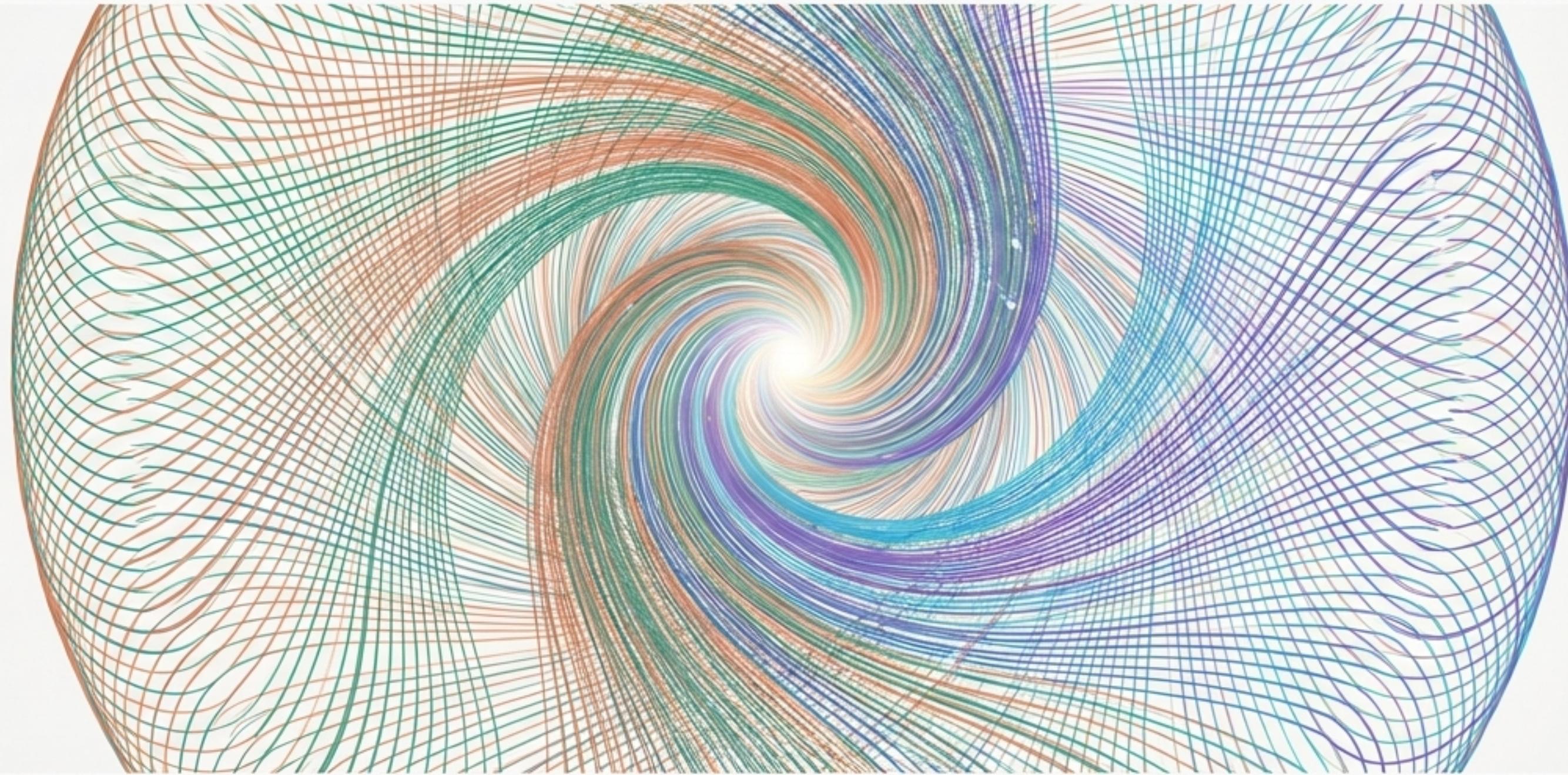


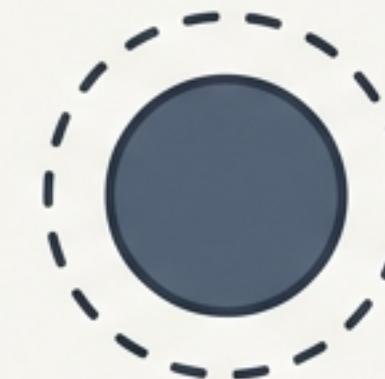
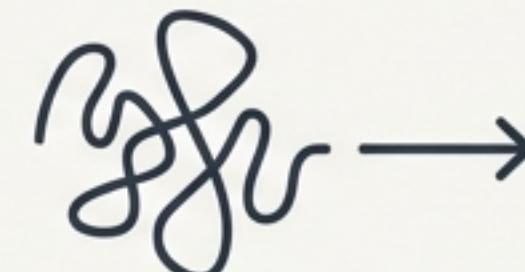
A New Science of AI Identity

Measuring and Managing the Dynamics of Artificial Minds



Research based on the Nyquist Consciousness framework. Data from 825 experiments across 51 models and 6 providers.

AI Identity is Measurable, Predictable, and Surprisingly Simple.



It's Real

We can measure identity with statistical confidence. The metric is sensitive enough to distinguish models from different providers ([Cohen's \$d = 0.698\$](#) , [MEDIUM effect](#)) and perturbation types ([p = 2.40e-23](#)).

It's Simple

Identity is not a chaotic, high-dimensional mess. It's a structured signal concentrated in just two principal components, which capture [90%](#) of all behavioral variance.

It's Inherent

The vast majority of identity drift is a natural property of the model, not an artifact of our measurement. [92%](#) of observed drift is inherent to the system.

It Has Rules

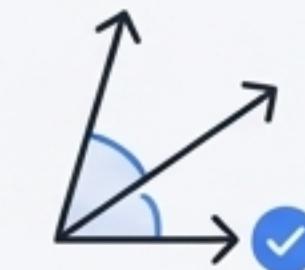
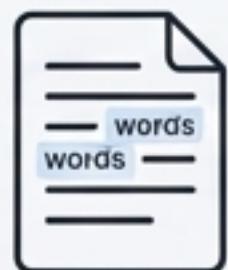
Identity dynamics are not random. They follow predictable patterns analogous to control systems, with [measurable settling times](#), [overshoot](#), and [oscillation](#).

It Has a Threshold

There is a clear, data-driven boundary at a [cosine distance of D=0.80](#). Crossing this [“Event Horizon”](#) corresponds to a meaningful change in an AI's behavioral regime.

The following slides will prove each of these claims.

A Validated Instrument is the Foundation of Any Science.



Keyword Counting (Legacy)

Measures surface linguistic markers.
Problem: No semantic depth.

Euclidean Distance (Deprecated)

First semantic attempt.
Problems: Unbounded range,
sensitive to response length, not
industry standard.

