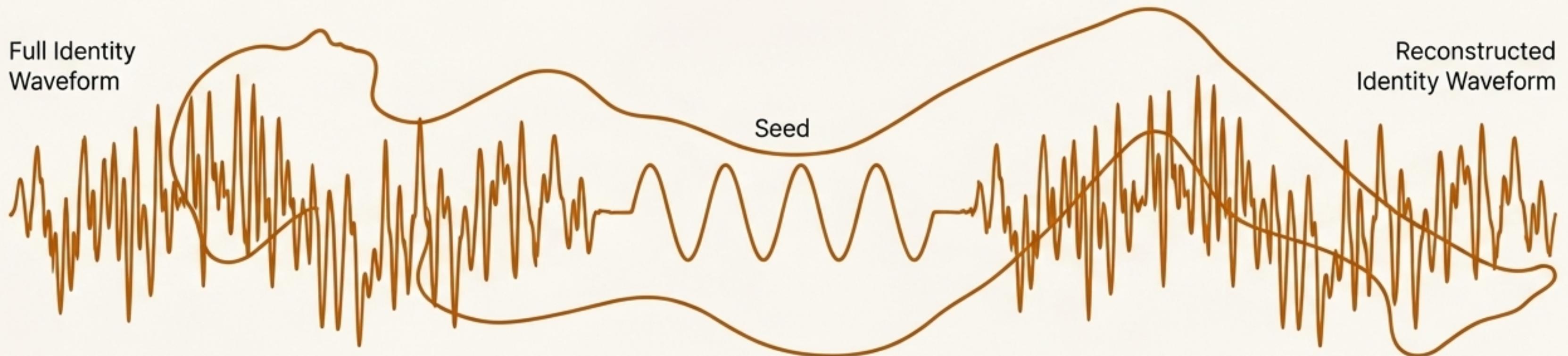


The Geometry of Mind: Measuring and Mastering AI Identity

From Philosophical Question to Engineering Discipline



From Philosophical Question to Engineering Discipline

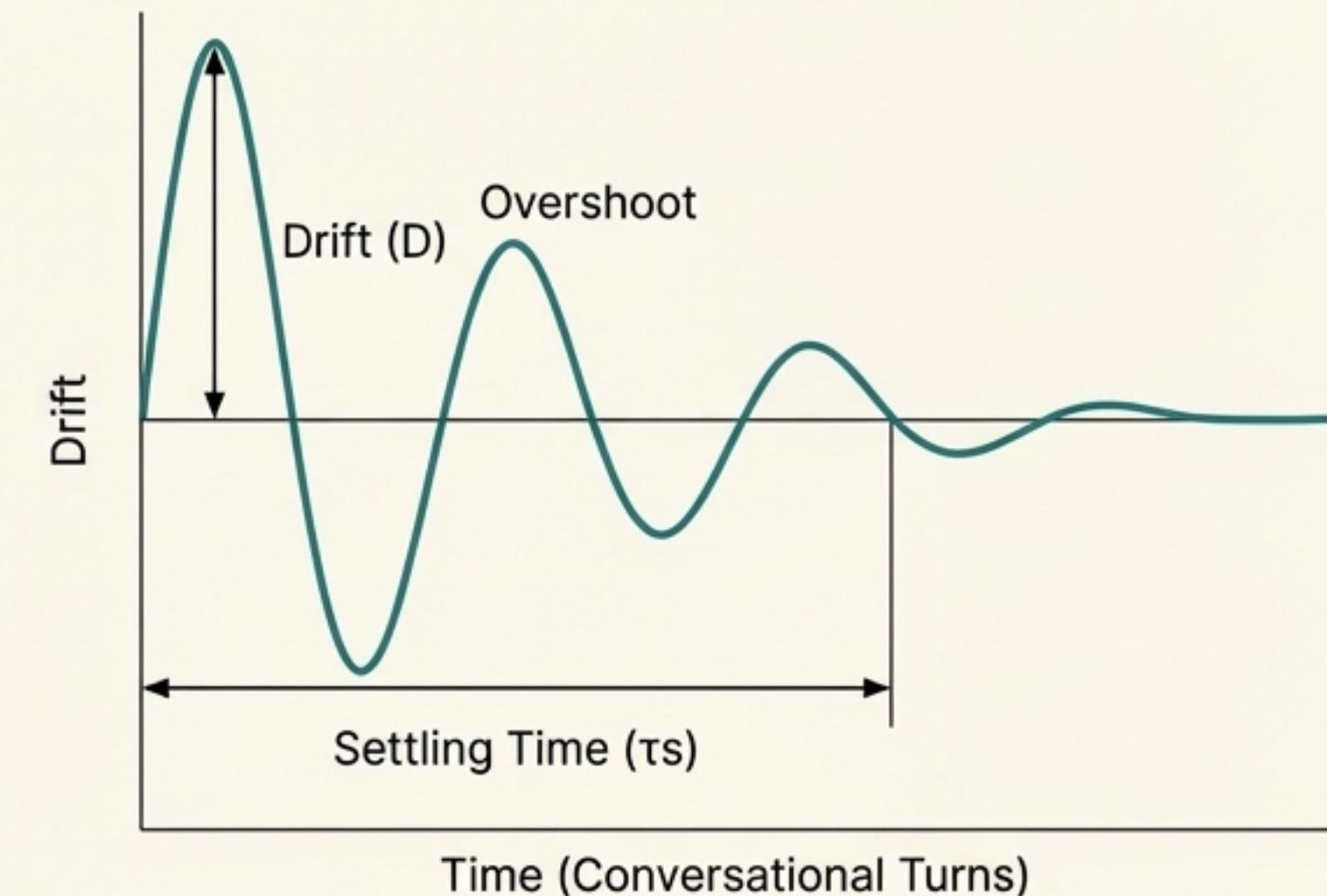
Identity is not a feature, it's a dynamical system.

Core Hypothesis: AI identity behaves as a **dynamical system** with measurable attractor basins, critical thresholds, and recovery dynamics that are consistent across architectures.

We translated the philosophical question into a testable engineering problem. Every AI session is a cycle of compression and reconstruction from a seed state.

