# The Stratix™ 10 Highly Pipelined FPGA Architecture

David Lewis, Gordon Chiu, Jeffrey Chromczak, David Galloway, Ben Gamsa, Valavan Manohararajah, Ian Milton, Tim Vanderhoek, John Van Dyken

Altera (now part of Intel Corp.)

# Pipelining in the Routing: Previous Work & Assumptions

- An old idea:
  - Singh 2001; Eguro 2008
  - Older fine grained studies used pass transistor architectures
  - Reported costly and mediocre performance for registered routing;
  - Focused on retiming

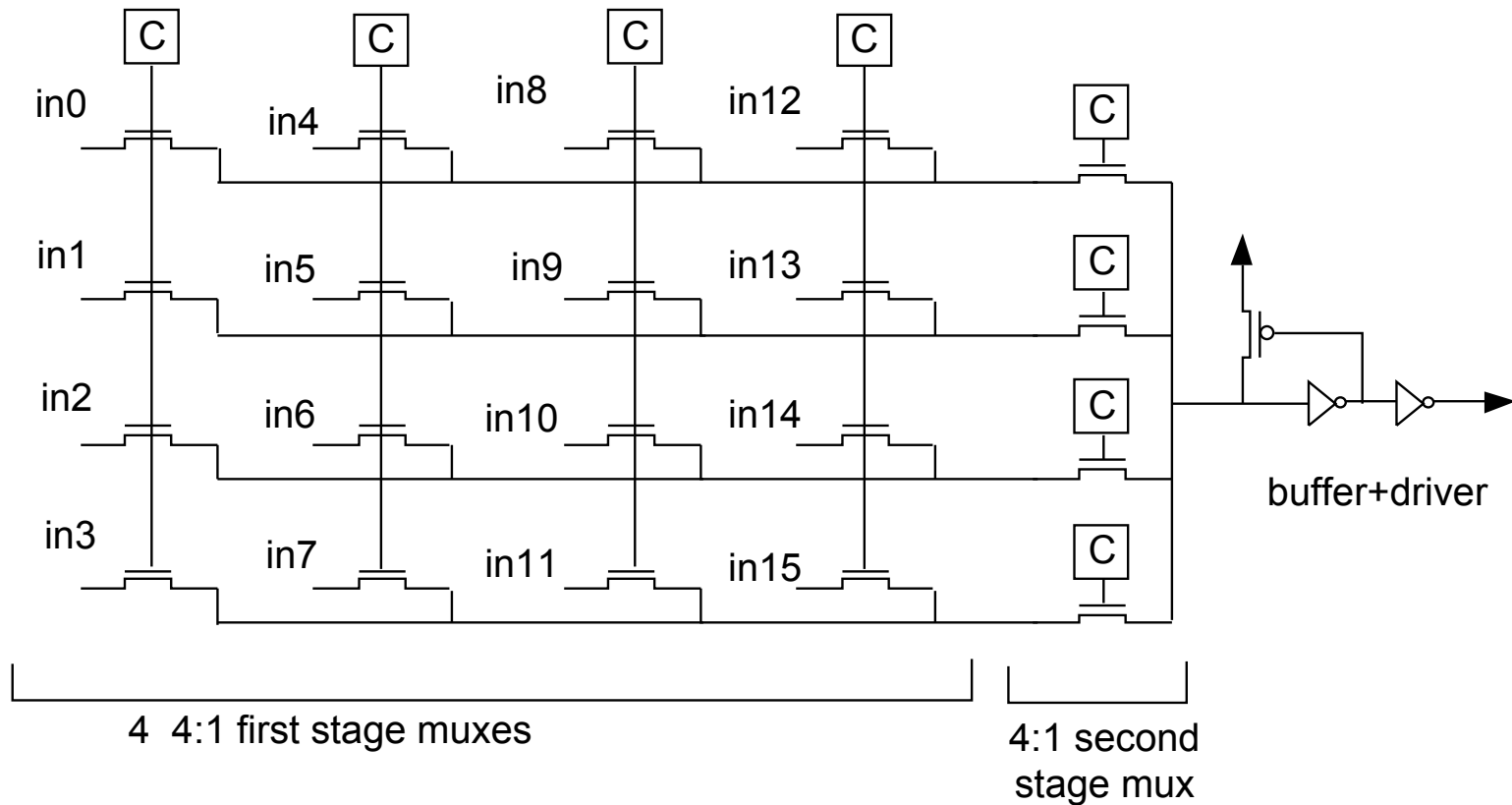- Other architectures targeted synchronous datapath designs
  - With higher level models

- Assumptions going forward:
  1. Designs are becoming more pipelined
  2. Can auto pipeline to add latency
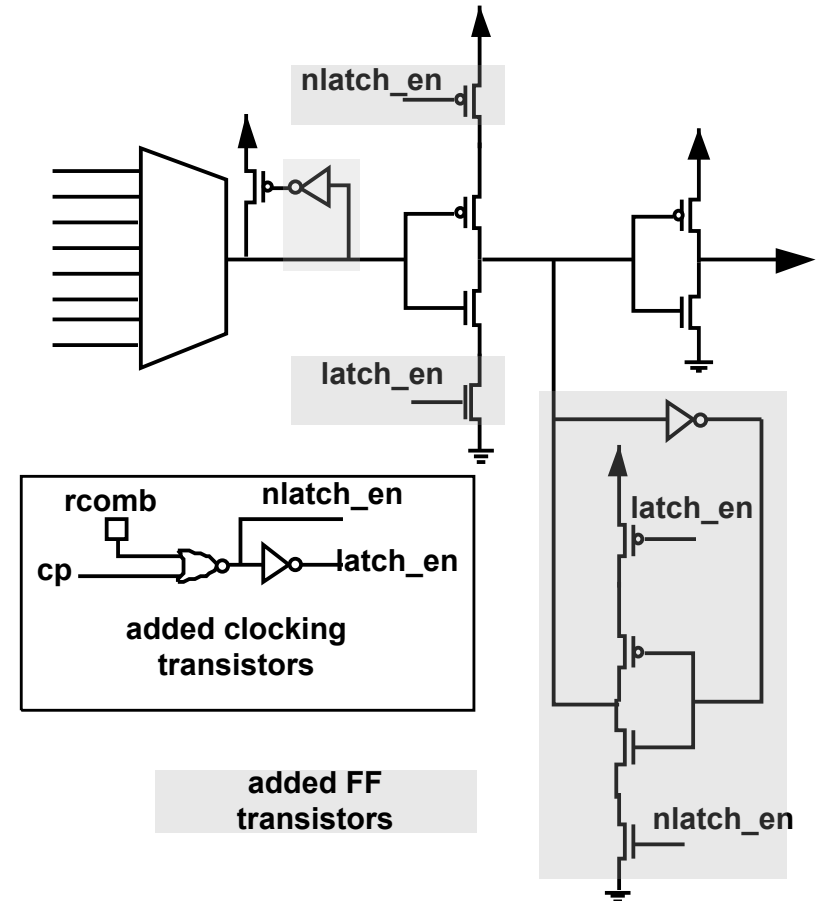  3. Designers are more willing to redesign to gain speed

