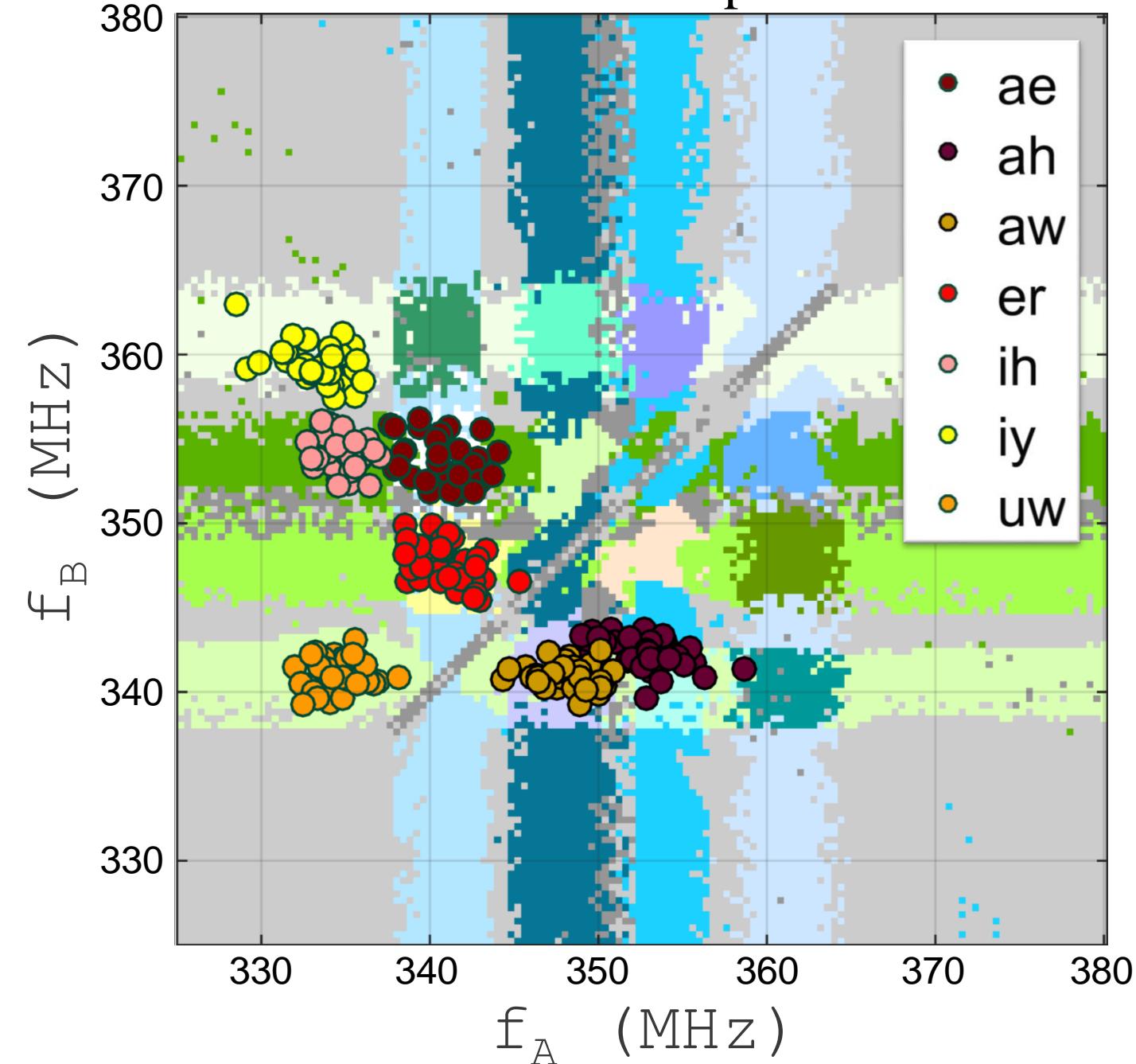
After 40 steps:



