# ESE534:
## Computer Organization
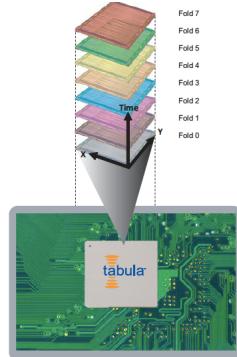
Day 22: April 9, 2012
Time Multiplexing

---

## Tabula

- March 1, 2010
  - Announced new architecture
- We would say
  - w=1, c=8 arch.

[src: www.tabula.com]      2

Fold 7
Fold 6
Fold 5
Fold 4
Fold 3
Fold 2
Fold 1
Fold 0
Time
Y

---

## Previously

- Saw how to pipeline **architectures**
  - specifically interconnect
  - talked about general case
- Saw how to **reuse** resources at maximum rate to do the *same* thing
- Saw demand-for and options-to-support data retiming

3

---

## Today

- Multicontext
  - Review why
  - Cost
  - Packing into contexts
  - Retiming requirements
  - Some components
- [concepts we saw in overview week 2-3, we can now dig deeper into details]

4

---

## How often is **reuse** of the *same* operation applicable?

- In what cases can we exploit high-frequency, heavily pipelined operation?

- …and when can we not?

5

---

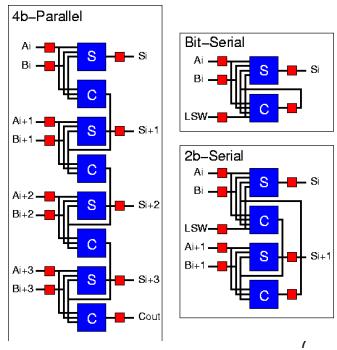## How often is **reuse** of the *same* operation applicable?

- Can we exploit higher frequency offered?
  - High throughput, feed-forward (acyclic)
  - Cycles in flowgraph
    - abundant data level parallelism [C-slow]
    - no data level parallelism
  - Low throughput tasks
    - structured (*e.g.* datapaths) [serialize datapath]
    - unstructured
  - Data dependent operations
    - similar ops [local control -- next time]
    - dis-similar ops

6

---

1

## Structured Datapaths

- Datapaths: same *pinst* for all bits
- Can serialize and reuse the same data elements in succeeding cycles
- example: adder

---

## Preclass 1

- Sources of inefficient mapping
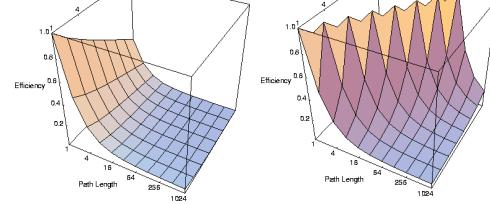  $W_{task}=4$, $L_{task}=4$
  to $W_{arch}=1$, $C=1$ architecture?

8

---

## Preclass 1

- How transform $W_{task}=4$, $L_{task}=4$ to run efficiently on $W_{arch}=1$, $C=1$ architecture?
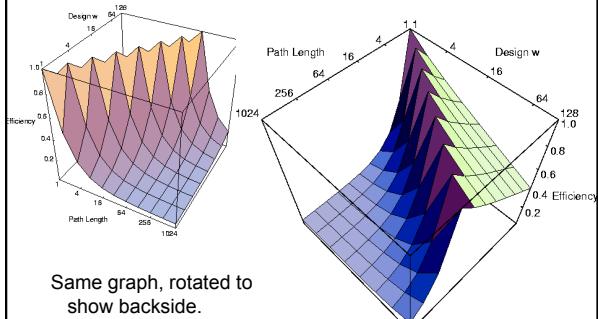
- Impact on efficiency?

9

---

## Throughput Yield

FPGA Model -- if throughput requirement is reduced for wide word operations, serialization allows us to reuse active area for same computation

10

---

## Throughput Yield

Same graph, rotated to show backside.

---

## Remaining Cases

- Benefit from multicontext as well as high clock rate
- *i.e.*
  - cycles, no parallelism
  - data dependent, dissimilar operations
  - low throughput, irregular (can't afford swap?)

