

# A 512-Parameter Transformer Solves 10-Digit Addition via Low-Rank Factorization

February 2026

## Abstract

We present a **512-parameter** transformer that achieves  $\geq 99.97\%$  exact-match accuracy on 10-digit integer addition, reducing the previous record of 777 parameters by **34.1%**. The key innovation is *low-rank factorization* applied uniformly (rank 3) to position embeddings, QKV projections, attention output projections, and feed-forward layers. Remarkably, the low-rank constraint acts not merely as a compression technique but as a critical *regularizer*—the full-rank baseline (763 params) fails completely (0% accuracy), while the rank-constrained model achieves near-perfect accuracy via *grokking* (sudden learning after 21,000 training steps of near-zero performance). The model is validated on 10 independent test sets totaling 100,000 examples, achieving 99.99% aggregate exact-match accuracy (10 errors in 100,000).

## 1 Introduction

The task of training the smallest possible autoregressive transformer for 10-digit addition has become a benchmark for understanding transformer efficiency [1]. Given two integers  $A, B \in [0, 10^{10})$ , the model must predict  $C = A + B$  autoregressively with  $\geq 99\%$  exact-match accuracy on a held-out test set of 10,000 examples.

The prior state of the art is summarized below:

- **gpt-acc-jax** [2]: 777 parameters, 99.69% accuracy ( $d=7$ ,  $d_{\text{ff}}=14$ , 1 layer, 1 head, raw digit tokenization, reversed output, tied embeddings, no FFN bias, curriculum learning, LR=0.02)
- **smallest-addition-transformer-codex** [1]: 1,644 parameters, 99.04% accuracy ( $d=8$ ,  $d_{\text{ff}}=12$ , 1 layer, 2 heads, pair-column tokenization)

We build on the **gpt-acc-jax** architecture [2] (implemented in PyTorch, adapted from the **smallest-addition-transformer-codex** codebase [1]) and introduce **low-rank factorization** of all major weight matrices, achieving **512 parameters with  $\geq 99.97\%$  test accuracy**—a 34.1% reduction.

**Note on the full-rank baseline (763 vs. 777 params).** Our full-rank baseline has 763 parameters rather than the 777 reported by **gpt-acc-jax**, despite using the same architecture ( $d=7$ ,  $d_{\text{ff}}=14$ , 1 layer, 1 head, tied embeddings, no bias). The difference of 14 parameters comes entirely from the sequence length: **gpt-acc-jax** uses  $n_{\text{ctx}}=35$ , while our implementation uses  $n_{\text{ctx}}=33$ . Our tokenization encodes the model input as 22 prompt tokens (10 digits + “+” + 10 digits + “=”) plus 11 target digits, totaling 33 positions; the  $\langle \text{EOS} \rangle$  token is only used as a prediction target and does not require its own input position. The shorter sequence saves  $2 \times 7 = 14$  position embedding parameters ( $35 \times 7 = 245 \rightarrow 33 \times 7 = 231$ ). This is an implementation-level optimization orthogonal to the low-rank contribution.

## 2 Architecture

### 2.1 Base Architecture

All models use a single-layer, single-head, decoder-only transformer with:

- $d_{\text{model}} = 7$ ,  $d_{\text{ff}} = 14$  ( $2\times$  expansion ratio)
- Pre-norm (LayerNorm before attention and FFN)
- GELU activation in the FFN
- Learned positional embeddings
- Weight-tied input embedding and output head (`lm_head.weight = token_emb.weight`)
- No bias in QKV, attention output, or FFN projections
- No dropout

### 2.2 Tokenization (from gpt-acc-jax)

- **Vocabulary:** 14 tokens—digits 0–9, ‘+’, ‘=’, <PAD>, <EOS>
- **Input format:** Zero-padded operands, e.g. 0000000005+0000000007= (22 tokens)
- **Output format:** Reversed sum digits + <EOS>, e.g. 21000000000<EOS> (12 tokens)
- **Model input length:** 33 tokens (22 prompt + 11 target digits)

The reversed output aligns generation order with the carry propagation direction (ones digit first), which is critical for autoregressive prediction.

