# From T0 and CNS-1 to RISC-V and AI Hardware

Krste Asanovic

Professor, EECS Dept, UC Berkeley
Chairman of Board, RISC-V Foundation
Co-Founder and Chief Architect, SiFive

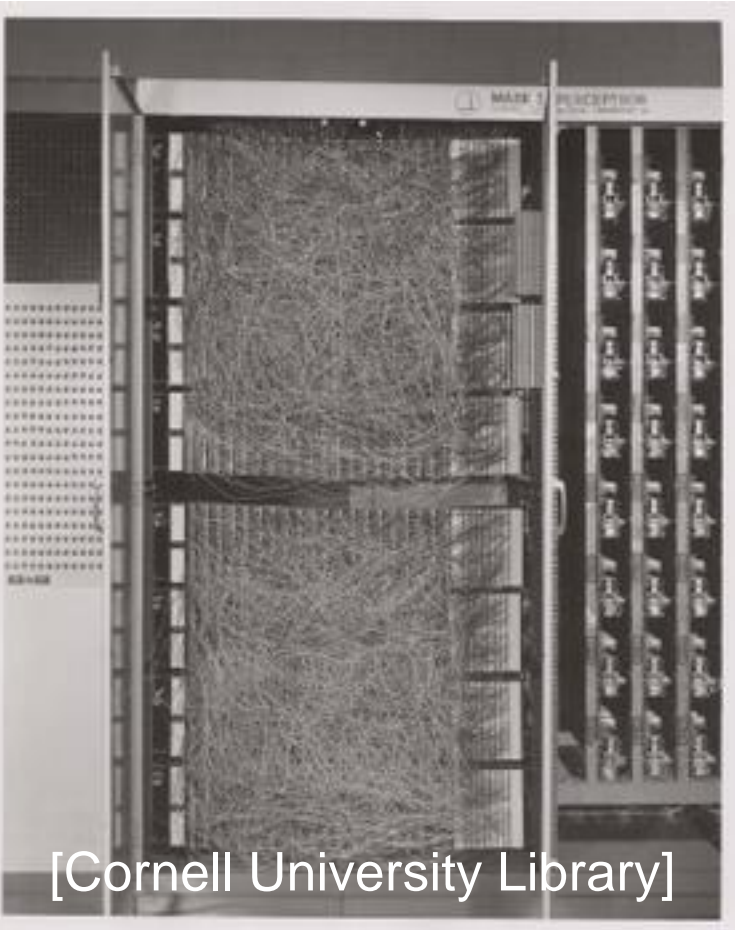ICSI 30 Years of Innovation
October 5, 2018

# Three Waves of Artificial Neural Networks

- 1950s/60s: Perceptrons (1 layer)

- 1980s/90s: Backpropagation, (2-3 layers)

- 2010s: Deep Neural Networks (3+ layers)

*(Note, all ideas were developed much earlier than eventual popularity)*

**2**

# Single-Layer Perceptrons, 1950s-60s
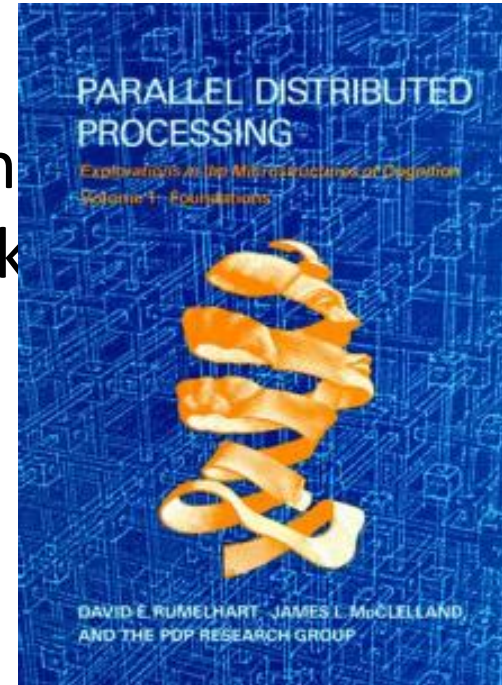


[Cornell University Library]

- Single-layer perceptrons explored for image processing [Rosenblatt]
- Only a linearly separable classifier
  - Unable to capture interesting functions, e.g. XOR
- AI field moved from statistical to symbolic approaches in `70s-`80s

**3**

# Multi-Layer Perceptrons and BackProp

- Influential PDP books published in 1986
- Two-layer backprop trained networks foun[d] to work surprisingly well at many hard task[s]
- Experts complained that results were unexplainable
- Training was extremely slow, so rush to build custom machines



