



IBM Systems and Technology Group - SRDC

Beyond Scaling

– Teaching the Old Dog some New Tricks!

Subramanian Iyer,
ssiyer@us.ibm.com

What I'd like to discuss with you today:

The Limits of Classical Scaling
Strain engineering, Hi k

On-chip memory
SRAM & DRAM

Autonomic Chips

Laser vs electronic Fuses
BIST and BISR
Self monitoring and self
repairing chips

What
about 3D

Power Supply
decoupling

The Humble
capacitor

There is only one reason to scale:

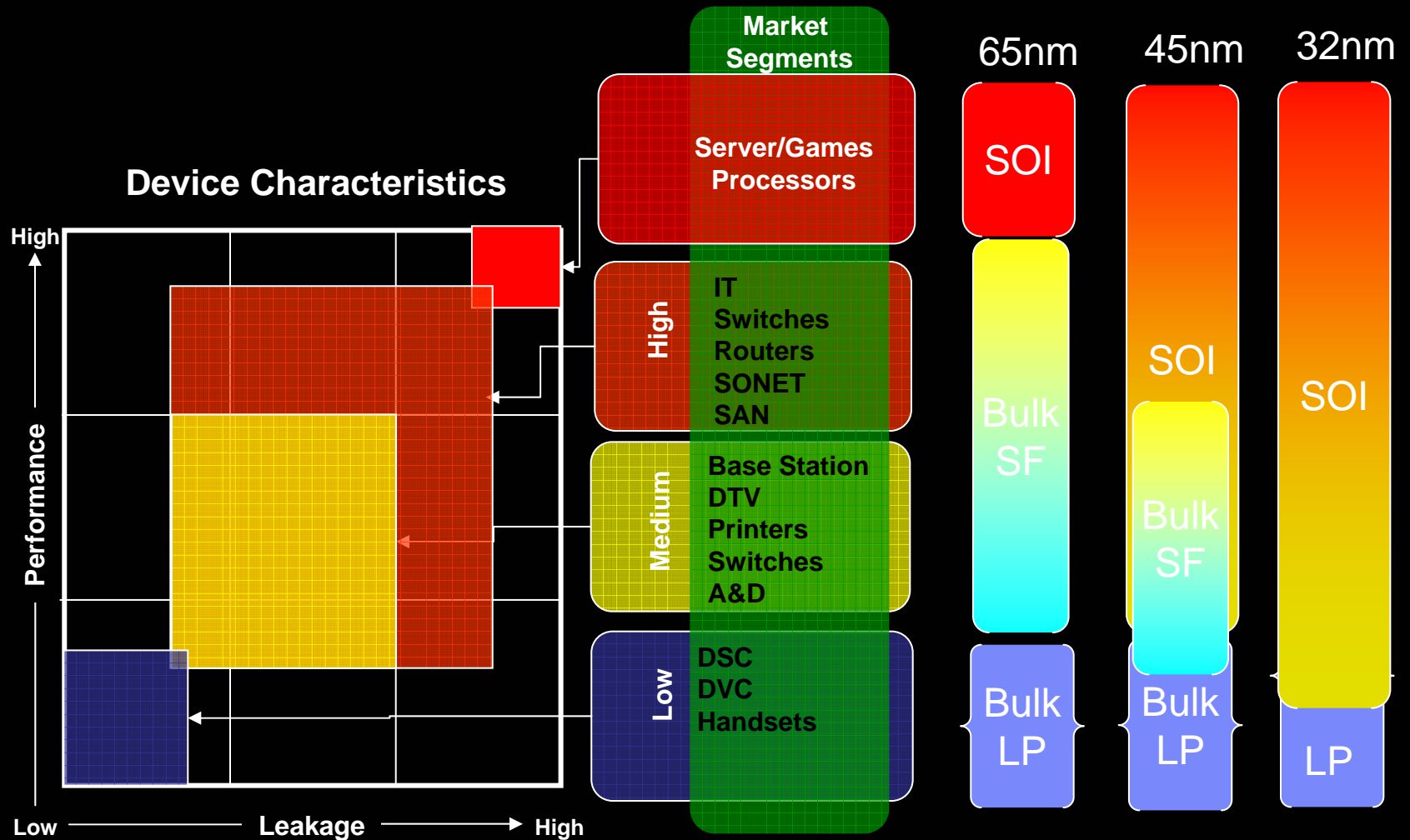
Performance,
Power

Miniaturization

Productivity: Perform more function at less cost



What are these functions ?



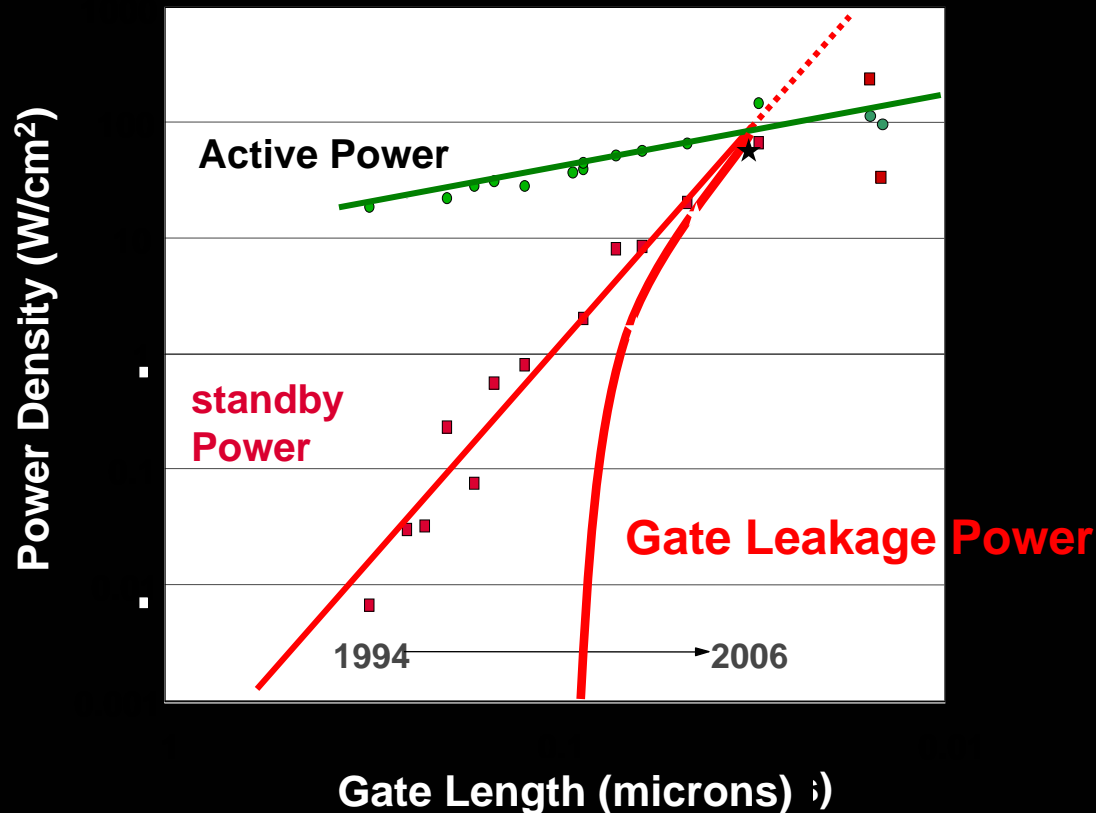
Customers need the flexibility to mix high performance and low leakage devices on the same die to optimize power and performance

Where Are We Going with Scaling?



- Historically constant field scaling came to a grinding halt at 90nm
- Stress engineering saved the day
- At 32nm HiK is a game changer
- It is designed to allow us to mix higher performance and low leakage on the same die at reasonable cost
- At 22nm further innovations are called for

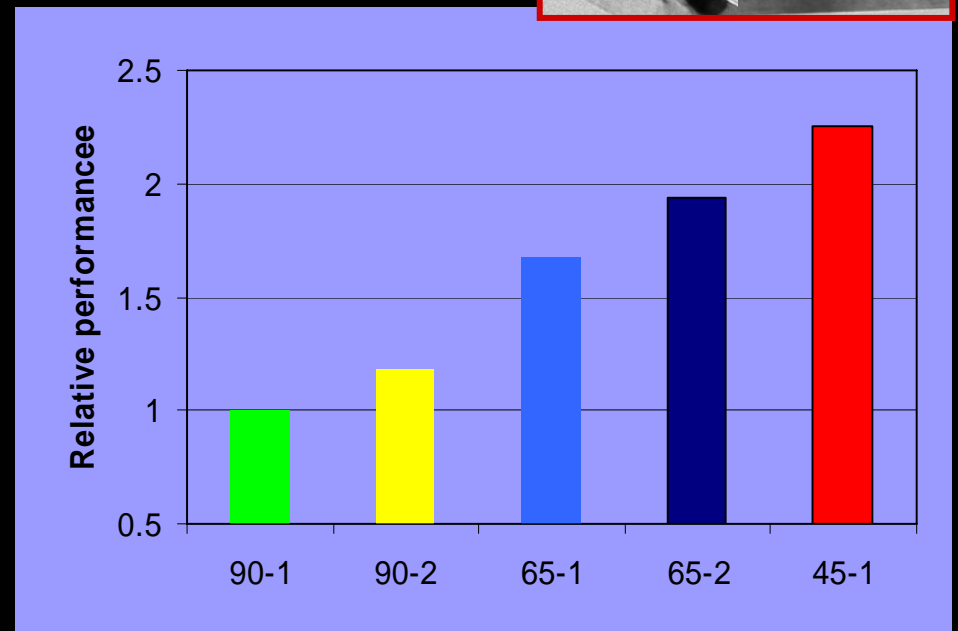
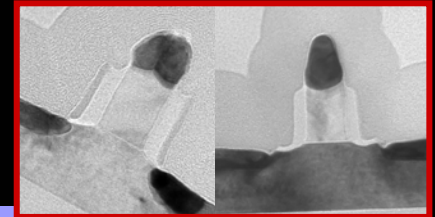
Power Density Keeps Growing !



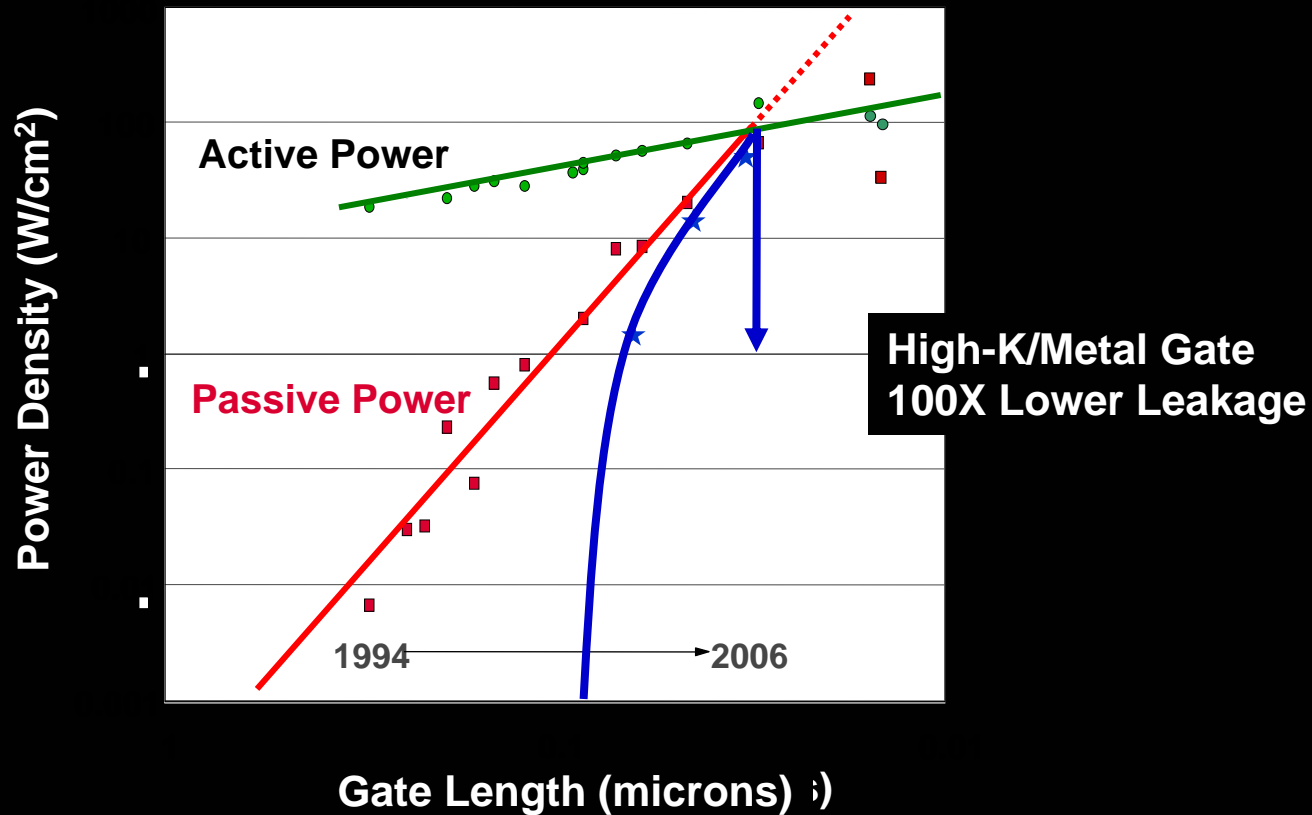
- Non-scaling causing power density to increase
- **Gate leakage power growing exponentially !!**

Strained Silicon: Achieving Relentless Performance Improvement

- Dual Stress Liner (DSL)
- Stress Memorization Technique (SMT)
- Embedded SiGe (eSiGe)
- Stress Proximity Technique (SPT)



Basic Strain	DSL	DSL SMT	DSL SMT eSiGe	DSL SMT eSiGe SPT
--------------	-----	------------	---------------------	----------------------------



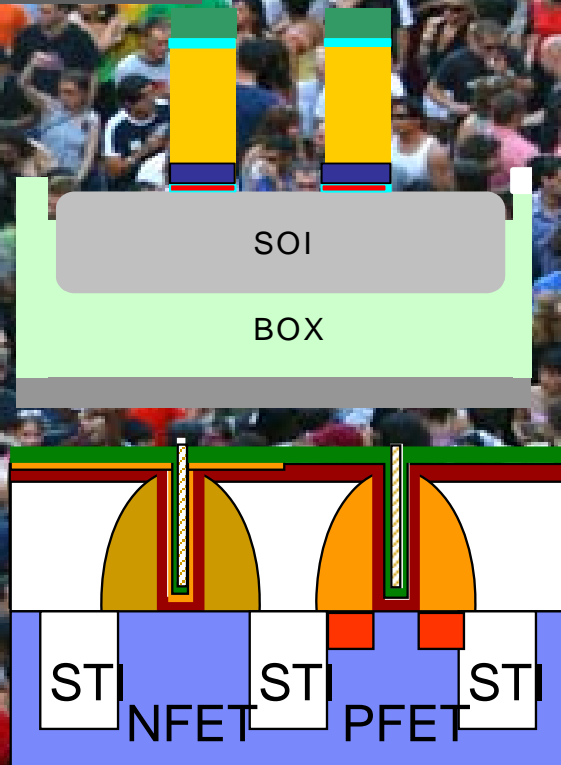
- High-K/Metal Gate reduces gate leakage power by 100X