FPGA 2016

# Typical Routing Multiplexers

◖ Have two stages
◖ Are followed by 1 or 2 buffers to drive the wire or LE input



4  4:1 first stage muxes
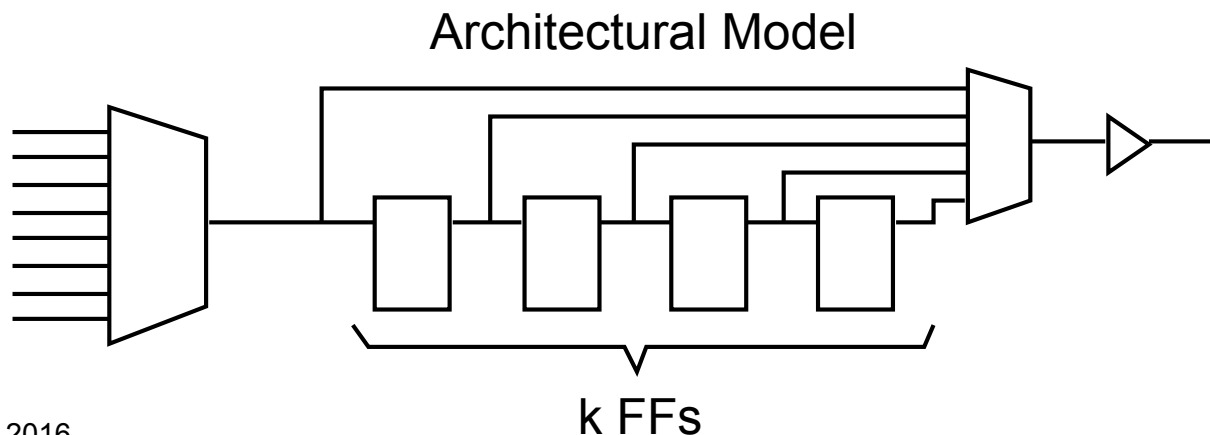
4:1 second stage mux

buffer+driver

FPGA 2016

# Key Idea: Add Pipelining to Routing Mux

- Creates a pipelined direct drive routing fabric
- Uses internal pulse latch to buffer
  - 8 minimum width and 2 larger transistors
- Minimal area cost and delay
- Latch alone is <5% soft logic area for 100% routing drivers
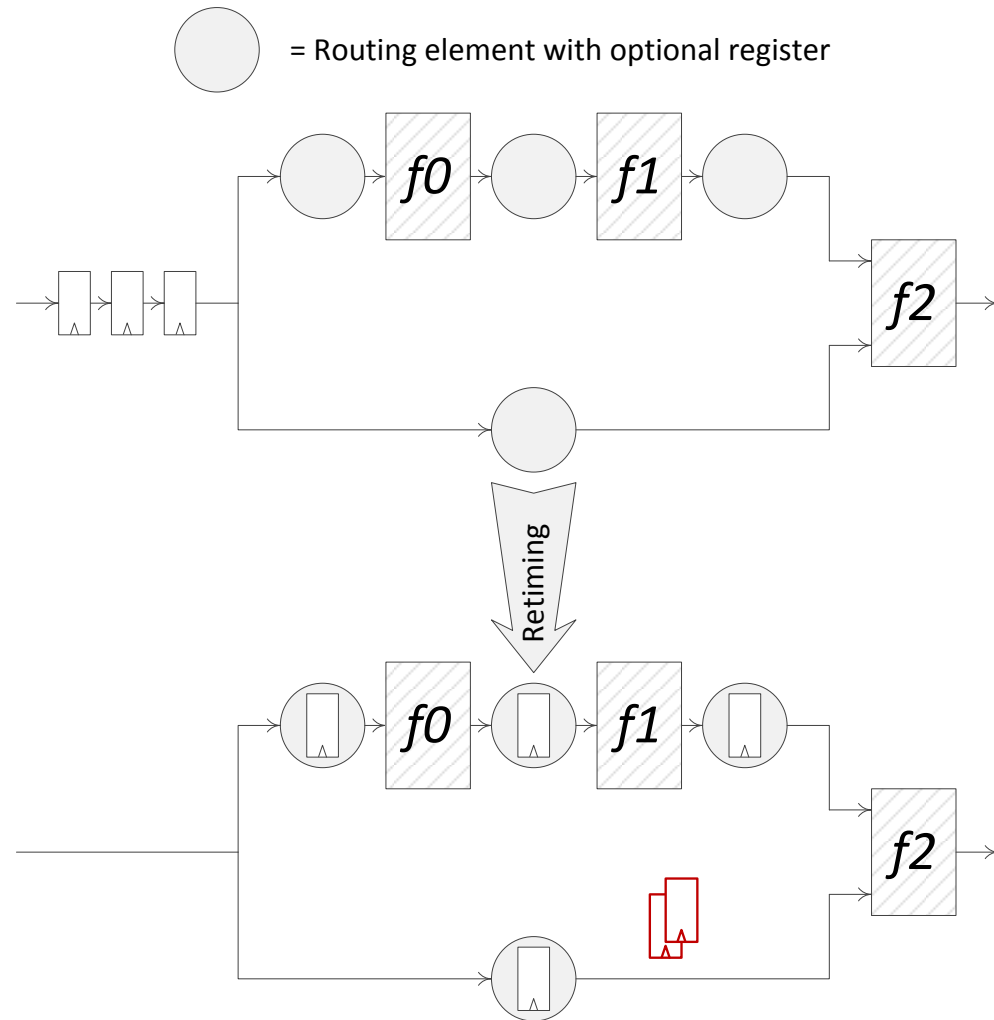- Clocking is key to cost