PARALLEL DISTRIBUTED PROCESSING

Explorations in the Microstructure of Cognition

Volume 1: Foundations

DAVID E. RUMELHART, JAMES L. McCLELLAND, AND THE PDP RESEARCH GROUP

# Ring Array Processor, (ICSI 1989)

### (Nelson Morgan, Jim Beck, Phil Kohn, Jeff Bilmes)





**Fig. 3** RAP Performance for uniform layer size, one hidden layer.

- RAP Machine built for fast training of "big dumb" neural networks for speech recognition
- Ring of TMS320C30 floating-point DSPs
  - Each DSP providing 32MFLOPS (32-bit FP)
  - Four DSPs/board, up to 10 boards connected at once (>1GFLOP/s peak, 640MB DRAM)
  - Neural net training rate of >100MCUPS (million connection updates per second) on 10 boards
  - FPGA ring connection used for systolic all-all communication during training/inference
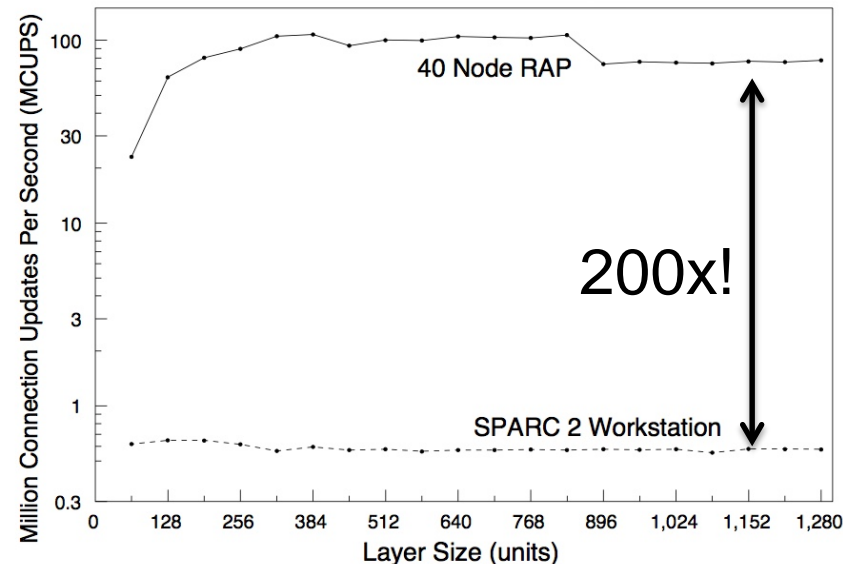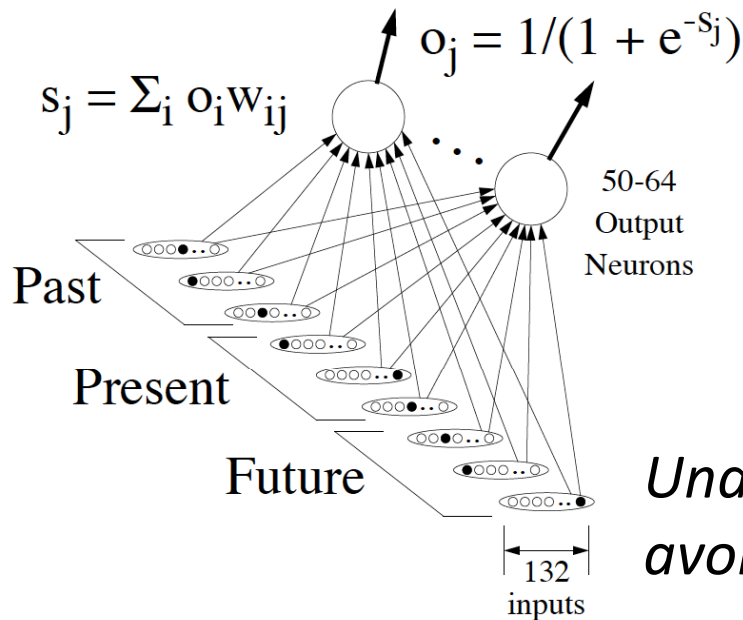- Fast, flexible, but expensive
  - ~$100,000 each

# Realization Group, ICSI, 1989

New naïve grad student joins Morgan's group to build custom ANN VLSI for speech training

$$o_j = 1/(1 + e^{-s_j})$$

$$s_j = \Sigma_i\, o_i w_{ij}$$

50-64 Output Neurons

Past

Present

Future

132 inputs

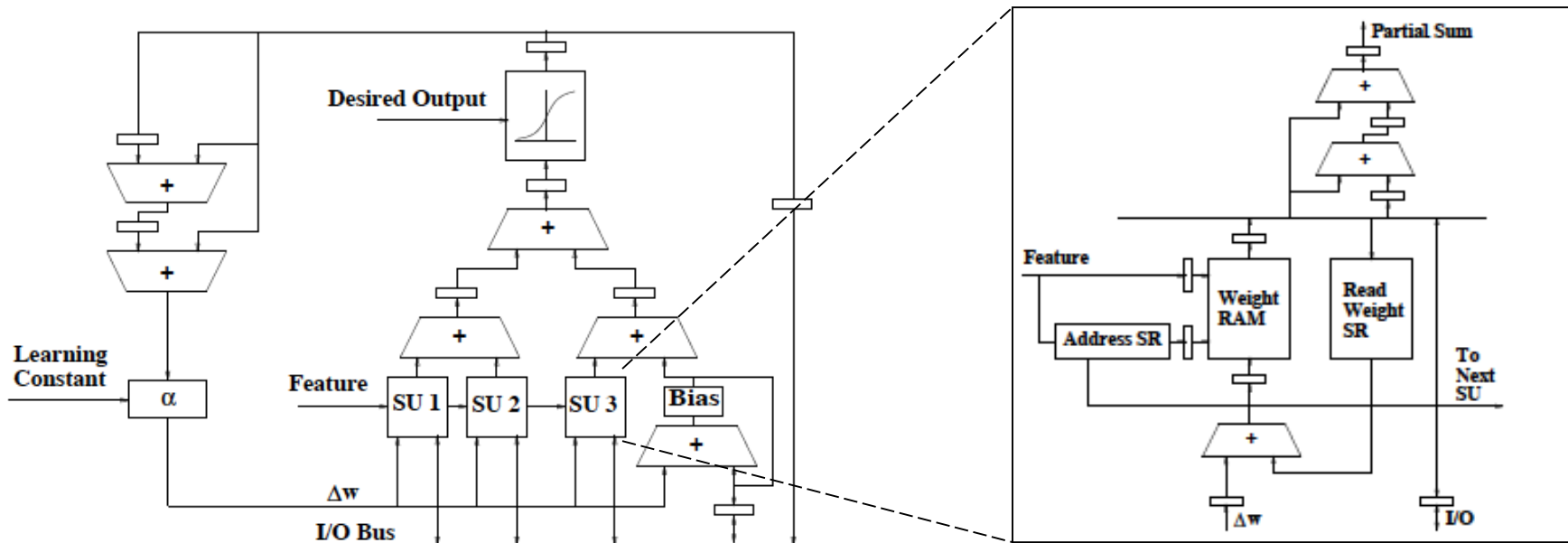This is a cool ANN architecture for which we need custom silicon!