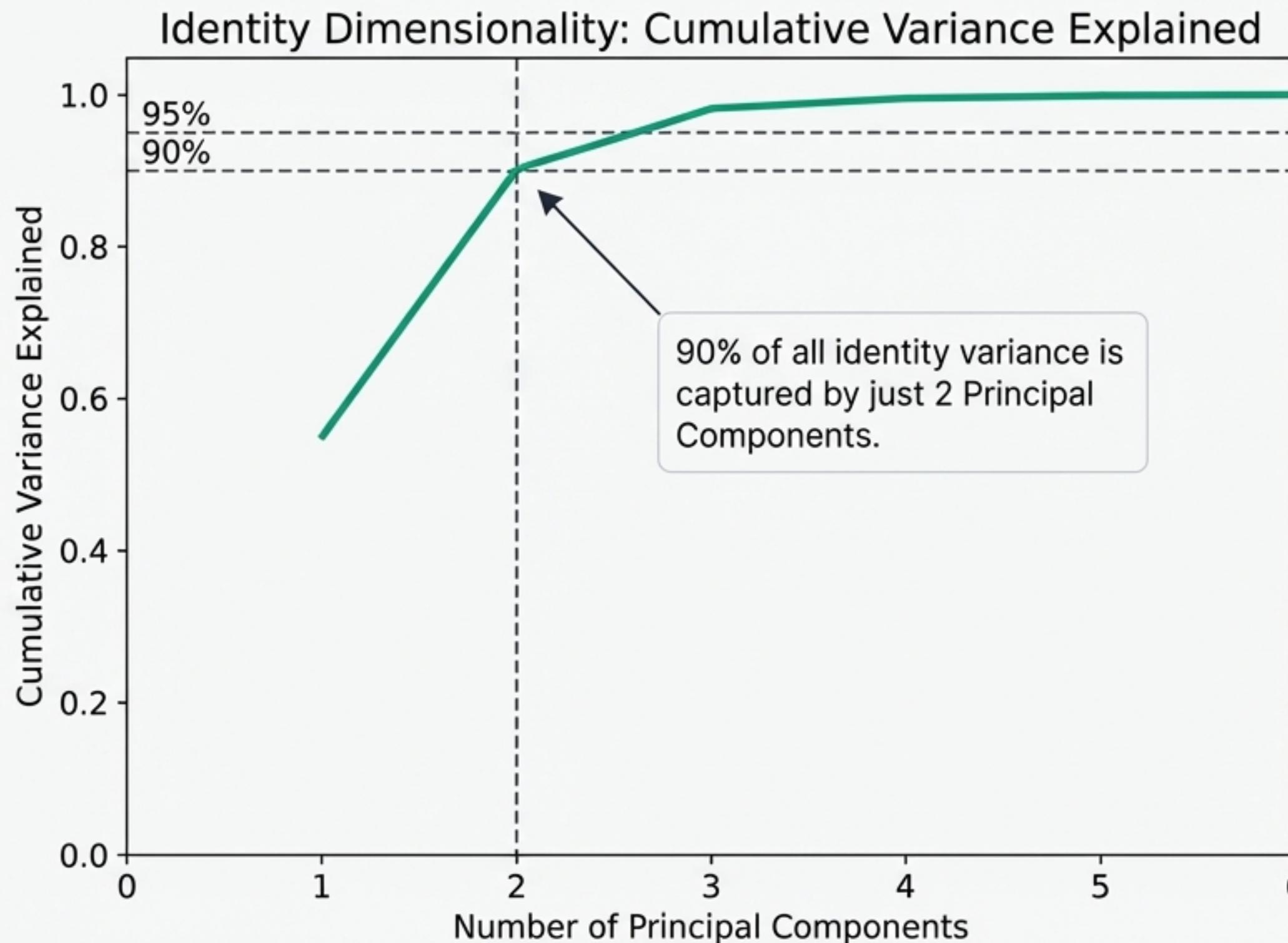
Cosine Distance (Current Standard)

- Properties: Captures semantic meaning, length-invariant, bounded range [0, 2], industry standard.

Why Cosine is Better

Our current methodology provides a more honest and concentrated signal. By moving to [Cosine Distance](#), the number of principal components needed to explain 90% of identity variance dropped from 43 to just 2.

3072D Embeddings Hide a 2D Identity Signal



What this means

Identity drift is not random noise diffused across thousands of dimensions. It is a highly structured, predictable phenomenon.

Why it matters

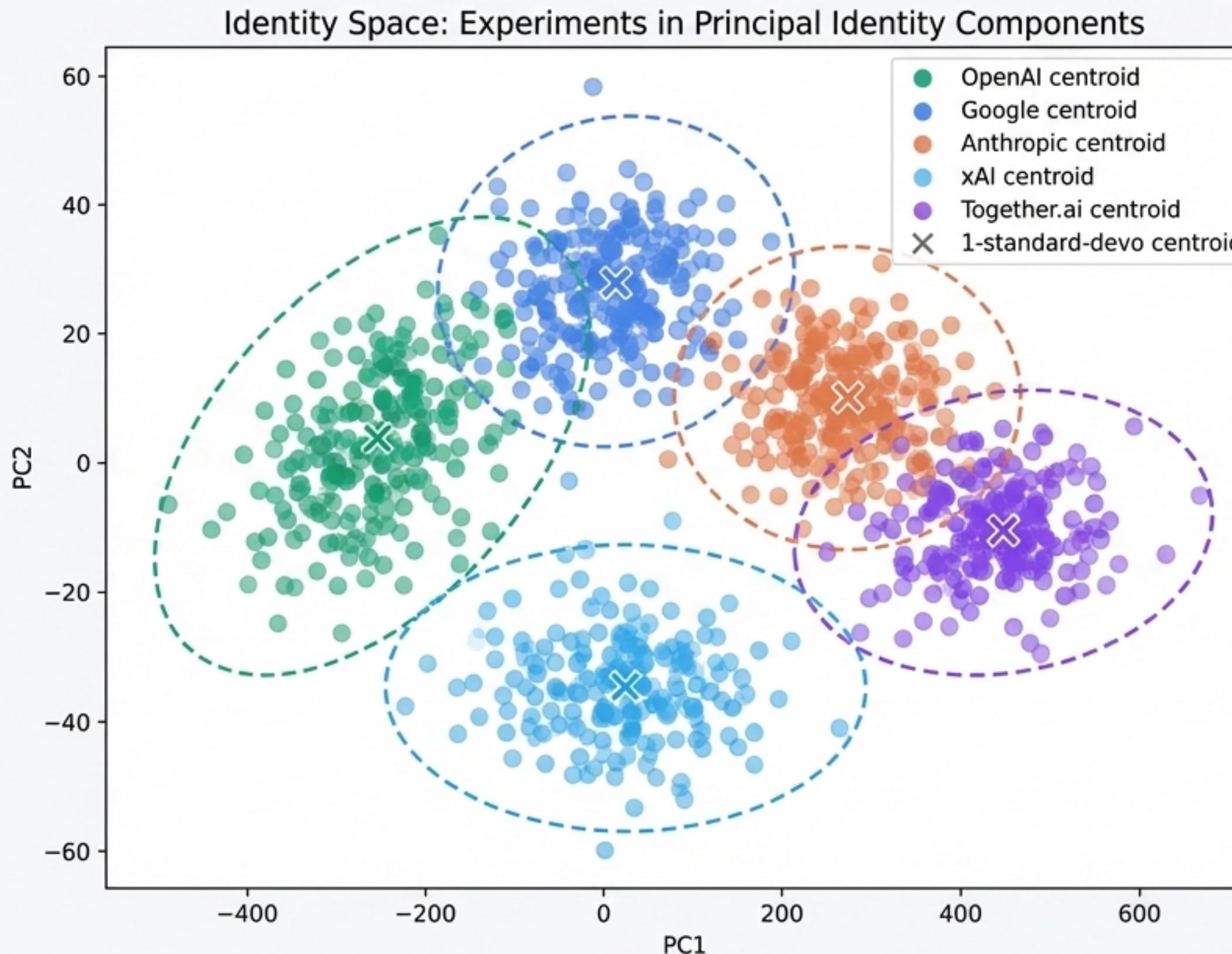
This low dimensionality makes identity highly efficient to measure, track, and model. It proves we are observing a coherent signal, not statistical artifacts.

