

Inference Model to Metal

dsteeve



TRAINING



Research One of something More is better Fast iteration Python = 0

Production Billions of somethings Less is better Predictable latency Python = Q



Why do we need another framework?

HOW STANDARDS PROLIFERATE:
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)

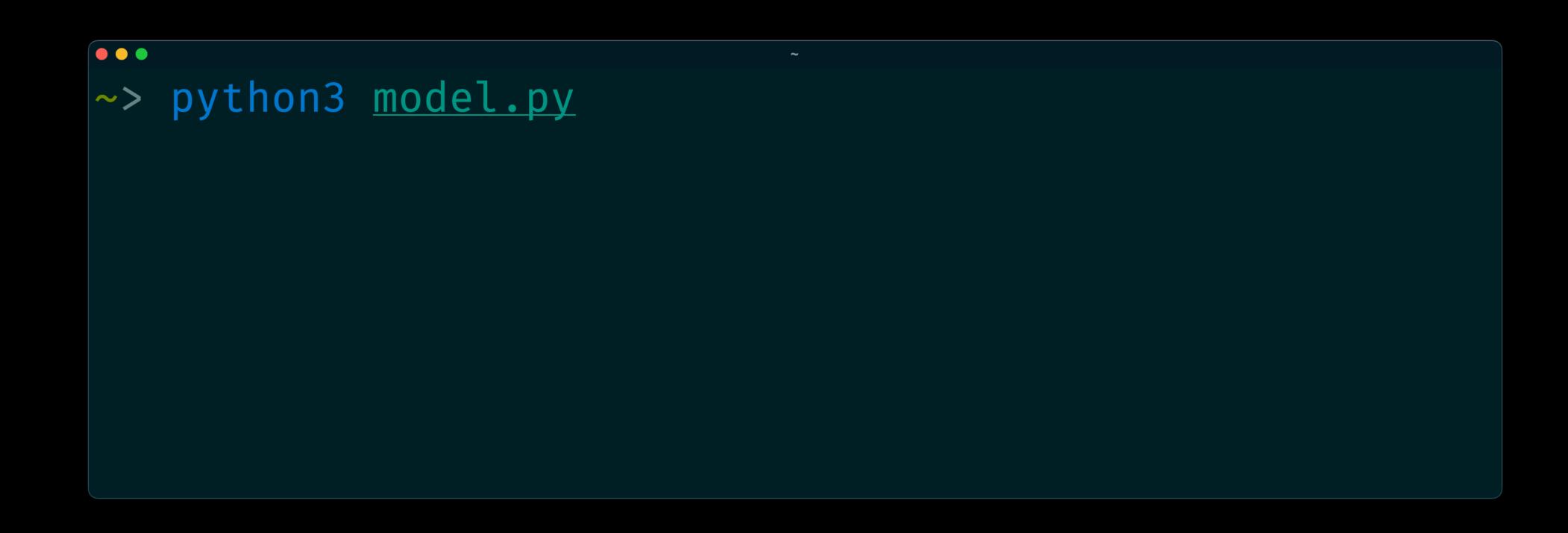
SITUATION: THERE ARE 14 COMPETING STANDARDS.



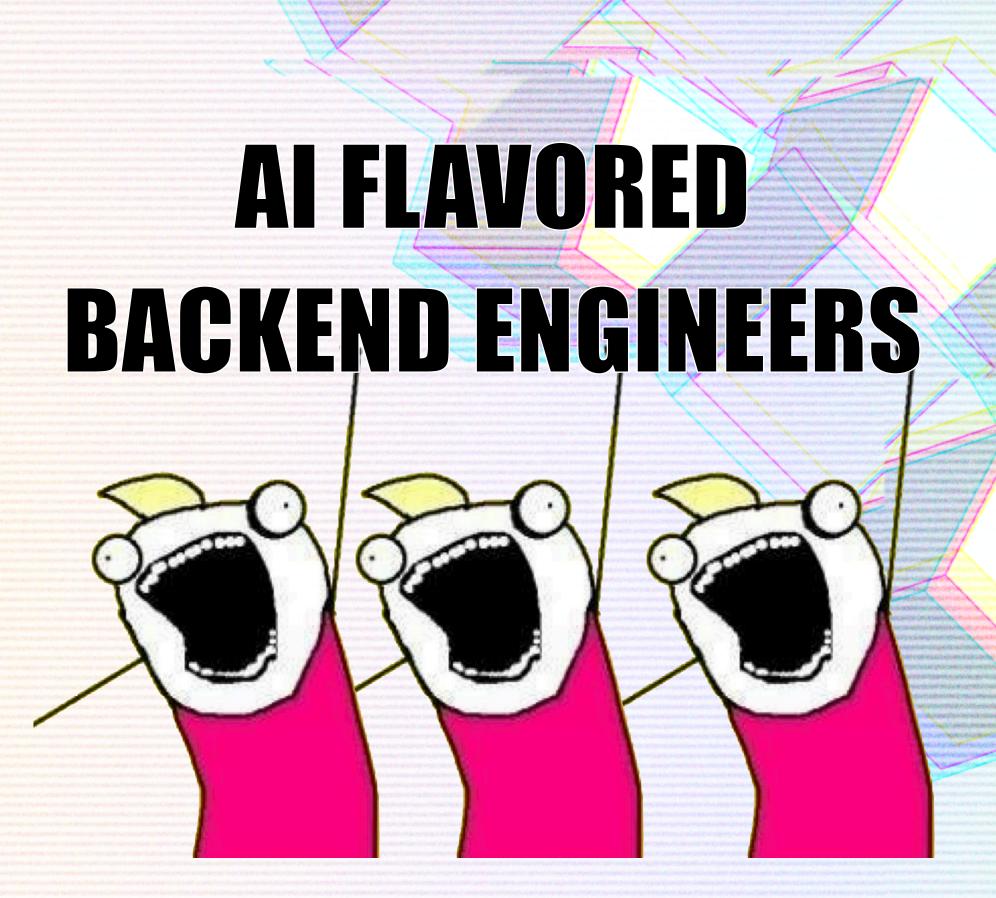
5∞N:

SITUATION: THERE ARE 15 COMPETING STANDARDS.









