# Pipeline Frequency Boosting: Hiding Dual-Ported Block RAM Latency using Intentional Clock Skew

Alexander Brant
alexb@ece.ubc.ca

Ameer Abdelhadi
ameer@ece.ubc.ca

Aaron Severance
aaronsev@ece.ubc.ca

Guy G.F. Lemieux
lemieux@ece.ubc.ca

Department of Electrical and Computer Engineering
University of British Columbia
Vancouver, Canada

## Objectives

Hide dual-ported block RAM access latencies without additional pipeline stages or architectural changes.
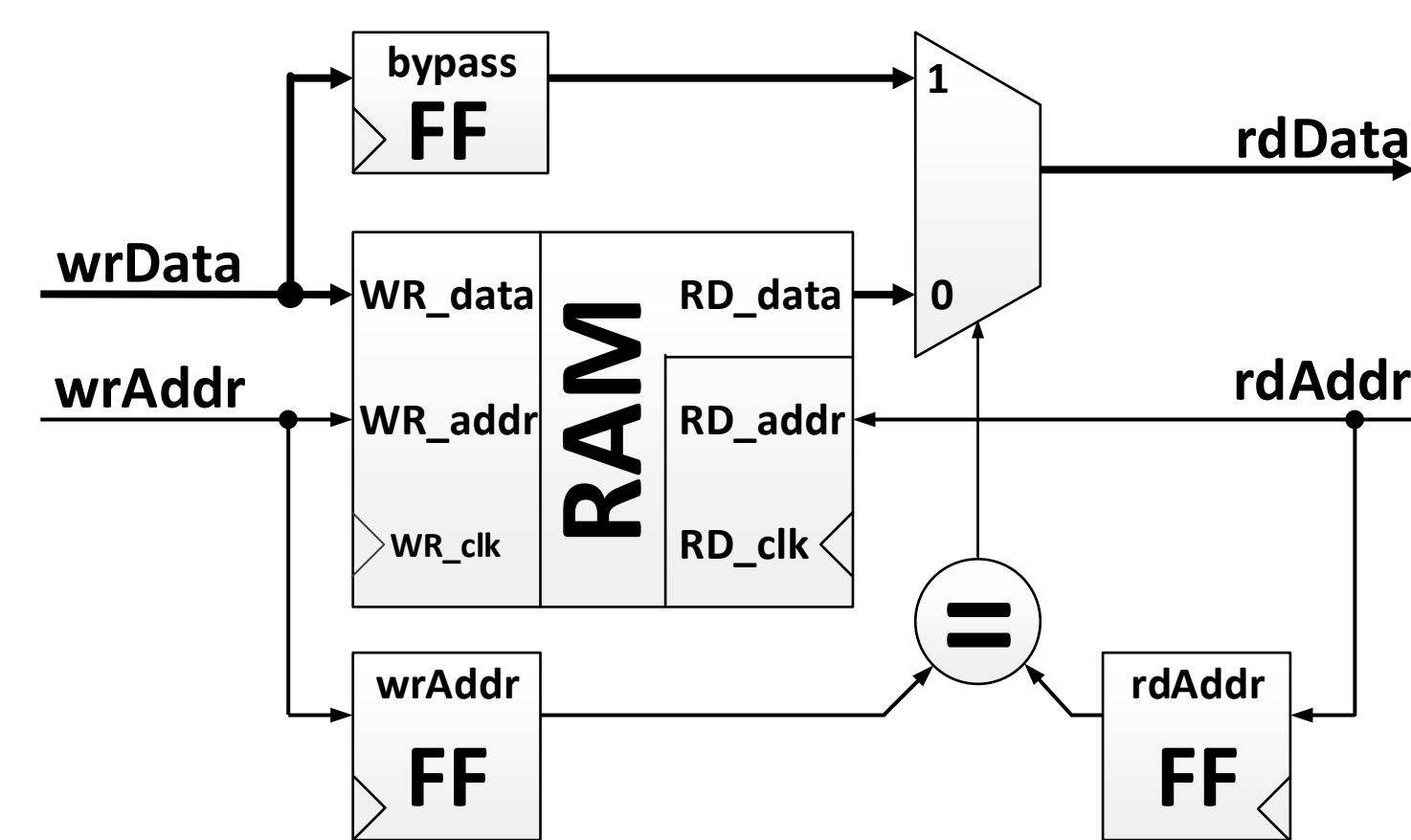
## Method

