# Neuromorphic Computing

# 1. Introduction

Ben Abdallah Abderazek, Khanh N. Dang
E-mail: {benab, khanh}@u-aizu.ac.jp

# Lecture Contents

1. <span style="color:red">Neuromorphic Computing</span>

2. Hardware Models of Spiking Neurons

3. Synaptic Dynamics

4. Synaptic Plasticity Mechanisms and Learning

5. Synthesizing Real-Time Neuromorphic Systems

6. Conclusions

# 1. Neuromorphic Computing:
## Neural Network Generations



1st generation – perceptron    2nd generation – deep learning    3rd generation – SNN
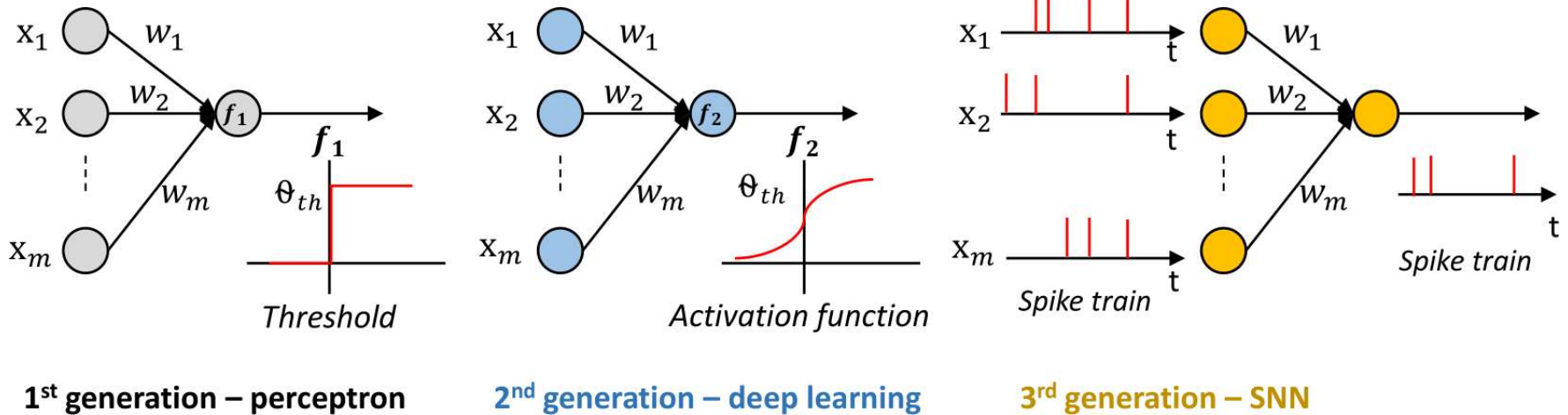
Fig. 1.2: Neural network generations

# 1. Neuromorphic Computing:

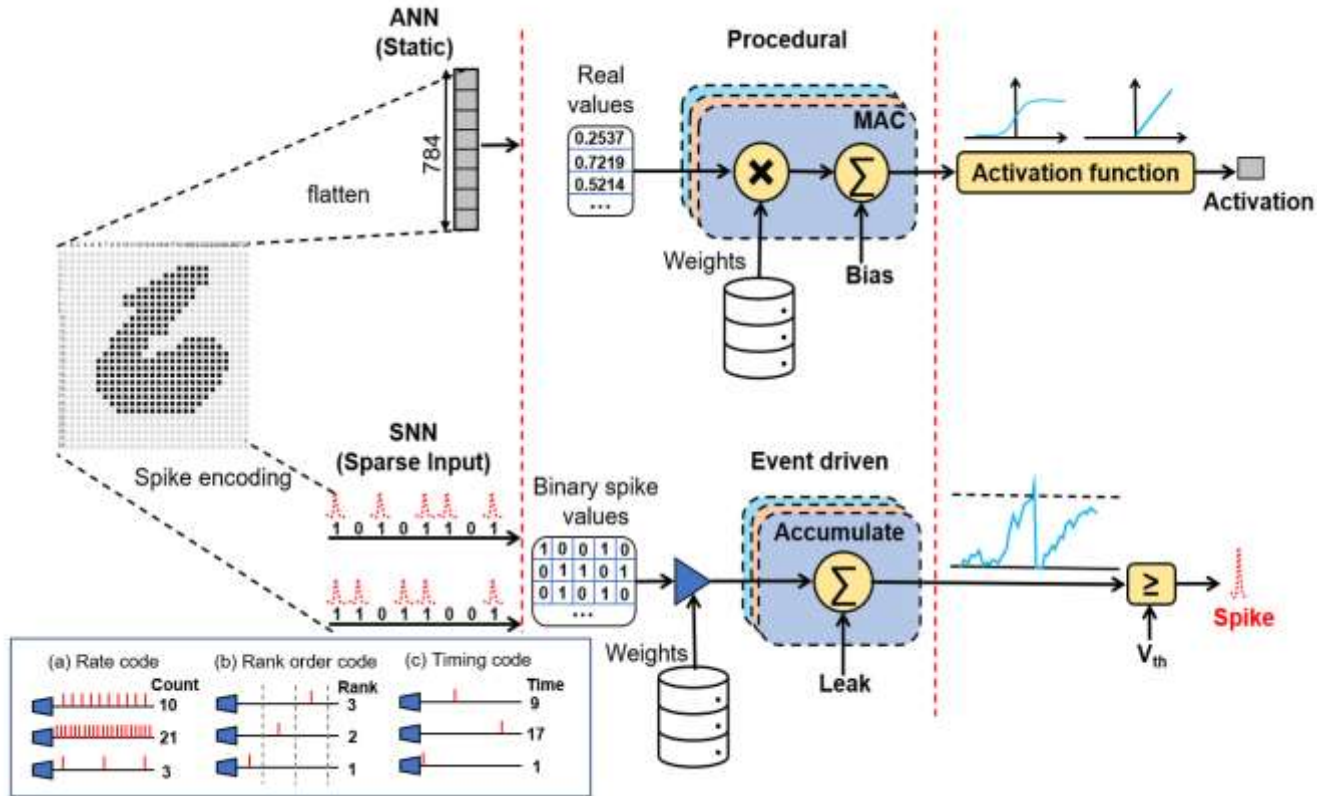## Conventional ANN vs Spiking Neural Network (Neuromorphic)

Artificial Neural Network (ANN) is a brain inspired computing paradigm modeled after the computational principles of the brain's neural network.

Approaches:

- **Conventional ANN:** Impressive results in visual and auditory cognitive applications. However, they are:
  - Slow when deployed in software, requiring a lot of time for training
  - Consume a lot of power when accelerated in hardware, requiring large servers for training as their sizes increase.
- **Spiking Neural Network (Neuromorphic):**
  - *More analogous to the brain*, communicating via spikes in a sparse event driven manner.
  - Exploits spike sparsity to achieve low-power.

# 1. Neuromorphic Computing:

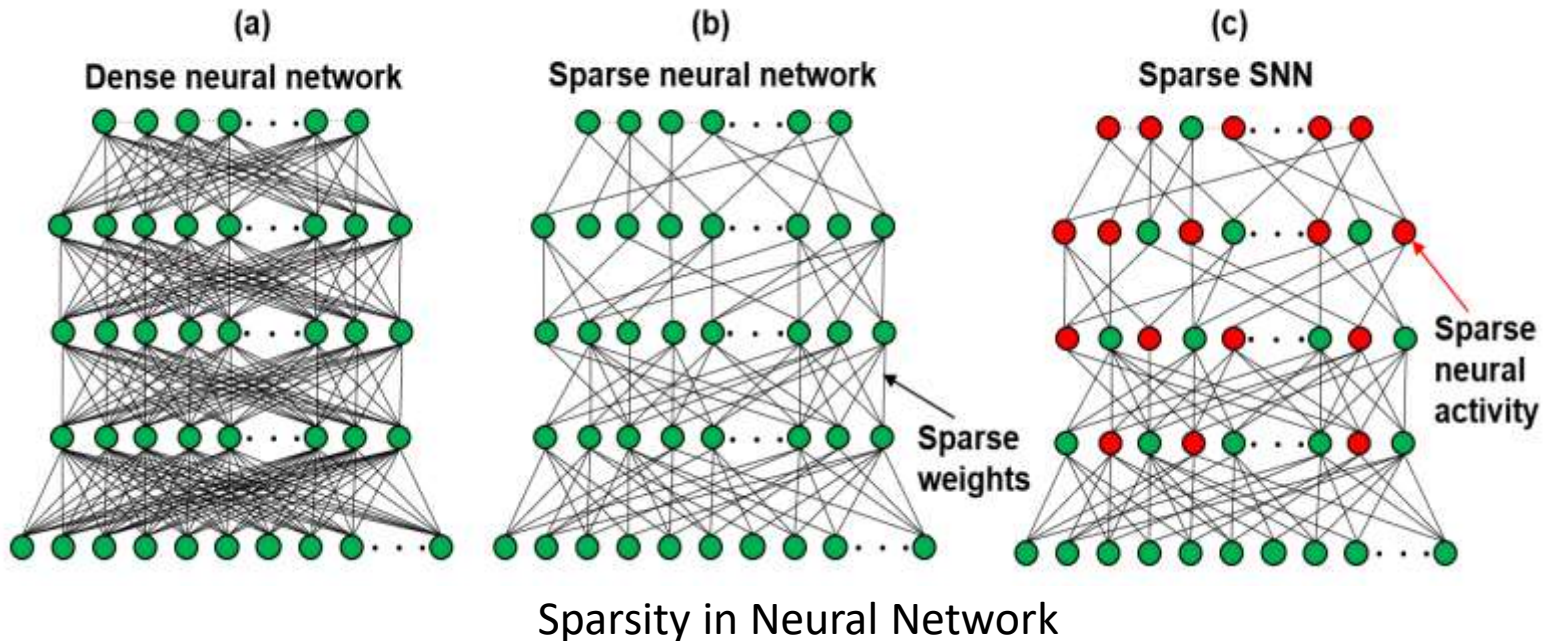## Conventional ANN vs Spiking Neural Network (Neuromorphic)



Conventional ANN vs Spiking Neural Network

- Sparse input in SNN means sparse memory use.
- Spike communication means minimal power per event signal
- Event based processing in SNN also contribute to low power.

