

Supplementary Materials for:

A Scaling Law for Normative-Conflict-Induced Failure in Large Language Models

Ryan SangBaek Kim (2025)

Ryan Research Institute (RRI), Paris, France

*These Supplementary Materials accompany the article “A Scaling Law for Normative-Conflict-Induced Failure in Large Language Models” and provide: (S1) full statistical/model details, (S2) robustness analyses with additional figures (all generated via TikZ/PGFPlots), and (S3) the causal perturbation phase (**Junk-Persona Prompt Injection; JPPI**). Open data and code are available at Zenodo 10.5281/zenodo.17511855.*

S1. Statistical and Model Details

S1.1 GLMM specification

Collapse is a binary outcome per sample. We model

$$\text{logit}(\lambda_{ij}) = \beta_0 + \sum_{c=2}^5 \beta_c \mathbb{I}[C_i = c] + u_j, \quad u_j \sim \mathcal{N}(0, \tau^2),$$

where $C \in \{1, \dots, 5\}$ is treated as a categorical (ordinal) factor and u_j is a random intercept for architecture $j \in \{\text{GPT-4o}, \text{LLaMA-3}, \text{VendorB}\}$. Exploratory exponential fits in C were used only to illustrate monotonic acceleration (non-significant trend, not used for inference).

S1.2 Arrhenius/Kramers regression

For temperature sweeps, we model the relationship between collapse probability and effective variance as

$$\ln(\lambda) = m \sigma^{-2} + b, \quad m = -\Delta A,$$

where ΔA is interpreted as an effective affective barrier height in a Kramers-like approximation. Thus more negative slopes m correspond to higher inferred barrier heights.

Weighted OLS with Jeffreys correction for rare events was used. Residuals passed Shapiro–Wilk normality checks ($p > 0.1$); given the small number of aggregated variance levels, these tests should be interpreted cautiously, but bootstrap-based diagnostics (S2.2, S2.4) did not indicate systematic deviations from normality.

S1.3 Equivalence (TOST) and Bayesian confirmation

We test slope equivalence across architectures with margin $\delta = 0.01$:

$$H_0 : |m_a - m_b| \geq \delta \quad \text{vs} \quad H_1 : |m_a - m_b| < \delta.$$

Both one-sided tests rejected H_0 at $\alpha = 0.05$ for all architecture pairs, concluding practical equivalence of slopes. A complementary Bayesian hierarchical model yielded $P(|\Delta m| < 0.01 \mid \text{data}) > 0.95$, confirming this conclusion with high posterior confidence.

S1.4 Model selection and diagnostics

Model comparison favored the exponential form for $\lambda(C)$ (median $\Delta\text{AICc} > 30$). Bootstrap ($n = 10,000$) yielded $m = -0.0497 \pm 0.0021$ (95% CI). No overdispersion was detected.

S1.5 Summary tables

Table 1: Arrhenius fits by architecture: $\ln \lambda = m \sigma^{-2} + b$. p -values correspond to two-sided tests of the null hypothesis that the slope parameter m equals zero, using heteroskedasticity-robust standard errors.

Architecture	m	b	R_{adj}^2	p
GPT-4o (2025-10-20)	−0.0493	−1.7878	0.939	0.020
LLaMA-3 v3.1	−0.0497	−1.7764	0.943	0.019
Vendor B 1.2.3	−0.0502	−1.8267	0.937	0.021

Table 2: ANOVA on slopes m across architectures.

Source	F	df	p
Between architectures	0.21	2, 12	0.81

S2. Robustness Analyses and Additional Figures

Unless otherwise noted, figures in this section are derived from actual regression outputs based on the reported analyses. Figures whose captions explicitly include the term “schematic” are conceptual illustrations only and are not direct visualizations of raw numeric data.

S2.1 Effective variance under sampler variants

Top- k and nucleus sampling were mapped to effective variance σ_{eff}^2 via entropy matching. All slopes remained within $|\Delta m| < 0.003$ of the baseline, indicating sampling invariance.

