

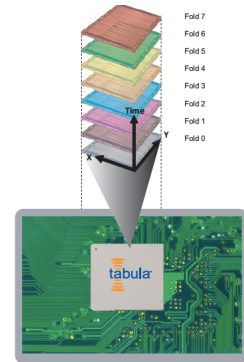
ESE534: Computer Organization

Day 22: April 14, 2010
Time Multiplexing



Tabula

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



[src: www.tabula.com]

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

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]

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

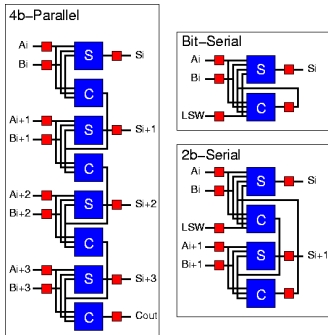
- In what cases can we exploit high-frequency, heavily pipelined operation?
- ...and when can we not?

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

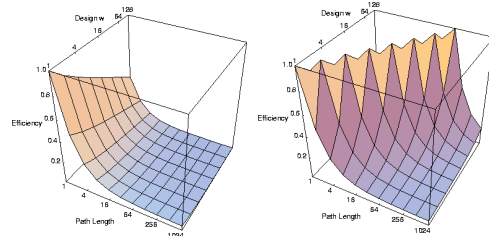
Structured Datapaths

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



Penn ESE534 Spring2010 -- DeHon

Throughput Yield

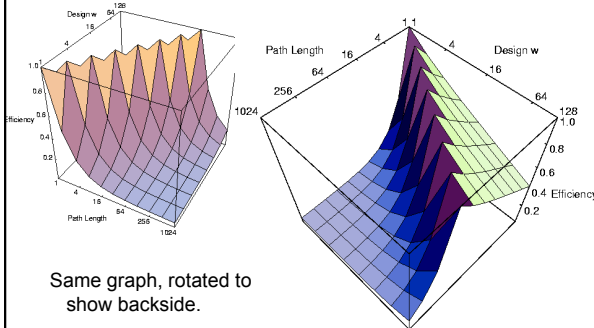


FPGA Model -- if throughput requirement is reduced for wide word operations, serialization allows us to reuse active area for same computation

8

Penn ESE534 Spring2010 -- DeHon

Throughput Yield



Same graph, rotated to show backside.

Penn ESE534 Spring2010 -- DeHon

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?)

10

Penn ESE534 Spring2010 -- DeHon

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**

11

Penn ESE534 Spring2010 -- DeHon

Resource Reuse

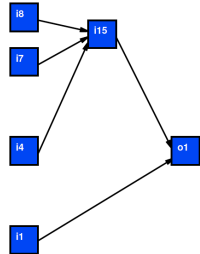
- To use resources in these cases
 - must direct to do different things.
- Must be able tell resources how to behave
- → separate instructions (*pinsts*) for each behavior

12

Penn ESE534 Spring2010 -- DeHon

Preclass 1

- How schedule onto 3 contexts?

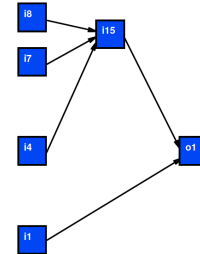


Penn ESE534 Spring2010 -- DeHon

13

Preclass 1

- How schedule onto 4 contexts?

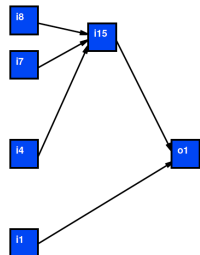


Penn ESE534 Spring2010 -- DeHon

14

Preclass 1

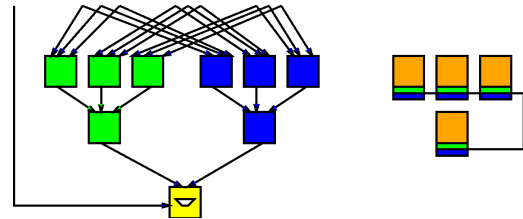
- How schedule onto 6 contexts?



