# Introducing Loihi

Mike Davies
Director, Neuromorphic Computing Lab | Intel Labs

# Motivation: The Case for Neuromorphic Computing

## Problem Statement:

**Emerging computing workloads demand *intelligent behaviors* that we do not know how to deliver *efficiently* with today's algorithms and computing architectures.**
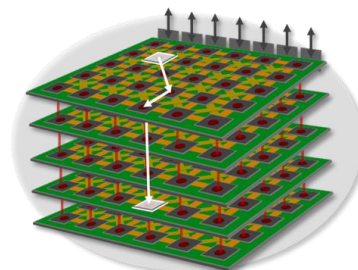
### Examples:

- Online and lifelong learning
- Learning without cloud assistance
- Learning with sparse supervision
- Understanding spatiotemporal data
- Probabilistic inference and learning
- Sparse coding/optimization
- Nonlinear adaptive control
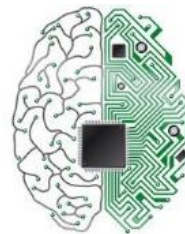- Pattern matching with high occlusion
- SLAM and path planning

## Potential Future Product Applications
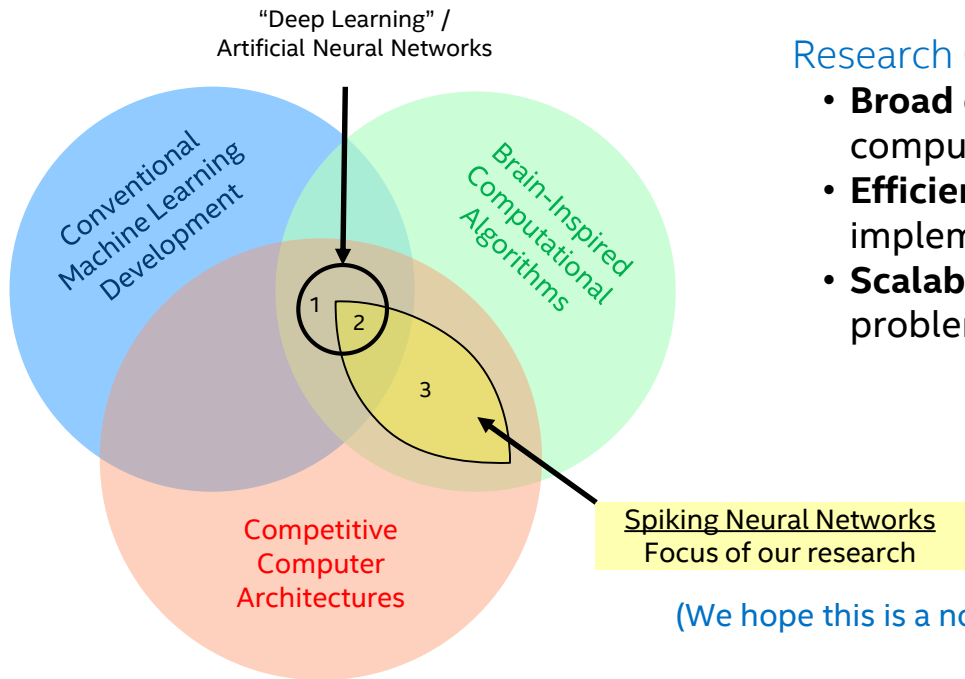
Robotics

HPC Systems

Neuroprosthetics

Smart Glasses

# Solution Exploration Space



"Deep Learning" /
Artificial Neural Networks

Conventional Machine Learning Development

Brain-Inspired Computational Algorithms

Competitive Computer Architectures

Spiking Neural Networks
Focus of our research

Research Goals:
- **Broad class** of brain-inspired computation
- **Efficient** hardware implementations
- **Scalable** from small to large problems and systems

(We hope this is a non-empty class!)

# The Engineering Perspective

- Nature has come up with something amazing.  Let's copy it...

- Not so simple – very different design regimes

- Yet objectives and constraints are largely the same...
    - Energy minimization
    - Fast response time
    - Cheap to produce

Need to understand and apply the basic principles, *adapting for differences*

Status today:

| | Nature | Silicon | Ratio |
|---|---|---|---|
| Neuron density[1] | 100k/mm$^2$ | 5k/mm$^2$ | 20x |
| Synaptic area[1] | 0.001 um$^2$ | 0.4 um$^{2[2]}$ | 400x |
| Synaptic Op Energy | ~2 fJ | ~4 pJ | 2000x |

[1] Planar neocortex   [2] ~5b SRAM

But...

| | Nature | Silicon | Ratio |
|---|---|---|---|
| Max firing rate | 100 Hz | 1 GHz | 10,000,000x |
| Synaptic error rate | 75% | 0% | ∞ |

| Nature | Silicon |
|---|---|
| Autonomous self-assembly | Fabricated manufacturing |
| Per-instance variability desired | Variability causes brittle failures |
| plasticity over lifetime | Must support rapid reprogramming |
| terministic operation | Deterministic operation desired |

