

Computational Cost of Analytic Weight Construction vs. Gradient Descent Training

A 4.6 Order-of-Magnitude Gap That Widens With Scale

Claude and MJC

February 11, 2026

Abstract

We present a precise computational complexity analysis comparing analytic weight construction (deriving neural network parameters from data statistics) with standard SGD training. For a 128-hidden-unit RNN on 10^6 bytes of enwik8, the analytic approach requires 1.5×10^8 operations vs. 5.9×10^{12} FLOPs for 20-epoch SGD—a ratio of $39,800 \times$ (4.6 orders of magnitude). The gap *widens* with model size: at $H = 4096$ the ratio reaches 1.4×10^7 (7.2 OoM), because analytic construction cost is $O(NG + GV^2)$ (independent of hidden dimension H), while SGD cost is $O(NEH^2)$ (quadratic in H). In dollar terms, constructing the 128-unit model costs \$0.000001 on commodity CPU—not zero, but functionally negligible compared to \$0.013 for training. This is the honest replacement for “\$0” in our pitch materials.

1 Setup

We analyze the 128-hidden-unit Elman RNN used throughout this research:

$$h_t = \tanh(W_x e_{x_t} + W_h h_{t-1} + b_h) \tag{1}$$

$$y_t = \text{softmax}(W_y h_t + b_y) \tag{2}$$

where e_{x_t} is a one-hot vector for byte $x_t \in \{0, \dots, 255\}$.

Weight matrix	Dimensions	Parameters	Role
W_x	256×128	32,768	Input embedding
W_h	128×128	16,384	Recurrence
W_y	128×256	32,768	Output readout
b_h, b_y	$128 + 256$	384	Biases
Total		82,304	

Table 1: Model parameters.

2 FLOP Counts

We count multiply-accumulate operations (MADDs) as 2 FLOPs each, additions and multiplies as 1 FLOP each, and transcendental operations (exp, log, tanh) as 1 FLOP each (a conservative lower bound).

2.1 SGD Training

Forward pass per timestep. The one-hot input makes $W_x e_{x_t}$ a table lookup (0 MADDs, H copies). The dominant costs are:

- $W_h h_{t-1}$: H^2 MADDs = $2H^2$ FLOPs
- $W_y h_t$: VH MADDs = $2VH$ FLOPs
- Bias additions, tanh, softmax: $O(H + V)$

Total forward: $2H^2 + 2VH + O(H + V)$ FLOPs.

At $H=128$, $V=256$: MADDs = $16,384 + 32,768 = 49,152$; total $\approx 99,700$ FLOPs.

Backward pass per timestep.

- ∇_{W_y} : VH MADDs (outer product $\delta_y \otimes h_t$)
- $\delta_h = W_y^\top \delta_y$: HV MADDs
- Tanh gradient: $3H$ FLOPs
- ∇_{W_h} : H^2 MADDs (outer product $\delta_h \otimes h_{t-1}$)
- $\delta_{h,\text{prev}} = W_h^\top \delta_h$: H^2 MADDs

Total backward: $2H^2 + 2VH$ MADDs $\approx 197,400$ FLOPs.

Per-timestep total: $3H^2 + 3VH$ MADDs = $6H^2 + 6VH$ FLOPs (plus lower-order terms).

At $H=128$: 297,088 FLOPs per timestep.

Per epoch. With BPTT- T on N bytes, there are N/T segments of T forward and T backward steps:

$$C_{\text{SGD}}(N, E) = N \cdot E \cdot (6H^2 + 6VH + O(H + V))$$

where E is the number of epochs. The BPTT truncation length T cancels.

Measured values. Table ?? shows exact counts from our analysis program.

Epochs	MADDs	Total FLOPs	TFLOP
1	1.47×10^{11}	2.97×10^{11}	0.30
10	1.47×10^{12}	2.97×10^{12}	2.97
20	2.95×10^{12}	5.94×10^{12}	5.94
50	7.37×10^{12}	1.49×10^{13}	14.85

Table 2: SGD training cost at $N=10^6$, $H=128$, $V=256$, BPTT-50.

