ASP-DAC 2011 26th January 2011 Yokohama, Japan

Non-Volatile Memory and Normally-Off Computing

T. Kawahara Central Research Laboratory, Hitachi, Ltd.

Original works cited in this presentation were supported in part by; • The "High-Performance Low-Power Consumption Spin Devices and Storage Systems" program (headed by Professor Hideo Ohno of Tohoku University) under Research and Development for Next-Generation Information Technology of MEXT,

• One of projects on "Fundamental Technologies for the Next Generation Supercomputing" by MEXT, and,

• The Japan Society for the Promotion of Science (JSPS) through its "Funding Program for World-Leading Innovative R&D on Science and Technology (FIRST Program). "

The author would like to thank to;

Tohoku University: S. Ikeda, T. Meguro, R. Sasaki, M. Yamanouchi, I. Morita, T. Hirata, T. Hanyu, and H. Ohno.

Hitachi: R. Takemura, K. Ono, T. Ishigaki, T. Yamada, A. Kotabe, S. Hanzawa, S. Yamaguchi, K. Miura, H. Yamamoto, J. Hayakawa, N. Matsuzaki, Y. Mouri, K. Ito, H. Takahashi, H. Hasegawa, Y. Goto, N. Osakabe, and H. Matsuoka.



- Introduction
- Non-volatile memories
- Spin-transfer torque RAM scalability
- Normally-off instant-on computing
- Conclusion

4 Social system needs large power

• World is power hungry, but...



5 **Necessity of low power IT**

- In 2025, five times larger power consumption is expected in Japan than that in 2005.
- Will reach nine times in the world.

(BkWh/year)



6 Sustainable and innovative society

 Solutions are not only by LSI technology, but LSI can contribute in many for sustainable world.



7 Beyond Low Voltage Operation

- Already made extensive progress for low voltage operation with high performance in last twenty years.
- Complementarily, a new kind of innovation for further.



8 Memory and Computing

- Memory systems are constructed according to a deep hierarchy based on speed and capacity. Volatile memory and nonvolatile memory are often combined
- Leads to slow start-up, long idle times, and low power efficiency.



9 Normally OFF, Instant ON Computing

- New innovation in addition to low-voltage.
- Normally OFF, Instant ON
 - -Turn off anytime when not in use,
 - but operate instantly with full performance when needed.
- Need data stored in the state of operation before any turned off, with zero power for retaining data.
 Nonvolatile RAM (NV-RAM)
- NV-RAM: Infinite number of fast write and read operations with non-volatility.
- Adaptive solutions in wide time domain, and wide area domain. LSI level to System level.

10 Example (but not ideal yet)

- Power ON time is improved to 1/9 with NV-RAM.
- Nonvolatility achieves low power also.



11 Outline

Introduction

Non-volatile memories

- Spin-transfer torque RAM scalability
- Normally-off instant-on computing
- Conclusion

12 RAM Comparison

- NV-RAM: Infinite # of write cycle, non-volatile, fast R/W.
- Magnetic RAMs are good candidates for NV-RAM.



13 Origin of Non-volatility

 Various phenomena is applicable: insulator barrier, change of structure (bistable), and alignment of electron condition (bistable).

Insulator barrier



Floating gate





Change of structure



Phase change



Ferroelectrics

Electrode

Electrode



Ionization

Alignment of electron condition



Ferromagnetism (+ Magneto resistive effect for memory)

-> Easier to achieve infinite write cycle (endurance).

14 Phase Change Memory

- Resistance change between amorphous and poly of chalcogenide material. Large resistance change.
- Good scalability.



15 Development of Phase Change Memory

- Principle was announced in the 70's.
- Production as NOR flash replacement in recent years.



16 Phase Change Memory Chips



17 Stacked Phase Change Memory

- Selecting element and the memory device are stacked.
- Selecting element by non-substrate Si, formed on wiring.



Cross sectional view Memory cell Polysilicon diode formed on wiring. Current drivability: 160uA@30nm.

18 TMR Device and Memory cell

- Resistance changes between parallel and anti-parallel.
- Use this two-state as an information bit.



TMR ratio = $(R_{AP} - R_P) / R_P$

19 Innovation in TMR device: MgO

- MgO barrier provides two breakthroughs.
- High TMR ratio and small write current.



20 SPRAM (SPin-transfer torque RAM)

- MRAM: Matured in product. Inferior in scalability.
- SPRAM: Good potential in scalability.



21 Thermal Stability Factor: *E/k_BT*

- Thermal field affects retention and disturbance.
- High E/k_BT attainable by perpendicular magnetization.
- Perpendicular TMR with typical CoFeB was reported.

(S. Ikeda et. al, Nature Materials, Volume 9, pp.721-724 (2010))



22 Development of Magnetoresistive Memory

- MRAM: Matured in product.
- R&D moved to SPRAM.



