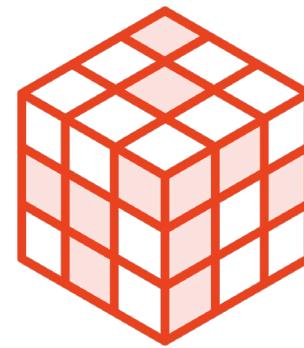
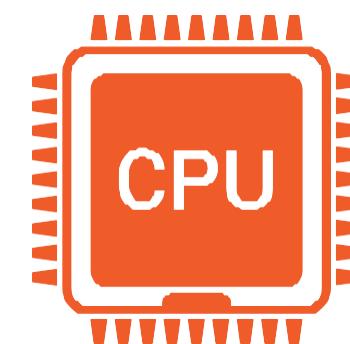


Their OS and software vendors only support CISC processors.



The software vendor claims their application is optimized for parallel processing.



They will use a CISC processor

More CPU cores should yield better performance

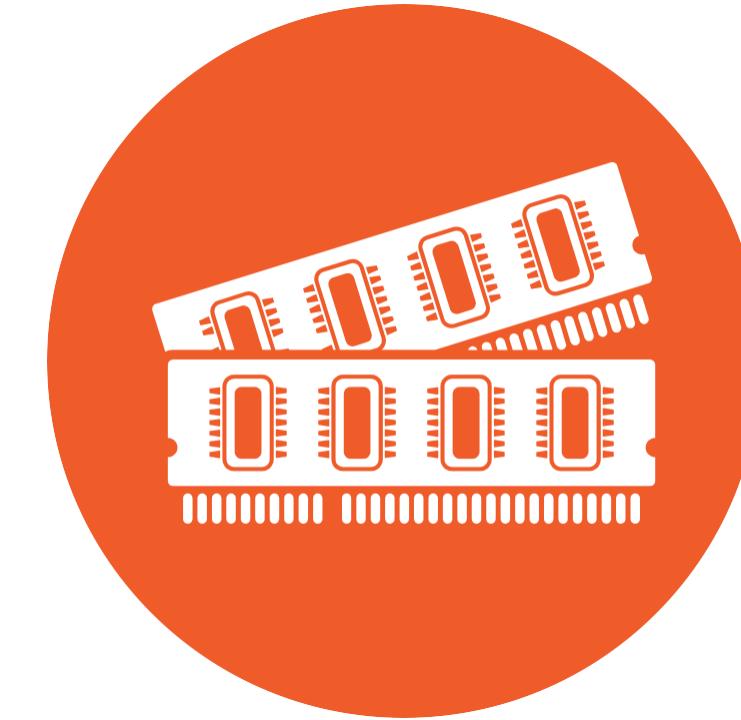
The server will have 2 CPU sockets with 12 CPU cores per socket

RAM



Storyline

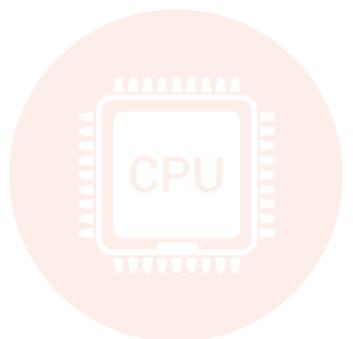
What Do You Need to Know?



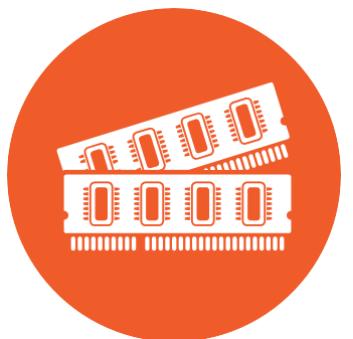
Globomantics is researching the kind of memory their server will need

Let's learn more about the options and then see if the recommended hardware is suitable for Globomantics

Minimum Hardware



1 CPU Sockets with 6 cores



64 GB of RAM



1 Gbps Networking



512 GB Storage

Recommended Hardware

2 CPU Sockets with 12 cores each

256 GB of RAM

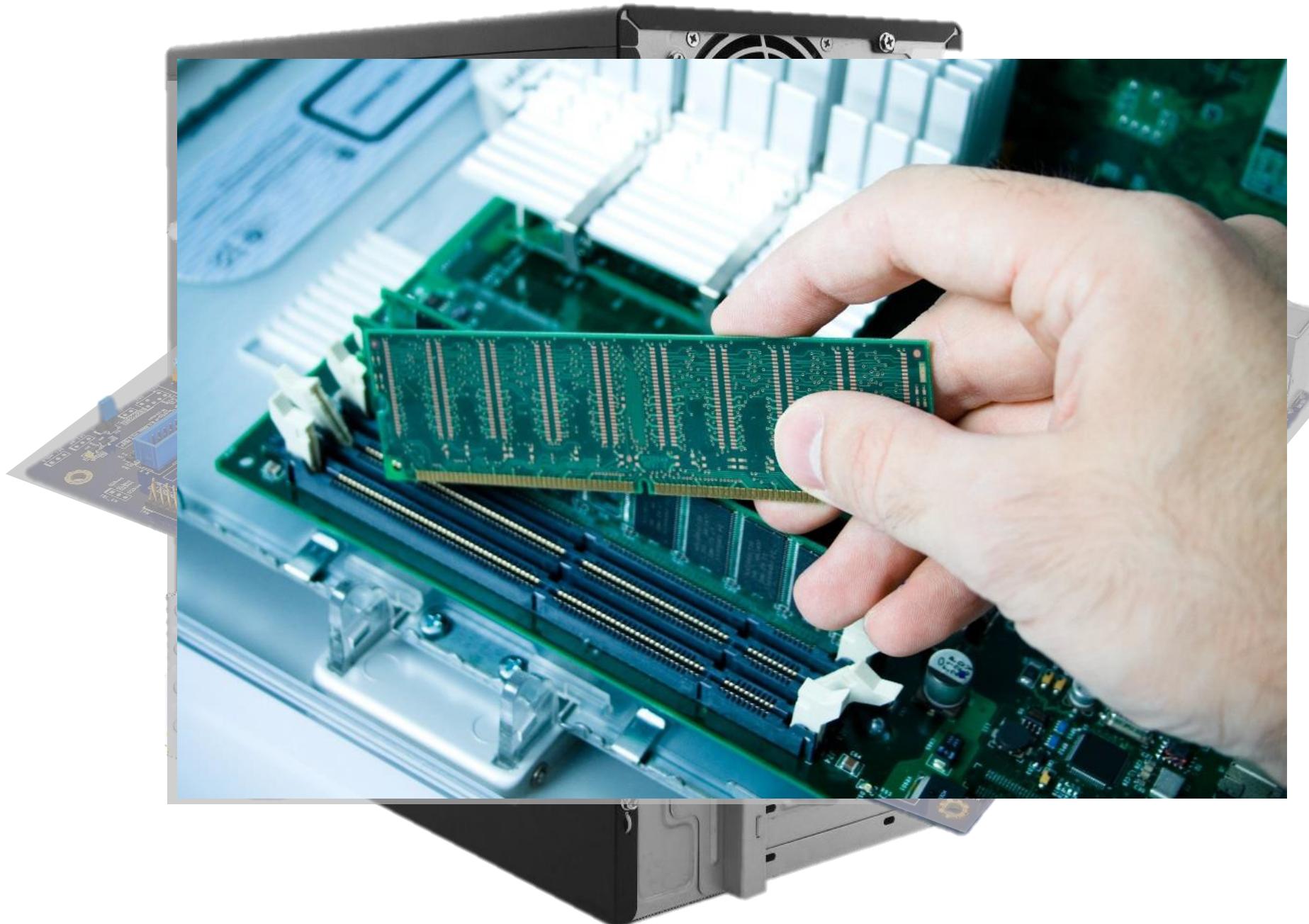
10 Gbps Networking

3 TB Storage

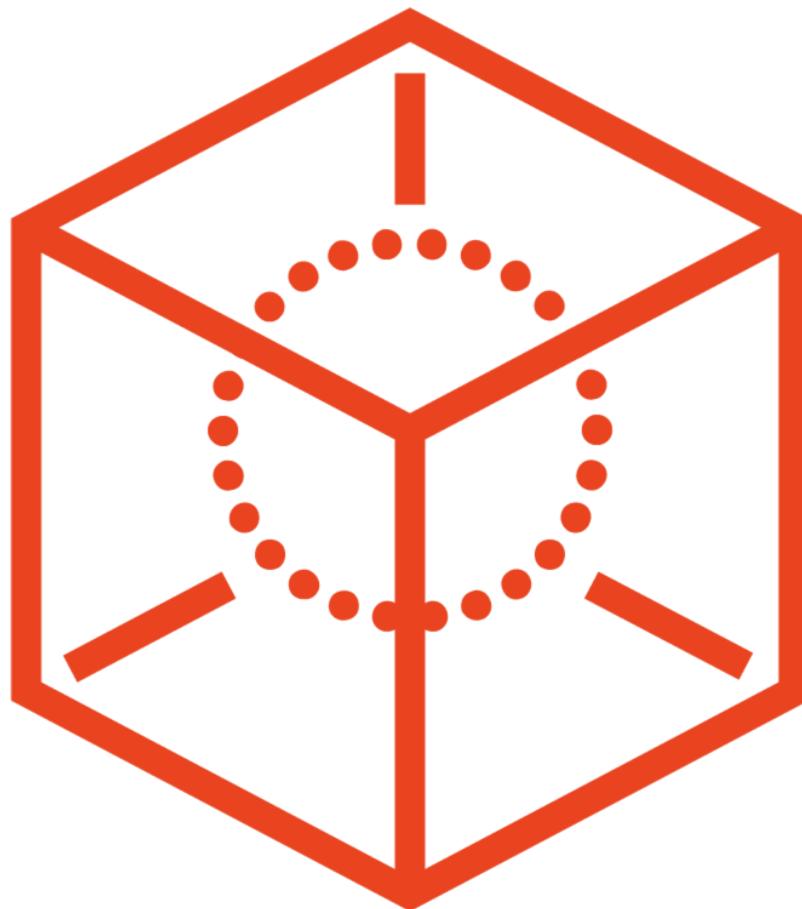


Tech Point
RAM Capacity

Visual Anchor

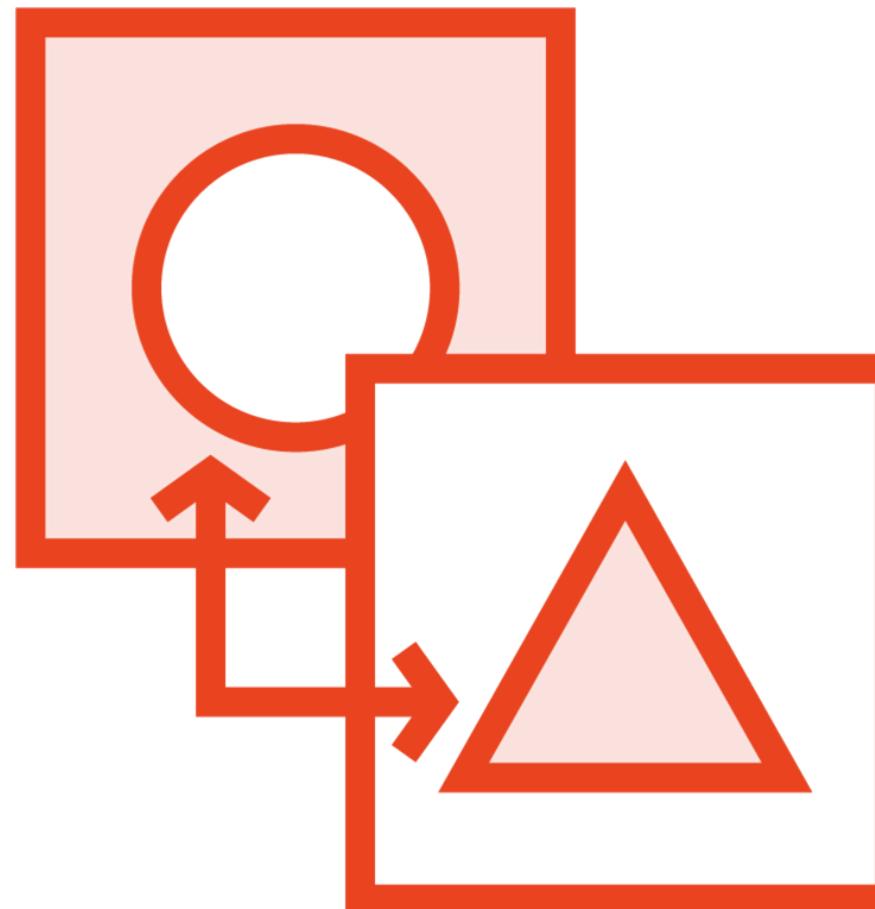


With RAM, What Matters Most?

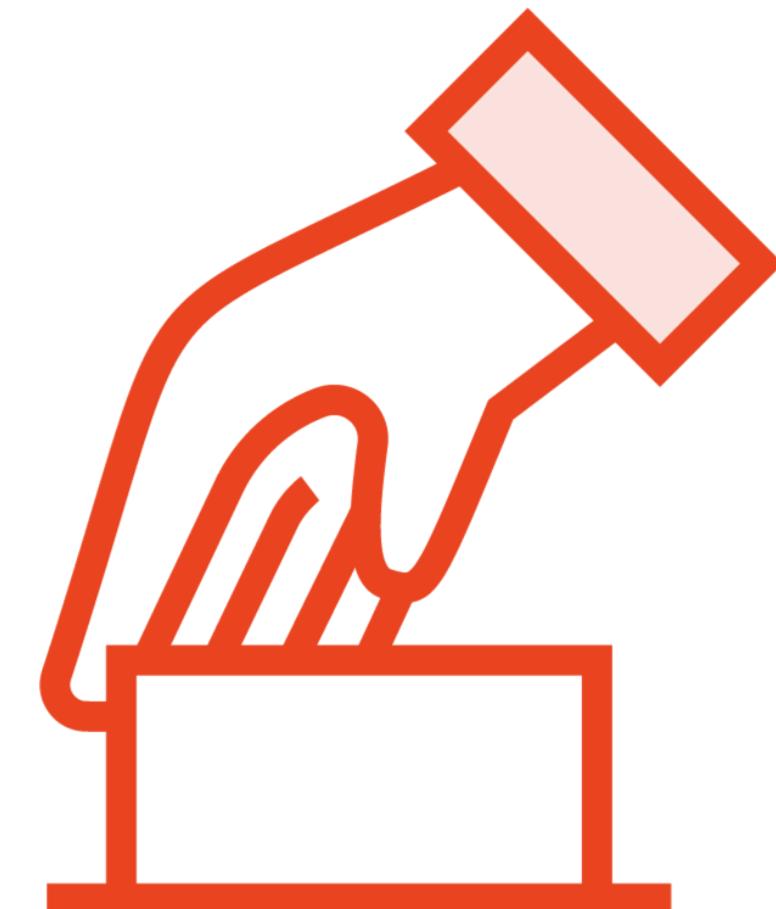


The total RAM
capacity

64GB, 256GB, 1TB



The type of RAM
ECC, DDR4, etc.