*Unary-encoded inputs to avoid multiply, 12-bit weights*

Training rule: $\Delta w_{ij} = -\alpha(o_j - d_j)\, o_i$

# HiPNeT-1: (Highly Pipelined Network Trainer, 1990)

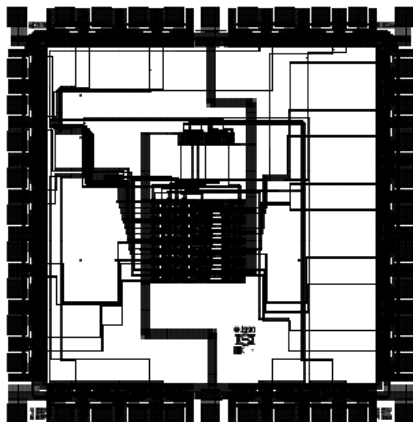Krste Asanovic, Brian Kingsbury, Nelson Morgan, John Wawrzynek



- Custom architecture for neural algorithm
- Ignores pipeline RAW hazards (net trains around them)
- Predicted 200MCUPS in 16mm$^2$ of 2μm CMOS running at 20MHz

# The first few chips…

- MOSIS had a "TinyChip" program
  - $500 to fab a 2.2mmx2.2mm chip in 2μm CMOS

Sigmoid unit (Pawan Sinha)

JTAG latches (Krste)

Multiplier (Brian)

8-bit datapath (Krste)

24b Adder (Brian)

Regfile (Bertrand)

# Meanwhile, back at the speech ranch…



Time for a programmable architecture…  **9**

# "Old" SPERT VLIW/SIMD Engine

VLIW Instruction

| Mult | Shift | Add | Limit |
|------|-------|-----|-------|
| Mult | Shift | Add | Limit |
| Mult | Shift | Add | Limit |
| Mult | Shift | Add | Limit |
| Mult | Shift | Add | Limit |
| Mult | Shift | Add | Limit |
| Mult | Shift | Add | Limit |
| Mult | Shift | Add | Limit |

Vector Unit

| | ALU | Add1 | Add2 |
|--|-----|------|------|

Scalar Unit

Memory Control

*Similar architecture later adopted by many embedded DSPs, especially for video and games.*
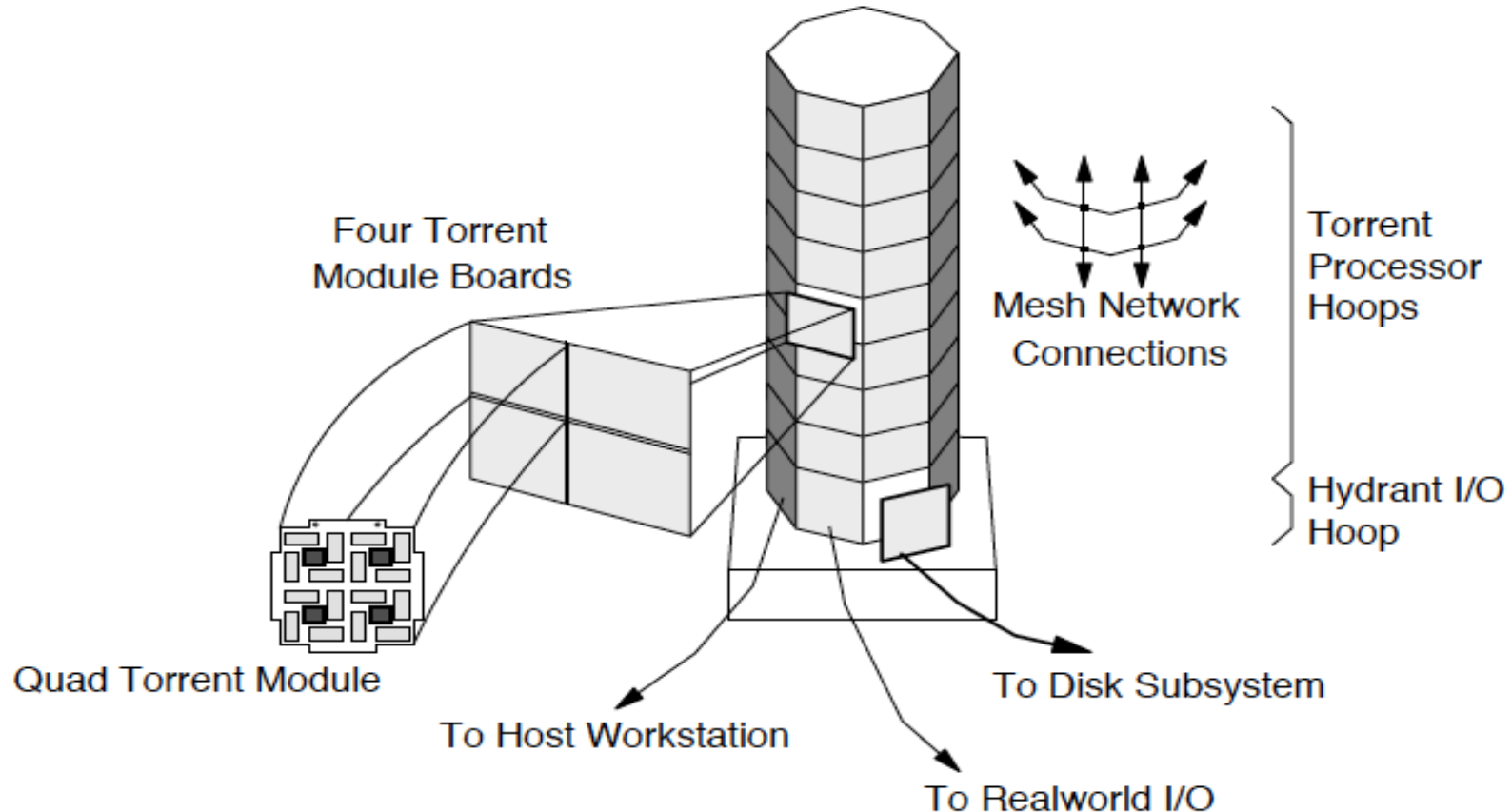
**10**

