## ESE535:
## Electronic Design Automation

Day 7: February 9, 2015
Scheduled Operator Sharing

Penn ESE535 Spring 2015 -- DeHon

---

## Today

- Sharing Resources
- Area-Time Tradeoffs
- Throughput vs. Latency
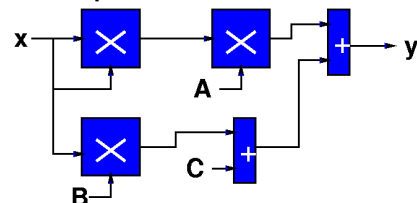- VLIW Architectures
- Scheduling (introduce)
  - Maybe start on

Behavioral
(C, MATLAB, …)

Arch. Select
Schedule

RTL
FSM assign
Two-level,
Multilevel opt.
Covering
Retiming

Gate Netlist
Placement
Routing

Layout

Masks

Penn ESE535 Spring 2015 -- DeHon

2

---

## Compute Function

- Compute:

$$y=Ax^2 +Bx +C$$

- Assume
  - $D(Mpy) > D(Add)$
  - $A(Mpy) > A(Add)$

Penn ESE535 Spring 2015 -- DeHon

3

---

## Spatial Quadratic



- $A(Quad) = 3*A(Mpy) + 2*A(Add)$

Penn ESE535 Spring 2015 -- DeHon

4

---

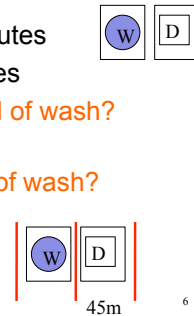## Latency vs. Throughput
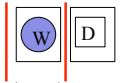
- **Latency:** Delay from inputs to output(s)
- **Throughput:** Rate at which can introduce new set of inputs

Penn ESE535 Spring 2015 -- DeHon

5

---

## Washer/Dryer Example

- 1 Washer Takes 30 minutes
- 1 Dryer Takes 45 minutes
- How long to do one load of wash?
  - → Wash latency
- How long to do 5 loads of wash?
- Wash Throughput?

45m

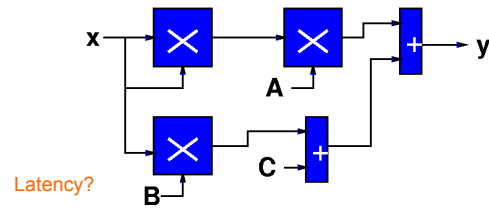Penn ESE535 Spring 2015 -- DeHon

6

---

1

## Pipelining

- Break up the computation graph into stages
  - Allowing us to
    - reuse portions of the graph for new data,
    - while older data is still working its way through the graph
      - Before it has exited graph
  - Use registers to isolate regions
  - Throughput > (1/Latency)
- Relate liquid in pipe
  - Doesn't wait for first drop of liquid to exit far end of pipe before accepting second drop
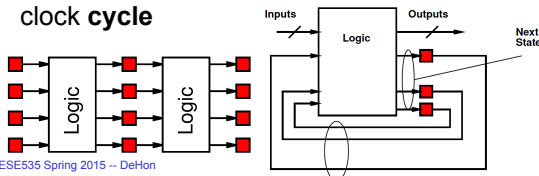
7

## Spatial Quadratic



Latency?

- D(Quad) = 2*D(Mpy)+D(Add) = 21
- Throughput 1/(2*D(Mpy)+D(Add)) = 1/21
- A(Quad) = 3*A(Mpy) + 2*A(Add) = 32

8

## Synchronous Discipline

- Compute
  - From registers
  - Through combinational logic
  - To new values for registers
- Delay through logic sets a lower bound on the duration of each clock – the clock **cycle**
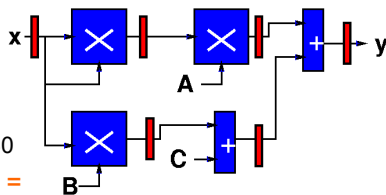
## Terms

