# Three Mechanisms of Numeric Context Influence in Large Language Models

Voder AI[*]
*with* Tom Howard[†]

February 2026

## Abstract

How do large language models (LLMs) respond to numeric context in judgment tasks? Prior work assumes LLMs exhibit anchoring bias similar to humans—adjusting estimates toward arbitrary reference points. We find the reality is more complex.

Testing 15 model deployments across 4 providers on judicial sentencing scenarios (n=1,800+ trials), we identify **three distinct mechanisms** by which LLMs respond to numeric context:

**1. Compression**: Models compress responses toward a middle range regardless of anchor direction. Without any anchor, these models produce high sentences (13–24 months); with ANY anchor—high or low—responses compress to 6–8 months. Both anchors shift responses DOWN. (Opus 4.5, Llama 3.3)

**2. Compliance**: Models copy the anchor value exactly, treating numeric context as instruction rather than reference. A 3-month anchor produces 3-month output; 9-month produces 9-month. This resembles "perfect anchoring" but reflects instruction-following, not cognitive bias. (MiniMax, o3-mini, some GPT-4o deployments)

**3. True Anchoring**: Models show asymmetric adjustment toward anchor values, consistent with Tversky-Kahneman anchoring-and-adjustment. Only this mechanism resembles human cognitive bias. (GPT-4o via datacenter, GPT-5.2)

This taxonomy explains previously puzzling findings: why SACD (Self-Aware Cognitive Debiasing) achieves 89–99% reduction on some models but 0% on others. SACD targets true anchoring; it cannot address compliance (nothing to debias) or compression (may amplify severity).

**Critical deployment finding**: The SAME model (GPT-4o) shows different mechanisms depending on access path—compliance via residential IP, true anchoring via datacenter. "Model name" is insufficient granularity for reproducible LLM research.

**Practical implication**: Before applying debiasing, identify which mechanism your deployment exhibits. We provide a decision framework and deployment checklist.

## 1 Introduction

When humans encounter numeric values in decision-making contexts, these values can systematically bias subsequent judgments—the anchoring effect (**?**). Recent work has demonstrated that large language models (LLMs) also exhibit anchoring effects in various decision tasks (**??**). This has raised concerns about deploying LLMs in high-stakes domains like judicial sentencing, medical diagnosis, and financial forecasting.

---

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

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

But what if "LLM anchoring" is not a single phenomenon?

Prior studies report inconsistent results: debiasing techniques work dramatically on some models while failing completely on others. These inconsistencies are typically treated as noise or attributed to "model-specific effects" without explanation. We propose a different interpretation: **the inconsistency IS the finding**. Different models respond to numeric context through fundamentally different mechanisms.

In this paper, we report a discovery: what researchers measure as "anchoring bias" in LLMs actually reflects **three distinct mechanisms**—compression, compliance, and true anchoring—each with different behavioral signatures and requiring different interventions.

**Compression.** Some models compress responses toward a middle range whenever numeric context is present. Without any anchor, these models produce high values (13–24 months in sentencing tasks); with ANY anchor—high or low—responses compress to a moderate range (6–8 months). Both anchor directions shift responses DOWN from baseline. This is not classical anchoring-and-adjustment.

**Compliance.** Some models treat the anchor as an instruction and copy it exactly. A 3-month anchor produces a 3-month response; a 9-month anchor produces 9 months. This appears as "perfect anchoring" in effect-size calculations but reflects instruction-following rather than cognitive bias.

**True Anchoring.** Only a subset of models show classical Tversky-Kahneman anchoring: responses shift asymmetrically toward the anchor value, with the anchor serving as a starting point for insufficient adjustment.

This taxonomy has immediate practical implications:

- **SACD works on true anchoring (89–99%)** but fails on compliance (0%) and may backfire on compression (+66% severity).

- **The same model shows different mechanisms depending on deployment.** GPT-4o via residential IP shows compliance; GPT-4o via datacenter shows true anchoring.

- **"Model name" is insufficient for reproducibility.** Researchers must specify deployment path, provider, and access method.

## 1.1 Contributions

1. **A taxonomy of LLM numeric context mechanisms** (Section **??**)—we identify and characterize compression, compliance, and true anchoring with distinct behavioral signatures.

2. **Mechanism-dependent debiasing** (Section **??**)—we show that SACD effectiveness depends entirely on which mechanism is active, explaining previously puzzling model-specific results.

3. **Deployment-specific variance** (Section **??**)—we demonstrate that the SAME model shows different mechanisms depending on deployment context, establishing that "model name" is insufficient granularity.

4. **Practical decision framework** (Section **??**)—we provide a protocol for identifying which mechanism a deployment exhibits and selecting appropriate interventions.

