

ESE532: System-on-a-Chip Architecture

Day 1: September 1, 2010
Introduction and Overview
(lecture start target 10:20am)

Note: both preclass and feedback linked to web (also slides)

www.seas.upenn.edu/~ese532/fall2021/fall2021.html

- Preclass
- Feedback form



Penn ESE532 Fall 2021 -- DeHon

Today

- Part 1: Case for Programmable SoC
- Part 2
 - Course Goals
 - Outcomes
 - Evolving Course, Risks, Tools
- Part 3: Sample Optimization
- Part 4: This course
 - (incl. policies, logistics)

Penn ESE532 Fall 2021 -- DeHon

2

Apple A14 Bionic

- 88mm², 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?

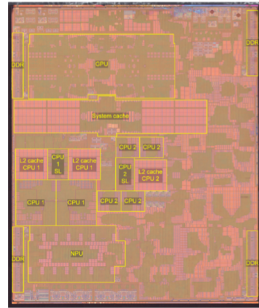


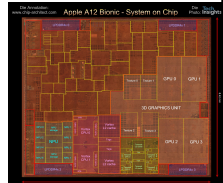
Image from <https://www.extremetech.com/computing/318715-comparison-of-apple-m1-a14-shows-differences-in-soc-design>
details: <https://www.tomshardware.com/news/apple-a14-bionic-revealed>
<https://www.anandtech.com/show/16226/apple-silicon-m1-a14-deep-dive-2>

Penn ESE532 Fall 2021 -- 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) ?

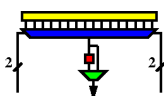


Penn ESE532 Fall 2021 -- DeHon

FPGA

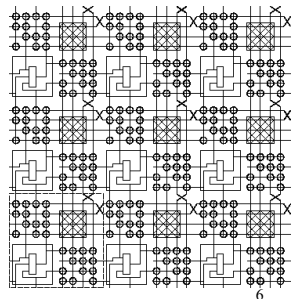
Field-Programmable Gate Array

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



ESE150, CIS571

Penn ESE532 Fall 2021 -- DeHon



Case for Programmable SoC

Penn ESE532 Fall 2021 -- DeHon

7

End of uProcessor Scaling

Old

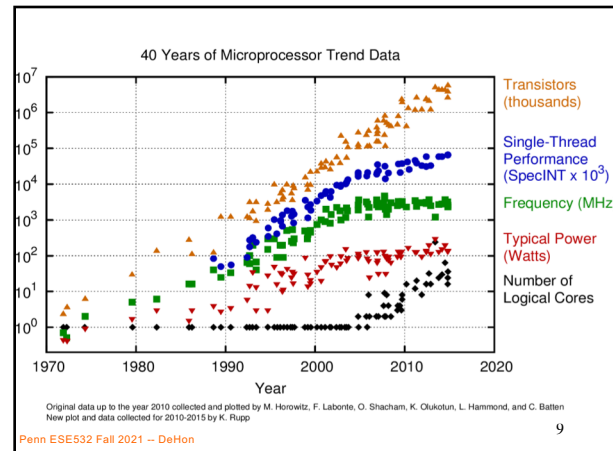
- Moore's Law scaling delivered faster transistors
- Processors rode Moore's Law
 - Turning transistors into performance
- Could wait and ride technology curve

Now

- Dennard's Law kicked in
- uP were burning more power
- Lost ability to scale down voltage
- Processor performance stalled

Penn ESE532 Fall 2021 -- DeHon

8



Penn ESE532 Fall 2021 -- DeHon

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 IC
- Wanted product differentiation
 - Got it at the board level
 - Select which ICs and how wired
- Build a custom IC
 - It was about gates and logic

