

ESE534: Computer Organization

Day 8: February 8, 2012
Operator Sharing, Virtualization,
Programmable Architectures



Penn ESE534 Spring2012 -- DeHon

Preclass Parity

- How many gates?
- Draw solutions

Penn ESE534 Spring2012 -- DeHon

2

Previously

- Pipelining – reuse in time for **same operation**
- Memory
- Memories pack state compactly
 - densely

Penn ESE534 Spring2012 -- DeHon

3

Day 7

What is Importance of Memory?

- **Radical Hypothesis:**
 - Memory is simply a very efficient organization which allows us to store data compactly
 - (at least, in the technologies we've seen to date)
 - A great engineering **trick** to optimize resources
- **Alternative:**
 - memory is a **primary**

Penn ESE534 Spring2012 -- DeHon

4

Today

- Operator Sharing (from Day 4)
- Datapath Operation
- Virtualization
- Memory
 - ...continue unpacking the role of memory...

Penn ESE534 Spring2012 -- DeHon

5

Design Space for Computation

(from Day 4)

Penn ESE534 Spring2012 -- DeHon

6

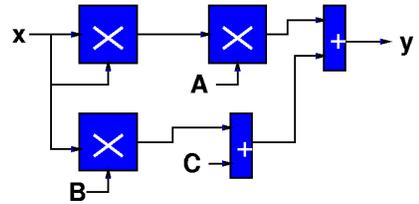
Compute Function

- Compute:
 $y = Ax^2 + Bx + C$
- Assume
 - $D(\text{Mpy}) > D(\text{Add})$
 - E.g. $D(\text{Mpy})=24, D(\text{Add})=8$
 - $A(\text{Mpy}) > A(\text{Add})$
 - E.g. $A(\text{Mpy})=64, A(\text{Add})=8$

Penn ESE534 Spring2012 -- DeHon

7

Spatial Quadratic

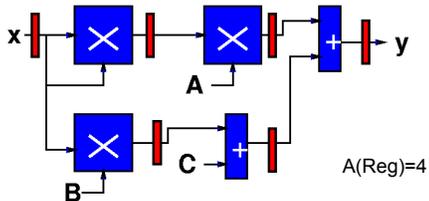


- $D(\text{Quad}) = 2 * D(\text{Mpy}) + D(\text{Add}) = 56$
- Throughput $1 / (2 * D(\text{Mpy}) + D(\text{Add})) = 1 / 56$
- $A(\text{Quad}) = 3 * A(\text{Mpy}) + 2 * A(\text{Add}) = 208$

Penn ESE534 Spring2012 -- DeHon

8

Pipelined Spatial Quadratic



- $D(\text{Quad}) = 3 * D(\text{Mpy}) = 72$
- Throughput $1 / D(\text{Mpy}) = 1 / 24$
- $A(\text{Quad}) = 3 * A(\text{Mpy}) + 2 * A(\text{Add}) + 6A(\text{Reg}) = 232$

Penn ESE534 Spring2012 -- DeHon

9

Quadratic with Single Multiplier and Adder?

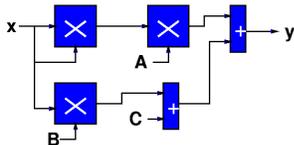
- We've seen reuse to perform the **same** operation
 - pipelining
 - bit-serial, homogeneous datapath
- We can also reuse a resource in time to perform a **different** role.

Penn ESE534 Spring2012 -- DeHon

10

Quadratic Datapath

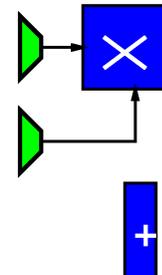
- Start with one of each operation
- HW2.4 showed could just use adders



11

Quadratic Datapath

- Multiplier serves multiple roles
 - $x * x$
 - $A * (x * x)$
 - $B * x$
- Will need to be able to steer data (switch interconnections)