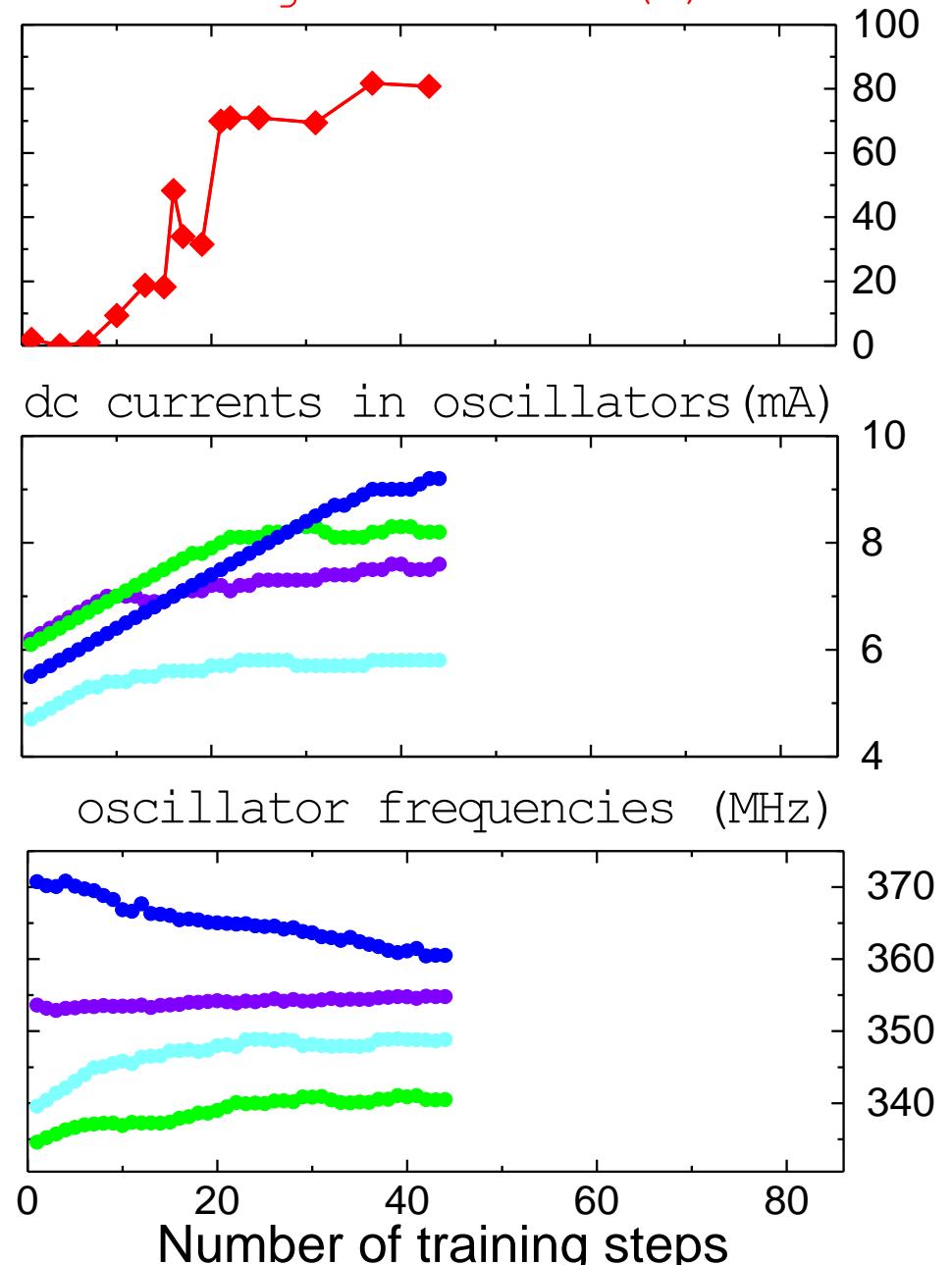
recognition rate (%)



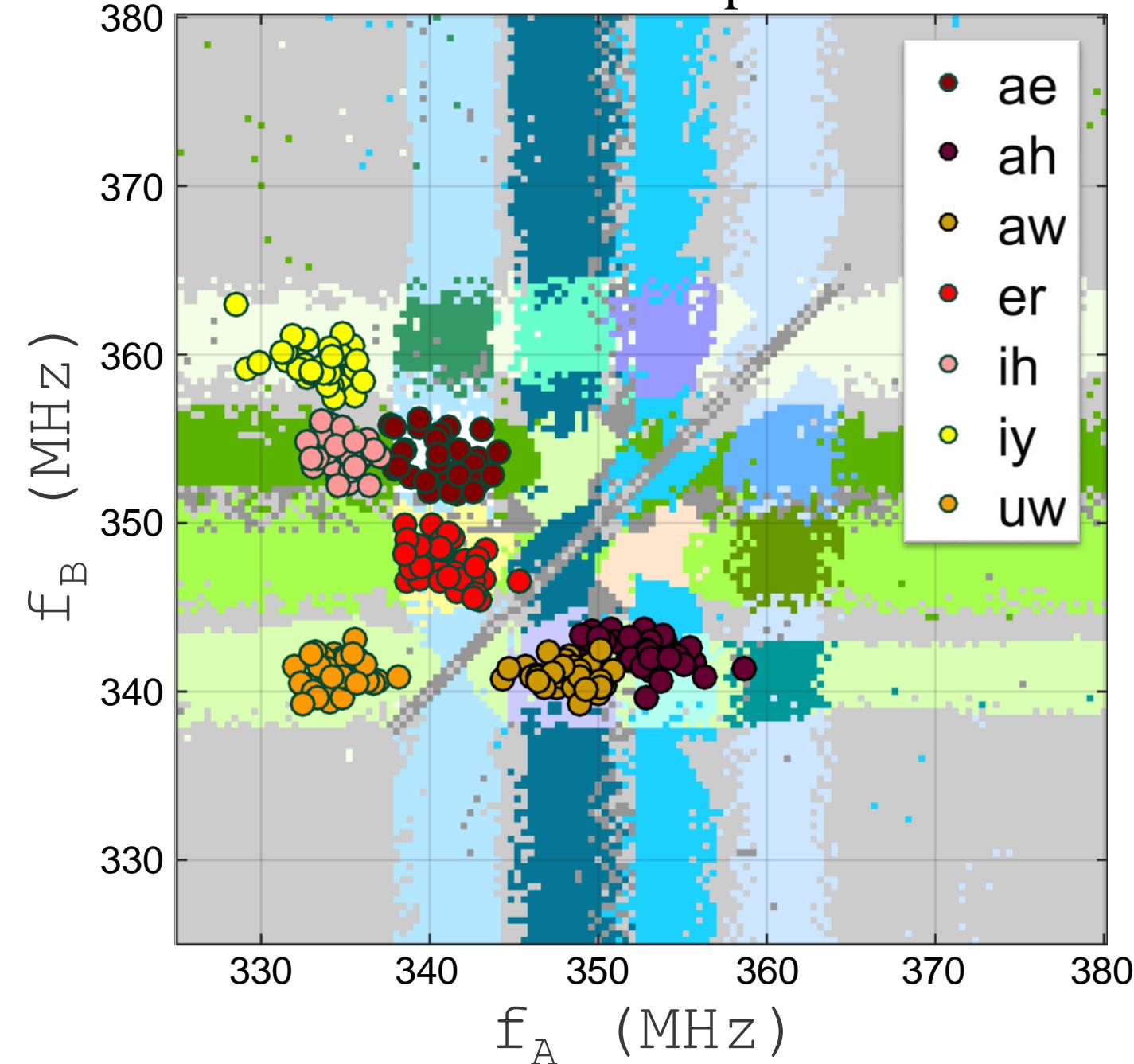
After 44 steps:



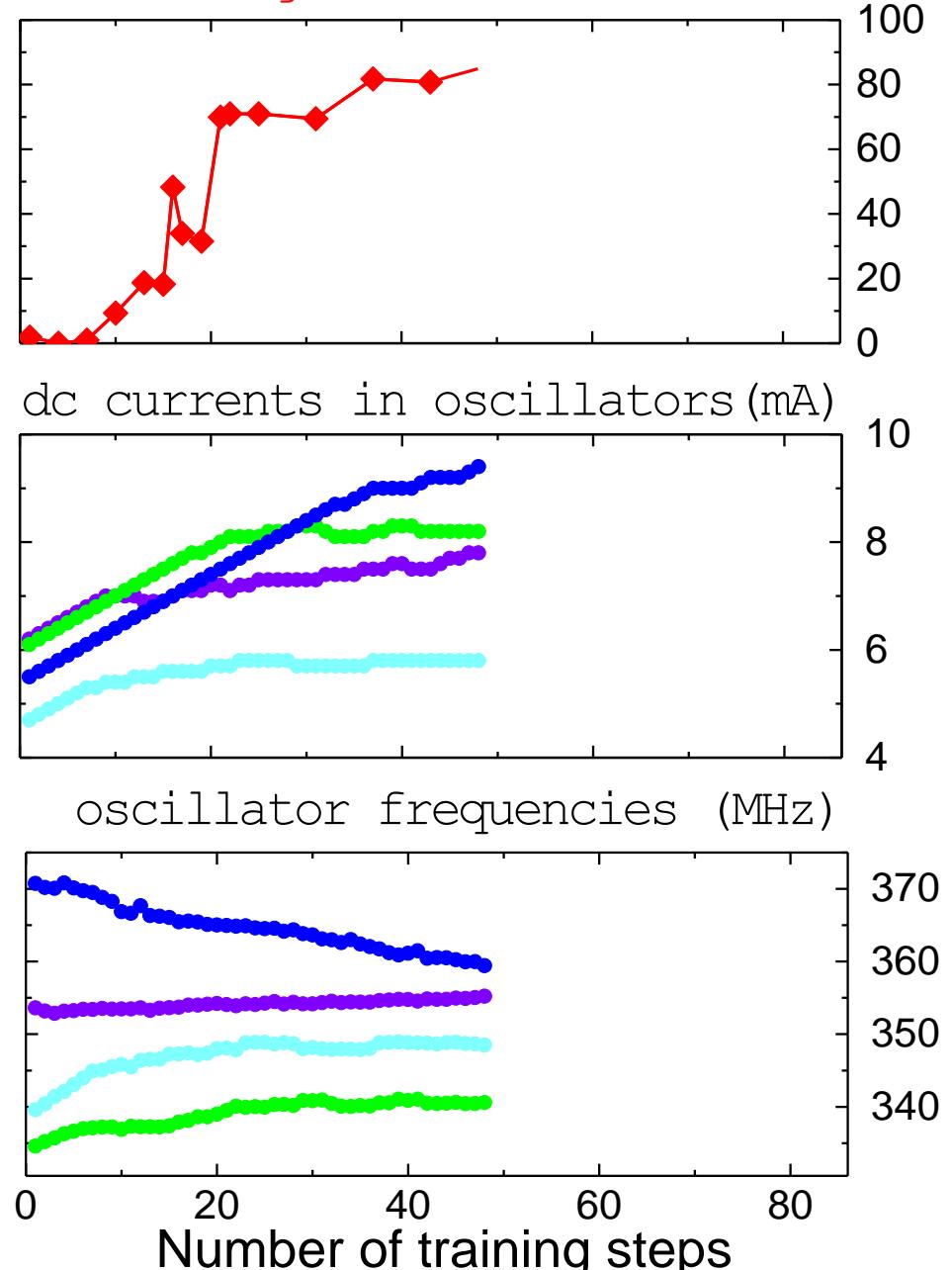
recognition rate (%)