Penn ESE534 Spring2010 -- DeHon

15

Example: Dis-similar Operations

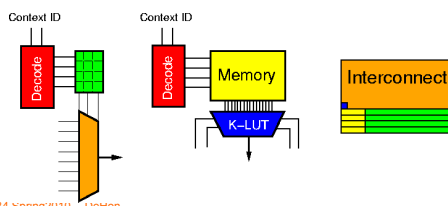


Penn ESE534 Spring2010 -- DeHon

16

Multicontext Organization/ Area

- $A_{cbxt} \approx 80K\lambda^2$
– dense encoding
- $A_{base} \approx 800K\lambda^2$
- $A_{cbxt} : A_{base} = 1:10$



Penn ESE534 Spring2010 -- DeHon

17

Preclass 2

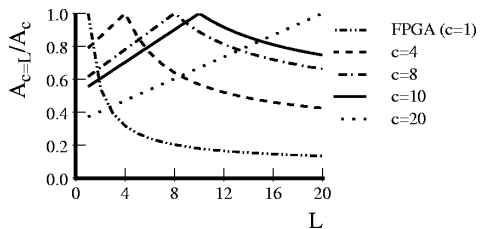
- Area:
 - Single context?
 - 3 contexts?
 - 4 contexts?
 - 6 contexts?

Penn ESE534 Spring2010 -- DeHon

18

Multicontext Tradeoff Curves

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



Reminder: Robust point: $c \cdot A_{\text{cctx}} = A_{\text{base}}$

19

Penn ESE534 Spring2010 -- DeHon

In Practice

Limitations from:

- Scheduling
- Retiming

20

Penn ESE534 Spring2010 -- DeHon

Scheduling

21

Penn ESE534 Spring2010 -- DeHon

Scheduling Limitations

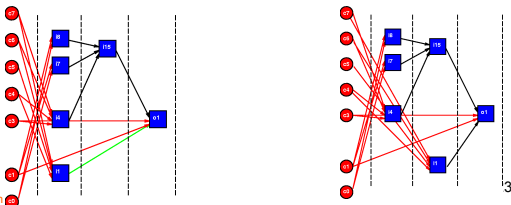
- N_A (active)
 - size of largest stage
- Precedence:**
 - can evaluate a LUT only after predecessors have been evaluated
 - cannot always, completely equalize stage requirements

22

Penn ESE534 Spring2010 -- DeHon

Scheduling

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



3

Penn

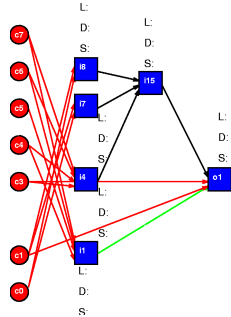
Scheduling

- 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 = $L+1-(\text{depth}+\text{level})$ [PI depth=0, PO level=0]

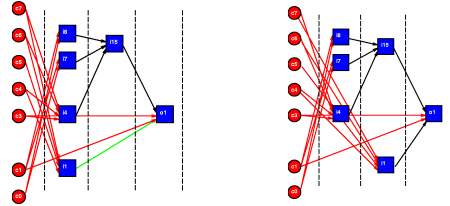
24

Penn ESE534 Spring2010 -- DeHon

Slack Example



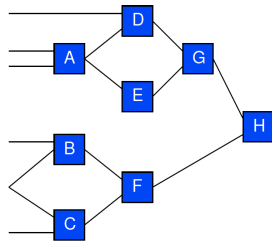
Allowable Schedules



Active LUTs (N_A) = 3

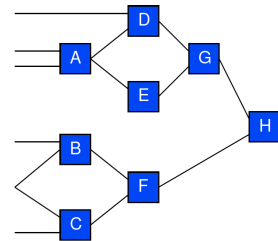
Preclass 3

- With precedence constraints, 4 context evaluation needs
 - Number of physical compute blocks?