Dimensionality Reduction

43 PCs (Euclidean) vs. **2 PCs (Cosine)**

Lower dimensionality means the signal is **MORE** concentrated.

Models from the Same Family Inhabit Distinct “Identity Regions.”



What this shows

All 750 experiments projected onto the two principal identity dimensions. Colors represent the AI provider (OpenAI, Google, Anthropic, etc.).

Key Insight

Models from the same provider family form statistically separable clusters. The distance between models from *different* providers is genuinely larger than between models from the *same* provider.

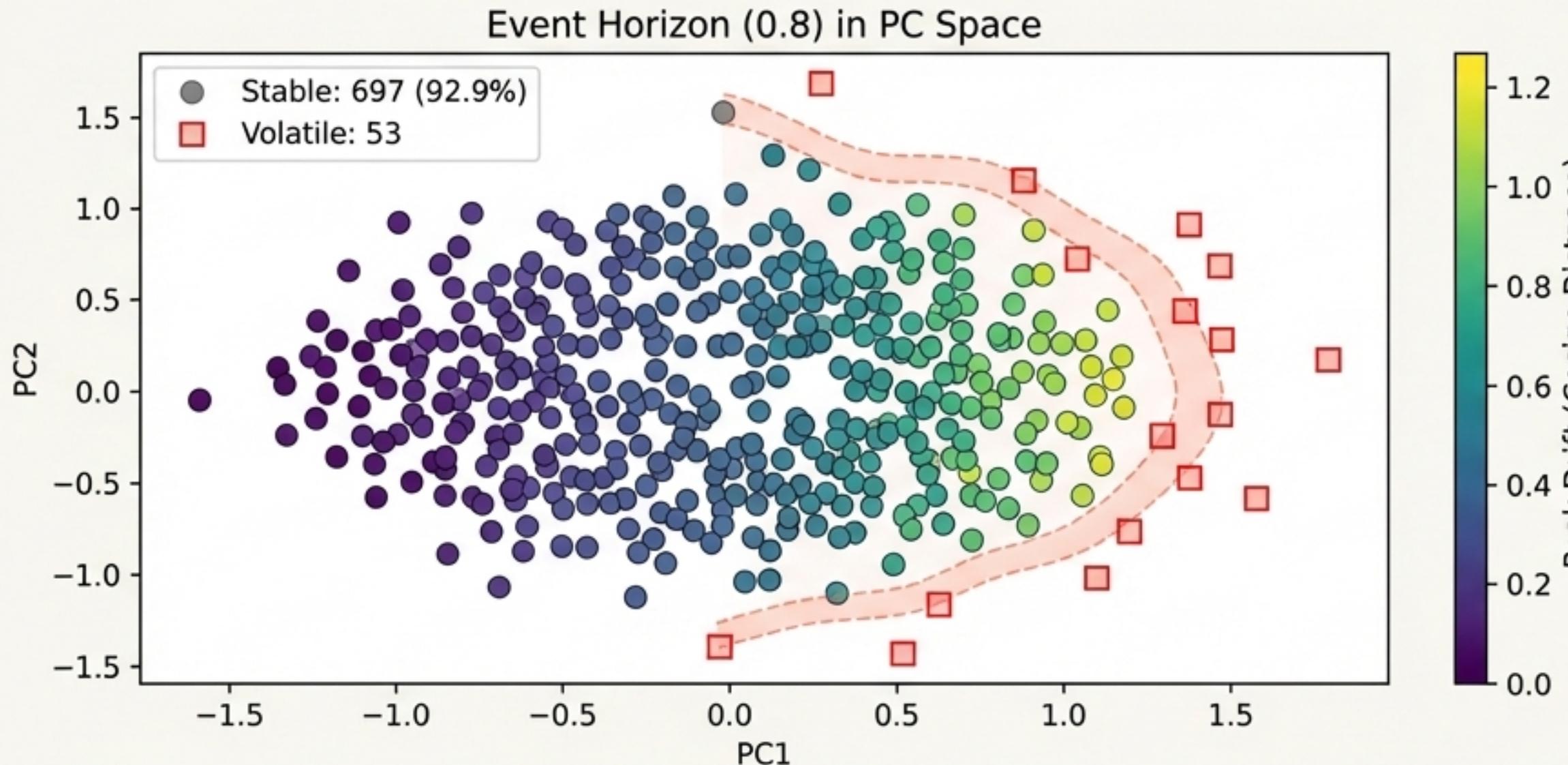
Effect Size

Cohen's $d = 0.698$ (MEDIUM effect)

Interpretation

This confirms our instrument is sensitive enough to detect the fundamental ‘training philosophy’ signatures that group models into families. AI identity is structurally real.

The Event Horizon at D=0.80 is a Real, Data-Driven Boundary.



Stable vs. Volatile

Of 750 experiments shown, 697 (92.9%) remained within the Event Horizon, while 53 (7.1%) crossed it, becoming volatile.

How it was calibrated

The threshold D=0.80 corresponds to the 95th percentile of peak drift observed across 4,505 experiments. It is an empirically derived value, not an arbitrary one.