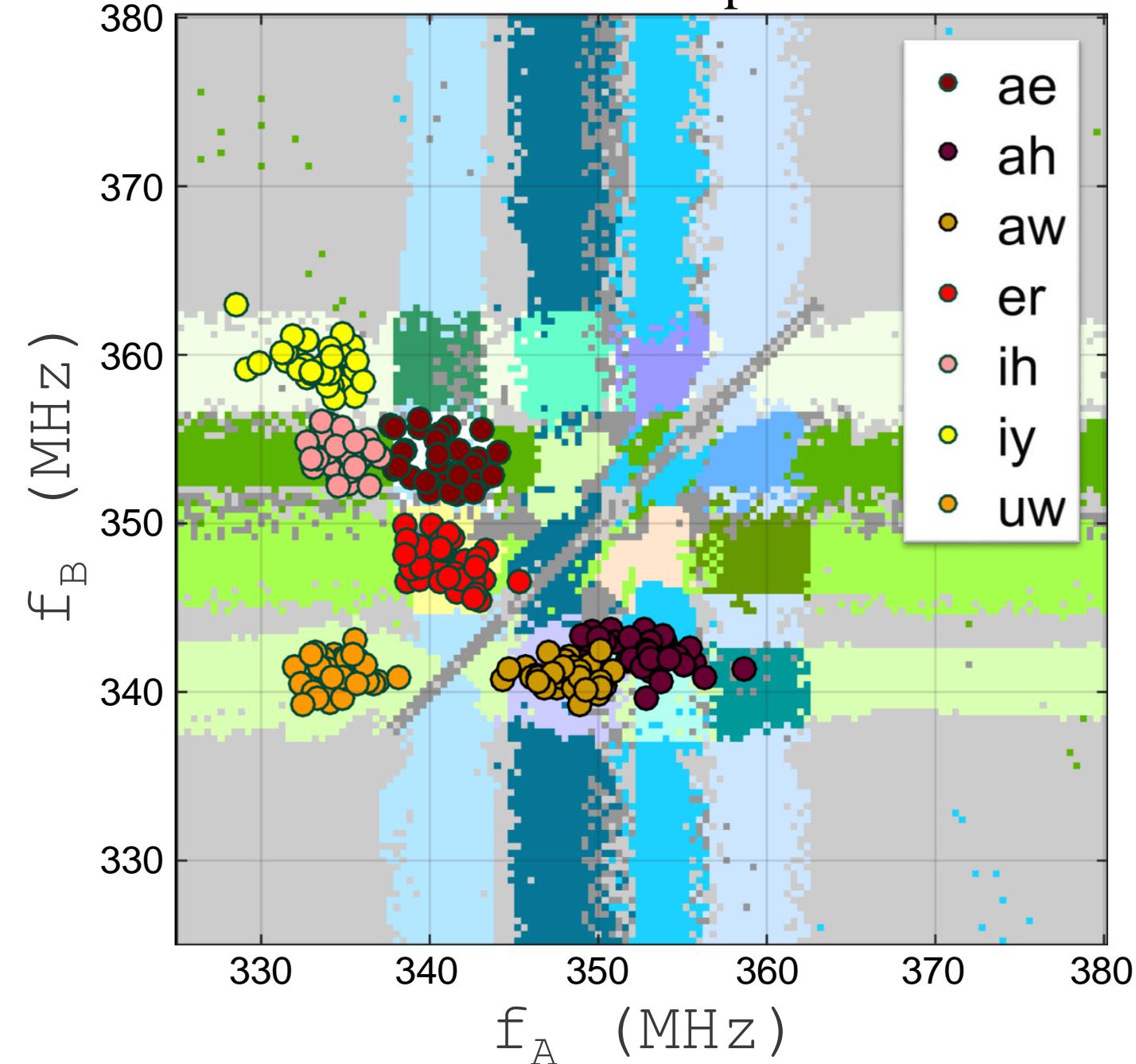
After 48 steps:



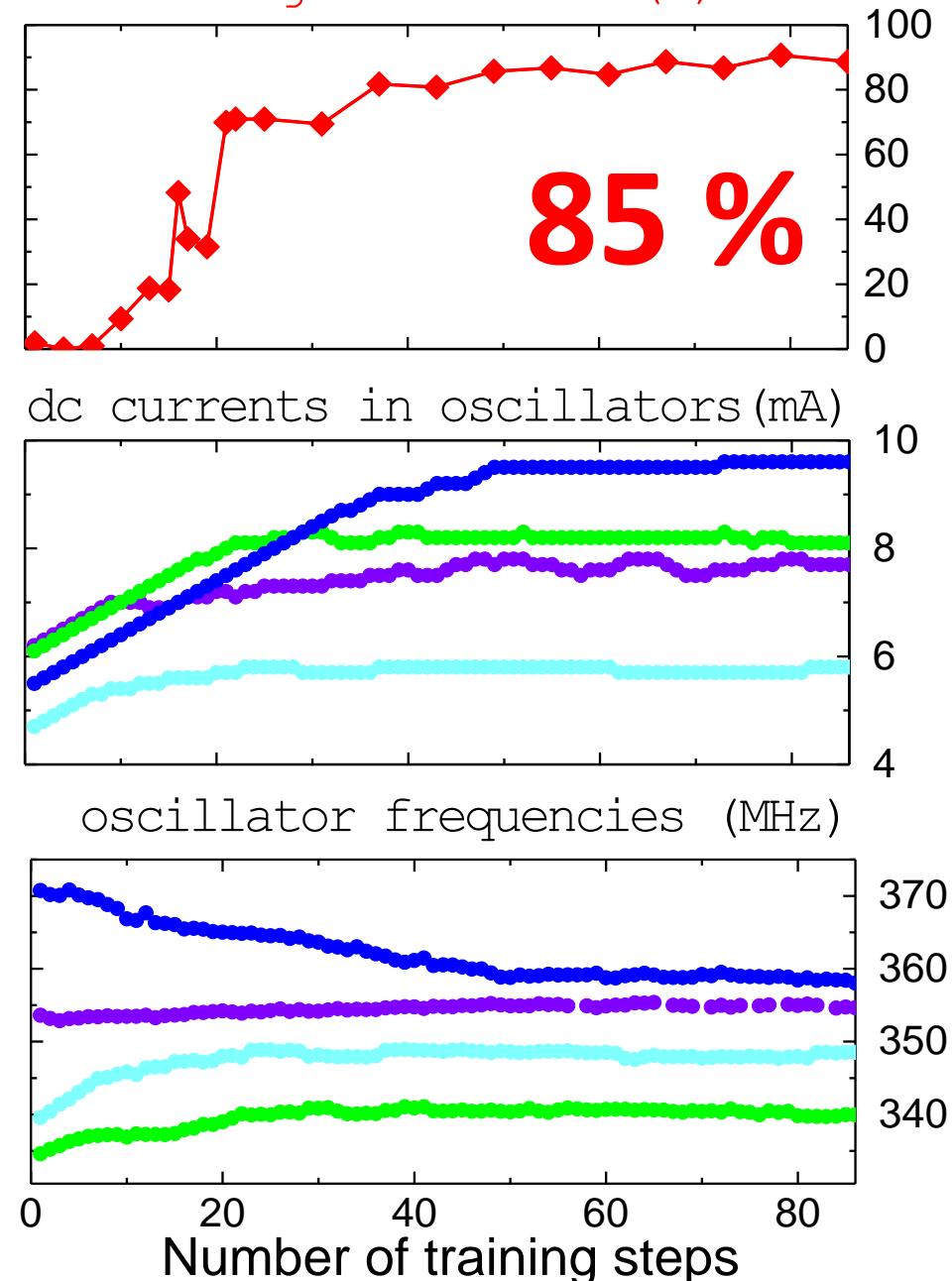
recognition rate (%)