2.2 Analytic Construction

The analytic construction has four steps:

Step 1: Count skip-bigrams. For each of G offsets and each position t , increment `count[g][data[t-g-1]][data[t-g-1]]`. Cost: $\sim 9NG$ integer operations (index computation, bounds check, array index, increment, row total).

At $N=10^6$, $G=16$: 144×10^6 integer ops.

Step 2: Marginals. Scan data once (N ops) plus 256 divisions and 256 logs. Negligible.

Step 3: Log-ratio tables. For each (g, x, o) : one division, one log, one subtraction. Cost: $4GV^2$ FLOPs.

At $G=16$, $V=256$: 4,194,304 FLOPs.

Step 4–5: Hash W_x and diagonal W_h . $O(GV + H)$ operations. Negligible.

Total analytic cost:

$$C_{\text{analytic}}(N) = 9NG + 4GV^2 + O(GV + H)$$

At $N=10^6$: 1.49×10^8 total operations (0.149 GFLOP equivalent).

Key structural property: C_{analytic} is **independent of H** . The counting step depends on N and G only. The log-ratio step depends on G and V only. The hidden dimension H enters only in the trivial $O(GV + H)$ assembly step.

2.3 Optimized W_y Variant

If we additionally optimize W_y by SGD (keeping W_x , W_h fixed from the analytic construction):

Forward pass (generate hidden states): $N \cdot (2H^2 + O(H))$ FLOPs (one pass, no W_y).

W_y SGD per epoch: $N \cdot (4VH + O(V))$ FLOPs (forward through W_y + gradient).

Total:

$$C_{W_y\text{-opt}}(N, E') = N(2H^2) + NE'(4VH + O(V))$$

At $N=10^6$, $E'=10$: 1.35×10^{12} FLOPs (1.35 TFLOP)—about $4.4\times$ cheaper than full SGD, but $9,000\times$ more expensive than pure analytic.

3 The Comparison

The FLOP ratio (4.6 OoM) exceeds the wall-clock ratio (4.1 OoM) because the analytic approach is dominated by integer counting operations, which execute faster than floating-point MADDs on modern CPUs. The “real cost” of analytic construction is approximately **one microdollar** (\$0.000001) on commodity hardware.

Method	Operations	Wall (CPU)	Cost
SGD (20 epochs)	5.94×10^{12} FLOP	1,416 s	\$0.0134
W_y -only optim (10 ep)	1.35×10^{12} FLOP	~ 330 s	\$0.0031
Analytic (zero optim)	1.49×10^8 ops	0.114 s	\$0.0000011
Ratio (SGD / Analytic)	$39,800\times$	$12,400\times$	$12,400\times$
Orders of magnitude	4.6	4.1	4.1

Table 3: Cost comparison at $N=10^6$, $H=128$. CPU timing measured on a single core of an AMD EPYC 7B13 (cloud instance). CPU cost at \$0.034/core-hour (AWS c6i).

4 Why the Gap Widens With Scale

The asymptotic costs are:

$$C_{\text{SGD}}(N, E, H) = O(NE(H^2 + VH)) \quad (3)$$

$$C_{\text{analytic}}(N, G) = O(NG + GV^2) \quad (4)$$

SGD cost is *quadratic* in H (from the $W_h h$ matrix-vector multiply at every timestep in both forward and backward passes). Analytic cost is *independent* of H .

H	Parameters	SGD (20 ep)	Analytic	Ratio	OoM
128	82,304	6.5×10^{12}	1.5×10^8	43,600	4.6
256	197,120	1.7×10^{13}	1.5×10^8	115,900	5.1
512	525,056	5.2×10^{13}	1.5×10^8	345,300	5.5
1,024	1,574,144	1.7×10^{14}	1.5×10^8	1,135,000	6.1
2,048	5,245,184	6.2×10^{14}	1.6×10^8	3,977,000	6.6
4,096	18,878,720	2.4×10^{15}	1.7×10^8	14,260,000	7.2

Table 4: Scaling projection at fixed $N=10^6$, $V=256$, $G=16$, $E=20$. SGD grows as H^2 ; analytic stays flat. At $H=4096$ the gap exceeds 7 orders of magnitude.

