

ESE532: System-on-a-Chip Architecture

Day 5: September 16, 2018
Dataflow Process Model



Today

Dataflow Process Model

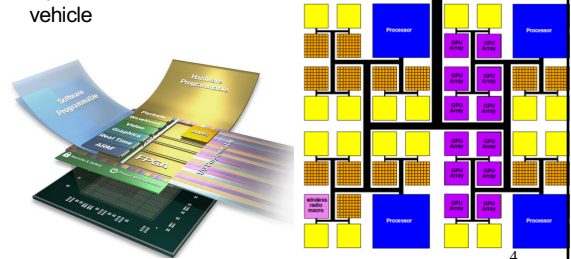
- Terms
- Issues
- Abstraction
- Performance Prospects
- Basic Approach
- If time permits
 - Dataflow variants
 - Motivations/demands for variants

Message

- Parallelism can be natural
- Expression can be agnostic to substrate
 - Abstract out implementation details
 - Tolerate variable delays may arise in implementation
- Divide-and-conquer
 - Start with coarse-grain streaming dataflow
- Basis for performance optimization and parallelism exploitation

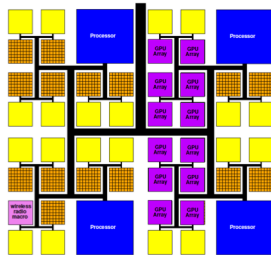
Programmable SoC

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



Reminder

- Goal: exploit parallelism on heterogeneous PSoC to achieve desired performance (energy)



Term: Process

- Abstraction of a processor
- Looks like each process is running on a separate processor
- Has own state, including
 - Program Counter (PC)
 - Memory
 - Input/output
- **May not actually run on processor**
 - Could be specialized hardware block
 - May share a processor

Thread

- Has a separate locus of control (PC)
- May share memory (contrast process)
 - Run in common address space with other threads

Penn ESE532 Fall 2019 -- DeHon

7

FIFO

- First In First Out
- Delivers inputs to outputs in order
- Data presence (empty signal)
 - Consumer knows when data available
- Back Pressure (full signal)
 - Producer knows when at capacity
 - Typically stalls
- Decouples producer and consumer processes
 - Hardware: maybe even different clocks

Penn ESE532 Fall 2019 -- DeHon

8

Process

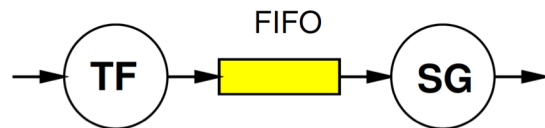
- Processes allow *expression* of independent control
- Convenient for things that advance independently
- Process (thread) is the easiest way to express some behaviors
 - Easier than trying to describe as a single process
- Can be used for performance optimization to improve resource utilization

Penn ESE532 Fall 2019 -- DeHon

9

Preclass 2

- Throughput with FIFO
 - (for concrete example on preclass)
- Correlation in delays?
- What benefit from FIFO and processes?



Penn ESE532 Fall 2019 -- DeHon

10

Preclass 2

- Independent probability of miss
 - P_f, P_g
- Concretely
 - 1 cycle in map
 - 100 run function and put in map
- If each runs independently (in isolation)
 - $T_{\sim} = 1 * (1 - P) + P * 100$
- If run together in lock step
 - Either can stall: $P = P_f + P_g - P_f P_g$
 - $T_{\sim} = 1 * (1 - P) + (P) * 100$

Penn ESE532 Fall 2019 -- DeHon

11

Model (from Day 4) Communicating Threads

- Computation is a collection of sequential/control-flow “threads”
- Threads may communicate
 - Through dataflow I/O
 - (Through shared variables)
- View as hybrid or generalization
- CSP – Communicating Sequential Processes → canonical model example

Penn ESE532 Fall 2019 -- DeHon

12

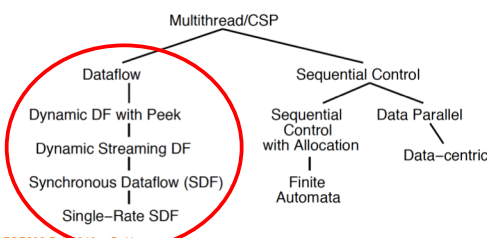
Issues

- **Communication** – how move data between processes?
 - What *latency* does this add?
 - *Throughput* achievable?
- **Synchronization** – how define how processes advance relative to each other?
- **Determinism** – for the same inputs, do we get the same outputs?

