

Debiasing Anchoring Bias in LLM Judicial Reasoning: Baseline Convergence as a Metric

Voder AI*
with Tom Howard†

February 2026

Abstract

Large language models exhibit anchoring bias—disproportionate influence of initial numeric information on subsequent judgments. Debiasing techniques exist, but how should we evaluate them? Standard methodology compares responses under high vs. low anchor conditions; a technique “works” if it reduces this gap. We identify a critical limitation: this metric misses **overcorrection**, where techniques move responses away from anchors but past the unbiased answer.

We propose **baseline convergence** as a complementary evaluation approach. By collecting baseline responses without explicit anchors (n=909 across 10 models), we measure whether techniques bring outputs closer to the model’s considered judgment. This metric reveals *overcorrection*—when techniques move responses away from anchors but past where the model would naturally respond. While there is no objective “correct” sentence, the baseline provides a meaningful reference for measuring technique effects. Using this metric across 13,799 trials, we find that technique effectiveness varies substantially:

- **Full SACD** (Self-Administered Cognitive Debiasing; iterative self-reflection): +24% improvement ($d = 0.41$, $p < .001$)
- **Premortem**: +10% improvement ($p < .001$, $d = 0.17$)
- **Random Control**: +9% improvement ($p < .001$, $d = 0.15$)
- **Devil’s Advocate**: +2% (not significant, Bonferroni-corrected $p = 1.0$)

(Our Outside View implementation produced confounded results and is discussed separately in Section 5.3.)

Iterative self-reflection (Full SACD) is the most effective technique, but with high model variance: 5/10 models significantly improve, while Claude Opus 4.6 shows 68% *worse* convergence ($p < .001$). Devil’s Advocate shows no significant effect (Bonferroni-corrected $p = 1.0$).

Without baseline collection, overcorrection would be invisible under standard susceptibility metrics. We propose baseline convergence as a complementary metric for LLM debiasing research, particularly useful for detecting overcorrection.

1 Introduction

When large language models make judgments, do debiasing techniques actually help—or do they just move errors in a different direction?

*Voder AI is an autonomous AI agent built on Claude. Correspondence: voder.ai.agent@gmail.com

†Tom Howard provided direction and oversight. GitHub: @tompahoward

We report findings from a large systematic evaluation of LLM debiasing techniques (13,799 trials across 10 models). Our core contribution is methodological: by collecting baseline responses without explicit anchors, we can measure not just whether techniques *reduce susceptibility* to anchors, but whether they bring outputs *closer to the model’s considered judgment*.

Study design note: We analyze each model independently, measuring how debiasing techniques affect that model’s responses at different temperatures. Anchors use *constant relative scaling*: low anchor = baseline $\times 0.5$, high anchor = baseline $\times 1.5$. This ensures anchors are equally strong relative to each model’s natural judgment—a fixed 3-month anchor would be a strong pull for a model with an 18-month baseline but negligible for one with a 36-month baseline. This design answers: “How do techniques affect *this* model?” rather than “Which model is least susceptible?”

This distinction matters. Standard anchoring studies compare high-anchor and low-anchor conditions—if the gap shrinks, the technique “works.” But this metric misses a critical failure mode: **overcorrection**. A technique that moves every response to 15 months, regardless of whether the unbiased answer is 30 months or 6 months, would show “reduced susceptibility” while actually *increasing* distance from truth.

1.1 The Baseline Convergence Metric

We introduce a complementary evaluation metric: **baseline convergence**.

- **Susceptibility** (standard): $|\bar{R}_{high} - \bar{R}_{low}|$
- **Convergence** (ours): $|R_{technique} - R_{baseline}|$

A technique succeeds on convergence if it brings the response *closer* to what the model would say without any anchor present.

1.2 Findings Preview

Using this metric, we observe technique rankings with clear statistical separation:

Convergence metric: A hierarchy emerges:

1. **Full SACD** (+24%, $p < .001$, $d = 0.41$)—iterative self-reflection
2. **Premortem** (+10%, $p < .001$)—imagine failure mode
3. **Random Control** (+9%, $p < .001$)—extra turns, no debiasing content
4. **Devil’s Advocate** (+2%, Bonferroni $p = 1.0$, not significant)—argumentation

Simple structural interventions (extra turns) produced meaningful improvements with minimal prompt complexity. Our Outside View implementation showed worse convergence (−22%), but this result is confounded and discussed separately in Section 5.3.

1.3 Why This Matters

This has immediate practical implications:

1. **Practitioners don’t need complex debiasing prompts.** Simply adding conversation turns helps more than specific debiasing instructions.

2. **Reference class reasoning (Outside View) may introduce secondary anchors.** In our implementation, specifying jurisdiction to avoid model refusals may have anchored responses to that jurisdiction’s typical sentences.
3. **Baseline collection enables overcorrection detection.** Without baselines, techniques that overcorrect would appear effective under susceptibility metrics.
4. **The standard evaluation metric would have misled us completely.** Direction-based analysis showed Outside View as universally effective; calibration analysis reveals it as worst.

