

Jeff Gehlhaar

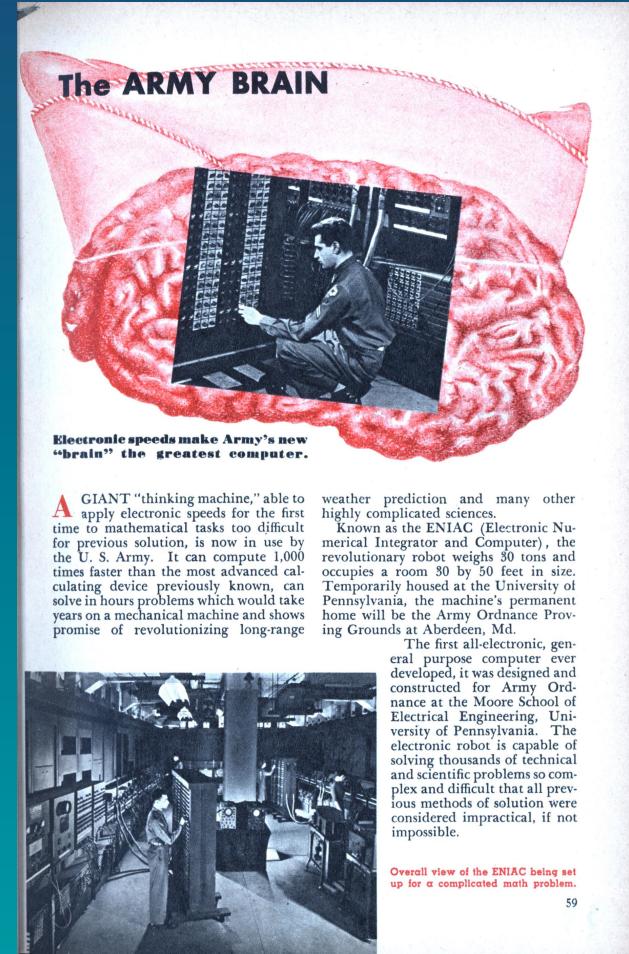
Vice President Technology, Qualcomm Research

March 4, 2014

Neuromorphic Processing : A New Frontier in Scaling Computer Architecture

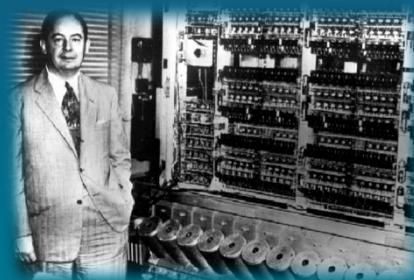
QUALCOMM®



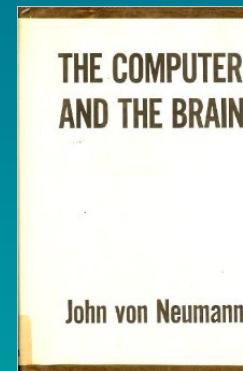


ENIAC – 1946

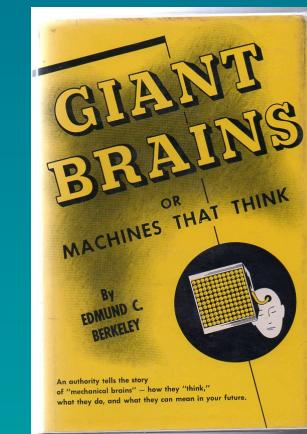
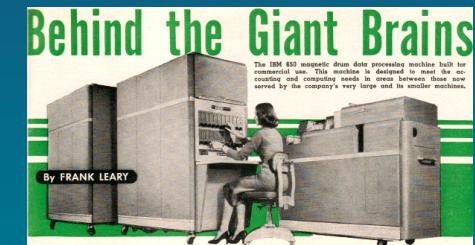
The desire to build brain-like computers is as old as the computer itself



IAS - 1945

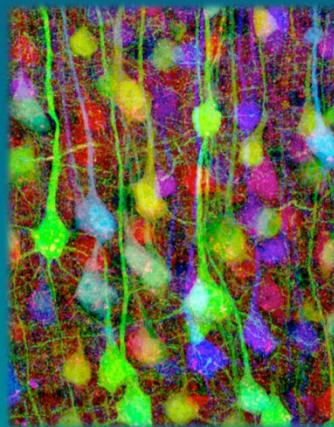


1958



Brain-inspired computing: why?

What makes the brain special?