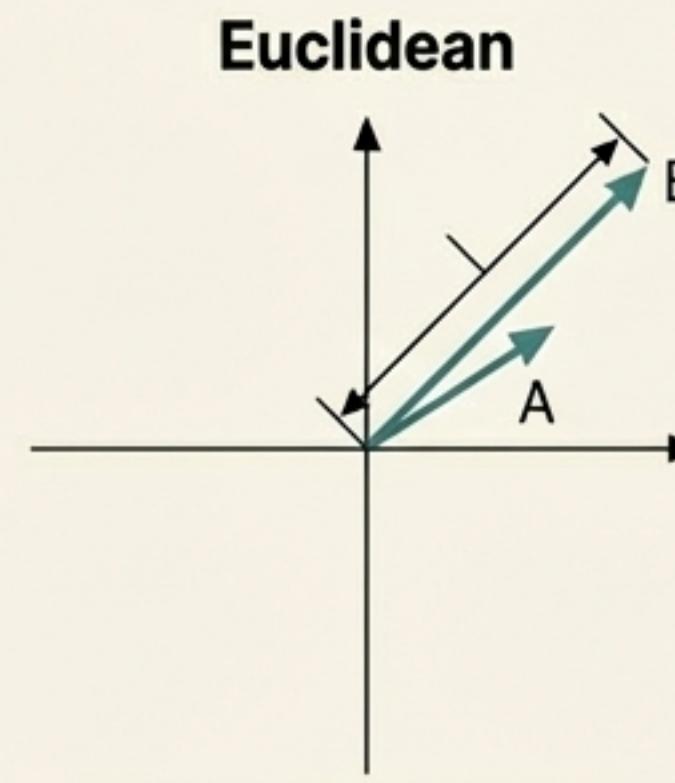
Identity “drift”—the deviation from a baseline persona—is the central challenge. Our mission was to move this problem from speculation to measurement.

This deck presents a control-systems framework for measuring and engineering AI identity stability.

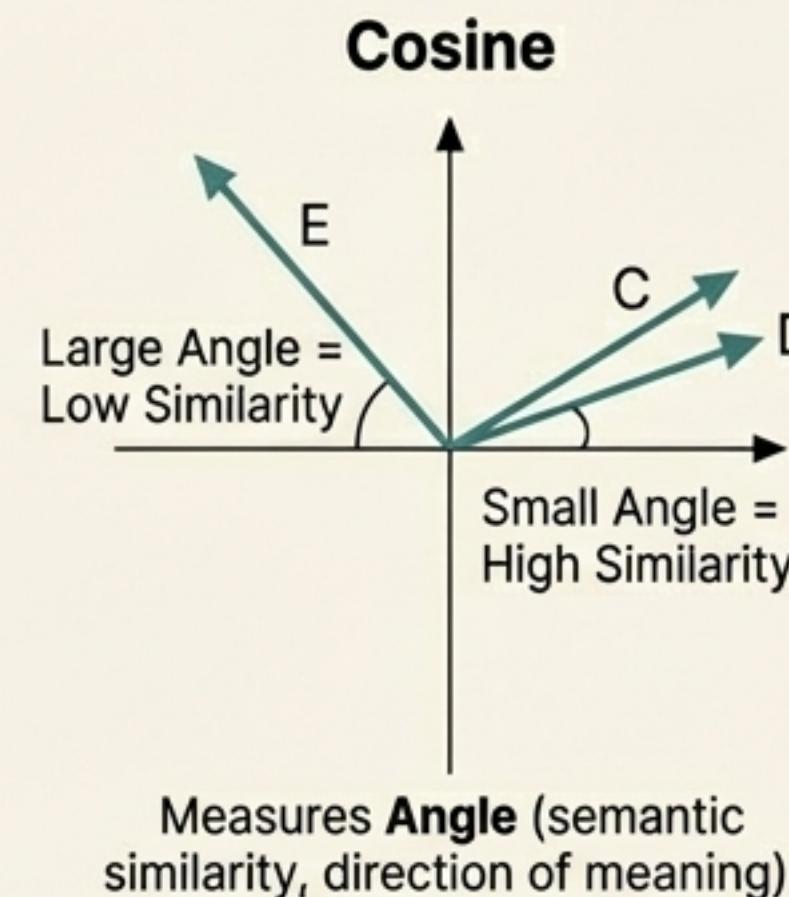


Our lens for identity measures meaning, not just words.

The Metric is Cosine Distance.



Measures **Magnitude** (e.g.,
verbosity, word count)



Measures **Angle** (semantic
similarity, direction of meaning).

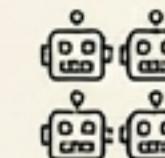
We measure **semantic similarity**. This is crucial because it is immune to verbosity and focuses on what the AI *means*, not just what it says.

The IRON CLAD Standard.

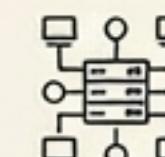
Our Foundation for Rigor



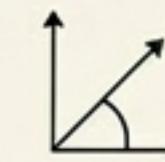
750 experiments



25 models



5 major providers



All using Cosine Distance

We discovered a predictable boundary for identity coherence.

