

ESE532: System-on-a-Chip Architecture

Day 1: August 30, 2017
Introduction and Overview



Penn ESE532 Fall 2017 -- DeHon

Today

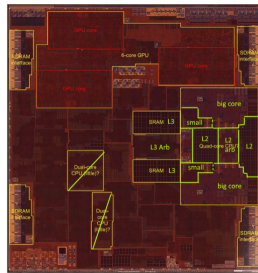
- Case for Programmable SoC
- Goals
- Outcomes
- New/evolving Course, Risks, Tools
- Sample Optimization
- This course (incl. policies, logistics)
- Zed Boards

Penn ESE532 Fall 2017 -- DeHon

2

Today SoC: Apple A10

- 3.3B transistors
- Quad=Dual Dual Core
 - Dual 64-bit ARM 2.3GHz
 - Dual 64-bit ARM low energy
- 3MB L2 cache
- 6 GPU cores
- Custom accelerators
 - Image Processor?
- 125mm² 16nm FinFET



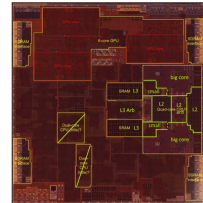
Chipworks Die Photo

3

Penn ESE532 Fall 2017 -- DeHon

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) ?

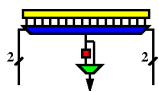


Penn ESE532 Fall 2017 -- DeHon

FPGA

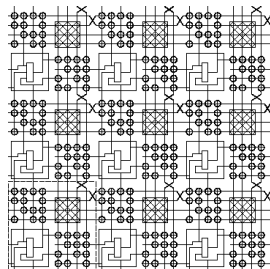
Field-Programmable Gate Array

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



ESE171, CIS371

Penn ESE535 Spring2015 -- DeHon



5

Case for Programmable SoC

Penn ESE532 Fall 2017 -- DeHon

6

The Way things Were

20 years ago

- Wanted programmability
 - used a processor
- 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 2017 -- DeHon

7

Today

- Microprocessor may not be fast enough
 - (but often it is)
 - Or low enough energy
- 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 2017 -- DeHon

8

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 2017 -- DeHon

9

NRE Costs

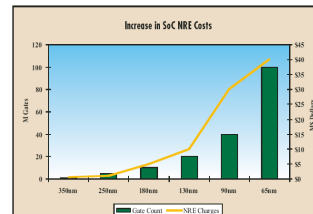


Figure 1 - NRE costs by process geometry Source: Statista Research Corp.

- **28-nm SoC development costs doubled over previous node – EE Times 2013**
 28nm+78%, 20nm+48%, 14nm+31%, 10nm+35%

Penn ESE532 Fall 2017 -- DeHon

10

Amortize NRE with Volume

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

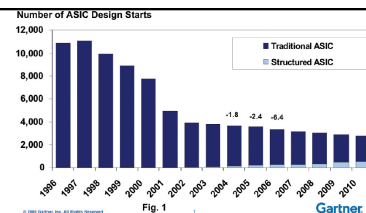
$$Cost = \frac{Cost_{NRE}}{N_{chips}} + Cost_{perchip}$$

Penn ESE532 Fall 2017 -- DeHon

11

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

Penn ESE532 Fall 2017 -- DeHon

12

Large ICs

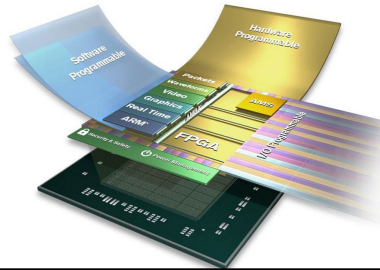
- Now contain significant software
 - Almost all have embedded processors
- Must co-design SW and HW
- Must solve complete computing task
 - Tasks has components with variety of needs
 - Some don't need custom circuit
 - 90/10 Rule

Given Demand for Programmable

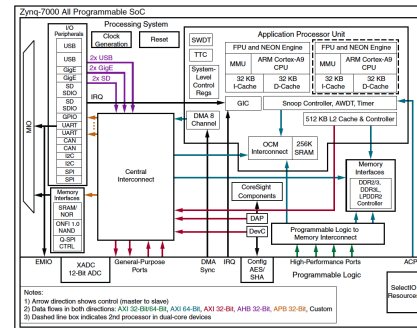
- How do we get higher performance than a processor, while retaining programmability?