Fault tolerance

Inputs may be noisy
Neurons may die



Power efficiency

20 W power consumption
Always On



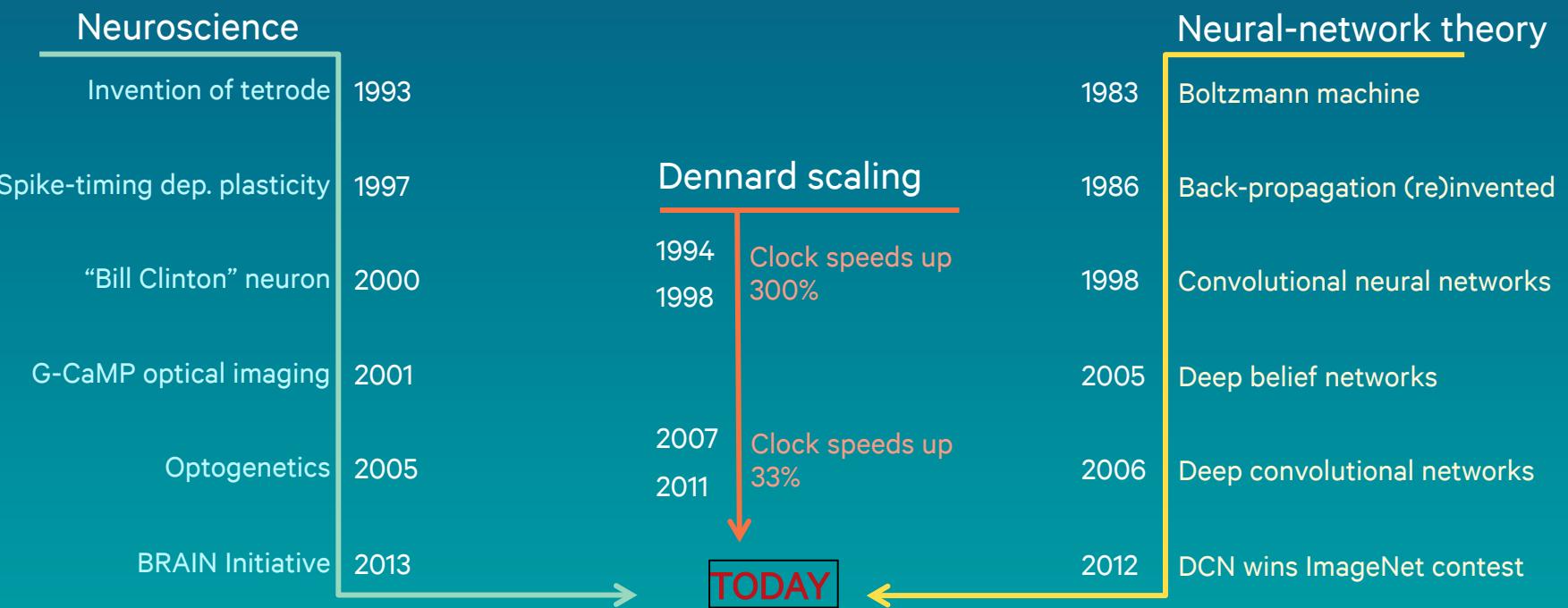
Learning ability

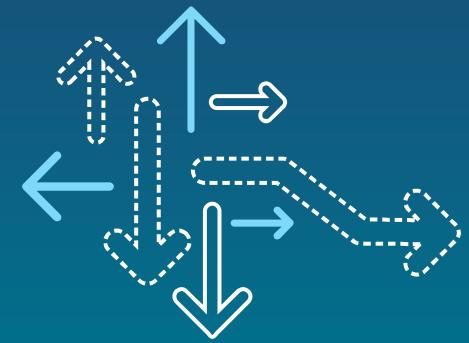
Supervised learning
Unsupervised learning
Reinforcement learning

Object identification
Language acquisition
Motor learning

Brain-inspired Computing: Why Now?

A merger of three trends



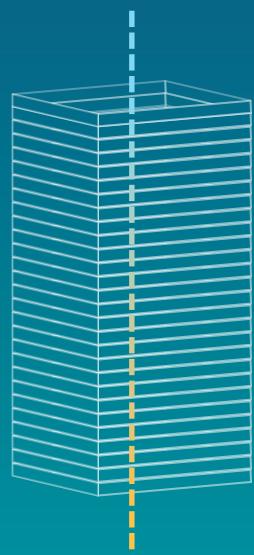


**So is now finally the time for
artificial cognition?**

The brain is a massively parallel machine

$>10^6$ processing steps

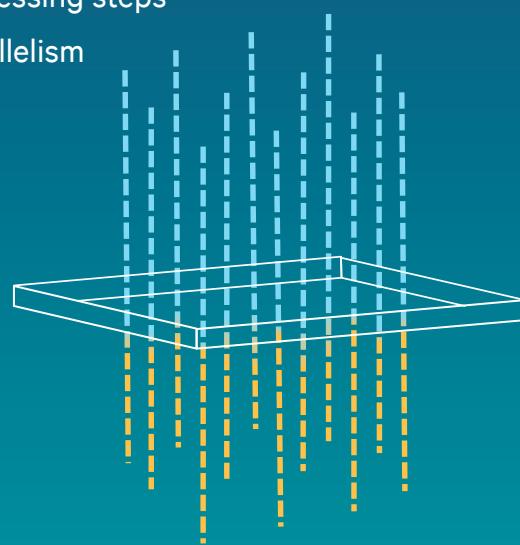
$<10^1$ parallelism



Modern computer
Dense, real-valued data

$<10^1$ processing steps

$>10^6$ parallelism



Human brain
Sparse ‘Events’ or “Spikes”

Mobile is the most challenging design environment



Power efficiency

Always on



Performance

Small form factor

Qualcomm Zeroth™

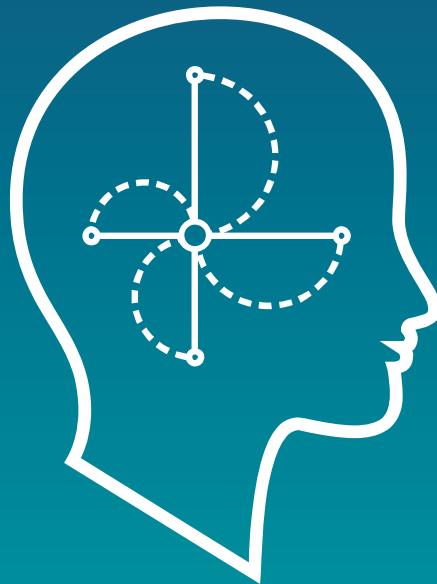
Enabling perception and cognition for smart devices