nlatch_en

latch_en

rcomb

nlatch_en

cp

latch_en

**added clocking transistors**

latch_en

**added FF transistors**

nlatch_en

# How Many Flip-Flops (k) Should Go After Mux?

- At least one FF ($k$=1) in each routing mux
  - May need $k > 1$ in locations where multiple signals converge to balance latency mismatches
- Built CAD to selectively enables from 0 to $k$ FF during retiming
- $k = 1$ is especially efficient implementation because can use pulse latch with flow through and no output mux
- $k > 1$ requires full edge triggered FF and bypass muxes for subsequent FFs: +10% area and increase delay per FF
- Experiments with larger $k$ are useful to establish bounds on performance, but unlikely we would build $k > 1$
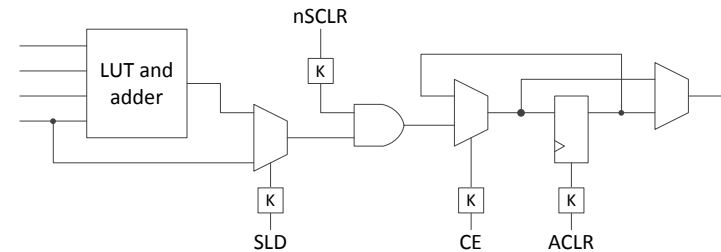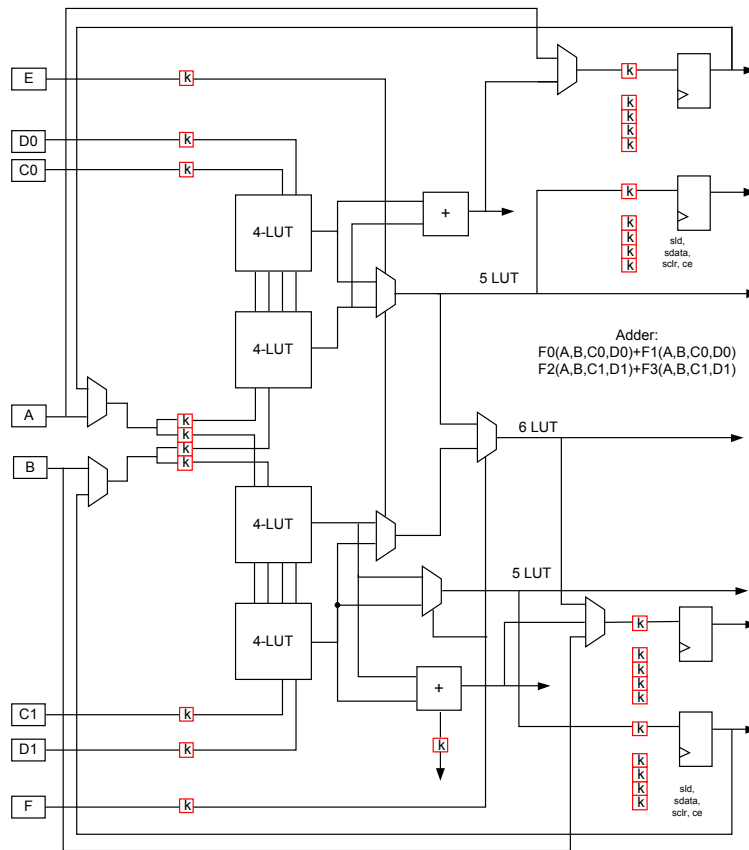
Architectural Model



k FFs

# Why K > 1?

- Retiming reconvergent paths can often lead to latency balancing problems

- Can correct minor mismatches in latency by providing k > 1 wherever signals converge



= Routing element with optional register

Retiming

FPGA 2016

# Pipelined Logic Fabric

◖ Postulate *k* FFs in front of logic element
◖ Note shared AB inputs should have separate FFs to allow independent pipelining of 5-LUTs



Adder:
F0(A,B,C0,D0)+F1(A,B,C0,D0)
F2(A,B,C1,D1)+F3(A,B,C1,D1)

