# Don't Forget the Memory: Automatic Block RAM Modelling, Optimization, and Architecture Exploration

**S. Yazdanshenas,  K. Tatsumura[*], and V. Betz**

**University of Toronto, Canada**
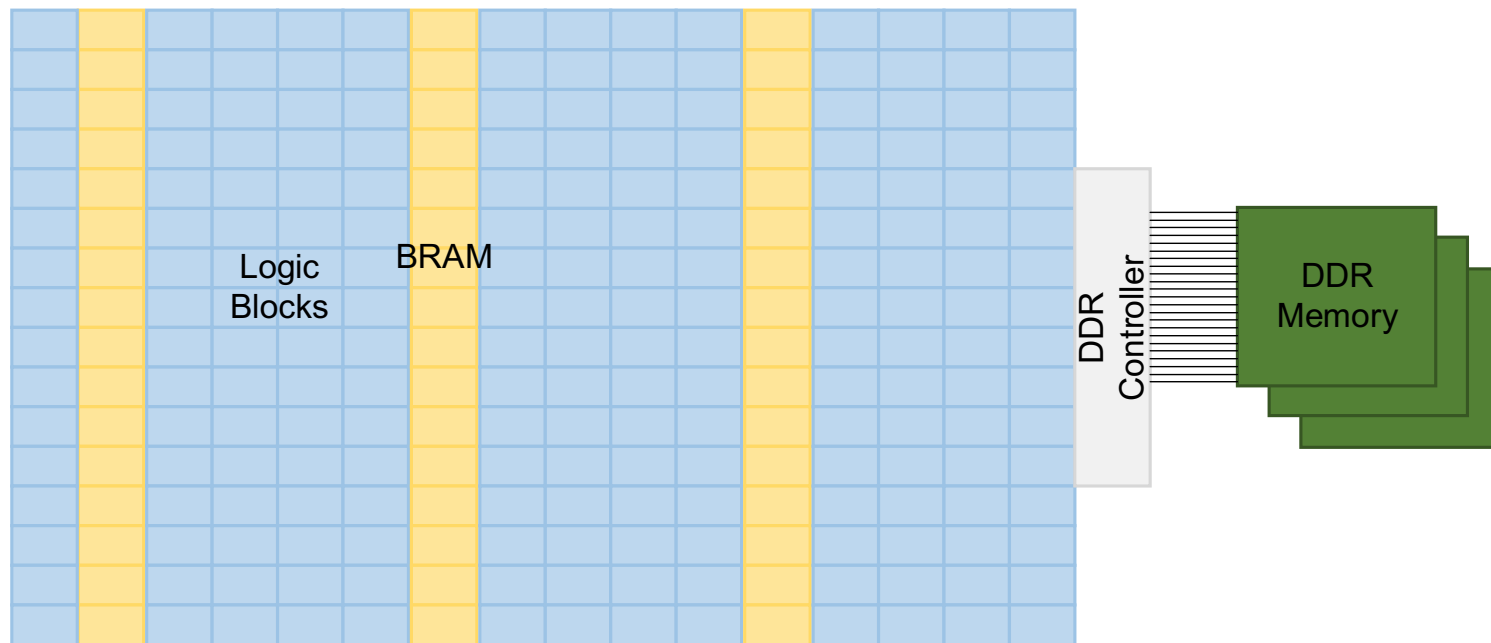
**[*]Toshiba Corporation, Japan**

The Edward S. Rogers Sr. Department