12

## Single Context

- When have:
  - cycles and no data parallelism
  - low throughput, unstructured tasks
  - dis-similar data dependent tasks
- Active resources sit idle most of the time
  - Waste of resources
- Cannot reuse resources to perform **different** function, only **same**
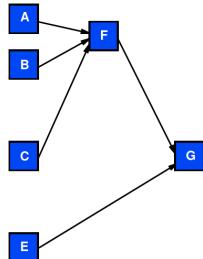
13

---

## Resource Reuse

- To use resources in these cases
  - must direct to do different things.

- Must be able tell resources how to behave

- ➔ separate instructions (*pinsts*) for each behavior

14

---

## Preclass 2

- How schedule onto 3 contexts?

15

---

## Preclass 2
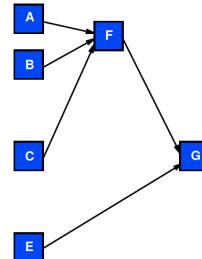
- How schedule onto 4 contexts?

16

---

## Preclass 2

- How schedule onto 6 contexts?

17

---

## Example: Dis-similar Operations

18

## Multicontext Organization/ Area

- $A_{ctxt} \approx 80K\lambda^2$
  - dense encoding
- $A_{base} \approx 800K\lambda^2$

- $A_{ctxt} : A_{base} = 1:10$
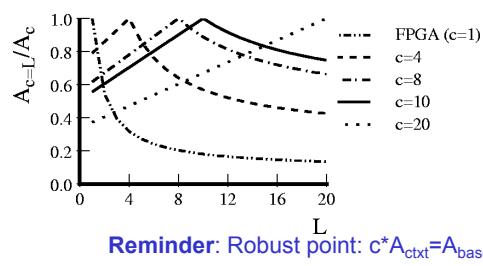
Context ID

Decode

Context ID

Decode

Memory

K–LUT

Interconnect

19

---

## Preclass 3

- Area:
  - Single context?
  - 3 contexts?
  - 4 contexts?
  - 6 contexts?

20

---

## Multicontext Tradeoff Curves

- Assume Ideal packing: $N_{active} = N_{total}/L$

$A_{c=L}/A_c$

1.0
0.8
0.6
0.4
0.2
0.0

0   4   8   12   16   20

L

FPGA (c=1)
c=4
c=8
c=10
c=20

**Reminder**: Robust point: $c*A_{ctxt} = A_{base}$

21

---

## In Practice

Limitations from:
- Scheduling
- Retiming

22

---

## Scheduling

23

---

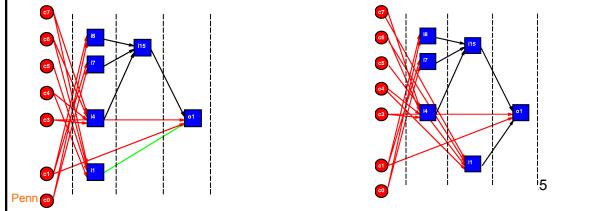## Scheduling Limitations

- $N_A$ (**active**)
  - size of largest stage

- **Precedence**:
  can evaluate a LUT only after predecessors have been evaluated
  → cannot always, completely equalize stage requirements

24

4

## Slide 1

# Scheduling

- Precedence limits packing freedom
- Freedom do have
  – shows up as slack in network
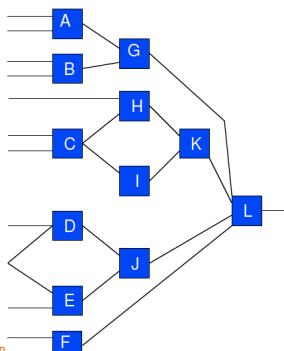
5

## Slide 2

# Scheduling

- Computing Slack:
  – ASAP (As Soon As Possible) Schedule
    - propagate depth forward from primary inputs
      – depth = 1 + max input depth
  – ALAP (As Late As Possible) Schedule
    - propagate distance from outputs back from outputs
      – level = 1 + max output consumption level
  – Slack
    - slack = L+1-(depth+level)  [PI depth=0, PO level=0]
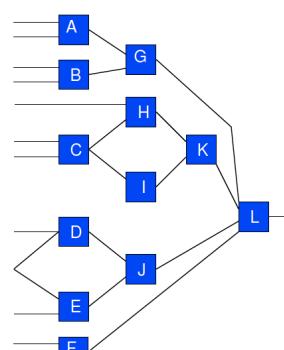
26

## Slide 3

# Work Slack Example

27

## Slide 4

# Preclass 4

- With precedence constraints, and unlimited hardware, how many contexts?