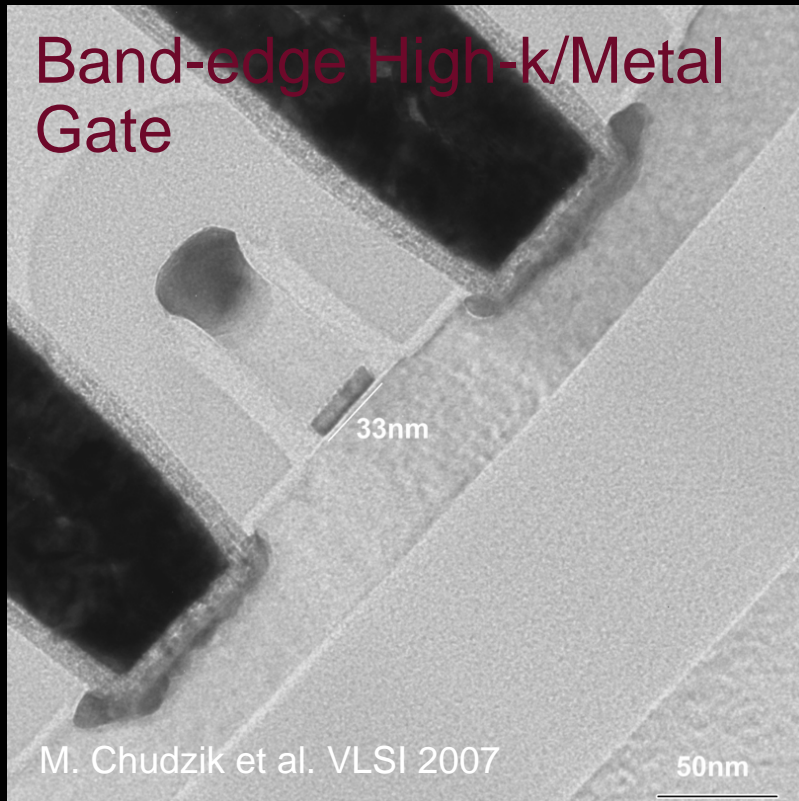
Where are Hi K materials when you need them ?

• Material Issues

- Compatible gate materials
 - Band edge vs Near Bandedge
- Performance vs Power Tradeoffs
- Temperature Stability
- Integration Schemes
- Scaling channel length



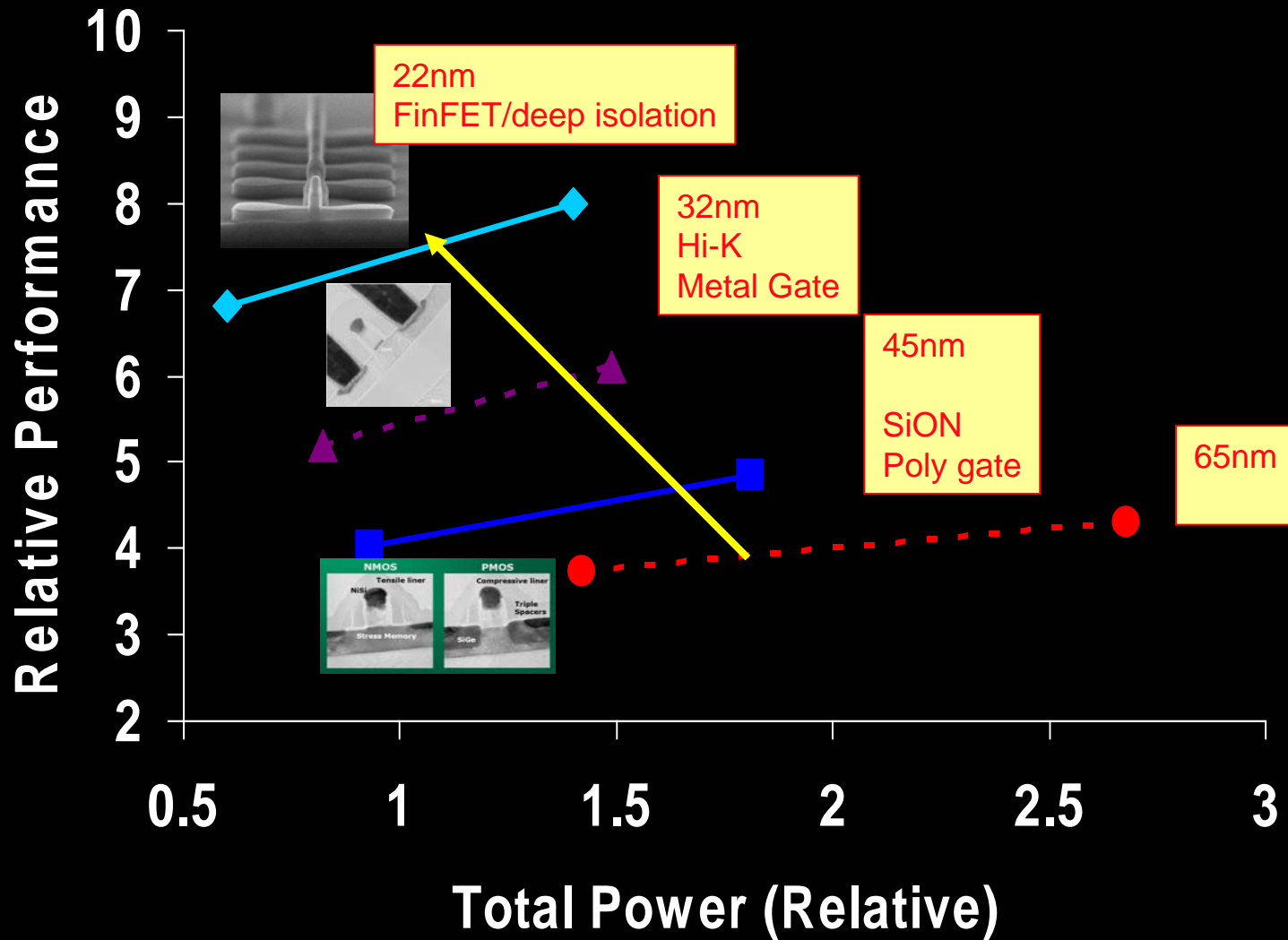
Scaling to 32 nm and beyond



- **TEM cross-sectional images of band-edge NFET HKMG stacks after full build showing a physical gate length of <40nm**

- Hi K allows us to return to scaling as we knew it
- Lower leakage devices will be possible by adjusting threshold
- But there are better ways of adjusting threshold
- So we will be able to offer both high performance as well as low leakage devices with a single gate dielectric without suffering from gate leakage limitations

But you must be able to scale the channel length !!



Key Issues Facing our Industry

- Scaling for performance is getting more difficult and very expensive
- Chip Power – both standby and active are getting to be un-manageable
- Manufacturing cost is becoming huge while more and more parts are being commoditized
- Technology development is getting very expensive

Net: Value needs to come from system concepts applied to ICs

So how did you like the ride ?

455 HP! This car moves!!!

But the gas tank is too small !! And my Kids don't fit in the back seat !!

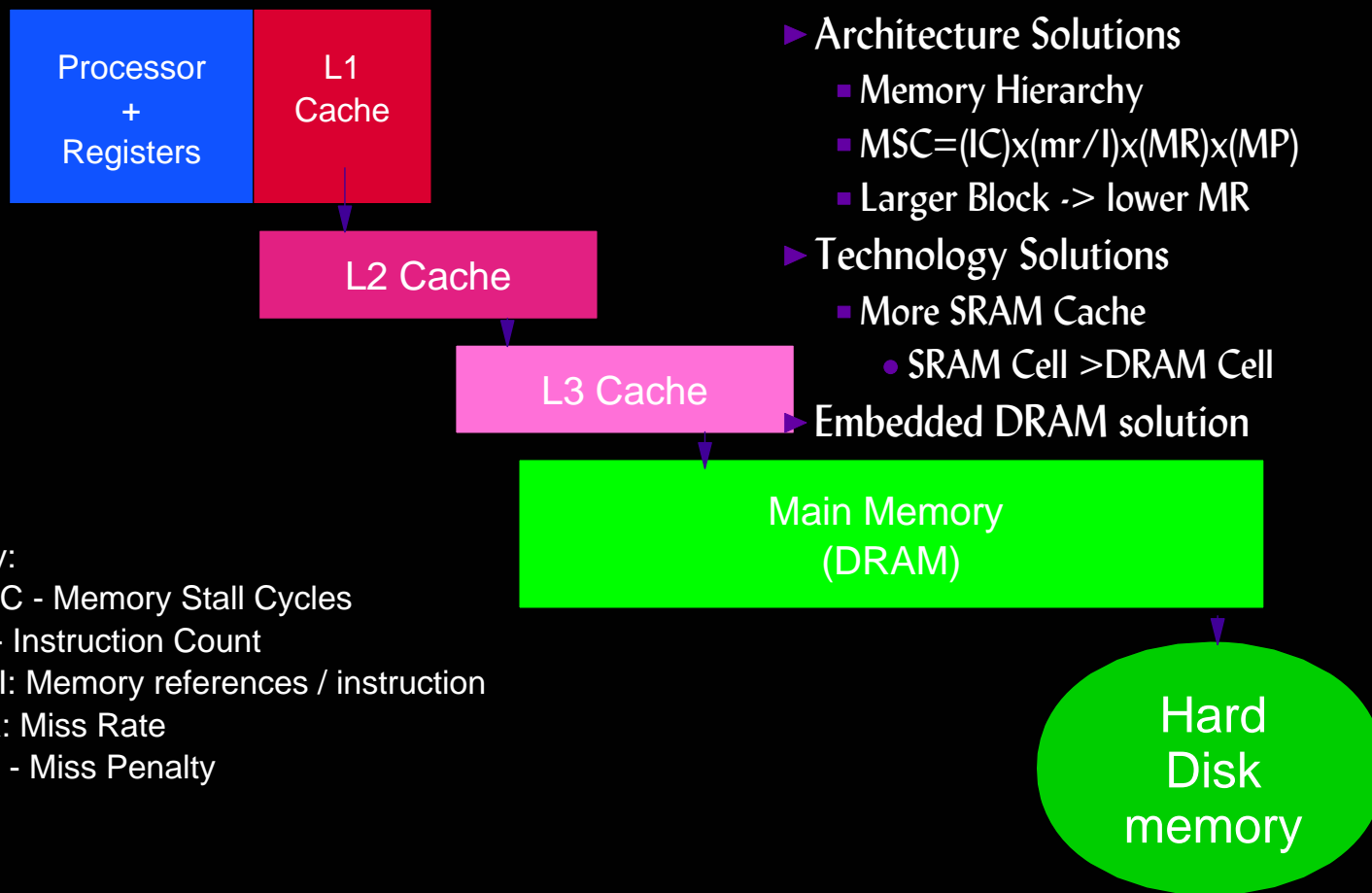
Performance Auto



Going beyond Scaling and New Materials

- Scaling, strain engineering, and improved materials will continue to improve performance , though at diminishing rates and certainly with diminishing returns
- A combination of voltage supply reduction, power budget constraints and design IP migration suggest that the days of dramatic raw performance gains are over
- **Performance must come from elsewhere**

Memory Hierarchy in a Modern Processor

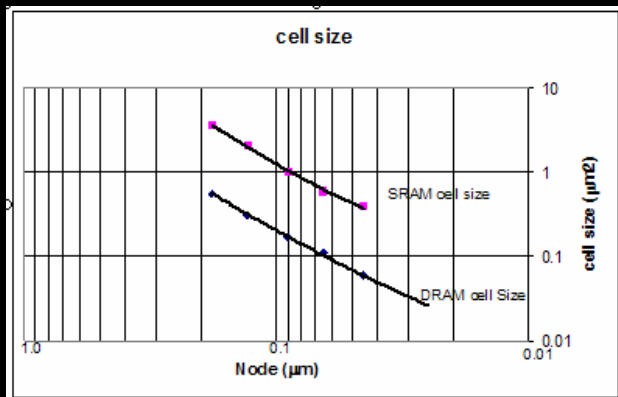
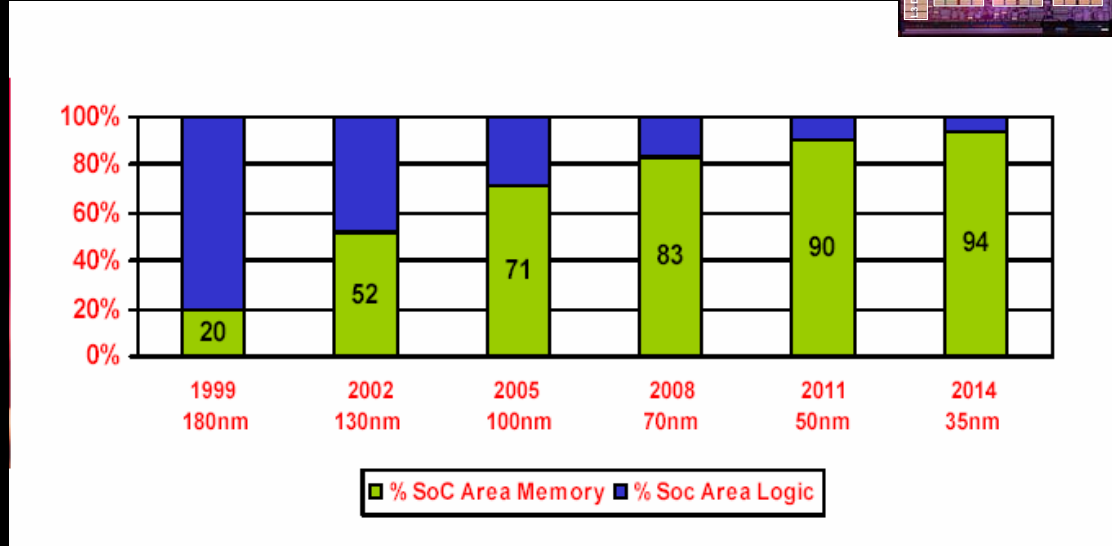
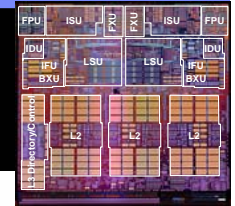


Key:
 MSC - Memory Stall Cycles
 IC - Instruction Count
 mr/l: Memory references / instruction
 MR: Miss Rate
 MP - Miss Penalty

At the system level, you never have enough memory!