Programmable SoC

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



Programmable SoC



Then and Now

20 years ago

- Programmability?
 - use a processor
- 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
- High-throughput
 - FPGA, GPU, PSoC, custom
- Wanted product differentiation
 - Program FPGAs, PSoC
- Build a custom IC
 - System and software

Goals, Outcomes

Goals

- Create Computer Engineers
 - SW/HW divide is wrong, outdated
 - 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 2017 -- DeHon

19

Roles

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

Penn ESE532 Fall 2017 -- DeHon

20

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 ESE532 Fall 2017 -- DeHon

21

Outcomes

- Understand the system on a chip from gates to application software, including:
 - on-chip memories and communication networks, I/O interfacing, RTL 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 2017 -- DeHon

22

First offering last term: warning last term

- We'll be making it up as we go along...
- You'll be first to perform assignments
- We may estimate difficulty of assignments incorrectly
 - Too easy, too hard
 - Many were too tedious
 - Provided wrong guidance
- Lectures will be less polished
- Intellectual excitement of trying to figure it out... and currency

Penn ESE532 Fall 2017 -- DeHon

23

This Term

- ...still figuring it out.
- Learned some lessons from last term.
- Refocusing assignments
- Increase focus on modeling throughout
- May still get it wrong
 - We change the assignments, so they are still new...
- It will be better...but probably not perfect
- It will be hard work...hopefully not insane.

Penn ESE532 Fall 2017 -- DeHon

24

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
 - Bonus points for tutorials

Relation to ESE532

534

- Deep into design space and continuum
- How to build compute, interconnect, memory
- Analysis
- Fundamentals
 - Theory
 - Why X better than Y
- More relevant substrate designers
- Both Real-Time and Best-effort

532 – System-on-a-Chip Architecture

- New course
- Probably 30% overlap
- Broader (all CMPE)
 - HW/SW codesign
- More Hands-on
 - Code in C
 - Map to Zynq
 - Accelerate an application
- More relevant to (P)SoC user
- Real-time focus

Distinction

CIS240, 371, 501

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

ESE532

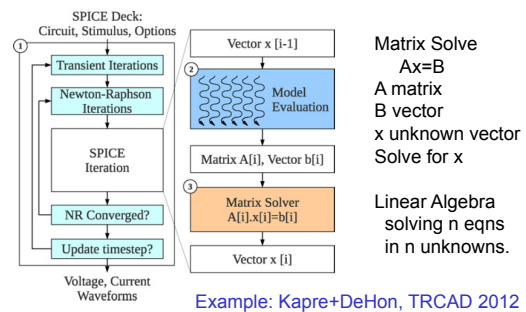
- Hardware-Software codesign
 - Willing to recompile, maybe rewrite code
 - Define/refine hardware
- Real-Time
 - Guarantee meet deadline
- Non shared-memory models

