

# RESISTIVE COMPUTATION: AVOIDING THE POWER WALL WITH LOW-LEAKAGE, STT-MRAM BASED COMPUTING

Xiaochen Guo, Engin Ipek, and Tolga Soyata

Rochester Computer Systems Architecture Laboratory

# Multicore Scaling Limited by Power

2

- Traditional MOSFET scaling theory relies on reducing  $V_{DD}$  in proportion to device dimensions

$$P = P_{\text{dynamic}} + P_{\text{leak}} = N \cdot (C_{\text{eff}} \cdot V_{DD}^2 \cdot f) + I_{\text{leak}} \cdot V_{DD}$$

2x ↑      1.4x ↑      
 $I_{\text{leak}} \propto e^{-V_{th}}$

1.4x ↓      2x ↓      1.4x ↓

- $V_{DD}$  has scaled very slowly since 90nm
- Multicore scaling severely challenged by power

# Our Approach: Resistive Computation

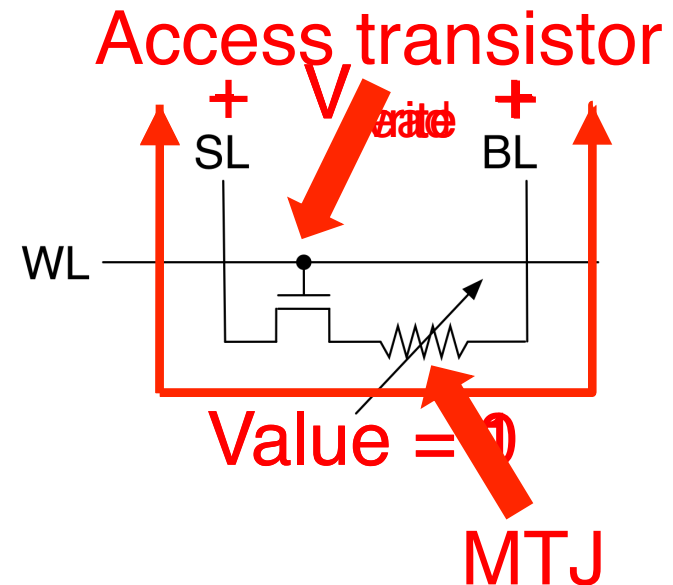
3

- Opportunity: spin-torque transfer magnetoresistive RAM (STT-MRAM)
  - ▣ Near-zero leakage power
  - ▣ Low-energy read operation
  
- Goal: selectively migrate on-chip storage and combinational logic to STT-MRAM to reduce power
  - ▣ On-chip storage
    - Caches, TLBs, RF, queues
  - ▣ Combinational logic
    - Lookup-table (LUT) based computing

# STT-MRAM

4

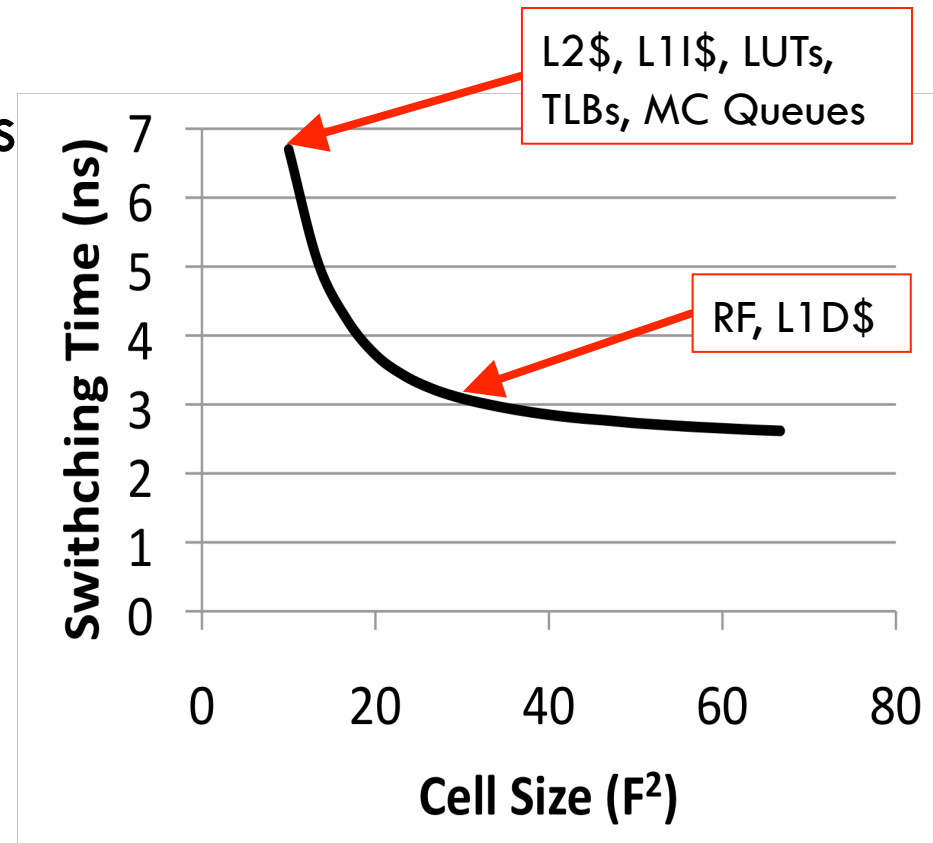
- Desirable properties
  - ▣ CMOS compatibility
  - ▣ Read speed as fast as SRAM
  - ▣ Density comparable to DRAM
  - ▣ Unlimited write endurance
  