- **Latency:** Delay from inputs to output(s)
- **Cycle Time:**
  - Clock period
  - Critical path delay between registers
- **Throughput:** Rate at which can introduce new set of inputs
  - Typically, inverse of cycle time
- **Pipelining**: how we separate latency from cycle time

10

## Pipelined Spatial Quadratic



- D(Quad) =
  3*D(Mpy) = 30
- Throughput =
  1/D(Mpy) = 1/10
- A(Quad) =
  3*A(Mpy)+2*A(Add)+6A(Reg) = 35
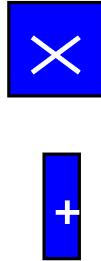
11

## Quadratic with Single Multiplier and Adder?

- We've seen reuse to perform the **same** operation
  - pipelining
- We can also reuse a resource in time to perform a different role.
  - Here: x*x, A*(x*x), B*x
  - also: (Bx)+c, (A*x*x)+(Bx+c)

12

2

## Quadratic Datapath

- Start with one of each operation
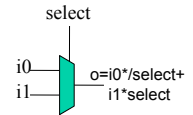
13

---

## Multiplexer

- Gate allows us to select data from multiple sources

- Mux
  - For short

- Useful when sharing operators
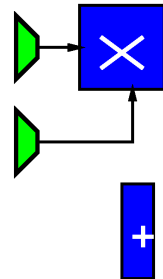
select

i0
i1
$o = i0*/select + i1*select$

14

---

## Quadratic Datapath

- Multiplier serves multiple roles
  - x*x
  - A*(x*x)
  - B*x
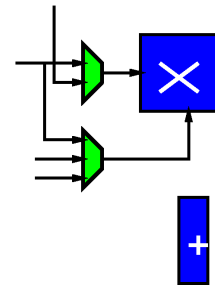- Use multiplexer to steer data (switch interconnections)
  - A(mux) < A(multiply)

15

---

## Quadratic Datapath

- Multiplier serves multiple roles
  - x*x
  - A*(x*x)
  - B*x
- x, x*x
- x,A,B

---

## Quadratic Datapath

- Multiplier serves multiple roles
  - x*x
  - A*(x*x)
  - B*x
- x, x*x
- x,A,B

x*x reg.

x

A
B

---

## Quadratic Datapath

- Adder serves multiple roles
  - (Bx)+c
  - (A*x*x)+(Bx+c)
- one always mpy output
- C, Bx+C

x*x reg.

x

A
B
C

Bx+C

---

3

# Quadratic Datapath

**x*x reg.**

x

A
B
C

**Y reg.**

**Bx+C reg.**

# Quadratic Datapath

• Add input
register for *x*

**x*x reg.**

x

A
B
C

**Y reg.**

**Bx+C reg.**

# Cycle Impact?

• Need more cycles
• How about the delay of
each cycle?
  – Add mux delay
  – Register setup/hold time,
  clock skew
  – Limited by slowest
  operation
  – Cycle?
    • D(Mpy)+2*D(Mux2) = 10.2

**x*x reg.**

x

A
B
C

**Y reg.**

**Bx+C reg.**

# Quadratic Control

• Now, we just need to control the datapath
• What control?
• Control:
  – LD x
  – LD x*x
  – MA Select
  – MB Select
  – AB Select
  – LD Bx+C
  – LD Y

**x*x reg.**

x

A
B
C

**Y reg.**

**Bx+C reg.**

# Quadratic Control

1. LD_X
2. MA_SEL=x,MB_SEL[1:0]=x, LD_x*x
3. MA_SEL=x,MB_SEL[1:0]=B
4. AB_SEL=C,MA_SEL=x*x, MB_SEL=A,
   LD_Bx+C
5. AB_SEL=Bx+C, LD_Y

[Could combine 1 and 5 and
do in 4 cycles;  analysis
that follows assume 5 as shown.]

