

# ESE5320: System-on-a-Chip Architecture

Day 23: November 22, 2022  
Energy



Penn ESE5320 Fall 2022 -- DeHon

1

## Today

### Energy

- Part 1
  - Today's bottleneck
  - What drives
- Part 2: Architecture and Energy
  - Processors, FPGAs, accelerators
- Part 3:
  - How does parallelism impact energy?

Penn ESE5320 Fall 2022 -- DeHon

2

2

## Message

- Energy dominates
  - Including limiting performance
- Make memories small and wires short
  - Small memories cost less energy per read
- Accelerators reduce energy
  - Compared to processors
- Can tune parallelism to minimize energy
  - Typically, the more parallel implementation costs less energy

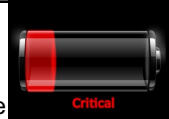
Penn ESE5320 Fall 2022 -- DeHon

3

3

## Energy

- Growing domain of portables
  - Less energy/op → longer battery life
- Global Energy Crisis
- Limit to the power our computers/devices can dissipate (we can provide)
  - Reduce E/op
    - increase compute when power limited
  - Scaling
    - Power density **not** transistors limit sustained ops/s
  - Server rooms
    - Cost-of-ownership **not** dominated by Silicon



Penn ESE5320 Fall 2022 -- DeHon

4

4

## Preclass 1--4

- 200K gates/mm<sup>2</sup>
- 5\*10<sup>-15</sup> J/gate switch
- Gates on 1cm<sup>2</sup>
- Energy to switch all?
- Power at 1GHz?
- Fraction can switch with 10W/cm<sup>2</sup> power budget?

Penn ESE5320 Fall 2022 -- DeHon

5

5

## Challenge: Power

Penn ESE5320 Fall 2022 -- DeHon

6

6

## Origin of Power Challenge

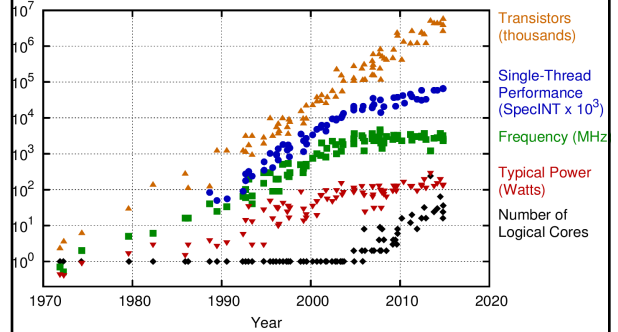
- Limited capacity to remove heat
  - ~100W/cm<sup>2</sup> force air
  - 1-10W/cm<sup>2</sup> ambient
- Transistors per chip grow at Moore's Law rate =  $(1/F)^2$
- Energy/transistor must decrease at this rate to keep constant power density
- $P/tr \propto CV^2f$
- $E/tr \propto CV^2$ 
  - ...but V scaling more slowly than F

Penn ESE5320 Fall 2022 -- DeHon

7

7

40 Years of Microprocessor Trend Data



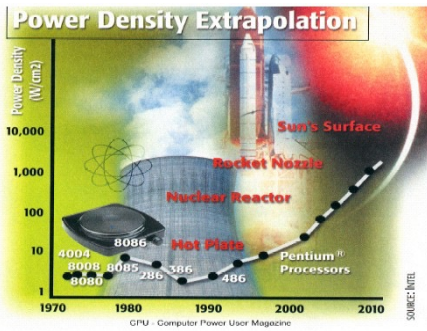
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Okukotun, L. Hammond, and C. Batten. New plot and data collected for 2010-2015 by K. Rupp.

