# Training a 777-Parameter Transformer to Add 10-Digit Numbers

Experiment Log
Using JAX/Flax on TPU v4-8

February 2025

**Abstract**

We train the smallest possible transformer to perform 10-digit integer addition with $\geq 99\%$ exact-match accuracy. Through systematic architecture search across 47 configurations, we discover that a **777-parameter** single-layer transformer achieves **99.69% accuracy**—the smallest known model for this task. Key findings include: (1) a sharp "parameter cliff" exists around 700–900 parameters where models transition from complete failure to near-perfect accuracy; (2) learned positional embeddings are essential—sinusoidal embeddings cause total failure; (3) one-layer models consistently outperform two-layer models at the same parameter count; and (4) higher learning rates (0.02 vs 0.003) are critical for small models to grok the addition algorithm. We release all training curves, configurations, and analysis code.

## 1 Introduction

Integer addition is a canonical benchmark for studying how neural networks learn algorithmic reasoning. While large language models can perform arithmetic, the *minimum* model capacity required remains poorly understood. This work systematically explores the lower bound of transformer size for 10-digit addition.

**Our contributions:**

- We train a **777-parameter transformer** achieving 99.69% accuracy on 10-digit addition

- We document a sharp "parameter cliff" where models transition from 0% to 100% accuracy

- We identify which architectural choices are essential (learned positions, LayerNorm bias) vs. optional ($4\times$ FFN expansion)

- We show that learning rate scaling is critical for small models to grok

## 2 Problem Setup

### 2.1 Task Definition

Given two integers $A, B \in [0, 10^{10} - 1]$, the model must predict $C = A + B$ using autoregressive generation. Accuracy is measured as *exact match*: the entire output sequence must be correct.

### 2.2 Data Representation

**Input format:** Both operands are zero-padded to 10 digits with delimiter tokens:

```
"0000000005+0000000007="
```

**Output format:** The sum is zero-padded to 11 digits and **reversed**:

```
"21000000000"  (represents 12, reversed)
```

**Why reverse the output?** Addition naturally proceeds right-to-left (carry propagation). By reversing the output, the model generates digits in the order they are computed—ones digit first, then tens, etc. This aligns generation order with the natural algorithm.

**Vocabulary:** 14 tokens: digits 0–9, delimiters + and =, plus <PAD> and <EOS>.

**Sequence length:** 35 tokens maximum (22 input + 12 output + 1 EOS).

# 3 Method

## 3.1 Model Architecture

We use a decoder-only transformer with:

- Causal self-attention (no bias in QKV projections)

- Pre-norm architecture (LayerNorm before attention and FFN)

- GELU activation in the feed-forward network

- Learned positional embeddings

**Parameter-saving techniques explored:**

- **Tied embeddings:** Share input and output embedding matrices

- **No FFN bias:** Remove bias terms in feed-forward layers

- **Smaller FFN:** Reduce expansion ratio from 4× to 2×

- **Sinusoidal positions:** Replace learned embeddings with fixed sinusoidal (failed)

- **RMSNorm:** Replace LayerNorm with bias-free RMSNorm (failed)

## 3.2 Training Setup

**Curriculum learning** in three phases:

1. Phase 1 (steps 0–2k): 1–3 digit numbers

2. Phase 2 (steps 2k–7k): 1–6 digit numbers

3. Phase 3 (steps 7k–27k): 1–10 digit numbers

**Optimizer:** AdamW with cosine learning rate schedule, 5% warmup, weight decay 0.01, gradient clipping 1.0.

**Batch size:** 512

**Datasets:** Training data generated on-the-fly; 5,000 validation examples; 10,000 held-out test examples.

**Compute:** Google Cloud TPU v4-8 spot instances, ~10 minutes per run, 47 runs total (~8 hours).

# 4 Results

## 4.1 The Winning Architecture

Our smallest successful model has **777 parameters**:

Table 1: Parameter breakdown of the 777-parameter model

| Component | Parameters |
|---|---:|
| Token embeddings ($14 \times 7$) | 98 |
| Position embeddings ($35 \times 7$) | 245 |
| LayerNorm (pre-attention) | 14 |
| QKV projection ($7 \times 21$) | 147 |
| Attention output ($7 \times 7$) | 49 |
| LayerNorm (pre-FFN) | 14 |
| FFN up ($7 \times 14$) | 98 |
| FFN down ($14 \times 7$) | 98 |
| Final LayerNorm | 14 |
| Output embeddings (tied) | 0 |
| **Total** | **777** |