Today's Stand

- Communication – FIFO-like channels
- Synchronization – dataflow with FIFOs
- Determinism – how to achieve
 - ...until you must give it up.

Dataflow Process Model



Operation/Operator

- **Operation** – logic computation to be performed
 - A *process* that communicates through dataflow inputs and outputs
- **Operator** – physical block that performs an Operation
 - E.g. processor, hardware block

Dataflow / Control Flow

Day 4

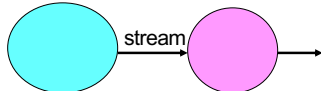
- | | |
|---|---|
| <p>Dataflow</p> <ul style="list-style-type: none"> • Program is a graph of operations • Operation consumes tokens and produces tokens • All operations run concurrently <ul style="list-style-type: none"> – All processes | <p>Control flow (e.g. C)</p> <ul style="list-style-type: none"> • Program is a sequence of operations • Operation reads inputs and writes outputs into common store • One operation runs at a time <ul style="list-style-type: none"> – defines successor |
|---|---|

Token

- Data value with presence indication
 - May be conceptual
 - Only exist in high-level model
 - Not kept around at runtime
 - Or may be physically represented
 - One bit represents presence/absence of data

Stream

- Logical abstraction of a persistent point-to-point communication link between operations (processes)
 - Has a (single) source and sink
 - Carries data presence / flow control
 - Provides in-order (FIFO) delivery of data from source to sink



Penn ESE532 Fall 2019 -- DeHon

19

Streams

- Captures communications structure
 - Explicit producer→consumer link up
- Abstract communications
 - Physical resources or implementation
 - Delay from source to sink
- Contrast
 - C: producer->consumer implicit through memory
 - Verilog/VHDL: cycles visible in implementation
 - (can add **on top of** either C or Verilog)

Penn ESE532 Fall 2019 -- DeHon

20

Variable Delay Source to Sink

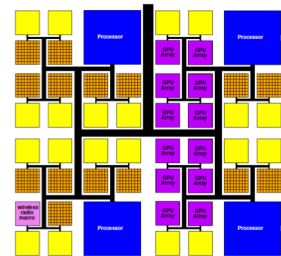
- Why might delay be variable between source and sink?

Penn ESE532 Fall 2019 -- DeHon

21

Communication Latency

- Once map to multiple processors
- Need to move data between processors
- That costs time



Penn ESE532 Fall 2019 -- DeHon

22

On-Chip Delay

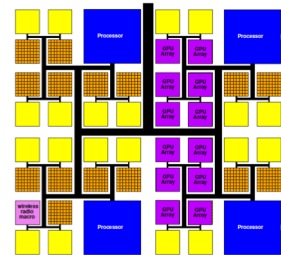
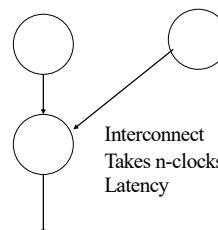
Day 3

- Delay is proportional to distance travelled
- Make a wire twice the length
 - Takes twice the latency to traverse
 - (can pipeline)
- Modern chips
 - Run at 100s of MHz to GHz
 - Take 10s of ns to cross the chip

Penn ESE532 Fall 2019 -- DeHon

23

Dataflow gives Clock Independent Semantics

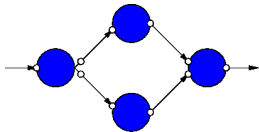


Penn ESE532 Fall 2019 -- DeHon

24

Dataflow Process Network

- Collection of Operations
- Connected by Streams
- Communicating with Data Tokens
- (CSP restricted to stream communication)



Penn ESE532 Fall 2019 -- DeHon

25

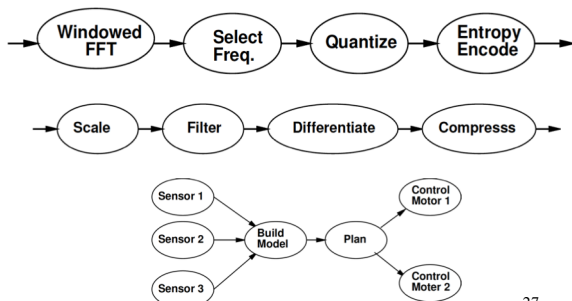
Dataflow Abstracts Timing

- Doesn't say
 - on which cycle calculation occurs
- Does say
 - What order operations occur in
 - How data interacts
 - i.e. which inputs get mixed together
- Permits
 - Scheduling on different # and types of resources
 - Operators with variable delay
 - Variable delay in interconnect

