

# ESE5320: System-on-a-Chip Architecture

Day 1: August 28, 2024  
Introduction and Overview  
(lecture start target 10:20am)

Note: slides, preclass linked to web  
[www.seas.upenn.edu/~ese5320/fall2024/fall2024.html](http://www.seas.upenn.edu/~ese5320/fall2024/fall2024.html)

- Work/finish preclass while waiting for lecture start
- Feedback form (turn in end of lecture)



Penn ESE5320 Fall 2024 -- DeHon

1

# Today

- Part 1: Case for Programmable SoC (**motivation**)
- Part 2: Course Goals, Outcomes, Tools (**philosophy?**)
- Part 3: Sample Optimization (**fast, flavor**)
- Part 4: This course (**operational details**)
  - (including policies, logistics)

Penn ESE5320 Fall 2024 -- DeHon

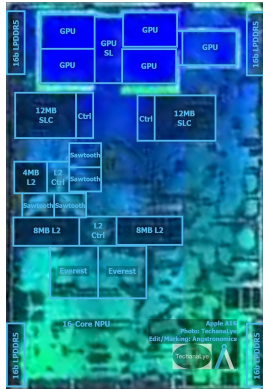
2

2

# Apple A16 Bionic

- ? 110+mm<sup>2</sup>, 4nm
- 16 Billion Tr.
- iPhone 14
- 6 ARM cores
  - 2 fast (3.5GHz)
  - 4 low energy (2GHz)
- 5 custom GPUs (1.4GHz)
- 16 Neural Engines
  - 17 Trillion ops/s?

Penn ESE5320 Fall 2024 -- DeHon

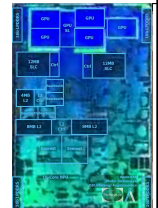


3

# Questions

- Why do today's SoC look like they do?
- How approach programming modern SoCs?
- How design a custom SoC?
- When building a System-on-a-Chip (SoC)
  - How much area should go into:
    - Processor cores, GPUs, FPGA logic, memory, interconnect, custom functions (which) .... ?

Penn ESE5320 Fall 2024 -- DeHon

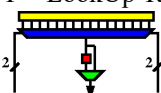


5

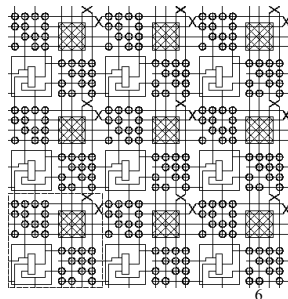
# FPGA

## Field-Programmable Gate Array

K-LUT (typical k=4 or 6)  
Compute block  
w/ optional  
output Flip-Flop  
(LUT = LookUp Table)



ESE1500, CIS5710  
Penn ESE5320 Fall 2024 -- DeHon



6

# Case for Programmable SoC

Penn ESE5320 Fall 2024 -- DeHon

7

7

## End of Microprocessor Scaling

Old

- Moore's Law scaling delivered faster transistors
- Processors rode Moore's Law
  - Turning transistors into performance
- Could wait and ride technology curve

Now

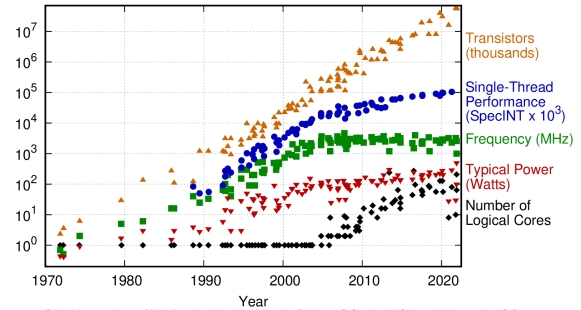
- Dennard's Law kicked in
  - How need to scale voltage with size
- Microprocessors were burning more power
- Lost ability to scale down voltage
- Processor performance stalled

Penn ESE5320 Fall 2024 -- DeHon

8

8

50 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Laborte, O. Shacham, K. Okukotun, L. Hammond, and C. Batten. New plot and data collected for 2010-2021 by K. Rupp.

Penn ESE5320 Fall 2024 -- DeHon

9

9

## The Way things Were

30 years ago

- Wanted programmability
  - used a processor
- Wanted it a little faster
  - Next year's processor would run faster...
- Wanted high-throughput
  - used a custom Integrated Circuit (IC) -- chip
- Wanted product differentiation
  - Got it at the board level
  - Select which ICs and how wired
- Build a custom IC (chip)
  - It was about gates and logic

Penn ESE5320 Fall 2024 -- DeHon

10

10

## Today

- Microprocessor may not be fast enough
  - (but often it is)
  - Or low enough energy
- Single core processor scaling has ended
- Time and Cost of a custom IC is too high
  - \$100M's of dollars for development, Years
- FPGAs promising
  - But build everything from prog. gates?
- Premium for small part count
  - And avoid chip crossing
  - ICs with 10–100 Billions of Transistors