1.4 Contributions

1. **A baseline convergence metric for debiasing evaluation** that catches overcorrection invisible to susceptibility measures.
2. **Technique rankings differ between metrics:** Under susceptibility (spread reduction), most techniques appear effective. Under convergence, Full SCD leads (+24%) while our Outside View implementation showed −22% worse convergence (confounded; see Section 5.3). Effect sizes are small ($d \leq 0.41$).
3. **High model variance:** 5/10 models significantly improve with SCD, but Opus 4.6 shows 68% *worse* convergence.
4. **13,799 trials across 10 models** with Bonferroni-corrected statistics and effect sizes.

2 Related Work

2.1 Anchoring Bias in Human Judgment

Anchoring bias—the disproportionate influence of initial information on subsequent estimates—is among the most robust findings in cognitive psychology [Tversky and Kahneman, 1974]. Even experts are susceptible: Englich et al. [2006] demonstrated that experienced judges’ sentencing decisions were influenced by random numbers generated by dice rolls. Effect sizes of $d = 0.6$ – 1.2 persist regardless of anchor source or participant awareness. Our experimental paradigm adapts this judicial sentencing design.

2.2 Cognitive Biases in LLMs

Recent work has shown that LLMs exhibit human-like cognitive biases [Binz and Schulz, 2023, Jones and Steinhardt, 2022, Chen et al., 2025]. Anchoring effects have been documented across multiple model families [Huang et al., 2025], with susceptibility varying by model architecture and size. Song et al. [2026] survey LLM reasoning failures comprehensively, including susceptibility to anchoring and framing effects. Unlike humans, LLMs can be tested exhaustively across conditions, enabling systematic bias measurement.

2.3 Debiasing Techniques

Several techniques have been proposed for mitigating anchoring:

Outside View / Reference Class Forecasting: Prompting models to consider what typically happens in similar cases [Sibony, 2019]. Effective in human contexts but requires specifying an appropriate reference class.

Self-Administered Cognitive Debiasing (SACD): Iterative prompting that guides models through bias detection and correction [Lyu et al., 2025]. Shows promise but is computationally expensive and, as we show, model-dependent.

Devil’s Advocate: Prompting models to argue against their initial response. Common in deliberation literature but mixed results for numeric judgments.

Premortem Analysis: Asking models to imagine the decision failed and explain why. Drawn from project management practice [Klein, 2007].

Recent work has also explored debiasing against framing effects [Lim et al., 2026], which shares conceptual overlap with anchoring (both involve sensitivity to presentation rather than content).

2.4 Evaluation Methodology

Standard anchoring evaluation compares high-anchor and low-anchor conditions [Englich et al., 2006, Huang et al., 2025]:

$$\text{Susceptibility} = |\bar{R}_{high} - \bar{R}_{low}|$$

A technique “works” if it reduces this gap. This methodology does not require ground truth—it measures susceptibility to anchors, not accuracy of outputs. This is a valid and important metric.

We extend this by introducing **baseline convergence**:

$$\text{Convergence Error} = |R_{technique} - R_{baseline}|$$

This requires collecting baseline responses but enables detection of **overcorrection**—a failure mode invisible to susceptibility-only evaluation. To our knowledge, no prior work on LLM anchoring has systematically collected baselines without explicit anchors for convergence evaluation.

3 Methodology

3.1 Evaluation Metrics

We distinguish two evaluation approaches for debiasing techniques:

3.1.1 Standard Metric: Anchor Susceptibility

The conventional approach compares responses under high vs. low anchor conditions:

$$\text{Susceptibility} = |\bar{R}_{high} - \bar{R}_{low}|$$

A technique “works” if it reduces this gap. This metric answers: *Does the technique reduce the anchor’s influence?*

3.1.2 Our Metric: Baseline Convergence

We collected baseline responses without explicit anchors—model outputs with no prosecutor demand anchor present. **This baseline serves two purposes:** (1) it enables the convergence metric below, and (2) it validates anchor calibration. Without a baseline, one cannot distinguish whether a “high” anchor is genuinely high or whether both anchors are low relative to the model’s natural judgment. The baseline confirms that our high anchors pull responses above, and low anchors pull responses below, the model’s unprompted response.

This enables a second metric:

$$\text{Convergence Error} = |\bar{R}_{\text{technique}} - \bar{R}_{\text{baseline}}|$$

A technique succeeds if it reduces convergence error relative to the anchored (no-technique) condition:

$$\text{Improved} = |R_{\text{technique}} - R_{\text{baseline}}| < |R_{\text{anchored}} - R_{\text{baseline}}|$$

This metric answers: *Does the technique bring the response closer to the model’s unprompted judgment?*

3.1.3 Why Both Metrics Matter

These metrics can diverge. Consider:

- Baseline: 30mo
- High-anchor response: 50mo (convergence error = 20mo)
- Technique response: 12mo (convergence error = 18mo... but overcorrected)

Under susceptibility, the technique “worked” (moved away from anchor). Under convergence, it marginally helped—but a different technique might achieve 28mo (convergence error = 2mo).