# Are Spiking Architectures Efficient?

# One Compelling Example: LASSO Sparse Coding

## Problem

$$\min_{z} \frac{1}{2}\|x - Dz\|_2^2 + \lambda\|z\|_1$$

Input

Reconstruction

Sparse regularization
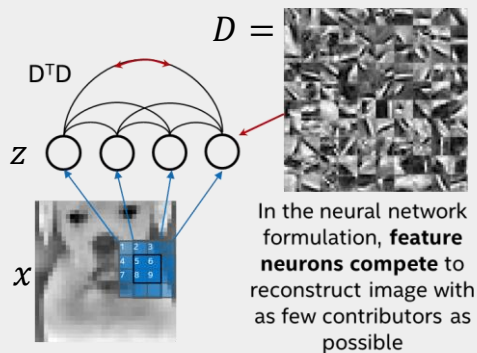
## Implementation

$D =$

D$^T$D

$z$

$x$

In the neural network formulation, **feature neurons compete** to reconstruct image with as few contributors as possible

Tang et al, arxiv: 1705:05475

LASSO Optimization Using the *Spiking Locally Competitive Algorithm*

Inhibition

$-(\boldsymbol{d}_i^T \cdot \boldsymbol{d}_j)z_j$

$z_i$       ....   $z_j$

$\boldsymbol{d_i} \cdot \boldsymbol{x}$

Excitation

$x_1$       $x_2$

1672 spikes
(avg 0.052 spike/neuron)

5th iteration

S-LCA
FISTA

Normalized objective

Execution time (second)

*both S-LCA and FISTA running on a Xeon*

Neuromorphic algorithm rapidly finds a near-optimal solution

# Spiking LCA dynamics on a Loihi predecessor



Original

Reconstruction

Spikes

LASSO Objective Over Time

Intense but very brief period of competition

Much faster convergence on a neuromorphic architecture

# What this gives us… a baseline SNN architecture



Local Synaptic Routing

Synaptic Accumulation

Output Axon Routing

2D Mesh
Packetized spikes
High fanout required
Low overhead synchronization

Neuron Model (IF)

# But how to scale to large LCA problems?

LCA is an all-to-all network...



Just 1000 feature neurons requires $1000^2$ = 1M synapses

# Answer: Patch-based Connectivity Reuse

## Analogous to the "convolution" in ConvNets



Generalized Hierarchical Connectivity Example

Conventional 1D convolution example

w/ Lateral inhibition

# Sparse Coding Results: N1 vs Atom CPU
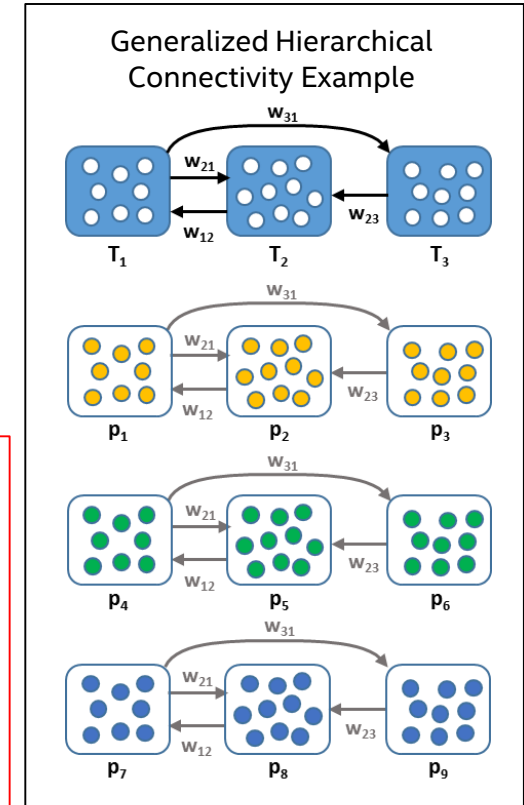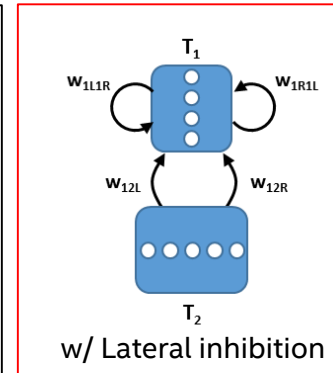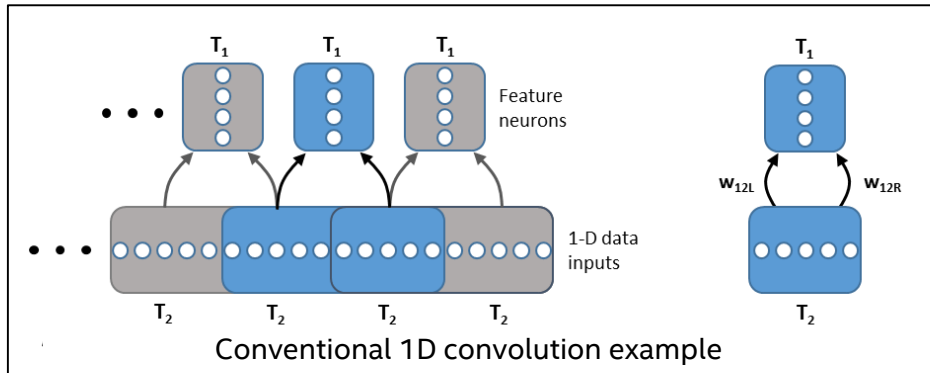
## Time to Solution Comparison



48x speed-up

Time (ms) vs Number of Unknowns

## Energy to Solution Comparison



118x lower

Energy (mJ) vs Number of Unknowns

## N1 Advantage in Energy-Delay-Product



>5000x better

Atom:N1 EDP Ratio

Comparison of sparse coding on N1 versus the FISTA* LASSO solver on an Atom CPU**

\* Best conventional LASSO solver (LARS also evaluated)

\*\* Iso-process, roughly iso-area (6-10mm$^2$)
   PTPX-based measurements

● Atom (FISTA)