Penn ESE532 Fall 2021 -- DeHon

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 Billions of Transistors

Penn ESE532 Fall 2021 -- DeHon

11

Non-Recurring Engineering (NRE) Costs

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

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

Penn ESE532 Fall 2021 -- DeHon

12

NRE Costs

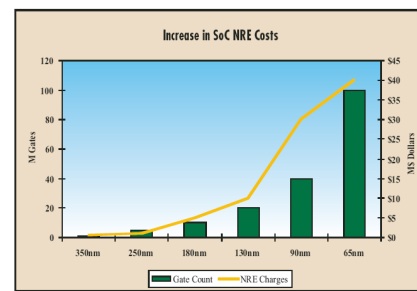


Figure 1 - NRE costs by process geometry

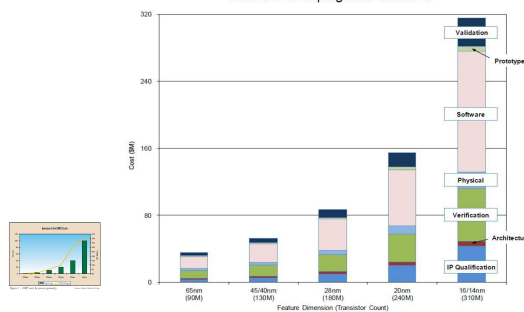
Source: Semico Research Corp.

Penn ESE532 Fall 2021 -- DeHon

13

NRE Cost (continued)

Cost of Developing New Products



Penn ESE532 Fall 2021 -- DeHon

<https://semiengineering.com/how-much-will-that-chip-cost/>

14

Amortize NRE with Volume

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

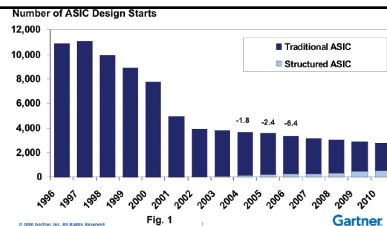
$$Cost = \frac{Cost_{NRE}}{N_{chips}} + Cost_{perchip}$$

Penn ESE532 Fall 2021 -- DeHon

15

Economics

Forcing fewer, more customizable chips



- Economics force fewer, more customizable chips
 - Mask costs in the millions of dollars
 - Custom IC design NRE 10s—100s of millions of dollars
 - Need market of billions of dollars to recoup investment
 - With fixed or slowly growing total IC industry revenues
 - Number of unique chips must decrease
 - Chips must be programmable

Penn ESE532 Fall 2021 -- DeHon

16

Large ICs

- 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 ESE532 Fall 2021 -- DeHon

17

Given Demand for Programmable

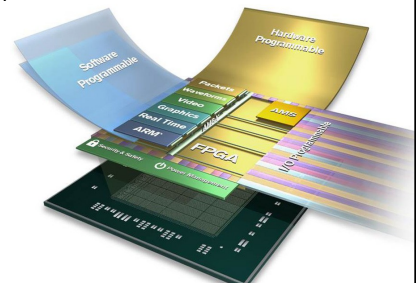
- How do we get higher performance than a processor, while retaining programmability?

Penn ESE532 Fall 2021 -- DeHon

18

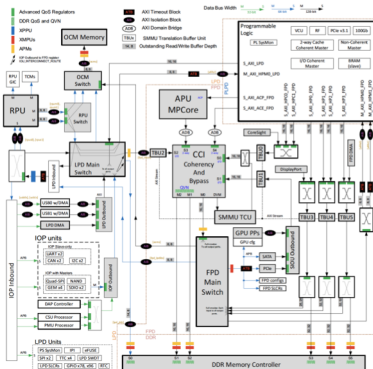
Programmable SoC

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



Penn ESE532 Fall 2021 -- DeHon

Programmable SoC



UG1085
Xilinx
UltraScale
Zynq
TRM
(p27)

20

Penn ESE532 I

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
 - 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
- Wanted product differentiation
 - Program FPGAs, PSoC
- Build a custom IC
 - System and software

Penn ESE532 Fall 2021 -- DeHon

21

Part 2: Course Goals, Outcomes

22

Penn ESE532 Fall 2021 -- DeHon

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 ESE532 Fall 2021 -- DeHon

Roles

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

24

Penn ESE532 Fall 2021 -- DeHon

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

25

Penn ESE532 Fall 2021 -- DeHon

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 ESE532 Fall 2021 -- DeHon

26

Evolving Course

- Spring 2017 – first offering
 - Raw, all assignments new, buggy, too tedious, long
- Fall 2017 – second offering
 - Refine assignments, project; increased explicit modeling emphasis
 - Hard, not insane
- Fall 2018 – third offering (similar 2017)
 - Added real-time ethernet data handling; project groups of 3
 - Many students challenged with C and software engineering
 - Stream debug and performance challenging
- Fall 2019 – fourth (structure same)
 - Try front-load more C, better introduce Stream optimization and debug
 - Group writeup on projects

