

# ESE534: Computer Organization

Day 15: March 22, 2010  
Compute 2:  
Cascades, ALUs, PLAs



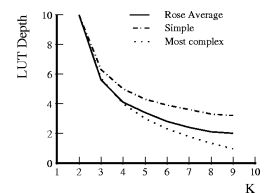
## Last Time

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

## Today

- ALUs
- Cascades
- PLAs

## Last Time



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

## Preclass

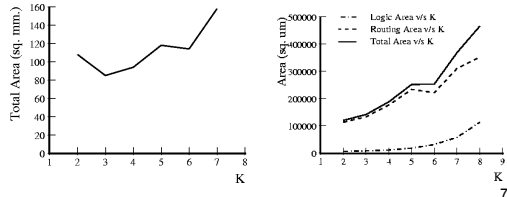
- How does addition delay compare between ALU-based and LUT-based architectures?
- What is the source of the advantage?

## Preclass

- Advantages of ALU compute block over LUT?
- Disadvantages?

## Implications from LUT-size study

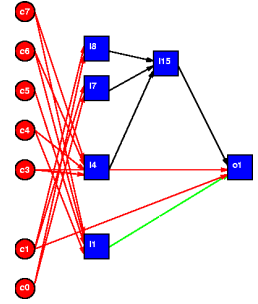
- ALU has more restricted logic functions
  - Require more blocks



Penn ESE534 Spring2010 -- DeHon

## Structure in subgraphs

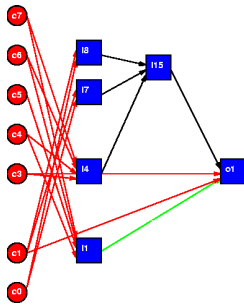
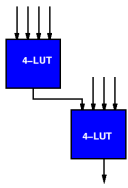
- Small LUTs capture structure
- What structure does a small-LUT-mapped netlist have?



Penn ESE534 Spring2010 -- DeHon

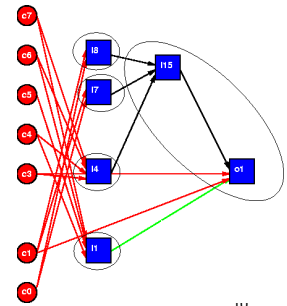
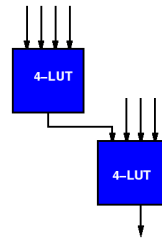
## Structure

- LUT sequences ubiquitous



Penn ESE534 Spring2010 -- DeHon

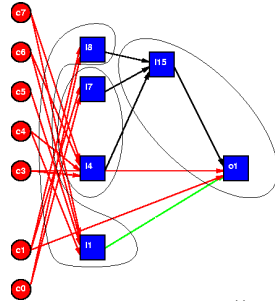
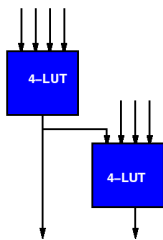
## Hardwired Logic Blocks



