# ESE532:
## System-on-a-Chip Architecture

Day 19: November 6, 2017
Design Space Exploration

Penn

---

# Today

- Design-Space Exploration
  - Generic
  - Fast Fourier Transform (FFT)

---

# Message

- The universe of possible implementations (design space) is large
  - Many dimensions to explore
- Formulate carefully
- Approach systematically
- Use modeling along the way for guidance

---

# Design-Space Exploration

Generic

---

# Design Space

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

---

# Design Space

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

1

## Preclass 1

- What choices (design-space axes) can we explore in mapping a task to an SoC?

- What showed up in homework so far?

## 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
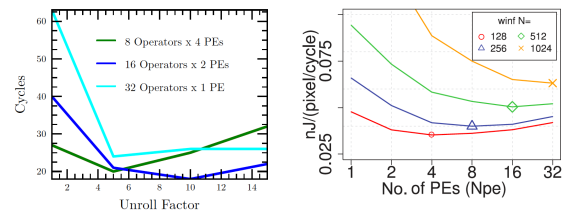- Where place data; how manage and move
- What precision use in computations

## 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?

## Finding Optima



- Kapre, FPL 2009
- Kadric, TRETS 2016

## 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

## 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

## Tools

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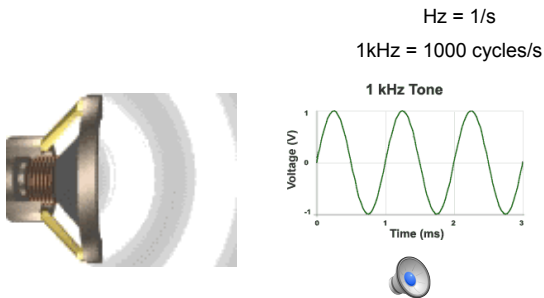  - What might you want that does not support?

13

## Design-Space Exploration

Example FFT

14

## Sound Waves

Hz = 1/s

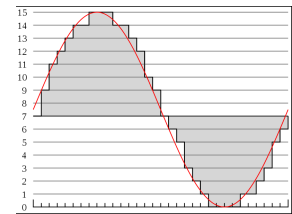1kHz = 1000 cycles/s

**1 kHz Tone**



Source: http://www.mediacollege.com/audio/01/sound-waves.html
15

## Discrete Sampling

- Represent as time sequence
- Discretely sample in time
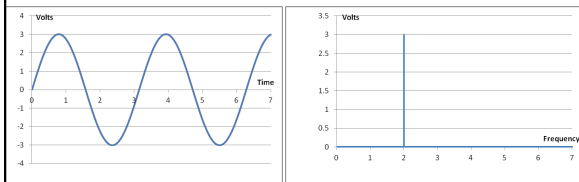- What we can do directly with an Analog-to-Digital (A2D) converter.



http://en.wikipedia.org/wiki/File:Pcm.svg

16

## Frequency-domain

- $T = \pi$, $A = 3$: $s(t) = A*\sin(2\pi*f*t) = 3*\sin(2*t)$

17

## Frequency-domain

- Can represent sound wave as linear sum of frequencies

18

3

# Time vs. Frequency

19

# Fourier Series – Why does it



sin(t)

sin(3t)

sin(2t)

The cos(nx) and sin(nx) functions form an orthogonal basis: they allow us to represent any periodic signal by taking a linear combination of the basis components without interfering with one another

20

# Fourier Transform

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

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

21

# FT as Matching

- Fourier Transform is essentially performing a dot product with a frequency
  - How much like a sine wave of freq. f is this?