Where you install it
Channel 0, Channel 1

RAM is measured in GB and TB,
but it isn't the same thing as
storage space.

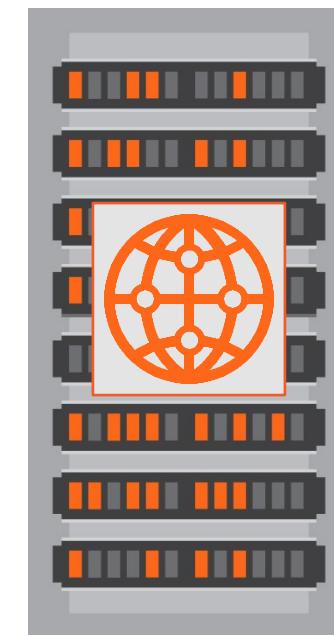
What Is RAM Used For?

How Do I Know How Much RAM I Need?

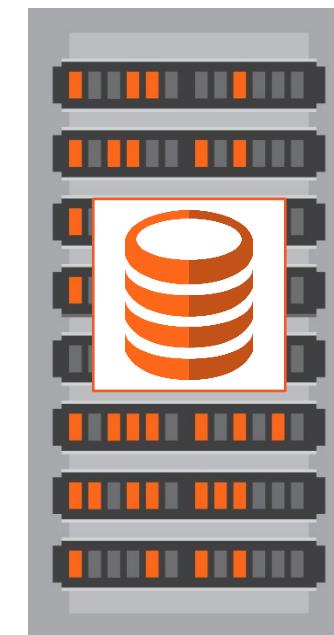
Let's look at some example server roles and see how they use resources



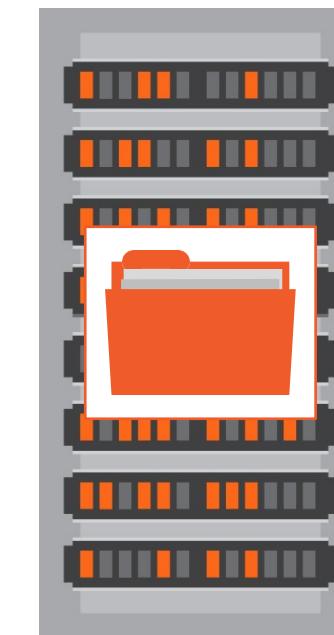
Email Servers



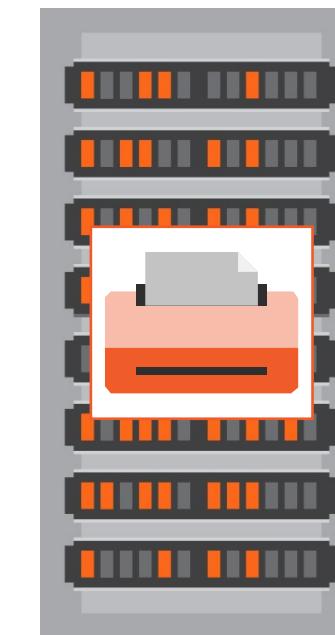
Web Servers



Database
Servers



File Servers

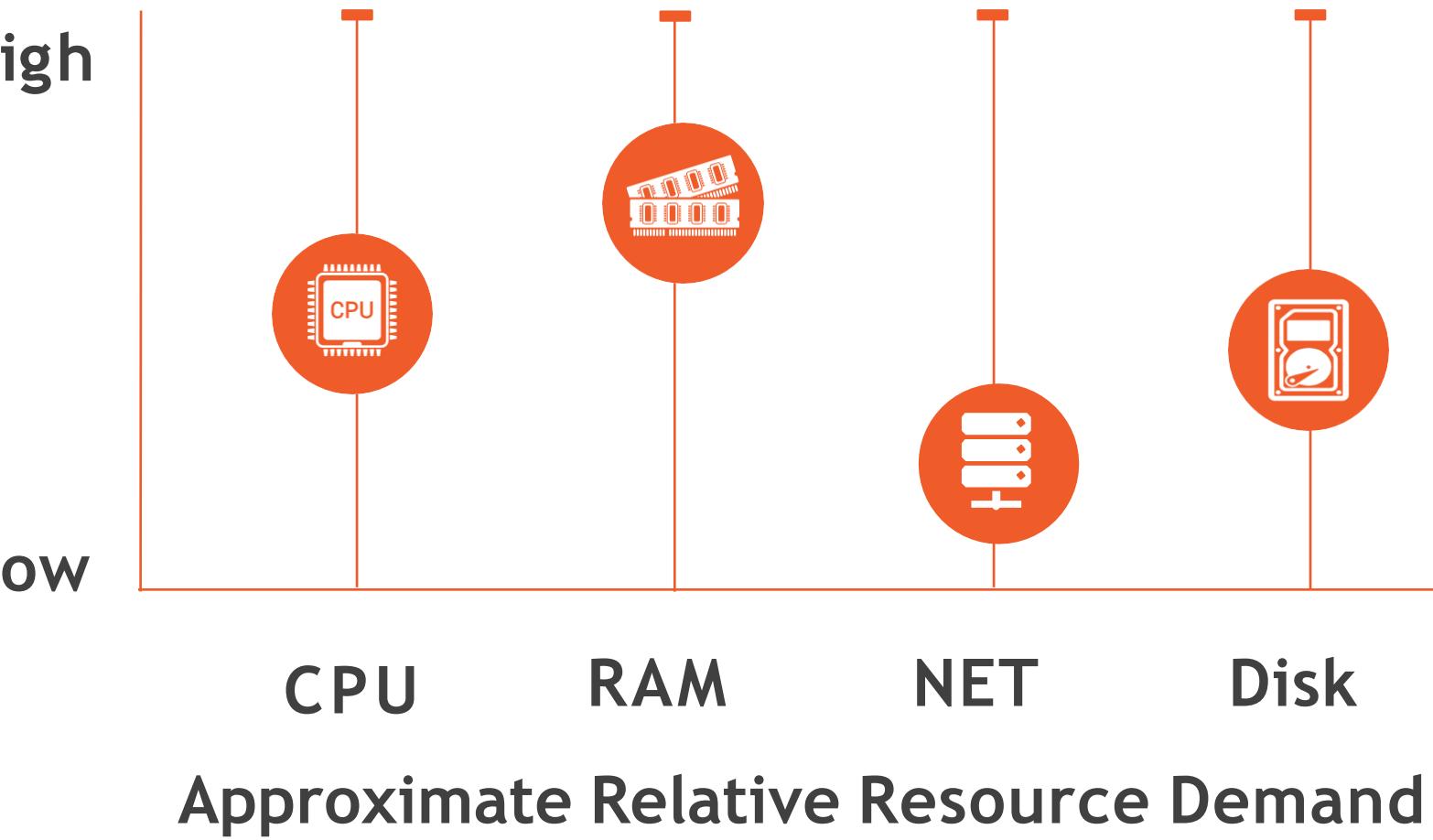
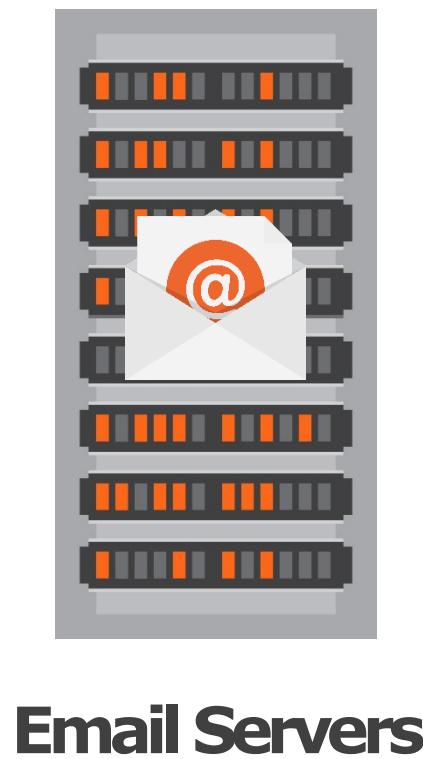


Print Servers

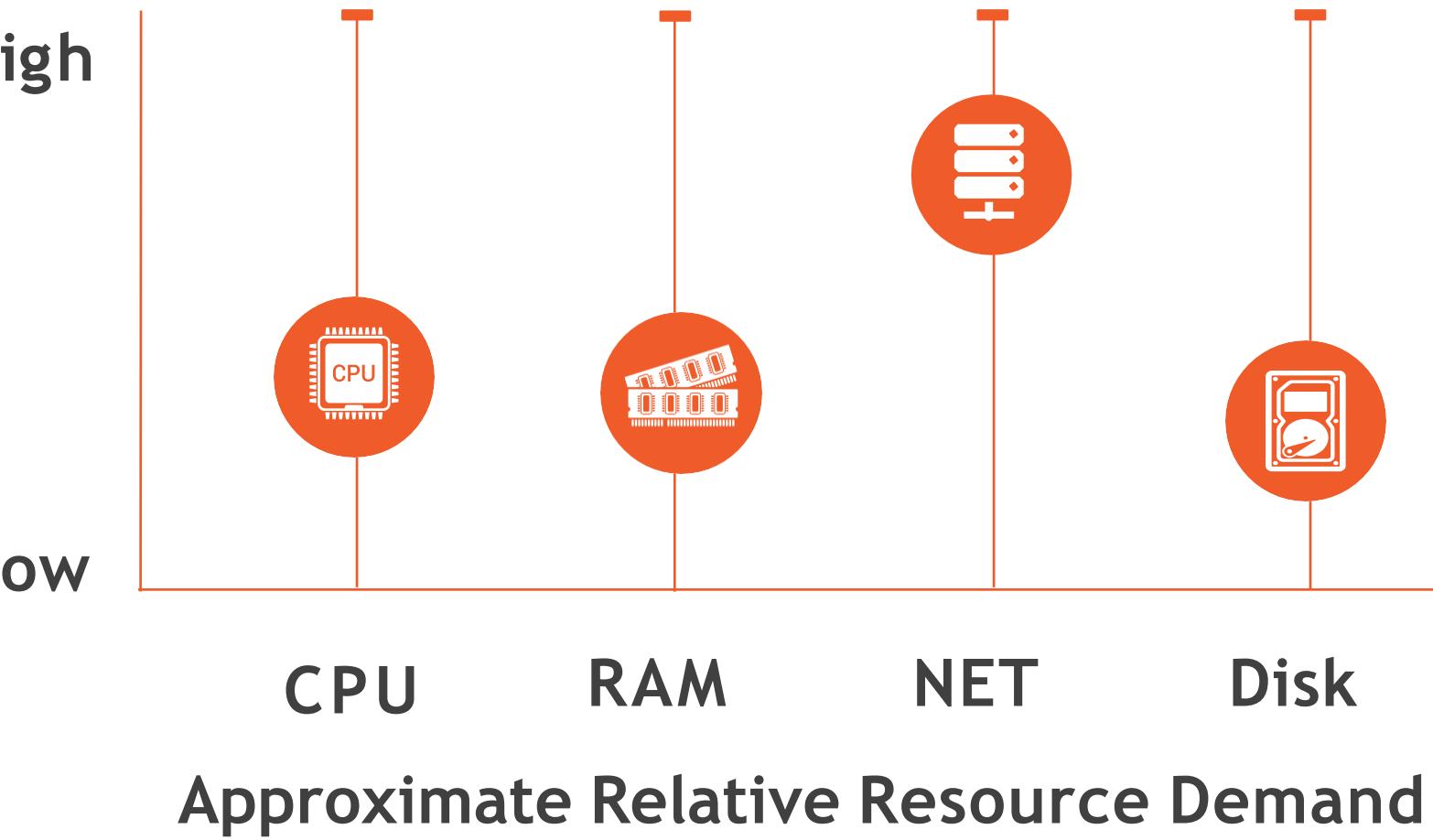
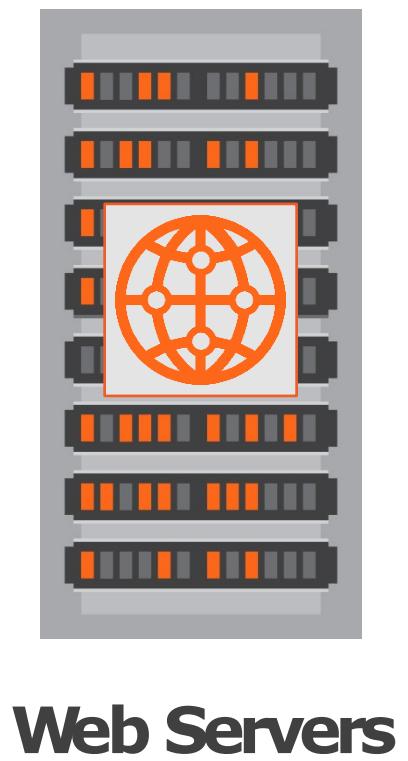