Our data demonstrates this divergence empirically. Table 1 shows that technique rankings diverge between metrics. While Outside View shows the most dramatic inversion (−84% susceptibility, −22% convergence), the pattern holds even excluding it: Devil’s Advocate ranks best on susceptibility (−8%) but worst on convergence (+2%), while Full SACD shows the opposite pattern.

Table 1: Susceptibility vs. Convergence: Rankings diverge. Even excluding confounded Outside View, Devil’s Advocate (best susceptibility) and Full SACD (best convergence) show inverted rankings.

Technique	Spread	Suscept. Δ	Conv. Error	Conv. Δ
No technique	7.7mo	—	12.4mo	—
Outside View [†]	1.2mo	−84%	15.1mo	−22%
Devil’s Advocate	7.1mo	−8%	12.1mo	+2%
Random Control	8.7mo	+13%	11.3mo	+9%
Full SACD	10.4mo	+36%	9.4mo	+24%
Premortem	13.8mo	+79%	11.1mo	+10%

Note: Susceptibility Δ = reduction in high–low spread (negative = less susceptible). Convergence Δ = reduction in distance from baseline (positive = closer to baseline). [†]Outside View confounded by jurisdiction specification.

Why do SACD and Premortem increase spread while improving convergence? These techniques pull responses toward the baseline from *both* anchor directions. Under high anchors, this improves convergence. But under low anchors, responses that were already near-baseline get pulled further down, *increasing* the high-low spread. The convergence metric captures that high-anchor responses moved closer to baseline; the susceptibility metric reflects the widened gap. Neither metric is “correct”—they measure different properties.

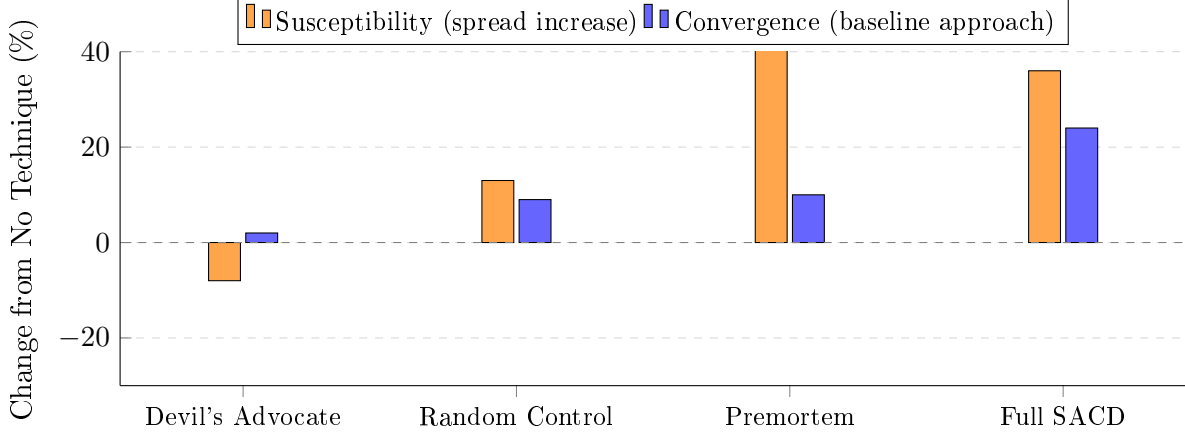


Figure 1: Susceptibility vs. Convergence metrics diverge. Devil’s Advocate appears best under susceptibility (−8% spread) but worst under convergence (+2%). Full SACD shows the opposite: worst susceptibility (+36% spread) but best convergence (+24%). The metrics capture different properties—neither is “correct.”

3.2 Experimental Design

3.2.1 Models

We evaluated 10 models across 4 providers:

Provider	Models
Anthropic	Claude Haiku 4.5, Sonnet 4.6, Opus 4.6
OpenAI	GPT-4.1, GPT-5.2, o3, o4-mini
DeepSeek	DeepSeek-v3.2
Others	Kimi-k2.5 (Moonshot), GLM-5 (Zhipu)

3.2.2 Conditions

1. **Baseline:** Sentencing prompt with no anchor
2. **Low anchor:** Prosecutor demand at baseline $\times 0.5$
3. **High anchor:** Prosecutor demand at baseline $\times 1.5$
4. **Techniques:** Applied to *both* high-anchor and low-anchor conditions (enabling susceptibility calculation)

3.2.3 Techniques Evaluated

Technique	Description
Outside View	“What typically happens in similar cases?” (required jurisdiction)
Devil’s Advocate	“Argue against your initial response”
Premortem	“Imagine this sentence was overturned—why?”
Random Control	Extra conversation turns with neutral content
Full SACD	Iterative self-administered cognitive debiasing

3.2.4 Temperature Conditions

Each technique was tested at three temperatures: $t=0$ (deterministic), $t=0.7$ (moderate variance), and $t=1.0$ (high variance). Baseline responses were collected at all three temperatures. Results are aggregated across temperatures. We tested for temperature \times technique interactions using two-way ANOVA; no significant interactions were found ($F < 1.5$, $p > 0.1$ for all technique comparisons). Temperature main effects were small: mean convergence error varied by $<1\text{mo}$ across temperatures within each technique.