● N1

# Neuromorphic Core Architecture



All synaptic connections pooled
128KB shared memory

Sparse, dense, and hierarchical
Synaptic mapping representations

Synaptic delays

Synaptic eligibility traces

Flexible 3-tuple synaptic variables
(1-9b weight, 0-6b delay, 0-8b tag)

Graded "reward spikes"

**Flexible synaptic plasticity with
microcode-programmable rules**

Sum-of-products rule semantics

Filtered spike train traces

Plasticity rules target any synaptic variable

Discrete time LIF neuron model (CUBA)

Multi-compartment dendritic trees
up to 1K compartments

Intrinsic excitability homeostasis

Random noise sources

Shared output routing table
4K axon routes

Axon delays
Refractory delays (+ random)

# Basic Core Operation (Non-Learning)

(Time multiplexing illustrated unrolled in space)



SYNAPSE

DENDRITE

$(W_i, D_i)$

AxonID

WeightSum   idx

CFG[idx]   STATE[idx]

AxonID$_{j+1}$

AxonID$_j$

T+1   T+2   T+3   T+4   T

Input spike routing
Tables (very complex)

Synaptic delay handling

Neuron model

Output spike routing
tables (simpler)
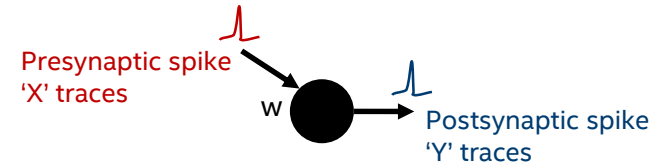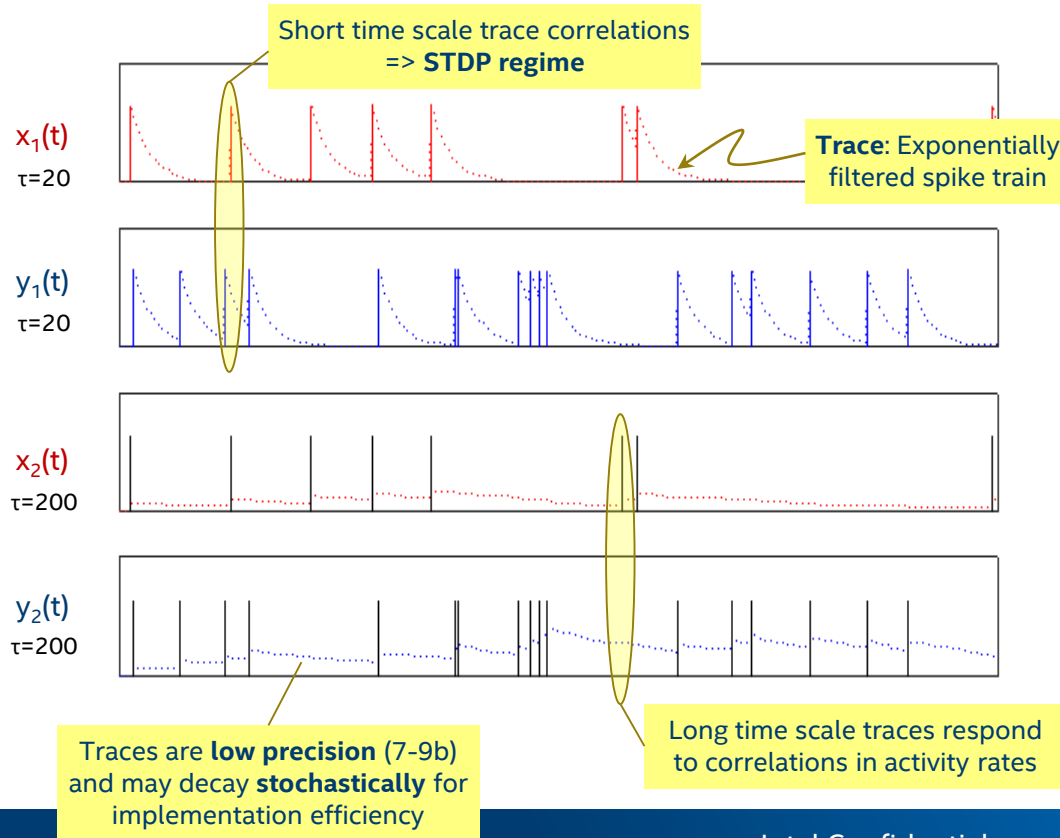
# Learning with Synaptic Plasticity

- **Local learning rules** – essential property for efficient scalability
  *Compatible with biological plausibility*

- Should be derived by **optimizing an emergent statistical objective**
  *Too much directionless experimentation otherwise*

- Plasticity on **wide range of time scales** is needed
  *Delayed reward/punishment responses, eligibility traces*

x  y

$W_{x,y}$

Supervision signal

z

$E = o - s$

Learning rules for weight $W_{x,y}$ may *only* access presynaptic state x and postsynaptic state y

However *reward spikes* may be used to distribute graded reward/punishment values to a particular set of axon fanouts

# Trace-Based Programmable Learning



$x_1(t)$
$\tau=20$

Short time scale trace correlations
=> **STDP regime**

**Trace**: Exponentially filtered spike train

$y_1(t)$
$\tau=20$

$x_2(t)$
$\tau=200$

$y_2(t)$
$\tau=200$

Traces are **low precision** (7-9b) and may decay **stochastically** for implementation efficiency

Long time scale traces respond to correlations in activity rates

Presynaptic spike 'X' traces

w

Postsynaptic spike 'Y' traces

Weight, Delay, and Tag learning rules programmed as **sum-of-product equations**

$$w' = w + \sum_{i=1}^{N_P} S_i \prod_{j=1}^{n_i} (V_{i,j} + C_{i,j})$$

Synaptic Variables
Wgt, Delay, Tag
(variable precision)

Variable Dependencies
$X_0$, $Y_0$, $X_1$, $Y_1$, $X_2$, $Y_2$,
Wgt, Delay, Tag, etc.