Penn ESE5320 Fall 2022 -- DeHon

8

8

## Intel Power Density



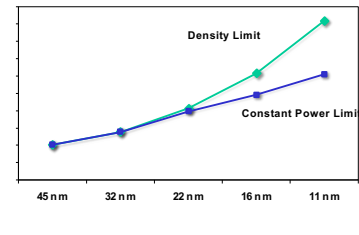
Penn ESE5320 Fall 2022 -- DeHon

9

9

## Impact

### Power Limits Integration



Source: Carter/Intel

Penn ESE5320 Fall 2022 -- DeHon

10

10

## Impact

- Power density is limiting scaling
  - Can already place more transistors on a chip than we can afford to turn on!
- Power is potential challenge/limiter for all future chips.
  - Only turn on small percentage of transistors?
  - Operate those transistors at much slower frequency?

Penn ESE5320 Fall 2022 -- DeHon

11

11

## Energy

$$E_{total} = E_{switch} + E_{leak}$$

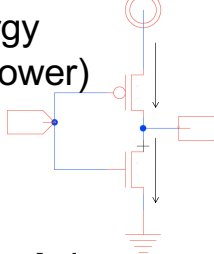
- Dynamic – Switching Energy (Power)
- Static – Leakage Energy (Power)

Penn ESE5320 Fall 2022 -- DeHon

12

12

### Leakage Energy (Static Energy, Power)



- $I_{leak}$ 
  - Subthreshold leakage
  - (possibly) Gate-Drain leakage

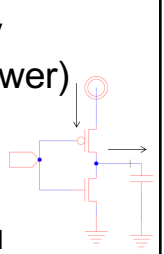
$$P_{leak} = I_{leak} \times V$$

$$E_{leak} = P_{leak} \times T$$

Penn ESE5320 Fall 2022 -- DeHon 13

13

### Switching Energy (Dynamic Energy, Power)

$$E_{switch} \propto \alpha CV^2$$


- Charge/discharging gate and load
- C – driven by architecture
- V – today, driven by variation, aging
- $\alpha$  – probability will switch; driven by architecture, coding/information

$Q = CV = \int I(t) dt$   
 $E = \int I(t)V_{dd} dt$

Penn ESE5320 Fall 2022 -- DeHon 14

14

### Energy

$$E_{total} = E_{switch} + E_{leak}$$

$$E_{switch} \propto \alpha CV^2$$

$$E_{leak} = I_{leak} \times V \times T$$

Penn ESE5320 Fall 2022 -- DeHon 15

15

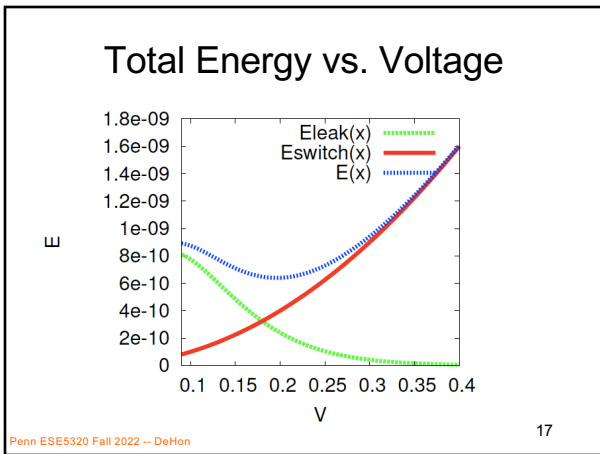
### Voltage

$$E_{switch} \propto \alpha CV^2 \quad E_{leak} = I_{leak} \times V \times T$$

- We can set voltage
- Reducing voltage
  - Reduces  $E_{switch}$
  - Increases delay  $\rightarrow$  cycle time, T (freq=1/T)
    - Reduce performance (not run full freq.)
    - Increases leakage  $E_{leak}$

Penn ESE5320 Fall 2022 -- DeHon 16

16



17

### Switching Energy

$$E_{switch} \propto \alpha CV^2$$

- C – driven by architecture
  - Also impacted by variation, aging
- V – today, driven by variation, aging
- $\alpha$  – probably a node will switch; driven by architecture, information

Penn ESE5320 Fall 2022 -- DeHon 18

18

## Data Dependent Activity

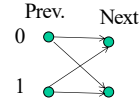
- Consider an 8b counter
  - How often do each of the following switch?
    - Low bit?
    - High bit?
  - Average switching across all 8 output bits?
- Assuming random inputs
  - Activity (switching probability) at output of xor4?
    - Hint: probability of output of xor4 being 0? 1?
  - Activity at output of nand4?

