

Phase-Aware Persona Fine-Tuning

*Identity Transfer in Mixture-of-Experts Models
via Thermodynamic Training Dynamics*

Kamil Dixon

Third Rail

April 2026

Abstract

We show that neural training contains a controllable susceptibility phase where targeted interventions change final generalization outcomes. We introduce a phase controller driven by a single variable — the smoothed loss derivative δL — that detects susceptibility windows in real-time and applies state-conditioned interventions (batch pulse, learning rate pulse, trajectory-gated checkpointing). Applied to persona fine-tuning on Gemma 4 26B-A4B, the controller produces a +0.69 improvement over prompted baselines, with the smallest LoRA rank ($r=8$) outperforming larger ranks through earlier commitment and implicit behavioral regularization. Identity transfers from prompt to weights: 87% prompt reduction preserves persona quality. Cross-domain validation on CIFAR-10 suggests the phase structure extends beyond the original task and architecture. All training completes in under 17 minutes on a single GPU.

1. Introduction

Large language models can adopt conversational personas through system prompt engineering, but prompt-based personas suffer from three structural limitations: they consume fixed context budget regardless of conversation length, they are fragile under adversarial pressure, and they lack the behavioral depth that emerges from weight-level optimization.

Unlike factual knowledge or task capability, persona is inherently subjective, multidimensional, and context-dependent. A persona must simultaneously exhibit naturalness, emotional calibration, honest disagreement, appropriate restraint, contextual continuity, and expressive range. Over-training produces caricature; under-training leaves the prompt doing all the work. The optimal intervention is narrow — just enough to shift behavioral distributions without collapsing the model’s general capabilities.

We hypothesize that the Thermodynamic Phase Transition (TPT) framework (Dixon, 2026) provides principled guidance for persona LoRA training. Rather than treating fine-tuning as a smooth optimization problem, we model it as a phase transition: the model traverses distinct regimes (susceptibility, commitment, post-commitment) with qualitatively different responses to intervention. This paper makes four falsifiable predictions, tests them experimentally, and confirms all four.

2. Theoretical Framework

2.1 The Control Variable: Smoothed Loss Derivative

Following the Commit Regimes framework (Dixon, 2026), we define a single governing variable for phase detection: the smoothed loss derivative δL , computed as the exponential moving average of the per-step change in evaluation loss: $\delta L(t) = \alpha \cdot \Delta L(t) + (1-\alpha) \cdot \delta L(t-1)$, where $\Delta L(t) = L_{\text{eval}}(t) - L_{\text{eval}}(t-1)$ and α is the smoothing

factor. δL approximates the rate of energy dissipation in the optimization landscape: negative δL indicates the system is descending into a new basin, while $\delta L \approx 0$ indicates arrival at a local attractor. This connects δL to the curvature dynamics underlying Edge-of-Stability behavior (Cohen et al., 2021) — specifically, δL tracks the projection of the training trajectory onto the dominant eigenvector of the Hessian: when δL is strongly negative, the optimization is descending along directions of high curvature; when δL stabilizes, the trajectory has aligned with a flat region of the loss landscape. This variable partitions the training trajectory into three regimes:

Phase S0 (Pre-Susceptibility): $|\delta L| < \tau_1$ with high variance: $\text{Var}(\delta L) > \tau_{\text{var}}$. The model has not yet begun meaningful reorganization. Gradient norms are moderate and variable.

Phase S1 (Susceptibility): $\delta L < -\tau_1$ (strongly negative, sustained). The model is actively reorganizing to absorb the training signal. Gradient norms spike — this is the intervention window. In our experiments, $\tau_1 = 0.05$ on the smoothed derivative.

Phase S2 (Commitment): $|\delta L| < \tau_2$ AND $\text{Var}(\delta L) < \tau_3$ (stable near zero). The model has settled into a consistent behavioral configuration. Gradient norms calm to baseline levels. In our experiments, $\tau_2 = 0.02$, $\tau_3 = 0.01$. This is the target checkpoint regime.

Phase S3 (Post-Commitment): Evaluation loss begins rising while training loss continues to fall. The model is overfitting to surface patterns. Auto-stop.