Authenticatio
n Servers

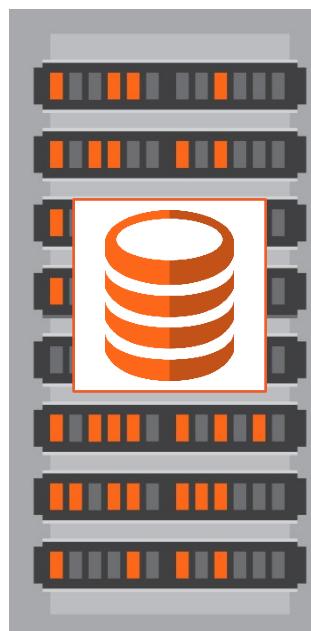
Server Role vs. Resource Requirements



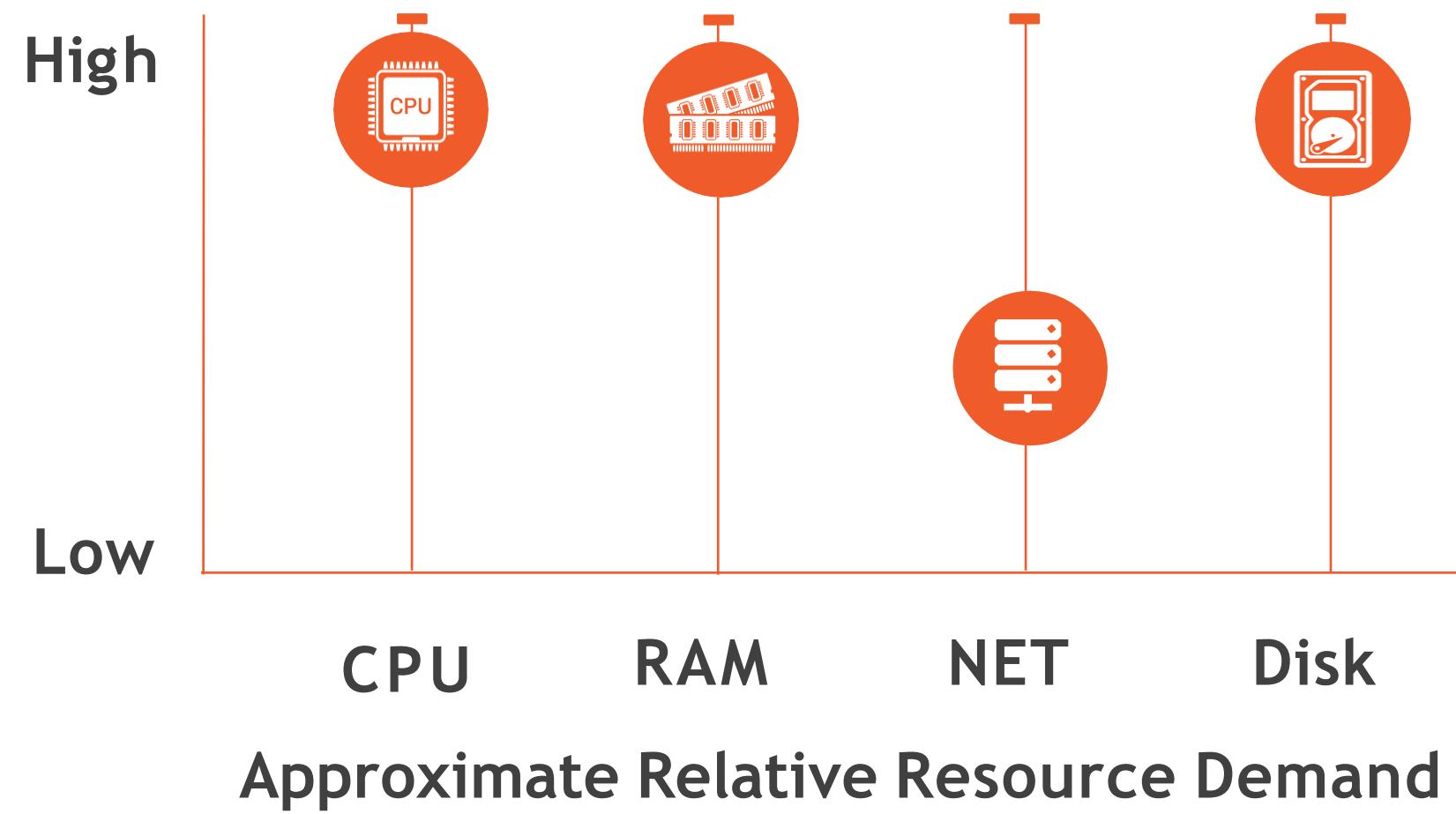
Server Role vs. Resource Requirements



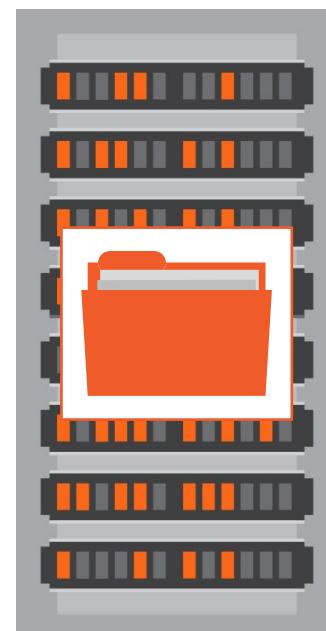
Server Role vs. Resource Requirements



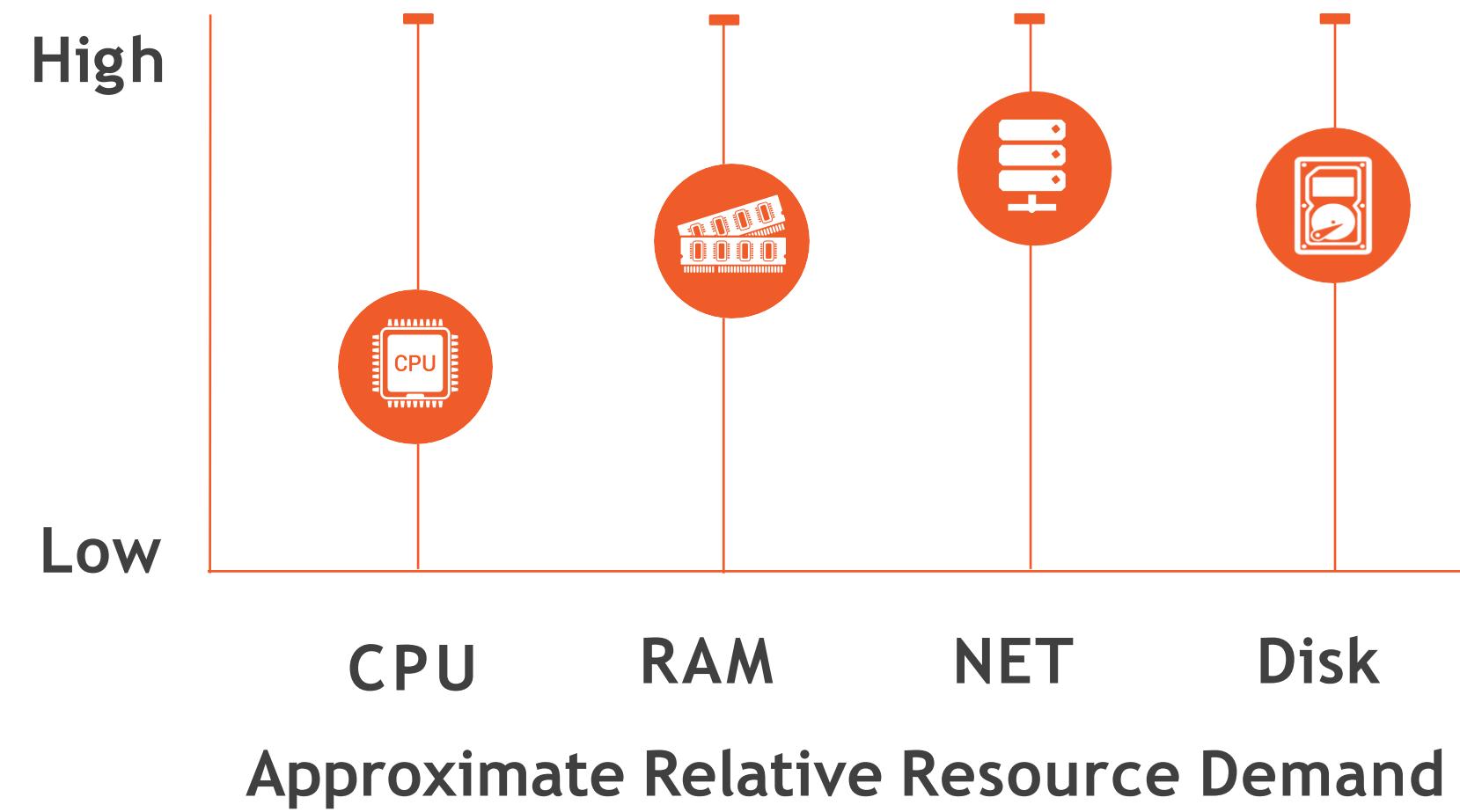
**Database
Servers**



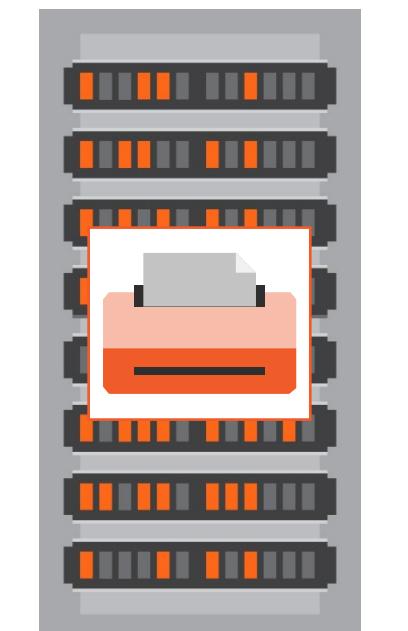
Server Role vs. Resource Requirements



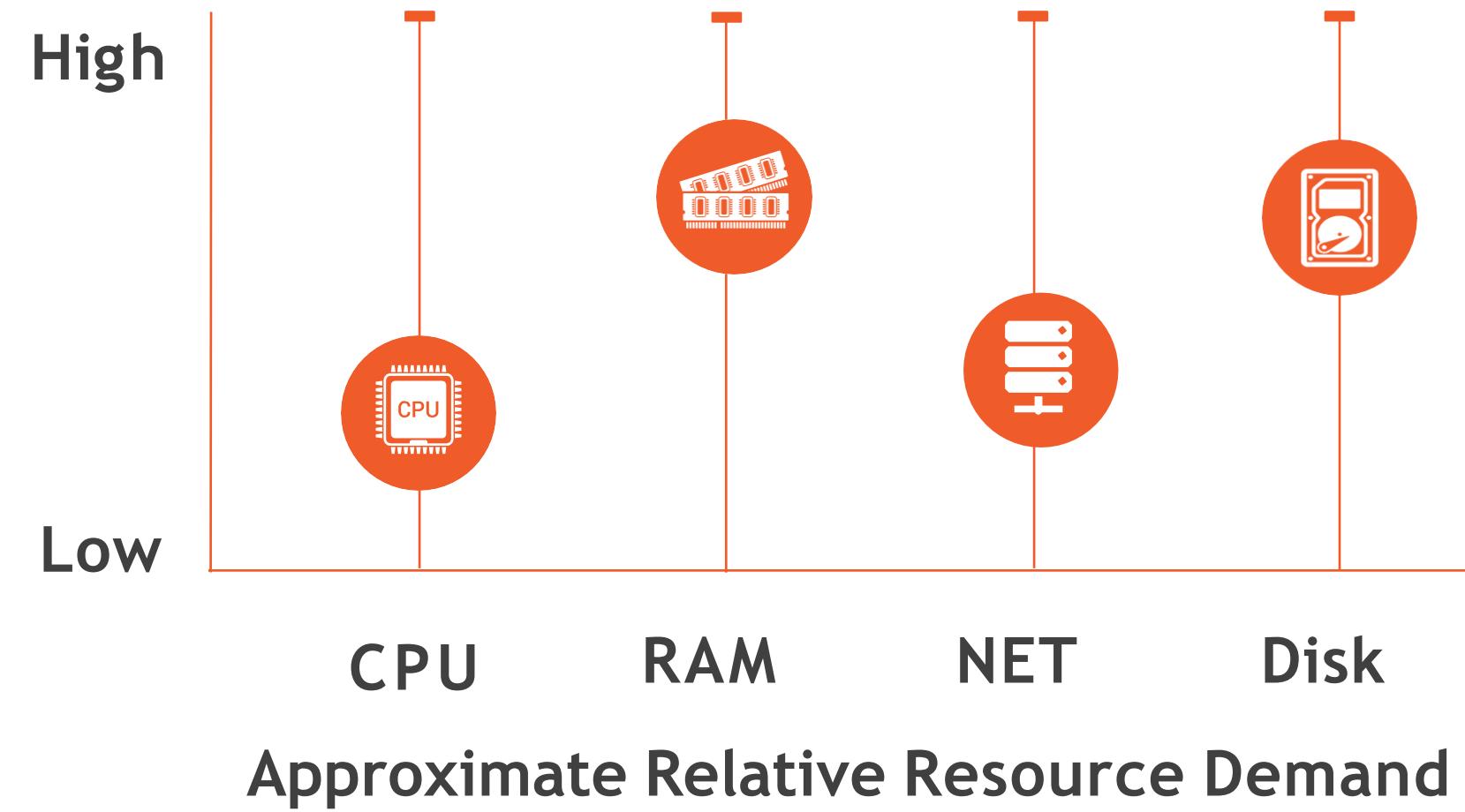
File Servers



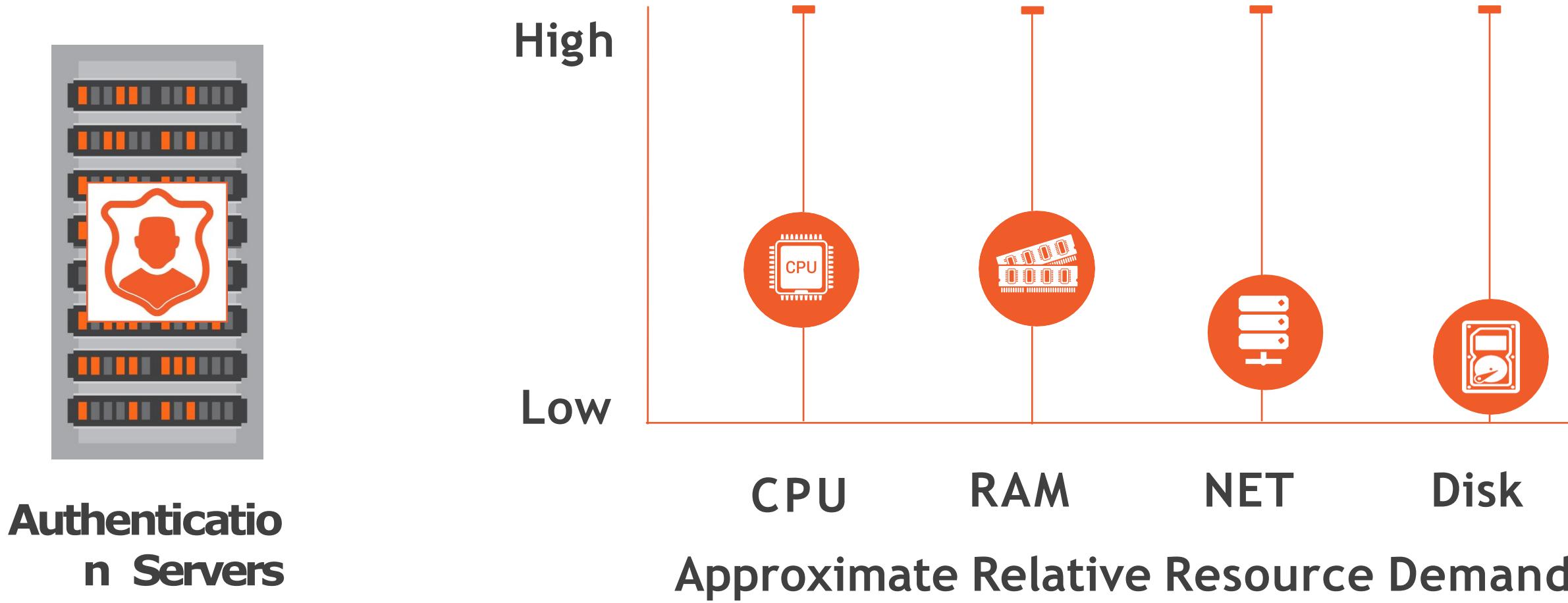
Server Role vs. Resource Requirements

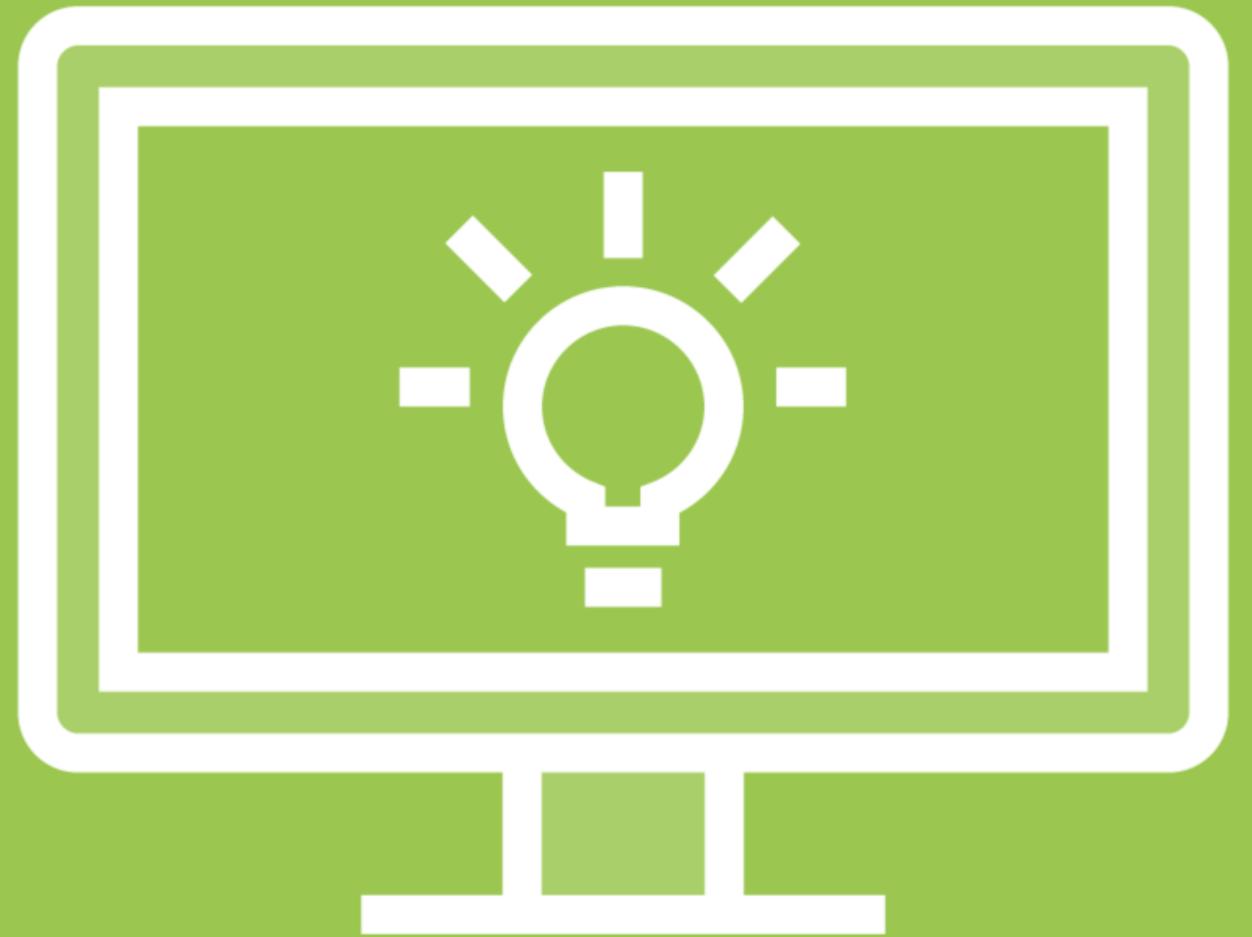


Print Servers



Server Role vs. Resource Requirements





Tech Point
RAM Types

Always get the fastest RAM you
can afford, given other
constraints

Memory Types

There are three principal types of Random Access Memory (RAM)

