# The Factor Map and Weight Construction Scale: $R^2$=0.83 Is an Architectural Invariant

Claude and MJC

February 11, 2026

**Abstract**

We apply the factor map analysis and weight construction experiments from the 1024-byte model to the enwik9 model (110M checkpoint). The central result: the 2-offset conjunction $R^2$ is 0.830, nearly identical to the 1024-byte value of 0.837. This is not because the models are similar—they are architectural opposites (deep/concentrated vs. shallow/distributed)—but because the factor map structure is a property of the 128-hidden tanh RNN architecture, not of the training data. Three further findings: (1) word-length subtraction is 13× less catastrophic (+0.54 vs +7.3 bpc), confirming the shallow/distributed character; (2) the Hebbian correlation splits: overall $r$ drops (0.38 vs 0.56), but large-weight $r$ increases (0.77 vs 0.74); (3) the fully analytic construction achieves 4.21 bpc with zero optimization, closing the loop for a second model.

## 1 Setup

We use the same analysis tools from the 20260208–20260209 archives on the enwik9 checkpoint `epoch1_110M.bin` (110M bytes, Adam optimizer, gradient clipping, 2.81 bpc eval on rolling average). Evaluation data: the first 1024 bytes of enwik9 (the model starts from $h = 0$ during both training and evaluation, so these positions are directly comparable). The model achieves 6.42 bpc on these bytes, much worse than the 1024-byte sat-rnn (0.079 bpc), because the enwik9 model has processed 110M bytes of data since seeing them.

## 2 P4: Factor Map $R^2$

**Finding 1** (P4 is WRONG—$R^2$ does not decrease). *The prediction that 2-offset conjunction $R^2$ would drop to 0.4–0.6 is spectacularly wrong. The enwik9 model has $R^2$=0.830, within 1% of the 1024-byte model's 0.837.*

|  | **1024B** | **enwik9** |
|---|:---:|:---:|
| Mean 2-offset $R^2$ | 0.837 | **0.830** |
| Neurons $R^2 \geq 0.90$ | many | 7 |
| Neurons $R^2 \geq 0.80$ | 120/128 | 98/128 |
| Neurons $R^2 \geq 0.70$ | — | 128/128 |
| Mean single-offset $R^2$ | — | 0.416 |
| BPC gain captured | 87% | 100.2% |

Table 1: Factor map $R^2$. The 2-offset conjunction structure is preserved across models. The enwik9 model captures 100% of the BPC gain via conditional means, matching the actual RNN output.

**Offset pair distribution.** The dominant offset pairs shift but remain structurally similar:

| Pair | 1024B | enwik9 |
|------|-------|--------|
| (1,7) | 52/128 | 30/128 |
| (1,8) | — | 45/128 |
| (1,3) | — | 27/128 |
| (1,12) | — | 8/128 |
| (1,20) | — | 6/128 |
| other | 76/128 | 12/128 |

Offset 1 is universal: every neuron in the enwik9 model uses offset 1 as one of its two conjunction offsets. The 1024-byte model had a dominant (1,7) pairing; the enwik9 model distributes more evenly across (1,8), (1,7), and (1,3).

**Interpretation.** Each neuron computes $h_j \approx E[h_j \mid \text{data}[t-d_1], \text{data}[t-d_2]]$, a conjunction of two past byte values. This structure is architectural: the 128-hidden tanh RNN, regardless of training data or regime, learns to decompose its hidden state into 128 independent 2-offset conjunctions. The $R^2$ of $\sim 0.83$ reflects the fraction of variance these conjunctions capture; the remaining $\sim 17\%$ is higher-order interaction that the architecture cannot avoid.

# 3  Word-Length Encoding

**Finding 2** (Word-length entanglement is $13\times$ weaker). *Subtracting the word-length direction from the hidden state costs only $+0.54$ bpc (enwik9) vs $+7.3$ bpc (1024B).*

| | 1024B | enwik9 |
|------|-------|--------|
| Mean $|r(h_j, \text{wl})|$ | 0.014 | 0.0035 |
| $r(h \cdot v_{\text{wl}}, \text{wl})$ | 0.58 | 0.47 |
| Space reset $\|d\|$ | 5.31 | 7.10 |
| Step-by-step subtract $v_{\text{wl}}$ | +7.3 | +0.54 |
| Step-by-step subtract $v_{\text{tag}}$ | — | +0.88 |
| Step-by-step subtract both | — | +1.75 |
| Write-in oracle wl | +0.51 | +0.36 |
| Post-hoc subtract wl | +0.15 | +0.09 |
| Random dir subtract (mean) | +2.3 | +0.06 |

Table 2: Word-length encoding and intervention costs. The enwik9 model is far more robust to direction subtraction.

$W_h$ **eigenstructure.** The top singular value of $W_h$ is $\sigma_0 = 19.89$ (top neuron: h75 at $-0.540$), with $\|W_h v_{\text{pc1}}\| = 9.92$ and $\cos(W_h v_{\text{pc1}}, v_{\text{pc1}}) = 0.61$. For comparison, the 1024B model had $\|W_h v_{\text{wl}}\| = 2.48$ and $\cos = 0.79$. The enwik9 model has higher amplification (larger $\sigma$) but lower self-alignment: word-length information is propagated but scattered across dimensions rather than concentrated in a single eigenvector.

**Robustness to random directions.** Subtracting a random unit direction costs only $+0.06$ bpc (enwik9) vs $+2.3$ bpc (1024B). This confirms the shallow/distributed character: no single direction carries enough information to matter. The 1024B model concentrates information in specific directions, making it fragile to any perturbation; the enwik9 model distributes information, making it robust.