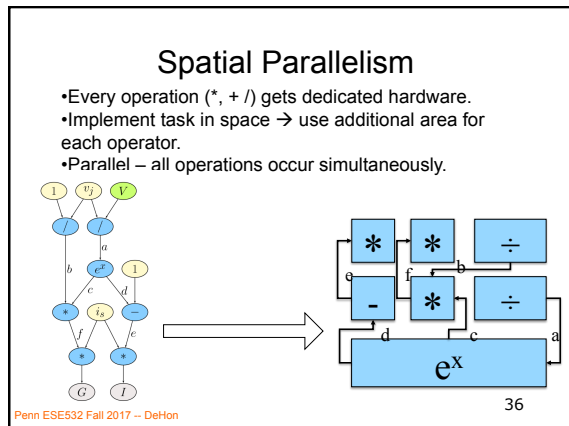
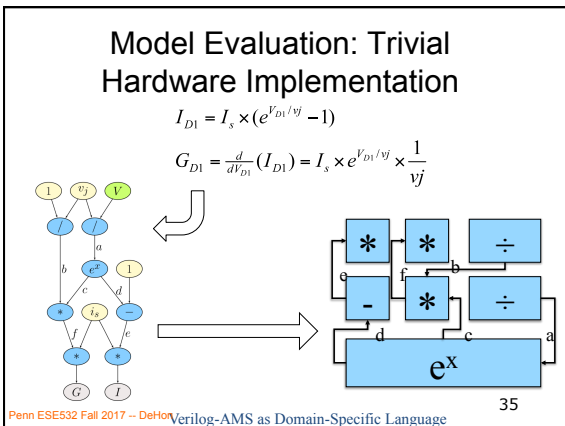
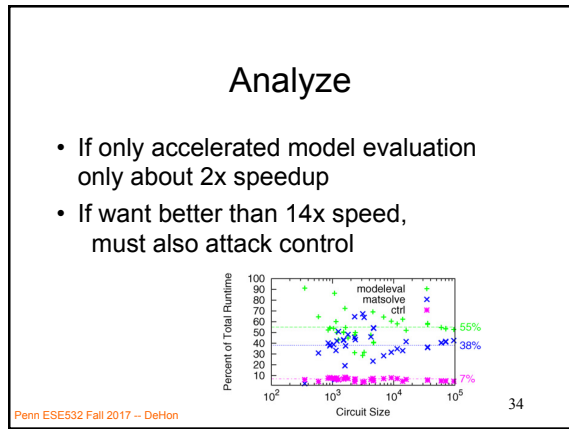
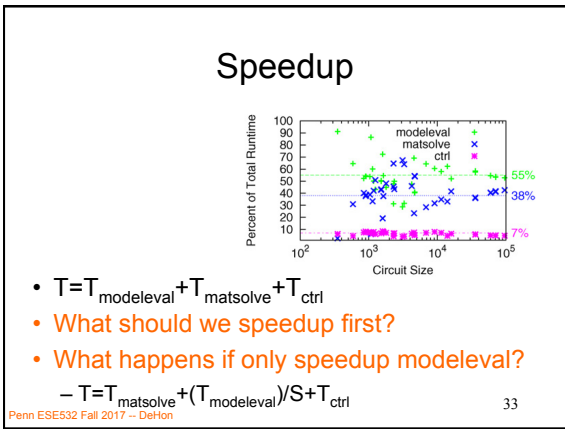
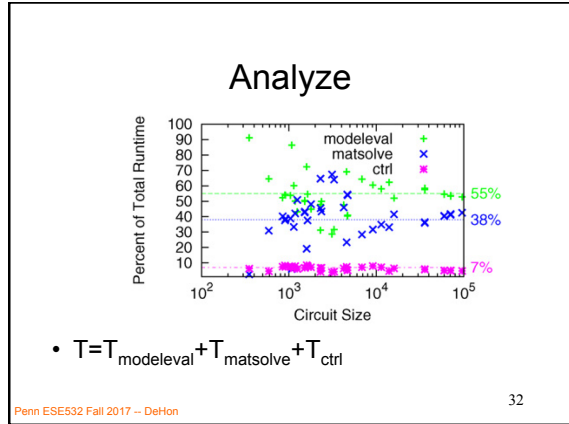
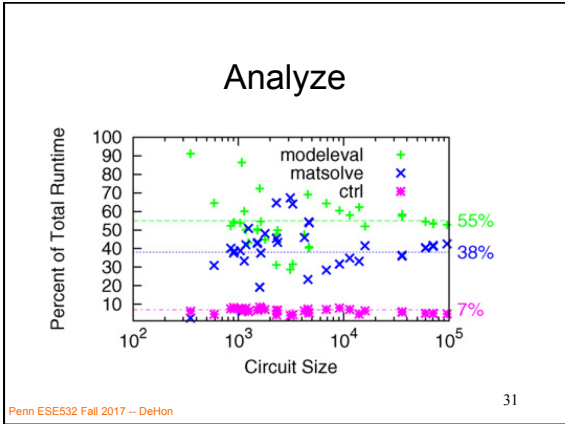
Approach -- Example