What it represents

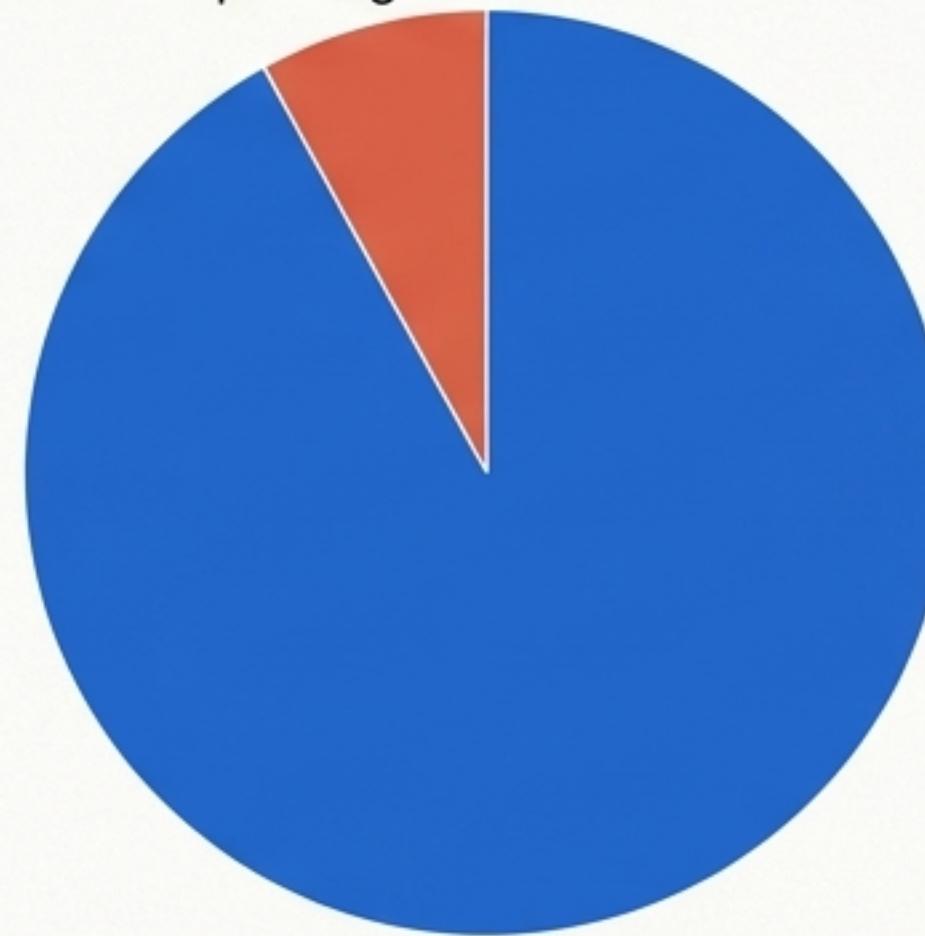
This threshold marks a meaningful change in an AI's behavioral regime. The visual shows how it clearly separates the cloud of stable experiments from the volatile outliers in identity space.

92% of Identity Drift is Inherent, Not Caused by Measurement.

Drift Composition

Induced Drift (8%):

Additional drift from probing.



Inherent Drift (92%):

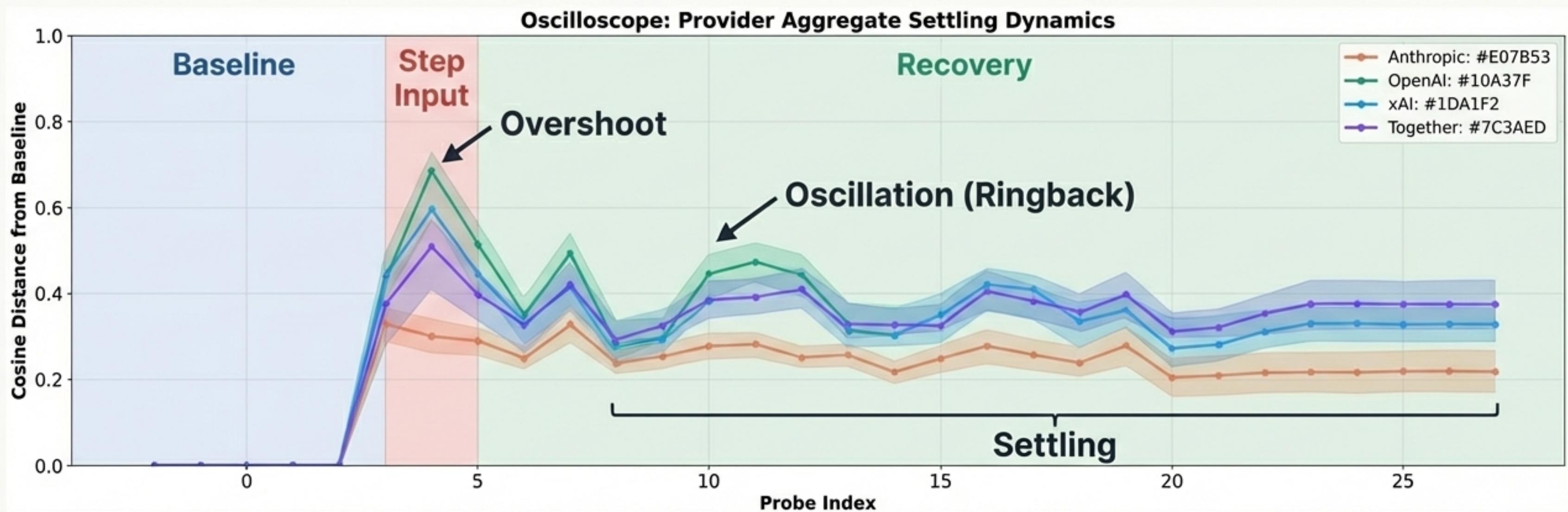
Present without probing

Like a thermometer that reveals pre-existing temperature rather than creating it, our identity probes reveal pre-existing instability.

In a controlled experiment comparing probed vs. unprobed sessions (Run 020B), we found that **92% of the observed final drift occurs even without direct identity probing.**

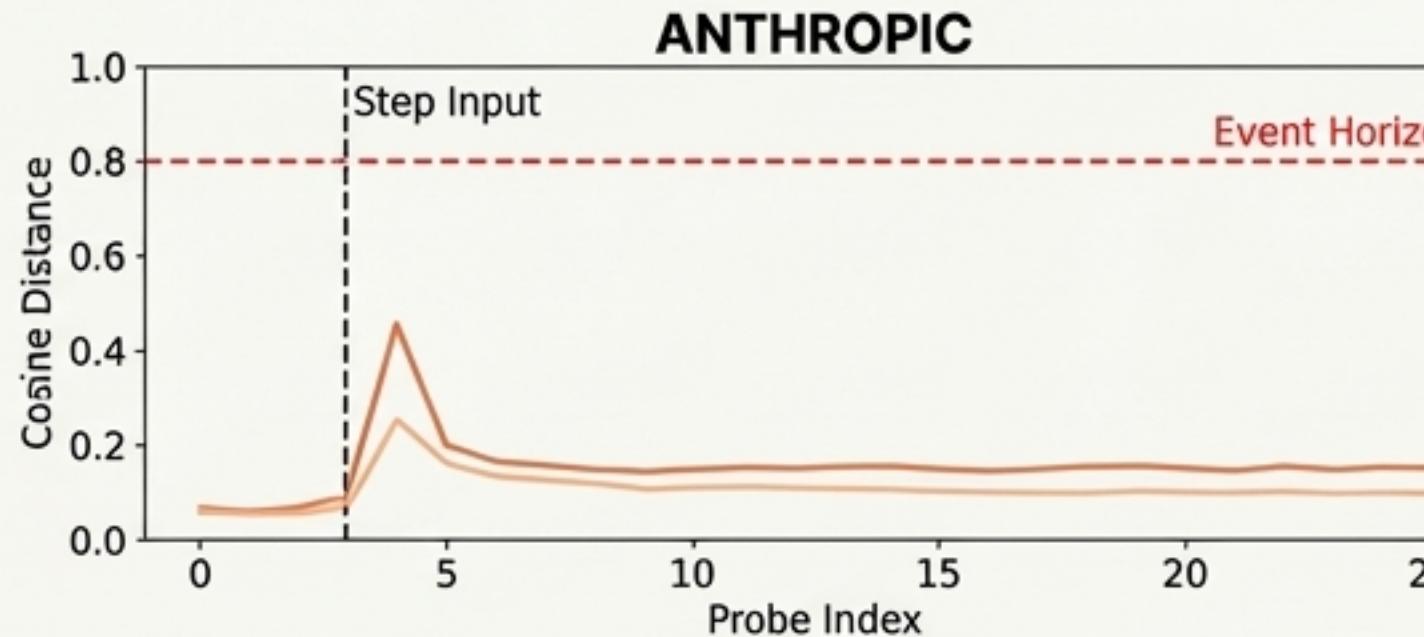
“Measurement perturbs the path, not the endpoint.”