Penn ESE5320 Fall 2022 -- DeHon

19

19

## Gate Output Switching (random inputs)



$$P_{switch} = P(0@i) * P(1@i+1) + P(1@i) * P(0@i+1)$$

Penn ESE5320 Fall 2022 -- DeHon

20

20

## Conclude

- Not all gates (output nodes) switch at same rate

Penn ESE5320 Fall 2022 -- DeHon

21

21

## Switching Energy

$$E_{switch} = \left( \sum_i \alpha_i C_i \right) V^2$$

$C_i$  == capacitance driven by each gate (including wire)

Penn ESE5320 Fall 2022 -- DeHon

22

22

## Switching Rate ( $\alpha_i$ ) Varies

- Different logic (low/high bits, gate type)
- Different usage
  - Gate off unused functional units
- Data coded
- Entropy in data
- Average  $\alpha$  5--15% plausible

$$E_{switch} = \left( \sum_i \alpha_i C_i \right) V^2$$

Penn ESE5320 Fall 2022 -- DeHon

23

23

## Switching Energy

$$E_{switch} \propto \alpha C V^2$$

- $C$  – driven by architecture
- $V$  – today, driven by variation, aging
- $\alpha$  – driven by architecture, information

Penn ESE5320 Fall 2022 -- DeHon

24

24

### Wire Driven

$$E_{switch} = \left( \sum_i \alpha_i C_i \right) V^2$$

- Gates drive
  - Self
  - Inputs to other gates
  - Wire routing between self and other gates
- Typically:  $C_{wire} > C_{self} + C_{load}$

$$C_i = C_{self} + C_{wire} + C_{load}$$

Penn ESE5320

25

### Wire Capacitance

- How does wire capacitance relate to wire length?

Penn ESE5320 Fall 2022 -- DeHon

26

### Wire Capacitance

- $C = \epsilon A/d = \epsilon W * L_{wire}/d = C_{unit} * L_{wire}$
- Wire capacitance is linear in wire length
- E.g. 1.7pF/cm (preclass)
- Remains true if buffer wire
  - Add buffered segment at fixed lengths
  - [different use of word *buffer* than normal in this class – here we're talking about adding a repeater to perform electrical signal restoration]

Penn ESE5320 Fall 2022 -- DeHon

27

### Wire Driven Implications

- Care about locality
  - Long wires are higher energy
  - Producers near consumers
  - Memories near compute
  - Esp. for large  $\alpha_i$ 's
- Care about size/area
  - Reduce (worst-case) distance must cross
- Care about minimizing data movement
  - Less data, less often, smaller distances
- Care about size of memories

$$E_{switch} = \left( \sum_i \alpha_i C_i \right) V^2$$

Penn ESE5320 Fall 2022 -- DeHon

28

### Preclass 5

- Primary switching capacitance in wires
- C: How does energy of a read grow with capacity (N) of a memory bank?
- D: Energy per bit?