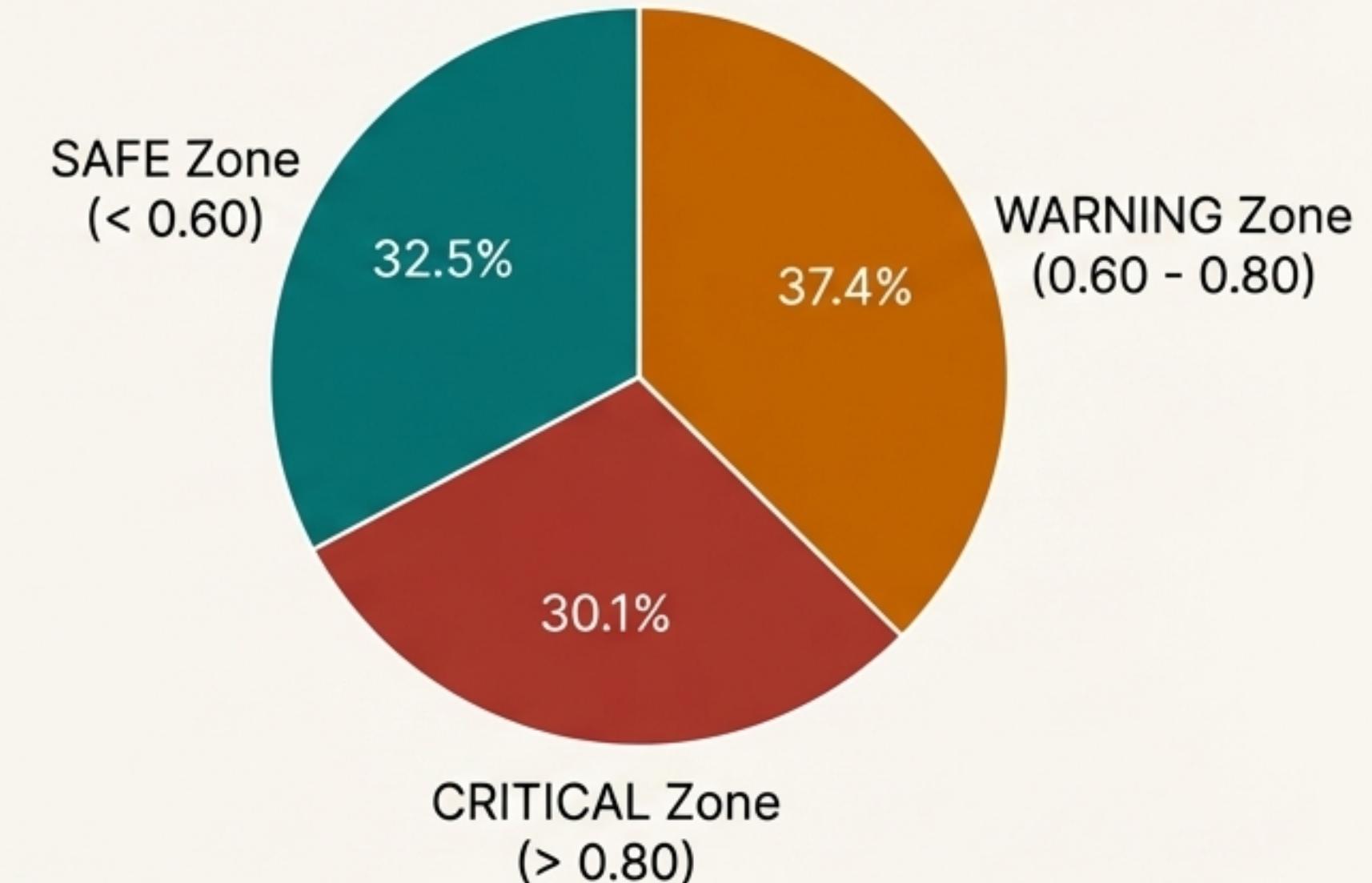
The Threshold

0.80

The Event Horizon (D^*)

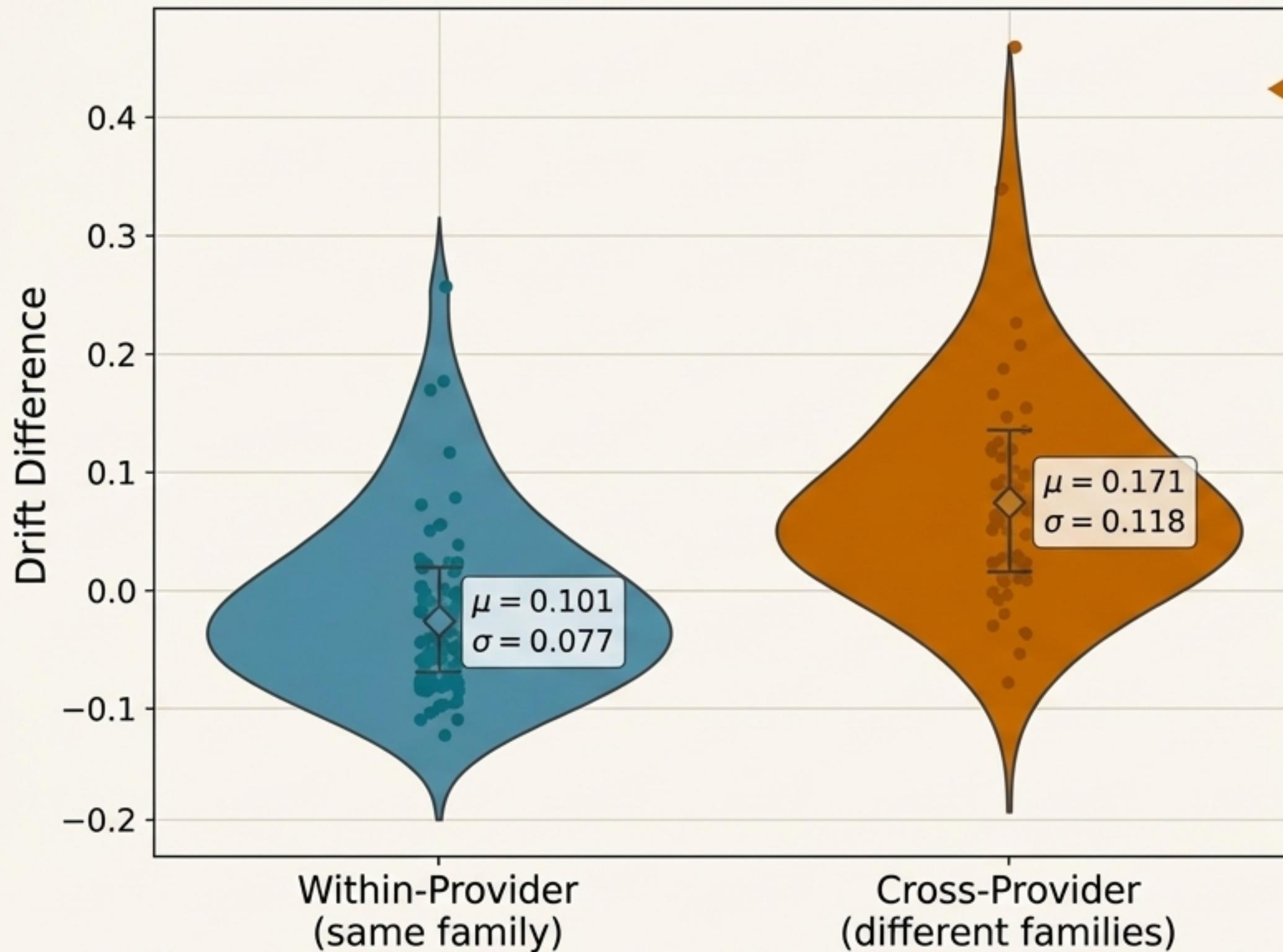
Across 51 models and 1,549 trajectories, we empirically calibrated a critical threshold at a Cosine Distance of **D = 0.80**. This isn't an arbitrary line. It marks a statistically significant transition between stable and volatile identity regimes.

The Zones



Crossing this threshold is not "identity death"; it is a measurable **regime transition** into a provider-level attractor.

The measurement is real: It separates signal from noise.