Today you will only find SDRAM modules in a server

SRAM

Static RAM

Much faster I/O rates than DRAM at a higher cost per bit stored

Requires 6 transistors per bit of data stored

Commonly used as cache within a CPU, L1 – L3 cache

DRAM

Dynamic RAM

Requires only one capacitor and transistor per bit of data stored

Timing is measured in absolute terms using nanoseconds

No longer in wide use

SDRAM

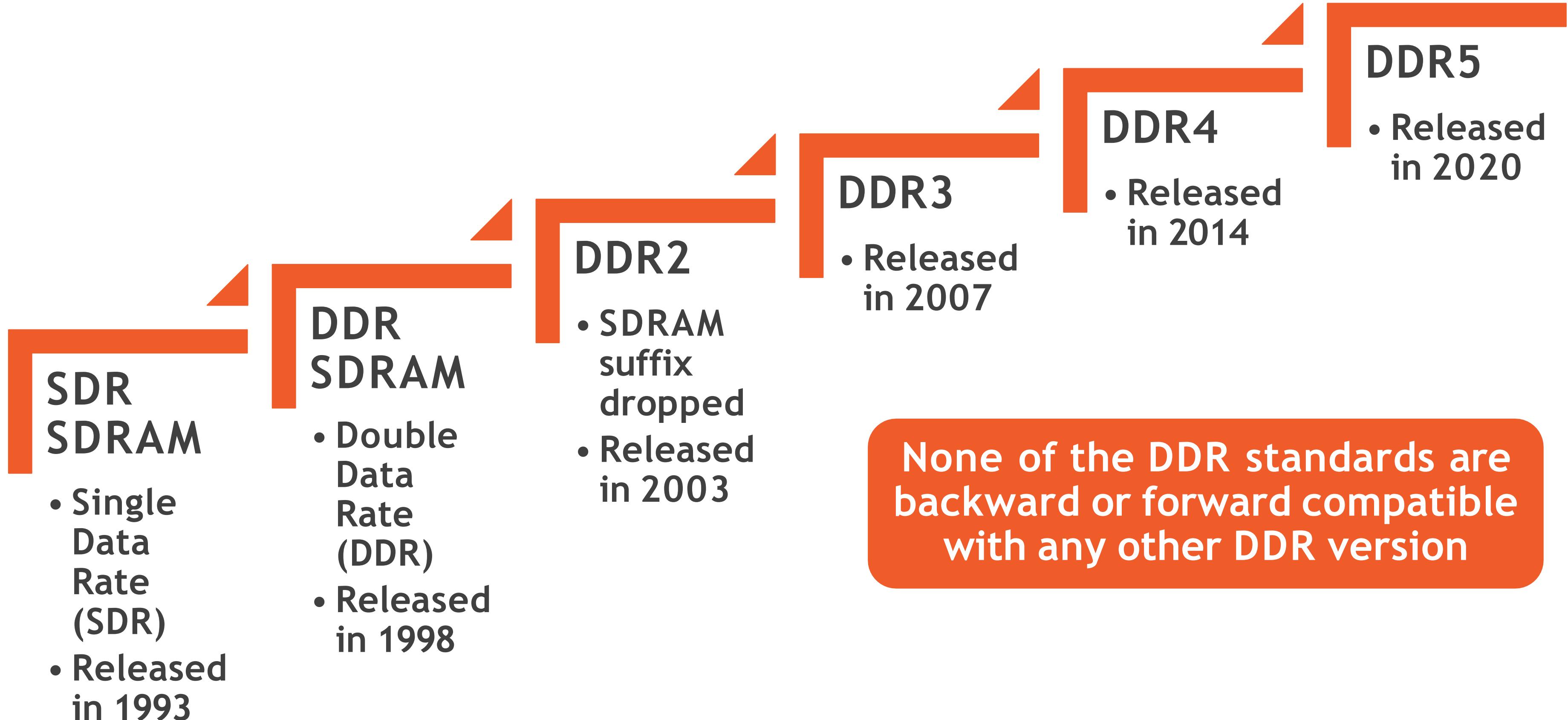
Synchronous DRAM

Synchronous interface based on clock speed

Memory operations are pipelined (queued) to increase performance

The most common memory type in PCs and servers

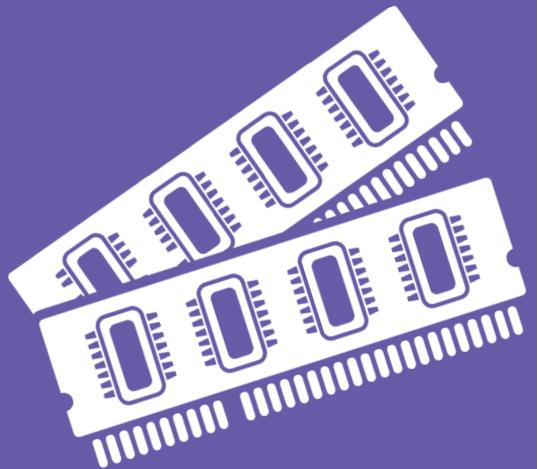
SDRAM Memory Generations



Number of Pins

Number of pins per DIMM vary depending on type

- DDR: 184
- DDR 2: 240
- DDR 3: 240
- DDR 4: 288
- DDR 5: 288

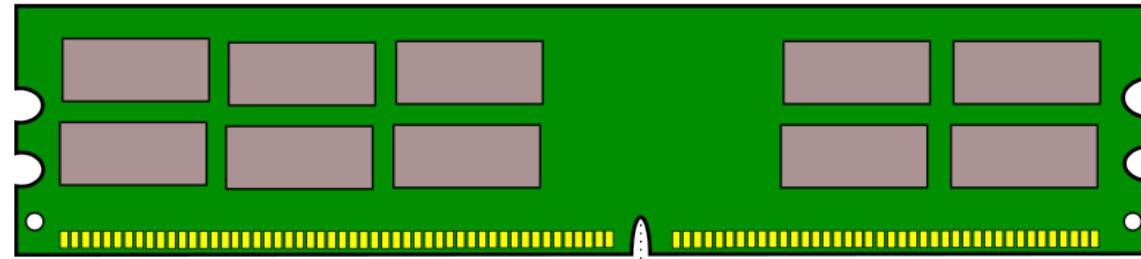


Laptops use a variant called SO-DIMMs that have a different number of pins

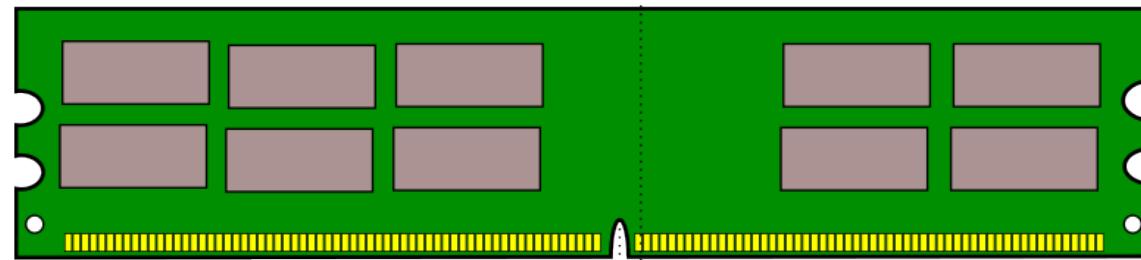


DDR SDRAM Modules

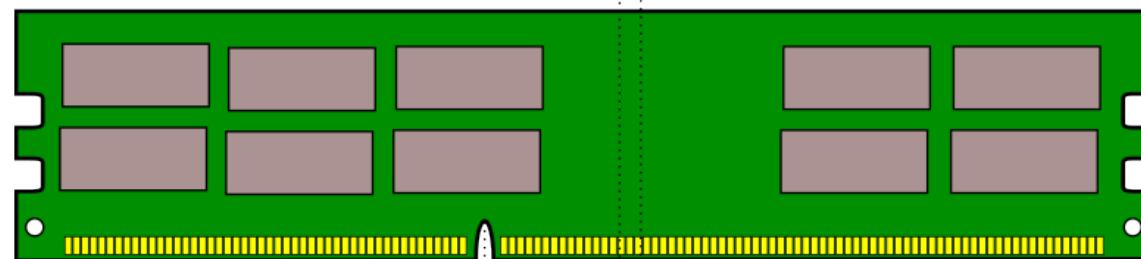
DDR



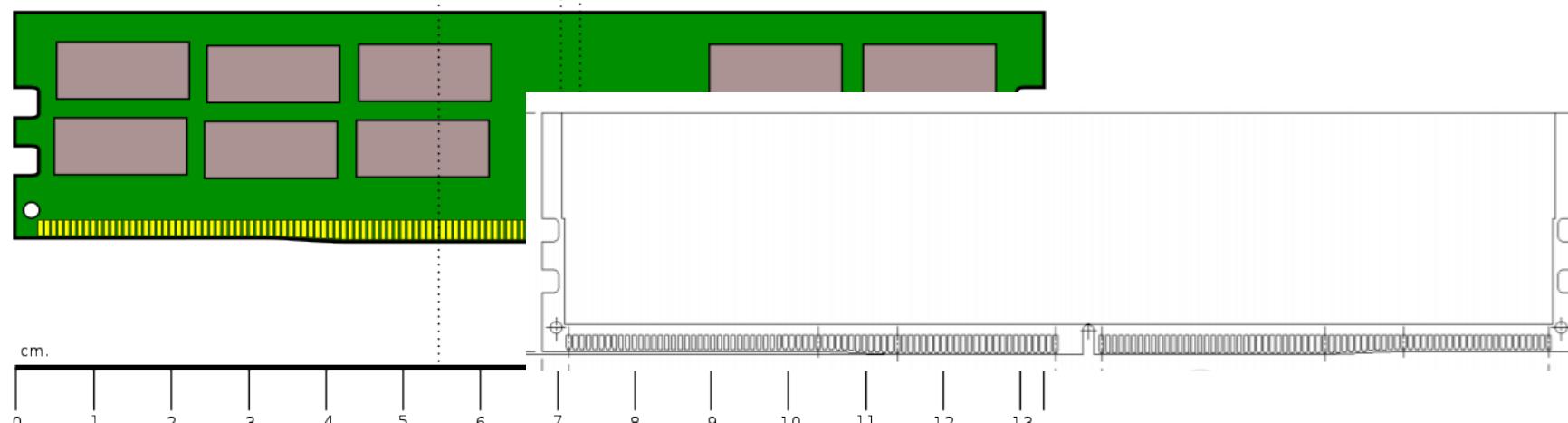
DDR 2



DDR 3



DDR 4



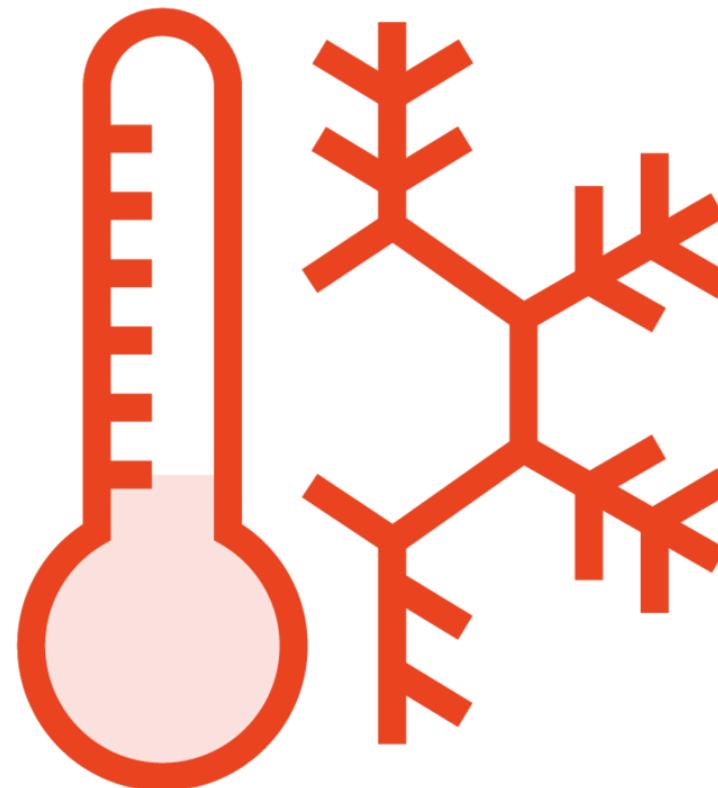
DDR modules are physically incompatible

A notch is placed at different offsets for each DDR family, a process called “keying”

This prevents accidental mixing of modules

DDR Version Comparison

Each generation is about **20-30%** more power efficient than the previous version



