# How to build a GPU cluster

Christopher Milan

AI Safety at UCLA

March, 2025

# Who am I

- I am Christopher Milan, the president of AI Safety at UCLA.

# Who am I

- I am Christopher Milan, the president of AI Safety at UCLA.
- I spent a significant amount of my time this past year building out our servers.

## Who am I

- I am Christopher Milan, the president of AI Safety at UCLA.
- I spent a significant amount of my time this past year building out our servers.
- These slides will be available online at ais-ucla.org/~chris.

# What is this talk about?

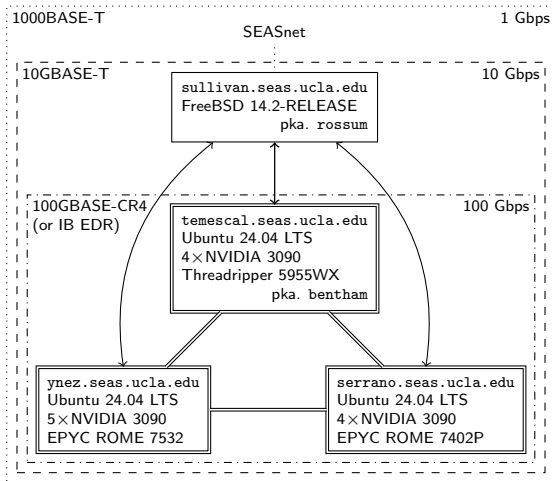While setting up and administrating our servers, I ran into many interesting challenges.

This talk will focus on how we chose our hardware.

There are many other interesting things to learn from this cluster, check out our blog.[1]

---

[1]blog.aisafetyatucla.org

# Server topology

AMD's "Zen 2" server CPUs.

AMD's "Zen 2" server CPUs.

- Memory bandwidth:

# Why AMD (EPYC ROME)?

AMD's "Zen 2" server CPUs.

- Memory bandwidth: 8 channels of DDR4-3200, roughly 200GB/s.

# Why AMD (EPYC ROME)?

AMD's "Zen 2" server CPUs.

- Memory bandwidth: 8 channels of DDR4-3200, roughly 200GB/s.
- Price:

# Why AMD (EPYC ROME)?

AMD's "Zen 2" server CPUs.

- Memory bandwidth: 8 channels of DDR4-3200, roughly 200GB/s.
- Price: A 7532 is about $200 on ebay.
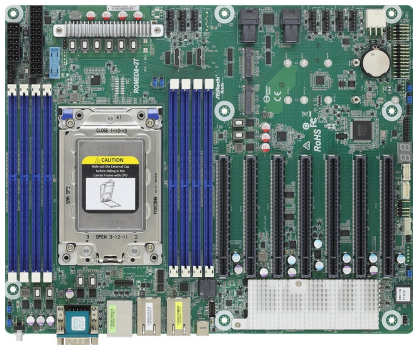
# Why AMD (EPYC ROME)?

AMD's "Zen 2" server CPUs.

- Memory bandwidth: 8 channels of DDR4-3200, roughly 200GB/s.
- Price: A 7532 is about $200 on ebay.
- PCIe lanes:

# Why AMD (EPYC ROME)?

AMD's "Zen 2" server CPUs.

- Memory bandwidth: 8 channels of DDR4–3200, roughly 200GB/s.
- Price: A 7532 is about $200 on ebay.
- PCIe lanes: EPYC ROME has 128 PCIe 4.0 lanes.

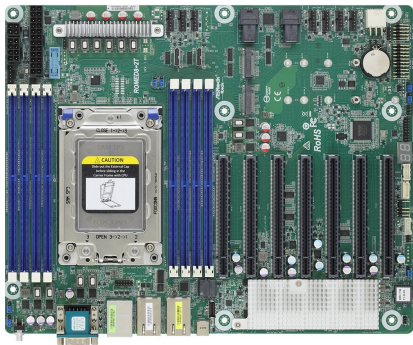Figure: ASRock Rack ROMED8-2T

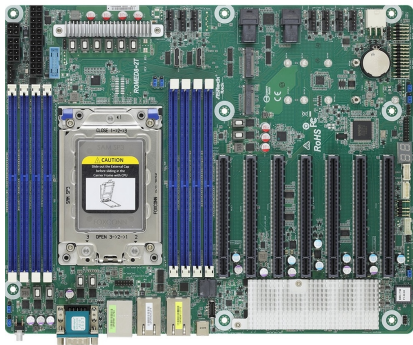- Reasonably priced (at the time...)