Clock skewing is employed to effectively eliminate the read and write latency of memories, while preserving functionality, and using fewer resources than conventional bypass designs.
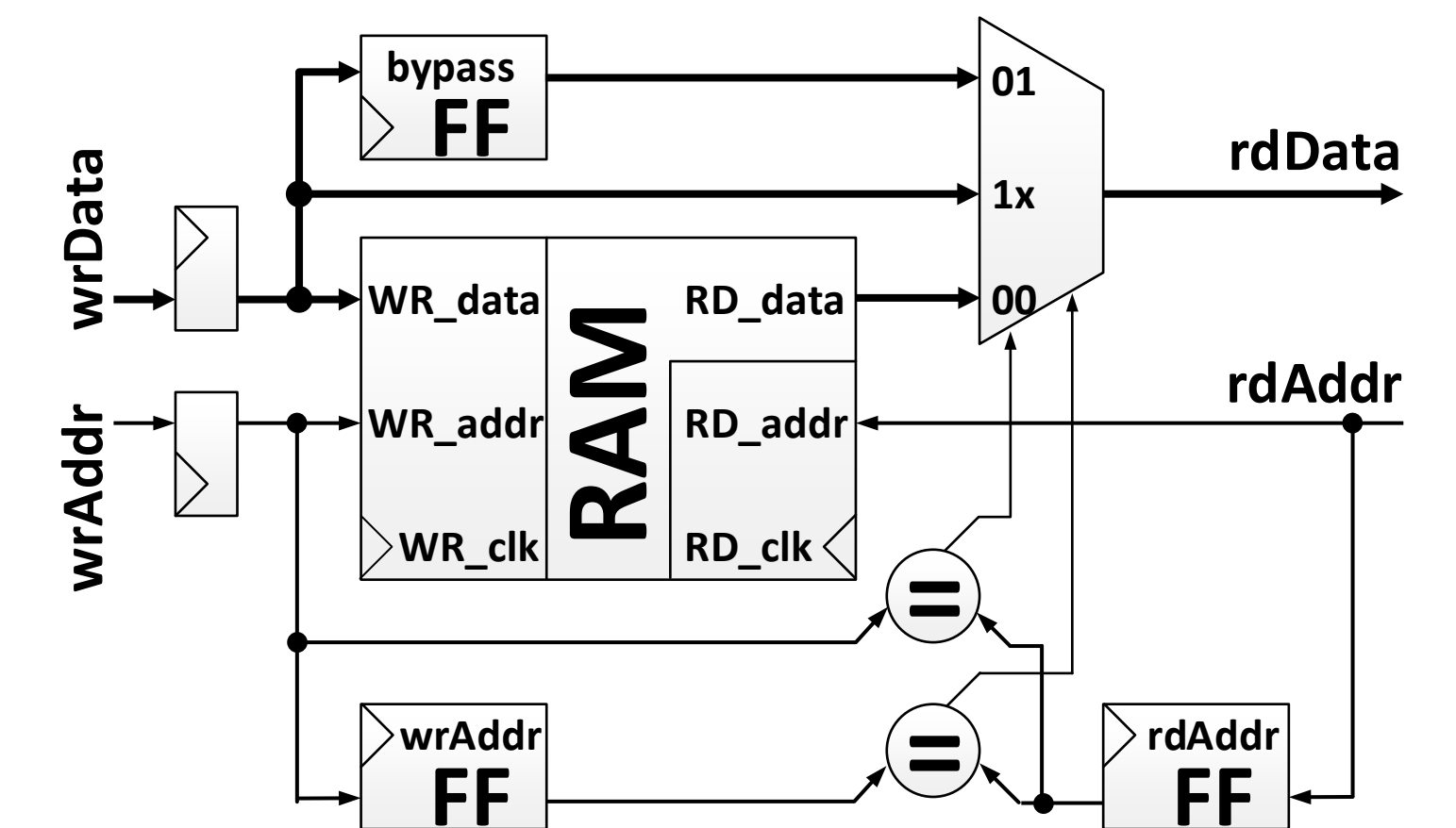
## Pipeline Bypassing/Forwarding

Written data is passed forward through a bypass register, skipping any additional cycles incurred by RAM writing process.

