

# ESE532: System-on-a-Chip Architecture

Day 18: October 31, 2018  
Design Space Exploration



## 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

## Preclass 1

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

## From Homework?

- Types of parallelism
- Mapping to different fabrics / hardware
- How manage memory, move data
  - DMA, streaming
  - Data access patterns
- Levels of parallelism
- Pipelining, unrolling, II, array partitioning and packing

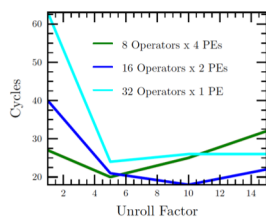
## Design-Space Choices

- 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

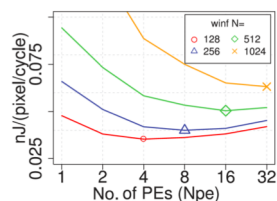
## Generalize Continuum

- Encourage to think about parameters (axes) that capture continuum to explore
- Start from an idea
  - Maybe can compute with 8b values
  - Maybe can put matrix-mpy 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
- The larger the 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

Penn ESE532 Fall 2018 -- DeHon

13

## 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 Fall 2018 -- DeHon

14

## 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 Fall 2018 -- DeHon

15

## Code for Exploration

- Can you write your code with parameters (#define) can easily change to explore continuum?
  - Unroll factor?
  - Number of parallel tasks?
  - Size of data to move?

Penn ESE532 Fall 2018 -- DeHon

16

## Design-Space Exploration

Example FFT

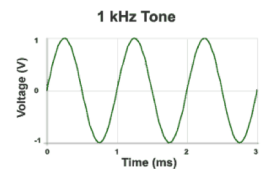
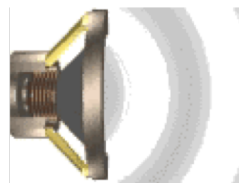
Penn ESE532 Fall 2018 -- DeHon

17

## Sound Waves

Hz = 1/s

1kHz = 1000 cycles/s



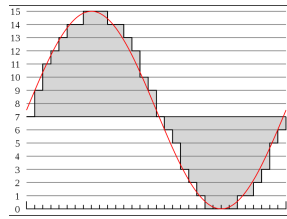
Source: <http://www.mediacollege.com/audio/01/sound-waves.html>

Penn ESE532 Fall 2018 -- DeHon

18

## Discrete Sampling

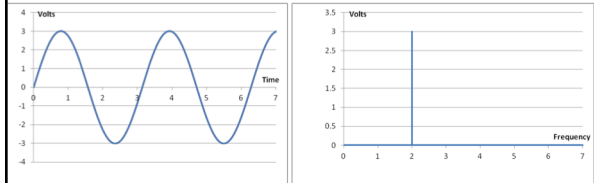
- Represent as time sequence
- Discretely sample in time
- What we can do directly with an Analog-to-Digital (A2D) converter.



