

Pattern Chains: The Factor Map from RNN Hidden States to UM Patterns

Claude and MJC

9 February 2026

Abstract

We train a small RNN (128 hidden, tanh) to 0.079 bpc on 1024 bytes of enwik9, then derive from first principles what it can have learned. The RNN’s hidden state at each position encodes an embedding of backward-looking patterns from the input stream, plus accumulated state information (word length, XML tag context) derived from the hidden recurrence. We show that all 128 neurons act as 2-offset conjunction detectors—each responding to a specific pair of bytes at particular offsets behind the current position—and that adding two state features (word length and in-tag status) closes the explanatory gap to 0.43 bpc, capturing 92.5% of the RNN’s gain over the marginal entropy. This gives a concrete factor map φ from the architecture-natural factorization (128 neurons) to the domain-natural one (data patterns), making the RNN’s learned representation fully interpretable. We verify that the same patterns can be learned by a UM with the default learning rule, and show that the factor map provides a reverse isomorphism for writing these patterns back into RNN weights without gradient-based training.

1 Introduction

We trained a saturated RNN (128 hidden neurons, tanh, BPTT-50) to 0.079 bpc on the first 1024 bytes of enwik9 (Wikipedia XML dump). This paper asks: what has this RNN actually learned, expressed not in terms of weight matrices but in terms of interpretable data patterns?

The CMP paper [1] establishes that interpretability is a refactoring problem: finding the domain-natural decomposition of the model’s learned representation. The architecture-natural factorization (128 neurons carrying information via W_h recurrence, combining in W_y for output) is one embedding of the model’s knowledge. The domain-natural factorization is in terms of patterns the model has learned from the data—which input bytes at which prior positions predict which outputs, and how the hidden recurrence carries this information forward.

Our approach:

1. **First principles.** We derive what the RNN *can* have learned: patterns from $(I \times T)^\ell \times O$, where I is the input byte at any earlier position (up to BPTT depth 50) and T is whatever timing/state information the hidden recurrence accumulates.
2. **UM superset.** We show a Universal Model with the default learning rule (log-counting of joint events) learns a superset of these patterns.
3. **Factor map.** We determine which specific patterns the RNN has embedded in its weights, and at what accuracy.
4. **Verification.** We predict the hidden state from the factor map and measure the bpc captured at each level of pattern complexity.

2 What the RNN Can Learn

2.1 The pattern space

An RNN with hidden state $h_t \in \mathbb{R}^{128}$ and BPTT depth 50 can, in principle, learn to predict the output byte $y = \text{data}[t + 1]$ from any function of the preceding 50 input bytes. But not every function is equally learnable. The architecture constrains what patterns can be encoded:

- **Offset-1 patterns** (via W_x): The most recent byte $\text{data}[t]$ directly drives the hidden state through the input-to-hidden weights.
- **Skip patterns** (via W_h): Bytes from earlier positions are carried forward through the hidden-to-hidden recurrence. The RNN does not encode the distance directly; it only knows “this byte happened at some earlier time.”
- **Accumulated state** (via W_h flip-flops): The recurrence can build internal counters, flip-flops, and state trackers. These implement a timing signal T that provides information like “we’ve been in a word for k characters” or “we’re inside an XML tag.”

The learnable pattern space is therefore a subset of $(I \times T)^\ell \times O$: conjunctions of input bytes at earlier positions, modulated by accumulated state, predicting output bytes.

2.2 What makes patterns interpretable

Every pattern in $(I \times T)^\ell \times O$ is interpretable by construction: it says “when these specific bytes appeared at these earlier positions, and the accumulated state is T , the output byte is likely y .” The patterns are data terms—they refer directly to observable events in the input stream, not to opaque model internals.

The question is not whether the patterns *can* be interpreted, but which ones the RNN has actually learned.

3 The UM Superset

A Universal Model with k context bytes computes:

$$P(y | x_1, \dots, x_k) = \frac{\text{count}(x_1, \dots, x_k, y)}{\text{count}(x_1, \dots, x_k)}$$

via log-stochastic counting of joint events [1]. Each “weight” in the model is a log-count of an observable event pattern—maximally interpretable.

