

# ESE534: Computer Organization

Day 12: February 27, 2012  
Compute 1: LUTs



# Previously

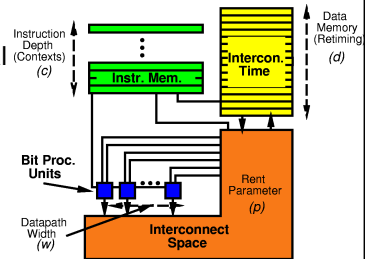
- Instruction Space Modeling
  - huge range of densities
  - huge range of efficiencies
  - large architecture space
  - modeling to understand design space

# Today

- Look at Programmable Compute Blocks
- Specifically LUTs
- Introduce recurring theme (methodology):
  - define parameterized space
  - identify costs and benefits
  - look at typical application requirements
  - compose results, try to find best point

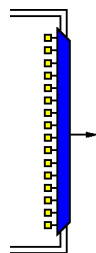
# Compute Function

- What do we use for “compute” function?
- Any Universal
  - NANDx
  - ALU
  - LUT

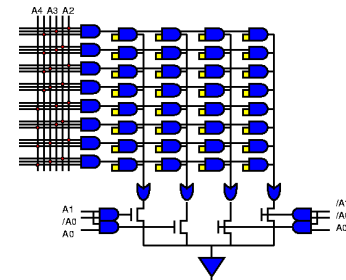


# Lookup Table

- Load bits into table
  - $2^N$  bits to describe
  - $\rightarrow 2^{2^N}$  different functions
- Table translation
  - performs logic transform



# Lookup Table



## We could...

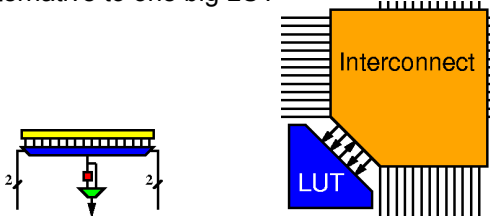
- Just build a large memory = large LUT
- Put our function in there
- What's wrong with that?

## How big is a k-LUT?

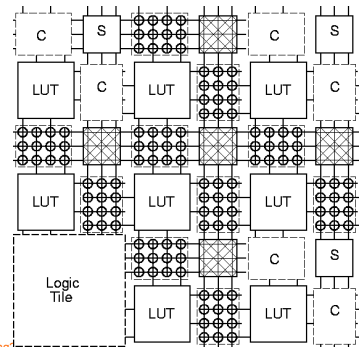
- k-input, 1-output?
- k-input, m-output?

## FPGA = Many small LUTs

Alternative to one big LUT



## Toronto FPGA Model



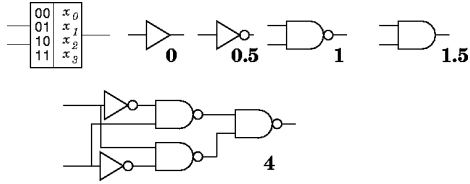
## What's best to use?

- Small LUTs
- Large Memories
- ...small LUTs or large LUTs
- **Continuum question:** how big should our memory blocks used to perform computation be?

## Start to Sort Out: Big vs. Small Luts

- Establish equivalence
  - how many small LUTs equal one big LUT?

## “gates” in 2-LUT ?

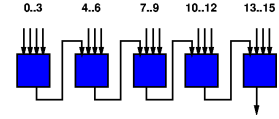


Penn ESE534 Spring2012 -- DeHon

13

## How Much Logic in a LUT?

- Lower Bound?
  - Concrete: 4-LUTs to implement M-LUT?
- Not use all inputs?
  - 0 ... maybe 1
- Use all inputs?
  - $(M-1)/3$



$(M-1)/(k-1)$  for K-lut

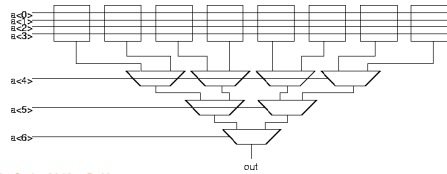
example M-input AND  
 • cover 4 ins w/ first 4-LUT,  
 • 3 more and cascade input  
 with each additional

14

Penn ESE534 Spring2012 -- DeHon

## How much logic in a LUT?

- Upper Upper Bound?:
  - M-LUT implemented w/ 4-LUTs
  - $M\text{-LUT} \leq 2^{M-4} + (2^{M-4} - 1) \leq 2^{M-3}$  4-LUTs