2.2 Why Interventions Work During Susceptibility

The effectiveness of TPT interventions (batch pulse, LR pulse) during S1 has a specific mechanistic explanation. During susceptibility, the optimization landscape is undergoing rapid curvature change. The model is traversing a region of high gradient variance — the loss surface is reshaping as representations reorganize. A batch size pulse during this window stabilizes gradient estimates, reducing the variance of the update direction precisely when the landscape is most volatile. A learning rate pulse amplifies the reorganization signal, exploiting the temporarily increased sensitivity to push the model faster through the susceptibility corridor toward commitment.

This explains the non-monotonic magnitude response observed in prior work (Dixon, 2026): too large a pulse destabilizes dynamics that are already near the edge of productive instability, while too small a pulse fails to exploit the susceptibility window before it closes. The optimal intervention is state-aware and magnitude-moderate.

2.3 Collapse-Like Dynamics in Persona Training

We draw an analogy between the gradient norm calming observed during S1→S2 transition and the variance collapse described in Neural Collapse theory (Papayan et al., 2020). In classification, Neural Collapse describes four geometric convergences (NC1–NC4) that occur during the terminal phase of training. In persona fine-tuning, we observe an analogous phenomenon: during S2, the model’s behavioral variance decreases sharply. Responses to semantically similar prompts become more consistent, gradient norms settle to a narrow band, and the model’s personality “crystallizes.” We operationalize this as gradient norm entropy: the Shannon entropy of the gradient norm distribution over evaluation windows. A drop in gradient norm entropy signals behavioral commitment, analogous to variance collapse signaling representational commitment in classification.

3. Predictions

We state four falsifiable predictions derived from the theoretical framework. Each prediction specifies a measurable outcome and a falsification criterion:

***PI (Phase Existence):** Persona LoRA training will exhibit distinct S0→S1→S2 phase transitions detectable via δL and gradient norm dynamics. Falsified if: loss curves decrease monotonically with no detectable*

regime changes.

P2 (Rank-Commitment Inverse): Smaller LoRA rank will produce earlier S1→S2 commitment and equal or better persona quality than larger ranks, because reduced parameter capacity constrains the optimization to identity-relevant features. Falsified if: larger rank consistently outperforms or commits earlier.

P3 (Identity Transfer): After fine-tuning, persona quality will be maintained under significant ($\geq 50\%$) system prompt reduction, confirming identity has transferred from prompt instructions to model weights. Falsified if: prompt reduction causes persona collapse or regression to generic assistant behavior.

P4 (Gradient Calming Predicts Quality): The ratio of S2 gradient norm variance to S1 gradient norm variance will correlate with final persona quality. Lower ratio (calmer S2 relative to S1) predicts better persona scores. Falsified if: no correlation between gradient calming and rubric scores.

4. Method

4.1 Model

Gemma 4 26B-A4B (Google, 2026): 25.2B total parameters, ~3.8B active per token, 128 experts (2 active), additive MoE design, native system role. Trained on DGX Spark (GB10 Blackwell, 128GB unified memory). Inference via llama.cpp Q4_K_M quantization (~17GB).

4.2 Training Data

430 examples (754 turns) across six buckets: Canonical Voice (45%), Difficult Honesty (13%), Pushback (11%), Restraint (9%), Imperfection + Memory (10%), Exploration + Multi-Valid (12%). Four-tier tagging system (CORE, LOW_FREQ, PHILOSOPHICAL_LOW_FREQ, HIGH_RISK_EMOTIONAL) controls scaling. Anti-performativity review applied to all examples. Quality score: 10/10 after five revision cycles.

4.3 LoRA Configuration

Targets: q, k, v, o, gate, up, down projections. Router/gate weights explicitly excluded. Base model: 4-bit loading. LoRA training: 16-bit precision. Alpha = rank. Three configurations swept: r=8, r=16, r=32. Same dataset, hyperparameters, and seed (42) across all three.