- Key challenge: expensive writes
  - ▣ Long switching latency (6.7ns @ 32nm)
  - ▣ High switching energy (0.3pJ/bit @ 32nm)



# Switching Time vs. Cell Size

5

- Faster switching with wider access transistors
  - + Faster writes
  - Slower reads
  - Lower density
  - Higher read energy



# Fundamental Building Blocks

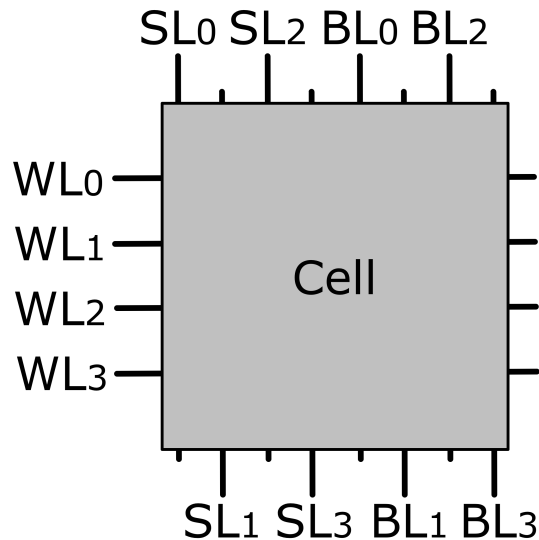
RAM Arrays and Lookup Tables

# STT-MRAM Arrays

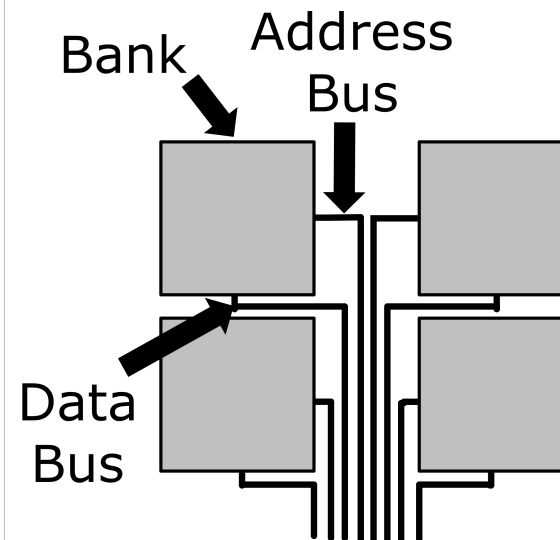
7

- Problem: low write throughput

## Multiplexing



## Banking

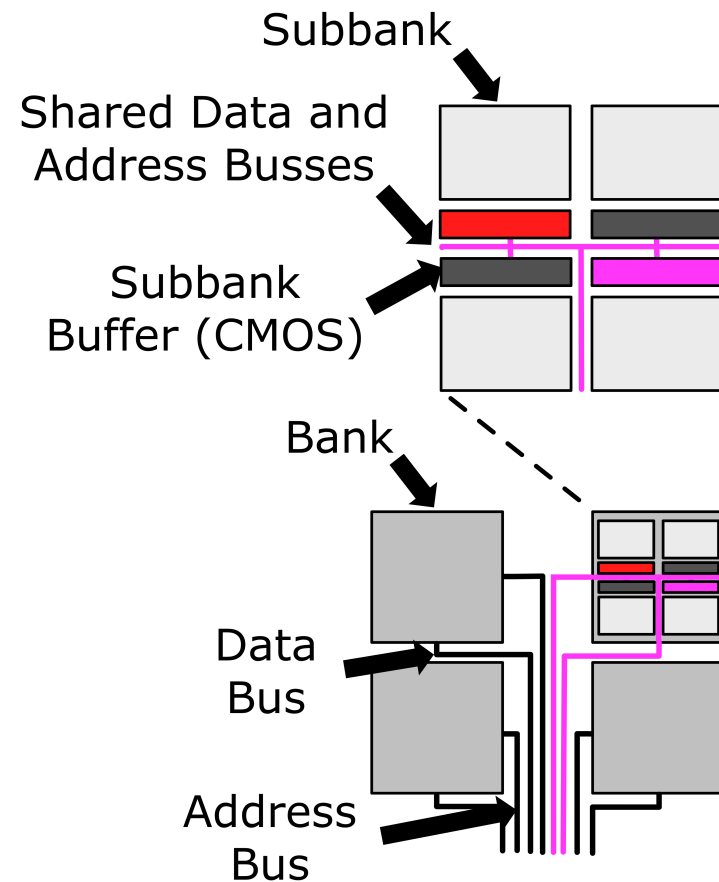


- Existing solutions incur high overheads to sustain adequate write throughput in STT-MRAM arrays

