
When Is an Attention Head a Computational Unit? Shared Circuits Despite Orthogonal Representations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Mechanistic analyses of transformers often proceed head-by-head, looking for at-
2 tention maps that correspond to computational roles. But when should an attention
3 head be interpreted as a computational unit? In principle, multi-head attention need
4 not organize computation in such a clean way. A single computational function
5 may be distributed across heads, and a single head may participate in several
6 computations, with invariants appearing only after OV routing and residual-stream
7 summation. In natural language models, this is hard to resolve because the under-
8 lying features and computations are unknown, and an apparently uninterpretable
9 head may be polysemantic, part of a distributed circuit, or a sign that the anal-
10 yst chose the wrong decomposition. To get a handle on this issue, we study a
11 controlled setting where the target computation is analytically specified: transfor-
12 mers trained on factored sequence processes. These processes require independent
13 predictive updates for multiple latent factors, represented in orthogonal residual-
14 stream subspaces. This lets us ask whether the attention circuits implementing
15 these independent updates are themselves factorized. We find that they need not
16 be. Per-factor updates constrain only the aggregate routed contribution across
17 heads, leaving substantial freedom in how computation is distributed. Heads may
18 specialize to individual factors, compose with other heads to implement one factor,
19 or contribute polysemantically across factor boundaries. The regime that emerges
20 depends on the generator’s spectrum, the head budget, and training dynamics. We
21 introduce *effective subspace attention*, a scalar quantity that combines attention
22 patterns with OV routing to recover invariant subspace-level contributions when
23 individual maps are illegible. Our results show that independent computations
24 can be implemented by shared attention circuits, and that individual heads may
25 look uninterpretable even when the collective routed circuit implements a precise,
26 theoretically predicted computation.

27 1 Introduction

28 When should an attention head be interpreted as a computational unit? Mechanistic analyses often
29 proceed head-by-head, looking for attention maps that correspond to distinct roles. But multi-head
30 attention need not organize computation this cleanly: a single computation may be distributed across
31 heads Jermyn et al. [2023], and a single head may contribute to multiple computations Janiak
32 et al. [2023]. In such cases, the invariant object is not an individual attention map, but the aggregate
33 residual-stream update produced after OV routing and summation.

34 In natural-language models we cannot easily tell these cases apart, because we lack ground truth for
35 what the underlying features are and what computations the model is meant to implement. To study
36 the question precisely, we need a setting where the functional units of the computation are known in

37 advance and are controllable. We can then ask if and how those known functions distribute across
 38 heads, and why.

39 We use transformers trained on *factored sequence processes*: sequences generated by independent
 40 computational units, called factors, that evolve independently. Prior work shows that transformers
 41 trained on such processes represent each factor in its own orthogonal subspace of the residual
 42 stream Shai et al. [2026]. In this setting the computational units are known (each factor), the subspaces
 43 they occupy are known (orthogonal subspaces of the residual stream), and the per-factor update that
 44 attention must implement at each layer can be derived analytically from the data-generating process.
 45 We then ask: when the computation and corresponding representations of a network decompose into
 46 separable computational units, do the attention circuits naturally factor in the same way?

47 We find that they do not. Even with independent computational units assigned to orthogonal subspaces
 48 of the residual stream, individual heads can be illegible while the collective routed contribution
 49 matches theory exactly. A computation downstream from an attention head simply receives the full
 50 residual stream, and whether the information present within was written by one or multiple heads is
 51 irrelevant.

52 Our main contributions are:

- 53 1. We derive a **predictive theory for the collective update an attention layer must write** into each
 54 independently updated subspace of the residual stream, along with the QK/OV division of labor
 55 it requires: QK supplies the temporal weighting, OV routes token-dependent contributions into
 56 the right subspace. From this we predict the minimum number of heads needed for any given
 57 configuration.
- 58 2. We show that this subspace update admits **many head-level decompositions**. Heads can be
 59 specialized, compositional, and polysemantic, all while achieving the same loss.
- 60 3. We introduce **effective subspace attention**, a quantity which decomposes the collective con-
 61 tribution of multiple attention heads into each subspace when individual attention maps are
 62 illegible.

63 2 Background

64 2.1 Data-generating processes and belief state geometry

65 We consider training data generated by edge-emitting hidden Markov models (HMMs). An HMM is
 66 defined by a token alphabet \mathcal{X} , a set of hidden states \mathcal{S} , an initial distribution $\eta^{(\emptyset)}$ over hidden states,
 67 and a collection of token-labeled transition matrices $(T^{(x)})_{x \in \mathcal{X}}$, where $T_{i,j}^{(x)} = \Pr(x, s_j \mid s_i)$ gives
 68 the joint probability of emitting token x and transitioning to state s_j given current state s_i .

69 An optimal Bayesian observer who knows the HMM structure but not the current hidden state
 70 maintains a *belief state* η over hidden states, updated upon observing each token x via:

$$\eta' = \frac{\eta T^{(x)}}{\eta T^{(x)} \mathbf{1}}, \quad (1)$$

71 Starting from an initial state $\eta^{(\emptyset)}$, the belief state after observing context $x_{1:\ell}$ is:

$$\eta^{(x_{1:\ell})} = \frac{\eta^{(\emptyset)} T^{(x_1)} \dots T^{(x_\ell)}}{\eta^{(\emptyset)} T^{(x_1)} \dots T^{(x_\ell)} \mathbf{1}}. \quad (2)$$

72 For stationary processes, the optimal initial belief state is the stationary distribution $\eta^{(\emptyset)} = \pi$, the
 73 left eigenvector of $T = \sum_x T^{(x)}$ associated with eigenvalue 1.

74 Prior work shows that transformers trained on next-token prediction linearly represent the belief
 75 states induced by the data-generating process in the residual stream [Shai et al., 2024, Riechers et al.,
 76 2025].

77 **Factored representations** The Factored World Hypothesis [Shai et al., 2026] formalizes generators
 78 with independently evolving latent factors, and shows that transformers represent the corresponding
 79 per-factor beliefs in approximately orthogonal residual-stream subspaces. We call a process *inde-*
 80 *pendent* if each token-labeled transition operator factorizes as $T^{(x)} = \bigotimes_{n=1}^N T_n^{(x)}$, where factor n

81 has hidden state space \mathcal{S}_n and transition dynamics $T_n^{(x)}$. The observed token x is determined by
 82 the joint output of latent subtokens $z^{(n)}$, though this decomposition is not exposed in the training
 83 data. Under independence, product beliefs remain product beliefs under Bayesian updating: if
 84 $\eta^{(x_{1:\ell})} = \bigotimes_n \eta_n^{(x_{1:\ell})}$, then $\eta^{(x_{1:\ell+1})} = \bigotimes_n \eta_n^{(x_{1:\ell+1})}$. Thus, the joint belief can be represented
 85 losslessly as per-factor beliefs with factor n occupying its own $(d_n - 1)$ -dimensional subspace, where
 86 $d_n = |\mathcal{S}_n|$.

87 2.2 An architectural division of labor

88 The controlled setting we work in has a structurally manifest decomposition in the data generator
 89 itself, and that decomposition maps cleanly onto distinct components of the attention architecture.
 90 This lets us state up front what the architecture *must* implement, and which parts of the architecture
 91 are capable of implementing it. We describe the intuition here.

92 For a factored generator, each factor transition matrix $T_n^{(x)}$ has its own independent state space, so
 93 optimal beliefs evolve independently by factor. Each update requires three pieces of information:
 94 (i) *which factor* to update, set by the generator’s tensor-product structure; (ii) *how strongly* to weight
 95 each temporal offset, set by the factor spectrum; and (iii) *what token-conditional vector* to write into
 96 each factor’s subspace. This maps directly onto attention: QK chooses the temporal weights, while
 97 OV chooses what to write into each factor subspace. This yields two predictions:

98 **OV carries factor identity.** For any attention pattern, head h writes only in the column space of
 99 $W_O^{(h)} W_V^{(h)}$: attention reweights value vectors but cannot change their span. Thus, if the residual
 100 stream decomposes into orthogonal factor subspaces $\bigoplus_n \mathbf{V}_{F_n}$ [Shai et al., 2026], the OV maps must
 101 route contributions into these subspaces. Factor identity is therefore enforced by OV, not by the
 102 attention weights.

103 **QK reflects temporal structure, not token identity.** In our toy model, each source contribution
 104 decomposes into a token-dependent direction and a factor-specific scalar ζ_n^{d-s} . Since the direction
 105 is routed by OV, QK is left to implement only the scalar: a lag-dependent weight from source s to
 106 destination d . Thus QK should be primarily offset-shaped rather than content-routed, with different
 107 factors inducing different lag profiles.

108 These predictions constrain the *aggregate* layer-level circuit, not individual heads. They specify the
 109 QK/OV division of labor, but leave open how a fixed head budget distributes that work: heads may
 110 specialize, collaborate on one factor, or contribute polysemantically across factors. Thus orthogonal
 111 factor representations need not imply specialized heads. The remainder of this section makes this
 112 mechanism explicit.

113 2.3 Constrained belief updates

114 Optimal prediction requires sequential Bayesian updating: each new token refines the current
 115 predictive vector, which itself reflects all previous tokens (Eq. (1)). But attention is fundamentally
 116 parallel — at each position, it computes a weighted sum of value vectors from all source positions
 117 simultaneously. This creates a tension: how can a parallel mechanism approximate an inherently
 118 sequential computation?

119 The key insight of Piotrowski et al. [2025] is that attention resolves this by treating each source
 120 token as *independent evidence*. The BOS token at position 0 provides the prior π through the OV
 121 circuit — the baseline prediction before any evidence is observed. Each subsequent token x_s at
 122 offset $k = d - s$ then contributes an independent correction: “given the prior, how does observing
 123 x_s exactly k steps ago change my current predictions?” This correction is the geometric embedding
 124 of the difference between the prior $\Pr(S_d) = \pi$ and the pointwise update from the observed token
 125 $\Pr(S_d | X_s = x_s) = \pi T^{|x_s} T^{d-s}$. $T^{|x_s}$ captures the immediate evidence from the token and T^{d-s}
 126 propagates that evidence forward through $d - s$ steps of latent dynamics. Summing the prior from
 127 BOS with independent corrections from all source tokens gives the *constrained belief update*:

$$\mathbf{r}_1^{(x_{1:d})} = \underbrace{\pi}_{\text{from BOS}} + \sum_{s=1}^d \underbrace{\left(\pi T^{|x_s} T^{d-s} - \pi \right)}_{\text{correction from source } s}. \quad (3)$$

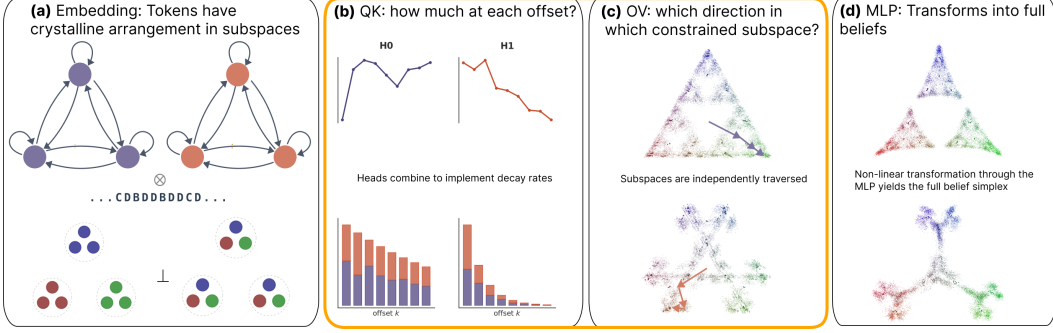


Figure 1: **Complete mechanistic description of a transformer trained on a factored process.** (a) The token embedding reflects the crystalline structure of the data generator: tokens arrange in factor-specific clusters within orthogonal subspaces of the residual stream. (b) The QK circuit determines how much each source token contributes at each offset, with heads combining to implement the temporal decay rates predicted by the spectral decomposition. (c) The OV circuit determines the direction and subspace of each contribution: displacement vectors within each factor’s constrained predictive geometry point toward vertices of the belief simplex, with each factor’s subspace traversed independently. (d) The MLP applies a nonlinear correction, transforming the constrained predictive geometry (intermediate fractal) into the full predictive geometry. In this work, we focus on the attention circuit (b–c, highlighted in orange): how the QK and OV circuits jointly implement belief updates across a factored world, the degrees of freedom in how this computation is distributed across heads, and what can be recovered when individual attention maps are not directly interpretable.

128 This is the best approximation to Bayesian updating achievable by summing independent pairwise
 129 contributions — precisely the functional form that attention can implement.