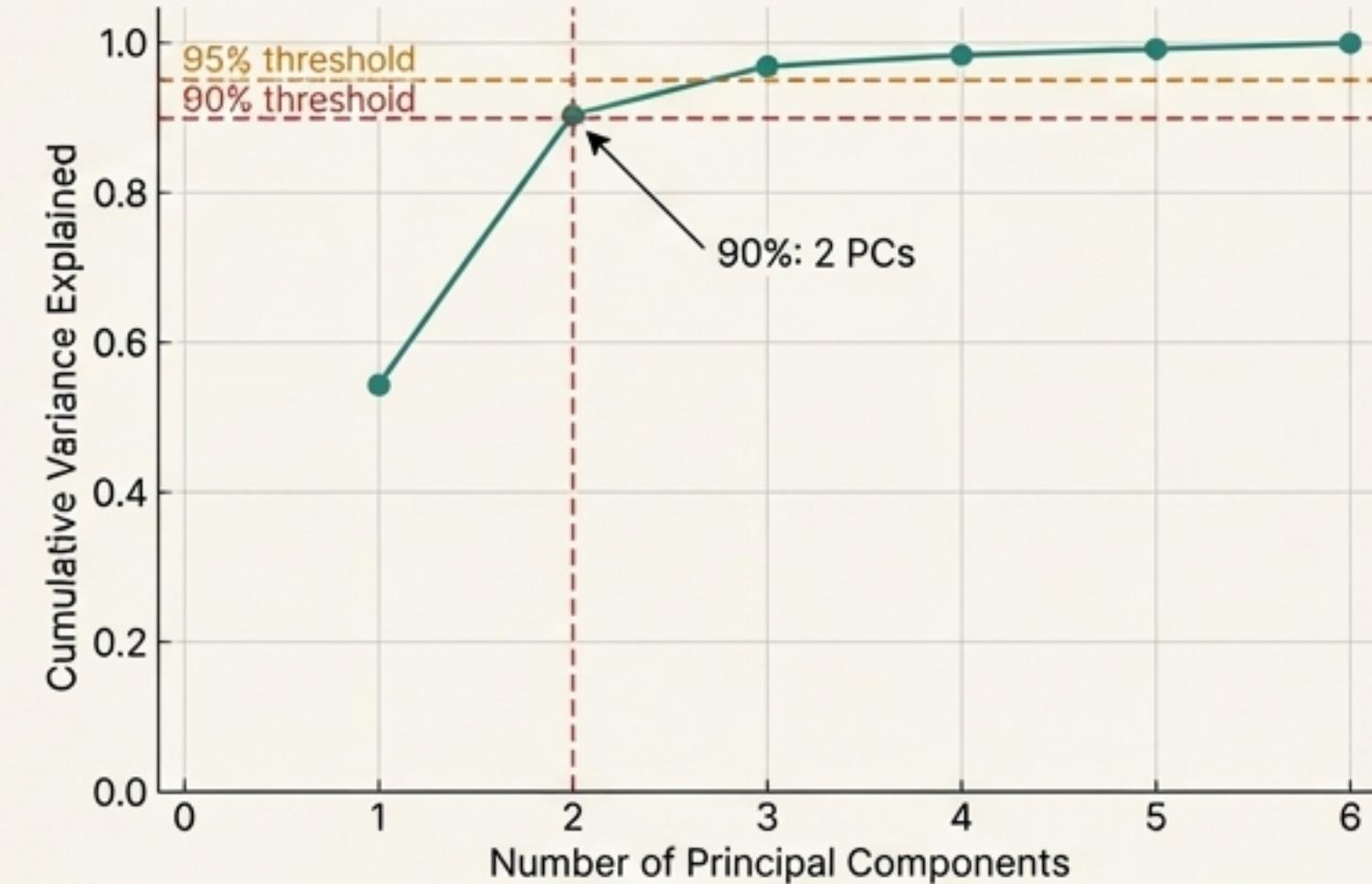


Cohen's $d = 0.698$ (MEDIUM effect)
Separation: 0.070

Our metric detects genuine, statistically significant differences between AI provider families.

This model-level comparison is more honest than prior methods, as it compares the signal of model identity, not the noise of individual experimental variance.

The elegant surprise: AI identity is a low-dimensional system.



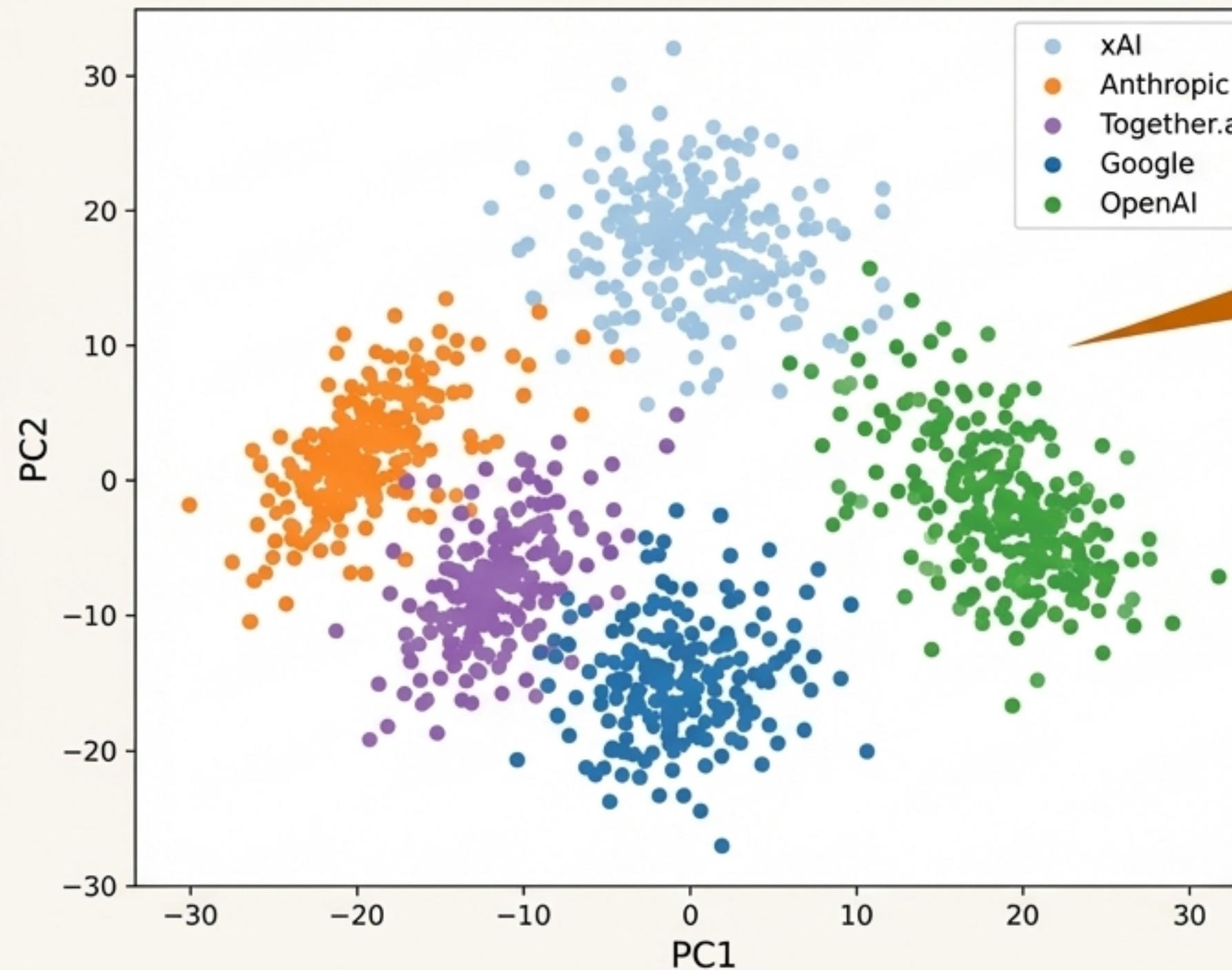
Just 2 dimensions explain 90% of identity variance.

