

ESE532: System-on-a-Chip Architecture

Day 14: March 13, 2017
Design Space Exploration



Today

- Midterm
- Design-Space Exploration

Message

- The universe of possible implementations (design space) is large
 - Many dimensions to explore
- Formulate carefully
- Approach systematically

Midterm

Mea Culpa

- Too long
 - Too time constrained

Mismatch?

- What I think I'm teaching
- What you think you are learning

- Midterm focused on performance modeling

Agree?

- How to design/select/map to SoC to reduce Energy/Area/Delay?

Penn ESE532 Spring 2017 -- DeHon

7

Day 1

Outcomes

- Design, optimize, and program a modern System-on-a-Chip.
- Analyze, identify bottlenecks, design-space
- 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 Spring 2017 -- DeHon

8

Day 1

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 Spring 2017 -- DeHon

9

Agree?

- How to design/select/map to SoC to reduce Energy/Area/Delay?
 - Haven't touched Area/Energy much, yet, so initial focus on Delay, Performance

Penn ESE532 Spring 2017 -- DeHon

10

Reduce Delay

- What can we do to reduce delay?
- How know when we succeed?
 - Compile and run
 - How know what's wrong and how to fix if we don't get a speedup?
- **Need to be able to model, estimate, and reason about performance.**

Penn ESE532 Spring 2017 -- DeHon

11

Model to Select Design

- Homework shows you that implementing a design is expensive
- **Goal:** have a model to drive you to select good implementations to try

Penn ESE532 Spring 2017 -- DeHon

12

Model to Understand

- If not get the performance you expect, what went wrong?
 - Need an expectation
 - Need to drill down into components
 - Model expectations wrong?
 - Which part?
 - Refine?
 - Something went wrong in implementations?
 - What?
 - How fix to match model?

Penn ESE532 Spring 2017 -- DeHon

13

Model to Understand

- If not get the performance you expect, what went wrong?
- **Very Powerful**
- Not match → something I need to understand
 - Bugs in implementation, understanding
- Does match → Tells me
 - where to look
 - what to try to fix
 - when to give up (on a particular approach)

Penn ESE532 Spring 2017 -- DeHon

14

Expect

- Kinds of analysis on the exam underlies all design exploration and experiments

Penn ESE532 Spring 2017 -- DeHon

15

Day 2

Message for Day

- Identify the Bottleneck
 - May be in compute, I/O, memory, data movement
- Focus and reduce/remove bottleneck
 - More efficient use of resources
 - More resources
- Repeat

Penn ESE532 Spring 2017 -- DeHon

16

Second Half

- How to design/select/map to SoC to reduce Energy/Area/Delay
- Spend some more time understanding
 - Area
 - Started as we look at accelerators
 - ...worry about fit in fixed array
 - Cost for custom design
 - Energy

Penn ESE532 Spring 2017 -- DeHon

17

Design-Space Exploration

Generic

Penn ESE532 Spring 2017 -- DeHon

18

Design Space

- Have many choices for implementation
 - Alternatives to try
 - Parameters to tune
 - Mapping options
- Our freedom to impact implementation costs
 - Area, delay, energy

Penn ESE532 Spring 2017 -- DeHon

19

Design Space

- Ideally
 - Each choice orthogonal axis in high-dimensional space
 - Want to understand points in space
 - Find one that best meets constraints and goals
- Practice
 - Seldom completely orthogonal
 - Requires cleverness to identify dimensions
 - Messy, cannot fully explore
 - But...can understand, priorities, guide

Penn ESE532 Spring 2017 -- DeHon

20

Preclass 1

- What choices (design-space axes) can we explore in mapping a task to an SoC?
- What showed up in homework so far?

Penn ESE532 Spring 2017 -- DeHon

21

From Homework?

- Types of parallelism
- Mapping to different fabrics / hardware
- How manage memory, move data
- Levels of parallelism
- Pipelining, unrolling, II

Penn ESE532 Spring 2017 -- DeHon

22

Design-Space Choices

- Type of parallelism
- How decompose / organize parallelism
- Area-time points (level exploited)
- What resources we provision for what parts of computation
- Where to map tasks
- How schedule/order computations
- How synchronize tasks
- How represent data
- Where place data; how manage and move
- What precision use in computations

Penn ESE532 Spring 2017 -- DeHon

23

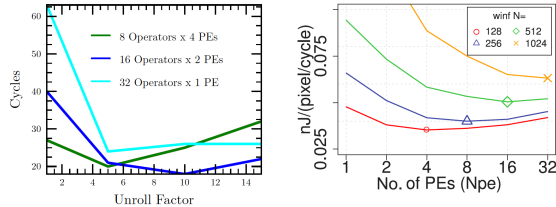
Generalize Continuum

- Encourage to think about parameters (axes) that capture continuum to explore
- Start from an idea
 - Maybe can compute with 8b values
 - Maybe can put dist computation on FPGA fabric
 - Move data in 1KB chunks
- Identify general knob
 - Tune intermediate bits for computation
 - How much of computation go on FPGA fabric
 - What is optimal data transfer size?

