# Energy Efficiency of Computing Architectures
## A Deep Dive into Processors and Emerging Computing Machines

Olivier Sentieys
Inria
Univ. Rennes, Irisa        olivier.sentieys@inria.fr

*informatiques mathématiques* **Inria**

**UMR IRISA**

**ENS rennes**   **UNIVERSITÉ DE RENNES 1**

http://people.rennes.inria.fr/Olivier.Sentieys

---

# Cairn Team at a Glance

Inria Rennes

Enssat Lannion

- *Energy-Efficient Computing Architectures*

- ~35 people, Rennes and Lannion campuses
- INRIA, Univ. Rennes 1, ENS Rennes
- Electrical Engineering & Computer Science

- Domain-specific computing architectures
- Design tools and compilers
- Wireless, signal, image, security

2

---

# Energy Efficiency Challenges

- Teraops/Watt?
  - $10^{12}$ op./s/W $\equiv$ 1 pJ/op
  - Several orders of magnitude from current processors and multicores
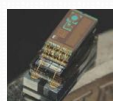- From Sensors to Clouds

  - 1 TOPS @ 1W
    - Clouds, embedded systems

  - 1 GOPS @ 1mW
    - IoT sensor nodes
  - 1 MOPS @ $100\mu W$
    - Energy harvesting

Intel CPU+FPGA     Micro Mote

3

---

# Improving Energy Efficiency

- Technology?
  - What can advanced technology nodes bring
- Accelerate
  - Energy advantages of specialized hardware
- Approximate
  - Playing with accuracy to reduce energy
- Manage the Power
  - Dynamic Voltage/frequency (Over-)Scaling
  - Energy Harvesting sensor nodes

4

# Key Questions

- A deep dive into processors… *(I hope not too deep)*
- Basics on transistors, logic gates, registers, memory
- Energy consumption of processor core/uncore
- Computers are parallel
  - Billions of transistors doing the job at the same time
  - Are multicore processors the solution?
- Specializing the computer
  - Reconfigurable computing
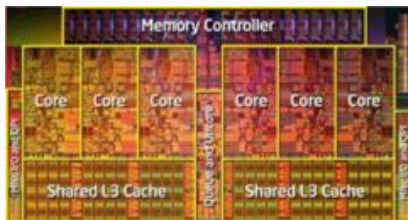- Emerging paradigms
  - Neuromorphic, approximate, stochastic

5

# Outline

- Part I: From Transistors to Logic Gates
  - Basic Element, Delay, Power Consumption
  - The Issue of Synchronization
- Part II: Inside a Processor
  - Von-Neumann Architecture, Instruction Set Architecture, Operating Systems
  - Multicore Processors, Power and Utilization Walls
- Part III: Pushing the Accelerator!
  - Hardware Accelerators
  - Reconfigurable Computing
- Part IV: Emerging Computing Paradigms
  - Neuromorphic Computing
  - Approximate Computing
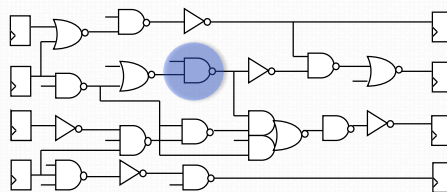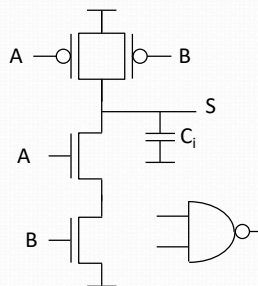  - Chips are going 3D