130 The spectral decomposition of the transition matrix T reveals a clean separation of roles. In our case,
 131 each factor is represented by a “Mess3” generator which has a single nontrivial decaying eigenvalue
 132 ζ of T . In this case, the corrections can be re-written through spectral decomposition as

$$\Delta \mathbf{r}_1^{(x_{1:d})} = \sum_{s=1}^d \zeta^{d-s} \pi T^{|x_s} P_\zeta, \quad (4)$$

133 where $P_\zeta T = T P_\zeta = \zeta P_\zeta$ (see Appendix I.1 for additional details).

134 Equation (4) realizes the division of labor above: lag appears only in the scalar ζ^{d-s} , while token
 135 identity appears only in the displacement $\pi T^{|x_s} P_\zeta$. The current-token term ($s = d, k = 0$) is the only
 136 exception, since it can also enter through the residual stream via the token embedding; see App. D.

137 Piotrowski et al. [2025] confirmed this correspondence in transformers trained on the Mess3 process
 138 (Appendix I) and further showed that when the eigenvalue ζ is negative, the alternating pattern ζ^{d-s}
 139 cannot be implemented by a single attention head, whose weights must be non-negative.

140 2.4 Per-factor constrained belief updates

141 For a factored process with transition operators $T^{(x)} = \bigotimes_{n=1}^N T_n^{(x)}$, each Mess3 factor n has its
 142 own transition matrix $T_n = \sum_{z^{(n)}} T_n^{(z^{(n)})}$, stationary distribution π_n , and eigenvalue ζ_n . Since the
 143 factors are represented in orthogonal subspaces, the constrained belief update (Eq. (3)) decomposes
 144 into independent per-factor updates:

$$\Delta \mathbf{r}_n^{(x_{1:d})} = \sum_{s=1}^d \zeta_n^{d-s} \cdot \mathbf{g}_n(z_s^{(n)}), \quad (5)$$

145 where $z_s^{(n)}$ is the sub-token for factor n at position s and $\mathbf{g}_n(z) = \pi_n T_n^{|z} P_{\zeta_n, n}$ is a token-dependent
 146 displacement direction that is independent of offset. For a given destination d , each correction term
 147 depends only on factor n ’s transition dynamics and the corresponding sub-token at source s .

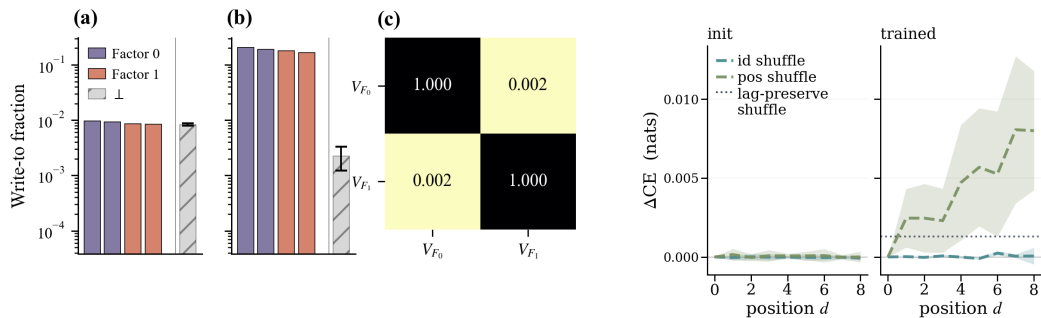


Figure 2: **OV routes factor-specific information into subspaces while QK specializes to temporal over semantic information.** **Left:** OV-circuit energy projected into the residual-stream subspaces associated with each factor, comparing (a) early training with the (b) end of training. Note that the bar for the perpendicular component represents an average over the remaining $d_{\text{model}} - 4$ dimensions not accounted for by $\mathbf{V}_{F_0} \oplus \mathbf{V}_{F_1}$. (c): The two factor-associated subspaces are each orthogonal to one another. This is averaged over 8 different models and 2 different heads per model. **Right:** holding all else fixed, we perform causal interventions on the QK-circuit and measure the impact on cross-entropy loss (ΔCE). Interventions include shuffling token identity and position information of the embeddings entering the circuit, as well as shuffling attention weights within lag bins.

148 3 From specialized heads to aggregate computation

149 All experiments use single-layer pre-norm decoder-only transformers trained on two-factor Mess3
150 sequences; full hyperparameters are in Appendix L.1.

151 3.1 When heads are interpretable: distinct positive eigenvalues

152 We first train a two-head transformer on a two-factor Mess3 process with eigenvalues $\zeta_0 = 0.7$
153 and $\zeta_1 = 0.4$. The natural ansatz is head specialization: head n implements the temporal decay
154 $A^{(n)}(d, s) \propto \zeta_n^{d-s}$ through QK, independent of subtoken identity and absolute position, while OV
155 routes its output into factor subspace \mathbf{V}_{F_n} .

156 First, we use a behavioral method to find $\mathbf{V}_{F_n} \subset \mathbb{R}^{d_{\text{model}}}$, 2-d orthogonal subspaces that correspond
157 to each factor. Once identified, we can define a write-to fraction that measures how much of the
158 image of the OV circuit lands inside or outside of each subspace (see Appendix E for a detailed
159 description of finding and evaluating these subspaces). Figure 2 shows that OV matrices initially
160 write no more to the factor subspaces \mathbf{V}_{F_n} than to other directions, but by the end of training place
161 most of their write-to energy in $\mathbf{V}_{F_0} \oplus \mathbf{V}_{F_1}$, a 4-dimensional subspace of the 120-dimensional
162 residual stream. Thus OV learns to route updates into the orthogonal factor subspaces. Separating
163 by head shows near-specialization. Recovering the intermediate activation geometry in Figure 13
164 also suggests strong alignment between a specific head and a factor’s belief states. However, it is not
165 exact specialization: in Figure 3, 11% of one factor’s write-to energy comes from the head primarily
166 associated with the other factor. This leakage is consistent with the training trajectory. Early in
167 training, both heads reduce loss on the slower-decaying factor (largest eigenvalue), which has the
168 larger available loss reduction; only later does one head re-specialize to the faster-decaying factor.
169 This kind of staged learning is consistent with prior work on learning dynamics in neural networks
170 and transformers [Kunin et al., 2025, Saxe et al., 2014, Yüksel et al., 2026]. Thus, even in a setting
171 where a specialized solution is available, gradient descent can transiently pass through collaborative
172 decompositions.

173 QK shows the complementary structure predicted by the theory. The rightmost panel in Figure
174 2 shows token-shuffle interventions that indicate that attention is primarily lag dependent: loss
175 increases when attention weights are shuffled across source positions, but not when shuffled only
176 among tokens at the same position. When the position shuffles are constrained to preserve lag, the
177 effect is significantly reduced. This indicates the QK circuits’ sensitivity to the position shuffle is

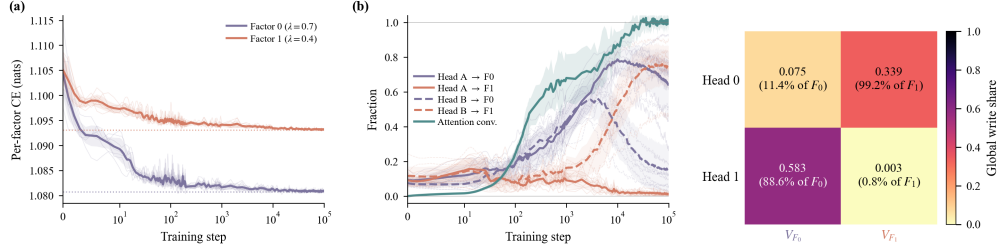


Figure 3: We measure how much the loss is reduced per-factor during training **(a)**, and for each factor, how much of each head’s energy is being written into that factor as it progresses towards its final attention maps **(b)**. For $\zeta_0 = 0.7$ and $\zeta_1 = 0.4$, during training, we see that the model prefers reducing loss for the factor with the larger eigenvalue (as there is more loss to reduce). Initially, this admits a collaborative solution where both heads focus on factor 0. Eventually, the second head re-specializes to factor 1. Plots are averaged across 8 heads, and heads are aligned such that head A is the one primarily responsible for factor 0 (and vice versa). In **(right)** we examine a single seed and see that the heads are largely specialized, but there is some leakage from head 0 into factor 0, explained by the fact that both heads initially focused on that factor.

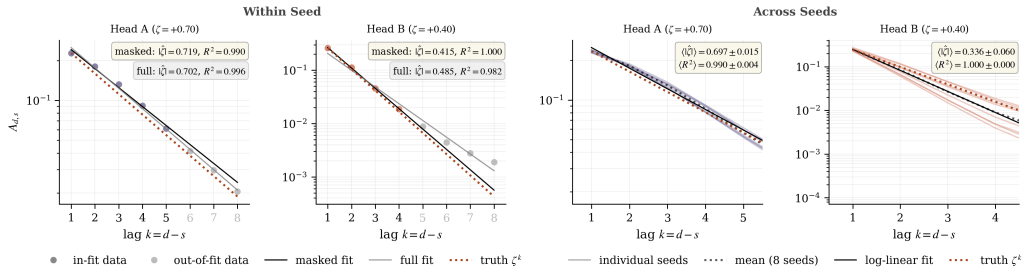


Figure 4: **Per-head attention patterns can be interpretable on their own.** In this case, for two factors with positive but different magnitude eigenvalues ($\zeta_0 = 0.70$ and $\zeta_1 = 0.40$), the decay rates of the attention profiles of each head uniquely match the temporal dynamics of one of the two factors – indicating we are in a regime of head *specialization*.

178 largely due to a sensitivity to lag and not the absolute position. Figure 4 shows attention weights
 179 binned by lag $k = d - s$. Since the $k = 0$ term shares a degree of freedom with the residual
 180 skip connection, we fit only off-diagonal attention weights; see Appendix D. The $\zeta_0 = 0.7$ head
 181 matches the predicted exponential decay, while the $\zeta_1 = 0.4$ head drifts at large lags, where the
 182 required corrections are small. Masking attention beyond lag τ confirms this: the slower-decaying
 183 factor remains sensitive up to $\tau = 5$, while the faster-decaying factor is sensitive only up to $\tau = 4$.
 184 We determine these thresholds using a noise floor. Lags whose ablation changes loss by less than
 185 10^{-5} nats are deemed irrelevant. See Appendix F for details. Fitting only behaviorally relevant lags
 186 yields good agreement for both heads.

187 This example cleanly separates OV routing from QK temporal weighting, but also shows why
 188 individual heads are not the invariant unit of analysis. The model is constrained only by the aggregate
 189 routed update written into the residual stream, not by which head supplies it. This motivates effective
 190 *subspace* attention.

191 3.2 Effective subspace attention

192 To extract the invariant factor-level computation, let f_n map the residual-stream contribution onto the
 193 belief update for factor n in Eq. (5). We define the *effective subspace attention*

$$\alpha_n(d, s) = \frac{\left\langle f_n \left(\sum_h A_{d,s}^{(h)} \mathbf{v}_s^{(h)} \right), \mathbf{g}_n(z_s^{(n)}) \right\rangle}{\left\| \mathbf{g}_n(z_s^{(n)}) \right\|^2}. \quad (6)$$

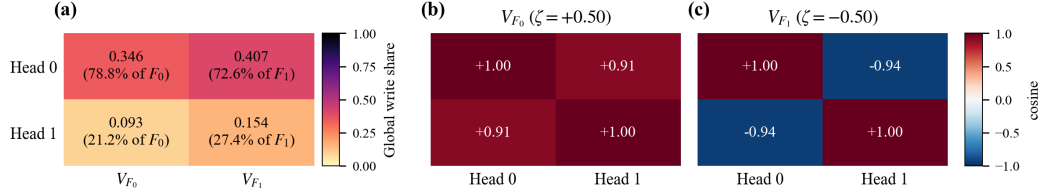


Figure 5: **In the mixed-sign case, heads are not specialized by write magnitude.** In (a) both heads write substantially into both factor subspaces. In (b) the heads write in nearly the same direction in the positive-eigenvalue subspace, but anti-parallel in the negative-eigenvalue subspace (c), providing the signed OV-coupling needed for alternating decay.

194 This projects the summed routed contribution from source s to destination d onto the token-dependent
 195 displacement direction for factor n . Under the constrained belief-update theory,

$$\alpha_n(d, s) = \zeta_n^{d-s}.$$

196 By linearity, α_n decomposes across heads as $\alpha_n(d, s) = \sum_h \alpha_n^{(h)}(d, s)$, so the theory constrains
 197 only the sum:

$$\sum_h \alpha_n^{(h)}(d, s) = \zeta_n^{d-s}. \quad (7)$$

198 Specialization is the case where one head supplies this quantity for a factor. Collaboration is the
 199 case where several heads sum to it. Polysemanticity is the case where a head contributes to multiple
 200 factors. We quantify these regimes with the coupling matrix induced by per-head effective subspace
 201 attention: polysemy measures spread across subspaces per head, and collaboration measures spread
 202 across heads per subspace; see Appendix C. Effective subspace attention lets us distinguish head-level
 203 structure from the invariant computation implemented by the full attention layer.

204 3.3 When heads are not interpretable: positive/negative eigenvalues.