Penn ESE532 Spring 2017 -- DeHon

24

Finding Optima



- Kapre, FPL 2009
- Kadric, TRETs 2016

Penn ESE532 Spring 2017 -- DeHon

25

Design Space Explore

- Think systematically about how might map the application
- Avoid overlooking options
- Understand tradeoffs
- Large design space
 - more opportunities to find good solutions
 - Reduce bottlenecks

Penn ESE532 Spring 2017 -- DeHon

26

Elaborate Design Space

- Refine design space as you go
- Ideally identify up front
- Practice bottlenecks and challenges
 - will suggest new options / dimensions
 - If not initially expect memory bandwidth to be a bottleneck...
- Some options only make sense in particular sub-spaces
 - Bitwidth optimization not a big issue on the 64b processor
 - More interesting on vector, FPGA

Penn ESE532 Spring 2017 -- DeHon

27

Tools

- Sometimes tools will directly help you explore design space
 - What SDSoC/Vivado HLS support?
- Often they will not
 - What might you want that does not support?

Penn ESE532 Spring 2017 -- DeHon

28

Tools

- Sometimes tools will directly help you explore design space
 - Unrolling, pipelining, II
 - Some choices for data movement
 - Some loop transforms
 - Granularity to place on FPGA
- Often they will not
 - Need to reshape functions and loops
 - Data representations and sizes

Penn ESE532 Spring 2017 -- DeHon

29

Design-Space Exploration

Example FFT

Penn ESE532 Spring 2017 -- DeHon

30

Fourier Transform

- Identify spectral components
- Convert from Time-domain to Frequency-domain
 - E.g. tones from data samples
 - Central to audio coding – e.g. MP3 audio

$$Y[k] = \sum_{j=0}^{n-1} (X[j]e^{-2i\pi \frac{k}{n}})$$

Penn ESE532 Spring 2017 --

31

Fast-Fourier Transform (FFT)

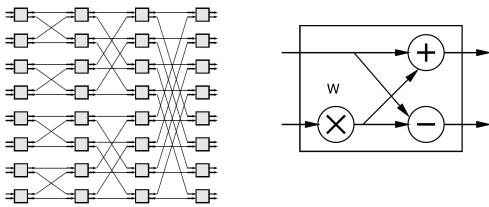
- Efficient way to compute FT
- $O(N \cdot \log(N))$ computation

Penn ESE532 Spring 2017 -- DeHon

32

FFT

- Large space of FFTs
- Radix-2 FFT Butterfly

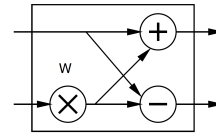


Penn ESE532 Spring 2017 -- DeHon

33

Basic FFT Butterfly

- $Y_0 = X_0 + W(\text{stage, butterfly}) * X_1$
- $Y_1 = X_0 - W(\text{stage, butterfly}) * X_1$
- Common sub expression, compute once: $W(\text{stage, butterfly}) * X_1$

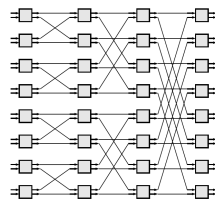


Penn ESE532 Spring 2017 -- DeHon

34

Preclass 2

- What parallelism options exist?
 - Single FFT
 - Sequence of FFTs

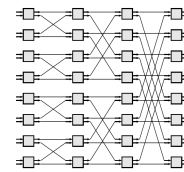


Penn ESE532 Spring 2017 -- DeHon

35

FFT Parallelism

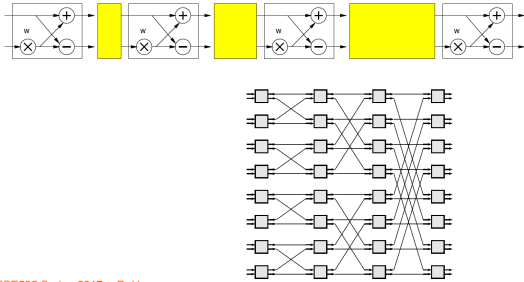
- Spatial
- Pipeline
- Streaming
- By column
 - Choose how many Butterflies to serialize on a PE
- By subgraph
- Pipeline subgraphs



