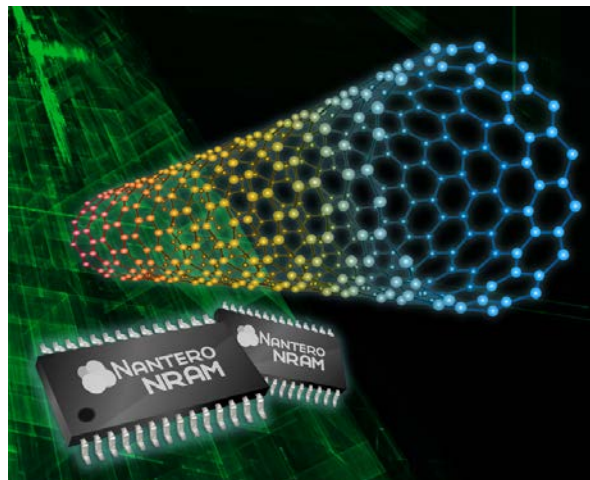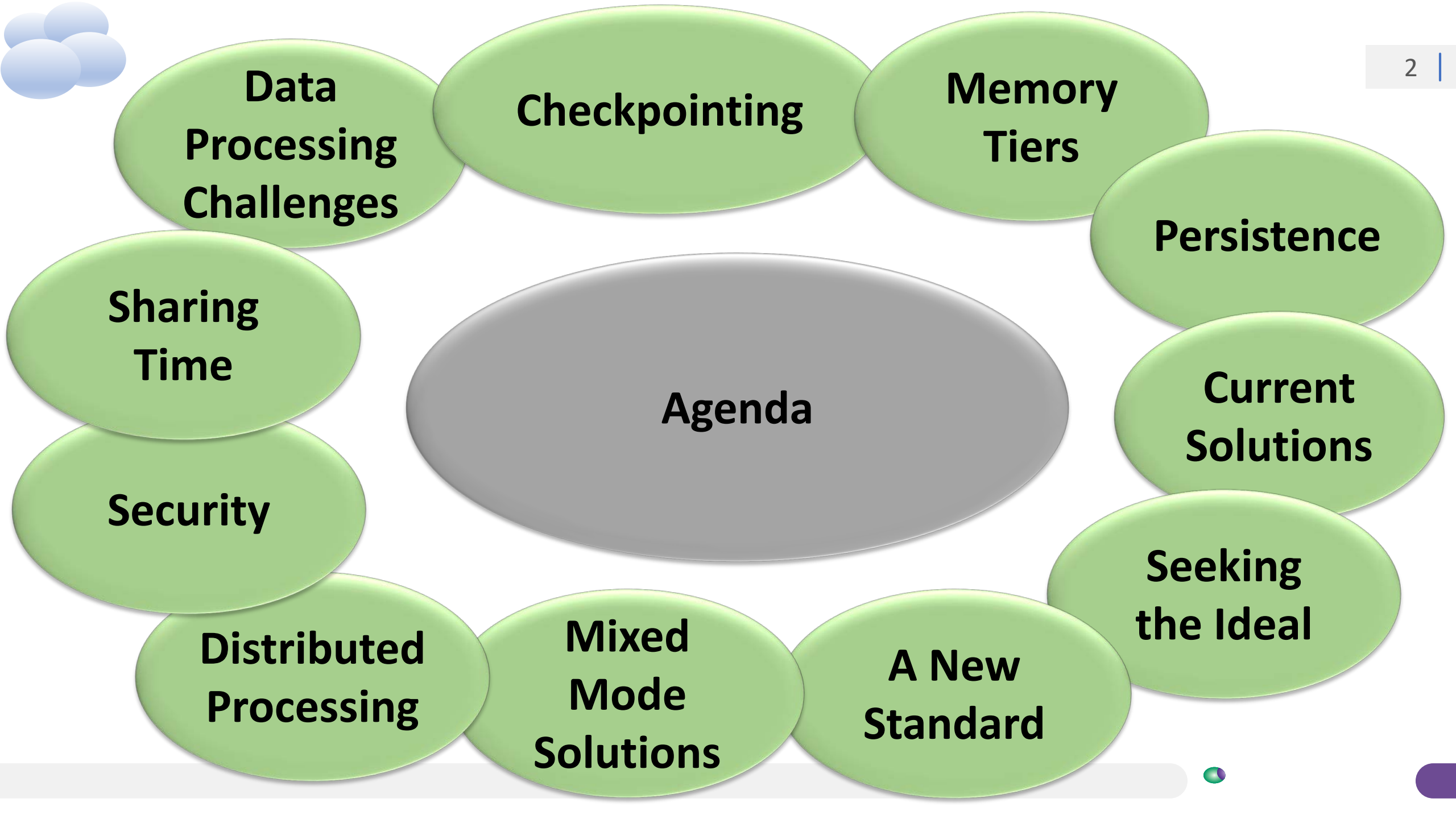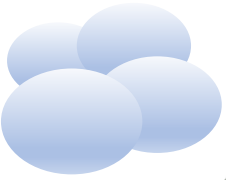# Expanding the World of Heterogenous Memory Hierarchies

## The Evolving Non-Volatile Memory Story



Bill Gervasi
Principal Systems Architect

Data Processing Challenges

Checkpointing

Memory Tiers

Persistence

Sharing Time

Agenda

Current Solutions

Security

Seeking the Ideal

Distributed Processing

Mixed Mode Solutions

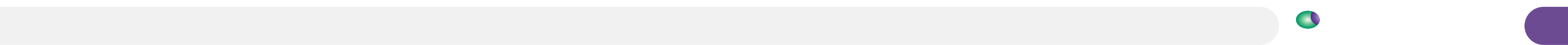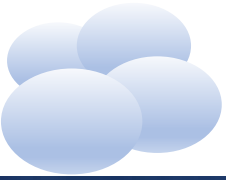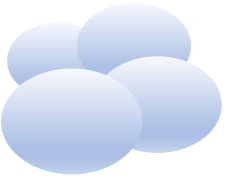A New Standard

# Data processing is great
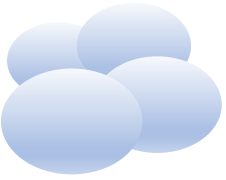
Data processing is great

Until something goes wrong

# The Cost of Power Failure

According to Gartner, the average cost of IT downtime is **$5,600** per minute. Because there are so many differences in how businesses operate, downtime, at the low end, can be as much as $140,000 per hour, **$300,000** per hour on average, and as much as $540,000 per hour at the higher end. Jun 18, 2018

The 20 | The Cost of IT Downtime | The 20
https://www.the20.com/blog/the-cost-of-it-downtime/

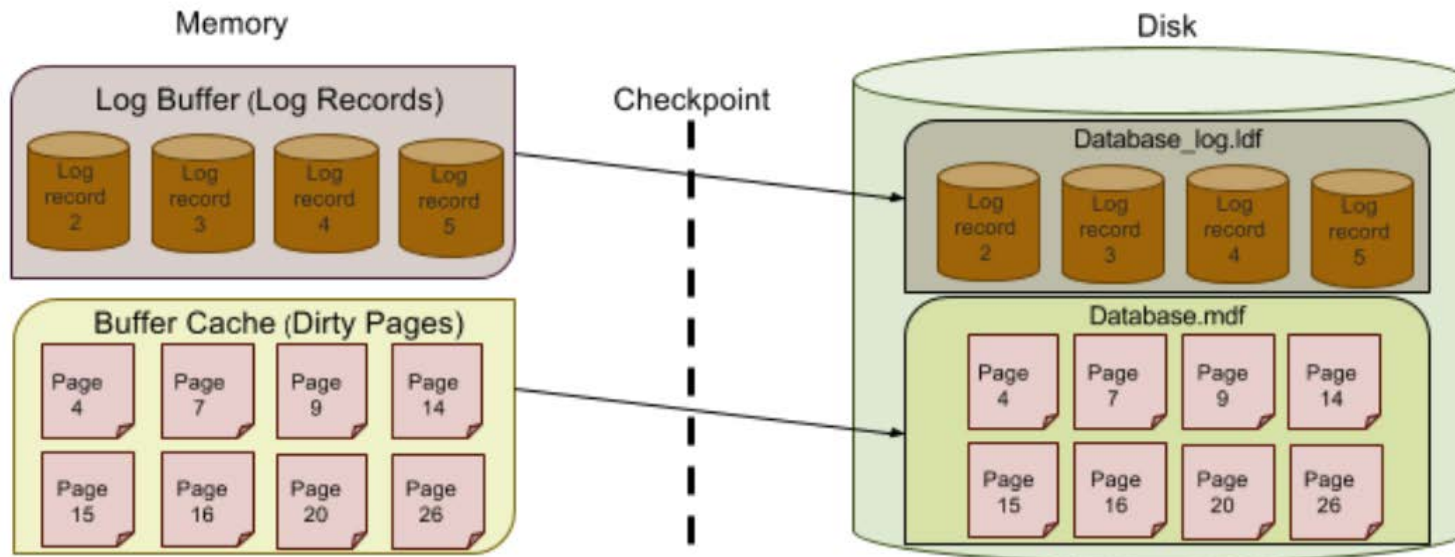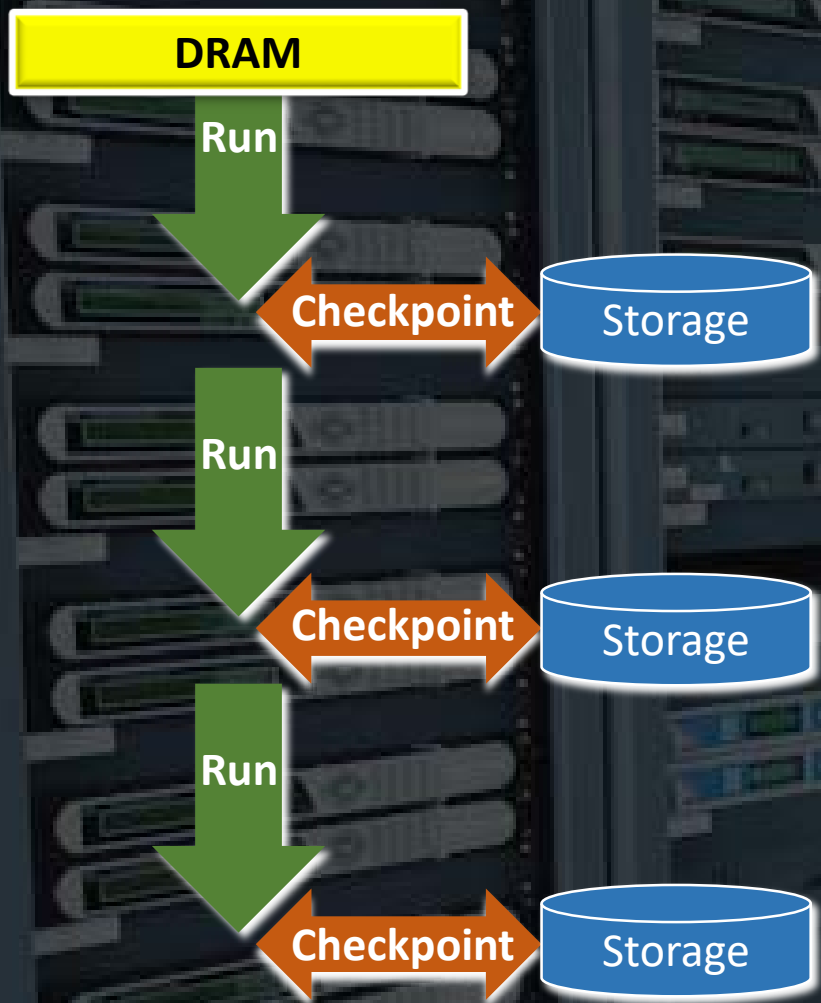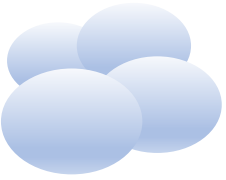## Amazon.com Goes Down, Loses $66,240 Per Minute

# Checkpoint

📅 November 12, 2015  👤 Alexandr Omelchenko  📁 Glossary

⭐⭐⭐⭐⭐ 📊 [Total: 21   Average: 4.2/5]

Checkpoint is a process that writes current in-memory dirty pages (modified pages) and transaction log records to physical disk. In SQL Server checkpoints are used to reduce the time required for recovery in the event of system failure. Checkpoint is regularly issued for each database. The following set of operations starts when checkpoint occurs:

1. Log records from log buffer (including the last log record) are written to the disk.
2. All dirty data file pages (pages that have been modified since the last checkpoint or since they were read from disk) are written into the data file from the buffer cache.
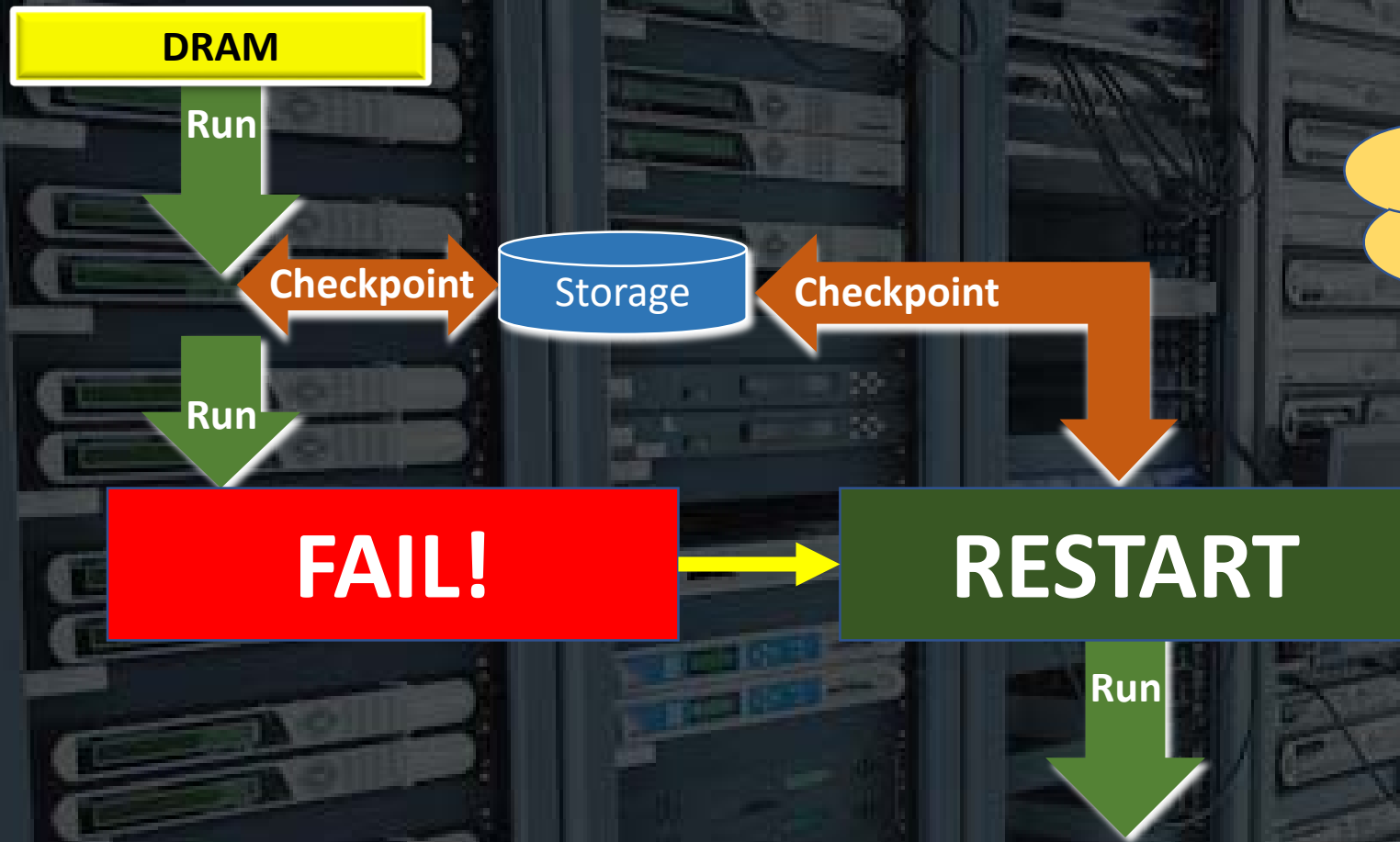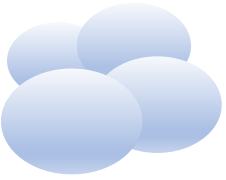3. Checkpoint LSN is recorded in the database boot page.
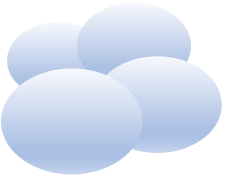
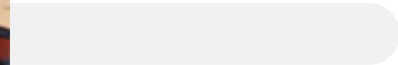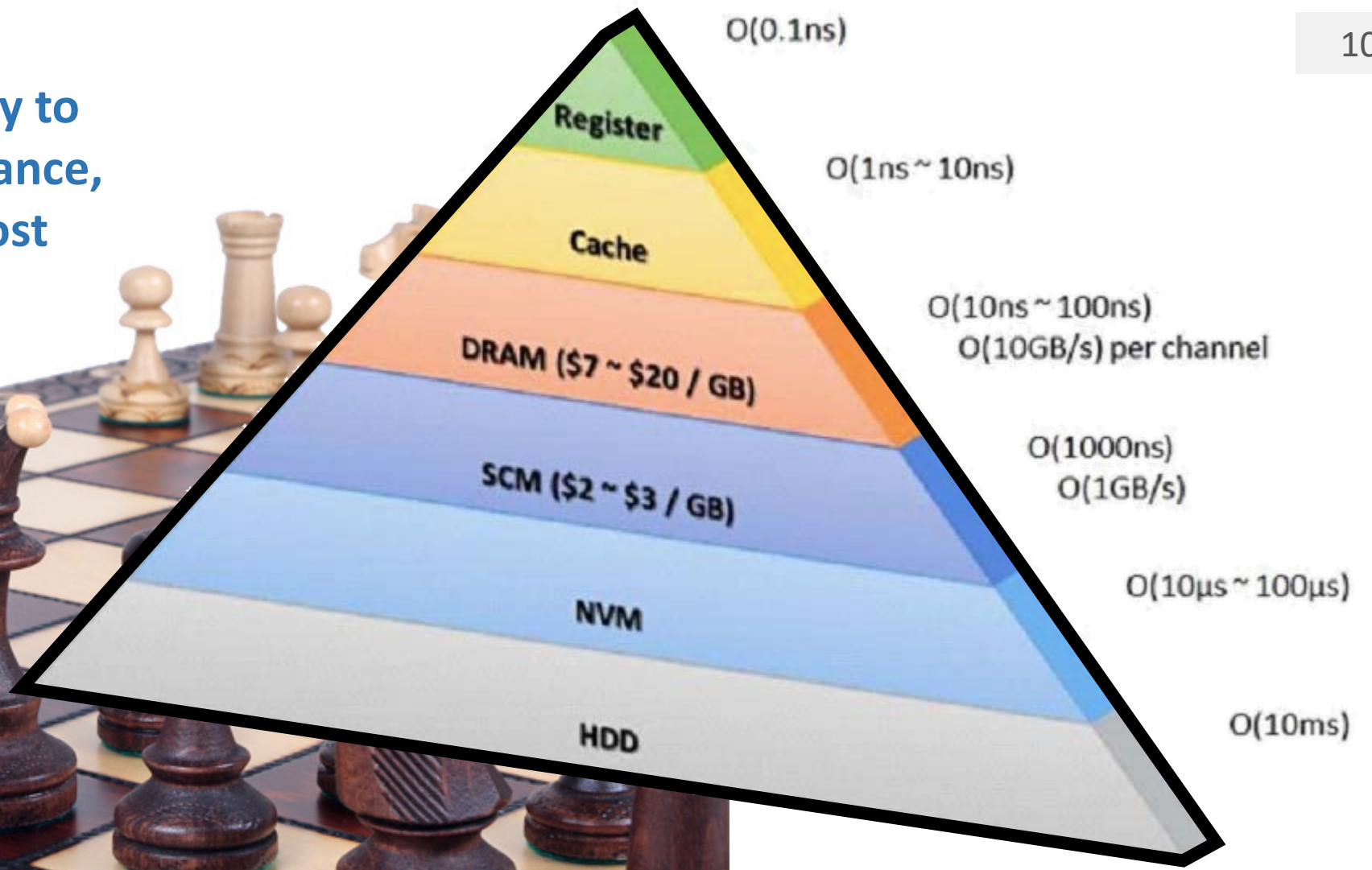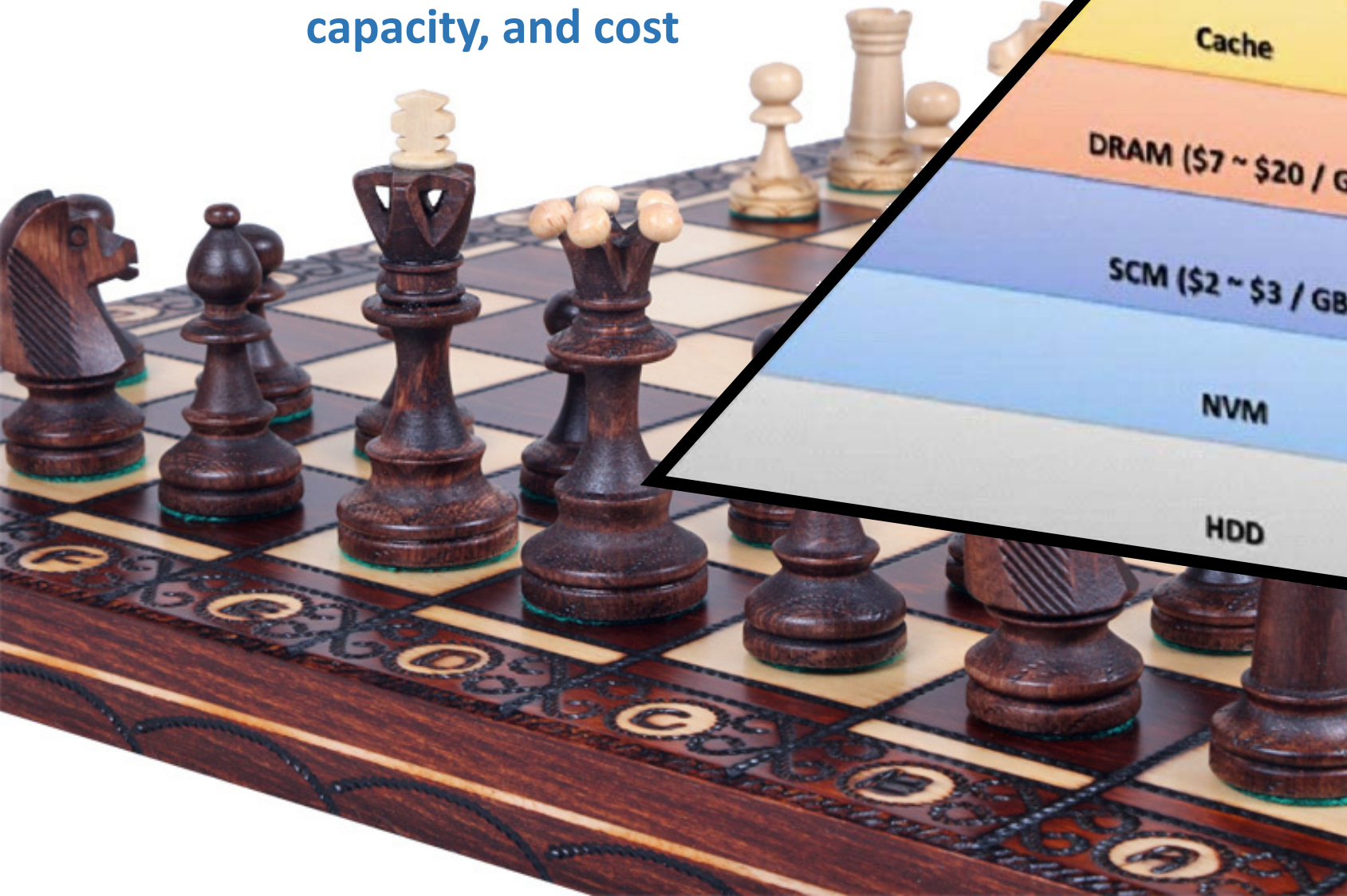**System failure is a key factor in server software design**
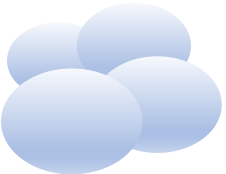
**Data persistence is essential**

**Storage access time impacts transaction granularity**

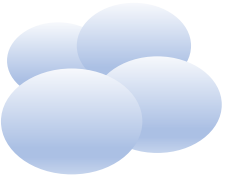**The game we play to trade off performance, capacity, and cost**



Register — O(0.1ns)

O(1ns~10ns)

Cache

O(10ns~100ns)
O(10GB/s) per channel

DRAM ($7 ~ $20 / GB)

O(1000ns)
O(1GB/s)

SCM ($2 ~ $3 / GB)

O(10μs~100μs)

NVM

HDD — O(10ms)

**To reduce the penalties
from checkpointing…**

**…move non-volatile
storage closer to the CPU**

# Traditional Server Architecture Review

Network

**I/O**

**CPU** **$**

**Memory Control**

Storage ... Storage ... Storage

**Faster, lower latency**

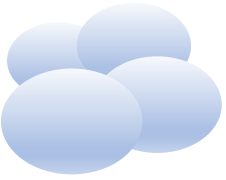Memory Memory : Memory Memory : Memory Memory : Memory Memory

*The Search for*

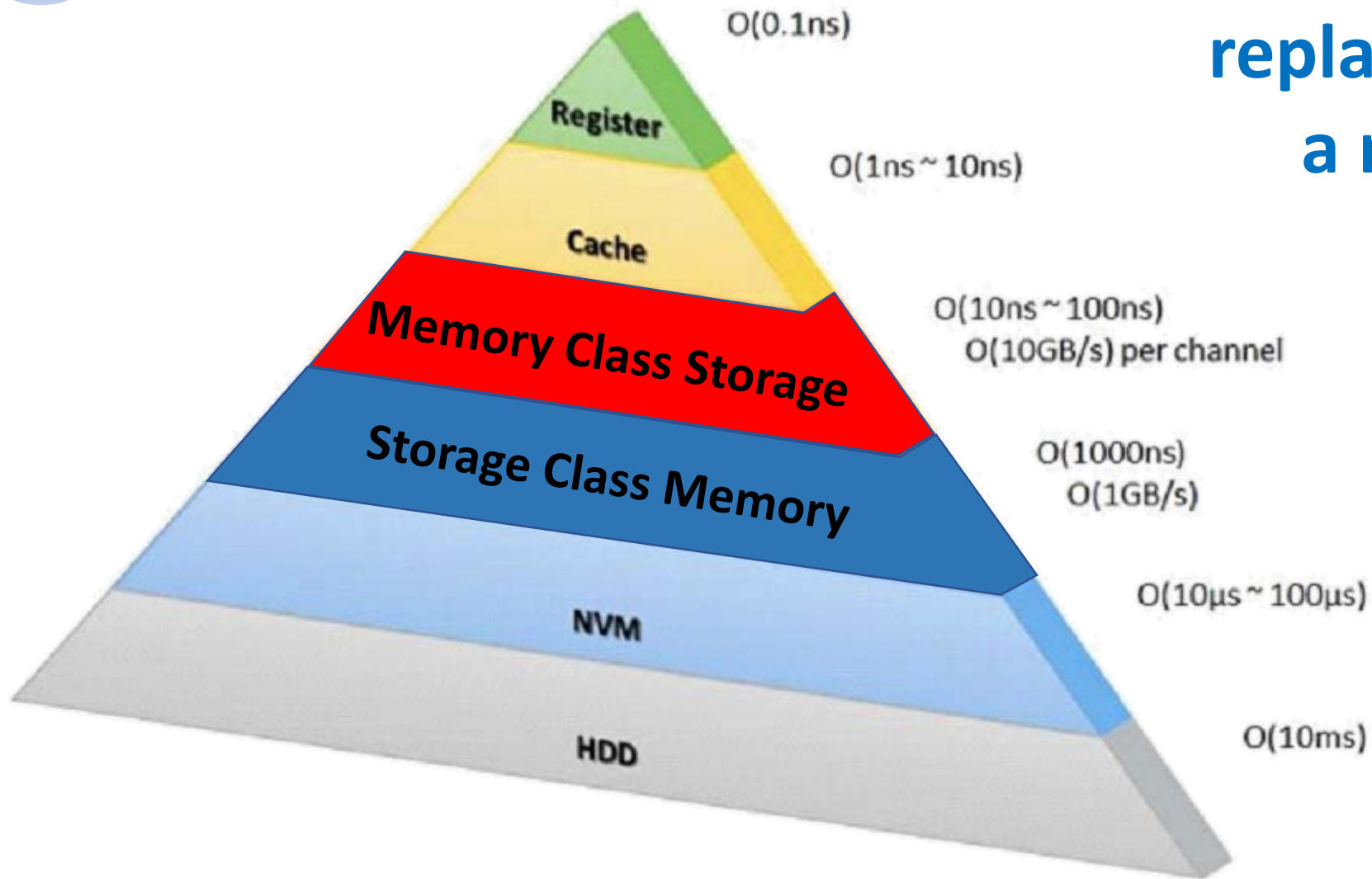# THE HOLY GRAIL
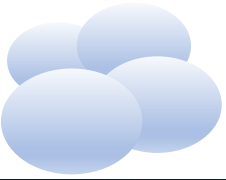
When we no longer
fear power failure…

DATA PERSISTENCE

What if you could replace DRAM with a non-volatile memory?

You'd call it **Memory Class Storage**

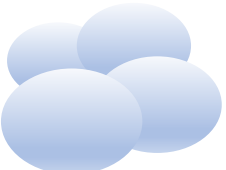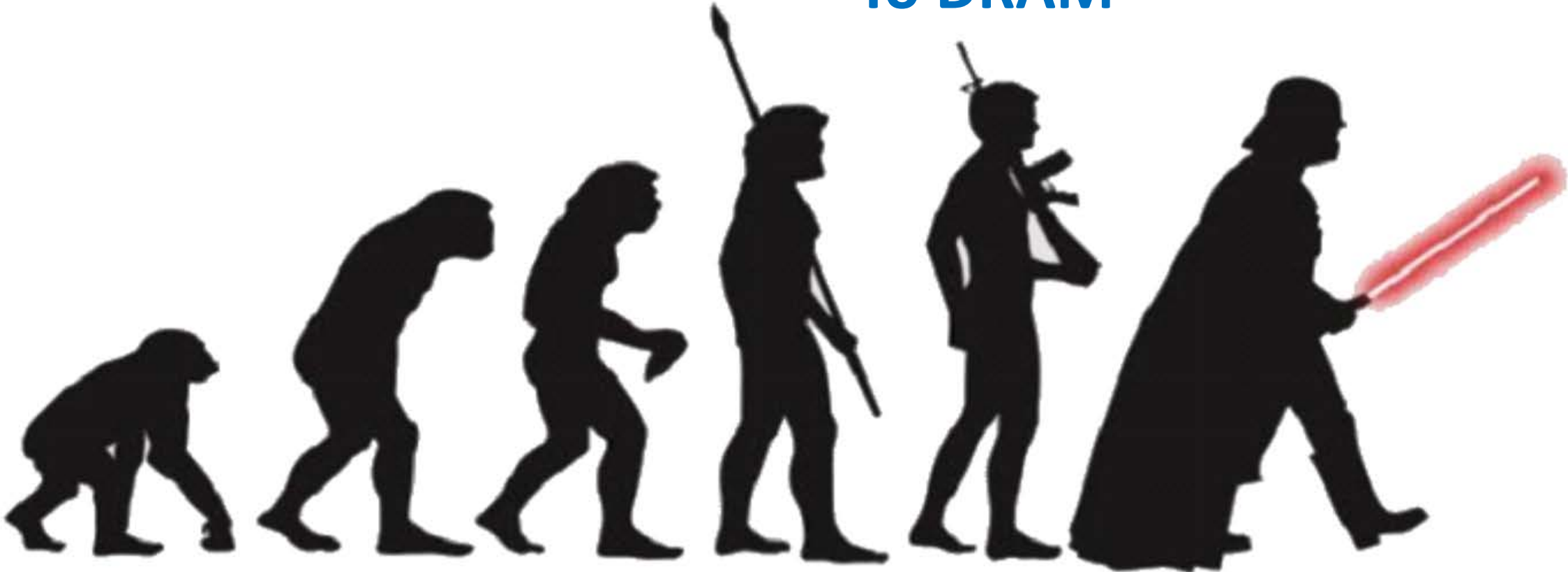**When was the last time you read about a new volatile memory?**

**NRAM™**

**MRAM**

**PCM**

**3DXP**

**ReRAM**

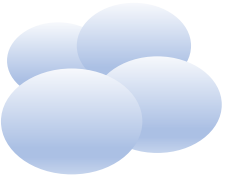**The non-volatile memory revolution is under way**