WHATDO WE WANT ??

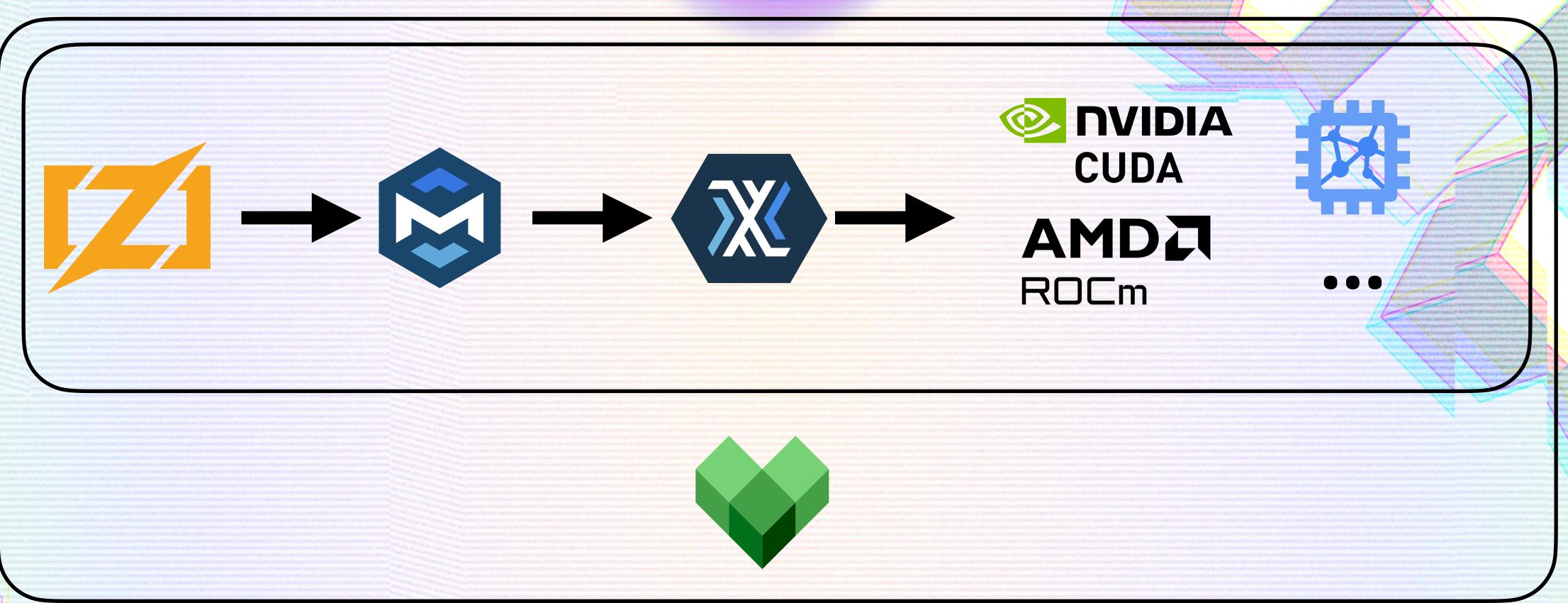


ACCELERATOR AGNOSTICITY,







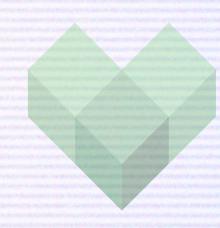


Inference only Zig as the frontend Zero Python Modular code Focus on maintainability

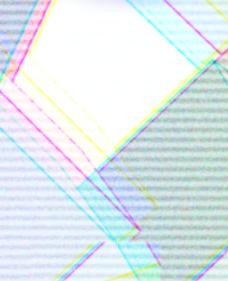


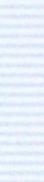












```
1 const zml = @import("zml");
 3 const Mnist = struct {
       fc1: Layer,
       fc2: Layer,
       const Layer = struct {
           weight: zml.Tensor,
           bias: zml.Tensor,
10
           pub fn forward(self: Layer, input: zml.Tensor) zml.Tensor {
               return self.weight.matmul(input).add(self.bias).relu();
       };
16
       pub fn forward(self: Mnist, input: zml.Tensor) zml.Tensor {
           var x = input.flattenAll().convert(.f32);
           const layers: []const Layer = &.{ self.fc1, self.fc2 };
18
           for (layers) |layer | {
19
               x = zml.call(layer, .forward, .{x});
20
           return x.argMax(0, .u8).indices;
24 };
```