Penn ESE5320 Fall 2024 -- DeHon

11

11

## Non-Recurring Engineering (NRE) Costs

- Costs spent up front on development
  - Engineering Design Time
  - Design Verification
  - Prototypes
  - Mask costs
- Recurring Engineering
  - Costs to produce each chip

$$Cost(N_{chips}) = Cost_{NRE} + N_{chips} \times Cost_{perchip}$$

Penn ESE5320 Fall 2024 -- DeHon

12

12

## NRE Costs

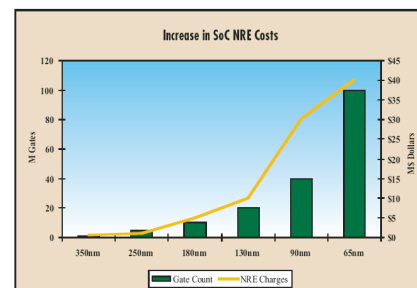


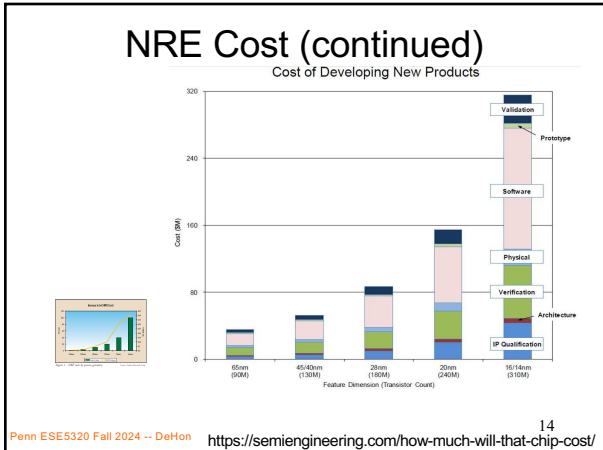
Figure 1 - NRE costs by process geometry

Source: Semico Research Corp.

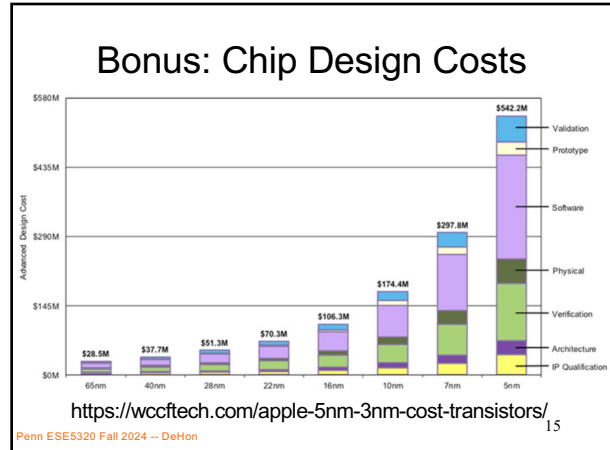
Penn ESE5320 Fall 2024 -- DeHon

13

13



14



15

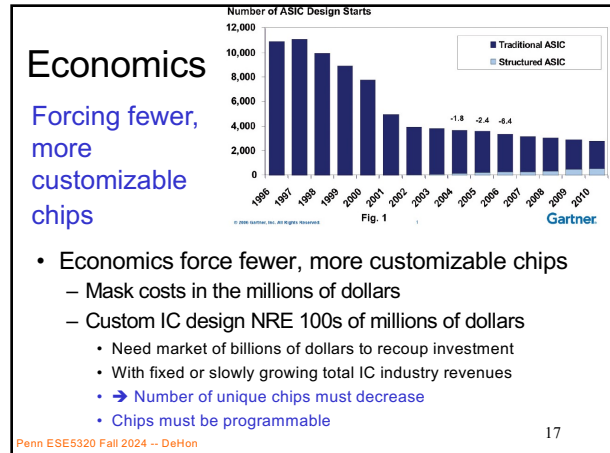
### Amortize NRE with Volume

$$Cost(N_{chips}) = Cost_{NRE} + N_{chips} \times Cost_{perchip}$$

$$Cost = \frac{Cost_{NRE}}{N_{chips}} + Cost_{perchip}$$

Penn ESE5320 Fall 2024 -- DeHon 16

16



17

- ### Large ICs (Chips)
- Now contain significant software
    - Almost all have embedded processors
  - Must co-design SW and HW
  - Must solve complete computing task
    - Tasks has components with variety of needs
    - Some don't need custom circuit
    - 90/10 Rule
- Penn ESE5320 Fall 2024 -- DeHon 18

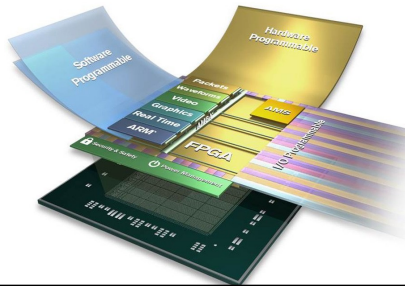
18

- ### Given Demand for Programmable
- How do we get higher performance than a processor, while retaining programmability?
    - Programmability – don't have to spend 100s of millions of dollars and months for fabrication?
- Penn ESE5320 Fall 2024 -- DeHon 19

