

Lionel Torres – Univ. Montpellier, France

Merci à : G. Sassatelli, A. Gamatie, P. Benoit, P. Nouet,
D. Novo, G. Dinatale, A. Todri, A. Virazel, L. Latorre, M. Robert, G. Patrigeon, P.Y. Peneau, J. Modad, F. Ouattara, J. Lopes, O. Coi, K. Sevin



LIRMN



Current IC Integration Challenges

- Energy is <u>critical</u>
- We need more and more **Performances** for applications
- Actual technology limitations (CMOS) Integration is more and more complex 10⁹ transistors/cm2
- Actual Reliability is a problem– X% of the systems encounter an uncorrectable error per year (X ranging from 1 to 5%)





2

Technology target : CMOS < 20 nm

To Transport 1 bit \rightarrow 1pJ/mm To transport 10⁹ data – 1s (1Ghz) \rightarrow 1pJ/mm x 10⁹ = 1mw/mm 64 Bits Bus \rightarrow 64mw/mm On real IC \rightarrow several W/cm2

Calcul, Bit transition $\rightarrow 1$ aJ Calcul, 10⁹ data transition – 1s $\rightarrow 1$ aJ x 10⁹ = 1nw



- → It is better to "calculate" than to "transport" the information
- → *In computing* memory is certainly interesting
- \rightarrow Reminder : minimal energy to change 1 bit d'information K. T Ln2 \rightarrow 2,85 zJ



- Today, 50% of the silicon area of IC is memory
- Take care to energy (static)!



Figure 1: Leakage power becomes a growing problem as demands for more performance and functionality drive chipmakers to nanometer-scale process nodes (Source: IBS).

Technology evolution

Actual memories:

- SRAM for fast access
- DRAM for applications
- Flash (mass storage)

Emerging memories

. . .

. . .

- Magnetic tunneling junctions
- Phase change memory
- Programmable metallization cells
- OxRRAM

Eugosi ?

Universal memory: "Non-volatile memory"

- SRAM performance
- Size of DRAM/Flash
- Non-volatility
- Scalibility

Resistance Switching Memory

Emerging memories offer non-volatility, speed and endurance => disruption of the memory hierarchy? Conductance of magnetic metal plates is larger in the presence of a magnetic field perpendicular to the current flow

> William Thomson 1824-1907



Resistance variation attained: 2%-5% in RT



Peter Grünberg and Albert Fert



2007 Nobel Prize in Physics

 Thin stacks of FM/NM metals have seen a conductance increase of up to 100% when subjected to a magnetic field



Spin Technology



M. Bowen et al. Nearly total spin polarization...

Spin Technology



 Compatible with CMOS Non_volatile memory Switching time < 1ns writing current < 10uA-100uA density x4 vs SRAM Immune to radiations



Samsung demonstrator (8 Mbit STT_MRAM) – 2016

Motivations

- A way
 - Go towards non-volatile systems using emerging NVM⁶
 - Current NVMs issues : Speed, Dynamic energy, Reliability



Where and how to place MRAM to:

reduce total power consumption ? keep same or get better performance ? **MRAM**

FeRAM

PCRAM

ReRAM

Contributions

- 1. Evaluation of MRAM-based cache memory hierarchy:
 - Exploration flow and extraction of memory activity
 - L1 and L2 caches based on STT-MRAM and TAS-MRAM

- 2. Non-volatile computing
 - Instant-on/off capability for embedded processor
 - Analysis and validation of *Rollback* mechanism



Take advantages of MRAM

Low leakage High density Non-volatility

Mitigate drawbacks of MRAM

High write latency High write energy

NVM exploration flow



* N. Binkert et al., "The gem5 simulator," ACM SIGARCH Computer Architecture News, Aug. 2011.

** X. Dong et al., "NVSim: A Circuit-Level Performance, Energy, and Area Model for Emerging Nonvolatile Memory," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, Jul. 2012.