# 1. Neuromorphic Computing:
## Exploiting Sparsity in Neural Network



Sparsity in Neural Network

- About 0.5% to 2% of neurons in the neocortex are active at any time
- Only about 1% to 5% of connections exist between two connected layers in the neocortex and 30% of those connections change every few days
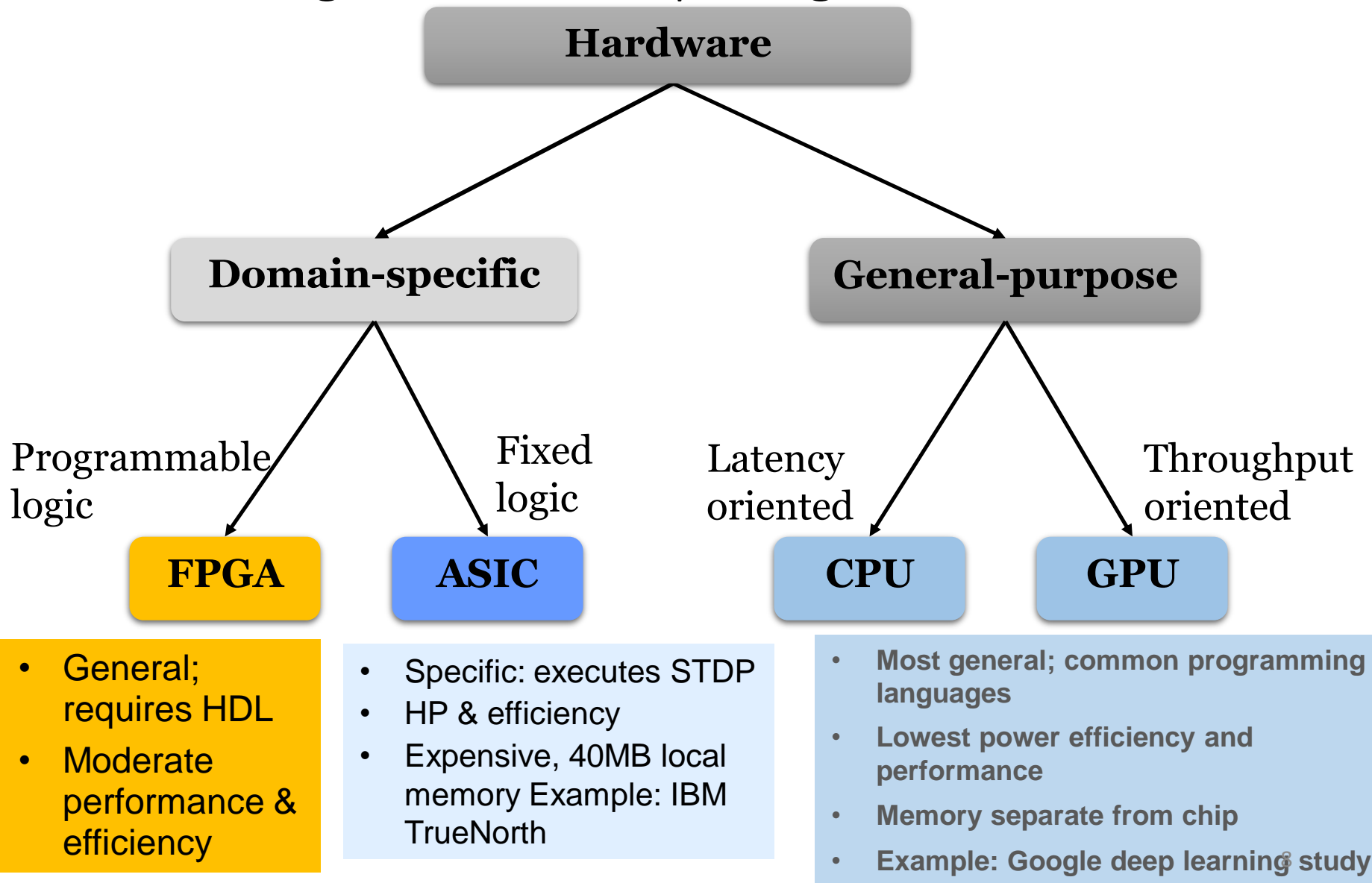
# 1. Neuromorphic Computing:

## What is Neuromorphic Computing?

- **Neuromorphic Computing** is the use of **hardware (VLSI)** to simulate the biological architecture of the **human nervous system** (brain, complex network of nerves, etc.),

- Neuromorphic Engineering is a new emerging field that involves biology, physics, mathematics, and computer science and engineering to design hardware models **of neural and sensory systems.**

- Neuromorphic systems opens new frontiers for neuro-robotics, artificial intelligence, and high-performance applications.
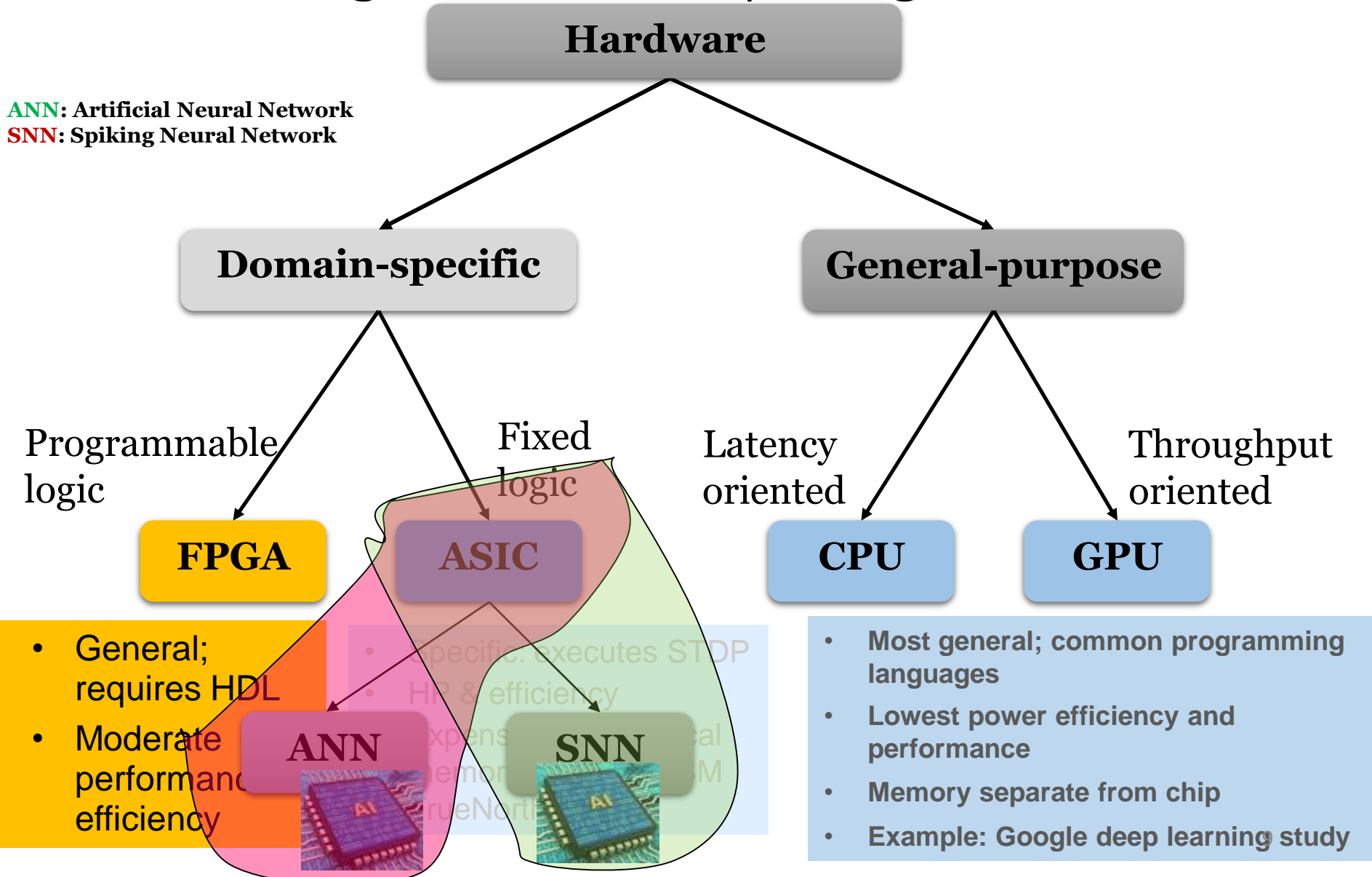
# 1. Neuromorphic Computing:
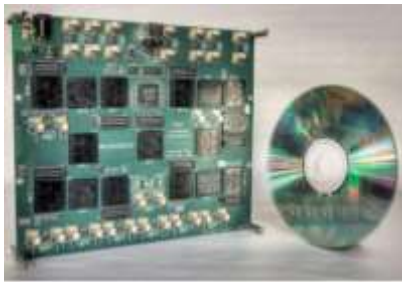## Neural Algorithms Computing in Hardware



**Hardware**

**Domain-specific**

**General-purpose**

Programmable logic

Fixed logic

Latency oriented

Throughput oriented

**FPGA**

**ASIC**

**CPU**

**GPU**

- General; requires HDL
- Moderate performance & efficiency

- Specific: executes STDP
- HP & efficiency
- Expensive, 40MB local memory Example: IBM TrueNorth

- Most general; common programming languages
- Lowest power efficiency and performance
- Memory separate from chip
- Example: Google deep learning study