19

## Programmable SoC

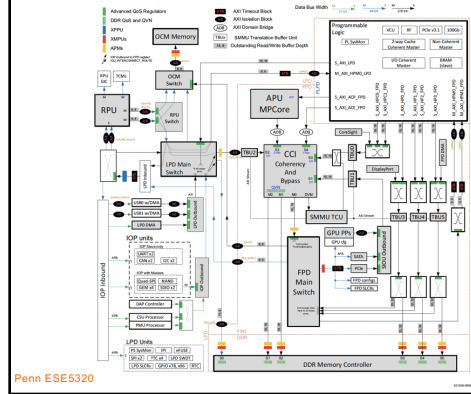
- Implementation Platform for innovation
  - This is what you target (avoid NRE)
  - Implementation vehicle



Penn ESE5320 Fall 2024 -- DeHon

20

## Programmable SoC



Penn ESE5320

UG1085  
Xilinx  
UltraScale  
Zynq  
TRM  
(p27)

21

21

## Then and Now

30 years ago

- Programmability?
  - use a processor
- Faster
  - Processors scaled
- High-throughput
  - used a custom IC
- Wanted product differentiation
  - board level
  - Select & wired IC
- Build a custom IC (Chip)
  - It was about gates and logic

Today

- Programmability?
  - uP, FPGA, GPU, PSoC
- Faster
  - Can't get with single core
- High-throughput
  - FPGA, GPU, PSoC, custom IC
- Wanted product differentiation
  - Program FPGAs, PSoC
- Build a custom IC (Chip)
  - System and software

Penn ESE5320 Fall 2024 -- DeHon

22

22

## Part 2: Course Goals, Outcomes

Penn ESE5320 Fall 2024 -- DeHon

23

23

## Goals

- Create Computer Engineers
  - SW/HW divide is wrong, outdated
  - Computer engineers understand computation
    - HW and SW are just tools and design options
  - Parallelism, data movement, resource management, abstractions
  - Cannot build a **chip** without software
- SoC user – know how to exploit
- SoC designer – architecture space, hw/sw codesign
- Project experience – design and optimization

Penn ESE5320 Fall 2024 -- DeHon

24

## Roles

- PhD Qualifier
  - One broad Computer Engineering
- CMPE Concurrency Lab
- Hands-on Project course

Penn ESE5320 Fall 2024 -- DeHon

25

25

## Outcomes

- Design, optimize, and program a modern System-on-a-Chip.
- Analyze, identify bottlenecks, design-space
  - Modeling → write equations to estimate
- Decompose into parallel components
- Characterize and develop real-time solutions
- Implement both hardware and software solutions
- Formulate hardware/software tradeoffs, and perform hardware/software codesign

Penn ESE5320 Fall 2024 -- DeHon

26

26

## Outcomes

- Understand the system on a chip from gates to application software, including:
  - on-chip memories and communication networks, I/O interfacing, design of accelerators, processors, firmware and OS/infrastructure software.
- Understand and *estimate* key design metrics and requirements including:
  - area, latency, throughput, energy, power, predictability, and reliability.

Penn ESE5320 Fall 2024 -- DeHon

27

27

## Course Programming

- Write *everything* in C
  - including for hardware (FPGA, spatial) operators
- Avoid learning separate language
  - Don't require or teach Verilog/VHDL
- Do focus on how tailor C for hardware
  - Focus on what's unique about specifying and guiding hardware
- Code → CHIPS

Penn ESE5320 Fall 2024 -- DeHon

28

28

## Tools

- Are complex
- Will be challenging, but good for you to build confidence can understand and master
- Tool runtimes can be long
- Learning and sharing experience will be part of assignments

Penn ESE5320 Fall 2024 -- DeHon

29

29

## Distinction

### CIS2400, 4710, 5710

- Best Effort Computing
  - Run as fast as you can
- Binary compatible
- ISA separation
- Shared memory parallelism

### ESE5320

- Real-Time
  - Guarantee meet deadline
- Hardware-Software codesign
  - Willing to recompile, maybe rewrite code
  - Define/refine hardware
- Non shared-memory parallelism models

Penn ESE5320 Fall 2024 -- DeHon

30

30

## Distinction

### ESE5390:

Hardware/Software Co-Design for Machine Learning

- Deep on Application (ML)
- More accessible to CS
  - Less previous experience with circuits and architecture
- Won't be as deep on understanding HW and optimization
- Program in Pytorch, OpenCL