of Electrical & Computer Engineering
UNIVERSITY OF TORONTO
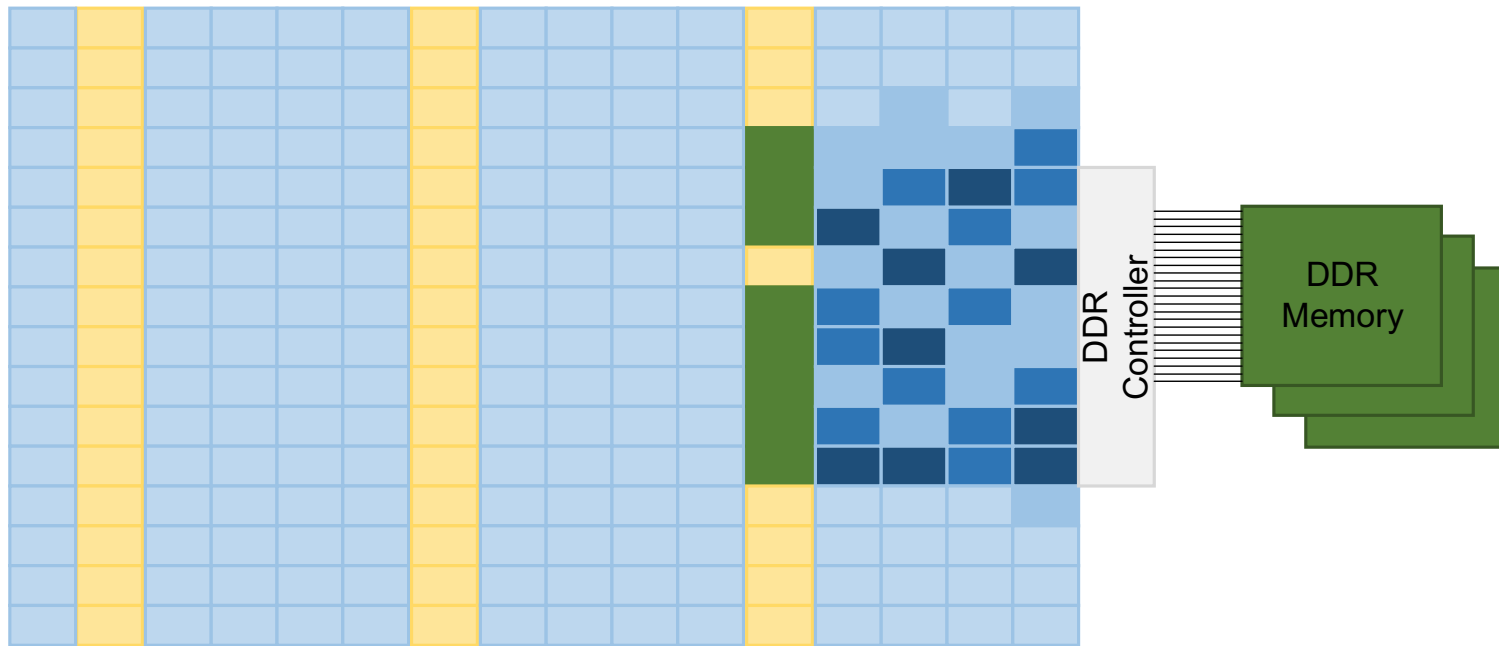
**TOSHIBA**
**Leading Innovation ≫**

# Memory: An Indispensable FPGA Component
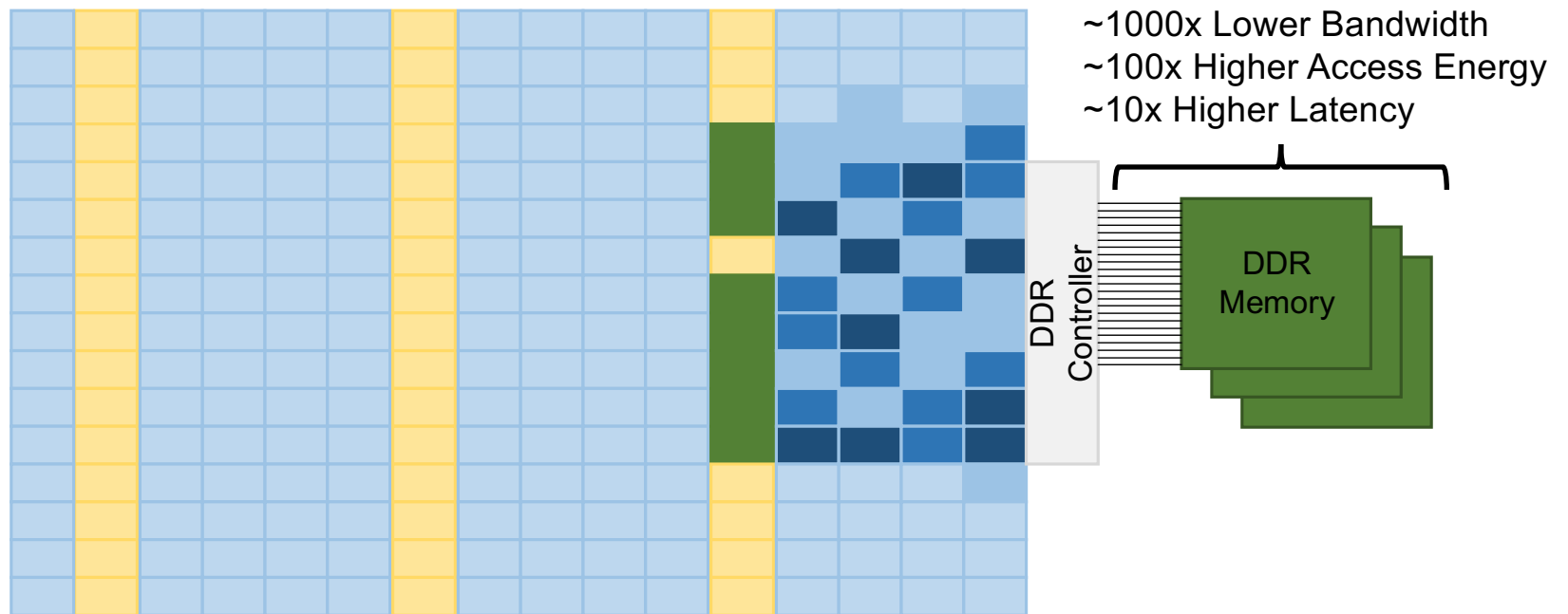
# Memory: An Indispensable FPGA Component