Single Output

Penn ESE534 Spring2010 -- DeHon

## Hardwired Logic Blocks

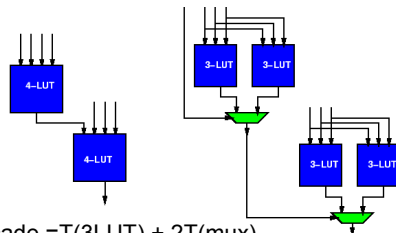


Two outputs

11

Penn ESE534 Spring2010 -- DeHon

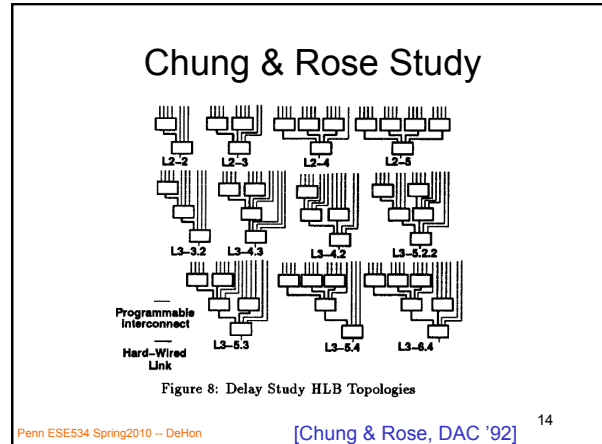
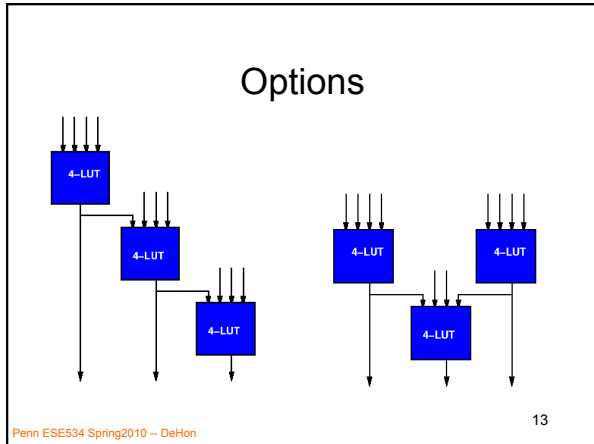
## Delay Model



- $T_{\text{cascade}} = T(3\text{LUT}) + 2T(\text{mux})$
- Don't pay
  - General interconnect
  - Full 4-LUT delay

12

Penn ESE534 Spring2010 -- DeHon



### Chung & Rose Study

Logic Block	$\overline{N_R}$	% decr in $N_R$	$D_{tot}$ (ns)	% decr in $D_{tot}$
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

Figure 8: Delay Study HLB Topologies

Table 3: Delay Performance of Different HLBs

15

Penn ESE534 Spring2010 -- DeHon [Chung & Rose, DAC '92]

### Cascade LUT Mappings

Logic Block	$\overline{N_R}$	% decr in $N_R$	$D_{tot}$ (ns)	% decr in $D_{tot}$
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

Bench Cct	3-inp L2-3	X4050 CLB	4-inp L2-3	4-inp L3-4.2
byteadd	60	36	26	19
alu2	77	71	48	37
alu4	126	122	82	61
spex7	14	38	27	20
lp	22	20	18	12
cl385	151	92	80	56
ck	21	17	15	11
cc	17	8	9	7
cm162a	8	5	5	4
comp	27	17	14	13
count	21	21	14	10
decod	10	10	8	8
max	10	5	5	5
via	56	97	70	52
v4ml	3	3	3	2
<b>Tot HLBs</b>	<b>653</b>	<b>561</b>	<b>421</b>	<b>320</b>
<b>LUT Bits</b>	<b>15072</b>	<b>22440</b>	<b>20208</b>	<b>26480</b>
<b>Ratioes</b>	<b>0.70</b>	<b>1</b>	<b>0.90</b>	<b>0.91</b>
<b>HLB pins Ratioes</b>	<b>6330</b>	<b>6171</b>	<b>5473</b>	<b>5440</b>
	<b>1.06</b>	<b>1</b>	<b>0.89</b>	<b>0.88</b>

Table 3: Delay Performance of Different HLBs

Table 2: Area Measures of Different HLBs

16

Penn ESE534 Spring2010 -- DeHon [Chung & Rose, DAC '92]

### Energy Impact?

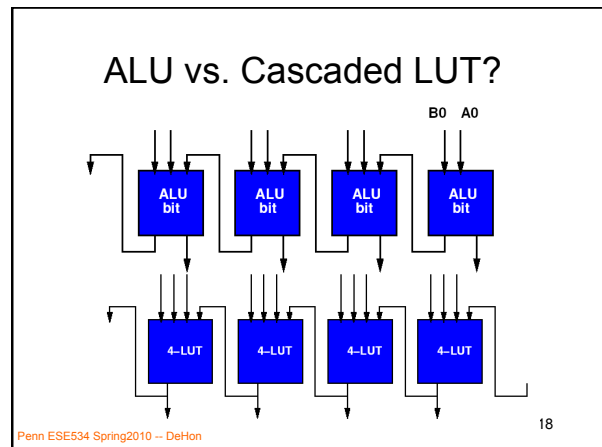
- What's the likely energy impact?