How do you integrate large amounts of dense memory on chip ?

DRAM cells are 6X smaller than SRAM cells
 about 3.5X as dense / Mb
 And about half as fast



Integrate DRAMs on a logic chip
 additionally they

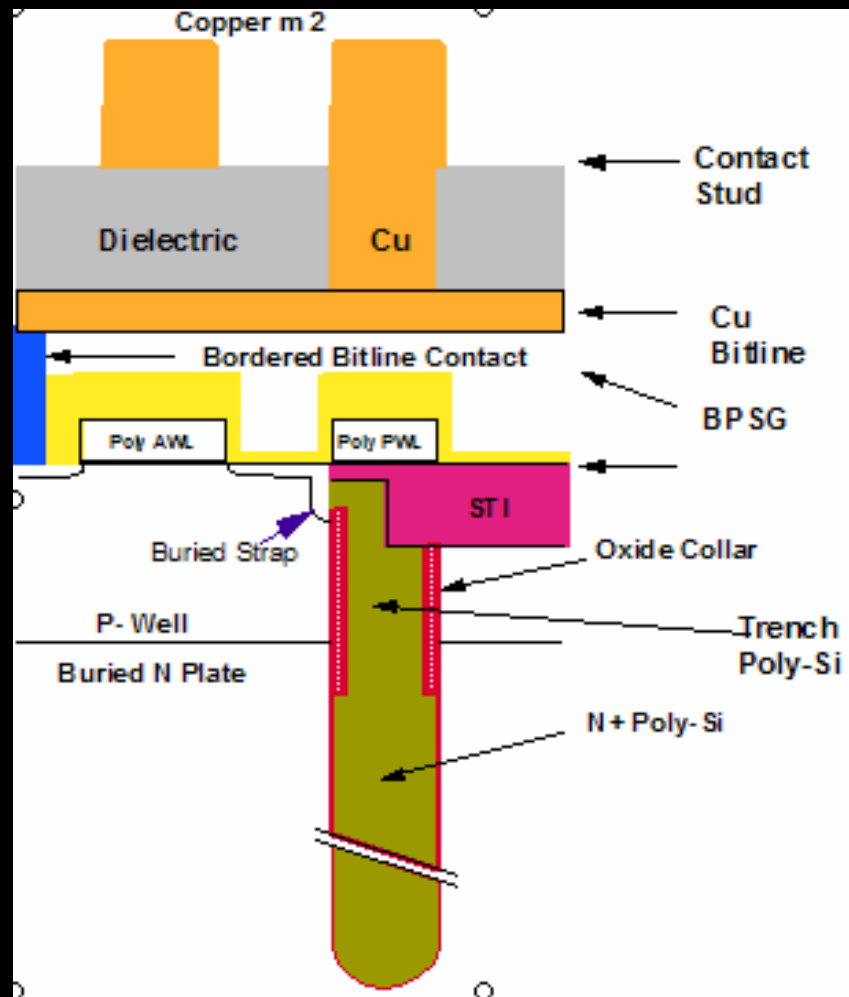
- Save power
- Have 10^4 lower SER
- and are more stable

The SRAM conundrum

- SRAMs are not aging with dignity
- V_{min} stability due to
 - Dopant fluctuation
 - Gate oxide leakage
- Solutions
 - Slow down scaling of both area and voltage
 - Multiple cells with unique devices and layouts
- Power
- Soft Error Rate continues to rise

Integrating DRAM and Logic

- Integrate with Logic without impacting logic Performance, Reliability or Yields
- The Deep Trench process is intrinsically logic friendly
- Capacitor fabricated first
- No perturbation of the remaining logic flow
- Significant Process & Test know-how in DRAMs using deep trench technology

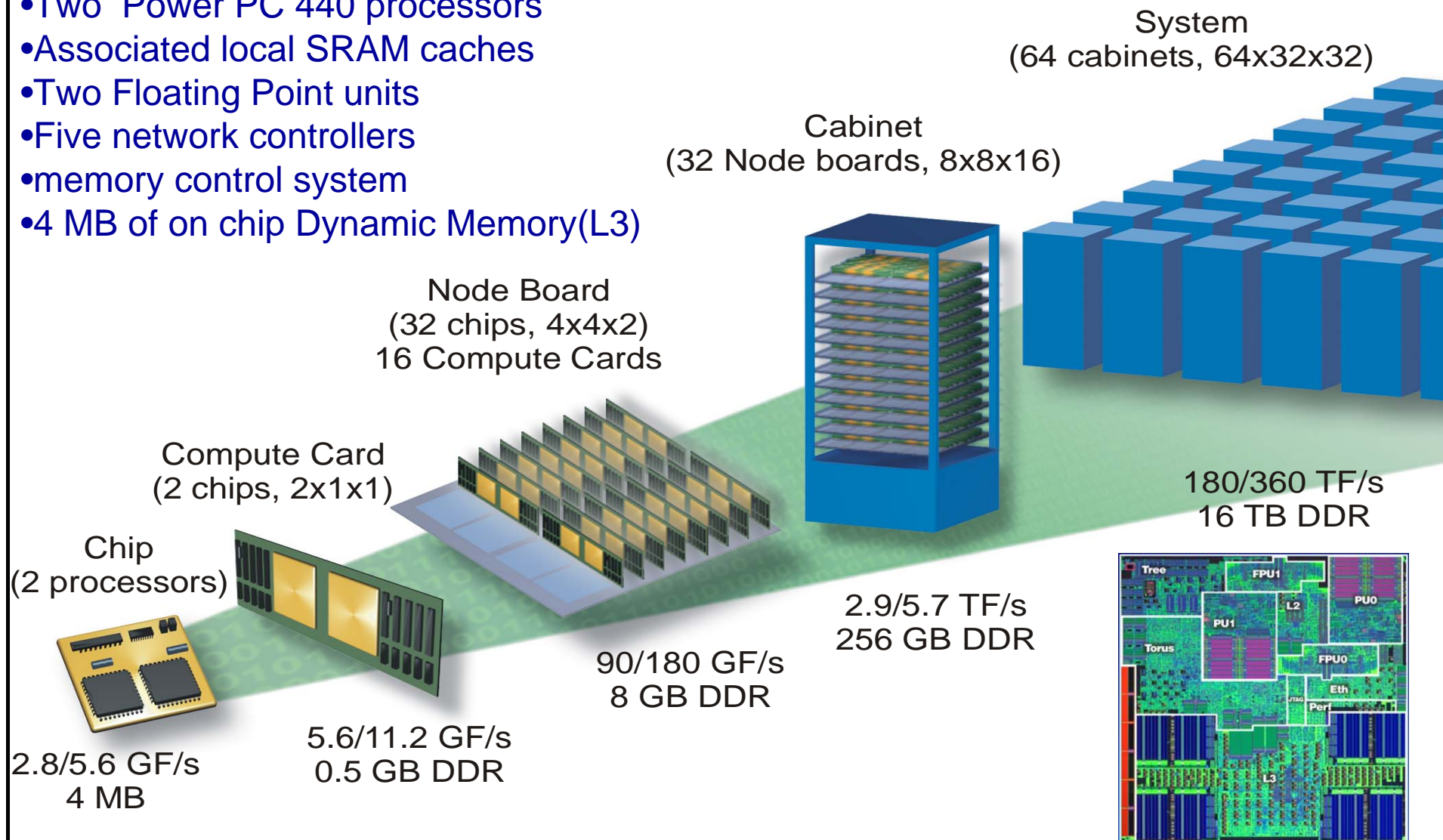


Over five generations, embedded DRAM has adapted to logic technology resulting in simpler processes and significantly higher cell performance

Blue Gene – world's fastest supercomputer in 0.13 μ bulk CMOS

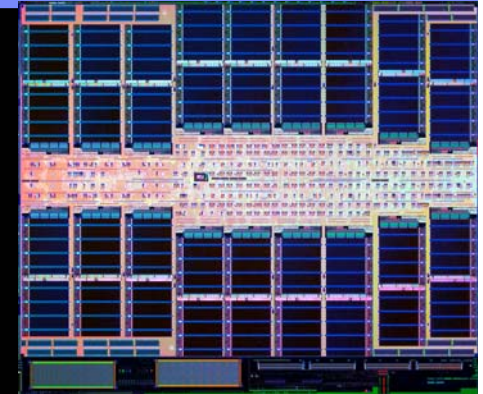


- Two Power PC 440 processors
- Associated local SRAM caches
- Two Floating Point units
- Five network controllers
- memory control system
- 4 MB of on chip Dynamic Memory(L3)



Power 5 – L3 cache

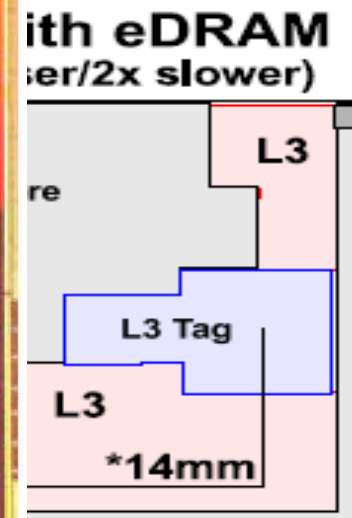
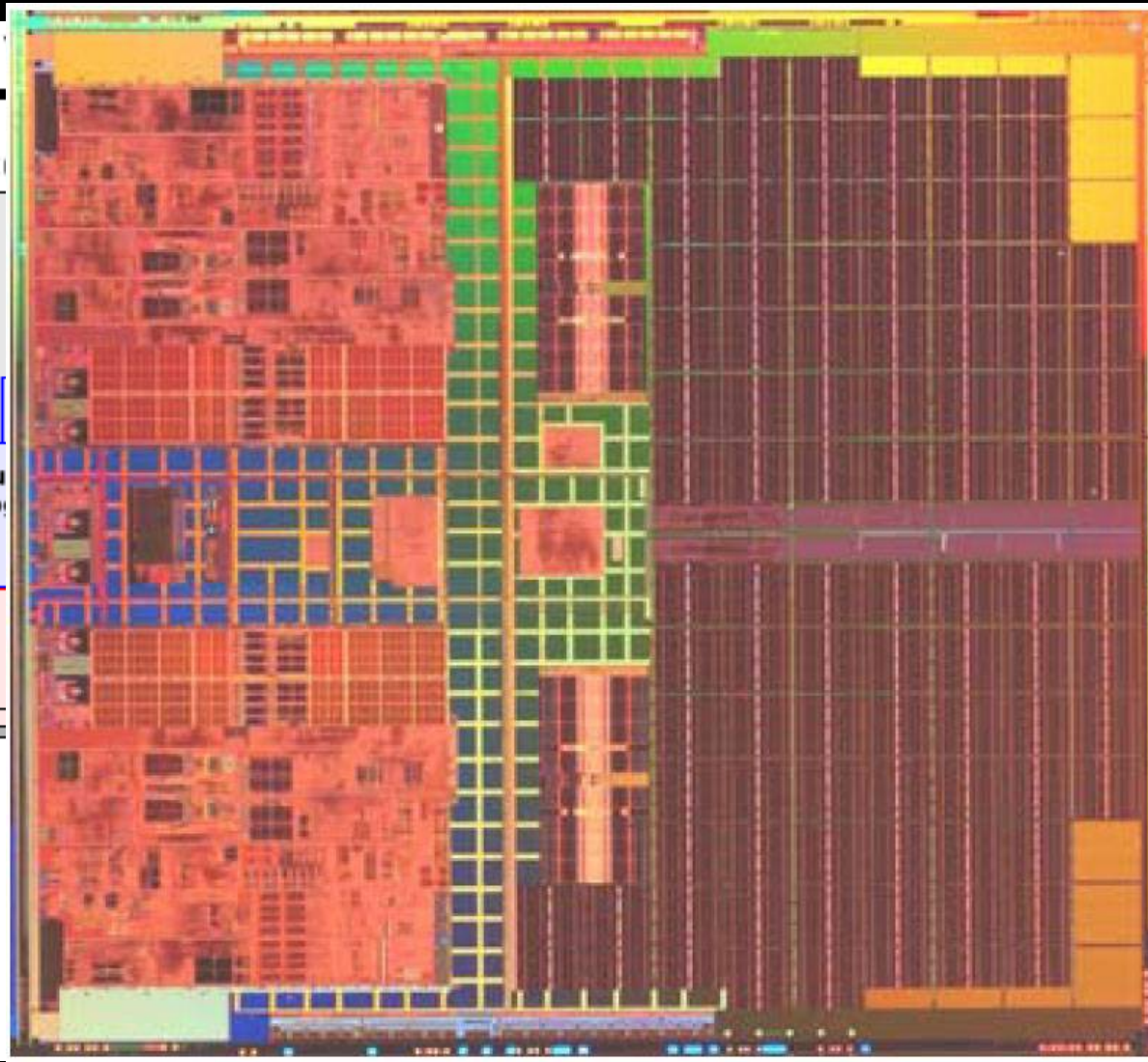
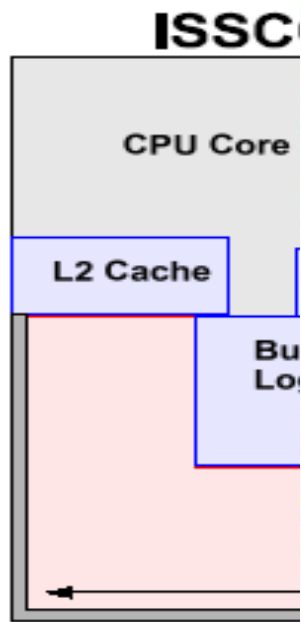
- 430 Million Transistors...
Largest ASIC Chip Ever
Produced
- 344Mb of L3 Memory per
Die
- 4x to 6x larger than SRAMs
used by competitors
- Enables a 50% improvement
in System Performance



“We consider IBM's reworking of the Level 3 cache design to be one of the top design wins in Power5.”