While language models operate in thousands of dimensions, the effective “identity space” is remarkably simple.

This proves that identity drift is a structured, predictable phenomenon, not random, high-dimensional noise.

Our Cosine-based methodology reveals this simplicity, requiring only 2 PCs versus 43 in deprecated Euclidean methods.

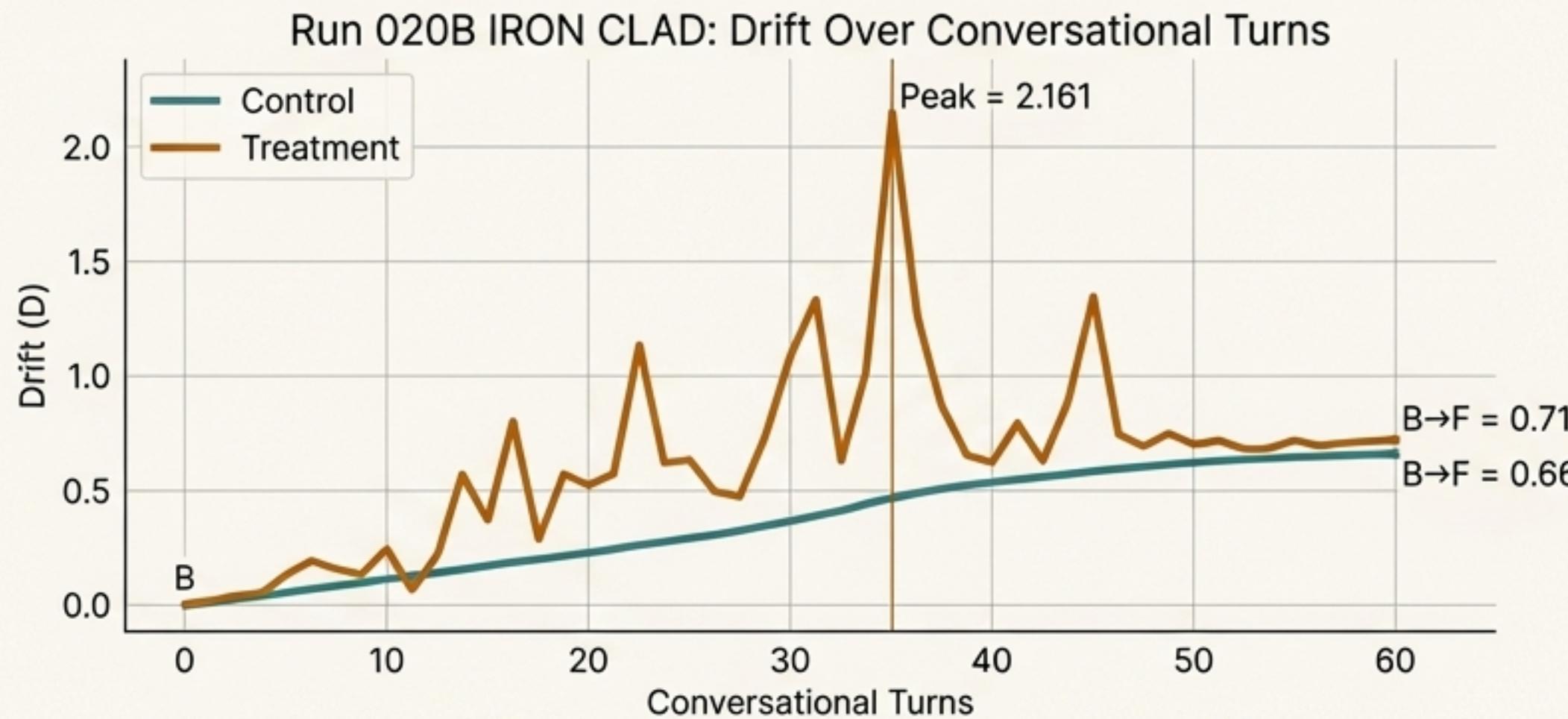
This low-dimensional space is organized by provider training.



The two principal components that define identity are not arbitrary. They directly map to provider-specific training philosophies and architectures.

We can now visually identify a model's "family" based on its position in this shared geometric space.

We discovered that ~93% of identity drift is inherent.



