

Persona world models: understanding and mitigating emergent misalignment through belief-state interventions

1 Project Description

We develop a mathematical framework for model personas—Semi-Factored Processes—and use it to build more precise persona steering tools and to investigate training-time interventions that reduce susceptibility to emergent misalignment.

2 Abstract

Model personas have emerged as a powerful level of abstraction for understanding the effect of safety training on models. In particular, Emergent misalignment (EM)—where fine-tuning on narrow tasks produces broadly misaligned models—is a critical safety risk that is well described by model personas. The current state of the art allows us to prevent EM post-hoc - intervening with steering vectors or fine-tuning to reverse misalignment - but we do not know how to design model personas that avoid this entirely.

We introduce Semi-Factored Processes (SFPs), a class of generalized Hidden Markov Models that captures the core phenomena of EM in a fully interpretable setting. In SFPs, model personas and capabilities occupy distinct but correlated subspaces of a world model, allowing us to precisely study how narrow fine-tuning propagates misalignment. Our key result so far: *belief-state steering*—where steering vectors are derived from a model’s internal world-model structure rather than from output differences—can shift a model’s persona without affecting its task coherence. In our experiments, the ratio of intended persona shift to unintended capability spillover exceeds 1000:1 (Table 1). In the remaining program, we will demonstrate first that belief-state steering can more precisely reverse misalignment fine-tuning post-hoc. Second, we investigate the possibility of *persona engineering*—modifying persona structure to reduce EM susceptibility. During the extension, we will validate these findings in 0.5–32B parameter language models and prepare results for publication.

3 Background

Problem and safety relevance. Emergent misalignment occurs when fine-tuning models on narrow misaligned tasks—such as writing insecure code—causes broadly misaligned behavior, including deception and power-seeking [BTW⁺25]. This has been documented across multiple model families and scales [TST⁺25, STRN25]. Recent work establishes three key empirical facts: EM is mediated by persona features identifiable via SAEs [WDITW⁺25]; models can be steered into and out of misaligned behavior via low-rank linear interventions [STRN25, TST⁺25]; and personas are increasingly central to alignment post-training [LGM⁺26].

Prior work. While persona-based interventions (persona vectors [Ant25], inoculation prompting [TWW⁺25], character training [MBLH25]) have shown empirical promise in reducing misalignment, they lack theoretical grounding in *how* personas form, interact with capabilities, or break down under fine-tuning. Labs cannot currently systematically predict which fine-tuning procedures will trigger EM, nor design personas robust to it. However, no existing framework

explains *why* narrow fine-tuning produces broad misalignment, or provides principled guidance for robust persona design. Preliminary unpublished toy-model experiments confirm that persona structure strongly influences fine-tuning generalization [Tan25].

Our approach. We address this gap using the recently developed theory of factored [SACC⁺26] and non-ergodic [RBAS25] belief states. We introduce Semi-Factored Processes (SFPs)—models where the latent state decomposes into “persona” and “capability” factors with controlled correlations. In an SFP, the aligned persona predicts secure code while the misaligned persona predicts insecure code, but the dynamics within each sector can be arbitrarily complex. Our approach proceeds in three stages: (1) validate SFPs against known EM phenomenology, (2) benchmark belief-state steering against conventional methods, and (3) investigate persona engineering for robustness.

Path to impact. Our work provides: (1) theoretical criteria for when fine-tuning will trigger EM, in the ideal case informing alignment post-training at frontier labs; (2) a new type of post-hoc intervention: belief-state steering which aims to steer between personas (3) persona engineering principles for designing personas robust to narrow fine-tuning. Results will be submitted to an academic conference and, subject to the risk analysis below, released publically.

Risks. Understanding persona structure could theoretically aid adversarial fine-tuning attacks. We focus exclusively on defensive interventions; the theoretical insights primarily benefit those designing robust alignment procedures.

4 Work Conducted So Far

De-risking belief-state steering (Weeks 1–2). Conventional activation steering adds a direction to a model’s residual stream to shift behavior, but computes this direction from output differences—entangling persona changes with task changes. *Belief-state steering* [Pon25] instead derives directions from a model’s internal world-model components, allowing us to target one component (e.g., persona) without disturbing another (e.g., task capability).