Penn ESE534 Spring2012 -- DeHon

12

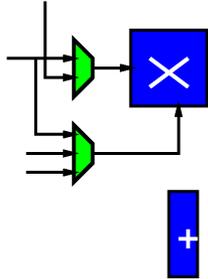
Quadratic Datapath

- Multiplier serves multiple roles

- x^2
- $A \cdot (x^2)$
- $B \cdot x$

- Inputs

- a) x, x^2
- b) x, A, B



Penn ESE534 Spring2012 -- DeHon

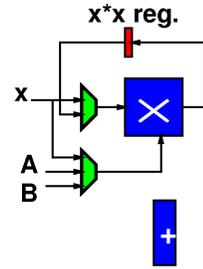
Quadratic Datapath

- Multiplier serves multiple roles

- x^2
- $A \cdot (x^2)$
- $B \cdot x$

- Inputs

- a) x, x^2
- b) x, A, B



Penn ESE534 Spring2012 -- DeHon

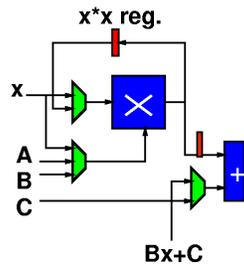
Quadratic Datapath

- Adder serves multiple roles

- $(Bx) + C$
- $(A \cdot x^2) + (Bx + C)$

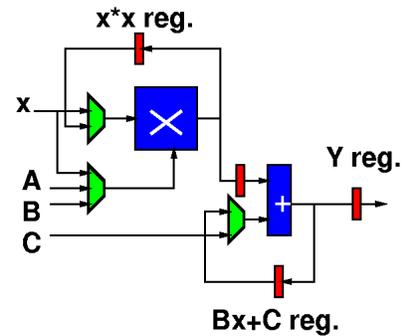
- Inputs

- one always mpy output
- $C, Bx + C$



Penn ESE534 Spring2012 -- DeHon

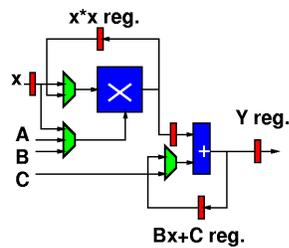
Quadratic Datapath



Penn ESE534 Spring2012 -- DeHon

Quadratic Datapath

- Add input register for x



Penn ESE534 Spring2012 -- DeHon

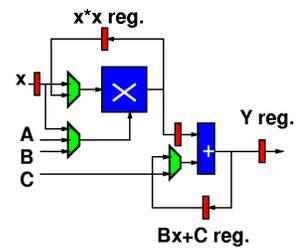
Quadratic Control

- Now, we just need to control the datapath

- What control?

- Control:

- LD x
- LD x^2
- MA Select
- MB Select
- AB Select
- LD $Bx + C$
- LD Y



Penn ESE534 Spring2012 -- DeHon

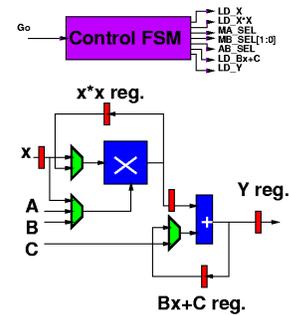
FSMD

- FSMD = FSM + Datapath
- Stylization for building controlled datapaths such as this (a **pattern**)
- Of course, an FSMD is just an FSM
 - it's often easier to think about as a datapath
 - synthesis, place and route tools have been notoriously bad about discovering/exploiting datapath structure

Penn ESE534 Spring2012 -- DeHon

19

Quadratic FSMD



Penn ESE534 Spring2012 --

20

Quadratic FSMD Control

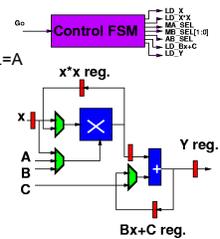
- S0: if (go) LD_X; goto S1
 - else goto S0
- S1: MA_SEL=x, MB_SEL[1:0]=x, LD_x*x
 - goto S2
- S2: MA_SEL=x, MB_SEL[1:0]=B
 - goto S3
- S3: AB_SEL=C, MA_SEL=x*x, MB_SEL=A
 - goto S4
- S4: AB_SEL=Bx+C, LD_Y
 - goto S0

Penn ESE534 Spring2012 -- DeHon

21

Quadratic FSMD Control

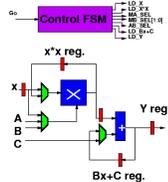
- S0: if (go) LD_X; goto S1
 - else goto S0
- S1: MA_SEL=x, MB_SEL[1:0]=x, LD_x*x
 - goto S2
- S2: MA_SEL=x, MB_SEL[1:0]=B
 - goto S3
- S3: AB_SEL=C, MA_SEL=x*x, MB_SEL=A
 - goto S4
- S4: AB_SEL=Bx+C, LD_Y
 - goto S0



Penn ESE534 Spring2012 -- DeHon

Quadratic FSM

- $D(\text{mux3}) = D(\text{mux2}) = 1$
- $A(\text{mux2}) = 2$
- $A(\text{mux3}) = 3$
- $A(\text{QFSM}) \approx 10$
- Latency/Throughput/Area?
- Latency: $5 \cdot (D(\text{MPY}) + D(\text{mux3})) = 125$
- Throughput: $1/\text{Latency} = 1/125$
- Area: $A(\text{Mpy}) + A(\text{Add}) + 5 \cdot A(\text{Reg}) + 2 \cdot A(\text{Mux2}) + A(\text{Mux3}) + A(\text{QFSM}) = 109$