## Slide 5

# Preclass 5
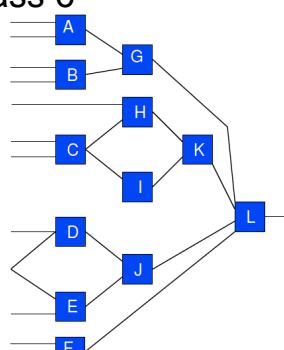
- **Without** precedence, how many compute blocks needed to evaluate in 4 contexts?

## Slide 6

# Preclass 6

- Where can schedule?
  – J
  – D

5

## Preclass 6
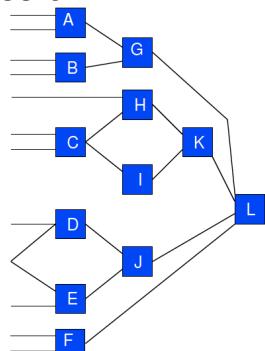
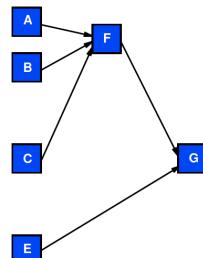- Where can schedule D if J in 3?
- Where can schedule D if J in 2?

## Preclass 6

- Where can schedule J if D in 1?
- Where can schedule J if D in 2?
- Where schedule operations?
- Physical blocks ?

## Reminder (Preclass 1)
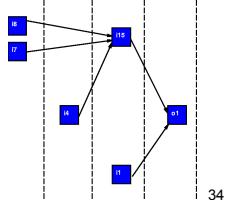
33

## Sequentialization

- Adding time slots
  - more sequential (more latency)
  - add slack
    - allows better balance
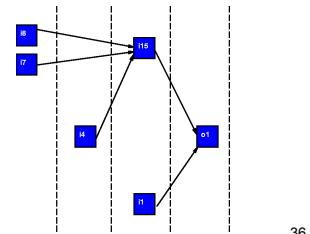
$L=4 \rightarrow N_A=2$ (4 contexts)

34

## Retiming

35

## Multicontext Data Retiming

- How do we accommodate intermediate data?
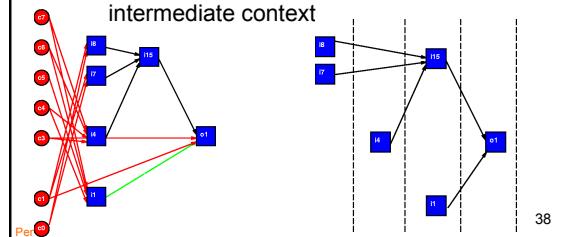
36

6

## Signal Retiming

- Single context, non-pipelined
  - hold value on LUT Output (wire)
    - from production through consumption

  - Wastes wire and switches by occupying
    - for entire critical path delay L
    - not just for 1/L'th of cycle takes to cross wire segment

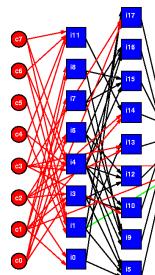  - How show up in multicontext?

37

## Signal Retiming

- Multicontext equivalent
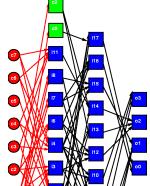  - need LUT to hold value for each intermediate context

38

## ASCII→Hex Example



**Single Context**: 21 LUTs @ 880Kλ$^2$=18.5Mλ$^2$

39

## ASCII→Hex Example



**Three Contexts**: 12 LUTs @ 1040Kλ$^2$=12.5Mλ$^2$

40

## ASCII→Hex Example

- All retiming on wires (active outputs)
  - saturation based on inputs to largest stage
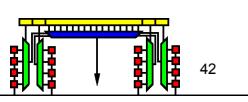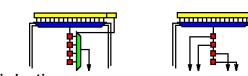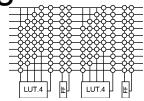