Penn ESE532 Fall 2021 -- DeHon

27

Evolving Course

- Fall 2018 – third offering (similar 2017)
 - Added real-time ethernet data handling; project groups of 3
 - Many students challenged with C and software engineering
 - Stream debug and performance challenging
- Fall 2019 – fourth (structure same)
 - Try front-load more C, better introduce Stream optimization and debug
 - Group writeup on projects
- Fall 2020 – fifth (structure same)
 - Move to Vitis (from SDSoc)
 - Use Amazon cloud for first half; F1 instance for FPGA access HW
 - Then transition to Ultra96 (SoC FPGA) for projects
- Fall 2021 – sixth
 - Stay with Vitis; use DFX (see next slide)
 - Ultra96; no Amazon cloud
 - Introduce project components earlier

Penn ESE532 Fall 2021 -- DeHon

28

Tools

- Are complex
- Will be challenging, but good for you to build confidence can understand and master
- Tool runtimes can be long
 - Maybe DFX will help
 - DFX – Dynamic Function Exchange
 - Partial reconfiguration
- Learning and sharing experience will be part of assignments

Penn ESE532 Fall 2021 -- DeHon

29

Distinction

CIS240, 371, 471, 571

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

ESE532

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

Penn ESE532 Fall 2021 -- DeHon

30

Distinction

ESE539:

Hardware/Software Co-Design for Machine Learning

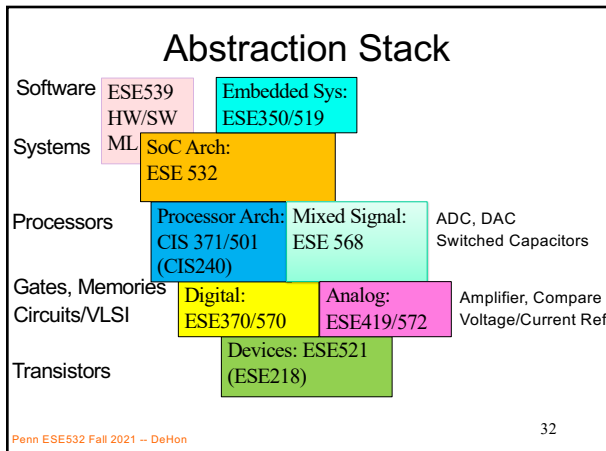
- Deep on Application (ML)
- More accessible to CS
 - Less previous experience with circuits and architecture
- Won't be as deep on understanding HW and optimization
- Program in Pytorch, OpenCL

ESE532:

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

Penn ESE532 Fall 2021 -- DeHon

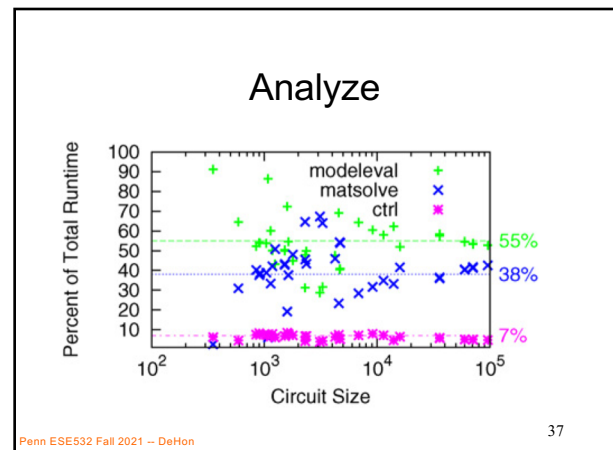
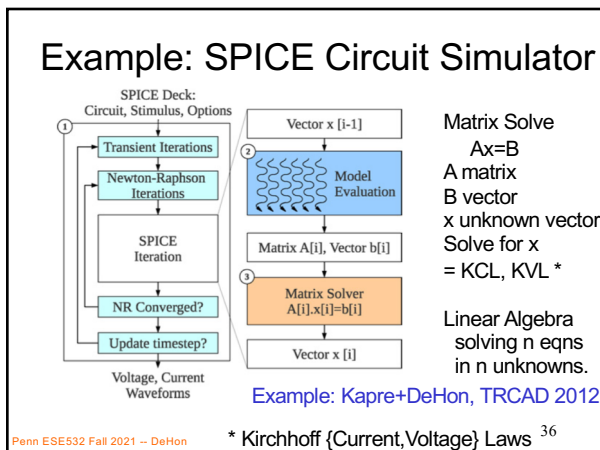
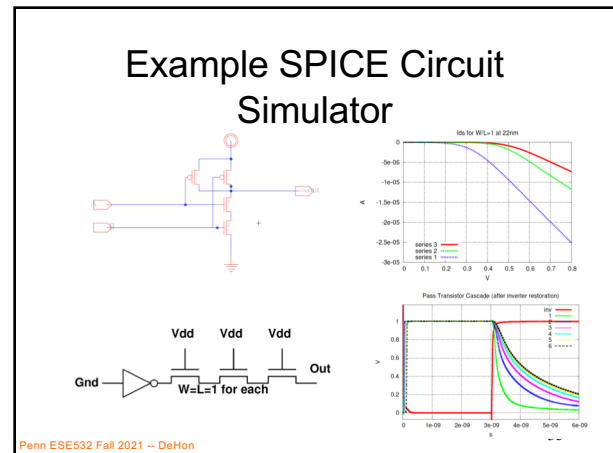
31