6

# Integrated Circuit Design

- Chips, logic gates and transistors



Intel's Xeon Chip

```
#pragma hls_design top
void my_design (int *a, int *o) {
    static int i,j;
for(i=1; i<=n-1; i++)
  for(j=1; j<=n-1; j++)
    a[i][j] = (a[i-1][j]+a[i][j]+a[i][j-1])/3.0;
}
```
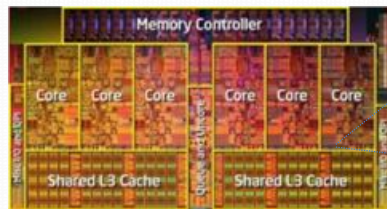```
pr
beg
 i
 end if;
end process;
```

# Part I: From Transistors to Logic Gates

- The Fundamental Element: MOSFET Transistor
- Design of CMOS Cells: Combinatorial Logic
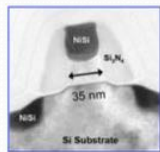- Memory Cells
- Delay and Power Consumption
- Synchronous Design
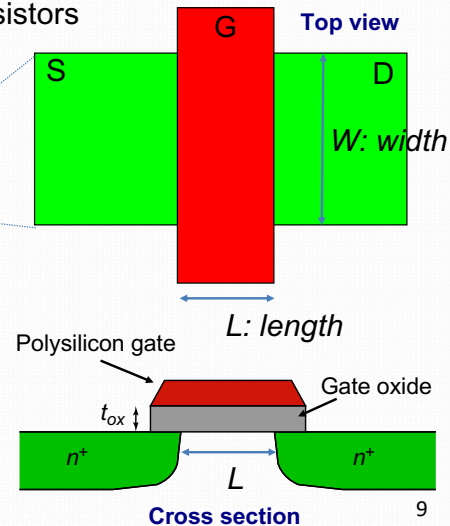
8

# Fundamental Building Block: MOSFET Transistor

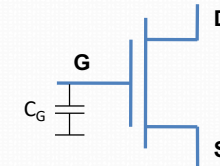Now several billions or transistors


Intel's Xeon Chip

**Top view**

S    G    D

*W: width*

*L: length*

Polysilicon gate

$t_{ox}$

Gate oxide

$n^+$    $L$    $n^+$

**Cross section**

MOSFET: Metal Oxide Silicium Field Effect

# The Basic Element: Transistor

- Transistor as a switch

- Vgs > Vt: NMOS on
  — Resistance $R_{DS}$

  D ———○———○——— S

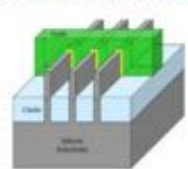- Vgs < Vt: NMOS off
  — Leakage $I_{off}$

  D ———○——/——○——— S

$C_G$   G

D

S

Ids

$I_{off}$

Vgs

**Vt**: threshold voltage

- Gate: capacitance $C_G$
- Switch: resistance $R_{DS}$

# Transistors Nowadays

- Intel FinFET: transistors go 3D



- Fully Depleted SOI[1]
  — Low-power

[1]Silicon on Insulator
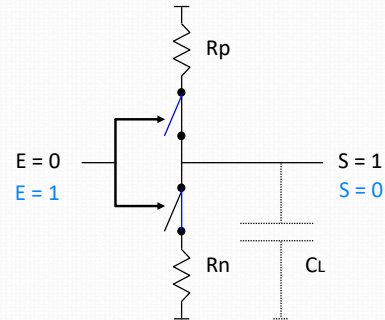
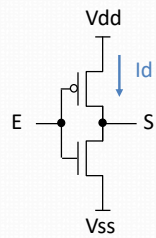# Part I: From Transistors to Logic Gates

- The Fundamental Element: MOSFET Transistor

- Design of CMOS Cells: Combinatorial Logic
- Memory Cells
- Delay and Power Consumption
- Synchronous Design
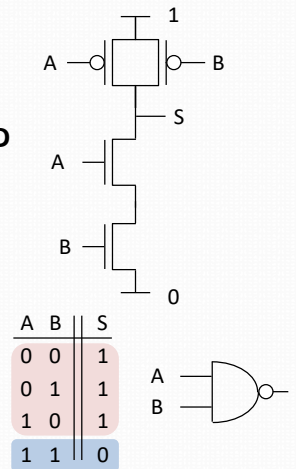
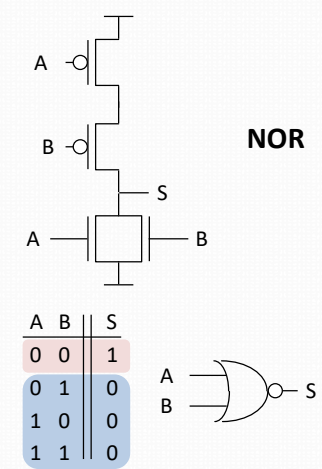# Combinatorial Logic Cells
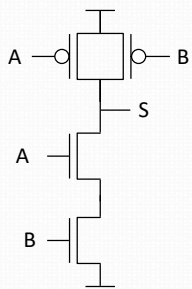
- Complementary Logic (CMOS)

CMOS Inverter

Vdd

Id

E — S

Vss

Rp

E = 0
E = 1

S = 1
S = 0

Rn    $C_L$

# NAND and NOR

**NAND**

1

A — B

S

A

B

0

| A | B | S |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

A
B — S

**NOR**

A

B

S

A — B

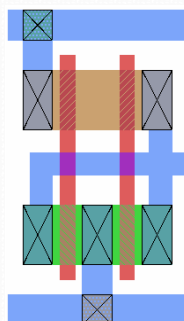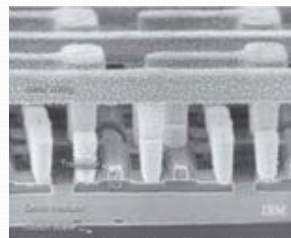| A | B | S |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 0 |

A
B — S

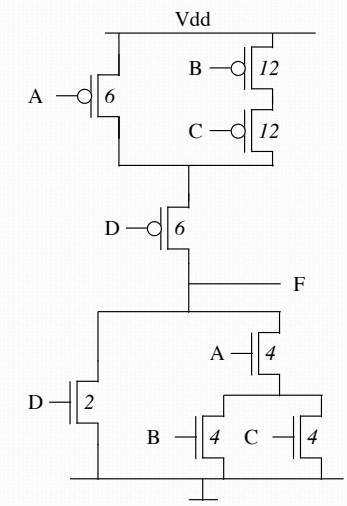# Layout Design

Transistor
Schematic

A — B

S

A

B

Layout

Silicon

# Complex Gates

- $F = \overline{A.(B+C) + D}$

- The art of transistor sizing

  – Equilibrate delay for $0 \rightarrow 1$ and $1 \rightarrow 0$ output transitions
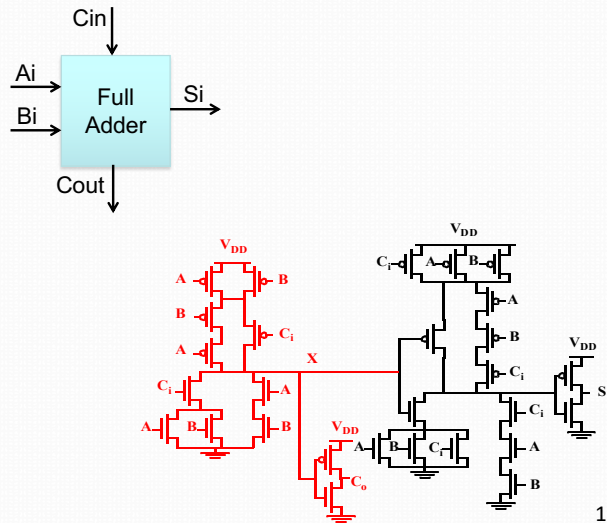  – Minimize cell area

Vdd

B — *12*

A — *6*

C — *12*

D — *6*

F

A — *4*

D — *2*

B — *4*   C — *4*

# Complex Gates: Full Adder

- Full Adder



Full Adder block: inputs Ai, Bi, Cin; outputs Si, Cout

| Ai | Bi | Ci | Co | Si |
|----|----|----|----|----|
| 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 1 | 1 | 1 |

17

# Complex Functions
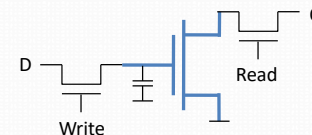
- 16-bit Adder (integer)



18

# Part I: From Transistors to Logic Gates

- The Fundamental Element: MOSFET Transistor
- Design of CMOS Cells: Combinatorial Logic
- Memory Cells
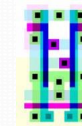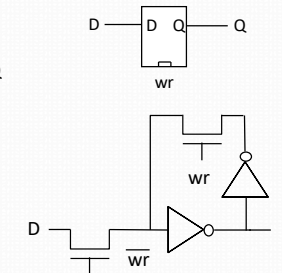- Delay and Power Consumption
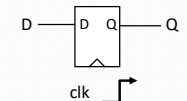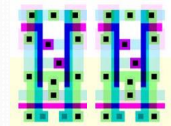- Synchronous Design

19

# Storing Values

| Capacitor (DRAM) | Latch (SRAM) | Flip-Flop (Register) |
|---|---|---|


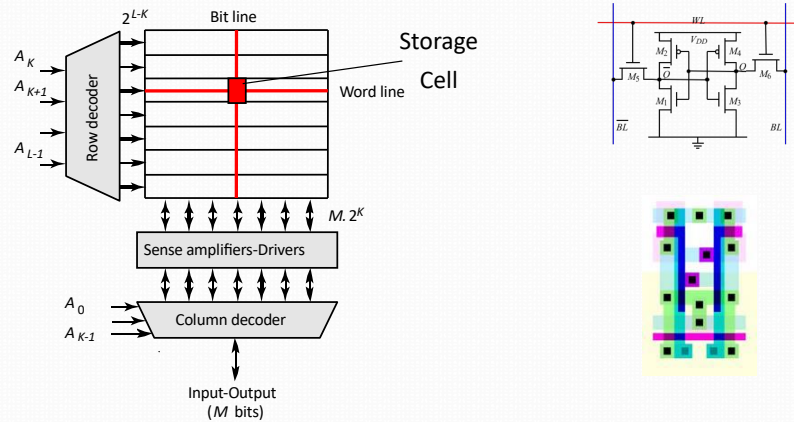
- Setup Time: Tsetup
- Hold Time: Thold
- Propagation Time: Tp

20

## Memory

- L2 Cache contains 4 Millions SRAM cells
  - Raw/column of 2000 cells



Bit line

Storage Cell

Word line

Row decoder

Sense amplifiers-Drivers

Column decoder

Input-Output ($M$ bits)

$2^{L-K}$

$M.2^K$

$A_K$, $A_{K+1}$, $A_{L-1}$, $A_0$, $A_{K-1}$

---