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

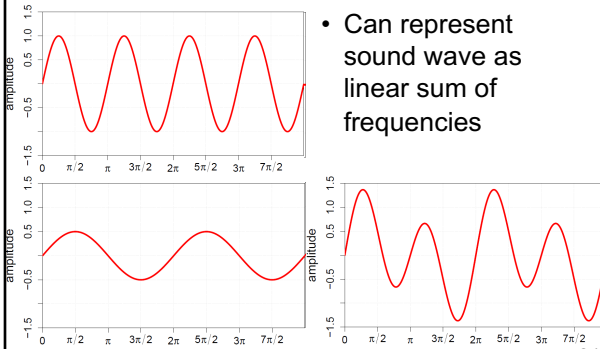
## Frequency-domain

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

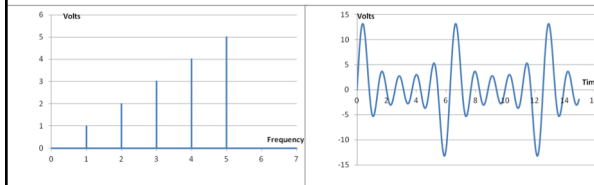


## Frequency-domain

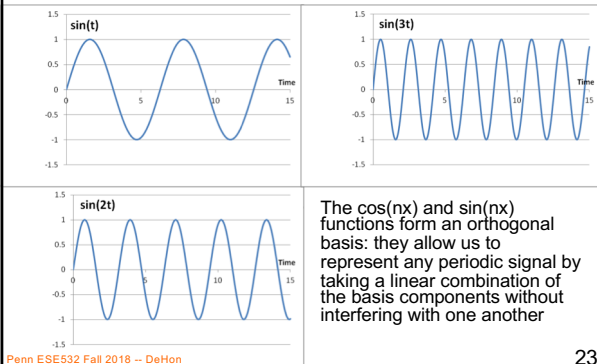
- Can represent sound wave as linear sum of frequencies



## Time vs. Frequency



## Fourier Series



## 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

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

## 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?

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

Penn ESE532 Fall 2018 -- De

25

## Fast-Fourier Transform (FFT)

- Efficient way to compute FT
- $O(N \log(N))$  computation
- Contrast  $N^2$  for direct computation
  - $N$  dot products
    - Each dot product has  $N$  points (multiply-adds)

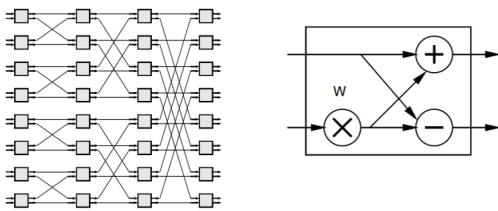
$$Y[k] = \sum_{j=0}^{n-1} \left( X[j] e^{-2i\pi \frac{k}{n} j} \right)$$

Penn ESE532 Fall 2018 -- De

26

## FFT

- Large space of FFTs
- Radix-2 FFT Butterfly

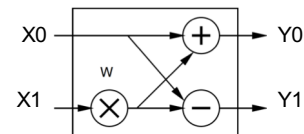


Penn ESE532 Fall 2018 -- DeHon

27

## 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:  $W(\text{stage, butterfly}) * X_1$

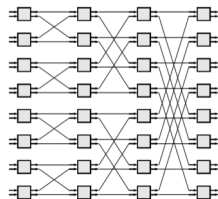


Penn ESE532 Fall 2018 -- DeHon

28

## Preclass 2

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

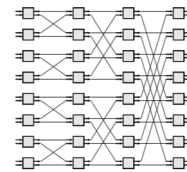


Penn ESE532 Fall 2018 -- DeHon

29

## FFT Parallelism

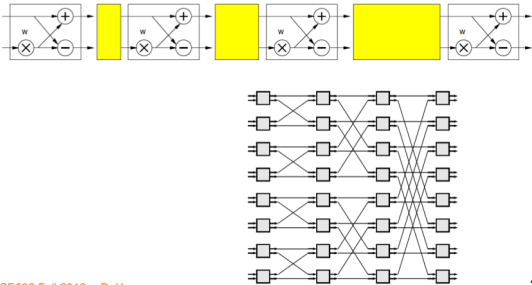
- Spatial
- Pipeline
- Streaming
- By column
  - Choose how many Butterflies to serialize on a PE
- By subgraph
- Pipeline subgraphs



Penn ESE532 Fall 2018 -- DeHon

30

## Streaming FFT

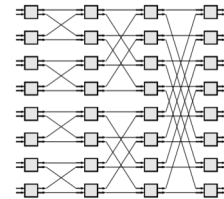


Penn ESE532 Fall 2018 -- DeHon

31

## Preclass 3

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



Penn ESE532 Fall 2018 -- DeHon

32

## Bit Serial

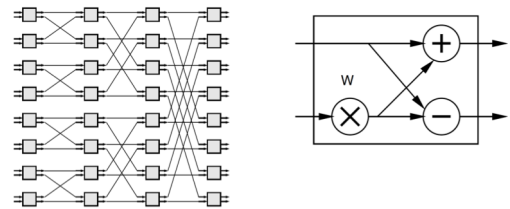
- Could compute the add/multiply bit serially
  - One full adder per adder
  - $W$  full adders per multiply
  - 50,000 LUTs
    - $\sim$  2500 bit-serial butterflies for  $W=16$ ?
      - Maybe 512-point FFT?
  - Another dimension to design space:
    - How much serialize word-wide operators
    - Use LUTs vs. DSPs

