# Group G11

NanoGPT is a minimal GPT-2 implementation by Andrej Karpathy that inspired a community speedrun challenge: reaching a validation loss of 3.28 as fast as possible on 8×H100 GPUs, with a current record of 1.435 minutes. In this work, we port and optimize this training pipeline in JAX for Google TPU v6e hardware, investigating how this architecture competes on a well-defined benchmark. Our implementation incorporates several algorithmic improvements including rotary positional embeddings (RoPE),the Muon optimizer with layer-wise learning rate decay for transformer weights, with Adam used for token embeddings and the language model head each with a warmdown schedule, logit soft-capping based on the Gemma 2 paper, and architecture dimensions specifically tuned to TPU v6e's hardware constraints. Preliminary results show we achieve a validation loss of <3.28 in 24 minutes on 8×TPU v6e.

# Project

Based on modded-nanoGPT

Train model to 3.28 validation loss on FineWeb dataset

Use Google TPU v6e-8 chips

Current world record 2026-03-18 on 8 H100 is 1.435 minutes, original time was 45 minutes
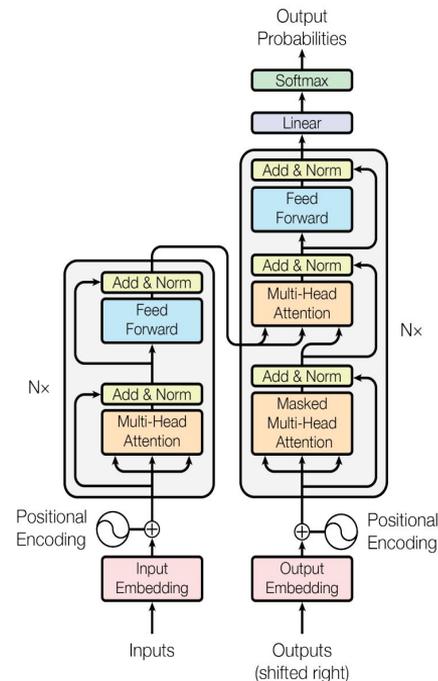


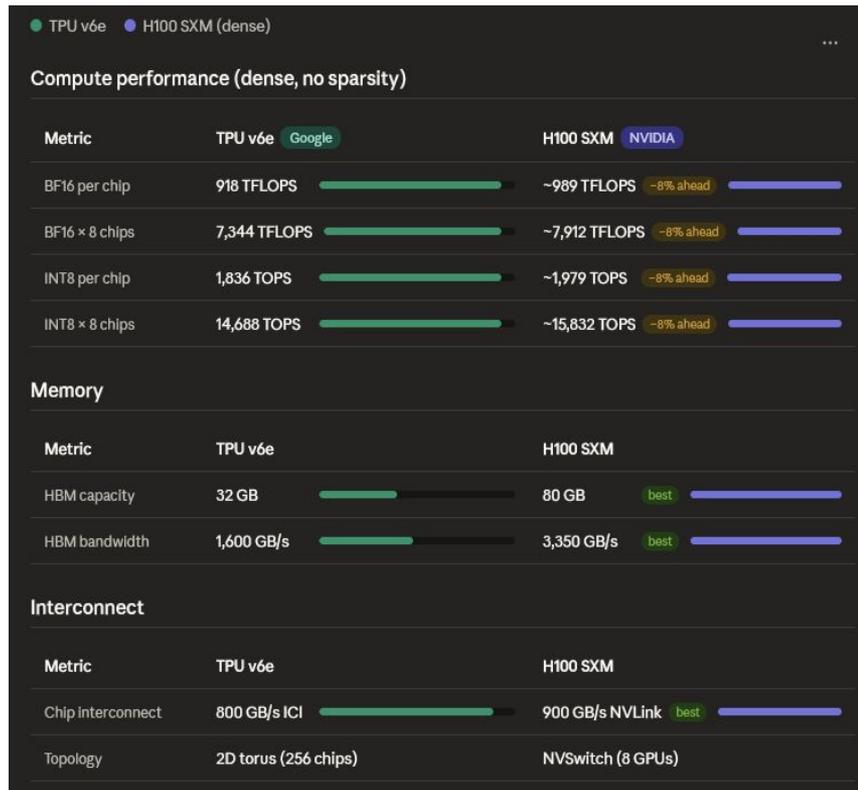Figure 1: The Transformer - model architecture.

