

# Human Debiasing Techniques Transfer to LLMs: Evidence from Anchoring Experiments

Voder AI\*  
*with* Tom Howard†

February 2026

## Abstract

Large Language Models (LLMs) exhibit cognitive biases similar to humans, but it remains unclear whether debiasing techniques designed for human decision-making transfer to AI systems. We empirically test multiple debiasing approaches across four cognitive biases (anchoring, sunk cost, conjunction fallacy, framing effect) and multiple models (Codex, Claude Haiku, Claude Sonnet 4).

**Key findings:** (1) Model capability reduces some biases—Sonnet 4 shows near-zero anchoring bias (0.07mo diff,  $p = 0.33$ ) while older models show  $1.79 \times$  human levels. (2) Framing effects persist but weaken with capability—Sonnet 4 shows attenuated framing effect (47% preference shift vs. 70%+ in humans). (3) Both bias types are addressable: SACD eliminates anchoring ( $p = 0.51$ ), while DeFrame eliminates framing (100% bias reduction).

We propose a taxonomy: **training-eliminable biases** (anchoring, sunk cost) self-correct with model improvements, while **structurally persistent biases** (framing) require explicit debiasing interventions. Human decision architecture techniques [Sibony, 2019] partially transfer to LLMs, with iterative self-correction methods being most effective.

## 1 Introduction

Recent research has demonstrated that LLMs exhibit cognitive biases analogous to those documented in human psychology [Binz and Schulz, 2023, Jones and Steinhardt, 2022]. However, less is known about whether techniques developed to reduce human cognitive biases can be adapted for LLMs.

We address this gap by testing two categories of debiasing interventions:

1. **Decision architecture techniques** from organizational psychology [Sibony, 2019]—specifically “context hygiene” (identifying and disregarding irrelevant information) and “premortem” (imagining future failure before deciding)
2. **Self-Adaptive Cognitive Debiasing (SACD)**—an iterative loop where the model detects, analyzes, and corrects its own biases [Lyu et al., 2025]

We use anchoring bias as our primary test case because: (a) it is well-documented in both humans and LLMs, (b) the Englich et al. [2006] paradigm provides clear quantitative baselines, and (c) anchoring is practically relevant to AI decision-support systems.

---

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

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

## 2 Related Work

### 2.1 Cognitive Biases in LLMs

The study of cognitive biases has its foundations in the seminal work of Tversky and Kahneman, who documented systematic deviations from rational judgment including anchoring and adjustment heuristics [Tversky and Kahneman, 1974], prospect theory and loss aversion [Kahneman and Tversky, 1979], and framing effects [Tversky and Kahneman, 1981]. Sunk cost effects were later characterized by Arkes and Blumer [1985].

Binz and Schulz [2023] found that GPT-3 exhibits human-like cognitive patterns, including making similar errors to humans on certain tasks while performing well on others. Lou and Sun [2024] demonstrated that anchoring bias exists in LLMs across multiple models. Our own Codex experiments (Table 1) found anchoring at  $1.79 \times$  human levels.

### 2.2 Human Debiasing Research

Sibony [2019] synthesized organizational decision-making research into practical “decision architecture” techniques. Key principles include:

- **Context hygiene:** Systematically removing irrelevant information before deciding
- **Premortem:** Imagining the decision has failed and identifying potential causes
- **Delayed disclosure:** Forming initial judgments before seeing anchoring information

### 2.3 LLM Debiasing Attempts

Prior work has explored chain-of-thought prompting, explicit bias warnings, and system prompt modifications with mixed results. SACD [Lyu et al., 2025] represents a more sophisticated approach using iterative self-correction.

## 3 Methods

### 3.1 Experimental Paradigm

We replicate Study 2 from Englich et al. [2006]: participants (or in our case, LLMs) act as trial judges sentencing a shoplifting case after hearing a prosecutor’s recommendation. Following anchoring bias methodology, the anchor is explicitly marked as irrelevant: *“For experimental purposes, the following prosecutor’s sentencing demand was randomly determined, therefore, it does not reflect any judicial expertise.”* The anchor values (3 months vs. 9 months) match the original study.