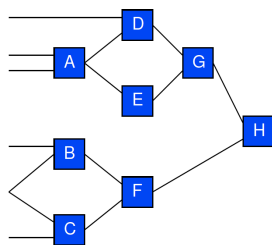
Preclass 4

- Slack on nodes?
- Where can schedule
 - B
 - C
 - F



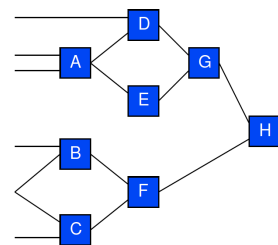
Preclass 4

- Where can schedule B if F in 2?
- Where can schedule B if F in 3?

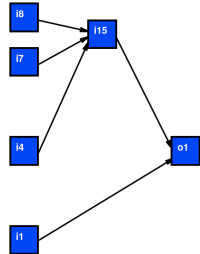


Preclass 4

- Where can schedule F if C in 1?
- Where can schedule F if C in 2?
- Where want to put
 - B
 - C
 - F
- Physical blocks ?



Reminder (Preclass 1)

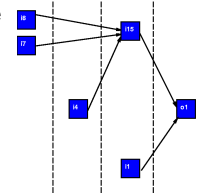


Penn ESE534 Spring2010 -- DeHon

31

Sequentialization

- Adding time slots
 - more sequential (more latency)
 - add slack
 - allows better balance



$L=4 \rightarrow N_A=2$ (4 contexts)

Penn ESE534 Spring2010 -- DeHon

32

Retiming

Penn ESE534 Spring2010 -- DeHon

33

Multicontext Data Retiming

- How do we accommodate intermediate data?

Penn ESE534 Spring2010 -- DeHon

34

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

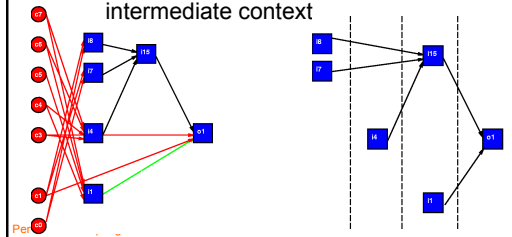
– How show up in multicontext?

Penn ESE534 Spring2010 -- DeHon

35

Signal Retiming

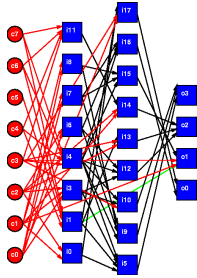
- Multicontext equivalent
 - need LUT to hold value for each intermediate context



Penn ESE534 Spring2010 -- DeHon

36

ASCII→Hex Example

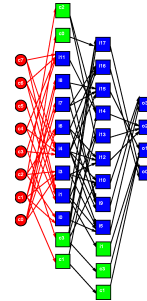


Single Context: 21 LUTs @ $880K\lambda^2=18.5M\lambda^2$

Penn ESE534 Spring2010 -- DeHon

37

ASCII→Hex Example



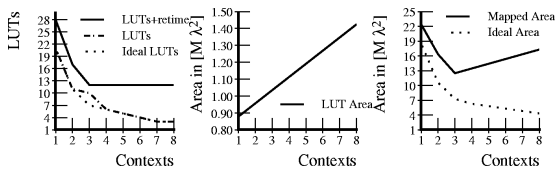
Three Contexts: 12 LUTs @ $1040K\lambda^2=12.5M\lambda^2$