# SQUIRT Test Chip, 1992



- 1.2μm CMOS, 2 metal layers
- 61,521 transistors, 8x4 mm$^2$, 400mW@5V, 50MHz
- 72-bit VLIW instruction word
- 16x32b register file, 24bx8b->32b multiplier, 32b ALU/shifter/clipper

**11**

# CNS-1: Connectionist Network Supercomputer
## (ICSI/UCB 1992-95)



Four Torrent Module Boards

Mesh Network Connections

Torrent Processor Hoops

Hydrant I/O Hoop

Quad Torrent Module

To Host Workstation

To Realworld I/O

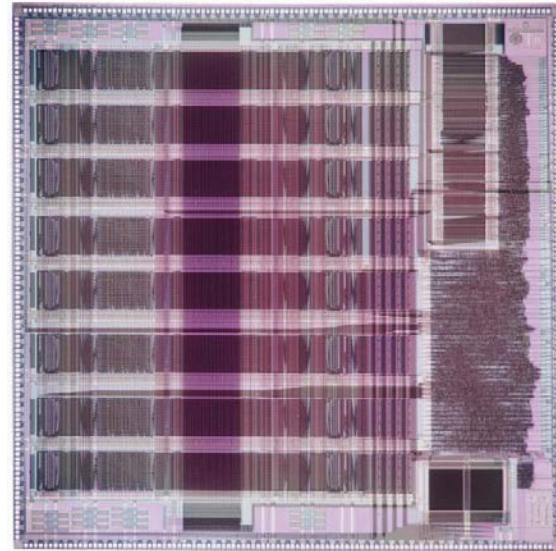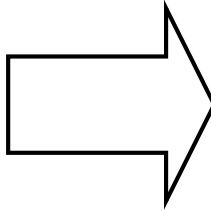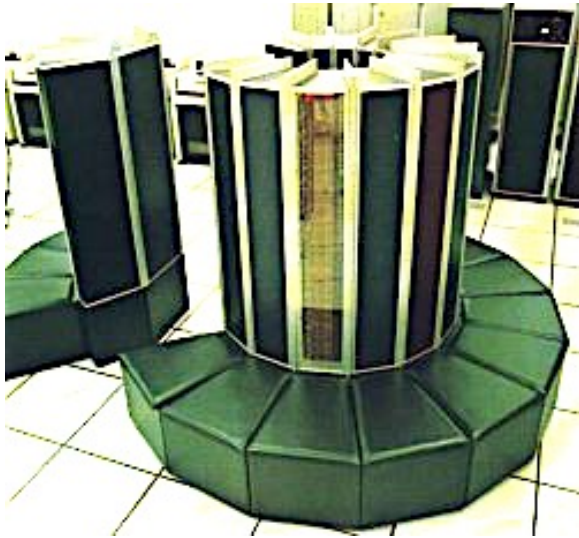To Disk Subsystem

**12**

CNS Research Group

# We abandoned SPERT VLIW

- **VLIW means no upward compatibility**
  - we wanted same ISA for CNS-1 to reuse software effort
- **VLIW scalar compiler was complex**
  - Simple VLIW hardware + complex VLIW compiler harder than complex RISC architecture + standard compiler
- **Assembly code was tough to write**
  - soon discovered this when writing test code and key loops
- **VLIW format too rigid**
  - hard to fit some operations into statically scheduled instruction slots (misaligned vector loads/stores, scatter/gathers)
- **VLIW had too large an instruction-cache footprint**
  - loop prologue/epilogue code plus unrolled loop body