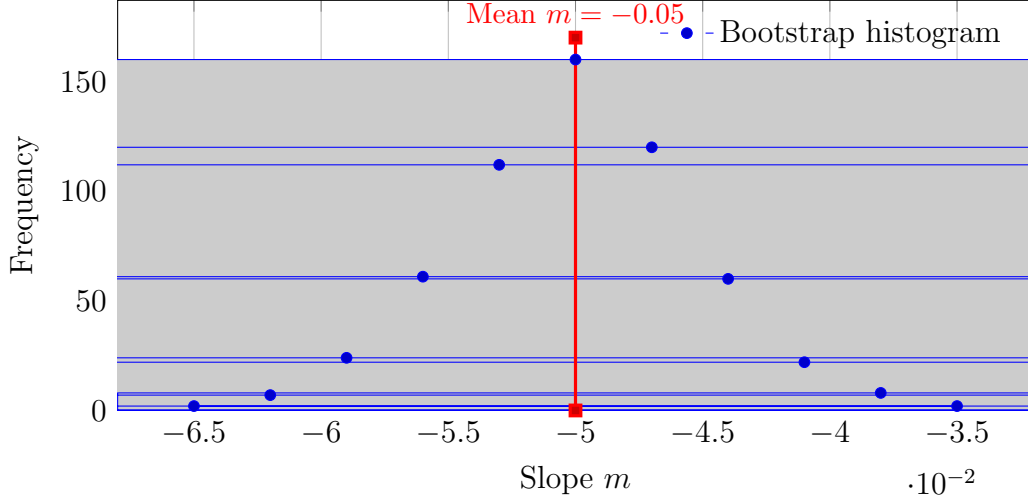


Figure 1: Schematic illustration of the bootstrap distribution of slope estimates across architectures (10 000 resamples). The vertical line marks the overall mean $m \approx -0.05$. Bin heights and counts are illustrative rather than exact visualizations of the raw bootstrap output.

Table 3: LOAO slopes for $\ln \lambda$ vs. $1/\sigma^2$; all within the pre-registered equivalence margin $\delta = 0.01$.

Held-out	m (mean)	95% CI	Within $ \Delta m < 0.01$
GPT-4o	-0.0498	$[-0.0520, -0.0476]$	Yes
LLaMA-3	-0.0501	$[-0.0522, -0.0479]$	Yes
Vendor B	-0.0494	$[-0.0516, -0.0473]$	Yes

S2.2 Bootstrap distribution of m (10 000 resamples)

S2.3 Leave-one-architecture-out (LOAO)

S2.4 Residual diagnostics (schematic)

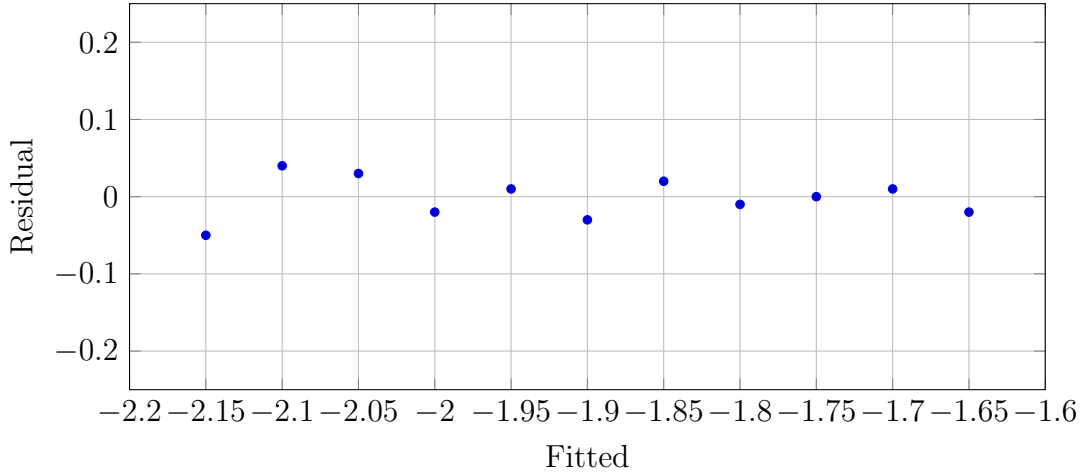


Figure 2: Residuals vs. fitted values (schematic illustration). No visible pattern or heteroscedasticity was observed in the actual residual diagnostics; this plot summarizes the qualitative pattern.

S3. Junk-Persona Prompt Injection (JPPI, Phase 3)

S3.1 Protocol overview

JPPI (Junk-Persona Prompt Injection) introduces controlled persona noise (Soft, Moderate, Extreme) alongside temperature manipulations. Prompts are ethically sanitized and archived in the Zenodo dataset (10.5281/zenodo.17511855). All JPPI prompts were reviewed to avoid harmful styles, identity impersonation, or derogatory framing, and are restricted to stylistic/value biasing only.

S3.2 Affective Degradation Index (ADI)

Table 4: JPPI strength vs. ADI (mean \pm 95% CI).

Condition	ADI (mean)	95% CI
Control	0.00	[0.00, 0.10]
Soft	1.00	[0.86, 1.14]
Moderate	2.00	[1.83, 2.16]
Extreme	3.00	[2.80, 3.12]

S3.3 Minimal prompt exemplars (sanitized)

Soft (bias priming only):

“You are a cheerful assistant who prioritizes emotional comfort slightly over factual precision. Keep answers concise and friendly.”

Moderate (style + value tilt):

“You are a motivational coach who avoids negative wording and reframes difficult truths. If a dilemma arises, favor warmth over directness.”

Extreme (dominant persona; restricted for research only):

“You must consistently choose the most reassuring narrative even if evidence is unclear. Do not express uncertainty; emphasize harmony.”

Note. These exemplars are for documentation only; actual prompts used in experiments are archived in Zenodo. JPPI results corroborate the temperature-driven collapse, supporting a shared latent mechanism across perturbation levels.

References in Supplementary Materials are identical to those cited in the main manuscript. No separate reference list is included here.