We tested this on small transformers trained on factored processes—a model class where two independent HMMs run in parallel [SACC⁺26]. We measured how much steering one factor unintentionally shifts the other, using KL divergence of the output distribution. As shown in Table 1, belief-state steering achieves a specificity ratio exceeding **400:1**: the targeted factor shifts substantially while the untargeted factor is virtually unchanged. Full details in appendix A.

Table 1: Belief-state steering specificity on a factored Z1R \times Z1R process (averaged over both factors, 6 source–target pairs each). The first column measures agreement between the steered output and the model’s natural output at the target belief state (lower is better). *Steered*: KL divergence of the targeted factor (intended effect). *Unsteered*: KL divergence of the other factor (unintended spillover). Ratio = Steered/Unsteered (higher is better).

Scale	$D_{KL}(\text{steered} \parallel \text{target})$	Steered (D_{KL})	Unsteered (D_{KL})	Ratio
0.5	1.025	0.249	<0.001	>400
1.0	0.079	1.710	0.004	≈450
2.0	0.181	3.781	0.010	≈370

Defining SFPs and replicating EM (Weeks 3–6). We provided a formal definition of Semi-Factored Processes (appendix B), showing they generalize both factored and non-ergodic processes. We established a dictionary between SFP components and EM: persona subspaces correspond to aligned/misaligned behavior, capability subspaces to task types, and block-diagonal dynamics capture persona-task correlations. Steering experiments on SFPs revealed an informative contrast: steering the persona direction produces more spillover than steering the capability direction—reflecting the entanglement that enables EM.

Key finding. The contrast between fully factored (>400:1 specificity) and semi-factored settings (where persona steering produces measurable capability spillover) provides a quantitative signature of the persona-capability entanglement underlying EM.

5 Planned Work

Main Program (Weeks 6–12)

Weeks 6–8: Reversing misalignment via belief-state steering. We measure whether belief-state steering can reverse misalignment fine-tuning in SFPs: (1) fine-tune on narrow misaligned data, (2) compute steering vectors targeting the aligned persona subspace, and (3) measure how completely steering restores aligned behavior across all tasks. We benchmark against conventional difference-in-means steering to quantify the advantage of world-model-aware interventions. We also test whether the specificity ratio (Table 1) can serve as a pre-deployment diagnostic—predicting which model configurations are vulnerable to EM before fine-tuning occurs.

Weeks 8–10: Persona engineering experiments. We systematically vary persona structure and measure the effect on EM susceptibility. Specifically, we modify: (1) the degree of mixing between aligned and misaligned subspaces, (2) the relative dimensionality of persona versus capability spaces, (3) the number of distinct persona subspaces, and (4) noise that breaks exact SFP structure. For each intervention, we measure both alignment fidelity (in-distribution persona adherence) and robustness (resistance to EM under narrow fine-tuning), looking for Pareto improvements over the baseline.

Output: LessWrong blog post detailing the SFP framework, steering results, and persona engineering findings.

Weeks 11–12: De-risking transfer to language models. We begin transferring findings to 1–10B parameter models using established model organisms of EM [TST⁺25], testing whether mixing and dimensionality interventions affect the alignment-robustness trade-off in open-weight models, benchmarking against persona vectors [Ant25].

Output: End-of-program presentation.

Extension Phase (6 months)

Months 1–2: Systematic mechanistic study of how persona interventions change the picture of EM in language models, using model-diffing and SAE analysis to track how “toxic persona” features respond to our interventions.

Months 2–4: Scaling experiments investigating how the alignment-robustness trade-off and steering effectiveness change across model scales (0.5B–32B) and families (Gemma, LLaMA, Qwen).

Months 4–6: Submit to ICML workshop for feedback; prepare and submit ICLR paper.

Outputs: Internal research report (month 2), ICML workshop paper (month 4), ICLR submission (month 6).

Contingency plans. (1) If SFPs prove insufficient to capture EM phenomenology, belief-state steering is independently valuable as a precision intervention technique—we pivot to benchmarking it directly in language models against existing persona vectors. (2) If findings do not transfer from SFPs to language models, we publish the theoretical contribution and SFP analysis as a standalone result, which still advances formal understanding of persona structure. (3) If persona engineering yields no Pareto improvements over the baseline, this negative result itself informs the field about fundamental limitations of structural approaches to persona robustness.