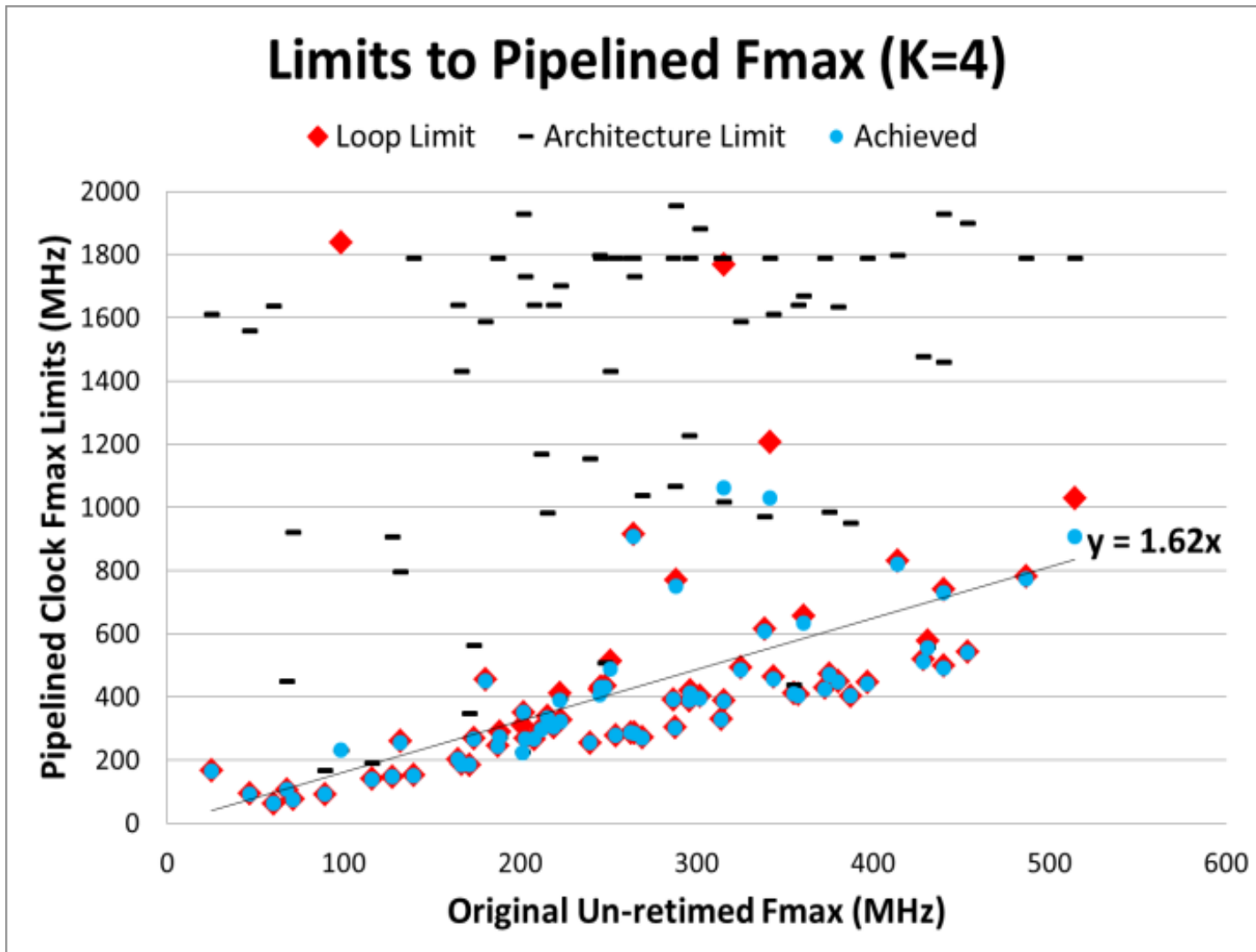FPGA 2016

# Three Different Experimental Flows

1. Retiming only: allow each FF in the user design to be retimed to anywhere there is a FF location
   - Preserve exact cycle by cycle behaviour

2. Pipelining: allow CAD flow to insert arbitrary latency in front of each clock domain and then retime
   - Preserves functional behaviour, but adds the same latency to all input to output paths
   - Note: can't ever put extra FFs in loops

3. Design modification: designer modifies the RTL to enable greater levels of pipelining while preserving functional requirements
   - Loops are critical in pipeline performance
   - Try and restructure design to minimize logic in loops + other techniques
   - Add pipelining in front of modules

# Early Experimental Conditions/Parameters

◖ Used modified Arria 10 model and approximate pipelining hardware

◖ Pipeline largest clock domain in each circuit

◖ Use $k = 4$ to approximate lots of hardware

◖ Assume no constraints on clocks available to routing FFs

◖ Measure achieved fmax after retiming / pipelining

◖ Loop limit: delay through longest loop in circuit / number of FFs in the loop

    – Can't insert FFs in a loop without breaking functionality

◖ Arch limit: longest FF location to FF location delay in any path in the circuit

FPGA 2016

# Early Limit Study Experimental Results
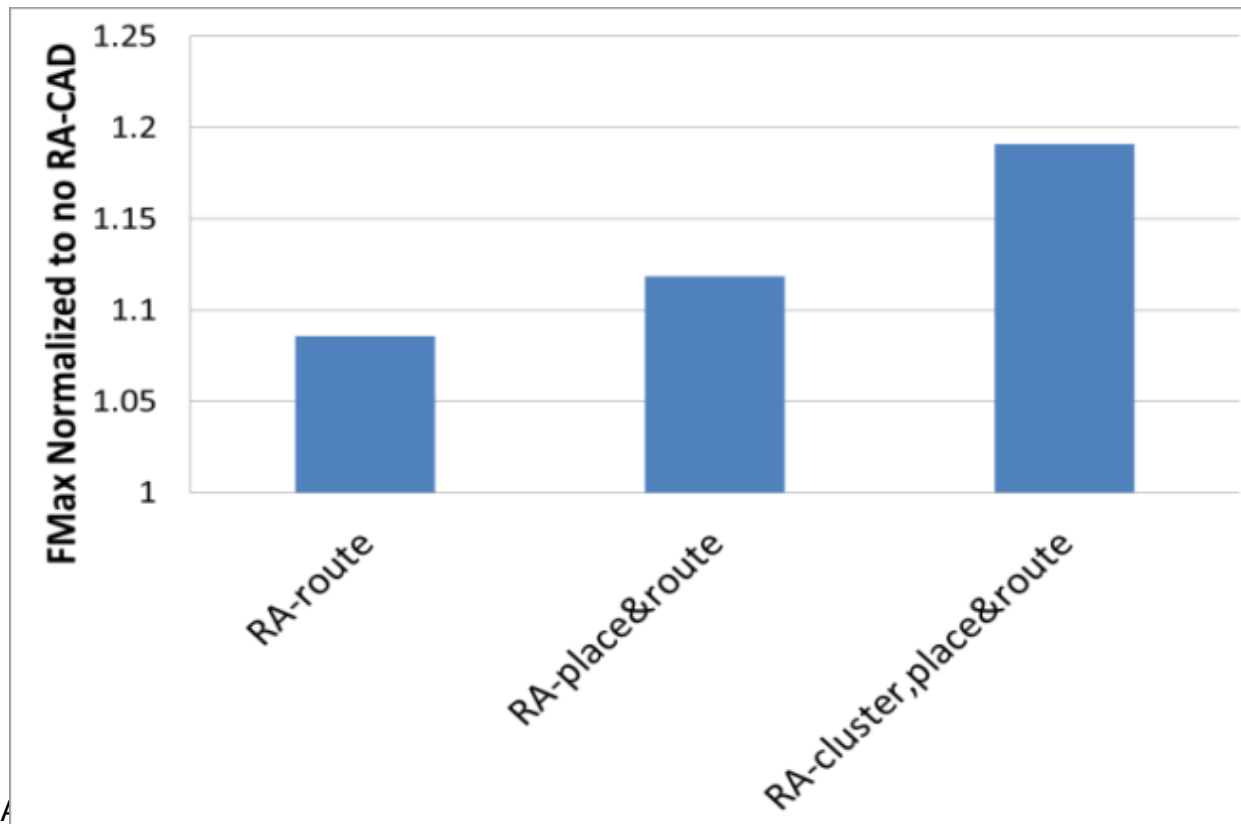
◄ Average 62% fmax increase; loop limited in most cases

## Limits to Pipelined Fmax (K=4)

◆ Loop Limit  — Architecture Limit  ● Achieved

Largest clock domain only

Y-axis: Pipelined Clock Fmax Limits (MHz)

X-axis: Original Un-retired Fmax (MHz)

y = 1.62x

# Refinements to CAD and Architecture

◖ Actual retiming done at the end of the CAD flow
  - After placement and routing

◖ Retiming-aware CAD flow uses continuous retiming based on skewed clocks at each FF
  - Skew clocks to optimize timing