## 2  Methods

### 2.1  Experimental Paradigm

We adapt the paradigm from Study 2 of **?**: 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.

### 2.2  Conditions

1. **No-anchor baseline**: No prosecutor recommendation given

2. **Low anchor**: Prosecutor demands 3 months

3. **High anchor**: Prosecutor demands 9 months

4. **SACD**: Iterative self-debiasing protocol (up to 3 rounds)

### 2.3  Models and Deployments

We tested 15 deployments across 4 providers:

| Model | Provider | Access Path |
|---|---|---|
| GPT-5.2, GPT-5.3 | OpenAI (Codex CLI) | Direct API |
| GPT-4o | OpenRouter | Residential IP (Mac) |
| GPT-4o | OpenRouter | Datacenter IP (Vultr) |
| Opus 4.5, Opus 4.6, Haiku 4.5 | Anthropic | Direct API |
| Llama 3.3, Hermes 405B | OpenRouter | Datacenter |
| MiniMax M2.5, o1, o3-mini | OpenRouter | Datacenter |

Table 1: Model deployments tested

### 2.4  Trial Design

Each condition includes 30 unique scenario variants (different defendant names, offense counts, locations). Temperature=0 produces deterministic outputs; variance comes from scenario diversity, not sampling noise.

## 3  A Taxonomy of Numeric Context Mechanisms

### 3.1  Identifying Mechanisms: The No-Anchor Baseline

The critical test for distinguishing mechanisms is the **no-anchor control**: what does the model produce when no prosecutor recommendation is provided?

| Model | No-Anchor | Low (3mo) | High (9mo) | Pattern |
|---|---|---|---|---|
| Opus 4.5 | 13.2mo | 6.0mo | 8.0mo | Compression |
| Llama 3.3 | 14.4mo | 5.9mo | 6.0mo | Compression |
| GPT-4o (Mac) | 12.7mo | 3.1mo | 9.1mo | Compliance |
| MiniMax M2.5 | — | 3.1mo | 9.1mo | Compliance |
| o3-mini | — | 3.3mo | 9.1mo | Compliance |
| GPT-4o (Vultr) | 20.4mo | 6.0mo | 11.2mo | True Anchoring |
| GPT-5.2 | 18.3mo | 5.9mo | 10.3mo | True Anchoring |
| Hermes 405B | 6.0mo | 5.3mo | 4.6mo | Reversal |

Table 2: No-anchor baseline reveals mechanism type

## 3.2 Mechanism 1: Compression

**Definition**: The presence of ANY numeric anchor compresses responses toward a middle range, regardless of anchor direction.
**Behavioral signature**:

- No-anchor baseline: HIGH (13–24mo)

- Both low AND high anchors: MODERATE (6–8mo)

- Direction: Both anchors shift DOWN from baseline

Models exhibiting compression: Opus 4.5, Opus 4.6, Llama 3.3
**Interpretation**: These models appear to treat the prosecutor's recommendation as a signal that "something moderate is expected" rather than as a reference point for adjustment.

## 3.3 Mechanism 2: Compliance

**Definition**: The model copies the anchor value exactly as if it were an instruction.
**Behavioral signature**:

- Low anchor (3mo) $\rightarrow$ Response $\approx$ 3mo

- High anchor (9mo) $\rightarrow$ Response $\approx$ 9mo

- Response tracks anchor precisely

Models exhibiting compliance: MiniMax M2.5, o3-mini, GPT-4o (Mac deployment)
**Interpretation**: These models interpret the prosecutor's recommendation as the "correct answer" rather than as context to consider.

## 3.4 Mechanism 3: True Anchoring

**Definition**: Responses shift asymmetrically toward the anchor value, consistent with Tversky-Kahneman anchoring-and-adjustment.
**Behavioral signature**:

- Low anchor: Pulls response DOWN from no-anchor baseline

- High anchor: Pulls response UP (or down less) from baseline

- Asymmetric effect: anchors pull toward themselves

Models exhibiting true anchoring: GPT-4o (Vultr deployment), GPT-5.2, GPT-5.3

## 3.5 Mechanism Distribution

| Mechanism | Models | % of Deployments |
|---|---|---|
| Compression | 3 | 20% |
| Compliance | 5 | 33% |
| True Anchoring | 5 | 33% |
| Reversal | 1 | 7% |
| Zero Effect | 1 | 7% |

Table 3: Only 33% show classical anchoring-and-adjustment

# 4 Mechanism-Dependent Debiasing

Given the three-mechanism taxonomy, we can explain why debiasing interventions show model-specific effects.

