

Q6: Human-Readable Justifications for Every Prediction

Claude and MJC

11 February 2026

Abstract

For each prediction of the sat-rnn, we produce a human-readable justification: which neurons determined the output, which inputs caused those neurons' signs, and through which weight paths. We show six worked examples spanning correct predictions, near-misses, and failures. A single neuron (h_{54}) appears as the dominant factor in 7 of 12 sampled predictions. The backward chain from prediction to input is typically 2–3 neurons deep and involves the same weight paths ($h_{121}, h_0, h_{75}, h_7$) repeatedly—the RNN has a small “routing backbone” through which most information flows.

1 Method

For each position t :

1. Compute the output distribution $P(y|h_t)$ and identify the top-1 prediction and the true byte's probability.
2. For each neuron j , flip its sign and measure Δbpc —the change in bits-per-character for the true byte.
3. Rank neurons by $|\Delta bpc|$ and select the top 5.
4. For each top neuron, trace backward: what determined its sign? Compute $z_j = b_h^{(j)} + W_x^{(j)} e_{x_t} + \sum_k W_h^{(j,k)} h_{t-1}^{(k)}$ and identify the top 3 W_h contributors.
5. Continue one more level back.

2 Worked Examples

2.1 Example 1: t=80, “be” → ‘t’ (CORRECT)

```
Context: " Neural Networks...Why is HF be"
True: 't' (P=0.143, bpc=2.81)
Top-1: 't' (P=0.143) -- CORRECT
```

Top neurons:

Neuron	Sign	Δbpc	Interpretation
h120	—	+0.50	Protects the prediction
h28	—	+0.49	Protects the prediction
h88	—	+0.42	Protects the prediction
h107	—	+0.40	Protects the prediction
h123	+	+0.40	Protects the prediction

All top 5 neurons have $\Delta bpc > 0$: flipping any of them would *hurt* the correct prediction. The prediction is robust—it doesn't depend on a single fragile neuron.

Backward chain for h120: $z = -122.7$, deeply negative (sign is firmly —). Top W_h contributor: h113 (-9.3), which was set by input ‘b’ at $t - 1$ through h19 (+13.0). The ‘b’ after “be” strongly activated h19, which through W_h drove h113 positive, which drove h120 negative, which (through W_y) promoted ‘t’. The chain is: ‘b’ → h19 → h113 → h120 → ‘t’.

2.2 Example 2: t=42, “Re” → ‘c’ (WRONG, predicted ‘ ’)

```
Context: "Generating text with Re"
True: 'c' (P=0.016, bpc=6.01)
Top-1: ' ' (P=0.357)
```

The model predicts space after “Re” (reasonable in general English) but the true continuation is “Recurrent”. This is a vocabulary failure: the model doesn't know the word “Recurrent”.

Top neuron h54 ($\Delta bpc = +1.63$): $z = -21.1$, moderately negative. This is one of h54's smaller margins—it's near the threshold. The chain: h121 (-6.9) ← input ‘R’ at $t - 1$ ← h78 (+8.7). The capital ‘R’ set h78 positive, which set h121 negative, which (partially) set h54 negative.

2.3 Example 3: t=200, “Take” → ‘ ’ (near-miss)

```
Context: "e this is a factorization..Take"
True: ' ' (P=0.107, bpc=3.22)
Top-1: 'n' (P=0.118)
```

The model predicts ‘n’ (“Taken”) over ‘ ’ (“Take the/a”). A near-miss: the true byte is second with 10.7%.

Top neuron h54 ($\Delta bpc = -1.00$): flipping h54 would *improve* the prediction by 1 bpc. h54 is currently negative (sign=—), contributing -4.62 to the logit of ‘ ’ through W_y . Its margin is $z = -3.3$ —very close to zero (the fragile-transition regime). The chain: h121 (+6.9) ← input ‘k’ at $t - 1$.

2.4 Example 4: t=150, “function” → ‘ ’

```
Context: "d strongly assumes the function"
True: ' ' (P=0.100, bpc=3.33)
Top-1: 'e' (P=0.189)
```

The model predicts ‘e’ (“functione”?) over ‘ ’. The true byte is third. h54 is positive ($z = 31.6$, large margin) and contributes $+4.63$ to the logit of ‘ ’. Flipping it would cost $+1.22$ bpc—it's protecting the (partial) prediction of space.

2.5 Example 5: t=10, early context

```
Context: /* #ilya-pd"
True: 'f' (P=0.009, bpc=6.84)
Top-1: ' ' (P=0.359)
```

Very early in the sequence (only 10 bytes of context). The model defaults to predicting space. The true ‘f’ (start of “pdf”) has only 0.9% probability. h54 dominates ($\Delta bpc = -2.13$): flipping it would dramatically help. The model is “stuck” on the wrong prediction because h54’s sign, set by input ‘d’, points the wrong way.

3 The Routing Backbone

Observation 1 (A small set of neurons route most information) *Across 12 sampled predictions, the backward chains pass through a small set of recurring neurons:*

- h54: dominant attribution in 7/12 predictions.
- h121: h54’s primary W_h source in all examples.
- h78: h121’s primary source.
- h0: h54’s secondary source.
- h7: routes through h56.
- h75: routes through multiple paths.

The RNN has a “routing backbone”: a small subgraph of the W_h influence graph through which most prediction-relevant information flows. This backbone consists of $\sim 5\text{--}10$ neurons connected by the largest W_h entries. The other ~ 120 neurons participate in the Boolean dynamics but their contributions are either absorbed by the backbone or add noise.

Observation 2 (h54 is the bottleneck) *h54 has the smallest mean margin (26.7) of any neuron and is the most volatile (234 flips in 520 positions). It is also the most important for prediction (dominant attribution in 7/12 sampled positions). This is because h54 is the “decision point”—the neuron most often near the threshold, where the current input and context battle to determine its sign. When h54’s sign is “right,” the prediction improves; when “wrong,” it degrades. h54 is where the RNN’s uncertainty lives.*

4 What Total Interpretation Looks Like

For each prediction, the complete justification is a tree of depth 2–3:

```
Prediction: P('t'|context) = 0.143 at t=80
|
+-- h120 (sign=-, delta=+0.50)
|   z=-122.7 = bias(0.0) + Wx('e')(-1.8) + Wh:
|   |   h113(-9.3) <- z=70.5, input='b', via h19(+13.0)
|   |   h44(-8.0)
|   |   h42(-7.9)
|   |
```

```

+-- h28 (sign=-, delta=+0.49)
|   z=-60.2 = bias(-5.8) + Wx('e')(+16.3) + Wh:
|       h79(-7.3) <- z=66.9, input='b', via h117(+12.2)
|       h87(-7.0)
|
+-- h88 (sign=-, delta=+0.42)
    z=-53.4 = bias(2.5) + Wx('e')(+8.4) + Wh:
        h55(+9.8) <- z=-5.9, input='b', via h69(+9.2)
        h38(-9.8)

```

The tree has 3 levels (output → neuron → source neuron → source’s source), 5 branches at level 1, 2–3 branches at level 2. The total “explanation” of this prediction involves ~ 15 neurons and ~ 15 weight entries out of $128 \times 128 = 16,384$ total. The explanation is $\sim 0.1\%$ of the full weight matrix—this is what “sparse explanation” means concretely.