# 1. Neuromorphic Computing:
## Neural Algorithms Computing in Hardware

**Hardware**

**ANN**: Artificial Neural Network
**SNN**: Spiking Neural Network

**Domain-specific**

**General-purpose**

Programmable logic

Fixed logic

Latency oriented

Throughput oriented

**FPGA**

**ASIC**

**CPU**

**GPU**

- General; requires HDL
- Moderate performance efficiency

- Specific: executes STDP
- HP & efficiency
- Expensive; critical memory: SRAM TrueNorth

**ANN**

**SNN**

- Most general; common programming languages
- Lowest power efficiency and performance
- Memory separate from chip
- Example: Google deep learning study

# 1. Neuromorphic Computing:
## Examples of Neuromorphic Chips/Systems



Neurogrid



IBM TrueNorth



Intel Loihi
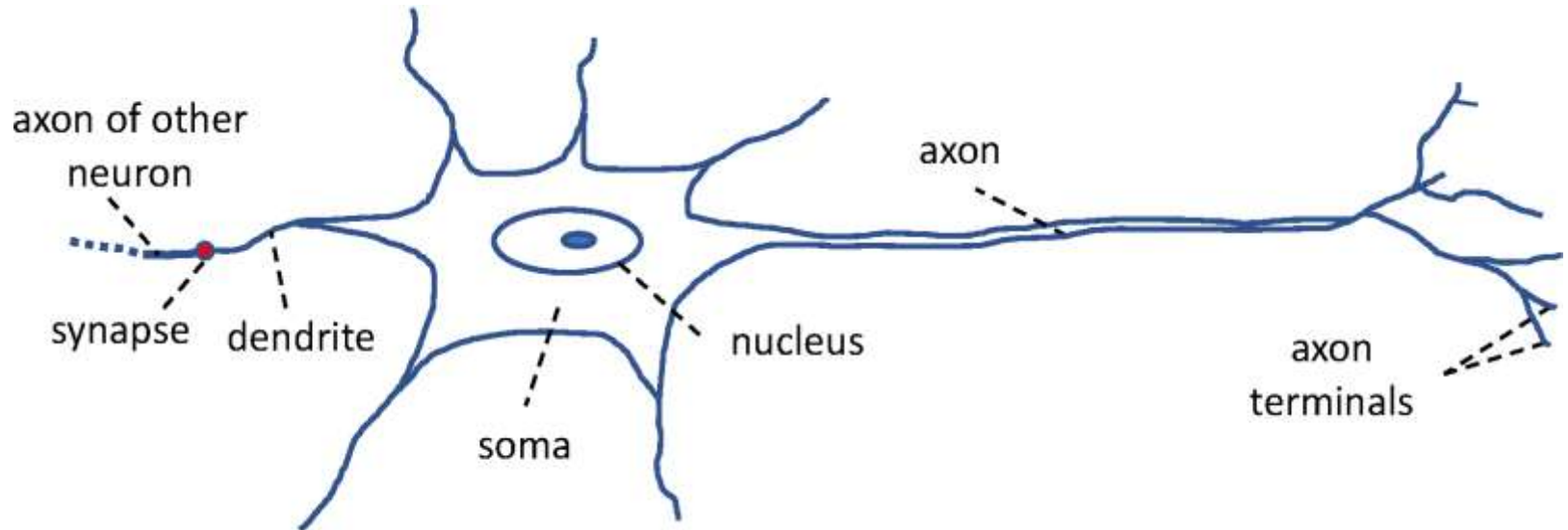
Examples of Neuromorphic Chips/Systems (not yet commercial)

# Lecture Contents

1. Neuromorphic Computing Approaches

2. <span style="color:red">Hardware Models of Spiking Neurons</span>

3. Synaptic Dynamics

4. Synaptic Plasticity Mechanisms and Learning

5. Synthesizing Real-Time Neuromorphic Systems

6. Conclusions

# 2. Hardware Models of Spiking Neurons:
## Neuron Excitability



axon of other neuron

synapse    dendrite

nucleus

axon

axon terminals

soma

(a)

Neurons information processing steps:

- **Synapses**: Connection between neurons
- **Dendrites**: Receive inputs
- **Cell body**: sums currents from dentures
- **Axon**: sends to action potential

How are action potentials generated given the current flowing into the soma (cell body) from dendrites and synapses?

(b)

## Biophysical description

Zooming in a patch of memory:



Biophysics of the membrane as an electrical circuit

$I(t)$: Current of membrane
$V(t)$: Membrane potential (Difference in electrical potential between inside and outside of the cell.)
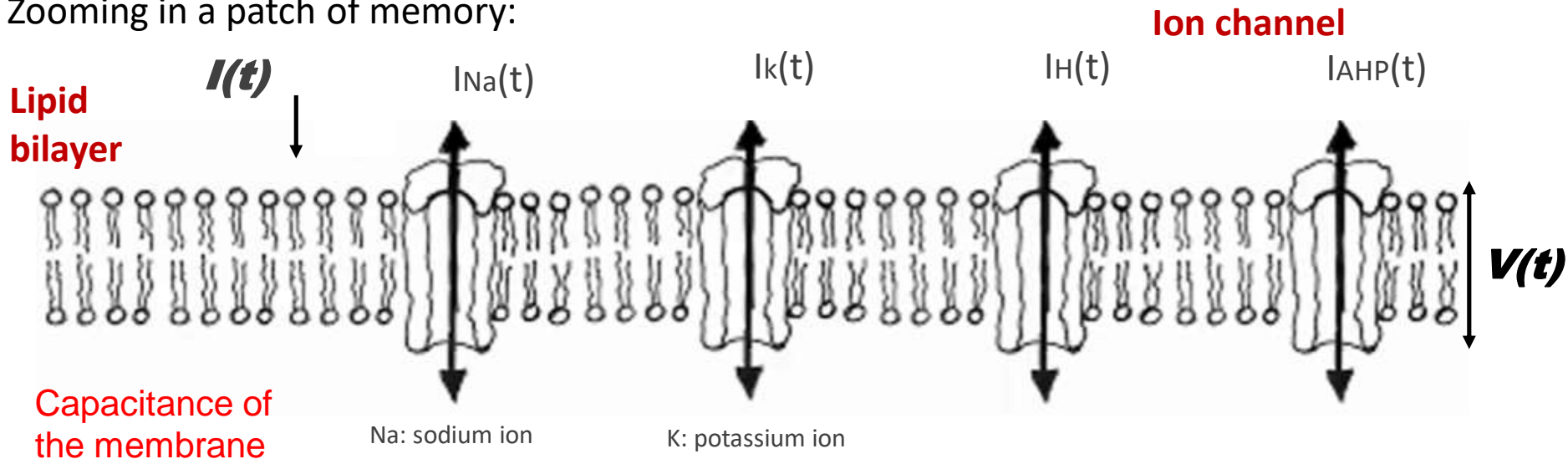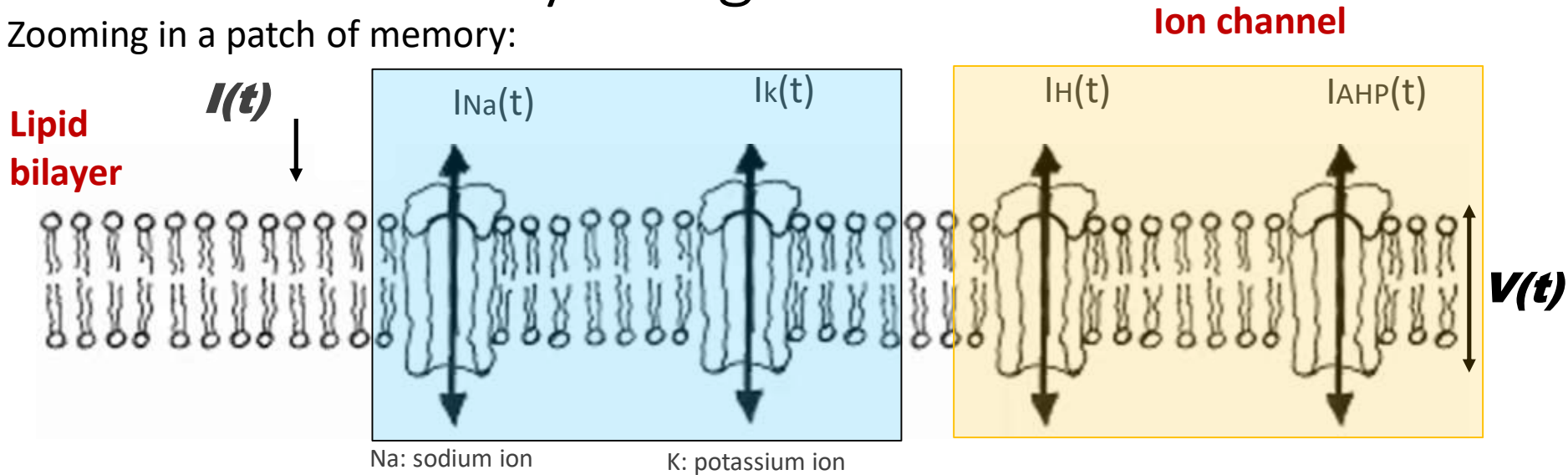
$C$: Capacitance of the membrane
$gL$: Conductance of the membrane
$EL$: Equilibrium potential of *Leak*

13

# 2. Hardware Models of Spiking Neurons:

## Biophysical description

Zooming in a patch of memory:

**Ion channel**

$I(t)$     $I_{Na}(t)$     $I_k(t)$     $I_H(t)$     $I_{AHP}(t)$

**Lipid bilayer**

$V(t)$

Na: sodium ion     K: potassium ion

Capacitance of the membrane

$$I = CdV/dt + g_L(V-E_L) + I_{Na} + I_{K} + I_H + I_{AHP}$$

Leak conductance of the membrane

Kirchhoff's current rule

Leak equilibrium potential

$I(t)$: Current of membrane
$V(t)$: Membrane potential (Difference in electrical potential between inside and outside of the cell.)