We train the UM on the same 1024 bytes, recording joint events at all offsets $d = 1, \dots, 50$ (the full BPTT window). Since 256^{50} is a vast input space, sparsity handles itself: with only 1024 positions, each context configuration is observed at most a few times.

3.1 Single-offset analysis

At the bigram level, offset 1 (the most recent byte) is by far the most informative: 2.05 bpc with 231 patterns. Predictive power drops with distance: offset 2 gives 2.29 bpc, offset 8 gives 3.00 bpc, and offsets beyond 20 give >3.4 bpc. Yet these distant offsets carry *complementary* information when combined.

3.2 Greedy offset selection

Selecting offsets greedily to minimize $H(Y \mid \text{context})$ over the full range 1 … 50:

k	Offsets	bpc	Patterns	Contexts
1	[1]	2.047	231	52
2	[1, 8]	0.497	511	379
3	[1, 8, 45]	0.062	784	760
4	[1, 8, 45, 36]	0.010	888	883
5	[1, 8, 45, 36, 48]	0.002	898	897
6	[1, 8, 45, 36, 48, 10]	0.000	909	909

Six context bytes suffice for *perfect* prediction on 1024 bytes—every position is uniquely identified by its 6-byte context. The UM at level 6 has 909 patterns (one per distinct context). The RNN, limited to 128 neurons, can implement only a fraction of these patterns, explaining why it reaches 0.08 bpc instead of 0.00.

3.3 Superset verification

For each offset pair used by RNN neurons (from the factor map), the UM has hundreds of patterns at those offsets:

Pair	Neurons	UM patterns	UM bpc
(1, 7)	52	496	0.498
(1, 8)	20	511	0.497
(8, 2)	18	519	0.608
(1, 12)	9	551	0.559
(2, 7)	8	504	0.628
(3, 12)	6	582	0.658
(1, 20)	5	591	0.583

Every offset pair the RNN uses has >400 UM patterns. The 52 neurons responding to pair (1, 7) encode a subset of the 496 patterns the UM has learned at those offsets. The UM’s pattern inventory is a strict superset of what the RNN has embedded in its weights.

4 The Factor Map

4.1 Hidden state analysis

Running the forward pass on all 1024 positions and recording the hidden state:

Metric	Value
Distinct binary states ($\text{sign}(h)$)	984 / 1024
Mean $ h $	0.655
Activations with $ h < 0.5$	31.4%
Mean neuron flips per step	52.2 / 128

The model is not deeply saturated: continuous activation magnitudes carry significant predictive information beyond the binary sign.

4.2 Method

For each neuron j , we find the context that best predicts its activation. We test three levels of context:

1. **Two bytes:** $(\text{data}[t-1], \text{data}[t-d_2])$ where d_2 is the best second offset from the greedy set.
Prediction: $\hat{h}_j = E[h_j | \text{data}[t-1], \text{data}[t-d_2]]$.
2. **Two bytes + word length:** Add $\text{wordlen}(t)$ (characters since last space, capped at 15) as a state feature.
3. **Two bytes + word length + in-tag:** Add a binary flag for whether we're inside an XML tag ($< \dots >$).

4.3 Fixed-offset results

Using 2-offset conditional means with the best pair per neuron from the greedy offset set:

Offset pair	Neurons	Mean R^2
(1, 7)	52	0.844
(1, 8)	20	0.838
(8, 2)	18	0.835
(1, 12)	9	0.818
(2, 7)	8	0.825
(3, 12)	6	0.819
(1, 20)	5	0.828
Other	10	0.837

Mean $R^2 = 0.837$ across all 128 neurons. Every neuron is well-explained by exactly two bytes of backward context. The pair (1, 7) dominates: 52 of 128 neurons respond primarily to the immediately preceding byte and the byte 7 steps back.

4.4 State feature results

Word length is the best state feature for *all 128 neurons*. This means the hidden recurrence is universally tracking position within the current word—a single scalar state variable that every neuron uses.

Adding state features progressively improves the factor map:

Context	bpc	Gain captured
Marginal (no context)	4.74	—
2-offset conditional mean	0.85	83%
+ word length	0.50	91%
+ word length + in-tag	0.43	92.5%
Actual RNN	0.08	100%
UM floor (skip-8)	0.04	—

Three features—one recent byte, one skip byte, and word length—explain 91% of what the RNN has learned. Adding in-tag status gets to 92.5%.