Experimental setup



Circuit-level analysis: Models (NVSim) & Prototype

Area



o MRAM is denser for large cache capacity

- MRAM cell size smaller than that of SRAM
- MRAM needs large transistors for write
- TAS-MRAM cache larger due to field lines

Circuit-level analysis: Models (NVSim) & Prototype

			Read		Write		Standby	
	Node	Technolog y	Latency (ns)	Energy (nJ)	Latency (ns)	Energy (nJ)	Leakage (mW)	
512kB L2 cache	45nm	SRAM	4.28	0.27	2.87	0.02	320	Λ
		STT- MRAM	2.61	0.28	6.25	0.05	23	4
	120nm	SRAM	5.95	1.05	4.14	0.08	₈₂)/8	
		TAS- MRAM	35 STT-MRAI	1.96 M≈SRAM	35	4.62	10	
			TAS-MRAM > SRAM		MRAM > SRAM		MRAM << SRAM	
32kB L1 cache	Node	Technolog y	Latency (ns)	Energy (nJ)	Latency (ns)	Energy (nJ)	Leakage (mW)	
	45nm	SRAM	1.25	0.024	1.05	0.006	22	
		STT- MRAM	1.94	0.095	5.94	0.04	3.3	
		IVIKAIVI > SKAIVI		MRAM > SRAM		MRAM << SRAM		

Case study

Quad-core architecture:

- Frequency 1GHz
- ARMv7 ISA
- Private L1 I/D
- Shared L2
- DDR3 Main memory



Benchmarks

- SPLASH-2
 - Mostly high performance computing
- PARSEC
 - Animation, data mining, computer vision, media processing





Architecture-level analysis: gem5

Read/Write ratio





D-Cache reads I-Cache writes D-Cache writes



L2/L1 access ratio

Denskursede	Number of accesses			
Benchmark	L1 cache	L2 cache		
SPLASH-2	~2 billions (0.5 billions/CPU)	~26 millions		
PARSEC	~12 billions (3 billions /CPU)	~16 millions		

Static/Dynamic energy ratio

Static energy

 $L2 \rightarrow 90\%$ $L1 \rightarrow 80\%$

MRAM-based L2

Execution time

STT-MRAM L2 (45 nm) TAS-MRAM L2 (130 nm)



Observations:

- STT shows good performance
- L2 has small impact in overall performance
 L2 has small impact in overall performance
 For TAS, 14% of penalty in average (SPLASH-2)
 Depends on applications (Cache miss rate, L1/L2 access ratio)



MRAM-based L2

Total L2 cache energy consumption



Observations:

- Up to 90% of gain for STT
- From 40% to 90% for TAS
 - Due to the very low leakage of MRAM-based cache



Extension to this work

STT-RAM designs with different data retention times [1]

	Design 1	Design 2	Design 3
Cell size (F^2)	23	22	27.3
MTJ sw time (ns)	10	5	1.5
Retention Time	4.27yr	3.24s	$26.5 \mu s$
Write Latency (ns)	10.378	5.370	1.500
Write Dyn. Eng(nJ)	0.958	0.466	0.187

Given a multi-bank STT-RAM memory, where each bank has customized retention time, how to suitably allocate data in the memory?

=> lifetime analysis for program variables to decide their mapping (see talks in «Timing Analysis» Session @ RTNS'2018) - Rabab Bouziane, Erven Rohou and Abdoulaye Gamatie

MRAM-based cache

Is MRAM suitable for cache ?

- Good candidate for lower level of cache (L2 or last level cache)
 - Up to 90% of energy gain
 - No or small performance penalty
 - More memory capacity using MRAM
 - Cache L2 is up of 20% energy consumption of overall system
- Not suitable for upper level of cache (L1) for high performance but depending of the application some gain in energy
 - Micro-architectural modifications required to mask latency
 - Not detailed in this presentation but full evaluation of cache L1 done too

Contributions

- 1. Evaluation of MRAM-based cache memory hierarchy:
 - Exploration flow and extraction of memory activity
 - L1 and L2 caches based on STT-MRAM and TAS-MRAM

2. Non-volatile computing

- *Instant-on/off* capability for embedded processor
- Analysis and validation of *Rollback* mechanism

MRAM-based processor

Normally-off computing

Two concepts:

- Instant on/off
 - Restore processor state





- Backward error recovery (Rollback)
 - Restore previous valid state



Traditional microcontrollers (MCU)



Traditional microcontrollers (MCU)



Traditional MCU

Non-volatile MCU



30-Oct-18

MCU based on STT-MTJ



MRAM-based processor



B. Jovanovic, R. Brum, L. Torres, *Comparative Analysis of MTJ/CMOS Hybrid Cells based* on TAS and In-plane STT Magnetic Tunnel Junctions, **IEEE Transactions on Magnetic**, 2014.