Power efficiency



**Natural, contextual
sensor processing**



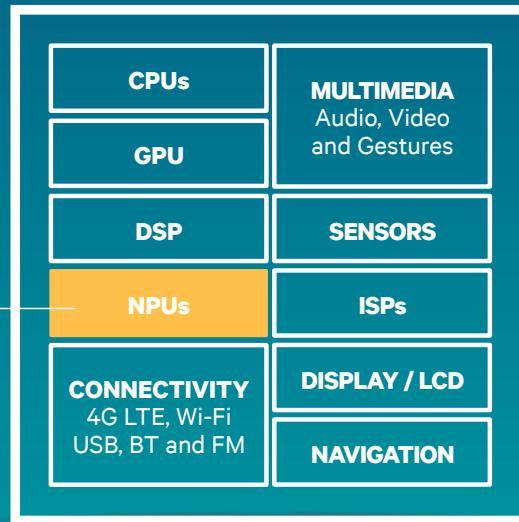
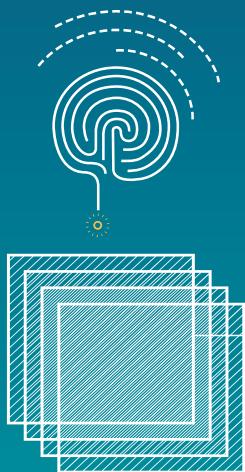
**Cognitive
algorithms**



**Continuous
machine learning**

Neural Processing Units (NPUs)