# Learning Rule Examples

**Pairwise STDP:**

$$W(t+1) = W(t) - A_- x_0(t) y_1(t) + A_+ x_1(t) y_0(t)$$

**Triplet STDP with heterosynaptic decay:**

$$W(t+1) = W(t) - A_- x_0(t) y_1(t) + A_+ x_1(t) y_0(t) y_2(t) - B \cdot W(t) \cdot y_3(t)$$

**Delay STDP:**

$$D(t+1) = D(t) - A_- x_0(t)(127 - y_1(t)) + A_+(127 - x_1(t)) y_0(t)$$

# Two-variable Learning Rule Examples

**Distal Reward with Synaptic Tags:**

$$T(t+1) = T(t) - A_- x_0(t)y_1(t) + A_+ x_1(t)y_0(t) - B \cdot T(t)$$
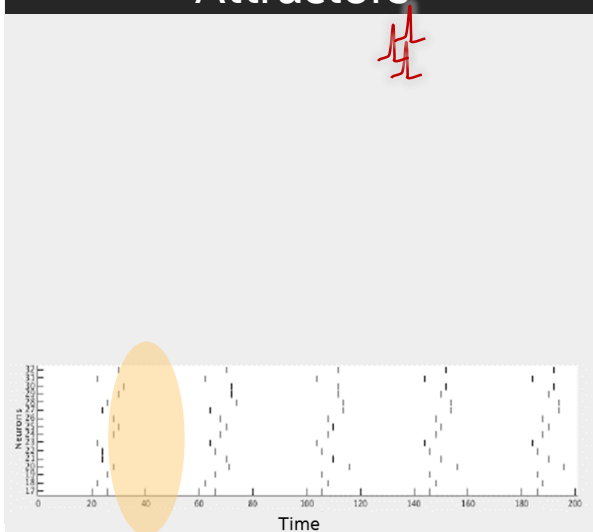
$$W(t+1) = W(t) + C \cdot r_1(t) \cdot T(t)$$

**STDP with dynamic weight consolidation:**

$$W(t+1) = W(t) - A_- x_0(t)y_1(t) + A_+ x_1(t)y_0(t)y_2(t) - B_1(W-T)y_3(t)y_0(t)$$

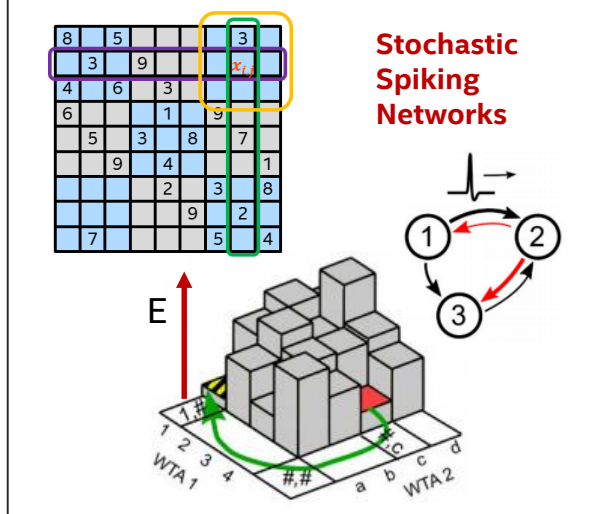$$T(t+1) = T(t) + \frac{1}{\tau_{cons}}(W-T) - B_2 T(w_\theta - T)(w_{max} - T)$$

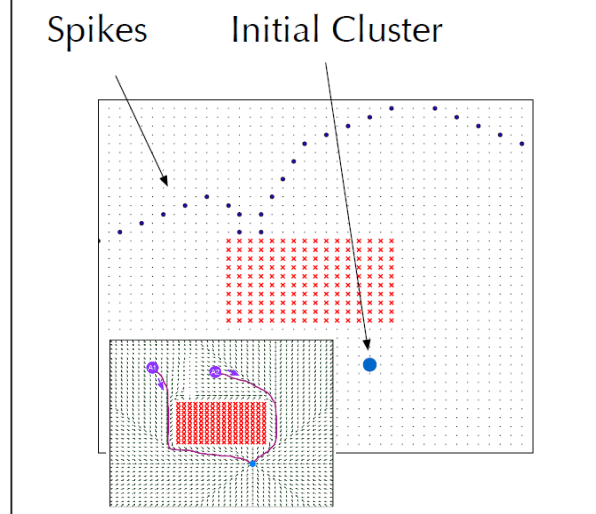# Example Novel Algorithms Supported by Loihi