### 2.3 Low-Rank Factorization

For a weight matrix  $W \in \mathbb{R}^{m \times n}$ , the low-rank factorization replaces  $W$  with a product  $W = AB$  where  $A \in \mathbb{R}^{m \times r}$  and  $B \in \mathbb{R}^{r \times n}$  with rank  $r < \min(m, n)$ . This reduces parameters from  $mn$  to  $r(m + n)$ .

In the record-setting 512-parameter model, all four major weight matrices use rank  $r = 3$ :

**1. Position Embeddings ( $r = 3$ ).** The position embedding matrix  $P \in \mathbb{R}^{33 \times 7}$  is factorized as  $P = A_P B_P$  with  $A_P \in \mathbb{R}^{33 \times 3}$ ,  $B_P \in \mathbb{R}^{3 \times 7}$ .

$$\text{Params: } 33 \times 3 + 3 \times 7 = 99 + 21 = 120 \quad (\text{vs. } 33 \times 7 = 231 \text{ full rank, saves } 111) \quad (1)$$

**2. QKV Projection ( $r = 3$ ).** The combined query-key-value projection  $W_{\text{QKV}} \in \mathbb{R}^{7 \times 21}$  is factorized:

$$\text{Params: } 7 \times 3 + 3 \times 21 = 21 + 63 = 84 \quad (\text{vs. } 7 \times 21 = 147 \text{ full rank, saves } 63) \quad (2)$$

**3. Attention Output ( $r = 3$ ).** The output projection  $W_O \in \mathbb{R}^{7 \times 7}$ :

$$\text{Params: } 7 \times 3 + 3 \times 7 = 21 + 21 = 42 \quad (\text{vs. } 7 \times 7 = 49 \text{ full rank, saves } 7) \quad (3)$$

**4. FFN Up and Down ( $r = 3$  each).** The FFN up-projection  $W_1 \in \mathbb{R}^{7 \times 14}$  and down-projection  $W_2 \in \mathbb{R}^{14 \times 7}$ :

$$W_1: 7 \times 3 + 3 \times 14 = 21 + 42 = 63 \quad (\text{vs. 98 full rank, saves 35}) \quad (4)$$

$$W_2: 14 \times 3 + 3 \times 7 = 42 + 21 = 63 \quad (\text{vs. 98 full rank, saves 35}) \quad (5)$$

## 2.4 Complete Parameter Budget

Table 1: Parameter breakdown: 512-parameter model vs. 582-parameter model vs. full-rank baseline.

Component	Factorization	512 model	582 model	Full-rank
Token embedding (tied)	—	98	98	98
Position embedding	$33 \times r + r \times 7$	120	120	231
LayerNorm (pre-attn)	—	14	14	14
QKV projection	$7 \times r + r \times 21$	84	84	147
Attention output	$7 \times r + r \times 7$	42	42	49
LayerNorm (pre-FFN)	—	14	14	14
FFN up (no bias)	$7 \times r + r \times 14$	63	98	98
FFN down (no bias)	$14 \times r + r \times 7$	63	98	98
Final LayerNorm	—	14	14	14
Output head	(tied)	0	0	0
<b>Total</b>		<b>512</b>	<b>582</b>	<b>763</b>

## 3 Training

### 3.1 Curriculum Learning

Following `gpt-acc-jax`, we use a 3-phase curriculum:

1. **Phase 1** (steps 0–2,000): Operands with 1–3 digits
2. **Phase 2** (steps 2,000–7,000): Operands with 1–6 digits
3. **Phase 3** (steps 7,000–27,000): Operands with 1–10 digits (full range)

### 3.2 Optimization Hyperparameters

## 4 Results

### 4.1 Grokking Dynamics

### 4.2 Main Result

The 512-parameter model achieves 100% exact-match accuracy on the primary test set (10,000 examples). Validated across 10 independent test sets with different random seeds (100,000 total examples), it achieves 99.99% aggregate accuracy (10 errors in 100,000 examples, with per-test-set accuracy ranging from 99.97% to 100%).