# STT-MRAM Arrays

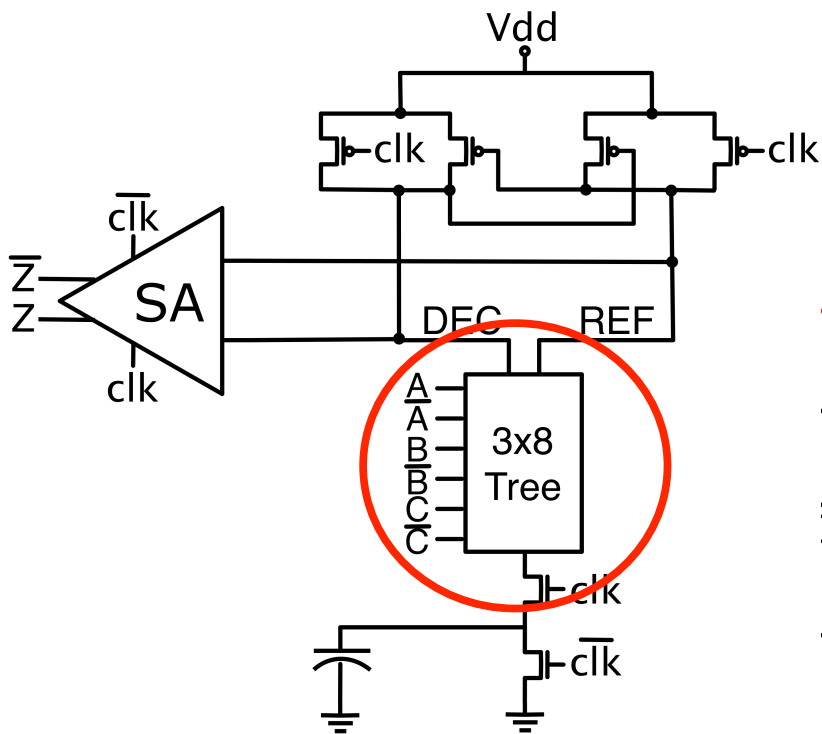
8

- CMOS subbank buffers
  - ▣ Latch in addr/data and release H-tree; complete write locally
  - ▣ Allow forwarding from ongoing writes
  - ▣ Facilitate local differential writes
  
- Reads access subbank via exclusive read port





# STT-MRAM LUTs [Suzuki09, Matsunaga08]



- Store truth tables of logic functions directly in STT-MRAM

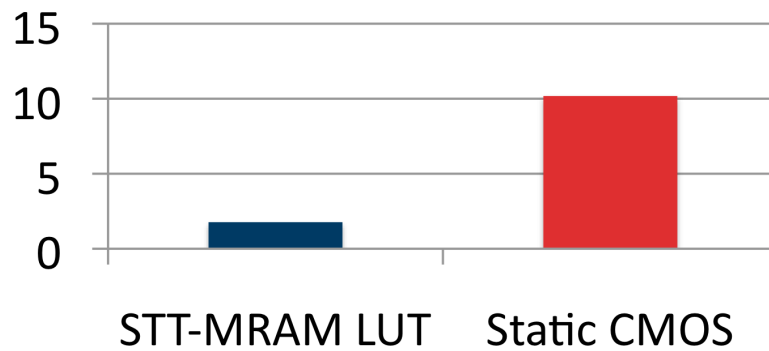
## Benefits

- ▣ Leakage confined to peripheral circuitry
  - ▣ Low-power (low swing) lookups
  - ▣ Fast lookups using sense amp
- Logic functions with many minterms can utilize LUTs effectively

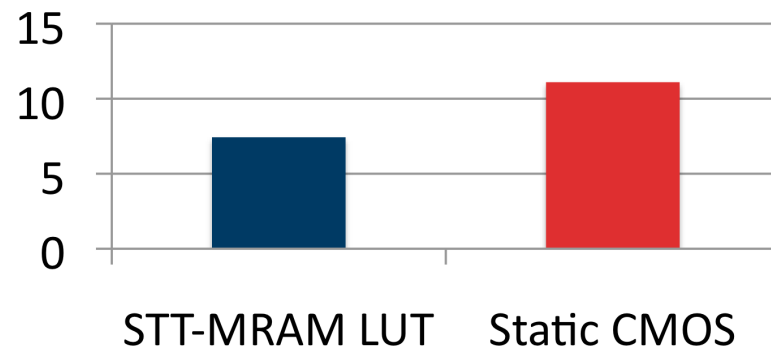
# Case Study: 3-bit Adder

10

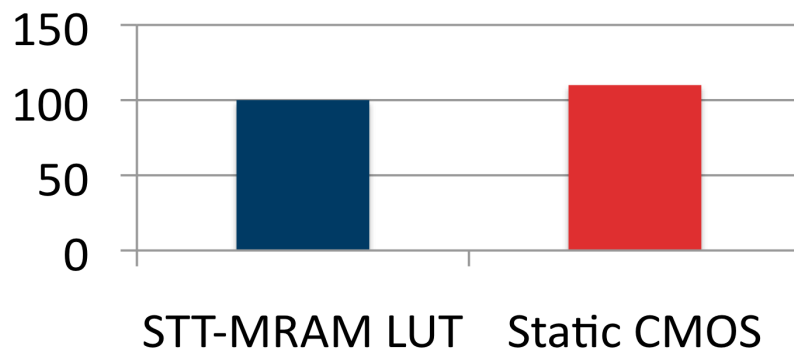
### Leakage Power (nW)



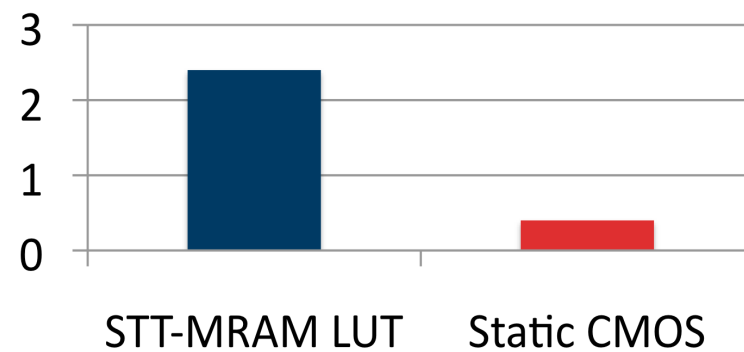
### Access Energy (fJ)