205 We apply effective subspace attention to a two-head transformer trained on a factored process with
 206 eigenvalues $\zeta_0 = 0.5$ and $\zeta_1 = -0.5$. Piotrowski et al. [2025] showed that a negative eigenvalue
 207 requires two heads: because attention weights are nonnegative, a single head cannot implement the
 208 sign-alternating decay $(-0.5)^{d-s}$. In this setting, we find that the two heads collaborate to implement
 209 the negative-eigenvalue factor, but also collaborate on the positive-eigenvalue factor.

210 Let us begin by validating our predictions on the OV circuit. Just like the pattern shown in Figure
 211 2, we find that the subspaces \mathbf{V}_{F_0} and \mathbf{V}_{F_1} are *nearly* orthogonal to each other and carry the lion's
 212 share of variation in the residual stream of the trained network. Unlike the previous case, however,
 213 we find no evidence of head specialization when looking at the write-to fraction. Figure 5 shows
 214 that instead, both heads contribute to both factors in approximate proportion to their total write-to
 215 magnitude. However, the *direction* in which they write differs. We establish this by looking at the
 216 cosine between the two heads' write operators into \mathbf{V}_{F_n} , computed in matrix-Frobenius geometry
 217 (details can be found in Appendix E).

218 The right panel in Figure 5 reveals that while the global write share doesn't distinguish factor 0 and 1,
 219 the signed value does. The two heads write to \mathbf{V}_{F_0} in nearly the same direction, but they are almost
 220 exactly anti-parallel in how they write to \mathbf{V}_{F_1} , providing the mechanism to fit the alternating decay
 221 rate. The top panel of Figure 6 shows a naive attempt at fitting a decay rate for each head individually
 222 for a trained model. We can see each head implements an alternating attention pattern, with one head
 223 specializing to odd k and one specializing to even k . Fits to the expected exponential decays are very
 224 poor when we look only at a single head, but the effective subspace attention allows contributions
 225 from both heads, and recovers the predicted monotone decay for ζ_0 and the predicted alternating
 226 decay for ζ_1 . Recovering the geometry associated with each factor in Figure 14 also reveals that no
 227 single head is sufficient to represent the full geometry; both are required.

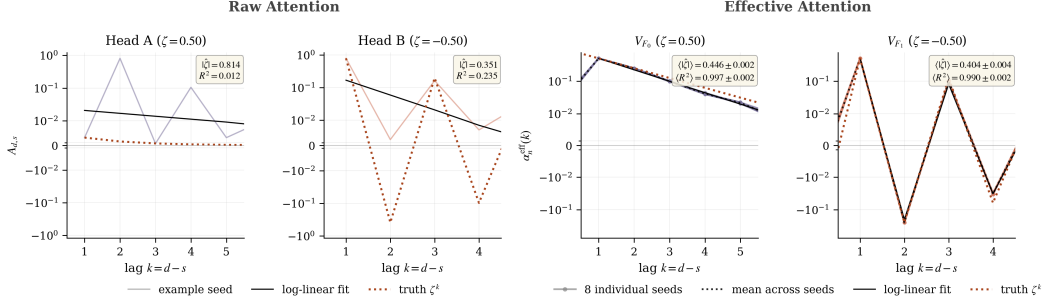


Figure 6: **Raw per-head attention profiles fail to recover the ground truth dynamics in the mixed-sign case, while effective subspace attention succeeds.** **Left:** neither individual head matches the predicted decay for either factor. **Right:** considering the relevant subspaces routed to by the OV vectors lets us recover the monotone positive decay and the alternating negative decay.

228 3.4 Minimum architecture

229 For a factored process, an attention head supplies a non-negative lag-dependent attention pattern,
 230 while its OV circuit determines how that pattern is signed and routed into factor subspaces. Thus the
 231 relevant architectural resource is not one head per eigenvalue or one head per factor, but the number
 232 of non-negative temporal rays needed to express the required signed lag-dependent updates after OV
 233 coupling. We refer to this as a *conic decomposition* of the update geometry.

234 This distinction matters in mixed-sign settings. A negative eigenvalue cannot be implemented by a
 235 single head, since attention weights are non-negative; two non-negative rays are needed to compose
 236 the alternating signed update. However, the same two non-negative rays can also be coupled with the
 237 same sign through OV to produce a monotone positive decay. Therefore, when a positive eigenvalue
 238 is present, it can sometimes be implemented "for free" by reusing the rays already required for the
 239 negative mode, reducing the head count below a naive per-eigenmode assignment. Appendix J gives
 240 the full conic argument. Table 1 lists the resulting predictions, which we test in the next section.

Table 1: Minimum head count under different two-eigenvalue cases.

Configuration	Values	Identifier	H_{\min}	H_{\min} (modes)
$+\zeta_1, +\zeta_2, \zeta_1 \neq \zeta_2 $	+0.7, +0.4	$D_{h=2}^+$	2	2
$+\zeta_1, +\zeta_1$	+0.7, +0.7	$I_{h=1}^+$	1	1
$+\zeta_1, -\zeta_1$	+0.5, -0.5	$S_{h=2}^\pm$	2	2
$+\zeta_1, -\zeta_2, \zeta_1 \neq \zeta_2 $	+0.7, -0.5	$D_{h=2}^\pm / D_{h=3}^\pm$	2	3
$-\zeta_1, -\zeta_1$	-0.5, -0.5	$I_{h=2}^-$	2	2

241 4 Effective attention recovers the invariant computation

242 Table 1 outlines the 5 different configurations of 2-factor processes we generate sequences from.
 243 Note that we omit the case of $-\zeta_1, -\zeta_2$ with $|\zeta_1| \neq |\zeta_2|$. This is because the narrow range of negative
 244 values for Mess3, $[-0.5, 0)$, means that the eigenvalues will either be so close together that the
 245 process can be approximated with fewer heads, or that one of the values will be close to 0 and
 246 therefore almost IID. The identifiers will be used in subsequent figures.

247 For each configuration in Table 1, we train models with $n_{\text{heads}} \in \{1, \dots, 6\}$ across eight random
 248 seeds. Figure 7 summarizes the central empirical result: the conic head-count prediction matches the
 249 empirical loss plateau, while effective subspace attention recovers the predicted factor-level updates
 250 even when raw attention maps are not individually interpretable. First, the minimum head count
 251 predicted by our conic framework aligns with the point at which loss plateaus: adding heads below this
 252 threshold substantially improves performance, while additional heads provide little benefit. Second,
 253 raw attention maps are not reliable indicators of whether the correct computation has been learned.
 254 In several configurations, individual heads do not match any predicted decay pattern. However,

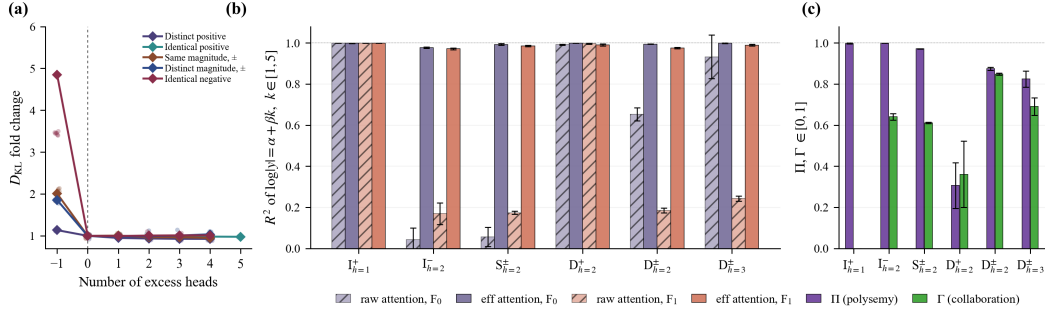


Figure 7: (a) D_{KL} versus the number of excess heads evaluated on 10,240 sequences. The conic minimum head count aligns with the empirical loss plateau across all five configurations. (b) Raw attention maps are not reliable indicators of whether the correct factor-level computation has been learned: some heads match predicted decay patterns, but many do not. Effective subspace attention, which incorporates OV routing and sums across heads, consistently recovers the predicted decay rates. (c) Polysemy and collaboration are common whenever multiple heads are available: heads often contribute to multiple factors, and factors are often implemented by multiple heads. Thus the invariant computation is generally the routed aggregate update, not an individual attention head.

255 effective subspace attention, which incorporates OV routing and sums contributions across heads,
 256 consistently recovers the theoretically predicted factor-level decay rates. Finally, polysemy and
 257 collaboration show that distributed implementations are common. Except in the single-head case,
 258 heads often contribute to multiple factors and factors are often implemented by multiple heads. Thus,
 259 the invariant computation is recovered at the routed aggregate level, not necessarily at the level of
 260 individual attention heads. In appendix C, we see that as we make more heads available than is
 261 needed, signatures of collaboration become even stronger.

262 5 Conclusion

263 We asked when an attention head should be treated as a computational unit. In our setting, the
 264 answer is: the computational unit is at the level of the routed aggregate, not the individual head. The
 265 data-generating spectrum sets a minimum head count that trained models reliably reach across all
 266 five configurations we tested. Above that floor, the same predicted update admits many head-level
 267 decompositions that loss treats as equivalent and per-head inspection cannot distinguish, regardless
 268 of whether they are specialized, collaborative, or polysemantic.

269 What recovers the invariant they share is effective subspace attention: it routes the post-OV, summed
 270 contribution into each factor’s subspace and returns the lag-dependent scalar predicted by the con-
 271 strained belief update, regardless of which heads supplied it. The unit of computation, in our setting,
 272 is the residual-stream subspace into which a layer collectively writes rather than the heads that write
 273 into it. This opens a question for natural-data interpretability: how can we find such subspaces in
 274 scaled models, and in doing so uncover their natural units of computation?

275 6 Limitations and future work.

276 Our analysis is restricted to single-layer transformers trained on processes with known ground truth,
 277 where the computation depends entirely on lag. Extending to multi-layer architectures, richer spectral
 278 structure, and processes beyond conditional independence and exclusively temporal dynamics remains
 279 important future work. More broadly, our current theory is combinatorial in the spectrum—counting
 280 distinct magnitudes and sign-changing modes—whereas the optimization problem faced by a finite
 281 model is also metric. A natural next step is to characterize how approximation quality degrades as
 282 eigenvalues move closer together or toward zero, and to replace a binary notion of “enough heads”
 283 with a quantitative tradeoff between spectral separation, exact implementability, and loss.

284 **References**

- 285 Clémentine C. J. Dominé, Nicolas Anguita, Alexandra M. Proca, Lukas Braun, Daniel Kunin, Pedro
286 A. M. Mediano, and Andrew M. Saxe. From lazy to rich: Exact learning dynamics in deep linear
287 networks. 2025. URL <https://arxiv.org/abs/2409.14623>.
- 288 Joakim Edin, Róbert Csordás, Tuukka Ruotsalo, Zhengxuan Wu, Maria Maistro, Casper L. Chris-
289 tensen, Jing Huang, and Lars Maaløe. Gim: Improved interpretability for large language models,
290 2025. URL <https://arxiv.org/abs/2505.17630>.
- 291 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
292 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli,
293 Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal
294 Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris
295 Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
296 <https://transformer-circuits.pub/2021/framework/index.html>.
- 297 Mario Geiger, Stefano Spigler, Arthur Jacot, and Matthieu Wyart. Disentangling feature and lazy
298 training in deep neural networks. *Journal of Statistical Mechanics: Theory and Experiment*,
299 2020(11):113301, November 2020. ISSN 1742-5468. doi: 10.1088/1742-5468/abc4de. URL
300 <http://dx.doi.org/10.1088/1742-5468/abc4de>.
- 301 Michael Hanna, Ollie Liu, and Alexandre Variengien. How does GPT-2 compute greater-than?:
302 Interpreting mathematical abilities in a pre-trained language model. In *Thirty-seventh Conference*
303 *on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=p4PckNQR8k>.
- 304
- 305 Sarthak Jain and Byron C. Wallace. Attention is not explanation. In *Proceedings of the 2019*
306 *Conference of the North American Chapter of the Association for Computational Linguistics:*
307 *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3543–3556, 2019. URL
308 <https://aclanthology.org/N19-1357/>.
- 309 Jett Janiak, Chris Mathwin, and Stefan Heimersheim. Polysemantic attention head in a 4-layer
310 transformer. *AI Alignment Forum*, 2023.
- 311 Adam Jermyn, Chris Olah, and Tom Henighan. Attention head superposition. *Transformer Circuits*
312 *Thread, Circuits Updates – May 2023*, 2023.
- 313 Harish Kamath, Emmanuel Ameisen, Isaac Kauvar, Rodrigo Luger, Wes Gurnee, Adam Pearce, Sam
314 Zimmerman, Joshua Batson, Thomas Conerly, Chris Olah, and Jack Lindsey. Tracing attention
315 computation through feature interactions. *Transformer Circuits Thread*, 2025.
- 316 Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. Attention is not only a weight:
317 Analyzing transformers with vector norms. 2020. URL <https://arxiv.org/abs/2004.10102>.
- 318 Daniel Kunin, Giovanni Luca Marchetti, Feng Chen, Dhruva Karkada, James B Simon, Michael R
319 DeWeese, Surya Ganguli, and Nina Miolane. Alternating gradient flows: A theory of feature learn-
320 ing in two-layer neural networks. In *The Thirty-ninth Annual Conference on Neural Information*
321 *Processing Systems*, 2025. URL <https://openreview.net/forum?id=t7LKc0MMW6>.
- 322 Maximilian Li and Lucas Janson. Optimal ablation for interpretability. In *The Thirty-eighth Annual*
323 *Conference on Neural Information Processing Systems*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=opt72TYzwZ)
324 [forum?id=opt72TYzwZ](https://openreview.net/forum?id=opt72TYzwZ).
- 325 S. E. Marzen and J. P. Crutchfield. Nearly maximally predictive features and their dimensions. *Phys.*
326 *Rev. E*, 95(5):051301(R), 2017. doi: 10.1103/PhysRevE.95.051301. SFI Working Paper 17-02-007;
327 [arxiv.org:1702.08565](https://arxiv.org/abs/1702.08565) [cond-mat.stat-mech].
- 328 Callum Stuart McDougall, Arthur Conmy, Cody Rushing, Thomas McGrath, and Neel Nanda. Copy
329 suppression: Comprehensively understanding a motif in language model attention heads. In *The*
330 *7th BlackboxNLP Workshop*, 2024. URL <https://openreview.net/forum?id=5Hd6813x3U>.