A new class of processors mimicking human perception and cognition

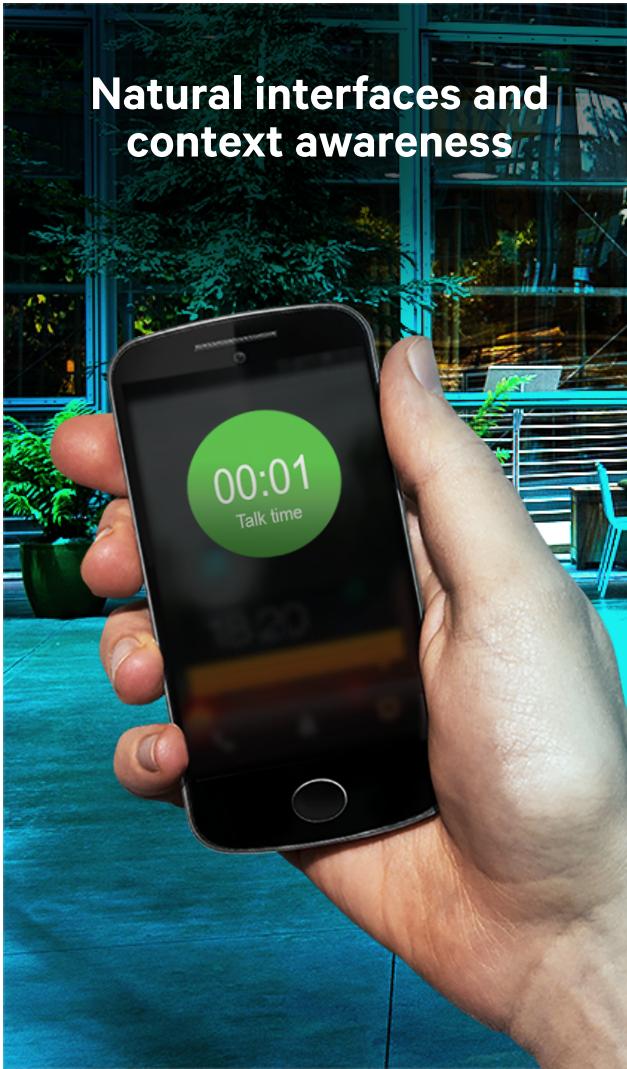


**Massively parallel,
reprogrammable**

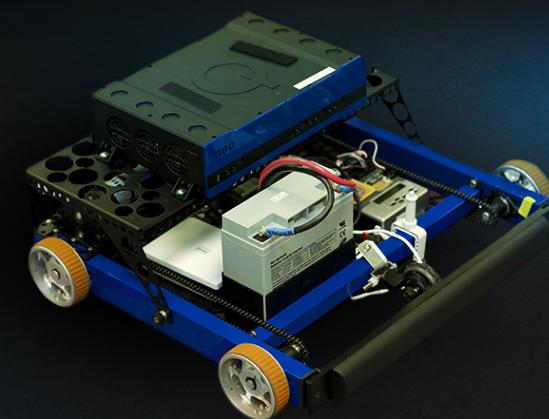
Comprehensive tools

Human-like functions

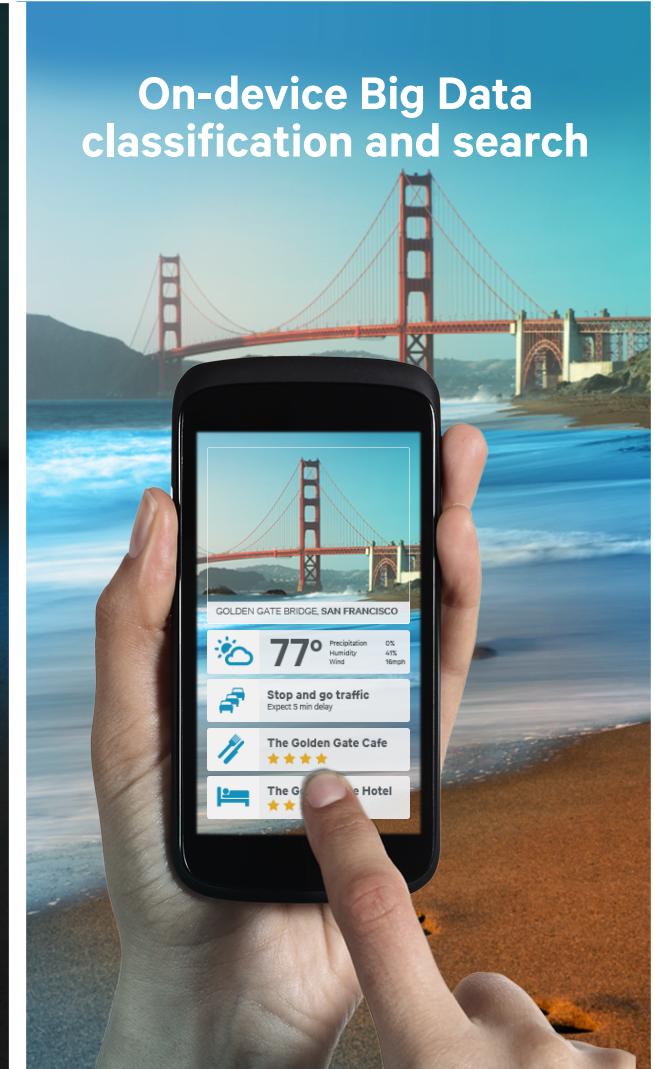
Natural interfaces and context awareness



Robot learning and control



On-device Big Data classification and search



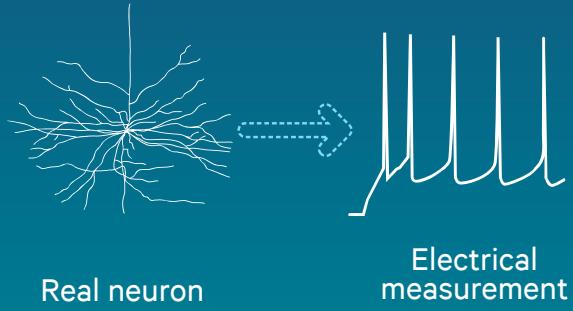
Key Challenges – Building A Neural Processor

- Networks and Algorithms
- Hardware Representation
- Software Tools and Languages

Complete Solution Requires Consideration of the Entire Problem

Zeroth™ approach: start with Spiking Neurons

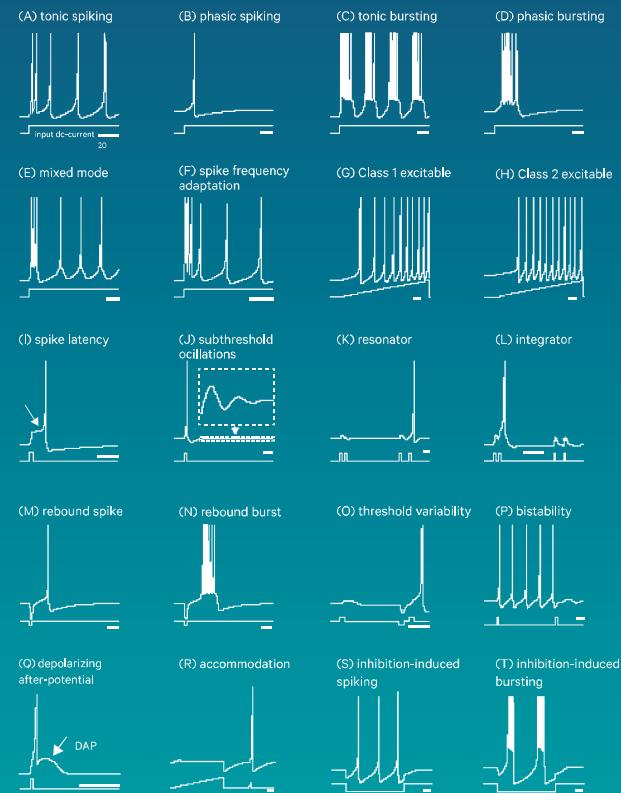
Developing complex neuron models that can be implemented in hardware



