

Emerging Memory Technologies

KARL-MAGNUS PERSSON



Karl-Magnus Persson | Nanoelectronics



Future CMOS Development Irrelevant?



Technology Benchmarking





Technology Benchmarking







Technology Benchmarking







IRDS Review





Research Trend I – Stacking Circuits and Mem

- Performance predictions of ReRAM with large scale circuit simulations using calibrated models
- Study shows benchmarks of a contemporary Intel Xeon Phi system VS a system with CNT-cores and STT-MRAM + 3D RRAM
- Proposed system shows up to 1000x gains in combined power and speed









Research Trend I – Stacking Circuits and Mem



Conventional CPU is idle 97% of the time!!

M. Aly et al – IEEE computer 2015



Machine Learning Hardware Implications

- Iterative re-programming of memory
- Large data sets → Limited by read/write of none-volatilememory (NVM)

Hardware Challenges

- Component improvement stagnated Moore's law has halted
- NVM technologies 10,000x slower than computing
- Separate compute and memory circuitry infer large inefficiencies

Possible solutions

- New methods → co-integrated circuits and memory in 3D, introducing new materials
- True neuromorphic hardware → synaptic networks using computational units



NANO

ELECTRONICS

GROUP



Limitations in NAND Flash



- The upside
 - None volatile
 - 3D integrated with 128 layers
 - Minimal feature size (F) down to 5 nm
- The downside
 - Read in ns but write in ms
 - Further scaling of the dielectric leads to electron leakage





Bit Cost Scalable NAND Flash

Schematic structure of planar NAND Flash

https://www.storagenewsletter.com

Resistive Memory





RRAM Mechanics

NANO ELECTRONICS GROUP

- Unipolar Switching
 - Joule heating
 - Ion movement due to diffusion
 - Vacancies dissolved into oxide
 - Occur at the center of the dielectric (hottest)
 - Large reset current required (most cases)
- Bipolar Switching
 - Bi-polar-bias requirement
 - Ion movement drift due applied field
 - Vacancies absorbed into metal
 - Occur at the oxide/metal interface



Overview of different ReRAM Technologies



- Anode filament, oxygen vacancies form conductive path
- High endurance, 1^{12} cycles at device level
- 3D compatible
- CBRAM
 - Cathode filament, bridging with metal ion movement
 - Similar structure to RRAM
 - Endurance questionable (finite number of switches)
- PCRAM
 - Phase change memory, a flash heating switches dielectric film between amorphous and crystalline state
 - Endurance questionable (finite number of switches)
- STT-MRAM
 - Spin-transfer-torque magnetic RAM
 - Changing orientation of spin changes the conductivity
 - Advanced material stack, 3D compatibility unlikely





NANO

ELECTRONICS

GROUP

Overview of different ReRAM Technologies

- Anode filament, oxygen vacancies form conductive path
- High endurance, 1^{12} cycles at device level
- 3D compatible
- CBRAM
 - Cathode filament, bridging with metal ion movement
 - Similar structure to RRAM
 - Endurance questionable (finite number of switches)
- PCRAM
 - Phase change memory, a flash heating switches dielectric film between amorphous and crystalline state
 - Endurance questionable (finite number of switches)
- STT-MRAM
 - Spin-transfer-torque magnetic RAM
 - Changing orientation of spin changes the conductivity
 - Advanced material stack, 3D compatibility unlikely



NANO

ELECTRONICS

GROUP



Stanford Memory Trends, H.-S. P. Wong et al <u>https://nano.stanford.edu/stanford-memory-trends</u>

Current (A)

RRAM – BE Metal Considerations





RRAM – BE Metal Considerations







16







RRAM Oxide and TE Properties



RRAM – Area Scaling



- RRAM switching is ideally independent of device area as only one filament forms
- Area dependence is instead partly coupled to self-capacitance, and a reduction in parasitic current discharge
- However, the probability to form a filament increase with area
- To reduce increased forming voltage and spread in the distribution, surface roughness and material quality at the interfaces are of crucial importance

Y. Y. Chen – ME 2013

Ann Chen – Globalfoundires 2013



RRAM – Area Scaling







H.S. Philip Wong – IEEE Proceedings 2012

3D RRAM - Architectural Concepts





- Bit performance improves with # layers
- Most scalable/cost efficient

VRAM type II

HRRAM – Commercialized



HRRAM

- Most simplistic
- Superior performance due to low RC interconnects
- Least cost efficient

Intel 3D Xpoint

- 3x faster than NAND flash
- 5x more expensive
- Switching mechanism unknown
- Limited stacking potential due to diode selector



Vertical RRAM – Lithography Free

Litho-free formation of a stair-case structure

Tanaka et al – VLSI symposium 2007

Karl-Magnus Persson Nanoelectronics

RRAM - Oxide Thickness Scaling

- Scaling the dielectric necessary to reduce minimum feature size
- Surface roughness affect the spread of the performance distribution
- Etched out vertical pillar have a smoother surface than deposited metal

Considerations on Scaled 3D Arrays

- Large arrays require MOSFET selectors to reduce leakage
- Vertical geometry allows for more aggressive thickness scaling as it reduces roughness
- Simulations show metal plane thickness will limit array size due to resistance of the vias, sub 6-nm metal is highly resistive
- Graphene and other 2D materials way become a viable way forward for large scaled arrays

Dense, fast, and energy efficient memory opens for computations directly in the memory

- Permutations are made along the vertical pillars
- NOR and NAND operations can be accomplished with propagation of specific pulse trains
- In memory computations implementations could be task specific, drastically reduces integration area

Research Trend II – In Memory Computations

TiN/Ti (50 nm)