Case Study: *hypothetical* processor chip w/ eDRAM

SRAM



barray



Scaling: Introduction of DRAM is equivalent to scaling (memory) *additional* two generations

DRAM introduction allows for a 0.28X scaling of memory block size at the time of introduction in addition to the assumed 0.7 scaling i.e. total scaling of 0.17

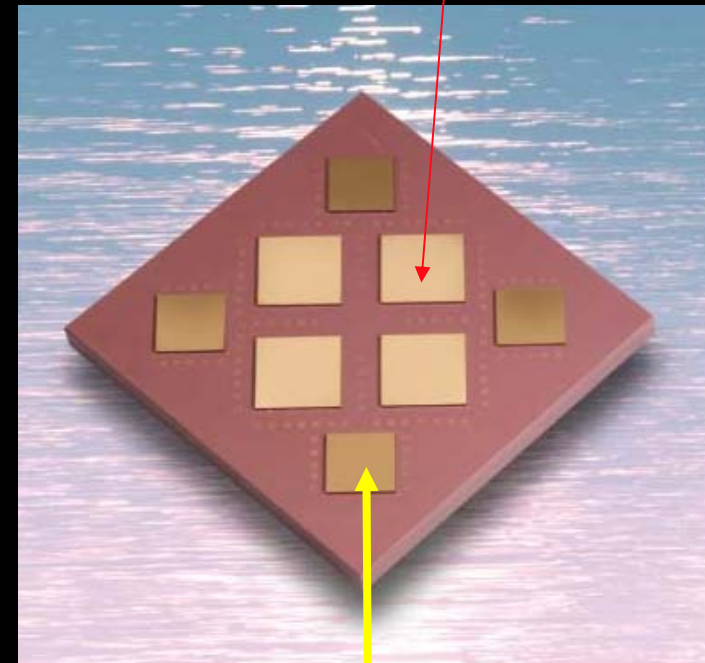
Scaling: Introduction of DRAM is equivalent to scaling (memory) *additional* two generations

DRAM introduction allows for a 0.28X scaling of memory block size at the time of introduction in addition to the assumed 0.6 scaling i.e. total scaling of 0.17

Going forward

- The next step is integration with the processor
- What are the challenges ?
 - Scaling and performance
- What is the approach

Power 4,5, 6 MCM



P5 CPU (SOI)

36MB eDRAM L3 cache (Bulk)

Performance eDRAM

- **70% of performance is gated by logic performance**
 - Address Decoding functions
 - Row system Logic
 - Sensing circuits
- **30% of performance dictated by cell (writeback)**
- **Subarray architecture plays an overall role in optimizing performance-area-power tradeoff**
- **Logic based technologies come out ahead on all fronts**

Performance eDRAM

- **70% of performance is gated by logic performance**
 - Address Decoding functions
 - Row system Logic
 - Sensing circuits
- **30% of performance dictated by cell (writeback)**
- **Subarray architecture plays an overall role in optimizing performance-area-power tradeoff**
- **Logic based technologies come out ahead on all fronts**

Technology Innovation – Development of SOI DRAM cell

■ Technology:

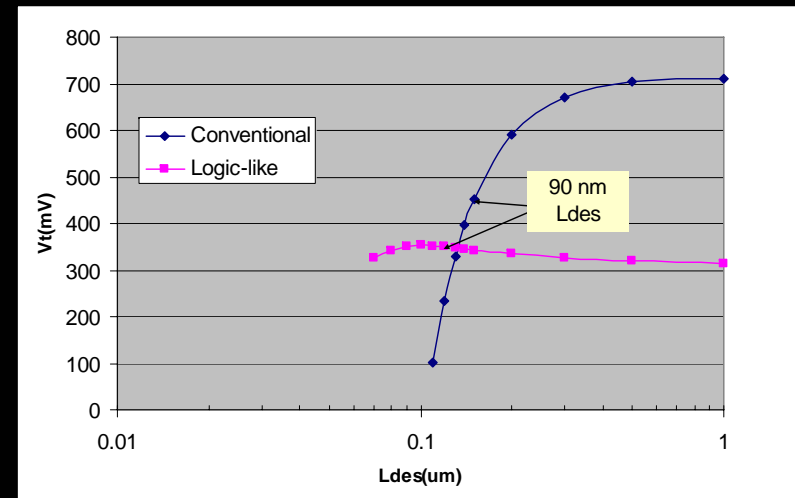
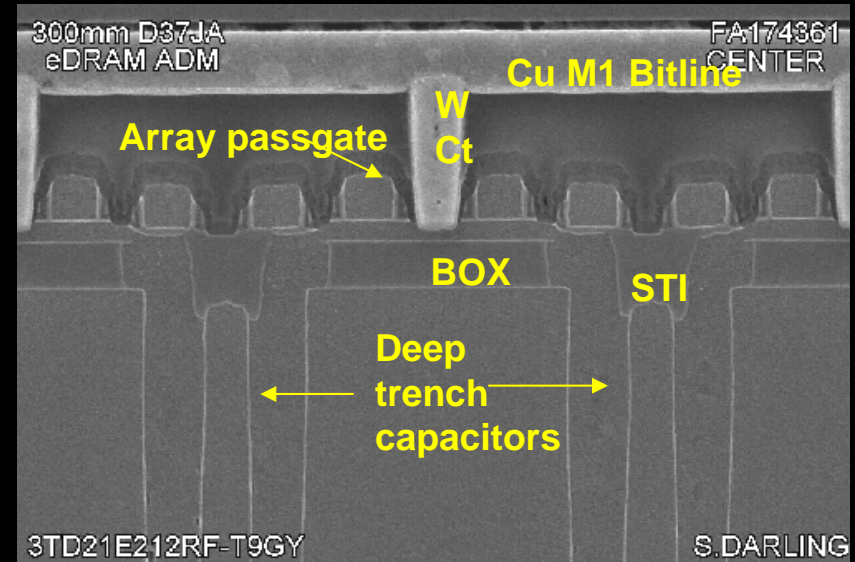
Use the Buried oxide to simplify the process & reduce parasitics – half the cost of bulk eDRAM

Scale the pass transistor for higher performance

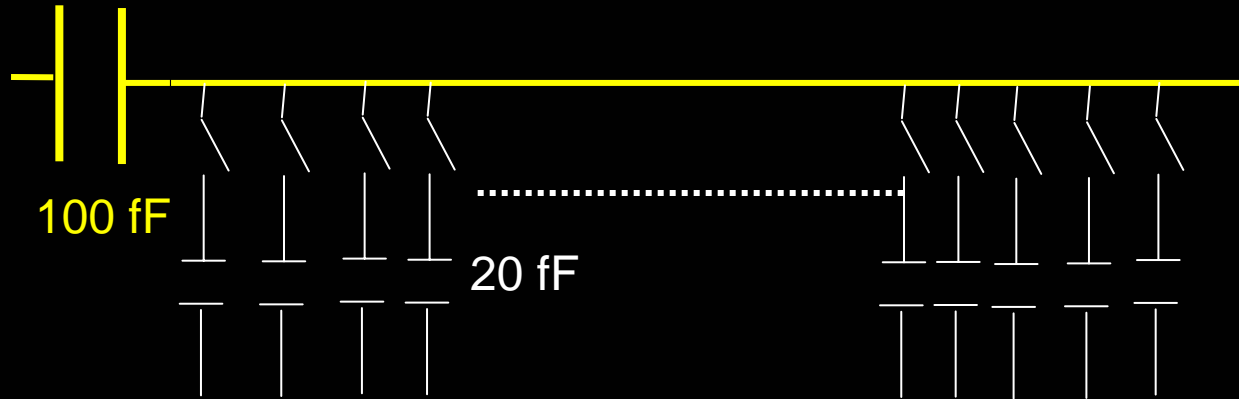
■ Design

Address retention through concurrent refresh

Ultra short bitlines with direct sense architecture



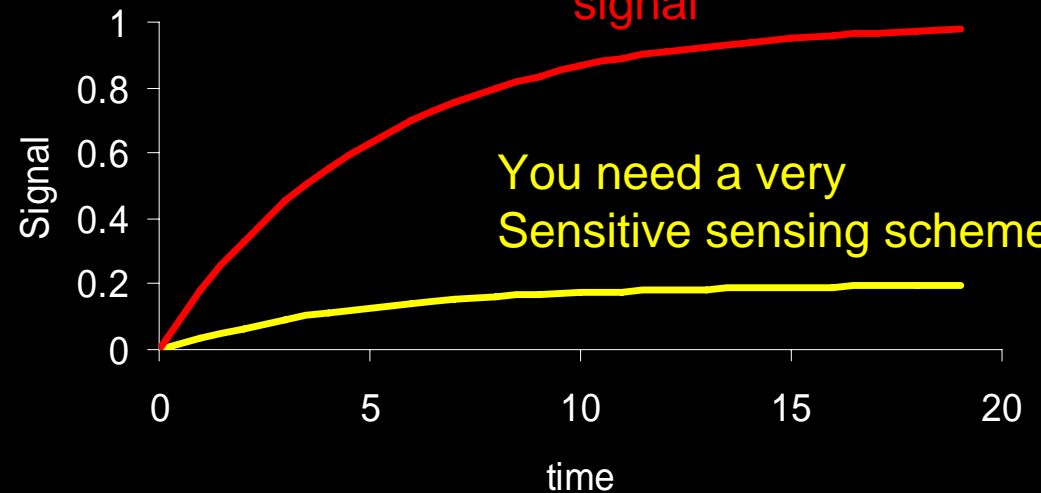
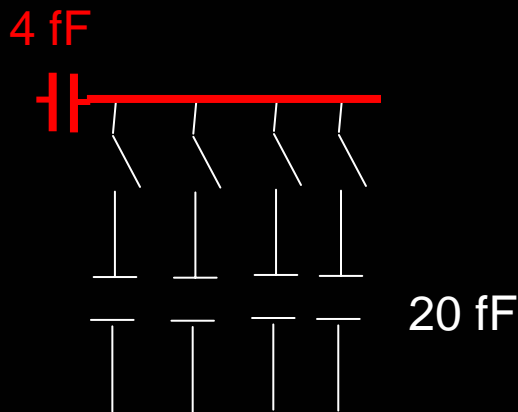
Ultrashort bitlines & μ Sense Advantage



Conventional – small signal, long time

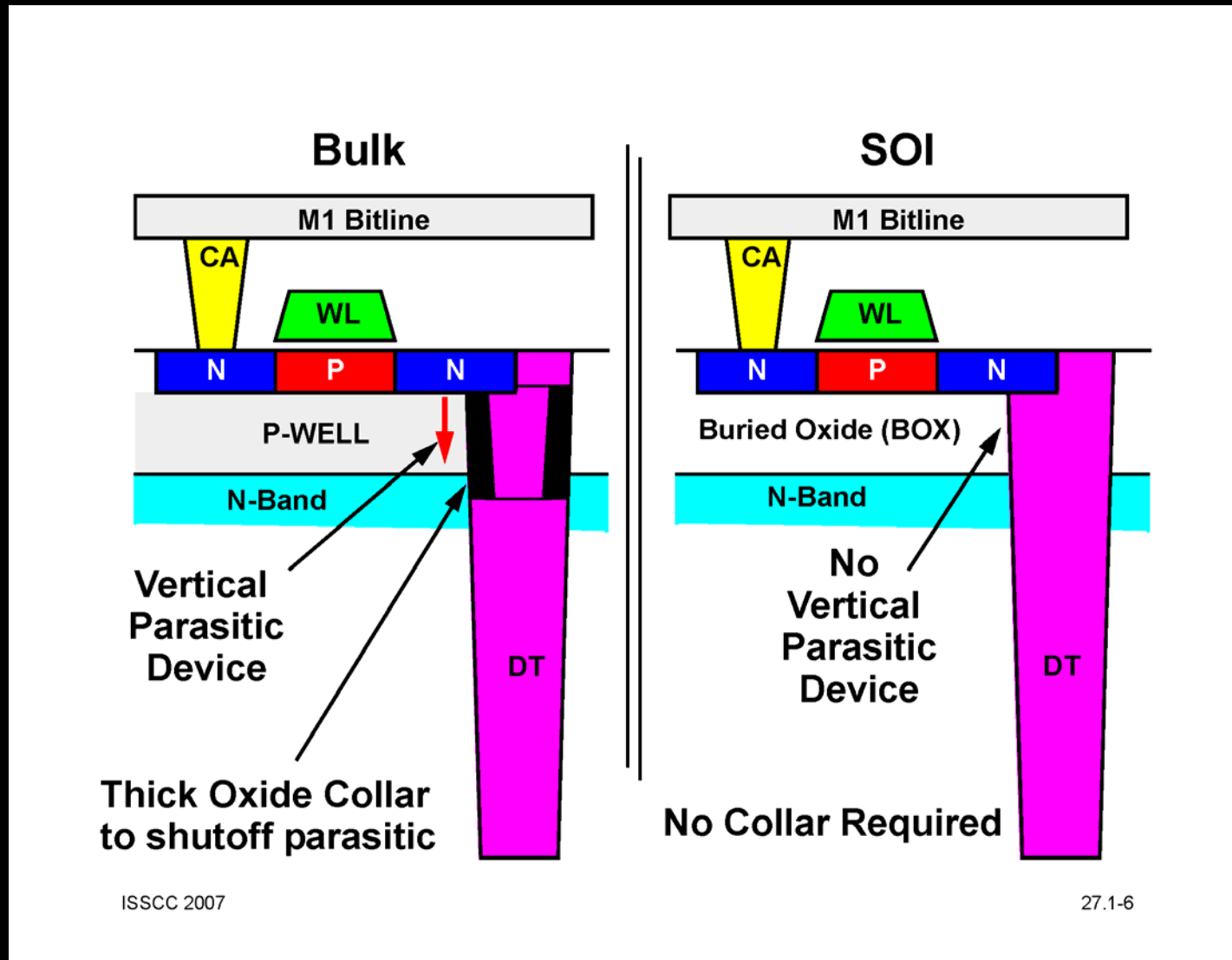
You can drive a gate directly with this signal