### Single-Stage Bypass



Load bypass register if read-after-write
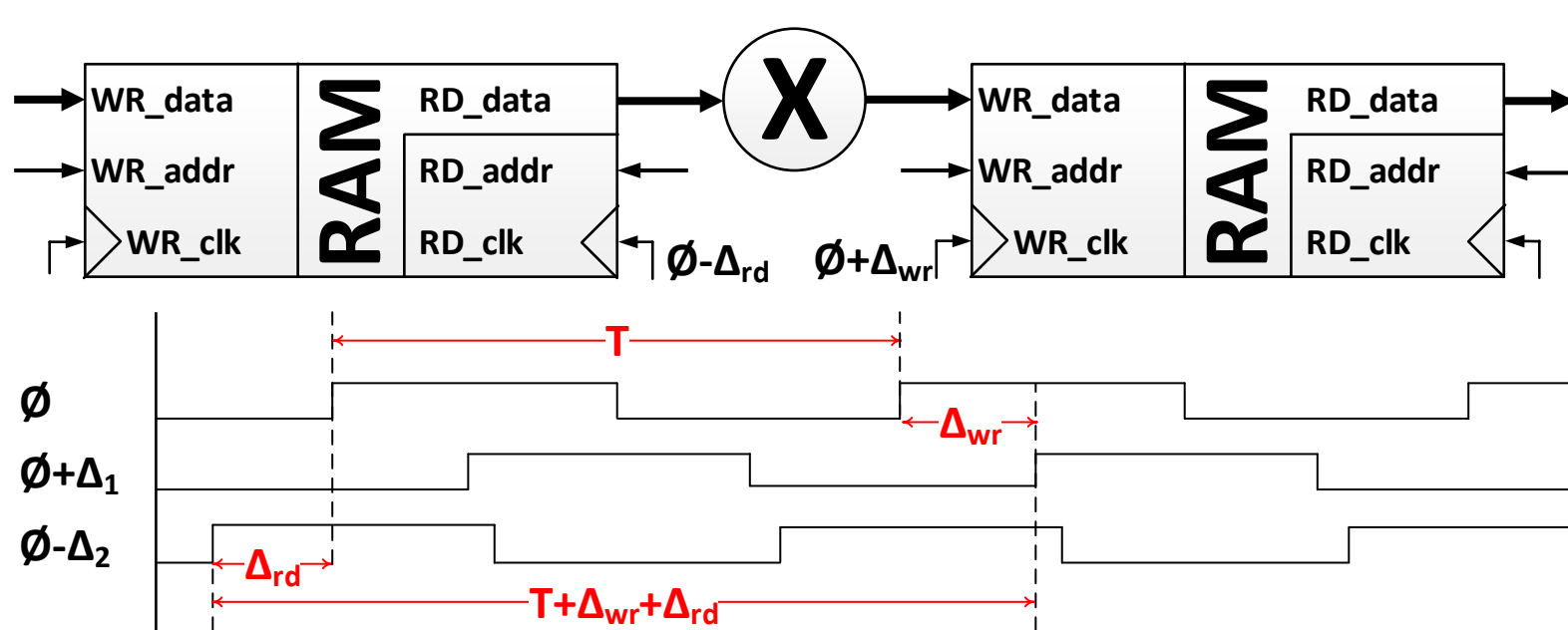
### Fully Pipelined Bypass