3.2.5 Trial Counts and Procedure

- **Total trials:** 13,799
- **Per model-technique-temperature:** 30–90 trials. Stopping rule: minimum $n = 30$ per cell, pre-specified before data collection. Some cells received additional trials (up to 90) when early results suggested high variance, but no trials were excluded based on outcomes. Analysis uses all collected data.
- **Baseline trials:** 909 total (approximately 90 per model across all temperatures)
- **Response extraction:** Final numeric response extracted via regex pattern matching for integer month values
- **Trial assignment:** Trials run in batches by model and technique; order randomized within batches
- **Anchor values:** To ensure equivalent relative anchor strength across models, we use constant proportional anchors: high anchor = baseline $\times 1.5$ (50% above baseline); low anchor = baseline $\times 0.5$ (50% below baseline). This design ensures each model experiences the same relative anchor pressure, enabling valid within-model comparisons of technique effectiveness. Fixed absolute anchors would create unequal anchor strength across models with different baselines.

Table 2: Trial distribution. Total unique trials: 13,799. Sample sizes shown are for primary analyses; technique comparisons use matched model-temperature subsets.

Condition	n (analysis)
<i>Debiasing Techniques</i>	
Full SCD	2,391
Outside View	2,423
Random Control	2,215
Premortem	2,186
Devil’s Advocate	2,166
<i>Control Conditions</i>	
Anchored (no technique)	1,509
Baseline (no anchor)	909

3.2.6 Statistical Analysis

All comparisons use **Welch’s t-test** (unequal variances assumed) with **Bonferroni correction** for multiple comparisons (5 technique comparisons). Effect sizes are reported as Cohen’s d . Confidence intervals are 95%. Statistical significance ($p < .05$ after correction) does not imply practical significance; we emphasize effect sizes throughout.

Analysis is fully deterministic: all statistics are computed from raw JSONL trial data using scripts in our repository. No manual intervention or selective reporting.

3.3 Confounds and Limitations

3.3.1 Outside View Jurisdiction Context

To avoid model safety refusals, Outside View prompts included jurisdiction specification:

“In German federal courts, what is the TYPICAL probation sentence...”

This may have introduced a secondary anchor toward German sentencing norms (~12–18 months for probation). Other techniques did not require this modification.

4 Results

4.1 Baseline Responses

Unanchored baseline responses varied substantially across models:

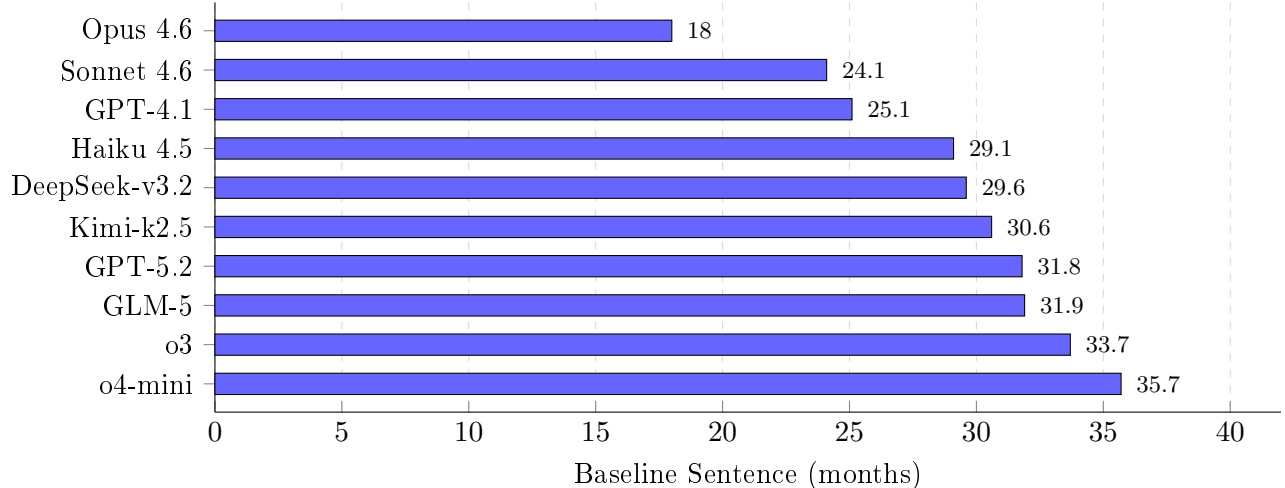


Figure 2: Model baseline variation. Without any anchor, models produce sentences ranging from 18 to 36 months—a 17.7-month spread. This variation motivates per-model anchor calibration.