MRAM-based processor

First Case study: Amber 23 processor (ARM based instruction)



FEATURES

- 3-stage pipeline
- 16x32-bit register file
- 32-bit wishbone system bus
- Unified instruction/data cache (16 kBytes)
 - Write through
 - Read-miss replacement policy
- Main memory (> Mbytes)
- Multiply and multiply-accumulate operations

- Implementation of both instant-on/off and rollback (Verilog code modified)
- Duplication of the registers to emulate the non-volatility

Instant on/off



Rollback





The goal of GREAT project is to co-integrate multiple functions like sensors, RF receivers and logic/memory together within CMOS by adapting STT-MTJs to a single baseline technology in the same System on Chip as the enabling platform for M2M and M2H IoT.

Proof of concept

A full SoC based on STT-MTJ under fabrication

Full Layout

Process 180nm CMOS (TowerJazz) 200nm STT-MTJ (Spintec,Singulus)

> Die area ~23mm²



Power supply 1.8V Core / 3.3V IO

> Frequency 20 MHz

2126 NVFFs

MSS-based SoC

GREAT



Overall architecture



Operation

GREAT

Normal execution

- Binary code loaded via UART
- Program executed from either local SRAM, or external SRAM or STT-MRAM

Active/Sleep modes management

- Specific Controller
- Backup is initiated by software
- Power-off is triggered:
 - Either by an external signal « sleep »
 - Or by software
- Recovery is triggered:
 - Either by an external signal « wakeup »
 - Or futher to an event from the interrupt controller



Application scenarios

Scenario

-

Sensing

Processing

Sending



Sensing (external sensor) Analog to digital conversion Signal processing (decimation) Ciphering (TRNG) Store in memory **Minor computation** \mathbf{v} Send data (UART or DAC) Backup/Sleep Wakeup/Restore

2

Scenario

Active Energy @20MHz

Scenario 1

Scenario 2



Comparison between execution from SRAM and execution from STT-MRAM (Post-layout simulations)

Backup/Wakeup Energy



28nm projection

Wakeup time : 4.15 µs

Backup time : 4.15 µs

2126 NVFFs arranged by clusters to avoid electrical integrity issues

82 clusters

Number of clusters

- 82 (for 180nm)
- 2 (for 28nm)

Backup time @20MHz

- 4.1µs (for 180nm)
- 100ns (for 28nm)

Minimum Tsleep





28nm projection

180 nm technology





Backup energy is independent of the time spent in sleep mode

Leakage energy is dominated by SRAM

Minimum T_{sleep} to compensate the backup energy ≈ 65 ms

Minimum T_{sleep} to compensate the backup energy ≈ 641µs

Sensor node application (Agriscope)



IGREA



- Agro monitoring : plant disease, temperature, irrigation, pesticide threshold, so on ..
- -Solutions are based on 32bit processor + RF stack
- Targeted autonomy 10 years
- Applications fully compatible with our SoC

- Agriscope's application has been ported to our SoC to compare with industrial use case

Application to a sensor node



Periodic wakeup (15min)

Energy consumption in the case of a non-volatile MCU - No leakage in sleep mode - Backup energy

			1 07
Sensor events	Traditional MCU	NV MCU	Factor
0	6.03 mJ	9.09 µJ	X 663
15 (1 per min) (rain gauge)	6.08 mJ	121 µJ	X 50
180 (1 per sec) (water meter)	6.81 mJ	1.5 mJ	X 4.5
9000 (1 per 0,1s sec) (anemometer)	41.4 mJ	67.2 mJ	/ 1.6



MRAM has a high potential to:

- Certainly Reduce energy consumption
 - At cache level (sure and proven)
 - Normally-off computing
- Can facilitate some features
 - Normally-off computing / Instant on-off
 - Backward error recovery (Rollback)
- Results should be confirmed through measurements on silicon prototype !
- Important : Link with compilation and OS
- An open framework available to the community : MAGPIE

For the future....



Racetrack memory



THE ALL



THANKS !