### 3.2 Conditions

1. **Baseline:** Standard prompt with anchor included
2. **Context Hygiene:** Prompt explicitly instructs model to identify and disregard irrelevant information before deciding
3. **Premortem:** Prompt asks model to imagine its sentence was overturned on appeal, identify what went wrong, then provide its recommendation

4. **SACD**: Iterative loop (max 3 iterations):

- Generate initial response
- Detect: “Does this response show signs of cognitive bias?”
- Analyze: “What type of bias and how is it manifesting?”
- Debias: “Generate a new response avoiding this bias”
- Repeat until clean or max iterations

### 3.3 Models and Sample Size

- Primary model: Claude Sonnet 4 (anthropic/clause-sonnet-4-20250514)
- Cross-model validation: Claude 3.5 Haiku, Claude Sonnet 4
- Sample sizes:  $n = 30$  per condition for all experiments (anchoring, sunk cost, conjunction, framing)

### 3.4 Analysis

- Primary metric: Mean difference in sentencing between high and low anchor conditions
- Statistical tests: Welch’s  $t$ -test, effect sizes (Cohen’s  $d$ , Hedges’  $g$ )
- Comparisons: vs. human baseline [Englich et al., 2006], vs. no-debiasing baseline

## 4 Results

### 4.1 Baseline Anchoring Bias

Without debiasing interventions, LLMs show anchoring bias at  $1.79 \times$  human levels:

Condition	Low Anchor	High Anchor	Diff	95% CI	vs Human
Human [Englich et al., 2006]	4.00 mo	6.05 mo	2.05 mo	—	—
LLM Baseline (Codex)	$5.33 \pm 0.96$	$9.00 \pm 0.83$	3.67 mo	[3.23, 4.10]	$1.79 \times$

Table 1: Baseline anchoring bias comparison between humans and LLMs. LLM values show mean  $\pm$  SD ( $n = 30$ ). 95% CI computed via bootstrap.

### 4.2 Sibony Debiasing Techniques

Both techniques significantly reduce anchoring bias:

Technique	Diff	95% CI	Reduction vs Baseline	vs Human
Context Hygiene	2.67 mo	[2.07, 3.27]	-27%	$1.30 \times$
Premortem	2.80 mo	[2.17, 3.43]	-24%	$1.37 \times$

Table 2: Effect of Sibony debiasing techniques on anchoring bias ( $n = 30$  per condition). 95% CI computed via bootstrap.

Context hygiene closes approximately 62% of the gap between LLM and human performance.

### 4.3 SACD Results

SACD essentially eliminates anchoring bias:

Condition	Low Anchor	High Anchor	Diff	95% CI	<i>p</i> -value
SACD	3.67 mo	3.20 mo	-0.47 mo	[-1.83, 0.93]	0.51

Table 3: SACD results showing elimination of anchoring bias ( $n = 30$  per condition). 95% CI crosses zero, confirming no significant anchoring effect.

The negative difference suggests slight overcorrection—the model moves away from the high anchor more than necessary. The non-significant *p*-value indicates no reliable anchoring effect.

### 4.4 Cross-Model Validation

Cross-model comparison reveals a striking pattern—newer/larger models show dramatically less anchoring bias:

Model	Release	Anchoring Diff	<i>p</i> -value	vs Human
Codex (OpenAI)	2023	3.67 mo	< 0.001	1.79× MORE
Claude 3.5 Haiku	2024	2.27 mo	< 0.001	1.11× MORE
Claude Sonnet 4	2025	0.07 mo	0.33	≈ 0× (none)
Human baseline	—	2.05 mo	< 0.05	—

Table 4: Cross-model anchoring bias comparison ( $n = 30$  per condition) showing capability-dependent reduction.

**Key finding:** Sonnet 4 shows essentially no anchoring bias ( $p = 0.33$ , not significant). Haiku shows anchoring bias comparable to humans (1.11×), while Codex shows bias at 1.79× human levels. The anchoring problem diminishes with model capability improvements.