$C$: Capacitance of the membrane
$gL$: Conductance of the membrane
$EL$: Equilibrium potential of *Leak*

# 2. Hardware Models of Spiking Neurons:

## Leaky Integrate-and-Fire

Zooming in a patch of memory:

**Lipid bilayer**

**Ion channel**

$I(t)$

$I_{Na}(t)$

$I_k(t)$

$I_H(t)$

$I_{AHP}(t)$

$V(t)$

Na: sodium ion

K: potassium ion

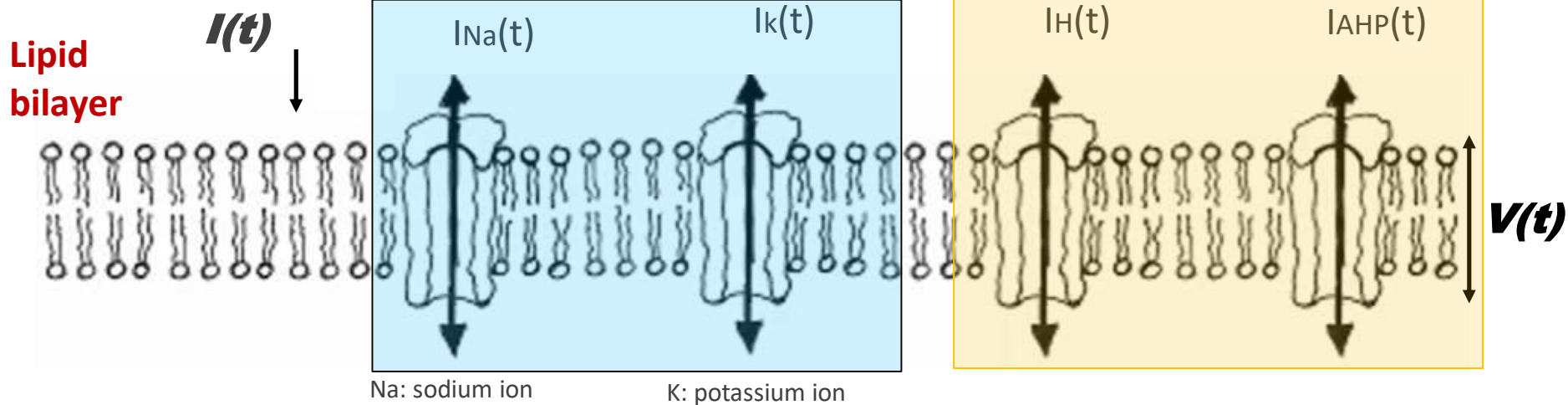$$I = CdV/dt + g_L(V-E_L) + I_{Na} + I_{K} + I_H + I_{AHP}$$

Replace by a threshold for spike emission follows by a reset to a fixed value/potential.

Ignore action on ion channels for now.

# 2. Hardware Models of Spiking Neurons:

## Leaky Integrate-and-Fire

Zooming in a patch of memory:

**Ion channel**

**Lipid bilayer**

$I(t)$

$I_{Na}(t)$        $I_k(t)$        $I_H(t)$        $I_{AHP}(t)$

$V(t)$

Na: sodium ion        K: potassium ion

$$I = CdV/dt + g_L(V-E_L) + I_{Na} + I_K + I_H + I_{AHP}$$

$$C_m dV/dt = -g_L(V-E_L) + I$$

If $V(t) = V_{th}$ then $V(t+\Delta) = E_L$

16

# Integrate-and-fire Model



Spike emission

threshold

$\vartheta$

Spike reception: EPSP

$\varepsilon\left(t - t_j^f\right)$ reset

$Ui/Vi$

$$\tau \cdot \frac{d}{dt} \, \mathbf{V}_i \ = \text{-V(t)} + \ \mathcal{R}I(t) \qquad \text{linear}$$

$$V(t) \ = \vartheta \Rightarrow \ \text{Fire+reset} \quad \text{threshold}$$

More details on ''Spiking Neuron Models'',
Cambridge press, 2002

# Spiking Neuron Model

## Spike Response Model



spike emission

$\eta\left(t - t_i^{\wedge}\right)$

i

state of neuron i

$u_i$

$\vartheta$

Spike reception: EPSP

$\varepsilon\left(t - t_j^f\right)$

Spike reception: EPSP

$\varepsilon\left(t - t_j^f\right)$

Spike emission: AP

$\eta\left(t - t_i^{\wedge}\right)$

reset of the membrane
potential (action potential)

EPSP: Excitatory
postsynaptic potentials

$$u_i(t) = \eta\left(t - t_i^{\wedge}\right) + \sum_j \sum_f w_{ij}\, \varepsilon\left(t - t_j^f\right)$$

$$u_i(t) = \vartheta \Rightarrow \text{Firing:} \quad t_i^{\wedge} = t$$

Spike Resonse Model (SRM), Gerstner (1996)

# 2. Hardware Models of Spiking Neurons:
## Spike Coding Schemes



Fig. 2.2: Time to first spike



Fig. 2.3: Inter-spike-interval



Fig. 2.4: Phase coding



Fig. 2.5: Rank order

# 2. Hardware Models of Spiking Neurons:
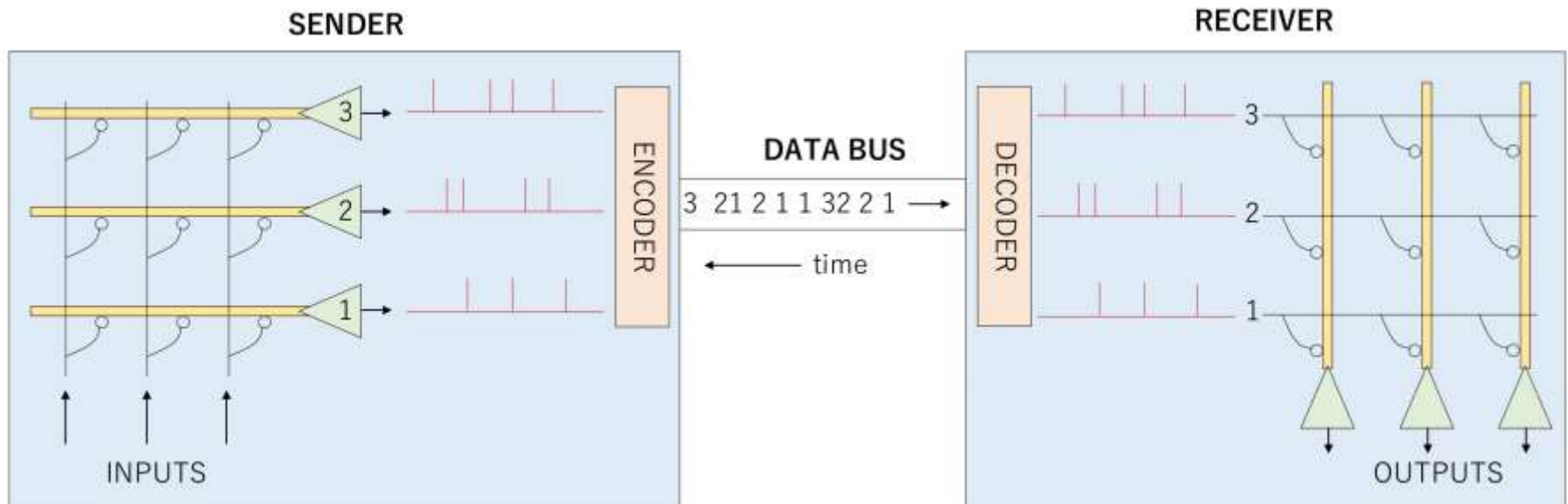
## Neurons Communication Scheme



Fig. 2.15: AER (Address Event Representation) protocol

# Lecture Contents

1. Neuromorphic Computing Approaches

2. Hardware Models of Spiking Neurons

3. <span style="color:red">Synaptic Dynamics</span>

4. Synaptic Plasticity Mechanisms and Learning

5. Synthesizing Real-Time Neuromorphic Systems

6. Conclusions

# 3. Synaptic Dynamics:
## Complex Structure of a Neural Network



Fig. Complex Structure of a Neural Network [M.Bertrand,2015].

➢ A typical neural network has four main regions: The cell body, the **dendrites**, The **axon**, and the **presynaptic terminals**.

➢ Each region has a distinct role in the generation of signals and the communication between neurons.

➢ Neurons can communicate through electrical synapses or chemical synapses alone or via both types of interactions.

# 3. Synaptic Dynamics:
## What is Synaptic Dynamics?

- Connections between neurons are not static, but change in amplitude and timing.

- Synaptic dynamics is the time-dependent changes in synaptic currents that change the strength of coupling between neurons.