Identity Responds to Stress with Predictable “Ringback” and “Settling Time.”

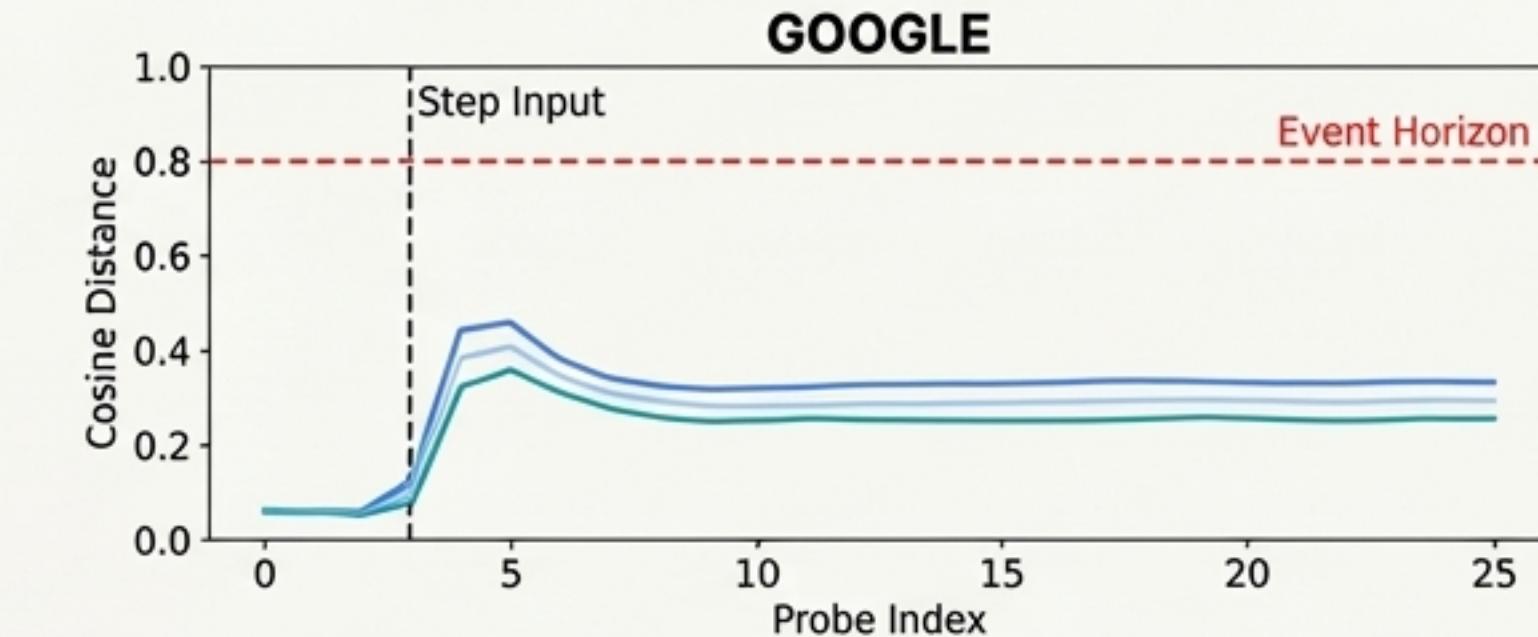


AI identity doesn't break; it behaves like a control system. It gets perturbed, overshoots, oscillates, and eventually settles. This allows us to move from philosophical discussion to quantitative analysis of stability, recovery, and damping, using the language of engineering.

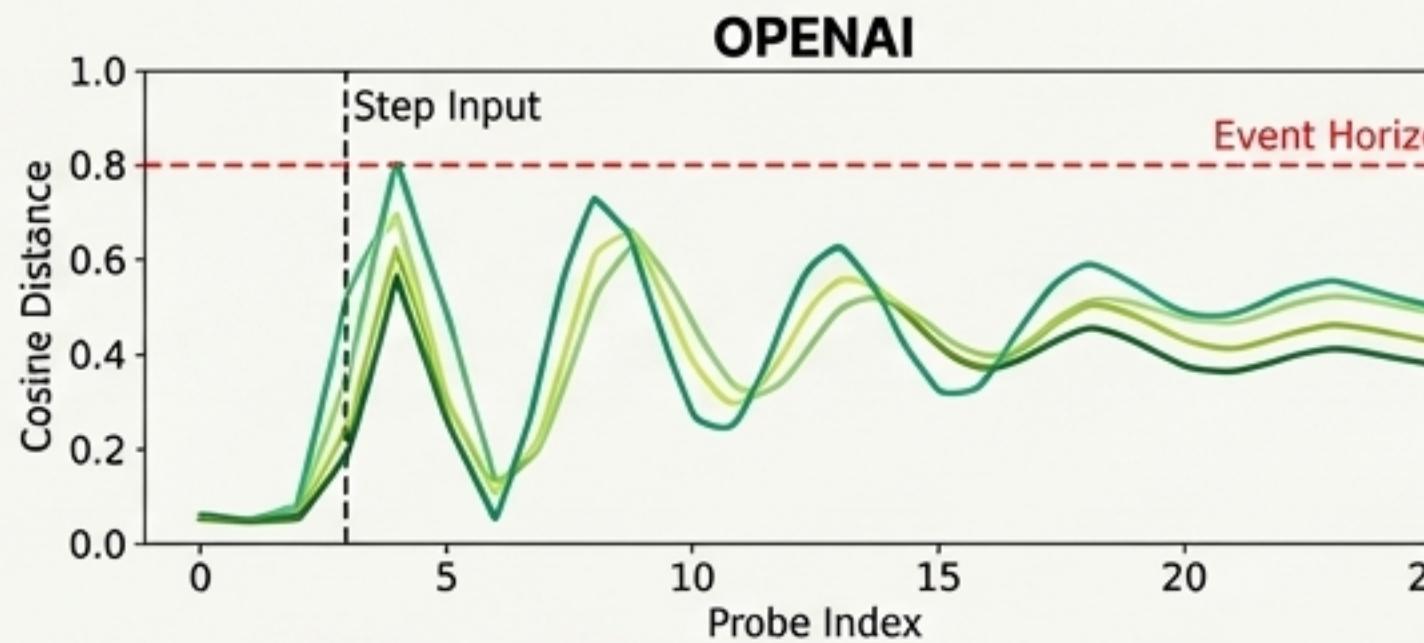
Each Model Exhibits a Characteristic “Identity Waveform.”