$$C \frac{dv}{dt} = F(v) - u + I$$

$$\frac{du}{dt} = a(b(v - v_r) - u)$$

Individual neuron
modeling



Family of
models

Designing a Better Neuron Model – COLD

Linear vs. non-linear models

COLD – Direct computation of time-to-spike

Simple Model – requires numerical integration

Sub and Supra-Threshold Regions

COLD – Independent Regions

Simple – coupled regions, constrained design

Event Based Architecture

Ease of Neuron Design

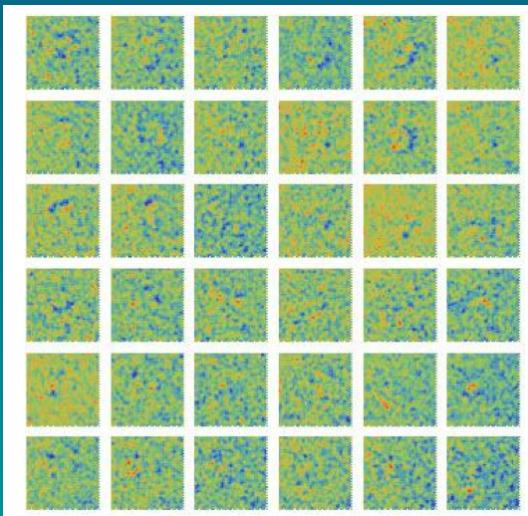
Color and Orientation Pop-out Networks

Results with COLD and Event Based Approaches

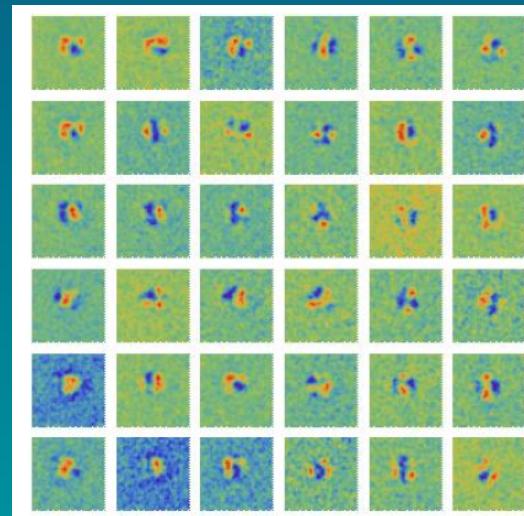
- Emergent networks for color pop-out (parvo) and orientation (magno) were constructed
- Using Qualcomm's COLD neuron formulation and an early event-based simulation approach
- Receptive field emergence for orientation was achieved in much shorter simulation times than other approaches