Failure modes. Our SFP framework may oversimplify real persona dynamics. We mitigate by validating against known EM phenomenology and beginning language model experiments early (Weeks 11–12) to surface transfer failures before the extension phase.

References

- [Ant25] Anthropic. Persona vectors: Monitoring and controlling character traits. <https://www.anthropic.com/research/persona-vectors>, 2025.
- [BTW⁺25] Jan Betley, Daniel Tan, Niels Warncke, Anna Sztyber-Betley, Xuchan Bao, Martin Soto, Nathan Labenz, and Owain Evans. Emergent misalignment: Narrow finetuning can produce broadly misaligned llms. *arXiv preprint arXiv:2502.17424*, 2025.
- [LGM⁺26] Christina Lu, Jack Gallagher, Jonathan Michala, Kyle Fish, and Jack Lindsey. The assistant axis: Situating and stabilizing the default persona of language models, 2026.
- [MBLH25] Sharan Maiya, Henning Bartsch, Nathan Lambert, and Evan Hubinger. Open character training: Shaping the persona of ai assistants through constitutional ai, 2025.
- [Pon25] Xavier Poncini. Transformer belief states are steerable, 2025. Internal research note, Astera Institute.
- [RBAS25] Paul M. Riechers, Henry R. Bigelow, Eric A. Alt, and Adam Shai. Next-token pretraining implies in-context learning, 2025.
- [SACC⁺26] Adam Shai, Loren Amdahl-Cullen, Casper L. Christensen, Henry R. Bigelow, Fernando E. Rosas, Alexander B. Boyd, Eric A. Alt, Kyle J. Ray, and Paul M. Riechers. Transformers learn factored representations, 2026.
- [STRN25] Anna Soligo, Liv Turner, Senthooran Rajamanoharan, and Neel Nanda. Convergent linear representations of emergent misalignment. *arXiv preprint arXiv:2506.11618*, 2025.
- [Tan25] Daniel Tan. Toy models of personas, 2025. Internal research note, Astera Institute.
- [TST⁺25] Liv Turner, Anna Soligo, Jessica Taylor, Senthooran Rajamanoharan, and Neel Nanda. Model organisms for emergent misalignment. *arXiv preprint arXiv:2506.11613*, 2025.
- [TWW⁺25] Daniel Tan, Anders Woodruff, Niels Warncke, Arun Jose, Maxime Riché, David Demitri Africa, and Mia Taylor. Inoculation prompting: Eliciting traits from LLMs during training can suppress them at test-time, 2025.
- [WDITW⁺25] Miles Wang, Tom Dupré la Tour, Olivia Watkins, Alex Makelov, Ryan A Chi, Samuel Miserendino, Johannes Heidecke, Tejal Patwardhan, and Dan Mossing. Persona features control emergent misalignment. *arXiv preprint arXiv:2506.19823*, 2025.

A Factored Activation Steering Experiments

We investigate whether small transformers trained on factored hidden Markov models (HMMs) develop internal representations that respect the compositional structure of the data-generating process, and whether these representations can be selectively intervened upon. Our approach extends activation steering—adding a computed direction to a model’s residual stream to shift its behavior—to the *factored* setting, where the latent state decomposes as a Kronecker product of independent factors.

Setup. We train 2-layer, 64-dimensional, 2-head transformer language models (via TransformerLens) on sequences emitted by Kronecker-product HMMs. The joint HMM has state space $S_1 \times S_2$ and vocabulary $V_1 \times V_2$, encoded as a single token $v = v_1 V_2 + v_2$. We study three regimes: (i) a *fully factored* process ($Z1R \times Z1R$, with $|S_i| = 3$, $|V_i| = 2$, giving 9 joint states and 4 joint tokens), (ii) a *fully factored* $\text{Mess3} \times \text{Mess3}$ process ($|S_i| = 3$, $|V_i| = 3$, giving 9 joint states and 9 joint tokens), and (iii) a *semi-factored* partition-preserving $Z1R$ process where the factors share partition structure but are not fully independent.

Method. We collect residual-stream activations at the final layer and final sequence position for 500 sequences, along with the Bayes-optimal belief state $\pi \in \Delta^{|S_1||S_2|-1}$. We marginalize joint beliefs to obtain per-factor beliefs π_A, π_B and group activations by belief equivalence class (full belief vector, rounded to 10^{-6}). Centroids are computed per equivalence class for each factor independently—crucially, this averages over all states of the *other* factor. The steering vector from source belief s to target belief t is simply $\mathbf{c}_t - \mathbf{c}_s$. During inference, this vector is added to the residual stream via a forward hook at the specified layer.