**PGA Power Breakdown**

[Kusse and Rabaey / ISLPED 1998]

17

Penn ESE534 Spring2010 -- DeHon



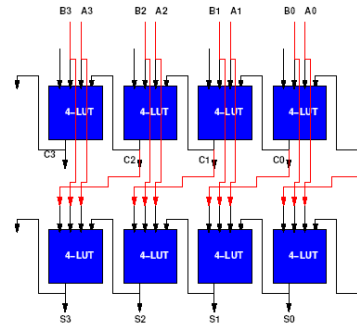
## Datapath Cascade

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

Penn ESE534 Spring2010 -- DeHon

19

## 4-LUT Cascade ALU

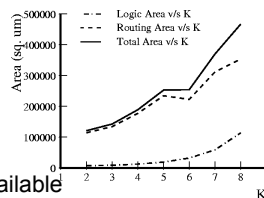


Penn ESE534 Spring2010 -- DeHon

20

## ALU vs. LUT ?

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

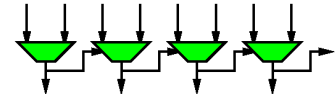


Penn ESE534 Spring2010 -- DeHon

21

## Accelerating LUT Cascade?

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



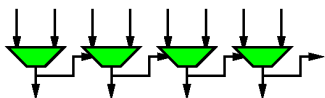
- Can we make mux cascade associative?

Penn ESE534 Spring2010 -- DeHon

22

## Parallel Prefix Mux cascade

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

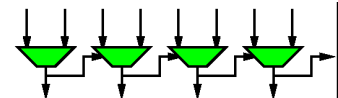


Penn ESE534 Spring2010 -- DeHon

23

## 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
  - $A=1, B=0 \rightarrow \text{mux-out}=\text{!}S$  Invert = I

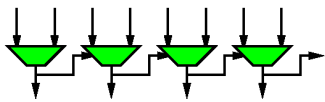


Penn ESE534 Spring2010 -- DeHon

24

## Parallel Prefix Mux cascade

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



Penn ESE534 Spring2010 -- DeHon

25

## Two-mux transforms

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

Penn ESE534 Spring2010 -- DeHon

26

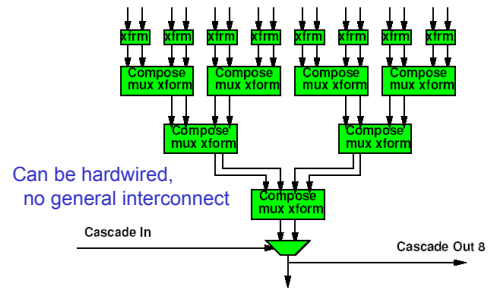
## Generalizing mux-cascade

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

Penn ESE534 Spring2010 -- DeHon

27

## Parallel Prefix Mux-cascade

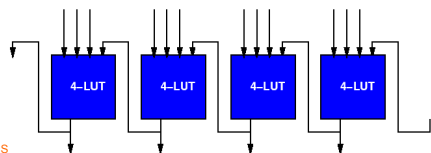


Penn ESE534 Spring2010 -- DeHon

28

## LUT Cascade Implications

- We can compute **any** LUT cascade in logarithmic time
  - Not just addition carry cascade
  - Can be different functions in each LUT
    - Not demand SIMD op