- 331 Thomas McGrath, Matthew Rahtz, Janos Kramar, Vladimir Mikulik, and Shane Legg. The hydra
332 effect: Emergent self-repair in language model computations. 2023. URL [https://arxiv.org/
333 abs/2307.15771](https://arxiv.org/abs/2307.15771).
- 334 Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? In
335 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, ed-
336 itors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates,
337 Inc., 2019. URL [https://proceedings.neurips.cc/paper_files/paper/2019/file/
338 2c601ad9d2ff9bc8b282670cdd54f69f-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2019/file/2c601ad9d2ff9bc8b282670cdd54f69f-Paper.pdf).
- 339 Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for
340 grokking via mechanistic interpretability. In *The Eleventh International Conference on Learning
341 Representations*, 2023. URL <https://openreview.net/forum?id=9XFSbDPmdW>.
- 342 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan,
343 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli,
344 Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane
345 Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish,
346 and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022.
347 <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- 348 Mateusz Piotrowski, Paul M Riechers, Daniel Filan, and Adam S Shai. Constrained belief updates
349 explain geometric structures in transformer representations. *arXiv preprint arXiv:2502.01954*,
350 2025.
- 351 Paul M Riechers, Thomas J Elliott, and Adam S Shai. Neural networks leverage nominally quantum
352 and post-quantum representations. *arXiv:2507.07432*, 2025.
- 353 Cody Rushing and Neel Nanda. Explorations of self-repair in language models. 2024. URL
354 <https://arxiv.org/abs/2402.15390>.
- 355 Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics
356 of learning in deep linear neural networks. 2014. URL <https://arxiv.org/abs/1312.6120>.
- 357 Adam Shai, Loren Amdahl-Culleton, Casper L. Christensen, Henry R. Bigelow, Fernando E. Rosas,
358 Alexander B. Boyd, Eric A. Alt, Kyle J. Ray, and Paul M. Riechers. Transformers learn factored
359 representations. *arXiv:2602.02385, ICML*, 2026.
- 360 Adam S Shai, Sarah E Marzen, Lucas Teixeira, Alexander Gietelink Oldenziel, and Paul M Riechers.
361 Transformers represent belief state geometry in their residual stream. *NeurIPS*, 2024.
- 362 Aaditya K Singh, Ted Moskowitz, Felix Hill, Stephanie C.Y. Chan, and Andrew M Saxe. What
363 needs to go right for an induction head? a mechanistic study of in-context learning circuits
364 and their formation. In *Forty-first International Conference on Machine Learning*, 2024. URL
365 <https://openreview.net/forum?id=08rrXl71D5>.
- 366 Eric Todd, Millicent Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau.
367 Function vectors in large language models. In *The Twelfth International Conference on Learning
368 Representations*, 2024. URL <https://openreview.net/forum?id=AwyxtyMwaG>.
- 369 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing Multi-Head
370 Self-Attention: Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned. In *Proceedings
371 of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5797–5808,
372 2019.
- 373 Kevin Ro Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.
374 Interpretability in the wild: a circuit for indirect object identification in GPT-2 small. In
375 *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=NpsVSN6o4u1>.
- 377 Sarah Wiegrefe and Yuval Pinter. Attention is not not Explanation. In *Proceedings of the 2019
378 Conference on Empirical Methods in Natural Language Processing and the 9th International
379 Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 11–20, 2019. URL
380 <https://aclanthology.org/D19-1002>.

- 381 Keith Wynroe and Lee Sharkey. Decomposing the qk circuit with bilinear sparse dictionary learning.
382 *AI Alignment Forum*, 2024.
- 383 Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. Hierarchical
384 attention networks for document classification. In *Proceedings of the 2016 Conference of the*
385 *North American Chapter of the Association for Computational Linguistics: Human Language*
386 *Technologies*, pages 1480–1489, 2016. URL <https://aclanthology.org/N16-1174/>.
- 387 Oğuz Kaan Yüksel, Rodrigo Alvarez Lucendo, and Nicolas Flammarion. Incremental learning of
388 sparse attention patterns in transformers. 2026. URL <https://arxiv.org/abs/2602.19143>.
- 389 Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:
390 Metrics and methods. In *The Twelfth International Conference on Learning Representations*, 2024.
391 URL <https://openreview.net/forum?id=Hf17y6u9BC>.

392 A Related Work

393 A.1 Interpretability of Attention

394 How interpretability relates to attention depends strongly on choosing *what* attention is supposed to
395 explain. Early work often drew equivalences between attention weights and token-level explanations:
396 if a token was attended strongly to, that token was considered important for prediction Yang et al.
397 [2016]. This interpretation has been heavily contested. Jain and Wallace [2019] argue that attention
398 weights often do not correlate with other measures of feature importance, and that alternative attention
399 distributions can often produce similar outputs. Wiegrefe and Pinter [2019] respond that attention
400 *can* still be explanatory under some definitions, but only when the relationship between attention
401 and the rest of the model is respected. A reasonable conclusion is that attention maps are useful
402 diagnostic objects, but not faithful explanations by default, and that interpretation should take the rest
403 of the model into account.

404 Mechanistic interpretability shifts the question away from focusing on individual tokens to instead
405 whether attention heads implement specific, identifiable computations. This is a more complete
406 framing because a head is not just a pattern: its behavior depends also on value vectors, output
407 projections, and downstream effects in the residual stream [Elhage et al., 2021, Kobayashi et al.,
408 2020]. Work on induction heads Olsson et al. [2022], the IOI circuit Wang et al. [2023], greater-than
409 circuits Hanna et al. [2023], modular addition Nanda et al. [2023], copy suppression McDougall et al.
410 [2024], and function vectors Todd et al. [2024] shows that some heads participate in relatively clean
411 mechanisms. In these cases, heads are interpretable not as isolated components but as participants
412 in larger circuits. Individual heads can be polysemantic, exhibiting multiple task-relevant attention
413 patterns in different contexts Janiak et al. [2023]. More generally, attentional features may be
414 represented in superposition across heads rather than localized to a single head Jermyn et al. [2023].
415 Recent work therefore decomposes attention computation at a finer granularity, either by tracing
416 interactions between attention features Kamath et al. [2025] or by decomposing the QK circuit into
417 sparse query-key feature pairs Wynroe and Sharkey [2024]. This suggests that heads are useful units
418 of analysis, but not always the most fundamental units of explanation.

419 However, even successful circuit analyses complicate the claim that heads are clean. Many heads can
420 be pruned with little effect, suggesting redundancy [Michel et al., 2019, Li and Janson, 2024, Voita
421 et al., 2019]. In more detailed analyses, heads often collaborate, compensate, or inhibit one another
422 Zhang and Nanda [2024]. The IOI circuit, for example, includes name-mover heads, backup heads,
423 and negative heads, rather than a single decisive component. Similarly, work on self-repair [Rushing
424 and Nanda, 2024, McGrath et al., 2023, Edin et al., 2025] shows that ablating one head can change
425 the behavior of other components which then compensate, and research within training dynamics
426 suggests that a transformer may consist of many such circuits Singh et al. [2024].

427 The emerging view is therefore mixed. Attention weights alone are weak explanations, and individual
428 heads should not be presumed to correspond to clean, independent mechanisms. However, attention
429 heads can still be meaningful units when the full circuit is accounted for.

430 **B Effective Attention**

431 The constrained belief update for factor n at destination position d is:

$$\mathbf{r}_n^{(z_{1:d})} = \boldsymbol{\pi}_n + \sum_{s=1}^d \left(\boldsymbol{\pi}_n T_n^{z_s^{(n)}} T_n^{d-s} - \boldsymbol{\pi}_n \right). \quad (8)$$

432 Each displacement can be written as:

$$\boldsymbol{\pi}_n T_n^{z_s^{(n)}} T_n^{d-s} - \boldsymbol{\pi}_n = \zeta_n^{d-s} \cdot \mathbf{g}_n \left(z_s^{(n)} \right), \quad (9)$$

433 where $\mathbf{g}_n(z) = \left(\boldsymbol{\pi}_n T_n^z - \boldsymbol{\pi}_n \right) P_{\zeta_n, n}$ is the token-dependent displacement direction (independent
434 of offset), and ζ_n is the non-trivial eigenvalue of T_n . The eigenvalue decay ζ_n^{d-s} is the scalar we wish
435 to recover.

436 The displacement map f_n maps the summed per-source OV contribution into factor n 's constrained
437 belief update. For a given (destination d , source s) pair:

$$f_n \left(\sum_h A_{d,s}^{(h)} \mathbf{v}_s^{(h)} \right) \approx \zeta_n^{d-s} \cdot \mathbf{g}_n \left(z_s^{(n)} \right). \quad (10)$$

438 To isolate the decay, we project onto the known $k=0$ displacement direction $\mathbf{g}_n \left(z_s^{(n)} \right)$, which
439 corresponds to offset zero where $\zeta_n^0 = 1$:

$$\alpha_n(d, s) = \frac{\left\langle f_n \left(\sum_h A_{d,s}^{(h)} \mathbf{v}_s^{(h)} \right), \mathbf{g}_n \left(z_s^{(n)} \right) \right\rangle}{\left\| \mathbf{g}_n \left(z_s^{(n)} \right) \right\|^2}. \quad (11)$$

440 This is a scalar-valued function of (d, s) . Under the constrained belief update theory:

$$\alpha_n(d, s) = \zeta_n^{d-s}. \quad (12)$$

441 The denominator normalizes out the token-dependent magnitude, leaving only the offset-dependent
442 decay. This is well-defined whenever $\mathbf{g}_n \left(z_s^{(n)} \right) \neq \mathbf{0}$, which holds for any token z that is informative
443 about factor n .

444 By linearity of f_n , the effective attention decomposes additively across heads:

$$\alpha_n(d, s) = \sum_h \alpha_n^{(h)}(d, s), \quad (13)$$

445 where

$$\alpha_n^{(h)}(d, s) = \frac{\left\langle f_n \left(A_{d,s}^{(h)} \mathbf{v}_s^{(h)} \right), \mathbf{g}_n \left(z_s^{(n)} \right) \right\rangle}{\left\| \mathbf{g}_n \left(z_s^{(n)} \right) \right\|^2}. \quad (14)$$

446 Under head specialization, $\alpha_n^{(h)}(d, s) \approx \zeta_n^{d-s}$ for exactly one head h and ≈ 0 for all others. In the
447 compositional specialized case, multiple heads may contribute to the same factor while remaining
448 factor-specific. Under polysemantic attention, no individual $\alpha_n^{(h)}$ matches ζ_n^{d-s} , but their sum does.