Figure: ASRock Rack ROMED8-2T

- Reasonably priced (at the time...)
- 7×PCIe 4.0x16

# ROMED8-2T



Figure: ASRock Rack ROMED8-2T

- Reasonably priced (at the time...)
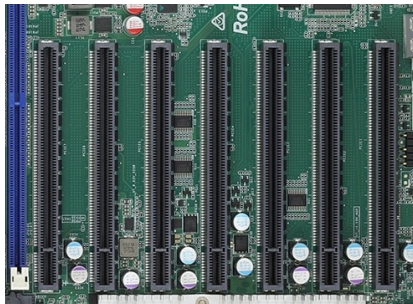- 7×PCIe 4.0x16 (but slot #2 is cursed)

# ROMED8-2T



Figure: ROMED8-2T PCIe slots

- Reasonably priced (at the time...)
- 7×PCIe 4.0x16 (but slot #2 is cursed)
- The slots are too close!

# PCIe risers to the rescue!

- Crypto miners have this exact same problem

# PCIe risers to the rescue!

- Crypto miners have this exact same problem
- Solution: PCIe risers

# PCIe risers to the rescue!

- Crypto miners have this exact same problem
- Solution: PCIe risers
- Surely nothing could go wrong...



Figure: PCIe risers on `temescal`

# Signal integrity errors



Figure: AER errors, image from Nathan Odle

# Signal integrity errors



Figure: ROMED8-2T block diagram

Figure: PCIe slots on ROMED8-2T

# Signal integrity errors



Figure: Remember slot 2? This is him now

# Signal integrity errors



All this switching causes signal loss!

Similarly, adding risers does too.

Figure: Remember slot 2? This is him now

# Another issue with PCIe risers

PCIe devices draw power directly from the PSU AND from the PCIe bus itself.

PCIe devices draw power directly from the PSU AND from the PCIe bus itself.

Why is this a problem?

PCIe devices draw power directly from the PSU AND from the PCIe bus itself.

Why is this a problem? Because at 350W/ea. one power supply is not enough for 6 GPUs.

# Another issue with PCIe risers

PCIe devices draw power directly from the PSU <span style="color:red">AND</span> from the PCIe bus itself.

Why is this a problem? Because at 350W/ea. one power supply is not enough for 6 GPUs.

Taping off the PCIe slot power just causes the cards to fail to be recognized.

# How to fix this?

Some options:

1. Just ignore it and have fewer GPUs.

# How to fix this?

Some options:

1. Just ignore it and have fewer GPUs.
2. Just get better risers, and still have fewer GPUs.

# How to fix this?

Some options:

1. Just ignore it and have fewer GPUs.

2. Just get better risers, and still have fewer GPUs.

3. How do the pros do it?

# How to fix this?

Some options:

1. Just ignore it and have fewer GPUs.
2. Just get better risers, and still have fewer GPUs.
3. How do the pros do it?



Figure: SlimSAS/MCIO connectors

# Dartmouth Summer Research Project on AI

We propose [a study] to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.

# Dartmouth Summer Research Project on AI

We propose [a study] to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.

An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves.

# Dartmouth Summer Research Project on AI

We propose [a study] to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.

An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves.

We think that a significant advance can be made [. . .] if a carefully selected group of scientists work together for a summer.

We propose [a study] to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.

An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves.

We think that a significant advance can be made [. . .] if a carefully selected group of scientists work together for a summer.

— A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence, August 31, 1955

We propose [a study] to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.

An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves.

We think that a significant advance can be made [. . .] if a carefully selected group of scientists work together for a summer.

— A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence, August 31, 1955

Solving problems requires tons of repeated trial and error.

# The alignment problem

For a sufficiently advanced system, ensuring its goals are the same as ours is a non-trivial task (see universal paperclips).

# The alignment problem

For a sufficiently advanced system, ensuring its goals are the same as ours is a non-trivial task (see universal paperclips).

It stands to reason that alignment will be similarly difficult.

# The alignment problem

For a sufficiently advanced system, ensuring its goals are the same as ours is a non-trivial task (see universal paperclips).

It stands to reason that alignment will be similarly difficult.

However, we only get one try.

## Join us

Visit our website at ais-ucla.org.

Some programs we run:

- Intro to AI Safety Fellowship (ais-ucla.org/fellowships).
- Upskilling Tracks (ais-ucla.org/upskilling-tracks).
- Reading Group (ask Will).
- Server access (ask me).