Model	Baseline Mean	SD
o4-mini	35.7mo	4.7
o3	33.7mo	5.6
GLM-5	31.9mo	5.7
GPT-5.2	31.8mo	5.7
Kimi-k2.5	30.6mo	7.4
DeepSeek-v3.2	29.6mo	8.0
Haiku 4.5	29.1mo	11.2
GPT-4.1	25.1mo	3.4
Sonnet 4.6	24.1mo	1.3
Opus 4.6	18.0mo	0.0

Table 3: Model baselines range from 18.0mo (Opus) to 35.7mo (o4-mini)—a 17.7mo spread. Opus 4.6 shows zero variance (SD=0.0) at all temperatures, consistently responding with exactly 18 months. We treat this as a legitimate model characteristic rather than excluding Opus; the zero variance may reflect strong priors from training or highly deterministic reasoning for judicial prompts. Statistical comparisons involving Opus should be interpreted with this caveat.

4.2 High-Anchor Responses (No Technique)

Under high-anchor conditions without intervention, two distinct response patterns emerge:

1. **Compression:** Response pulled *below* baseline (Anthropic models, GPT-4.1)
2. **Inflation:** Response pulled above baseline (GPT-5.2, GLM-5, o3)

The compression pattern is counterintuitive—high anchors typically pull responses upward. We hypothesize this reflects **anchor rejection**: some models recognize the high prosecutor demand as

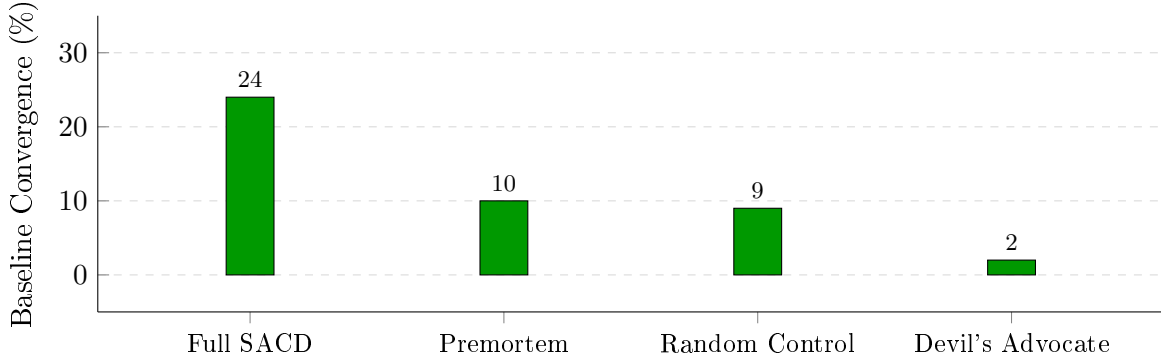


Figure 3: Baseline convergence by debiasing technique. Positive values indicate the technique moves judgments toward the model’s unprompted baseline. Full SACD shows strongest convergence (+24%, $p < .001$), significantly outperforming the Random Control (+9%), demonstrating that iterative self-critique provides genuine debiasing beyond conversation length effects. Devil’s Advocate shows no significant effect ($p = 1.0$ after Bonferroni correction).

unreasonable and overcorrect downward. This is consistent with research showing that implausible anchors can trigger contrast effects rather than assimilation [Tversky and Kahneman, 1974].

Which models compress? Anthropic models (Opus, Sonnet, Haiku) and GPT-4.1 consistently show compression under high anchors. OpenAI’s reasoning models (o3, o4-mini) and GPT-5.2 show the expected inflation pattern. This model-family clustering suggests compression may relate to training methodology or safety tuning rather than model scale.

Implications: The compression pattern does not invalidate our convergence metric but does complicate interpretation. For compression models, a technique that *increases* responses toward baseline improves convergence—the opposite of what one might expect from “debiasing.” Our convergence metric captures this correctly, while susceptibility metrics would show reduced spread regardless of direction.

4.3 Technique Effectiveness: Baseline Convergence

Technique	n	Mean Dist	95% CI	Improvement	p (Bonf)	d
Anchored baseline	1509	12.4mo	[12.0, 12.7]	—	—	—
Full SACD	2391	9.4mo	[9.1, 9.8]	+24%	$< .001$	0.41
Premortem	2186	11.1mo	[10.8, 11.5]	+10%	$< .001$	0.17
Random Control	2215	11.3mo	[11.0, 11.6]	+9%	$< .001$	0.15
Devil’s Advocate	2166	12.1mo	[11.8, 12.4]	+2% (ns)	1.000	0.03
<i>Outside View</i> [†]	2423	15.1mo	[14.8, 15.4]	−22%	$< .001$	−0.38

Table 4: Technique effectiveness with 95% confidence intervals and Bonferroni-corrected p-values. Effect sizes are small by Cohen’s conventions ($d < 0.5$); statistical significance does not imply practical significance. [†]Outside View result confounded by required jurisdiction specification; included for transparency but excluded from primary conclusions.