449 **C Measuring polysemy and collaboration**

450 The per-head effective attention (14) also defines a natural coupling matrix $C \in \mathbb{R}^{H \times N}$ that can be
451 used to quantify how factor information is distributed across heads:

$$C_{h,n} \propto \mathbb{E}_{z,k} \left[\left| \bar{\alpha}_n^{(h)}(z^{(n)}, k) \right|^2 \right]. \quad (15)$$

452 where

$$\bar{\alpha}_n^{(h)}(z^{(n)}, k) = \mathbb{E}_{(d,s):d-s=k, z_s^{(n)}=z} [\alpha_n^{(h)}(d, s)] \quad (16)$$

453 The average effective attention $\bar{\alpha}_n^{(h)}(z, k)$ describes the average attention head h pays to subtoken
 454 $z^{(n)}$ at a position k steps back from the query token. The average is over positions consistent with this
 455 lag, and sequences conditioned on a particular subtoken identity at this position. Since the effective
 456 attention can be negative, the coupling matrix measures the average attention power relevant to factor
 457 n within head h . From this matrix, we can define two probability distributions:

$$P(n | h) = \frac{C_{hn}}{\sum_{n'} C_{hn'}} \quad (17)$$

$$P(h | n) = \frac{C_{hn}}{\sum_{h'} C_{h'n}} \quad (18)$$

458 which characterize the amount a given head is responsible for a factor and the amount a given factor
 459 is handled by a particular head, respectively. We can also define relative weightings

$$w_n = \frac{\sum_h C_{hn}}{\sum_{n'h'} C_{hn'}} \quad (19)$$

$$w_h = \frac{\sum_n C_{hn}}{\sum_{n'h'} C_{h'n}} \quad (20)$$

460 that quantify the overall magnitude of a given head or factor’s contribution to the coupling matrix.
 461 The entropy $\mathbf{H}[P]$ (not to be confused with number of heads H) of distribution P quantifies how
 462 peaked or spread a distribution is. Using the distributions in (17)- (18), we define

$$\Pi = \frac{1}{\log N} \sum_{h=1}^H w_h \mathbf{H}[P(n | h)] \quad (\text{polysemy}) \quad (21)$$

$$\Gamma = \frac{1}{\log H} \sum_{n=1}^N w_n \mathbf{H}[P(h | n)] \quad (\text{collaboration}), \quad (22)$$

463 The polysemy score Π quantifies the extent to which a given head attends to multiple factors on
 464 average. It is zero when heads are responsible for one factor only and one when heads attend to all
 465 factors evenly. Similarly, the collaboration score Γ quantifies the fraction of total heads responsible
 466 for a particular factor on average. It is zero when each factor is handled by a single head and one
 467 when each factor is handled evenly by all heads. Finally, we define the effective number of heads for
 468 a particular factor

$$N_{\text{eff}}(n) = \exp(\mathbf{H}[P(h | n)]) \quad (23)$$

469 and its weighted average

$$\bar{N}_{\text{eff}} = \sum_n w_n N_{\text{eff}}(n) \quad (24)$$

470 C.1 Π , Γ , and N_{eff} as a function of H

471 Figure C.1 plots the metrics Π , Γ , and N_{eff} as for a varying number of total heads.

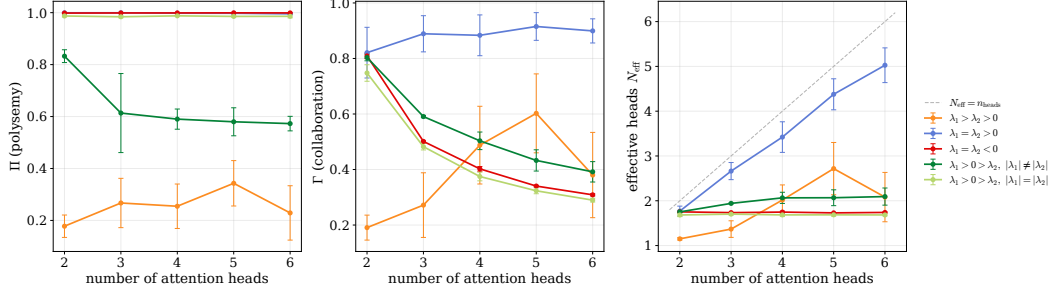


Figure 8: Polysemy Π , collaboration Γ , and effective number of heads as a function of total number of heads in the model. In each run, we fix d_{model} and vary n_{heads} which determines d_{head} . Error bars reflect the 95% confidence intervals over 8 random seeds.

472 D Skip-connection shares degrees of freedom with diagonal attention

473 The residual stream at position d after the attention layer is:

$$474 \quad \mathbf{x}_{\text{mid}}^{(d)} = \underbrace{\mathbf{x}_{\text{pre}}^{(d)}}_{\text{skip connection}} + \sum_h \left[A_{d,0}^{(h)} \mathbf{v}_{\text{BOS}}^{(h)} + A_{d,d}^{(h)} \mathbf{v}_d^{(h)} + \sum_{s=1}^{d-1} A_{d,s}^{(h)} \mathbf{v}_s^{(h)} \right]. \quad (25)$$

474 Projecting into factor n 's subspace, this must equal the constrained belief update (Eq. 4). The three
475 terms on the right-hand side of Eq. 25 correspond to three groups of contributions:

476 **Non-local sources** ($1 \leq s < d$). These contribute the offset $k \geq 1$ displacements and enter
477 exclusively through attention. The spectral theory applies unambiguously:

$$\sum_h A_{d,s}^{(h)} f_n(\mathbf{v}_s^{(h)}) = \zeta_n^{d-s} \cdot \mathbf{g}_n(z_s^{(n)}). \quad (26)$$

478 **BOS token** ($s = 0$). The BOS token encodes a share of the prior π_n through attention to position 0:

$$\sum_h A_{d,0}^{(h)} f_n(\mathbf{v}_{\text{BOS}}^{(h)}) = \gamma_n \cdot \pi_n \quad (27)$$

479 for some scalar γ_n that may vary with position d but not with token identity.

480 **Current token** ($s = d$) **and skip connection**. The offset $k = 0$ displacement $\mathbf{g}_n(z_d^{(n)})$ and the
481 remaining share of the prior $(1 - \gamma_n) \cdot \pi_n$ must be jointly delivered by two pathways:

$$f_n(\mathbf{x}_{\text{pre}}^{(d)}) + \sum_h A_{d,d}^{(h)} f_n(\mathbf{v}_d^{(h)}) = (1 - \gamma_n) \cdot \pi_n + \mathbf{g}_n(z_d^{(n)}). \quad (28)$$

482 The left-hand side has two terms but only their sum is constrained. The model is free to distribute the
483 $k=0$ displacement across these two pathways without affecting the total. This is a degree of freedom
484 that does not exist at any other offset.

485 D.1 Experimental test: one-hot embeddings

486 If the token embedding carries part of $\mathbf{g}_n(z_d^{(n)})$ through the skip connection, then removing this
487 information pathway should force the model to compensate through diagonal attention. We test this
488 by comparing two training conditions:

- 489 • **Learned embeddings** (standard): the embedding matrix W_E is trained jointly with all
490 other parameters. The skip connection can carry arbitrary token-dependent information into
491 `resid_mid`.

492 • **One-hot embeddings:** W_E is frozen as a one-hot encoding (identity matrix, zero-padded to
 493 d_{model}). The skip connection passes through a vector that identifies the token but contains no
 494 learned displacement structure. The $k=0$ displacement must be delivered entirely through
 495 diagonal attention.

496 All other parameters are trained normally in both conditions. We train 8 seeds per condition on the
 497 $\zeta_1 = 0.5$, $\zeta_2 = -0.5$ case (as we observe this setting empirically having little diagonal-attention),
 498 holding all hyperparameters fixed.

499 D.2 Results

500 Figure 9 shows the mean attention weights across lag for both conditions. With learned embeddings,
 501 $\bar{A}_{d,d}$ is lower across lag 0, but all other positions match or are higher, reflecting the degree of freedom:
 502 different seeds distribute the $k=0$ displacement differently between the embedding and diagonal
 503 attention. With one-hot embeddings, $\bar{A}_{d,d}$ is consistently higher.

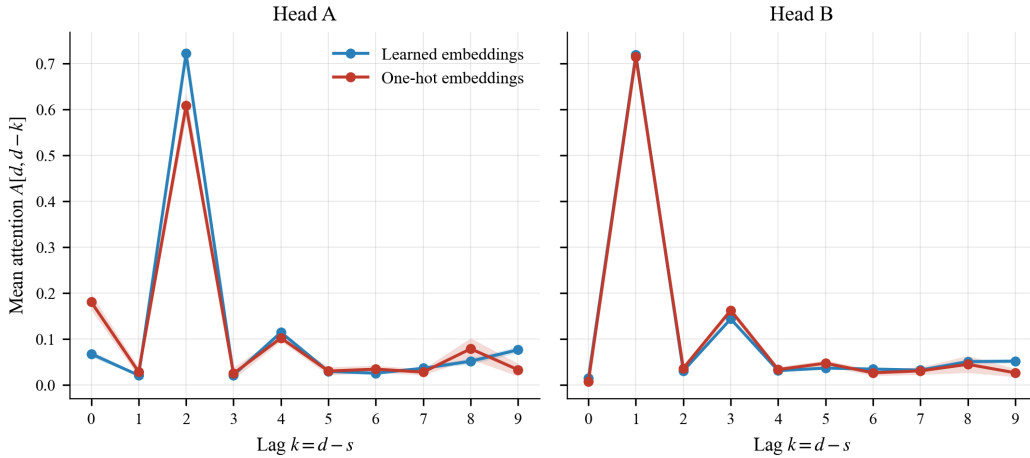


Figure 9: Mean attention-weights grouped by lag across seeds for learned vs. one-hot embeddings. Learned embeddings show lower diagonal attention across all seeds for head A, which we take to be the head implementing the $\zeta_1 = 0.5$ and the even offsets of $\zeta_2 = -0.5$. Notably, scores at all other lags are essentially an exact match. The only difference is in the diagonal attention. One-hot embeddings force higher diagonal attention because the model cannot use the embedding to store additional information. This isolates the $k = 0$ degree of freedom: the current-token displacement can be shared between the learned embedding/skip pathway and diagonal attention, while non-local attention patterns remain unchanged.

504 These results support Eq. 28: the skip connection and diagonal attention are interchangeable pathways
 505 for the $k=0$ displacement, and the model may exploit this freedom when learned embeddings are
 506 available.

507 E Factor Subspace Identification and Quantification

508 We use the vary-one routine from Shai et al. [2026] to identify subspaces.

509 Because our data generator is a product of two independent processes, we can sample datasets that
 510 vary one factor’s subsequence while keeping the other fixed. For each factor, we fix the other to a
 511 single realization, collect activations across many realisations of the varied factor, and mean-center
 512 them per position. The top 2 principal components then form an orthonormal basis for that factor’s
 513 subspace: per-position centering cancels everything that is constant across the batch leaving only the
 514 variation induced by the varied factor, and since Mess3 beliefs are 2-dimensional, two components
 515 span the belief simplex.

516 We define these bases as $\mathbf{V}_{F_n} \subset \mathbb{R}^{d_{\text{model}}}$. We find that these subspaces account for the vast majority of
 517 the variance in the residual stream, explaining on average 97.2% of the variation, and the subspaces

518 \mathbf{V}_{F_1} and \mathbf{V}_{F_2} are effectively orthogonal to each other, with a near-zero overlap. Once identified, we
 519 can define a write-to-fraction that measures how much of the image of the OV circuit lands inside or
 520 outside of each subspace.

521 Once equipped with the subspaces, $\mathbf{V}_{F_n} \subset \mathbb{R}^{d_{\text{model}}}$ we build the orthogonal-complement basis as the
 522 null space of the stacked factor bases $\mathbf{V}_{\perp} \in \mathbb{R}^{d_{\text{model}} \times (d_{\text{model}} - 4)}$ since $[\mathbf{V}_{F_0} \mid \mathbf{V}_{F_1}]$ is rank 4. Finally,
 523 we stack everything into a change-of-basis matrix:

$$\mathbf{B} = [\mathbf{V}_{F_0} \mid \mathbf{V}_{F_1} \mid \mathbf{V}_{\perp}] \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$$

524 When \mathbf{V}_{F_0} and \mathbf{V}_{F_1} are exactly orthogonal, $\mathbf{B} \in O(d_{\text{model}})$ is an orthogonal matrix and we can
 525 decompose any $x \in \mathbb{R}^{d_{\text{model}}}$ into coordinate vectors relative to those subspaces

$$y = \mathbf{B}^{\top} x = [y_0; y_1; y_{\perp}], \quad y_0 \in \mathbb{R}^2, y_1 \in \mathbb{R}^2, y_{\perp} \in \mathbb{R}^{d_{\text{model}} - 4}. \quad (29)$$

526 In d_{model} -dim form, the orthogonal projectors are:

$$P_n = \mathbf{V}_{F_n} \mathbf{V}_{F_n}^{\top}, \quad P_{\perp} = \mathbf{V}_{\perp} \mathbf{V}_{\perp}^{\top} \approx I - P_0 - P_1$$

We note that the \mathbf{V}_{F_n} subspaces are not completely orthogonal to each other, which is why $P_1 + P_2 + P_{\perp}$ is only approximately the identity. This allows the following approximate decomposition.

$$x \approx P_0 x + P_1 x + P_{\perp} x = \mathbf{V}_{F_0} y_0 + \mathbf{V}_{F_1} y_1 + \mathbf{V}_{\perp} y_{\perp}$$

527 With these definitions in hand, we can clearly measure how much each head’s OV circuit writes to
 528 different interpretable subspaces of the residual stream.

$$f_{\text{out}}(h, j) \equiv \frac{|\text{OV}_h \mathbf{b}_j|_F^2}{|\text{OV}_h|_F^2} \quad (30)$$

529 where $\text{OV}_h \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ is head h ’s OV matrix and $\mathbf{b}_j \in \mathbb{R}^{d_{\text{model}}}$ is the j -th column of \mathbf{B} , lying
 530 in $\mathbf{V}_{F_0}, \mathbf{V}_{F_1}$, or \mathbf{V}_{\perp} depending on which block index j falls into. Summing per-direction within
 531 a block recovers the subspace fraction: $f_{\text{out}}(h, n) = \sum_{j \in \mathbf{V}_{F_n}} f_{\text{out}}(h, j)$ which we interpret as the
 532 percent of head h ’s total write-to budget that lands in factor n ’s subspace.