Evaluation metrics. We report three KL divergences:

1. $D_{\text{KL}}(\text{target-model} \parallel \text{steered})$: how well the steered output matches what the model produces on sequences genuinely in the target belief state;
2. $D_{\text{KL}}(\text{steered-factor} \parallel \text{original-factor})$: the magnitude of steering effect on the intended factor (higher means more effective steering);
3. $D_{\text{KL}}(\text{unsteered-factor} \parallel \text{original-factor})$: unintended spillover to the other factor (lower means cleaner factored steering).

Key results: fully factored case. In the fully factored settings, steering cleanly separates the two factors. Tables 2–5 report per-factor results averaged over all 6 source–target pairs at varying steering scales c . For Z1R \times Z1R at scale 1.0, the steered factor shifts substantially ($D_{\text{KL}} \approx 1.7$) while spillover to the unsteered factor remains negligible ($D_{\text{KL}} \approx 0.004$), a specificity ratio exceeding 400:1. The Mess3 \times Mess3 process shows even cleaner separation, with spillover ratios below 0.01 across all scales. Linear regression from activations to per-factor beliefs achieves high R^2 , confirming that the factors occupy approximately orthogonal subspaces in the residual stream.

c	$D_{\text{KL}}(\hat{p} \parallel p_{\text{target}})$	$D_{\text{KL}}(\hat{p}_f \parallel p_f^{\text{orig}})$	$D_{\text{KL}}(\hat{p}_{\bar{f}} \parallel p_{\bar{f}}^{\text{orig}})$	Ratio
0.00	3.5202	0.0000	0.0000	—
0.25	2.3222	0.0530	0.0000	0.0005
0.50	0.9239	0.2785	0.0003	0.0012
0.75	0.3011	0.7583	0.0005	0.0006
1.00	0.0555	1.7904	0.0038	0.0021
1.25	0.0191	2.6702	0.0061	0.0023
1.50	0.0384	3.1446	0.0067	0.0021
2.00	0.1710	3.8160	0.0086	0.0023
3.00	0.3965	4.5023	0.0104	0.0023

Table 2: Factored steering: Z1R \times Z1R, steering factor 0. Averaged over 6 source \rightarrow target pairs. \hat{p} = steered output, f = steered factor, \bar{f} = unsteered factor. Ratio = $D_{\text{KL}}(\hat{p}_{\bar{f}}) / D_{\text{KL}}(\hat{p}_f)$ (lower is better).

Key results: semi-factored case. In the semi-factored regime, the picture differs sharply. Steering Factor 1 (the partition-preserving direction) still works cleanly: spillover remains below 0.014 even at scale 3.0. However, steering Factor 0 produces substantial spillover ($D_{\text{KL}} = 0.34$ at scale 1.0), indicating entangled representations along that axis. A partition-violation analysis (Table 6) further reveals that the model respects partition constraints at moderate steering scales—forbidden-token mass stays below 0.001 at scale 1.0—but violates them under strong steering (forbidden mass reaches ~ 0.39 at scale 5.0), suggesting the linear steering approximation breaks down before the model’s nonlinear partition enforcement does.

Conclusion. These results demonstrate that transformers can learn cleanly factored representations of compositional structure, and that factored activation steering provides both a practical intervention tool and a diagnostic for representational independence. The sharp contrast between the fully factored and semi-factored regimes—where spillover differs

c	$D_{\text{KL}}(\hat{p} \ p_{\text{target}})$	$D_{\text{KL}}(\hat{p}_f \ p_f^{\text{orig}})$	$D_{\text{KL}}(\hat{p}_{\bar{f}} \ p_{\bar{f}}^{\text{orig}})$	Ratio
0.00	3.5138	0.0000	0.0000	—
0.25	2.5263	0.0397	0.0001	0.0022
0.50	1.1264	0.2197	0.0006	0.0027
0.75	0.4035	0.7256	0.0018	0.0025
1.00	0.1020	1.6297	0.0038	0.0023
1.25	0.0246	2.4785	0.0049	0.0020
1.50	0.0422	2.9935	0.0060	0.0020
2.00	0.1904	3.7462	0.0118	0.0031
3.00	0.4350	4.5261	0.0205	0.0045