Penn ESE534 Spring2012 -- DeHon

15

## How Much?

- Lower Upper Bound:
  - $2^{2^M}$  functions realizable by M-LUT
  - Say Need  $n$  4-LUTs to cover; compute  $n$ :
    - strategy count functions realizable by each
    - $(2^{2^4})^n \geq 2^{2^M}$
    - $n \log(2^{2^4}) \geq \log(2^{2^M})$
    - $n 2^4 \log(2) \geq 2^M \log(2)$
    - $n 2^4 \geq 2^M$
    - $n \geq 2^{M-4}$

Penn ESE534 Spring2012 -- DeHon

16

## How Much?

- Combine
  - Lower Upper Bound
  - Upper Lower Bound
  - (number of 4-LUTs in M-LUT)

$$2^{M-4} \leq n \leq 2^{M-3}$$

Penn ESE534 Spring2012 -- DeHon

17

## Memories and 4-LUTs

- For the **most complex** functions
  - an M-LUT has  $\sim 2^{M-4}$  4-LUTs
- ◇ SRAM 32Kx8  $\lambda=0.6\mu\text{m}$ 
  - $170M\lambda^2$  (21ns latency)
  - $8 \cdot 2^{11} = 16K$  4-LUTs
- ◇ XC3042  $\lambda=0.6\mu\text{m}$ 
  - $180M\lambda^2$  (13ns delay per CLB)
  - 288 4-LUTs
- Memory is 50+x denser than FPGA
  - ... and faster

Penn ESE534 Spring2012 -- DeHon

18

## Memory and 4-LUTs

- For “regular” functions?
  - ◇ 15-bit parity
    - entire 32Kx8 SRAM
    - How many 4-LUTs?
    - 5 4-LUTs
      - (2% of XC3042  $\sim 3.2M\lambda^2 \sim 1/50$ th Memory)

## Preclass: 16-bit Adder from Memory and 3-LUTs

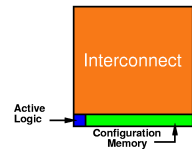
- How many inputs? outputs?
- Area for single large LUT?
- How many 3-LUTs?
- Area per 3-LUT?
- LUT area to implement adder with 3-LUTs?
  - Not include interconnect
- Ratio?

## Memory and 4-LUTs

- Same 32Kx8 SRAM
  - ◇ 7b Add
    - entire 32Kx8 SRAM (largest will support)
    - 14 4-LUTs
      - (5% of XC3042,  $8.8M\lambda^2 \sim 1/20$ th Memory)

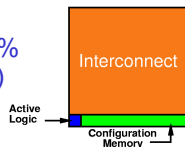
## LUT + Interconnect

- Interconnect allows us to exploit **structure** in computation
- Consider addition:
  - N-input add takes
    - 2N 3-LUTs
    - one N-output (2N)-LUT
  - $N \times 2^{(2N)} \gg 2N \times 2^3$
  - $N=16: 16 \times 2^{32} \gg 2 \times 16 \times 2^3$
  - $2^{36} \gg 2^8 \rightarrow$  factor of  $2^{28} = 256$  Million