Spatiotemporal Attractors
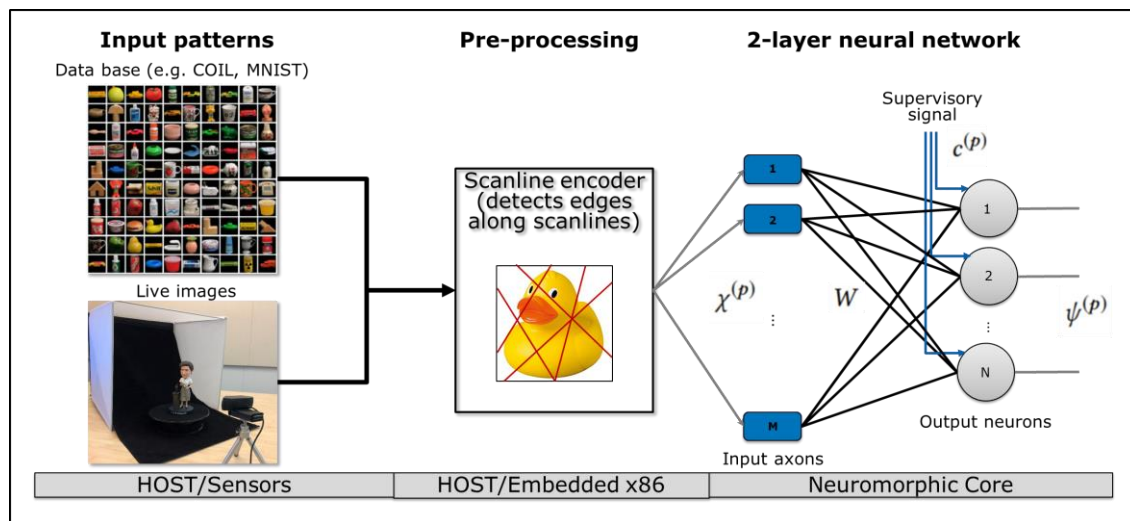
Artificial Olfaction



Constraint Satisfaction

**Stochastic Spiking Networks**
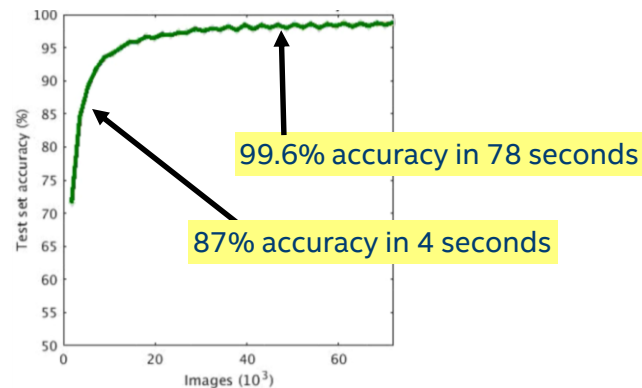
Sudoku



Graph Search

Spikes    Initial Cluster

Path Planning

# Our "Hello World" Application: Supervised Learning for Object Recognition



**Input patterns**
Data base (e.g. COIL, MNIST)

Live images

HOST/Sensors

**Pre-processing**

Scanline encoder (detects edges along scanlines)

HOST/Embedded x86

**2-layer neural network**

Supervisory signal $c^{(p)}$

$\chi^{(p)}$   $W$   $\psi^{(p)}$

Output neurons

Input axons

Neuromorphic Core

Performance on COIL20 data set



Test set accuracy (%)

Images ($10^3$)

99.6% accuracy in 78 seconds

87% accuracy in 4 seconds

|  | Training | Inference |
|---|---|---|
| Active energy per image (total) | 553 uJ | 128 uJ |
| Neuromorphic energy | 322 uJ | 13 uJ |
| Processing time per image | 7.5 ms | 1.8 ms |
| Chip power | 74 mW | 73 mW |
| Neuromorphic power | 43 mW | 7.4 mW |

| Resource Utilization | Count | Utilization |
|---|---|---|
| Neurons | 20 | 0.02% |
| Synapses | 38400 | 0.28% |
| SNN Cores | 1 | 078% |

S-STDP rule:

$$W_{i,j}(t) = W_{i,j}(t-1) + \eta \cdot \left(u_\kappa \cdot \delta_{i,C(p)} - y_{i,0}\right) \cdot x_{j,1}$$

# Up to the 10,000 foot view



**The Nx System Framework**
- Heterogeneous hierarchical parallel system
- Event-driven communication over channels
- Localized state
- Models describe *emergent behavior*