23 SPRAM Chips

Perpendicular TMR (Toshiba) 65 nm

- **Process:**
- **Power Supply:**
- **Density:**
- Chip size:
- Cell size:
- Write Current:
- TMR device:

50 uA **Perpendicular**

0.3584 μm² (84.8F²)

1.2 V

64 Mb

47.12 mm²

7.14mm SWLEGC # BLOM2 CSL/HDQ/0H4 6.6mm VCLMP/VPCT Germanatory Replica Amay 11 Suppose Arrange Poweri0M4 Meshed Power-line

K. Tsuchida, et al., **ISSCC 2010**



Modified DRAM Process (Hynix, Grandis)

- **Process:** •
- **Power Supply:**
- **Density**:
- Chip size:
- Cell size:
- Write Current:
- TMR device:

- 54 nm
- 1.8 V
- 64 Mb
 - 23.36 mm²
 - 0.041 µm² (14F²)
 - 140 uÅ

In-plane



S. Chung, et al., **IEDM 2010**



24 Progress in memory development

- Memory tied up closely to deep material and device tech.
- Modeling is indispensable to quantitative cooperation with circuitry.



25 Ex. Modeling of TMR Memory Cell (1/2)

- Physics was built into analog and mixed-signal (AMS) simulation, and simulated with CMOS circuits.
- Dependence of basic physical characteristic to memory function can be directory simulated.





Non-linear *I-R* behavior *T*-dependent *I-R* behavior

•P-state: $G_P \propto V^2$ (Simmons) •AP-state: $\Delta R/R \propto \exp(-|V|/Vh)$ • $\Delta R/R(V=0)$: function of *T*dependent polarization •Non-linear parameter *V*h : *T*independent •Numerical solution of stochastic-LLG equation •Include thermal fluctuation of macro-spin

$$\frac{\partial \mathbf{s}}{\partial \tau} = \left[\mathbf{s} \times \left(\mathbf{h}_{eff} + \mathbf{h}_{fl} \right) - \alpha \mathbf{s} \times \left(\mathbf{s} \times \left(\mathbf{h}_{eff} + \mathbf{h}_{fl} \right) \right) \right]$$
$$\mathbf{h}_{fl} = \sqrt{\frac{\alpha}{(1 + \alpha^2)}} \frac{2k_{\mathrm{B}}T}{M_{\mathrm{S}}H_{\mathrm{K}}Vo} \boldsymbol{\varsigma}_{fl}$$
$$\Delta \tau = \frac{\gamma_0 H_{\mathrm{K}}}{1 + \alpha^2} \Delta t$$
$$\Delta \tau = \frac{\gamma_0 H_{\mathrm{K}}}{1 + \alpha^2} \Delta t$$
$$\langle \boldsymbol{\varsigma}_{fl}(t) \rangle = 0, \ \left\langle \boldsymbol{\varsigma}_{\mathfrak{g}}^{\varphi}(t) \boldsymbol{\varsigma}_{\mathfrak{g}}^{\varphi}(t') \right\rangle = 2\delta_{\varphi \theta} \delta(t - t')$$

- s: normalized spin vector $H_{\rm K}$: easy-plane anisotropy field $\gamma_{\rm C}$: gyromagnetic ratio $h_{\rm eff}$: effective field $h_{\rm fl}$: stochastic field
- *M*s: saturation magnetization α : damping coefficient *Vo*: free-layer volume $k_{\rm B}$: Boltzmann's constant *T*: temperature

26 Ex. Modeling of TMR Memory Cell (2/2)

- Physics was built into analog and mixed-signal (AMS) simulation, and simulated with CMOS circuits.
- Dependence of basic physical characteristic to memory function can be directory simulated.

Transient write current and stochastic switching_{κ. Ono, et al.,}



27 For emerging memory simulation

- Need path to generate physical model with new materials. Memory tech. should handle physical phenomena.
- Estimation of phenomena between adjacent cells and the size effect are important.



"Visualization": important to""intuitive understanding of peculiarTophenomenon to the material andmgrasping the disturbance betweenofadjacent cells in the array.in

"Size dependent phenomena": TCAD necessary to handle materials more than current Si, otherwise size dependence, that is important, cannot be estimated.

28 Outline

Introduction

- Non-volatile memories
- Spin-transfer torque RAM scalability
- Normally-off instant-on computing
- Conclusion

29 Memory cell size scalability

- Transistor drivability \propto F, but TMR write current \propto F².
- Cell becomes easy to shrink with advancing feature size.



TMR size	F ²	F ²	F ²	F ²
Gate width	10F	4F	2F	F
Cell size	40F ²	16F ² -	8-6F ²	4F ²



30 Scalable 4*F*² Memory-cell configurations

• TMR memory cell has better scalability than DRAM cell due to low aspect ratio. R. Takemura, et al., IMW 2010



31 Scalable write current vs read current

Good scalability causes small read current
Disruptive reading by using dead region in write.



32 DDRx-SDRAM-compatible operation

Restore operation can be performed with PRE. NV-RAM or long-retention DRAM for super LP.