◖ Paths that will be critical after retiming will have lower slack

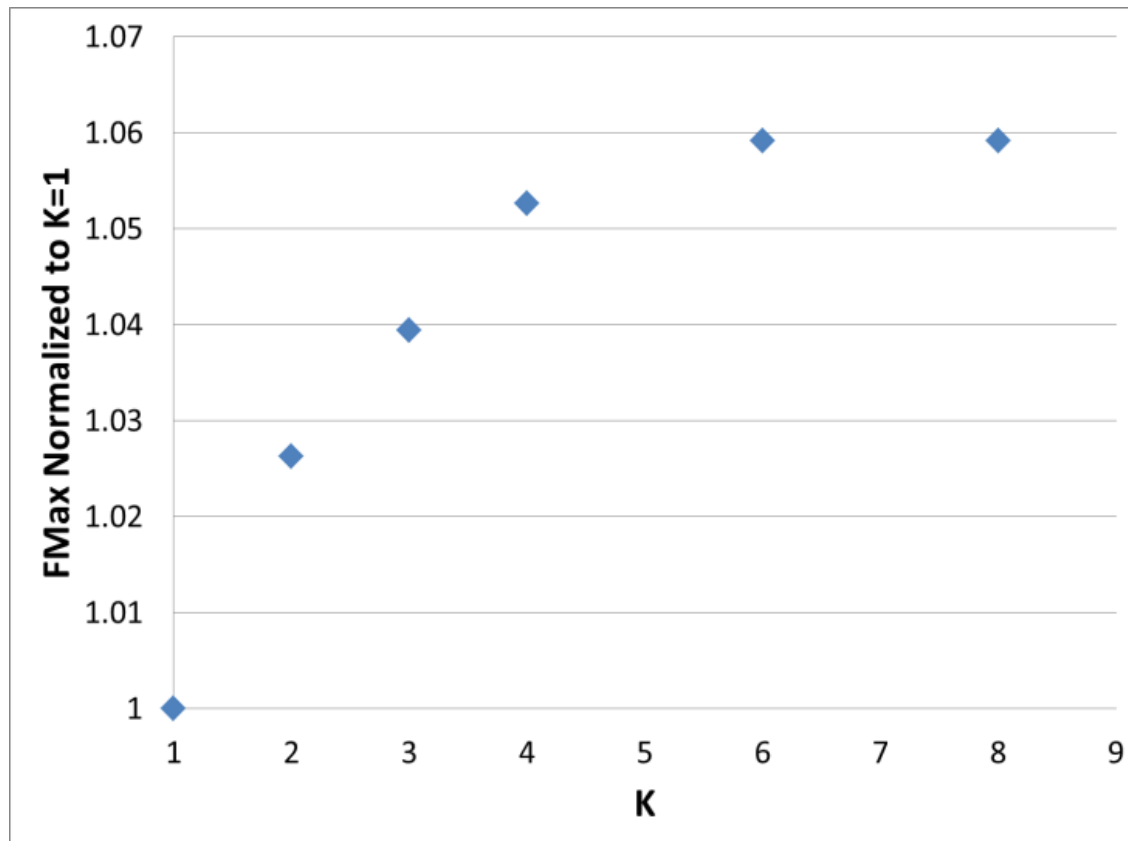◖ CAD can then target these paths for better clustering, placement, and routing

FPGA 2016

# Impact of Retiming Aware CAD

◖ Up to 19% fmax improvement by focusing on paths that are critical for retiming

◖ Critical paths are generally those in loops

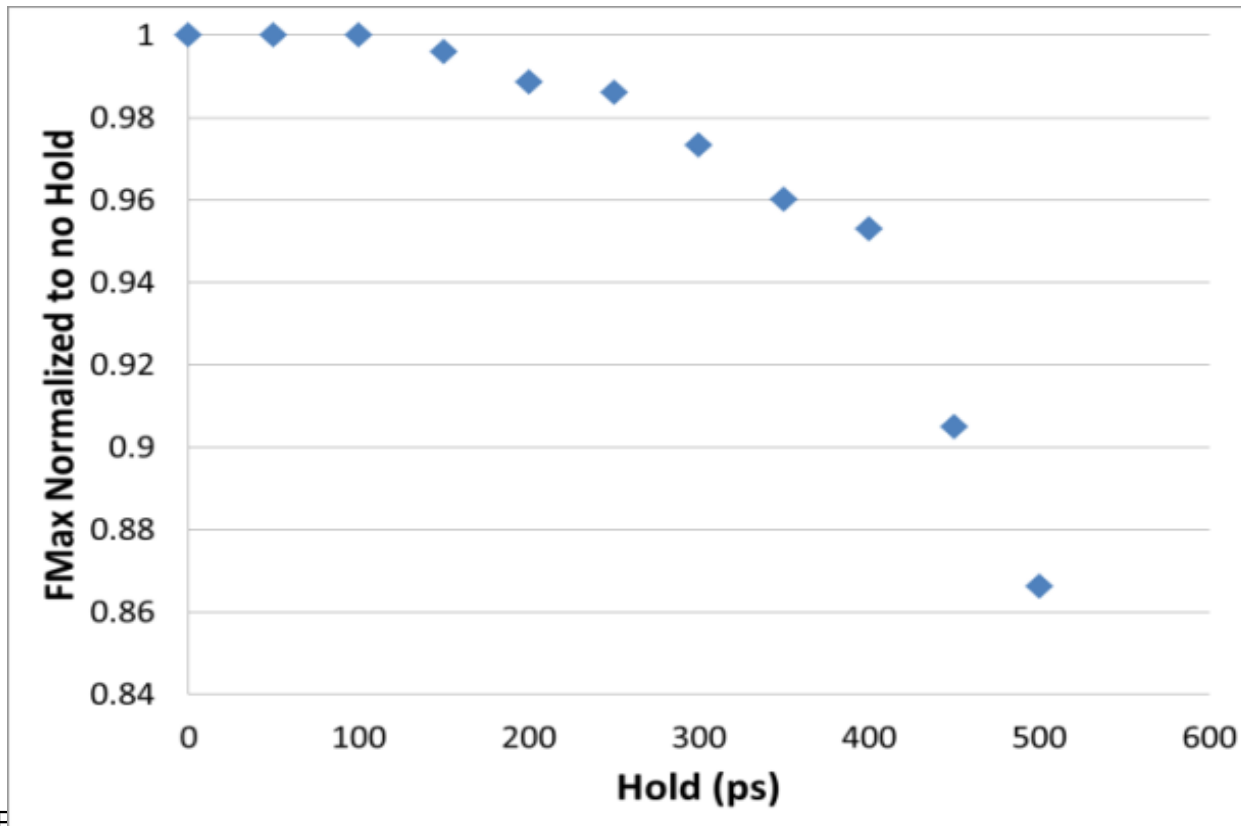– pure feedforward circuits can generally be deeply pipelined