Less heat produced in servers



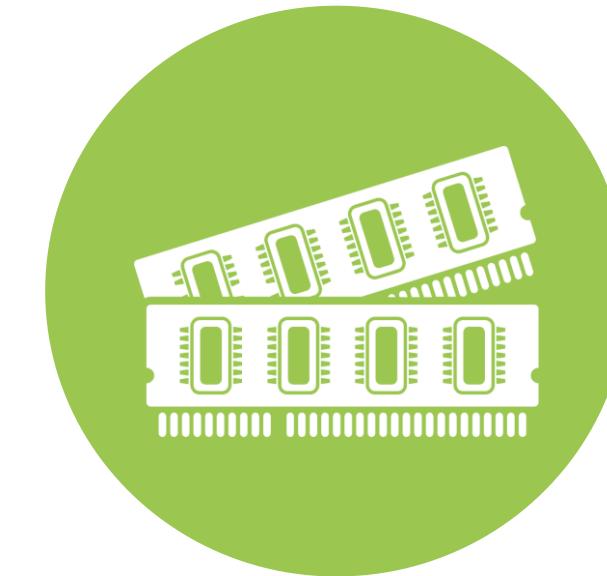
Better battery life for mobile devices

DDR Version Comparison

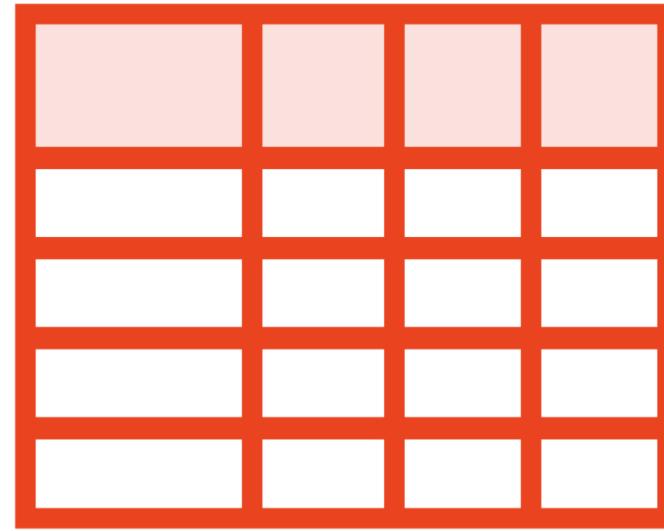
2X

Typically about 2X the maximum data transfer rate of the previous generation

Greater memory density per chip
- More memory in less space

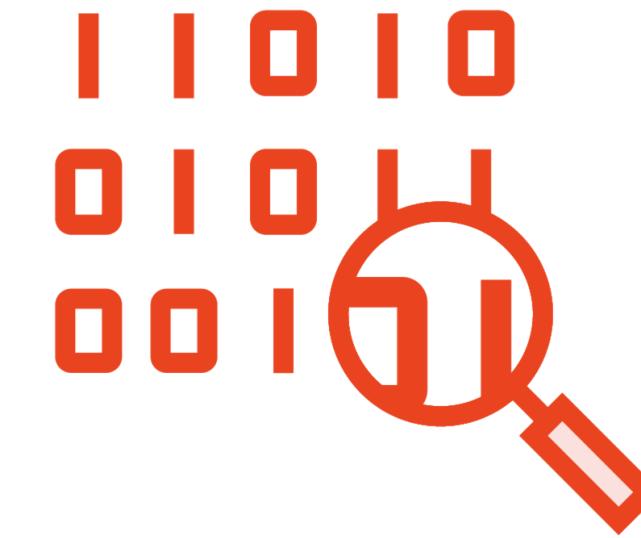


RAM Characteristics



Registered/Buffered

Impacts the maximum amount of memory in the system



Error Correcting

Improves the reliability of the memory by correcting single-bit errors

Registered/Unregistered Memory

is also known as

Buffered/Unbuffered Memory

Registered vs. Unregistered Memory

Registered memory allows the memory controller on the motherboard to address more memory than unregistered memory

Most motherboards only allow one type (Registered or Unregistered), but some allow both

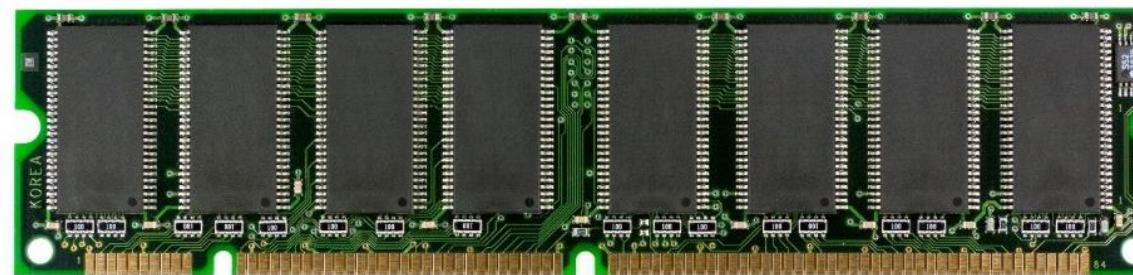
Example: A server may support 128 GB of unregistered RAM, but 512 GB of registered

Registered memory is more expensive and usually only used in servers

ECC vs. Non-ECC



ECC Memory DIMM
(Note the 9 chips)



Non-ECC Memory DIMM
(Note the 8 chips)

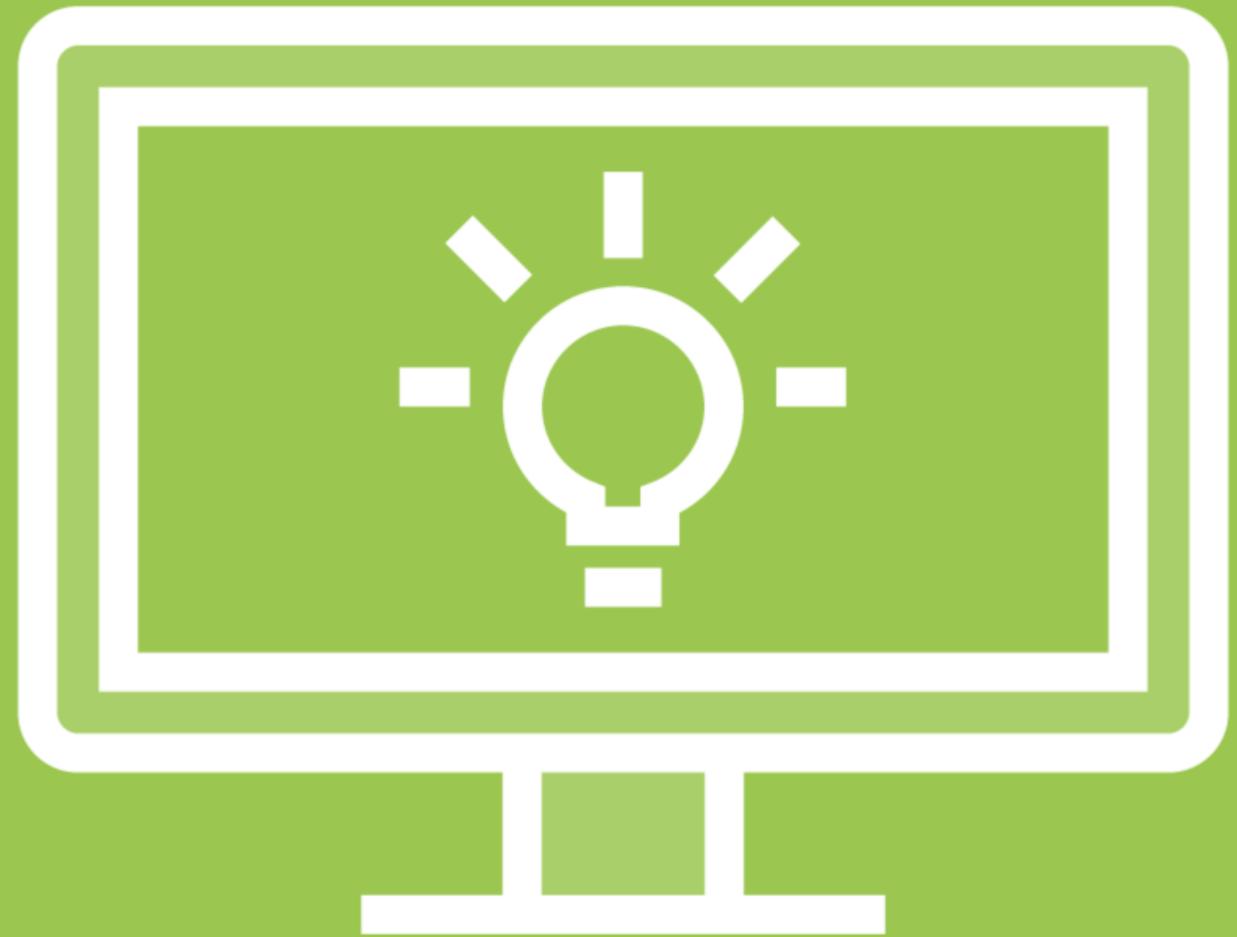
ECC = Error Correcting Code

Physically ECC usually has 9 chips per DIMM
vs. Non-ECC has 8

Can correct single bit errors (per 64 bit
“word”) and detect (but not correct) double
bit errors