Penn ESE532 Fall 2019 -- DeHon

26

Some Task Graphs



Penn ESE532 Fall 2019 -- DeHon

27

Synchronous Dataflow (SDF) with fixed operators

- Particular, restricted form of dataflow
- Each operation
 - Consumes a **fixed** number of input tokens
 - Produces a **fixed** number of output tokens
 - Operator performs **fixed number of operations (in fixed time)**
 - When full set of inputs are available
 - Can produce output
 - Can fire any (all) operations with inputs available at any point in time

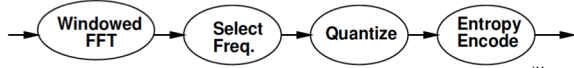
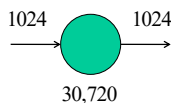
Penn ESE532 Fall 2019 -- DeHon

28

SDF Operator

FFT

- 1024 inputs
- 1024 outputs
- 10,240 multiplies
- 20,480 adds
- (or 30,720 primitive operations)



Penn ESE532 Fall 2019 -- DeHon

29

Processor Model

- Simple (for today's lecture)
 - Assume one primitive operation per cycle
- Could embellish
 - Different time per operation type
 - E.g. adds: 1 cycle, multiply: 3 cycles
 - Multiple memories with different timings

Penn ESE532 Fall 2019 -- DeHon

30

Time for Graph Iteration on Processors

- Single processor $T_{one} = \sum_i Nops_i$

- One processor per Operator
 - $T_{each} = \max(Nop_1, Nop_2, Nop_3, \dots)$

- General

$$T_{map} = \max \left(\sum_i c(1, i) \times Nops_i, \sum_i c(2, i) \times Nops_i, \sum_i c(3, i) \times Nops_i, \dots \right)$$

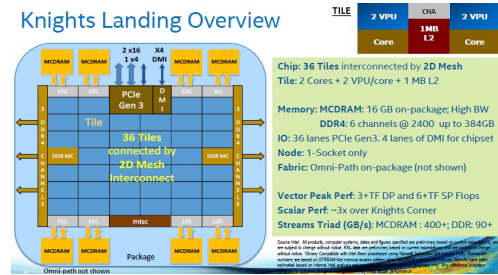
$c(x, y) = 1$ if Processor x runs task y

Penn ESE532 Fall 2019 -- DeHon

31

Intel Knights Landing

Knights Landing Overview



<https://www.nextplatform.com/2016/06/20/intel-knights-landing-yields-big-bang-buck-jump/>
 Penn ESE532 Fall 2019 -- DeHon [Intel, Micro 2016]

32

GRVI/Phallanx

- Puts 1680 RISC-V32b Integer cores
- On XCVU9P FPGA
- <http://fpga.org/2017/01/12/grvi-phalanx-joins-the-kilocore-club/>

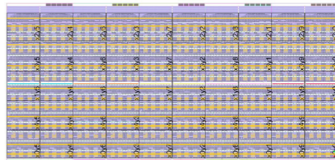
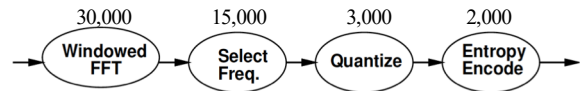


Fig 6: A 400 GRVI Phalanx. 10x5 clusters of 8 PEs (KU040)

[Gray, FCCM 2016] 33

Penn ESE532 Fall 2019 -- DeHon

Map to different processors



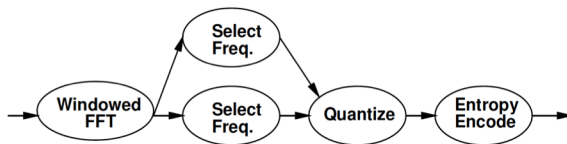
- Map to (preclass 1)
 - One processor performance?
 - One process per processor performance?
 - Two processors
 - How?
 - Performance?
 - Bottleneck?

Penn ESE532 Fall 2019 -- DeHon

34

Refine Data Parallel

- If component is data parallel, can split out parallel tasks

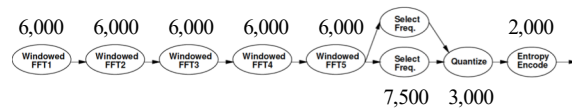


Penn ESE532 Fall 2019 -- DeHon

35

Refine Pipeline

- If operation internally pipelineable, break out pipeline into separate tasks



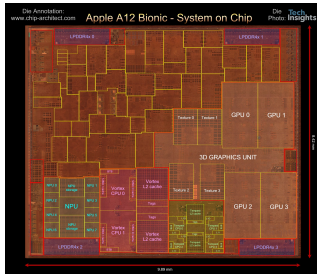
- Performance with one processor per operator?
- Achieve same performance with how many processors?