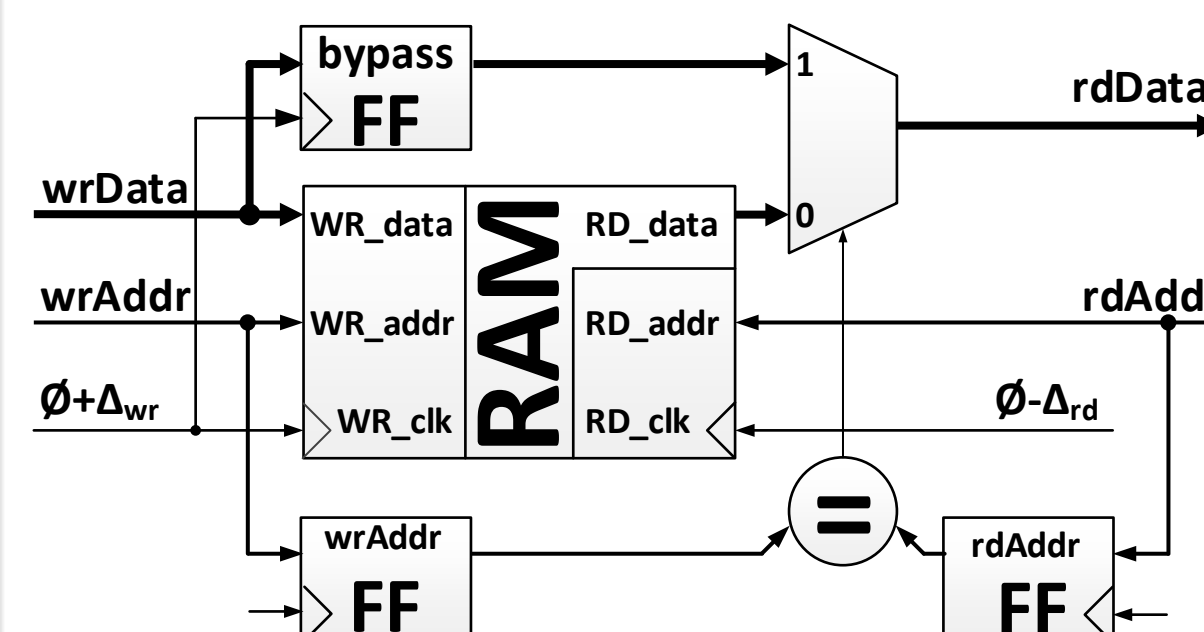Earlier read-after-write is also bypassed

## Clock Skew Scheduling

Useful or intentional skewing introduces skews to clocked elements allowing critical paths longer periods.
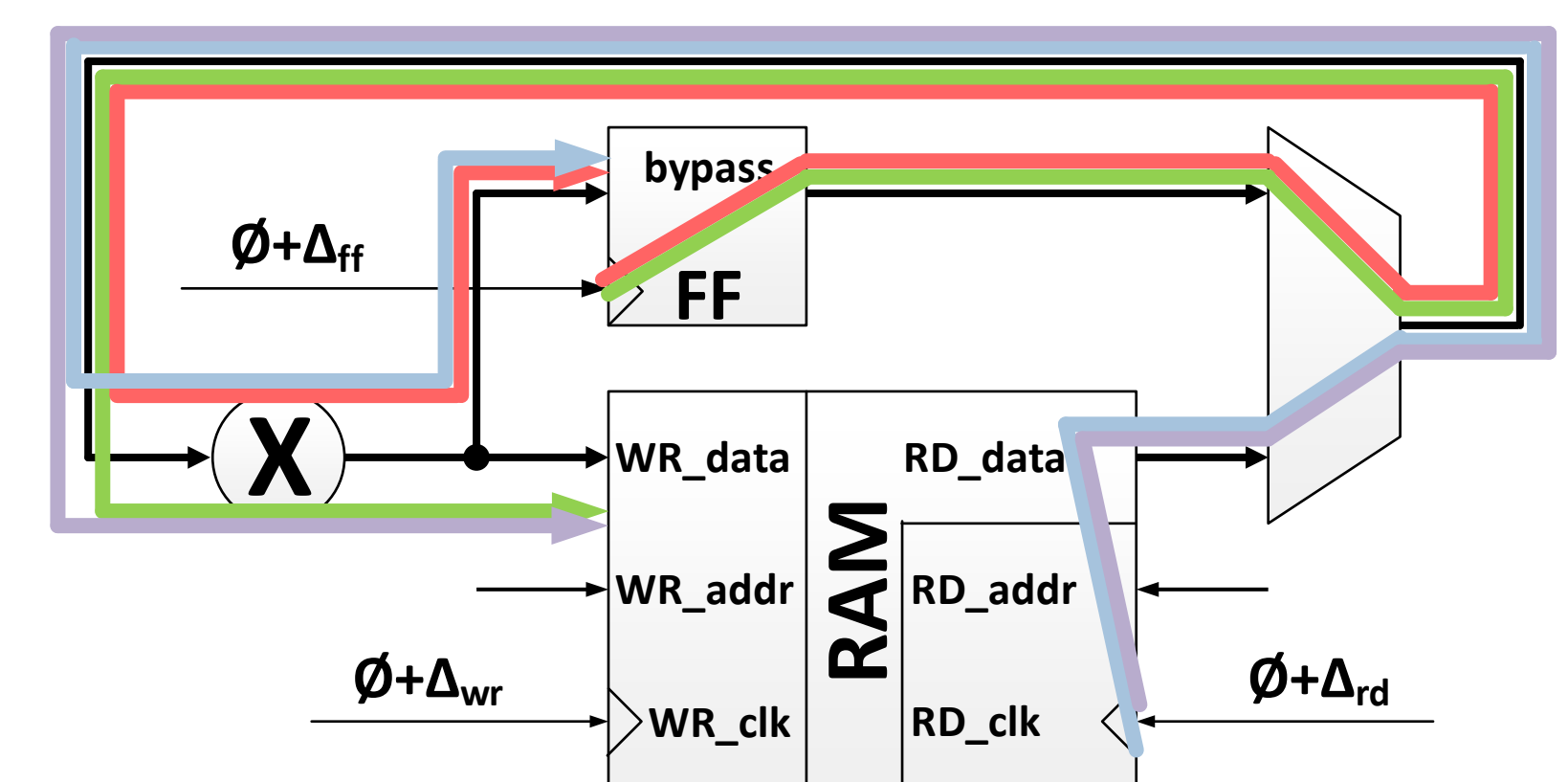
### Pipelined Block RAM



