

# ESE534: Computer Organization

Day 13: February 29, 2012  
Compute 2:  
Cascades, ALUs, PLAs



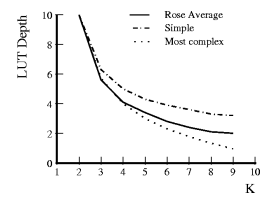
## Last Time

- 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

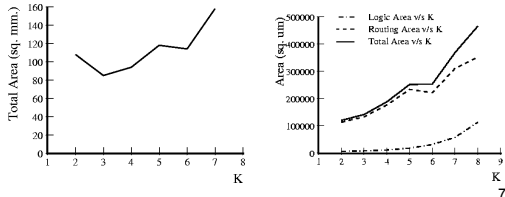
- How does addition delay compare between ALU-based and LUT-based architectures?
- What is the source of the advantage?

## Preclass

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

## Implications from LUT-size study

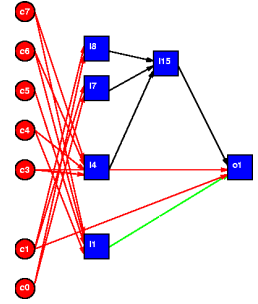
- ALU has more restricted logic functions
  - Will require more blocks



Penn ESE534 Spring2012 -- DeHon

## Structure in subgraphs

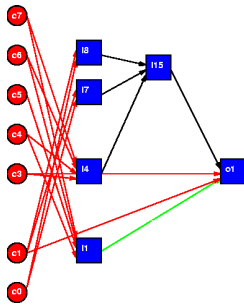
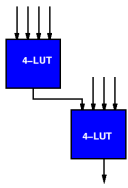
- Small LUTs capture structure
- What structure does a small-LUT-mapped netlist have?



Penn ESE534 Spring2012 -- DeHon

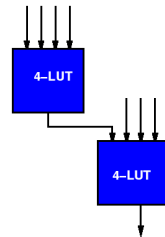
## Structure

- LUT sequences ubiquitous

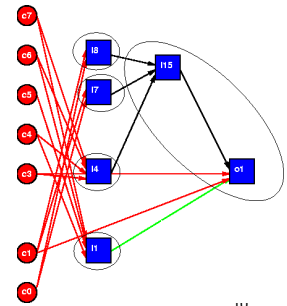


Penn ESE534 Spring2012 -- DeHon

## Hardwired Logic Blocks

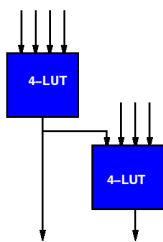


Single Output

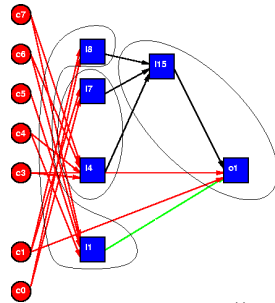


Penn ESE534 Spring2012 -- DeHon

## Hardwired Logic Blocks



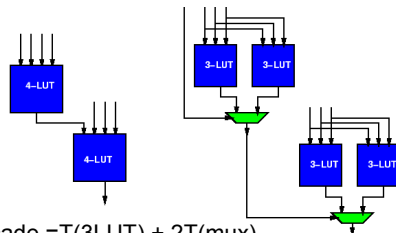
Two outputs



11

Penn ESE534 Spring2012 -- 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 Spring2012 -- DeHon

## Options

13

Penn ESE534 Spring2012 -- DeHon

## Chung & Rose Study

- Which most help:
  - Ripple adder?
  - Associative reduce tree?
  - Multiplier?

Figure 8: Delay Study HLB Topologies

14

Penn ESE534 Spring2012 -- DeHon [Chung & Rose, DAC '92]

## Chung & Rose Study

Figure 8: Delay Study HLB Topologies

Logic Block	$\bar{N}_R$	% decr in $\bar{N}_R$	$\bar{D}_{tot}$ (ns)	% decr in $\bar{D}_{tot}$
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

Table 3: Delay Performance of Different HLBs

15

Penn ESE534 Spring2012 -- DeHon [Chung & Rose, DAC '92]

## Cascade LUT Mappings

Logic Block	$\bar{N}_R$	% decr in $\bar{N}_R$	$\bar{D}_{tot}$ (ns)	% decr in $\bar{D}_{tot}$
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

Table 3: Delay Performance of Different HLBs

Bench Cct	3-inp L2-3	X4000 CLB	4-inp L2-3	4-inp L3-4.2
byteadd	60	36	26	19
aha2	77	71	48	37
aha4	126	122	82	61
npw7	14	38	27	20
lp	22	20	18	13
cl355	151	91	80	56
cc	21	17	15	11
cc	17	8	9	7
cm16a	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>22460</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>6630</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 2: Area Measures of Different HLBs

16

Penn ESE534 Spring2012 -- DeHon [Chung & Rose, DAC '92]

## Energy Impact?

- What's the likely energy impact?