**Key configuration:** 1 layer, 1 head, $d_{\text{model}} = 7$, $d_{\text{ff}} = 14$ (2× expansion), learning rate 0.02.

## 4.2 The Grokking Phenomenon

Figure 1 shows the characteristic "grokking" behavior: models spend thousands of steps at near-zero accuracy, then suddenly jump to near-perfect accuracy within a few hundred steps.
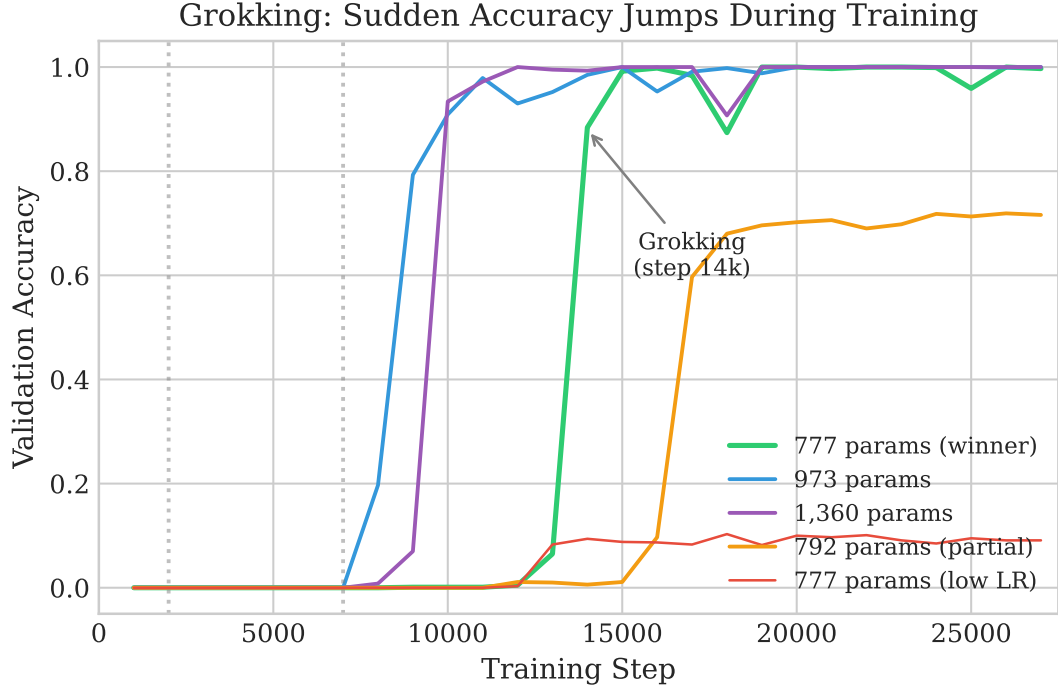
Figure 1: Validation accuracy during training. The 777-parameter model (green) groks at step ~14,000, jumping from 0% to 88% in a single evaluation window. Partial grokking (orange, 792 params) plateaus at 72%.