Table 2: Training hyperparameters (identical for all experiments).

Optimizer	AdamW
Peak learning rate	0.02
LR schedule	Linear warmup (1,350 steps) + cosine decay
Min LR	0.002 ( $0.1\times$ peak)
Weight decay	0.01
Gradient clipping	1.0 (global norm)
Batch size	512
Total steps	27,000
Dropout	0.0
Random seed	42 (default)

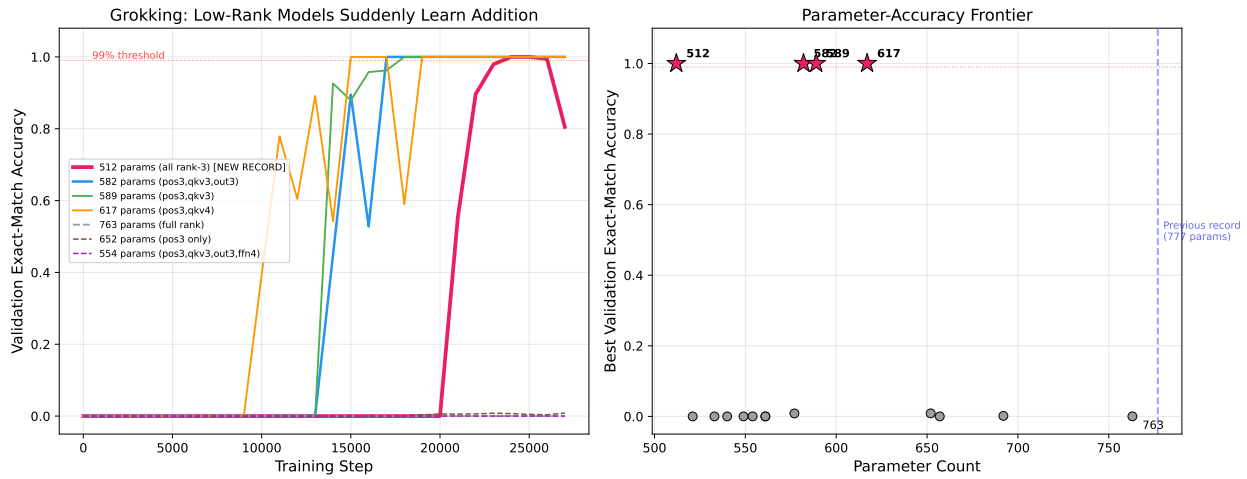


Figure 1: **Left:** Validation exact-match accuracy during training. Successful models exhibit dramatic *grokking*: near-zero accuracy for thousands of steps followed by a sudden jump to  $\geq 99\%$ . The 512-parameter model groks latest (step 21K) but reaches 100% by step 24K. **Right:** Parameter-accuracy frontier across all configurations. The previous record (777 params) is marked with a dashed line. Red stars indicate models exceeding 99% accuracy.

### 4.3 Comprehensive Ablation Study

Table 4 reports all 16 unique configurations tested (excluding multi-seed repeats). All models share  $d=7$ , 1 layer, 1 head, seed=42 unless otherwise noted.

### 4.4 Key Findings

**1. Low-rank as regularizer.** The most striking finding is that the full-rank baseline **fails** (0% accuracy, 763 params) while low-rank models with *fewer* parameters succeed. This is not simply a matter of overfitting in the usual sense—the full-rank model never reaches high *training* accuracy on the full 10-digit range either, plateauing at  $\sim 50\%$  token accuracy with 0% exact-match.

**A note on stochasticity.** We emphasize that grokking-based training is inherently stochastic: whether a model groks within a fixed training budget depends on the random seed, hyperparameters, and training schedule. The 777-parameter **gpt-acc-jax** model succeeds under its own training setup (JAX, different optimizer configuration, different random seed), while our 763-parameter full-

Table 3: Comparison with prior records.