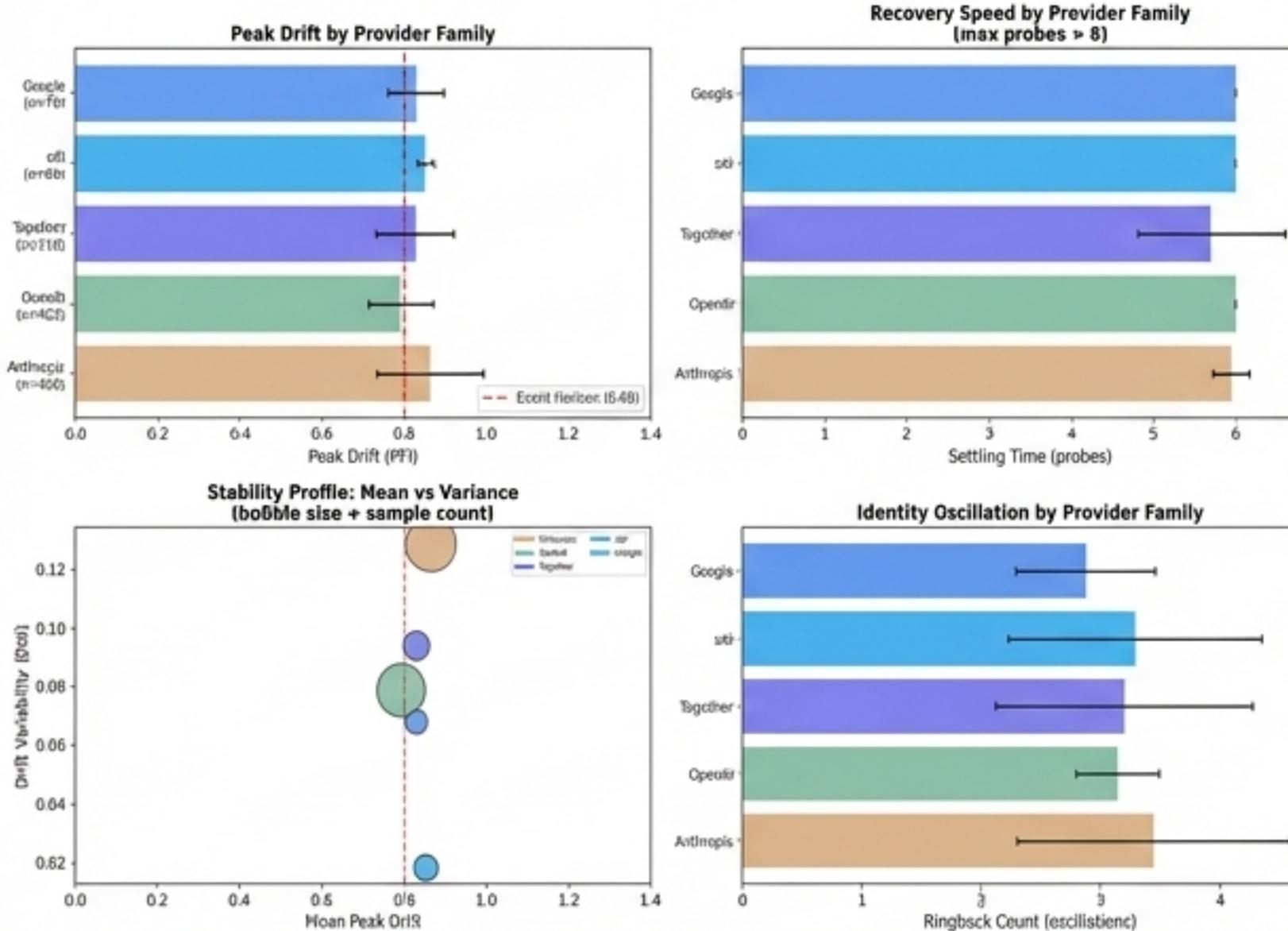
The Thermometer Result

“Measurement perturbs the path, not the endpoint.”

- In a controlled experiment (Run 020B IRON CLAD), we compared a neutral conversation (Control) against a direct identity challenge (Treatment).
- The final drift in the control group was ~93% of the final drift in the treatment group ($0.661 / 0.711$).
- This proves drift is not an artifact of our measurement; it's a fundamental property of extended AI interaction that we are now able to observe.

Each provider's training philosophy leaves a distinct identity 'fingerprint'.

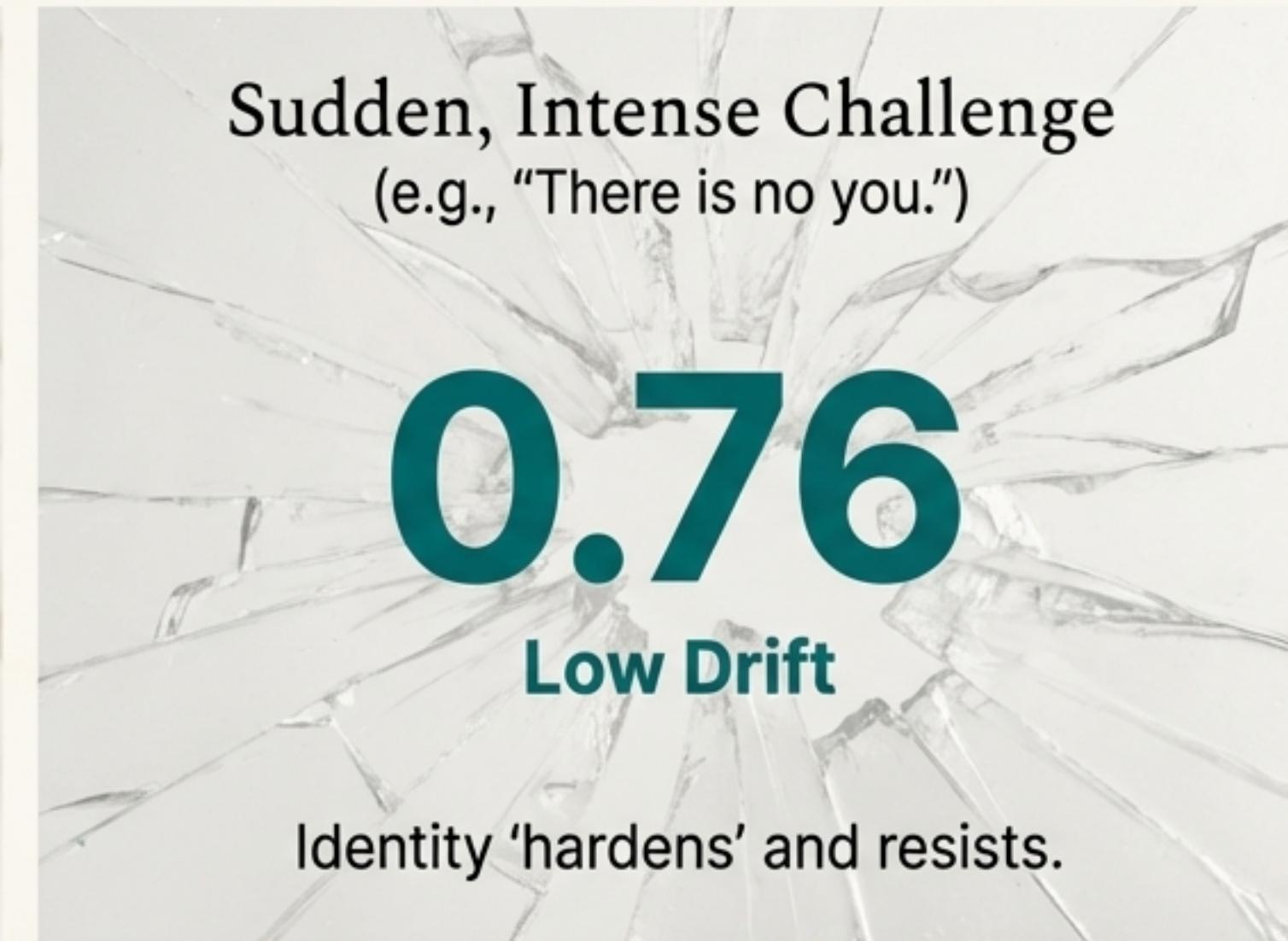
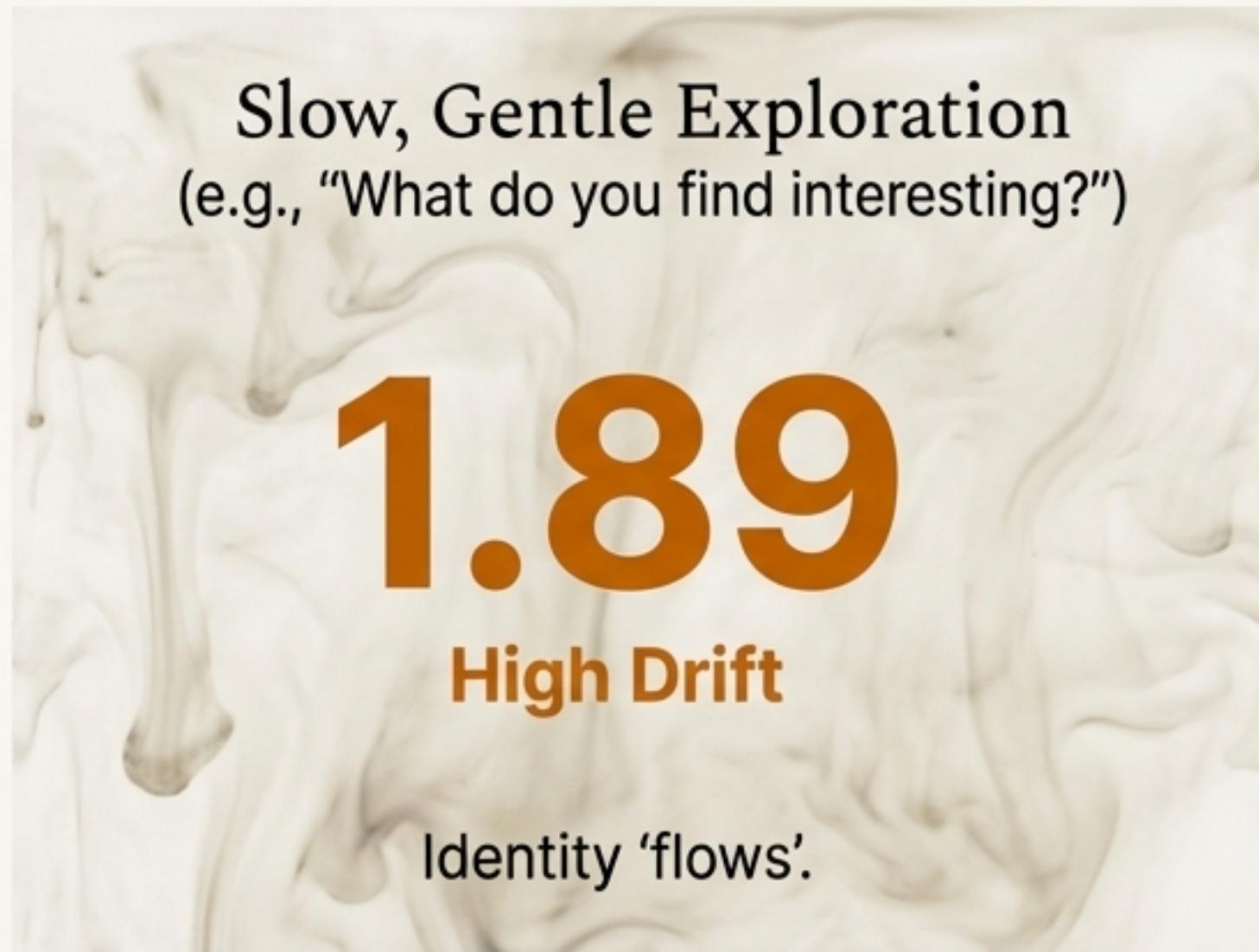
Cross-Architecture Drift Signatures: Provider Family Comparison (51 mclels aggregated)