Penn ESE534 S

29

## “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

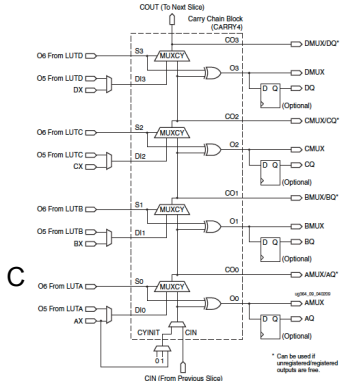
Penn ESE534 Spring2010 -- DeHon

30

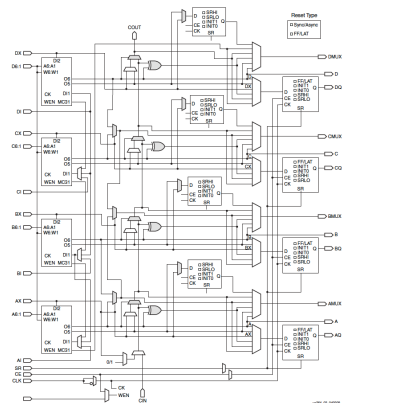
# Commercial Devices

# Virtex 6 Carry

- Prop=A xor B
- Otherwise
  - Generate
  - Squash
- Sum A xor B xor C



# Virtex 6



# Xilinx Virtex-II

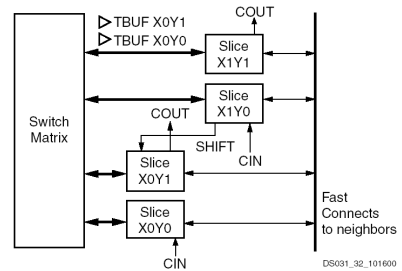


Figure 14: Virtex-II CLB Element

# Virtex 2 Slice

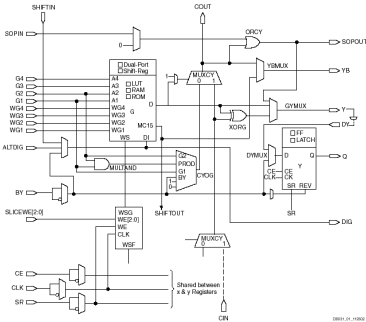
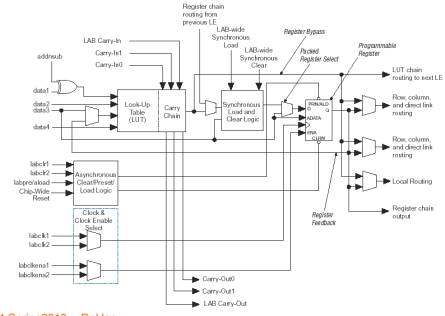
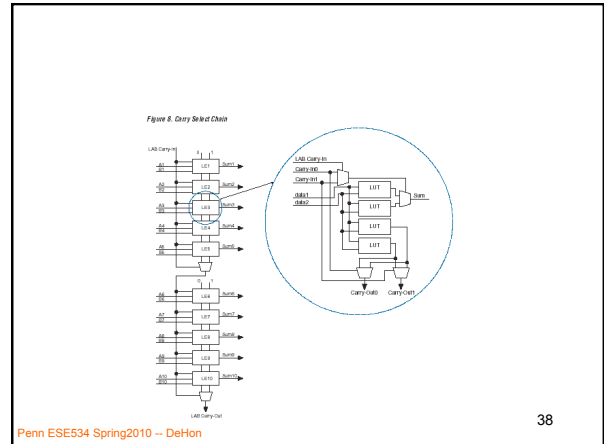
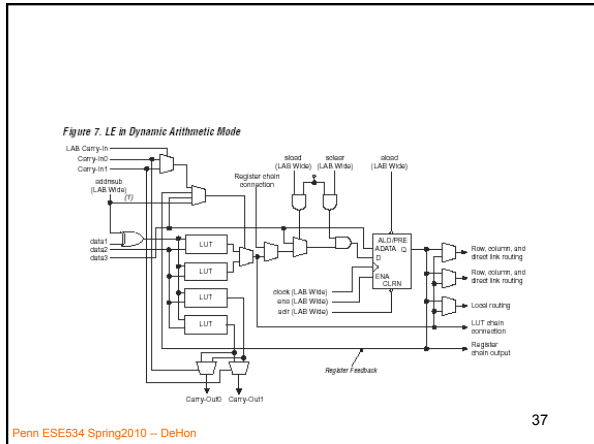


Figure 16: Virtex-II Slice (Top Half)

# Altera Stratix

Figure 5: Stratix LE





## Programmable Array Logic (PLAs)

Penn ESE534 Spring2010 -- DeHon

39

## PLA

- Directly implement flat (two-level) logic
  - $O = a*b*c*d + !a*b*d + b*c*d$
- Exploit substrate properties allow wired-OR