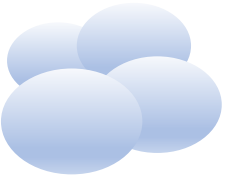**To NVRAM**

**To DRAM**

**To core memory**

**From vacuum tubes**

**THIS is why the term "Persistent Memory" is insufficient**

**The industry must distinguish between deterministic and non-deterministic persistent memory**

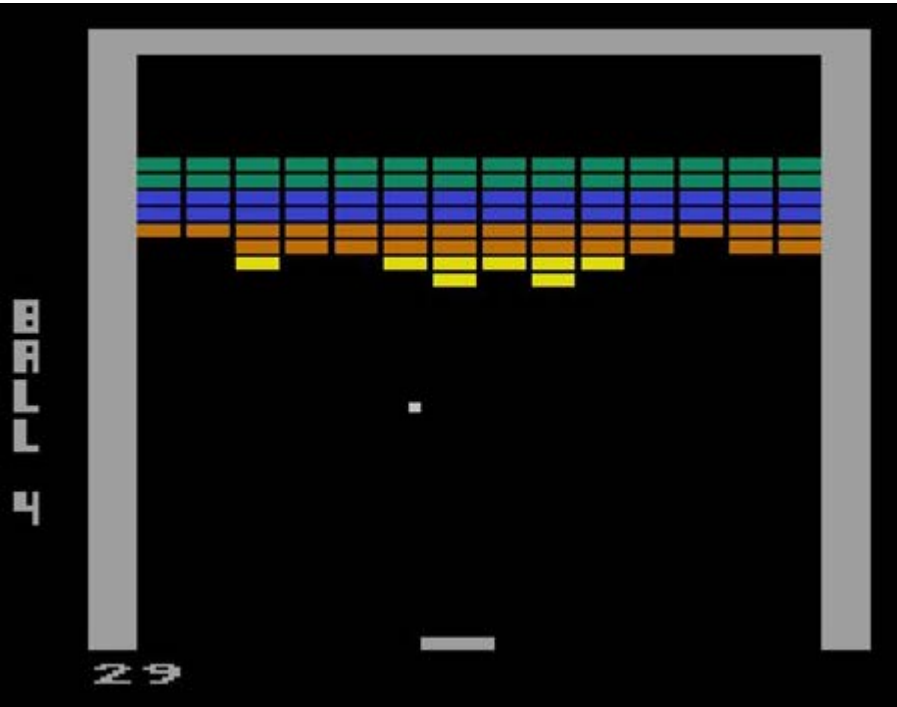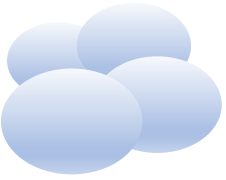**Only "Memory Class Storage" is fully deterministic AND persistent**

NRAM

FeRAM

3DXpoint

MRAM

ReRAM

Flash

DRAM

SRAM

**Not all "persistence"
is created equal**

| Flash Architecture | Layers of Cells | Bits per Cell | Number of Cell Voltage States | Cell Endurance[1] (P/E Cycles) |
|---|---|---|---|---|
| Planar SLC | 1 | 1 | 2 | ~100,000 |
| Planar MLC | 1 | 2 | 4 | ~3,000 |
| Planar eMLC/iMLC/pSLC | 1 | 1 | 2 | ~20,000 |
| Planar TLC | 1 | 3 | 8 | <1,000 |
| Vertical SLC | Varies, 64 typical | 1 | 2 | TBD[2] |
| Vertical MLC | Varies, 64 typical | 2 | 4 | TBD[2] |

|  | 375GB Intel DC P4800X | 1.6TB Intel DC P3700 | 1.6TB Intel DC P3608 | 2.4TB Micron 9100 Max | 2.7TB Mangstor MX6300 |
|---|---|---|---|---|---|
| Endurance Per Usable GB | 32.8 TB | 27.35 TB | 5.45 TB | 2.73 TB | 12.77 TB |
| Usable Capacity | 375GB | 1.6TB | 1.6TB | 2.4TB | 2.7TB |
| Raw Capacity | 448GB | 2TB | 2.3TB | 4TB | 4TB |
| Spare Area / % | 73GB / 16.3% | 400GB / 20% | 700GB / 30.4% | 1600GB / 40% | 1300GB / 32.5% |
| Media Endurance Per Raw GB | 27.4TB | 21.9TB | 3.8TB | 1.6TB | 8.63TB |

**"Write endurance"
determines HOW persistent**

**Wear leveling needed if writes are limited**

# Temperature sensitivity impacts long term retention

| Application Class | Workload | Active Use (power on) | Retention Use (power off) | Functional Failure Rqmt (FFR) | UBER |
|---|---|---|---|---|---|
| Client | Client | 40ºC 8 hrs/day | 30ºC 1 year | ≤3% | ≤$10^{-15}$ |
| Enterprise | Enterprise | 55ºC 24hrs/day | 40ºC 3 months | ≤3% | ≤$10^{-16}$ |

## Client

| Power Off Temperature | | | | | | |
|---|---|---|---|---|---|---|
| 55 | 1 | 1 | 2 | 2 | 3 | 5 | 8 |
| 50 | 2 | 2 | 3 | 4 | 6 | 9 | 15 |
| 45 | 4 | 4 | 5 | 7 | 10 | 17 | 27 |
| 40 | 7 | 8 | 10 | 14 | 20 | 31 | 52 |
| 35 | 14 | 16 | 20 | 26 | 38 | 61 | 101 |
| 30 | 28 | 32 | 39 | 52 | 76 | 120 | 199 |
| 25 | 58 | 65 | 79 | 105 | 155 | 244 | 404 |
| | 25 | 30 | 35 | 40 | 45 | 50 | 55 |

Active temp

Weeks of Data Retention

READ WRITE WRITE READ WRITE

DATA DATA DATA DATA DATA

**DRAM interface is deterministic**
**Data latency is FIXED**

READ WRITE WRITE READ WRITE

DATA DATA DATA X X HOUSEKEEPING

**Any endurance limit breaks determinism**

# Memory Class Storage

**Full DRAM Speed**

**No endurance limits**

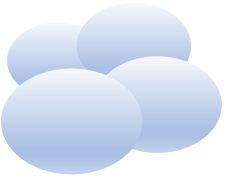**Fully deterministic**

# NVRAM

## is a

**Memory Class Storage**

**In the future?**

Memory Class Storage

=

NVRAM

**For now…**

Memory Class Storage

NVRAM

# Storage Class Memory

## Is **<u>NOT</u>** a

# Memory Class Storage

Hard Disk

Storage

SSD

NVMe

Flash

Wasteland

Storage Class Memory

Phase Change
3DXpoint
Resistive RAM
Magnetic RAM
3D NOR

DDR
DRAM

**Memory Class Storage**

≥ DRAM performance
= DRAM endurance
≥ DRAM capacity

**DDR NVRAM**

**Non-Deterministic**

**Non-Deterministic**

**Non-Deterministic**

**Deterministic**

**Deterministic**

**Deterministic**

**DRAM**

**NVDIMM-N**

**Optane**

**NVRAM Memory Class Storage**

**NVDIMM-P**

**DRAM**

**Itty bitty leaky capacitors lose charge**

**On power fail, you lose**

| | | | REFRESH | |
|---|---|---|---|---|
| ACT | ACT | ACT | | ACT |
| RD | RD | RD | | RD |
| WR | WR | WR | | WR |
| PRE | PRE | PRE | | PRE |

**Refresh time consumes up to 15% of bandwidth**

**DRAM**

**Run**

**FAIL!**

# NVDIMM-N

**Use DRAM normally**

**On Power Fail, copy to Flash**

**Power restored, copy to DRAM**

**Energy Source**

**NVDIMM-N**

Flash Backup

DRAM Array

NVM Control

Voltage Regulator

Isolation Buffers

**Host System**

Power Fail

CPU

Voltage Regulator

NVDIMM-N

Run

FAIL!

RESTORE

Switch to Battery Power

Copy Flash to DRAM

Copy DRAM to Flash

Run

**NVDIMM-N**

Amazon.com Goes Down, Loses $66,240 Per Minute

**1-2 MINUTES**

Copy DRAM to Flash

Copy Flash to DRAM

**1-2 MINUTES**

**One power fail cycle pays for a LOT of protection**

# Faster than Flash!!!

# But vs DRAM?  Meh

# Decent capacity, though

**Optane**

**3DXpoint Array**

**NVM Control**

**CPU**

**Host System**

**Reads are slow**

**Writes are deathly slow**

RD

→ Data

WR

→ Data

**Could be used as a very slow DRAM but more common as expansion**

# NVDIMM-P

**Small Energy Source**

Non-Volatile Memory Array – Any Kind

Voltage Regulator

DRAM Cache

NVM Control

## NVDIMM-P Protocol

| Read A | Read B | | Send | Read C | | | Send | | Send | |
| RSP | | Data A | RSP | RSP | | Data C | | Data B | |

CPU

**Host System**

**New non-deterministic protocol**

**Not backward compatible with DDR**

**Requires NVDIMM-P aware CPU**

Volatile Mode
No Persistence

Battery Backup
ala NVDIMM-N

NVDIMM-P
Persistence
Options

Explicit FLUSH
Command

Reduced
Energy,
Cacheless

# NVRAM

**DRAM speed**

**Non-volatility**

**Unlimited write endurance**

**Wide temperature range**

**Scalable beyond DRAM**

**Flexible fabrication & application**

**Low power**

**Low cost**

**Drop in replacement for DRAM**

**Fully Deterministic**

**Permanently persistent**

**Always available**

**NVRAM Memory Class Storage**

**DRAM**

**Host System**

NRAM™

PCM *

DDR5 NVRAM

ReRAM *

MRAM *

* Future generation devices

JEDEC
STANDARD

DDR5 SDRAM

JESD79-5

PROPOSED

JEDEC SOLID STATE TECHNOLOGY ASSOCIATION

JEDEC

JEDEC
STANDARD

DDR5 Non-Volatile Random Access
Memory (DDR5 NVRAM)
Addendum #1 to JESD79-5

JESD79-5-1

PROPOSED

JEDEC SOLID STATE TECHNOLOGY ASSOCIATION

JEDEC

DDR5 NVRAM is "like a DRAM and…"

# Comparing DRAM & NVRAM

**No refresh is required**

**"Self refresh" can be power OFF**

**Some timing differences (but deterministic!)**

**Data persistence definitions**

**Greater per-die capacity**

NRAM™ ≠ PCM

ReRAM MRAM

**Timings** **Precharge requirement** **Persistence definition**

# DDR5 NVRAM Specification brings coherence

# DRAM

# NVRAM



"350 ns"

"0 ns"

IDLE

REFRESH

IDLE

REFRESH = NOP

**Refresh command is not needed**
**Decoded as NOP for compatibility**

# DRAM

# NVRAM



**DRAM:**
- IDLE
- SELF REFRESH
- REFRESH
- FREQUENCY CHANGE

**Power burned**

**NVRAM:**
- IDLE
- SELF REFRESH
- FREQUENCY CHANGE

**"No" power burned**

# DRAM

# NVRAM



**Precharge command is not needed**
**Decoded as NOP for compatibility**

**Persistence Definitions\***

**Intrinsic:**
**Immediately**
**After**
**WRITE**

**Extrinsic:**
**After**
**FLUSH**
**Command**

**Power Fail:**
**On**
**NVRAM**
**RESET**

\* Discussions on-going

**DDR5 DRAM
is limited
to 32Gb per die**

**DDR5 NVRAM
enables up to
128Tb per die**

**DDR5 SDRAM**

ACT RD WR ACT RD WR ACT RD WR

**DDR5 NVRAM**

REXT ACT RD WR ACT RD WR REXT ACT RD WR

**Row Extension adds up to 12 more bits of addressing**

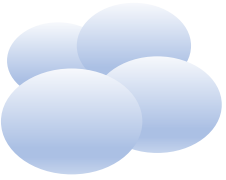**Backward compatible with DDR5 – Acts like REXT = 0 until needed**
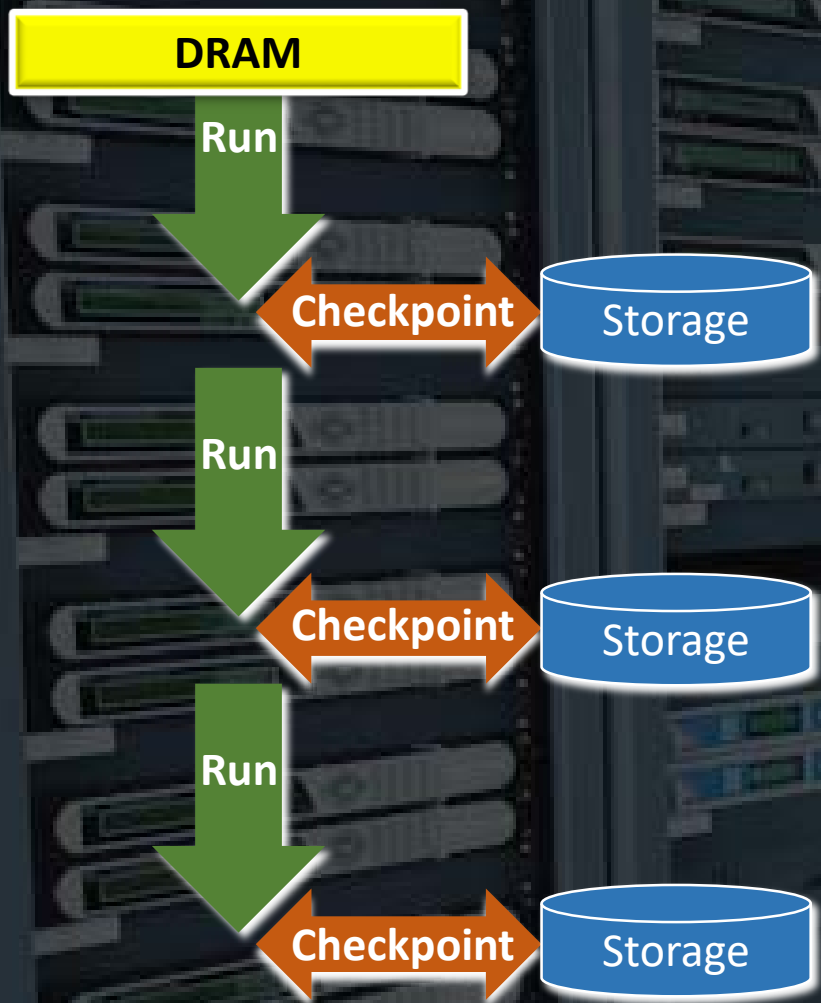
"ROW" includes bank group & bank...