# Memory: An Indispensable FPGA Component



~1000x Lower Bandwidth
~100x Higher Access Energy
~10x Higher Latency

DDR Controller

DDR Memory

# Memory: An Indispensable FPGA Component

- BRAM growing in importance
  - Many applications (search engine, CNNs, …) BRAM-intensive
  - Can't fully utilize FPGA's computation capacity without on-chip memory

**Memory Richness Increase**

18kb +
288kb

Ultrascale+
16nm

18kb

Xilinx RAM bit/LE

4kb
Virtex 220nm

200

150

100

50

0

1.E+03    1.E+04    1.E+05    1.E+06    1.E+07

**Number of Logic Elements (4-input-LUT equivalent)**

# BRAM's Evolution

- Memory-richness growth
- Organization changes



Evolution of On-Chip Memory in Intel Stratix

# The Key Point

- BRAMs can't be neglected!
  - ~25% of area
  - Should respond to application demands
- Need BRAM models
  - How efficient is an architecture?
  - What's the best architecture?

# BRAM design is Difficult

- BRAM design is challenging!
  - Analog nature of some components



Sense Amplifier

# BRAM design is Difficult

- BRAM design is challenging!
  - Analog nature of some components
  - Variability of memory cells

# BRAM design is Difficult

- BRAM design is challenging!
  - Analog nature of some components
  - Variability of memory cells
  - Custom layout style
  - Significant FPGA-specific peripheral circuitry
- **Hand design of each candidate BRAM infeasible**

# Use existing tools?

# BRAM design is Different

- CACTI underestimates area

# BRAM design is Different

- Also underestimates read energy



13

# BRAM design is Different

- Overestimates operating frequency

# Emerging Memory Technologies

- Model promising emerging memory technologies
    - Magnetic Tunnel Junction (MTJ)
    - Phase Change Memory (PCM)
    - Resistive RAM (RRAM)
- **Ideal: model any technology with SPICE support**



**Magnetic Tunnel Junction (MTJ)**

Parallel    Anti-parallel

Free Layer
Insulator
Fixed Layer

Low resistance state (LRS)    High resistance state (HRS)

# BRAM Design Tool

# What do we need?

# COFFE: Logic & Routing



C. Chiasson et al. "COFFE: Fully-Automated Transistor Sizing for FPGAs," FPT 2013

# COFFE BRAM Flow

# 2-Bank SRAM-Based BRAM

**Memory Cells**

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Width-Configurable Decoders

Write Drivers and Sense Amplifiers

Write Drivers and Sense Amplifiers

Column Decoders

Input Crossbar and Level Shifter

Precharge and Equalizer

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

# 2-Bank SRAM-Based BRAM



| Memory Cells |
| --- |
| Traditional Decoders |
| Read/Write Circuitry |
| Width-Configurability Circuitry |
| To/From Routing |
| Pipeline Registers |

Output Crossbar

Width-Configurable Decoders

Write Drivers and Sense Amplifiers

Write Drivers and Sense Amplifiers

Column Decoders

Input Crossbar and Level Shifter

Precharge and Equalizer

Precharge and Equalizer

Row Decoders and Wordline Drivers

Bank #1 Memory Cell Array

Bank #2 Memory Cell Array

WL Port B          Vdd          WL Port B

$\overline{BL}$ Port B          BL Port B

$\overline{BL}$ Port A          BL Port A

WL Port A          WL Port A

SRAM Cell

# 2-Bank SRAM-Based BRAM

Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Width-Configurable Decoders

Write Drivers and Sense Amplifiers

Write Drivers and Sense Amplifiers

Input Crossbar and Level Shifter

Column Decoders

Precharge and Equalizer

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

# 2-Bank SRAM-Based BRAM

**Memory Cells**

**Traditional Decoders**

**Read/Write Circuitry**

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Width-Configurable Decoders

Input Crossbar and Level Shifter

Write Drivers and Sense Amplifiers

Write Drivers and Sense Amplifiers

Column Decoders

Precharge and Equalizer

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

# 2-Bank SRAM-Based BRAM

| |
|---|
| Memory Cells |
| Traditional Decoders |
| Read/Write Circuitry |
| Width-Configurability Circuitry |
| To/From Routing |
| Pipeline Registers |

Output Crossbar

Width-Configurable Decoders

Input Crossbar and Level Shifter

Write Drivers and Sense Amplifiers

Column Decoders

Write Drivers and Sense Amplifiers

Precharge and Equalizer

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

Vdd    Vdd

$\overline{BL}$    BL

Precharge and equalizer

# 2-Bank SRAM-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Width-Configurable Decoders

Input Crossbar and Level Shifter

Write Drivers and Sense Amplifiers

Column Decoders

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Write Drivers and Sense Amplifiers