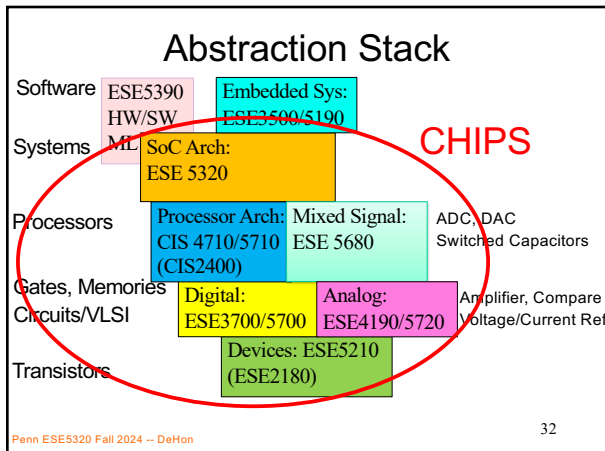
### ESE5320:

- Deep computer engineering
- Broad application
- Program in C
- Suitable followup if want to dig deeper

Penn ESE5320 Fall 2024 -- DeHon

31

31



32

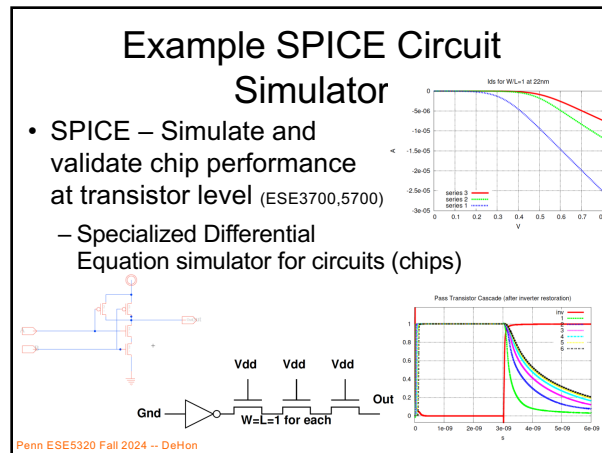
## Part 3: Approach -- Example

Penn ESE5320 Fall 2024 -- DeHon 33

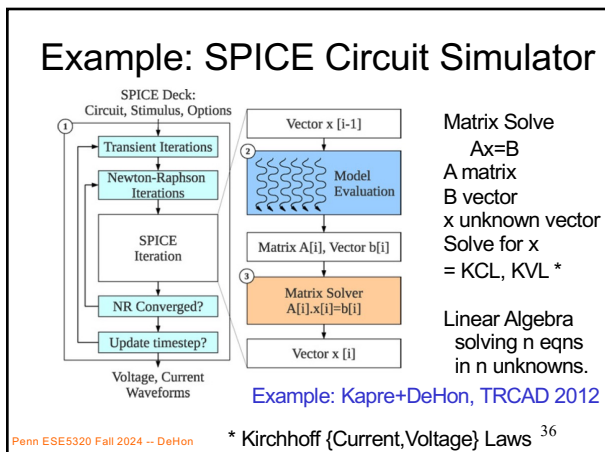
33

- ### Abstract Approach
- Identify requirements, bottlenecks
  - Decompose Parallel Opportunities
    - At extreme, how parallel could make it?
    - What forms of parallelism exist?
      - Thread-level, data parallel, instruction-level
  - Design space of mapping
    - Choices of where to map, area-time tradeoffs
  - Map, analyze, refine
    - Write equations to understand, predict
- Penn ESE5320 Fall 2024 -- DeHon 34

34



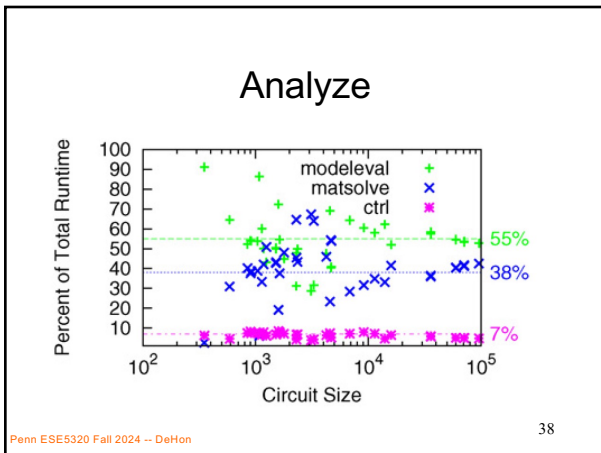
35



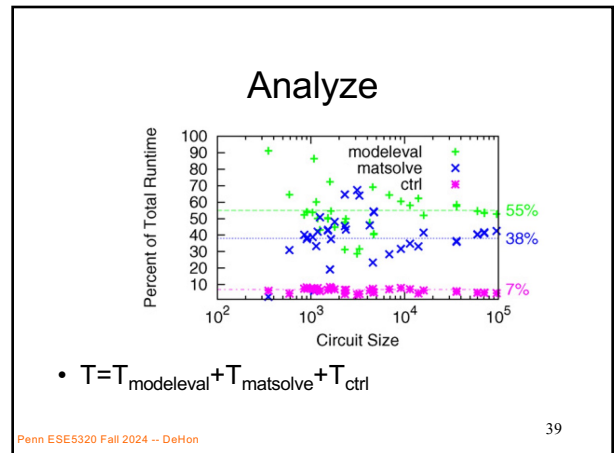
36

- ### Abstract Approach
- Identify requirements, bottlenecks
  - Decompose Parallel Opportunities
    - At extreme, how parallel could make it?
    - What forms of parallelism exist?
      - Thread-level, data parallel, instruction-level
  - Design space of mapping
    - Choices of where to map, area-time tradeoffs
  - Map, analyze, refine
    - Write equations to understand, predict
- Penn ESE5320 Fall 2024 -- DeHon 37