## Delay: Parasitic Elements



- Drain-Source Resistance: $R_{DS} = \dfrac{L}{W}\dfrac{1}{k(V_{dd}-V_t)}$
- Gate Capacitance: $C_g = \dfrac{\epsilon W.L}{t_{ox}} = W.L.C_{ox}$

$$\text{Delay} \propto R_{DS}.C_g \propto \frac{L^2}{Vdd - Vt}$$

---

## Power and Energy Consumption

- Dynamic power
  - Charge and discharge of node capacitance
- Energy = $C.Vdd^2$
- Power

$$P_{dyn_i} = C.Vdd^2.f.Prob_{0\to1}$$

- Static power: *Ps*
  - Sub-threshold and junction leakage current

$$P_{stat_i} = N.I_{off}.Vdd$$



$Vdd$
$Idd = Isc + Ic$
$Isc$ $Ic$
$C$

Id
$I_{off}$
Vg
Vt (low)

---

## Power at Higher Level

- Propagating activity



$$P = \sum_i \left[\alpha_i.f_i.C_i.Vdd^2 + I_{leak_i}.Vdd\right]$$

## Activity

- Activity $\alpha_i$ is the probability to have a *0→1* transitions at the output of a gate
- Example: AND gate
  - $P_S = P(S=1) = P_A P_B$
  - $\alpha_i = P_S(1 - P_S)$

| A | B | S |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

$\alpha = 3/16$

- Activity propagation

25

## Propagating Activity is not So Simple

- Conditional probabilities

- Glitches: gate delay
  - Significant in arithmetic

Glitches

26

## Dynamic Power vs. Performance

- Decreasing Vdd reduces power but increases delay

$$P_{dyn_i} = \alpha_i . f_{clk} . C_i . Vdd^2$$

$$Delay \propto \frac{1}{V_{dd} - V_t}$$

27

## Leakage vs. performance

- High performance

- Low leakage

$$P_{stat_i} = N.I_{off}.Vdd$$

$$Delay \propto \frac{1}{V_{dd} - V_t}$$

Ioff:
- Exponential in inverse of Vt
- Exponential in temperature
- Linear in device count

28

## Minimum Energy per Operation

- Putting all together

## On-Chip Interconnect?

- Gate delay decreases but... wire delay increases
- Crossing chip in 5-10 clock cycles
- Also affected by noise...



- Metal layers to reduce wire delay
- Repeaters

- Towards network-on-chip

## Conclusion: Power in CMOS

$$P = \sum_i \left[ \alpha_i . f_i . C_i . Vdd^2 + I_{leak_i} . Vdd \right]$$

- Dynamic power
  - 40-70% today
  - Decreasing relatively
  - DVFS becomes more and more difficult

- Leakage power
  - 20-50 % today
  - Increasing rapidly
    - number of transistors
    - Vdd/Vt scaling
  - Critical for memory

$$P = \frac{energy}{operation} \times rate + static\ power$$

# Inside (Simple) Processor Architecture

# Von Neumann Computers

- Processing address, data, control, on the same resources
- Single memory for data and program
- Sequential behavior

- Practically, most processors use Harvard model: separated data and program memory

CPU (Central Processing Unit)

Memory

# Instruction Set Architecture (ISA)

- ISA defines a programmer's interface
- Each instruction is defined by coding (binary) and semantics

Select Register in RF

| 0 | 1 | 0 | 1 | 0 | 1 | B = READ(@A) |
| 1 | 1 | 0 | 0 | 0 | 1 | WRITE(B,@A) |
| 0 | 1 | 1 | 0 | 1 | 0 | C = A + B |
| 1 | 1 | 0 | 1 | 0 | 1 | ... |

| 31 | 30 | 29 | 28 | 27 | 26 | 25 | 24 | 23 | 22 | 21 | 20 | 19 | 18 | 17 | 16 | 15 | 14 | 13 | 12 | 11 | 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A[4:0] | | | | | B[4:0] | | | | | C[4:0] | | | | | OPX[10:0] | | | | | | | | | | | OPCOD[5:0] | | | | | |

# Microarchitecture Pipeline

- Microarchitecture defines how instructions are executed (not unique)

# Execution of an instruction involves

1. Instruction fetch
2. Decode and register fetch
3. ALU operation
4. Memory operation (optional)
5. Write back (optional)

and compute address of next instruction

## Slide (top-left)

Fetch Instruction

Decode Instr.
Load Registers

Read/Write Memory

Execute Instr.

Write Back to
Registers

Compute next address



---

## Achieving Higher Performance

- Branch/value prediction
- Cache memory
- In-core parallelism
  - Multiple Fus
  - Out of order execution
  - VLIW+good compilers

- Multiple cores on a single chip

---

## Abstraction in Computer Systems

- Maximum of an array T