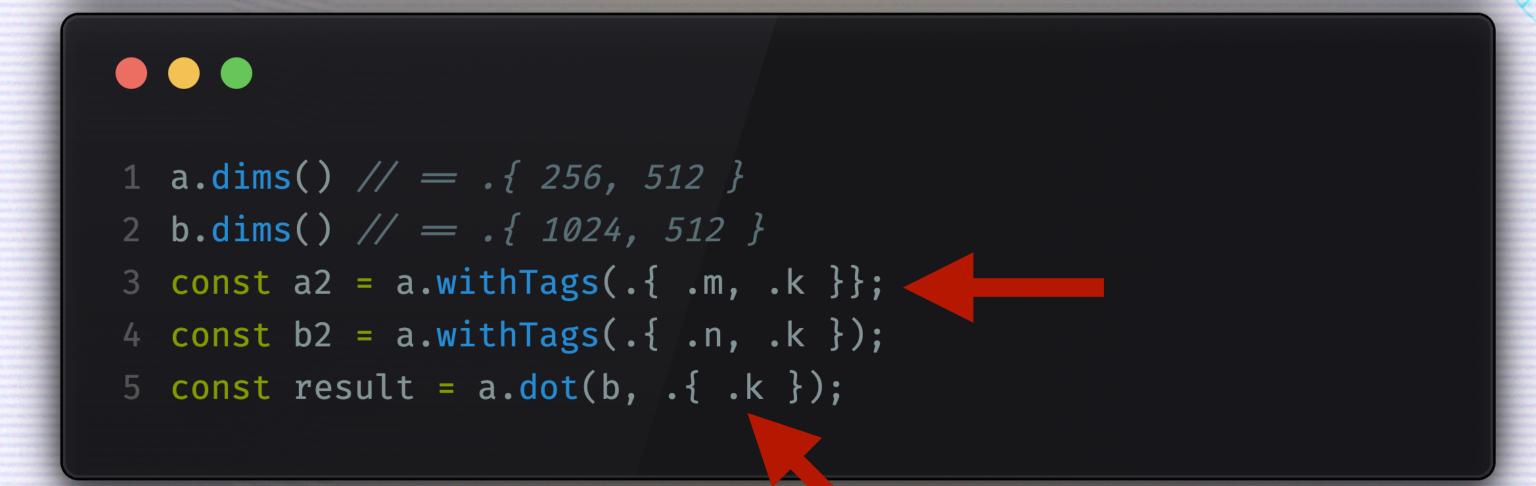








```
const sh = zml.Shape.init(.{ .w = 1920, .h = 1080, .c = 3 });
sh.dim(0) // = 1920
sh.dim(.w) // = 1920
const sh2 = sh.transpose(.{ .h, .w, .c });
sh.dim(0) // = 1080
sh.dim(.w) // = 1920
```









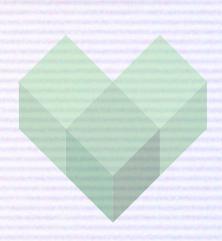


```
1 pub fn sdpa(q: Tensor, k: Tensor, v: Tensor, opts: SdpaOpts) Tensor {
       /// ... error checking ...
       const attn_mask = if (opts.attn_mask) | m | m else null;
       const sqrtHeadDim: f16 = 1.0 / std.math.sqrt(@as(f16, @floatFromInt(dims.hd)));
       const scale_logit = if (opts.scale) |s| s else Tensor.scalar(sqrtHeadDim, .f16);
       k = k.mul(scale_logit.convert(k.dtype()));
       var attn_weights = q.dot(k, .{.hd});
10
       if (attn_mask) | mask | {
           attn_weights = attn_weights.add(mask.broadcastLeft(attn_weights.shape()));
       attn_weights = attn_weights.convert(.f32);
       if (opts.bias) | bias | {
16
           attn_weights = attn_weights.add(bias);
18
       attn_weights = attn_weights.softmax(.k).convert(q.dtype());
19
20
       var attn = attn_weights.dot(v, .{.k});
       return attn.transpose(q.shape());
23 }
```