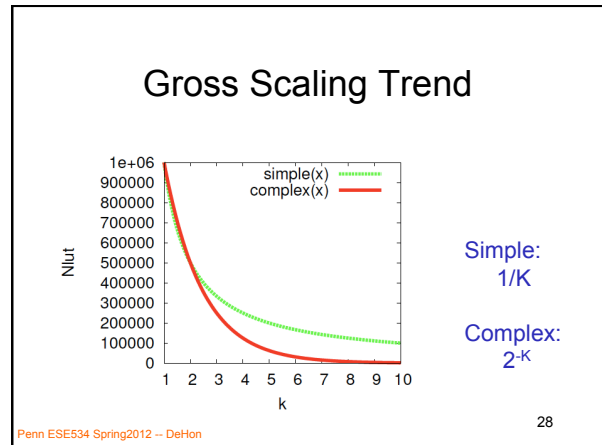
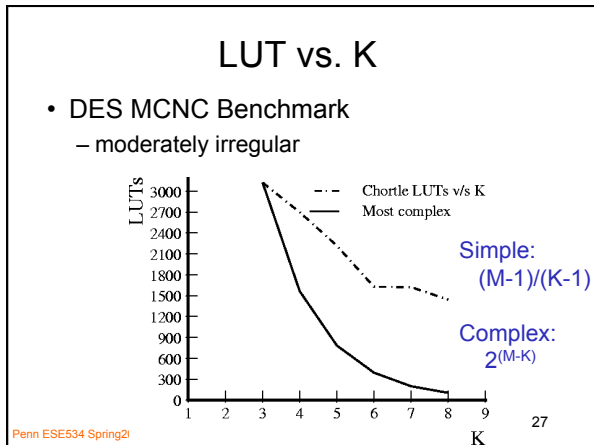
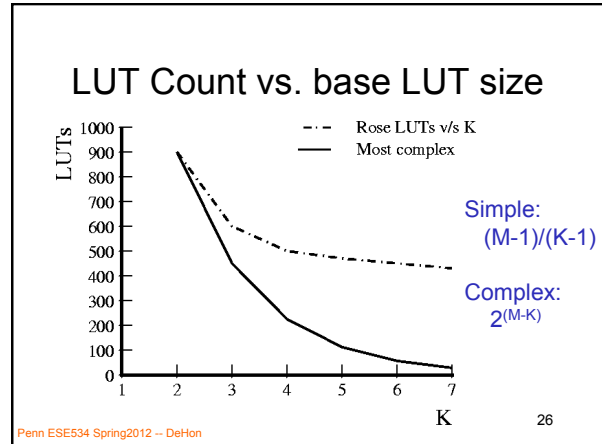
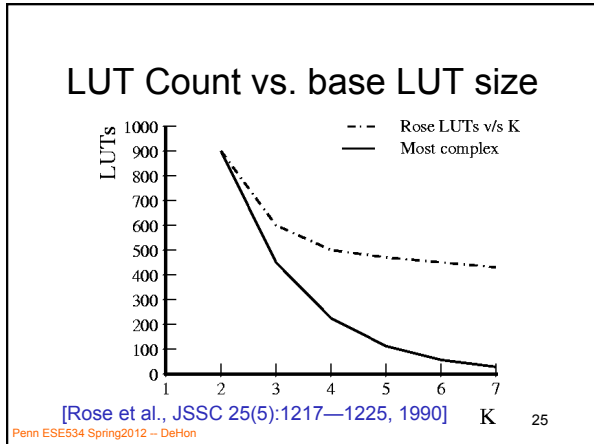
## LUT + Interconnect

- Interconnect allows us to exploit **structure** in computation
- Even if Interconnect was 99% of the area (100x logic area)
  - Would still be worth paying!
  - Add:  $N \times 2^{(2N)} \gg 2N \times (2^3 \times 128)$
  - $N=16: 16 \times 2^{32} \gg 2 \times 16 \times 2^{10} = 2^{15}$
  - $\rightarrow$  factor of  $2^{21} = 2$  Million
- Structure exploitation to avoid exponential costs is worth it!



## Different Instance of a Familiar Concept

- The most general functions are huge
- Applications exhibit **structure**
  - Typical functions not so complex
- Exploit structure to optimize “common” case

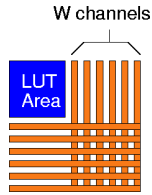


- ### Toronto Experiments
- Want to determine best K for LUTs
  - Bigger LUTs
    - handle complicated functions efficiently
    - less interconnect overhead
  - Smaller LUTs
    - handle regular functions efficiently
    - interconnect allows exploitation of compute structure
  - What's the typical complexity/structure?   
 [Rose et al., JSSC 25(5):1217–1225, 1990] 29
- Penn ESE534 Spring2012 -- DeHon

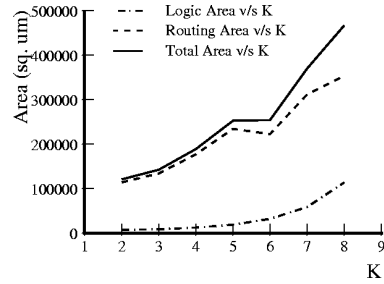
- ### Standard Systematization
- Define a design/optimization space
    - pick key parameters
    - e.g. K = number of LUT inputs
  - Build a cost model
  - Map designs
  - Look at resource costs at each point
  - Compose:
    - Logical Resources ⊕ Resource Cost
  - Look for best design points
- Penn ESE534 Spring2012 -- DeHon 30