Chip configuration

Timing diagram

33 Concept of MLC-SPRAM

Series 2-TMRs with different *Ic* and *∆R* provide digital levels of resistance, possible <3F²/bit.
 Proven basic operation by measurement.



Memory cell configuration

Characteristic

T. Ishigaki, et al., VLSI Tech. 2010

34 Toward 2F²/bit and More...

2ⁿ value level can achieve n-bits per cell.
 Needs high *E/k_BT* and sufficiently small *RA*.



Memory cell circuit

3-bit/cell state transition scheme

35 Future tech.: Voltage-driven magnetic RAM

Basic research is progressing for voltage-driven magnetic RAM.

Current-driven device should face power consumption concern, voltage drop in line, and ability for current supply switch. -> Solved by Voltage-driven.



※This work was supported by Research and Development for Next-Generation Information Technology of MEXT.



Introduction

- Non-volatile memories
- Spin-transfer torque RAM scalability
- Normally-off instant-on computing
- Conclusion

37 Non-Volatile Architecture

Merits for non-volatile architecture (non-volatile computing)

- Disappear the difference between suspend and hibernation.
 - Easy to control On/Off of server due to the load.
- Assured write back operation.
 - Fast and quick restarting.

Issues

- New OS for the best use of non-volatile memory.
 - When restarting, activate the contents of memory.
 - For security issue, handle the erase of contents.
 - Same as uninterruptible computation.
- Cooperation with processor vender and OS vender.
 - Ultimately, compiler schedules On/Off.



38 In CPU and Basic Circuitry Layer

- Achieving a normally-off state as much as possible by detecting any standby operation, when the power needed is preferably lower than required for transition.
- Eliminating the power required for communication between the memory and logic circuits.
- High energy efficiency in operation and no DC power consumed in CPU's long standby state, as occurs often in most business applications.



39 Nonvolatile Logic-in-Memory Architecture

• <u>Logic-in-Memory Architecture</u> (proposed in 1969): Storage elements are distributed over a logic-circuit plane.



 Storage is nonvolatile: (Leakage current is cut off)
 MTJ devices are put on the CMOS layer
 Storage/logic are merged: (global-wire count is reduced)

-----> Static power is cut off.

Chip area is reduced.

Wire delay is reduced.

Dynamic power is reduced.

40 Design of a Nonvolatile Full Adder

- Demonstrated quick on/off operation.
- Dynamic power reduce to 1/4 with keeping performance.



T.Kawahara, Hitachi

41 High-Density TCAM Cell Design

- Halves active power-delay product with 1/100 stand-by.
- 2bit TMR storage merged into CMOS logic.



T.Kawahara, Hitachi

42 Design of a Compact Nonvolatile FPGA

• TMR device/logic direct combination by differential current-mode tech. simplifies the circuit, reduces power.



T.Kawahara, Hitachi

43 In Main Memory and System Layer

- For instant-on, keeping the status in main memory in standby without power dissipation.
 - Simplified control circuit because the refresh operation is minimized.
 - Battery backup in cache logging area is eliminated.
- Hot standby configuration (speed) with deep standby mode (power) is attainable.
 - Achieving system reliability for mission-critical systems (high SLA) with low-power.



44 Ex. In Parallel Processing System



45 Power-down during Parallel Processing?



46 **Power-down Controllability**



47 Power Control in Parallel Processing

- High-performance processor capacity is underutilized.
- Can reduce power. Need automatic implementation tools.



48 Non-Volatile Architecture

- Operation of sub-processor has some patterns.
 Need tool that optimizes these to reduce power in the actual flow of system.
- Advancement going on for logic level normallyoff and instant-on operation.
 - Need tool that automatically optimized for logic synthesis and layout included power supply line
- Need merged with event-driven, asynchronous computing.

49 Normally-off, Instant-on Computing

• IT equipment to be turned on instantly and be resumed the state prior to the interruption.

50 Sustainable and innovative society

- Low power innovation is indispensable for sustainability.
- Normally-off, Instant-on computing: Harmonizing the power control tech. and the information tech. for our social innovation.

51 Non-Volatile RAM: a Key Component

NV-RAM: infinite number of fast write and read operations as well as non-volatility.

Spin-transfer torque RAM (SPRAM)

32Mb Chip (VLSI 2009, Hitachi/Tohoku U)

2Mb Chip (ISSCC 2007, Hitachi/Tohoku U)

52 Outline

Introduction

- Non-volatile memories
- Spin-transfer torque RAM scalability
- Normally-off instant-on computing
- Conclusion

53 Conclusion

Normally-off, Instant-on Computing: a way to allow computing equipment to normally be turned off when not in use but be able to turn on instantly with full performance when needed.

Non-volatile memory: any internal status of computation is memorized at any time before the power is turned off without consuming power. NV-RAM: Infinite number of fast write and read operations as well as non-volatility.

Tools for physical analysis and implementation with circuits are important for the deign of emerging memory.

Tools that clarify internal computation status and manage fine grain power control combined with NVRAM are necessary.

Harmonizing power control tech. and information tech. for the social innovation.