```
numpy.amax(T)
```

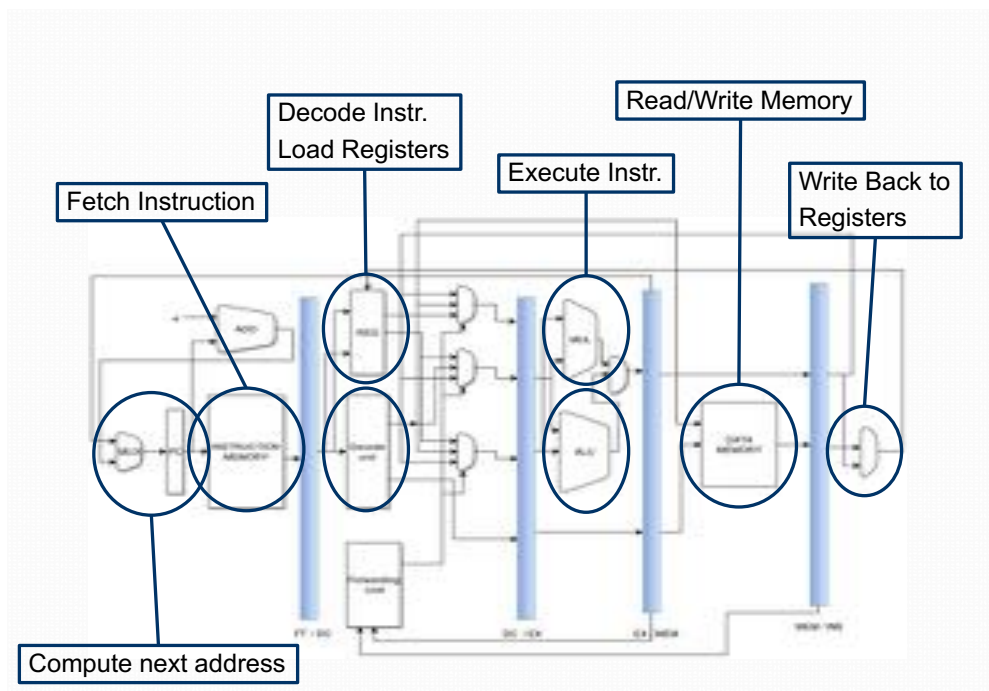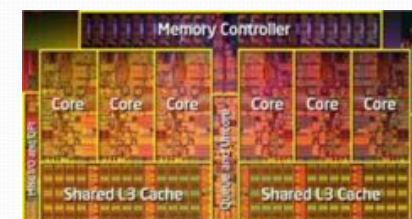| 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 |
| 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 |

```
int largest(int T[], int length) {
  int max = T[0];
  for(i=1; i<length; i++) {
      if (max < T[i]) {
            max = T[i];
      }
  }
  return max;
}
```