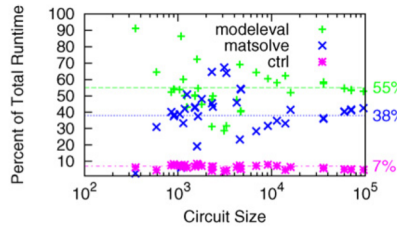
Part 3: Approach -- Example

Penn ESE532 Fall 2021 -- DeHon 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 ESE532 Fall 2021 -- DeHon 34



Analyze

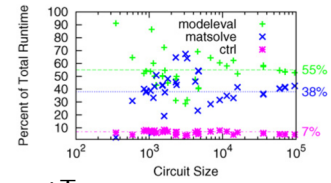


$$T = T_{\text{modelval}} + T_{\text{matsolve}} + T_{\text{ctrl}}$$

Penn ESE532 Fall 2021 -- DeHon

38

Speedup



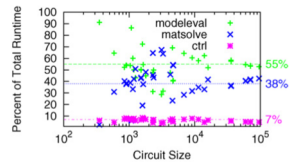
- $T = T_{\text{modelval}} + T_{\text{matsolve}} + T_{\text{ctrl}}$
- What should we speedup first?
- What happens if only speedup modelval?
 - $T = T_{\text{matsolve}} + (T_{\text{modelval}})/S + T_{\text{ctrl}}$

Penn ESE532 Fall 2021 -- DeHon

39

Analyze

- If only accelerated model evaluation only about 2x speedup
- If want better than 14x speed, must also attack control



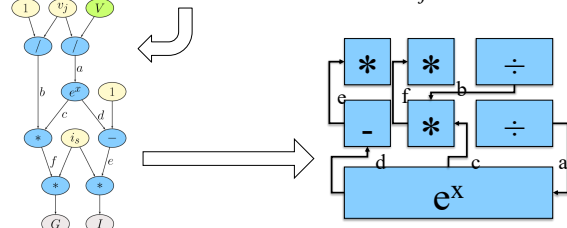
Penn ESE532 Fall 2021 -- DeHon

40

Model Evaluation: Trivial Hardware Implementation

$$I_{D1} = I_s \times (e^{V_{D1}/V_T} - 1)$$

$$G_{D1} = \frac{d}{dV_{D1}}(I_{D1}) = I_s \times e^{V_{D1}/V_T} \times \frac{1}{V_T}$$

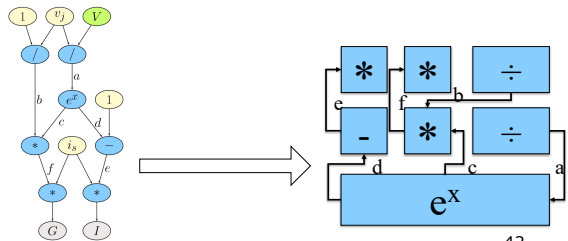


Penn ESE532 Fall 2021 -- DeHon Verilog-AMS as Domain-Specific Language

41

Spatial, Pipelined Parallelism

- Every operation (*, + /) gets dedicated hardware.
- Implement task in space \rightarrow use additional area for each operator.
- Parallel – all operations occur simultaneously.