XC4003A data from Eric Kusse (UCB MS 1997)

17

Penn ESE534 Spring2012 -- DeHon [Virtex II, Shang et al., FPGA 2002]

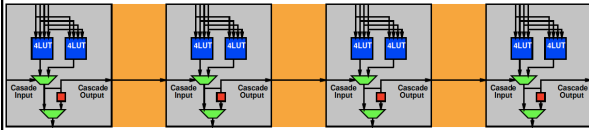
## ALU vs. Cascaded LUT?

18

Penn ESE534 Spring2012 -- DeHon

## Datapath Cascade

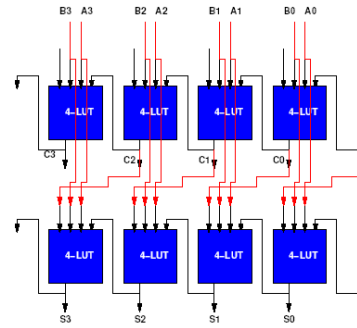
- ALU/LUT (datapath) Cascade
  - Long “serial” path w/out general interconnect
  - Pay only Tmux and nearest-neighbor interconnect



Penn ESE534 Spring2012 -- DeHon

19

## 4-LUT Cascade ALU

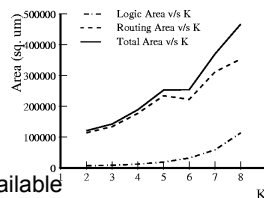


Penn ESE534 Spring2012 -- DeHon

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]

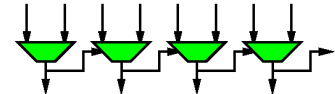


Penn ESE534 Spring2012 -- DeHon

21

## Accelerating LUT Cascade?

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



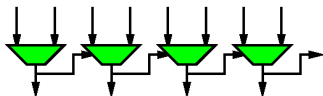
- Can we make mux cascade associative?

Penn ESE534 Spring2012 -- DeHon

22

## Parallel Prefix Mux cascade

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

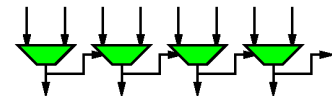


Penn ESE534 Spring2012 -- DeHon

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

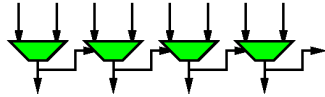


Penn ESE534 Spring2012 -- DeHon

24

## Parallel Prefix Mux cascade

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



Penn ESE534 Spring2012 -- 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 Spring2012 -- DeHon