Penn ESE532 Fall 2019 -- DeHon

36

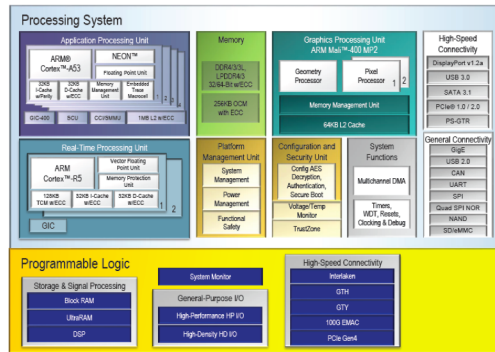
Apple A12 Bionic

- 84mm², 7nm
- 7 Billion Tr.
- iPhone XS, XR
- 6 ARM cores
 - iPad 2019
- 6 ARM cores
 - 2 fast
 - 4 low energy
- 4 custom GPUs
- Neural Engine
 - 5 Trillion ops/s?



A13 8.5B tr; still 6 ARM cores 37

Zynq® UltraScale+™ MPSoCs: EG Block Diagram



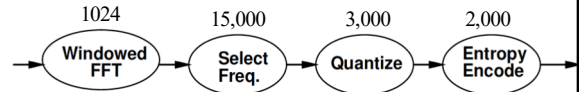
Heterogeneous Processor



- GPU perform 10 primitive FFT Ops per cycle
- Fast CPU can perform 2 ops/cycle
- Slow CPU 1 op/cycle
- Map: FFT to GPU, Select to 2 Fast CPUs, quantize and Entropy each to own Slow CPU
- Cycles/graph iteration?

Custom Accelerator

- Dataflow Process doesn't need to be mapped to a processor
- Map FFT to custom datapath on FPGA logic
 - Read and produce one element per cycle
 - 1024 cycles to process 1024-point FFT



Operations

- Can be implemented on different operators with different characteristics
 - Small or large processor
 - Hardware unit
 - Different levels of internal
 - Data-level parallelism
 - Instruction-level parallelism
 - Pipeline parallelism
- May itself be described as
 - Dataflow process network, sequential, hardware register transfer language

Streams

- Stream: logical communication link
- How might we implement:
 - Two threads running on a single processor (sharing common memory)?
 - Two processes running on different processors on the same die?
 - Two processes running on different hosts
 - E.g. one at Penn, one on Amazon cloud

Add Delay

- What do to computation if add an operation that copies inputs to outputs with some latency?
 - Impact on function?
 - Throughput impact if Identity operation has
 - Latency 10, throughput 1 value per cycle?
 - (reminder 1024 values between FFT and Select Freq.)



Penn ESE532 Fall 2019 -- DeHon

43

Semantics (meaning)

- Need to implement semantics
 - *i.e.* get same result as if computed as indicated
- But can implement any way we want
 - That preserves the semantics
 - Exploit freedom of implementation

Penn ESE532 Fall 2019 -- DeHon

44

Basic Approach

Penn ESE532 Fall 2019 -- DeHon

45

Approach (1)

- Identify natural parallelism
- Convert to streaming flow
 - Initially leave operations in software
 - Focus on correctness
- Identify flow rates, computation per operator, parallelism needed
- Refine operations
 - Decompose further parallelism?
 - E.g. data parallel split, ILP implementations
 - model potential hardware

Penn ESE532 Fall 2019 -- DeHon

46

Approach (2)

- Refine coordination as necessary for implementation
- Map operations and streams to resources
 - Provision hardware
 - Scheduling: Map operations to operators
 - Memories, interconnect
- Profile and tune
- Refine

Penn ESE532 Fall 2019 -- DeHon

47

Dataflow Variants

Time Permitting

Penn ESE532 Fall 2019 -- DeHon

48

Turing Complete

- Can implement any computation describable with a Turing Machine
 - (theoretical model of computing by Alan Turing)
- Turing Machine – captures our notion of what is computable
 - If it cannot be computed by a Turing Machine, we don't know how to compute it