Penn ESE534 Spring2012 -- DeHon

23

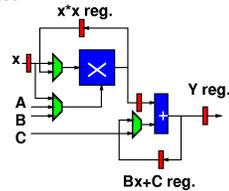
Universal Sharing

Penn ESE534 Spring2012 -- DeHon

24

Review

- Given a task: $y = Ax^2 + Bx + C$
- Saw how to share primitive operators
- Got down to one of each



Penn ESE534 Spring2012 -- DeHon

Very naively

- Might seem we need one of each different type of operator

Penn ESE534 Spring2012 -- DeHon

26

..But

- Doesn't fool us
- We already know that **and** gate (and many other things—HW1.3) are universal
- So, we know, we can build a universal compute operator

Penn ESE534 Spring2012 -- DeHon

27

Temporal Composition

Penn ESE534 Spring2012 -- DeHon

28

Temporal

- Don't have to implement all the gates *at once*
- Can *reuse* one gate over time

Penn ESE534 Spring2012 -- DeHon

29

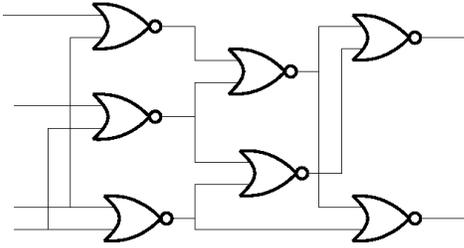
Temporal Decomposition

- Take Set of gates
- Sort topologically
 - All predecessors before successors
- Give a unique number to each gate
 - Hold value of its outputs
- Use a memory to hold the gate values
- Sequence through gates

Penn ESE534 Spring2012 -- DeHon

30

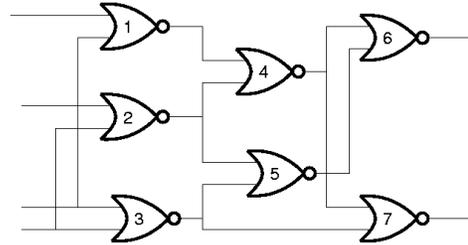
Example Logic



Penn ESE534 Spring2012 -- DeHon

31

Numbered Gates



Penn ESE534 Spring2012 -- DeHon

32

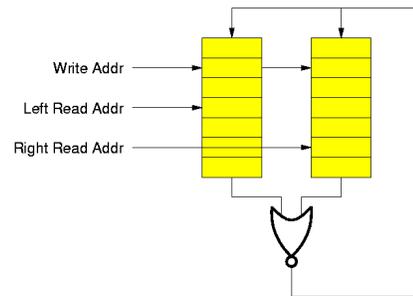
Preclass

- Number gates

Penn ESE534 Spring2012 -- DeHon

33

nor2 Memory/Datapath



Penn ESE534 Spring2012 -- DeHon

34

Programming?

- How do we program this netlist?

Penn ESE534 Spring2012 -- DeHon

35

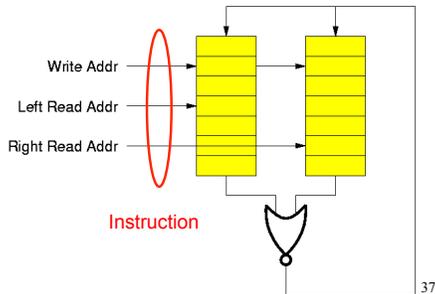
Programming?

- Program gates
 - Tell each gate where to get its input
 - Tell gate n where its two inputs come from
 - Specify the memory location for the output of the associated gate
 - Each gate operation specified with
 - two addresses (the input sources for gate)
 - This is the *instruction* for the gate

Penn ESE534 Spring2012 -- DeHon

36

nor2 Memory/Datapath



Penn ESE534 Spring2012 -- DeHon

37

Supply Instruction

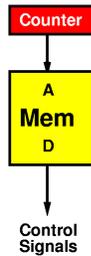
- How can we supply the sequence of instructions to program this operation?

Penn ESE534 Spring2012 -- DeHon

38

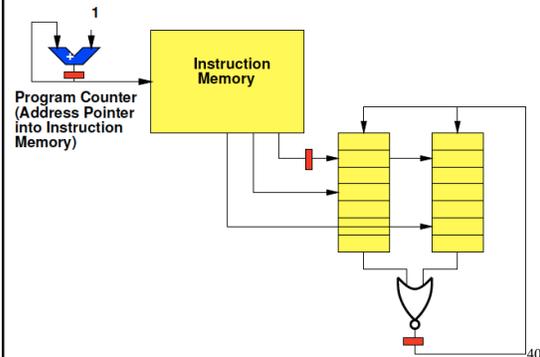
Simplest Programmable Control

- Use a memory to “record” control instructions
- “Play” control with sequence



Penn ESE534 Spring2012 -- DeHon

Temporal Gate Architecture



Penn ESE534 Spring2012 -- DeHon

40

How program preclass computation?