## Toronto LUT Size

- Map to K-LUT
  - use Chortle
- Route to determine wiring tracks
  - global route
  - different channel width  $W$  for each benchmark
- Area Model for  $K$  and  $W$ 
  - $A_{lut}$  exponential in  $K$
  - Interconnect area based on switch count

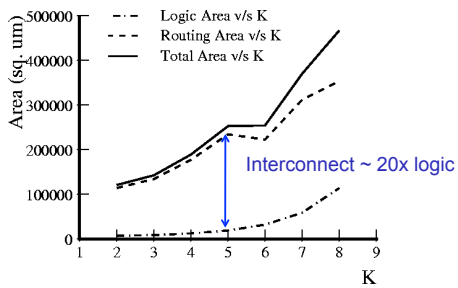


## LUT Area vs. K



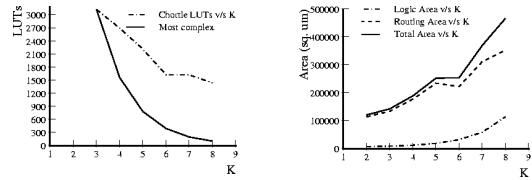
- Routing Area roughly linear in  $K$  ?

## LUT Area vs. K



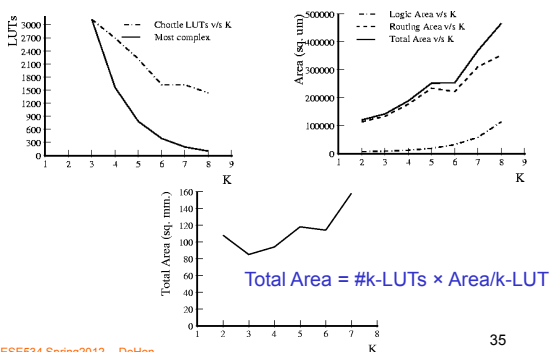
## Mapped LUT Area

- Compose Mapped LUTs and Area Model

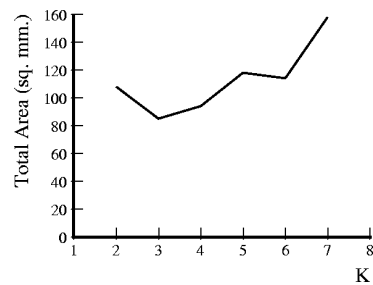


$$\text{Total Area} = \#k\text{-LUTs} \times \text{Area}/k\text{-LUT}$$

## Mapped LUT Area

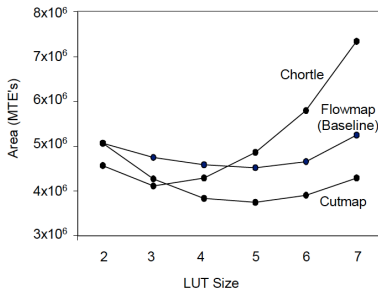


## Mapped Area vs. LUT K



N.B. unusual case minimum area at  $K=3$

## Area vs. K (different tools)



Penn ESE534 Spring2012 -- DeHon

[Yan et al., FPGA 2002]

37

## Toronto Result

- Minimum LUT Area
  - at K=4
  - robust for different switch sizes
    - (wire widths)
    - [see graphs in paper]

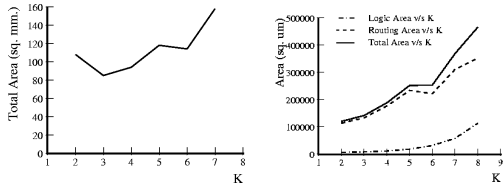
Penn ESE534 Spring2012 -- DeHon

38

## Implications

Can we make more general conclusions?

- More restricted logic functions than LUTs?



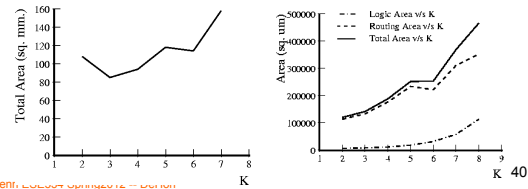
Penn ESE534 Spring2012 -- DeHon

39

## Implications (Deep)

In the range that minimizes area:

- LUT area negligible compared to interconnect
- Anything less flexible than LUT will require more interconnect



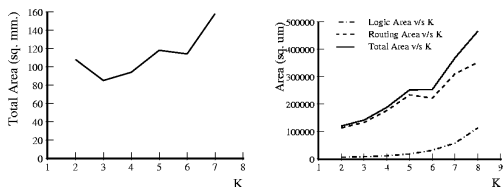
Penn ESE534 Spring2012 -- DeHon

40

## Implications

Can we make more general conclusions?