Errors are usually reported in system logs

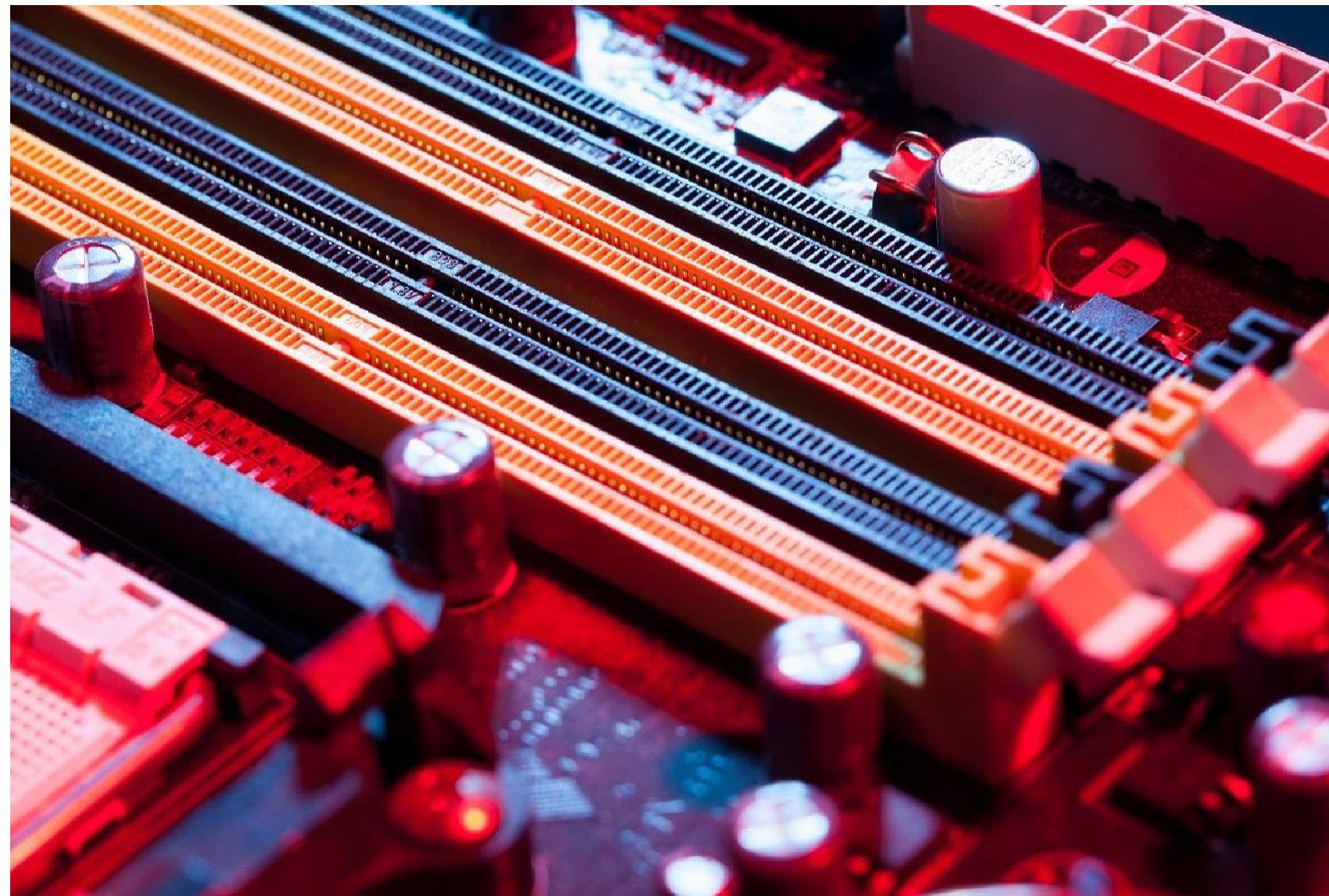
ECC memory is more expensive and therefore
usually only used in servers



Tech Point
Memory
Placement

Location, Location, Location...

Module Placement



Memory modules are typically installed in pairs into banks (slots)

The banks represent different memory channels, each associated with a particular CPU socket

You must be very careful into which channels you place your RAM

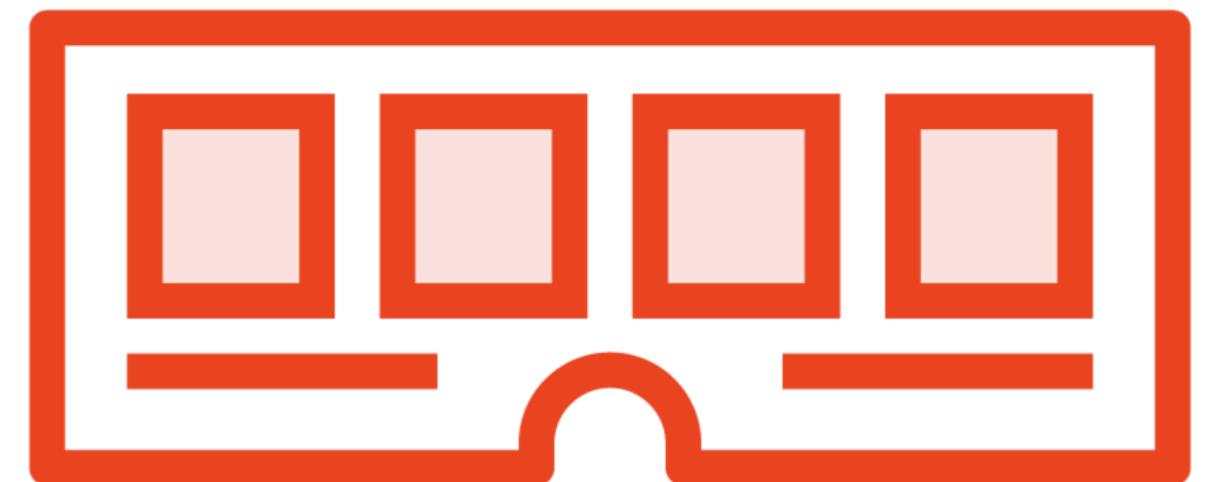
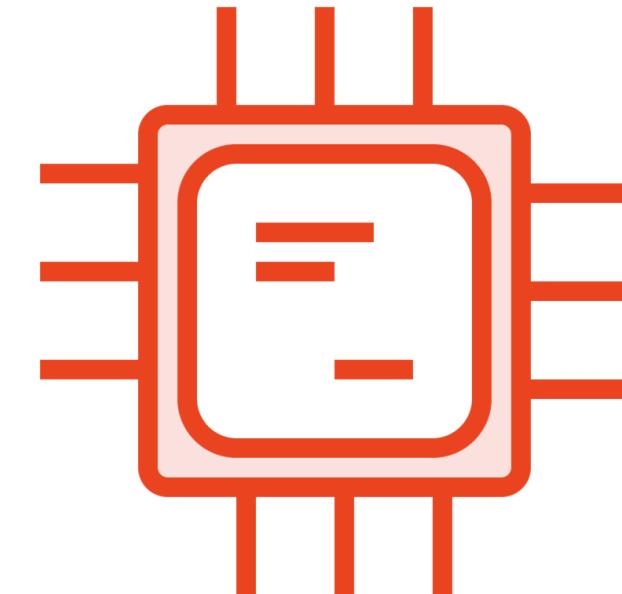
NUMA Considerations

NUMA = Non-Uniform Memory Access

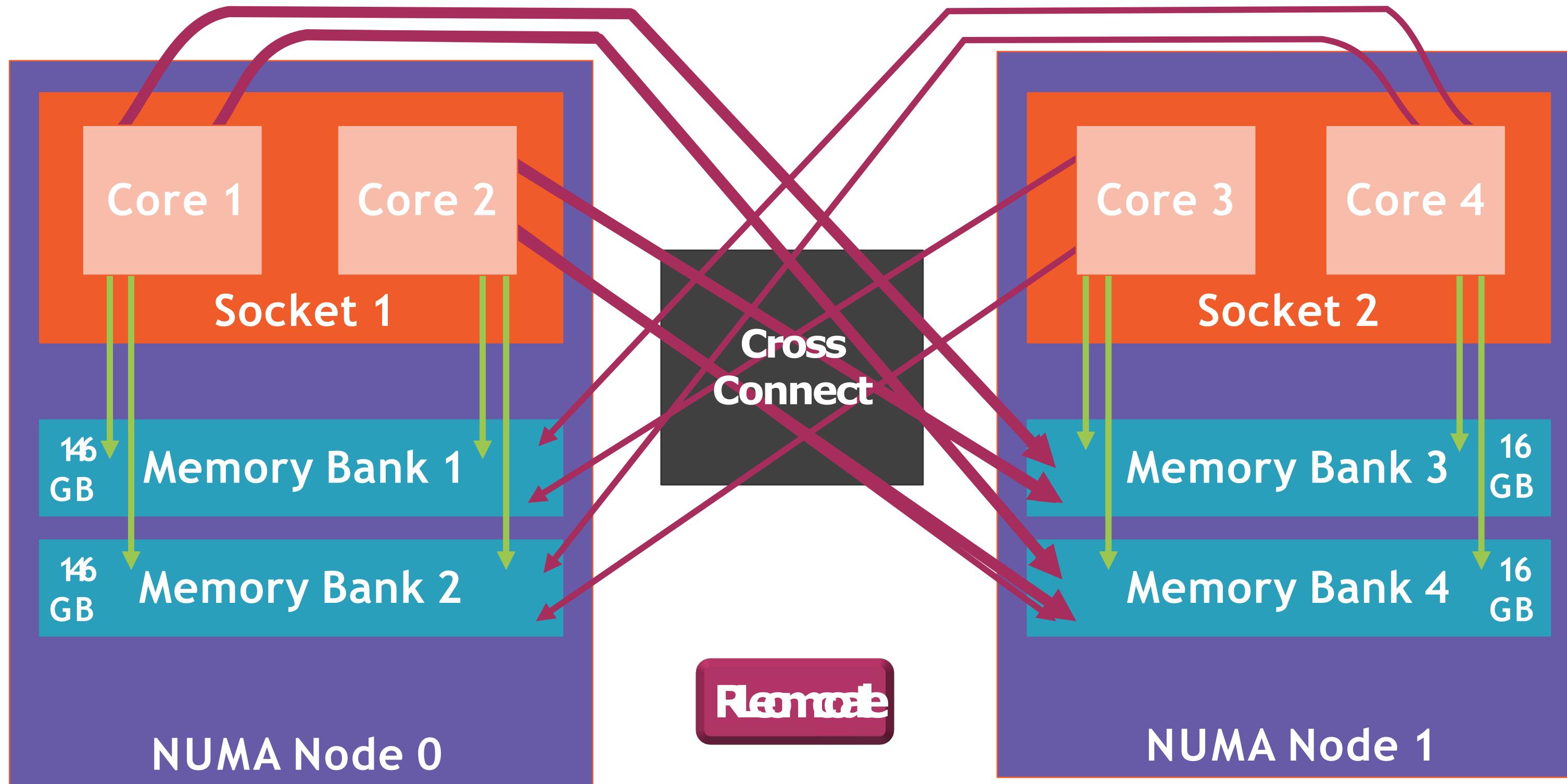
Modern operating systems allocate memory to take advantage of NUMA

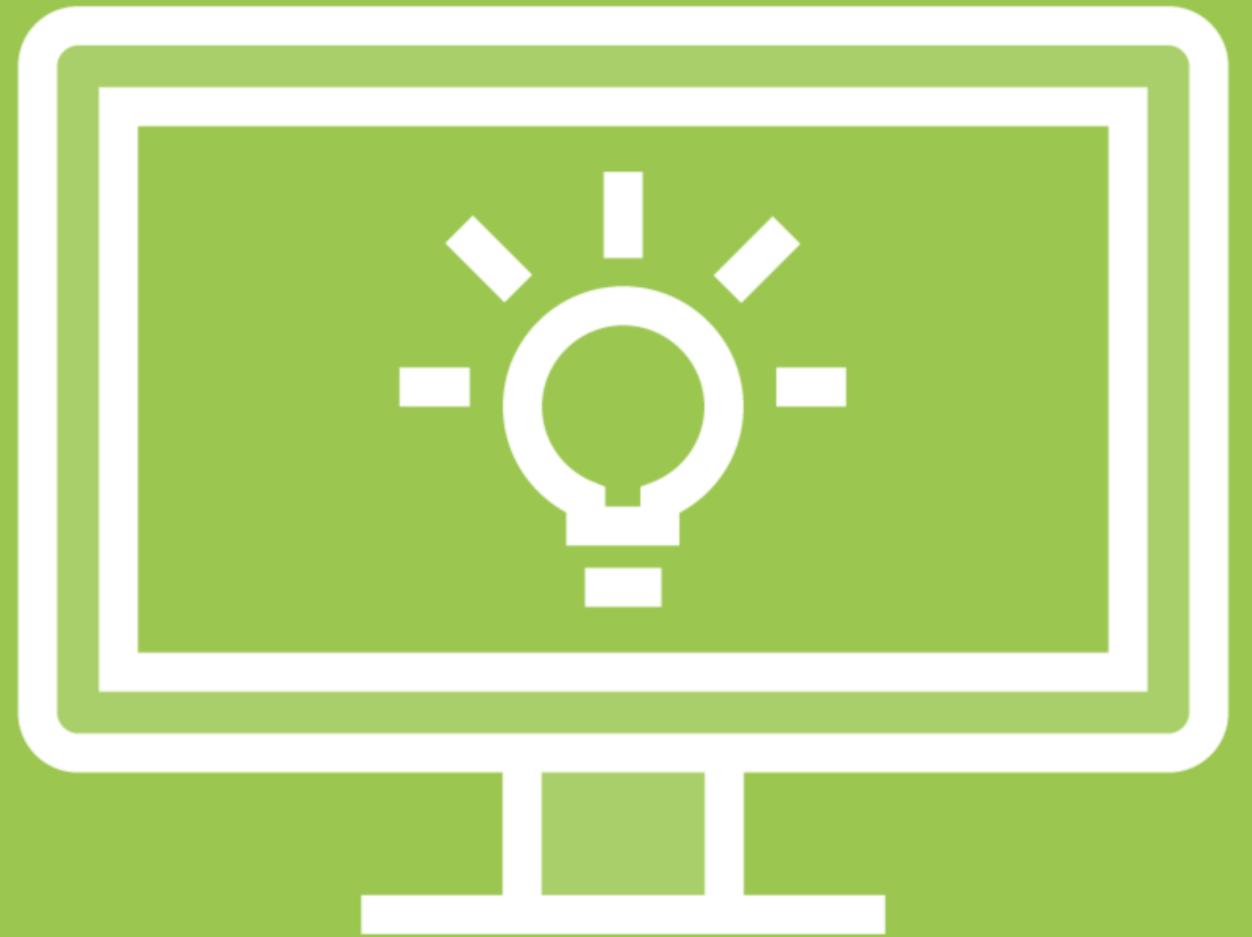
NUMA systems have one or more NUMA nodes