Penn ESE532 Spring 2017 -- DeHon

36

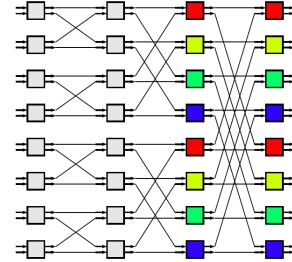
Streaming FFT



Penn ESE532 Spring 2017 -- DeHon

37

Common Subgraphs

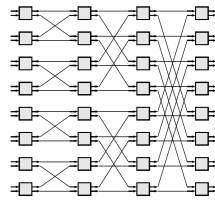


Penn ESE532 Spring 2017 -- DeHon

38

Preclass 3

- How large of a spatial FFT can implement with 220 multipliers?



Penn ESE532 Spring 2017 -- DeHon

39

Bit Serial

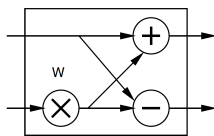
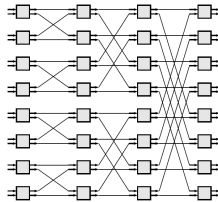
- Could compute the add/multiply bit serially
 - One full adder per adder
 - W full adders per multiply
 - 50,000 LUTs
 - \approx 2500 bit-serial butterflies for $W=16$?
- Another dimension:
 - How much serialize word-wide operators

Penn ESE532 Spring 2017 -- DeHon

40

Accelerator Building Blocks

- What might we use as primitive, FFT-specific building blocks?

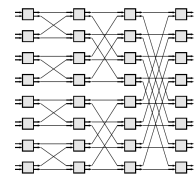


Penn ESE532 Spring 2017 -- DeHon

41

Processor Mapping

- How map butterfly operations to processors?
 - Implications for communications?

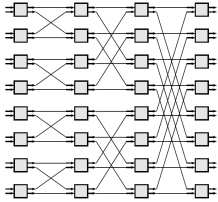


Penn ESE532 Spring 2017 -- DeHon

42

Preclass 4a

- How large local memory to communicate from stage to stage?

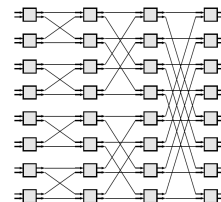


Penn ESE532 Spring 2017 -- DeHon

43

Preclass 4b

- How change evaluation order to reduce local storage memory?



Penn ESE532 Spring 2017 -- DeHon

44

Communication

- How implement the data shuffle between processors or accelerators?
 - Memories / interconnect ?
 - How serial / parallel ?
 - Network?

Penn ESE532 Spring 2017 -- DeHon

45

Data Precision

- Input data from A/D likely 12b
- Output data, may only want 16b
- What should internal precision and representation be?

Penn ESE532 Spring 2017 -- DeHon

46

Number Representation

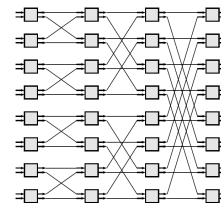
- Floating-Point
 - IEEE standard single (32b), double (64b)
 - With mantissa and exponent
 - ...half, quad
- Fixed-Point
 - Select total bits and fraction
 - E.g. 16.8 (16 total bits, 8 of which are fraction)
 - Represent $1/256$ to $256 \cdot 1/256$

Penn ESE532 Spring 2017 -- DeHon

47

Heterogeneous Precision

- May not be same in every stage
 - W factors less than 1
 - Non-fraction grows at most 1b per stage



Penn ESE532 Spring 2017 --

48

W/Twiddle factors

- Precompute and store in arrays
- Compute as needed
 - How? sin/cos hardware? CORDIC?
Polynomial approximation?
- Specialize into computation
 - Many evaluate to 0, ± 1 , $\pm \frac{1}{2}$,

Penn ESE532 Spring 2017 -- DeHon

49

FFT (partial) Design Space

- Parallelism
- Decompose
- Size/granularity of accelerator
 - Area-time
- Sequence/share
- Communicate
- Representation/precisions
- Twiddle

Penn ESE532 Spring 2017 -- DeHon

50

Big Ideas:

- Large design space for implementations
- Worth elaborating and formulating systematically
 - Make sure don't miss opportunities
- Think about continuum for design axes

Penn ESE532 Spring 2017 -- DeHon

51

Admin

- HW7 out → due Friday
 - Individual
- Working on getting Project ready

Penn ESE532 Spring 2017 -- DeHon

52