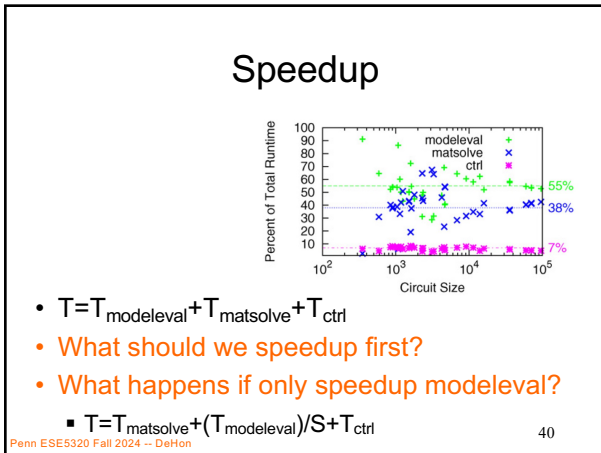
37



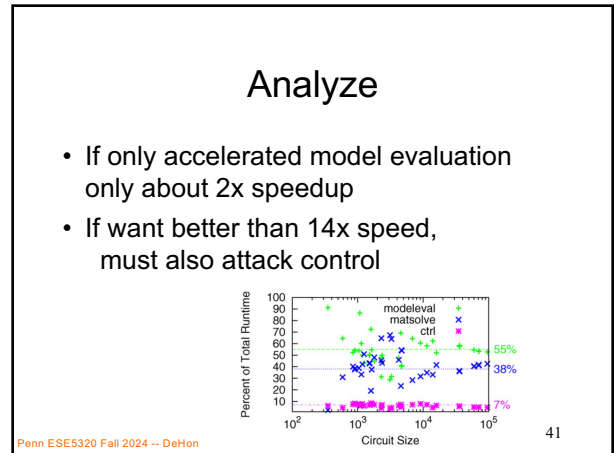
38



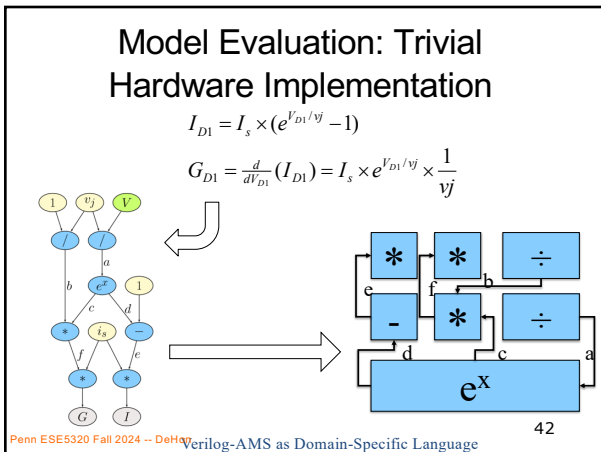
39



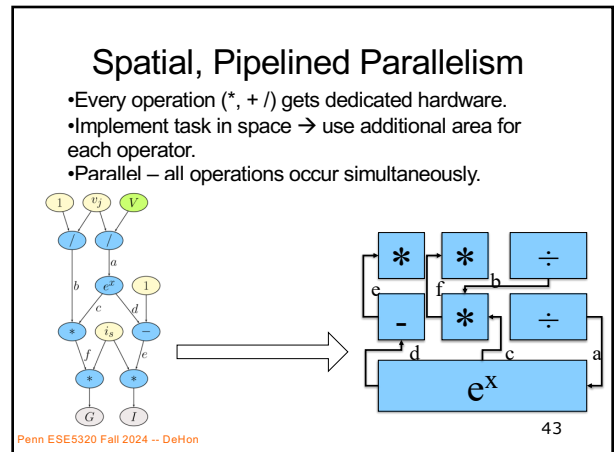
40



41



42

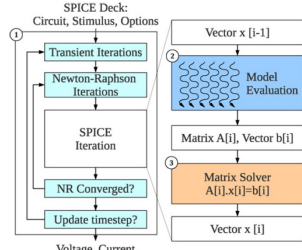


43



## Parallelism: Model Evaluation Data Parallel

- Every device independent
- Many of each type of device
- Can evaluate in parallel
  - $T = T_{seq} / N_{proc}$
- Build pipelined circuit for model
  - $T_{seq} = N_{comp} * T_{cycle}$
  - vs.  $T_{pipe} = T_{cycle}$