### Delay (ps)



### Active Area ( $\mu\text{m}^2$ )



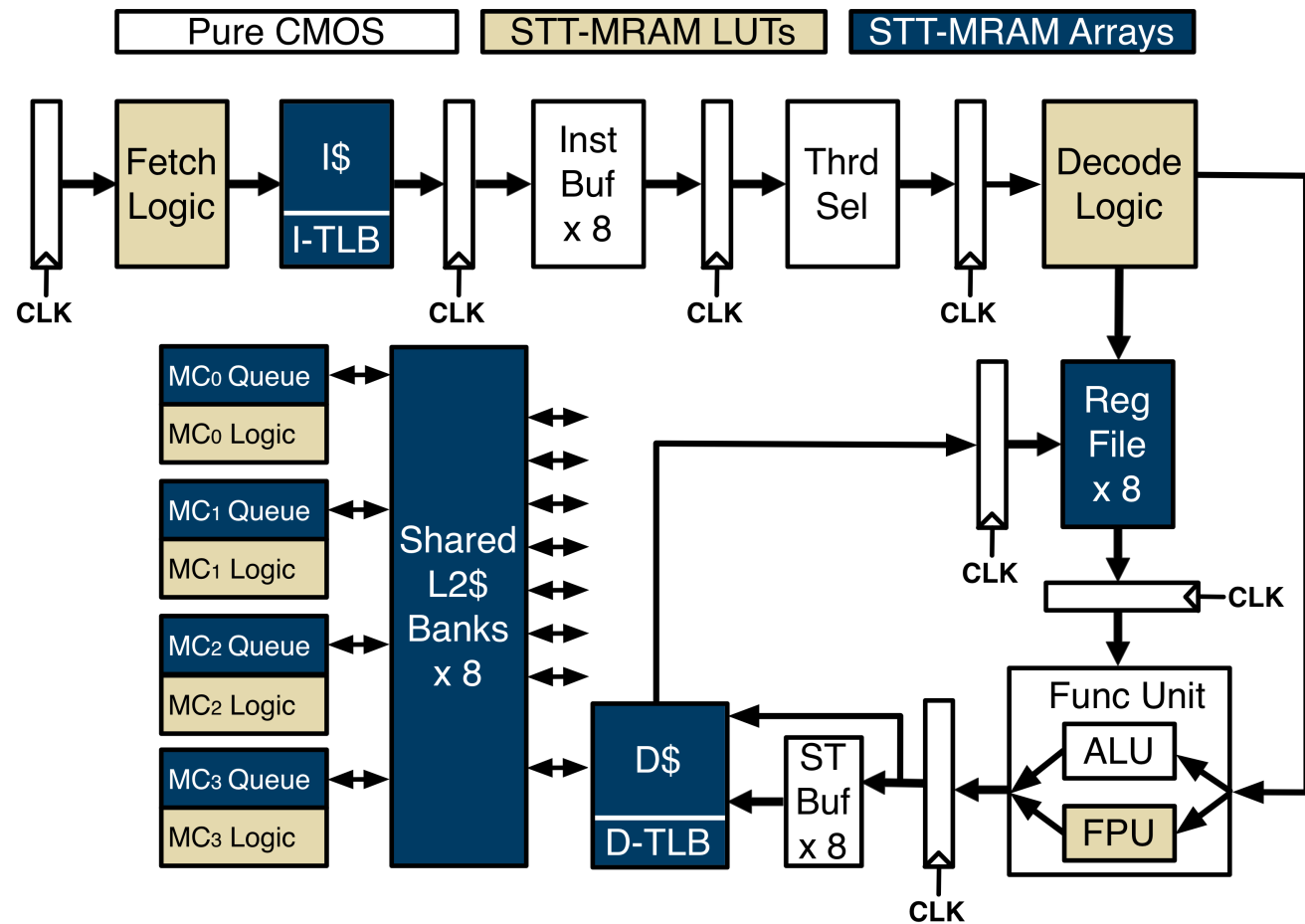
# Pipeline Organization

# Hybrid CMT Pipeline

12

Small arrays and simple logic in CMOS

Large arrays and complex logic in STT-MRAM



# Front End

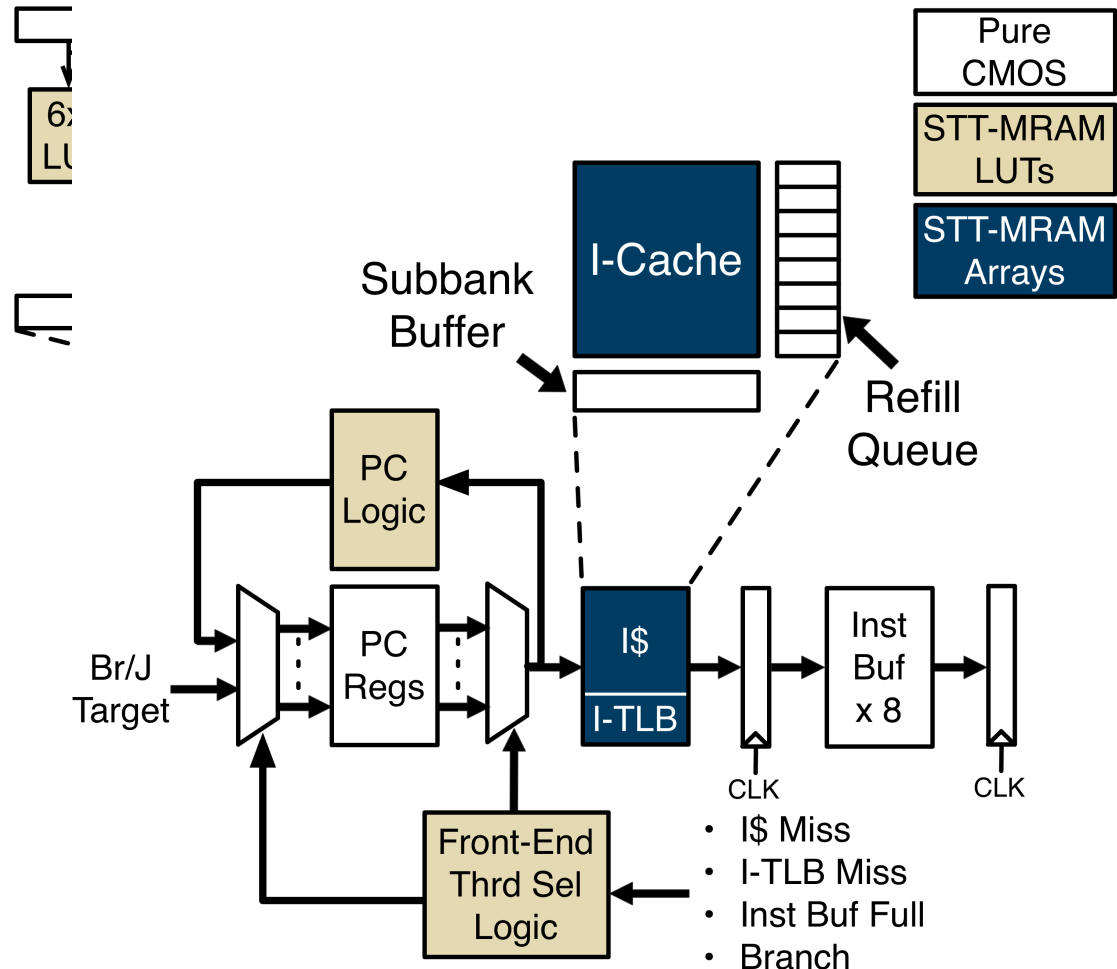
13

LUT-based carry-select adder to compute PC+4