The write port is clocked late and the read port is clocked early to provide the processing logic a wider clock period.

### Single-Stage Bypass Block RAM



Clock scheduling for a Single-stage Bypass: The write lags by $\Delta_{wr}$ and the read leads by $\Delta_{rd}$.

## Timing Paths Analysis



$T \geq \Delta_{ff} + t_{c \to o}(ff) + t_d(ff \to mux) + t_d(mux) + t_d(mul) + t_d(mul \to ff) + t_{su}(ff) - \Delta_{ff}$

$T \geq \Delta_{ff} + t_{c \to o}(ff) + t_d(ff \to mux) + t_d(mux) + t_d(mul) + t_d(mul \to ram_{wr}) + t_{su}(ram_{wr}) - \Delta_{wr}$

$T \geq \Delta_{rd} + t_{c \to o}(ram_{rd}) + t_d(ram_{rd} \to mux) + t_d(mux) + t_d(mul) + t_d(mul \to ff) + t_{su}(ff) - \Delta_{ff}$

$T \geq \Delta_{rd} + t_{c \to o}(ram_{rd}) + t_d(ram_{rd} \to mux) + t_d(mux) + t_d(mul) + t_d(mul \to ram_{wr}) + t_{su}(ram_{wr}) - \Delta_{wr}$