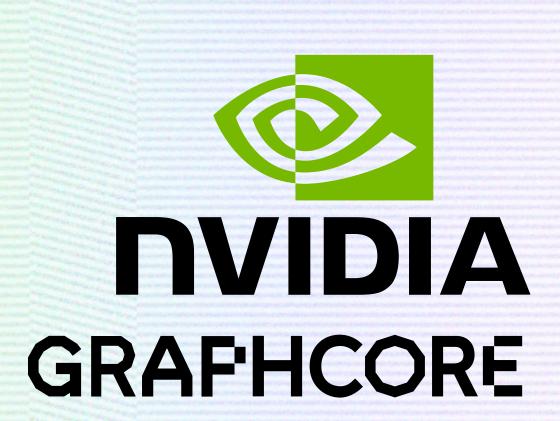




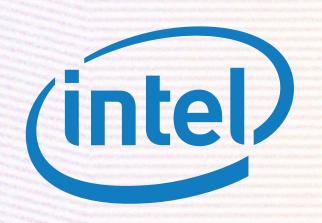
OpenXLA Founding Members



CIrm

















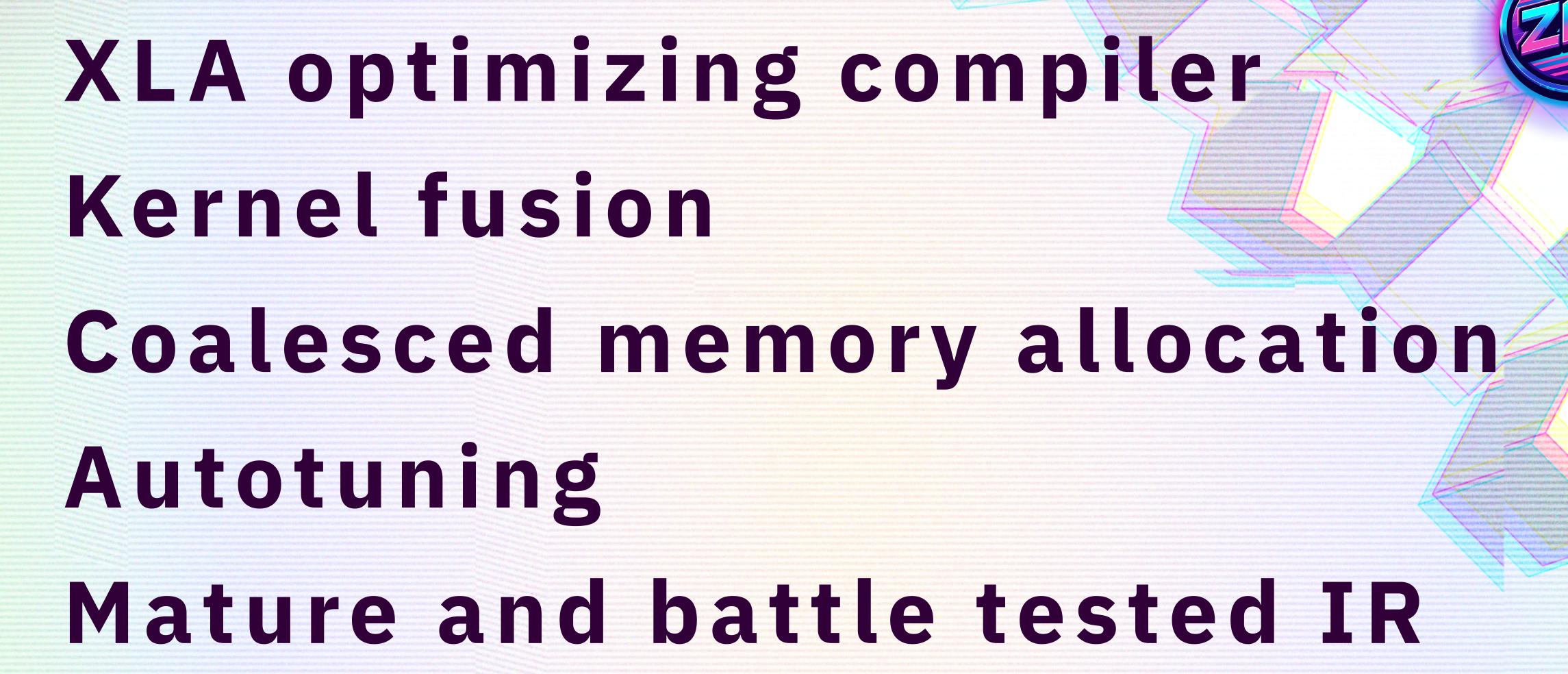
Apple







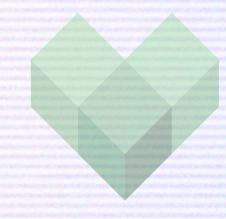


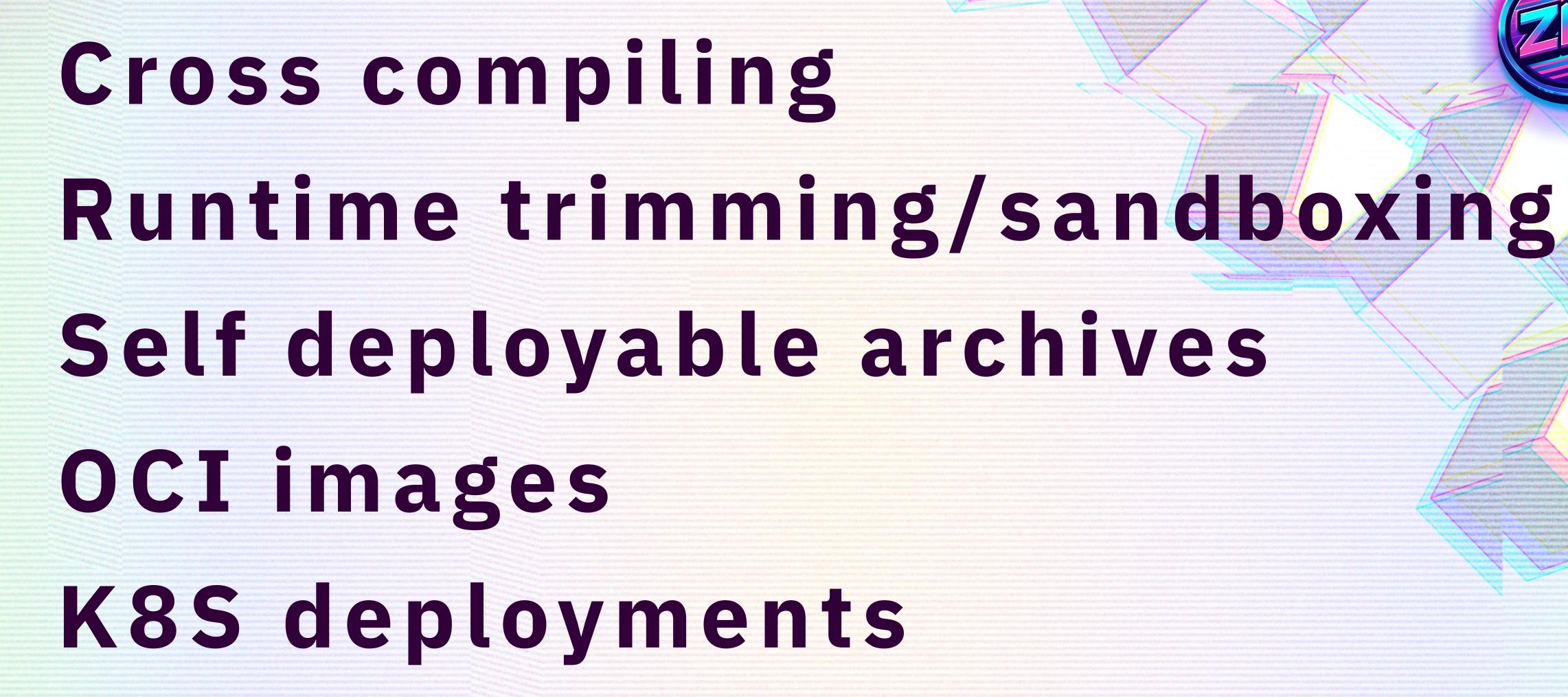