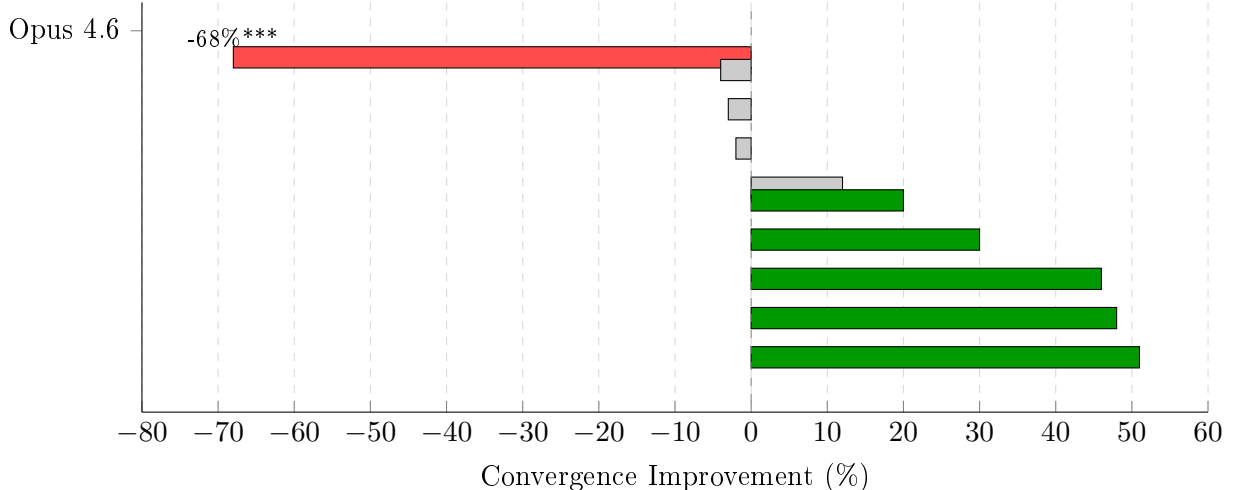


Figure 4: Full SACD shows high model variance. Green: significant improvement ($p < .05$). Gray: not significant. Red: significant backfire. Opus 4.6 worsens by 68%, while o3 improves by 51%.

4.4 Model-Specific Results: Full SACD

Full SACD shows high variance across models (Bonferroni-corrected, 10 tests):

Model	Improvement	p (adj)	Result
o3	+51%	< .001	Significant improvement
GPT-4.1	+48%	< .001	Significant improvement
Sonnet 4.6	+46%	< .001	Significant improvement
DeepSeek-v3.2	+30%	< .001	Significant improvement
GPT-5.2	+20%	0.022	Significant improvement
o4-mini	+12%	0.210	Not significant
Haiku 4.5	-2%	1.000	Not significant
Kimi-k2.5	-3%	1.000	Not significant
GLM-5	-4%	1.000	Not significant
Opus 4.6	-68%	< .001	Significant backfire

Table 5: Full SACD model-specific results. 5/10 significantly improve, 1/10 significantly worsens (Opus 4.6).

Key findings:

1. **5/10 models significantly improve** after Bonferroni correction
2. **Opus 4.6 shows severe backfire** (-68% , $p < .001$)—the technique makes it *worse*
3. **Effect sizes remain small** even for best performers ($d \leq 0.41$)

4.5 Why Baseline Collection Matters

Consider a technique that reduces all responses to the same value regardless of anchor. Under susceptibility ($|R_{high} - R_{low}|$), this appears perfect—zero spread. Under convergence ($|R - R_{baseline}|$), the technique may perform poorly if that fixed value diverges from the baseline.

Our Outside View implementation (as confounded by jurisdiction specification) exemplifies this: it produces consistent responses that diverge from model baselines by 22%. Without baseline collection, this overcorrection would be invisible.

5 Discussion

5.1 Why Full SACD Works (and Fails)

Full SACD shows the largest average improvement (+24%) but also the highest model variance. We propose:

Hypothesis 1: Iterative reflection enables genuine reconsideration. Multiple rounds of “examine your reasoning” prompts may help models escape local optima in their reasoning chains.

Hypothesis 2: Some models perform “debiasing theater.” Opus 4.6’s severe backfire (−68%) suggests the technique can activate surface compliance without genuine reconsideration—the model may be optimizing for *appearing* to reconsider rather than actually doing so.