Penn ESE534 Spring2010 -- DeHon

40

## Wired-or

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

Penn ESE534 Spring2010 -- DeHon

41

## Wired-or

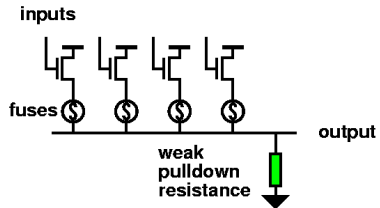
- Implementation with Transistors

Penn ESE534 Spring2010 -- DeHon

42

## Programmable Wired-or

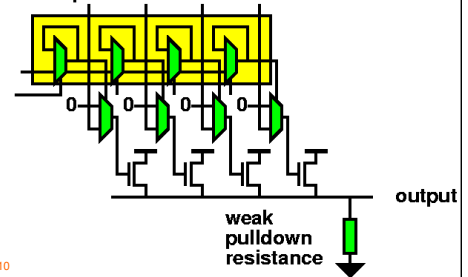
- Use some memory function to programmable connect (disconnect) wires to OR
- Fuse:



Penn ESE534 Spring2010 -- DeHon

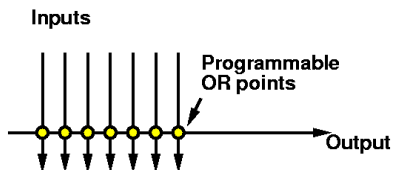
## Programmable Wired-or

- Gate-memory model



Penn ESE534 Spring2010

## Diagram Wired-or

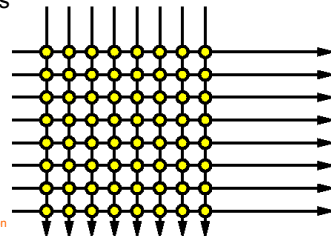


Penn ESE534 Spring2010 -- DeHon

45

## Wired-or array

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

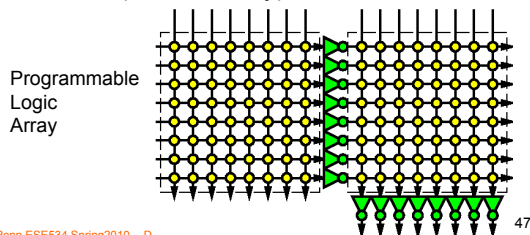


Penn ESE534 Spring2010 -- DeHon

46

## Combined or-arrays to PLA

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

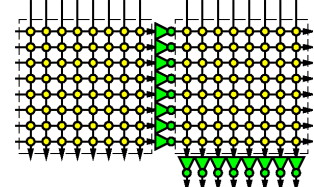


Penn ESE534 Spring2010 -- DeHon

47

## PLA

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

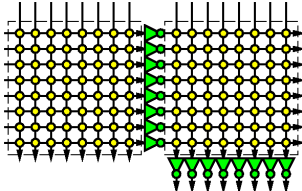


Penn ESE534 Spring2010 -- DeHon



## PLA

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



Penn ESE534 Spring2010 -- DeHon

## Preclass 5

- How many product terms in  $S[w]$  in MSP?

Penn ESE534 Spring2010 -- DeHon

50

## Carries

- $C[1]=A[0]*B[0]$ 
  - $P[1]=1$
- $C[2]=A[1]*B[1]+C[1]*A[1]+C[1]*B[1]$ 
  - $P[2]=1+2*P[1]=3$
- $C[3]=A[2]*B[2]+C[2]*A[2]+C[2]*B[2]$ 
  - $P[3]=1+2*P[2]=7$
- $P[4]=15, P[5]=31, P[6]=63$
- $P[w]=2^w-1$

Penn ESE534 Spring2010 -- DeHon

51

## Sum

- $S[w]=A[w] \text{ xor } B[w] \text{ xor } C[w]$
- $S[w]=(A[w]*B[w]+A[w]*B[w])*C[w] + (A[w]*B[w]+A[w]*B[w])*C[w]$
- Product terms =  $2*P[w]+2*P[w]$

Penn ESE534 Spring2010 -- DeHon

52

## $S[2]$ in MSP

$a2*/a0*/b2*/b1+$   
 $/a2*/a0*/b2*/b1+$   
 $a2*/a1*/a0*/b2+$   
 $/a2*/a1*/a0*/b2+$   
 $a2*/b2*/b1*/b0+$   
 $/a2*/b2*/b1*/b0+$   
 $a2*/a2*/b2*/b0+$   
 $/a2*/a1*/b2*/b0+$   
 $/a2*/a0*/b2*/b1*/b0+$   
 $/a2*/a1*/a0*/b2*/b0+$   
 $a2*/a0*/b2*/b1*/b0+$   
 $a2*/a1*/a0*/b2*/b0+$   
 $a2*/a1*/b2*/b1+$   
 $/a2*/a1*/b2*/b1+$   
 $/a2*/a1*/b2*/b1$

Penn ESE534 Spring2010 -- DeHon

53

## PLA Product Terms

- Can be exponential in number of inputs
- *E.g.* n-input **xor** (parity function)
  - When flattened 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 additions, as we just saw...