26

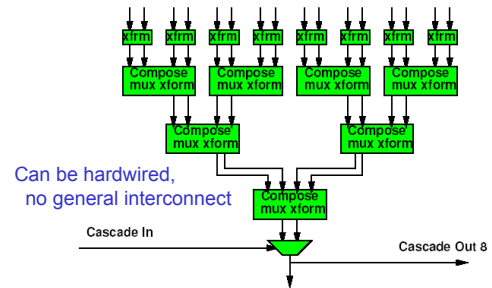
## Generalizing mux-cascade

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

Penn ESE534 Spring2012 -- DeHon

27

## Parallel Prefix Mux-cascade

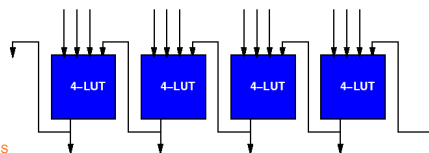


Penn ESE534 Spring2012 -- DeHon

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 combine 3 ideas

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

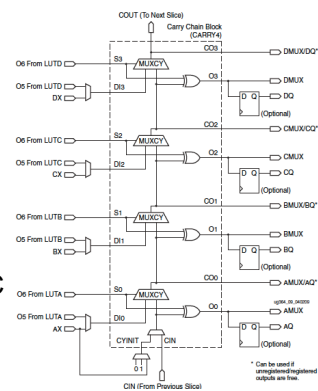
Penn ESE534 Spring2012 -- DeHon

30

## Commercial Devices

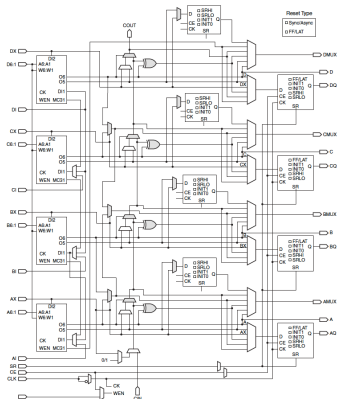
## Virtex 6 & 7 Carry

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



## Virtex 6

- V7 similar



## Programmable Array Logic (PLAs)

## PLA

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

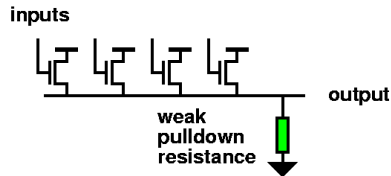
## Wired-or

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



## Wired-or

- Implementation with Transistors

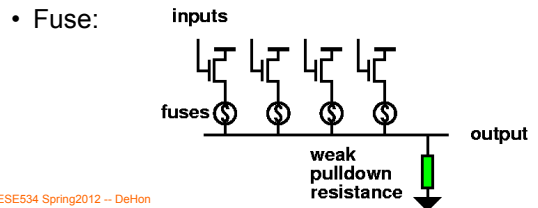


Penn ESE534 Spring2012 -- DeHon

37

## Programmable Wired-or

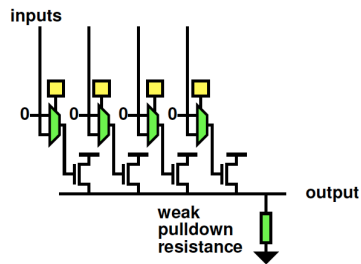
- Use some memory function to programmable connect (disconnect) wires to OR



Penn ESE534 Spring2012 -- DeHon

## Programmable Wired-or

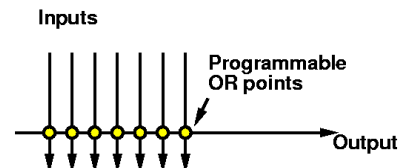
- Gate-memory model



Penn ESE534 Spring2012 -- DeHon

39

## Diagram Wired-or

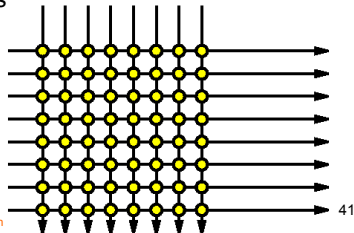


Penn ESE534 Spring2012 -- DeHon

40

## Wired-or array

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

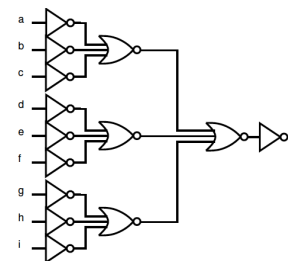


Penn ESE534 Spring2012 -- DeHon

41

## Preclass 6

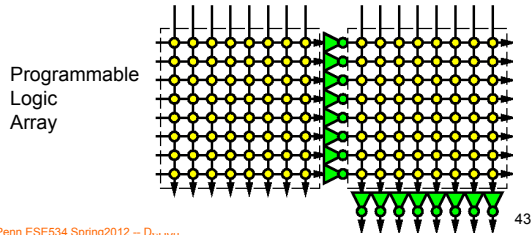
- What function (MSP) ?



Penn ESE534 Spring2012 -- DeHon

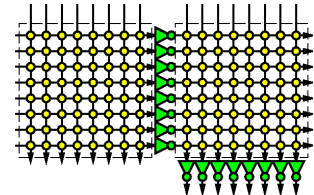
## Combined or-arrays to PLA

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



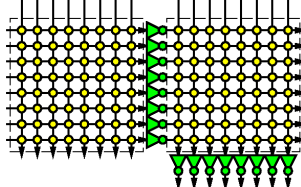
## PLA

- 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?