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

22

# Fast-Fourier Transform (FFT)

- Efficient way to compute FT
- O(N*log(N)) computation

23

# FFT
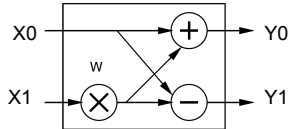
- Large space of FFTs
- Radix-2 FFT Butterfly

24

4

## Basic FFT Butterfly

- $Y0 = X0 + W(stage, butterfly) * X1$
- $Y1 = X0 - W(stage, butterfly) * X1$
- Common sub expression, compute once: $W(stage, butterfly) * X1$

25

## Preclass 2

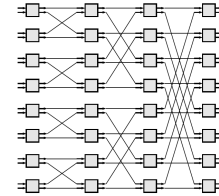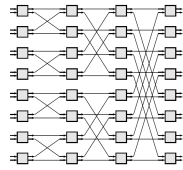- What parallelism options exist?
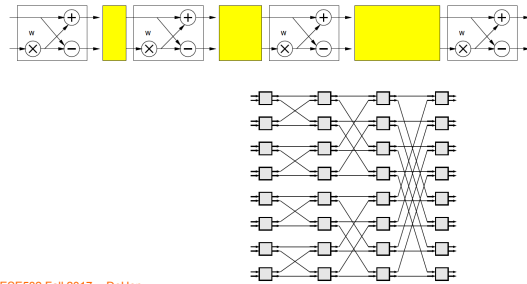  - Single FFT
  - Sequence of FFTs

26

## FFT Parallelism



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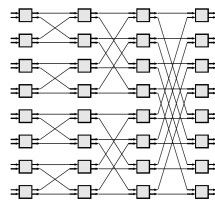- By subgraph
- Pipeline subgraphs

27

## Streaming FFT

28

## Preclass 3

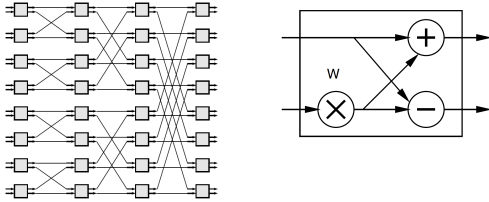- How large of a spatial FFT can implement with 220 multipliers?

29

## Bit Serial

- Could compute the add/multiply bit serially
  - One full adder per adder
  - W full adders per multiply
  - 50,000 LUTs
    - ~= 2500 bit-serial butterflies for W=16?
      - Maybe 512-point FFT?
- Another dimension:
  - How much serialize word-wide operators

30

5

## Accelerator Building Blocks

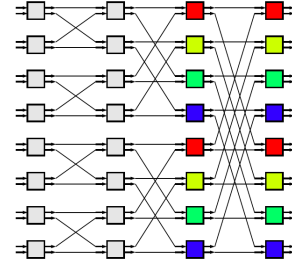- What might we use as primitive, FFT-specific building blocks?

## Common Subgraphs

## Processor Mapping

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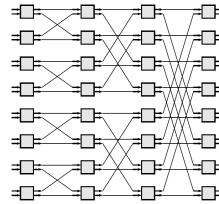  - Implications for communications?

## Preclass 4a

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

## Preclass 4b

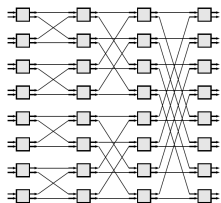- How change evaluation order to reduce local storage memory?

## Communication

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

## Data Precision

- Input data from A2D likely 12b
- Output data, may only want 16b
- What should internal precision and representation be?
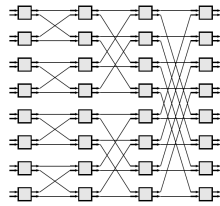
## Number Representation

- 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 256-1/256

## Heterogeneous Precision

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

## W/Twiddle factors

- Precompute and store in arrays
- Compute as needed
  - How?  sin/cos hardware? CORDIC? Polynominal approximation?
- Specialize into computation
  - Many evaluate to 0, ±1, ±½, ….

## FFT (partial) Design Space

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

## Big Ideas:

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

# Admin

- 1st milestone parallelism feedback
  – Should have seen Friday evening
- 2nd milestone due Friday
  – Asks you to identify design space

8