Precharge and Equalizer

Bank #2 Memory Cell Array

Vdd    Vdd

Data

Write Driver

# 2-Bank SRAM-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Width-Configurable Decoders

Input Crossbar and Level Shifter

Write Drivers and Sense Amplifiers

Column Decoders

Write Drivers and Sense Amplifiers

Precharge and Equalizer

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

Sense Amplifier

26

# 2-Bank SRAM-Based BRAM



**Legend:**
- Memory Cells
- Traditional Decoders
- Read/Write Circuitry
- Width-Configurability Circuitry
- To/From Routing
- Pipeline Registers

Output Crossbar

Input Crossbar and Level Shifter

Width-Configurable Decoders

Write Drivers and Sense Amplifiers

Column Decoders

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

# 2-Bank SRAM-Based BRAM



Legend:
- Memory Cells
- Traditional Decoders
- Read/Write Circuitry
- Width-Configurability Circuitry
- To/From Routing
- Pipeline Registers

Diagram labels:
- Output Crossbar
- Width-Configurable Decoders
- Input Crossbar and Level Shifter
- Write Drivers and Sense Amplifiers
- Column Decoders
- Precharge and Equalizer
- Bank #1 Memory Cell Array
- Bank #2 Memory Cell Array

28

# 2-Bank SRAM-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Input Crossbar and Level Shifter

Width-Configurable Decoders

Write Drivers and Sense Amplifiers

Column Decoders

Precharge and Equalizer

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

# 2-Bank MTJ-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Input Crossbar and Level Shifter

Sense Amplifier

Write Driver

Write Driver

Column Selector and Predischarger

Reference Cells

Bank #1 Memory Cell Array

Width-Configurable Decoders

Column Decoders

Row Decoders and Wordline Drivers

Sense Amplifier

Write Driver

Write Driver

Column Selector and Predischarger

Reference Cells

Bank #2 Memory Cell Array

MTJ-based Memory Cell

30

# 2-Bank MTJ-Based BRAM

| Memory Cells |
| --- |
| Traditional Decoders |
| Read/Write Circuitry |
| Width-Configurability Circuitry |
| To/From Routing |
| Pipeline Registers |

Output Crossbar

Input Crossbar and Level Shifter

Sense Amplifier

Write Driver

Write Driver

Column Selector and Predischarger

Reference Cells

Bank #1 Memory Cell Array

Width-Configurable Decoders

Column Decoders

Row Decoders and Wordline Drivers

Sense Amplifier

Write Driver

Write Driver

Column Selector and Predischarger

Reference Cells

Bank #2 Memory Cell Array

# 2-Bank MTJ-Based BRAM



| Memory Cells |
| Traditional Decoders |
| Read/Write Circuitry |
| Width-Configurability Circuitry |
| To/From Routing |
| Pipeline Registers |

Output Crossbar

Input Crossbar and Level Shifter

Width-Configurable Decoders

Sense Amplifier

Sense Amplifier

Write Driver

Write Driver

Write Driver

Write Driver

Column Selector and Predischarger

Column Decoders

Column Selector and Predischarger

Reference Cells

Reference Cells

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

# 2-Bank MTJ-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Input Crossbar and Level Shifter

Width-Configurable Decoders

Sense Amplifier

Sense Amplifier

Write Driver

Write Driver

Write Driver

Write Driver

Column Selector and Predischarger

Column Decoders

Column Selector and Predischarger

Reference Cells

Reference Cells

Bank #1 Memory Cell Array

Row Decoders and Wordline Drivers

Bank #2 Memory Cell Array

Column Selector and Predischarger

33

# 2-Bank MTJ-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Input Crossbar and Level Shifter

Width-Configurable Decoders

Sense Amplifier

Write Driver

Column Selector and Predischarger

Column Decoders

Reference Cells

Row Decoders and Wordline Drivers

Bank #1 Memory Cell Array

Bank #2 Memory Cell Array

Write Driver

# 2-Bank MTJ-Based BRAM



Memory Cells

Traditional Decoders

Read/Write Circuitry

Width-Configurability Circuitry

To/From Routing

Pipeline Registers

Output Crossbar

Width-Configurable Decoders

Input Crossbar and Level Shifter

Sense Amplifier

Write Driver

Write Driver

Column Selector and Predischarger

Reference Cells

Column Decoders

Row Decoders and Wordline Drivers

Bank #1 Memory Cell Array

Sense Amplifier

Write Driver

Write Driver

Column Selector and Predischarger

Reference Cells

Bank #2 Memory Cell Array

Vdd

Output

Sense Amplifier

# 2-Bank MTJ-Based BRAM



**Legend:**
- Memory Cells
- Traditional Decoders
- Read/Write Circuitry
- Width-Configurability Circuitry
- To/From Routing
- Pipeline Registers

**Diagram labels:**
- Output Crossbar
- Input Crossbar and Level Shifter
- Sense Amplifier
- Width-Configurable Decoders
- Write Driver
- Column Decoders
- Column Selector and Predischarger
- Reference Cells
- Bank #1 Memory Cell Array
- Row Decoders and Wordline Drivers
- Bank #2 Memory Cell Array

# 2-Bank MTJ-Based BRAM

# Simulation In Context: SRAM Precharge

# Simulation In Context: SRAM Precharge

# Simulation In Context: SRAM Precharge

# Simulation In Context: Worst-case SRAM Cell

- Variation changes SRAM cell read/write currents significantly

# Simulation In Context: Worst-case SRAM Cell

- Variation changes SRAM cell read/write currents significantly
- Using the nominal memory cell will be inaccurate
  - BRAM energy will be underestimated
  - BRAM frequency will be overestimated

# Simulation In Context: Worst-case SRAM Cell

- Variation changes SRAM cell read/write currents significantly
- Using the nominal memory cell will be inaccurate
  - BRAM energy will be underestimated
  - BRAM frequency will be overestimated
- Don't have the Spice model for the worst-case cell!