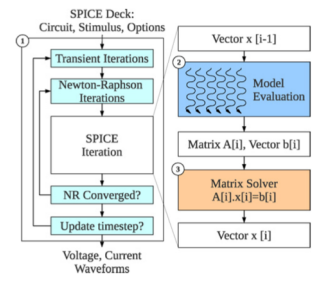


Penn ESE532 Fall 2021 -- DeHon

42

Parallelism: Model Evaluation Data Parallel

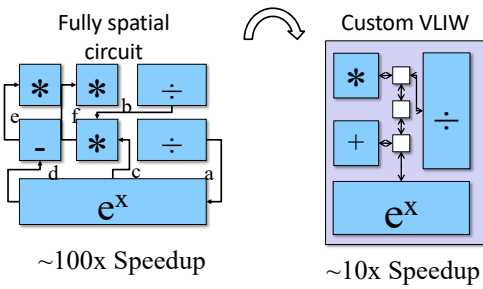
- Every device independent
- Many of each type of device
- Can evaluate in parallel
 - $T = T_{\text{seq}}/N_{\text{proc}}$
- Build pipelined circuit for model
 - $T_{\text{seq}} = N_{\text{comp}} \times T_{\text{cycle}}$
 - vs. $T_{\text{pipe}} = T_{\text{cycle}}$



Penn ESE532 Fall 2021 -- DeHon

43

Spatial Too Big?



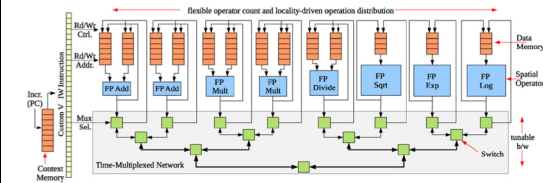
VLIW=Very Long Instruction Word
exploits Instruction-Level Parallelism

44

Penn ESE532 Fall 2021 -- DeHon

Parallelism: Model Evaluation

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

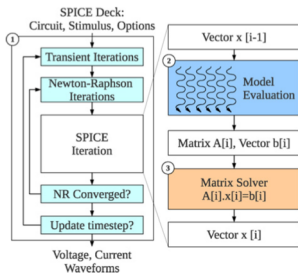


Penn ESE532 Fall 2021 -- DeHon

45

Parallelism: Matrix Solve

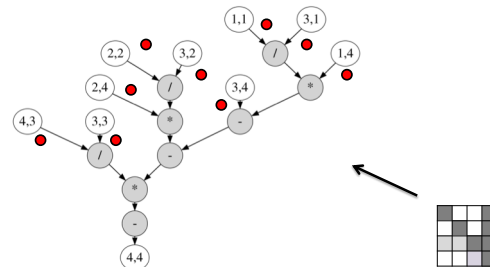
- Needed direct solver?
- E.g. Gaussian elimination
- Data dependence on previous reduce
 - Limited data parallelism
- Parallelism in subtractions
- Some row independence



