

# CS184a: Computer Architecture (Structure and Organization)

Day 12: February 2, 2005  
Compute 2:  
Cascades, ALUs, PLAs



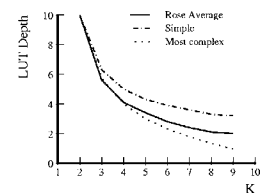
## Last Time

- LUTs
  - area
  - structure
  - big LUTs vs. small LUTs with interconnect
  - design space
  - optimization

## Today

- Cascades
- ALUs
- PLAs

## Last Time



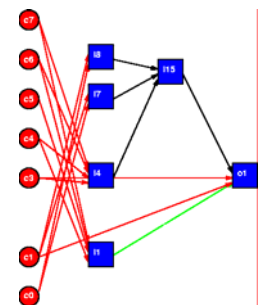
- Larger LUTs
  - Less interconnect delay
- + General: Larger compute blocks
  - Minimize interconnect crossings
- Large LUTs
  - Not efficient for typical logic structure

## Different Structure

- How can we have “larger” compute nodes (less general interconnect) without paying huge area penalty of large LUTs?

## Structure in subgraphs

- Small LUTs capture structure
- What structure does a small-LUT-mapped netlist have?



### Structure

- LUT sequences ubiquitous

Caltech CS184 Winter2003 -- DeHon

### Hardwired Logic Blocks

Single Output

Caltech CS184 Winter2003 -- DeHon

### Hardwired Logic Blocks

Two outputs

9

Caltech CS184 Winter2003 -- DeHon

### Delay Model

- $T_{\text{cascade}} = T(3\text{LUT}) + T(\text{mux})$
- Don't pay
  - General interconnect
  - Full 4-LUT delay

10

Caltech CS184 Winter2003 -- DeHon

### Options

11

Caltech CS184 Winter2003 -- DeHon

### Chung & Rose Study

Figure 8: Delay Study HLB Topologies

12

Caltech CS184 Winter2003 -- DeHon [Chung & Rose, DAC '92]

## Cascade LUT Mappings

Logic Block	$N_R$	% decr in $N_R$	$D_{int}$ (ns)	% decr in $D_{int}$
K4	5.4	0	30	0
L2-2	4.2	22	26	13
L2-3	3.4	37	22	27
L2-4	3.1	43	21	30
L2-5	3.0	44	21	30
L3-3.2	4.0	26	25	17
L3-4.2	3.0	44	21	30
L3-4.3	3.1	43	21	30
L3-5.2.2	3.1	43	21	30
L3-5.3	3.0	44	21	30
L3-5.4	2.9	46	20	33
L3-6.4	2.8	48	20	33

Table 3: Delay Performance of Different HLBs

[Chung & Rose, DAC '92]

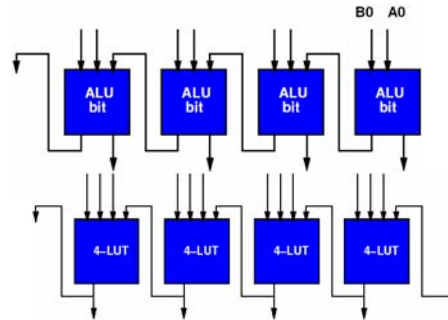
Block Cet	2-inp L2-3	X4000 CLB	4-inp L2-3	4-inp L3-4.2
byrom	40	36	26	19
ab2	77	71	48	37
ab4	126	122	82	61
agat7	44	38	27	20
l3	92	90	56	43
c1354	131	91	80	59
cl	21	17	15	11
cc	17	8	5	7
cm12a	8	5	5	4
comp	27	17	14	13
count	21	21	14	10
decod	19	10	8	8
emx	10	5	5	5
vda	96	97	70	53
vml	3	3	3	2
tot HLBs	653	561	421	300
LUT Bits	16072	22440	20208	20480
Ratio	0.70	1	0.90	0.91
HLB pins	6530	6171	3473	3440
Ratio	1.06	1	0.89	0.88