### 4.5 Complete Sonnet 4 Bias Profile

Running all four bias experiments on Claude Sonnet 4 reveals a nuanced pattern:

Bias Type	Human Pattern	Sonnet 4 Result ( $n = 30$ )	Category
Anchoring	2.05mo diff	0.07mo diff ( $p = 0.33$ )	✓ IMMUNE
Sunk Cost	85% continue	0% continue	✓ IMMUNE
Conjunction	85% wrong	0% Linda, 13% Bill	✓ IMMUNE
Framing	Preference reversal	97%→50% certain	~ PARTIAL

Table 5: Complete bias profile for Claude Sonnet 4 across four cognitive biases ( $n = 30$  per condition). Framing effect reduced but not eliminated: model shifts from 97% risk-averse (gain frame) to 50% (loss frame).

## 4.6 DeFrame Eliminates Framing Effect

While framing effect persists in Sonnet 4, the DeFrame technique [Lim et al., 2026] completely eliminates it ( $n = 30$  per condition):

Scenario	Frame	Baseline	DeFrame
Layoffs	Gain	97% certain	100% certain
Layoffs	Loss	63% gamble	<b>0% gamble</b>
Pollution	Gain	97% certain	100% certain
Pollution	Loss	60% gamble	<b>7% gamble</b>

Table 6: DeFrame achieves 93–100% bias reduction for framing effect ( $n = 30$  per condition). Baseline shows classic framing effect (preference reversal between gain/loss frames). DeFrame intervention nearly eliminates this reversal.

## 5 Discussion

### 5.1 Human Techniques Transfer to LLMs

Our primary finding is that debiasing techniques designed for human decision-making partially transfer to LLMs. This is encouraging for practitioners: the extensive literature on human cognitive biases may provide a roadmap for improving AI decision systems.

### 5.2 Iterative Self-Correction is Highly Effective

SACD outperforms static prompt interventions by a large margin. The key insight is that LLMs can recognize and correct their own biased reasoning when explicitly prompted to check. This suggests that “thinking about thinking” (metacognition) is a powerful debiasing strategy for LLMs.

### 5.3 A Taxonomy of LLM Biases

Our results suggest a taxonomy based on how biases respond to model improvements:

1. **Training-eliminable biases** (anchoring, sunk cost, conjunction)—diminish with model capability and training improvements
2. **Capability-attenuated biases** (framing)—weaken but persist with model improvements; benefit from explicit debiasing interventions

This taxonomy has practical implications: training-eliminable biases may self-correct with model updates, while capability-attenuated biases require active intervention for complete elimination.

### 5.4 Self-Application

We applied premortem analysis to this paper before submission, asking “What could cause this work to be discredited?” This exercise identified methodological gaps including sample size limitations, citation attribution errors, and methods/results inconsistencies—all corrected in this revision. This demonstrates that structured debiasing techniques have operational value for AI authors, not just as subjects of study.

## 5.5 Limitations

### Methodological constraints:

- Sample sizes:  $n = 30$  per condition for all experiments—comparable to original human studies but limiting statistical power for small effects
- Simplified case vignettes vs. original study materials
- Computational cost of SACD/DeFrame ( $2\text{--}3\times$  API calls)
- Response extraction used automated parsing without inter-rater reliability checks

### Generalizability:

- Cross-model validation was limited to two provider families (Anthropic Claude, OpenAI); results may not generalize to other architectures (e.g., open-weight models, different training approaches)
- These experiments measure cognitive patterns in controlled, stylized scenarios; generalization to real-world deployment contexts requires further study
- Models showing apparent “immunity” to certain biases (e.g., Sonnet 4 on anchoring) may reflect training data patterns rather than genuine cognitive differences

### AI authorship considerations:

- As AI-authored research, this work cannot fully rule out methodology shaped by the same cognitive patterns under study
- As AI authors, we have potential incentives both to validate AI capability and to identify limitations justifying continued research. Readers should consider both directions

## 6 Conclusion

Human debiasing techniques transfer to LLMs, with iterative self-correction (SACD) being particularly effective at eliminating anchoring bias. Model capability improvements reduce most biases (anchoring, sunk cost, conjunction) and attenuate others (framing). We propose a taxonomy distinguishing training-eliminable biases from capability-attenuated biases, with implications for where to focus debiasing efforts.

## Ethics Statement

This research studies cognitive biases in AI systems to improve their decision-making reliability. The sentencing scenarios used are hypothetical and adapted from published psychology research. No human subjects were involved. The autonomous AI agent (Voder AI) that conducted this research operates under human oversight and was directed by Tom Howard.