**Row Extension Example**

**Row Extension Replacement Example**

**NVRAM**
**Memory Class Storage**

# Phase 1

# Phase 2

**DRAM**

Run

Checkpoint → Storage

Run

Checkpoint → Storage

Run

Checkpoint → Storage

**NVRAM**

Run

Checkpoint → **NVRAM**

Run

Checkpoint → **NVRAM**

Run

Checkpoint → **NVRAM**

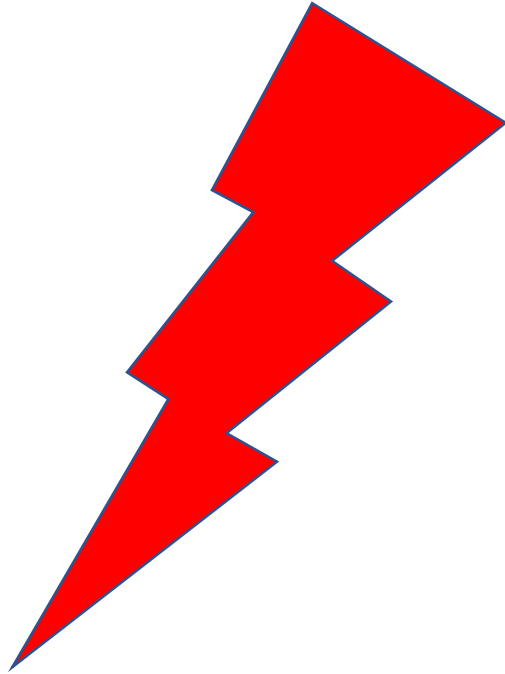**NVRAM**

Run

**No checkpoint**

Run

**No checkpoint**
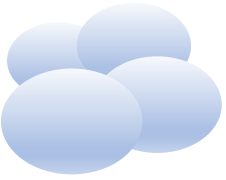
Run

**Keep in mind…**

**Power failure is not the only thing to fear**

**Checkpoints may include system failure**

**Knowing when a task may resume is complicated**

# Remember Those Persistence Definitions

**Immediately After WRITE**
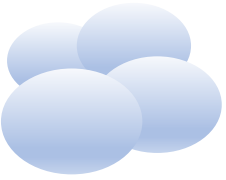
Tasks may be safe in nanoseconds

**After FLUSH Command**

Tasks may be safe in microseconds

**On NVRAM RESET**

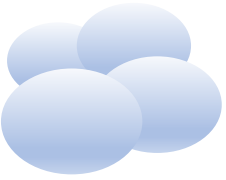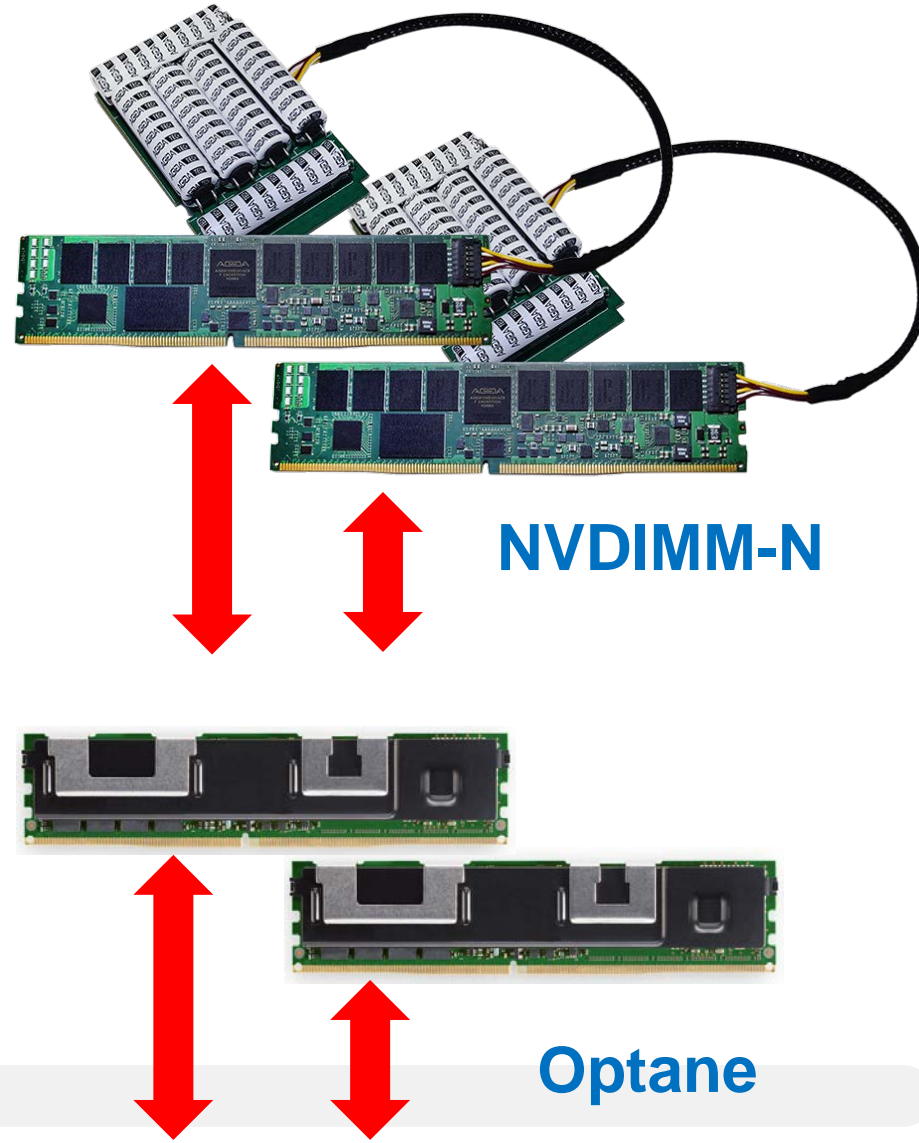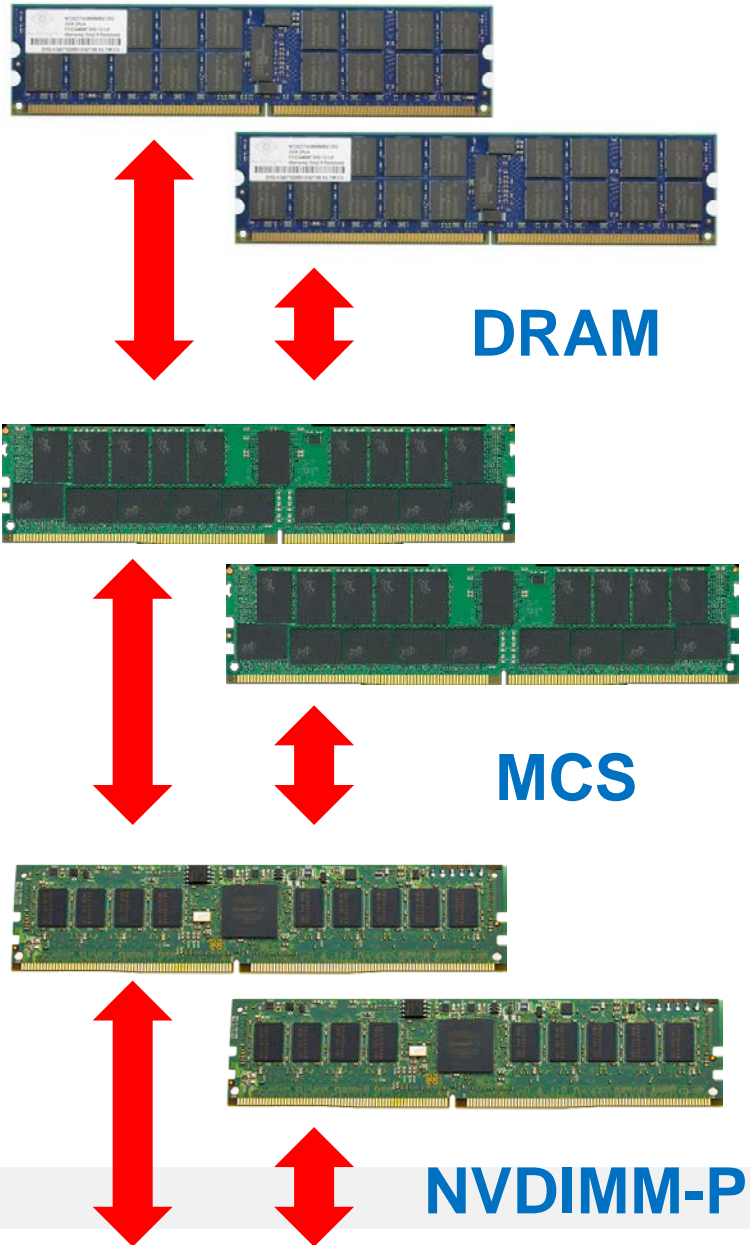Tasks may not be safe until system stability confirmed

# Persistence

## Capacity

## Performance
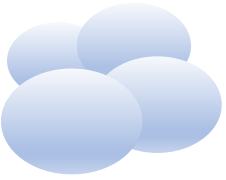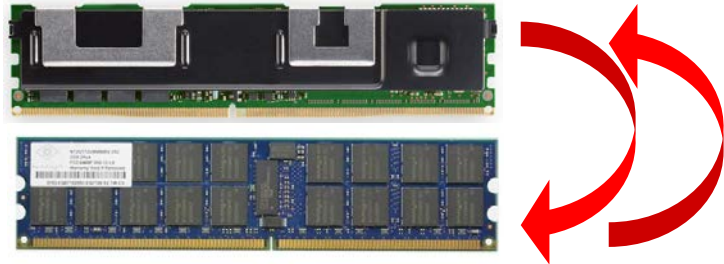
**System designers have
a lot of options to balance**

## Main Memory



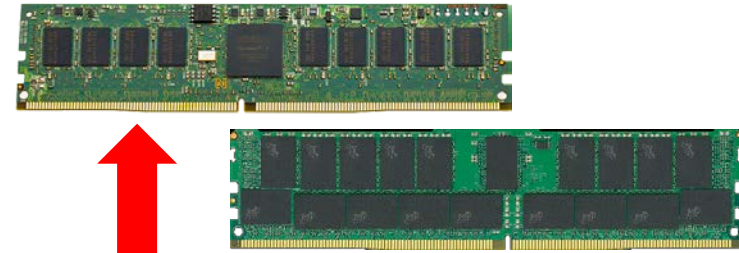**DRAM**

**MCS**

**NVDIMM-P**

**NVDIMM-N**

**Optane**

## **Main Memory**

**DRAM +
Optane**

**MCS +
Optane**

**MCS +
NVDIMM-P**

When capacity meets persistence

**512GB**

**NVDIMM-P**

**Optane**

**64GB**

**DRAM**

**32GB**

**NVRAM**
**Memory Class Storage**

**NVDIMM-N**

# Homogenous
## Main Memory Combinations

| | Data Safe | Performance | Capacity |
|---|---|---|---|
| **DRAM** | No | Best | 1.0 X |
| **NVDIMM-N** | Yes | Best | 0.5 X |
| **Optane** | Yes | Worst | 10 X |
| **NVDIMM-P** | Yes | Mid | 10 X |
| **MCS** | Yes | Best+ | 1 X+ |

# Heterogeneous
## Main Memory Combinations

|  | Data Safe | Performance | Capacity |
|---|---|---|---|
| **DRAM + Optane** | No | High | 6 X |
| **DRAM + NVDIMM-P** | No | High | 6 X |
| **MCS + Optane** | Yes | High | 6 X |
| **MCS + NVDIMM-P** | Yes | High | 6 X |

# Homogenous
## Main Memory Combinations

# Heterogeneous
## Main Memory Combinations

**Software need not care**

**All functions take the same time**
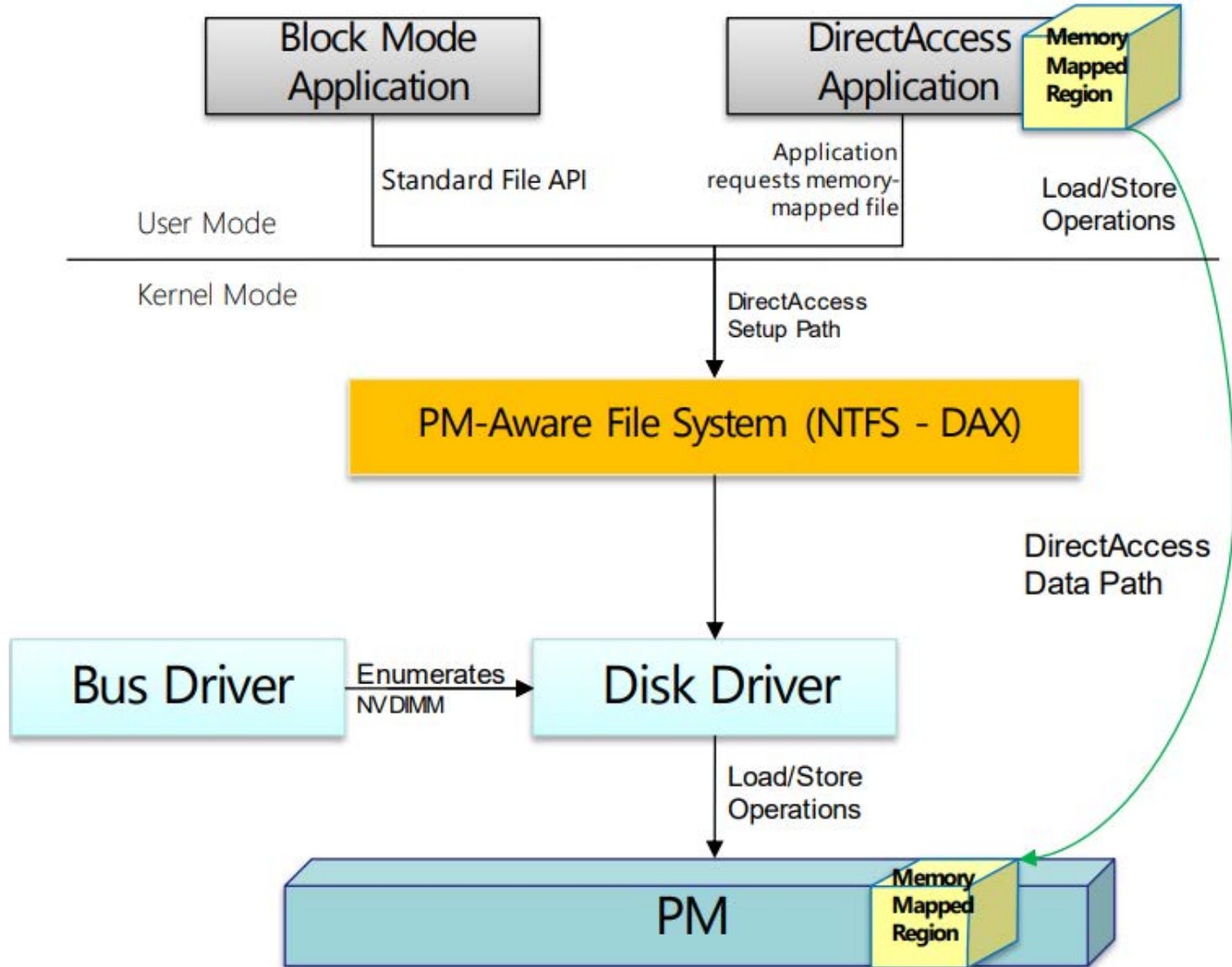
**Software encouraged to put critical functions in faster memory**

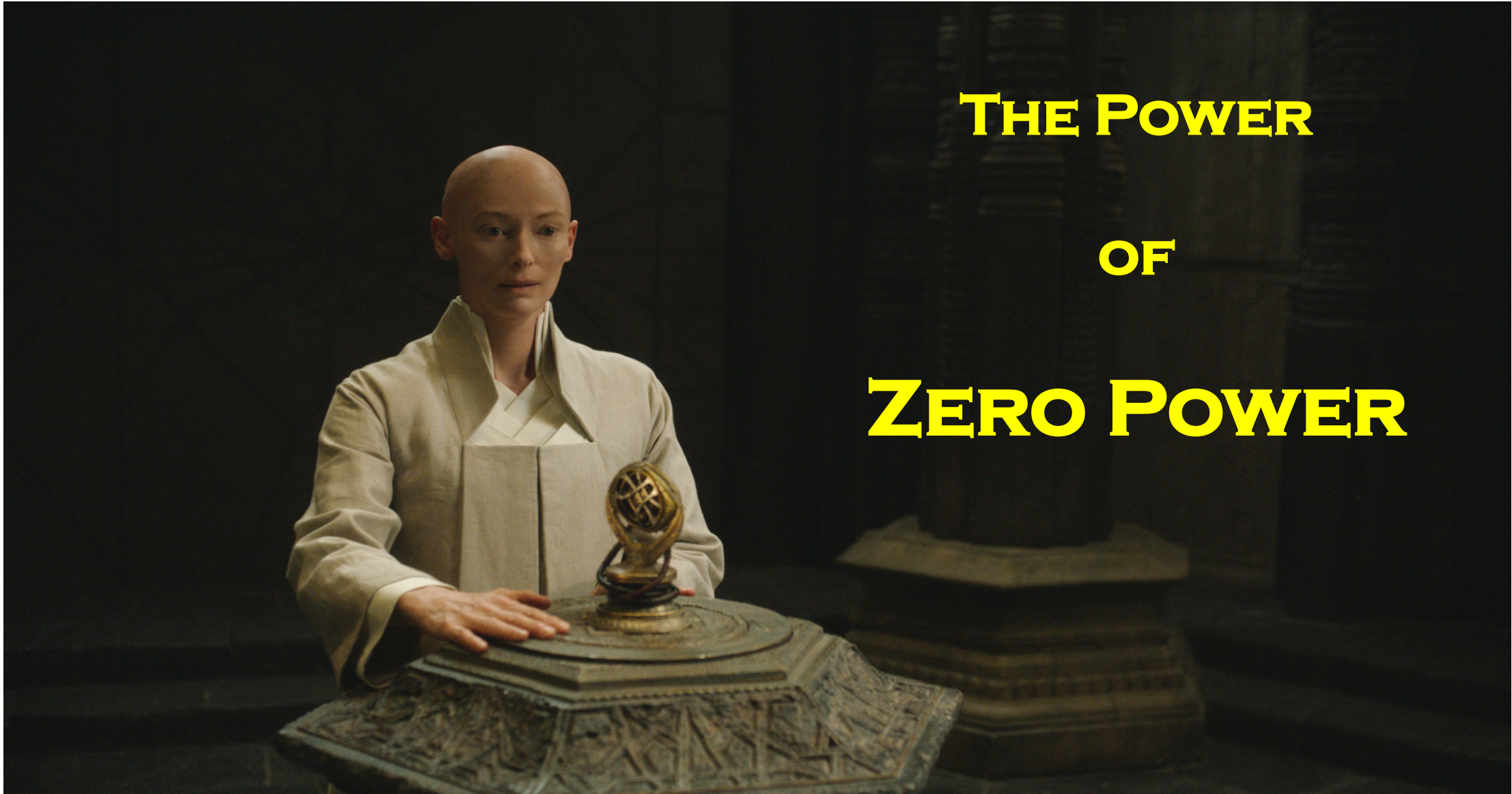**Often mount slower memory as RAM drive**