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

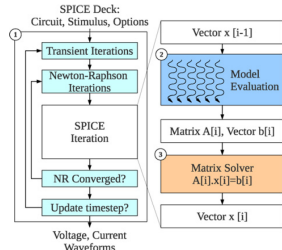
SPICE Circuit Simulator





Parallelism: Model Evaluation

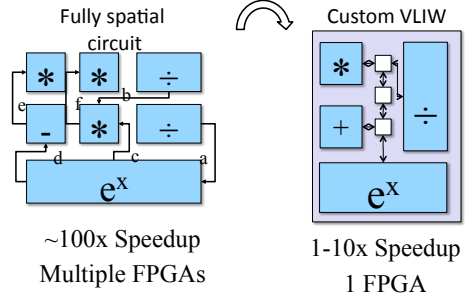
- Every device independent
- Many of each type of device
- Can evaluate in parallel
 - $T = T_{seq} / N_{proc}$
- Build pipelined circuit for model
 - $T_{seq} = N_{comp} * T_{cycle}$
 - vs. $T_{pipe} = T_{cycle}$



Penn ESE532 Fall 2017 -- DeHon

37

Single FPGA Mapping



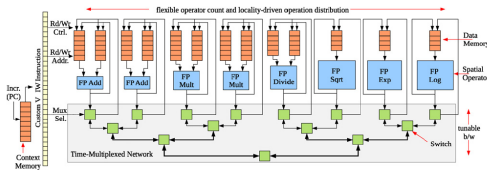
VLIW=Very Long Instruction Word
exploits Instruction-Level Parallelism

Penn ESE532 Fall 2017 -- DeHon

38

Parallelism: Model Evaluation

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

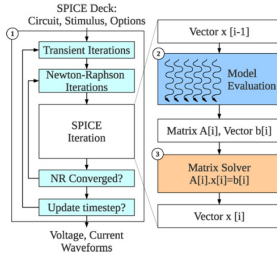


Penn ESE532 Fall 2017 -- DeHon

39

Parallelism: Matrix Solve

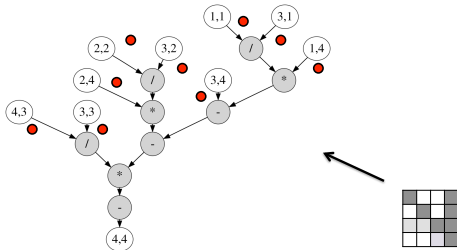
- Needed direct solver?
- E.g. Gaussian elimination
- Data dependence on previous reduce
- Parallelism in subtracts
- Some row independence