46

Penn ESE532 Fall 2021 -- DeHon

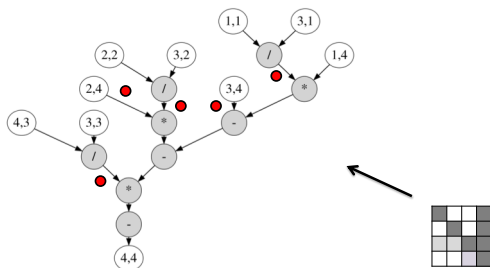
Example Matrix



Penn ESE532 Fall 2021 -- DeHon

47

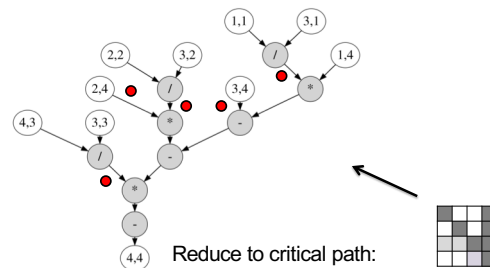
Example Matrix



48

Penn ESE532 Fall 2021 -- DeHon

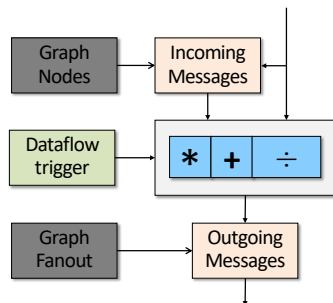
Example Matrix



Penn ESE532 Fall 2021 -- DeHon

49

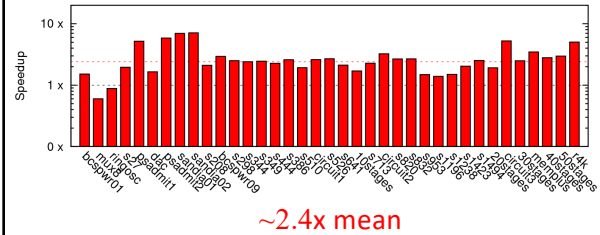
Dataflow Processing Element (PE)



Penn ESE532 Fall 2021 -- DeHon

50

Matrix Solve Only

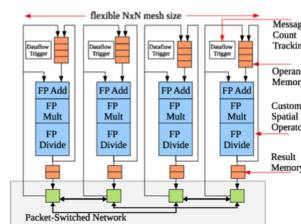


Penn ESE532 Fall 2021 -- DeHon

51

Parallelism: Matrix Solve

- Settled on constructing dataflow graph
- Graph can be iteration independent
 - Statically scheduled
 - (cheaper)
- This is bottleneck to further acceleration

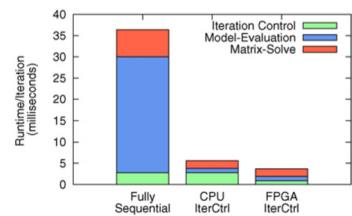


Penn ESE532 Fall 2021 -- DeHon

52

Parallelism Controller?

- Could leave sequential
- For some designs, becomes the bottleneck once others accelerated
- Has internal parallelism in condition evaluation



$$T = T_{\text{model eval}}/S_1 + (T_{\text{matrix solve}}/S_2) + T_{\text{ctrl}}$$

Penn ESE532 Fall 2021 -- DeHon

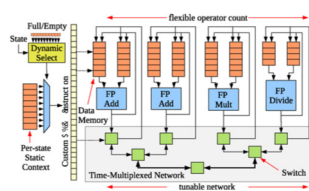
53

Parallelism Controller

- Customized datapath controller

$$T_{\text{seqctrl}} = N_{\text{add}} + N_{\text{mul}} + 10 \cdot N_{\text{divide}}$$

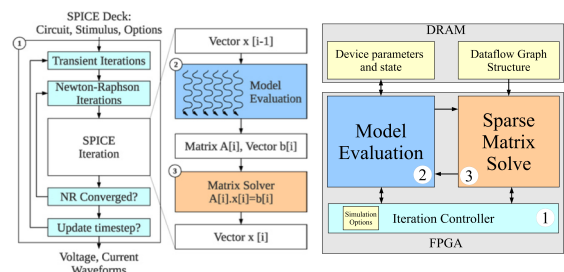
$$T_{\text{vliwctrl}} = \text{Max}(N_{\text{add}}/2, N_{\text{mul}}, 10 \cdot N_{\text{divide}})$$