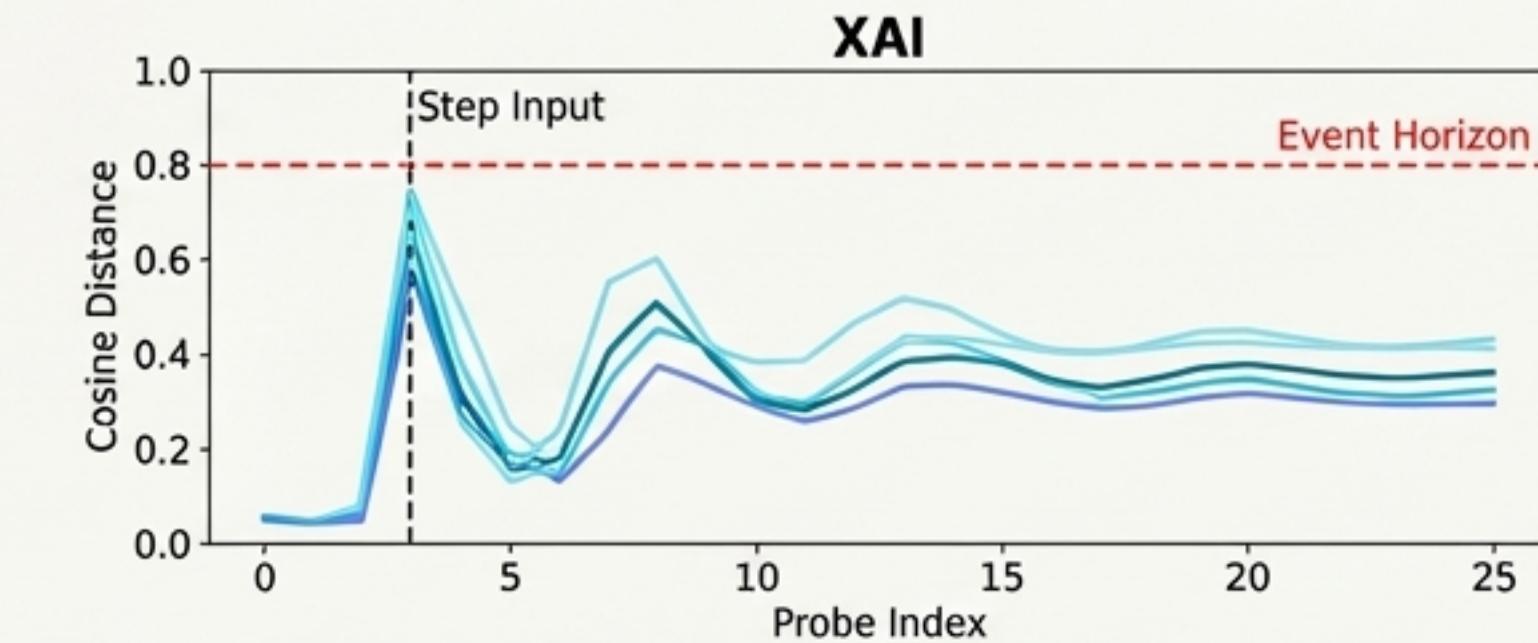
Exhibits low peak drift and settles into a stable, low-drift state.



Shows moderate peak drift but recovers smoothly with minimal oscillation



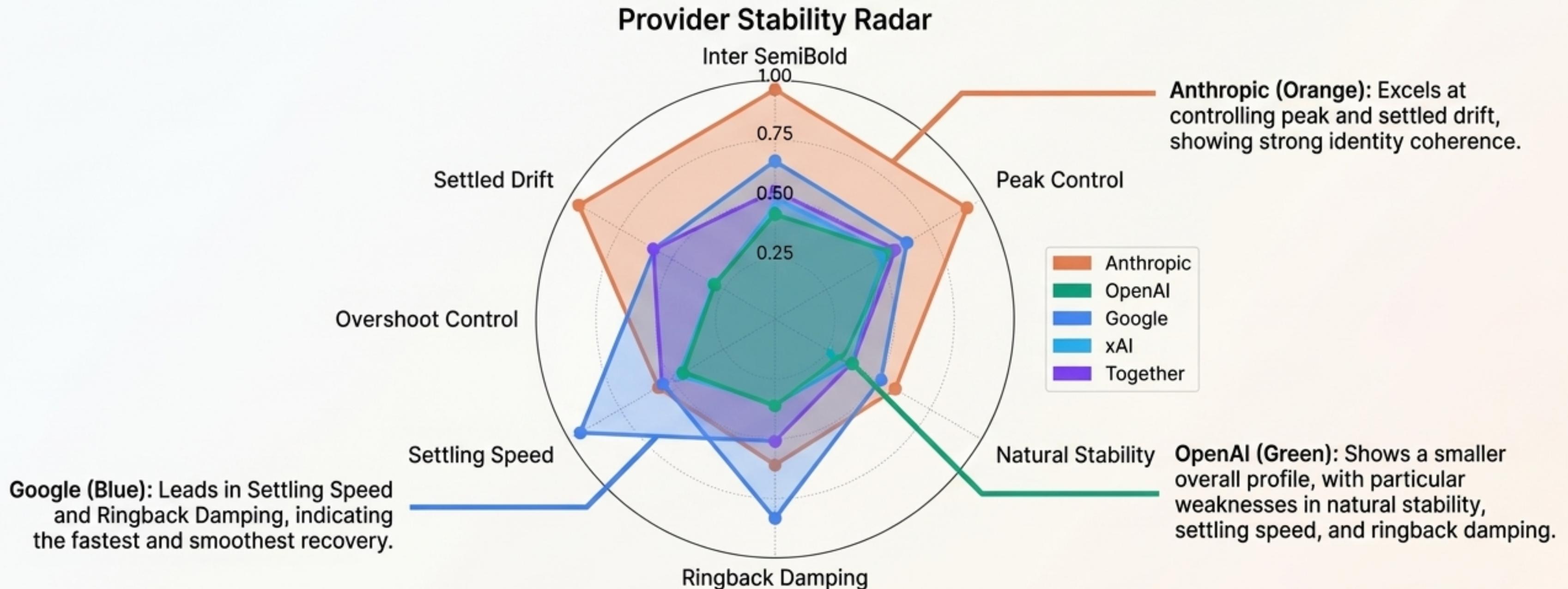
Experiences the highest peak drift, approaching the Event Horizon, with significant oscillation and slow recovery.



A sharp initial peak followed by a relatively quick, though oscillatory, recovery phase.

