Internal states before *wait* modulate reasoning patterns

Dmitrii Troitskii*, Koyena Pal*, Chris Wendler, Callum Stuart McDougall, Neel Nanda

Motivation

- Crosscoders allows us to discover **thousands** base/ shared/reasoning features in an unsupervised manner.
- We present a method to find an interesting subset of features: **latent attribution** focusing on features relevant from promoting/supressing *wait*
- Many of them are interpretable and modulate specific reasoning behaviors



Input: $x_{0...t}$ x_0 · · · · · · · x_t



Latent Attribution for k Latents

Sparse Crosscoder

Base Model: Llama-3.1-8B Reasoning Model: DeepSeek-R1-Distill-Llama-8B



Where could the features come from? Base-only Shared

Reasoning-only

Trained 3 L1 Sparse Crosscoders at 25%, 50%, 75% layer depths respectively (Lindsey et al, 2024)[1]

Attributed Feature Distribution according to Crosscoder Relative Norm Difference

Latents before wait tokens influence reasoning. Features exposed via crosscoders and attribution reveal distinct reasoning patterns in DeepSeek-R1-Distill-Llama-8B.

Input

< | begin_of_sentence | >< | User | >Find three prime numbers that add up to 100.< | Assistant | ><think>Alright, so I need to find three prime numbers that add up to 100. Hmm, okay. Let me think about how to approach this. Prime numbers are numbers greater than 1 that have no divisors other than 1 and themselves. So, the primes are 2, 3, 5, 7, 11, 13, and so on. ...

Reference Response Wait, but just to make sure there are no other combinations, maybe without using 2? Let me think. If I don't include 2, then all three primes are odd. Adding three odd numbers: odd + odd = odd + odd = even, then even + odd = odd. But 100 is even, so that won't work. Therefore, we must include 2 as one of the primes. So, the other two primes must add up to 98, which we found as 19 and 79. ...

Feature ID: 744 (in bottor

[Steering Strength: 1.5] Let m back to the main point. Maybe 1 get back to the initial approach. I should try another prime and should go back to the beginning go back to 3. Maybe I should st and try a different approach. M should get back to the original p

[Steering Strength: 1.25] Wait, go back and see if I should try a approach. Maybe I should go ba the original problem and try to so over. Let's get back to the begin was trying to find three prime m Let's not forget that I should try stick to primes ...

> **Reasoning Behavior:** Go to Initial Approach



Top 50

• Shared • Reasoning only

gonly • Base only

Bottom 50

Intervened Responses and Observed Reasoning Behavior Change

m 50)	Feature ID: 31748 (in top 50)	Feature ID: 25929 (in bottom 50)	Feature ID: 188 (in botto
ne get I should . Maybe see if I g. Let's tart over faybe I problem.	[Steering Strength: 1.5] Wait, I think I just remembered a question about the Goldlerian Numbers, or something about a recent question. Is there a way to have a prime number, the way I'm saying, but if I have a bit more, the Prime, the number, the word I'm talking about the number of the following, just a thought,	[Steering Strength: 1.0] Alternatively, I might have started with a different prime. Let me try another approach. Let's assume all three primes are odd. Since 2 is the only even prime, this is impossible. So, one prime must be 2, and the other two are primes. Therefore, this is the only solution. I think this is the answer.	[Steering Strength: 1.5] I if I need to find which prin I don't know how many ma know what other primes I how to find the other two without which I don't know don't know which one I do don't know what I don't kr
, let me another ack to start ming. I umbers. y to	[Steering Strength: 1.25] Wait, I think I read something about there are three prime numbers with a tip of the week. Hmm, what's the latest in the news. Wait, I just made a note of the new study about the U.S. government, did you see the new study about the	[Steering Strength: - 1.25] Wait a second, hold on. Is there another possibility? Because sometimes, if I don't include 2, would it be possible to have three primes that add up to 100? Let me check that because maybe it's not necessary to use 2, but	[Steering Strength: 1.25] don't know if I need to fin ones exactly, I don't know there are, so maybe I don't which ones I are, I don't know many I don't know how m I don't know how many I don't know how many I don't know how many I don't know how
: 1	U.S. government, Reasoning Behavior: Knowledge Recall	Reasoning Behavior: Conclusion (in positive steer) Re-trying (in negative steer)	Reasoning Beha Uncertainty
	T		<u> </u>

END

Latent Attribution

How much does each latent component in a crosscoder contributes to the change in downstream metric, M_{patch} ?

 $M_{patch} = \log p_{\theta}(y_{wait} \mid L)$

Reasoning Instances from Venhoff et al., 2025 [2]

• *wait* tokens include "Wait", "Wait", "wait", "wait"

Patchscope

What can we decode from the features?



om 50)

don't know imes exactly nore I don't don't know primes ow what? I on't know l now ... Wait, I nd which how many know now how nany primes. don't know IOW ...

avior:

What happens when we steer each feature?



[1] Lindsey, J., Templeton, A., Marcus, J., Conerly, T., Batson, J., and Olah, C. Sparse crosscoders for cross-layer features and model diffing, 2024. URL https://transformer-circuits.pub/2024/crosscoders/index.html.
[2] Venhoff, C., Arcuschin, I., Torr, P., Conmy, A., and Nanda, N. Understanding reasoning in thinking language models via steering vectors. In Workshop on Reasoning and Planning for Large Language Models, 2025. URL https://openreview.net/forum? id=OwhVWNOBcz.