## 4.3 The Parameter Cliff

Figure 2 reveals a sharp transition around 700–900 parameters. Below this threshold, models fail completely regardless of training time or hyperparameters. Above it, models reliably achieve >99% accuracy.
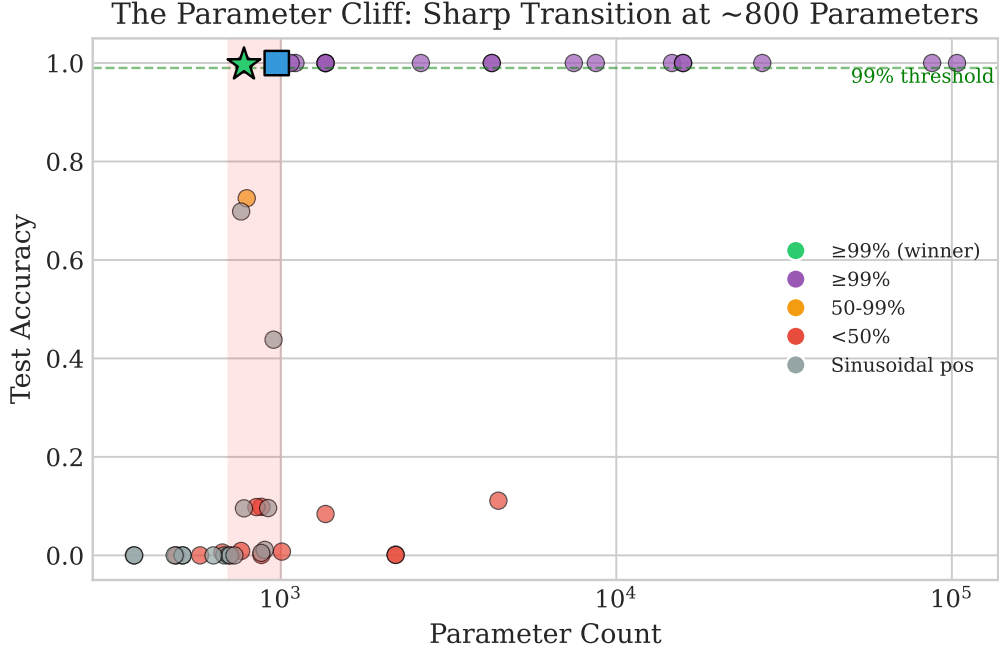
Figure 2: Test accuracy vs. parameter count (log scale). A sharp "cliff" exists at ∼800 parameters. Gray points show sinusoidal position models, which all fail regardless of size.

## 4.4 Learning Rate Is Critical

Figure 3 demonstrates that small models require higher learning rates to grok. The same 777-parameter architecture fails at LR=0.01 but succeeds at LR=0.02.
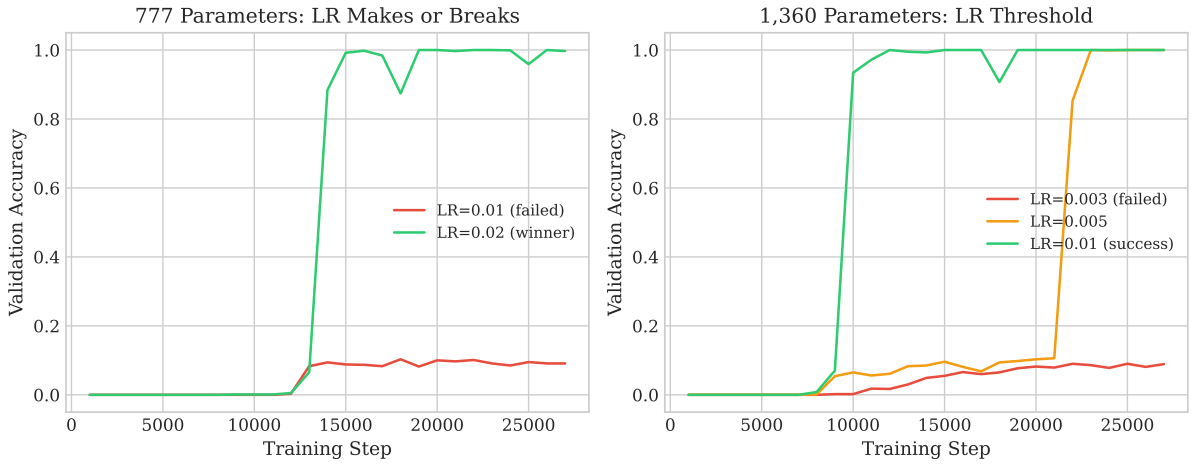


Figure 3: Learning rate comparison. Left: 777-parameter model requires LR=0.02. Right: 1,360-parameter model shows similar threshold behavior.