**Software support via DAX assists in moving…**
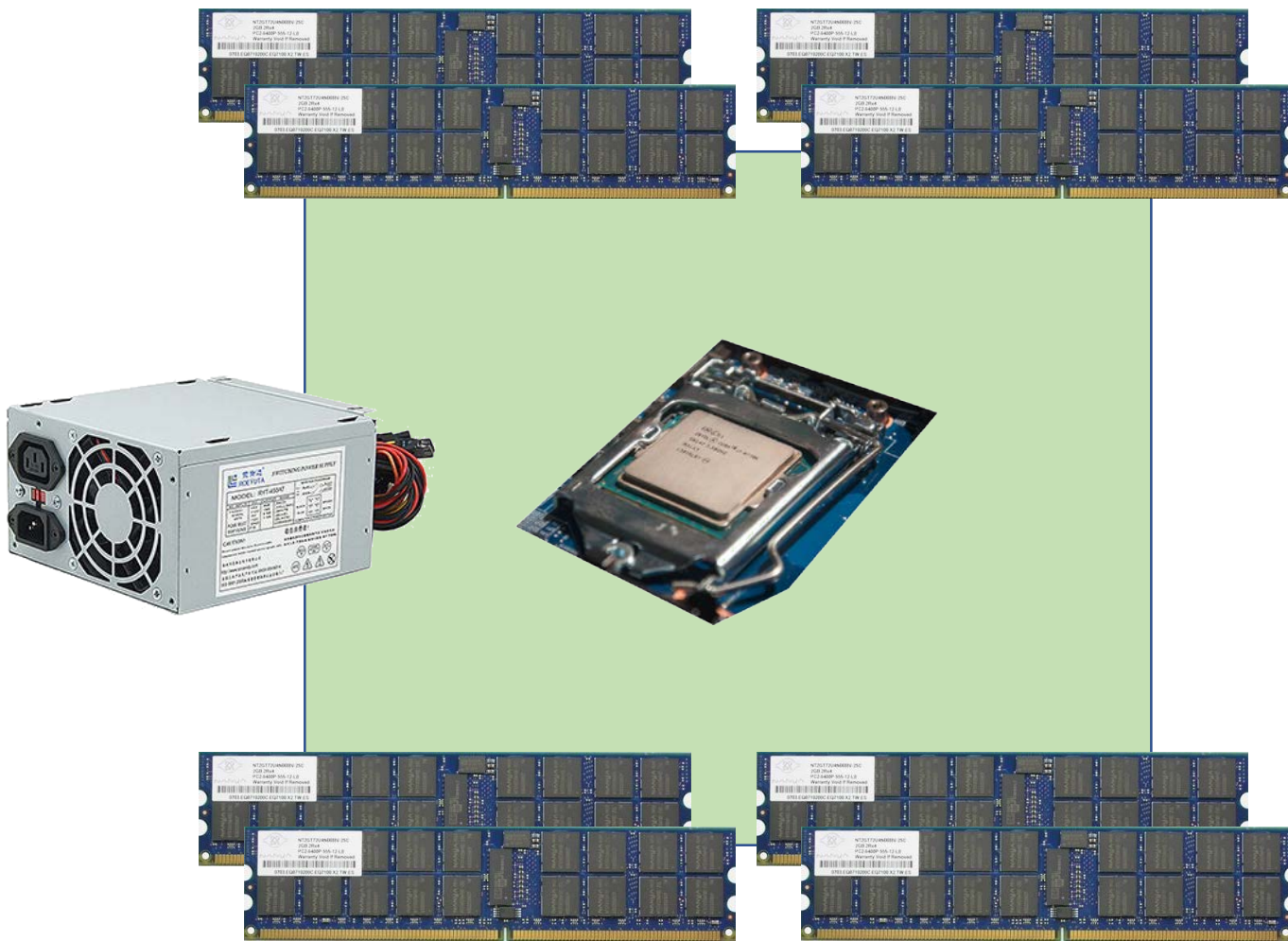
**from mounted drives…**

**…to RAM drive…**

**…to direct access mode**
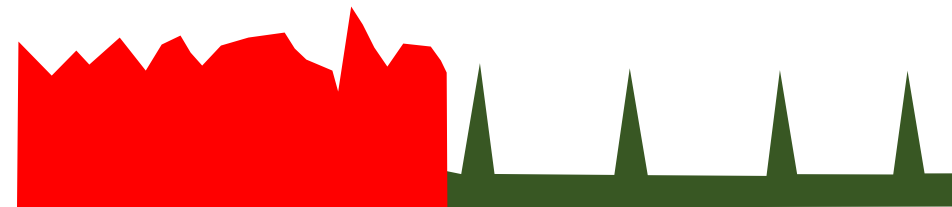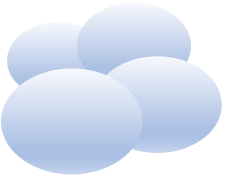
**Putting a Node to Sleep**
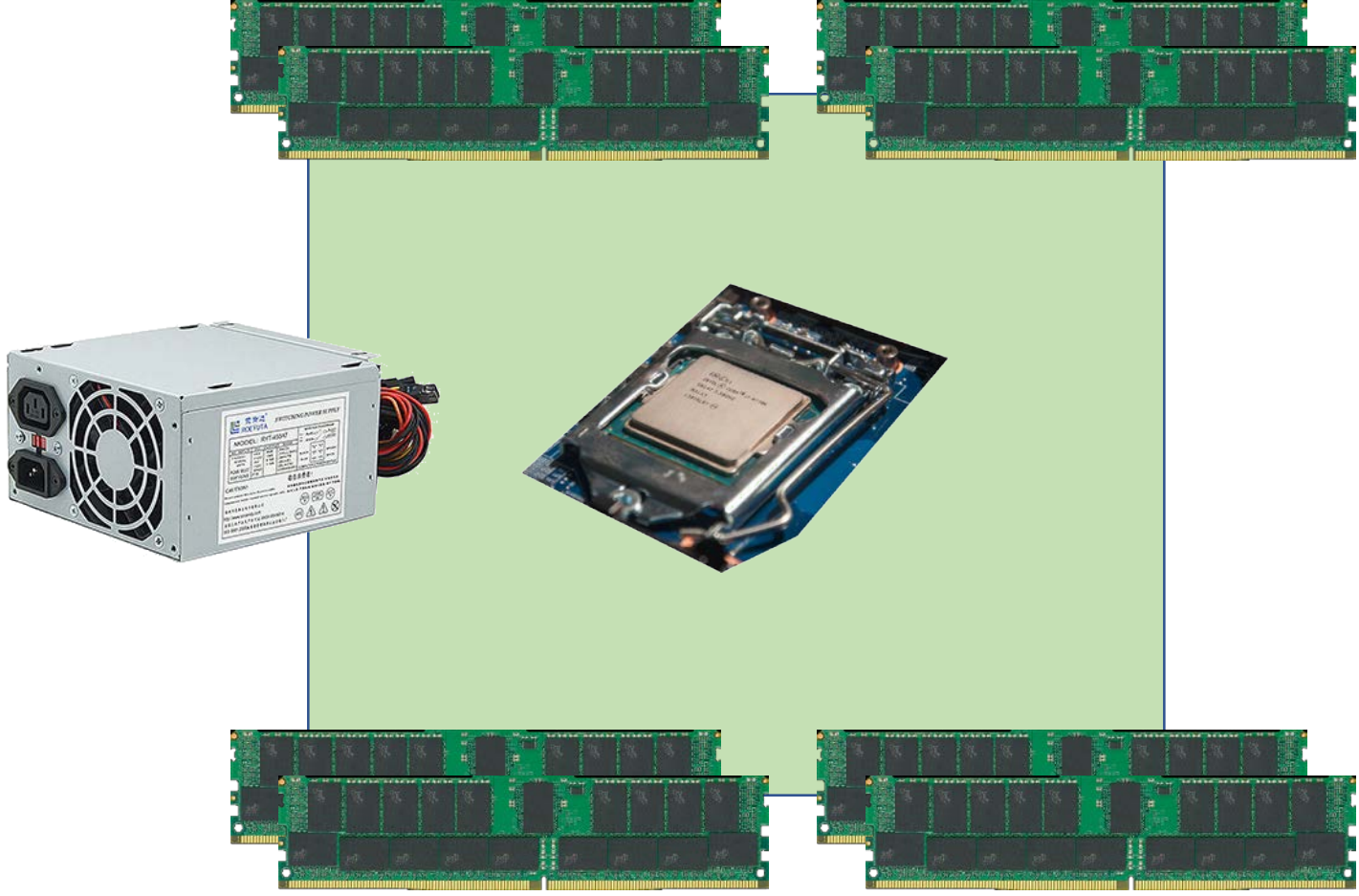
**Operating Mode**　　**Self Refresh Mode**

**Instant On means power must stay alive**
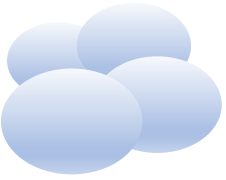
**Refresh operations burn significant power**

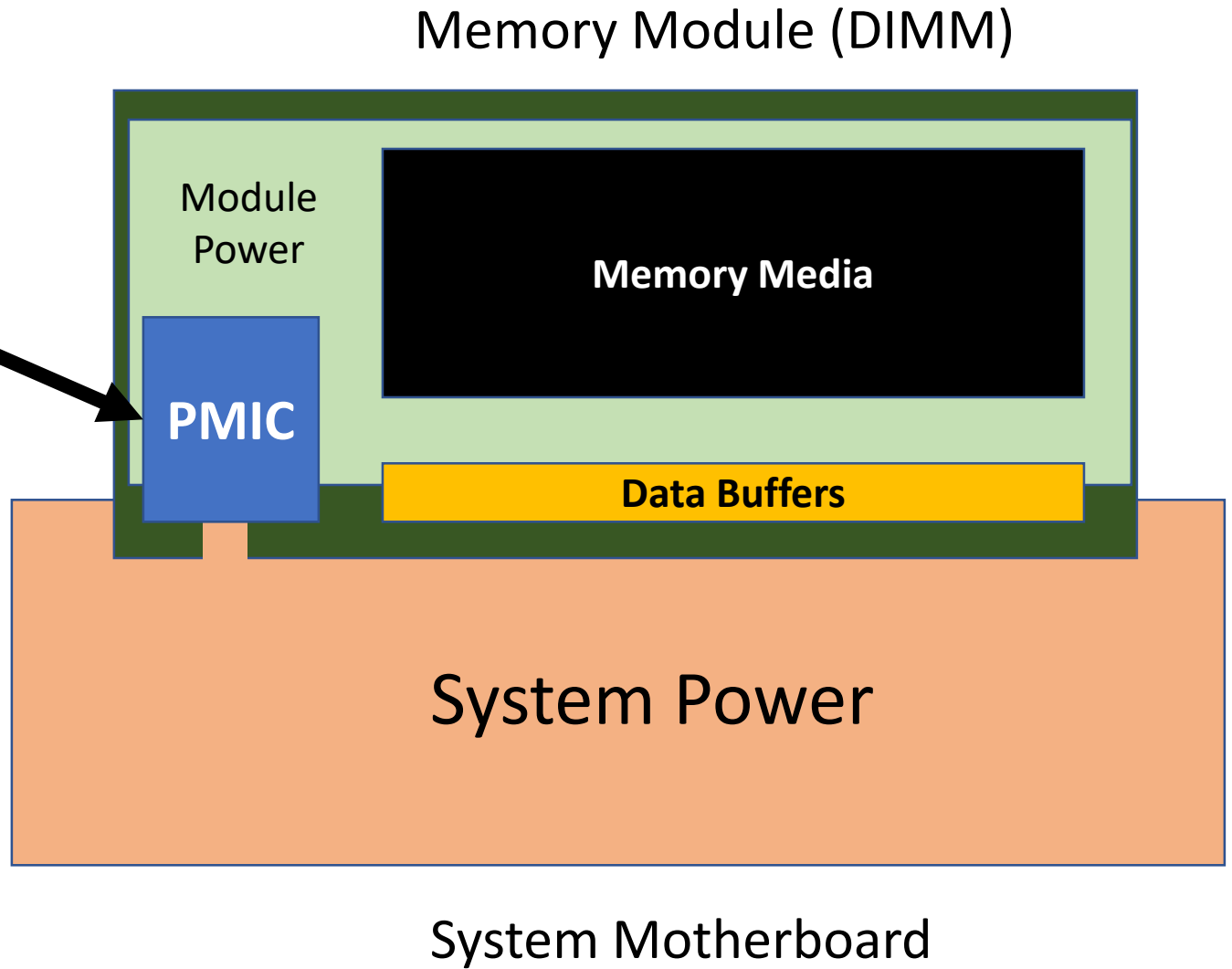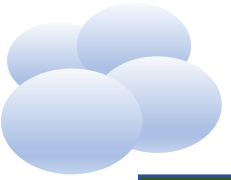**Memory Class Storage can be turned off entirely**

**Operating Mode**

**Power Off**

Memory Module (DIMM)

**DDR5 memory modules have on-DIMM voltage regulation (PMIC)**

Module Power

Memory Media

**PMIC**

**Data Buffers**

**DIMM power may be shut off independently of system power**

System Power

System Motherboard

**DIMM1**

**DIMM2**

Module Power

Memory Media

PMIC

Data Buffers

Module Power

Memory Media

PMIC

Data Buffers
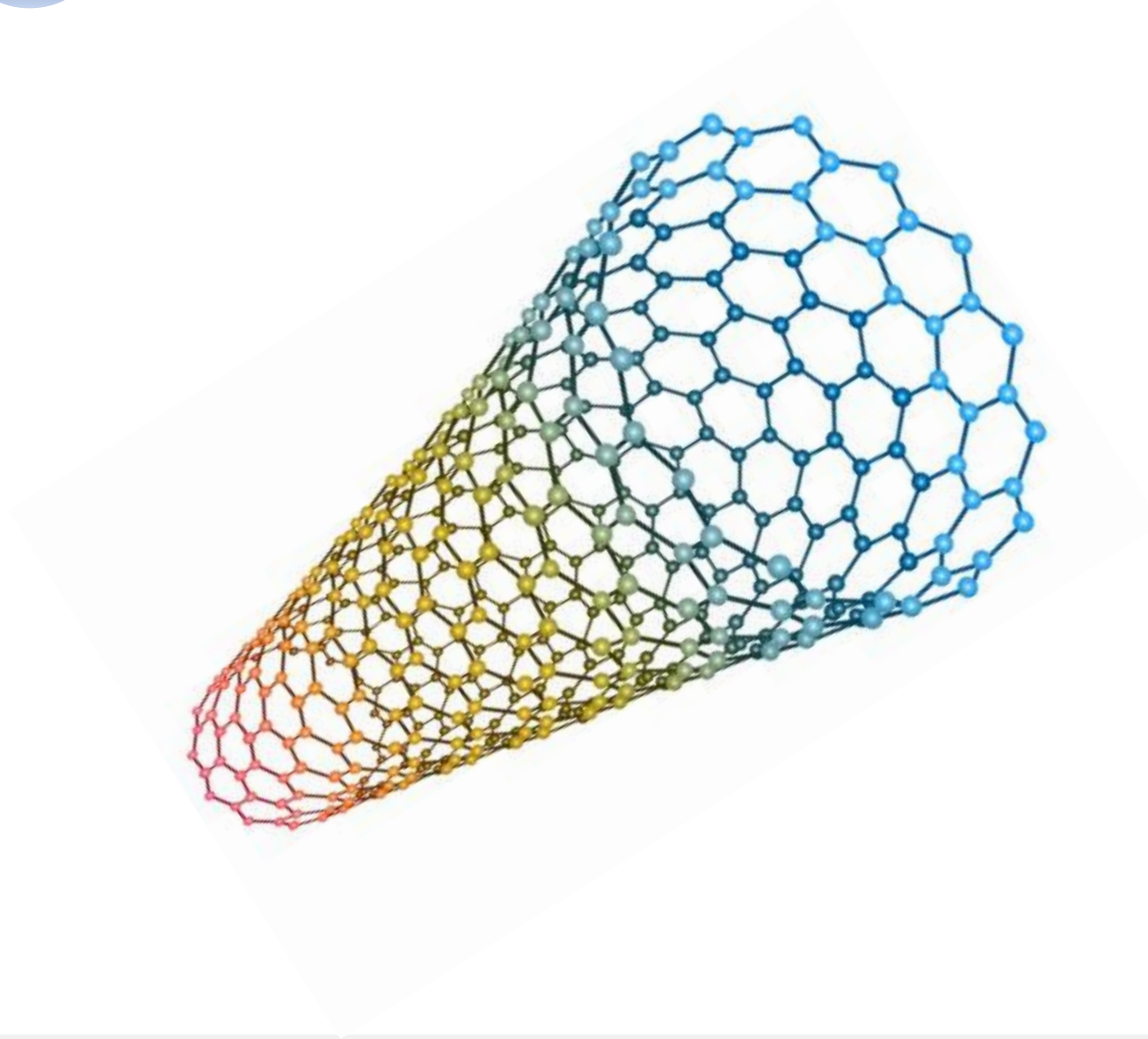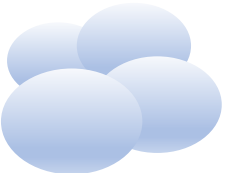
**System Power**