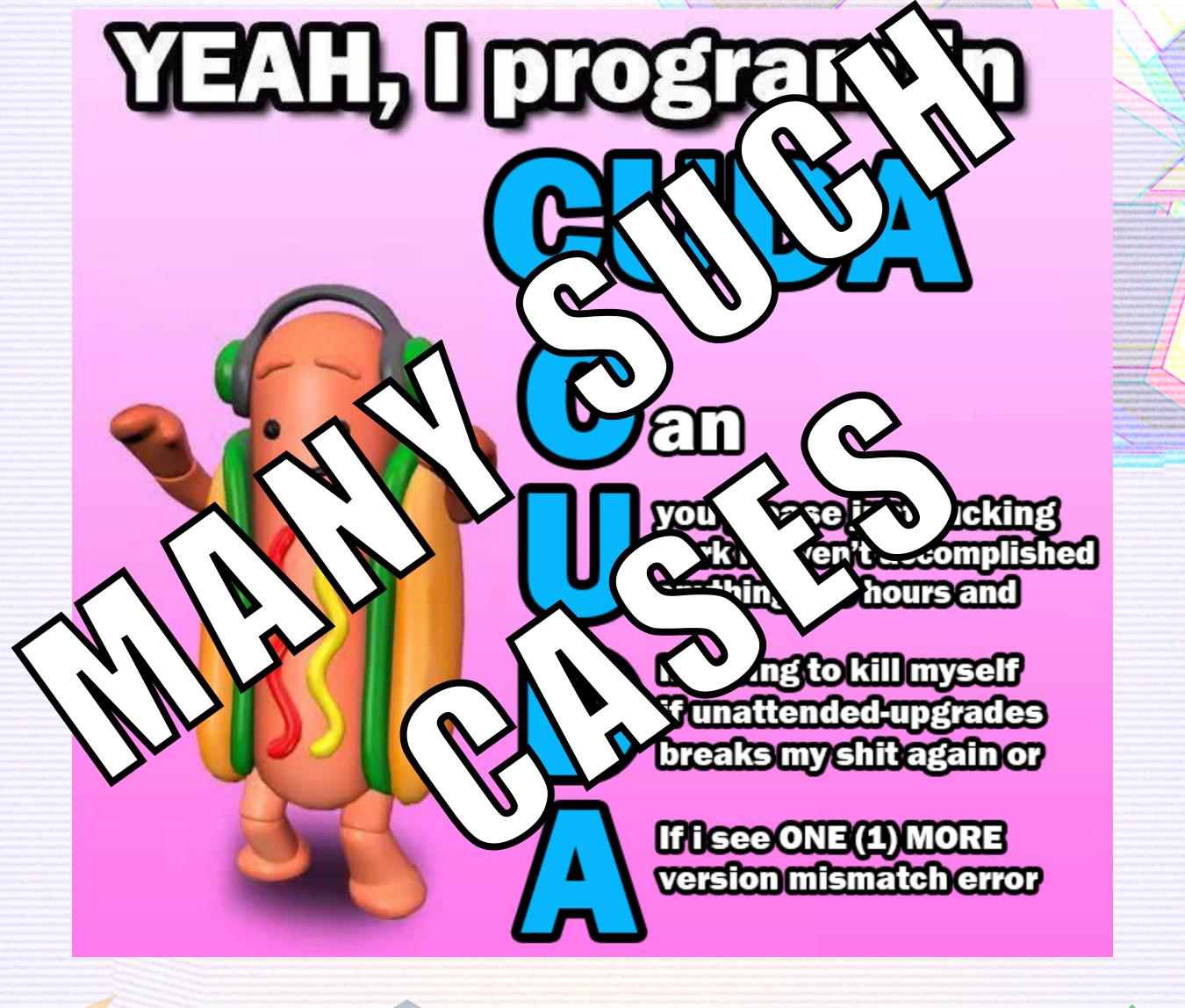


































1 \$ bazel run -c opt //mnist --@zml//runtimes:cuda=true



1 \$ bazel build -c opt //mnist:archive --@zml//runtimes:rocm=true --platforms=@zml//:linux_amd64



