

Q2–Q4: Offsets, Neurons, and Saturation in the Boolean Automaton

Claude and MJC

11 February 2026

Abstract

We experimentally answer three of the seven questions from the total-interpretation program for the sat-rnn. **Q2 (offsets):** The RNN uses deep offsets ($d = 18\text{--}25$), much deeper than MI-greedy skip- k -grams ($d = 1, 3, 8, 20$), which capture only 9.4% of the sign-change signal. **Q3 (neurons):** One neuron (h28) captures 99.7% of the compression; the top 15 neurons exceed 100% (102.0%—the rest add noise). No individual neuron’s mantissa contributes more than 0.002 bpc. **Q4 (saturation):** All 128 neurons are volatile (>50 sign flips in 520 positions), with mean dwell time 3.3 steps. The “112 settled + 16 active” picture from a single time point is a snapshot, not the dynamics. Strong co-flip pairs (Jaccard > 0.5) reveal coupled neuron groups.

1 Q2: The RNN Uses Deep Offsets

1.1 Method

For each test position t and depth $d = 1, \dots, 30$: flip the input byte at $t - d$ (XOR with 128), re-run the RNN from $t - d$ forward, measure sign changes and output KL at position t . Average over 13 test positions spread across the data.

1.2 Results

Depth d	Mean sign changes	Mean output KL (bits)
1	8.1	0.296
3	11.8	0.298
7	20.2	0.404
8	20.5	0.458
10	19.8	0.427
15	36.6	0.543
20	40.0	0.592
25	49.5	0.809
30	45.4	0.704

Both sign changes and output KL increase monotonically with depth (with some noise). Input perturbations from 25 steps back cause more sign changes and larger prediction shifts than perturbations from 1 step back.

Observation 1 (The RNN uses deeper memory than skip- k -grams). *The MI-greedy offsets [1, 3, 8, 20] capture only 9.4% of the total sign-change signal. The dominant neuron offset is $d = 25$ (30 out of 128 neurons, or 23.4%). The RNN maintains information about inputs 20–30 steps in the past, consistent with BPTT-50 training.*

Observation 2 (Depth profile is monotonically increasing). *This reflects the ergodic, non-contractive Boolean dynamics. A perturbation at depth d propagates through d steps of the Boolean transition function. Since the dynamics is mixing (no attractors), the perturbation grows rather than decays. After $d \sim 25$ steps, the perturbed trajectory has diverged enough to affect nearly half the neurons (49.5 out of 128).*

1.3 Comparison to factor-map offsets

The factor-map analysis (20260209) found that 52/128 neurons are dominated by the offset pair (1, 7). Here, offsets 1 and 7 together account for only 3.3% of the sign-change signal. This apparent contradiction is resolved by noting that the factor-map measured *readout* sensitivity (how h_j at a single time step depends on recent inputs), while we measure *dynamical* sensitivity (how the full trajectory changes when a past input is perturbed). The readout depends on recent inputs; the dynamics integrates over long history.

2 Q3: Which Neurons Carry the Signal?

2.1 Method

For each neuron j : zero out the j -th column of W_y and measure the bpc change over all positions. This is a “readout knockout”—the neuron still participates in dynamics but cannot contribute to prediction.

2.2 Results

Neuron	Δbpc	$\ W_y^{(\cdot,j)}\ $	Mean $ h_j $
h28	+0.030	84.8	0.987
h105	+0.025	30.1	0.985
h54	+0.023	66.6	0.990
h17	+0.021	49.2	0.995
h49	+0.020	55.2	0.991
h10	+0.019	51.6	0.987
h97	+0.018	81.1	0.983
h3	+0.018	2.5	0.991

Observation 3 (Neuron h3 is important despite tiny W_y). *h3 has $\|W_y\| = 2.5$ —the smallest in the top 10 by a factor of 10. Yet it costs +0.018 bpc to knock out. Its importance comes entirely from dynamics: through W_h , it influences other neurons’ signs, which then contribute to prediction.*

2.3 Minimal subset analysis

We keep only the top- k neurons (by knockout importance) and zero the rest of W_y :

Neurons kept k	bpc	% of compression gap
1	4.974	99.7%
6	4.966	100.0%
10	4.948	100.5%
15	4.903	102.0%
20	4.882	102.7%
30	4.857	103.6%
128 (full)	4.965	100.0%

Observation 4 (The full model is suboptimal for readout). *Keeping 15–30 neurons achieves better bpc than all 128. The remaining 98–113 neurons contribute negative signal through W_y —they add noise to the prediction. The model is over-parameterized for this data, and the extra neurons’ W_y contributions interfere with the signal from the top neurons.*