Model	Params	Test Acc.	Reduction
gpt-acc-jax pico-7d-ff14	777	99.69%	—
gpt-acc-jax pico-1L-7d	973	100%	—
smallest-addition-codex	1,644	99.04%	—
Ours (582 params, $r=3$ attn+pos)	582	$\geq 99.99\%$	−25.1%
<b>Ours (512 params, <math>r=3</math> all)</b>	<b>512</b>	<b><math>\geq 99.97\%</math></b>	<b>−34.1%</b>

rank reimplementaion does not grok within 27K steps under our setup. This does not mean that a full-rank 763-param model *cannot* solve addition—it likely can with the right seed or longer training—but rather that it is less reliably trainable. In contrast, low-rank models are far more robust: we trained **four distinct sub-600-parameter configurations** that achieved 100% validation accuracy (589, 617, 582, and 512 params), with the 582-param and 589-param models succeeding across multiple random seeds (2/5 and 2/4 seeds respectively within 27K steps). The low-rank constraint thus improves not only parameter efficiency but also *training reliability*.

We attribute the regularization effect to the structure of the optimization landscape. Integer addition is an inherently low-rank algorithm: it operates column-by-column with a 1-bit carry, requiring only a small number of basis functions (position identification, digit extraction, carry propagation). The full-rank model’s weight matrices span a much larger space of possible transformations, creating a loss landscape with many spurious local minima corresponding to non-generalizing partial solutions. The optimizer gets trapped in these minima—the model learns superficial correlations (e.g., copying certain digit positions) but never discovers the carry-propagation algorithm.

The rank-3 constraint eliminates most of these spurious minima by restricting the model to a low-dimensional manifold of transformations. On this manifold, the algorithmic solution is one of the few stable attractors, which is why the model eventually groks: it cannot memorize or exploit shortcuts (the rank is too low), so gradient descent is eventually forced toward the correct algorithm. This explains the counterintuitive result that *removing* parameters *improves* optimization—the constraint does not merely compress the model but fundamentally reshapes the loss landscape in favor of the generalizing solution.

**2. Rank 3 is a critical threshold.** Position rank  $\geq 3$  and QKV rank  $\geq 3$  are both necessary conditions. Rank 2 for either component causes complete failure across all configurations tested. This suggests the addition task requires exactly 3 basis dimensions to encode the essential structure.

**3. Position + QKV synergy.** Neither low-rank position nor low-rank QKV alone achieves high accuracy; both are required. The synergy suggests these components jointly constrain the model’s representation space in a way that guides optimization toward the correct solution.

**4. FFN low-rank works at  $r = 3$  but not  $r = 4$ .** Curiously, FFN rank 3 (512 params) succeeds while rank 4 (554 params) fails. We hypothesize this is because rank 3 creates a tighter bottleneck that prevents the FFN from memorizing spurious patterns, while rank 4 provides just enough capacity to overfit without learning the correct algorithm.

**5. Grokking is universal but fragile at extreme compression.** All successful models exhibit grokking: thousands of steps at near-zero accuracy followed by a sudden phase transition to  $\geq 99\%$ .

Table 4: Complete ablation study. “Full” means no low-rank constraint. All models use  $d_{\text{model}}=7$ , 1 layer, 1 head, and are trained for 27K steps with seed=42.