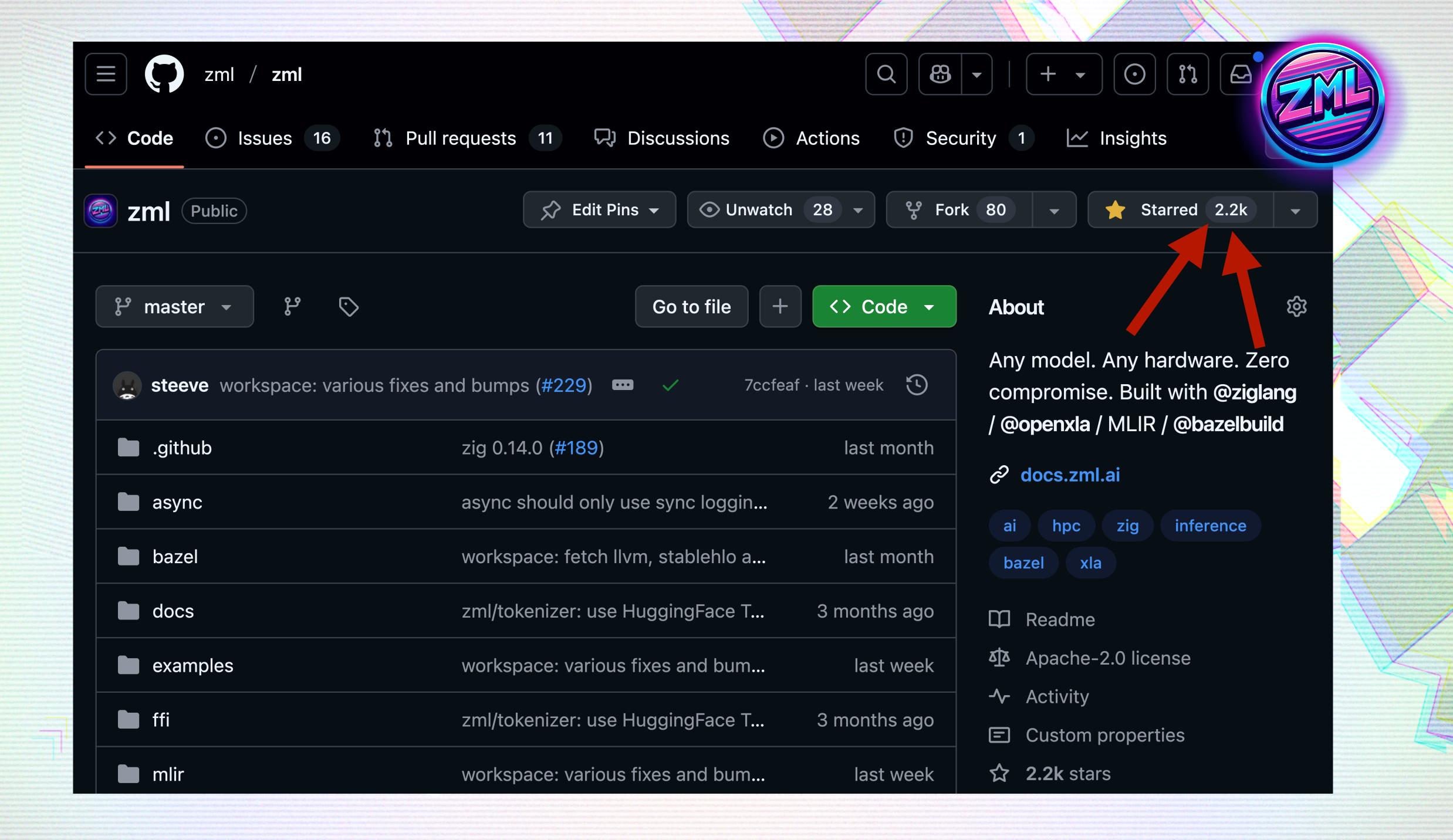
1 \$ bazel run -c opt //mnist:push --@zml//runtimes:cuda=true --@zml//runtimes:tpu=true