```
         R1 ← *R2      // max
loop     R2 ← R2+1     // T[]
         R3 ← *R2
         R1 < R3 ?
         BZ next
         R1 ← R3
next     B loop
```

---

## Abstraction and Performance?

- Matrix Multiply: relative speedup to a Python version (18 core Intel)



Matrix Multiply Speedup Over Native Python

["There's Plenty of Room at the Top," Leiserson, et. al.]

## Energy Cost in a Processor

- Operation:
  - 32-bit addition: 0.05pJ
  - 16-bit multiply: 0.25pJ
  - 64-bit FPU: 20pJ/op

- Instruction:
  - fetch, decode, read 2 operands from RF, execute, write back
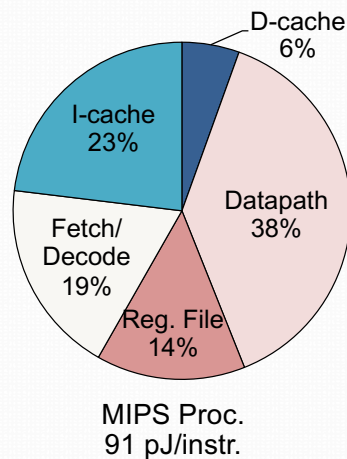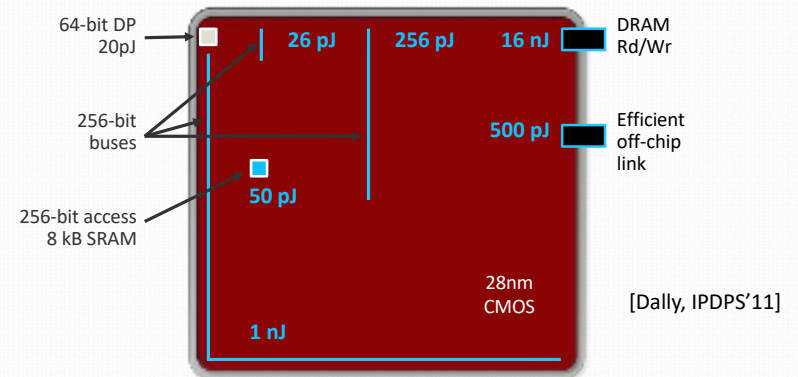


D-cache 6%
I-cache 23%
Datapath 38%
Fetch/Decode 19%
Reg. File 14%

MIPS Proc. 91 pJ/instr.

---

## Energy Cost in a Processor

- Fetching operands costs more than computing



64-bit DP 20pJ → 26 pJ   256 pJ   16 nJ   DRAM Rd/Wr
256-bit buses
256-bit access 8 kB SRAM   50 pJ   500 pJ   Efficient off-chip link
1 nJ   28nm CMOS   [Dally, IPDPS'11]

---

# Multicore: it's all a trick!
## Power and Utilization Walls

---

## And then came the "Power Wall"



Power Density: 100 W/chip (~25W/cm²) is a limit

Source: C. Batten, Cornell

## and the "Multicore Era"

- Increasing performance by increasing # of cores

**256 Cores**
- 4-way SIMD FMACs @ 2.5–5 GHz
- 5–10 TFlops on one chip
- Some apps require 1 byte/flop
- Need 5–10 TB/s of off-chip BW
- Need 5–10 TB/s of on-chip BW too!

Number of Cores axis: 512, 256, 128, 64, 32, 16, 8, 4, 2, 1

Manycore Era

Labels: Intel TFlops, Intel TFlops, Tilera TILE64, NVIDIA GT200, Cavium Octeon, MIT RAW, Raza XLR, Rock, Cell, Niagra, Nehalem, Barcelona, Nehalem, Power4, Opteron, XBox360, Core2, Power6, 286, 386, 486, Pentium, P2, P3, P4, Athalon, Itanium

Years axis: 1980 1985 1990 1995 2000 2005 2010 2015 2020

–Source: C. Batten, Cornell

56

---

## Moving to multicore

- 1 core@2GHz@1.2V@1W

  2GHz | 1W 1.2V

- 1 core@1GHz@0.8V@0.25W

  1GHz | 0.22W 0.8V

- 2 cores@1GHz@0.8V@0.5W

  1GHz

- But… twice area (and not so simple)

  1GHz

- **Advanced technology nodes?**

62

---

## Technology Scaling

Intel's Xeon Chip

$S$

**28 nm**    **20 nm**    **14 nm**

**Classical (Dennard's) scaling**

| Device count | $S^2$ |
|---|---|
| Device frequency | $S$ |
| Capacitance, Vdd | $1/S$ |
| Device power | $1/S^2$ |
| **Utilization** | **1** |

$Core_i$

$Core_i$

$50W@1.4.f$

$100W@f$

63

---

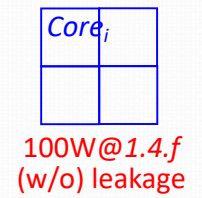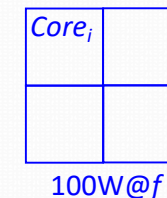## End of Dennard's Scaling

- Energy efficiency is not scaling along with integration capacity

**Leakage limited scaling**

| Device count | $S^2$ |
|---|---|
| Device frequency | $S$ |
| Device power (cap) | $1/S$ |
| **Device power ($V_{dd}$)** | **~1** |
| **Utilization** | **$1/S^2$** |

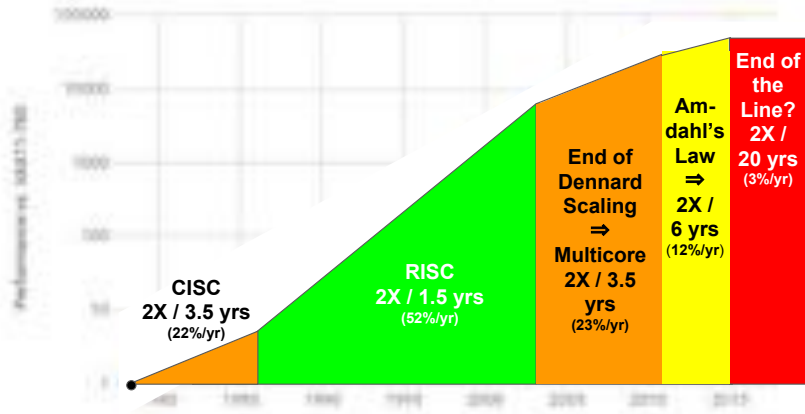$Core_i$

$Core_i$

$100W@f$

$100W@1.4.f$
(w/o) leakage