Penn ESE532 Fall 2021 -- DeHon

54

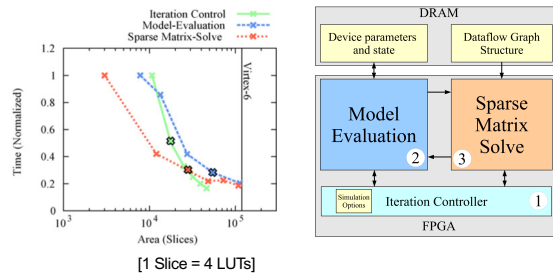
Single-Chip Solution



Penn ESE532 Fall 2021 -- DeHon

55

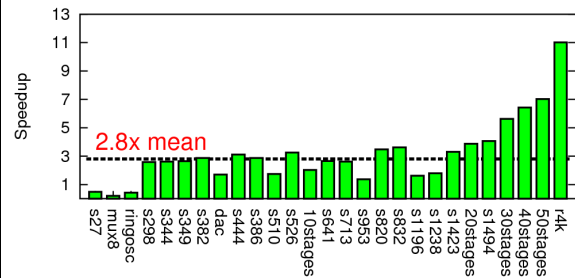
Area-Time for Each



Penn ESE532 Fall 2021 -- DeHon

56

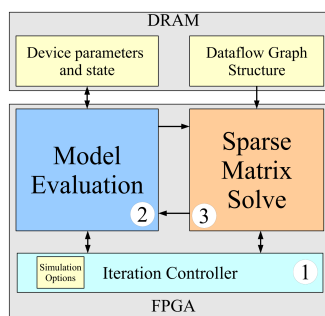
Composite Speedup



Penn ESE532 Fall 2021 -- DeHon

57

Modern SoC



Penn ESE532 Fall 2021 -- DeHon

58

Part 4: Class Components

Penn ESE532 Fall 2021 -- DeHon

59

Class Components

- Lecture (incl. preclass exercise)
 - Slides on web before class
 - (you can print if want a follow-along copy)
 - N.B. I encourage class participation
 - In class/synch recording; 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 F5pm Eastern)
- Project – open-ended (~6 weeks)



Penn ESE532 Fall 2021 -- DeHon

Note syllabus, course admin online

60

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 ESE532 Fall 2021 -- DeHon

61

Second Half

- Use everything on project
- Schedule more tentative
 - Adjust as experience and project demands
- Going deeper
- Real-time
- Reactive
- Memory
- Networking
- Energy
- Scaling
- Chip Cost
- Verification

Penn ESE532 Fall 2021 -- DeHon

62

Teaming

- HW in Groups of 2
- HW: we assign
- Individual assignment writeup
- Project in Groups of 3
- Project: you propose, we review
 - Most portions group writeup
 - Few components individual writeup

Penn ESE532 Fall 2021 -- DeHon

63

Office & Lab Hours

- Andre: T 4:15pm—5:30pm
 - Levine 270, Zoom
 - See canvas
- TAs -- Detkin
 - Tuesday 5 pm
 - Wednesday 5 pm (not today)
 - Thursday 5—7pm (first office hours tomorrow)

Penn ESE532 Fall 2021 -- DeHon

64

C Review

- Course will rely heavily on C
 - Program both hardware and software in C
- HW1 has some C warmup problems
- TAs will hold C review
 - on Sept. 7th, 5:00pm
 - (before our next class meeting since Monday 9/6 is Labor day)

Penn ESE532 Fall 2021 -- DeHon

65

Preclass Exercise

- Motivate the topic of the day
 - Introduce a problem
 - Introduce a design space, tradeoff, transform
- Available before lecture (and in lecture)
 - (only available for 24-48 hours; download)
 - Should work before lecture starts
- Do bring/use calculator
 - Will be numerical examples

Penn ESE532 Fall 2021 -- DeHon

66

Diagnostic Quiz

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

Penn ESE532 Fall 2021 -- DeHon

67

Lecture Timeline

- Preclass available before class
 - In class hardcopy circa 10:10am
- Start lecture at 10:20am
- Lecture until 11:40am
- (most days) stay for remaining questions
 - Pending course after us
- Post video to canvas later in day

Penn ESE532 Fall 2021 – DeHon

68

Feedback

- Will have anonymous feedback {paper, google forms} for each lecture
 - Clarity?
 - Speed?
 - Vocabulary?
 - General comments
- Paper hardcopy for in-person
- Linked on syllabus for not in-person:
 - <https://forms.gle/FWdNWjCsnv6F4pGf8>

Penn ESE532 Fall 2021 – DeHon

69

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 ESE532 Fall 2021 – DeHon

70

Hybrid

- Uncertain how hybrid zoom/in-person will work
 - Setup a bit of a challenge
 - Interactive work?

Penn ESE532 Fall 2021 – DeHon

71

- Your action: **Admin**
 - Find course web page
 - Read it, including the policies
 - Find Syllabus
 - Find 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 piazza group for course
 - Signup for detkin/ketterer access

Penn ESE532 Fall 2021 – DeHon

72

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 ESE532 Fall 2021 – DeHon

73

Questions?