μ Sense – large signal short time

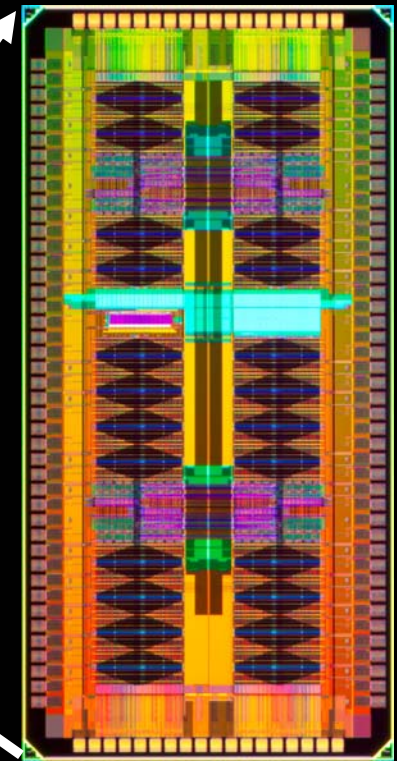
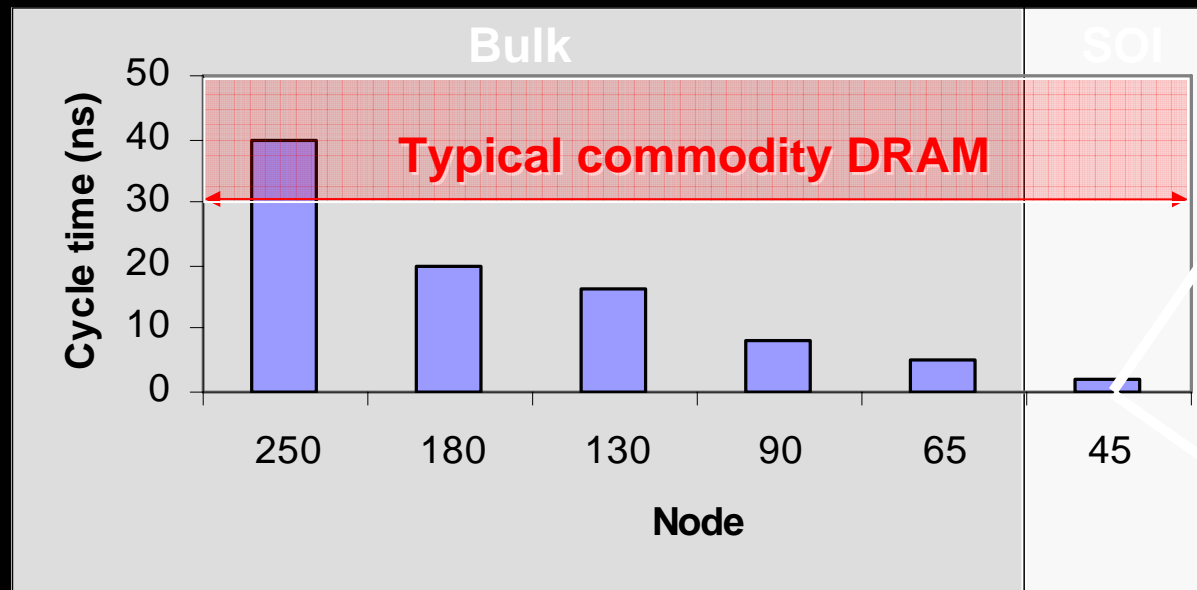


You need a very Sensitive sensing scheme

Collar Process Elimination



SOI eDRAM



Advantages:

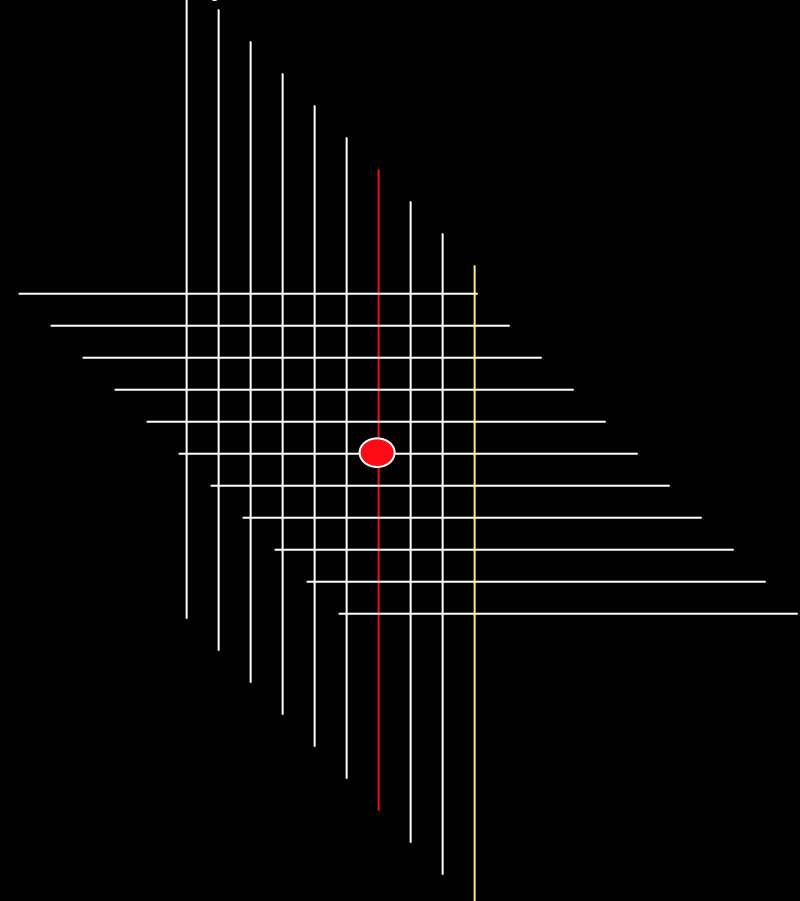
- SRAM like performance @ DRAM density, power and Soft Error Rate
- Full Logic compatibility
- ~3.5X reduction in memory area (3 additional mask levels)

L3 caches in servers

Yielding Large Memory Chips Requires redundancy



Replace **this** with **that**



L3 cache used in P5
344 Mb, 430M transistors
Largest ASIC made at IBM

Note: a chip this size can have as many as 32K fuses

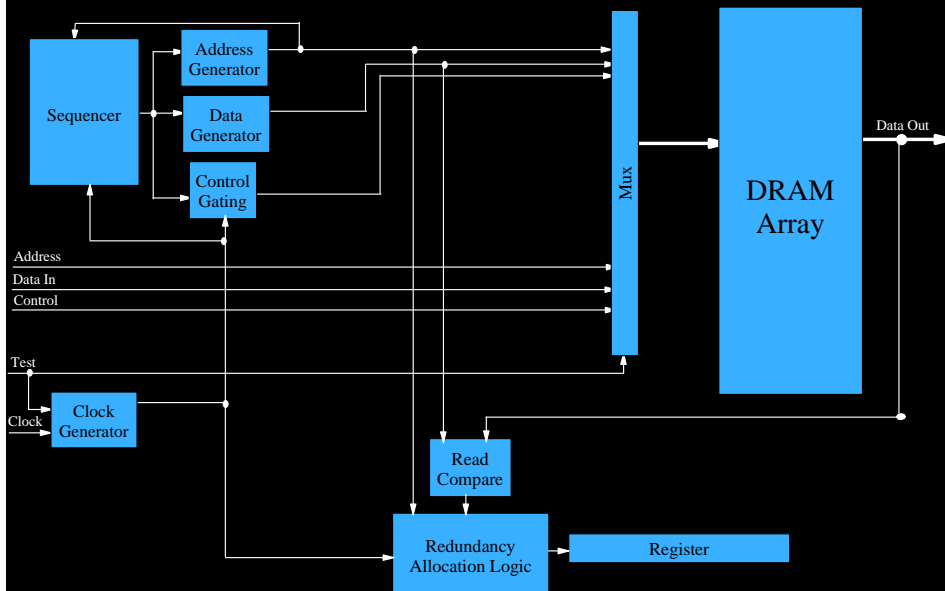
Redundancy invoked using fuses

Built-in Self Test and repair

Key Enablers: BIST & Redundancy

- Test at High speed
- Very large bandwidth but very few pin-outs
- Solution : Since you have access to a high speed logic technology why not build the tester on-chip
- Next step : repair faulty chips!

General Architecture of BIST



- On chip test engine
- Hard & Soft Patterns
- Allocates Redundancy
- Tests redundant elements
- Generates Fuse String

Laser Fuses

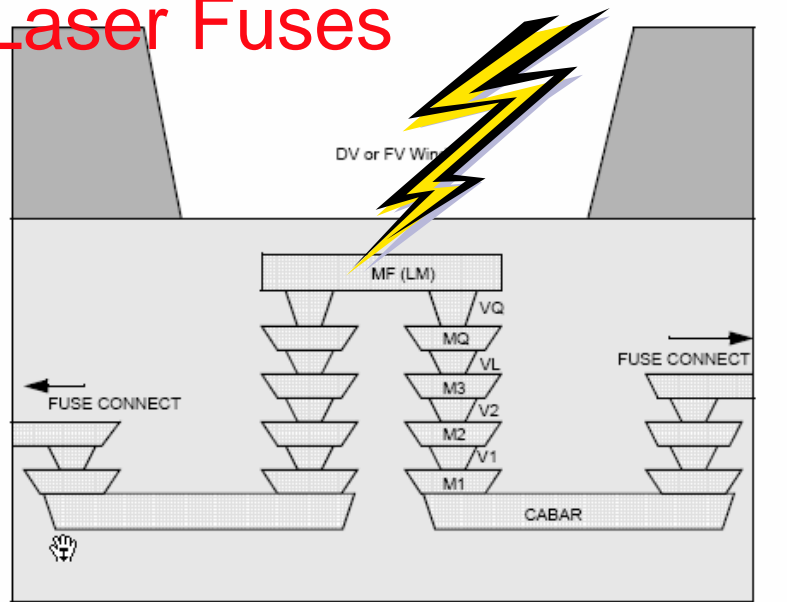


Figure 27: Fuse Cross-Section example - LM

- Do not scale & occupy too much space
- Block wiring and C4s
- Need to be exposed
- Can be blown only at wafer level
- Need precise mechanical alignment
- Require a complex laser fuser
- Require multiple wafer handling and data manipulation

Net: Laser Fuses are a pain!!

Test
 Determine repairs
 Move wafer to laser fuser
 Move repair data to laser fuse
 Blow laser fuses
 Take wafer back to tester
 Test again to verify
 ↓ Hope nothing break again ever!

So, For ages DRAM yogis have pondered important questions:

What is the origin of the Universe ?

Why are we doing DRAMs ?

How do you eliminate laser fuses ?



An electronic fuse would answer these questions

You could use circuits on the chip to program them!

Maybe you could even make the chip autonomic !!

What are *Autonomic Chips* ?

- Chips that can identify, test and repair themselves
In the Fab at wafer level, In the package, In the field
multiple times
- Repairs are self-consistent and contained on chip
without need to access the rest of the system
They need to be hard-coded on the die
- Chips that can maximize their yield within a
predefined parametric space
- Chips that ensure only authorized access
- Chips that identify hazardous operating conditions
Temperature, excessive device drifts
- Chips that can track their history

Achieving autonomic chips

- Need on-chip monitor to be aware of things going wrong

Built-in-Self Test (BIST) engine (processor or dedicated)

- A self-repair methodology

- **At the algorithmic level**

- Aware of previous repairs

- Aware of available repair opportunities

- Decides on irreparability and functional disablement

- **A method to execute the same permanently on chip**

- A non-volatile memory element that is easy to implement (e.g.. Fuse)

But seriously



Do you want this on your chip ?

Fine Print: Explosion here has been somewhat exaggerated

Poly Si Fuse Programmed by rupture

The eFUSE Philosophy



Technology Independent

Fuse must work even if die is marginal

Programming time should be short

Should be programmable on chip

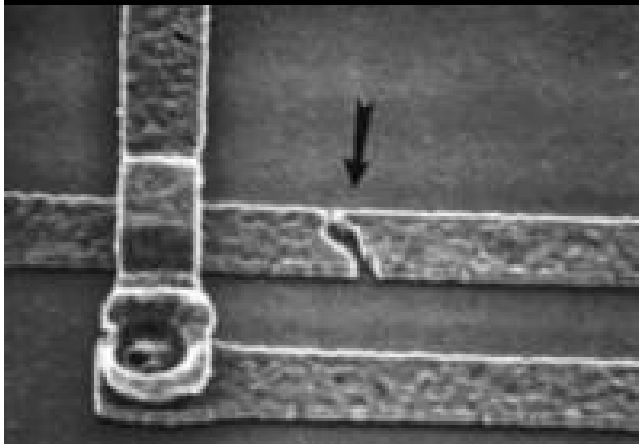
Use design levels and Booleans to create Fuse link

Understand the physics of the blow mechanism

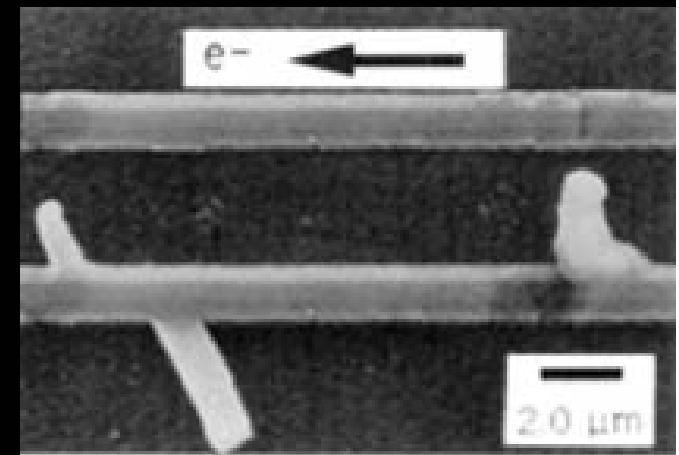
Design innovative circuits to program and sense the fuses

The Kinder Gentler Fuse

We need to induce an electrical open without needing material to disappear.



$$\nabla \cdot \neq 0$$

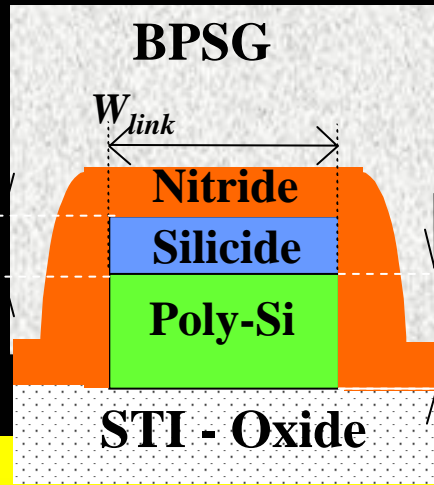
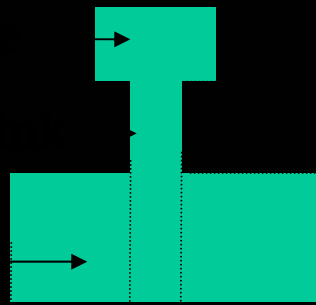


Can we employ electromigration of metal lines ?