- **Utilization Wall**: percentage of a chip that can switch at full frequency drops exponentially
- Replace dark cores with **specialized cores** (10-100x more energy efficient)

64

## End of Growth of Speed?
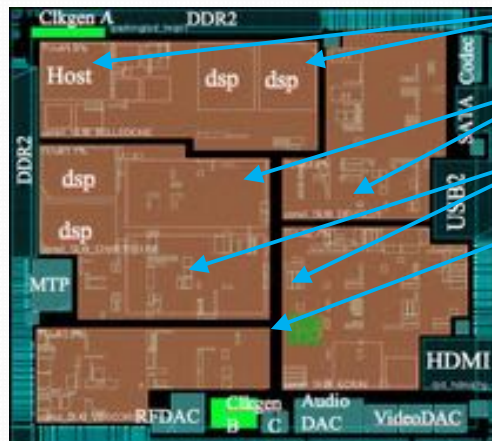
40 years of Processor Performance

End of
the
Line?
2X /
20 yrs
(3%/yr)

Am-
dahl's
Law ⇒
2X /
6 yrs
(12%/yr)

End of
Dennard
Scaling
⇒
Multicore
2X / 3.5
yrs
(23%/yr)

RISC
2X / 1.5 yrs
(52%/yr)

CISC
2X / 3.5 yrs
(22%/yr)

Based on SPECintCPU. Source: John Hennessy and David Patterson, Computer Architecture: A Quantitative Approach, 6/e. 2018

65

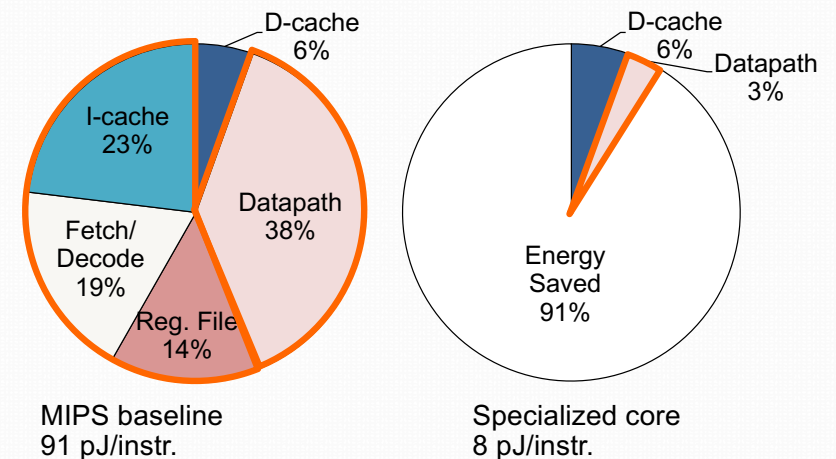## Part III: Pushing the Accelerator!

66

## What is a HW accelerator?

- 16 processors
- 38 HW blocks
- 140 memory blocks
- 5 Gbytes/s on-chip interconnection network

67

## Energy Savings in Specialized HW

D-cache
6%

I-cache
23%

Datapath
38%

Fetch/
Decode
19%

Reg. File
14%

MIPS baseline
91 pJ/instr.

D-cache
6%

Datapath
3%

Energy
Saved
91%

Specialized core
8 pJ/instr.

[Goulding et al., Hot Chips'10]

68

## An example: Bitcoin Mining

| Type | Model | Mhash/s | Mhash/J | Power (W) |
|------|-------|---------|---------|-----------|
| GPP | Intel Xeon X5355 (dual) | 22.76 | 0.09 | 120 |
| GPP | ARMCortex-A9 | 0.57 | 1.14 | 1.5 |
| GPP | Intel Core i7 3930k | 66.6 | 0.51 | 130 |
| GPU | AMD 7970x3 | 2050 | 2.41 | 850 |
| GPU | Nvidia GTX460 | 158 | 0.66 | 240 |
| ASIC | AntMiner S1 | 180.000 | **500** | 360 |
| ASIC | AntMiner S5 | 1.155.000 | **1957** | 590 |
| FPGA | Bitcoin Dominator X5000 | 100 | **14.7** | 6.8 |
| FPGA | Butterflylabs Mini Rig | 25.200 | **20.16** | 1250 |

69

## Making ANN Inference more Efficient

- Main motivation: AlphaGo consumes around 250,000 Watts!
- Bring Logic and Memory closer
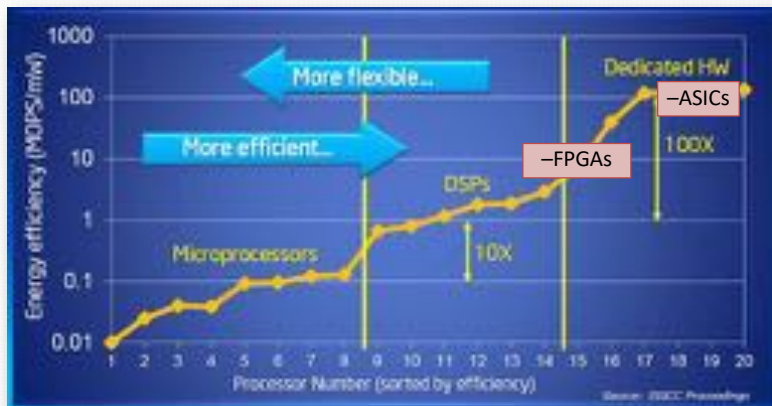- Compute less precisely



Tensor Processing Unit

ASIC for TensorFlow
Designed by Google
10x better perf / watt
latency and efficiency
bit quantization

- Google Tensor Processing Units (TPU)
  - Computations close to memory
  - 8 bit operations

70

## The Efficiency of Specialization



−* Source: Ning Zhang and Bob Brodersen, ISSCC data

−100-1000X Gap in Efficiency … but Specialization comes with Penalties in Programmability

71

# Reconfigurable Hardware Accelerators

75

# Field Programmable Gate Array (FPGA)
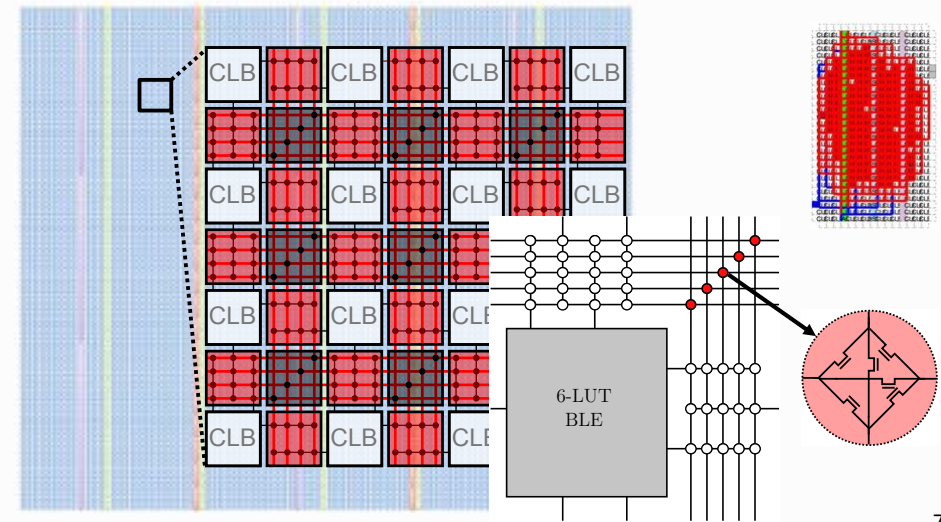


**>4K Multipliers/Adders**

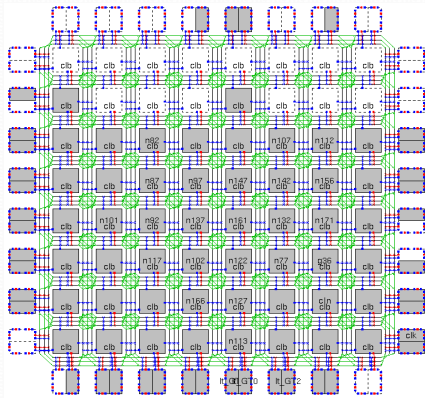**>50MB RAM blocks**

**>2M Configurable Logic Blocks**

76

# Field Programmable Gate Array (FPGA)
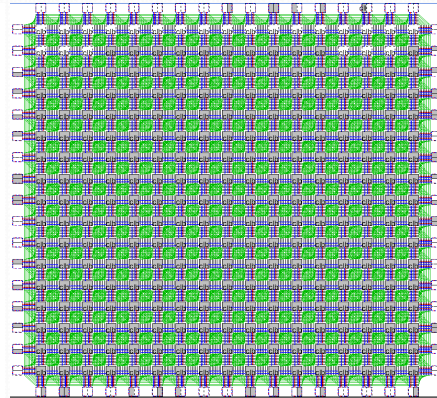


CLB CLB CLB CLB

6-LUT BLE
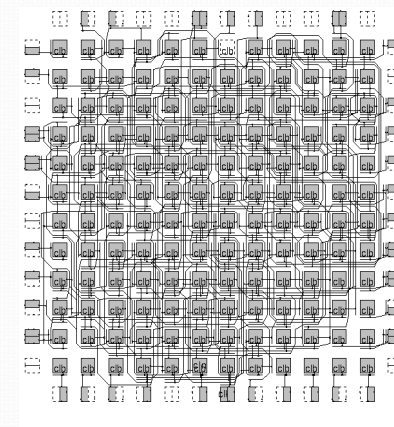
77

# The Program is the Configuration
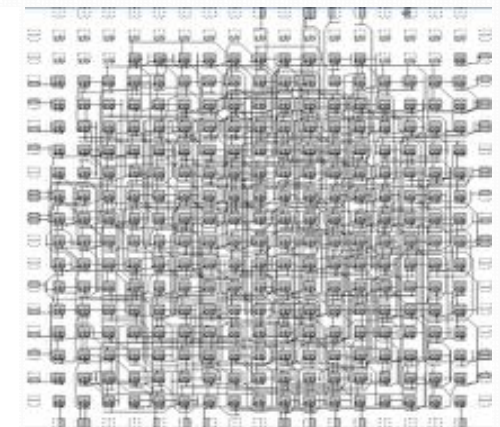


(a) abs    (b) calcNeighbor

78

# The Program is the Configuration



(a) Crc16    (b) calcNeighbor

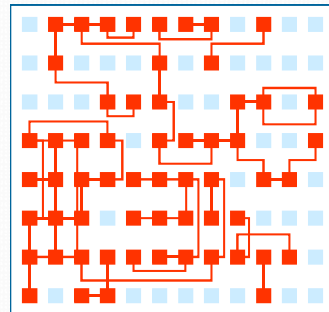79

# Space-Time Computation