# Simulation In Context: Worst-case SRAM Cell

- Variation changes SRAM cell read/write currents significantly
- Using the nominal memory cell will be inaccurate
  - BRAM energy will be underestimated
  - BRAM frequency will be overestimated
- We don't have the Spice model for the worst-case cell!
- Use Monte Carlo simulation to find distribution of cell properties

# Monte Carlo Simulation: Worst-case Cell



K. Tatsumura et al."High Density, Low Energy, Magnetic Tunnel Junction Based Block RAMs for Memory-Rich FPGAs," FPT 2016
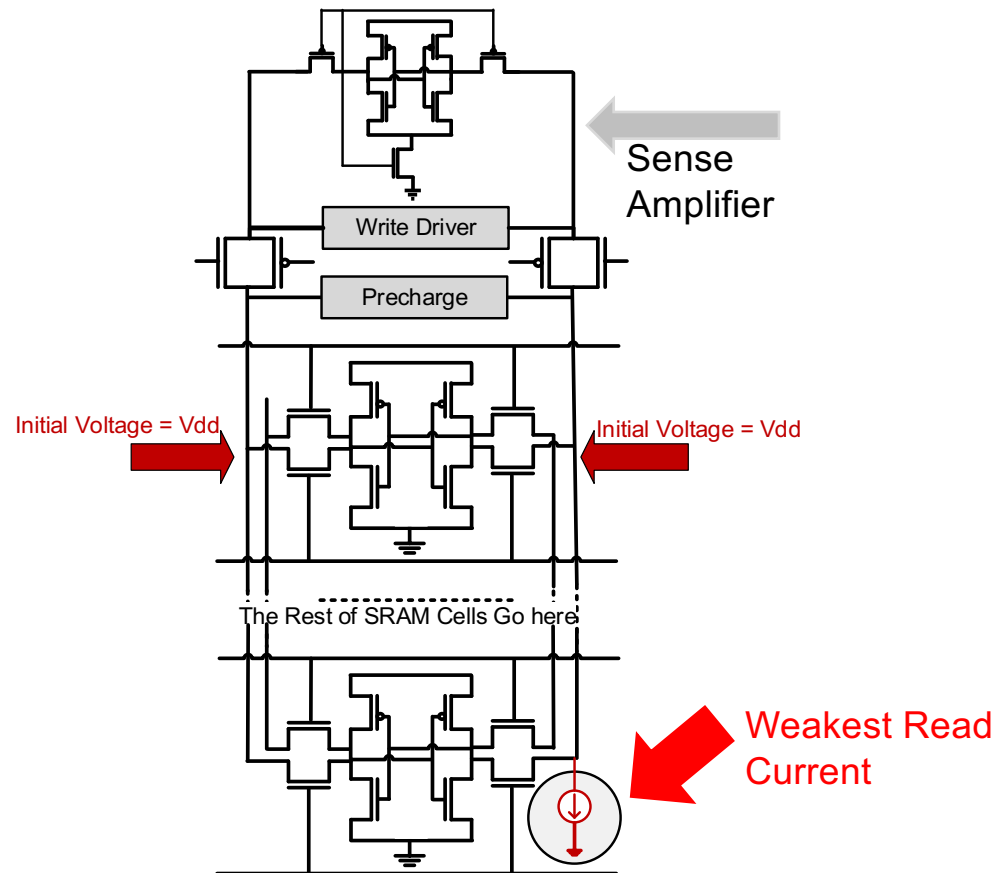
# Simulation In Context: Worst-Case Cell