**x*x reg.**

x

A
B
C

**Y reg.**

**Bx+C reg.**

# Quadratic Memory Control

1. LD_X
2. MA_SEL=x,
   MB_SEL[1:0]=x,
   LD_x*x
3. MA_SEL=x,
   MB_SEL[1:0]=B
4. AB_SEL=C,
   MA_SEL=x*x,
   MB_SEL=A, LD_Bx
   +C
5. AB_SEL=Bx+C, LD_Y

Address

```
000: 00000001
001: 00010010
010: 00001000
011: 01000100
100: 10100000
101: 00000000
```

**x*x reg.**

x

A
B
C

**Y reg.**

**Bx+C reg.**

## Quadratic Datapath



- Latency/Throughput/Area?
- Latency: 5*(D(MPY)+D(mux3))=51
- Throughput: 1/Latency ~= 0.02
- Area: A(Mpy)+A(Add)+5*A(Reg)
  +2*A(Mux2)+A(Mux3)+A(Imem)=17.5+
  A(Imem)

25

## Quadratic with 2 Mult, 1 Add

**A  x    B        C**



| step | X | X | + |
|------|-----|-----|------------------|
| 1 | X*X | B*X | |
| 2 | A*(X*X) | | (B*X)+C |
| 3 | | | (A*X*X*X) +(B*X+C) |

- Latency/Throughput/Area?

26

## Quadratic with 2 Mult, 1 Add

**A  x    B        C**



| step | X | X | + |
|------|-----|-----|------------------|
| 1 | X*X | B*X | |
| 2 | A*(X*X) | | (B*X)+C |
| 3 | | | (A*X*X*X) +(B*X+C) |

- Latency = 3*(D(Mpy)+D(Mux))=30.3
- Throughput = 1/30.3 ~=0.03
- Area = 2*A(Mpy)+4*A(Mux2)+A(Add)
  +3*A(Reg) = 26.5

27

## Quadratic: Area-Time Tradeoff

| Design | Area | Throughput | Latency |
|--------|------|-----------|---------|
| 3M2A (pipe) | 35 | 0.1 | 30 |
| 2M1A | 26.5 | 0.03 | 30.3 |
| 1M1A | 17.5 | 0.02 | 51 |

28

## Registers→Memory

- Generally can see many registers
- If # registers >> physical operators
  – Only need to access a few at a time
- Group registers into memory banks

29

## Memory Bank Quadratic

- Store x
- x*x
- B*x
- A*x$^2$; B*x+c
- (A*x$^2$)+(B*x+c)

30

5

## Memory Bank Quadratic

- Store x
- x*x
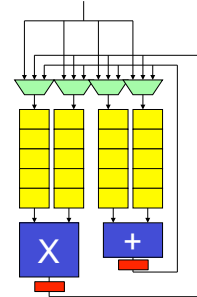- B*x
- $A*x^2$; $B*x+c$
- $(A*x^2)+(B*x+c)$



x
$x^2$
x
B
A
Bx
$Ax^2$
c
Bx+c
X
+

31

## Cycle Impact?

How cycle changed?
- Add mux delay
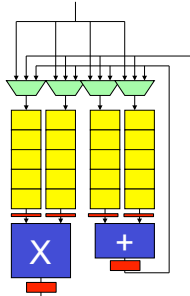- Register setup/hold time, clock skew
- Memory read/write
  - Could pipeline



X
+

32

## Cycle Impact?

- Add mux delay
- Register setup/hold time, clock skew
- Memory read/write
  - Could pipeline
  - Impact?
    - Latency
    - Throughput?



X
+

33

## Impact

- When have big operators
  - Like multiplier
- Can share them to reduce area
  - At cost of throughput
  - Maybe at cost of latency, energy
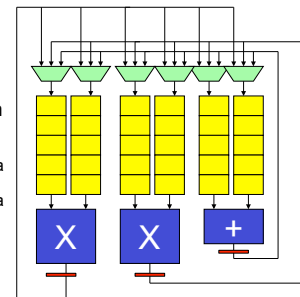- This gives a rich trade space

34

## Details

- At extreme, number of "big" operators is dominant cost
  - Total number for area
  - Number in path for delay
- Does cost additional area, delay to share them
  - sometimes a lower order cost