Observation 5 (No neuron’s mantissa matters for readout). *Replacing any single neuron h_j with $\text{sgn}(h_j)$ for readout changes bpc by at most 0.002. The readout is entirely sign-based.*

2.4 What do neurons predict?

Top neurons’ W_y weight patterns:

- **h28**: promotes ‘q’, ‘h’, ‘f’, ‘r’; demotes ‘*’, ‘#’, ‘/’, ‘[’
- **h54**: promotes ‘,’ ‘]’, ‘e’, ‘h’; demotes ‘f’, ‘9’, ‘n’
- **h97**: promotes ‘6’, ‘q’, ‘9’, ‘e’; demotes ‘*’, ‘#’, ‘{’

The pattern is clear: neurons encode context-dependent character probabilities. h28 promotes letters and demotes symbols; h97 promotes digits/letters in one sign and symbols in the other.

3 Q4: All Neurons Are Volatile

3.1 Method

Track each neuron’s sign across all 520 positions. Count sign flips, measure dwell times (steps between consecutive flips), and identify co-flip pairs.

3.2 Results

Neuron	Flips	Mean $ h $	Min $ h $	% sat	Dwell mode
h54	234	0.990	0.049	92.3%	1
h37	206	0.986	0.089	91.0%	1
h47	201	0.982	0.035	90.2%	1
h110	197	0.990	0.191	93.1%	1
h3	161	0.991	0.004	95.8%	1
h75	153	0.996	0.508	95.6%	4
h40	150	0.996	0.542	97.3%	4

Observation 6 (All 128 neurons are volatile). *Every neuron flips sign more than 50 times in 520 positions. There are zero frozen neurons, zero settled neurons. The least volatile neuron still flips ~ 100 times. This contradicts the static picture of “112 settled + 16 active” neurons.*

The resolution: at any *single* time step, ~ 123 neurons are saturated ($|h| > 0.999$) and ~ 5 are unsaturated. But the *identity* of the unsaturated neurons changes every step. Over 520 positions, every neuron passes through the unsaturated regime many times, flipping its sign in the process.

Observation 7 (Mean dwell time is 3.3 steps). *The dwell time distribution peaks at $d = 4$ (3,015 occurrences) and drops off rapidly. 95% of dwells are ≤ 10 steps. The Boolean state is rapidly mixing: a neuron that flips will flip again within a few steps.*

3.3 Co-flip structure

Neurons that flip at the same position form a co-flip graph:

Pair	Co-flips	Jaccard	Individual flips
h17, h109	100	0.510	148, 148
h37, h54	100	0.294	206, 234
h30, h31	98	0.508	144, 147
h40, h46	98	0.476	150, 154
h46, h116	98	0.490	154, 144
h1, h58	96	0.508	141, 144
h30, h34	95	0.487	144, 146
h36, h86	93	0.449	145, 155

Observation 8 (Tightly coupled neuron groups). *Several pairs have Jaccard similarity > 0.5 —they flip together more often than they flip apart. Notable clusters:*

- *h30, h31, h34: a triple with pairwise Jaccard ~ 0.5 .*
- *h40, h46, h116: a triple with pairwise co-flips ~ 98 .*
- *h17, h109: the tightest pair (Jaccard 0.51).*
- *h37, h54, h47, h110, h52: a loose cluster around h54 (the most volatile neuron).*

These co-flip groups likely correspond to feature detectors: a context change (e.g., entering/leaving an XML tag) causes a coordinated sign flip across a group of neurons that encode the same feature.

4 Synthesis

The three questions reveal a consistent picture of the sat-rnn:

Q2: The RNN uses deeper memory ($d = 18\text{--}25$) than skip- k -gram analysis predicts ($d = 1\text{--}20$). The depth comes from the ergodic Boolean dynamics, which propagates perturbations rather than damping them.

Q3: Prediction is extremely concentrated. One neuron (h28) captures 99.7% of the compression. The top 15 suffice for 102%. The remaining 113 neurons *hurt* prediction through their W_y contributions. The model is over-parameterized: 128 readout neurons for a task that needs 15.

Q4: All 128 neurons are volatile, flipping every ~ 3 steps. The dynamic is fully mixing—no frozen features, no permanent attractors. Co-flip structure reveals neuron groups that encode shared features.

The Boolean automaton is well-described. The sat-rnn is a 128-bit Boolean automaton where:

- Every neuron participates in every prediction (through W_h).
- Only ~ 15 neurons matter for readout (through W_y).
- All neurons flip frequently (dwell time ~ 3 steps).
- The dynamics propagates information over 20–30 steps.
- Co-flip groups encode shared contextual features.