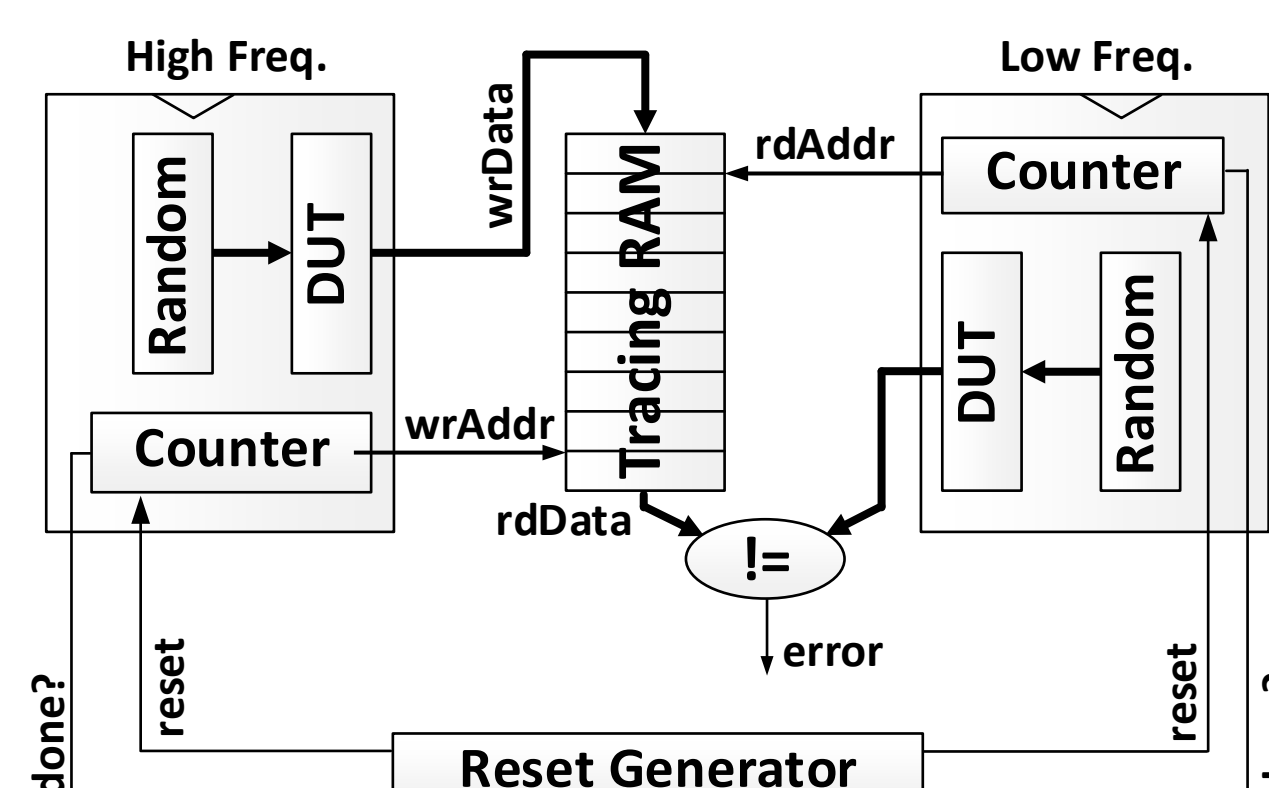For minimal clock period, $\Delta_{wr}$ and $\Delta_{rd}$ are

$$\Delta_{wr} = t_d(mul \to ram_{wr}) + t_{su}(ram_{wr}) - t_d(mul \to ff) - t_{su}(ff) + \Delta_{ff}$$

$$\Delta_{rd} = t_{c \to o}(ff) + t_d(ff \to mux) - t_{c \to o}(ram_{rd}) - t_d(ram_{rd} \to mux) + \Delta_{ff}$$

## Performance Testing Circuit

- Altera/Terasic DE2-115 board with a Cyclone IV-E
- Static timing analysis is derived using QuartusII
- A dual-clocked tracing RAM is written by a DUT at high freq. and read by a reference at low freq.
- Randomly initialized RAMs and random addresses
- Clock enables were avoided for high performance



## Experimental Results

| Design Method | | Fmax (MHz) (QuartusII STA) | Fmax (MHz) (Tested) | #LUTs | #FFs |
|---|---|---|---|---|---|
| Fully Pipelined bypass | | 179 | 208 | 37 | 34 |
| Single-Stage bypass | $\Delta_{wr} = \Delta_{rd} = 0$ | 113 | 153 | 21 | 17 |
| | Optimal $\Delta_{wr}, \Delta_{rd}$ | 183 $\Delta_{wr}=3.9ns; \Delta_{rd}=0$ | 238 $\Delta_{wr}=1.9ns; \Delta_{rd}=0$ | 21 | 17 |

A pipeline with 16-bit data width and 8-bit address width

50% less FFs    43% less LUTs    14% performance improvement