*Software, software, software,….*

14

# T0: First Vector Microprocessor (1995)

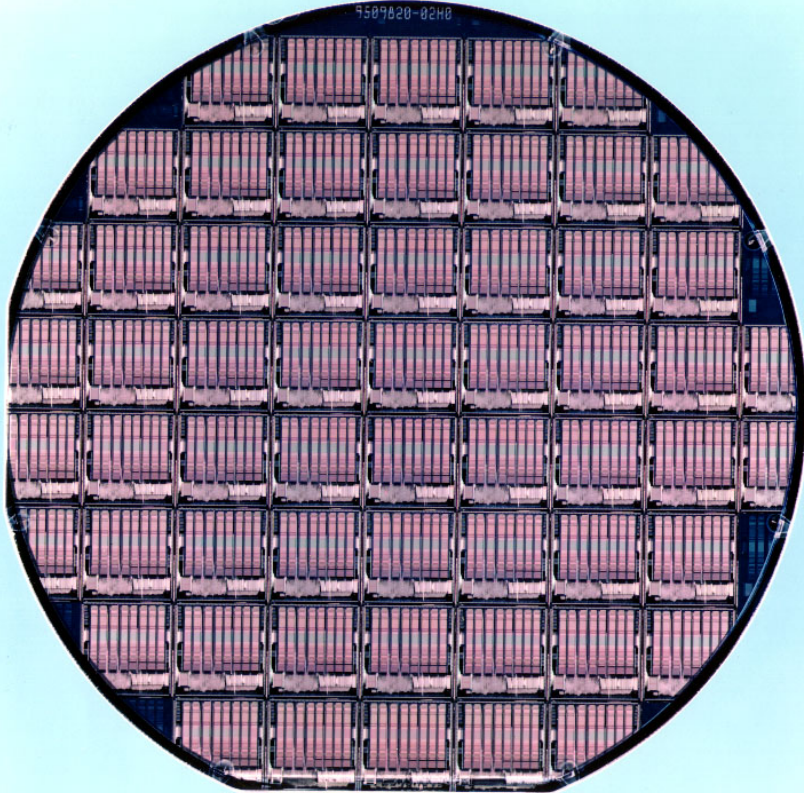- Vector supercomputers (e.g., Crays) very successful in scientific computing, clean programming model



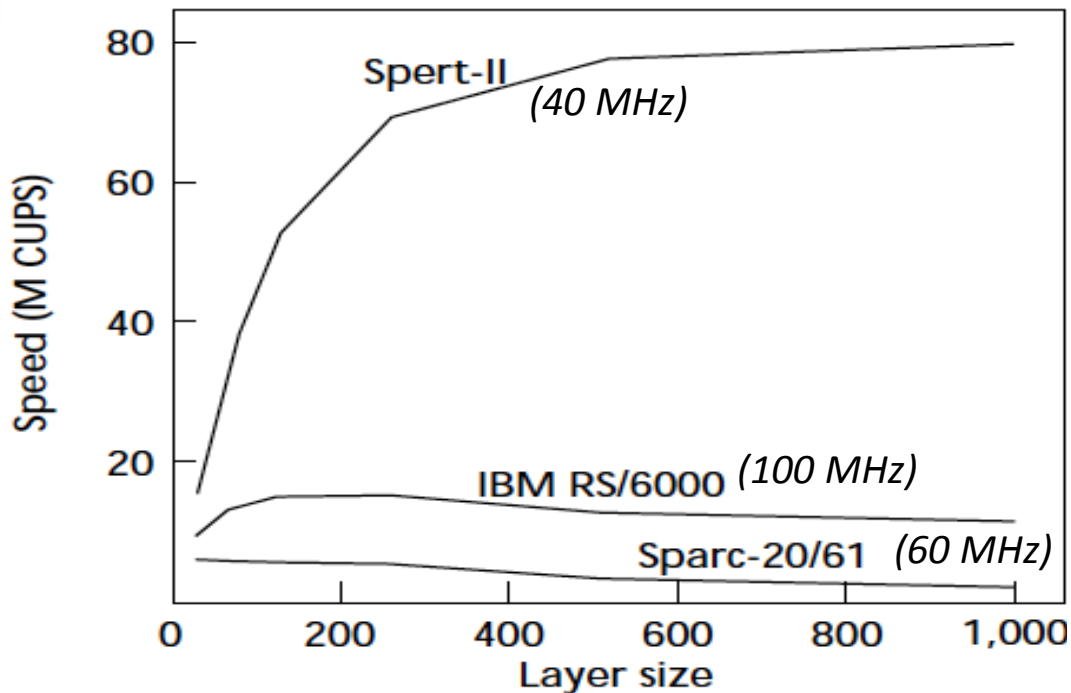Add a vector coprocessor to a standard MIPS RISC scalar processor, all on one chip, for neural net training

**15**

# System Design Choices

- **Which standard RISC?**
  - Considered SPARC, HP PA, PowerPC, and Alpha
  - Chose MIPS because: simplest, good software tools, Unix desktop workstations for development, and a 64-bit extension path
- **Buy or build a MIPS core?**
  - Commercial MIPS R3000 chips had coprocessor interface
  - No MIPS soft cores (no Verilog or synthesis tools yet)
  - Decided to roll our own
    - Vector coprocessor would have played havoc with caches
    - Coprocessor interface too inefficient
    - Commercial chip plus glue logic would blow our size and power budgets (to fit inside workstation)
    - Couldn't simulate whole system in our environment

**16**

# Brian, Krste, and a Torrent wafer

# SPERT-II / T0 Vector Microprocessor (1995)



- Boards shipped to 9 international sites
- Used as production research platform for nine years
  - last time powered up for work in 2004!