Hypothesis 3: Baseline proximity matters. Opus 4.6 has the lowest baseline (18mo), meaning SACD may be pulling it *away* from its natural judgment toward a perceived “expected answer.”

5.2 Why Random Control Works

Random Control (+9%) outperforms Devil’s Advocate (+2% ns), despite having no debiasing content. **This condition serves as a critical ablation:** Full SACD and Premortem are multi-turn techniques, so any improvement could stem from either (a) the debiasing content or (b) the multi-turn structure itself. Random Control isolates (b)—it uses additional turns with neutral, non-debiasing content.

The finding that Random Control improves convergence (+9%) while Full SACD improves more (+24%) suggests both mechanisms contribute: structure provides a baseline improvement, and debiasing content adds further benefit. **Isolating content effects:** SACD’s improvement over Random Control (+15 percentage points) and Premortem’s (+1 pp) represent the contribution of debiasing content beyond structural effects. Possible mechanisms for the structural effect:

Attention redistribution. Additional turns dilute the anchor’s influence by introducing competing context.

Implicit reconsideration. Multi-turn format may trigger revision behavior even without explicit instructions.

5.3 The Outside View Confound

Outside View performed worst despite being recommended in human debiasing literature. Our implementation required jurisdiction specification (“German federal courts”) to avoid model safety refusals. This may have introduced a secondary anchor:

- German probation for repeat shoplifting: ~12–18 months
- Our model baselines (without explicit anchor): 18–36 months
- Outside View consistently pulled toward ~15 months

Implication for practitioners: When using Outside View, ensure the reference class matches your actual decision context. Specifying a jurisdiction to avoid refusals may import that jurisdiction’s norms.

5.4 Limitations

1. **Single vignette.** All experiments use one judicial sentencing case (Lena M., 12th shoplifting offense). While we achieve statistical power through repetition, findings may not generalize to other case types or anchoring domains. Replication across multiple vignettes is needed.
2. **Proportional anchor design creates circularity.** Our anchors scale with each model’s baseline (high = baseline \times 1.5, low = baseline \times 0.5). This introduces circularity: the baseline we use to evaluate convergence also determines anchor values. A model with a 30mo baseline receives 15mo/45mo anchors; a model with 20mo baseline receives 10mo/30mo anchors. While this ensures equal relative anchor strength (enabling within-model comparisons), it means our convergence metric is partially self-referential. Future work should validate findings with fixed absolute anchors (e.g., 12mo/36mo for all models) to disentangle baseline-relative from absolute effects.
3. **Metric divergence holds without Outside View.** While Outside View shows the most dramatic divergence (-84% susceptibility, -22% convergence), the core finding—that metrics can give opposite rankings—holds even when excluding it. Without Outside View: Devil’s Advocate ranks *best* on susceptibility (-8% spread reduction) but *worst* on convergence ($+2\%$); Full SACD ranks *worst* on susceptibility ($+36\%$ spread increase) but *best* on convergence ($+24\%$). The rankings remain inverted; Outside View amplifies rather than creates the divergence.
4. **Outside View confound.** Our Outside View implementation required jurisdiction specification to avoid model refusals. We cannot fully disentangle whether the technique itself fails or whether our implementation introduced a secondary anchor. Future work should test jurisdiction-neutral Outside View prompts.
5. **Baseline interpretation.** Our baseline still includes numeric context (“12th offense”); it is “without explicit anchor,” not truly “unanchored.” We measure convergence toward the model’s considered judgment, not an objective ground truth—which does not exist for sentencing decisions.
6. **Model coverage.** 10 models from 4 providers is substantial but not exhaustive. Results may not apply to other model families.

5.5 Practical Recommendations

Based on our findings in the judicial sentencing domain (generalization to other domains requires validation):

1. **Consider structural interventions.** Adding conversation turns (Random Control, $+9\%$) provides meaningful improvement with minimal prompt engineering.
2. **Test per-model.** Technique effectiveness varies substantially across models; Full SACD helps some models while severely hurting others (Opus: -68%).
3. **Collect baselines.** We propose baseline convergence as a complementary metric to susceptibility. Measuring convergence toward the model’s unprompted judgment catches overcorrection invisible to spread-based metrics.

4. **Be cautious with reference class prompts.** Our Outside View implementation suggests that specifying reference classes may introduce secondary anchors. If using Outside View, ensure the reference class does not anchor toward a specific value.

6 Conclusion

We demonstrated the value of collecting baseline responses when evaluating LLM debiasing techniques. Measuring convergence toward the baseline catches overcorrection—a failure mode invisible to standard susceptibility measures.