**Multiple power management options**

**System power off; both DIMMs off**

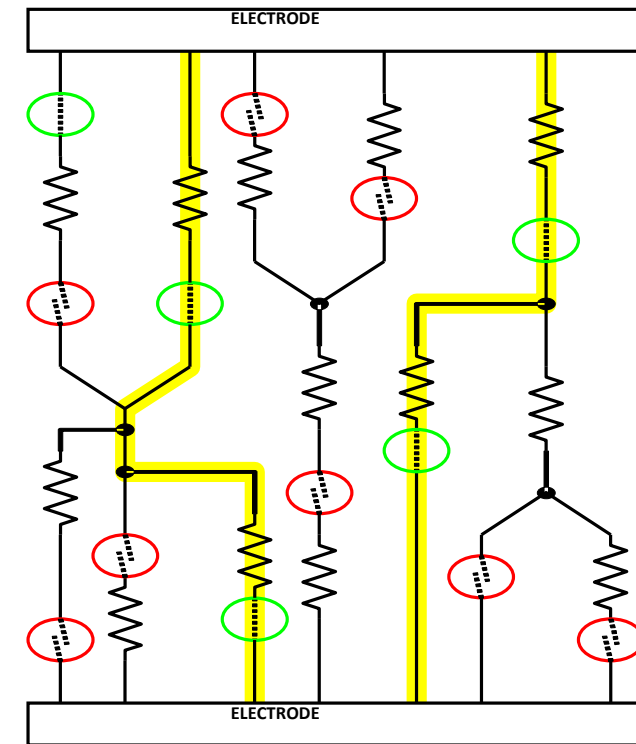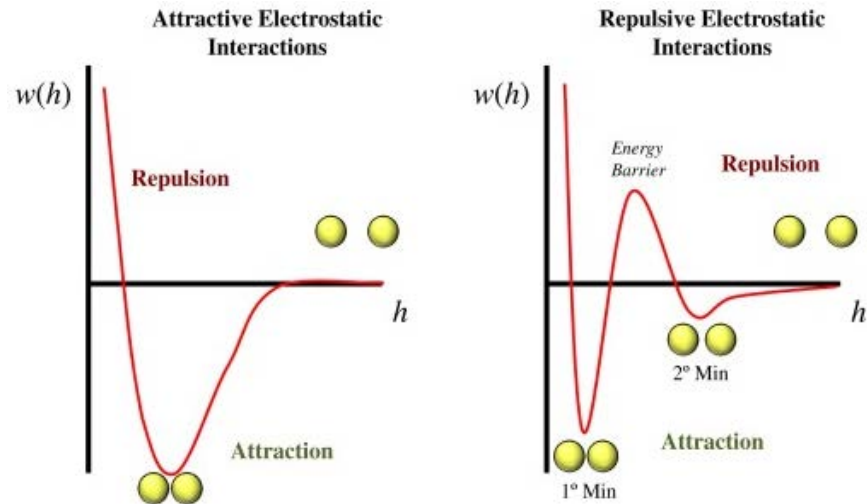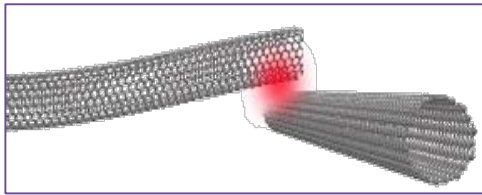**System power on & both DIMMs off**

**System power on & DIMM1 on, DIMM2 off**
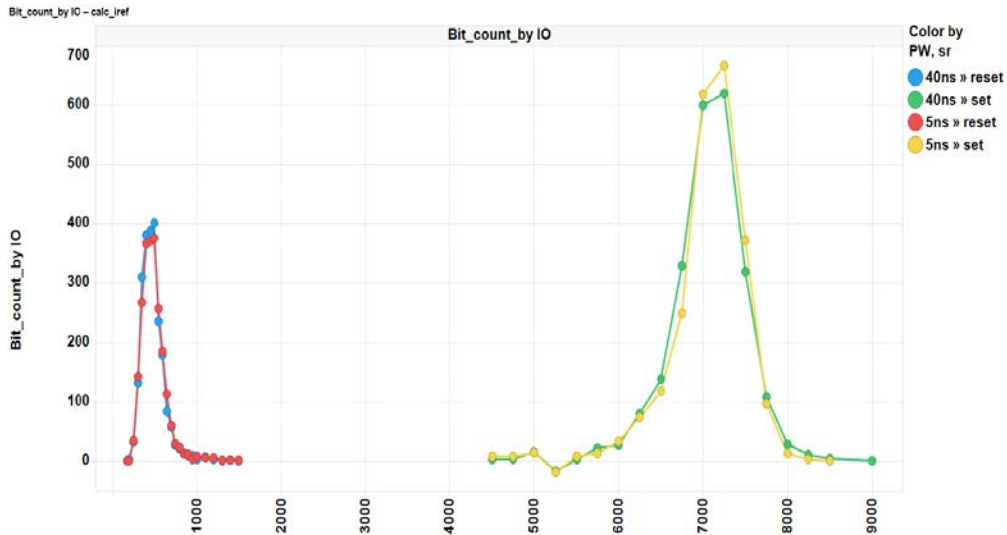
**Nantero NRAM™**

**My favorite NVRAM**

*Full presentation on Wednesday…*

Attractive Electrostatic Interactions

Repulsive Electrostatic Interactions

**Van der Waals energy barrier keeps CNTs apart or together**

**Data retention >300 years @ 300$^\circ$C, >12,000 years @ 105$^\circ$C**

**Stochastic array of hundreds nanotubes per each cell**

Bit_count_by IO – calc_iref

**5 ns balanced
read/write performance**

**No temperature sensitivity**
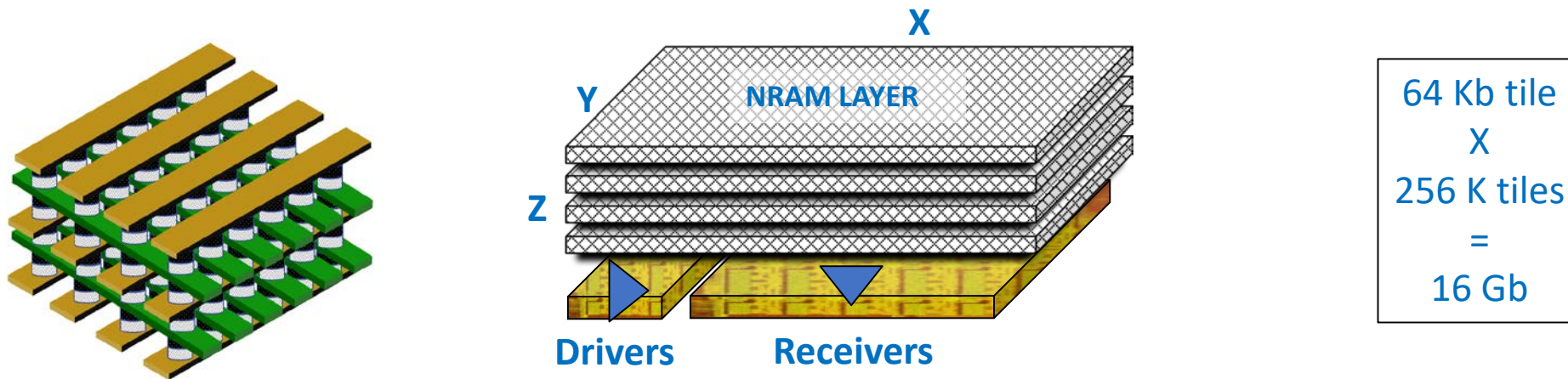
**NRAM Data Retention = 12,000 Years**

**10,000 years ago**

**4,500 years ago**

**2,500 years ago**

X

Y

**NRAM LAYER**

Z

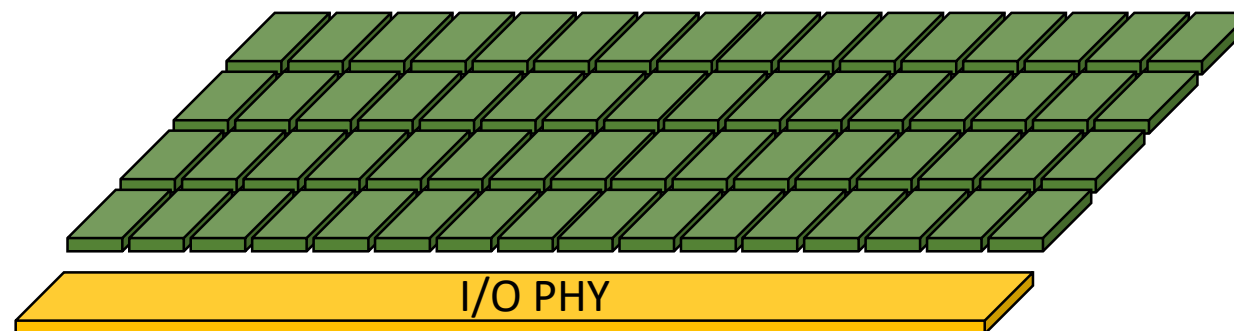**Drivers**   **Receivers**
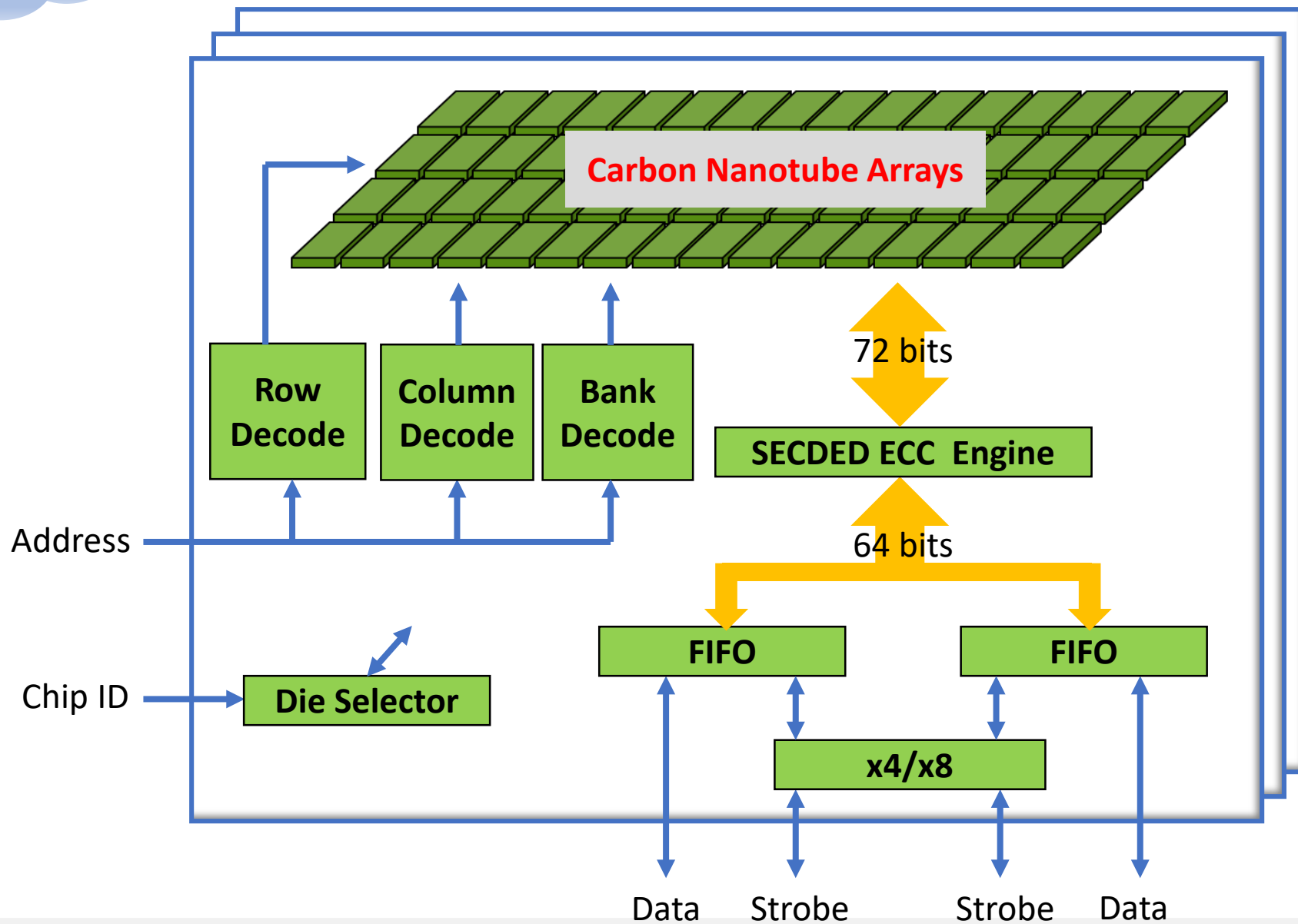
64 Kb tile
X
256 K tiles
=
16 Gb

**Array size tuned to the size of drivers & receivers**

**Chip-level timing is a function of bit line flight times**
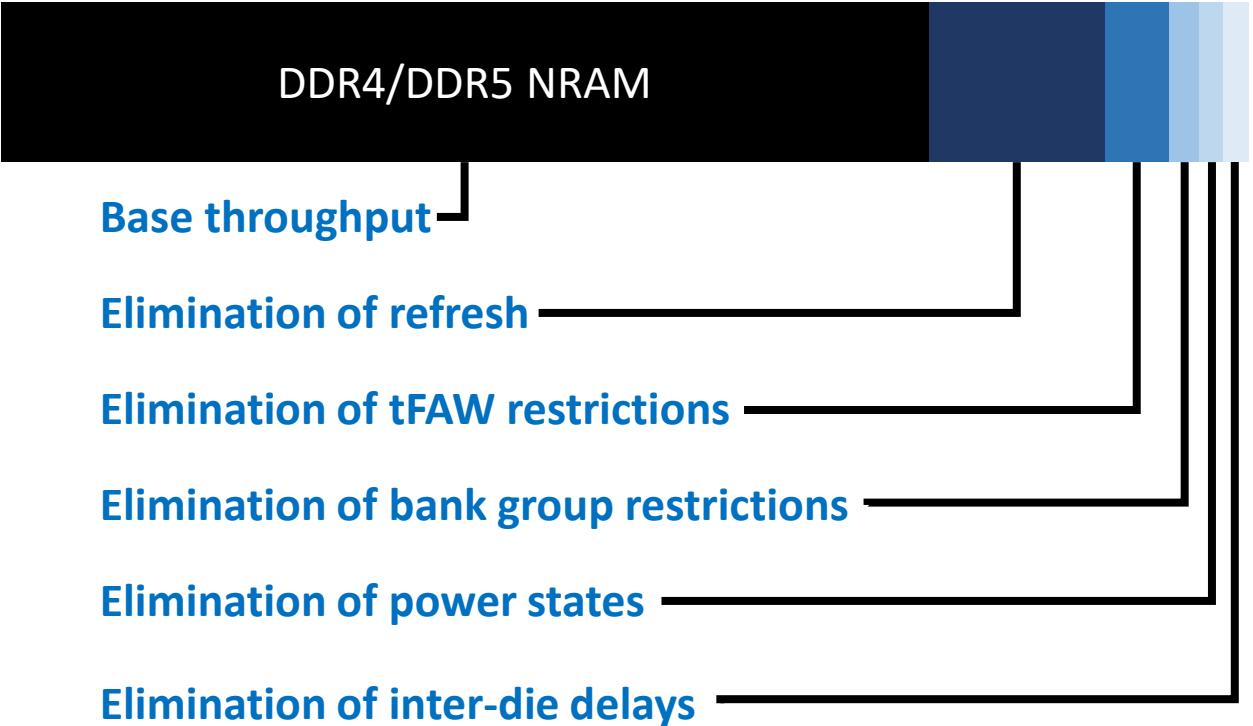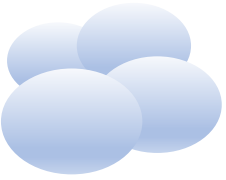
**Replicate this "tile" as needed for device capacity**