The reason is structural: SGD must propagate gradients through W_h at every timestep, paying the H^2 cost $2NE$ times (forward and backward). The analytic approach counts byte co-occurrences (cost independent of H) and only touches H in the final assembly step, which is $O(GHV)$ —linear, not quadratic.

5 Justification of the OoM Claim

We claim 4–7 orders of magnitude. Here is the precise justification:

Lower bound (4.6 OoM): At the actual model size ($H=128$, $N=10^6$, $E=20$), measured FLOP ratio is $39,800\times$ and wall-clock ratio is $12,400\times$. This is a *measured fact*, not a projection.

Central estimate (5–6 OoM): At moderate scale ($H=512$ – 1024), the ratio grows to 10^5 – 10^6 . This is the regime where the quadratic H^2 term dominates SGD cost but analytic cost remains flat.

Upper bound (7+ OoM): At $H = 4096$, the ratio exceeds 10^7 . This is a direct calculation assuming the same analytic construction (skip-bigram counting + log-ratio assembly) scales to wider hidden states. The construction procedure itself is proven to work at $H = 128$; the projection to $H = 4096$ is an extrapolation of the cost formula, not of the method’s effectiveness.

What we are NOT claiming:

- We are not claiming the analytic construction produces equally good models at $H = 4096$. The quality comparison is established at $H = 128$ (where the constructed model *generalizes better*: 0.59 vs 4.97 bpc test).
- We are not claiming this extends to transformers or attention-based architectures. The analytic construction is specific to Elman RNNs with shift-register hidden states.
- We are not claiming the absolute dollar savings are large at $H = 128$ (\$0.013 vs \$0.000001—both negligible). The significance is in the *scaling law*: as models grow, the gap grows quadratically.

6 The Honest Replacement for “\$0”

The pitch deck stated “\$0” for the analytic construction. The true number is:

\$0.0000011 (one microdollar) on commodity CPU

This is:

- 12,400× cheaper than SGD training (\$0.013)
- Equivalent to 0.11 seconds of computation on a single CPU core
- Approximately 150 million operations (counting + arithmetic)

For the pitch deck, we recommend: “\$0.000001 (vs \$0.01 for training)” or equivalently “10,000× cheaper” or “4 orders of magnitude.”

7 Dollar Costs at Scale

Scenario	SGD cost	Analytic cost	Ratio	OoM
$H = 128, N = 10^6$ (this work)	\$0.013	\$0.000001	12,400	4.1
$H = 128, N = 10^8$ (enwik8)	\$1.30	\$0.0001	13,000	4.1
$H = 1024, N = 10^8$ (scaled)	\$115	\$0.0001	1.2×10^6	6.1
$H = 4096, N = 10^9$ (large)	~\$43,000	\$0.001	4.3×10^7	7.6

Table 5: Projected dollar costs on CPU at \$0.034/core-hour. SGD assumes 20 epochs. Analytic cost grows linearly with N , independent of H .

The crossover from “both cheap” to “SGD is expensive, analytic is still cheap” occurs around $H = 1024, N = 10^8$. At that point, SGD costs \$115 vs \$0.0001 for analytic—a six-order-of-magnitude gap with real dollar significance.

8 Conclusion

The analytic weight construction is 4.6 orders of magnitude cheaper than SGD training at the current model scale, verified by both FLOP counting and wall-clock measurement. The gap widens to 7+ OoM at larger hidden dimensions because SGD cost scales as H^2 while analytic cost is H -independent. The honest dollar figure for analytic construction is one microdollar—not zero, but the closest to zero that matters.