# Simulation In Context: Worst-Case Cell

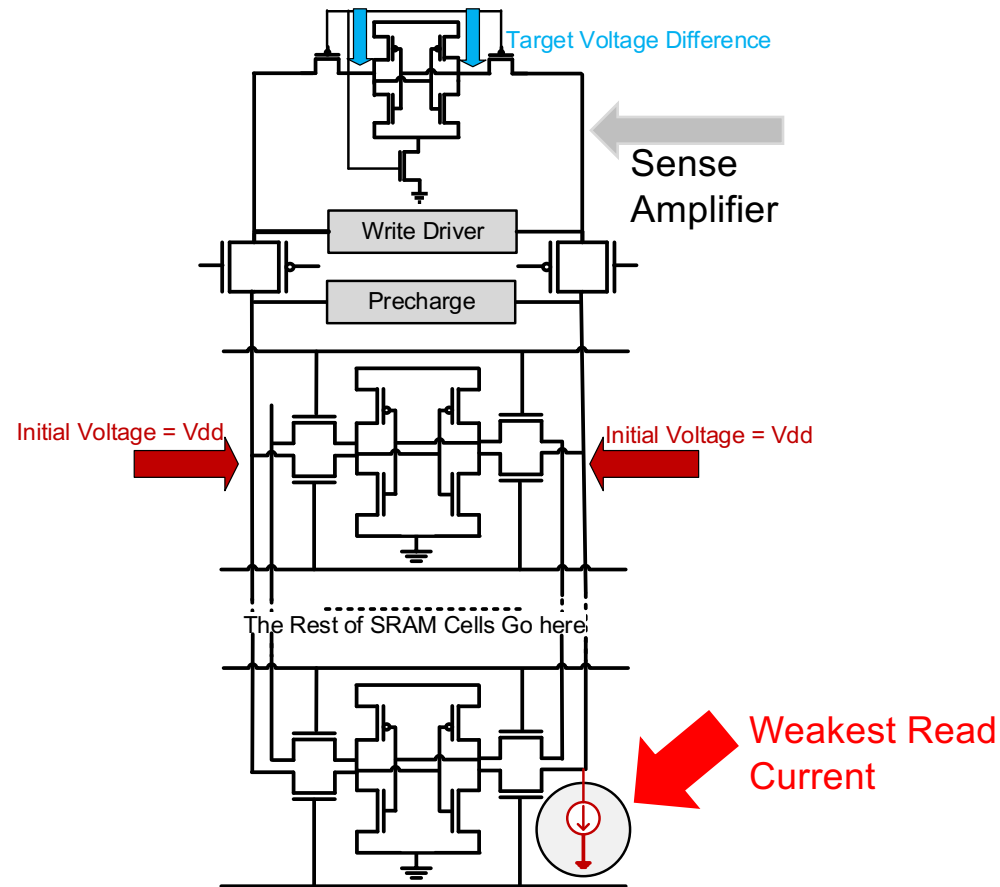# Simulation In Context: Worst-Case Cell



Sense Amplifier

Write Driver

Precharge

Initial Voltage = Vdd

Initial Voltage = Vdd

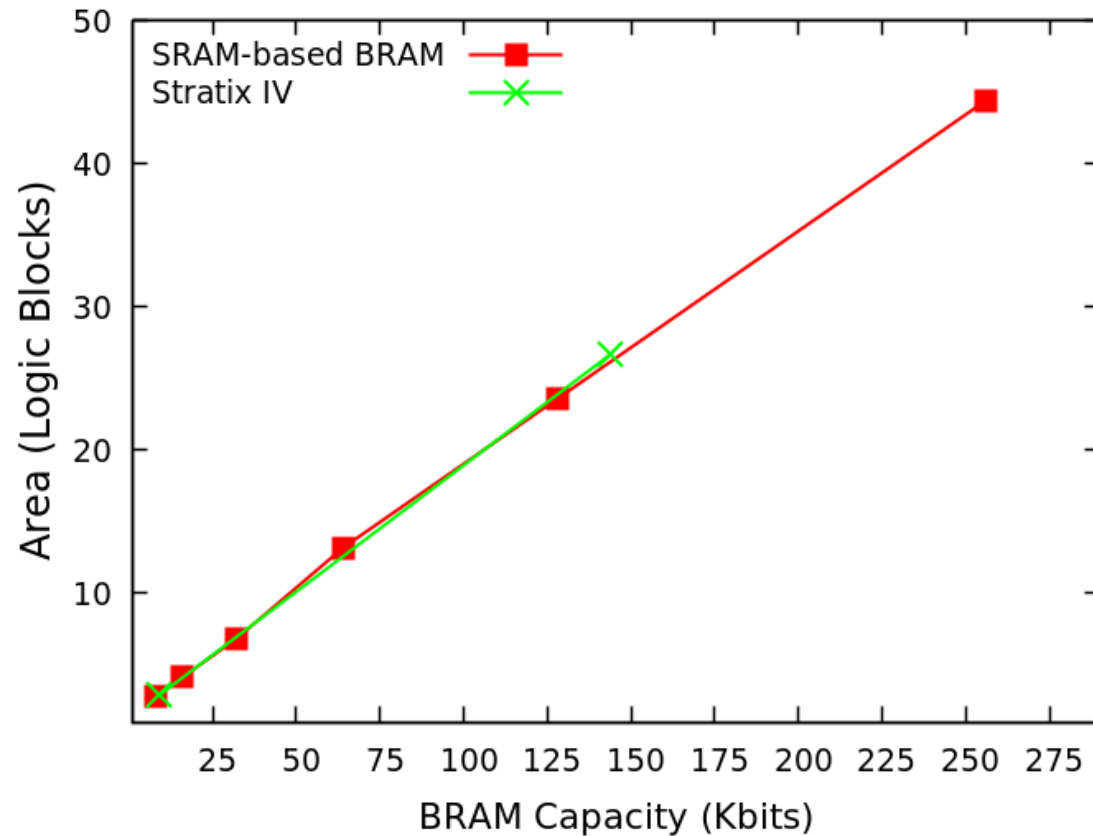The Rest of SRAM Cells Go here