LUT-based front-end thread selection logic

SRAM-based refill queue to avoid I\$ conflicts

Predecode and back-end thread selection with MRAM-related stall conditions

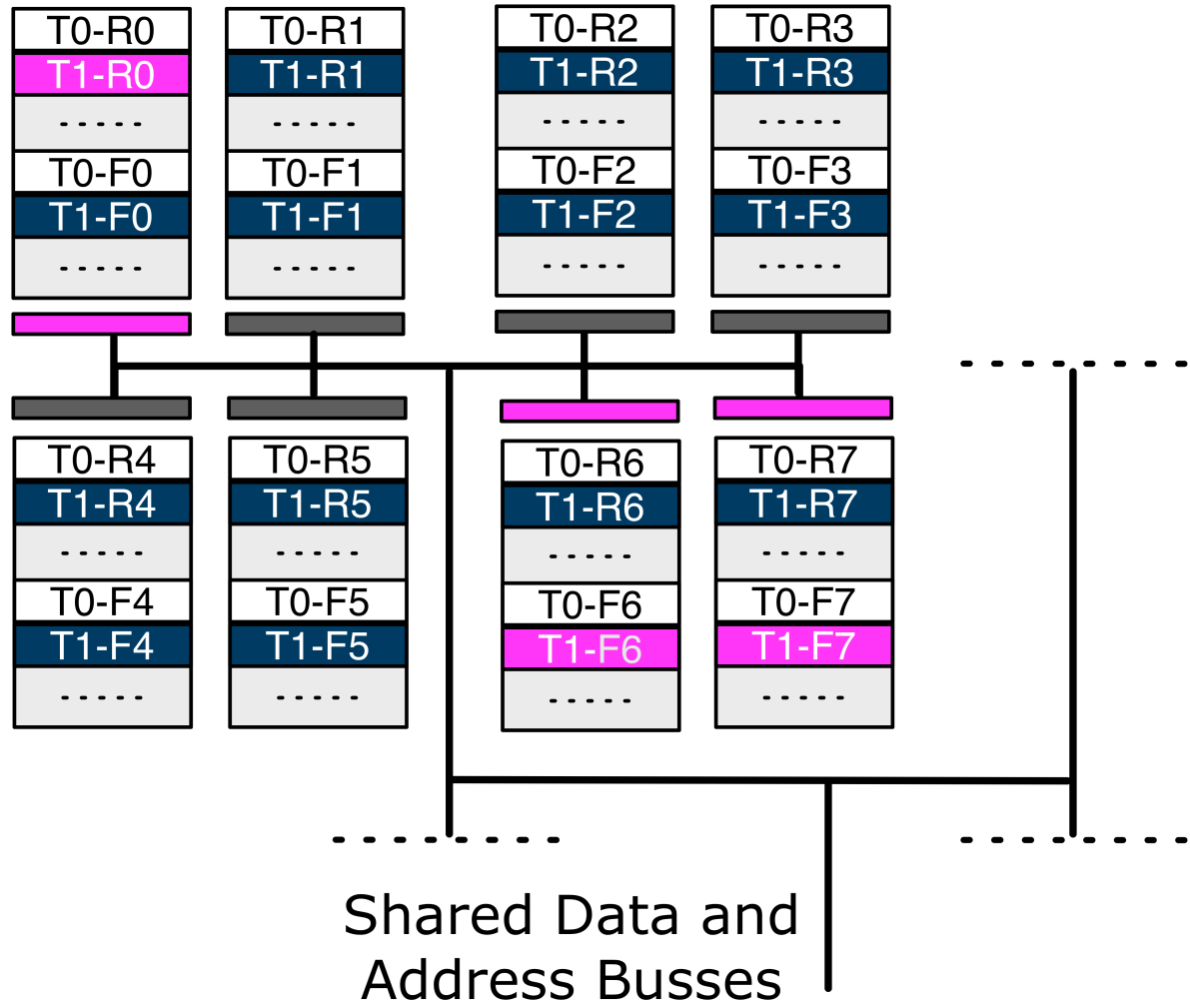


# Register File

14

Architectural registers of all threads aggregated in a unified STT-MRAM array to amortize subbank buffers