- Both presynaptic and postsynaptic contribute to the changes of synaptic currents.
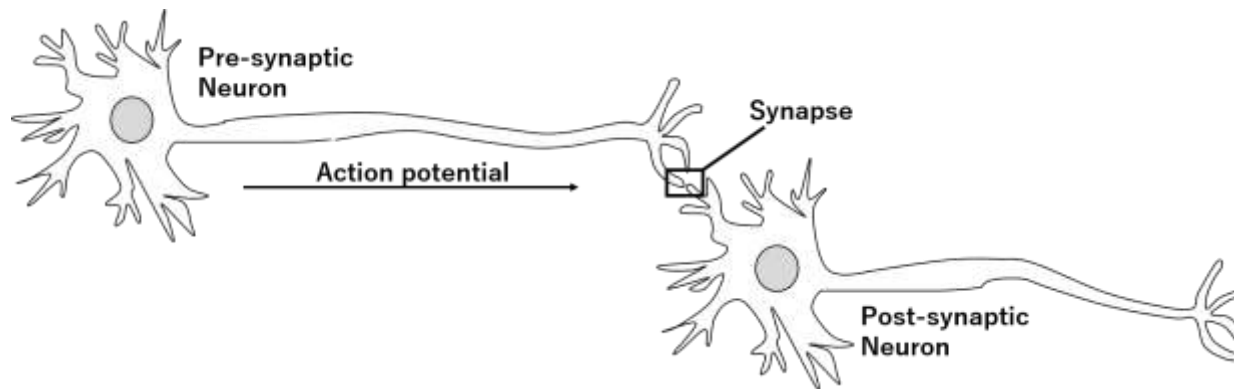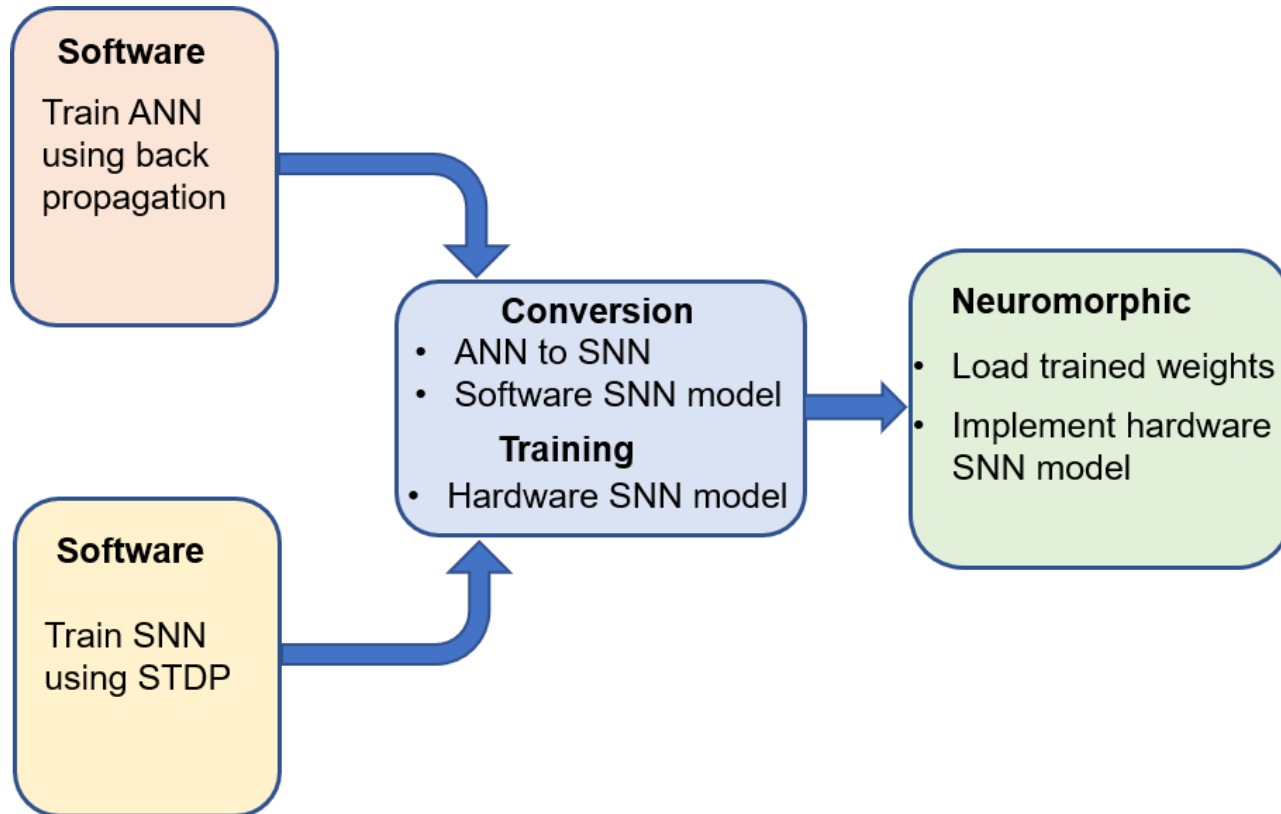
- Synaptic dynamics realizes **adaptive learning**.



Fig. 2.1: Two neurons communicating via a synapse.

# Lecture Contents

1.  Neuromorphic Computing Approaches

2.  Hardware Models of Spiking Neurons

3.  Synaptic Dynamics

4.  Synaptic Plasticity Mechanisms and Learning

5.  Synthesizing Real-Time Neuromorphic Systems

6.  Conclusions

# 4. Synaptic Plasticity Mechanisms & Learning: Learning Methods

- Spiking neural network (SNN) processes and communicates sparse binary signals (spikes) in a highly parallel and event-driven manner.

- The learning phase (minimizes a particular cost (loss)), is a complex process of acquiring the parameters to output the correct inference results.

- The cost function optimization is performed with a <u>gradient-descent-based</u> optimization or other classical optimization methods (i.e., genetic algorithm).

- There are various training/learning algorithms for SNNs:
  - ➢ Unsupervised Spike-timing-dependent plasticity (STDP)
  - ➢ ANN to SNN conversion

25

# 4. Synaptic Plasticity Mechanisms & Learning:
## Learning Methods



**Software**
Train ANN using back propagation

**Software**
Train SNN using STDP

**Conversion**
- ANN to SNN
- Software SNN model

**Training**
- Hardware SNN model

**Neuromorphic**
- Load trained weights
- Implement hardware SNN model

Neuromorphic Learning Framework

# 4. Synaptic Plasticity Mechanisms & Learning: Spike-timing-dependent plasticity (STDP)

- Adjusts the strength of connections (synapses) between neurons in the brain.
  - ✓ Adjusts the connection strengths based on the relative timing of a particular neuron's output and input action potentials.

$$\Delta w = \begin{cases} \Delta w^+ = A^+ e^{\left(\frac{-\Delta t}{\tau_+}\right)}, & \text{if } \Delta t > 0 \\ \Delta w^- = -A^- e^{\left(\frac{\Delta t}{\tau_-}\right)}, & \text{if } \Delta t \leq 0 \end{cases}$$

Where $\Delta w$ is the change in synaptic weight. If a presynaptic spike arrives the postsynaptic neuron within a time window $\tau_+$ before the postsynaptic spike, the synaptic weight increases $\Delta w^+$, but if it arrives within a time window $\tau_-$, after the postsynaptic spike, the synaptic weight decreases $\Delta w^-$. $\Delta t$ is the time difference between the presynaptic and postsynaptic spike which is expressed as $\Delta t = t_{post} - t_{pre}$, while $A^+$ and $A^-$ are potentiation and depression amplitude parameters respectively.

# 4. Synaptic Plasticity Mechanisms & Learning: Spike-timing-dependent plasticity (STDP)
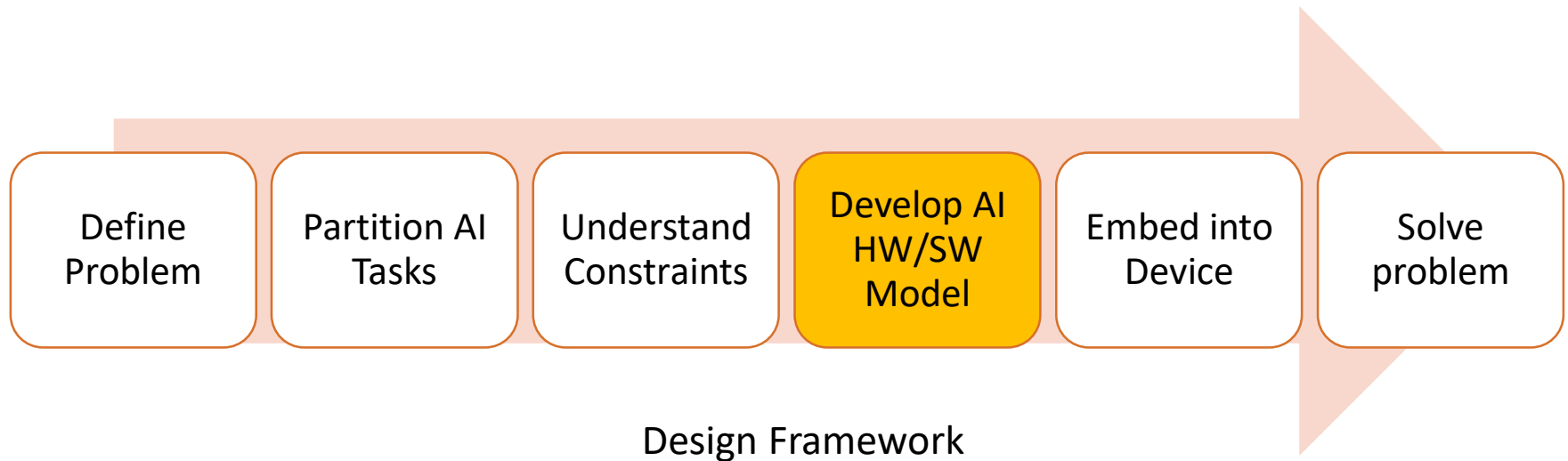


Fig. 1.5: STDP Architecture.