Weakest Read Current

48

# Simulation In Context: Worst-Case Cell

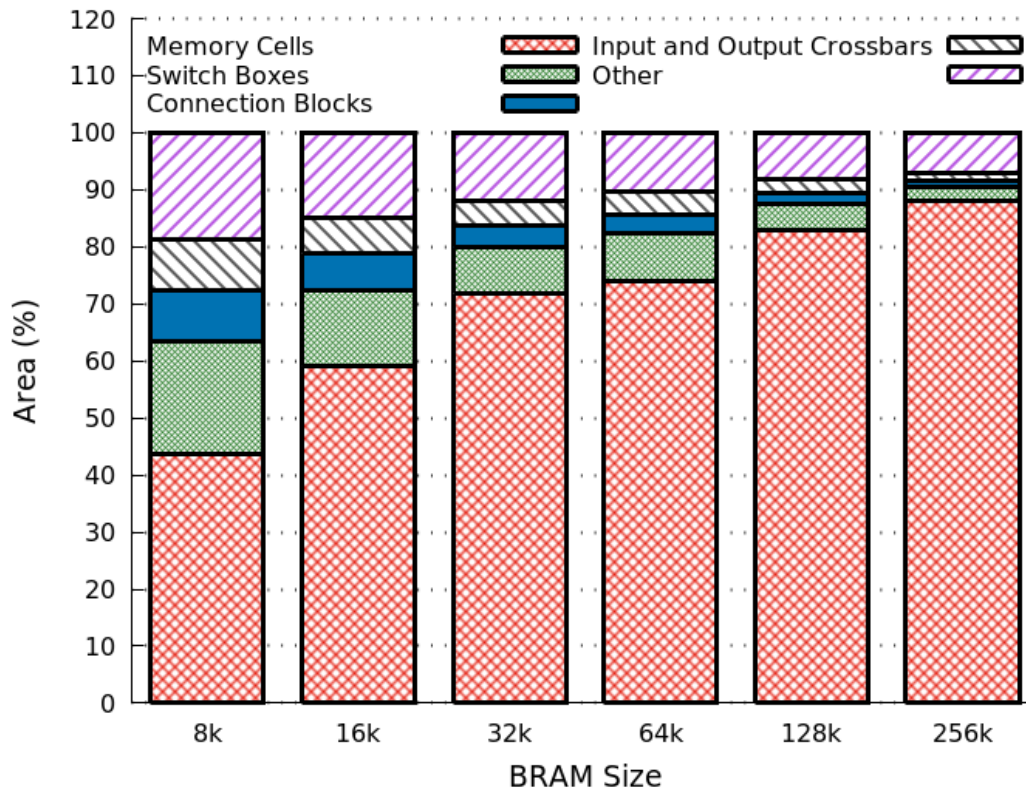

49

# Validation and Results

# Area Validation

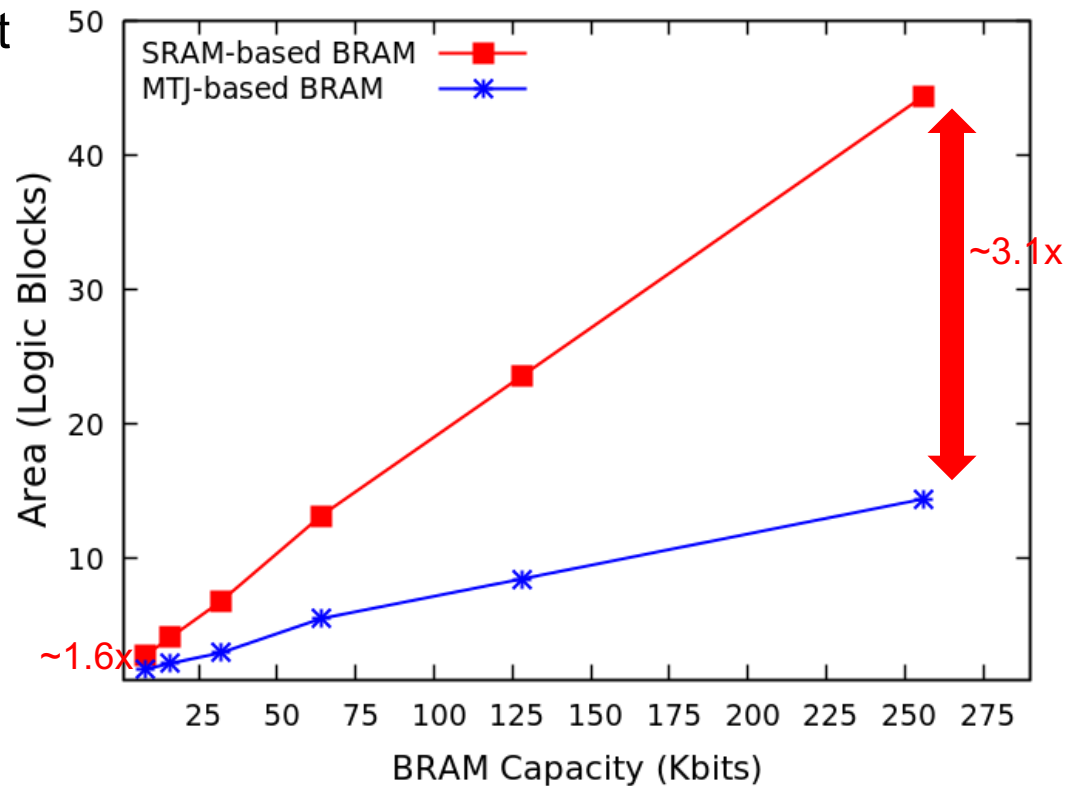- Area results align well with commercial data