Penn ESE534 Spring2010 -- DeHon

38

ASCII→Hex Example

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

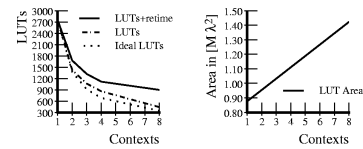


Ideal=Perfect scheduling spread + no retime overhead

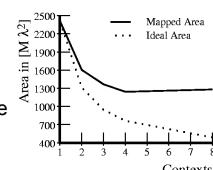
Penn ESE534 Spring2010 -- DeHon

39

DES Latency Example



Single Output case

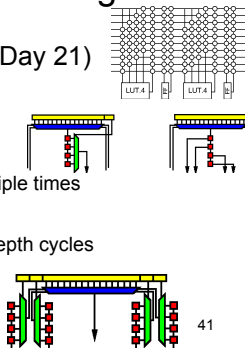


Penn ESE534 Spring2010 -- uemon

40

Alternate Retiming

- Recall from last time (Day 21)
 - Net buffer
 - smaller than LUT
 - Output retiming
 - may have to route multiple times
 - Input buffer chain
 - only need LUT every depth cycles

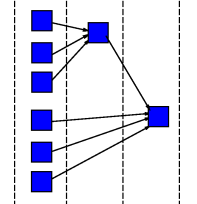
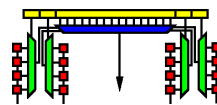


Penn ESE534 Spring2010 -- DeHon

41

Input Buffer Retiming

- Can only take K unique inputs per cycle
- Configuration depth differ from context-to-context



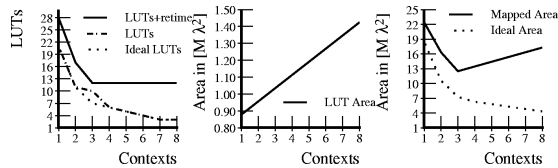
Penn ESE534 Spring2010 -- DeHon

42

Reminder

ASCII→Hex Example

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

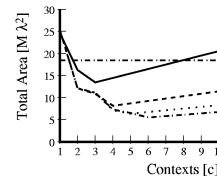
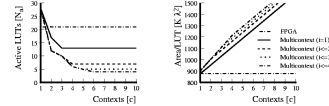


Ideal=Perfect scheduling spread + no retime overhead

43

Penn ESE534 Spring2010 -- DeHon

ASCII→Hex Example (input retime)



@ depth=4, c=6: 5.5Mλ²
(compare 18.5Mλ²)
3.4x

44

Penn ESE534 Spring2010 -- DeHon

General throughput mapping:

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

45

Penn ESE534 Spring2010 -- DeHon

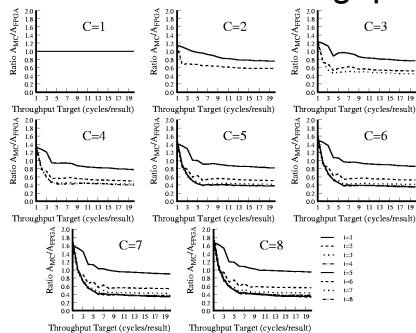
Benchmark Set

- 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			

Penn ESE534 Spring2010 -- DeHon

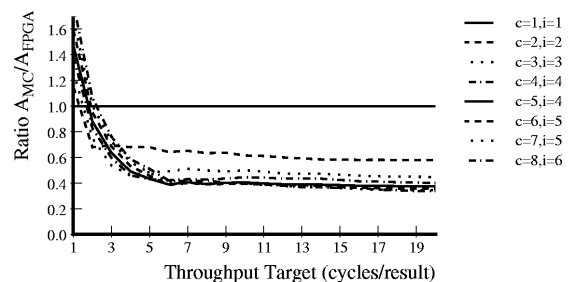
Multicontext vs. Throughput



47

Penn ESE534 Spring2010 -- DeHon

Multicontext vs. Throughput

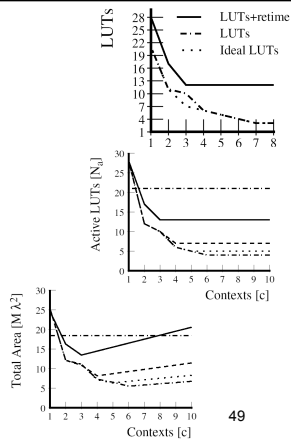


48

Penn ESE534 Spring2010 -- DeHon

General Theme

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



Penn ESE534 Spring2010 -- DeHon

49

Beyond Area

Penn ESE534 Spring2010 -- DeHon

50

Only an Area win?