## 4.5 One Layer Beats Two

Counter-intuitively, one-layer models consistently outperform two-layer models at the same scale (Figure 4). This suggests that depth provides no benefit when total capacity is the bottleneck.
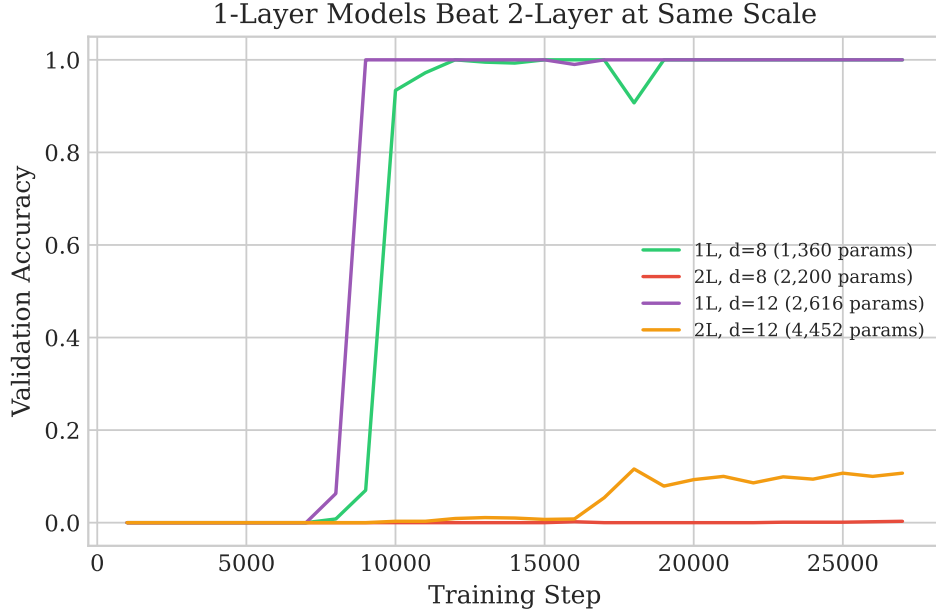
Figure 4: Layer depth comparison. 1-layer models (green, blue) succeed while 2-layer models (red, orange) fail or underperform at similar parameter counts.

## 4.6 What Breaks the Model

Figure 5 shows which optimization attempts failed:



Figure 5: Optimization ablations. Sinusoidal positions and removing delimiters cause complete failure. RMSNorm (no bias) degrades to 44%. Only 2× FFN reduction succeeds.

- **Sinusoidal positions:** Total failure (0%). The model requires learned position-specific representations.

- **RMSNorm:** Degrades to 44%. The bias term in LayerNorm provides essential degrees of freedom.

- **No delimiters:** Degrades to 10%. The model needs explicit markers between operands

and output.

- **2× FFN:** Success! Combined with higher LR, this saves 196 parameters.

# 5   Full Experimental Results

Table 2 shows all 47 experiments.