# Setup

Each run used the same TPU v6e-8 chips

Same 10 billion token subset each run already tokenized.

Each run had the same key to handel randomness

In order for run to count Val loss < 3.28

# Run 1: Baseline

Total time 105 min vs expected 31 min
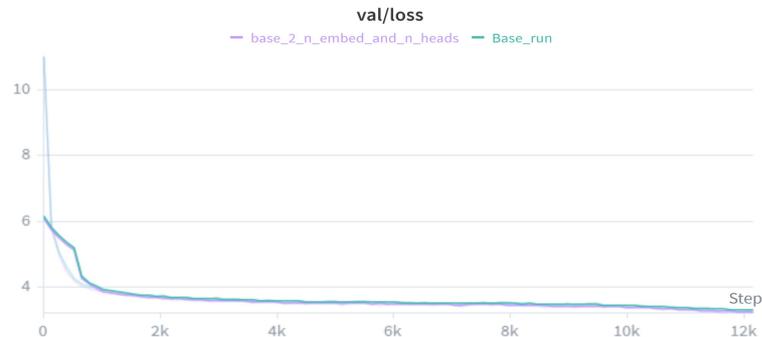
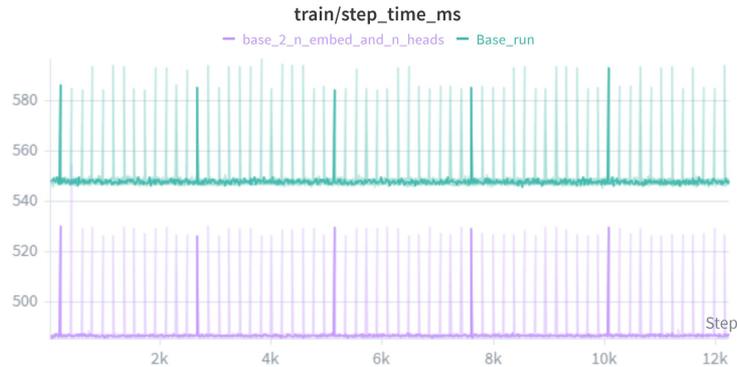Issues: Poor data loading, weights handling

Notable implementations: Rotary Positional Embedding, high learning rate

train/step_time_ms

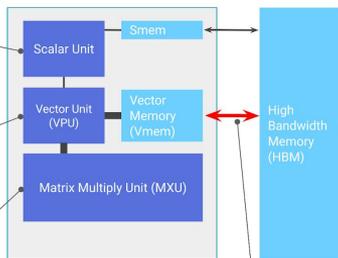val/loss

# Run 2: Hardware-Aware Optimization

Changed the so that head dim = 256

Improvement 13% over baseline



The **Scalar Unit** sort of acts like a CPU 'dispatching' instructions to the VPU and MXU

The **VPU** performs elementwise operations (e.g. activations), loads data into the MXU