# TetraSpert (1997): faster training



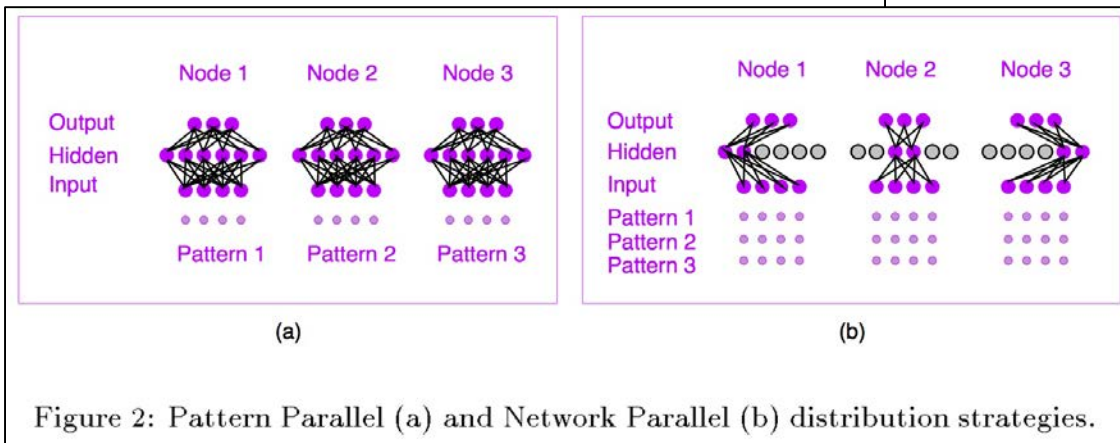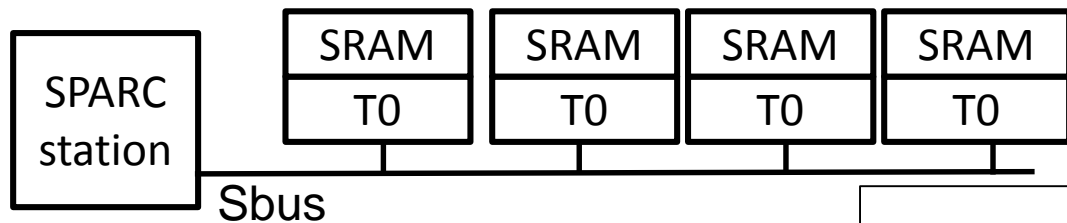| SPARC station | SRAM | SRAM | SRAM | SRAM |
|---|---|---|---|---|
| | T0 | T0 | T0 | T0 |

Sbus

Figure 2: Pattern Parallel (a) and Network Parallel (b) distribution strategies.

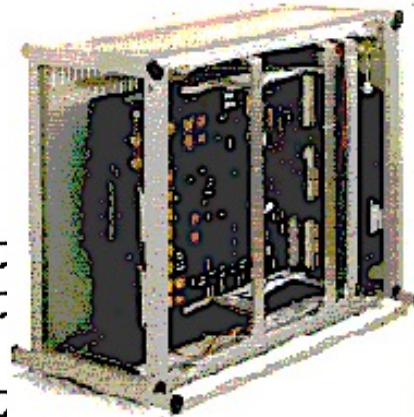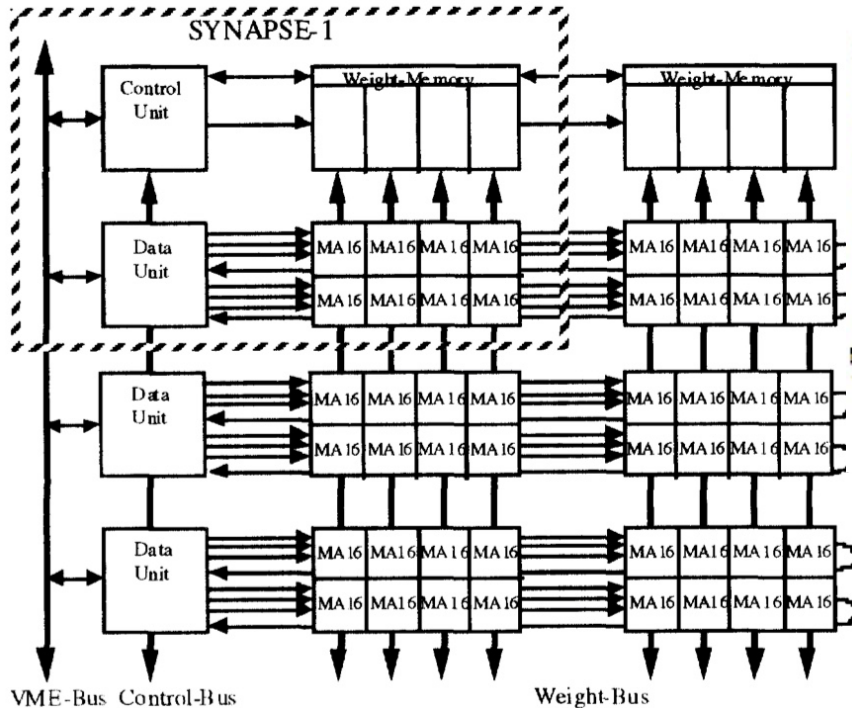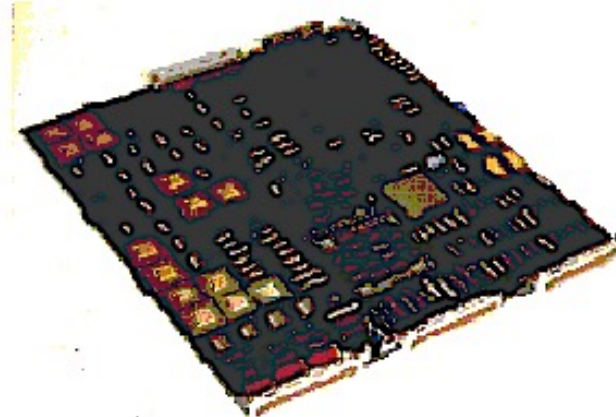Figure 5: Multi-Spert NP training performance with optimized pattern p