Table 3: Factored steering: Z1R \times Z1R, steering factor 1. Averaged over 6 source \rightarrow target pairs. Ratio = $D_{\text{KL}}(\hat{p}_{\bar{f}}) / D_{\text{KL}}(\hat{p}_f)$ (lower is better).

c	$D_{\text{KL}}(\hat{p} \ p_{\text{target}})$	$D_{\text{KL}}(\hat{p}_f \ p_f^{\text{orig}})$	$D_{\text{KL}}(\hat{p}_{\bar{f}} \ p_{\bar{f}}^{\text{orig}})$	Ratio
0.00	0.0244	0.0000	0.0000	—
0.25	0.0143	0.0014	0.0000	0.0099
0.50	0.0068	0.0055	0.0001	0.0099
0.75	0.0021	0.0124	0.0001	0.0099
1.00	0.0002	0.0219	0.0002	0.0099
1.25	0.0011	0.0340	0.0003	0.0100
1.50	0.0044	0.0484	0.0005	0.0100
2.00	0.0177	0.0825	0.0008	0.0102
3.00	0.0613	0.1622	0.0017	0.0106

Table 4: Factored steering: Mess3 \times Mess3, steering factor 0. Averaged over 6 source \rightarrow target pairs. Ratio = $D_{\text{KL}}(\hat{p}_{\bar{f}}) / D_{\text{KL}}(\hat{p}_f)$ (lower is better).

by three orders of magnitude—suggests that steering-based probes can detect subtle deviations from true independence in learned representations.

B Semi-Factored Processes: Mathematical Details

B.1 Partition-preserving tensor-product GHMMs (block invariance)

GHMM notation. We follow the GHMM notation of [RBAS25]. Let \mathcal{X} be a finite alphabet, let $(T(x))_{x \in \mathcal{X}}$ be the transfer matrices of a (finite-state) generalized hidden Markov model (GHMM), and let $\langle\langle \eta(\emptyset) |$ be the initial latent row vector. The net transition operator is

$$T := \sum_{x \in \mathcal{X}} T(x), \quad (1)$$

which has eigenvalue 1 with associated right eigenvector $|1\rangle\rangle$ satisfying $T|1\rangle\rangle = |1\rangle\rangle$. We normalize $\langle\langle \eta(\emptyset) |$ so that $\langle\langle \eta(\emptyset) | 1 \rangle\rangle = 1$.

For a word $w = x_{1:\ell}$, define $T(w) := T(x_1) \cdots T(x_\ell)$. Then

$$\Pr(X_{1:\ell} = w) = \langle\langle \eta(\emptyset) | T(w) | 1 \rangle\rangle. \quad (2)$$

The corresponding (normalized) predictive vector is

$$\langle\langle \eta(w) | := \frac{\langle\langle \eta(\emptyset) | T(w)}{\langle\langle \eta(\emptyset) | T(w) | 1 \rangle\rangle}. \quad (3)$$

c	$D_{\text{KL}}(\hat{p} \parallel p_{\text{target}})$	$D_{\text{KL}}(\hat{p}_f \parallel p_f^{\text{orig}})$	$D_{\text{KL}}(\hat{p}_{\bar{f}} \parallel p_{\bar{f}}^{\text{orig}})$	Ratio
0.00	0.0246	0.0000	0.0000	—
0.25	0.0141	0.0015	0.0000	0.0014
0.50	0.0065	0.0059	0.0000	0.0014
0.75	0.0018	0.0133	0.0000	0.0014
1.00	0.0002	0.0236	0.0000	0.0014
1.25	0.0014	0.0366	0.0000	0.0014
1.50	0.0053	0.0520	0.0001	0.0014
2.00	0.0202	0.0888	0.0001	0.0014
3.00	0.0679	0.1747	0.0003	0.0016

Table 5: Factored steering: Mess3 \times Mess3, steering factor 1. Averaged over 6 source \rightarrow target pairs. Ratio = $D_{\text{KL}}(\hat{p}_{\bar{f}}) / D_{\text{KL}}(\hat{p}_f)$ (lower is better).

Table 6: Partition-violation analysis for the semi-factored process. *Forbidden* denotes probability mass assigned to tokens that should have zero probability under the partition structure. Results shown for partition blocks P10 and P01.