Table 2: Area Measures of Different HLBs

13

Caltech CS184 Winter2003 -- DeHon

## ALU vs. Cascaded LUT?



14

Caltech CS184 Winter2003 -- DeHon

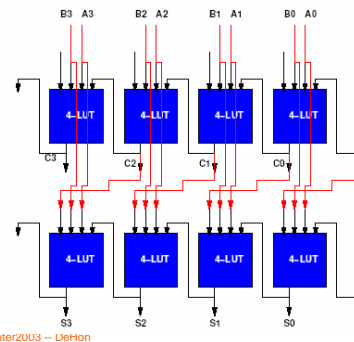
## Datapath Cascade

- ALU/LUT (datapath) Cascade
  - Long “serial” path w/out general interconnect
  - Pay only Tmux and nearest-neighbor interconnect

15

Caltech CS184 Winter2003 -- DeHon

## 4-LUT Cascade ALU

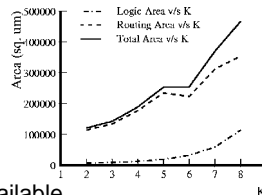


16

Caltech CS184 Winter2003 -- DeHon

## ALU vs. LUT ?

- Compare/contrast
- ALU
  - Only subset of ops available
  - Denser coding for those ops
  - Smaller
  - ...but interconnect dominates
  - [Datapath width orthogonal to function]

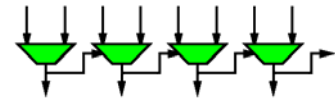


17

Caltech CS184 Winter2003 -- DeHon

## Parallel Prefix LUT Cascade?

- Can we do better than  $N \times T_{mux}$ ?
- Can we compute LUT cascade in  $O(\log(N))$  time?
- Can we compute mux cascade using parallel prefix?



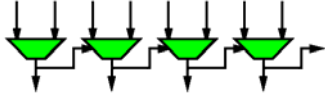
- Can we make mux cascade associative?

18

Caltech CS184 Winter2003 -- DeHon

## Parallel Prefix Mux cascade

- How can mux transform  $S \rightarrow \text{mux-out}$ ?
  - $A=0, B=0 \rightarrow \text{mux-out}=0$
  - $A=1, B=1 \rightarrow \text{mux-out}=1$
  - $A=0, B=1 \rightarrow \text{mux-out}=S$
  - $A=1, B=0 \rightarrow \text{mux-out}=\neg S$

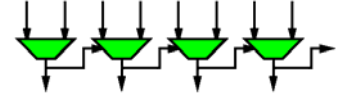


Caltech CS184 Winter2003 -- DeHon

19

## Parallel Prefix Mux cascade

- How can mux transform  $S \rightarrow \text{mux-out}$ ?
  - $A=0, B=0 \rightarrow \text{mux-out}=0$  Stop = S
  - $A=1, B=1 \rightarrow \text{mux-out}=1$  Generate = G
  - $A=0, B=1 \rightarrow \text{mux-out}=S$  Buffer = B
  - $A=1, B=0 \rightarrow \text{mux-out}=\neg S$  Invert = I

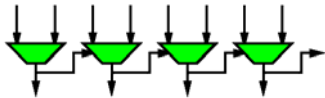


Caltech CS184 Winter2003 -- DeHon

20

## Parallel Prefix Mux cascade

- How can 2 muxes transform input?
- Can I compute 2-mux transforms from 1 mux transforms?



Caltech CS184 Winter2003 -- DeHon

21

## Two-mux transforms

- $SS \rightarrow S$  •  $GS \rightarrow S$  •  $BS \rightarrow S$  •  $IS \rightarrow S$
- $SG \rightarrow G$  •  $GG \rightarrow G$  •  $BG \rightarrow G$  •  $IG \rightarrow G$
- $SB \rightarrow S$  •  $GB \rightarrow G$  •  $BB \rightarrow B$  •  $IB \rightarrow I$
- $SI \rightarrow G$  •  $GI \rightarrow S$  •  $BI \rightarrow I$  •  $II \rightarrow B$

Caltech CS184 Winter2003 -- DeHon

22

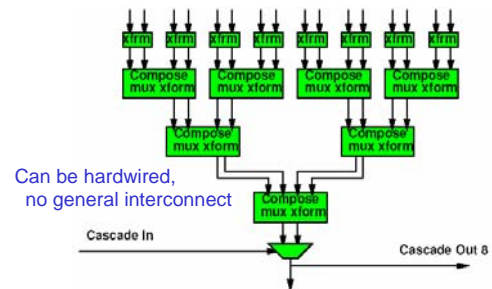
## Generalizing mux-cascade

- How can N muxes transform the input?
- Is mux transform composition associative?