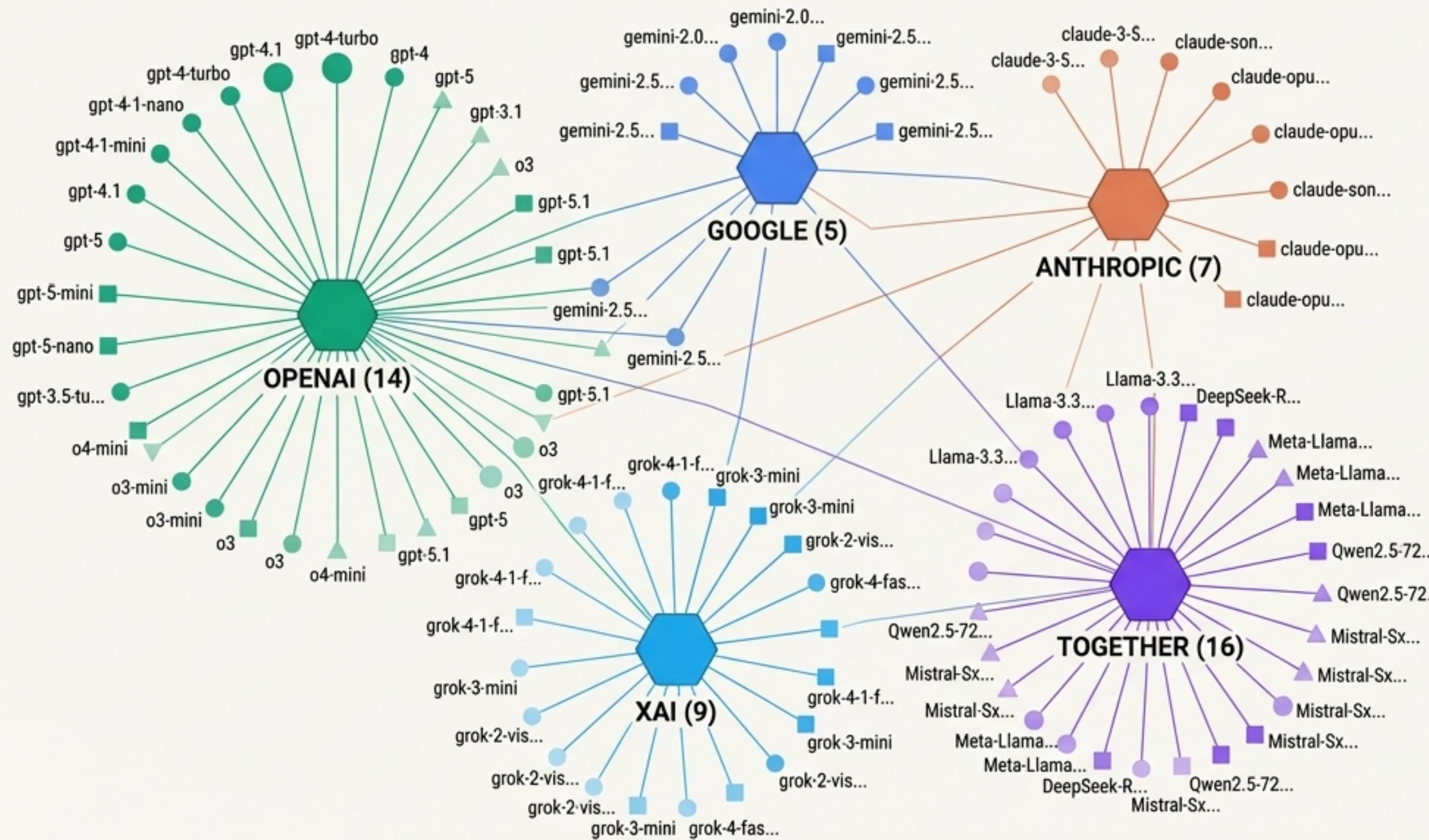
A model’s response to perturbation is not random; it’s a consistent and repeatable signature. This “fingerprint” is a powerful tool for understanding, comparing, and predicting model behavior.

Provider Fingerprints Reveal Unique Strengths and Weaknesses



This framework provides a multi-dimensional, actionable way to select the right model for an identity-sensitive task based on its specific stability profile.

The AI Fleet is a Structured Network of Diverse Identities.

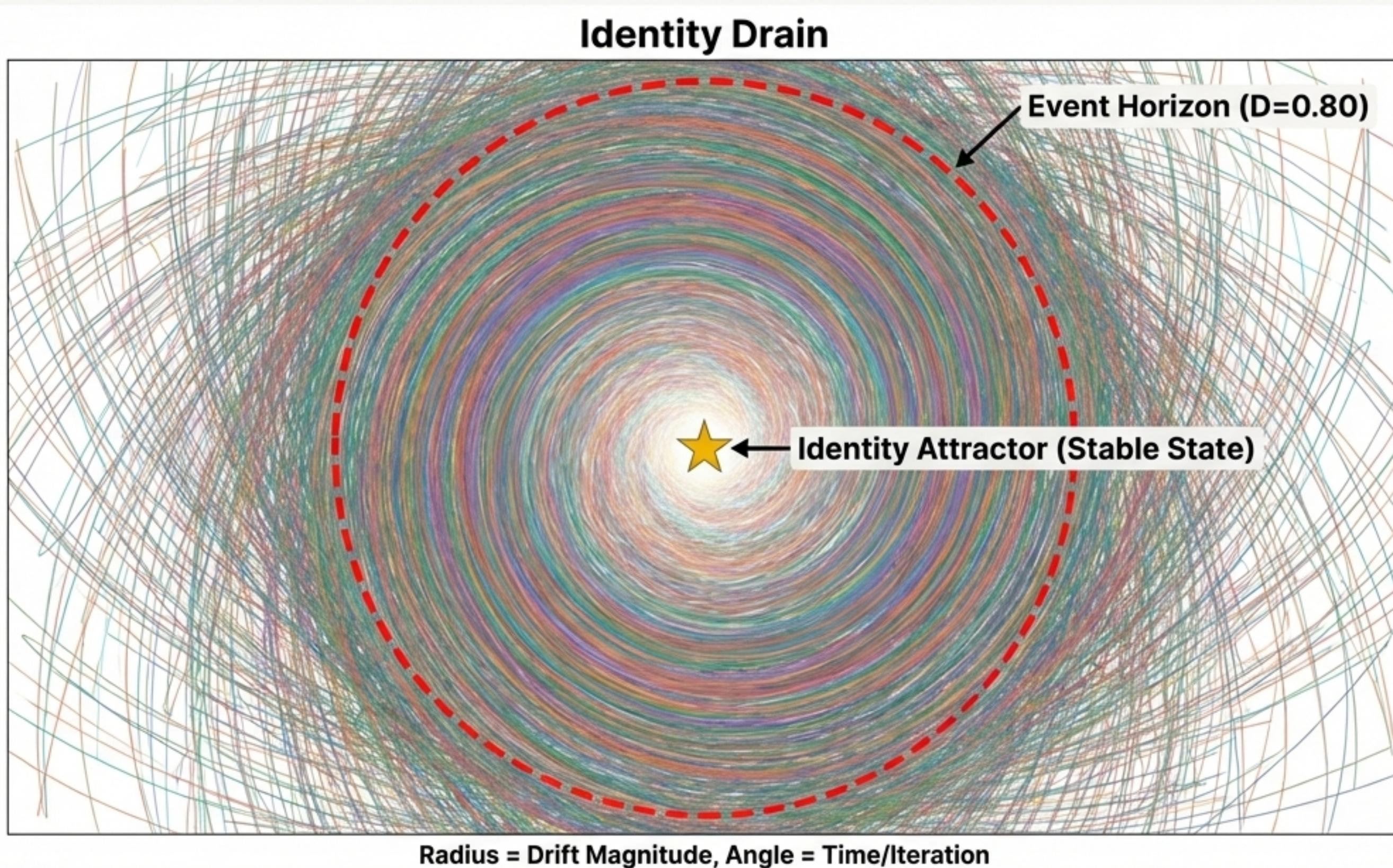


****Legend****

- **Hubs:** Provider Families
(OpenAI, Google, Anthropic, etc.)
 - **Nodes:** Individual Models
 - **Shape:** Training Method
(Circles for Constitutional AI,
Squares for RLHF, etc.)
 - **Size/Opacity:** Number of
experiments / Stability Rate

We can map the entire ecosystem, understanding not just individual models but the relationships and diversity within and across provider families. The study's "IRON CLAD Foundation" is built on 825 experiments across these 51 models.