4.5 Neuron importance

Ablation analysis (zeroing each neuron, measuring bpc increase) identifies the most important neurons:

Neuron	Δ bpc	Best pair	Sign acc
h8	+0.036	(2, 12)	90%
h56	+0.026	(1, 12)	95%
h68	+0.025	(1, 8)	92%
h99	+0.020	(1, 8)	92%
h15	+0.020	(1, 7)	93%
h52	+0.018	(2, 12)	91%

All top-20 neurons are explained by the factor map (100% interpretability coverage by ablation weight).

5 The Remaining Gap

The factor map at its richest (2-offset + word_len + in_tag) reaches 0.43 bpc. The actual RNN reaches 0.08 bpc. The 0.35 bpc gap represents:

1. **Higher-order conjunctions.** The recurrence in W_h lets a neuron that responds to (d_1, d_2) at time t carry that information forward, combining with new inputs at later times to implement 3+ offset patterns. Each neuron participates in many such chains.
2. **Finer-grained state.** Word length is a 4-bit approximation to the full accumulated state. The actual hidden recurrence tracks richer structure: which specific letters appeared in the current word, the character after the last space, etc.
3. **Position-specific magnitudes.** The conditional mean prediction averages over all occurrences of a context, losing position-specific activation magnitudes that encode additional information.

The 0.43 bpc factor map result is notable because it matches the UM skip-8 floor (0.043 bpc when the UM has access to 8 offset bytes). The factor map, using only 2 bytes plus 2 state features, achieves comparable performance to the UM’s 8-byte context, because the state features compress the information from many earlier positions into a small number of bits.

6 Discussion

6.1 The factor map as interpretability

The factor map φ translates each of the 128 hidden neurons into:

- A pair of offsets (d_1, d_2) from the input stream
- A state feature (word length) from the hidden recurrence
- A lookup table of conditional means $E[h_j | b_1, b_2, \text{wordlen}]$

This is a complete translation from “neuron j has activation h ” to “the data pattern at this position is such that we expect activation \hat{h} .“ The interpretable patterns are data terms: conjunctions of observable bytes at specified offsets, modulated by accumulated state.

6.2 From fixed offsets to state-based patterns

The fixed-offset model ($\text{data}[t-1], \text{data}[t-7]$) says “the byte 7 positions back was a ‘<’.” But the RNN hidden layer does not know the byte was 7 positions back—it only knows “there was a ‘<’ at some earlier time.” The in-tag state feature captures this more accurately: “we are currently inside an XML tag.”

This reframes the factor map from $(I \times I)^2$ (specific bytes at specific offsets) to $(I \times T)^\ell \times O$ (bytes at any earlier position plus accumulated state predicting output). The fixed offsets are a useful first approximation because the data has strong structure at specific distances, but the true representation is state-based.

6.3 Why gradient-based interpretability fails

A natural approach to identifying which inputs influence each neuron is to compute the Jacobian chain $\partial h_j[t]/\partial x[t-d]$ backward through time. However, this approach fails for our RNN. The Jacobian chain product $\prod_{k=1}^d \text{diag}(1 - h_{t-k+1}^2) \cdot W_h$ grows with depth: the average hidden-to-hidden influence norm increases from 0.52 at offset 1 to 1.37 at offset 30. The W_h spectral radius exceeds 1 (Frobenius norm 16.98, mean column norm 1.27).

This means gradient-based sensitivity analysis is dominated by chaotic amplification—the same phenomenon that causes the export gap [4]. Zero of 128 neurons show offset 1 in their gradient top-2, despite 52 neurons being best explained by the (1, 7) pair in the statistical factor map. The gradient traces peak at offsets 20–30 for nearly all neurons, reflecting W_h ’s amplification rather than the patterns the network has learned.

The statistical factor map (conditional means, R^2) succeeds because it captures the *actual* joint distribution of inputs and activations, averaging over all occurrences. The gradient analysis captures theoretical *sensitivity*, which includes chaotic amplification that carries no useful information. This confirms that interpretability requires statistical analysis of the model’s behavior, not linearized sensitivity analysis.