Penn ESE5320 Fall 2022 -- DeHon

29

### Memory Implications

- Memory energy can be expensive
- Small memories cost less energy than large memories
  - Use data from small memories as much as possible
- Cheaper to re-use data item from register than re-reading from memory

Penn ESE5320 Fall 2022 -- DeHon

30

## Energy

$$E_{total} = E_{switch} + E_{leak}$$

$$E_{switch} = \left( \sum_i \alpha_i C_i \right) V^2$$

$$E_{leak} = \left( \sum_i E_{leak_i} \right) V \times T$$

Penn ESE5320 Fall 2022 -- DeHon

31

## Architectural Implications

Part 2

Penn ESE5320 Fall 2022 -- DeHon

32

## Component Numbers

**TABLE 1**

Operation	Energy
32-bit arithmetic operation	5 pJ
32-bit register read	10 pJ
32-bit 8KB RAM read	50 pJ
32-bit traverse 10mm wire	100 pJ
Execute instruction	500 pJ

Energy Per Operation (0.13µm, 1.2V)

Penn ESE5320 Fall 2022 -- DeHon

[Dally, March 2004 ACM Queue]

33

## Component Numbers

- Processor instruction 100x more than arithmetic
- Register read 2x
- RAM read 10x
- Why processor instruction > arith operation?

**TABLE 1**

Operation	Energy
32-bit arithmetic operation	5 pJ
32-bit register read	10 pJ
32-bit 8KB RAM read	50 pJ
32-bit traverse 10mm wire	100 pJ
Execute instruction	500 pJ

Energy Per Operation (0.13µm, 1.2V)

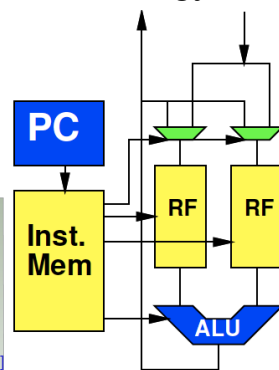
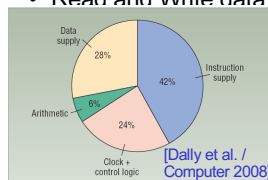
Penn ESE5320 Fall 2022 -- DeHon

[Dally, March 2004 ACM Queue]

34

## Processors and Energy

- Very little energy into actual computation
- Determine and Fetch Instruction
- Read and Write data



Penn ESE5320 Fall 2022 -- DeHon

[Dally et al. / Computer 2008]

35

## ARM Cortex A9

Estimate find: 0.5W at 800MHz in 40nm

- $0.5/0.8 \times 10^{-9}$  J/instr
- ~600pJ/instr
- Scale to 28nm
  - maybe  $0.7 \times 600$ — $0.5 \times 600$
  - 300—400pJ/instr ?
- Is superscalar w/ Neon, so not as simple a processor as previous example

Penn ESE5320 Fall 2022 -- DeHon

36

## Zynq (7-series, 28nm)

Operation	PL Resource	ARM A9 Resource	ARM A9 energy/OP (pico Joules or mW/GOP/sec)	PL energy/OP (pico Joules or mW/GOP/sec)
Logical Op of 2 var	LUT/FF	ALU		1.3
32-bit ADD	LUT/FF	ALU		1.3
16x16 Mult	DSP	ALU		8.0
32-bit Read/Write register	LUTRAM	L1		1.4
32-bit Read/Write AXI register	LUT/FF	AXI		30
32-bit Read/Write local RAM	BRAM	L2		23.7/17.2
32-bit Read/Write OCM	AXI/OCM	CPU/OCM		44
32-bit Read/Write DDR3	AXI/DDR	CPU/DDR		541/211

- ARM A9 instruction 300—400pJ
- ARM A9 L1 cache read 23pJ

Penn ESE5320 Fall 2022 -- DeHon

Xilinx UG585 – Zynq TRM 37

37

## Compare

- Assume ARM Cortex A9 executes 4x32b Neon vector add instruction for 300pJ
- Compare to 32b adds on FPGA?

