

swiss scientific initiative in health / security / environment systems

Approximate Computing with Unreliable Dynamic Memories

ICYSOC

S.Ganapathy*, A.Teman*, R.Giterman[†], A.Burg*, and G Karakonstantis[‡]

*Telecommunications Circuits Lab (TCL), Ecole Polytechnique Fédérale de Lausanne, Switzerland

⁷ENICS Lab, Faculty of Engineering, Bar Ilan University, Israel

⁷High Performance and Distributed Computing, Queen's University Belfast, United Kingdom





Large amount of on-chip memory required for enhancing system performance Reduced SRAM cell dimensions for higher density causes increase in parametric variations

Low Vdd Operation Tremendous power reduction Reduced Vdd lowering SRAM cell noise margin increasing susceptibility to failures

Low-Latency Access Faster access improves system performance Requires reduced timing margins increasing susceptibility to timing failures

	1TIC eDRAM	6T SRAM	Gain Cell eDRAM
Low-V _{dd} Margin	Poor	Poor (ratioed)	Good (ratioed, gain function)
Cell cap	<10 fF	Irrelevant	$\sim 1 \mathrm{fF}$
Cell Size	15F ²	130F ²	65F ²
Process	Trench cap	Logic Compatible	Logic Compatible
Access	2 ns	< 1 ns	2 ns



ÉCOLE POLYTECHNIQUE

FÉDÉRALE DE LAUSANNE

Bar-Ilan University

Belfast

Queen's University

FNSNF

Figure: Probability Density Function of Data Retention Time of 2T GC-eDRAM Bitcell in a 28 nm FD-SOI process. The plot was extracted for 100k Monte Carlo samples at 25C.(inset): Schematic and typical operational waveforms of the 2T Gain Cell.

MEMORIES FOR APPROXIMATE COMPUTING PARADIGM

Exploiting the inherent error-resilience of certain class of applications to trade-off quality of result for power and performance

Allows for imprecise computations as long as quality is acceptable

Tolerable Imprecision Exploited by:

"Good Enough"



Relax reliability requirements of memories and provide only sufficient amount of protection based on application requirements Tolerating failures also helps improve memory parametric yield



Impact of memory failures studied for 2 machine learning benchmarks - K-Nearest Neighbours (KNN) and ElasticNet.

RTD 2013

Measured R² and Score for ElasticNet and KNN respectively.

Elastic net

0.85

K–Nearest Nieghbors



Redundant Input Data Sets

Varying Human Approximations **Perception Levels**





MINIMUM REFRESH GC eDRAM FOR APPROXIMATE COMPUTING

eDRAM Refresh Strategies





Cumulative Distribution Functions of the data retention times of 2T GC-eDRAM bitcells for various process technologies: 180 nm CMOS, 65 nm CMOS, 28 nm FD-SOI, and a estimated median node.

Availability does not have much variation at mature nodes compared to advanced nodes. By tolerating as much as 177 errors in 16kB memory, availability can be increases between 18% and 50% when moving from 65nm to 28nm node

Memory availability of GC-eDRAM arrays in various technology nodes: 180 nm CMOS, 65 nm CMOS, and 28 nm FD-SOI, and an estimated median node. Availability is shown for traditional worst-case design ("No Failures") and for error-tolerance of Perr $< 10^{-3}$ for memory bank sizes and operating frequencies typical

Since failures can be tolerated, the refresh rate is adjusted to

For 65nm node, this potentially reduces the refresh power by as much as **5X** and for other technologies, there is a minimum



"Variation Aware Performance Analysis of Gain Cell Embedded DRAMs" W.Zhang et.al. ISLPED'10 "Approximate Computing with Unreliable Dynamic Memories" S.Ganapathy et.al. NEWCAS'15 "Energy versus data integrity trade-offs in embedded high-density logic compatible dynamic "Quality Programmable Vector Processors for Approximate Computing" S. Venkataramani et.al. MICRO'13 memories" A.Teman et.al. DATE'15 "Replica technique for adaptive refresh timing of gain-cell-embedded DRAM" A.Teman et.al. TCS II'14

Retention power savings during low activity periods for GC-eDRAM rrays at different nodes, operated with an error tolerance of Perr $< 10^{-3}$