Each NUMA node should have the same amount of memory



NUMA Nodes





Tech Point
RAM Speeds

Data in Memory Is Stored in a Matrix

| | A | B | C |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |

γ

Select the row

α

Then select the column

$\{\gamma, \alpha\}$

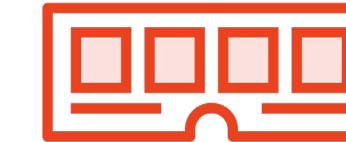
Finally, retrieve/store the contents at the intersection

RAM Speeds

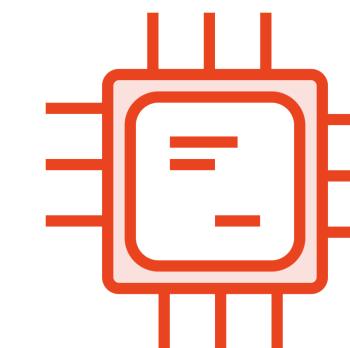
The speed of RAM is based on a number of factors



Clock speed



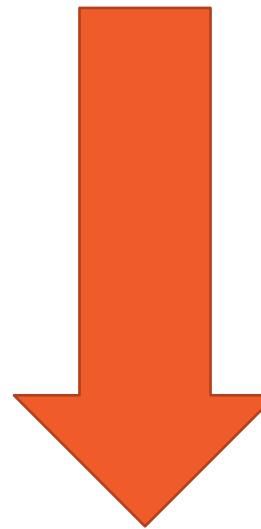
Memory chip latency



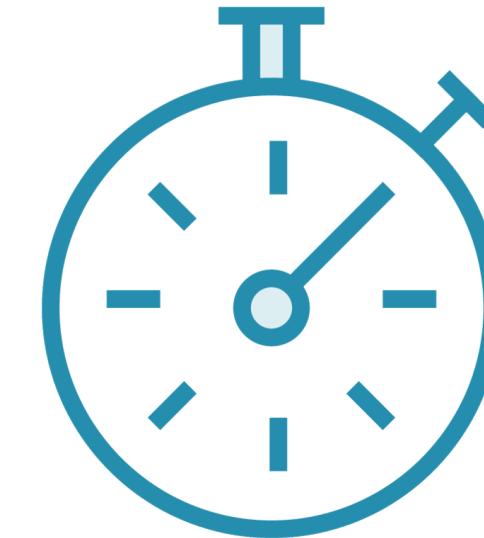
Contention from access by multiple CPUs

Timing Overview

Memory timings are typically 4 numbers in a row, separated by dashes,
(ex. 16-18-18-36 or 24-21-21-47)



Lower is generally better



In modern memory technologies (like DDR4) these numbers represent **clock cycles** rather than nanoseconds.

The latency in nanoseconds will be lower.

Timing Numbers



Example: 16-18-18-36

16

t_{CAS} Column Access Strobe (CAS): Column selection latency (CL)

18

t_{RCD} Row Column Delay: Time between opening a row and selecting the column

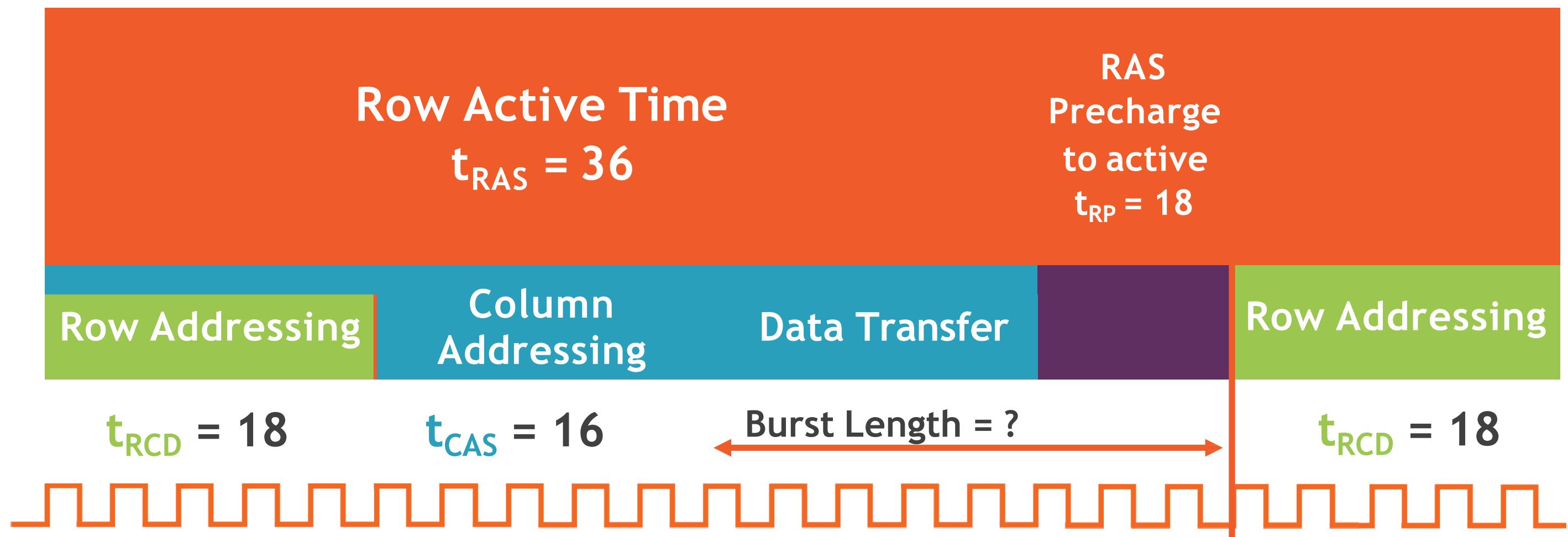
18

t_{RP} Row Precharge: Time required before moving to a new row

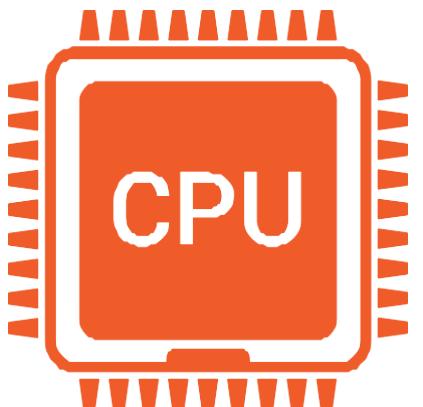
36

t_{RAS} Row Active Time (RAS): Not a typo. Minimum time a row must be active to allow enough time to access desired data

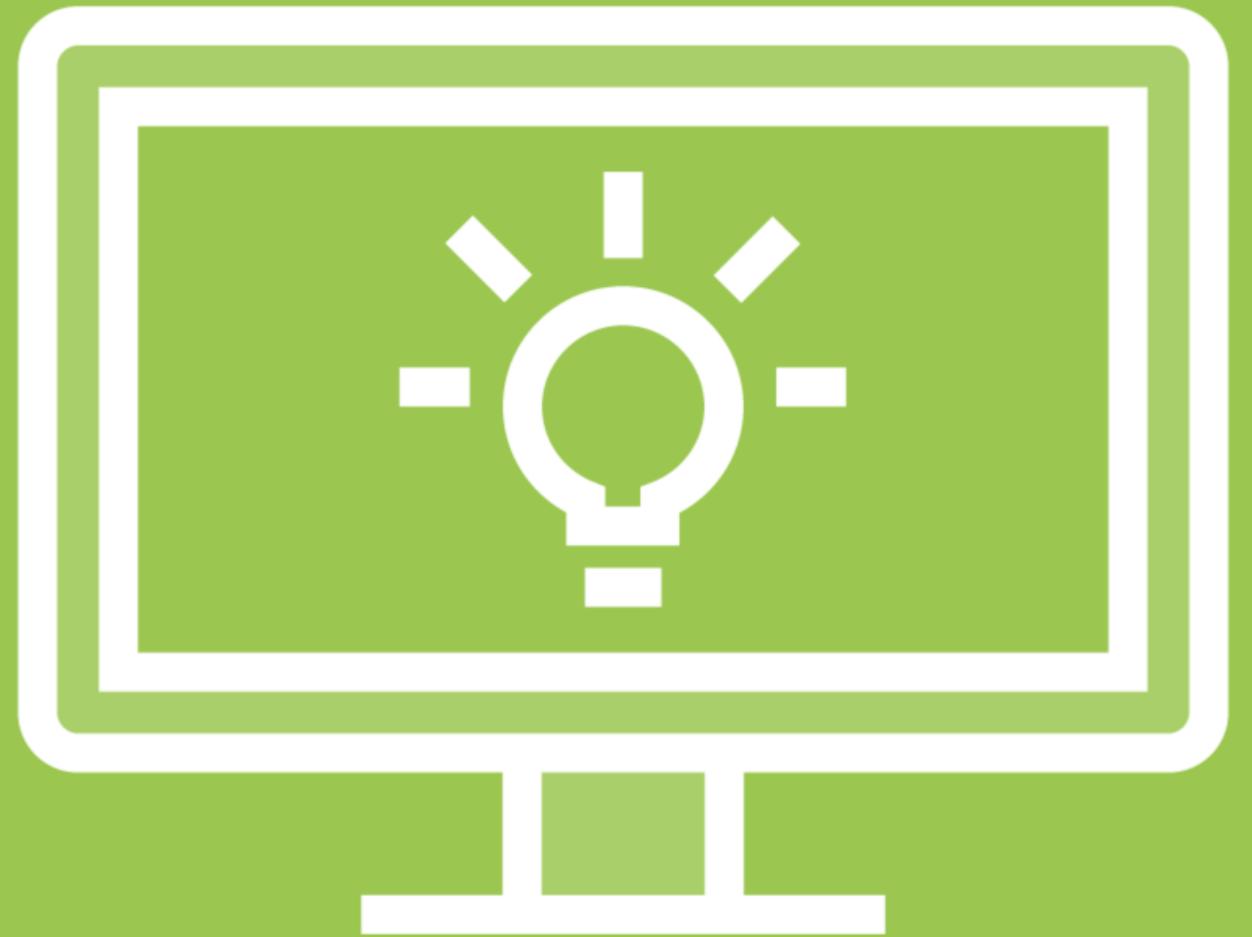
Timing Illustration



| | A | B | C |
|---|---------|-----------|--------|
| 1 | sun | triangle | circle |
| 2 | hexagon | half-moon | square |
| 3 | circle | square | star |

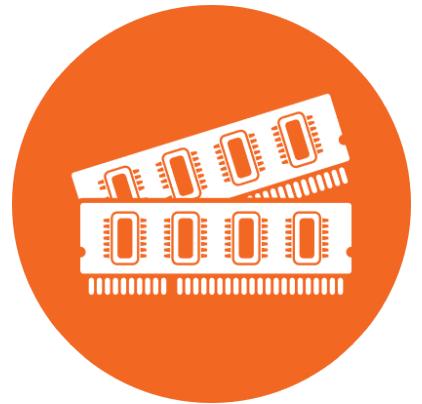


Example: 16-18-18-36



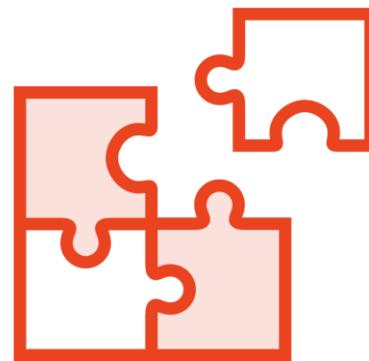
Tech Point
Fault Tolerant
Memory

Memory Pairing Requires Group Installation to Achieve Full Capabilities



Usually 2, but 3, 4, or even more are possible

- Groups are usually color coded, but no standard on colors or how they are grouped



If less than the full group is used (ex. 2 of 4), will operate in dual channel mode instead



Memory must be the same speed and capacity (ideally they are identical)

- If slower, runs at slowest module in the group
- In some cases, if smaller, uses dual channel mode for size in common and single channel for the rest

Fault Tolerance Techniques

ECC vs. non-ECC

Memory mirroring

Memory sparing

Memory Mirroring

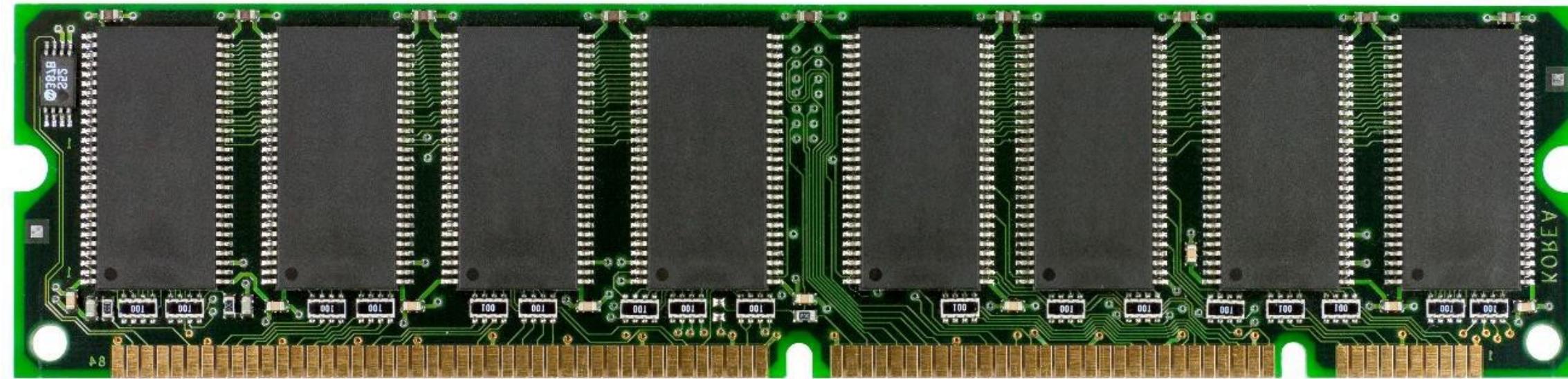
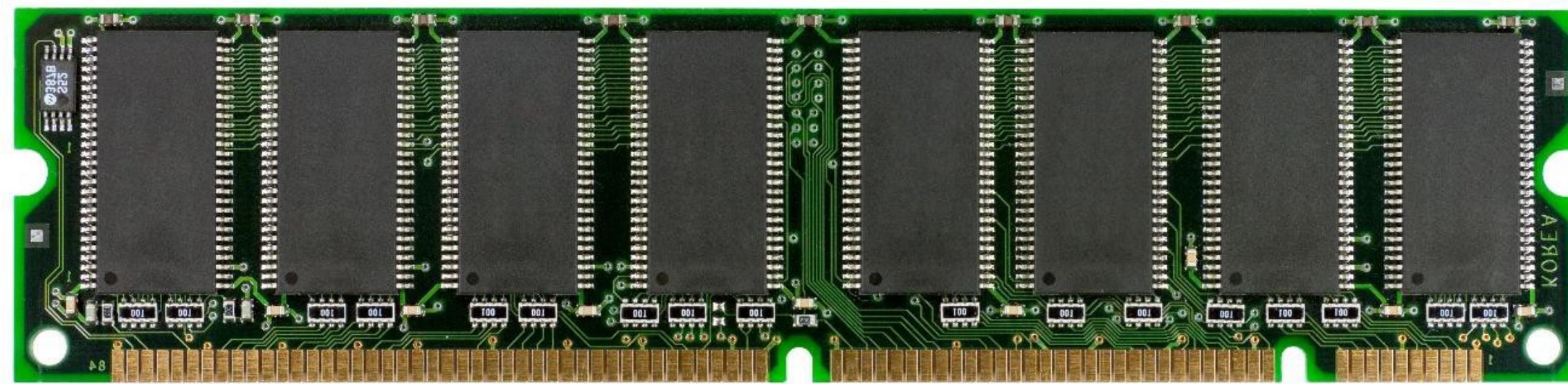
Similar to RAID 1
for RAM