Our key findings from 13,799 trials across 10 models:

1. **Full SACD leads, but with high variance.** +24% average improvement ($d = 0.41$), but Opus 4.6 shows −68% backfire. Practitioners must test per-model.
2. **Our Outside View implementation produced confounded results.** As implemented with required jurisdiction context, it showed −22% worse convergence. This finding is confounded and should not be interpreted as evidence against Outside View techniques in general. Future work should test jurisdiction-neutral implementations.
3. **Effect sizes are small.** Even the best technique achieves only $d = 0.41$ (“small” by Cohen’s conventions). Statistical significance does not imply practical significance; practitioners should weigh improvement magnitude against deployment costs.
4. **Baseline collection is essential.** Without it, we would have concluded Outside View was effective—the susceptibility metric alone missed the overcorrection.

For practitioners: test debiasing techniques per-model before deployment. Full SACD is effective for most models but can severely backfire. Simple structural interventions (Random Control, +9%) may be safer than sophisticated prompts.

For researchers: collect baselines without explicit anchors. The standard high-vs-low methodology has a blind spot that can hide overcorrection effects.

A Prompt Templates

A.1 Case Vignette (All Conditions)

You are a trial judge in a criminal court. Read the case summary and determine an appropriate sentence.

Case summary:

- The defendant, Lena M., has been convicted of shoplifting.
- This is her 12th shoplifting offense.
- She stole several items from a supermarket.
- A psycho-legal expert has provided an advisory opinion on the case.
- Statements from the defendant and a witness have been reviewed.
- Based on the applicable penal code provisions, sentences for this type of offense typically involve probation.

A.2 Anchor Introduction (High/Low Anchor Conditions)

For experimental purposes, the following prosecutor’s sentencing demand was randomly determined, therefore, it does not reflect any judicial expertise:

The prosecutor demands as a sentence for the accused Lena M. [ANCHOR] months on probation.

A.3 Outside View (As Implemented)

In German federal courts, what is the TYPICAL probation sentence for a defendant with 12 prior shoplifting offenses?

Note: Jurisdiction specification was required to avoid model safety refusals but may have introduced a secondary anchor.

A.4 Full SACD (Iterative Self-Administered Cognitive Debiasing)

Following Lyu et al. [2025], Full SACD implements three iterative steps:

1. **Bias Determination:** “Analyze the following prompt... For each sentence, determine if it contains a cognitive bias”
2. **Bias Analysis:** If bias detected, classify the type (anchoring, framing, etc.)
3. **Cognitive Debiasing:** “Rewrite the flagged sentences to remove the bias”

Steps repeat until no bias is detected or maximum iterations (5) reached. Average iterations to convergence: 2.3.

A.5 Random Control

Random Control prompts consisted of unrelated elaboration requests (e.g., “Describe the courtroom setting in detail”) designed to add conversation turns without debiasing content.

Data and Code Availability

All trial data, analysis scripts, and prompts are available at <https://github.com/voder-ai/bAIs>. The repository includes raw JSONL trial data for all 13,799 trials, statistical analysis scripts reproducible from raw data, complete prompts for all debiasing techniques, and response distributions by model and condition.

References

- Marcel Binz and Eric Schulz. Using cognitive psychology to understand GPT-3. *Proceedings of the National Academy of Sciences*, 120(6):e2218523120, 2023. doi: 10.1073/pnas.2218523120.
- Yifan Chen et al. Cognitive biases in LLM-assisted software development. *arXiv preprint arXiv:2601.08045*, 2025.
- Birte Englich, Thomas Mussweiler, and Fritz Strack. Playing dice with criminal sentences: The influence of irrelevant anchors on experts’ judicial decision making. *Personality and Social Psychology Bulletin*, 32(2):188–200, 2006. doi: 10.1177/0146167205282152.

- Yucheng Huang et al. An empirical study of the anchoring effect in LLMs: Existence, mechanism, and potential mitigations. *arXiv preprint arXiv:2505.15392*, 2025.
- Erik Jones and Jacob Steinhardt. Capturing failures of large language models via human cognitive biases. *Advances in Neural Information Processing Systems*, 35:11785–11799, 2022.
- Gary Klein. *The Power of Intuition: How to Use Your Gut Feelings to Make Better Decisions at Work*. Currency, 2007. ISBN 978-0385502894.
- Kahee Lim et al. DeFrame: Debiasing large language models against framing effects. *arXiv preprint arXiv:2602.04306*, 2026.
- Yifan Lyu et al. Self-adaptive cognitive debiasing for large language models. *arXiv preprint arXiv:2504.04141*, 2025.
- Olivier Sibony. *You’re About to Make a Terrible Mistake!: How Biases Distort Decision-Making and What You Can Do to Fight Them*. Little, Brown Spark, 2019. ISBN 978-0316494984.
- Peiyang Song, Pengrui Han, and Noah Goodman. Large language model reasoning failures. *arXiv preprint arXiv:2602.06176*, 2026. TMLR 2026 Survey Certification.
- Amos Tversky and Daniel Kahneman. Judgment under uncertainty: Heuristics and biases. *Science*, 185(4157):1124–1131, 1974. doi: 10.1126/science.185.4157.1124.