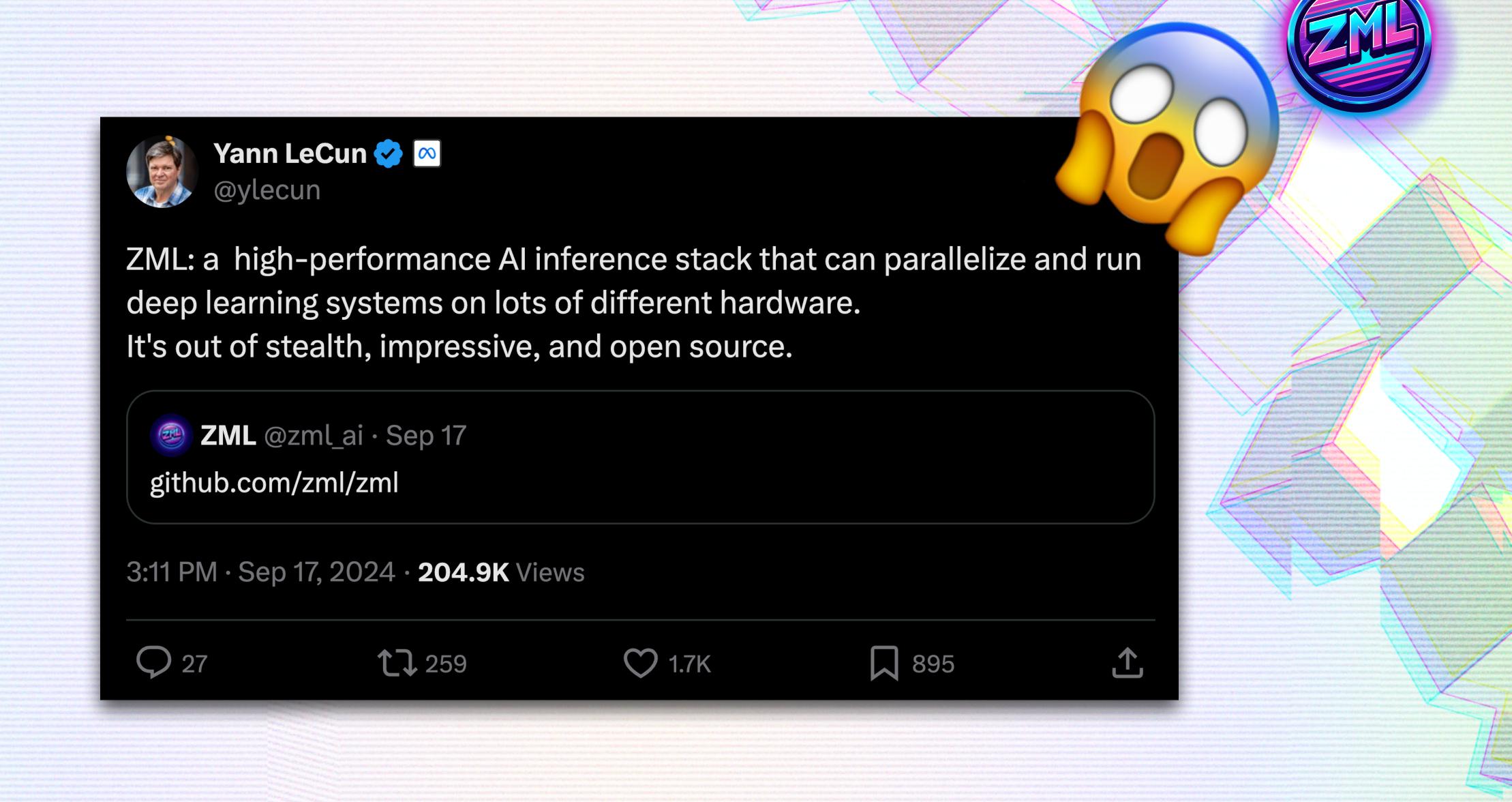








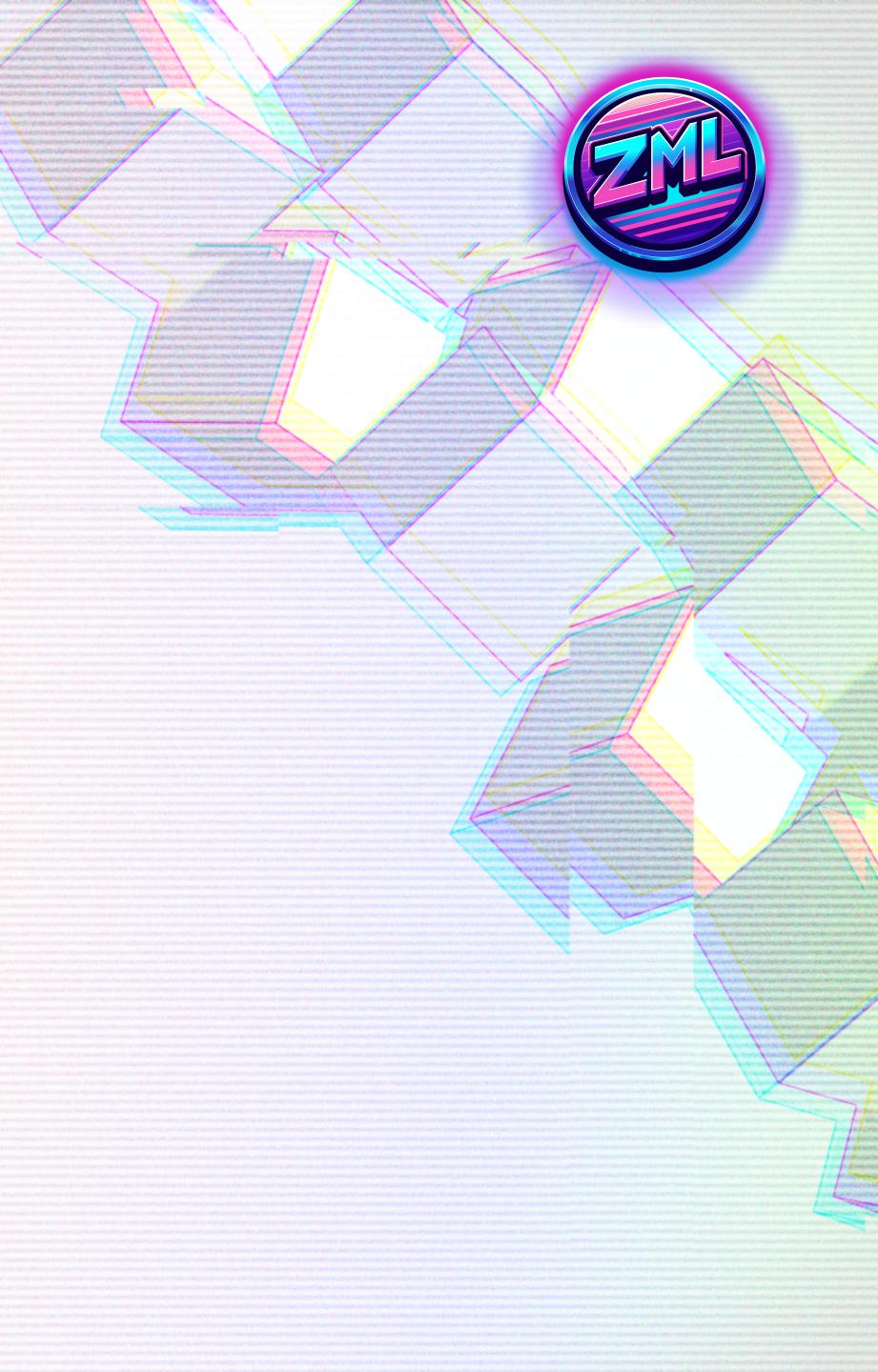








More chips More models More modalities More integrations





Serving Stack





```
info(pjrt): Loaded library: libpjrt_cpu.so
info(pjrt): Loaded library: libpjrt_rocm.so
2025-04-28 16:25:43.764975: I xla/service/service.cc:152] XLA service 0x1b3c60b0 initialized for platform ROCM (this do
2025-04-28 16:25:43.765008: I xla/service/service.cc:160] StreamExecutor device (0): AMD Radeon Graphics, AMDGPU ISA
info(zml/context): Available Platforms:
info(zml/context):
info(zml/context): • cuda
info(zml/context): <a>✓</a> rocm (AUTO-SELECTED)
info(zml/context):

    #0: AMD Radeon Graphics

info(zml/context):
                     • tpu
info(zml/context): • neuron
debug(zml/io): Loading shard: model-00004-of-00004.safetensors
debug(zml/io): Loading shard: model-00001-of-00004.safetensors
debug(zml/io): Loading shard: model-00002-of-00004.safetensors
debug(zml/io): Loading shard: model-00003-of-00004.safetensors
                Loading Llama weights from ../zml++huggingface+Meta-Llama-3.1-8B-Instruct/model.safetensors.index.json.
info(zml): Compiling llama.LlamaLM.forward with { Shape({s=64,u32}), Shape({s=64,u32}), Shape({s=64,u32}), meta.MapType
,hd=128,bf16}), .v = Shape({layer=32,b=64,k=256,h=8!,hd=128,bf16}), .layer_index = Shape({u32}) }, meta.MapType(tensor.