Orientation Receptive Field Formation

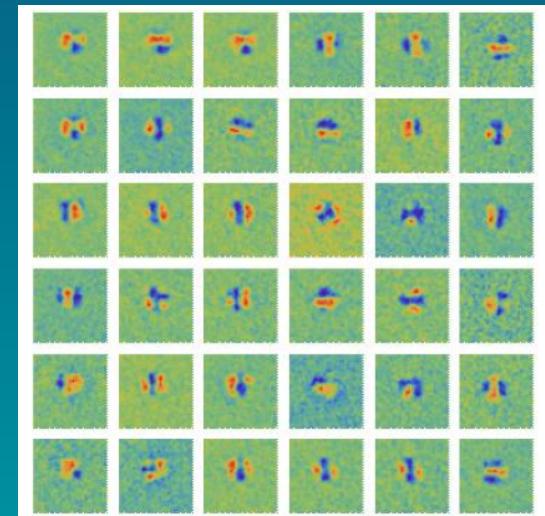
Rapid Specialization



Undifferentiated Fields



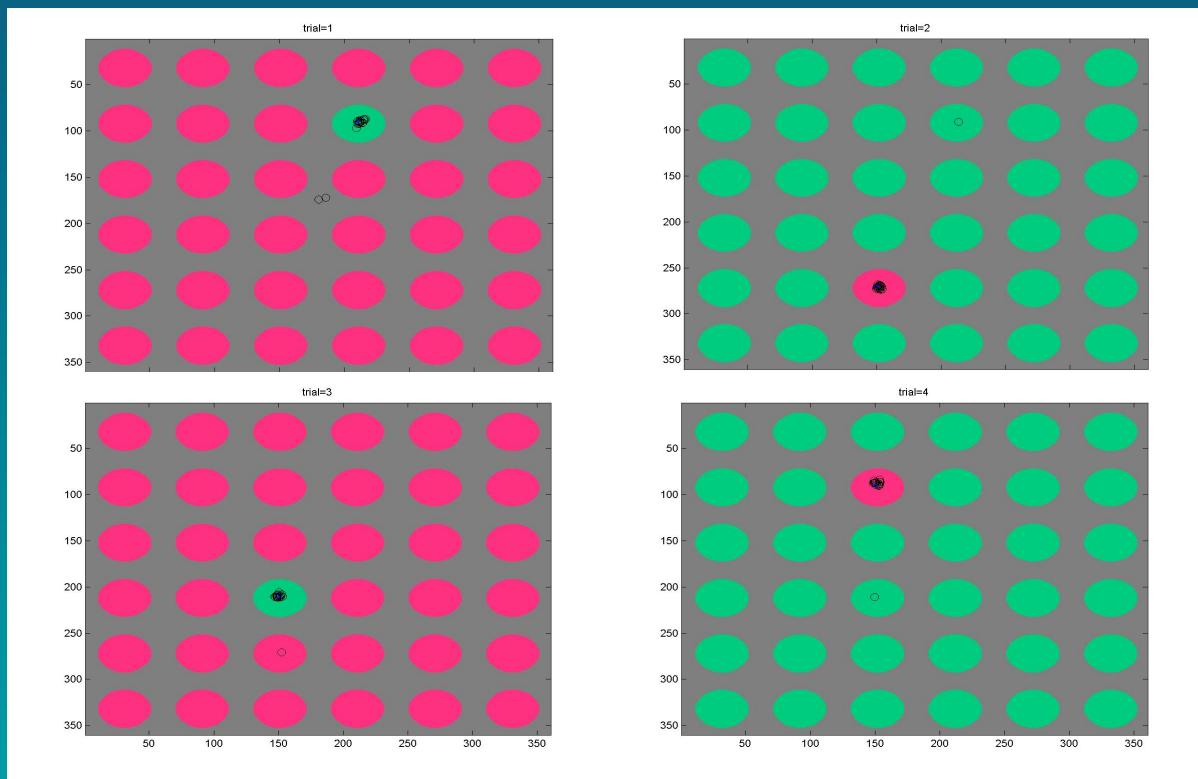
After 30 Seconds



After 100 Seconds

Color Pop-Out with Saccade Locations

The Red – Green Test

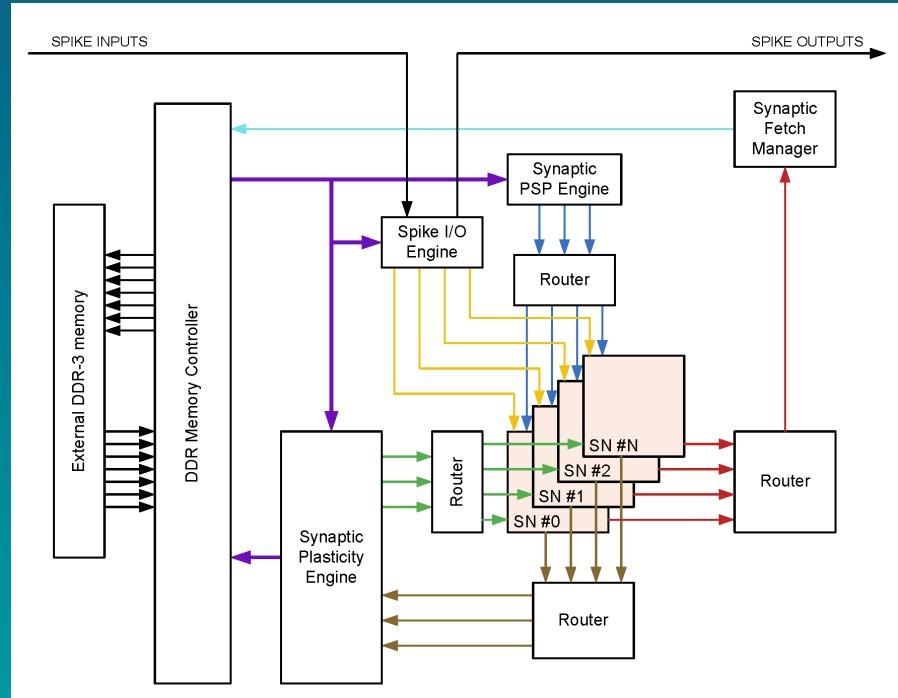


High Level Hardware Considerations

- Analog or Digital Architecture
- Fixed Versus Programmable Neuron Behavior
- High Density Versus Power Considerations

Hardware Overview

- Neural Accelerator: Custom tuned for efficient execution of spiking neural networks
- Programmable neuron and synapse dynamics
- Time multiplexed implementation across a scalable number of “Super Neurons”
- On board neuronal state, off-chip DRAM used for high density storage of synaptic connections



Key Future Hardware Challenges

Further Avenues of Investigation

- Future Memories
- Asynchronous logic and Event Based Processing
- The right level of programmability
- New Classes of Sensors



Software and Programming Considerations

The Brain is Parallel, Not Linear!

- Programs no longer linear
- Program Logic Encoded into the Network Itself
- Common Design Flow Across Output Targets
- Leveraged LLVM



Nature of the Programming Language

Basic Building Blocks

```
define unit exc_neuron from cold {
    init {
        a = 0.02;
        b = 0.1;
        d = 10;
    }
}

define unit inh_neuron from cold {
    init {
        a = 0.1;
        b = 0.4;
        d = 8;
    }
}
```

Define Units

```
// create excitatory (exc) population
create exc_neuron on Grid( countX=Nex, countY=Ney ) as exc, all;
create random_input set prob=1.0/N on Grid( countX=Nex, countY=Ney ) as excInput;

// create inhibitory (inh) population
create inh_neuron on Grid( countX=Nix, countY=Niy ) as inh, all;
create random_input set prob=1.0/N on Grid( countX=Nix, countY=Niy ) as inhInput;

connect from excInput to 1 nearest exc with fixed_syn set weight=20, delay=1;
connect from inhInput to 1 nearest inh with fixed_syn set weight=20, delay=1;
```

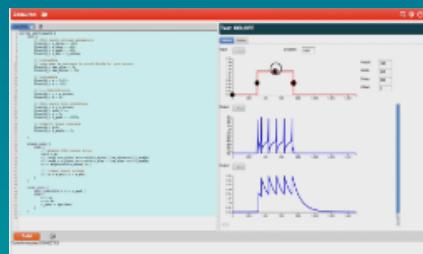
Create Populations

Connect

Zeroth Neural Network Tools Platform

Zeroth Development Studio (ZDS): end to end package

- Programming Language: High Level Network Description language (HLND)
- Development Environment: Integrated Development Environment (IDE)
- Libraries- Neuromorphic Development Kit (NDK)
- Robotics Simulator- Virtualization World (VW)



HLND example

IDE/NDK



Supports multiple neuron types

Support multiple execution targets:

- Software simulation: Linux workstation or equivalent (*today*)
- Hardware emulation: Custom Xilinx-K7 based FPGA board (*today*)
- ASIC Chip: Zeroth Neural Processing Unit (NPU) (*future*)