## Preclass 5

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

Penn ESE534 Spring2012 -- DeHon

46

## 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 Spring2012 -- DeHon

47

## 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 =  
 $2*P[w]+2*Pterms(/C[w])$

Penn ESE534 Spring2012 -- DeHon

48



## S[2] in MSP

$$\begin{aligned}
 &a^2/a^0/b^2/b^1+ \\
 &/a^2/a^0/b^2/b^1+ \\
 &a^2/a^1/a^0/b^2+ \\
 &/a^2/a^1/a^0/b^2+ \\
 &a^2/b^2/b^1/b^0+ \\
 &/a^2/b^2/b^1/b^0+ \\
 &a^2/a^2/b^2/b^0+ \\
 &/a^2/a^1/b^2/b^0+ \\
 &/a^2/a^0/b^2/b^1/b^0+ \\
 &/a^2/a^1/a^0/b^2/b^0+ \\
 &a^2/a^0/b^2/b^1/b^0+ \\
 &a^2/a^1/a^0/b^2/b^0+ \\
 &a^2/a^1/b^2/b^1+ \\
 &/a^2/a^1/b^2/b^1+ \\
 &/a^2/a^1/b^2/b^1
 \end{aligned}$$

49

Penn ESE534 Spring2012 -- DeHon

## 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...

50

Penn ESE534 Spring2012 -- DeHon

## 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

51

Penn ESE534 Spring2012 -- DeHon

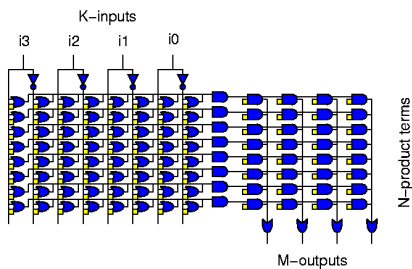
## PLAs vs. LUTs?

- Look at Inputs, Outputs, P-Terms
  - minimum area (one study, see paper)
  - $K=10, N=12, M=3$
- A(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)
    - Not comparing to hardwired cascades

[Kouloheris & El Gamal/CICC'92] 52

Penn ESE534 Spring2012 -- DeHon

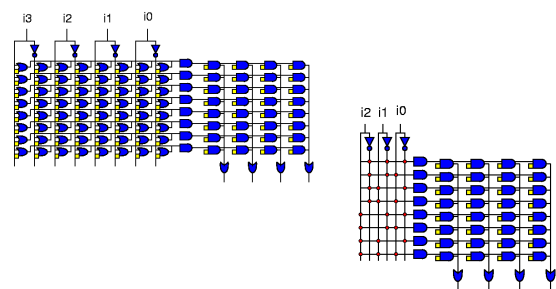
## PLA



53

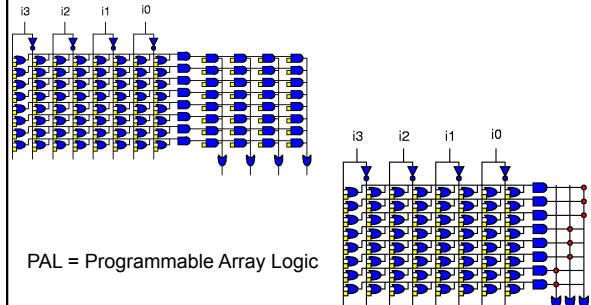
Penn ESE534 Spring2012 -- DeHon

## PLA and Memory



Penn ESE534 Spring2012 -- DeHon

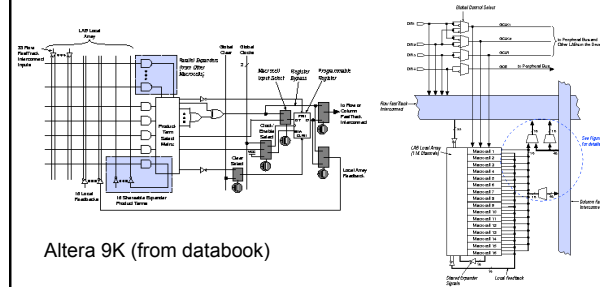
## PLA and PAL



PAL = Programmable Array Logic

Penn ESE534 Spring2012 -- DeHon

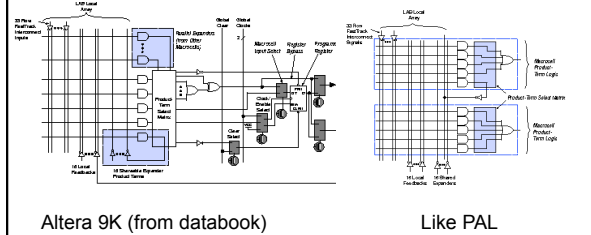
## Conventional/Commercial FPGA



Altera 9K (from databook)

Penn ESE534 Spring2012 -- DeHon

## Conventional/Commercial FPGA



Altera 9K (from databook)

Like PAL

Penn ESE534 Spring2012 -- DeHon

57

## Admin

- HW6.1-2 due Friday
- Spring Break next week
- Reading for Monday after break on web
- HW6.3-4 on Wed. after break

Penn ESE534 Spring2012 -- DeHon

58

## 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 Spring2012 -- DeHon

59