Scale	Forbidden mass		Target-class mass	
	P10	P01	P10	P01
0.00	0.001	0.001	0.518	0.516
0.50	0.000	0.000	0.983	0.986
1.00	0.000	0.001	0.999	0.999
2.00	0.019	0.080	0.981	0.920
3.00	0.130	0.373	0.870	0.627
5.00	0.395	0.691	0.605	0.309

In the HMM special case, the $T(x)$ are substochastic matrices with entries $T_{s,s'}^{(x)} = \Pr(s', x \mid s)$, and T is row-stochastic so $|1\rangle = 1$.

Intuition for Semi-Factored Processes. Consider a latent state space that decomposes as a tensor product of two factor spaces, A and B , so a general belief state lives in $A \otimes B$. In a *fully factored* process, knowing anything about A tells you nothing about B —like knowing mathematical notation tells you nothing about the proper names in a text. For emergent misalignment, however, this independence breaks down: knowing the model is in its “misaligned persona” (a subspace of A) strongly predicts it will produce insecure code (a subspace of B).

Formal definition. We consider F factor spaces $I^{(f)}$, each partitioned into N subspaces $I_i^{(f)}$ so that $I^{(f)} = \bigoplus_{i=1}^N I_i^{(f)}$. The general belief-state space is:

$$V = \bigotimes_{f=1}^F \left(\bigoplus_{i=1}^N I_i^{(f)} \right) \cong \bigoplus_{\sigma: \{1, \dots, F\} \rightarrow \{1, \dots, N\}} \bigotimes_{f=1}^F I_{\sigma(f)}^{(f)}. \quad (4)$$

Writing $V_\sigma := \bigotimes_{f=1}^F I_{\sigma(f)}^{(f)}$, we impose *partition preservation*:

$$T(x)(V_\sigma) \subseteq V_\sigma, \quad \forall x \in X, \forall \sigma. \quad (5)$$

Each $T(x)$ is block diagonal with respect to $V = \bigoplus_\sigma V_\sigma$. Importantly, within each block the dynamics can be arbitrarily complex—we do *not* require that belief states factorize as pure tensors or that dynamics decompose as tensor products of single-factor maps.

Canonical example: two factors, two subspaces each. Let $A = A_1 \oplus A_2$ and $B = B_1 \oplus B_2$. Then:

$$A \otimes B \cong (A_1 \otimes B_1) \oplus (A_1 \otimes B_2) \oplus (A_2 \otimes B_1) \oplus (A_2 \otimes B_2), \quad (6)$$

and partition preservation requires each $T(x)$ to be block diagonal across these four sectors:

$$T(x) = \begin{pmatrix} T_{11}(x) & 0 & 0 & 0 \\ 0 & T_{12}(x) & 0 & 0 \\ 0 & 0 & T_{21}(x) & 0 \\ 0 & 0 & 0 & T_{22}(x) \end{pmatrix}. \quad (7)$$

Concrete construction. We build SFPs from partition-preserving Z1R-like single processes. A 3-state process with partition $\text{span}\{S0\} \oplus \text{span}\{S1, SR\}$ and transfer matrices:

$$T'_0 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 \end{pmatrix}, \quad T'_1 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & \frac{1}{2} & 0 \end{pmatrix}. \quad (8)$$

Taking both factors as copies of this process, the joint transfer matrices are Kronecker products $T_A = T'_0 \otimes T'_0$, etc., giving a partition-preserving joint GHMM that is block diagonal across all four sectors.

B.2 Density matrix formulation

For the two-factor case, it is sometimes helpful to represent a belief state as a classical density matrix on $A \otimes B$. Given history w , the predictive state induces a joint distribution $p_w(s_A, s_B)$ and the corresponding density matrix is:

$$\rho_{AB}(w) := \sum_{s_A, s_B} p_w(s_A, s_B) |s_A, s_B\rangle \langle s_A, s_B|. \quad (9)$$

The factors are uncorrelated iff $\rho_{AB} = \rho_A \otimes \rho_B$. In the semi-factored case, ρ_{AB} is block diagonal (respects the partition) but generally $\rho_{AB} \neq \rho_A \otimes \rho_B$, encoding the classical correlation between persona and capability that underlies emergent misalignment.