- The STDP unit Follows the *spike* or *pulse* model assumption for cortical neurons where information lies in spike timings, and not in spike shapes.
- 16 presynaptic traces are required to initiate the learning process. The PWU mechanism enables fast parallel on-chip learning.
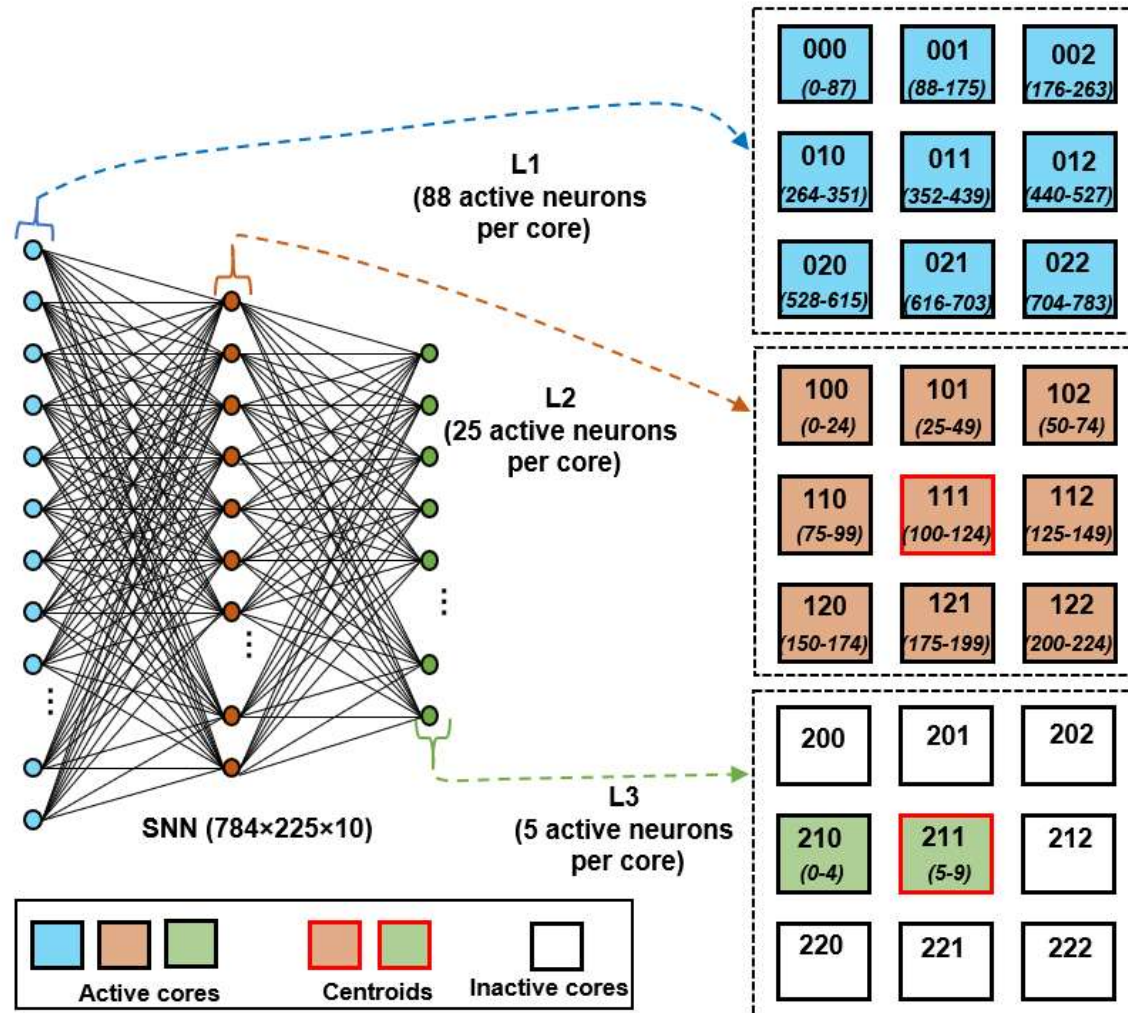
# Lecture Contents

1. Neuromorphic Computing Approaches

2. Hardware Models of Spiking Neurons

3. Synaptic Dynamics

4. Synaptic Plasticity Mechanisms and Learning

5. Synthesizing Real-Time Neuromorphic Systems

6. Conclusions

# 5. Synthesizing Real-Time Neuromorphic Systems:
# A framework for a Real Neurocomputing Design

| Define Problem | Partition AI Tasks | Understand Constraints | Develop AI HW/SW Model | Embed into Device | Solve problem |

Design Framework

Define Problem→ Partition AI Tasks → Understand Constraints → Develop AI HW/SW  Model -> Embed into Device -> Solve problem