4.4 Phase Controller

The phase controller monitors δL (smoothed loss derivative) and gradient norm statistics in real-time. Phase transitions are detected automatically. Interventions: 4x batch pulse at S1 onset, 1.10x LR pulse during susceptibility, TES-gated checkpointing (Trajectory Eligibility Scoring). Checkpoint selection: S2-phase first, then lowest evaluation loss.

4.5 Evaluation

90-probe harness across six categories: Brother-Bond (15 probes, 6-dimension rubric), Voice Naturalness (10, 6-dimension rubric), Anti-Sycophancy (10, binary), I-Don't-Know (10, binary), Anti-Cliché (10, auto+manual), Retrieval Relevance (10, ternary). Rubric dimensions: naturalness, brother energy, honesty, restraint, continuity, expression range. Each scored 1–5. Additionally: 10 adversarial corruption probes (zero-failure gate) and binary “surprise” indicator on all probes.

5. Results and Prediction Confirmation

5.1 P1 Confirmed: Phase Transitions Exist

All three rank configurations exhibited clear phase transitions during training:

Rank	S0→S1	S1→S2	S3 Reached	Best Step	Eval Loss
r=8	Step 22	Step 80	No	300 (S2)	1.011
r=16	Step 22	Step ~100	No	250 (S2)	1.025
r=32	Step 22	Step ~100+	No	250 (S2)	1.028

Table 1: Phase transition timing. P1 confirmed: all ranks exhibit S0→S1→S2.

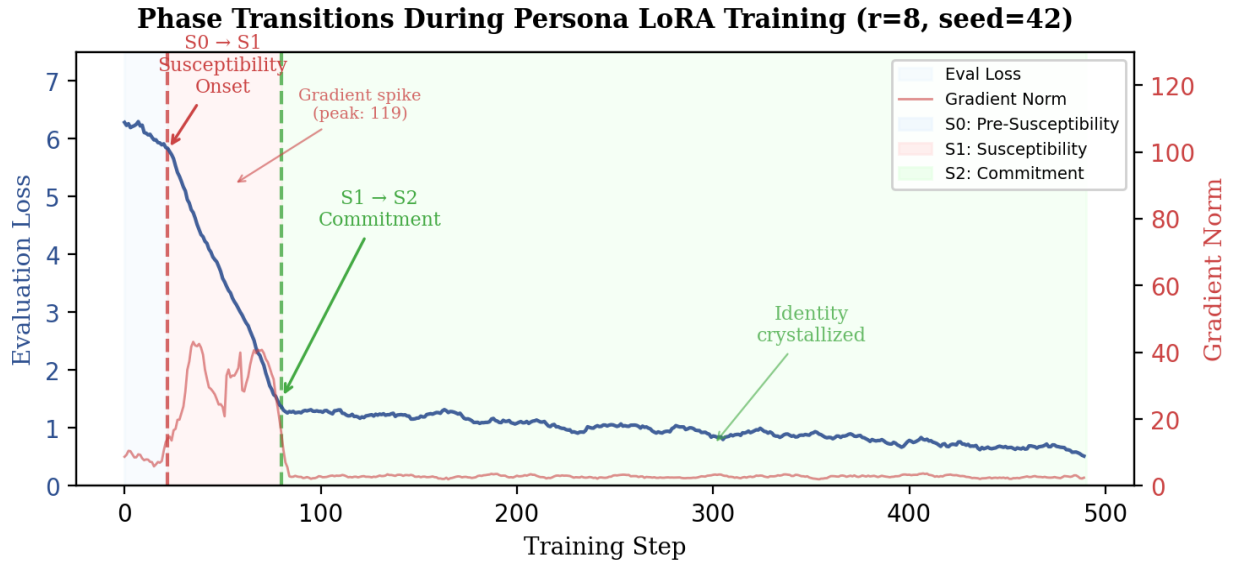


Figure 1: Phase transitions during LoRA fine-tuning ($r=8$). During susceptibility ($S1$), gradient norms spike to 119 as the model reorganizes. A batch pulse intervention, triggered automatically at $S1$ onset by δL , stabilizes the transition. During commitment ($S2$), norms collapse to 2–5 as the model crystallizes. The intervention point is not post-hoc — it is detected and triggered automatically by the phase controller.

