



```
`torch.OutOfMemoryError`
```

-- finally, an exception you'll never catch again!

The AI Memory Problem

**GPU Memory Capacity is
NOT ENOUGH**

**Agentic Context at
Scale (10-100TB)**

**GPU Memory
(288 GB)**

**Model
(50 GB)**

**Context
(100+ GB)**

The AI Memory Problem

**GPU Memory Capacity is
NOT ENOUGH**

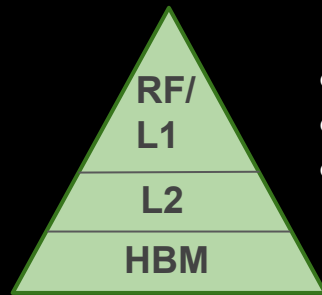
**Agentic Context at
Scale (10-100TB)**

**GPU Memory
(288 GB)**

**Model
(50 GB)**

**Context
(100+ GB)**

**AI Memory Hierarchy Goes
Beyond GPUs**



- model
- current context
- weights (training)

The AI Memory Problem

**GPU Memory Capacity is
NOT ENOUGH**

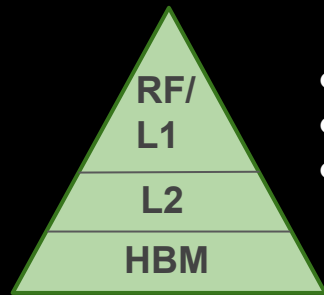
**AI Memory Hierarchy Goes
Beyond GPUs**

**Agentic Context at
Scale (10-100TB)**

**GPU Memory
(288 GB)**

**Model
(50 GB)**

**Context
(100+ GB)**



- model
- current context
- weights (training)

G2. System DRAM

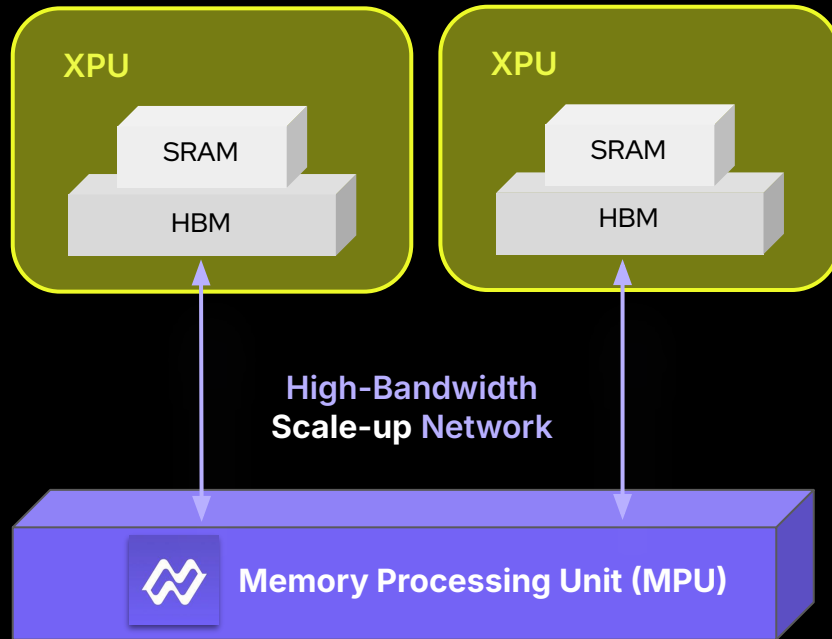
G3. NVMe (SSD)

G4. Remote Storage

- inactive ctx
- sparse attention
- serverless
- activations/
optimizer state

Netpreme™ X-Mem™ Technology

The world's fastest memory-to-compute fabric



US Patent App. 63/793,877 US Patent App. 63/838,102 US Patent App. 63/951,047

Now Integrated with




What you get

10x faster than RDMA/System Memory

High-Capacity near-memory computing

50% efficient at scale

Find us at

BOOTH #10

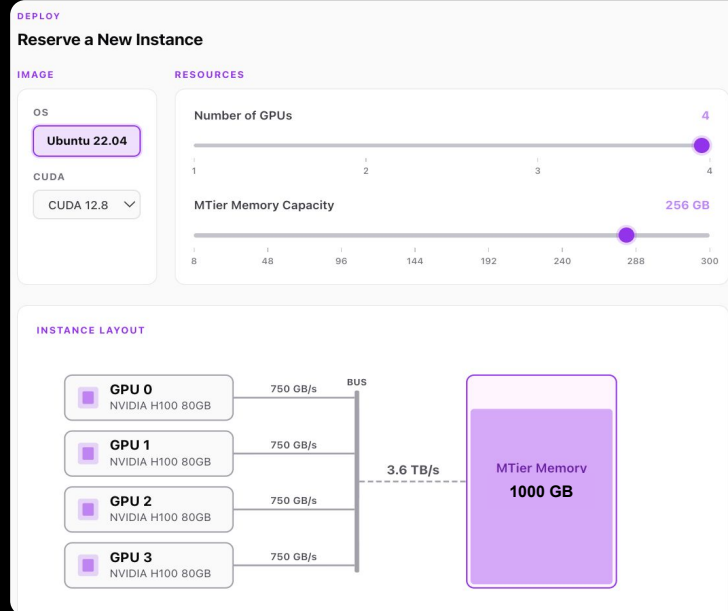
Try X-Mem™ Pilot

If you are....

ML Serving Engine Builder

AI Infrastructure Engineer

Struggling with long context?
with complex agentic flows?
with large batches?
Or you tired to see
``torch.OutOfMemoryError``



We're hiring

ML System

ML Kernels

ASIC and more.....