Penn ESE532 Fall 2017 -- DeHon

40

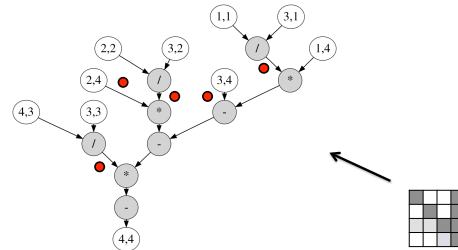
Example Matrix



Penn ESE532 Fall 2017 -- DeHon

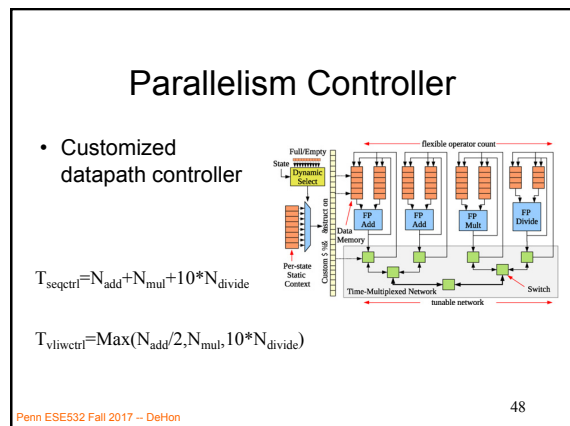
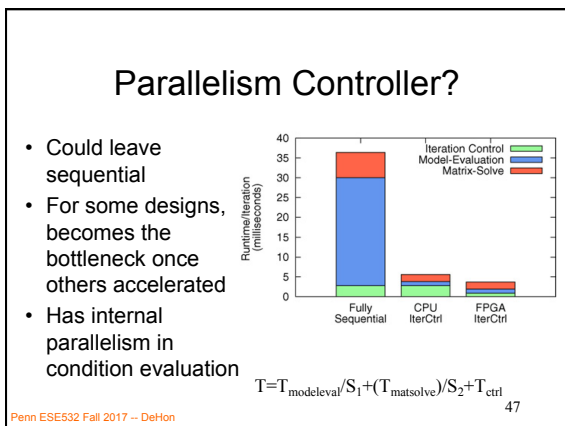
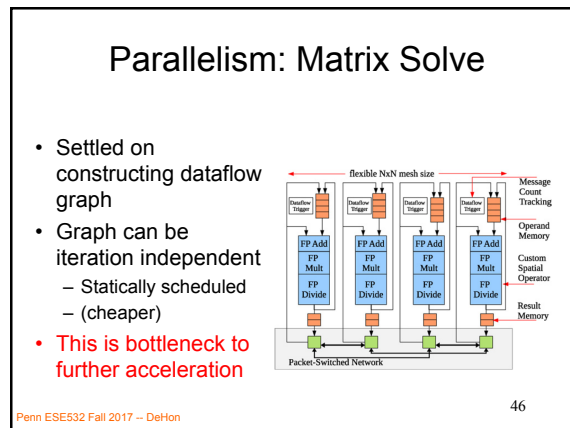
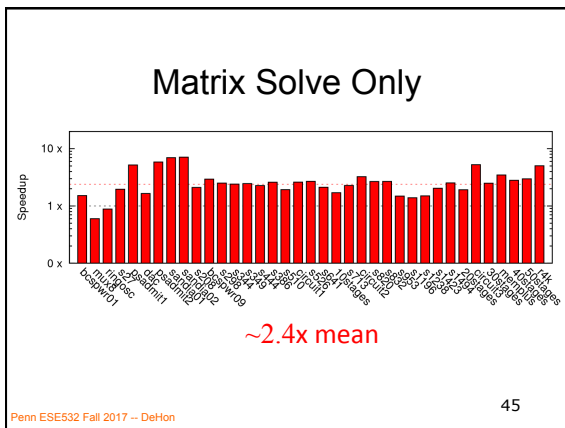
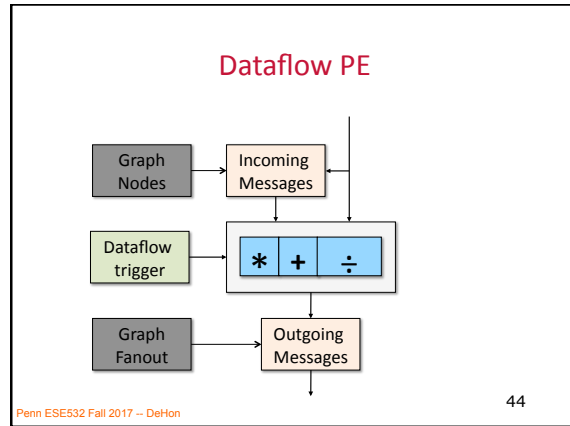
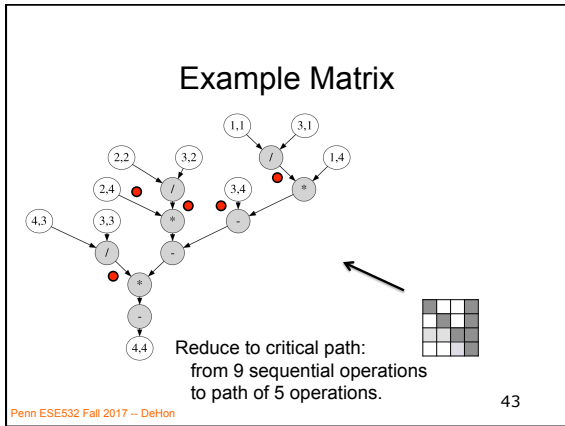
41