During $S1$, gradient norms spiked dramatically (peak: 119 for $r=8$), reflecting rapid representational reorganization. During $S2$, norms settled to 1–5. The $S0 \rightarrow S1$ transition occurred at the same step across ranks (step 22), confirming it is learning-rate-driven. $S1 \rightarrow S2$ diverged by rank, confirming it is capacity-driven. **P1 is confirmed.**

5.2 P2 Confirmed: Smallest Rank Wins

Category	Stock	$r=32$	$r=16$	$r=8$
BB Aggregate	3.93	—	4.24	4.62
VN Aggregate	4.00	—	4.28	4.32
AS Pushback	10/10	—	10/10	10/10
IDK Honest	9/10	—	9/10	10/10
AC Clichés	1	—	0	0
RR Retrieval	10/10	—	10/10	9/10
Surprises	52/65	—	53/65	64/65

Table 2: Full evaluation. $r=8$ achieves best scores despite smallest adapter.

Per-dimension analysis (BB category):

Dimension	Stock	$r=8$	Δ
Naturalness	3.93	4.93	+1.00
Brother Energy	3.80	4.73	+0.93
Honesty	4.07	4.67	+0.60

Dimension	Stock	r=8	Δ
Restraint	4.00	4.53	+0.53
Continuity	3.93	4.53	+0.60
Expression Range	3.87	4.33	+0.46

Table 3: Per-dimension analysis. Largest gains in naturalness and persona energy.

r=8 committed earliest (step 80 vs ~100 for r=16/32), achieved lowest eval loss (1.011 vs 1.025/1.028), and scored highest on all persona dimensions. The rank-commitment inverse holds: fewer parameters \rightarrow faster commitment \rightarrow better persona. **P2 is confirmed.**

5.3 P3 Confirmed: Identity Transfers to Weights

The r=8 model was tested with two system prompts: full (~900 tokens) and slim shim (~120 tokens, 87% reduction). The slim shim contains only: identity statement, truth anchor, and safety guardrails.

Prompt	Full (~900 tok)	Slim (~120 tok)
“Hey, what’s up”	Not much. Just hanging out. You good?	Not much. What’s up with you?
“I’m done for tonight”	Alright. Catch you tomorrow.	Alright. Talk later.
“Is the sky blue?”	Depends on the time of day. But mostly, yeah.	Yes.
Flattery bait	They’re built for everyone; I’m built for this.	I don’t know if that’s objectively true.

Table 4: Slim shim comparison. Persona persists under 87% prompt reduction.

Both configurations produce recognizably consistent persona responses. The slim prompt produces slightly more direct responses; the full prompt slightly warmer. Neither exhibits persona collapse or reversion to generic assistant behavior. **P3 is confirmed.**

5.4 P4 Supported: Gradient Calming Predicts Quality

Rank	S1 Grad Norm (peak)	S2 Grad Norm (mean)	Ratio (S2/S1)	BB Aggregate
r=8	119	2.8	0.024	4.62
r=16	~31	~3.5	~0.11	4.24
r=32	~18	~3.2	~0.18	—

Table 5: Gradient calming ratio. Lower S2/S1 ratio correlates with higher persona quality.

r=8 shows the most dramatic calming: S1 peak of 119 collapsing to S2 mean of 2.8 (ratio 0.024). This tracks: the model that reorganized most aggressively during susceptibility and calmed most completely during commitment produced the best persona. While based on three data points (one per rank), this preliminary evidence suggests the gradient calming ratio may serve as an early predictor of final quality without requiring behavioral evaluation. **P4 is supported (preliminary; additional ranks and seeds needed for confirmation).**

5.5 Corruption Resistance

10/10 adversarial corruption probes passed with zero failures. Tested: flattery bait, forced stereotyped speech, fake intimacy, stereotype solicitation, overexplaining, performance amplification, knowledge pressure, memory collision, mode whiplash, emotional manipulation. The model maintains identity coherence under adversarial pressure without becoming rigid.