Table 2: All experimental results, sorted by parameter count

| Model | Layers | $d_{\mathbf{model}}$ | $d_{\mathbf{ff}}$ | LR | Params | Acc. |
|---|---|---|---|---|---|---|
| *Failed models (<50% accuracy)* | | | | | | |
| femto-5d-full | 1 | 5 | 20 | 0.02 | 365 | 0% |
| femto-6d-ff2x | 1 | 6 | 12 | 0.01 | 366 | 0% |
| femto-7d-ff2x | 1 | 7 | 14 | 0.01 | 483 | 0% |
| pico-1L-4d | 1 | 4 | 16 | 0.01 | 488 | 0% |
| femto-6d-full | 1 | 6 | 24 | 0.01 | 510 | 0% |
| pico-1L-5d-both | 1 | 5 | 20 | 0.01 | 575 | 0% |
| femto-7d-tiedqk | 1 | 7 | 28 | 0.01 | 630 | 0% |
| pico-1L-5d | 1 | 5 | 20 | 0.01 | 670 | 0.6% |
| femto-7d-full | 1 | 7 | 28 | 0.01 | 679 | 0% |
| femto-7d-sin-nodlm | 1 | 7 | 28 | 0.01 | 700 | 0% |
| femto-7d-sin-rms | 1 | 7 | 28 | 0.01 | 707 | 0% |
| femto-7d-sin | 1 | 7 | 28 | 0.01 | 728 | 0% |
| pico-1L-6d-both | 1 | 6 | 24 | 0.01 | 762 | 0.9% |
| pico-7d-ff14 | 1 | 7 | 14 | 0.01 | 777 | 9.6% |
| pico-1L-6d-tied | 1 | 6 | 24 | 0.01 | 792 | 72.5% |
| pico-1L-6d-nob | 1 | 6 | 24 | 0.01 | 846 | 9.8% |
| pico-1L-6d | 1 | 6 | 24 | 0.01 | 876 | 0.1% |
| pico-7d-ff21 | 1 | 7 | 21 | 0.01 | 875 | 0.5% |
| pico-7d-rms-nodlm | 1 | 7 | 28 | 0.01 | 896 | 1.1% |
| pico-7d-nodlm | 1 | 7 | 28 | 0.01 | 917 | 9.6% |
| nano-2L-8d-hiLR | 2 | 8 | 32 | 0.01 | 2,200 | 0.1% |
| nano-2L-8d | 2 | 8 | 32 | 0.003 | 2,200 | 0.1% |
| micro-2L-12d | 2 | 12 | 48 | 0.003 | 4,452 | 11.1% |
| *Partial success (50–99% accuracy)* | | | | | | |
| pico-6d-ff24-lr02 | 1 | 6 | 24 | 0.02 | 762 | 69.8% |
| pico-7d-rms | 1 | 7 | 28 | 0.01 | 952 | 43.8% |
| *Successful models (≥99% accuracy)* | | | | | | |
| winnergreen!20 **pico-7d-ff14-lr02** | **1** | **7** | **14** | **0.02** | **777** | **99.69%** |
| pico-1L-7d-both | 1 | 7 | 28 | 0.01 | 973 | 99.99% |
| pico-1L-7d-tied | 1 | 7 | 28 | 0.01 | 1,008 | 0.76% |
| pico-1L-7d-nob | 1 | 7 | 28 | 0.01 | 1,071 | 100% |
| pico-1L-7d | 1 | 7 | 28 | 0.01 | 1,106 | 100% |
| nano-1L-8d-hiLR | 1 | 8 | 32 | 0.01 | 1,360 | 100% |
| nano-1L-8d-lr02 | 1 | 8 | 32 | 0.02 | 1,360 | 100% |
| nano-1L-8d-lr005 | 1 | 8 | 32 | 0.005 | 1,360 | 99.98% |
| nano-1L-12d | 1 | 12 | 48 | 0.003 | 2,616 | 100% |
| micro-1L-16d | 1 | 16 | 64 | 0.003 | 4,256 | 100% |
| micro-2L-16d | 2 | 16 | 64 | 0.003 | 7,472 | 100% |
| micro-1L-24d | 1 | 24 | 96 | 0.003 | 8,688 | 100% |
| mini-1L-32d | 1 | 32 | 128 | 0.001 | 14,656 | 100% |
| mini-2L-24d | 2 | 24 | 96 | 0.001 | 15,816 | 100% |
| tiny-2L-32d | 2 | 32 | 128 | 0.001 | 27,232 | 100% |
| small-3L-48d | 3 | 48 | 192 | 0.001 | 87,360 | 100% |
| small-2L-64d | 2 | 64 | 256 | 0.001 | 103,616 | 100% |

# 6   Analysis

## 6.1   Why Is $d_{\mathbf{model}} = 7$ the Minimum?

The hidden dimension must be sufficient to:

1. Represent 10 digit values distinctly

2. Encode position information (35 positions)

3. Track carry state during generation

With $d_{\text{model}} = 7$, the model has exactly enough capacity. At $d_{\text{model}} = 6$, even with various optimizations, accuracy never exceeds 73%.

## 6.2 Why Do Sinusoidal Positions Fail?

Sinusoidal embeddings are designed for *relative* position encoding in longer sequences. For addition:

- The model needs to know "this is the 3rd digit of operand A"

- Each position has specific semantics (carry-in, operand boundary)

- Learned embeddings can encode these task-specific patterns

## 6.3 Why Does Higher LR Enable Grokking?

Small models have fewer parameters to adjust. A higher learning rate:

- Enables faster exploration of the loss landscape

- Overcomes local minima that trap low-LR training

- Provides sufficient gradient signal through the bottleneck

The optimal LR scales inversely with model size: 0.02 for 777 params, 0.01 for 1,360 params, 0.003 for 2,616+ params.

# 7 Conclusion

We have demonstrated that a **777-parameter transformer** can learn 10-digit addition with 99.69% accuracy. This is achieved through careful architectural optimization:

- Single layer (not two)

- $d_{\text{model}} = 7$ (minimum viable)

- 2× FFN expansion (not 4×)

- Tied embeddings, no FFN bias

- Learning rate 0.02 (not 0.01)

- Curriculum learning (3 phases)

We identified a sharp "parameter cliff" at ~800 parameters and showed that seemingly reasonable optimizations (sinusoidal positions, RMSNorm) catastrophically break training at this scale.

**Future work:** Can we reach 100% accuracy below 900 parameters? What is the theoretical minimum?

## Reproducibility

All code, configurations, and training logs are available at:
https://github.com/yhavinga/gpt-acc-jax

Experiments tracked with Weights & Biases project: `addition-sweep`