Operation	PL Resource	ARM A9 Resource	ARM A9 energy/OP (pico Joules or mW/GOP/sec)	PL energy/OP (pico Joules or mW/GOP/sec)
Logical Op of 2 var	LUT/FF	ALU		1.3
32-bit ADD	LUT/FF	ALU		1.3
16x16 Mult	DSP	ALU		8.0
32-bit Read/Write register	LUTRAM	L1		1.4
32-bit Read/Write AXI register	LUT/FF	AXI		30
32-bit Read/Write local RAM	BRAM	L2		23.7/17.2
32-bit Read/Write OCM	AXI/OCM	CPU/OCM		44
32-bit Read/Write DDR3	AXI/DDR	CPU/DDR		541/211

38

## Compare

- Assume ARM Cortex A9 executes 4x16b Neon vector multiply instruction for 300pJ
- Compare to 16x16 multiplies on FPGA?

Operation	PL Resource	ARM A9 Resource	ARM A9 energy/OP (pico Joules or mW/GOP/sec)	PL energy/OP (pico Joules or mW/GOP/sec)
Logical Op of 2 var	LUT/FF	ALU		1.3
32-bit ADD	LUT/FF	ALU		1.3
16x16 Mult	DSP	ALU		8.0
32-bit Read/Write register	LUTRAM	L1		1.4
32-bit Read/Write AXI register	LUT/FF	AXI		30
32-bit Read/Write local RAM	BRAM	L2		23.7/17.2
32-bit Read/Write OCM	AXI/OCM	CPU/OCM		44
32-bit Read/Write DDR3	AXI/DDR	CPU/DDR		541/211

39

## Programmable Datapath

- Performing an operation in a pipelined datapath can be orders of magnitude less energy than on a processor
  - ARM 300pJ vs. 1.3pJ 32b add
  - Even neon 300pJ vs. 4x1.3pJ for 4x32b add
  - 300pJ vs. 4x8pJ for 4 16x16b multiplies

40

## Zynq

Operation	PL Resource	ARM A9 Resource	ARM A9 energy/OP (pico Joules or mW/GOP/sec)	PL energy/OP (pico Joules or mW/GOP/sec)
Logical Op of 2 var	LUT/FF	ALU		1.3
32-bit ADD	LUT/FF	ALU		1.3
16x16 Mult	DSP	ALU		8.0
32-bit Read/Write register	LUTRAM	L1		1.4
32-bit Read/Write AXI register	LUT/FF	AXI		30
32-bit Read/Write local RAM	BRAM	L2		23.7/17.2
32-bit Read/Write OCM	AXI/OCM	CPU/OCM		44
32-bit Read/Write DDR3	AXI/DDR	CPU/DDR		541/211

- Reading from OCM order of magnitude less than from DRAM
- ...and BRAM half that

Penn ESE5320 Fall 2022 -- DeHon

Xilinx UG585 – Zynq TRM 41

41

## FPGA vs. Std Cell Energy

Name	Method	Logic Only	Logic & DSP	Logic & Memory	Logic, Memory & DSP
booth	Sim	26			
rs.encoder	Sim	52			
cordic18	Const	6.3			
cordic8	Const	5.7			
des_area	Const	27			
des_perf	Const	9.3			
fir_restruct	Const	9.6			
mac1	Const	19			
aes192	Sim	12			
irc3	Const	12	7.5		
diffeq2	Const	15	12		
diffeq2	Const	16	12		
molecular	Const	15	16		
rs.decoder1	Const	13	16		
rs.decoder2	Const	11	11		
atm	Const			15	
aes	Sim			13	
aes_inv	Sim			12	
ethernet	Const			16	
serialproc	Const			16	
fir24	Const				5.3
pipe5proc	Const				8.2
raytracer	Const				8.3
Geomean		14	12	14	7.1

[Kuon/Rose TRCADv26n2p203--215 2007]

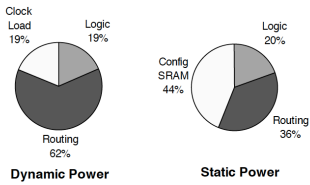
Penn ESE5320 Fall 2022 -- DeHon

42

42

## FPGA Disadvantage to Custom

- Interconnect Energy
  - Long wires → more capacitance → more E
  - Switch Energy is an overhead