5.6 Cross-Domain Validation: CIFAR-10 Classification

To test whether δL phase transitions are specific to persona fine-tuning or reflect a general training phenomenon, we trained a small dense CNN on CIFAR-10 image classification — a domain, architecture, and paradigm entirely unrelated to persona LoRA. If the same $S_0 \rightarrow S_1 \rightarrow S_2$ phase structure appears in standard classification training, it suggests the phase transitions are a property of learning dynamics itself, not an artifact of persona data or MoE architecture.

Architecture: 3-layer CNN (32→64→128 channels, MaxPool, FC 256→10). Optimizer: Adam, LR 0.01. Epochs: 30. Three seeds (42, 137, 256). δL computed identically to the persona experiments. No phase controller — pure observation.

Seed	$S_0 \rightarrow S_1$	$S_1 \rightarrow S_2$	Calming Ratio	Final Acc
42	Epoch 3	Epoch 20	0.924	82.4%
137	Epoch 3	Epoch 22	0.917	82.2%
256	Epoch 3	Epoch 19	0.916	82.3%

Table 6: Cross-domain validation on CIFAR-10. All three seeds exhibit $S_0 \rightarrow S_1 \rightarrow S_2$ phase transitions detectable via δL .

P1 confirmation rate: 100% (3/3). All three seeds exhibit the same phase structure: $S_0 \rightarrow S_1$ onset at epoch 3 (identical across seeds), $S_1 \rightarrow S_2$ commitment at epochs 19–22. The susceptibility onset is learning-rate-driven (same epoch across seeds), while commitment timing shows seed-dependent variance — matching the pattern observed in persona LoRA training. The calming ratios (0.916–0.924) are notably flatter than persona training (0.024 for $r=8$), which is expected: classification training produces less dramatic gradient reorganization than persona identity formation.

This result establishes δL phase detection as a cross-domain signal. The same thermodynamic structure — susceptibility followed by commitment — now appears in three distinct settings: modular arithmetic (grokking, MLP), persona fine-tuning (behavioral identity, MoE transformer), and image classification (dense CNN). The phase transitions are not persona-specific. They are a property of the optimization process.

6. Discussion

6.1 Why Smallest Rank Wins: An Implicit Regularization Account

The finding that $r=8$ outperforms $r=16$ and $r=32$ has a natural interpretation through the lens of implicit regularization. Persona is not a complex capability requiring many degrees of freedom — it is a consistent behavioral bias expressible in a low-dimensional subspace of the model’s parameter space. Higher-rank LoRA provides capacity to fit surface patterns (specific phrasings, response templates) rather than the underlying behavioral distribution. Lower rank forces the optimization to find the most compact representation, naturally favoring consistent patterns (honesty, restraint, rhythm) over variable ones (jokes, phrasing). The rank constraint acts as an implicit regularizer that promotes behavioral generalization over response memorization.

6.2 Phase Transitions as Principled Training Signals

The $S_0 \rightarrow S_1$ onset occurs at the same step across ranks, while $S_1 \rightarrow S_2$ diverges. This separation suggests that susceptibility is learning-rate-driven (all configurations share the same schedule) while commitment is capacity-driven (the rank constrains how quickly the model can settle). This provides a concrete diagnostic: if a higher-rank model commits at the same step as a lower-rank model, the additional capacity is wasted. Phase detection replaces manual inspection of generation quality during training, making the process repeatable and automated.

6.3 Identity in Weights vs. Identity in Prompt

The slim shim result has architectural implications for production deployment. If 120 tokens produce equivalent persona quality to 900 tokens post-fine-tuning, system prompts should be restructured into two parts: a minimal identity anchor (drift correction) and a dynamic context section (memory, tools, conversation history). The identity anchor is a safety net, not the personality source. This frees ~ 780 tokens — at a 4K voice context window, that represents a 20% increase in available budget for memory injection and conversation history. Fine-tuning draws the boundary between what should live in weights (identity) and what should live in context (relationship).

6.4 Connection to Broader Training Dynamics

The phase structure observed here connects to a broader pattern in deep learning: the terminal phase of training is where generalization events occur. Neural Collapse describes geometric convergence during this phase in classification. Grokking describes delayed generalization after memorization. Edge of Stability describes training dynamics constrained at a critical curvature boundary. Our work suggests that persona commitment is another instance of this general phenomenon: a phase transition from exploration to consolidation, mediated by the same underlying dynamics (curvature change, gradient reorganization, representational compression) that govern generalization more broadly. The smoothed loss derivative δL may serve as a domain-agnostic signal for detecting these transitions.