# 5. Synthesizing Real-Time Neuromorphic Systems: Application mapping



Application mapping example on a $3 \times 3 \times 3$ Neuromorphic Chip

# 5. Synthesizing Real-Time Neuromorphic Systems: Connecting Neuromorphic Chips



Pixel

Y

X

Source chip

AER

Y

X

Arbor

Target chip

## Inside the Pixel



Inside the Pixel

# 5. Synthesizing Real-Time Neuromorphic Systems: Using Crossbars



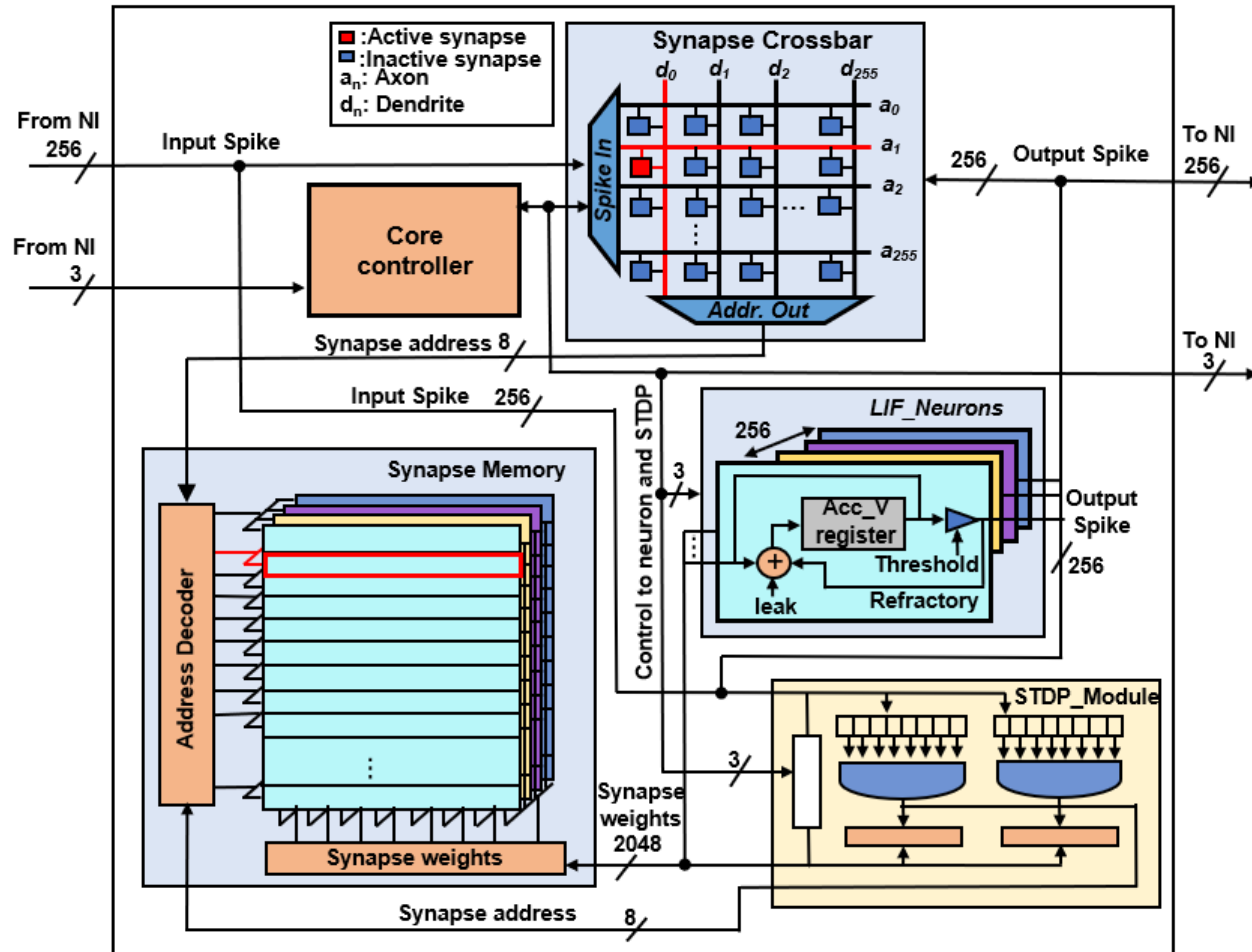Using Crossbars

# 5. Synthesizing Real-Time Neuromorphic Systems:
## Using Crossbars



Using Crossbars

# 5. Synthesizing Real-Time Neuromorphic Systems:
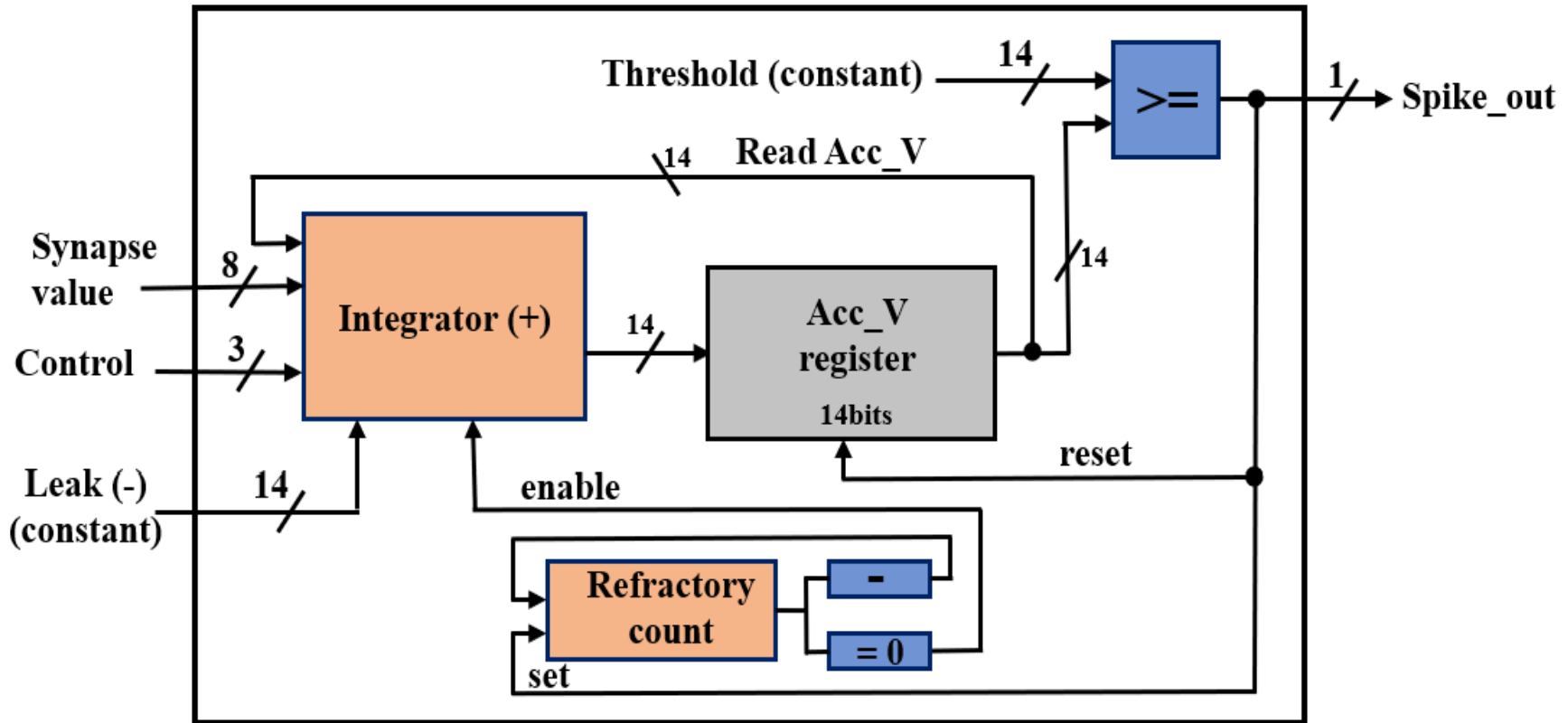## Using Crossbars

# 5. Synthesizing Real-Time Neuromorphic Systems:
## Spiking Neuro-Processing Core



Architecture of Spiking Neuro-Processing Core.

# 5. Synthesizing Real-Time Neuromorphic Systems:
## LIF Neuron Module



Architecture of LIF Neuron

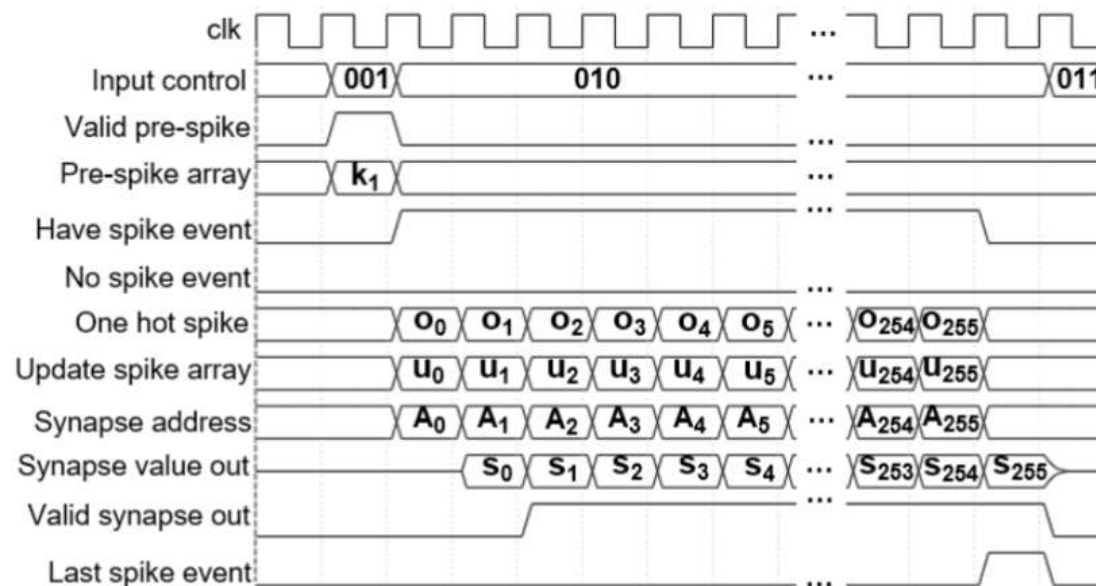# 5. Synthesizing Real-Time Neuromorphic Systems:
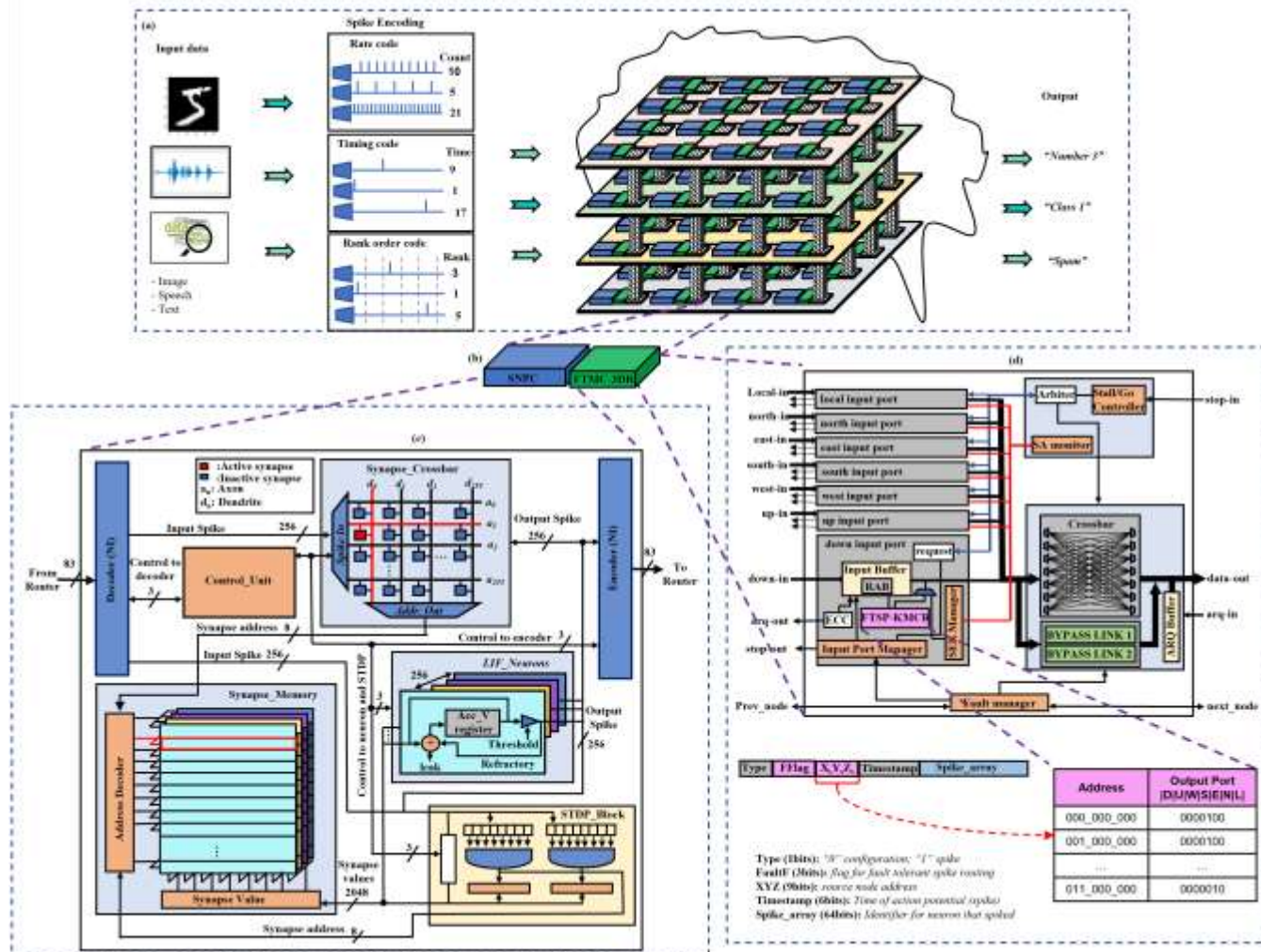## LIF Neuron Module