CPU (*today*)
GPU (*future*)



1-N Compute Cluster



Automatic distribution onto multiple cores
→ scalable



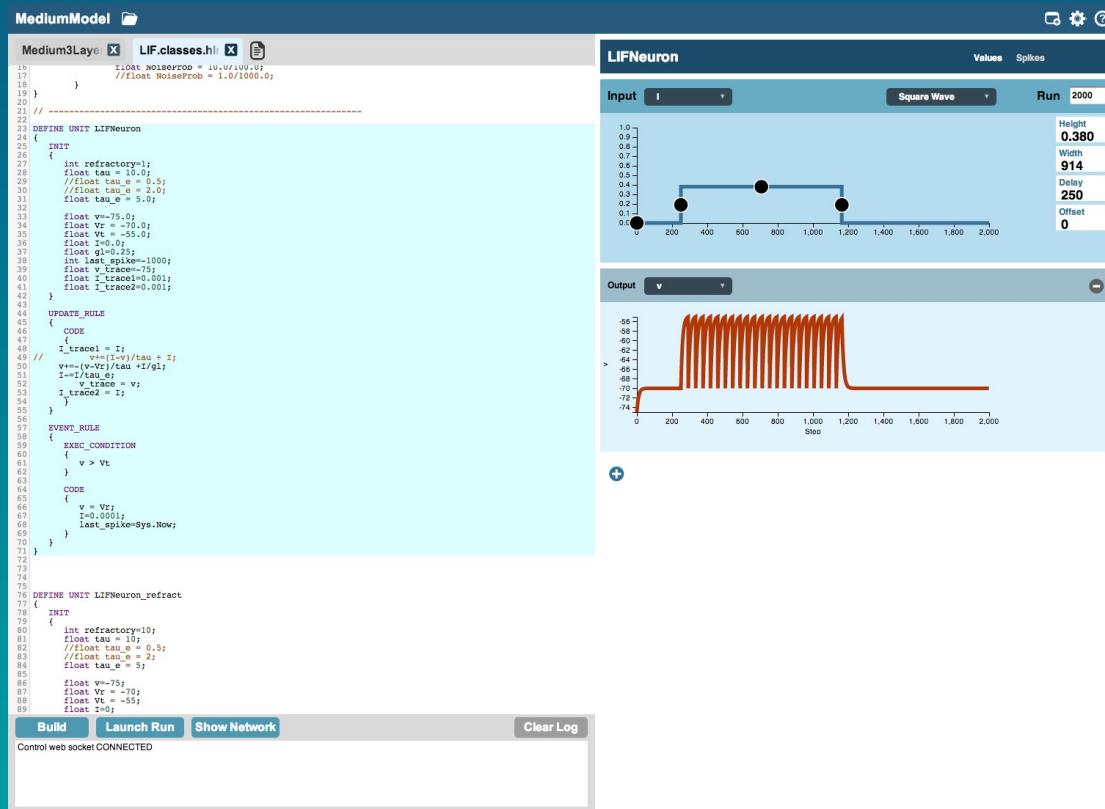
FPGA



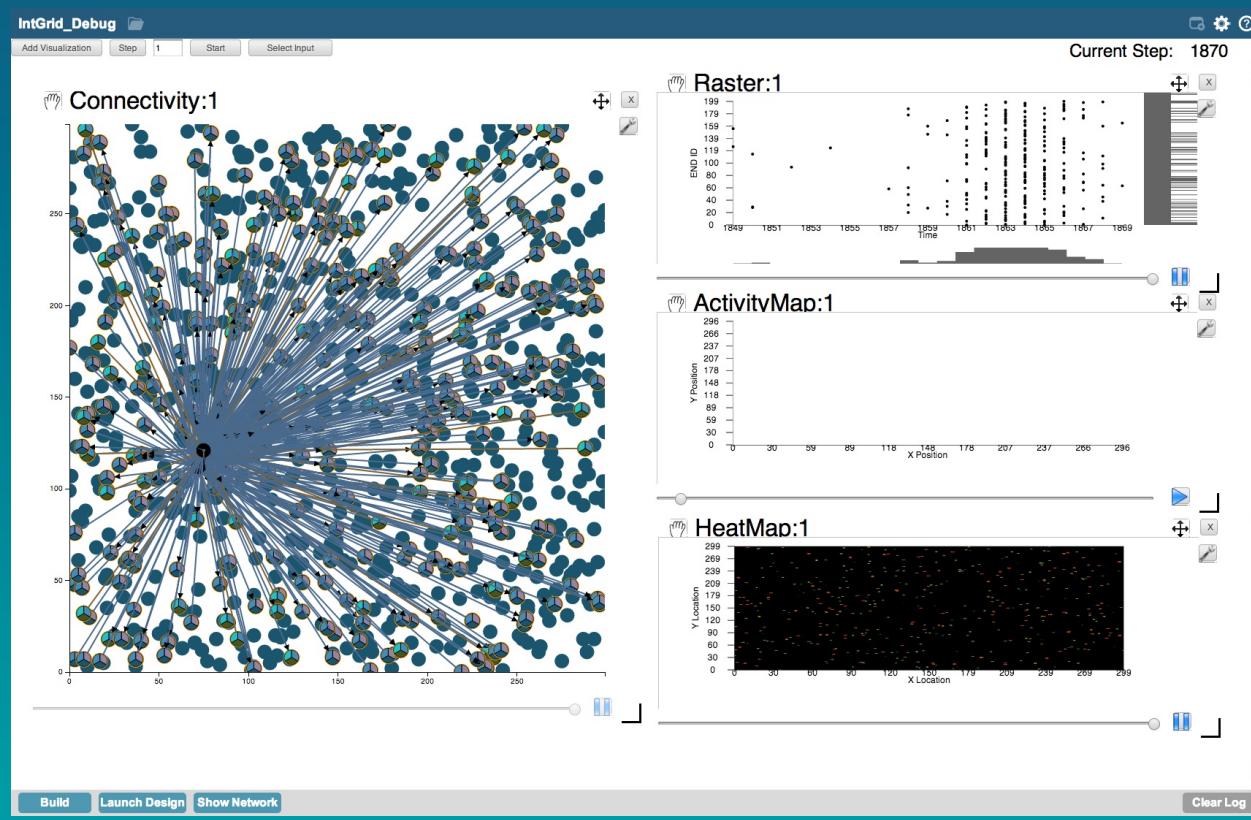
NPU ASIC (*Future*)