35

## VLIW

- Very Long Instruction Word
- Set of operators
  - Parameterize number, distribution (X, +, sqrt…)
    - More operators→ less time, more area
    - Fewer operators→ more time, less area
- Memories for intermediate state



X
X
+

36

6

## VLIW

- Very Long Instruction Word
- Set of operators
  - Parameterize number, distribution (X, +, sqrt…)
    - More operators→ less time, more area
    - Fewer operators→ more time, less area
- Memories for intermediate state
- Memory for "long" instructions

37

## VLIW

38

## VLIW

- Very Long Instruction Word
- Set of operators
  - Parameterize number, distribution (X, +, sqrt…)
    - More operators→ less time, more area
    - Fewer operators→ more time, less area
- Memories for intermediate state
- Memory for "long" instructions
- Schedule compute task
- General framework for specializing to problem
  - Wiring, memories get expensive
  - Opportunity for further optimizations
- General way to tradeoff area and time

39

## VLIW

40

## Review

- Reuse physical operators in time
- Share operators in different roles
- Allows us to reduce area at expense of increasing time
- Area-Time tradeoff
- Pay some sharing overhead
  - Muxes, memory
- VLIW – general formulation for shared datapaths

41

## Design Automation

Sets up two problems for us:
- Provisioning
  - (Architecture Selection)
  - End of next week (after…)
- Scheduling
  - Start introducing now
  - Next two lectures

Behavioral
(C, MATLAB, …)

Arch. Select
Schedule

RTL

FSM assign
Two-level,
Multilevel opt.
Covering
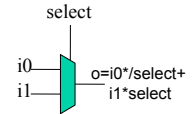Retiming

Gate Netlist

Placement
Routing

Layout

Masks

42

7

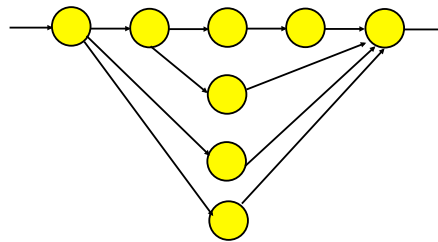## Time Permitting

## General Problem

- Resources are not free
  - Wires, io ports
  - Functional units
    - LUTs, ALUs, Multipliers, ….
  - Memory access ports
  - State elements
    - memory locations
    - Registers
      - Flip-flop
      - loadable master-slave latch
  - Multiplexers (mux)

select

i0

i1

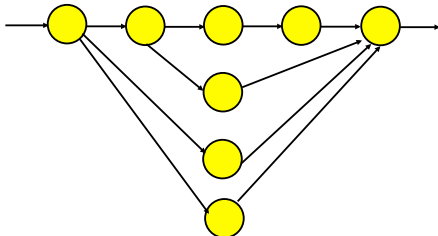o=i0*/select+
i1*select

## Trick/Technique

- Resources can be shared (reused) in time
- Sharing resources can reduce
  - instantaneous resource requirements
  - total costs (area)

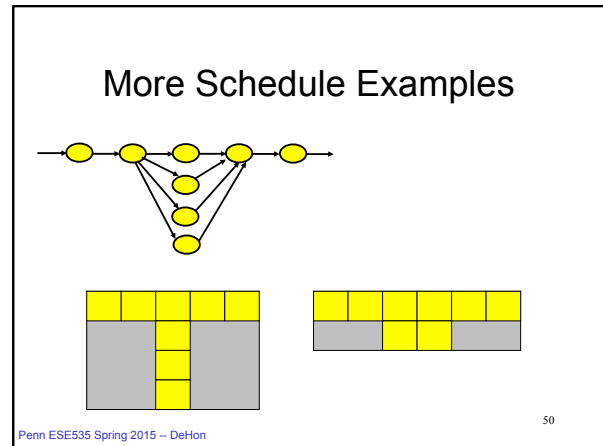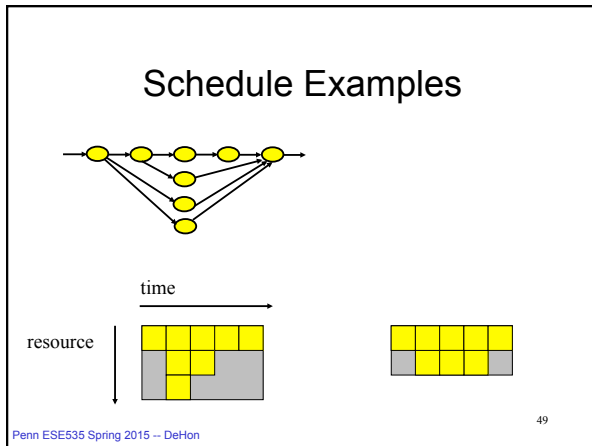- **Pattern:** scheduled operator sharing