50 nm

7 TiN (20 nm)

Research Trend III – Hyper-Dimensional Vectors

- 10,000 bits/vector, any type of object can be represented
- Orthogonality factor between vectors determine object resemblance
- Doing in-memory reduces energy consumption slightly, but lowers chip area drastically (660x for a 2k-vector)
- Error resilience of HD computing mean that a wide range of RRAM technologies can be used
- Tricks such as generating random vectors with a high percentage '0'-s may be a way to reduce the power

0.605 (%) E 0.6 28-nm node P Digital Accuracy Consumption 5.0 70 70 70 70 Encoding 0.335 0.318 60 n-VRRAM 52.2% HD (2 kb) Recognition 0.160 PERM HD (10 kb) 0.0 Dec Com 20 SA.MUX.PG 10 20 30 40 Hard Errors (%) 3D LP 3D Digital VRRAM Digital VRRAM HD = 2 kb0.1 HD = 1 kh1E-3 0.01 10 Proportion of Hard Errors (%) Accuracy (%) 5 Energ RRAMs having wide-range 60 endurance can be used ormalized Sparsity: % of '0 HD Vector Recognition -0-2 kb 40 Stuck-at errors after -O= 4 kb endurance failure -∆-6 kb 20 10 kb 80 85 70 75 1M 100k 10k 100 10M 1k Sparsity in 3D VRRAM Array (%) RRAM Endurance (cycles/cell)

26

IBM TrueNorth - 1 unit chip layout (2014)

ELECTRONICS

GROUP

Implications for Neuromorphic Computing

IBM TrueNorth

- SyNAPSE DARPA funded initiative to simulate the brain
- Dedicated neuromorphic hardware
- 4096 computational units (1 unit pictured)
- Memory (SRAM) occupies about 40% of the chip area

RRAM

- 3D RRAM with 128 layers
- 64 Tb per chip

1 unit chip layout (2014)

SRAM ~ 120-140 $F^2 \rightarrow$ 1T1R RRAM 20x smaller!!

Stanford Research: CNTs with 3D RRAM

~18 lithography steps + transfer and etch procedures

5 wafers in Lund

7nm-back-gate

2D Via-hole RRAM – Accepted at DRC

Fig. 1. (a) Schematic illustration of the RRAM structure post ALD of the electrical TE and the coule, as well spin-courting, baking, and patterning of the lifting layer that has been opened down to the active ERAM interface area and to the TE. (b) Four the electrical BE definition, and (i.), a magnification of the vishole that constitutes the RRAM device.

Fig. 2. 10 consecutive stratching-cycles for a RRAM device convergenzing to sample A. The writeking is performed without any current limiting device in series and achieves a 10x high to low revisionce ratio while seriesking at balaw.2500 mV.

Fig. 3. (a) Low-woltage switching surves showing set (negative voltage) and reset (positive voltage) for these different devices of the different types, and where the oursest is limited by the self-compliant TO electrode, (b) Devices biased at large positive voltages to the point where they from a reverse filament. Although they do not break due to the nature of the ITO, they have to go darough emanutive reresenting to be able to works the before

Fig. 4. (a) A baseled of the RFF voltage distribution, where each loss constrains the data from 10 different devices. The RFF is lower for the device B, indicating that the film has more varancies accusationted or the TO-4RO, interfaces (b) A baseled of the facturing voltage distribution fulfillences between the different samples.

	Fig. 5. (a) Schematic vacancy distributions for the different		TAJ	NAL O	AUDE DEPOS	TTEN.
	studied devices, where more vacancies at the electrical BE	Sample		T(C)	ALD rycles	tux (mett
		α.		,200	.27.102	3.2
2.4	are likely to reduce the RFF	31		300	27.100	32
	voltage.	E	•	300	22 FE/ 5 AL	12

Lund RRAM Roadmap – 5-Å-Al/22-Å-Hf

NANO

ELECTRONICS

GROUP

Karl-Magnus Persson | Nanoelectronics

3<u>3</u>

NWFET-RRAM 4x4 Cross-bar Array

FET Controlled RRAM Measurement – NW1

NANO

ELECTRONICS

GROUP

RRAM Oxide and TE Properties

Compound	Stable condition*	Structure	Heat of formation of compond (eV)	Formation energy of oxygen vacancy** (eV)		Band gap in perfect	Defect level (eV)		
				q = 0	q = +2	cystial (cv)	$\varepsilon(+2/0)$ (eV)	$\varepsilon(+2/+1)$ (eV)	$\varepsilon(+1/0)$ (eV)
Al ₂ O ₃	hypo. NP	bixbyite	-16.15	7.14	0.18	6.09	3,48	3.80	3.16
	normal	corundum	-16.35	7.08	1.32	6.90	2.88	3.07	2.69
In ₂ O ₃	normal	bixbyite	-7.81	1.53	0.29	1.05	0.62	0.56	0.67
	HP	corundum	-7.67	1.63	0.23	1.00	0.70	0.68	0.72
SnO ₂	normal	rutile	-6.04	3.49	2.41	1.26	0.54	0.50	0.58
	HP	fluorite	-5.05	-0.45	-0.41	0.29	-0.02	-0.05	0.01

ITO RRAM Considerations

💷 Oxygen vacancy 🔲 Conductor

Sze Group

ITO-NWFET-RRAM Memory Technology

NANO

ELECTRONICS

GROUP

Future of Computing

- Buzzwords + buzzwords
- Contemporary memory technologies are not on par with the development in CPU speed
- 3D stacking circuits and ReRAM could potentially improve efficiency for data intense computing by 1000x
- Rapidly growing research area due to the promise of neuromorphic computing
- Enormous space for innovation and new startups, specialized hardware to complement the software needed

39