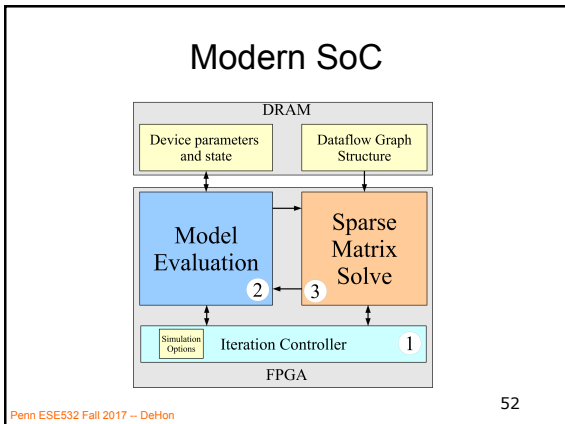
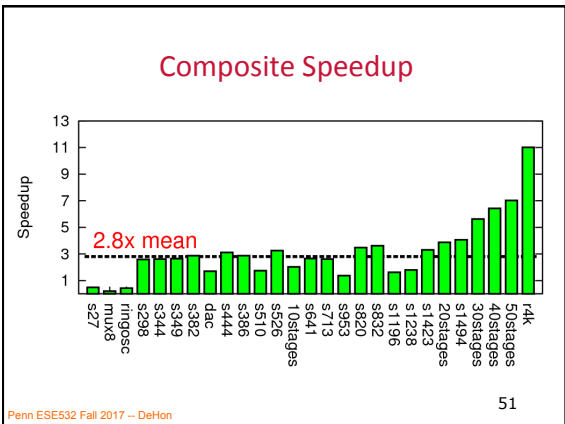
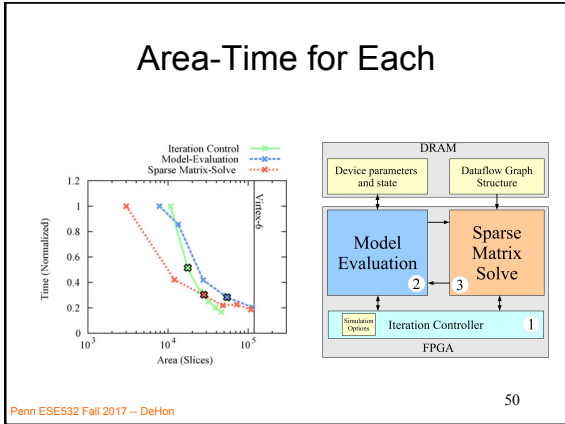
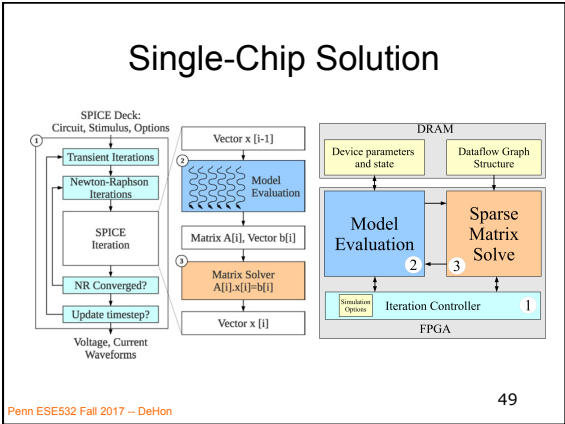
Example Matrix