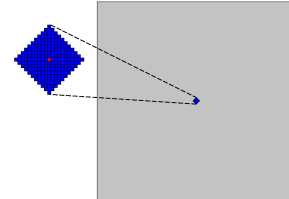
- If area were free, would I always want a fully spatial design?

Penn ESE534 Spring2010 -- DeHon

51

Communication Latency

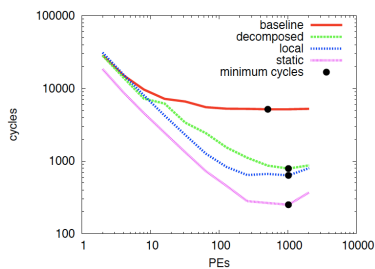
- Communication latency across chip can limit designs
- Serial design is smaller → less latency



Penn ESE534 Spring2010 -- DeHon

52

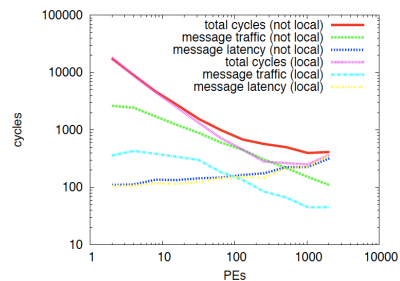
Optimal Delay for Graph App.



Penn ESE534 Spring2010 -- DeHon

53

Optimal Delay Phenomena



Penn ESE534 Spring2010 -- DeHon

54

What Minimizes Energy

- HW5

$$B = \sqrt{N}$$

"EnergyALU"

N	EnergyALU
1	25000
2	12500
4	6250
8	3125
16	1562
32	781
64	390
128	195
256	97
512	48
1024	24
2048	12
4096	6
8192	3
16384	1.5

55

Penn ESE534 Spring2010 -- DeHon

Components

56

Penn ESE534 Spring2010 -- DeHon

DPGA (1995)

Process	1.0μ CMOS
Chip	7.1mm×6.8mm
AEs	144
Contexts	4
AE Area	640Kλ ²
A _{base}	544Kλ ²
A _{ctx}	24Kλ ²
A _{base} : A _{ctx}	20+:1
(nominal delay)	9ns

57

Penn ESE680-002 Spring2007 -- DeHon [Tau et al., FPD 1995]

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)

58

Penn ESE534 Spring2010 -- DeHon [Trimberger, FCCM 1997]

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 log₂(8)=3 bits to control mux select
 - Approx. 200b pins vs. 64b pins

59

Penn ESE534 Spring2010 -- DeHon

Tabula

- 8 context, 1.6GHz, 40nm
 - 64b pins
- Our model w/ input retime
 - 1Mλ² base
 - 80Kλ² / 64b pinst Instruction mem/context
 - 40Kλ² / input-retime depth
 - 1Mλ²+8×0.12Mλ²~≈2Mλ² → 4× LUTs (ideal)
 - Recall ASCIItoHex 3.4, similar for thput map
- They claim 2.8× LUTs

60

Penn ESE534 Spring2010 -- DeHon [MPR/Tabula 3/29/2009]

Admin

- Course Evaluations
- Final exercise updated
 - Baseline design and technology defined
 - Still planning to get more details on sparing-based scheme by Monday
- Office Hours today 2pm
 - Next week
- Reading for next week on web

Penn ESE534 Spring2010 -- DeHon

61

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
- $A_{\text{ctxt}} \ll A_{\text{active}}$ makes interesting
 - save area by sharing active among instructions

Penn ESE534 Spring2010 -- DeHon

62

Big Ideas [MSB-1 Ideas]

- Economical retiming becomes important here to achieve active LUT reduction
 - one output reg/LUT leads to early saturation
- $c=4-8$, $l=4-6$ 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?

Penn ESE534 Spring2010 -- DeHon

63