```
for(i=1; i<length; i++) {
    if (max < T[i]) {
        max = T[i];
    }
}
```
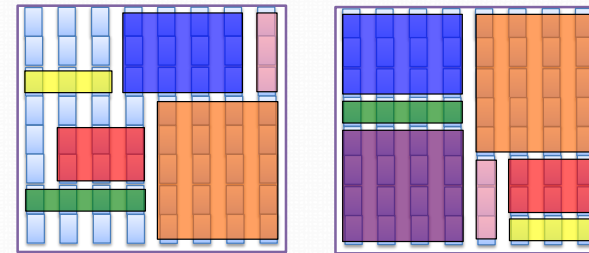
```
for(i=1; i<N; i++) {
    for(j=1; j<M; j++) {
        y[i][j]+=x[i][j]*h[j][i]
    }
}
```

# FPGA Acceleration

- FPGAs can run multiple tasks in parallel
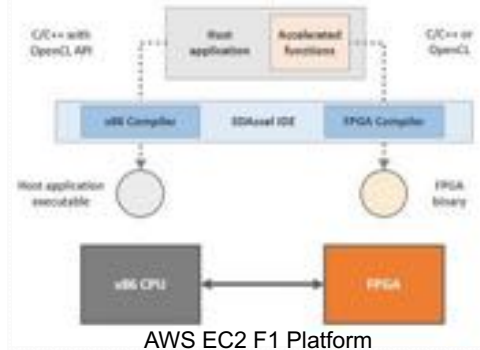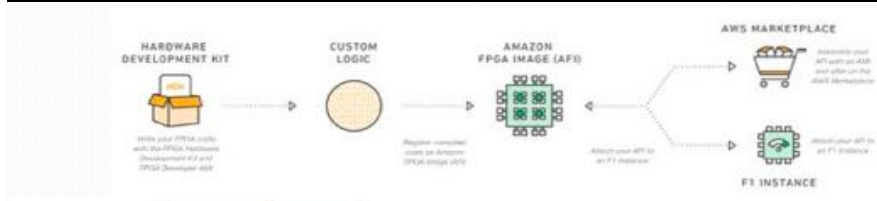
FPGA accelerators for HPC/Cloud

- Towards heterogeneous multicores

# Amazon AWS EC2 F1

aws
aws.amazon.com

| Instance Size | FPGAs | DDR-4 (GiB) | vCPUs | Instance Memory (GiB) | NVMe Instance Storage (GB) | Network Bandwidth |
|---|---|---|---|---|---|---|
| f1.2xlarge | 1 | 4 x 16 | 8 | 122 | 1 x 470 | Up to 10 Gbps |
| f1.16xlarge | 8 | 32 x 16 | 64 | 976 | 4 x 940 | 25 Gbps |

- Up to 8 Xilinx UltraScale+ FPGA devices in a single EC2/F1 instance

AWS EC2 F1 Platform

# Time has Come for Specialization

- Microsoft Unveils Catapult to Accelerate Bing
  - One FPGA per blade
  - 6 × 8 2-D torus topology
  - High-end Stratix V FPGAs
- Running Bing Kernels for feature extraction and machine learning
- Increase ranking throughput by 95% at comparable latency to software-only
- Increase power consumption by 10%
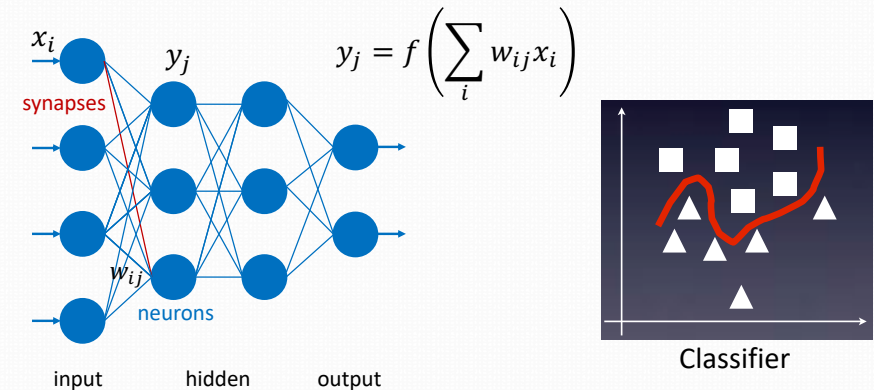- Increase total cost of ownership by less than 30%

# Part IV: Emerging Computing Paradigms

---
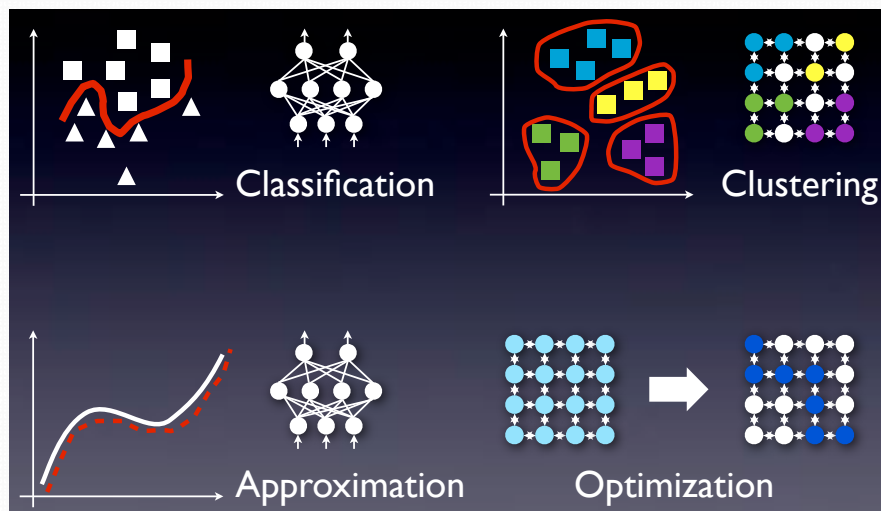
# How Do Artificial Neural Networks Work?

$$y_j = f\left(\sum_i w_{ij} x_i\right)$$

$x_i$ synapses  $y_j$  neurons

input   hidden   output

Classifier

- Neural networks are not fundamentally complicated
- The issue: finding the good weights with *learning*

---

# What ANNs Can Do

Classification   Clustering

Approximation   Optimization

[O. Temam, ISCA10]

---

# So What's New?

## Convergence of trends

- Computer performance (e.g. GPU) can train neural networks with millions of weights
- Access to gigantic datasets
  - Billions of images
  - Training can take weeks!
- More complex ANNs

Imagenet

  - Deep Convolutional Neural Networks (CNN)
  - Long Short-Term Memory (LSTM) Recurrent Neural Networks
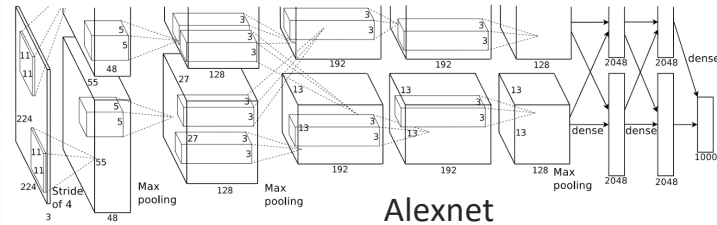- Trendy vision applications
- Emerging technologies offer opportunities

# So What's New?

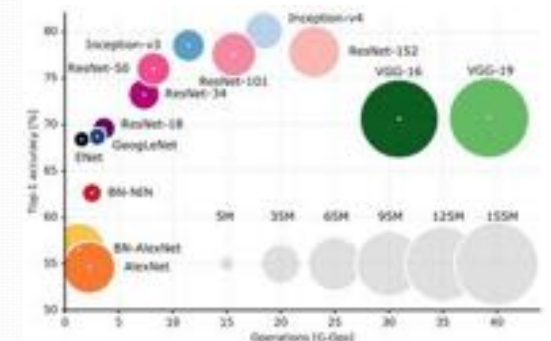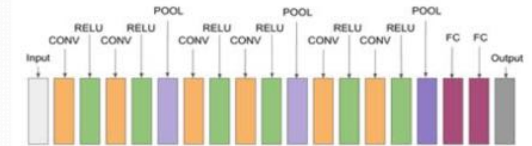- **Deeper Networks**



Alexnet



Alex Krizhevsky *et al.*, Imagenet classification with deep convolutional neural networks, 2012.

89

# Complexity of Deep CNNs

- 10-30 GOPS
  - Mainly convolutions
- 10-200 MB
  - Fully-connected layers



# And What About Energy?

- The brain seems to have something very special about energy efficiency

Lee Sedol (brain)
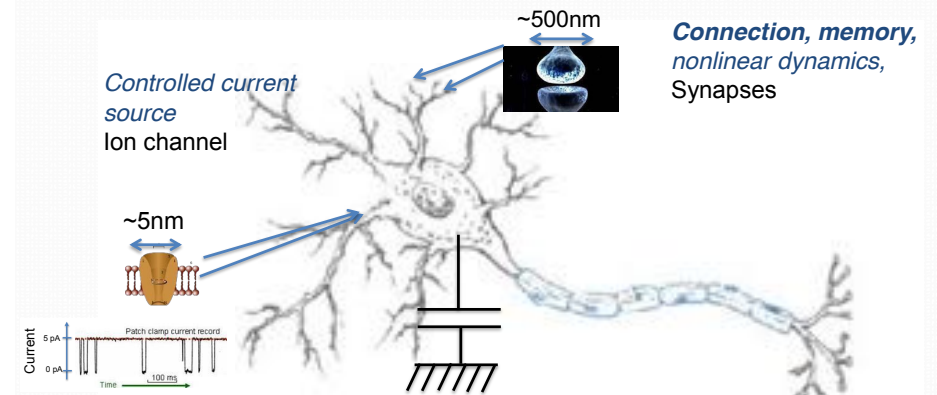
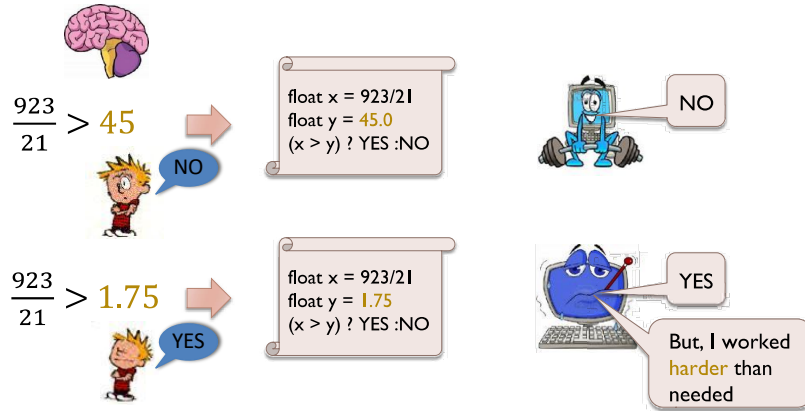AlphaGo (CPU+GPU with tree seach and deep neural networks)



**20 Watt**

**>250 000 Watt**

- Computers: arithmetic but chiefly memory transfers

# Real Biological Neurons

~500nm

***Connection, memory,*** *nonlinear dynamics,* Synapses

*Controlled current source*
Ion channel

~5nm



- Brain computes with strong approximations (mostly analog) based on low power, slow, noisy and variable nano-devices

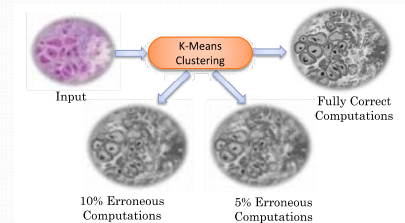# Humans Approximate....
# But Computers Do Not!

$$\frac{923}{21} > 45$$

NO