Registers of a single thread striped across subbanks to reduce subbank buffer conflicts

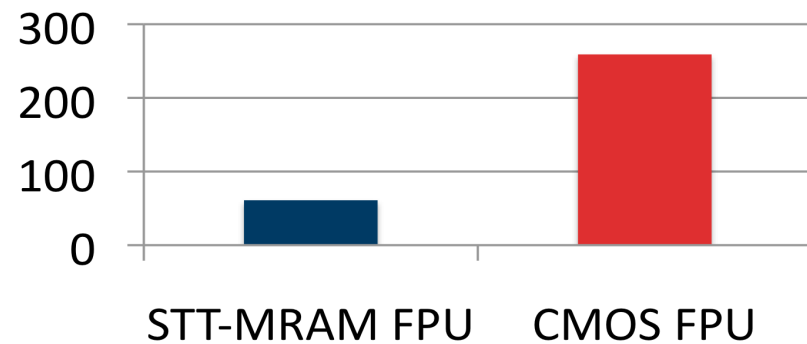


# Floating-Point Unit

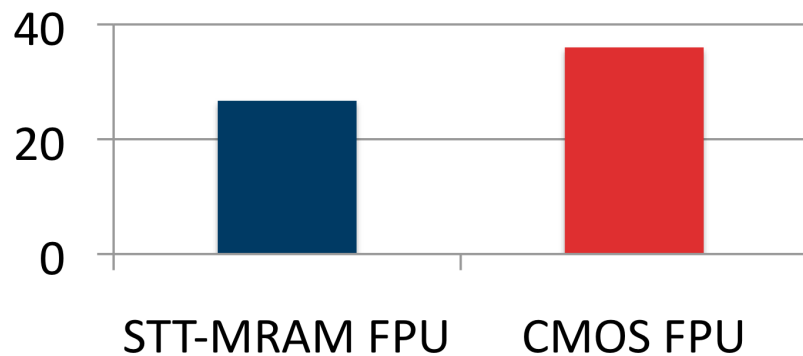
15

	STT-MRAM FPU	CMOS FPU
Add, Sub, Mult	24 cycles	12 cycles
Div	64 cycles	64 cycles

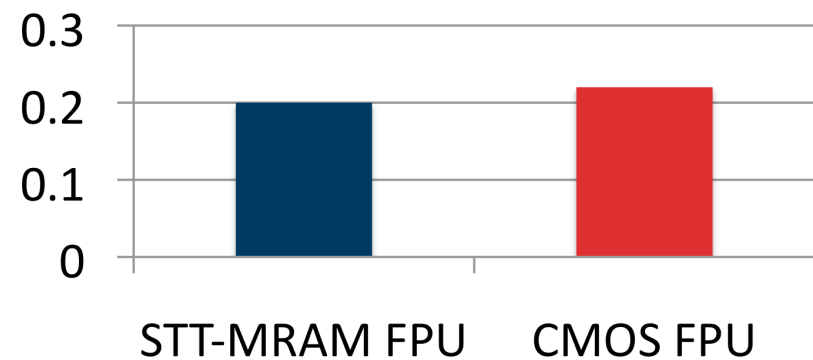
### Leakage Power (mW)