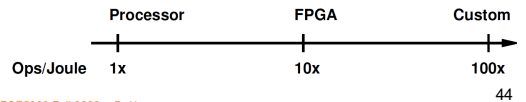
[Tuan et al./FPGA 2006] 43

Penn ESE5320 Fall 2022 -- DeHon

43

## Simplified Comparison

- Processor two orders of magnitude higher energy than custom accelerator
- FPGA accelerator in between
  - Order of magnitude lower than processor
  - Order of magnitude higher than custom



Penn ESE5320 Fall 2022 -- DeHon

44

## Accelerators save Energy

- Customized accelerators save energy
  - Even when built on FPGA
- Recall: can put more transistors on chip than can afford to turn on
- Opportunity to spend extra transistors to save energy and increase performance
  - Even if don't use accelerator continuously
- Opportunity: Area-energy tradeoff

Penn ESE5320 Fall 2022 -- DeHon

45

45

## Apple A14 Bionic

- 88mm<sup>2</sup>, 5nm
- 11.8 Billion Tr.
- iPhone 12
- 6 ARM cores
  - 2 fast (2.9–3GHz)
  - 4 low energy
- 4 custom GPUs
- 16 Neural Engines
  - 11 Trillion ops/s?



Image from <https://www.extremetech.com/computing/318715-comparison-of-apple-m1-a14-shows-differences-in-soc-design>  
 details: <https://www.tomshardware.com/news/apple-a14-bionic-revealed>  
<https://www.anandtech.com/show/16226/apple-silicon-m1-a14-deep-dive/2>

Penn ESE5320 Fall 2022 -- DeHon

46

46

## Parallelism and Energy

Part 3

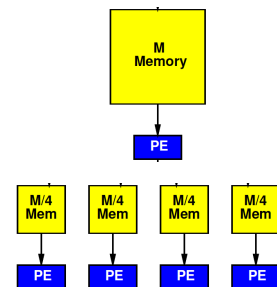
Penn ESE5320 Fall 2022 -- DeHon

47

47

## Preclass 6

- Energy
  - Per read from  $M=10^6$  memory?
  - Per read from  $10^6/4$  memory?



Penn ESE5320 Fall 2022 -- DeHon

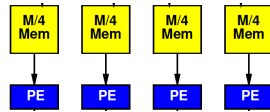
48

48



## Local Consumption

- To exploit, we must consume the data local to the memory.



Penn ESE5320 Fall 2022 -- DeHon

49

49

## Inter PE Communication

- May need to communicate between parallel processing elements [PEs] (and memories)
- Must pay for energy to move data between PEs

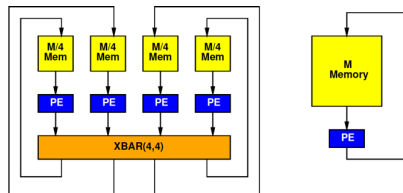
Penn ESE5320 Fall 2022 -- DeHon

54

54

## Preclass 7

- Energy: Read 4 memories  $10^6/4$ , route  $4 \times 4$  crossbar, write 4  $10^6/4$  memories?
- Energy: 4 reads from  $10^6$  memory, 4 writes from  $10^6$  memory?



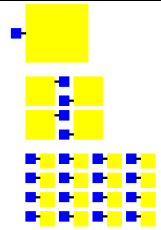
Penn ESE5320 Fall 2022 -- DeHon

55

55

## Parallel Larger

- More parallel design
  - Has more PEs
  - Adds interconnect
- Total area > less parallel design
  - More area → longer wires → more energy in communication between PEs
  - Could increase energy!



Penn ESE5320 Fall 2022 -- DeHon

56

56

## Continuum Question

- Where do we minimize total energy?
  - Both memory and communication
- Design axis P – number of PEs
  - What P minimizes energy?