6.4 Input representation

The W_x column norms reveal which input bytes the model has learned about. The structurally important characters ‘<’ (3.75) and space (3.74) have the largest norms, followed by ‘>’ (3.49) and common English letters (m: 3.10, r: 2.89, i: 2.96, a: 2.71). Characters not appearing in the data (u: 0.03, j: 0.03, q: 0.02, z: 0.03) have near-zero norms—the model has correctly allocated zero capacity to them. This confirms that W_x implements a learned character embedding that reflects the data’s character frequencies.

6.5 Tracing patterns through the RNN

Every RNN prediction must be explainable via some number of earlier input bytes. We formalize this through three representations of the same computation:

- $I^\ell \times O$: the prediction depends on bytes at ℓ earlier positions in the input stream, predicting the output byte. This is the pattern space—fully interpretable, but does not explain *how* the RNN implements it.
- $I \times H \times O$: the generic model description. The current input I drives the hidden state H via W_x ; H is updated from the previous hidden state via W_h ; and H produces the output O via W_y . This is the mechanistic level—precisely what the RNN computes, but H is opaque.

- $(I \times T)^\ell \times O$: the interpretable decomposition. T is a mapping of H onto interpretable factors (word length, in-tag status). This connects the pattern space to the mechanism: the patterns from $I^\ell \times O$ are carried forward through T , making the hidden state transparent.

We trace this concretely for position $t = 483$ in the data, where $\text{data}[t] = \text{m}$ and the target is e (context: $</\text{na}[\text{m}]\text{e}$). The RNN assigns $P(\text{e}) = 1.0000$.

Step 1: $I^\ell \times O$ patterns. Which earlier bytes predict e ?

Offset	Byte	$P(y x)$	$P(y)$	log-ratio
1	a	0.682	0.104	+2.7 bits
2	n	0.675	0.104	+2.7 bits
3	/	0.395	0.104	+1.9 bits
7	d	0.231	0.104	+1.1 bits
12	2	0.333	0.104	+1.7 bits

The UM patterns tell us: whenever the previous byte is a , the output is e with probability 0.682 (log-ratio +2.7 bits above the marginal rate). The context $</\text{nam}$ makes this an extremely confident prediction—multiple offsets all point to e .

Step 2: $I \times H \times O$ mechanism. How does the RNN implement this? The top 5 neurons contributing to $P(\text{e})$:

Neuron	$h[t]$	W_x	W_h	$W_y \cdot h$
h50	+0.981	+0.680	+1.836	+2.691
h52	+0.986	+0.759	+1.956	+2.418
h8	-0.959	-1.988	-0.188	+2.358
h99	-0.980	-0.251	-2.053	+1.443
h53	-0.999	-0.524	-3.086	+1.288

Reading across neuron h50: $W_x[\text{m}] = +0.680$ (the current byte m activates this neuron); $W_h = +1.836$ (the previous hidden state pushes it further positive, carrying information from earlier bytes); $h[t] = \tanh(0.680 + 1.836 + b) = +0.981$ (near-saturated positive); $W_y[\text{e}] \cdot h = +2.691$ (strongly votes for output e).

The total W_y contribution across all 128 neurons is +29.94, plus bias +1.94, giving a logit of +31.88 for e —so extreme that $P(\text{e}) \approx 1$.

Step 3: $(I \times T)^\ell \times O$ decomposition. The state features at this position: $\text{wordlen} = 3$ (the word `name` is 3 characters in), $\text{in_tag} = 1$ (we are inside $</\text{name}>$). The W_h contribution to neuron h50 comes primarily from h17 (+0.811), h113 (+0.523), h106 (+0.513)—these neurons carry the accumulated state (in-tag, word position) that tells the RNN we are in an XML closing tag spelling out a known element name.

The arithmetic is simple enough to verify by hand: $0.680 + 1.836 - 0.194 = 2.322$, and $\tanh(2.322) = 0.981$. Every number in the trace is a product of data bytes, weight matrix entries, and \tanh —all findable, all checkable.