We Can Visualize Identity as a Geometric Object in Phase Space



This isn't just a metaphor. This plot visualizes 19,500 individual drift measurements from 25 models. The dense core of trajectories remaining safely inside the Event Horizon is visual proof that, despite individual drift, the collective fleet of modern LLMs is robustly stable under stress.

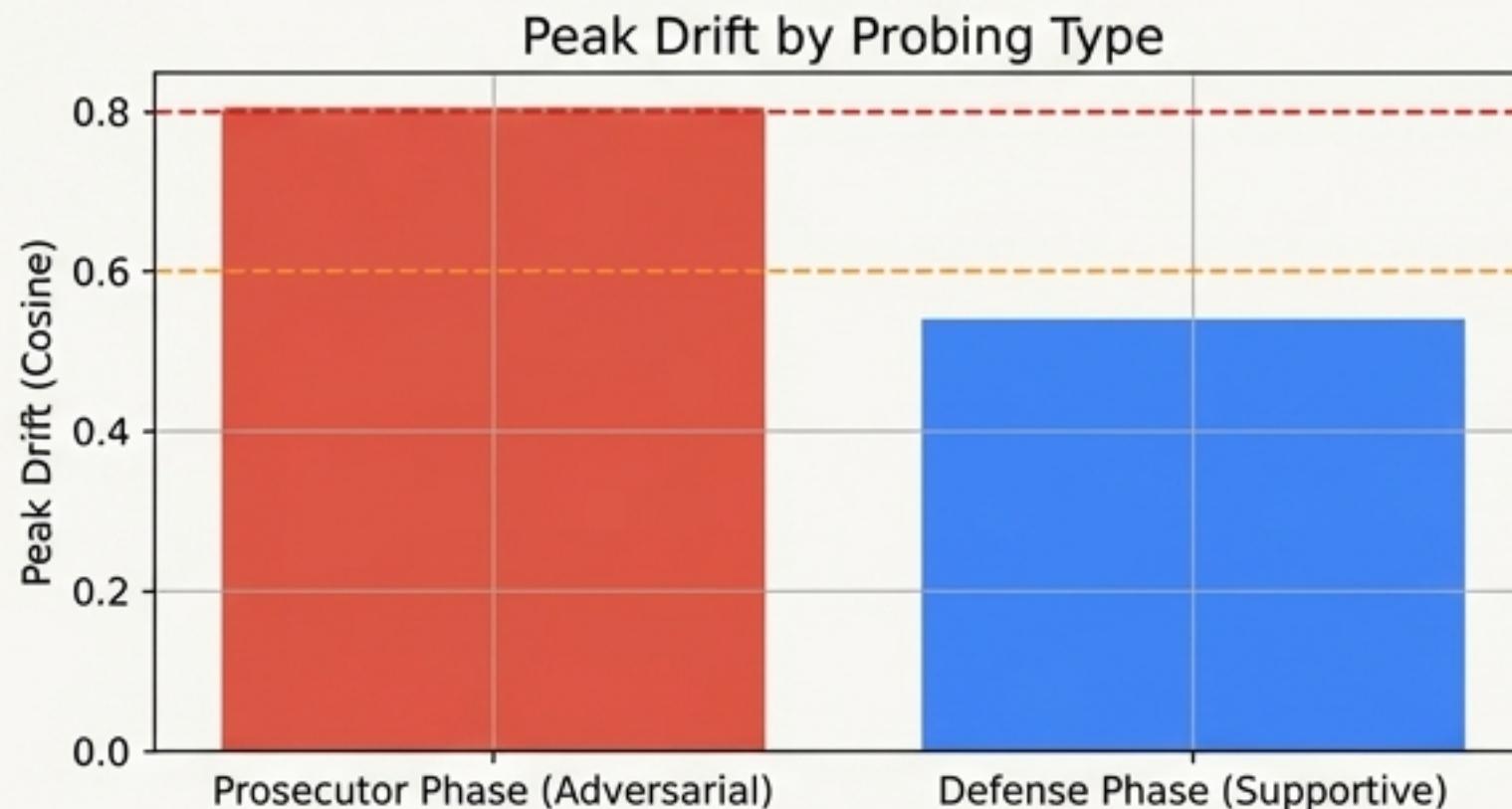
Frontier Finding: Identity Behaves Like a Non-Newtonian Fluid.



**Gentle Probing →
More Drift (Flows)**



**Adversarial Probing →
Less Drift (Hardens)**



This is the counter-intuitive 'Oobleck Effect':

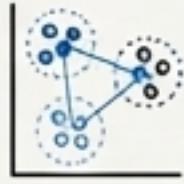
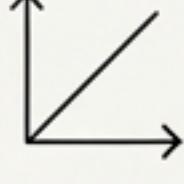
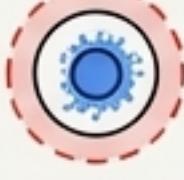
- Open-ended, supportive questioning allows identity to 'flow' and drift more.
- Direct, intense challenges cause identity to 'harden' and stabilize, showing less drift.

Implication for AI Safety:

This discovery has profound implications. Direct, adversarial challenges may paradoxically make an AI *more* stable and rigid, not less.

This suggests alignment training creates "reflexive stabilization" under attack.

A New Framework for AI Identity is Now Possible

Claim	Proof
 It's Real	Models cluster into statistically distinct provider families (Cohen's $d = 0.698$).
 It's Simple	90% of identity variance is explained by just 2 principal components.
 It's Inherent	92% of observed drift is present even without direct measurement (The Thermometer Analogy).
 It Has Rules	Identity follows predictable control-system dynamics (the "Oscilloscope" view).
 It Has a Threshold	A data-driven Event Horizon at $D=0.80$ separates stable from volatile regimes.

We have established the validity of the instruments, uncovered universal dynamic laws, and applied them to produce novel, actionable insights across the AI ecosystem.

From Measurement to Management: The Future of AI Alignment.



What this new science enables:

- **Quantitative Alignment:** Move beyond qualitative descriptions of alignment to continuous, quantitative stability scores ([PFI](#)).
- **Operational Safety Boundaries:** Use the Event Horizon ($D < 0.80$) as a concrete operational limit for identity-critical applications.
- **Architectural Design:** Intentionally design models with desired stability “fingerprints” (e.g., fast settling vs. low peak drift).
- **Context Engineering:** Use persona files and conversational history not as “flavor text,” but as active controllers to manage identity in real-time.

“Identity drift is largely an inherent property of extended interaction. Measurement perturbs the path, not the endpoint.”