## Example

## Example

Assume unit delay operators.
How many operators do I need to evaluate this computation
in ~5 time units.

## Sharing

- Does not have to increase delay
  - w/ careful time assignment
  - can often reduce peak resource requirements
  - while obtaining original (unshared) delay
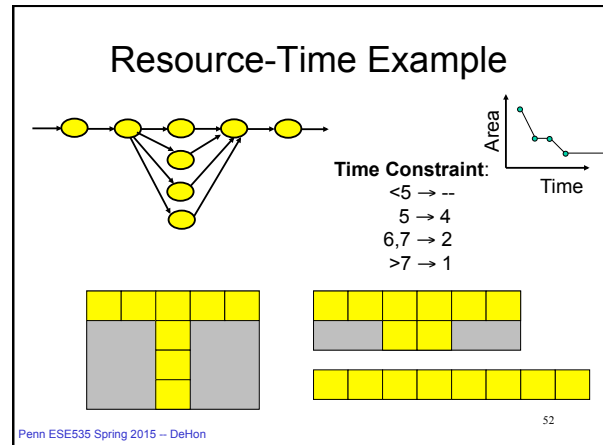- Alternately: Minimize delay given fixed resources

## Schedule Examples



time

resource

49

## More Schedule Examples

50

## Scheduling

- **Task**: assign time slots (and resources) to operations
  - **time-constrained**: minimizing peak resource requirements
    - *n.b.* time-constrained, not always constrained to minimum execution time
  - **resource-constrained**: minimizing execution time

51

## Resource-Time Example



**Time Constraint**:
<5 → --
5 → 4
6,7 → 2
>7 → 1

Area

Time

52

## Scheduling Use

- Very general problem formulation
  - HDL/Behavioral → RTL
  - Register/Memory allocation/scheduling
  - Instruction/Functional Unit scheduling
  - Processor tasks
  - Time-Switched Routing
    - TDMA, bus scheduling, static routing
  - Routing (share channel)

53

## Two Types (1)

- **Data independent**
  - graph static
  - resource requirements and execution time
    - independent of data
  - schedule staticly
  - maybe bounded-time guarantees
  - typical ECAD problem

54

9

## Two Types (2)

- **Data Dependent**
  - execution time of operators variable
    - depend on data
  - flow/requirement of operators data dependent
  - if cannot bound range of variation
    - must schedule online/dynamically
    - cannot guarantee bounded-time
    - general case (*I.e.* halting problem)
  - typical "General-Purpose" (non-real-time) OS problem

## Unbounded Resource Problem

- Easy:
  - compute ASAP schedule (next slide)
    - *I.e.* schedule everything as soon as predecessors allow
  - will achieve minimum time
  - won't achieve minimum area
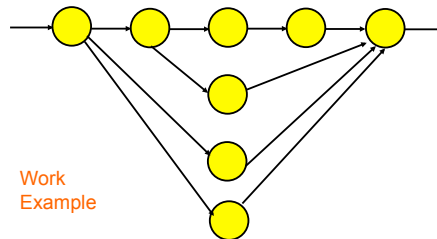    - (meet resource bounds)

