#### Exceptional service in the national interest







### Neural Computing for Scientific Computing Applications: More than Just Machine Learning

Neuromorphic Computing Workshop, Knoxville TN, 7/17/17 Brad Aimone (jbaimon@sandia.gov), Ojas Parekh, William Severa Sandia National Laboratories

Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. SAND NO. 2011-XXXXP

## Hardware Acceleration of Adaptive Neural Algorithms (HAANA) 2014-2017



Neural inspired computing lacks theoretical foundation to translate between fields





Neural inspired computing lacks theoretical foundation to translate between fields





What is the brain as inspiration?





#### Why?

- It is a challenge to separate brains (cognitive capability) from neurons (low-energy mechanism)
- Belief that neurons are noisy
- Moore's Law It has always been easier to wait for faster processors than to re-invent numerical computing on specialized parallel architecture

### The Brain as Computer: Bad at Math, Good at Everything Else

Modeling computers after the brain could revolutionize robotics and big data



#### A Unique Machine



#### nfographic:

From Macro to Micro: A Visual Guide to the Brain

#### Why We Should Copy the Brain

Trying to create consciousness may be the path to understanding this most deeply mysterious human attribute By Glenn Zorpette

#### In the Future, Machines Will Borrow Our Brain's Best Tricks

A researcher imagines how true artificial intelligence will change the world By Fred Rothganger

#### The Brain as Computer: Bad at Math, Good at Everything Else

Modeling computers after the brain could revolutionize robotics and big data By Karlheim Meler

## Theoretical models of the brain do not need to capture everything







Shallow Depth Inference Rapid, Stable Learning Context Modulated Decisions Memory Capacity Power Efficient Distributed Representations Not Consistently Logical Bad at Math = Implicit in model

= Not implicit in model



Neuroscience Systems Model

### Spiking neurons are a more powerful version of classic logic gates

Spiking threshold gates provide high degree of parallelism at very low power



### Are threshold gates and spiking neurons equivalent?





## HAANA has produced a number of spiking numerical algorithms

- Cross-correlation
  - Severa et al., ICRC 2016
- SpikeSort
  - Verzi et al., submitted
  - SpikeMin
  - SpikeMax
- SpikeOptimization
  - Verzi et al., IJCNN 2017
- Sub-cubic (i.e., Strassen) constant depth matrix multiplication

t = 1

Parekh et al., submitted









### **A Velocimetry Application**

- A motivating application is the determination of the local velocity in a flow field
- The maximal cross-correlation between two sample images provides a velocity estimate
- SNN algorithms are straightforward; exemplify core concepts
  - Highly parallel
  - Different neural representations
  - Modular, precise connectivity
  - Time/Neuron tradeoff



### **Time Multiplexed Cross Correlation**





Severa et al., ICRC 2016

## Cross-Correlation Exhibits Time/Neuron Tradeoff





Severa et al., ICRC 2016

- Exchange Time Cost ↔ Neuron Cost
- Complexity is unchanged
- Neurons:  $O(n^2) \leftrightarrow O(n)$
- Time:  $\boldsymbol{0}(1) \leftrightarrow \boldsymbol{0}(n)$



### "Neural" network for matrix multiplication



Strassen formulation of matrix multiply enables less than O(N<sup>3</sup>) neurons – resulting in less power consumption

Parekh et al., submitted



## Strassen multiplication in neural hardware may show powerful advantages

	Depth	# Gates	Value of $\epsilon$
Standard	3	O(N <sup>3</sup> )	-
"Direct" Strassen	d	Ο(N <sup>ω + ε</sup> )	1/d
Refined Strassen	d	Ο(Ν <sup>ω + ε</sup> )	O(1/c <sup>d</sup> )
Non-constant Depth	O(log log N)	Ο(Ν <sup>ω</sup> )	-

#### **Example: Triangle Counting in Graphs**



Output: does the graph have ≥ T triangles? Applications to social network analysis



Parekh et al., submitted

## Theoretical models of the brain do not need to capture everything





Shallow Depth Inference Rapid, Stable Learning Context Modulated Decisions Memory Capacity Power Efficient Distributed Representations Not Consistently Logical Bad at Math



### How do we take advantage of neuroscience?





Primate visual cortex Felleman and Van Essen, 1991

Hippocampus

### View of brain as computing system





### Cortex – hippocampus interaction can extend AI to more complete computing system





- Cortex learns to process sensory information at different levels of abstraction
  - Similar to deep learning, though more sophisticated in biology
- Hippocampus would be a content addressable memory
  - Provide context and retrieval cues to guide cortical processing

## A robust hippocampus abstraction can bring a complete neural system to AI

### Desired functions

- Learn associations between cortical modalities
- Encoding of temporal, contextual, and spatial information into associations
- Ability for "one-shot" learning
- Cue-based retrieval of information

#### Desired properties

- Compatible with spiking representations
- Network must be stable with adaptation
- Capacity should scale nicely
- Biologically plausible in context of extensive hippocampus literature
- Ability to formally quantify costs and performance



# Formalizing CAM function one hippocampus layer at a time



 Constraining EC inputs to have "grid cell" structure sets DG size to biological level of expansion (~10:1)



 Mixed code of broadtuned (immature) neurons and narrow tuned (mature) neurons confirms predicted ability to encode novel information

> William Severa, NICE 2016 Severa et al., Neural Computation, 2017



## Brain uses a different approach to processing in memory









### **Questions?**

Thanks to Sandia's LDRD HAANA Grand Challenge and the DOE NNSA Advanced Simulation and Computing program