### Dynamic Energy (pJ)



### Active Area (mm<sup>2</sup>)

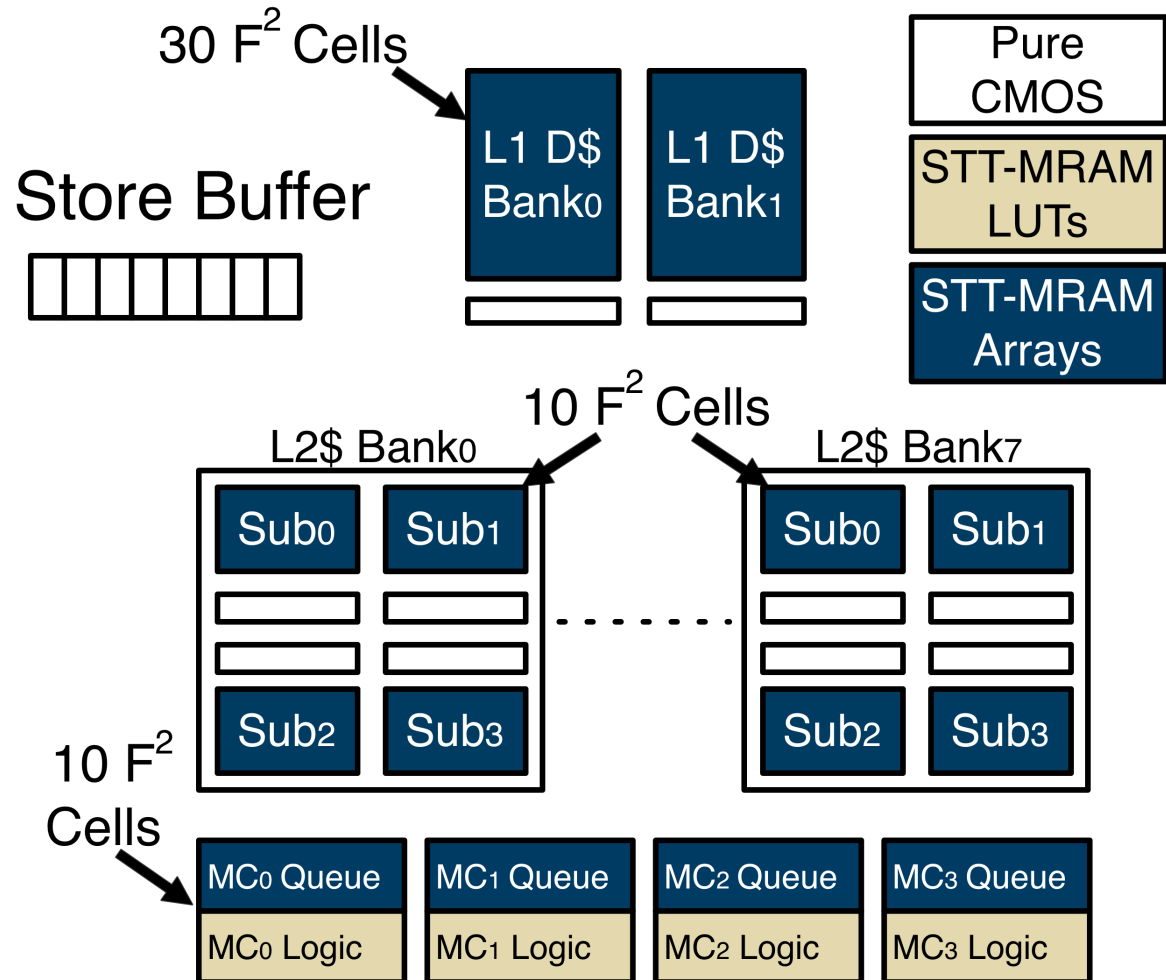


# Memory System

Use store buffers to avoid L1 D\$ subbank conflicts

L1s optimized for fast writes using  $30F^2$  cells

L2 and memory controllers optimized for density using  $10F^2$  cells

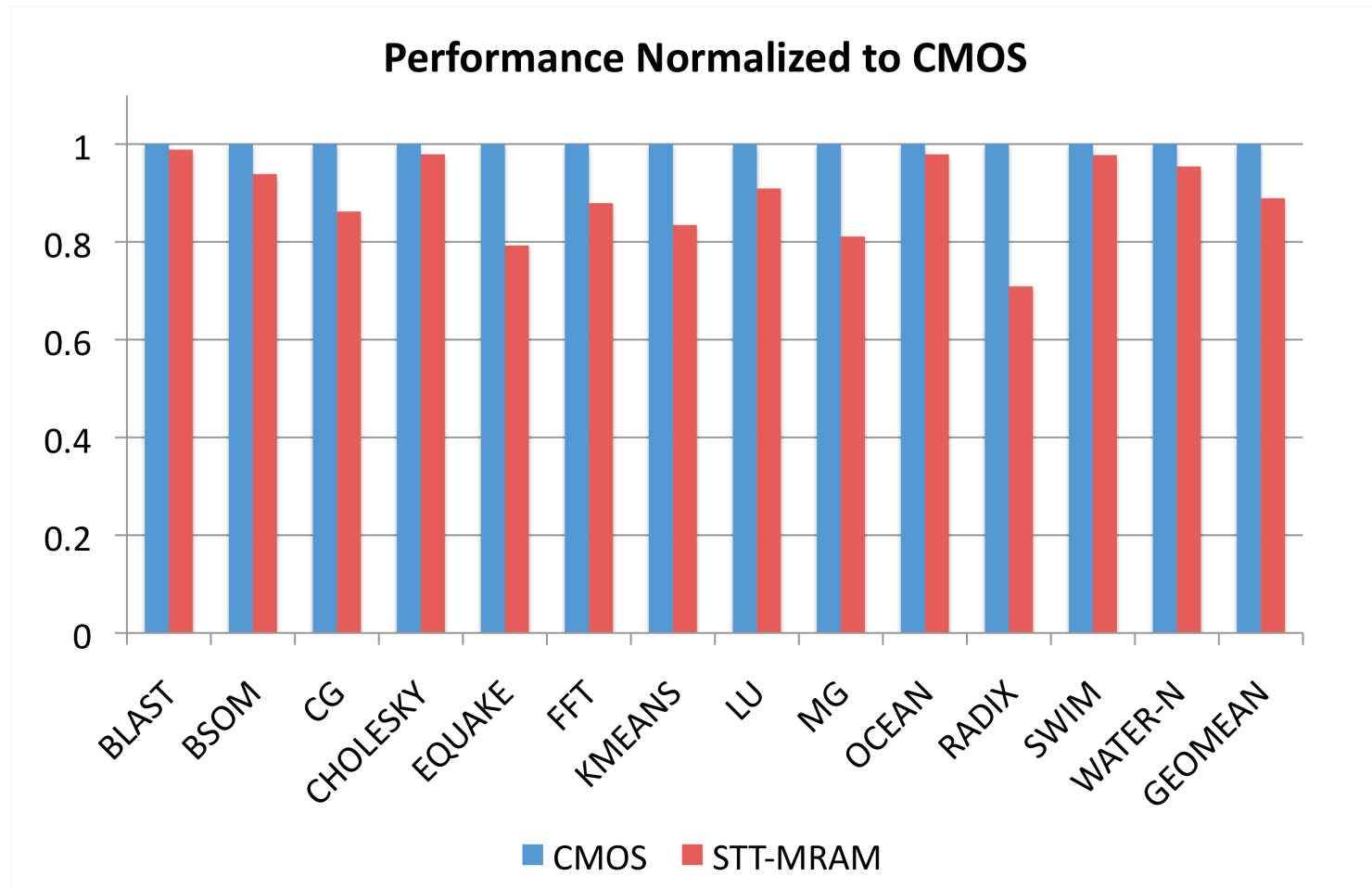




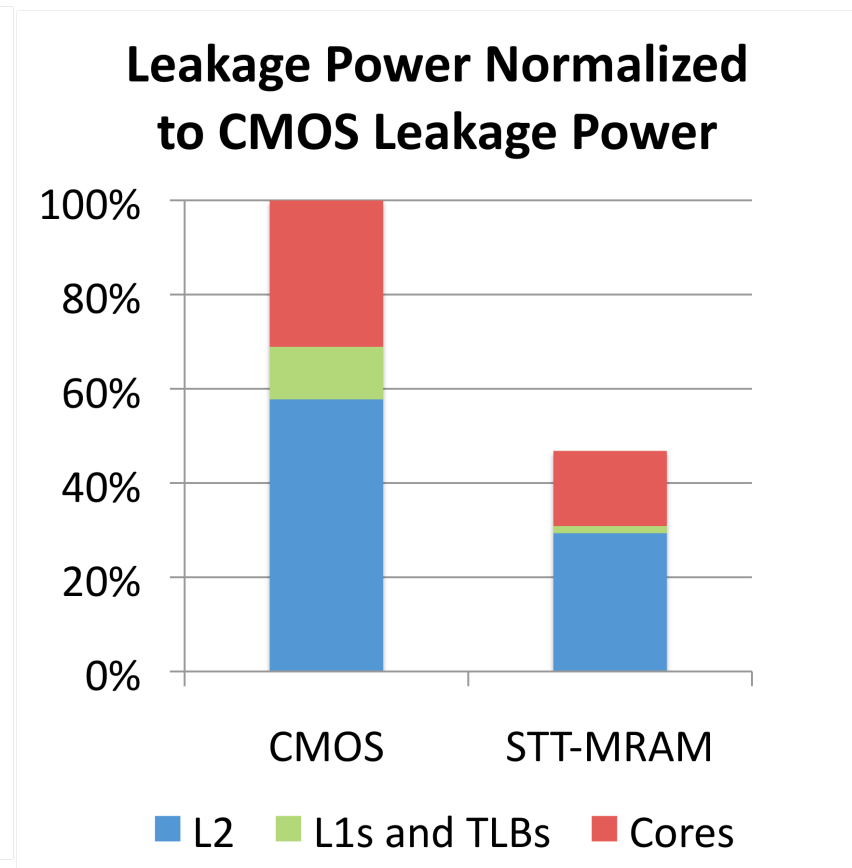
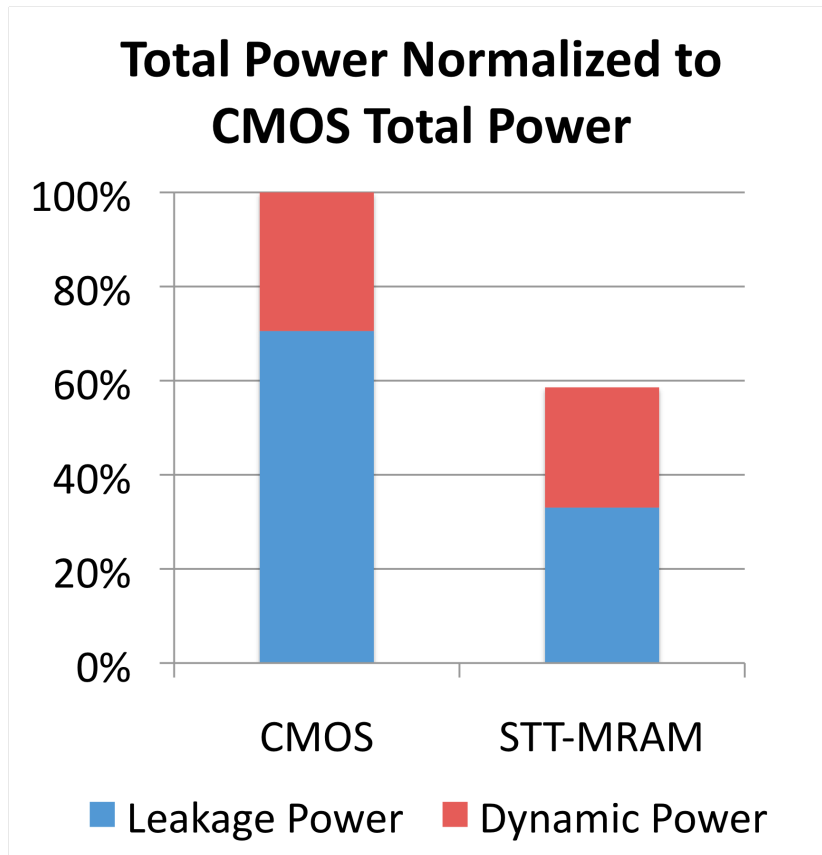
# Evaluation

# Performance

18



# Power



# Contributions and Findings

20

- New technique to reduce leakage and dynamic power in a deep-submicron microprocessor
  - ▣ Selectively migrate on-chip storage and combinational logic from CMOS to STT-MRAM
  - ▣ Use subbank buffers to alleviate long write latency
  
- STT-MRAM is an attractive low-power solution beyond 32nm
  - ▣ Dramatically lower leakage power
  - ▣ Modest loss in performance