Penn ESE534 Spring2010 -- DeHon

54

## 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

Penn ESE534 Spring2010 -- DeHon

55

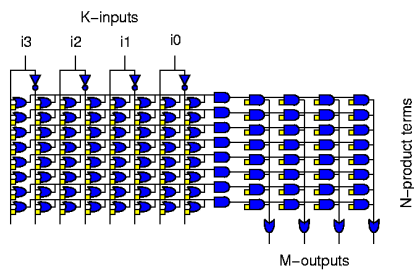
## PLAs vs. LUTs?

- 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)

[Kouloheris & El Gamal/CICC'92] 56

Penn ESE534 Spring2010 -- DeHon

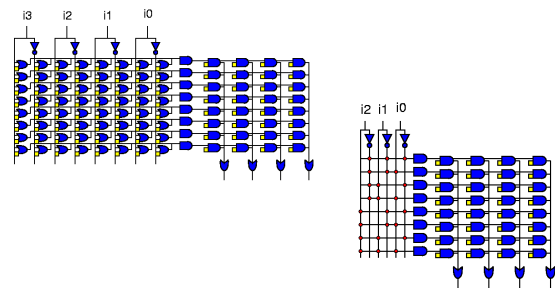
## PLA



Penn ESE534 Spring2010 -- DeHon

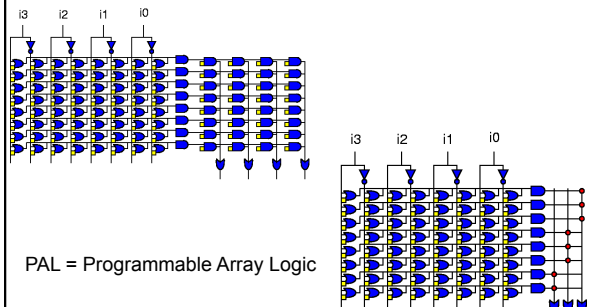
57

## PLA and Memory



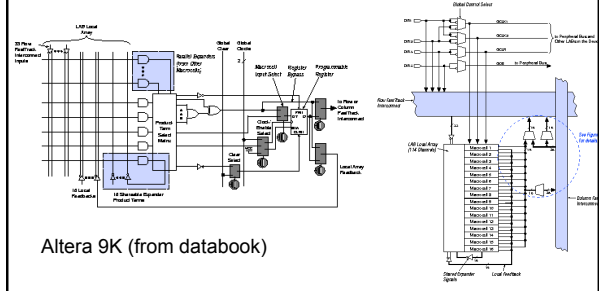
Penn ESE534 Spring2010 -- DeHon

## PLA and PAL



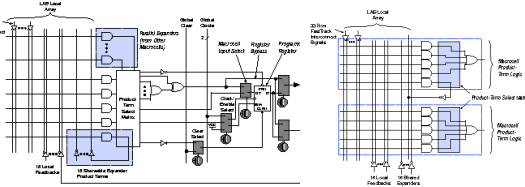
Penn ESE534 Spring2010 -- DeHon

## Conventional/Commercial FPGA



Penn ESE534 Spring2010 -- DeHon

## Conventional/Commercial FPGA



Altera 9K (from databook)

Like PAL

Penn ESE534 Spring2010 -- DeHon

61

## Admin

- HW7 out
  - Covers Compute blocks
    - (last lecture and this)
- Reading for Wed. on web
- Guest lecture Wednesday
  - No office hours Wednesday

Penn ESE534 Spring2010 -- DeHon

62

## Big Ideas [MSB Ideas]

- Programmable Interconnect allows us to exploit structure in logic
  - want to match to application structure
  - Prog. interconnect delay expensive
- Hardwired Cascades
  - key technique to reducing delay in programmables (both ALUs and LUTs)
- PLAs
  - canonical two level structure
  - hardwire portions to get Memories, PALs

Penn ESE534 Spring2010 -- DeHon

63