DataFlow (DF) Process Network Roundup

Model	Deterministic Result	Deterministic Timing	Turing Complete
SDF+fixed-delay operators	Y	Y	N
SDF+variable delay operators	Y	N	N
DDF blocking	Y	N	Y
DDF non-blocking	N	N	Y

Synchronous Dataflow (SDF) with fixed operators

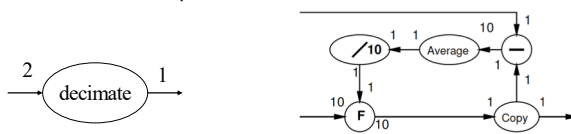
- Particular, restricted form of dataflow
- Each operation
 - Consumes a **fixed** number of input tokens
 - Produces a **fixed** number of output tokens
 - **Operator performs fixed number of operations (in fixed time)**
 - When full set of inputs are available
 - Can produce output
 - Can fire any (all) operations with inputs available at any point in time

Synchronous Dataflow (SDF)

- Particular, restricted form of dataflow
- Each operation
 - Consumes a **fixed** number of input tokens
 - Produces a **fixed** number of output tokens
 - (can take variable computation for operator)
 - When full set of inputs are available
 - Can produce output
 - Can fire any (all) operations with inputs available at any point in time

Multirate Synchronous Dataflow

- Rates can be different
 - Allow lower frequency operations
 - Communicates rates to tools
 - Use in scheduling, provisioning
 - Rates must be constant
 - Data independent



Dynamic Dataflow

- (Less) restricted form of dataflow
- Each operation
 - **Conditionally** consume input **based on data value**
 - **Conditionally** produce output **based on data value**
 - When full set of inputs are available
 - Can (optionally) produce output
 - Can fire any (all) operations with data-specified necessary inputs available at any point in time

Blocking

- Key to determinism: behavior doesn't depend on timing
 - Cannot ask if a token is present
- If (not_empty(in))
 - Out.put(3);
- Else
 - Out.put(2);

Penn ESE532 Fall 2019 -- DeHon

55

Process Network Roundup

Model	Deterministic Result	Deterministic Timing	Turing Complete
SDF+fixed-delay operators	Y	Y	N
SDF+variable delay operators	Y	N	N
DDF blocking	Y	N	Y
DDF non-blocking	N	N	Y

Penn ESE532 Fall 2019 -- DeHon

56

Motivations and Demands for Options

Time Permitting

Penn ESE532 Fall 2019 -- DeHon

57

Variable Delay Operators

- What are example of variable delay operators we might have?

Penn ESE532 Fall 2019 -- DeHon

58

Variable Delay Operators

- Operators with Variable Delay
 - Cached memory or computation
 - Shift-and-add multiply
 - Iterative divide or square-root

Penn ESE532 Fall 2019 -- DeHon

59

GCD (Preclass 3)

- What is delay of GCD computation?
 - while(a!=b)
 - t=max(a,b)-min(a,b)
 - a=min(a,b)
 - b=t
 - return(a);

Penn ESE532 Fall 2019 -- DeHon

60

Dynamic Rates?

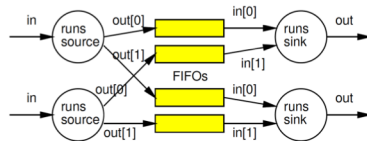
- When might static rates be limiting? (prevent useful optimizations?)

Dynamic Rates?

- Static Rates limiting
 - Compress/decompress
 - Lossless
 - Even Run-Length-Encoding
 - Filtering
 - Discard all packets from spamRus
 - Anything data dependent

When non-blocking necessary?

- What are cases where we need the ability to ask if a data item is present?
- Consider an IP packet router:



Non-Blocking

- Removed model restriction
 - Can ask if token present
- Gained expressive power
 - Can grab data as shows up
- Weaken our guarantees
 - Possible to get non-deterministic behavior

Process Network Roundup

Model	Deterministic Result	Deterministic Timing	Turing Complete
SDF+fixed-delay operators	Y	Y	N
SDF+variable delay operators	Y	N	N
DDF blocking	Y	N	Y
DDF non-blocking	N	N	Y

Big Ideas

- Capture gross parallel structure with Process Network
- Use dataflow synchronization for determinism
 - Abstract out timing of implementations
 - Give freedom of implementation
- Exploit freedom to refine mapping to optimize performance
- Minimally use non-determinism as necessary

Admin

- Reading for Day 6 on web
- HW3 due Friday
 - Implementing multiprocessor solutions on homogeneous (ARM A53) processor cores