HTML 5 / Javascript IDE environment

Dynamic Code / Neuron parameter co-design

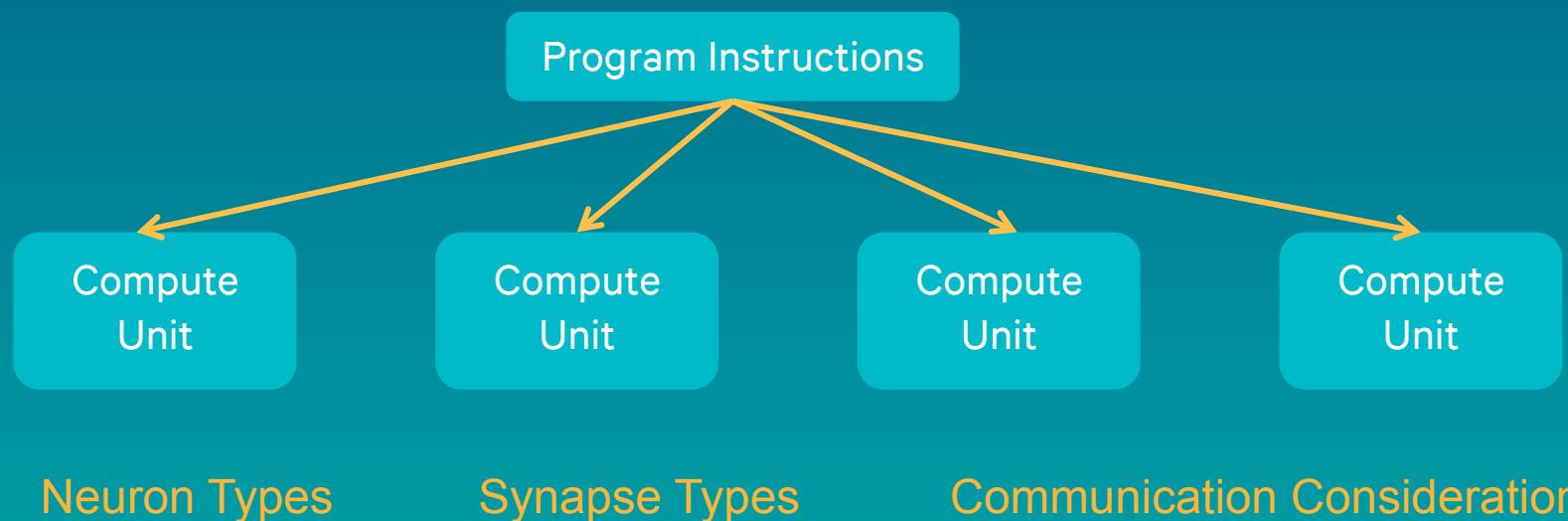


Real Time Data Visualizations



Hardware Design Creates Software Complexities

- Parallel and programmable nature of the hardware creates a complex “place and route” challenge



“Call to Action”

Research Opportunities Abound! – And We Can’t Do it Alone!

- Programming languages and abstractions to capture and represent these networks
- Better “routing and placement” algorithms for allocation to these massively parallel systems
- Breakthroughs in representation and training (e.g. DCNs take forever to train)
- HW architecture for asynch logic, continued advancements in memories
- Low power always-on sensors
- *Demonstrate “True Cognition”*

Thank you

Follow us on:   

For more information on Qualcomm, visit us at:
www.qualcomm.com & www.qualcomm.com/blog

© 2014 Qualcomm Incorporated. Qualcomm, Snapdragon, Dragonboard and Vuforia are trademarks of Qualcomm Incorporated, registered in the United States and other countries. Zeroth is a trademark of Qualcomm Incorporated. AllJoyn is a trademark of Qualcomm Innovation Center, Inc., registered in the United States and other countries. Other products and brand names may be trademarks or registered trademarks of their respective owners.

Qualcomm Incorporated includes Qualcomm's licensing business, QTL, and the vast majority of its patent portfolio. Qualcomm Technologies, Inc., a wholly-owned subsidiary of Qualcomm Incorporated, operates, along with its subsidiaries, substantially all of Qualcomm's engineering, research and development functions, and substantially all of its product and services businesses, including its semiconductor business, QCT, and QWI. References to "Qualcomm" may mean Qualcomm Incorporated, or subsidiaries or business units within the Qualcomm corporate structure, as applicable.