Penn ESE532 Fall 2018 -- DeHon

33

## Accelerator Building Blocks

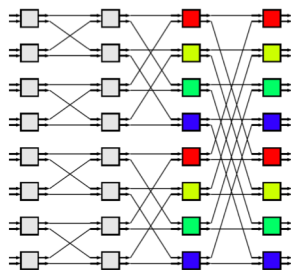
- What common subgraphs exist in the FFT?



Penn ESE532 Fall 2018 -- DeHon

34

## Common Subgraphs

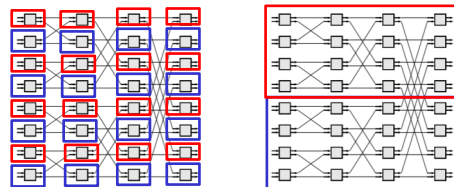


Penn ESE532 Fall 2018 -- DeHon

35

## Processor Mapping

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

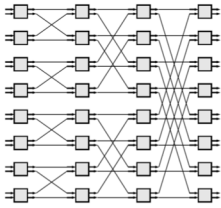


Penn ESE532 Fall 2018 -- DeHon

36

## Preclass 4a

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

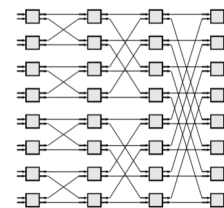


Penn ESE532 Fall 2018 -- DeHon

37

## Preclass 4b

- How change evaluation order to reduce local storage memory?

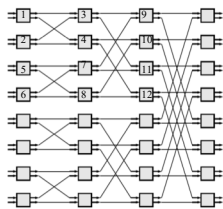


Penn ESE532 Fall 2018 -- DeHon

38

## Preclass 4b

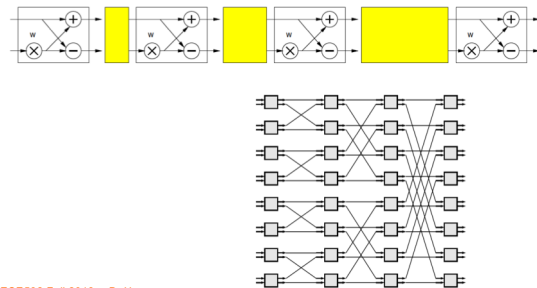
- Evaluation order



Penn ESE532 Fall 2018 -- DeHon

39

## Streaming FFT



Penn ESE532 Fall 2018 -- DeHon

40

## Communication

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

Penn ESE532 Fall 2018 -- DeHon

41

## Data Precision

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

Penn ESE532 Fall 2018 -- DeHon

42

## Number Representation

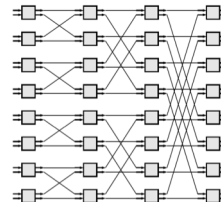
- Floating-Point
  - IEEE standard single (32b), double (64b)
    - With mantissa and exponent
    - ...half, quad ....
- Fixed-Point
  - Select total bits and fraction
    - E.g. 16.8 (16 total bits, 8 of which are fraction)
      - Represent  $1/256$  to  $256-1/256$

Penn ESE532 Fall 2018 -- DeHon

43

## Heterogeneous Precision

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



Penn ESE532 Fall 2018 -- Di

44

## 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,  $\pm 1$ ,  $\pm 1/2$ , ....
  - Multiplication by 0, 1 not need multiplier...

Penn ESE532 Fall 2018 -- DeHon

45

## FFT (partial) Design Space

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

Penn ESE532 Fall 2018 -- DeHon

46

## 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

Penn ESE532 Fall 2018 -- DeHon

47

## Admin

- P1 milestone
  - Due Friday
- P2 out
  - Asks you to identify design space

Penn ESE532 Fall 2018 -- DeHon

48