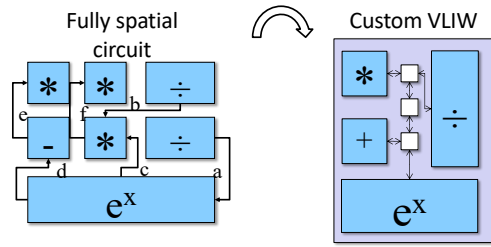


Penn ESE5320 Fall 2024 -- DeHon

44

44

## Spatial Too Big?



~100x Speedup  
Multiple FPGAs

~10x Speedup  
1 FPGA (2010)

VLIW=Very Long Instruction Word  
exploits Instruction-Level Parallelism

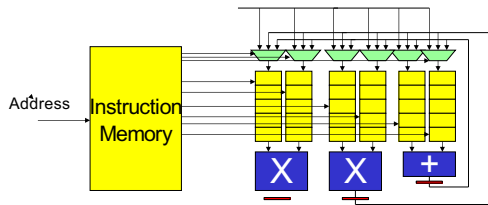
Penn ESE5320 Fall 2024 -- DeHon

45

45

## VLIW

- Very Long Instruction Word
- Supports Instruction/Operator-Level Parallelism
- Perform many primitive operations in parallel
  - Can parameterize and customize set of operations
- Using "long" instructions



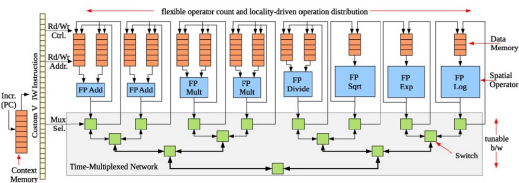
Penn ESE5320 Fall 2024 -- DeHon

46

46

## Parallelism: Model Evaluation

- Spatial end up bottlenecked by other components
- Use custom evaluation engines
- ...or GPUs



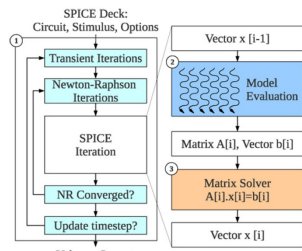
Penn ESE5320 Fall 2024 -- DeHon

47

47

## Parallelism: Matrix Solve

- Needed direct solver?
- E.g. Gaussian elimination
- Data dependence on previous reduce
  - Limited data parallelism
- Parallelism in subtractions
- Some row independence