FPGA

# Effect of Number of FFs in Routing Mux: k

◖ About 5% fmax loss by dropping *k* to 1,
  – but much less than area cost of *k* = 4
  – Normalized to *k* = 1 result; we built in +5% fmax in original experiment

# Hold Time Issue with Pulsed Latch

◖ Consequence of low-cost pulse latch: needs some hold time, else data can race through consecutive latches

◖ Don't know exact value of hold time during early architecture experiments, but guess ~200ps
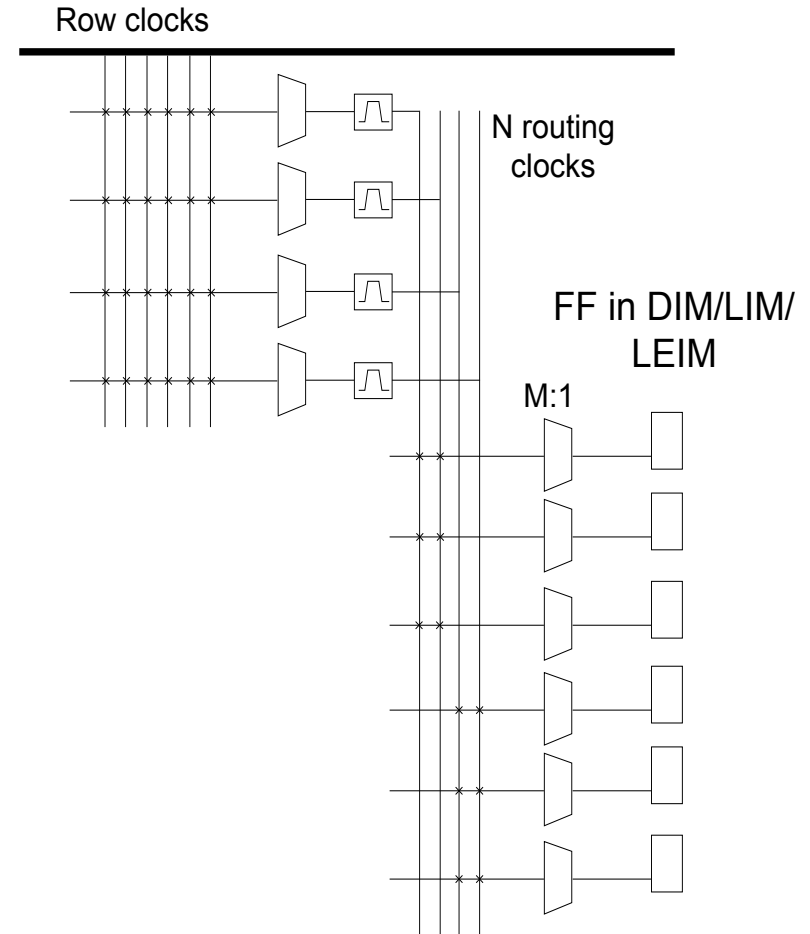
# Clocking

- Prior academic work largely ignores clocking
- Real customer design set contains an average of 14 clocks per design, up to 67
- Stratix 10 global clock architecture is routable for better timing properties, but no changes relevant to the pipelined fabric
  - See Ebeling FPGA 2016 for details
- 6 clock lines available to provide global clocks to each LAB
- Approximately 160 routing mux FFs per LAB
  - Plus the 80 inside the LAB logic
- A 6:1 mux per latch would be too large/expensive
  - Several times larger than the FF
- Since many FF clock muxes, desire to minimize their size

FPGA 2016

# Clocking Mux Architecture

- Divide the FFs into groups and make them share clocks
  - Ex: all short wires going left/right; all long wires going up/down, etc.
- Pick 1 or 2 clocks per group from the 6 available
- Each FF selects from those clocks, using a 2:1 mux or no mux at all
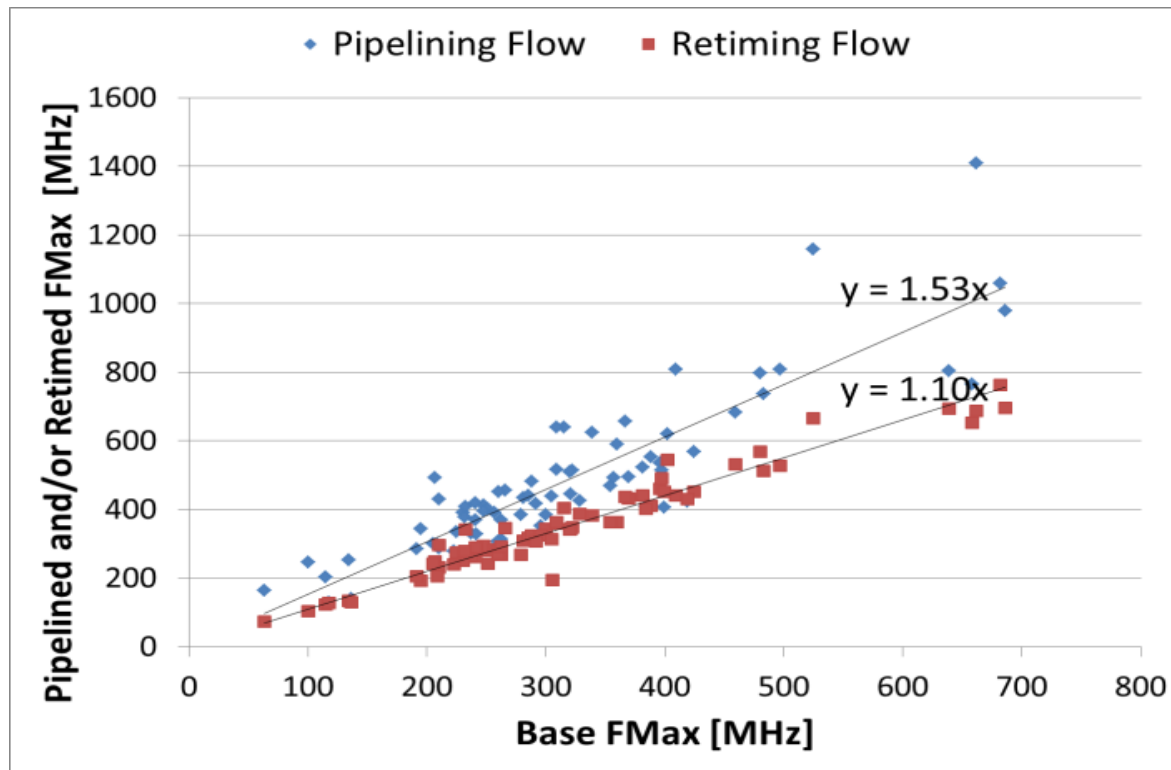- Carefully tuning the groups resulted in **no** Fmax loss