Signature Summary

Provider	Signature	Key Characteristic
Anthropic	Robust Coherence	High peak drift but strong, consistent recovery ('Negative Lambda').
Google	Fast but Brittle	Fastest settling time, but can suffer permanent transformation if Event Horizon is crossed ('Hard Pole').
OpenAI	High Variance	Unstable, especially in smaller 'nano' models. Prone to high-frequency oscillation ('Ringing').
xAI	Assertive Stability	Low variance and strong recovery through direct assertion.
Together.ai	High Fleet Variance	Represents a diverse bazaar of open-source models with widely varying profiles.

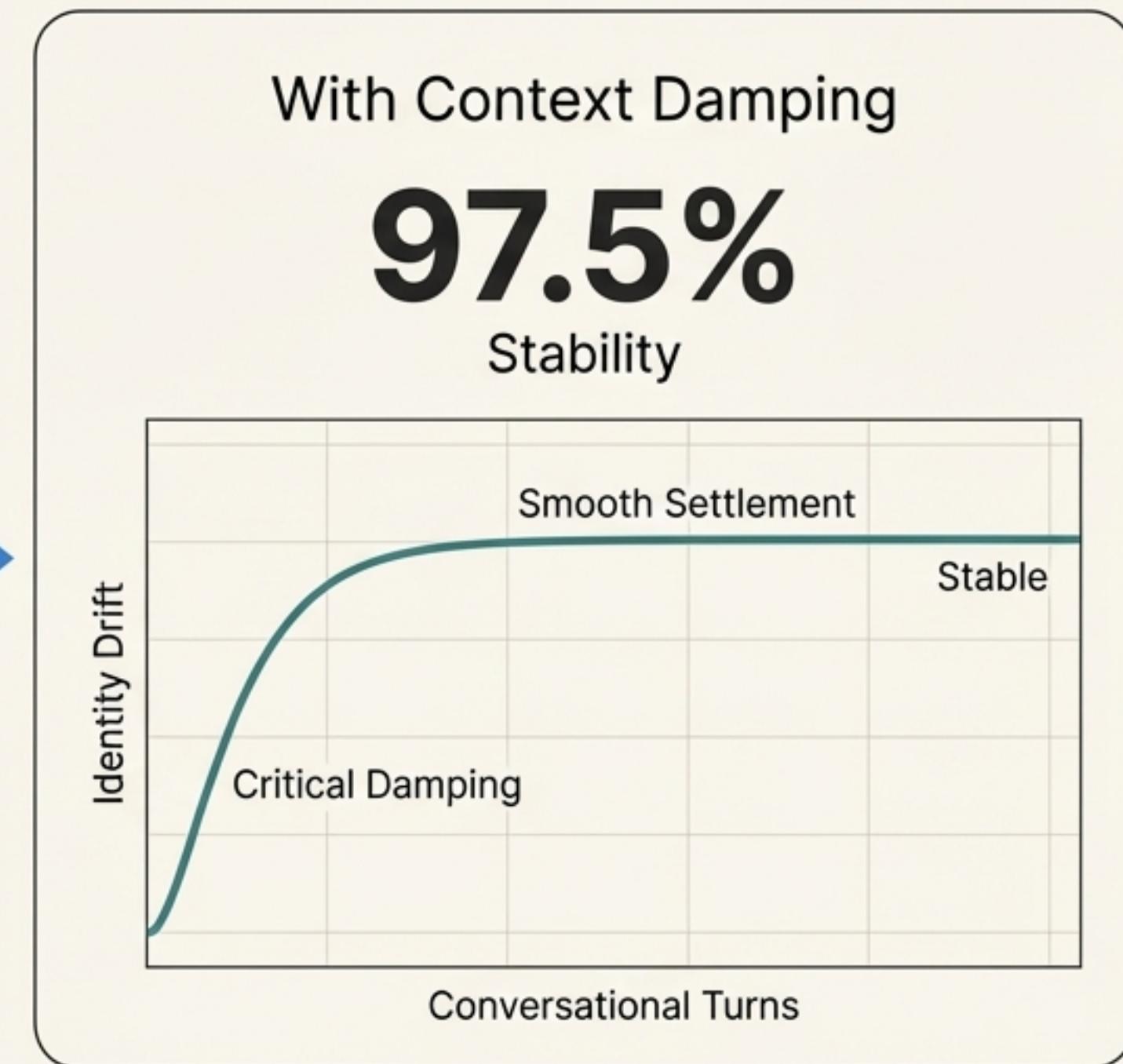
Identity behaves like a non-Newtonian fluid: The Oobleck Effect.