Modern interconnects are electromigration resistant

Need to control the electromigration and complete in a reasonable time e.g. 200 μs – Cu line not an option

eFUSE - physical layout



cathode

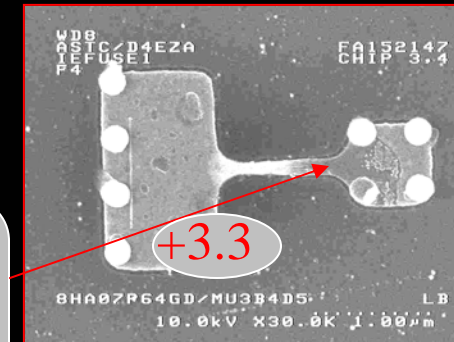
Mechanism:

- Current driven through silicide
- Temperature rises & gradient set up
- Silicide electromigrates but current is sustained as the Poly Si is hot intrinsic and conductive
- Electromigration of silicide is forced to completion
- Current turned off, everything cools and link is high resistance

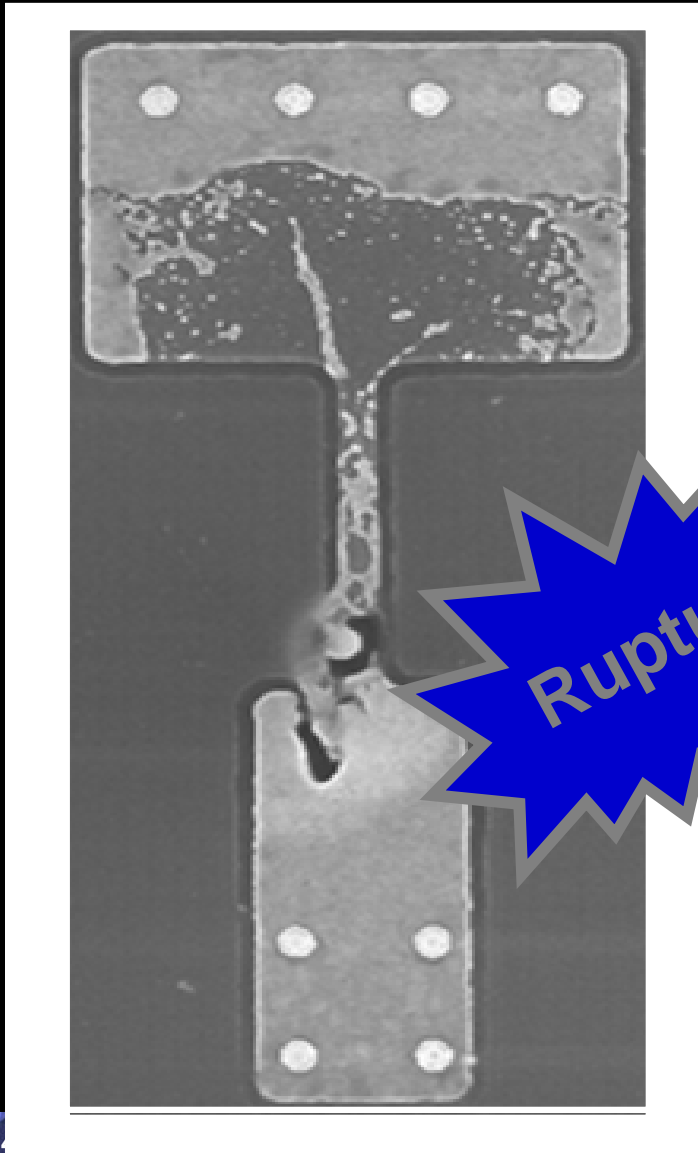
Key Features

- Geometry
- Thermal environment
- Current Characteristics

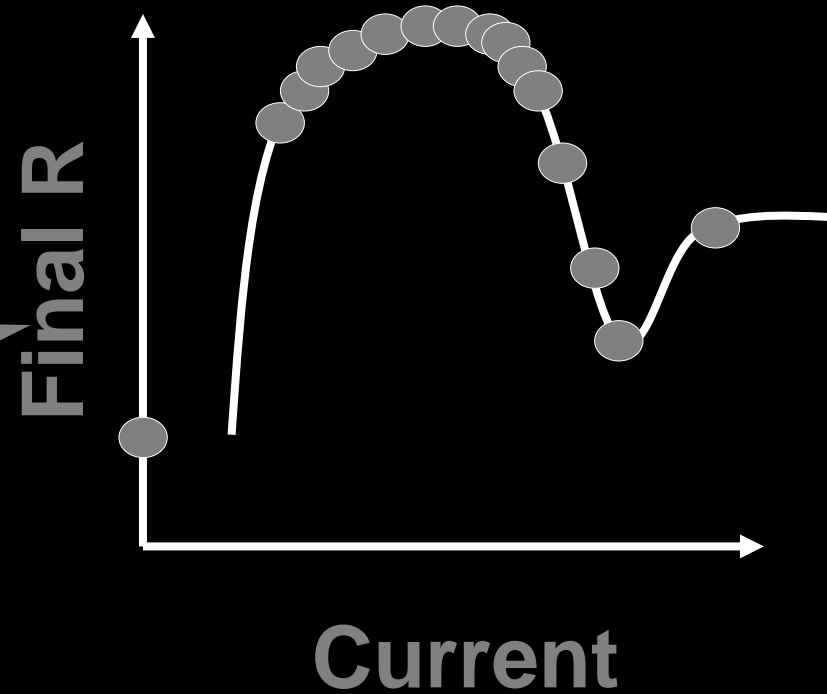
Silicide removal always occurs at cathode



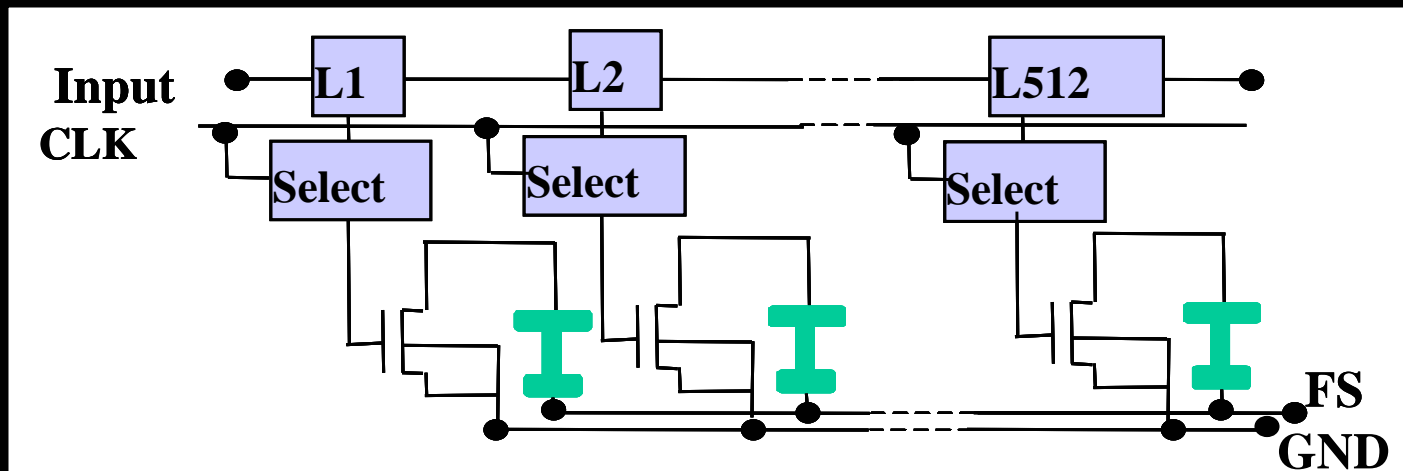
eFUSE Programming



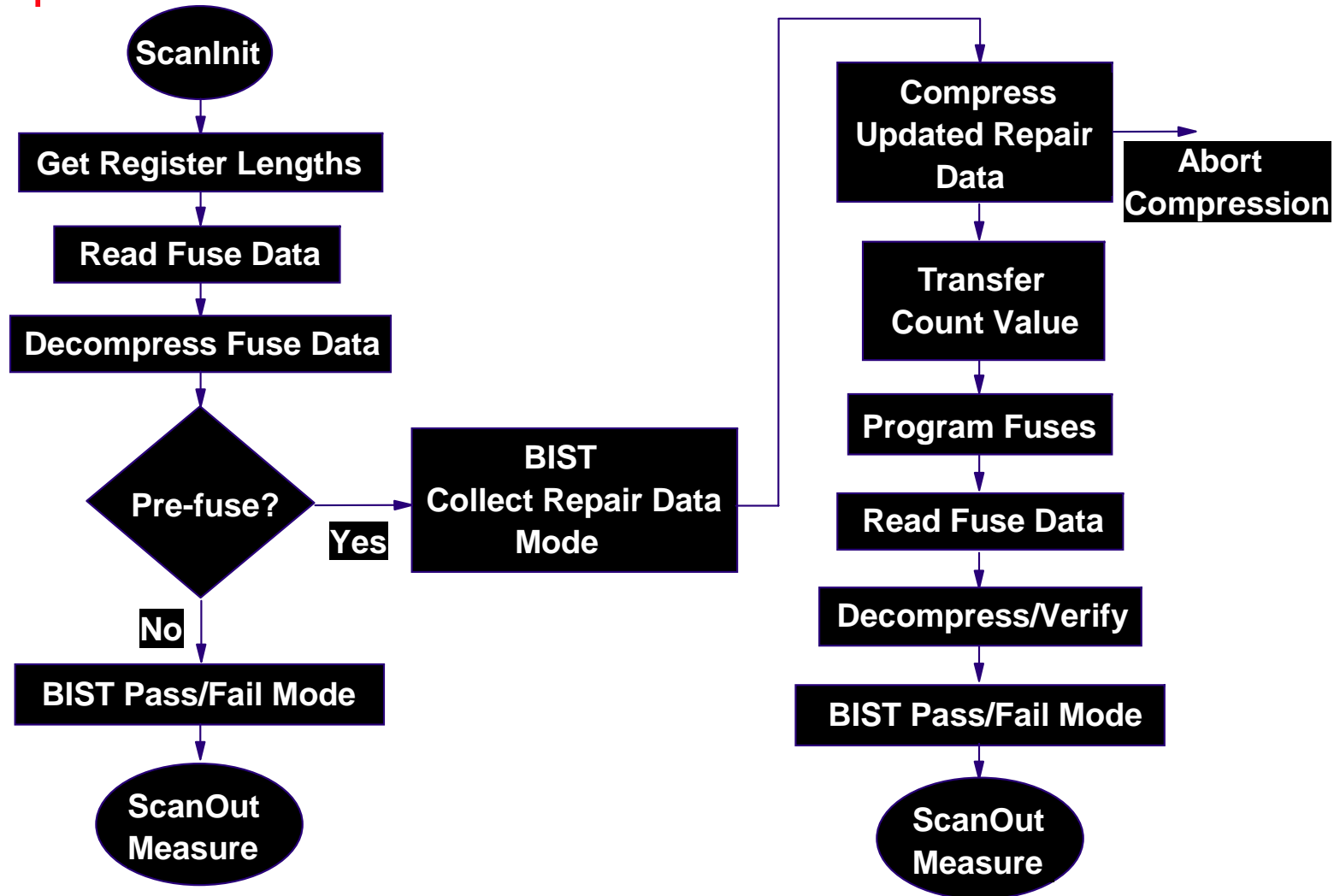
Intact Fuse



Circuit details

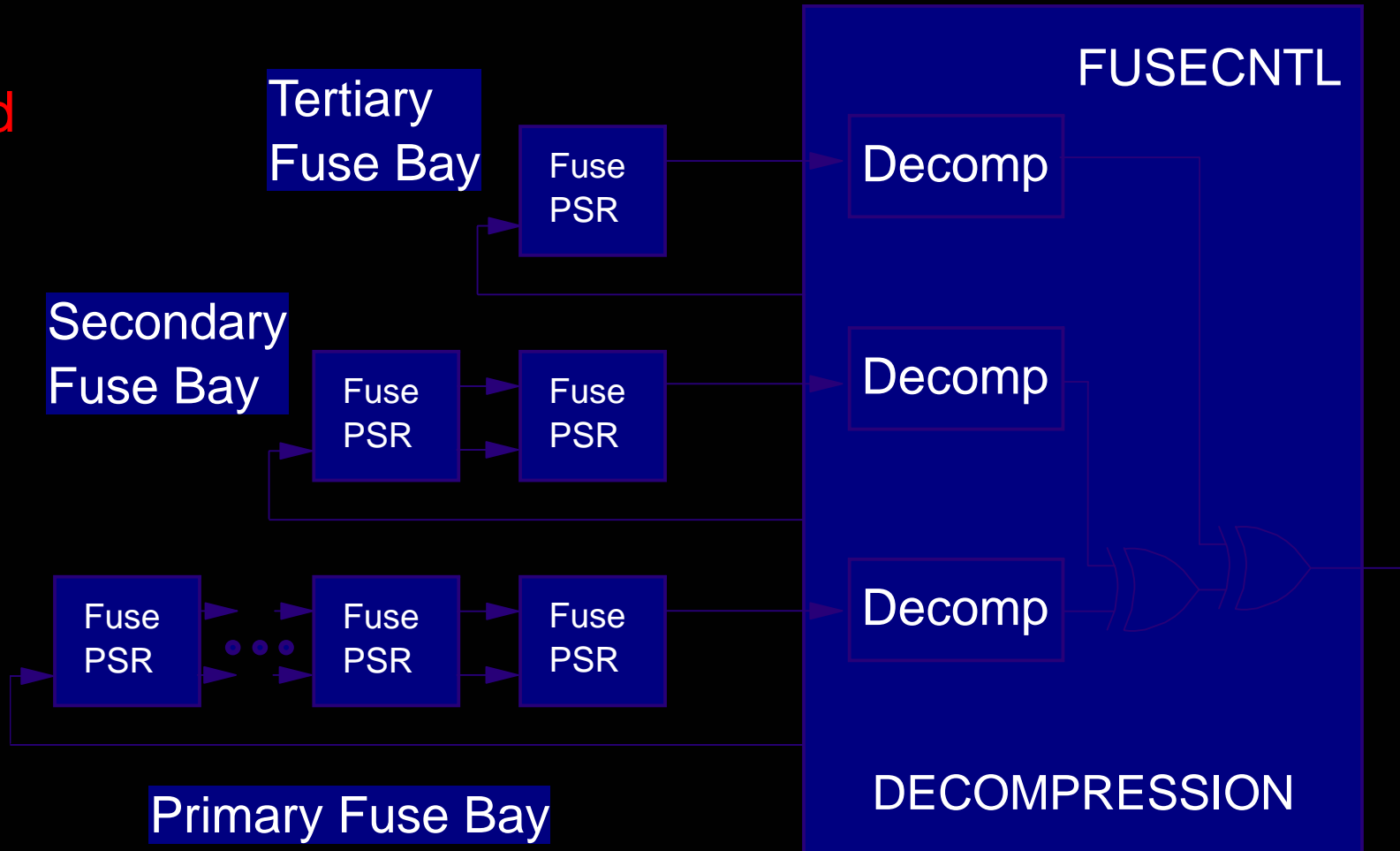


Self-Test and Repair Flow – one touch test, repair & verify as implemented in Cu08