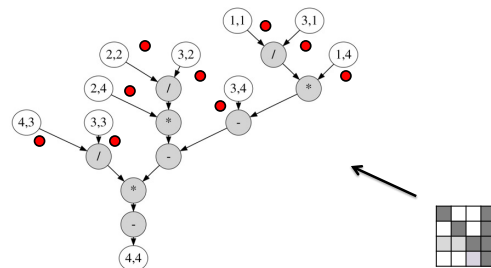


Penn ESE5320 Fall 2024 -- DeHon

48

48

## Example Matrix

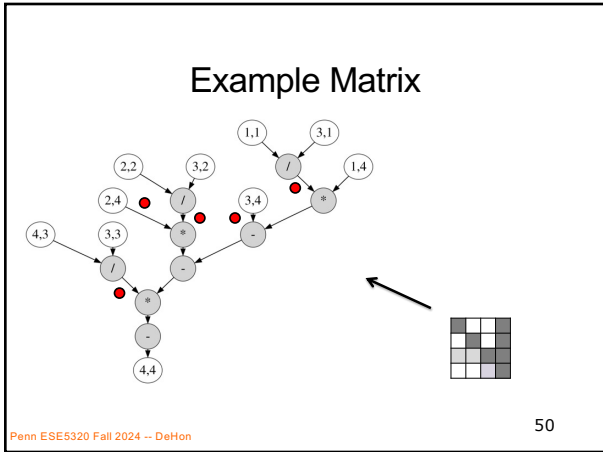


Penn ESE5320 Fall 2024 -- DeHon

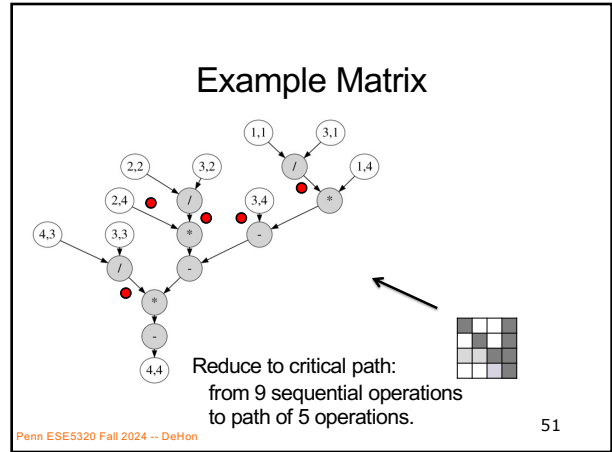
49

49

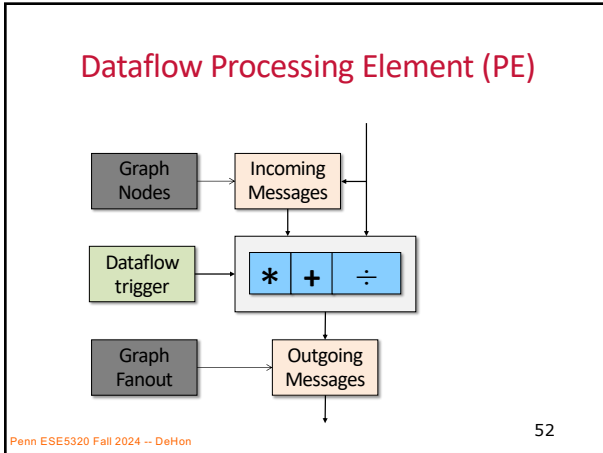




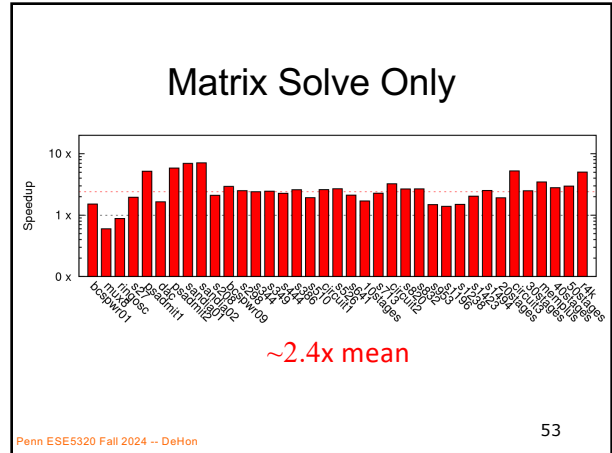
50



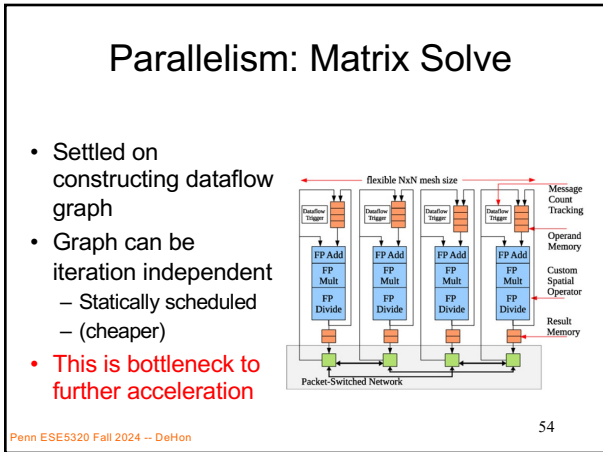
51



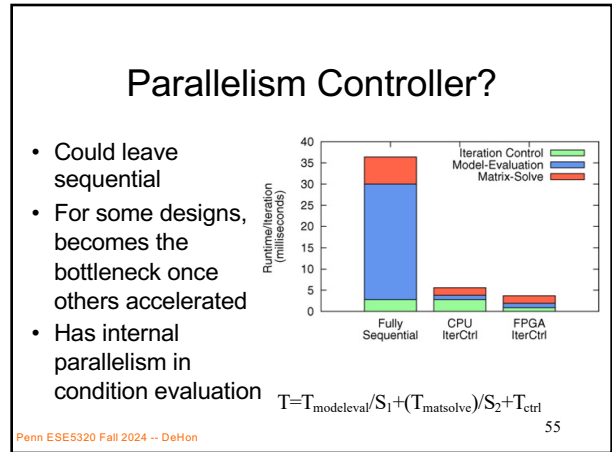
52



53



54



55

## Parallelism Controller

- Customized datapath controller

$$T_{seqctrl} = N_{add} + N_{mul} + 10 * N_{divide}$$

$$T_{vliwctrl} = \text{Max}(N_{add}/2, N_{mul}, 10 * N_{divide})$$

Penn ESE5320 Fall 2024 -- DeHon

56

56

## Single-Chip Solution

Penn ESE5320 Fall 2024 -- DeHon

57

57

## Area-Time for Each

Penn ESE5320 Fall 2024 -- DeHon

58

58

## Composite Speedup

Penn ESE5320 Fall 2024 -- DeHon

59

59

## Modern SoC

Penn ESE5320 Fall 2024 -- DeHon

60

60

## Part 4: Class Components

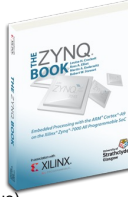
Penn ESE5320 Fall 2024 -- DeHon

61

61

## Class Components

- Lecture (incl. preclass exercise)
  - In-person (not hybrid, don't expect recordings)
  - Slides, preclass on web before class (print if you want)
  - N.B. I encourage class participation
    - In class; Questions ("warm" calls)
  - Daily Quiz
- Reading [~1 required paper/lecture]
  - online: Canvas, IEEE, ACM, also ZynqBook, Parallel Programming for FPGAs
- Homework (1 per week due 5pm)
- Project – open-ended (~6 weeks)
  - Oct. 23 – Dec. 9<sup>th</sup> (~ weekly milestones, details syllabus)