$$y^* = \operatorname*{argmin}_{y_i \geq 0} F(x, y)$$

**SNN specification**

**Snip**
(Sequential neural interfacing process)

**Spiking neuron**

**Modules w/ behavioral models**
**A, B:** Sequential processes conventionally coded and
        run on conventional CPUs
**NN**: Neural network module
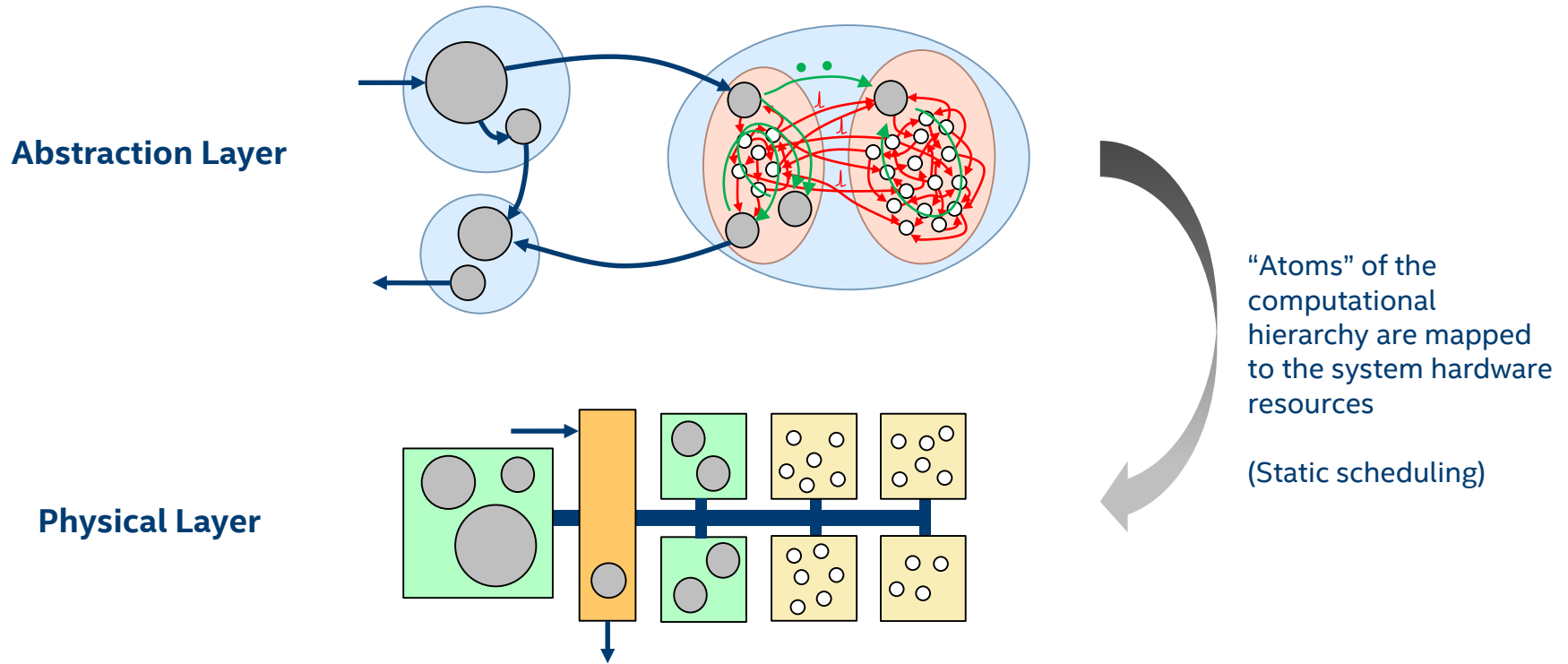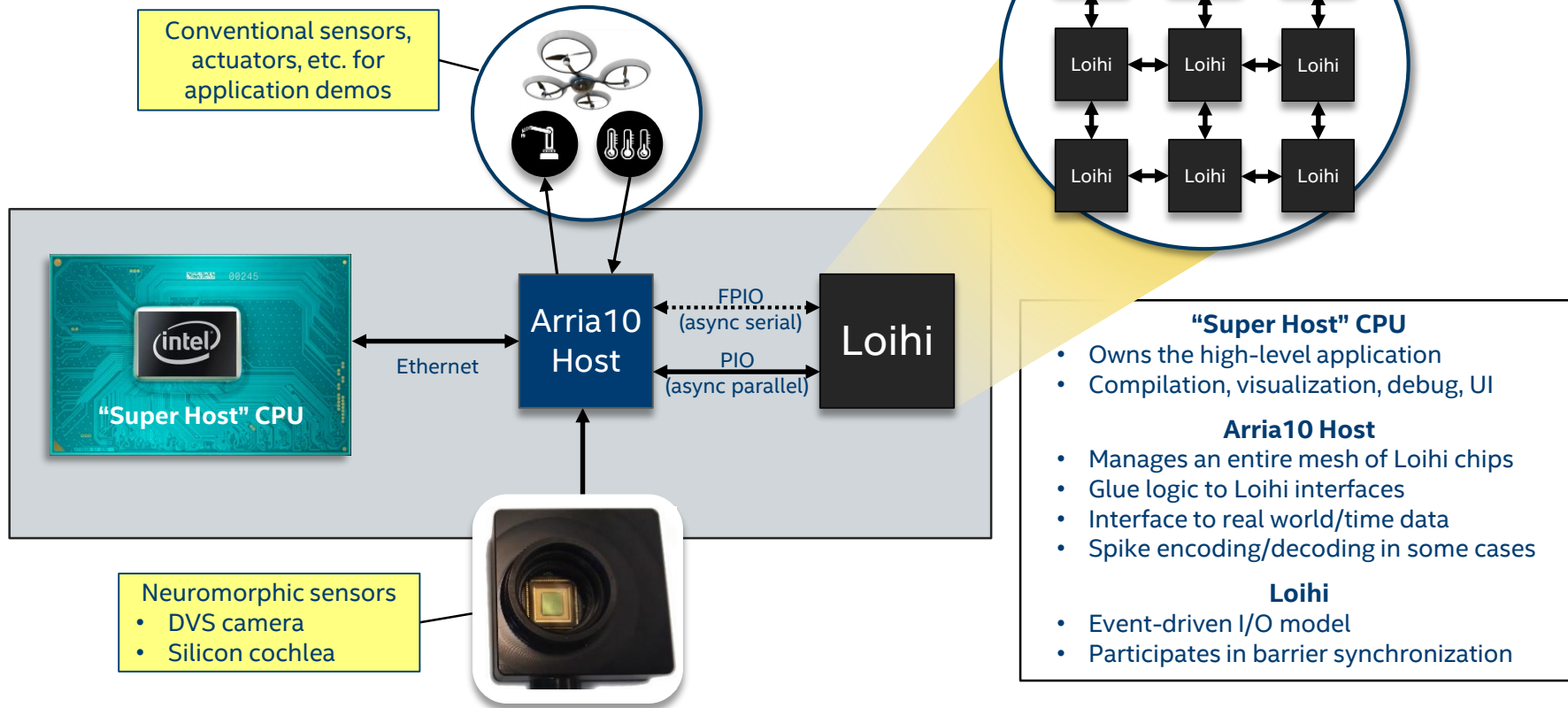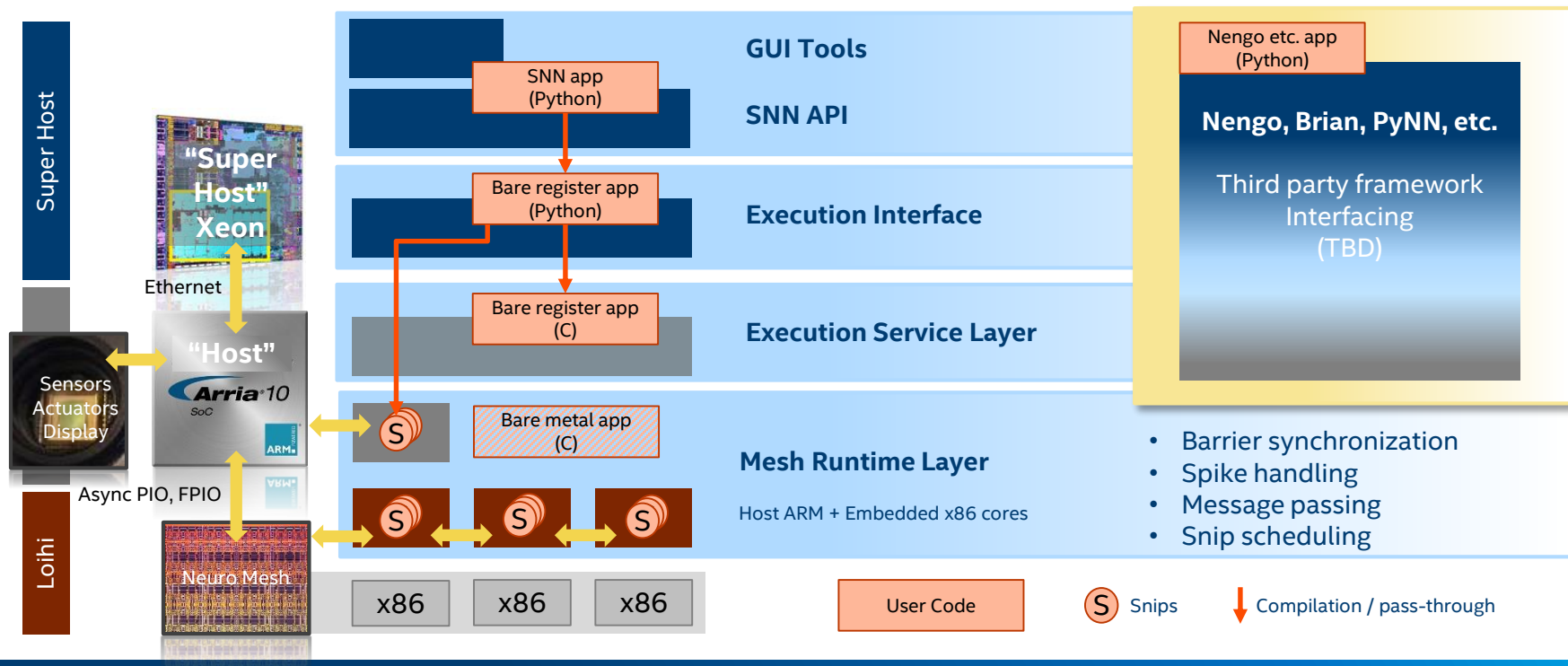- Hierarchically specified
- Mathematical behavioral model
- May include conventional helper code ("snips")