- Custom? / Gate Arrays?



Penn ESE534 Spring2012 -- DeHon

41

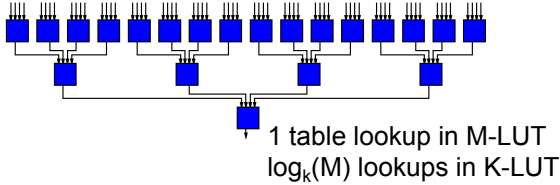
## Delay

Penn ESE534 Spring2012 -- DeHon

42

## Delay?

- Circuit Depth in LUTs?
- Lower bound?
  - (M-input fun using K-LUTs)
- “Simple Function” → M-input AND

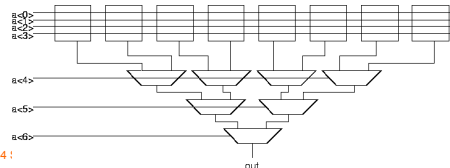


Penn ESE534 Spring2012 – DeHon

43

## Delay?

- M-input “Complex” function
  - Upper Bound:
    - use each k-lut as a  $k - \log_2(k)$  input mux
  - Upper Bound:  $\lceil (M-k)/\log_2(k - \log_2(k)) \rceil + 1$



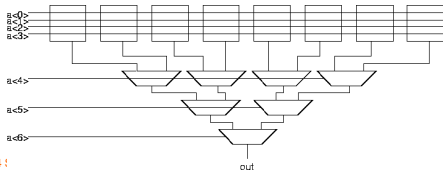
Penn ESE534 :

44

Will not cover in class, here if want to see additional details.

## Delay?

- M-input “Complex” function
  - 1 table lookup for M-LUT
  - Lower Upper bound:  $\lceil \log_k(2^{(M-k)}) \rceil + 1$
  - $\log_k(2^{(M-k)}) = (M-k)\log_k(2)$



Penn ESE534 !

45

Will not cover in class, here if want to see additional details.

## Some Math

- $Y = \log_k(2)$
- $k^Y = 2$
- $Y \log_2(k) = 1$
- $Y = 1/\log_2(k)$
- $\log_k(2) = 1/\log_2(k)$
- $(M-k)\log_k(2)$
- $(M-k)/\log_2(k)$

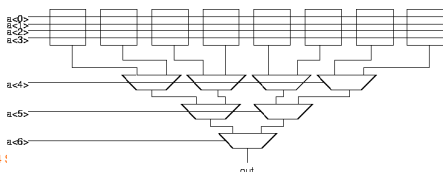
Penn ESE534 Spring2012 – DeHon

46

Will not cover in class, here if want to see additional details.

## Delay?

- M-input “Complex” function
  - Lower Upper bound:  $\lceil \log_k(2^{(M-k)}) \rceil + 1$
  - $\log_k(2^{(M-k)}) = (M-k)\log_k(2)$
  - Lower Upper Bound:  $\lceil (M-k)/\log_2(k) \rceil + 1$

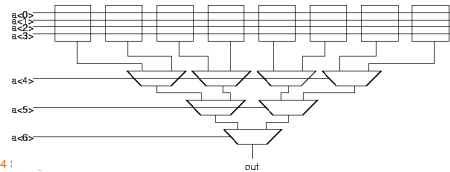


Penn ESE534 :

47

## Delay?

- M-input “Complex” function
  - 1 table lookup for M-LUT
  - between:  $\lceil (M-k)/\log_2(k) \rceil + 1$
  - and  $\lceil (M-k)/\log_2(k - \log_2(k)) \rceil + 1$



Penn ESE534 :