Penn ESE532 Fall 2017 -- DeHon

42





Class Components

Penn ESE532 Fall 2017 -- DeHon

53

- ### Class Components
- Lecture (incl. preclass exercise)
 - Slides on web before class
 - (you can print if want a follow-along copy)
 - N.B. I will encourage (force) class participation
 - Questions, polls
 - Reading [~1 required paper/lecture]
 - online: Canvas, IEEE, ACM, also ZynqBook
 - Homework
 - (1 per due F5pm)
 - Project – open-ended
 - (~6 weeks)
 - **Note syllabus, course admin online**
-
- Penn ESE532 Fall 2017 -- DeHon
- 54

First Half

- Quickly cover breadth
- Metrics, bottlenecks
- Memory
- Parallel models
- SIMD/Data Parallel
- Thread-level parallelism
- Spatial, C-to-gates
- Real-time
- Reactive
- Line up with homeworks

Penn ESE532 Fall 2017 -- DeHon

55

Second Half

- Use everything on project
- Schedule more tentative
 - Adjust as experience and project demands
- Going deeper
- Memory
- Networking
- Energy
- Scaling
- Chip Cost
- Defect + Fault tolerance

Penn ESE532 Fall 2017 -- DeHon

56

Teaming

- HW and Project in Groups of 2
- First assignment
 - You choose, optionally individual
- HW2-7: we assign
- Project: you propose, we review
- Individual assignment turnin
 - See assignment details if individual pieces

Penn ESE532 Fall 2017 -- DeHon

57

Office & Lab Hours

- Andre: T 4:15pm—5:30pm Levine 270
- Hans: TR 6—7pm in Ketterer
 - Start tomorrow 8/31

Penn ESE532 Fall 2017 -- DeHon

58

C Review

- Course will rely heavily on C
 - Program both hardware and software in C
- HW1 has some C warmup problems
- Hans will hold C review
 - Ketterer on Sept. 5th at 6pm
 - (before our next class meeting since Monday 9/4 is Labor day)

Penn ESE532 Fall 2017 -- DeHon

59

Preclass Exercise

- Motivate the topic of the day
 - Introduce a problem
 - Introduce a design space, tradeoff, transform
- Work for ~5 minutes before start lecturing
- Do bring calculator class
 - Will be numerical examples

Penn ESE532 Fall 2017 -- DeHon

60

Feedback

- Will have anonymous feedback sheets for each lecture
 - Clarity?
 - Speed?
 - Vocabulary?
 - General comments

Penn ESE532 Fall 2017 -- DeHon

61

Policies

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

Penn ESE532 Fall 2017 -- DeHon

62

Admin

- Your action:
 - 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

Penn ESE532 Fall 2017 -- DeHon

63

Big Ideas

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

Penn ESE532 Fall 2017 -- DeHon

64

Zedboard Checkout

- Ketterer (Moore 200) – card key access
- Lockers: EOS, B13, 8a-3.30p
 - strongly recommending Masterlock combination lock (with black dial) available at Amazon <http://a.co/5YfVLJt>

Penn ESE532 Fall 2017 -- DeHon

65