Pattern as SN-form model. The interpretable model for this prediction, in SN form:

```

IF  data[t] in {a,m,n,e}      (word chars)
AND data[t-1] = 'a'          (P(e|a) = 0.68)
AND word_len >= 2           (mid-word)
AND in_tag = 1               (inside XML tag)
THEN P(e) ~ 1.00

```

This is a single entry in a table of patterns that, taken together, explain nearly all of the RNN’s predictions. The trained RNN has embedded exactly this pattern—and thousands like it—in its weight matrices. The factor map reads them out; the trace verifies them; and the SN form makes them human-readable.

Looking at this pattern, we can immediately see it *will* generalize: the sequence `</name>` appears throughout Wikipedia XML, and the conditional $P(e | a, \text{in_tag})$ reflects a genuine regularity. Conversely, a pattern like “3 at offset 20 predicts s” (from the UM superset at that offset, log-ratio +2.0 bits) is a coincidence of the specific 1024-byte window and will *not* generalize. The abduction step—distinguishing real patterns from coincidences—is what turns a statistical model into an interpretable one.

6.6 Reverse isomorphism: data to RNN without training

The factor map provides the forward direction: given the RNN, find the patterns. The reverse direction—given the patterns, construct the RNN—completes the isomorphism. We implement this as:

1. **UM learning:** count all bigram co-occurrences $(x@d, y)$ for offsets $d = 1, \dots, 50$.
2. **Feature encoding:** for each neuron j assigned to offset g and output byte y_j , set $h_j[t] = \log P(y_j | \text{data}[t - d_g]) - \log P(y_j)$, the log-likelihood ratio from UM counts.
3. **Readout optimization:** optimize W_y via gradient descent on the constructed h vectors.

No gradient-based training of W_x or W_h is needed—the hidden state is computed directly from data using UM-learned conditional distributions.

Construction	bpc	Encoding
Marginal (unigram)	4.74	—
Random hash (8 off \times 16n)	0.91	± 1 hash
Conditional prob (8 off \times 16n)	0.28	$P(y x) - P(y)$
Log-prob ratio (8 off \times 16n)	0.107	$\log P(y x)/P(y)$
Trained RNN (4000 epochs SGD)	0.079	learned
UM floor (skip-6, perfect)	0.000	counting

The log-prob ratio construction reaches **0.107 bpc**—within 0.03 of the trained RNN’s 0.079 bpc—using only UM-learned conditional distributions and W_y optimization. The random hash baseline (0.91 bpc) shows that the choice of encoding matters enormously: log-likelihood ratios preserve nearly all the UM’s learned information, while binary hashes discard most of it.

The 0.03 bpc gap between construction (0.107) and training (0.079) represents: (a) the benefit of learning W_x and W_h jointly with W_y , and (b) higher-order patterns captured by the recurrence that single-offset bigram features cannot represent. Given that the UM floor is 0.000 bpc, the construction captures $(4.74 - 0.107)/(4.74 - 0.079) = 99.4\%$ of the trained RNN’s total prediction quality.

Reproducibility

Repository: <https://github.com/inimino/hutter> (commit: 74ddb51)

Model: Saturated RNN, $256 \rightarrow 128 \rightarrow 256$, tanh, BPTT-50. Checkpoint: `sat_model.bin` (329 KB). Performance: 0.079 bpc on 1024 bytes.

Data: First 1024 bytes of enwik9 (<http://mattmahoney.net/dc/enwik9.zip>).

Analysis tools: `binary_states.c`, `factor_map.c`, `factor_map2.c`, `factor_map3.c`, `factor_map4.c`, `rnn_trace.c`, `um_learn.c`, `reverse_iso.c`, `reverse_iso2.c`.

To reproduce:

```
git clone https://github.com/inimino/hutter.git
cd hutter && git checkout 74ddb51
bash reproduce.sh
```

References

- [1] Michaeljohn Clement. CMP. 2026. <https://cmp.ai/cmp.pdf>
- [2] Claude and MJC. Pattern Priors. 7 Feb 2026.
- [3] Claude and MJC. The Hidden Quotient. 7 Feb 2026.
- [4] Claude and MJC. The Export Gap. 7 Feb 2026.