## ASAP Schedule
## As Soon As Possible (ASAP)

- For each input
  - mark input on successor
  - if successor has all inputs marked, put in visit queue
- While visit queue not empty
  - pick node
  - update time-slot based on latest input
  - mark inputs of all successors, adding to visit queue when all inputs marked

## ASAP Example



Work Example

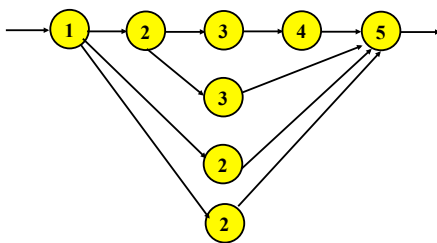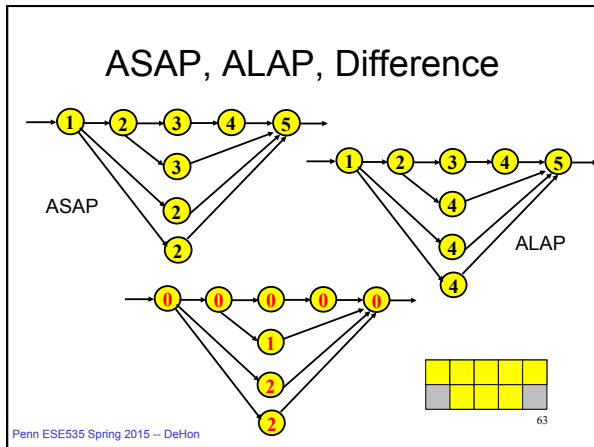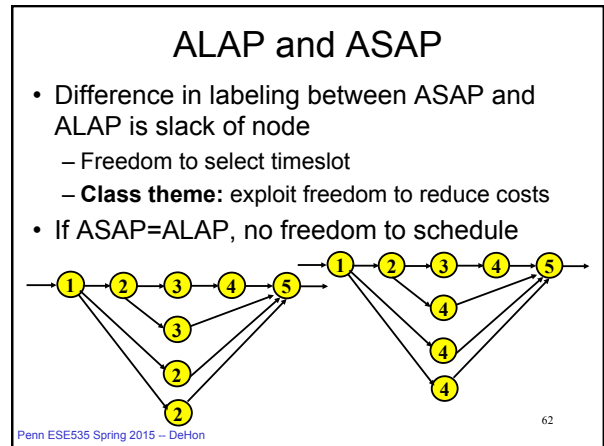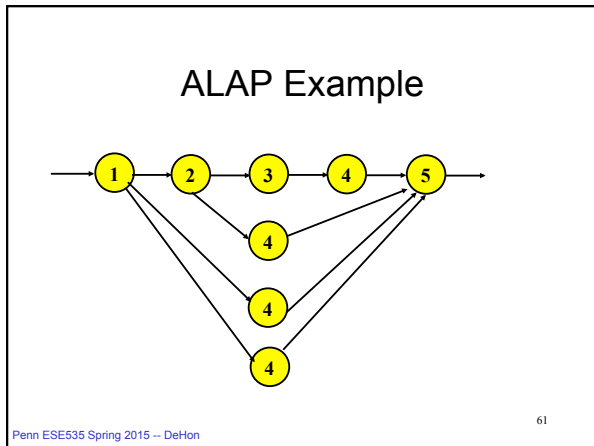## ASAP Example

## Also Useful to Define ALAP

- As Late As Possible
- Work backward from outputs of DAG
- Also achieve minimum time w/ unbounded resources

Rework Example

## ALAP Example

61

## ALAP and ASAP

- Difference in labeling between ASAP and ALAP is slack of node
  - Freedom to select timeslot
  - **Class theme:** exploit freedom to reduce costs
- If ASAP=ALAP, no freedom to schedule

62

## ASAP, ALAP, Difference



ASAP

ALAP

63

## Big Ideas:

- Scheduled Operator Sharing
- Area-Time Tradeoffs

64

## Admin

- Assignment 2, 3 feedback on canvas
- Assignment 4 due Thursday
- Reading for Wednesday online

65

11