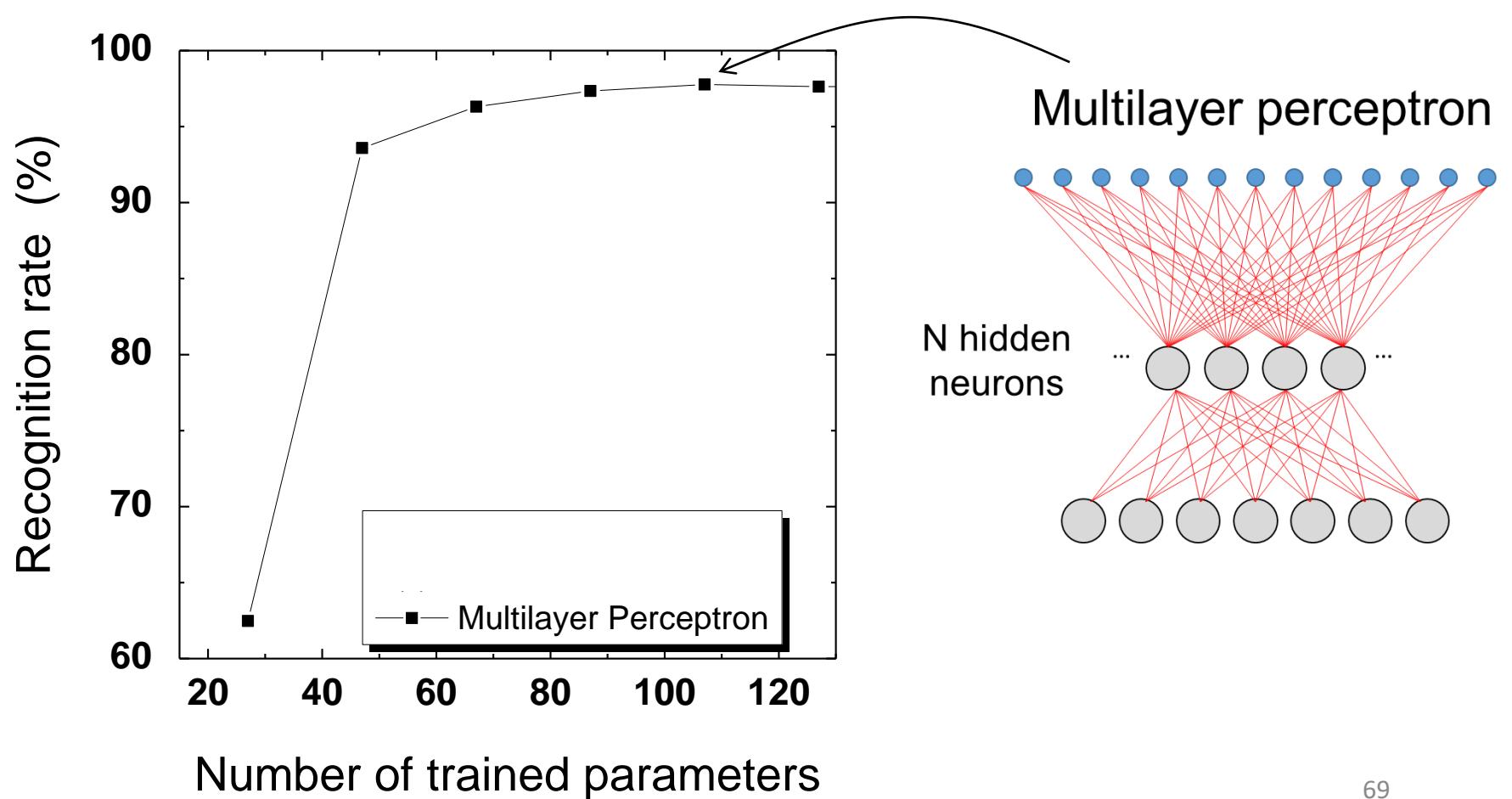
After 86 steps:



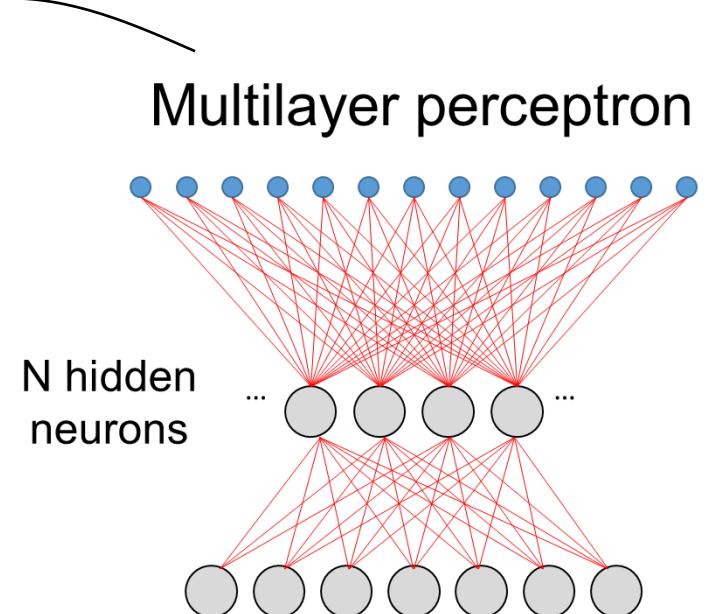
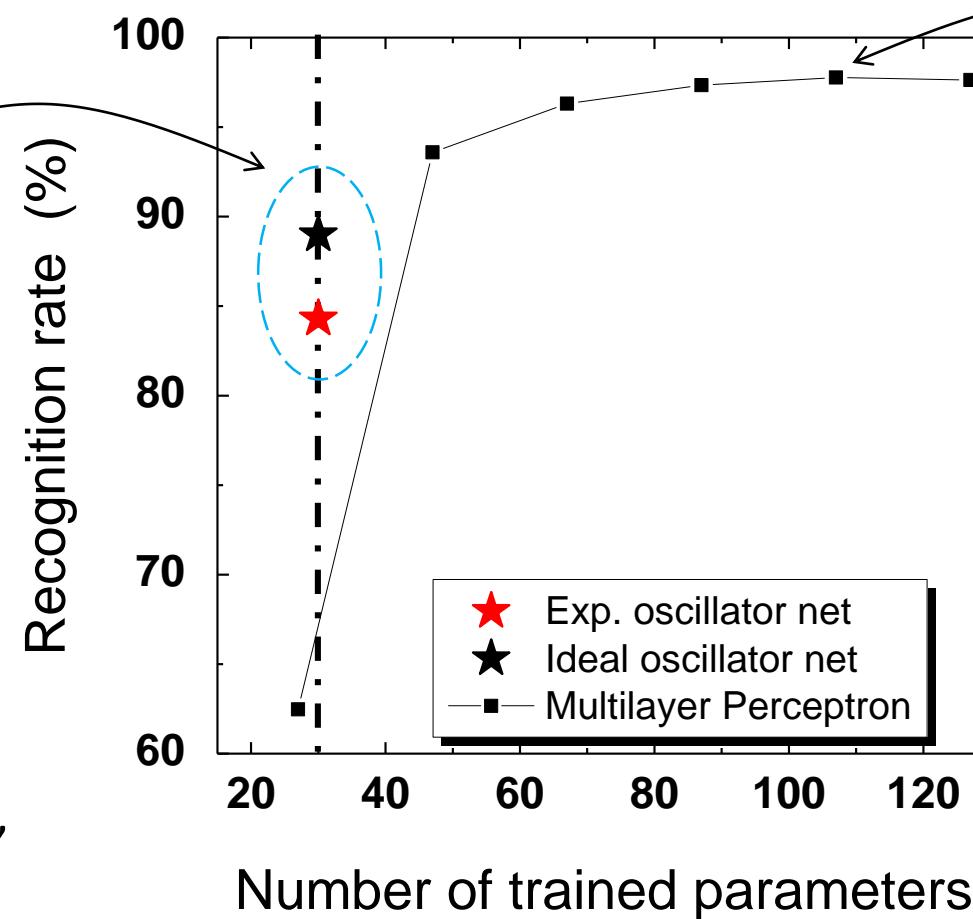
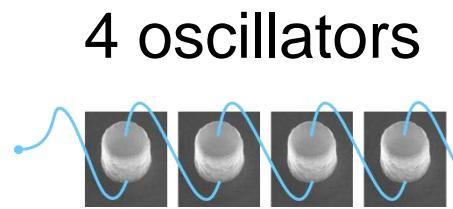
recognition rate (%)