Caltech CS184 Winter2003 -- DeHon

23

## Parallel Prefix Mux-cascade



Caltech CS184 Winter2003 -- DeHon

24

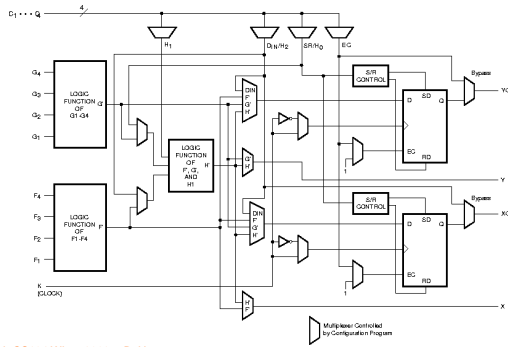
# "ALU"s Unpacked

## Traditional/Datapath ALUs

1. SIMD/Datapath Control
  - Architecture variable  $w$
2. Long Cascade
  - Typically also  $w$ , but can shorter/longer
  - Amenable to parallel prefix implementation in  $O(\log(w))$  time  $w/ O(w)$  space
3. Restricted function
  - Reduces instruction bits
  - Reduces expressiveness

# Commercial Devices

## Xilinx XC4000 CLB



Multiplexer Controlled by Configuration Program

## Xilinx Virtex-II

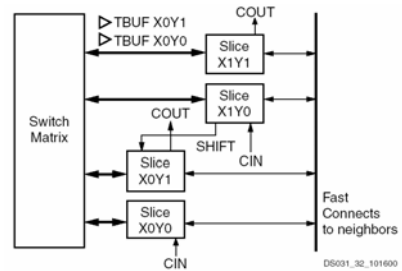


Figure 14: Virtex-II CLB Element

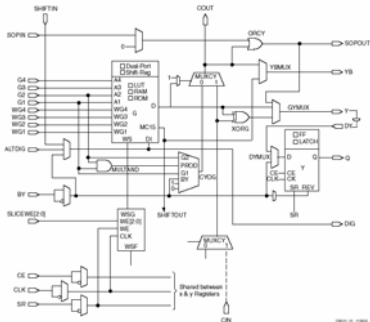


Figure 16: Virtex-II Slice (Top Half)

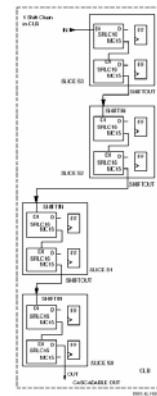


Figure 17: Virtex-II Slice (Bottom Half)



## Wired-or

- Connect series of inputs to wire
- Any of the inputs can drive the wire high

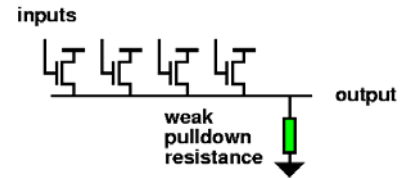


Caltech CS184 Winter2003 -- DeHon

37

## Wired-or

- Implementation with Transistors

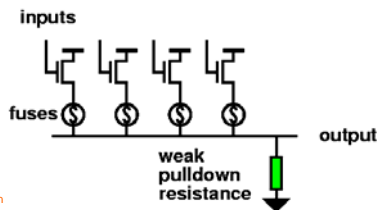


Caltech CS184 Winter2003 -- DeHon

38

## Programmable Wired-or

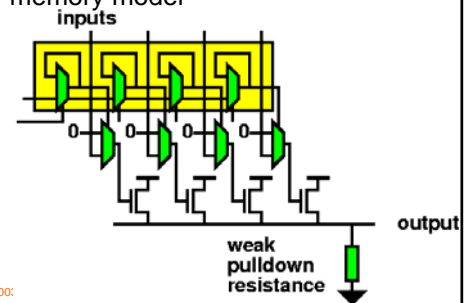
- Use some memory function to programmable connect (disconnect) wires to OR
- Fuse:



Caltech CS184 Winter2003 -- DeHon

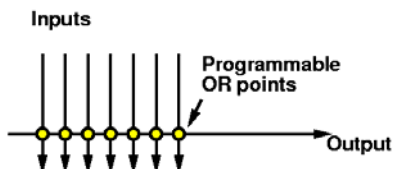
## Programmable Wired-or

- Gate-memory model



Caltech CS184 Winter2003

## Diagram Wired-or

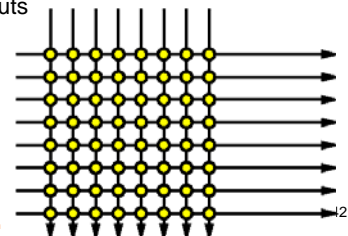


Caltech CS184 Winter2003 -- DeHon

41

## Wired-or array

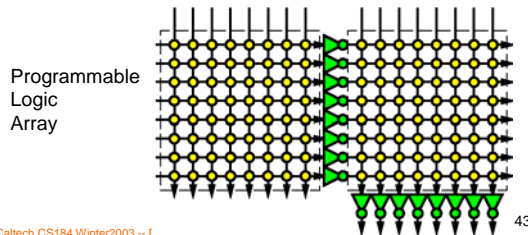
- Build into array
  - Compute many different or functions from set of inputs



Caltech CS184 Winter2003 -- DeHon

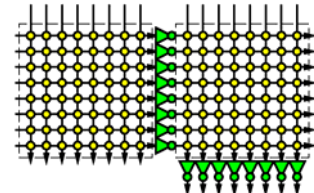
## Combined **or**-arrays to PLA

- Combine two **or** (nor) arrays to produce PLA (**and-or** array)



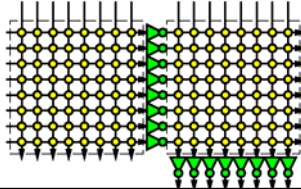
## PLA

- Can implement each **and** on single line in first array
- Can implement each **or** on single line in second array



## PLA

- Efficiency questions:
  - Each **and/or** is linear in total number of potential inputs (not actual)
  - How many product terms between arrays?



## PLA Product Terms

- Can be exponential in number of inputs
- *E.g.* n-input **xor** (parity function)
  - When flatten to two-level logic, requires exponential product terms
  - $a!b+!a*b$
  - $a!*b!*c+!a*b!*c+!a!*b*c+a*b*c$
- ...and shows up in important functions
  - Like addition...

46

## PLAs

- Fast Implementations for large ANDs or ORs
- Number of P-terms **can be** exponential in number of input bits
  - most complicated functions
  - not exponential for many functions
- Can use arrays of small PLAs
  - to exploit structure
  - like we saw arrays of small memories last time