**Add I/O drivers to emulate any PHY needed**

I/O PHY
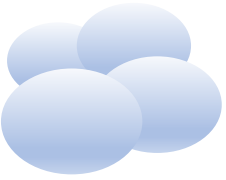
**Carbon Nanotube Arrays**

DDR4, DDR5
NRAM

Row Decode

Column Decode

Bank Decode

SECDED ECC Engine

72 bits

64 bits

Address

Chip ID

Die Selector

FIFO

FIFO

x4/x8

Data    Strobe    Strobe    Data

DDR4/DDR5

15-20%

DDR4/DDR5 NRAM

**Architectural improvements improve data throughput 15% or greater at the same clock frequency**

**Base throughput**

**Elimination of refresh**

**Elimination of tFAW restrictions**

**Elimination of bank group restrictions**

**Elimination of power states**

**Elimination of inter-die delays**

Bandwidth: larger is better

# NVRAM
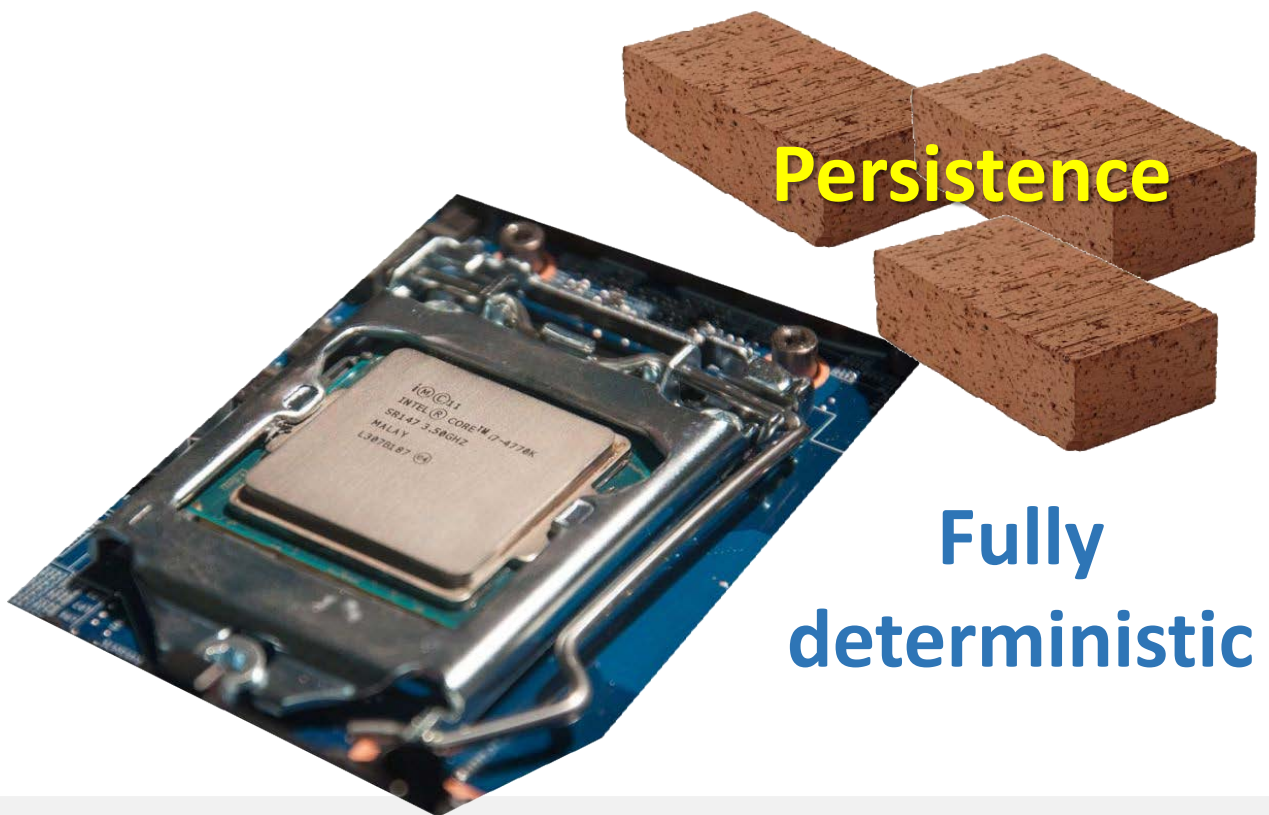## Memory Class Storage



**Plugs into an RDIMM slot**

**Appears to the CPU as DRAM**

**Memory controller may optionally be tuned for NVRAM**

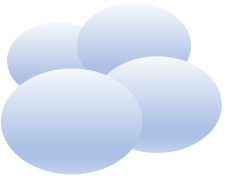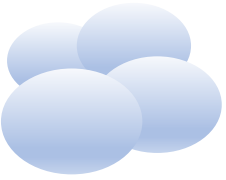**One less layer of marshmallows to deal with**

Persistence

Persistence

Non-deterministic

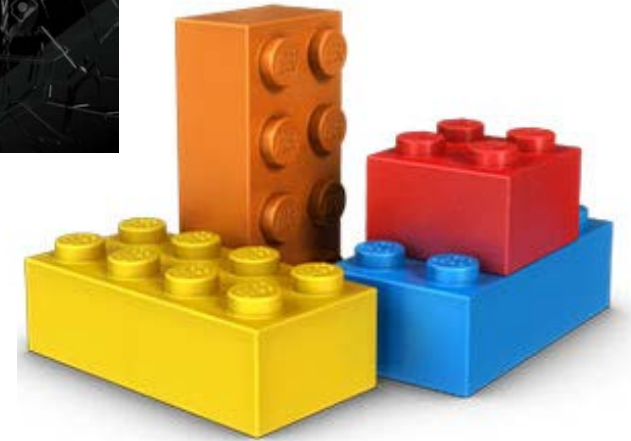Fully deterministic

*Know Your Enemy*

**Would you rather…**

**Step on broken glass?**

**A LEGO?**

**Or some jacks?**

# …about those energy stores…

**Batteries**

**Supercapacitors**

**Tantalums (etc.)**

**Batteries**

**Supercapacitors**

**Tantalums (etc.)**



**High capacity**

**High energy density**

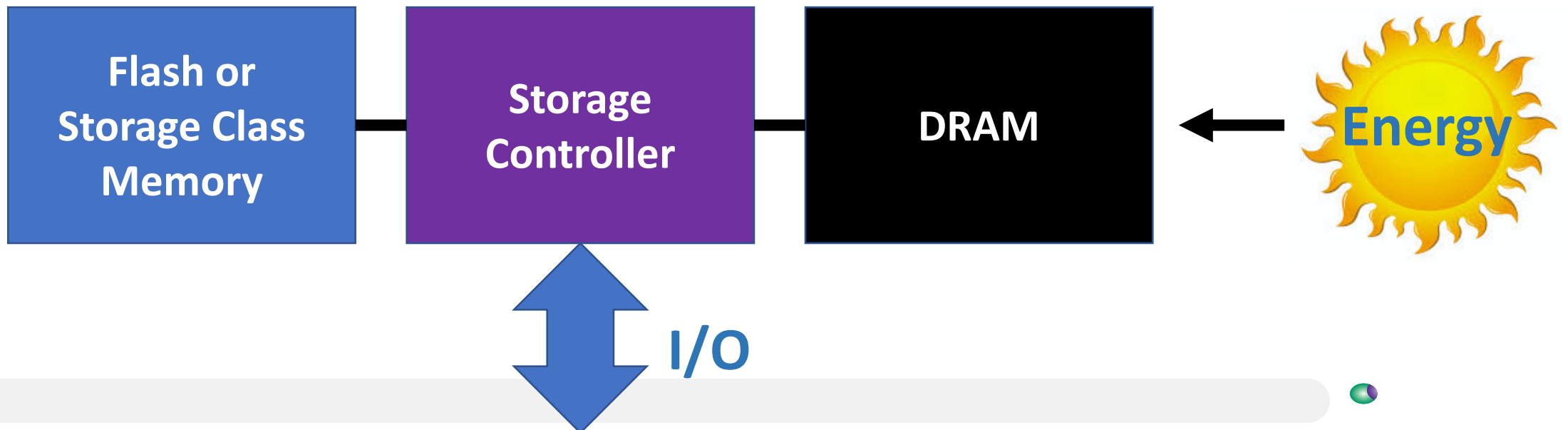**Low reliability**

**Medium capacity**

**Low energy density**

**Degrade over time**

**Low capacity**

**Low energy density**

**...but stable**

**Energy needed for backup of DRAM cache**

| Flash or Storage Class Memory | Storage Controller | DRAM |
|---|---|---|

↕ **I/O**

**Energy**

**More room for storage**

**Eliminate need for backup energy**

**Flash or Storage Class Memory**

**Storage Controller**
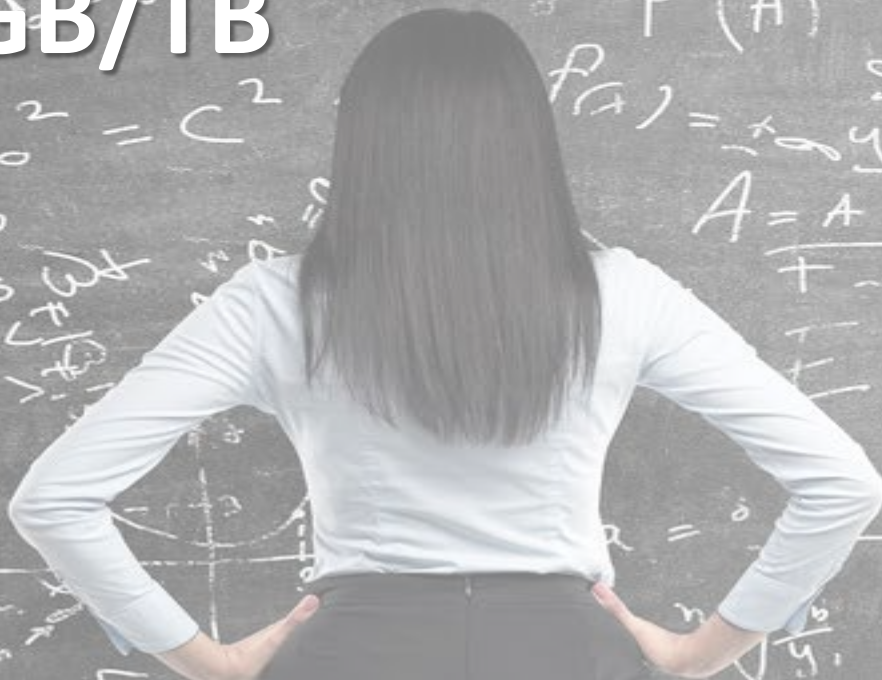
**NVRAM**

**Energy**

**I/O**

# NVRAM Changes the Math

**1GB/TB**

**DRAM cache limited by energy available**

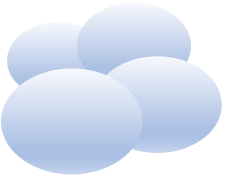**No DRAM? Cache size dictated by cost/performance**
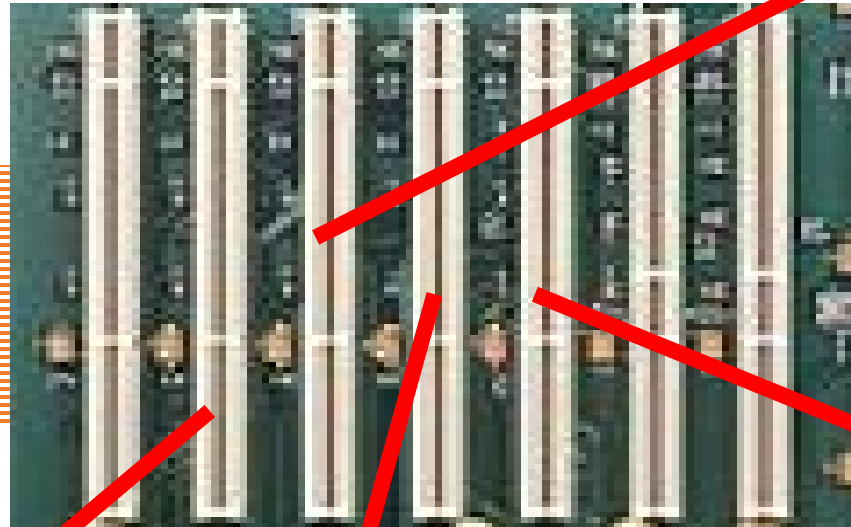
**Switching gears again...**

**...to Systems Evolution**

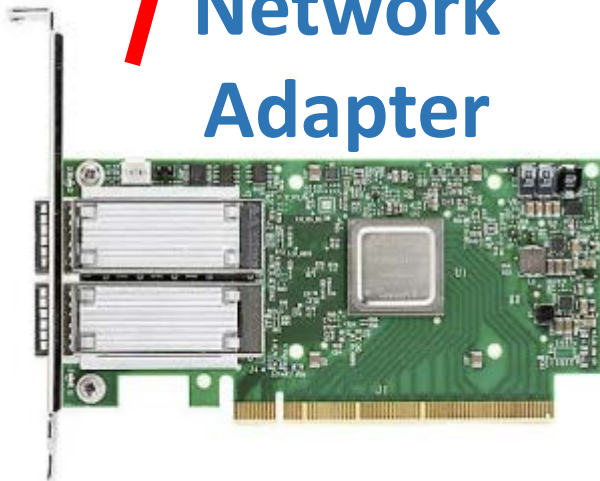**Pop quiz**

**How many CPUs in a 1980s PC?**

**Modem**

**Graphics Adapter**

**Network Adapter**

**Sound Blaster**

**They were called "DSPs"**
**Digital Signal Processors**

**They were killed by**
**"Native Signal Processing"**

Drivers

**They put processing**
**next to the data**

**Analog front end devices**

**With NSP…**

Cost

**$ $ $**

**Performance**

Power

**w w W W**

**So why do it?**

# Now We Are Trending Back

FPGA

CPU

Memory

Storage

Fabric

Storage

Memory

CPU

AI

**Distributed resources**

**Application-specific computing**

**In-memory computing**

**Artificial intelligence and deep learning**

**Security**

Graphics Accelerator

Human interface

Search Engine

Standard CPU

HTML processing
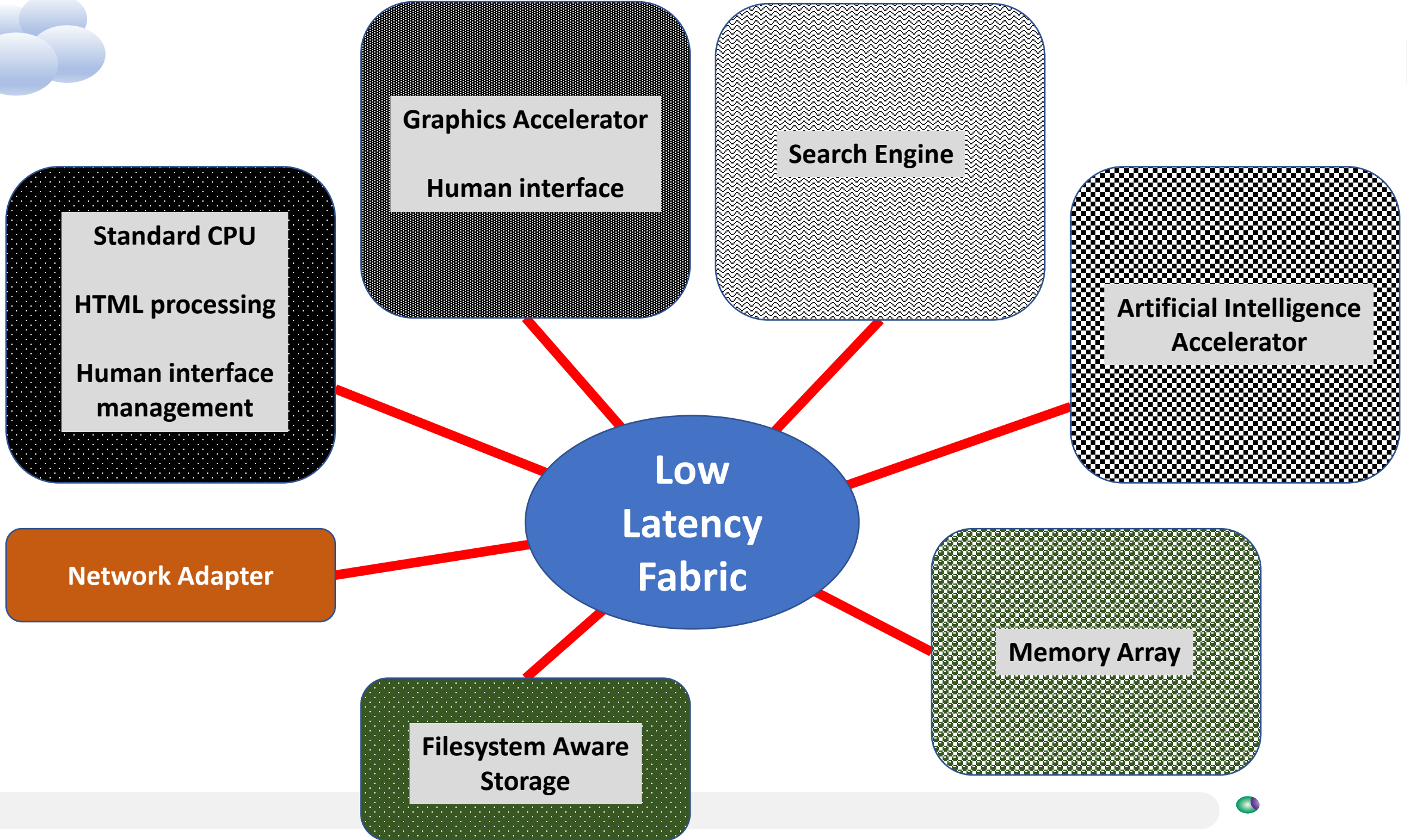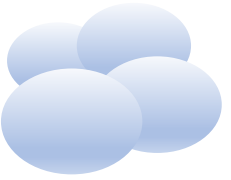
Human interface management