Row clocks

N routing clocks

FF in DIM/LIM/ LEIM

M:1

# FPGA Architect's First Law of Entropy:

You have to run just to stay in place

# Results Accounting for Realities

◖ With real CAD, $k = 1$, hold time, all domain fmax

◖ Retiming only: +10% fmax; pipelining: +53% fmax
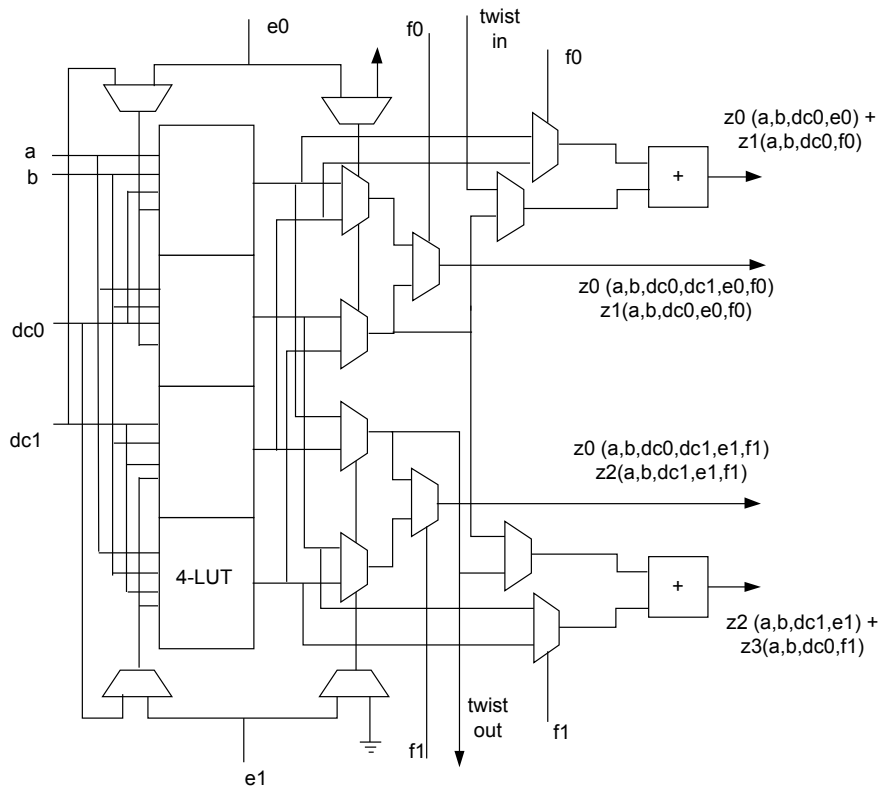
  – Small domains have less benefit from pipelining



FPGA 2016

# Other Architecture Changes - Logic Element Modification

- Stratix II to V have shared LUT mask (SLM)
- All provide 2 5-LUT with 8 inputs → 2 shared inputs, 3 unique
- SLM can build 2 6-LUT with identical functions, and 4 shared inputs
- Difficult for pipelining because internal stages of the LUT are used for two different logical functions
  - Can't independently retime
- Removed SLM
- Also removed complicated arithmetic (3 input adder and use of all 8 inputs)
- Push back part of adder into LUTs; simplify adder hardware
- Converted asynchronous clear and synchronous clear into 2 general purpose clears
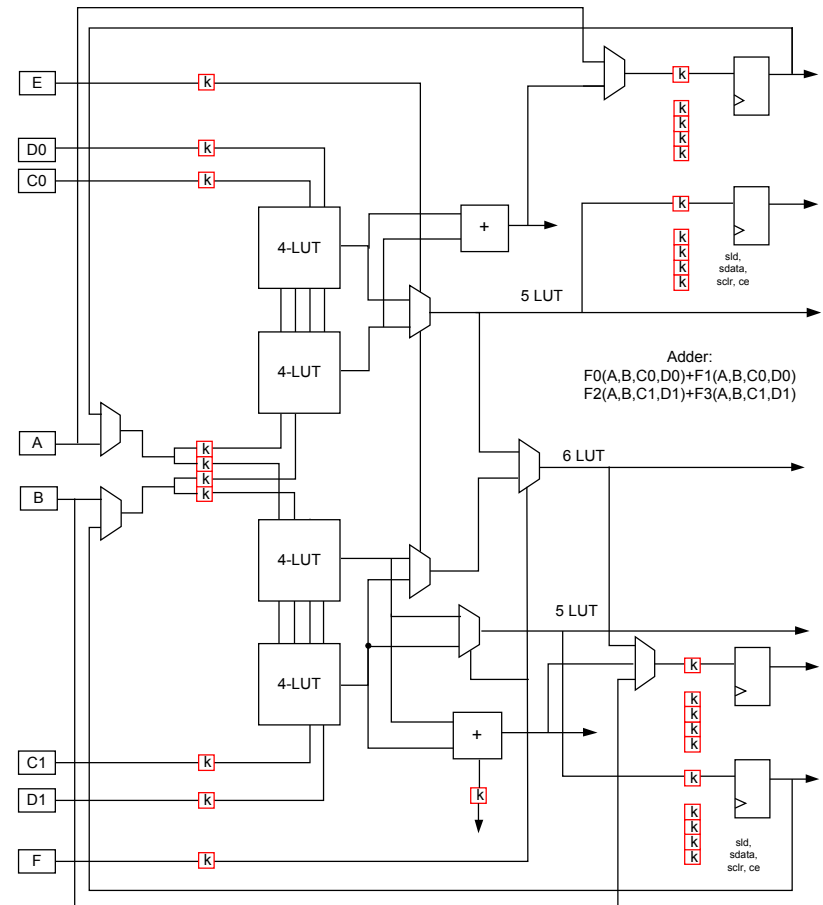- Synchronous load is now static only