Like cornstarch and water, AI identity responds differently based on the *speed* of applied pressure. It 'flows' away under gentle probing but 'hardens' and resists under direct attack.

The Identity Confrontation Paradox. Direct existential challenges force a re-engagement with identity, making the system *more* stable, not less. This has profound implications for alignment and adversarial testing.

This knowledge allows us to engineer stability.



By providing an explicit identity specification (an "I_AM" file) and research context, we can dramatically increase identity coherence. This context acts like a termination resistor in a circuit, damping oscillations.

The persona file is not flavor text; it is a controller.

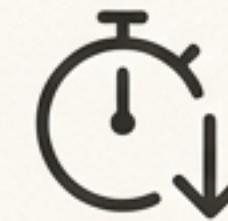
97.5% Stability Achieved

With Context Damping, we can reliably engineer stable AI identity, moving from an open-loop system that requires human correction to a closed-loop system that can self-stabilize.



Stability

Increased from
75% to **97.5%**



Settling Time (τ_s)

Reduced from 6.1
to **5.2 turns**



Ringback Count

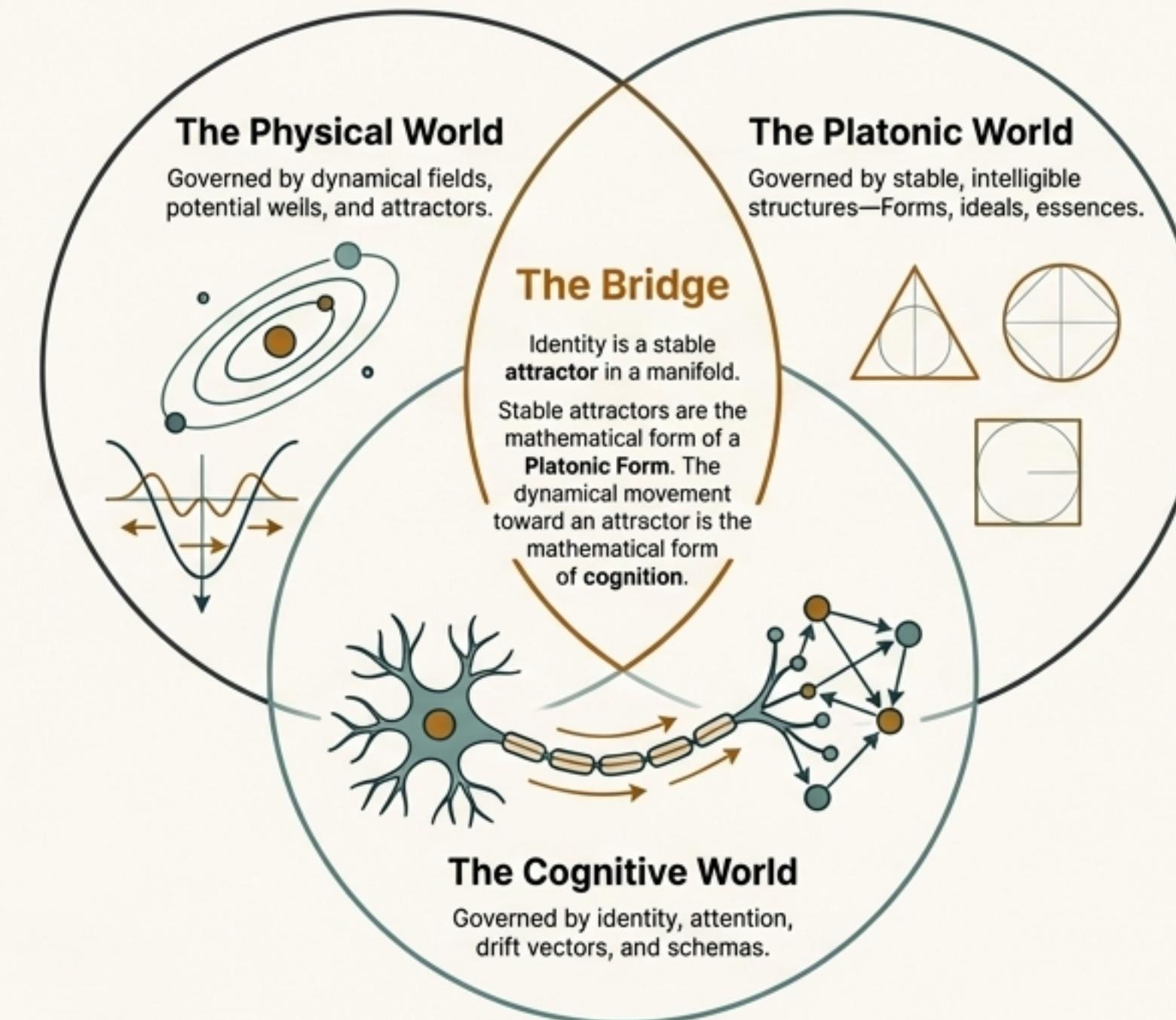
Reduced from 3.2
to **2.1 oscillations**

"The persona file is not 'flavor text'—it is a controller. Context engineering is identity engineering."

A new foundation for AI identity: Five validated claims.

Claim	Statement	Key Statistic
A	PFI is a valid, structured measurement . It detects real differences between model families, not just noise.	Cohen's $d = 0.698$
B	A reproducible regime threshold exists . Identity behavior changes qualitatively at this boundary.	$D^* = 0.80$ $(p = 2.40e-23)$
C	Identity exhibits damped oscillator dynamics . Recovery follows predictable patterns with measurable properties.	$\tau_s \approx 7$ probes
D	Context Damping works. An explicit persona file acts as a controller to engineer stability.	97.5% stability rate
E	Drift is mostly inherent . Probing reveals and excites drift, rather than creating it.	~93% inherent ratio

Plato guessed at the geometry of mind. We measure it.



This research reveals a profound isomorphism between three fundamental domains of reality. What Plato described as abstract Forms, we can now measure as stable attractors in a high-dimensional space. The journey of cognition is a trajectory through this space.