533 The above assumes a uniform-input distribution, but this can be replaced with an empirical average
 534 over real data that also folds in attention. We define head h ’s contribution to the residual stream for
 535 batch b at position d

$$r_h(b, d) = \sum_{s \leq d} A_{d,s}^{(h)}(b) v_{b,s} W_O^{(h)} \in \mathbb{R}^{d_{\text{model}}} \quad (31)$$

536 where $v_{b,s} \in \mathbb{R}^{d_{\text{head}}}$ are the post LayerNorm inputs from the residual stream projected through $W_V^{(h)}$.
 537 From this, we can define

$$f_{\text{realized}}(h, j) = \frac{\sum_{b,d} (\mathbf{b}_j^{\top} r_h(b, d))^2}{\sum_{b,d} |r_h(b, d)|_2^2} \quad (32)$$

538 and again, by summing over $j \in \mathbf{V}_{F_n}$, we can find $f_{\text{realized}}(h, n)$.

539 When looking for a measure that is sensitive to direction, we use $\mathbf{M}_h^{(n)} = \text{OV}_h \mathbf{V}_n \in \mathbb{R}^{d_{\text{model}} \times 2}$
 540 which is the OV map restricted to the rank-2 output subspace of \mathbf{V}_{F_n} . Its squared Frobenius norm is
 541 the unnormalized write-to- \mathbf{V}_{F_n} energy: $|\mathbf{M}_h^{(n)}|_F^2 = |\text{OV}_h \mathbf{V}_{F_n}|_F^2$ (the numerator of $f_{\text{out}}(h, n)$). We
 542 then look at the cosine between two heads’ write operators, computed in matrix-Frobenius geometry:

$$s(h, h', n) = \frac{\langle \mathbf{M}_h^{(n)} \mathbf{M}_{h'}^{(n)} \rangle_F}{|\mathbf{M}_h^{(n)}|_F \cdot |\mathbf{M}_{h'}^{(n)}|_F}; \quad (33)$$

543 where $\langle \mathbf{A}, \mathbf{B} \rangle_F$ is the inner product of the matrices viewed as vectors in $\mathbb{R}^{d_{\text{model}} \cdot 2}$.

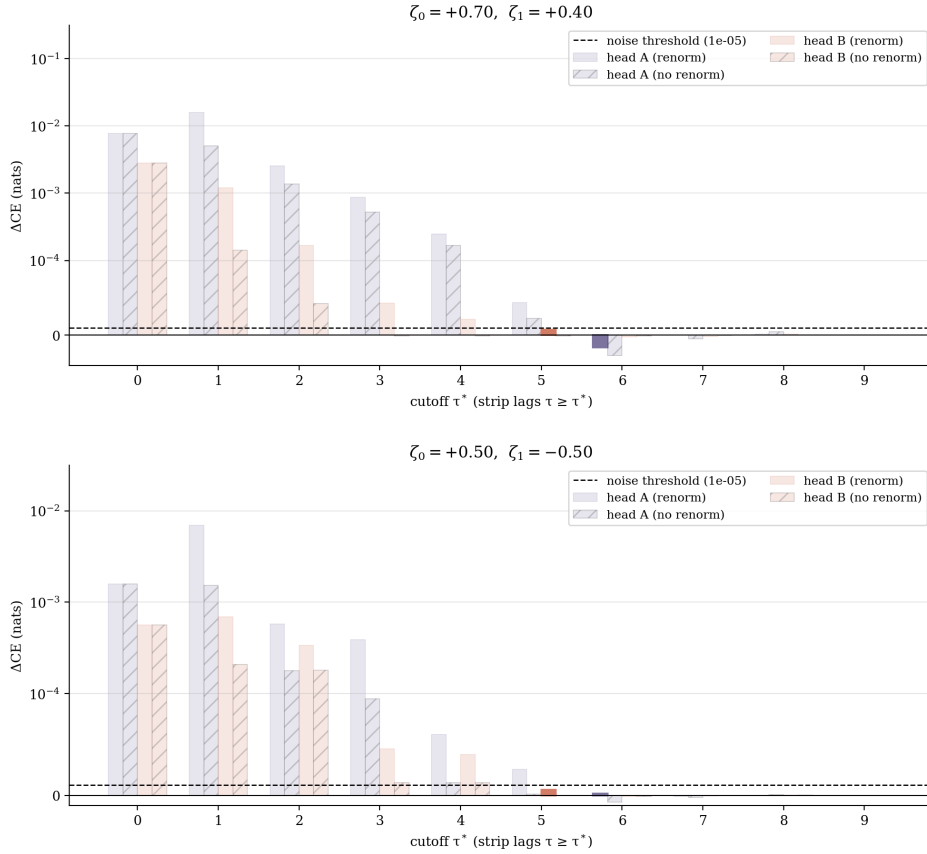


Figure 10: **The importance of attention weights to loss decay with lag.** We measure the cumulative impact on loss of ablating attention weights $A_{d,s}$ for all $\tau = d - s \geq \tau^*$. We then use this to determine the threshold τ for our attention decay fits. Bar heights are mean values over 8 model seeds.

544 **F Additional Experiment Details**

545 See Figs. 10 and 11 for details.

546 **G Weight initialization and attention structure**

547 The experiments in the main text use an initialization range of 0.02 (uniform $\mathcal{U}(-0.02, 0.02)$ for
 548 all weight matrices). Here we investigate how initialization scale affects both the final loss and the
 549 structure of the learned attention patterns.

550 **G.1 Background: rich vs. lazy learning**

551 The scale of weight initialization determines whether a neural network operates in the *rich* (feature
 552 learning) or *lazy* (kernel) regime Dominé et al. [2025] Geiger et al. [2020]. At small initialization,
 553 the network starts with negligible structure and the optimizer builds representations from scratch
 554 — weights move far from their initial values and the learned features bear no resemblance to the
 555 random initialization. At large initialization, the initial random structure is amplified: gradients refine
 556 rather than replace the existing weight configuration, and the network stays closer to its initialization
 557 throughout training.

558 For attention, this distinction has a concrete consequence. At small initialization, the QK dot products
 559 are near zero, so softmax produces *near-uniform attention patterns*. The optimizer has an easier time

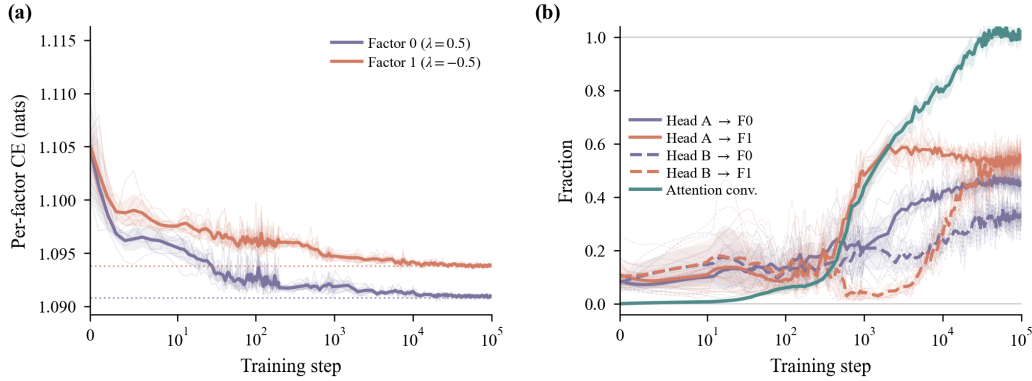


Figure 11: For $\zeta_1 = 0.5$ and $\zeta_2 = -0.5$, we see in (a) that loss initially reduces most for the positive factor. However, (b) shows that the attention heads converge to a polysemantic and collaborative solution quickly, which they must to further reduce loss for both factors.

560 changing these to match the spectral structure predicted by the constrained belief update framework.
 561 At large initialization, the QK dot products are large, so softmax produces *sharp, random attention*
 562 *patterns* from the start. The OV circuit adapts to make use of whatever temporal structure the
 563 random initialization provides, but the QK circuit — whose gradients through softmax are small
 564 when attention is already saturated — has limited ability to reshape these patterns.

565 G.2 Experimental setup

566 We compare two initialization scales across all five configurations from Table 1:

- 567 • **Small init** (init range = 0.02): the setting used throughout the main text.
- 568 • **Default init** (TransformerLens default): substantially larger initialization scale.

569 All other hyperparameters are held fixed. For each configuration and initialization scale, we train 8
 570 models with different random seeds.

571 G.3 Results: loss

572 Table 2 compares the final KL divergence between the two initialization scales. Models trained with
 573 small initialization consistently achieve lower loss across all five configurations. The gap is most
 574 pronounced for configurations requiring compositional attention, where the QK circuit must learn
 575 structured non-monotonic patterns that the random initialization does not provide.

Table 2: Final D_{KL} (nats, mean \pm std across seeds) for small vs. default initialization. H : number of heads. The ratio column shows how many times worse the large initialization is. Configurations requiring compositional attention (negative eigenvalues) are most affected. Not only is loss worse on average, the variance is much more extreme.

Configuration	H	Small init (0.02)	Default init	Ratio
Identical positive	1	0.00039 \pm 0.00001	0.00045 \pm 0.00001	1.14 \times
Distinct positive	2	0.00018 \pm 0.00001	0.00034 \pm 0.00004	1.96 \times
Distinct mag., \pm (2H)	2	0.00030 \pm 0.00001	0.00040 \pm 0.00001	1.36 \times
Distinct mag., \pm (3H)	3	0.00029 \pm 0.00001	0.00049 \pm 0.00008	1.69 \times
Same magnitude, \pm	2	0.00028 \pm 0.00001	0.00091 \pm 0.00026	3.29 \times
Identical negative	2	0.00016 \pm 0.00000	0.00106 \pm 0.00062	6.78 \times

576 **G.4 Results: raw and effective attention**

577 The loss difference is explained by the quality of the learned temporal dynamics. Figure 12 shows
 578 our predictive framework applied to higher initialization.

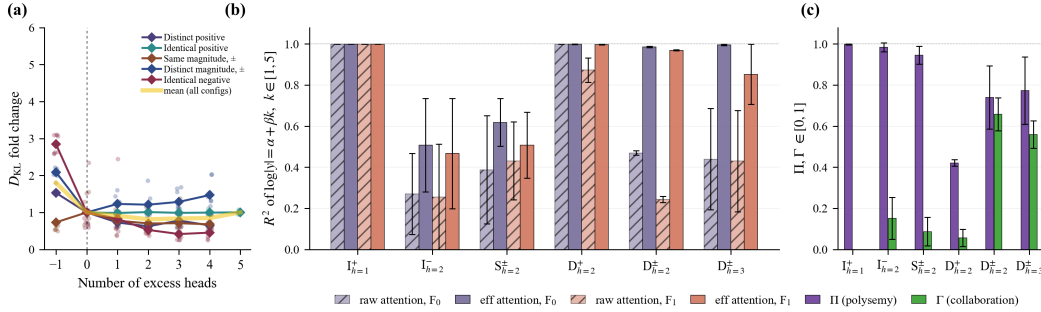


Figure 12: At larger initializations, the model fails to converge to optimal solutions. (a): Compared to the main text, we here see a much more inconsistent pattern, where sometimes adding additional heads will even worsen the model, as it gets stuck in a worse setting during optimization. (b): Effective-attention reveals that for some configurations, the models simply do not learn the correct patterns. We see that the cases with the worst matches are exactly the ones with the largest increase in loss from table 2. (c): Our metrics show evidence of much more degenerate solutions. For the 2 worst cases, we have high polysemy yet low collaboration, indicating a dead head in the model and therefore effectively an underparameterized architecture.

579 With small initialization, attention patterns achieve high R^2 to the predicted candidate patterns across
 580 all configurations (consistent with the main text results). With default initialization, the R^2 is lower —
 581 the patterns show exponential-like decay but at rates that deviate from the eigenvalues of the data
 582 generator. This is consistent with the rich/lazy distinction: small initialization allows the QK circuit
 583 to learn the precise spectral structure from scratch, while large initialization locks in approximate
 584 patterns early in training that the optimizer cannot fully correct.