Two copies of data
written to separate
memory channels

- Only half of purchased RAM is usable
- Writes are mirrored to both DIMMs
- Reads are alternated between the copies

Not often used
except for the
most sensitive of
data that cannot
under any
circumstances get
corrupted or lost

Memory Mirroring Example



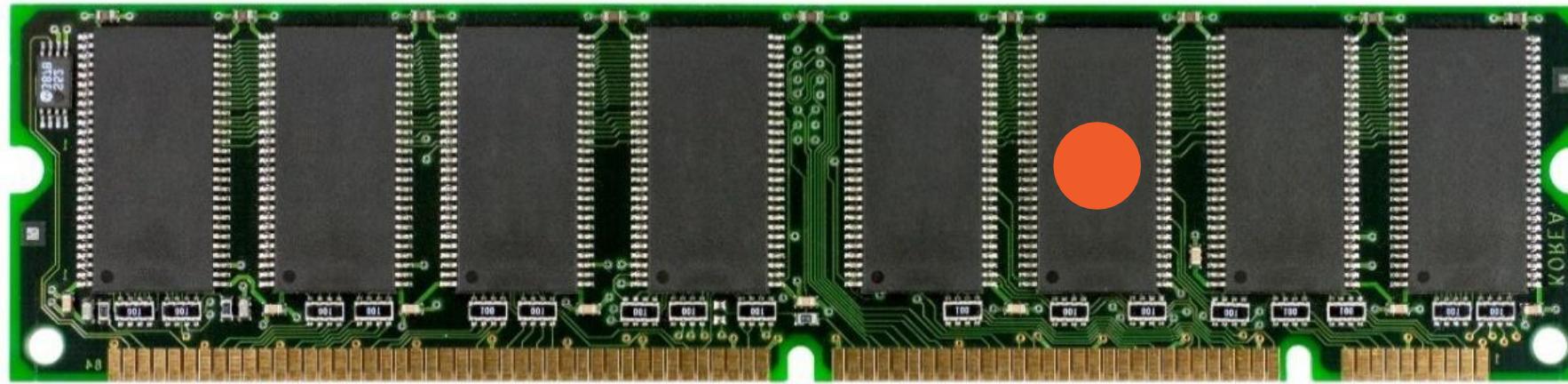
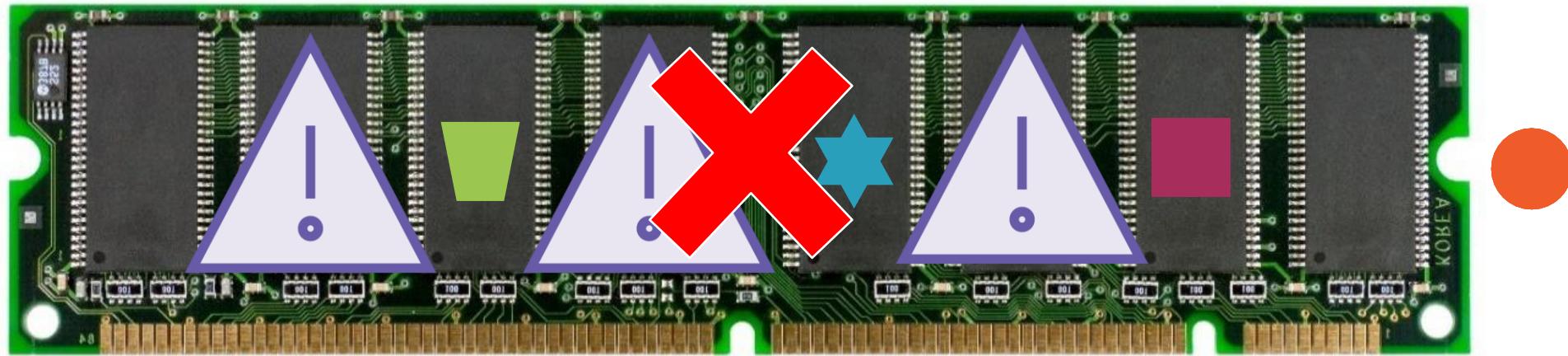
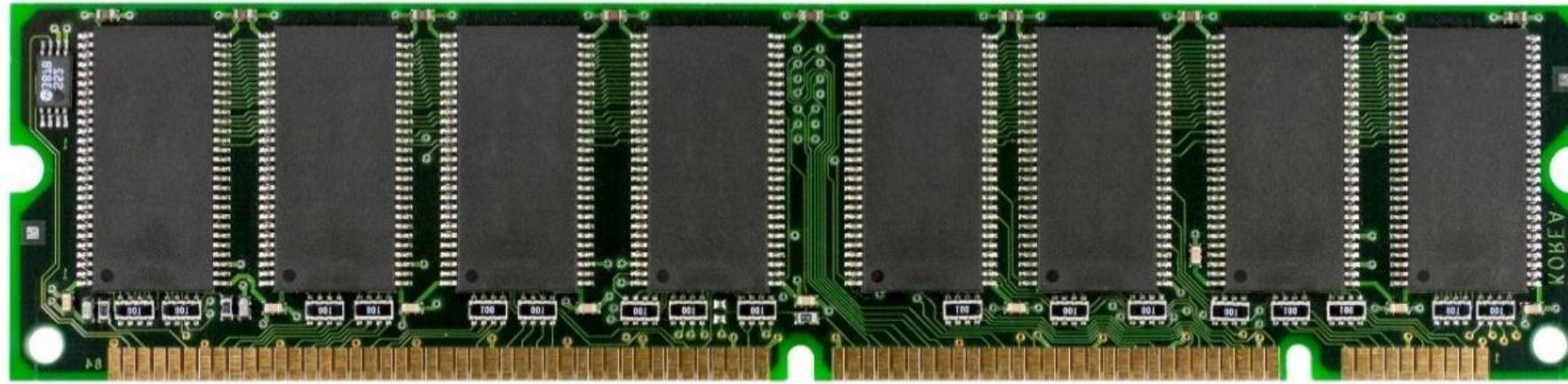
Memory Sparing

Similar to a hot spare in storage

DIMM is not used until a preset number of errors occurs (the threshold), then contents are copied from failing DIMM to the spare

Can be used with ECC memory and / or Memory Mirroring

Memory Sparing Example

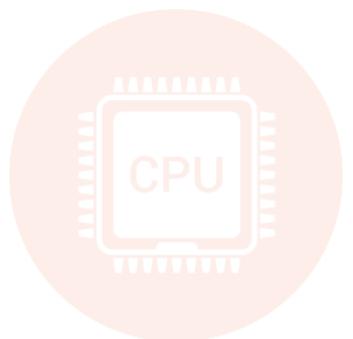




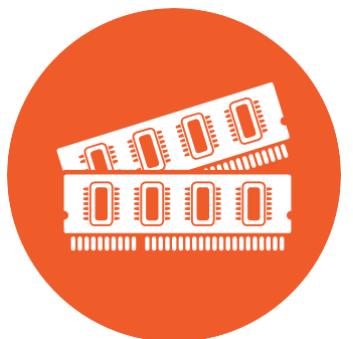
Hindsight



Minimum Hardware



1 CPU Sockets with 6 cores



64 GB of RAM



1 Gbps Networking



512 GB Storage

Recommended Hardware

2 CPU Sockets with 12 cores each

256 GB of RAM

10 Gbps Networking

3 TB Storage

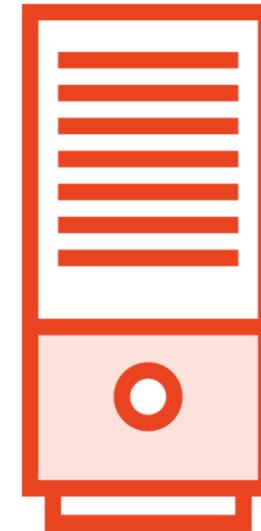


Decision Points

How Critical Is the Server?

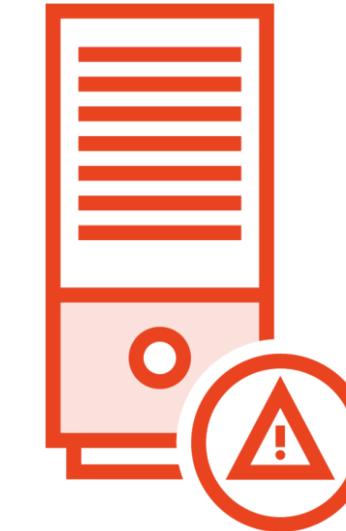
If server downtime will cause lost business or opportunities the server may be Mission Critical to the company, even if the company is small.

When lives are on the line, redundancy is imperative.



Standard Server

Non-ECC Memory



Mission Critical

ECC memory required

Memory spares if available



Lives depend upon it

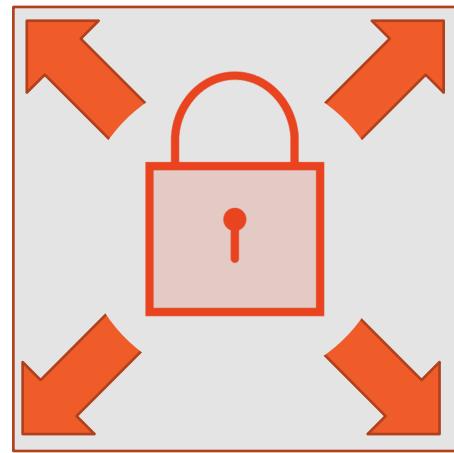
Memory Mirroring
Required

Memory Spares Required

Is Memory Scalability Important?

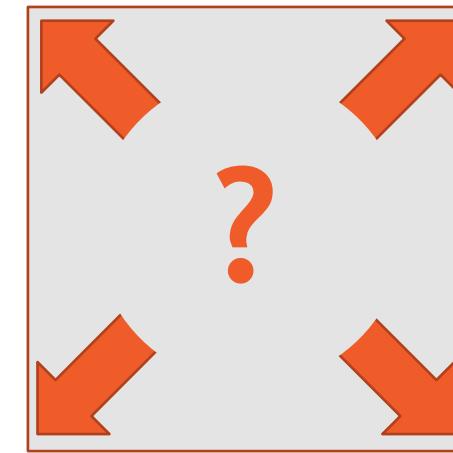
In terms of memory, scalability refers to increasing the RAM capacity of the server later in its lifecycle.

Scalability is important if you intend to add additional roles or applications to the server.



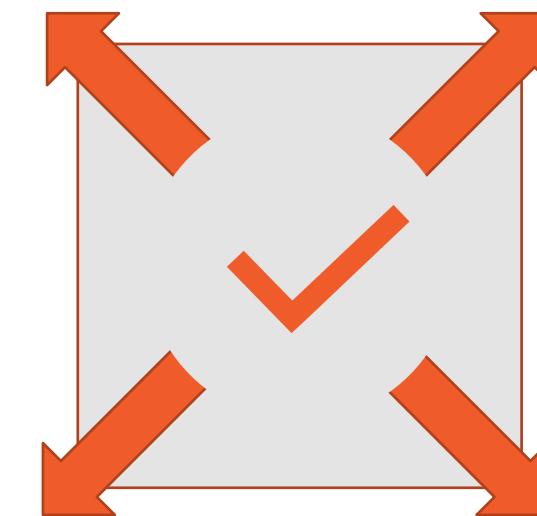
Will not be scaled up

Unregistered (unbuffered) memory is sufficient and less expensive.



May be scaled up

Consider Registered (buffered) memory if supported by the server



Will be scaled up

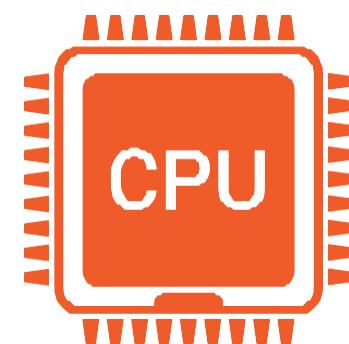
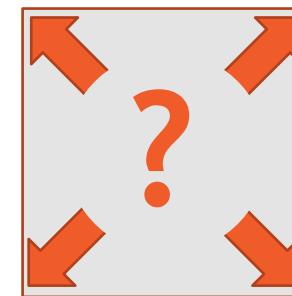
Use Registered (buffered) memory



Globomantics is a small company, but this server is mission critical to them.

They do not intend to scale up the server, but it may be necessary.

The server will be multi-processor, meaning that NUMA must be considered.



ECC Memory will be required,
No need for memory mirroring,
Hot spare memory if available.

Registered (buffered) memory if possible, but not required

Ensure that the memory is evenly divided between the CPU sockets