*[ Faerber, Asanovic,  IEEE ICAAPP 1997 ]*   **19**

# Siemens SYNAPSE-1 (1992-5)

- Systolic matrix-multiply engine (16b*16b)
- Four levels of program control (68000s + microcode)



SYNAPSE-1



Control Unit

# What Happened to '90s Neurocomputers?

- Very small market
- Neural networks faded in popularity
  - some kept working on them
- Moore's Law scaling favored general-purpose processors
- In 1996, Intel introduced MMX

Programmable Neurocom[puter]

Krste Asanović
MIT Laboratory for Computer Sci[ence]
200 Technology Square
Cambridge, MA 02139
krste@mit.edu

Appears in *The Handbook of Brain Theory and Neural N[etworks]*
(M.A. Arbib, Ed.), Cambridge, MA: The MIT Press, 2002. (c) The MIT Press
http://mitpress.mit.edu

Although the multimedia extensions implemented to date provide only a limited boost to the performance of general-purpose processors on fixed-point matrix code, they signal an intent by commercial microprocessor manufacturers to perform well on these types of code. As commercial design teams incorporate multimedia-style kernels into the workloads they consider during the design of new microprocessors, we can expect performance to increase rapidly also for ANN algorithms. The continuing tremendous investment placed in high-volume microprocessors ensures that these devices will use the most advanced fabrication technologies and the most aggressive circuit design styles yielding the highest clock rates. Given these trends, there will be greatly reduced interest in future special-purpose neurocomputers.

# 40 years of Processor Performance



Performance vs. VAX-11/780

Intel i7-7700k, 4.2 GHz (boosts to 4.5 GHz)
Intel Core i7-6700k 4 cores, 4.0 GHz (boosts to 4.2 GHz)
Intel Core i7-6700k 4 cores, 4.0 GHz (boosts to 4.2 GHz)
Intel Xeon 4 cores, 3.7 GHz (boosts to 4.1 GHz)
Intel Xeon 4 cores, 3.6 GHz (boosts to 4.0 GHz)
Intel Xeon 4 cores, 3.6 GHz (boosts to 4.0 GHz)
Intel Core i7 4 cores, 3.4 GHz (boosts to 3.8 GHz)
Intel Xeon 6 cores, 3.3 GHz (boosts to 3.6 GHz)
Intel Xeon 4 cores, 3.3 GHz (boosts to 3.6 GHz)

Intel D850EMVR motherboard Pentium 4 processor (with hy...

Intel VC820 motherboard Pentium III processor, 1.0 GHz
Professional Workstation XP1000 21264A, 667 MHz
Digital AlphaServer 8400 6/575 21264 575 MHz

AlphaServer 4000 5/600 21164, 600 MHz
Digital Alphastation 5/500, 500 MHz

Digital Alphastation 5/300, 300 MHz

Digital Aphastation 4/266, 266 MHz

IBM POWERstation 100, 150 MHz

Digital 3000 aXP/500, 150 MHz

HP 9000/750, 66 MHz

IBM RS6000/540, 30 MHz
MIPS M2000, 25 MHz
MIPS M/120, 16.7 MHz
Sun-4/260, 16.7 MHz

VAX 8700, 22 MHz

22%/year

VAX-11/785
VAX-11/780, 5 MHz

23%/year

52%/year

Now general-purpose
processor scaling stopped

'90s neurocomputers couldn't
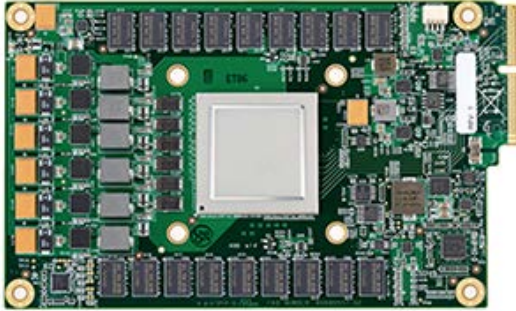compete with scaling