585 **G.5 Interpretation**

586 These results connect to the gauge freedom framework in two ways. First, they show that the degrees
 587 of freedom in how attention distributes computation are not merely theoretical — the initialization
 588 scale shifts which point in the solution space the optimizer finds, with small initialization favoring
 589 specialized solutions and large initialization favoring more entangled configurations where no head
 590 cleanly matches a single predicted pattern.

591 Second, the loss degradation at large initialization suggests that not all points in the solution space are
 592 equally accessible to gradient descent. The spectral theory predicts a family of equivalent solutions,
 593 but the QK parameterization’s inductive bias — shaped by the softmax saturation at large initial
 594 weights — restricts which solutions the optimizer can reach. The gauge freedom is real in the loss
 595 landscape but constrained in practice by the training dynamics.

596 **H Conic vs. eigenmode decomposition**

597 The conic analysis (Appendix J) establishes that the distinct magnitude \pm configuration ($\zeta_1 = +0.70$,
 598 $\zeta_2 = -0.50$) admits solutions at both $H = 2$ (H_{\min} , the conic minimum) and $H = 3$ (H_{\min} (mode),
 599 the eigenmode prediction). Both achieve equivalent loss (Figure 7a), but they implement the compu-
 600 tation in qualitatively different ways.

601 **H.1 Two heads: conic decomposition**

602 At $H = 2$, the model cannot assign one head per eigenvalue. There are two eigenvalues but one is
 603 negative, requiring sign-alternating contributions that a single non-negative pattern cannot produce.
 604 Instead, the model finds a *conic* solution: both heads learn non-standard attention patterns that are

605 individually difficult to interpret, but whose linear combination through the OV coupling recovers the
606 correct per-factor decay rates.

607 Figure 7b shows that for the 2-head case ($D_{h=2}^{\pm}$), the raw attention R^2 is low for both factors. No
608 head cleanly implements either 0.7^k or the even/odd pattern for $(-0.5)^k$. However, the effective
609 subspace attention R^2 is high: the OV coupling successfully unmixes the heads' contributions to
610 recover both ζ_1^{d-s} and ζ_2^{d-s} .

611 This is the conic decomposition in action. The two heads' patterns are non-negative elements in the
612 feasibility cone (Definition 49), and the coupling matrix decomposes them into the target eigenvalue
613 decays. This requires that no head is fully specialized; all heads must collaborate on all factors.

614 H.2 Three heads: eigenmode decomposition

615 At $H = 3$, the model has enough capacity to assign one head per spectral role: one head implementing
616 the positive decay, one implementing the even-offset pattern and one implementing the odd-offset
617 pattern.

618 Figure 7b confirms this: for the 3-head case ($D_{h=3}^{\pm}$), the raw attention R^2 for the positive 0.7 factor
619 is much higher, despite loss being equivalent. The model finds it useful to allocate a head primarily
620 to serving exponential decay. There is, however, also higher variance in this metric than for the
621 other cases. The model effectively has a head more than it needs, meaning there is a larger space
622 of equivalent solutions it can find. It has less pressure to implement the exact attention maps. The
623 hard requirement that 2 heads must implement the negative eigenvalue remains, and as such we need
624 effective attention to recover the contribution to the 2nd factor.

625 H.3 Implications

626 The crucial observation is that both solutions implement the same per-factor belief update. The
627 residual stream at `resid_mid` contains the same constrained predictive geometry in both cases, and
628 both achieve equivalent loss. The difference is purely in *legibility*: the 3-head eigenmode solution
629 admits a slightly more directly interpretable solution, while the 2-head conic solution is completely
630 opaque just from the attention maps.

631 These findings reveal that the opaqueness of raw attention maps can come from multiple sources:
632 in some cases, the model may be taking advantage of degrees of freedom and in others, there may
633 simply exist a minimal decomposition of the mechanism that admits no clearly interpretable map.

634 I The structure of our training data: A process with two non-Markovian 635 factors

636 Throughout this paper, we use non-Markovian training data sampled from an ergodic source composed
637 of two independent factors to assess how the transformer would learn to represent and predict tokens
638 from a simple "world" made of parts. Here each part is a particular three-state HMM from the
639 parametrized Mess3 family.

640 I.1 A Mess3 factor

641 The Mess3 process Marzen and Crutchfield [2017], Shai et al. [2024], Piotrowski et al. [2025] has
642 three hidden states and three observable tokens $\mathcal{Z} = \{0, 1, 2\}$. Since Mess3 is a 3-state HMM, its
643 local predictive vectors live in a 2-simplex. This 2-simplex corresponds to the set of probability
644 distributions over the three hidden states.

645 The Mess3 process is defined by two parameters, which can be interpreted as a latent-transition
646 probability $p \in (0, 1/2]$ and observation fidelity $q \in [0, 1]$, with dependent quantities $\beta = (1 - q)/2$
647 and $y = 1 - 2p$.¹

¹In previous publications, p and q were denoted by x and α respectively. We use p and q here to avoid a serious internal conflict of notation.

648 The labeled transition matrices are:

$$T^{(0)} = \begin{bmatrix} qy & \beta p & \beta p \\ qp & \beta y & \beta p \\ qp & \beta p & \beta y \end{bmatrix} \quad (34)$$

$$T^{(1)} = \begin{bmatrix} \beta y & qp & \beta p \\ \beta p & qy & \beta p \\ \beta p & qp & \beta y \end{bmatrix} \quad (35)$$

$$T^{(2)} = \begin{bmatrix} \beta y & \beta p & qp \\ \beta p & \beta y & qp \\ \beta p & \beta p & qy \end{bmatrix}. \quad (36)$$

649 The net transition matrix

$$T = T^{(0)} + T^{(1)} + T^{(2)} = (q + 2\beta) \begin{bmatrix} 1 - 2p & p & p \\ p & 1 - 2p & p \\ p & p & 1 - 2p \end{bmatrix} \quad (37)$$

650 has a uniform stationary distribution $\boldsymbol{\pi} = \boldsymbol{\pi}T = \frac{1}{3}[1 \ 1 \ 1]$, and a stationary right eigenstate of
651 $\mathbf{1} = T\mathbf{1} = [1 \ 1 \ 1]^\top$.

652 Besides the stationary eigenvalue of 1, the net transition matrix T has eigenvalue $\zeta = 1 - 3p$,
653 which has (algebraic and geometric) multiplicity of 2 for $p \in (0, 1/2] \setminus \{1/3\}$. (Notice that Mess3
654 degenerates into a memoryless fair coin at $p = 1/3$.)

655 The spectral projection operator associated with stationarity is $P_1 = \mathbf{1}\boldsymbol{\pi}$. The spectral projection
656 operator associated with ζ is $P_\zeta = I - P_1 = I - \mathbf{1}\boldsymbol{\pi}$. Note that these satisfy the eigen-relation
657 $P_\lambda T = TP_\lambda = \lambda P_\lambda$, are orthonormal $P_\lambda P_\xi = P_\lambda \delta_{\lambda,\xi}$, and form a decomposition of the identity
658 $\sum_{\lambda \in \Lambda_T} P_\lambda = I$. Accordingly, they form not only a decomposition of the net latent transition
659 operator $T = \sum_{\lambda \in \Lambda_T} \lambda P_\lambda$, but also an eigen-decomposition of powers of T , such that $T^n =$
660 $\sum_{\lambda \in \Lambda_T} \lambda^n P_\lambda = \mathbf{1}\boldsymbol{\pi} + \zeta^n P_\zeta$. Note that T^n corresponds to the cumulative latent transition dynamic
661 when marginalizing over observed symbols for n time steps.

662 I.2 The joint dynamic

663 The minimal HMM capable of generating the probabilistic structure of the training sequences is a
664 9-state HMM obtained via the tensor product of two constituent Mess3 transition matrices:

$$T^{(x)} = T_1^{(z^{(1)})} \otimes T_2^{(z^{(2)})} \quad \text{for } x = z^{(1)} + 3z^{(2)}. \quad (38)$$

665 Practically, training sequences are generated by running two HMMs in parallel and, based on the
666 outputs of the two constituent HMMs, producing a token $x \in \mathcal{X} = \{m\}_{m=0}^8$ at each timestep, where
667 each token corresponds to a unique $(z^{(1)}, z^{(2)})$ pair.

668 J Minimum Architectural Requirement

669 J.1 Minimum Architecture

670 For a factored process, let each factor m have stationary distribution $\boldsymbol{\pi}_m$, transition matrix T_m , and
671 non-stationary decay eigenvalue ζ_m . In Mess3, each factor has a two-dimensional non-stationary
672 displacement space—the plane of the 2-simplex. Since the two non-stationary eigenvalues are equal,
673 any vector in the plane of the simplex is an eigenstate of T_m associated with eigenvalue ζ_m . Thus lag
674 propagation acts as a scalar on the entire displacement space:

$$\boldsymbol{y}_m T_m^k = \zeta_m^k \boldsymbol{y}_m \quad (39)$$

675 for any non-stationary displacement $\boldsymbol{y}_m = \boldsymbol{\eta}_m - \boldsymbol{\pi}_m$ of factor m . Hence, for a source token z_s and
676 lag $k = d - s$, the contribution to factor m can be written as

$$\zeta_m^k \boldsymbol{g}_m(z_s^{(m)}), \quad (40)$$

677 where $\mathbf{g}_m(z_s^{(m)})$ is the token-dependent displacement direction in the factor- m subspace.

678 For two factors, write these two token-dependent displacement directions induced by the source token
679 as

$$\mathbf{u} = \mathbf{g}_1(z_s^{(1)}), \quad \mathbf{v} = \mathbf{g}_2(z_s^{(2)}). \quad (41)$$

680 Because the belief readout is factorized, the residual displacement decomposes as a direct sum of
681 factor-wise displacements. Therefore the required lag- k update has the idealized form

$$\mathbf{p}_k = \zeta_1^k \mathbf{u} + \zeta_2^k \mathbf{v}. \quad (42)$$

682 Using (\mathbf{u}, \mathbf{v}) as coordinates, this becomes the lag-update curve

$$\mathbf{p}_k = (\zeta_1^k, \zeta_2^k) \in \mathbb{R}^2. \quad (43)$$

683 We now ask how many attention heads are needed to express this family of updates under an idealized
684 OV coupling model. The OV circuit of each head defines how strongly that head writes into the
685 displacement subspace of each factor. In the two-factor case, collect these fixed OV couplings into a
686 heads-by-factors matrix

$$B = \begin{pmatrix} b_{1,1} & b_{1,2} \\ b_{2,1} & b_{2,2} \\ \vdots & \vdots \\ b_{H,1} & b_{H,2} \end{pmatrix}. \quad (44)$$

687 The h -th row

$$\mathbf{b}_h = (b_{h,1}, b_{h,2}) \quad (45)$$

688 is the factor-space write direction supplied by the OV circuit of head h . Equivalently, in the original
689 residual space, head h 's relevant OV contribution is idealized as

$$b_{h,1} \mathbf{u} + b_{h,2} \mathbf{v}. \quad (46)$$

690 The attention mechanism supplies a lag-dependent nonnegative scalar $A_h(k) \geq 0$ for each head.
691 Thus the total update expressible by H heads has factor coordinates

$$\mathbf{p}_k = B^\top A(k) = \sum_{h=1}^H A_h(k) \begin{pmatrix} b_{h,1} \\ b_{h,2} \end{pmatrix}, \quad A_h(k) \geq 0. \quad (47)$$

692 This separates the two roles: the OV circuits supply fixed factor-write directions, while the QK/atten-
693 tion circuits supply the lag-dependent nonnegative coefficients. A single head can vary how much it
694 writes as a function of lag, but its relative coupling to the two factors is fixed by its OV row.

695 Therefore, under this conic OV-coupling abstraction, the relevant head-count quantity is not simply the
696 number of eigenvalues or factors. It is the minimum number of OV coupling rays whose nonnegative
697 cone contains the finite lag-update set

$$\{(\zeta_1^k, \zeta_2^k) : k \geq 0\}. \quad (48)$$

698 Equivalently,

$$H_{\min} = \min \{H : \{(\zeta_1^k, \zeta_2^k) : k \geq 0\} \subseteq \text{cone}(b_1, \dots, b_H)\}. \quad (49)$$

699 This is a capacity calculation for the OV geometry. It assumes that the OV maps can realize
700 the required factor-level coupling directions, while the attention mechanism must still learn the
701 corresponding lag-dependent coefficients $A_h(k)$. Thus the transformer need not assign one head to
702 each factor or to each eigenvalue. Instead, it may learn a conic decomposition of the spectral update
703 geometry. Whether the transformer prefers this solution in practice is examined in section H.

704 Applying the rule above to our different configurations yields table 3 (partially copied from the
705 experiments section).

Table 3: Minimum head count under different two-eigenvalue configurations.