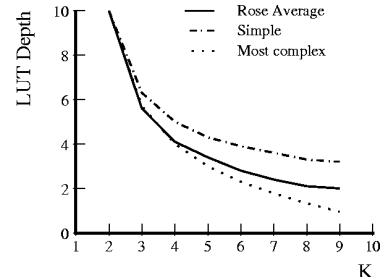
48



## Delay

- **Simple:**  $\log M$
- **Complex:** linear in  $M$
- Both scale with  $k$  as  $1/\log(k)$

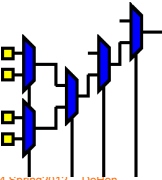
## Circuit Depth vs. K



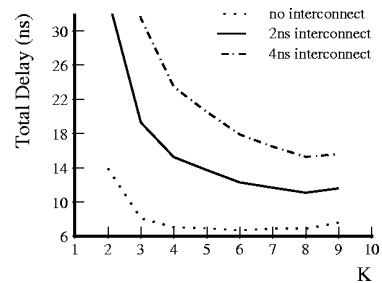
[Rose et al., JSSC 27(3):281–287, 1992] 50

## LUT Delay vs. K

- How LUT delay scale with  $k$  for small LUTs?
- Large LUTs:
  - add length term
  - $c_2 \times \sqrt{2^k}$
- Plus Wire Delay
  - $\sim \sqrt{\text{area}}$



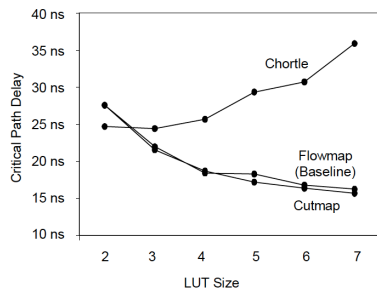
## Delay vs. K



Why not satisfied with this model?

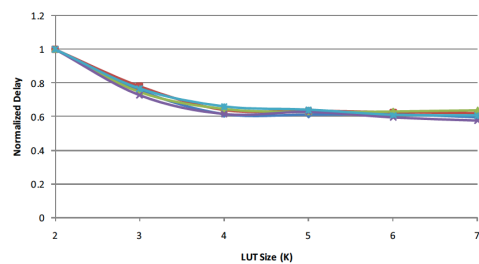
$$\text{Delay} = \text{Depth} \times (t_{\text{LUT}} + t_{\text{interconnect}})$$

## Delay vs. K (different tools)



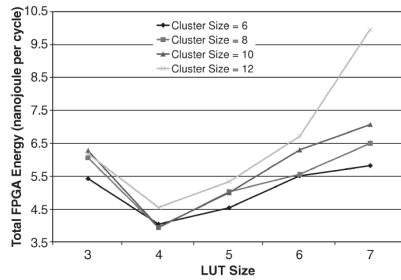
[Yan et al., FPGA 2002]

## Delay vs. K (proper critical path interconnect)



[Luu et al., FPGA 2009]

## Energy



[Li et al., TRCAD v24n11p1712 (2005)]

Penn ESE534 Spring2012 -- DeHon

55

## Observation

- General interconnect is expensive
- “Larger” logic blocks
  - ↳ fewer interconnect crossings
  - ↳ reduces interconnect delay
  - ↳ get larger
  - ↳ less area efficient
    - don't match structure in computation
  - ↳ get slower
    - Happens faster than modeled here due to area

Penn ESE534 Spring2012 -- DeHon

56

## Admin

- Reading
  - Today's: classic paper...**definitely read**
  - Wed. → no **required** reading
    - Are some suggestions
- Office hours Tuesday
  - Especially if still confused about HW6
- HW6.1-2 due on Friday

Penn ESE534 Spring2012 -- DeHon

57

## Big Ideas [MSB Ideas]

- Memory most dense programmable structure for the **most complex** functions
- Memory inefficient (scales poorly) for structured compute tasks
- **Most tasks have structure**
- Programmable interconnect allows us to exploit that structure

Penn ESE534 Spring2012 -- DeHon

58

## Big Ideas [MSB-1 Ideas]

- Area
  - LUT count decrease w/ K, but slower than exponential
  - LUT size increase w/ K
    - exponential LUT function
    - empirically linear routing area
  - Minimum area around K=4

Penn ESE534 Spring2012 -- DeHon

59

## Big Ideas [MSB-1 Ideas]

- Delay
  - LUT depth decreases with K
    - in practice closer to  $\log(K)$
  - Delay increases with K
    - small K linear + large fixed term

Penn ESE534 Spring2012 -- DeHon

60