- How would we program the preclass computation?
 - Complete the memory

Penn ESE534 Spring2012 -- DeHon

41

Simulate the Logic

- For Preclass
 - Go around the room calling out:
 - Identify PC
 - Identify instruction
 - Perform nor2 on slot ___ and slot ___
 - Result is ___
 - Store into slot ___

Penn ESE534 Spring2012 -- DeHon

42

What does this mean?

- With only one **active** component
 - **nor** gate
- Can implement **any** function
 - given appropriate
 - state (memory)
 - muxes (interconnect)
 - Control

Penn ESE534 Spring2012 -- DeHon

43

Defining Terms

Fixed Function:

- Computes one function (e.g. FP-multiply, divider, DCT)
- Function defined at fabrication time

Programmable:

- Computes “any” computable function (e.g. Processor, DSPs, FPGAs)
- Function defined after fabrication

Penn ESE534 Spring2012 -- DeHon

44

Result

- Can sequence together primitive operations in time
- **Communicating** state through memory
 - Memory as interconnect
- To perform “arbitrary” operations

Penn ESE534 Spring2012 -- DeHon

45

“Any” Computation? (Universality)

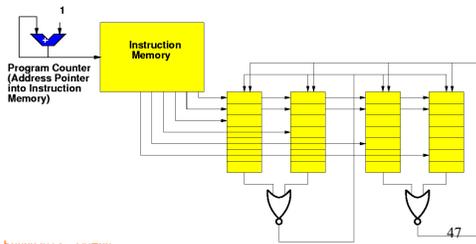
- Any computation which can “fit” on the programmable substrate
- **Limitations:** hold entire computation and intermediate data

Penn ESE534 Spring2012 -- DeHon

46

Temporal-Spatial Variation

- Can have any number of gates
 - Tradeoff Area for Reduce Time....



Penn ESE534 Spring2012 -- DeHon

47

Use of Memory?

- **What did we use memory for here?**
- State
- Instructions
- Interconnect

Penn ESE534 Spring2012 -- DeHon

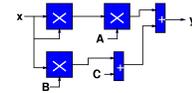
48

“Stored Program” Computer/ Processor

- Can build a datapath that can be *programmed* to perform **any** computation.
- Can be built with limited hardware that is *reused* in time.
- **Historically:** this was a key contribution from Penn’s Moore School
 - Computer Engineers: Eckert and Mauchly
 - ENIAC→EDVAC
 - (often credited to Von Neumann)

What have we done?

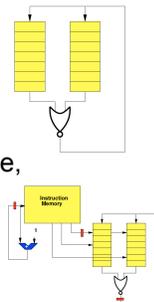
- Taken a computation: $y = Ax^2 + Bx + C$
- Turned it into operators and interconnect



- Decomposed operators into a basic primitive:
 - nor, adds

What have we done?

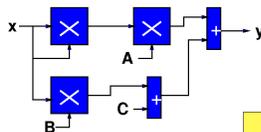
- Said we can implement it on as few as one of **compute unit** (nor2)
- Added a unit for **state**
- Added an **instruction** to tell single, universal unit how to act as each operator in original graph



Virtualization

- We’ve *virtualized* the computation
- No longer need one **physical** compute unit for each operator in original computation
- Can suffice with:
 1. shared operator(s)
 2. a **description** of how each operator behaved
 3. a place to store the intermediate data between operators

Virtualization



Why Interesting?

- Memory compactness
- This works and was interesting because
 - the area to describe a computation, its interconnect, and its state
 - is much smaller than the physical area to spatially implement the computation
- e.g. traded multiplier for
 - few memory slots to hold state
 - few memory slots to describe operation
 - time on a shared unit (adder, gate)

Questions?

Admin

- HW4 due Monday
- No new reading for Monday

Big Ideas [MSB Ideas]

- Memory: efficient way to hold state
 - ...and allows us to describe/implement computations of unbounded size
- State can be \ll computation [area]
- **Resource sharing**: key trick to reduce area
- Memory key tool for Area-Time tradeoffs
- “configuration” signals allow us to *generalize* the utility of a computational operator

Big Ideas [MSB-1 Ideas]

- First programmable computing unit
- Two key functions of memory
 - retiming (interconnect in time)
 - instructions
 - description of computation