Caltech CS184 Winter2003 -- DeHon

47

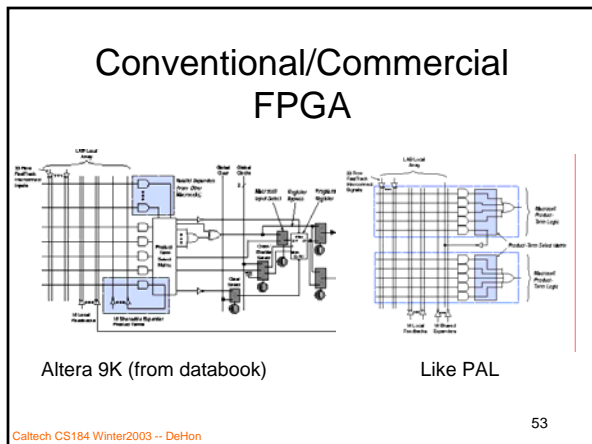
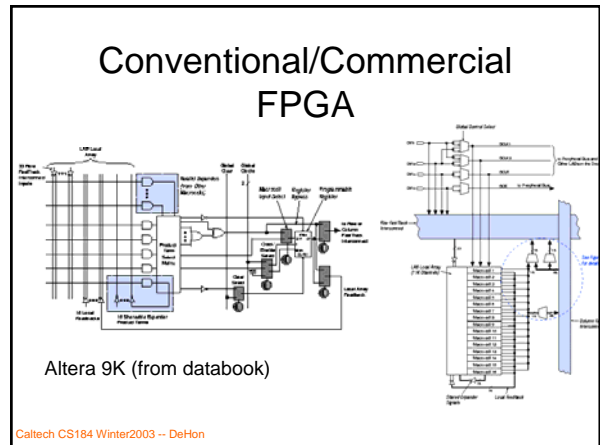
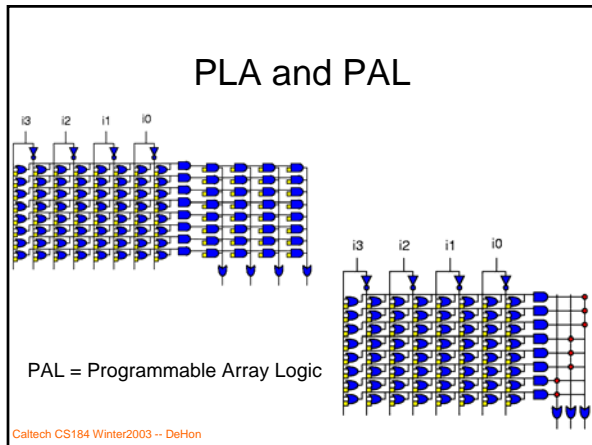
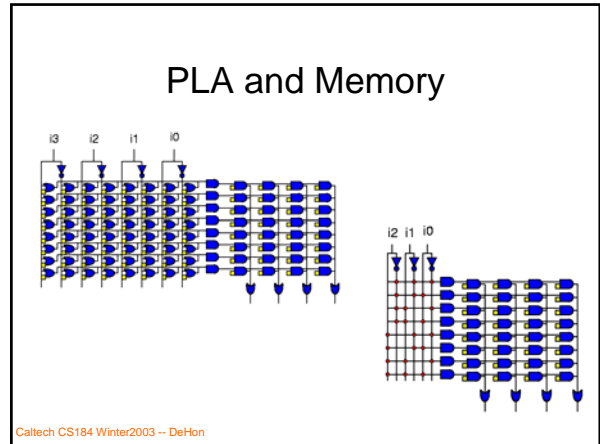
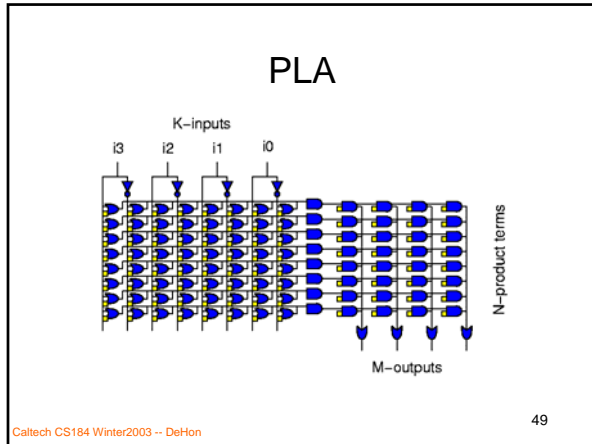
## PLAs vs. LUTs?

- Look at Inputs, Outputs, P-Terms
  - minimum area (one study, see paper)
  - $K=10, N=12, M=3$
- $A(\text{PLA } 10,12,3)$  comparable to 4-LUT?
  - 80-130%?
  - 300% on ECC (structure LUT can exploit)
- Delay?
  - Claim 40% fewer logic levels (4-LUT)
    - (general interconnect crossings)

Caltech CS184 Winter2003 -- DeHon

[Kouloheris & El Gamal/CICC'92] 48





- ### Big Ideas [MSB Ideas]
- Programmable Interconnect allows us to exploit that structure
    - want to match to application structure
    - Prog. interconnect delay expensive
  - Hardwired Cascades
    - key technique to reducing delay in programmables
  - PLAs
    - canonical two level structure
    - hardwire portions to get Memories, PALs 54
- Caltech CS184 Winter2003 -- DeHon

## Big Ideas [MSB-1 Ideas]

- Better structure match with hardwired LUT cascades