```
float x = 923/21
float y = 45.0
(x > y) ? YES :NO
```

NO

$$\frac{923}{21} > 1.75$$

YES

```
float x = 923/21
float y = 1.75
(x > y) ? YES :NO
```

YES

But, I worked harder than needed

- Leads to inefficiency
- Overkill (for many applications)

93

# Many Applications are Error Resilient

- Produce outputs of acceptable quality despite approximate computation
  - Perceptual limitations
  - Redundancy in data and/or computations
  - Noisy inputs

- Digital communications, media processing, data mining, machine learning, web search, …
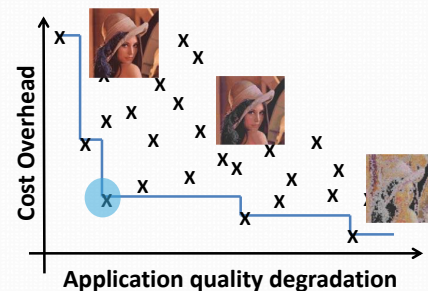
Input  →  K-Means Clustering  →  Fully Correct Computations

10% Erroneous Computations    5% Erroneous Computations

e.g. Image Segmentation        94

# Approximate Computing

- Play with **approximations** to reduce **energy** and increase execution speed while keeping **accuracy in acceptable limits**
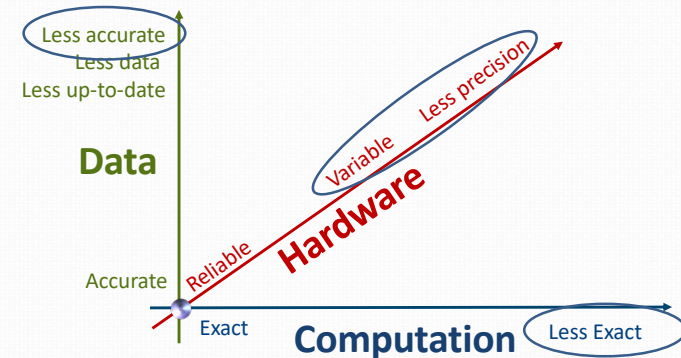  - Relaxing the need for fully precise operations

- Design-time/run-time
- Abstraction levels

Cost Overhead

Application quality degradation

95

# Approximate Computing

- Three dimensions to explore

Less accurate
Less data
Less up-to-date

**Data**

Less precision
Variable
**Hardware**
Reliable

Accurate

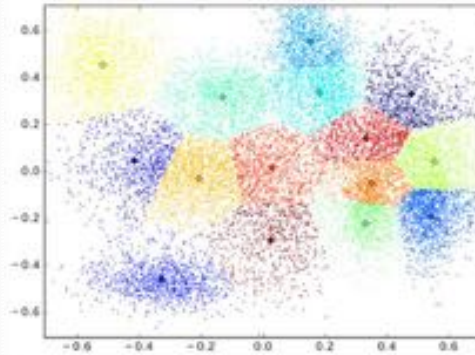Exact   **Computation**   Less Exact

*Note: Precision (#bits) ≠ Accuracy (quality)*

96

# K-Means Clustering

- Data mining, image classification, etc.
- A multidimensional space is organized as:
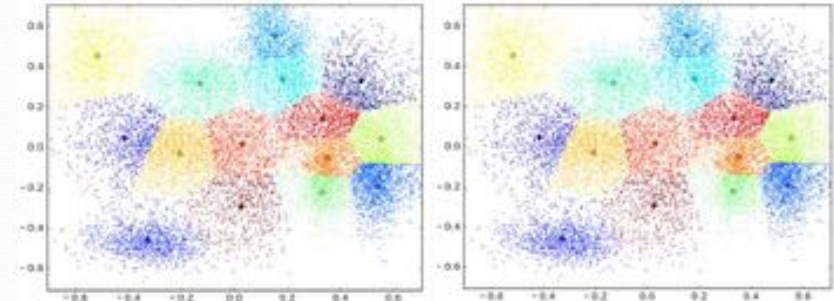  - $k$ clusters $S_i$,
  - $S_i$ defined by its centroid $\mu_i$



- Finding the set of clusters $S = \{S_i\}_{i \in [0, k-1]}$ satisfying $\arg\min_S \sum_{i=1}^{k} \sum_{x \in S_i} \|x - \mu_i\|^2$ is NP-hard (solved here by Lloyd's iterations)

---

# Approximate K-Means Clustering

- W = 16 bits, accuracy = $10^{-4}$
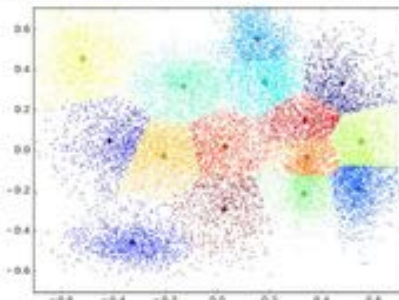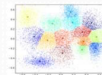- No major (visible) difference with reference



Reference: double

Floating-point: ct_float$_{16}$
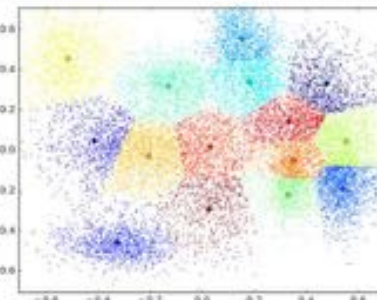5-bit exponent
11-bit mantissa

---

# Approximate K-Means Clustering

- W = 16 bits, accuracy = $10^{-4}$
- No major (visible) difference with reference



Fixed-Point: ac_fixed$_{16}$
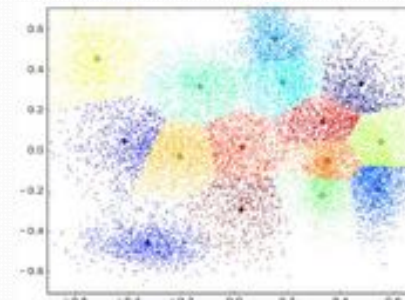3-bit integer part
13-bit fractional part

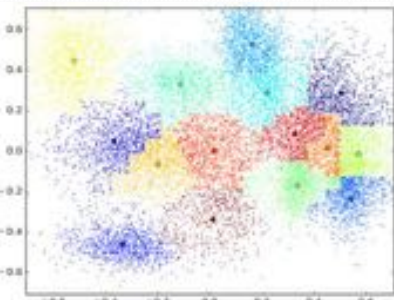Floating-point: ct_float$_{16}$
5-bit exponent
11-bit mantissa

---

# Approximate K-Means Clustering

- W = 8 bits, accuracy = $10^{-4}$
- 8-bit float is still practical
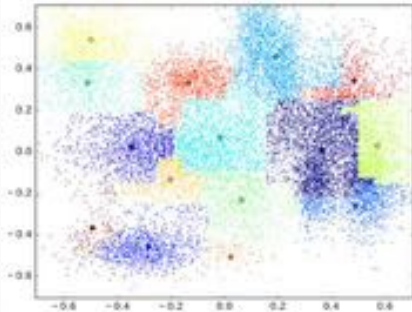


Reference: double

Floating-Point: ct_float$_8$
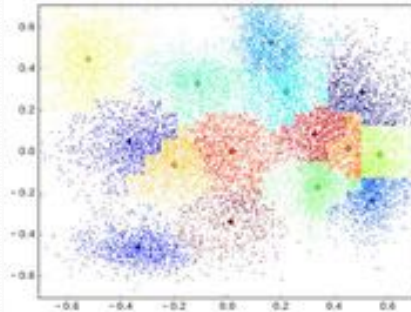5-bit exponent
3-bit mantissa

# Approximate K-Means Clustering

- $W = 8$ bits, accuracy $= 10^{-4}$
- 8-bit float is better and still practical



Fixed-Point: ac_fixed$_8$

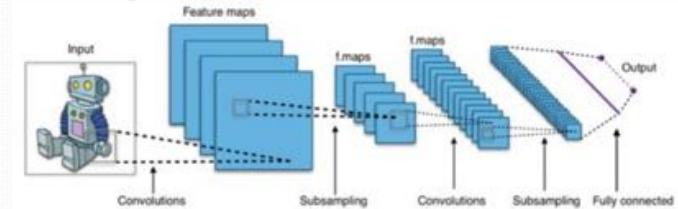3-bit integer part
5-bit fractional part

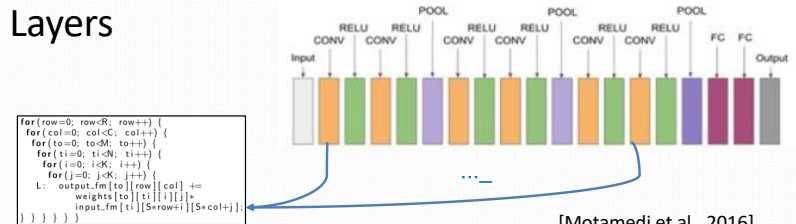Floating-Point: ct_float$_8$

5-bit exponent
3-bit mantissa

# Deep Convolutional Neural Networks

- General organization



- Layers
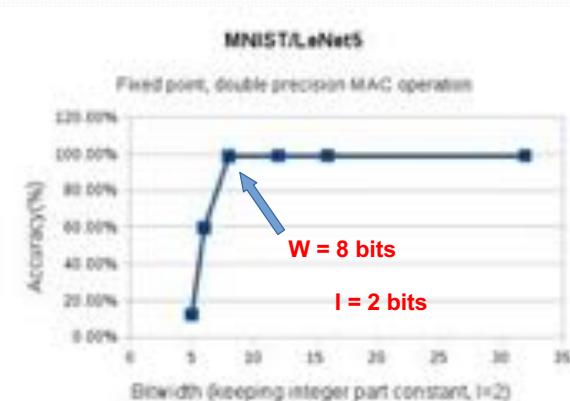


[Motamedi et al., 2016]

# Resilience

**Aoccdrnig to a rscheearch at Cmabrigde Uinervtisy, it deosn't mttaer in waht oredr the ltteers in a wrod are, the olny iprmoatnt tihng is taht the frist and lsat ltteer be at the rghit pclae. And we spnet hlaf our lfie larennig how to splel wrods. Amzanig, no!**

- Our biological neurons are fault tolerant to computing errors and noisy inputs

# Approximate CNNs
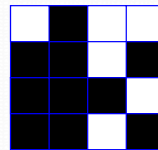
- 10k images, MNIST/Lenet
- Fixed-Point Arithmetic



MNIST/LeNet5

Fixed point, double precision MAC operation

$W = 8$ bits

$I = 2$ bits

# Summary



- Energy consum
  - True in embed
  - True in HPC, n

- End of Moore's law…

- Multicores but utilization wall
  - Percentage of a chip that can switch at full frequency drops exponentially
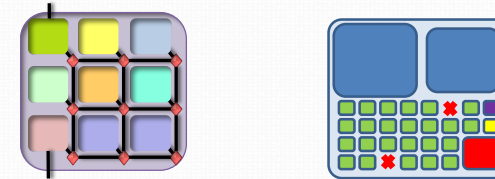


Dark Silicon

105

# What's next?

- Dark Silicon is also an opportunity
  - Heterogeneous manycore architectures



- Efficiency of hardware specialization
  - Domain-specific architectures and languages
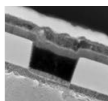- Computing just right
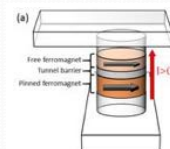  - @design-time or @run- time

106

# What's next?

- Emerging devices

- Cells, brain, neurons have "analog" behavior

- And compute with very low precision

- Making neuromorphic computing more efficient

 *Phase Change Memory*

 *Memristors, Oxide Resistive Memory*

 ***Spin Torque Magnetic Memory***

107