Penn ESE5320 Fall 2024 -- DeHon

62

• [Note syllabus, course admin online](#)

62

## First Half

Quickly cover breadth

- Metrics, bottlenecks
  - Memory
  - Parallel models
  - SIMD/Data Parallel
  - Thread-level parallelism
  - Spatial, C-to-gates
- Line up with homeworks

Penn ESE5320 Fall 2024 -- DeHon

63

63

## Second Half

- Use everything on project
- Going deeper
- Memory
- Verification
- VLIW
- Reduce
- Energy
- Chip Cost
- Real-time
- Reactive

Penn ESE5320 Fall 2024 -- DeHon

64

64

## Teaming

- HomeWorks (HW) in Groups of 2 (after 0, 1)
- HW: we assign
- Individual assignment writeup
- Project in Groups of 3
- Project: you propose team of 3, we review
  - Most portions group writeup
  - Maybe few components individual writeup

Penn ESE5320 Fall 2024 -- DeHon

65

65

## Office & Lab Hours

- Andre: TBD
  - Levine 270, Zoom
- TAs – Ketterer (starting next week)
  - Tuesday 8:30-9:30 pm
  - Wednesday 4:00-5:00 pm (not today)
  - Thursday 8:30pm-10:30pm (not tomorrow)

Penn ESE5320 Fall 2024 -- DeHon

66

66

## Diagnostic Assessment

- Course will rely heavily on C
  - Program both hardware and software in C
- If you cannot read/write code in C, this course will be a challenge
- Diagnostic Assessment intended as a quick indication if you aren't ready
  - Should be able to complete quickly
  - Better to find out now than after you're stuck in the course
  - Due next Wednesday (9/4)

Penn ESE5320 Fall 2024 -- DeHon

67

67

## C Review

- Course will rely heavily on C
  - Program both hardware and software in C
- HW1 has some C warmup problems
- TA will hold C review
  - on Sept. 4<sup>th</sup>, TBD (probably office hours)
  - (before our next class meeting since Monday 9/2 is Labor day)
  - See Ed Discuss for details

Penn ESE5320 Fall 2024 -- DeHon

68

68

## Preclass Exercise

- Motivate the topic of the day
  - Introduce a problem
  - Introduce a design space, tradeoff, transform
- Available on syllabus before lecture
  - May want to print a copy to bring to class
- Should work before lecture starts
  - Won't be available later
- Do bring/use calculator
  - Will be numerical examples

Penn ESE5320 Fall 2024 -- DeHon

69

69

## Daily Quiz

- Count for Engagement Points
- Only available until next lecture
- Incentive to keep up with material

Penn ESE5320 Fall 2024 -- DeHon

70

70

## Lecture Timeline

- Work on preclass before lecture start
- Start lecture at 10:20am
- Lecture until 11:40am
- (most days) stay for remaining questions
  - Pending course after us
  - Typically take questions in hall after clearing out for next course
  - [not today → CMPE Meet-and-Greet]

Penn ESE5320 Fall 2024 -- DeHon

71

71

## Feedback

- Will have anonymous feedback for each lecture
  - Clarity?
  - Speed?
  - Vocabulary?
  - General comments
    - Specificity most helpful
      - X was unclear because of Y
      - Subtopic Z went too fast
      - Need an example for Q

Penn ESE5320 Fall 2024 -- DeHon

72

72

## Policies

- Canvas turn-in of assignments
- No handwritten work
- Due on time
  - Individual assignments only
    - 3 free late days total
- Collaboration
  - Tools – allowed
  - Designs – limited to project teams as specified on assignments
- See web page

Penn ESE5320 Fall 2024 -- DeHon

73

73

- Your action: **Admin**
  - Feedback sheet for today
  - Find course web page
    - Read it, including the policies
    - Find Syllabus
      - Find diagnostic assessment, homework 1
      - Find lecture slides
        - » Will try to post before lecture
      - Find reading assignments
  - Find reading for lecture 2 on canvas and web
    - ...for this lecture if you haven't already
  - Find/join Ed Discussion group for course
  - Signup for detkin/ketterer access
  - ~~Complete/submit diagnostic assessment~~

Penn ESE5320 Fall 2024 -- DeHon

74

74

## Big Ideas

- Programmable Platforms
  - Key delivery vehicle for innovative computing applications
  - Reduce TTM (Time-to-Market), risk
  - More than a microprocessor
  - Heterogeneous, parallel
- Demand hardware-software codesign
  - Soft view of hardware
  - Resource-aware view of parallelism

Penn ESE5320 Fall 2024 -- DeHon

75

75

Questions?

Penn ESE5320 Fall 2024 -- DeHon

76

76