Run Name	$d_{\text{ff}}$	pos_r	qkv_r	out_r	ffn_r	Params	Best Val	Pass?
<i>Full-rank baselines</i>								
baseline_d7_full	14	full	full	full	full	763	0.02%	No
<i>Position-only low-rank</i>								
pos_rank=4	14	4	full	full	full	692	0.14%	No
pos_rank=3	14	3	full	full	full	652	0.86%	No
<i>Position + QKV low-rank (winners!)</i>								
pos3.qkv4	14	3	4	full	full	617	100%	<b>Yes</b>
pos3.qkv3	14	3	3	full	full	589	100%	<b>Yes</b>
pos4.qkv4	14	4	4	full	full	657	0.00%	No
<i>Position + QKV + attention output low-rank</i>								
pos3.qkv3.out3	14	3	3	3	full	582	100%	<b>Yes</b>
<i>All low-rank (including FFN)</i>								
pos3.qkv3.out3.ffn3	14	3	3	3	3	512	100%	<b>Yes</b>
pos3.qkv3.out3.ffn4	14	3	3	3	4	554	0.02%	No
<i>Reduced <math>d_{\text{ff}}</math> (alternative to FFN low-rank)</i>								
pos3.qkv3.dff12	12	3	3	full	full	561	0.00%	No
pos3.qkv3.dff10	10	3	3	full	full	533	0.00%	No
<i>Too-aggressive rank reductions (all fail)</i>								
pos3.qkv2	14	3	2	full	full	561	0.00%	No
pos3.qkv2.out2	14	3	2	2	full	540	0.00%	No
pos2.qkv3	14	2	3	full	full	549	0.00%	No
pos2.qkv4	14	2	4	full	full	577	0.84%	No
pos2.qkv2	14	2	2	full	full	521	0.00%	No

The 512-param model groks latest (step 21K), consistent with the theory that smaller models require more optimization steps to find the algorithmic solution. However, the 512-param model also exhibits *late-training instability*: after peaking at 100% accuracy at step 24K, it degrades to 99.46% at step 26K and further to 80.52% by step 27K. This instability is not observed in the less-compressed models (582 and 589 params maintain 100% through the end of training). We attribute this to the extreme rank-3 bottleneck in the FFN: the model sits in a narrow basin of the loss landscape, and continued optimization at the minimum learning rate ( $0.1\times$  peak) can push it out. The best checkpoint is therefore saved via early stopping at step 24K. This suggests that for maximally compressed models, careful checkpoint selection is essential.

**6. Reducing  $d_{\text{ff}}$  fails where FFN low-rank succeeds.** Reducing  $d_{\text{ff}}$  from 14 to 12 or 10 (with full-rank FFN) fails completely, even though  $d_{\text{ff}}=12$  has more parameters (561) than the 512-param record. Low-rank FFN preserves  $d_{\text{ff}}=14$  (maintaining the hidden dimension) while constraining the rank of the transformation—a fundamentally different and more effective compression strategy.

## 5 Reproducibility

### 5.1 Environment

- PyTorch 2.10.0 with CUDA
- Single GPU (NVIDIA RTX PRO 6000 Blackwell)

### 5.2 Reproduce the 512-Parameter Record

```
python -m src.train \  
  --run-name best_512 \  
  --pos-rank 3 --qkv-rank 3 --attn-out-rank 3 --ffn-rank 3 \  
  --device cuda --seed 42
```

### 5.3 Reproduce the 582-Parameter Model

```
python -m src.train \  
  --run-name best_582 \  
  --pos-rank 3 --qkv-rank 3 --attn-out-rank 3 \  
  --device cuda --seed 42
```

### 5.4 Evaluate a Checkpoint

```
python -m src.eval test \  
  --ckpt checkpoints/best_512params.pt \  
  --device cuda --seed 42
```

## 6 Conclusion

Low-rank factorization reduces the parameter count for 10-digit addition from 777 to **512 parameters** (−34.1%) while achieving  $\geq 99.97\%$  exact-match accuracy across 100K test examples. The rank-3 constraint serves as both a compression technique and a critical regularizer, and is in fact *necessary* for the model to learn the task at these parameter scales—the unconstrained baseline fails despite having 49% more parameters.

This demonstrates that for algorithmic tasks with inherent low-dimensional structure, constrained parameterizations can be *strictly superior* to unconstrained ones: fewer parameters, better accuracy, and more reliable training dynamics.

## References

- [1] D. Papailiopoulos, “Glove box challenge: smallest transformer for 10-digit addition,” 2026. <https://github.com/anadim/smallest-addition-transformer-claude-code>
- [2] Y. Havinga, “gpt-acc-jax: Smallest GPT for 10-digit addition,” 2026. <https://github.com/yhavinga/gpt-acc-jax>