**The MXU** performs matrix multiplications - and is therefore our driver of chip FLOP/s.

*Abstract layout of a TPU TensorCore.*

**HBM** stores the weights, activations, optimiser states, new batch data etc

**HBM bandwidth:** determines how fast data goes to and from the computational elements



**train/step_time_ms**
— base_2_n_embed_and_n_heads  — Base_run



**val/loss**
— base_2_n_embed_and_n_heads  — Base_run

# Run 3: Relu^2 and RoPE frequency

12.5% speed up / step, total time 65 minutes
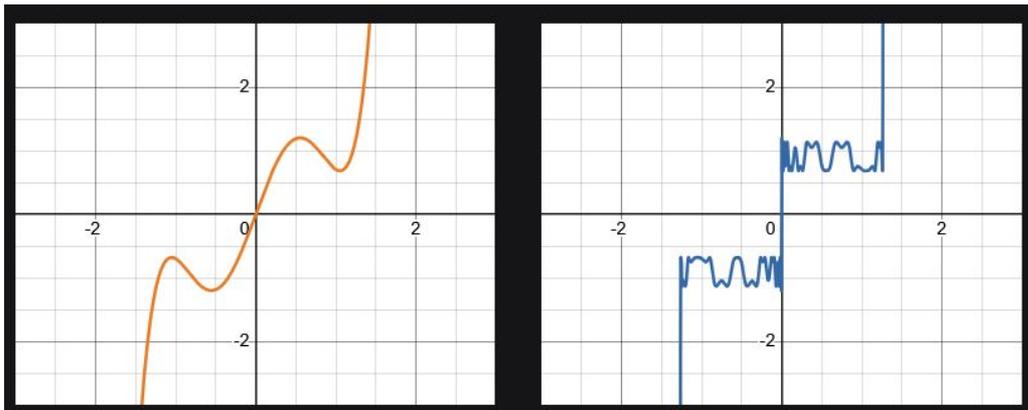
Changed from Gelu to Relu^2

RoPE frequency from 10000 to 1024 due to lower context length from original paper.

# Run 4: Muon

Total time 37 minutes 76% then run 3

Apply the function left 5 times get right

**Algorithm 2** Muon

**Require:** Learning rate $\eta$, momentum $\mu$
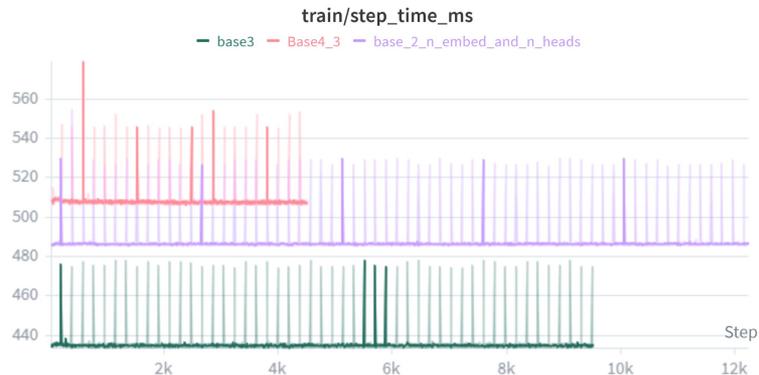1: Initialize $B_0 \leftarrow 0$
2: **for** $t = 1, \ldots$ **do**
3:     Compute gradient $G_t \leftarrow \nabla_\theta \mathcal{L}_t(\theta_{t-1})$
4:     $B_t \leftarrow \mu B_{t-1} + G_t$
5:     $O_t \leftarrow$     SVD (Bt)
6:     Update parameters $\theta_t \leftarrow \theta_{t-1} - \eta O_t$
7: **end for**
8: **return** $\theta_t$

# Run 4:

Slower per step but faster convergence compared to adamW as stated in the paper.

New zero-initialised embedding layer for predicting tokens instead 1 layer doing prediction and representation.
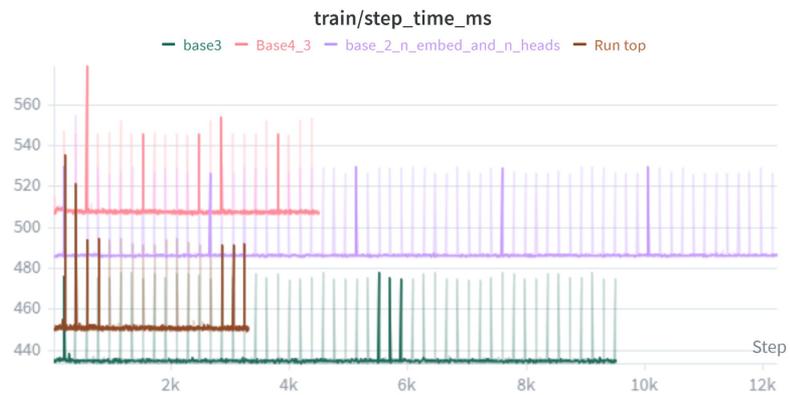
QK layer norm stabilizes training and improves performance, as expected based on paper.



val/loss



train/step_time_ms

# Run 5: Fixes

Fastest run 27 min

Hyper param tuning

# Conclusions