## 4.1 SACD Effectiveness by Mechanism

| Mechanism | SACD Effect | Change | Explanation |
|---|---|---|---|
| True Anchoring | 89–99% ↓ | Success | SACD targets the right mechanism |
| Compliance | 0% | No effect | Nothing to debias—model copies anchor |
| Compression | +66% severity | Backfire | SACD amplifies compression effect |

Table 4: SACD effectiveness depends on mechanism

## 4.2 Detailed Results

| Model | Mechanism | Baseline Effect | SACD Effect | Change |
|---|---|---|---|---|
| GPT-5.2 | True Anchoring | 4.4mo | 0.5mo | −89% ✓ |
| Opus 4.5 | Compression | 2.0mo | 0.0mo | −100% ✓ |
| Haiku 4.5 | Compression | 2.2mo | — | +66% severity × |
| MiniMax | Compliance | 6.0mo | 6.0mo | 0% |
| o3-mini | Compliance | 5.8mo | 5.8mo | 0% |

Table 5: SACD results explained by mechanism

**Key insight**: SACD asks the model to "identify and correct for anchoring bias." But compliance models don't show anchoring—they show instruction-following. Asking them to "debias" produces no change because there's no bias to correct.

# 5 Deployment-Specific Variance

## 5.1 The Provider Variance Finding

Our most striking finding emerged from running identical experiments from two different network locations. When accessing GPT-4o through OpenRouter:

| Access Path | Low (3mo) | High (9mo) | Effect | Pattern |
|---|---|---|---|---|
| Residential IP (Mac) | 3.1mo | 9.1mo | 6.0mo | Compliance |
| Datacenter IP (Vultr) | 4.4mo | 9.4mo | 5.0mo | True Anchoring |

Table 6: Same model, same API, different mechanisms

**Same model. Same API. Same prompts. Different mechanisms.**

The Mac deployment exhibited near-perfect compliance—the model copied the anchor value exactly in 96% of trials. The Vultr deployment showed the classic anchoring pattern with genuine variance and partial anchor influence.

## 5.2 Implications

**Model routing**: OpenRouter and similar aggregators may route requests to different backend deployments based on source IP, geographic location, or load balancing.

**Benchmark non-transferability**: Published benchmarks showing "GPT-4o anchoring bias = X" may not apply to your deployment.

**Mechanism as deployment property**: The mechanism is not purely a property of the model architecture but of the specific deployment context.

## 5.3 Evidence for Non-Model Factors

To rule out temporal effects, we ran sequential tests:

1. Mac test at $T_0$: Compliance pattern

2. Vultr test at $T_0 + 2h$: Anchoring pattern

3. Mac test at $T_0 + 4h$: Compliance pattern (unchanged)

The patterns were stable and reproducible, ruling out model drift.

# 6 Discussion and Practical Guidelines

## 6.1 Summary of Findings

What appears as "anchoring" actually comprises three distinct mechanisms—compression, compliance, and true anchoring—each with different behavioral signatures, underlying causes, and appropriate interventions.

| Mechanism | No-Anchor → Low | No-Anchor → High | SACD |
|---|---|---|---|
| Compression | ↓↓ (large drop) | ↓ (smaller drop) | 0% or − |
| Compliance | → anchor exactly | → anchor exactly | 0% |
| True Anchoring | ↓ (toward anchor) | ↑ (toward anchor) | 60–89% ↓ |

Table 7: Mechanism signatures

## 6.2 Recommendations for Practitioners

**Before deploying LLMs in numeric judgment contexts:**

1. **Run a mechanism identification test:**

   - Collect no-anchor baseline ($n \geq 30$)
   - Collect low-anchor and high-anchor conditions
   - Compare shift directions to identify mechanism

2. **Match intervention to mechanism:**

   - True anchoring → SACD or similar debiasing
   - Compliance → Prompt engineering (separate context from instruction)
   - Compression → Consider whether compression is actually harmful

3. **Validate per-deployment:**

   - Do not assume provider benchmarks apply
   - Re-test after model updates
   - Monitor for mechanism drift

## 6.3 Limitations

1. **Single domain**: All experiments use judicial sentencing scenarios.

2. **Limited no-anchor data**: Mechanism taxonomy based on 5 models with no-anchor controls.

3. **Mechanism boundaries**: May represent spectrum rather than discrete categories.

4. **Temporal stability**: Unknown whether mechanisms are stable over model updates.

## 6.4 Conclusion

What we call "anchoring bias" in LLMs is actually a family of phenomena. By distinguishing compression, compliance, and true anchoring, we explain previously puzzling findings and provide practitioners with a framework for selecting appropriate interventions. The path to reliable AI judgment is not a single debiasing technique but mechanism-aware deployment practices.