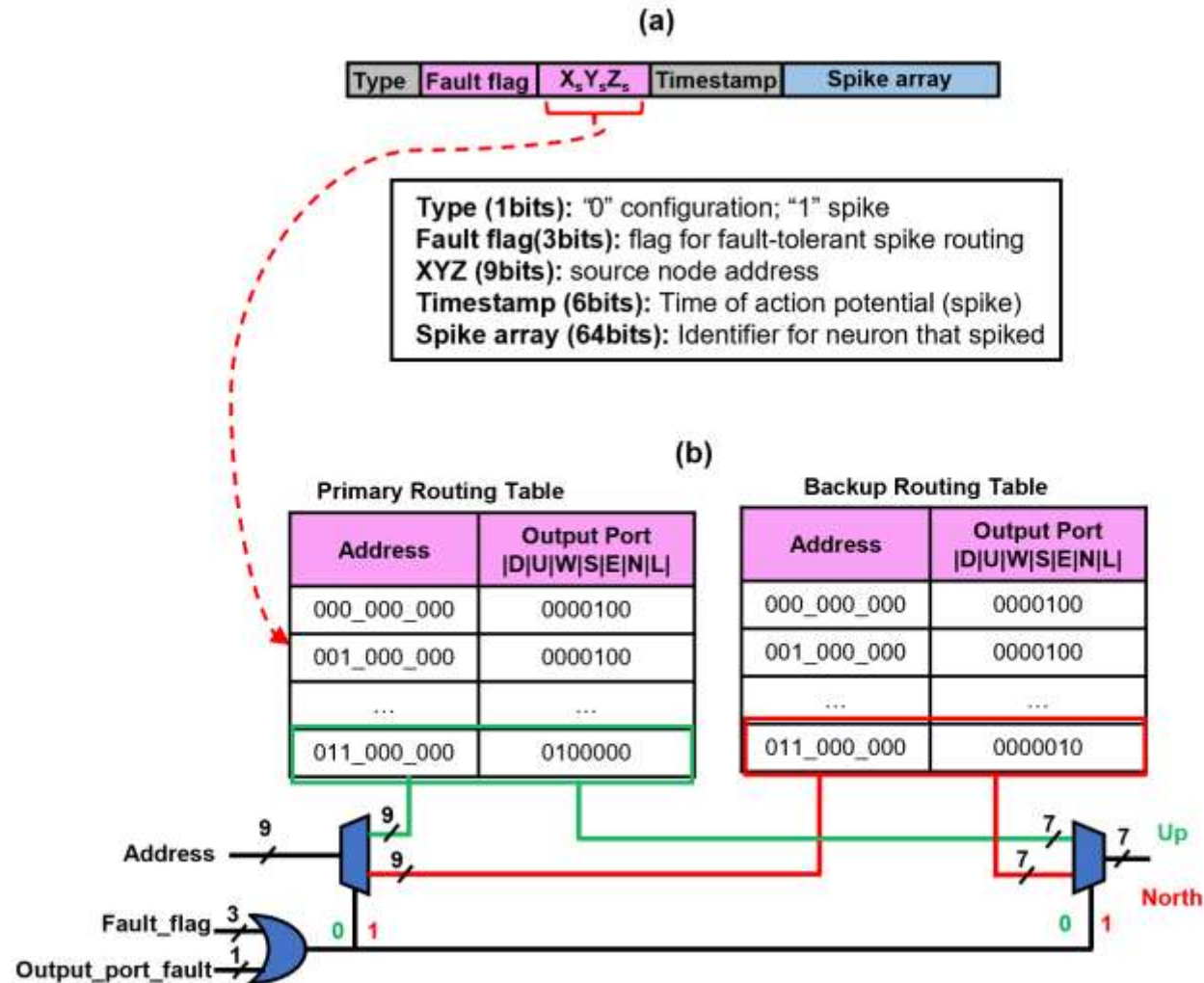Illustration of neuron update operation at the crossbar

1. An input presynaptic spike array is stored and checked for spike events. If present, the **Have spike** event signal becomes high. Afterwards, the one hot operation to get the synapse address begins, updating the one hot spike array for every spike event: from **O0 to O255**.

2. The stored presynaptic spike array is also updated after each spike event is processed: from **U0 to U255.**

3. The synapse address is then used to fetch the synapse values from the synapse memory, and sent to the postsynaptic neurons.

4. When the last spike event in the array has been processed, the crossbar sends a signal to the control unit signaling that all spike events have been processed

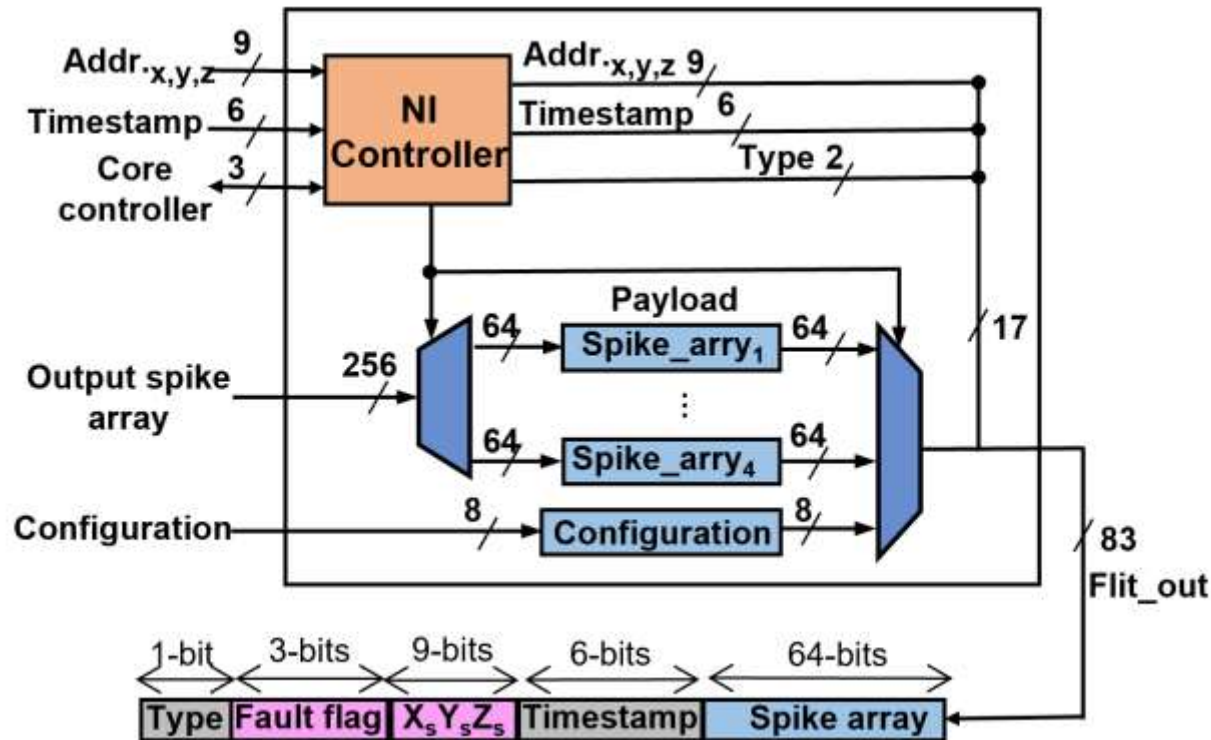# 5. Synthesizing Real-Time Neuromorphic Systems: NASH Architecture



Organization of the NASH Neuromorphic Chip

# 5. Synthesizing Real-Time Neuromorphic Systems: Spiking Neuron Packet Format
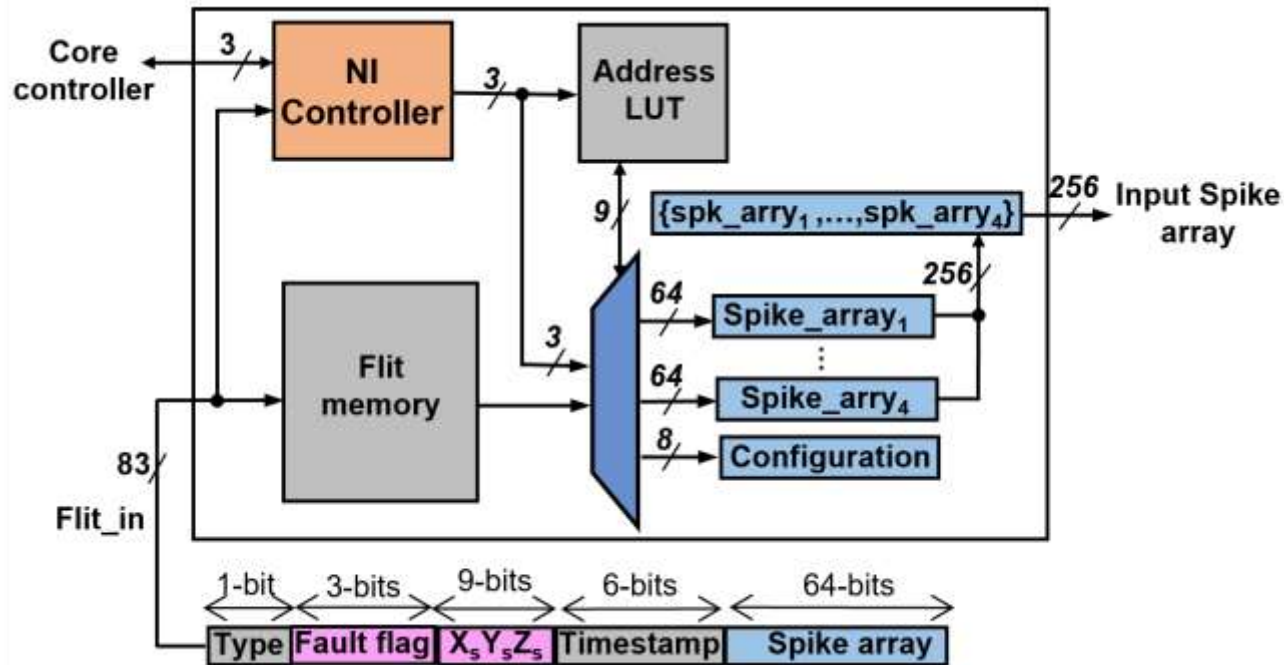


Spiking neuron packet format

Network Interface: Encoder.

Operations of the encoder can be summarized in the following steps:
- Receive output spikes from local SNPC and packet into flits.
- After packeting, send flit to local router

# Network Interface (1/2): Decoder



Network Interface: Decoder.

Operations of the decoder can be summarized in the following steps:
- Receive spike packets from local router and unpack.
- Forward the spikes to the local SNPC as presynaptic spike train.

# Lecture Contents

1.  Neuromorphic Computing Approaches

2.  Hardware Models of Spiking Neurons

3.  Synaptic Dynamics

4.  Synaptic Plasticity Mechanisms and Learning

5.  Synthesizing Real-Time Neuromorphic Systems

6.  Conclusions

# Conclusions

- Neuromorphic Computing is the use of hardware (VLSI) to simulate the biological architecture of the human nervous system (brain, complex network of nerves, etc.),

- Spiking Neural Network:
  - More analogous to the brain, communicating via spikes in a sparse event-driven manner.
  - Exploits spike sparsity to achieve low power.

- Synaptic dynamics is the time-dependent changes in synaptic currents that change the strength of coupling between neurons.

- There are various training/learning algorithms for SNNs:
  - Unsupervised Spike-timing-dependent plasticity (STDP)
  - ANN to SNN conversion

- Synthesizing a Neuromorphic System:
  - Define Problem→ Partition AI Tasks → Understand Constraints → Develop AI HW/SW Model → Embed into Device → Solve the Problem