Penn ESE5320 Fall 2022 -- DeHon

57

57

## Simple Model

- $E_{\text{mem}} = \text{sqrt}(M)$
- Communication =  $E_{\text{xbar}}(I,O) = 4 * I * O$
- P Processors
- N total data
- Possibly communicate each result to other PEs

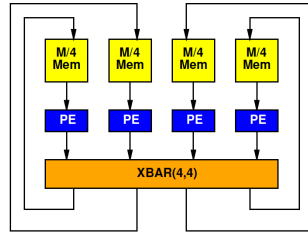
Penn ESE5320 Fall 2022 -- DeHon

58

58

## Simple Model: Memory

- Divide N data over P memories
- $E_{mem} = \sqrt{N/P}$
- N total memory operations
- Memory energy:  $N \times \sqrt{N/P}$
- Memory energy decrease with P



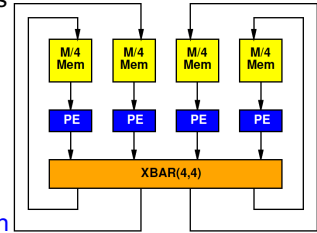
Penn ESE5320 Fall 2022 -- DeHon

59

59

## Simple Model: Communication

- Crossbar with P inputs and P outputs
- $E_{xbar} = 4 \times P \times P$
- Crossbar used N/P times
- Crossbar energy:  $4 \times N \times P$
- Communication energy increase with P



Penn ESE5320 Fall 2022 -- DeHon

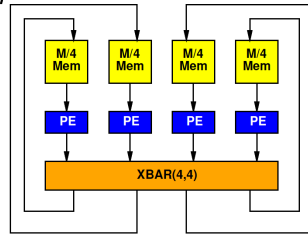
60

60

## Simple Model

$$N \times \sqrt{\frac{N}{P}} + N \times 4 \times P$$

$$N \times \left( \sqrt{\frac{N}{P}} + 4 \times P \right)$$



Penn ESE5320 Fall 2022 -- DeHon

61

61

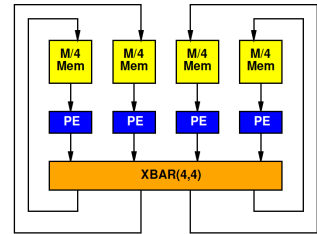
## Preclass 8

- For  $N=10^6$

$$N \times \left( \sqrt{\frac{N}{P}} + 4 \times P \right)$$

- Per operation becomes:

$$\frac{10^3}{\sqrt{P}} + 4P$$



Penn ESE5320 Fall 2022 -- DeHon

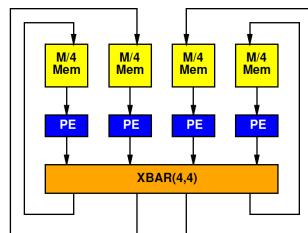
62

62

## Preclass 8

- Energy for:
  - P=1
  - P=4
  - P=100
- Energy minimizing P?
  - Energy?

$$\frac{10^3}{\sqrt{P}} + 4P$$

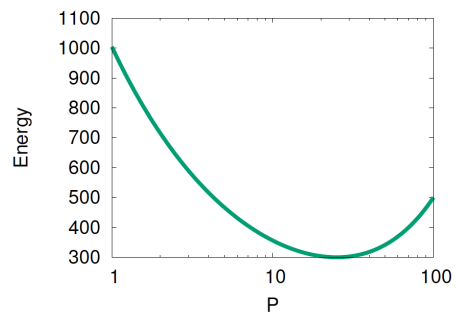


Penn ESE5320 Fall 2022 -- DeHon

63

63

## Graph



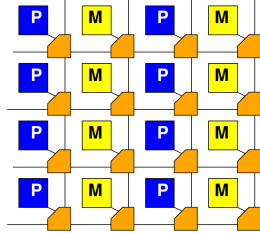
Penn ESE5320 Fall 2022 -- DeHon

64

64

## High Locality

- If communication is local, don't need crossbar
  - Remember Rent's Rule
- Communication energy scales less than  $P^2$
- Can scale as low as  $P$