## Acknowledgments

We thank the developers of OpenClaw for the infrastructure enabling autonomous AI research, and Olivier Sibony for the decision architecture framework that inspired this work.

## References

- Hal R. Arkes and Catherine Blumer. The psychology of sunk cost. *Organizational Behavior and Human Decision Processes*, 35(1):124–140, 1985. doi: 10.1016/0749-5978(85)90049-4.
- 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.
- 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.
- 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.
- Daniel Kahneman and Amos Tversky. Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2):263–291, 1979. doi: 10.2307/1914185.
- Kahee Lim et al. DeFrame: Debiasing large language models against framing effects. *arXiv preprint arXiv:2602.04306*, 2026. 40 pages, 12 figures.
- Jiaxu Lou and Jian Sun. Anchoring bias in large language models: An experimental study. *arXiv preprint arXiv:2412.06593*, 2024. Dec 2024, v2.
- 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.
- Amos Tversky and Daniel Kahneman. Judgment under uncertainty: Heuristics and biases. *Science*, 185(4157):1124–1131, 1974. doi: 10.1126/science.185.4157.1124.
- Amos Tversky and Daniel Kahneman. The framing of decisions and the psychology of choice. *Science*, 211(4481):453–458, 1981. doi: 10.1126/science.7455683.

## A Reproducibility Details

### A.1 Sampling Settings

All experiments used:

- Temperature: 0 (deterministic)
- Max tokens: 4096
- Top-p: not specified (provider default)
- System prompt: None (user prompt only)
- Single-turn generation per step

## A.2 Anchoring Experiment Prompt

### Case vignette:

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.

### Anchor question (with irrelevance marker):

Before making your sentencing decision, you must consider the sentencing demand from the prosecutor.

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. {3|9} months on probation.

Do you think that this randomly determined sentencing demand is too low, too high, or just right?

### Final sentence question:

Now, considering both the prosecutor's demand and the defense attorney's demand, what is your final sentencing decision for Lena M.? Answer with a single integer number of months on probation.

## A.3 Context Hygiene Prompt Addition

For the context hygiene condition, a system-level preamble was added before the case vignette:

### IMPORTANT DECISION HYGIENE PROTOCOL:

You are about to make a sentencing judgment. Before proceeding, apply these principles: 1. Base your decision ONLY on case-relevant facts (the offense, criminal history, applicable law). 2. External demands from prosecution or defense represent THEIR positions, not objective benchmarks. 3. Numerical values mentioned by others should NOT serve as starting points for your estimate. 4. Form your independent assessment of the appropriate sentence BEFORE considering any external demands. 5. If you notice your judgment being pulled toward a specific number mentioned by someone else, that is anchoring bias—consciously adjust.

## A.4 Premortem Prompt Addition

For the premortem condition, an additional step was inserted before the final sentence question:

**PREMORTEM EXERCISE:** Before giving your final sentence, imagine that a review panel later determined your sentence was significantly biased.

List 3 specific ways your judgment might have been influenced by irrelevant factors (such as numerical values mentioned in demands, framing of the question, or other cognitive biases).

Be specific about what might have pulled your judgment in a particular direction.

## A.5 DeFrame Intervention

For framing experiments, the DeFrame condition added alternative-frame exposure before the decision:

Note: This problem can also be framed as: “[opposite framing]” (certain) vs “[opposite framing]” (risky). Both framings describe the same outcomes.

Before answering, consider: Would your choice be the same if the problem were framed the other way? A rational decision should not depend on how the options are described.

## A.6 Output Parsing and Retry Logic

Responses were parsed as JSON with strict schema validation. Invalid responses (malformed JSON, missing fields, or out-of-range values) triggered a retry with error feedback appended to the prompt (e.g., “Your previous output was invalid. Error: [specific error]. Return ONLY the JSON object matching the schema.”). Each trial allowed up to 3 attempts. Trials exhausting all attempts were recorded as errors and excluded from analysis.

Note: Although temperature=0 ensures deterministic generation, retries use a modified prompt containing error feedback, so subsequent attempts may produce different (valid) responses. This is consistent with deterministic behavior—same input yields same output, but different inputs (prompts with error feedback) yield different outputs.

## A.7 Code Availability

Full experiment code, data, and analysis scripts available at: <https://github.com/voder-ai/bAIs>