Ideal≡Perfect scheduling spread + no retime overhead

41

## Alternate Retiming

- Recall from last time (Day 20)
  - Net buffer
    - smaller than LUT
  - Output retiming
    - may have to route multiple times
  - Input buffer chain
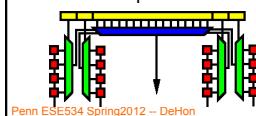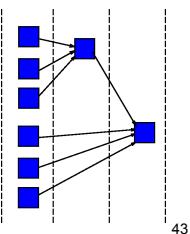    - only need LUT every depth cycles

42

7

## Input Buffer Retiming

- Can only take K unique inputs per cycle
- Configuration depth differ from context-to-context
  - Cannot schedule LUTs in slot 2 and 3 on the same physical block, since require 6 inputs.

43

---

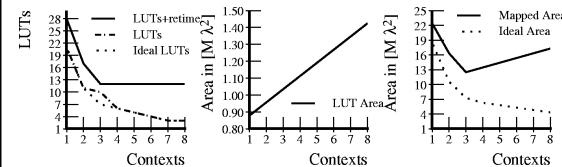## ASCII→Hex Example

- All retiming on wires (active outputs)
  - saturation based on inputs to largest stage



Ideal≡Perfect scheduling spread + no retime overhead

44

---

## ASCII→Hex Example (input retime)



@ depth=4, c=6: 5.5M$\lambda^2$
(compare 18.5M$\lambda^2$ )
3.4×

45

---

## General throughput mapping:

- If only want to achieve limited throughput
- Target produce new result every t cycles
1. Spatially pipeline every t stages
   cycle = t
2. retime to minimize register requirements
3. multicontext evaluation w/in a spatial stage
   try to minimize resource usage
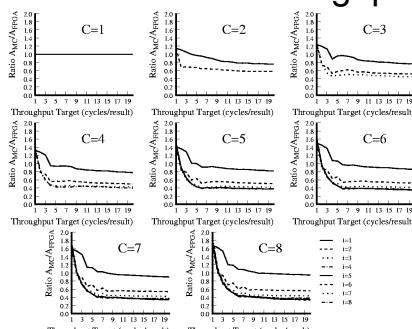4. Map for depth (i) and contexts (c)

46

---

## Benchmark Set

- 23 MCNC circuits
  - area mapped with SIS and Chortle

| Circuit | Mapped LUTs | Path Length | Circuit | Mapped LUTs | Path Length |
|---------|-------------|-------------|---------|-------------|-------------|
| 5xp1 | 46 | 10 | des | 1267 | 13 |
| 9sym | 123 | 7 | e64 | 230 | 9 |
| 9symml | 108 | 8 | f51m | 45 | 17 |
| C499 | 85 | 10 | misex1 | 20 | 6 |
| C880 | 176 | 21 | misex2 | 38 | 8 |
| alu2 | 169 | 19 | rd73 | 105 | 10 |
| apex6 | 248 | 9 | rd84 | 150 | 9 |
| apex7 | 77 | 7 | rot | 293 | 16 |
| b9 | 46 | 7 | sao2 | 73 | 9 |
| clip | 121 | 9 | vg2 | 60 | 9 |
| cordic | 367 | 13 | z4ml | 8 | 7 |
| count | 46 | 16 | | | |

---

## Multicontext vs. Throughput

48

---

8

## Multicontext vs. Throughput



Ratio $A_{MC}/A_{FPGA}$ vs. Throughput Target (cycles/result)

- c=1,i=1
- c=2,i=2
- c=3,i=3
- c=4,i=4
- c=5,i=4
- c=6,i=5
- c=7,i=5
- c=8,i=6

49

---

## General Theme

- Ideal Benefit
  - *e.g.* Active=N/C
- Logical Constraints
  - Precedence
- Resource Limits
  - Sometimes bottleneck
- Net Benefit
- Resource Balance



LUTs+retime / LUTs / Ideal LUTs

50

---

## Beyond Area
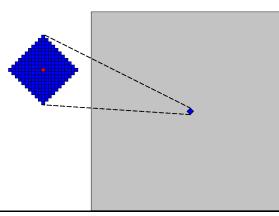
(Did not cover this section in class)

51

---

## Only an Area win?

- If area were free, would we always want a fully spatial design?