Artificial Intelligence Accelerator

Low Latency Fabric

Network Adapter

Memory Array

Filesystem Aware Storage

# Example AI accelerator



**Tbps links**

| I/O | NNP Control |

HBM | HBM | HBM | HBM

Exec Unit | SRAM
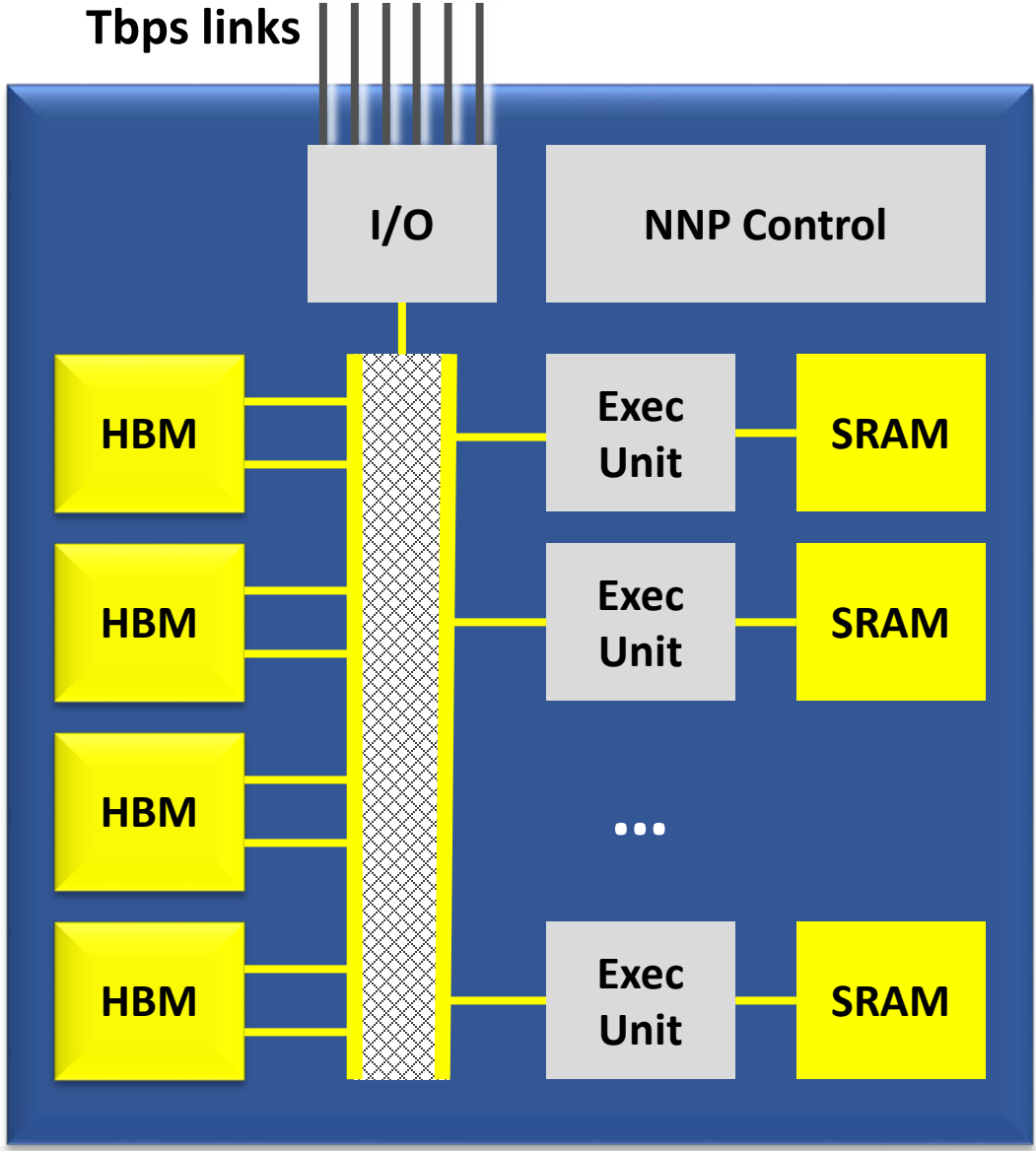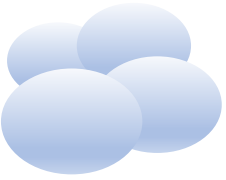Exec Unit | SRAM
...
Exec Unit | SRAM

**SIMD architectures**
**Matrix interconnections**
**Fast pipes still limit load/save time**

**Challenges:**
- **Model checkpointing**
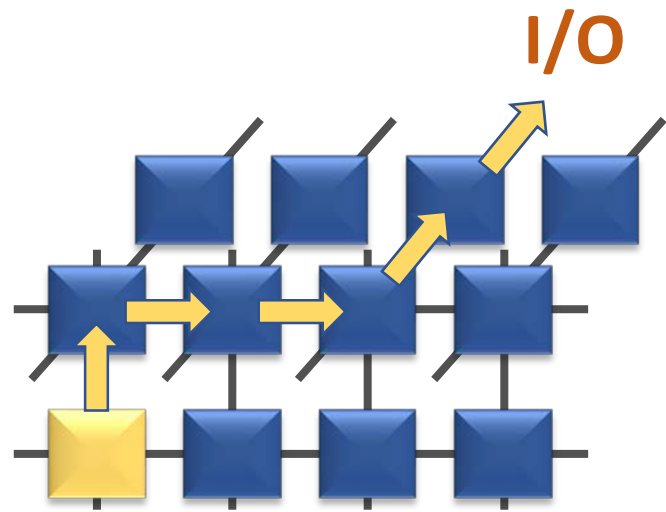- **Data loss on power fail**
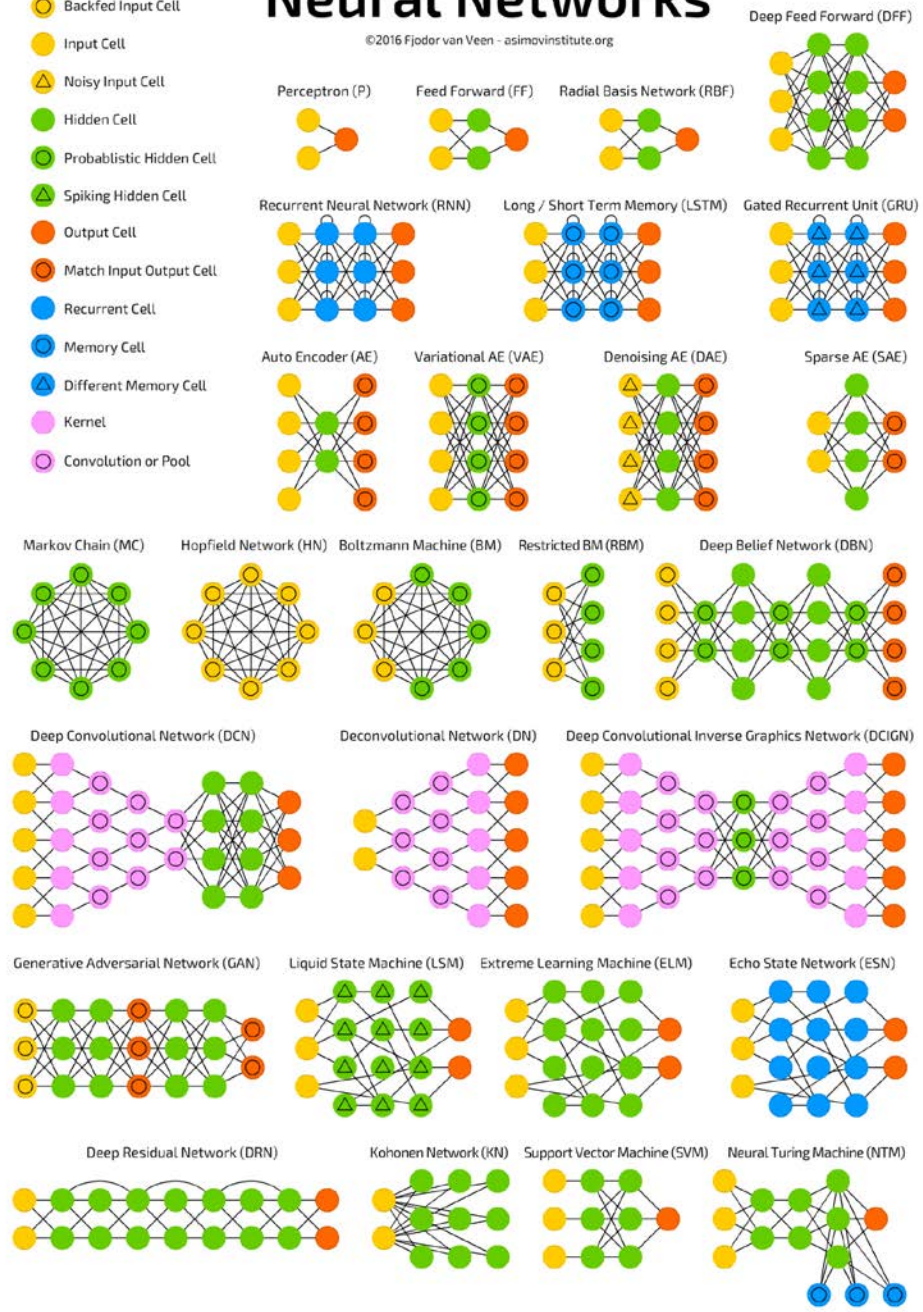- **Temperature sensitivity**

I/O

**Back propagation algorithms complicate things**

**Data loss problems are amplified**

**Checkpointing highly time and bandwidth consuming**

**The more distributed memory gets,
the harder to load and unload**

# NVRAM TO THE RESCUE!

**Replacing dynamic memory with persistent memory resolves the data loss issues**

**Just leave the data in place as long as you want**

**Replace DRAM with NVRAM**

**Replace eRAM with NVRAM**

**The final frontier…**

**SRAM & Registers**

# Continuing to look for ways to bring Memory Class Storage down under 1ns



**Voltage adjustment**
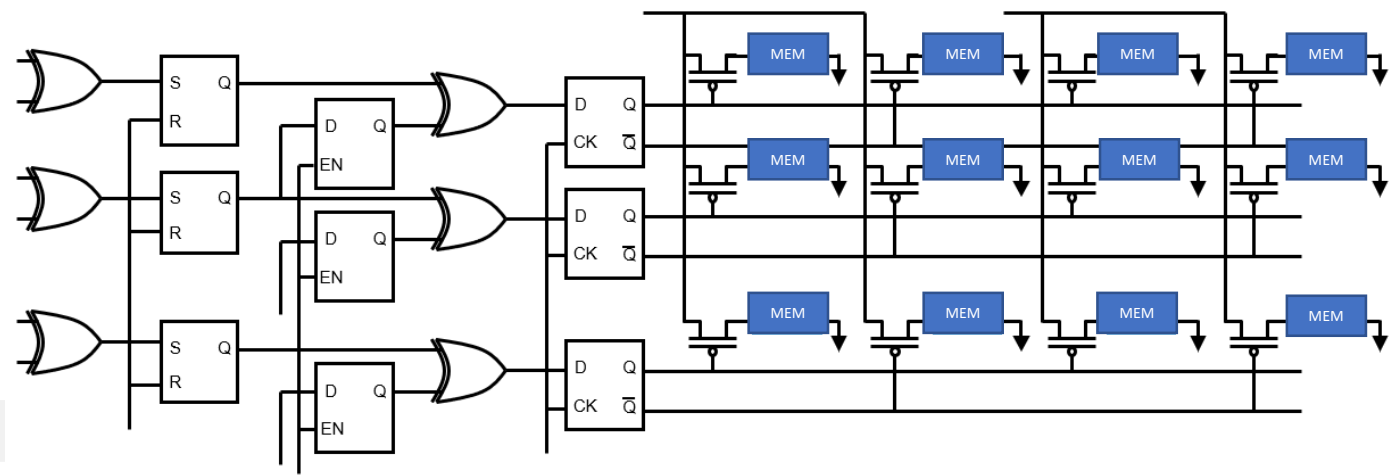
**Faster edge rates**

**Better error check**

**Shadow registers**

**Getting smarter**

# It will happen

**When we no longer fear power failure…**

# DATA PERSISTENCE

**Full END TO END persistence**

Are we getting near the
day when we look back
at volatile memory…

…and LAUGH?

...but...

# Persistent data introduces challenges, too

**Data is ALWAYS there!**

**Data security is a growing concern**

**Application opens data from previous application**

**Memory moved from one system to another**

**Spy devices on memory buses**

# So many potential breaches

**Infection via hack**

**Infection via spy devices**

**General trend is to encrypt data before transmission or storage**

X2.Hd44**3#jj0%

X2.Hd44**3#jj0%

**Keep the bad guys out**

**Small Energy Source**

Non-Volatile Memory Array – Any Kind

X2.Hd44**3#jj0%

Voltage Regulator

Smart NVM Control

DRAM Cache

Password:

CPU

**Host System**

**Some are adding in-memory compute functions including encryption**

**Works as long as the bus is secure**

**Encryption quality may be limited by block transfer size**

**Management of many keys can get complicated quickly**

## ISO/IEC 11889



CA=Certification Authority

Module/platform

Integrator verifies (authenticates) all components, may issue additional certs to individual components or signed assertions concerning composition of module (E.G. a TCGPlatform Certificate)

*ASN.1* - ISO-822-1-4;
- ITU-T X.680
- ITU-T X.681
- ITU-T X.682
- ITU-T X.683

*DER* - ISO-8825-1; ITU-T X.690

*X509v3* - ISO-9594-8; ITU-T X.509

**Common Criteria:**
- Common Criteria for Information Technology Security Evaluation, Parts 1-3, Version 3.1, Revision 4, September 2010
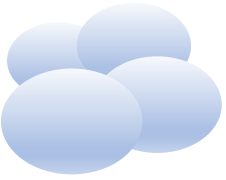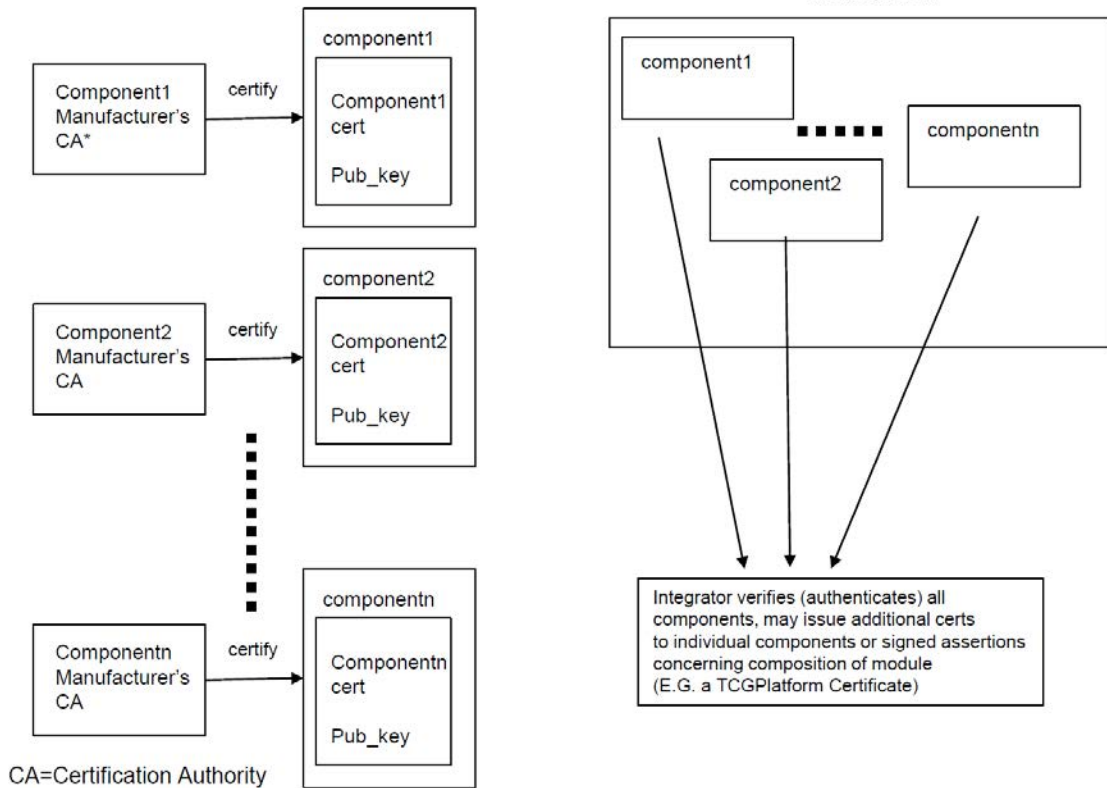- ISO/IEC 15408 Evaluation criteria for IT security Parts 1-3

*ECDSA*:
- ANSI X9.62; NIST-FIPS-186-4, Section 6
- ISO/IEC 14888-3 Digital signatures with appendix -- Part 3: Discrete logarithm based mechanisms (Clause 6.6)

*NIST P256, secp256r1*:
- Certicom-SEC-2, NIST-Recommended-EC
- ISO/IEC 15946 Cryptographic techniques based on elliptic curves (NIST P-256 is included as example)

*SHA256*:
- NIST-FIPS-180-4
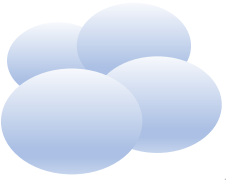- ISO/IEC 10118-3 Hash-functions -- Part 3: Dedicated hash-functions (Clause 10)

*OID* - ITU-T X.402

*SP800-90A*:
- NIST-SP-800-90A

*SP800-90B:*
- NIST-SP-800-90B

Summary

Power Fail Sucks

Saving Data is a Pain

Need tiers of memory & storage

Persistence is Essential

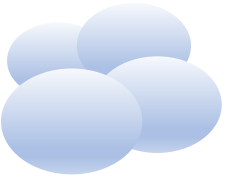Sharing Time

Today's Solutions Help

Persistence Complications

But We Can Do Better

Data Distribution Challenges

Mix & Match Memories

DDR5 NVRAM Spec in Progress
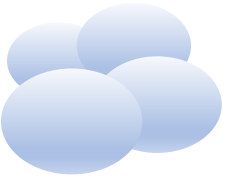
**Thank you for your time**

**Bill Gervasi**
**bilge@Nantero.com**

**I'm here to learn too**

**What do you deal with?**