# SRAM-based BRAM Area Breakdown

- SRAM area dominates for large BRAMs
- Smaller BRAMs: other components relevant

# Area: MTJ vs. SRAM

- MTJ is more area-efficient
- It gets increasingly more efficient with BRAM size

# Frequency validation

- Reasonable alignment with commercial data
- Less guardband
- No aggressive banking

# Operating Frequency: MTJ vs. SRAM

- SRAM is faster
- The gap narrows with increasing size

# Simulation Results: Energy

# SRAM-based BRAM Narrow Mode

- FPGA RAM configurable
- Often used in narrower modes
- Energy mostly unaffected

# Energy Per Bit: MTJ vs. SRAM

- MTJs are more efficient with large BRAM
- MTJ narrow mode more efficient

# Worst-Case Cell Modeling Crucial

- Use nominal memory cell?
  - Underestimates area and energy
  - Gets more severe with increasing memory size

| BRAM Capacity (Kbit) | | Change in Delay | Change in Energy per bit |
|---|---|---|---|
| 8 | | -21% | -9% |
| 16 | | -19% | -6% |
| 32 | | -27% | -15% |
| 64 | | -22% | -9% |
| 128 | | -30% | -20% |
| 256 | | -42% | -29% |

# Architecture Exploration

# Architecture Exploration: RAM-Mapping Flow

- BRAM models generated by COFFE can be used in architecture exploration
- Area-oriented RAM mapping
  - 69 industrial circuits
  - Used in development of Stratix V memory architecture
  - We have partial data:
    - Number of logic blocks used
    - Number, Sizes, and types of Logical RAMs
  - Gradually excluding less-memory-rich circuits

# SRAM-based BRAM

- 16K always the best
- Stratix V-like

# MTJ-based BRAM

- MTJ always saves area
- The best architecture changes
  - 32k or 64k best

# Architecture Exploration: VTR Flow

- Nine VTR benchmark circuits with memory
- MTJ vs. SRAM
- Architecture Parameters
  - 32kb BRAM, every 8 columns
    - MTJ BRAM is smaller → get more 2.3x more RAM blocks per column
  - Ten 6-luts per logic block

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

❑ **Changes by switching to MTJ-based BRAMs:**

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---------|---------------|------------|--------------|------------|-------------|---------------|-------------|--------------------|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---------|---------------|------------|--------------|------------|-------------|---------------|-------------|--------------------|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Architecture Exploration: VTR

- Changes by switching to MTJ-based BRAMs:

| Circuit | RAM/LUT Ratio | Block Area | Routing Area | Total Area | Block Delay | Routing Delay | Total Delay | Area-delay Product |
|---|---|---|---|---|---|---|---|---|
| mcml | 1% | -11% | 0 | -5% | 5% | -19% | -6% | -10% |
| LU32PEEng | 7% | -14% | 9% | -2% | 9% | -6% | 0 | -1% |
| LU8PEEng | 6% | -13% | 4% | -5% | -4% | 9% | 3% | -2% |
| ch_intrinsics | 2% | -15% | -2% | -10% | -1% | 6% | 3% | -7% |
| mkDelayWorker32B | 24% | -32% | -47% | -41% | 60% | -40% | 3% | -39% |
| mkPktMerge | 198% | -57% | -48% | -53% | 141% | -34% | 12% | -47% |
| mkSMAdapter4B | 8% | -16% | 0 | -10% | 92% | -32% | 9% | -2% |
| or1200 | 2% | -5% | -4% | 0 | 1% | -7% | -2% | -3% |
| raygentop | 1% | -3% | 2% | -1% | 9% | -17% | -4% | -5% |
| boundtop | 1% | -3% | 1% | -1% | 8% | -5% | 0 | -2% |
| Geometric Mean | 4% | -19% | -11% | -15% | 25% | -16% | 2% | -14% |

# Conclusion

- First transistor sizing tool capable of BRAM modeling
  - SRAM-based
  - MTJ-based
- Simulation results align well with available commercial data
- COFFE now enables BRAM architecture exploration!
  - RAM-Mapping
  - VTR