65

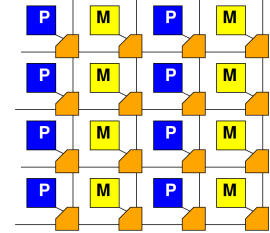
Penn ESE5320 Fall 2022 -- DeHon

65

## Model for High Locality

- $E_{\text{comm}} = \text{constant}$
- $E_{\text{comm}} = 10$
- Total comm:  $N * 10$

$$N \times \left( \sqrt{\frac{N}{P}} + 10 \right)$$



66

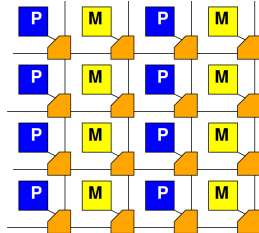
Penn ESE5320 Fall 2022 -- DeHon

66

## Preclass 9

- What is energy minimizing  $P$ ?

$$\frac{10^3}{\sqrt{P}} + 10$$



67

Penn ESE5320 Fall 2022 -- DeHon

67

## Task Locality Matters

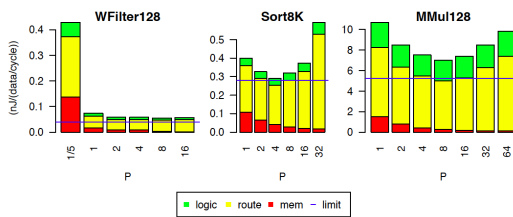
- Optimal  $P$  (processors) depends on communication locality (Rent Exponent  $p$ )
  - Very local problems ( $p < 0.5$ ) always benefit from parallelism
  - Highly interconnected problems ( $p > 0.5$ ) must balance energies  $\rightarrow$  intermediate parallelism

68

Penn ESE5320 Fall 2022 -- DeHon

68

## Tune Parallelism: Stratix-IV



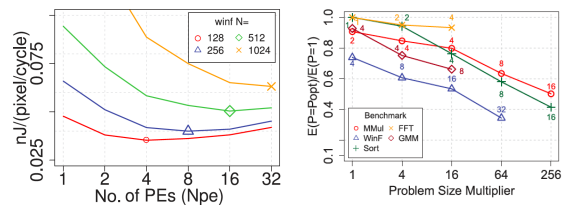
69

Penn ESE5320 Fall 2022 -- DeHon

[Kadic, FPGA 2015]

69

## PE Scaling with Problem Size



70

Penn ESE5320 Fall 2022 -- DeHon

[Kadic, TRET 2016]

70

## Power vs. Energy

- Notice: main comparisons are energy for operation or task
  - Not power
- Power is rate of energy use
  - Can slow down clock to reduce power
    - But leaves Energy unchanged
  - Power alone doesn't tell us energy requirement or energy efficiency
    - Need to also know time
      - Energy  $\approx$  Power \* Time
    - Sometimes see Gops/W  $\rightarrow$  Ops/J  $\rightarrow$  energy metric

Penn ESE5320 Fall 2022 -- DeHon

71

## Big Ideas

- Energy dominance
- With power-density budget
  - The most energy efficient architecture delivers the most performance
- Make memories small and wires short
- SoC, accelerators reduce energy by reducing processor instruction execution overhead
- Parallel design exploit locality  $\rightarrow$  reduce energy
- Optimal parallelism for problem
  - Driven by communication structure, size

Penn ESE5320 Fall 2022 -- DeHon

72

72

## Admin

- Feedback (including P3)
- Tuesday is Virtual Thursday
  - Thursday office hours
- No class on Wed.
  - This Wed. is a Virtual Friday
- Happy Thanksgiving!
- P4 due next Friday
- Reading for Monday on web

Penn ESE5320 Fall 2022 -- DeHon

73

73