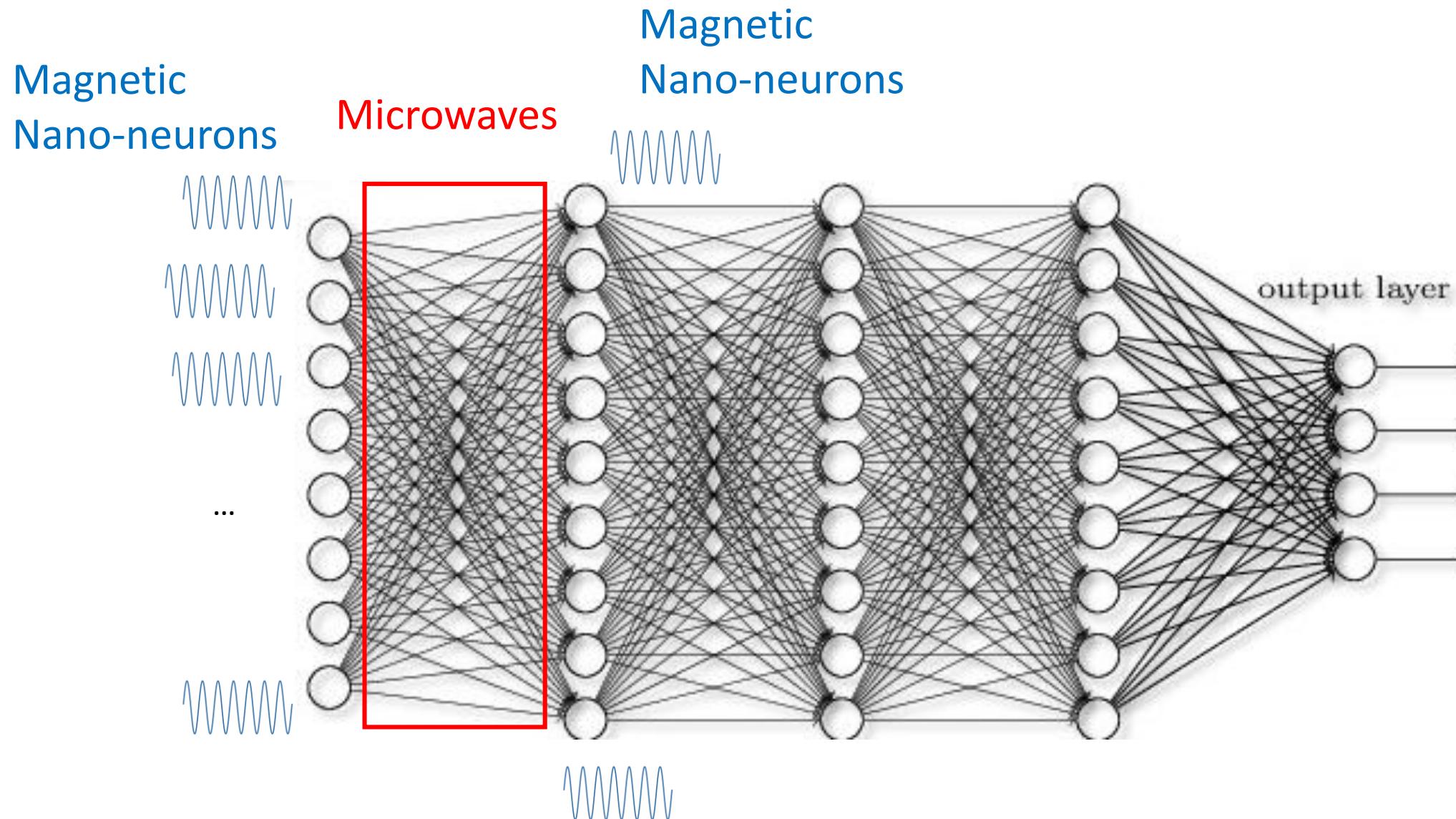
# For comparison, we trained a multilayer perceptron on the same database



# Coupled oscillators achieve similar recognition rates slightly higher than conventional neural networks with the same number of parameters



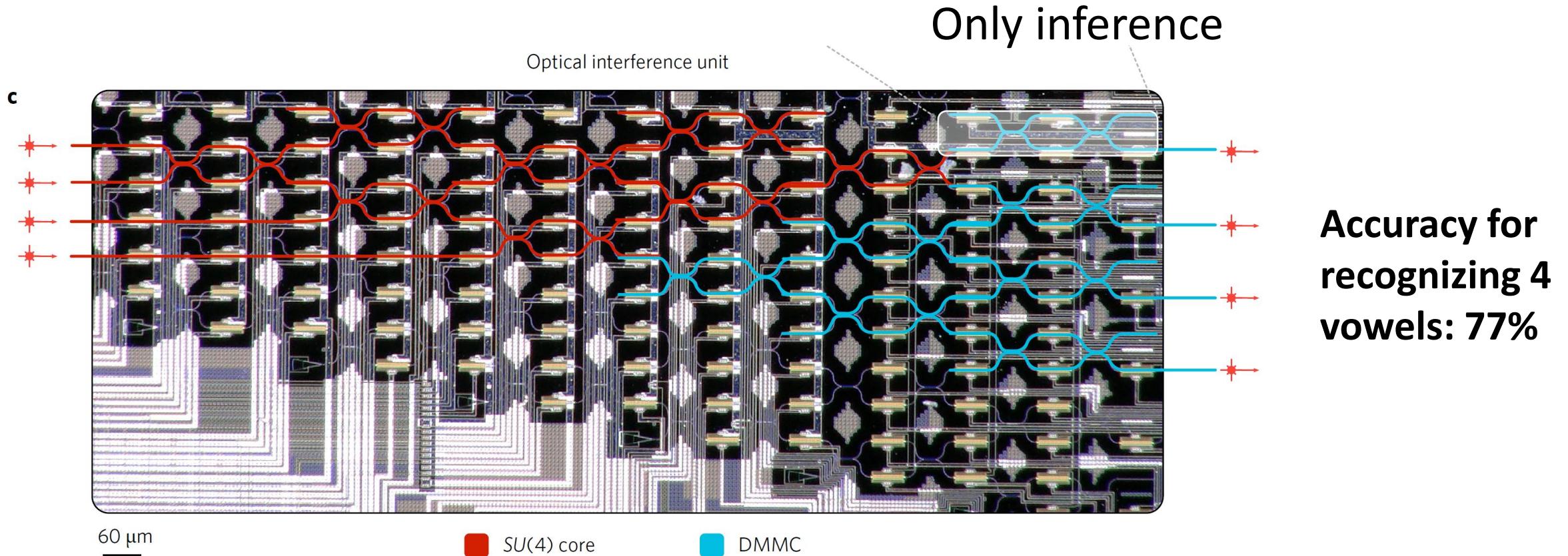
# Perspectives: deep learning with spintronic nanodevices ?