Configuration	Lag-update geometry	H_{\min}
$+\zeta_1, +\zeta_2, \zeta_1 \neq \zeta_2 $	Curve in positive quadrant	2
$+\zeta_1, +\zeta_1$	Single ray	1
$+\zeta_1, -\zeta_1$	Two parity rays in half-plane	2
$+\zeta_1, -\zeta_2, \zeta_1 \neq \zeta_2 $	Two parity branches in half-plane	2
$-\zeta_1, -\zeta_1$	One line through origin	2
$-\zeta_1, -\zeta_2, \zeta_1 \neq \zeta_2 $	Curve in alternating quadrants	3

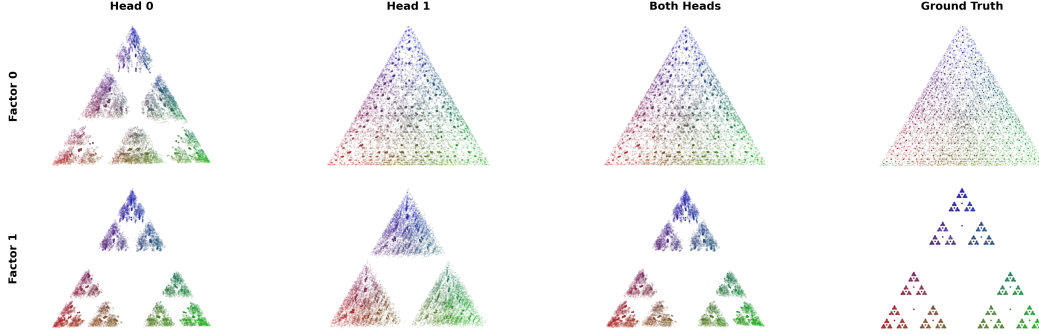


Figure 13: For $\zeta_1 = 0.7, \zeta_2 = 0.4$, regressing from the individual heads + the sum of both heads to the constrained belief simplex shows strong alignment between each head and a specific geometry.

706 K Belief State Geometry

707 To visualize how individual heads contribute to the constrained belief geometry, we regress from
 708 partial attention-layer outputs to the ground-truth factored constrained belief state. Let x_d^{pre} denote the
 709 residual stream before the attention layer at destination position d , and let $o_d^{(h)}$ denote the OV-routed
 710 output of head h at the same position. For a two-head model, we construct three activation sets:

$$x_d^{(1)} = x_d^{\text{pre}} + o_d^{(1)}, \quad x_d^{(2)} = x_d^{\text{pre}} + o_d^{(2)}, \quad x_d^{(1,2)} = x_d^{\text{pre}} + o_d^{(1)} + o_d^{(2)}.$$

711 These correspond to adding only head 1, only head 2, or both heads to the pre-attention residual
 712 stream.

713 For each activation set $S \in \{1, 2, (1, 2)\}$, we fit a linear map

$$\hat{r}_d^S = W_S x_d^S + b_S$$

714 to predict the factored constrained belief state. In the two-factor Mess3 setting, each factor has a
 715 three-dimensional constrained belief vector, so the regression target is

$$r_d^{\text{fact}} = [r_d^{(0)}; r_d^{(1)}] \in \mathbb{R}^6.$$

716 For visualization, we plot the recovered belief for each factor separately. Each three-dimensional
 717 factor belief is mapped to RGB coordinates, so points with similar inferred beliefs have similar colors.
 718 Comparing the regressions from $x_d^{(1)}$, $x_d^{(2)}$, and $x_d^{(1,2)}$ lets us see whether the constrained belief
 719 geometry is linearly recoverable from individual head contributions or only from their sum.

720 For the 2 cases we walked through in the main body, $\zeta_1 = 0.7, \zeta_2 = 0.4$ and $\zeta_1 = 0.5, \zeta_2 = -0.5$,
 721 we show here the per-head and combined regressions to the belief state geometry in figures 13 and 14.

722 L Experiment details

723 L.1 Model hyperparameters

724 We train pre-norm, decoder-only transformers on next-token prediction with cross-entropy loss.
 725 We use a single attention layer, sequence length of 11 (including BOS token), batch size of 128,

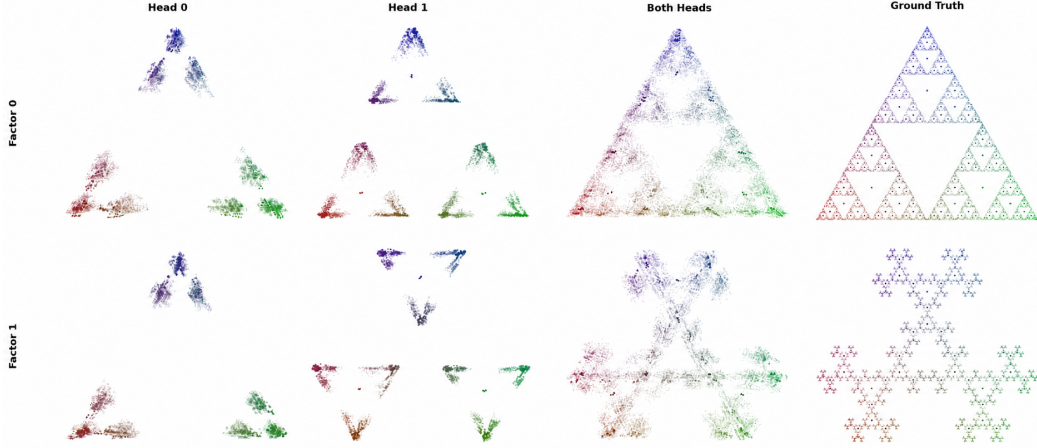


Figure 14: For $\zeta_1 = 0.5$, $\zeta_2 = -0.5$, regressing from the individual heads + the sum of both heads to the constrained belief simplex shows that no head is strongly responsible for specific ground truth geometry. In fact, both geometries require both heads.

726 $d_{\text{model}} = 120$, $d_{\text{head}} = d_{\text{model}} // n_{\text{heads}}$, and $d_{\text{MLP}} = 4 \cdot d_{\text{model}}$. We use the Adam optimizer with a
 727 learning rate of $5e - 4$ and train for 100,000 steps.

728 L.2 Compute

729 We trained on Nvidia H100 GPUs. No model at any point during training took up more than 10GB of
 730 total GPU memory, and every model required less than 20 minutes of training.

731 L.3 Mess3 parameters

Table 4: Mess3 parameters for each two-factor configuration. For each factor, the transition parameter is determined by $p = (1 - \zeta)/3$, and the emission parameter is fixed to $q = 0.6$.

Configuration	ζ_1	ζ_2	p_1	p_2	q_1	q_2
Distinct positive	+0.7	+0.4	0.1	0.2	0.6	0.6
Identical positive	+0.7	+0.7	0.1	0.1	0.6	0.6
Same magnitude, mixed sign	+0.5	-0.5	1/6	0.5	0.6	0.6
Distinct magnitude, mixed sign	+0.7	-0.5	0.1	0.5	0.6	0.6
Identical negative	-0.5	-0.5	0.5	0.5	0.6	0.6

732 **NeurIPS Paper Checklist**

733 **1. Claims**

734 Question: Do the main claims made in the abstract and introduction accurately reflect the
735 paper’s contributions and scope?

736 Answer: [Yes]

737 Justification: [TODO]

738 Guidelines:

- 739 • The answer [N/A] means that the abstract and introduction do not include the claims
740 made in the paper.
- 741 • The abstract and/or introduction should clearly state the claims made, including the
742 contributions made in the paper and important assumptions and limitations. A [No] or
743 [N/A] answer to this question will not be perceived well by the reviewers.
- 744 • The claims made should match theoretical and experimental results, and reflect how
745 much the results can be expected to generalize to other settings.
- 746 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
747 are not attained by the paper.

748 **2. Limitations**

749 Question: Does the paper discuss the limitations of the work performed by the authors?

750 Answer: [Yes]

751 Justification: [TODO]

752 Guidelines:

- 753 • The answer [N/A] means that the paper has no limitation while the answer [No] means
754 that the paper has limitations, but those are not discussed in the paper.
- 755 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 756 • The paper should point out any strong assumptions and how robust the results are to
757 violations of these assumptions (e.g., independence assumptions, noiseless settings,
758 model well-specification, asymptotic approximations only holding locally). The authors
759 should reflect on how these assumptions might be violated in practice and what the
760 implications would be.
- 761 • The authors should reflect on the scope of the claims made, e.g., if the approach was
762 only tested on a few datasets or with a few runs. In general, empirical results often
763 depend on implicit assumptions, which should be articulated.
- 764 • The authors should reflect on the factors that influence the performance of the approach.
765 For example, a facial recognition algorithm may perform poorly when image resolution
766 is low or images are taken in low lighting. Or a speech-to-text system might not be
767 used reliably to provide closed captions for online lectures because it fails to handle
768 technical jargon.
- 769 • The authors should discuss the computational efficiency of the proposed algorithms
770 and how they scale with dataset size.
- 771 • If applicable, the authors should discuss possible limitations of their approach to
772 address problems of privacy and fairness.
- 773 • While the authors might fear that complete honesty about limitations might be used by
774 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
775 limitations that aren’t acknowledged in the paper. The authors should use their best
776 judgment and recognize that individual actions in favor of transparency play an impor-
777 tant role in developing norms that preserve the integrity of the community. Reviewers
778 will be specifically instructed to not penalize honesty concerning limitations.

779 **3. Theory assumptions and proofs**

780 Question: For each theoretical result, does the paper provide the full set of assumptions and
781 a complete (and correct) proof?

782 Answer: [Yes]

783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: **[Yes]**

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887

Answer: [Yes]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: [TODO]

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- 888
- 889
- 890
- 891
- 892
- 893
- 894
- 895
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
 - For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
 - If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

896 8. Experiments compute resources

897 Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

900 Answer: [Yes]

901 Justification: [TODO]

902 Guidelines:

- 903
- 904
- 905
- 906
- 907
- 908
- 909
- 910
- The answer [N/A] means that the paper does not include experiments.
 - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
 - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
 - The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

911 9. Code of ethics

912 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

914 Answer: [Yes]

915 Justification: [TODO]

916 Guidelines:

- 917
- 918
- 919
- 920
- 921
- 922
- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
 - If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
 - The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

923 10. Broader impacts

924 Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

926 Answer: [N/A]

927 Justification: This paper is primarily theoretical and mechanistic. It studies transformers in a controlled synthetic setting where the data-generating process and target computations are analytically specified, with the goal of understanding when attention heads should or should not be interpreted as computational units. The work does not introduce a deployed system, application-specific model, dataset involving people, or new capability intended for real-world decision-making or content generation. We therefore do not identify a direct path to specific positive or negative societal impacts beyond the broad and indirect effects of foundational interpretability research.

935 Guidelines:

- 936
- 937
- 938
- The answer [N/A] means that there is no societal impact of the work performed.
 - If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.

- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

958 11. Safeguards

959 Question: Does the paper describe safeguards that have been put in place for responsible
960 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
961 image generators, or scraped datasets)?

962 Answer: [N/A]

963 Justification: **[TODO]**

964 Guidelines:

- 965
- 966
- 967
- 968
- 969
- 970
- 971
- 972
- 973
- 974
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

975 12. Licenses for existing assets

976 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
977 the paper, properly credited and are the license and terms of use explicitly mentioned and
978 properly respected?

979 Answer: [N/A]

980 Justification: **[TODO]**

981 Guidelines:

- 982
- 983
- 984
- 985
- 986
- 987
- 988
- 989
- 990
- 991
- 992
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 993
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- 994
- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.
- 995
- 996

997 **13. New assets**

998 Question: Are new assets introduced in the paper well documented and is the documentation
999 provided alongside the assets?

1000 Answer: [Yes]

1001 Justification: [TODO]

1002 Guidelines:

- The answer [N/A] means that the paper does not release new assets.
 - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
 - The paper should discuss whether and how consent was obtained from people whose asset is used.
 - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.
- 1003
- 1004
- 1005
- 1006
- 1007
- 1008
- 1009
- 1010

1011 **14. Crowdsourcing and research with human subjects**

1012 Question: For crowdsourcing experiments and research with human subjects, does the paper
1013 include the full text of instructions given to participants and screenshots, if applicable, as
1014 well as details about compensation (if any)?

1015 Answer: [N/A]

1016 Justification: [TODO]

1017 Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
 - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.
- 1018
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025

1026 **15. Institutional review board (IRB) approvals or equivalent for research with human
1027 subjects**

1028 Question: Does the paper describe potential risks incurred by study participants, whether
1029 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1030 approvals (or an equivalent approval/review based on the requirements of your country or
1031 institution) were obtained?

1032 Answer: [N/A]

1033 Justification: [TODO]

1034 Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
 - Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.
- 1035
- 1036
- 1037
- 1038
- 1039
- 1040
- 1041
- 1042
- 1043
- 1044

1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does *not* impact the core methodology, scientific rigor, or originality of the research, declaration is not required.

Answer: [N/A]

Justification: **[TODO]**

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.