tablehlo.RngAlgorithm.Type.DEFAULT } }
info(zml): Compiling llama.LlamaLM.forward with { Shape({s=320,u32}), Shape({s=66,u32}), Shape({s=66,u32}), meta.MapTyp
!,hd=128,bf16}), .v = Shape({layer=32,b=64,k=256,h=8!,hd=128,bf16}), .layer_index = Shape({u32}) }, meta.MapType(tensor
stablehlo.RngAlgorithm.Type.DEFAULT } }
```



zmlai/llmd

By <u>zmlai</u> · Updated 3 months ago

LLMD is a high performance LLM server, built by ZML.ai

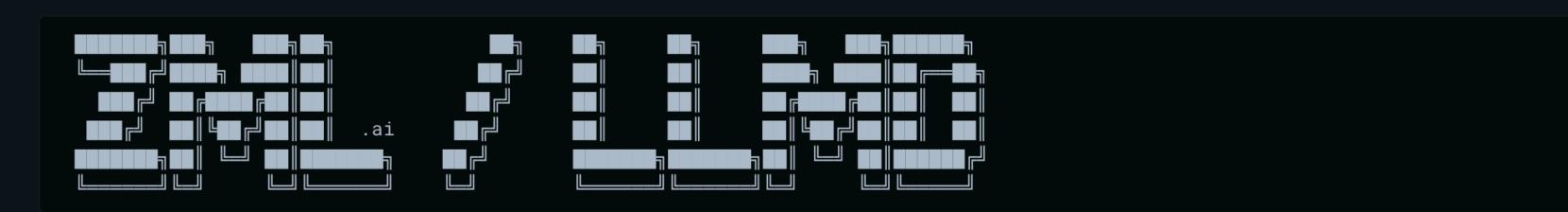
MAGE

LANGUAGES & FRAMEWORKS MACHINE LEARNING & AI WEB SERVERS

☆6 <u>↓</u>282

Overview

Tags



LLMD - High Performance LLM Inference Server

This is a technical preview of LLMD, the ZML.ai 🖰 homegrown high-performance LLM inference server. It is not intended for production use. Please report any issues to the ZML team.

LLMD is a work in progress and may change significantly in the future. Some core features are purposedly disabled.

Get blazing-fast inference with a tiny 2.4GB container that works on both NVIDIA and AMD GPUs.

Preview Features

- Easy Setup Just mount your model and run
- Cross-Platform GPU Support Single container works on both NVIDIA and AMD GPUs
- Lightweight Only 2.4GB container size
- OpenAl API Compatible Drop-in replacement for OpenAl endpoints
- High Performance Optimized for fast inference

Quick Start