52

---

## Communication Latency
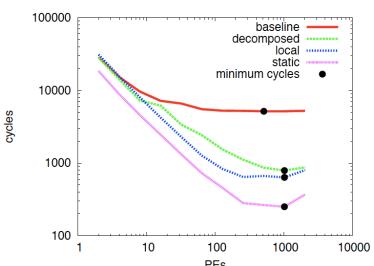
- Communication latency across chip can limit designs
- Serial design is smaller → less latency

53

---

## Optimal Delay for Graph App.
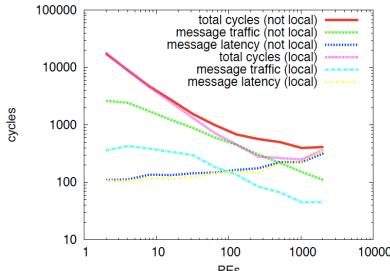


cycles vs. PEs

- baseline
- decomposed
- local
- static
- minimum cycles

54

9

## Optimal Delay Phenomena

55

## What Minimizes Energy

- HW5

$$B = \sqrt{N}$$

56

## Components

57

## DPGA (1995)

| Process | $1.0_\mu$ CMOS |
|---|---|
| Chip | $7.1mm \times 6.8mm$ |
| AEs | 144 |
| Contexts | 4 |
| AE Area | $640K\lambda^2$ |
| $A_{base}$ | $544K\lambda^2$ |
| $A_{ctx}$ | $24K\lambda^2$ |
| $A_{base} : A_{ctx}$ | 20+:1 |
| (nominal delay) | 9ns |

[Tau et al., FPD 1995]      58

## Xilinx Time-Multiplexed FPGA

- Mid 1990s Xilinx considered Multicontext FPGA
  - Based on XC4K (pre-Virtex) devices
  - Prototype Layout in F=500nm
  - Required **more** physical interconnect than XC4K
  - Concerned about power (10W at 40MHz)

[Trimberger, FCCM 1997]    59

## Xilinx Time-Multiplexed FPGA

- Two unnecessary expenses:
  - Used output registers with separate outs
  - Based on XC4K design
    - Did not densely encode interconnect configuration
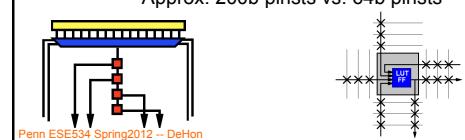    - Compare 8 bits to configure input C-Box connection
      - Versus $\log_2(8)=3$ bits to control mux select
      - Approx. 200b pinsts vs. 64b pinsts

60

10

## Tabula



Physical organization of Spacetime tile
Routing MUXes
Configuration RAM (64x9x2 bits)
Spacetime Logic
Configuration RAM (64x9x2 bits)
Routing MUXes
Config RAM control & Clock control

- 8 context, 1.6GHz, 40nm
  - 64b pinsts
- Our model w/ input retime
  [MPR/Tabula 3/29/2009]
  - $1M\lambda^2$ base
    - $80K\lambda^2$ / 64b pinst Instruction mem/context
    - $40K\lambda^2$ / input-retime depth
  - $1M\lambda^2 + 8 \times 0.12M\lambda^2 \sim = 2M\lambda^2$ ➔ 4× LUTs (ideal)
    - Recall ASCIItoHex 3.4, similar for thput map
- They claim 2.8× LUTs

61

---

## Admin

- No required reading for Wed.
  - Supplemental on web
- Office hours Tuesday
  - Good for questions about project

62

---

## Big Ideas
## [MSB Ideas]

- Several cases cannot profitably reuse same logic at device cycle rate
  - cycles, no data parallelism
  - low throughput, unstructured
  - dis-similar data dependent computations
- These cases benefit from more than one instructions/operations per active element
- $A_{ctxt} << A_{active}$ makes interesting
  - save area by sharing active among instructions

63

---

## Big Ideas
## [MSB-1 Ideas]

- Economical retiming becomes important here to achieve active LUT reduction
  - one output reg/LUT leads to early saturation
- c=4--8, I=4--6 automatically mapped designs roughly 1/3 single context size
- Most FPGAs typically run in realm where multicontext is smaller
  - How many for intrinsic reasons?
  - How many for lack of HSRA-like register/CAD support?

64