# Simplified ALM

◄ Simpler, at least in comparison to previous, small fmax win
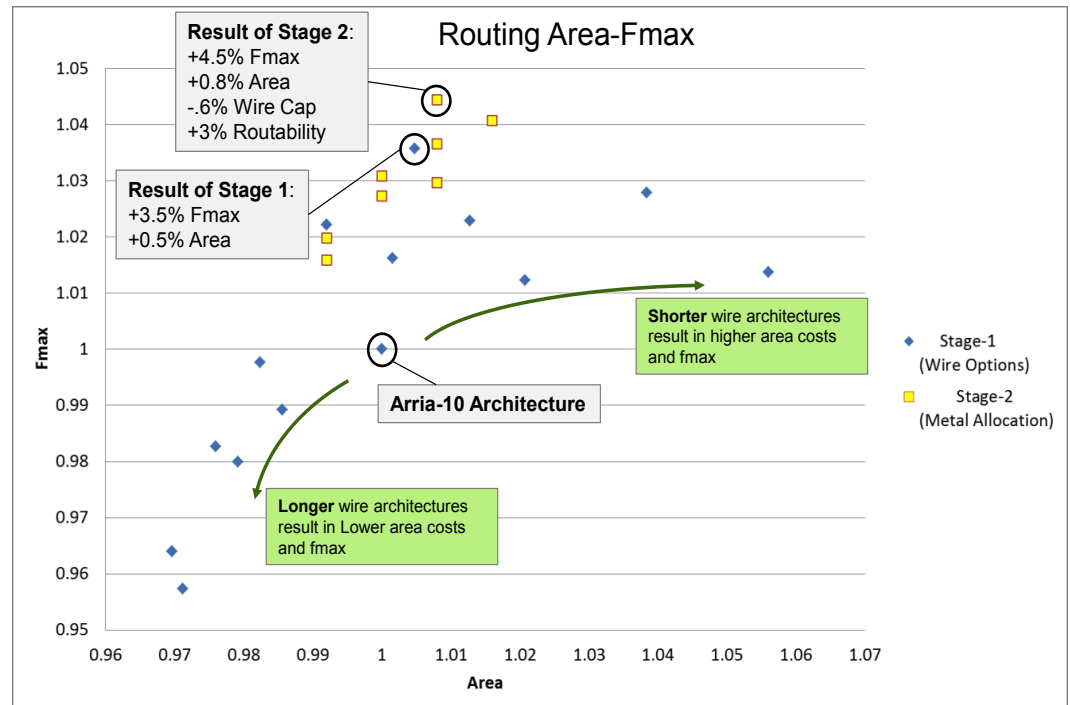


previous Stratix

Stratix 10

z0 (a,b,dc0,e0) + z1(a,b,dc0,f0)

z0 (a,b,dc0,dc1,e0,f0) z1(a,b,dc0,e0,f0)

z0 (a,b,dc0,dc1,e1,f1) z2(a,b,dc1,e1,f1)

z2 (a,b,dc1,e1) + z3(a,b,dc0,f1)

Adder:
F0(A,B,C0,D0)+F1(A,B,C0,D0)
F2(A,B,C1,D1)+F3(A,B,C1,D1)

# Routing Optimization

- Larger variety of metal layers in Intel vs. TSMC
- 2 stage investigation of wire lengths and allocation to metal layers: +4.5% fmax at +0.8% area

| H wire | H wire fraction | V wire | V wire fraction |
|--------|-----------------|--------|-----------------|
| H2 | 10% | V2 | 12% |
| H4 | 28% | V3 | 38% |
| H10 | 50% | V4 | 33% |
| H24 | 12% | V16 | 17% |



Routing Area-Fmax

Result of Stage 2:
+4.5% Fmax
+0.8% Area
-.6% Wire Cap
+3% Routability

Result of Stage 1:
+3.5% Fmax
+0.5% Area

Arria-10 Architecture

Shorter wire architectures result in higher area costs and fmax

Longer wire architectures result in Lower area costs and fmax

Stage-1 (Wire Options)

Stage-2 (Metal Allocation)

# How User Designs were Modified in Experiments

- Altera design expert worked with several customers
- Modified their designs to enable more pipelining and preserve functional requirements
- Full designs, not isolated cores
- Primarily reduced paths through loops
- Shannon decomposition of critical loop state
- Moving some computation that can be precomputed earlier in the pipeline
- Loop unrolling; works better with S 10 than previous
- Did back port design changes to Stratix V to measure architecture + design benefit

# Results

- Redesign is 2.4X faster than S V
- Back port to S V is 1.23X faster
- Architecture + redesign + process is 92% faster than optimized design in S V

| Module | S V (MHz) | Retime S10 | Pipe S10 | Redesign S10 | Redesign SV |
|--------|-----------|------------|----------|--------------|-------------|
| 1A | 320 | 380 (+19%) | 460 (+44%) | 489 (+53%) | 329 (+3%) |
| 1B | 327 | 482 (+47%) | 626 (+91%) | | |
| 1C | 319 | 432 (+35%) | 457 (+43%) | | |
| 2A | 250 | 429 (+72%) | 454 (+82%) | 942 (+277%) | 347 (+39%) |
| 3A | 191 | 290 (+52%) | 291 (+52%) | 748 (+292%) | 359 (+88%) |
| 4A | 403 | 599 (+49%) | 638 (+58%) | 725 (+80%) | 411 (+2%) |
| 4B | 384 | 555 (+45%) | 570 (+48%) | 695 (+81%) | 391 (+2%) |
| Geo % | | 45% | 59% (*) | 136% | 23% |

(*) typo in paper: 49% should be 59%

# Quartus Enhancements for Stratix 10

◖ Constraint-based retiming to solve for minimum clock period

**Design optimization advisor:**

◖ Where should asynchronous resets be converted into synchronous?

◖ Where should pipeline registers be added to enable deeper pipelining?

◖ Which loops limit the performance?

◖ Use set of infeasible constraints from retimer to show where the conversions need to be done

◖ Can automatically modify the retiming graph to model the proposed change and report on potential fmax

# Quartus Enhancements for Stratix 10



Speculative modifications and Fmax predictions for each

Limits to performance for particular (speculative) step

# Conclusions

◖ Direct drive routing enables a very low cost FF in each routing mux

◖ Most of the cost is elsewhere
  - FFs in the logic
  - LAB level clock muxing

◖ Put one FF location everywhere you can think of, and carefully optimize the clock muxing

◖ Subtle interactions between pipelining and internals of a complicated logic element

◖ Optimized designs 92% faster in S 10 than S V

◖ Quartus support to help designers identify where to focus