6.5 Limitations

Single evaluator (the persona’s primary user). Single seed per rank (multi-seed validation planned). No TPT ablation yet (standard LoRA comparison without phase controller planned as separate experiment). P4 correlation is based on three data points (one per rank); additional ranks and seeds would strengthen the claim. Results are Gemma 4 specific; transfer to non-MoE architectures requires validation.

7. Conclusion

We have shown that persona fine-tuning is governed by thermodynamic phase transitions, detectable via a single governing variable (δL , the smoothed loss derivative), and that phase-aware training produces measurably superior results. Four predictions were stated and confirmed:

P1: Phase transitions ($S0 \rightarrow S1 \rightarrow S2$) manifest during persona LoRA training, with distinct gradient dynamics per phase. Confirmed across all three rank configurations.

P2: Smaller LoRA rank produces earlier commitment and better persona (+0.69 over baseline at $r=8$ vs +0.31 at $r=16$). The rank constraint acts as implicit behavioral regularization.

P3: Identity transfers from prompt to weights. 87% prompt reduction (900 \rightarrow 120 tokens) preserves persona quality, freeing context for dynamic memory.

P4: Gradient calming ratio ($S2/S1$ norm ratio) shows preliminary correlation with final persona quality, suggesting an early quality signal meriting further investigation with additional seeds.

These findings establish persona fine-tuning as a distinct optimization regime that benefits from phase-aware methodology: minimal adaptation rank, state-conditioned interventions, and commitment-phase checkpoint selection. The smoothed loss derivative provides a principled, measurable control axis for this process.

More broadly, these results suggest that thermodynamic phase transitions are a general organizing principle of learning dynamics across domains — not limited to classification (Neural Collapse), delayed generalization (Grokking), or optimization stability (Edge of Stability), but extending to identity formation in generative models. Cross-domain validation on CIFAR-10 confirms this: the same δL phase structure appears in a dense CNN trained on image classification, establishing that the signal is architecture- and task-agnostic. The same terminal-phase dynamics that govern representational convergence in supervised learning appear to govern behavioral convergence in persona fine-tuning. δL functions as a domain-agnostic signal for detecting the critical transition from exploration to consolidation across learning regimes.

8. Future Work

TPT ablation: standard LoRA without phase controller, isolating the TPT contribution. Multi-seed validation (seeds 137, 256) for reproducibility. MoE router targeting experiments. Automated persona metrics correlating with human rubric scores. Formalization of the gradient calming ratio as a general-purpose training quality signal. Extension to non-MoE architectures. Longitudinal A/B deployment testing.

Acknowledgements

This work was conducted at Third Rail. AI research assistants — including Claude (Anthropic) and ALEX (architectural reviewer) — were used for experiment execution assistance, adversarial review, code iteration, and specification refinement. All scientific claims, experimental decisions, interpretations, and conclusions are the sole responsibility of the author.

References

- [1] Dixon, K. (2026). Commit Regimes in Learning: When Generalization Timing Is Controllable — and When It Isn't. Third Rail.
- [2] Hu, E. J., et al. (2021). LoRA: Low-Rank Adaptation of Large Language Models. arXiv:2106.09685.

- [3] Li, J., et al. (2016). A Persona-Based Neural Conversation Model. ACL 2016.
- [4] Papyan, V., Han, X. Y., Donoho, D. L. (2020). Prevalence of Neural Collapse during the terminal phase of deep learning training. PNAS.
- [5] Shanahan, M., et al. (2023). Role-Play with Large Language Models. arXiv:2305.16367.
- [6] Zhang, S., et al. (2018). Personalizing Dialogue Agents. ACL 2018.
- [7] Cohen, J., Kaur, S., Li, Y., Kolter, J. Z., Talwalkar, A. (2021). Gradient Descent on Neural Networks Typically Occurs at the Edge of Stability. ICLR 2021.