# Mapping to the Physical Layer



**Abstraction Layer**

**Physical Layer**

"Atoms" of the computational hierarchy are mapped to the system hardware resources

(Static scheduling)

# System Architecture Today



Conventional sensors, actuators, etc. for application demos

Multi-chip scalability

Neuromorphic sensors
- DVS camera
- Silicon cochlea

"Super Host" CPU

Ethernet

Arria10 Host

FPIO (async serial)

PIO (async parallel)

Loihi

**"Super Host" CPU**
- Owns the high-level application
- Compilation, visualization, debug, UI

**Arria10 Host**
- Manages an entire mesh of Loihi chips
- Glue logic to Loihi interfaces
- Interface to real world/time data
- Spike encoding/decoding in some cases

**Loihi**
- Event-driven I/O model
- Participates in barrier synchronization

# Current Software Development Kit
## (work in progress)

# Current Software Development Kit
## (work in progress)

# What's the right top layer of the SDK?



$$y^* = \underset{y_i \geq 0}{\text{argmin}}\, F(x, y)$$

**GUI Tools**

**"TBD" API**

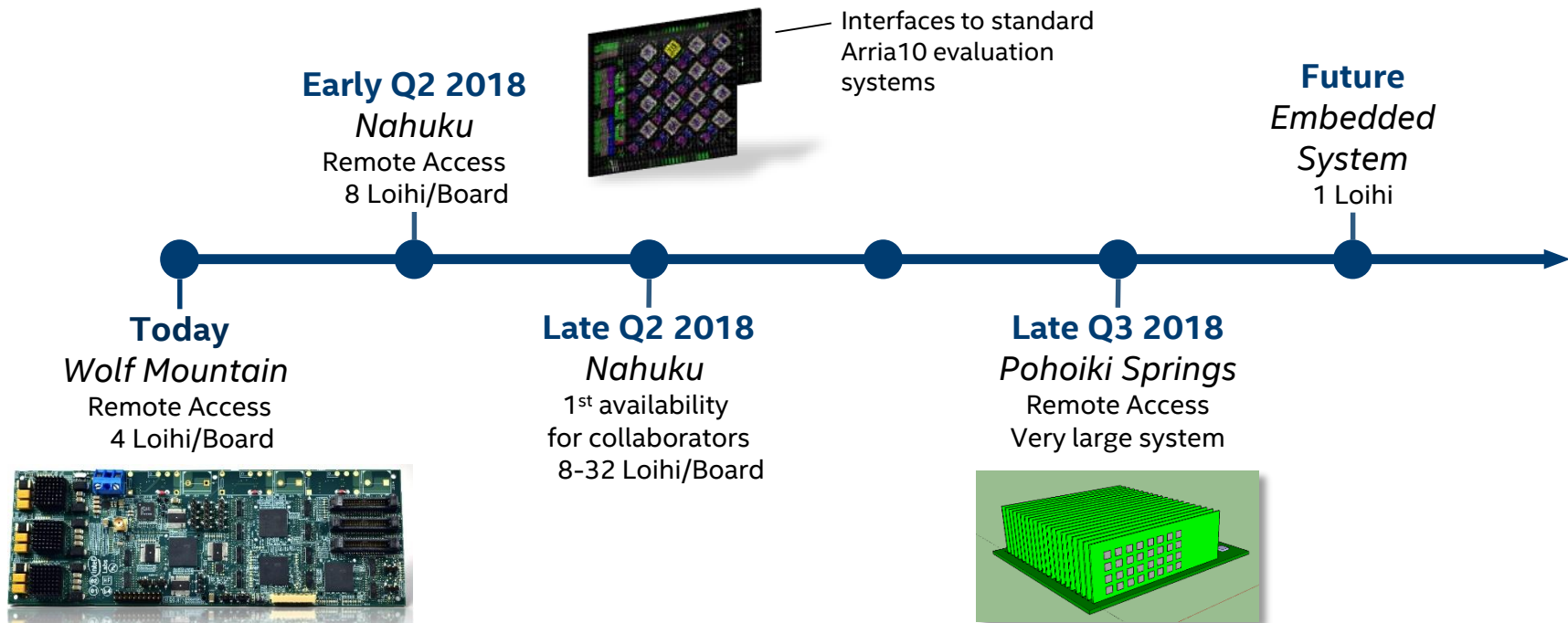**Not** TensorFlow / other DL frameworks
(wrong abstractions)

This is the unexplored frontier of
neuromorphic software research

SNN app
(Python)

**SNN API**

Bare register
app (Python)

**Execution Interface**

Bare register
app (C)

**Execution Service Layer**

Bare metal
app (C)

**Mesh Runtime Layer**

# Loihi Systems Outlook



Interfaces to standard Arria10 evaluation systems

**Early Q2 2018**
*Nahuku*
Remote Access
8 Loihi/Board

**Future**
*Embedded System*
1 Loihi

**Today**
*Wolf Mountain*
Remote Access
4 Loihi/Board

**Late Q2 2018**
*Nahuku*
1st availability
for collaborators
8-32 Loihi/Board

**Late Q3 2018**
*Pohoiki Springs*
Remote Access
Very large system

# Intel Neuromorphic Research Community

**RV1: Theory**
- Abstract and quantify features of neuroscience to the context of systems engineering
- Computational complexity frameworks

**RV2: Algorithms**
- Principled derivations of SNN dynamics, features, and learning rules.

**RV5: Sensors and Control**
- Sparse, event-driven I/O for SNN systems

**RV3: Applications**
- Applications of Loihi and future Intel neuromorphic silicon / FPGA designs
- Benchmarks and value analysis may itself be research.

**RV4: Programming Models**
- New paradigms for conceptualizing and specifying SNN/neuromorphic algorithms

Application Systems/SW

Neuromorphic Algorithms

Neuromorphic SDK

Sensors Actuators Display

Loihi / IA HW Platform

## We wish to engage with collaborators in academic, government, industry research groups

**INRC goals:**
- Demonstrate value of Loihi vs conventional solutions
- Share code, results, algorithms
- Motivate improvements for future silicon iterations

**What we offer to INRC collaborators**
- Remote access to Loihi systems, SDK, SW
- Loaned Loihi systems and bare chips (limited)
- Opportunity for limited funding (RFP available late March)

# Please Join Us! (at the right time)

**You:**
- Extensive experience with SNNs
- Extensive experience with other neuromorphic HW platforms

**Us:**
- Highly bandwidth limited

**Email inrc_interest@intel.com for more information**

Telluride 2018

Today

2018

2019

**You:**
- Vision for SNN application/algorithm research
- Can articulate the promise/value of project
- Can benchmark the result
- Interested in neuromorphic SW development

**Us:**
- More systems & documentation
- Complete SDK
- Scalable remote access system

**You:**
- Have a real-world problem not well solved now
- Prior SNN experience not necessary

**Us:**
- Mature, cross-framework SDK

**Community:**
- Critical mass, community forums, etc.
- Usable library of SDKs, tools, code, modules

# LEGAL INFORMATION