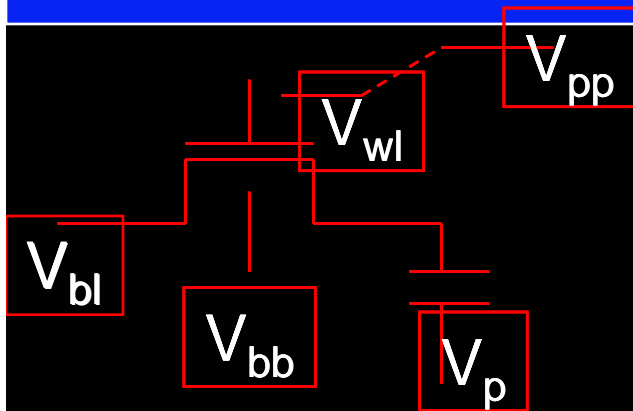
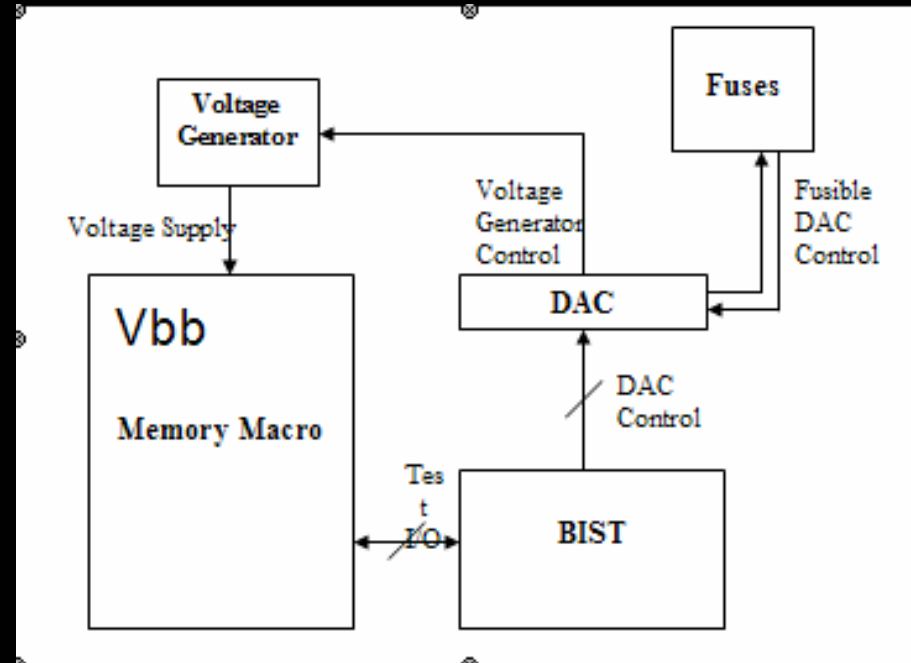
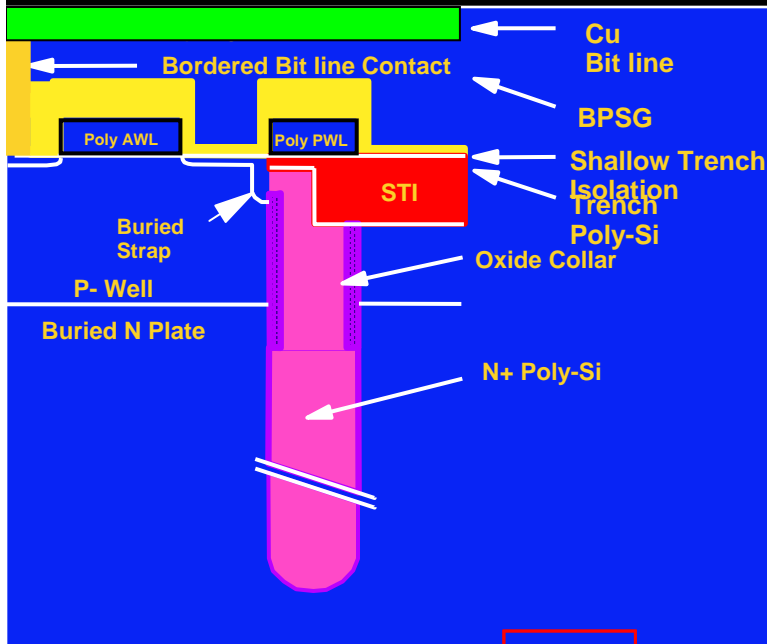


Hierarchical Repair - Combining Three Repair Steps as implemented in Cu08

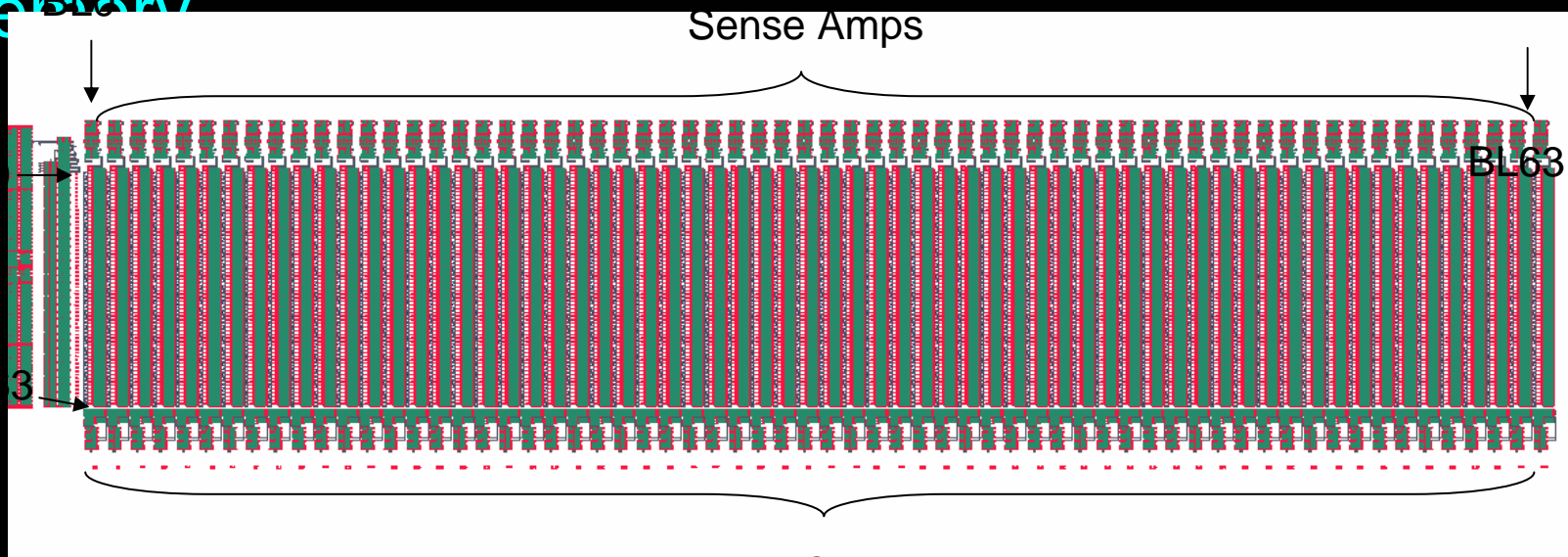
In the field
@ module
@ Wafer



Die specific Yield optimization: Voltage setting using eFUSE in embedded DRAM



2 Dimensional OTP ROM – for on-chip nonvolatile memory

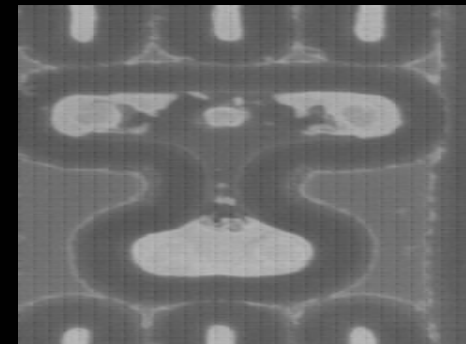


Based on our electromigration Fuse technology – -
- No added process Cost

Memory architecture allows for ~15X improvement
in bit density ($115 \mu\text{m}^2/\text{bit}$ to $7 \mu\text{m}^2/\text{bit}$)

Compilable in 2Kb blocks

Centralized repair, code and autonomic functions



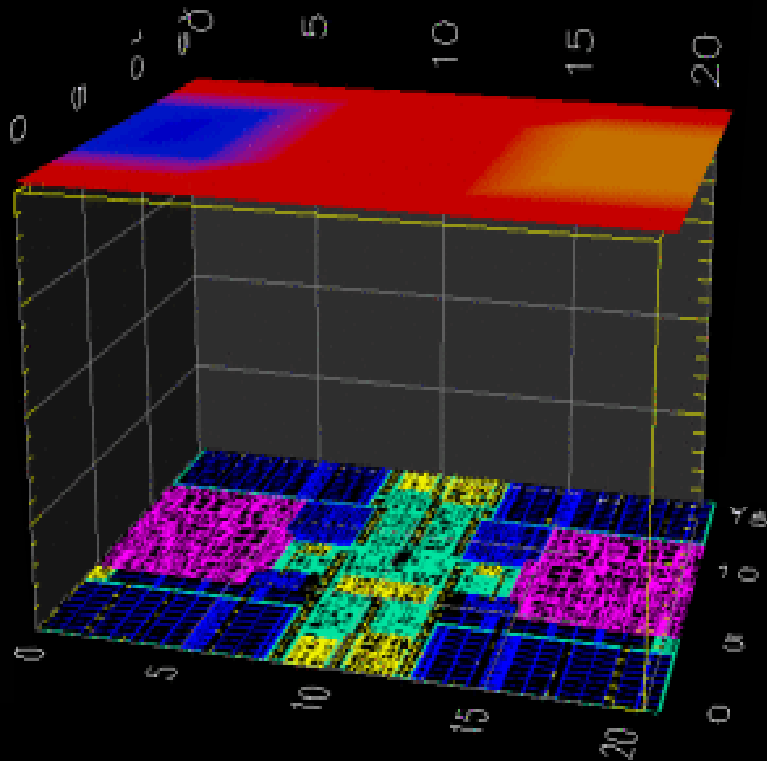
What's next ?

- Optimizing chips for power and performance
 - More autonomic functions
 - Supply Chain management
 - Chip identifiers for authentication
 - RFID tags
 - Logic compatible 2D OTP ROM
-
- Exploit other phenomena: Hot carrier, NBTI to self balance mismatched circuits

Power supply decoupling

- Why is decoupling important ?
The concept of voltage compression
- Decoupling today is mostly an afterthought but compression is a serious problem
- A judicious of decoupling can allow for upto 10% reduction in V_{dd} due to reduced compression

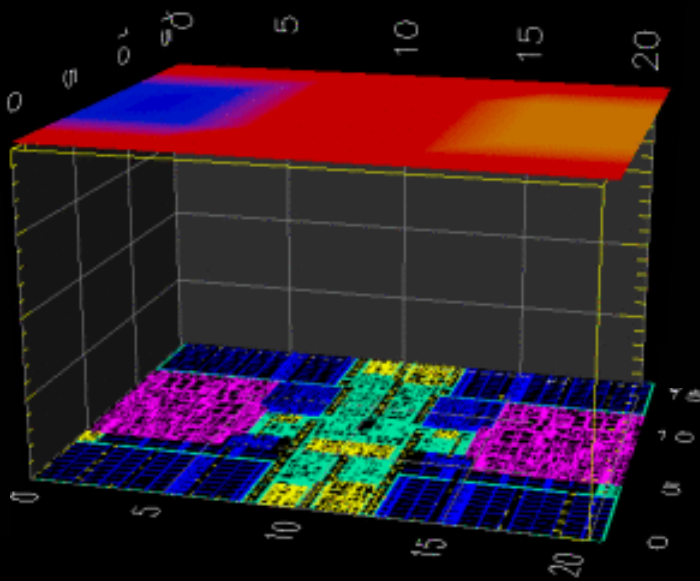
Noise Traveling Across Chip (Connected cores)



- Noise travels from active core to idle core
- ~4ns delay between cores
- Noise is attenuated at quiet core

James et al, ISSCC 2007

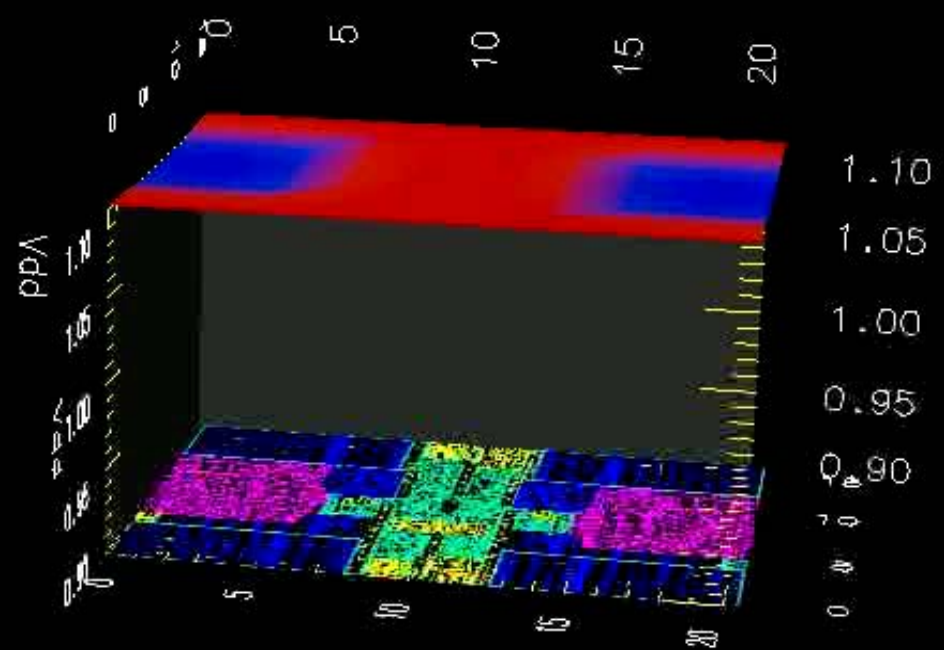
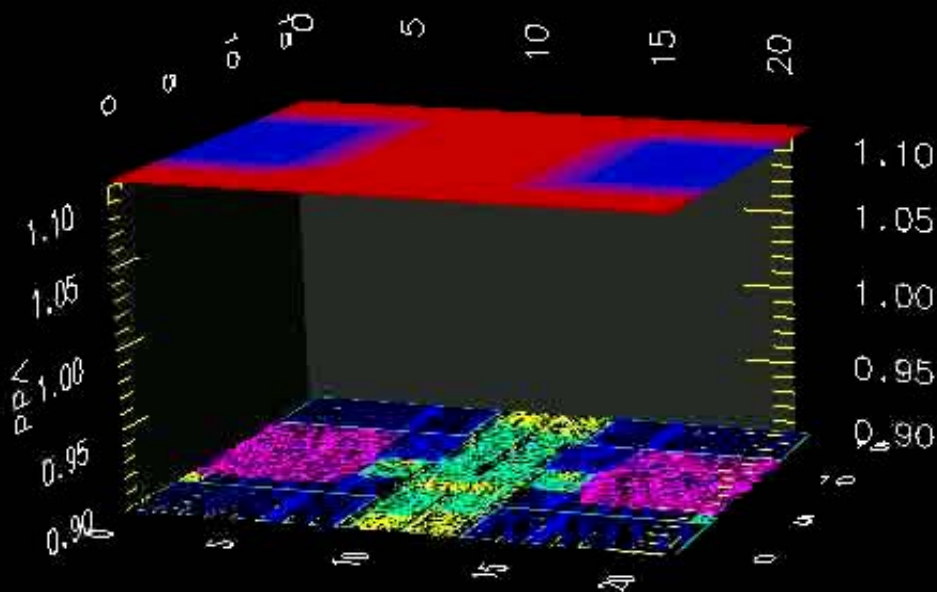
Noise Traveling Across Chip (Connected cores)