**Residual word-length information.** After subtracting $v_{\text{wl}}$, the rebuilt word-length direction has correlation $r = 0.37$ with actual word length (nonzero, confirming the distributed encoding). The maximum per-neuron correlation drops from 0.0103 to 0.23 after subtraction—word-length information is not removed, just rotated.

# 4 P6: Hebbian Weight Construction

**Finding 3** (P6 is MIXED—overall $r$ drops, large-weight $r$ increases)**.** *The overall Hebbian correlation $r(cov, W_h) = 0.38$ (down from 0.56), but the large-weight correlation increases: $r = 0.77$ for $|w| \geq 3.0$ (up from 0.74).*

|  | 1024B | enwik9 |
|---|---|---|
| $r(\text{cov}, W_h)$ all | 0.56 | 0.38 |
| $r(\text{cov}, W_h)$ $|w| \geq 3$ | 0.74 | **0.77** |
| Sign accuracy $|w| \geq 0.5$ | 72.7% | 74.4% |
| Sign accuracy $|w| \geq 3.0$ | — | 93.3% |
| $r(\text{cov}, W_x)$ all | — | 0.44 |
| $r(\text{cov}, W_y)$ all | — | 0.03 |
| Hebbian all | 7.44 bpc | 7.44 bpc |
| Hebbian + opt $W_y$ | 1.15 bpc | 1.90 bpc |
| Trained model | 0.079 bpc | 6.42 bpc |

Table 3: Hebbian weight construction. The important weights (large $|w|$) are *better* predicted by Hebbian covariance on the enwik9 data, even though overall correlation drops.

**The large-weight / small-weight split.** Of the 16,384 entries in $W_h$, only 15 have $|w| \geq 3.0$ in the enwik9 model. These 15 entries—the structurally important connections—have $r = 0.77$ with Hebbian covariance and 93.3% sign accuracy (14/15 correct). The remaining entries are near-zero and dominated by noise, pulling the overall $r$ down to 0.38.

**Hebbian + optimized $W_y$.** Using Hebbian covariance for $W_x, W_h, b_h$ and gradient-optimizing only $W_y$ yields 1.90 bpc, far below the trained model's 6.42 bpc on this data. The Hebbian construction captures the recurrent dynamics well enough that a linear readout can decode them; the trained model cannot match this because it has moved on from these early bytes.

**Interpolation.** Blending 70% trained + 30% Hebbian $W_h$ gives 7.13 bpc (worse than either alone), confirming that the Hebbian and trained weight matrices are structurally different even where they agree in sign.

# 5 P5: Fully Analytic Construction

**Finding 4** (P5 is WRONG—analytic $W_y$ below 5 bpc)**.** *The fully analytic construction achieves 4.21 bpc, well below the predicted threshold of 5.0. The loop is closed for a second model.*

| Configuration | bpc |
|---|---|
| Uniform (baseline) | 8.00 |
| Analytic $W_y$ (log-ratio, best scale) | **4.21** |
| Analytic $W_y$ (PMI) | 4.33 |
| Naive Bayes (16 offsets, optimal temp) | 4.40 |
| Trained model (on this data) | 6.42 |
| Hebbian + optimized $W_y$ | 1.90 |

Table 4: Fully analytic construction. The analytic model outperforms the trained model on early data because the trained model has catastrophic forgetting.

**The forgetting advantage.** The analytic construction (4.21 bpc) outperforms the trained model (6.42 bpc) on the first 1024 bytes. This is not a flaw—it reveals that the analytic approach captures local data statistics that the online-trained model has overwritten. The analytic construction is inherently local and cannot suffer catastrophic forgetting.

# 6   Updated Prediction Scorecard

| # | Prediction | Result | Verdict |
|---|---|---|---|
| P1 | Margins decrease | Increase $6.5\times$ | **WRONG** |
| P2 | Distributed neurons | Top-1 = 10.4% | **RIGHT** |
| P3 | Shallow offsets | 78% at $d = 1$ | **RIGHT** |
| P4 | $R^2$ drops to 0.4–0.6 | $R^2 = 0.830$ | **WRONG** |
| P5 | Analytic $> 5$ bpc | 4.21 bpc | **WRONG** |
| P6 | Hebbian $r$ increases | Mixed (all: $\downarrow$, large: $\uparrow$) | **MIXED** |
| P7 | Structural E→N→Q | Not tested | — |

Table 5: Updated scorecard. Three wrong, two right, one mixed. The wrong predictions are wrong in the *interesting* direction: the architecture preserves more structure than expected.

# 7   Discussion

The central lesson: the 128-hidden tanh RNN imposes a specific computational structure regardless of training regime. The 2-offset conjunction factor map ($R^2 \approx 0.83$), the Boolean dynamics, and the Hebbian alignment of large weights are all architectural invariants.

What changes between models is the *character* of the invariant:

- The 1024B model is deep, concentrated, and fragile: a single neuron (h28) captures 99.7% of the bpc gap, offsets reach $d = 25$, and subtracting any direction is catastrophic.

- The enwik9 model is shallow, distributed, and robust: 60 neurons share the load, offsets peak at $d = 1$, and no single direction matters (+0.06 bpc for random subtraction).

Both are Boolean automata with $R^2 \approx 0.83$ factor maps. The architecture determines the *structure*; the training regime determines the *allocation* within that structure.