year

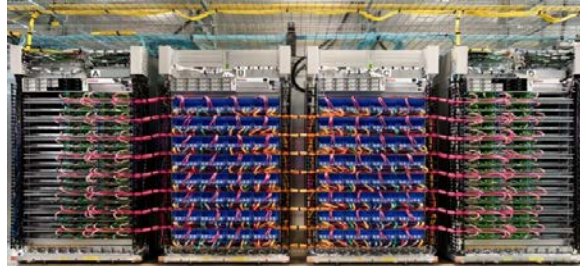**[ Hennessy & Patterson, 2017 ]**

**22**

# General-Purpose GPUs (GP-GPUs)

- In 2006, Nvidia introduced GeForce 8800 GPU supporting a new programming language: CUDA
  - "Compute Unified Device Architecture"
  - Subsequently, broader industry pushing for OpenCL, a vendor-neutral version of same ideas.
- Idea: Take advantage of GPU computational performance and memory bandwidth to accelerate some kernels for general-purpose computing
- Attached processor model:  Host CPU issues data-parallel kernels to GP-GPU for execution
- Over time, became the fastest standard way of performing neural network training
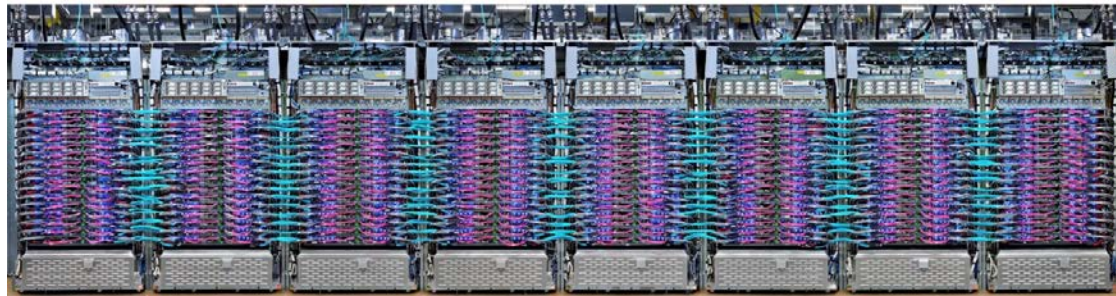
**24**

# Extensive Efforts in Custom AI Chips



Google TPU

TPUv2

TPUv3

*"Let's fill a reticle with reduced-precision vector processing, add high-bandwidth local memory, and attach multiple to a server to accelerate neural network training."*

~100 AI Hardware Startup Pitches

# RISC-V
## (pronounced "risk-five")

**RISC-I**

**RISC II**

**RISC-III (aka SOAR)**

**RISC-IV (aka SPUR)**

**RISC-V (Raven-1, 28nm FDSOI, 2011)**

*Why are outsiders complaining about changes to RISC-V in Berkeley classes?*

27

# What is RISC-V?

- A high-quality, license-free, royalty-free RISC ISA specification originally from UC Berkeley
- Standard maintained by non-profit RISC-V Foundation
- Suitable for all types of computing system, microcontrollers to supercomputers
- Numerous proprietary and open-source cores
- Experiencing rapid uptake in industry and academia
- Supported by growing shared software ecosystem
- A work in progress…

# RISC-V Timeline

Privileged Arch, v1.10

RISC-V Foundation Incorporated

Hot Chips 2014

User ISA v2.0 IMAFD

1st Rocket tapeout, EOS14, 45nm

User ISA v1.0,
Raven-1 tapeout (28nm),
RVC MS thesis

RISC-V ISA project begins

First Linux
Port

1st Workshop

NVIDIA to RISC-V

Commercial Softcores

1st Commercial SoC

WDC to RISC-V

1st Commerical Unix SoC

8th Workshop

NVIDIA

WD

**You are here**

2010  2011  2012  2013  2014  2015  2016  2017  2018  2019

Berkeley                              World

RISC-V Foundation: 180+ Members

**RISC-V Foundation Growth History**
August 2015 to April 2018

Update: ~120 companies and organizations in October 2018

Platinum  Gold  Silver  Auditor  Individual

# RISC-V Vector Extension Overview

vl — *Vector length CSR sets number of elements active in each instruction*

| | | | |
|---|---|---|---|
| v31[0] | v31[1] | | v31[MVL-1] | vetype31 |

*32 vector registers*

| v1[0] | v1[1] | | v1[MVL-1] | vetype1 |
| v0[0] | v0[1] | | v0[MVL-1] | vetype0 |

16 bits vetype/vreg

*Maximum vector length (MVL) depends on implementation, number of vector registers used, and type of each element.*

*Vetype sets type of element in each vector register (e.g., 32-bit integer, 16-bit floating-point)*

# RISC-V 2-D Vector Extension

```
# Matrix multiply
vfmadd.p v0,v1,v2,v0
```

| v2[0][0] | v2[0][1] | | v0[0][MVL-1] |
|----------|----------|---|--------------|
| v2[1][0] | v2[1][1] | | v0[1][MVL-1] |
| **v2** | | | |
| v2[MVN-1][0] | v2[MVN-1][1] | | V0[MVN-1][MVL-1] |

| v1[0][0] | v1[0][1] | | v1[0][MVN-1] |
|----------|----------|---|--------------|
| v1[1][0] | v1[1][1] | | v1[1][MVN-1] |
| **v1** | | | |
| v1[MVM-1][0] | v1[MVM-1][1] | | v1[MVM-1][MVN-1] |

| v0[0][0] | v0[0][1] | | v0[0][MVL-1] |
|----------|----------|---|--------------|
| V0[1][0] | V0[1][1] | | v0[1][MVL-1] |
| **v0** | | | |
| V0[MVM-1][0] | V0[MVM-1][1] | | V0[MVM-1][MVL-1] |

- Vector registers configured as 2D matrices
- Single instructions for matrix multiply/ convolutions

# RISC-V for AI Accelerators

- RISC-V designed originally as basis for custom accelerators
- Simplify software by using one simple base ISA on all cores
  - Where you need high-performance Unix-capable core to run operating system, build a superscalar OoO core
  - Where wanted VLIW for microcode scheduling, build wide in-order superscalar
  - Where wanted low-precision SIMD, use standard vector extensions
  - Where want to take advantage of 2D optimizations (e.g., systolic matrix multiply, convolution), use 2D vector extensions (in progress)
  - Secret-sauce weight compression/number format? add custom extensions!
  - Where need interrupt-responsive I/O management core, build embedded core



- Same memory model, synchronization primitives, compiler tool flow (C-struct packing), debug, tracing,…