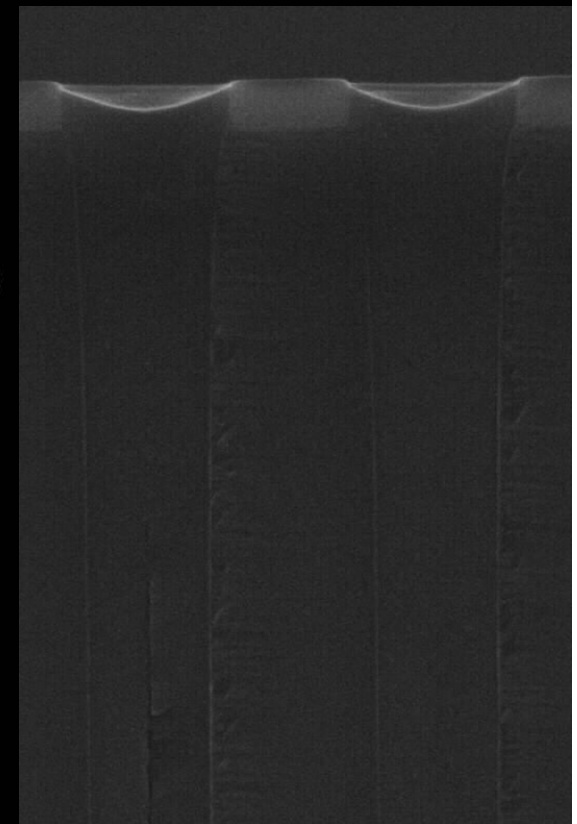
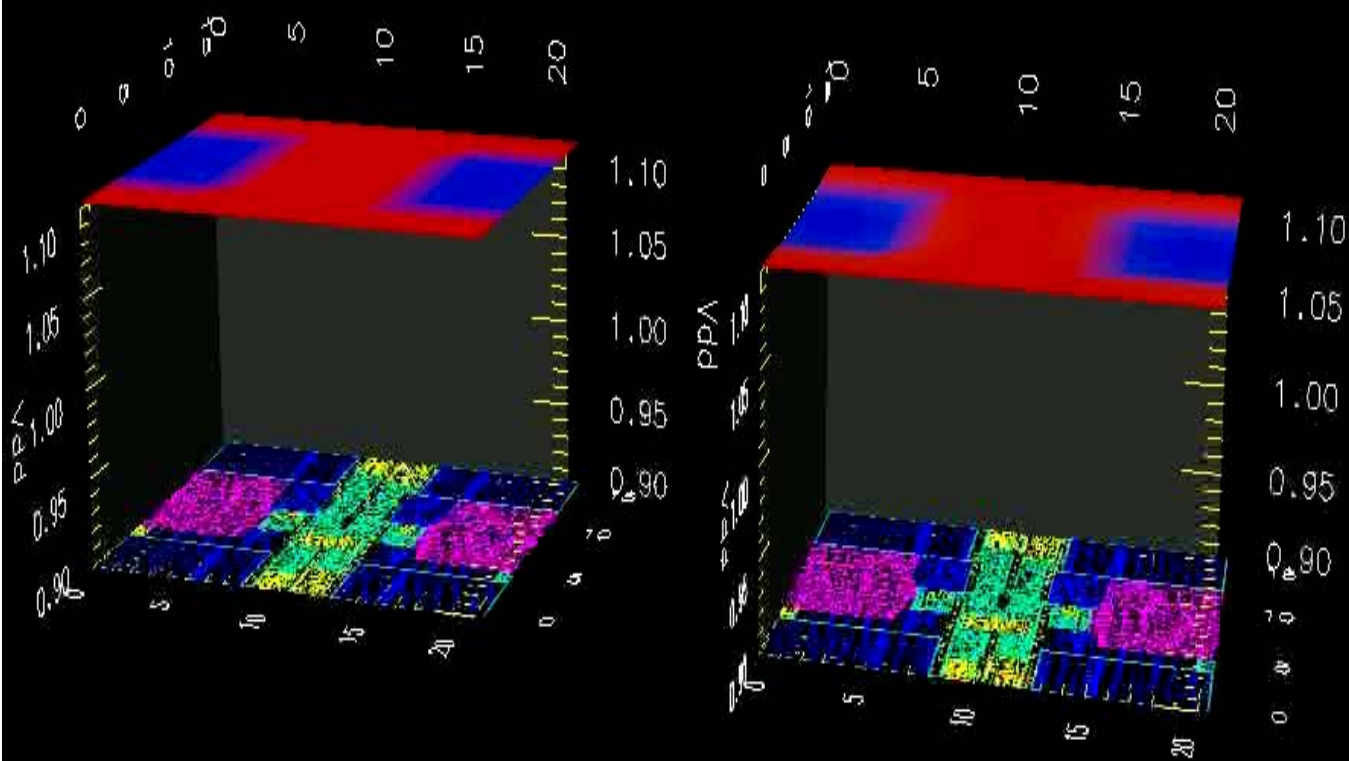
- Decoupling today is mostly an afterthought but Transient voltage droops due switching activity is a serious problem
- Planar Caps are large and leak tremendously
- So we cannot isolate this noise and must increase power supply to overcome
- Can be a serious issue in high performance multicore applications

Trench based Decoupling

- Deep trench decoupling reduces this by over 100mV and enables a 10% reduction of power supply voltage at same performance
- Comes free with our embedded DRAM



Trench based Decoupling



Current Solutions / Problems

Header / Footer based power gating popular

Lower array supply to save on leakage (reduce V_{DS})

Problem occurs during Wake Up (access cycle)

Most schemes have wake up cycles for DI/DT

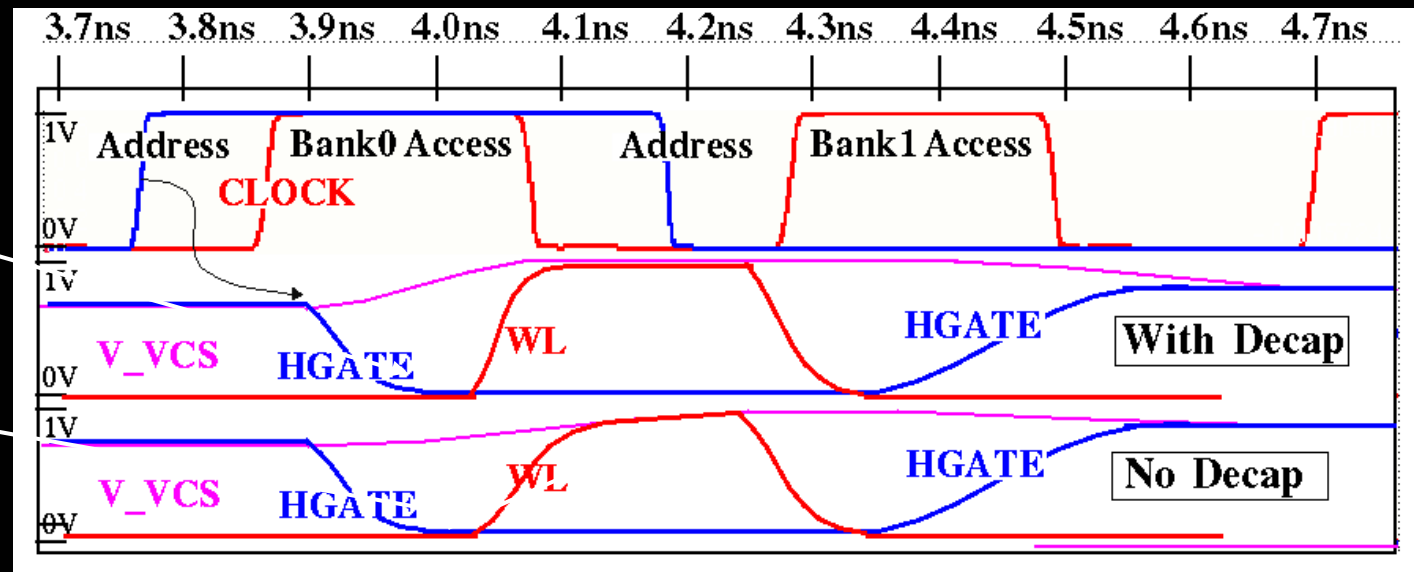
Power supply response time $> 1\text{ns}$

Need multiple wake cycles @ 3-4 GHz frequencies

Area efficient local charge reservoir can help !

Fast
Charge up
w/ decap

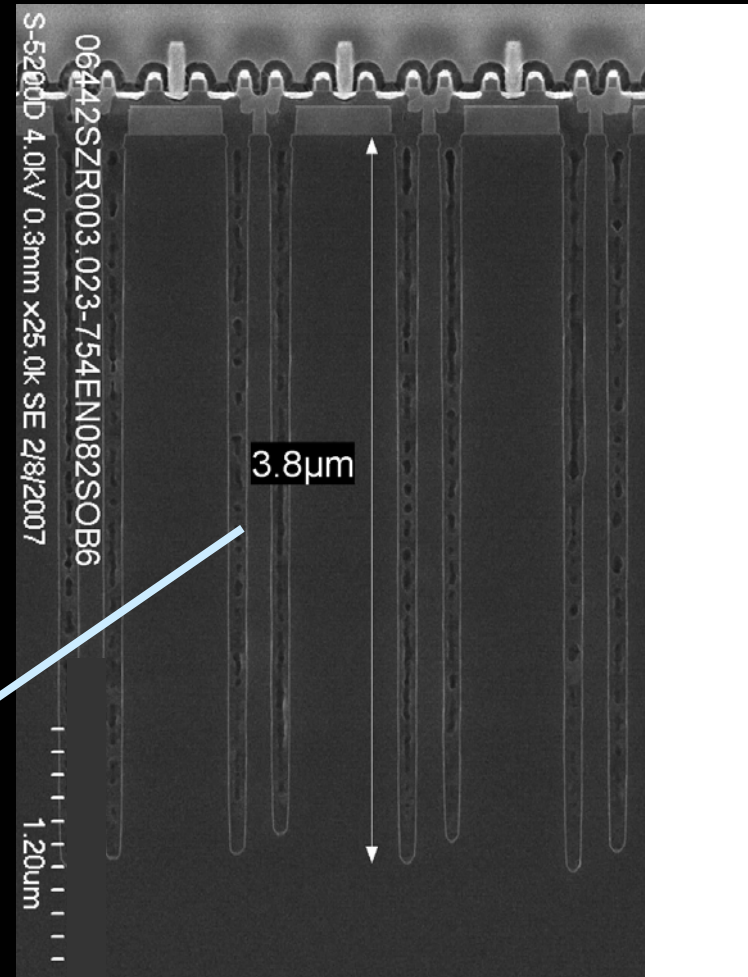
Slow RC
without
decap



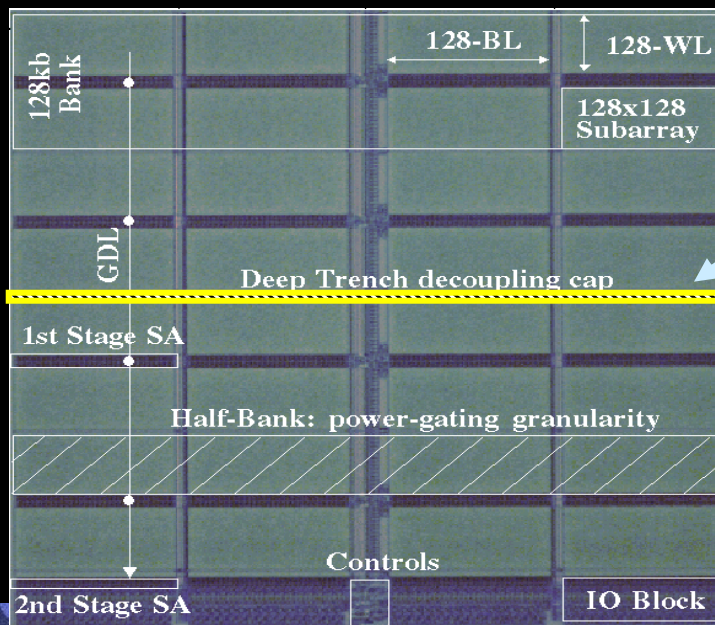
Performance penalty

The DT Decoupling Cap 20X Advantage

- DT decoupling capacitance: $200\text{fF}/\mu\text{m}^2$ vs. $11\text{fF}/\mu\text{m}^2$ for thick-oxide depletion device
- Incorporated in 45nm ASIC Methodology
- Local array-supply charge reservoir to enable dynamic leakage reduction with NO wake-up cycles for $> 2\text{GHz}$ operation



512Kb
ASIC
SRAM



**< 2% area impact for 250pF
of Embedded DT CAP**

Three Dimensional Integration – scaling in the third dimension



Summary

- Scaling provides some density benefit
- Technology performance comes from technology features rather than scaling: eg Strain Engineering
- Hi K if can get us back to scaling for a short time
- Memory integration is a very big thing
- Chips will have more autonomic system level like functions enabled by technology features such as eFUSE
- Innovative decoupling can provide for upto a 10% performance improvement
- Three Dimensional Integration is the next **huge** thing

I would like to acknowledge the contributions of the 45nm , Hi K, embedded DRAM and eFUSE teams

Yes..... The Old Dog does learn new tricks