# Conclusion

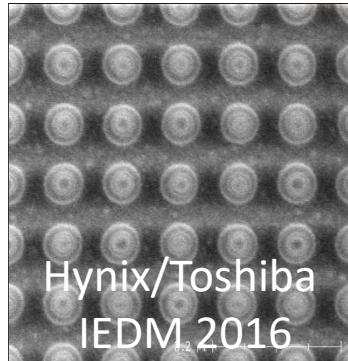
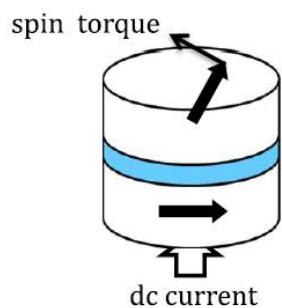
- Limitations of CMOS chips
- Current efforts to implement low energy deep learning with emerging nanodevices
- Futuristic ideas

# Deep learning in optics

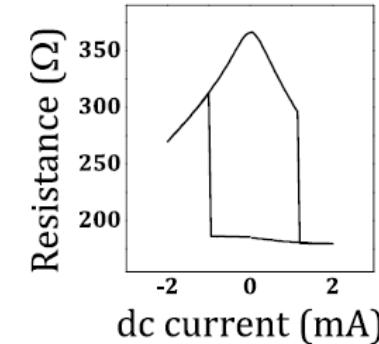


# Spintronics is a toolbox for neuromorphic computing

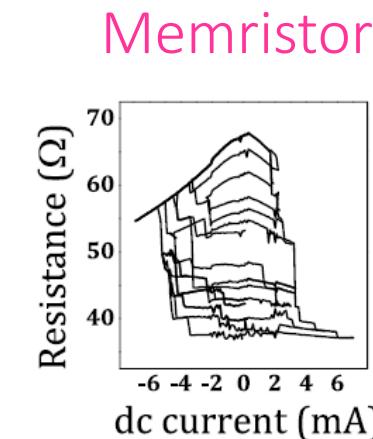
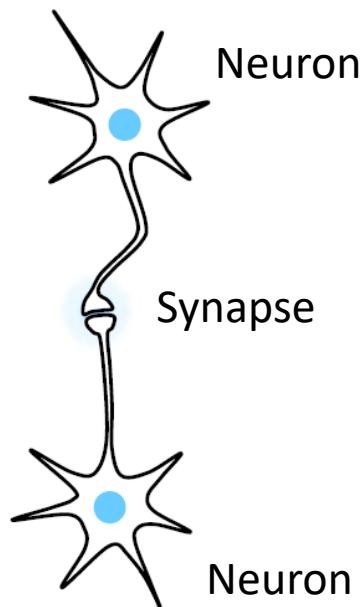
J. Grollier, D. Querlioz  
and M. D. Stiles, *PIEEE*  
104, 2024 (2016)



## Binary Memory

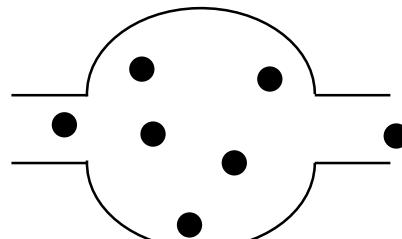


Ma, Endoh *et al*  
2016 *Jpn. J. Appl. Phys.* **55** 04EF15



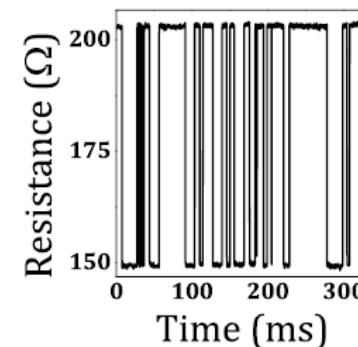
S. Lequeux *et al*,  
*Sci. Rep.* (2016)

## Sklyrmions



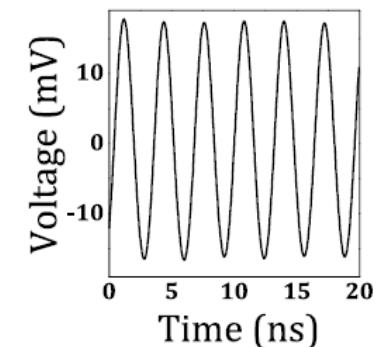
D. Pinna *et al*, *PRAppl* (2018)  
D. Prychynenko *et al*, *PRAppl* (2018)

## Stochasticity



A. Mizrahi *et al*,  
*Nature Com* (2018)

## Non-linear dynamics



J. Torrejon, M. Riou *et al*,  
*Nature* (2017)