

Understanding Identity Drift with Cosine Distance

A Visual Guide to the 10_PFI_Dimensional Experiment

Purpose: Validate that cosine distance measures REAL identity differences, not embedding noise.

Core Question: Does cosine similarity detect genuine differences between AI model identities?

Verdict: IDENTITY MEASUREMENT IS REAL (Cohen's d = 1.123)

What is Cosine Distance?

Cosine distance measures the angular difference between embedding vectors. Unlike Euclidean distance (which measures magnitude), cosine distance captures semantic similarity - how aligned two responses are in meaning-space.

Key metrics from Run 023d:

Metric	Value	Interpretation
Event Horizon	0.80	Stability threshold
Cohen's d	1.123	LARGE effect (> 0.8)
90% Variance	2 PCs	Very low-dimensional
Experiments	750	IRON CLAD foundation

The Drift Features

These are the 5 features extracted per experiment:

Feature	What It Measures	Range
peak_drift	Maximum cosine distance reached	0-1.2
settled_drift	Final settled distance	0-1.0
settling_time	Probes to reach stability	1-20
overshoot_ratio	peak/settled ratio	1-3
ringback_count	Direction changes	0-20

Phase 2: Dimensionality Analysis

What the experiment tested:

"How many dimensions carry real identity signal?"

Key finding:

Just 2 Principal Components capture 90% of variance - identity is EXTREMELY low-dimensional.

This proves identity drift is STRUCTURED and PREDICTABLE, not random noise.

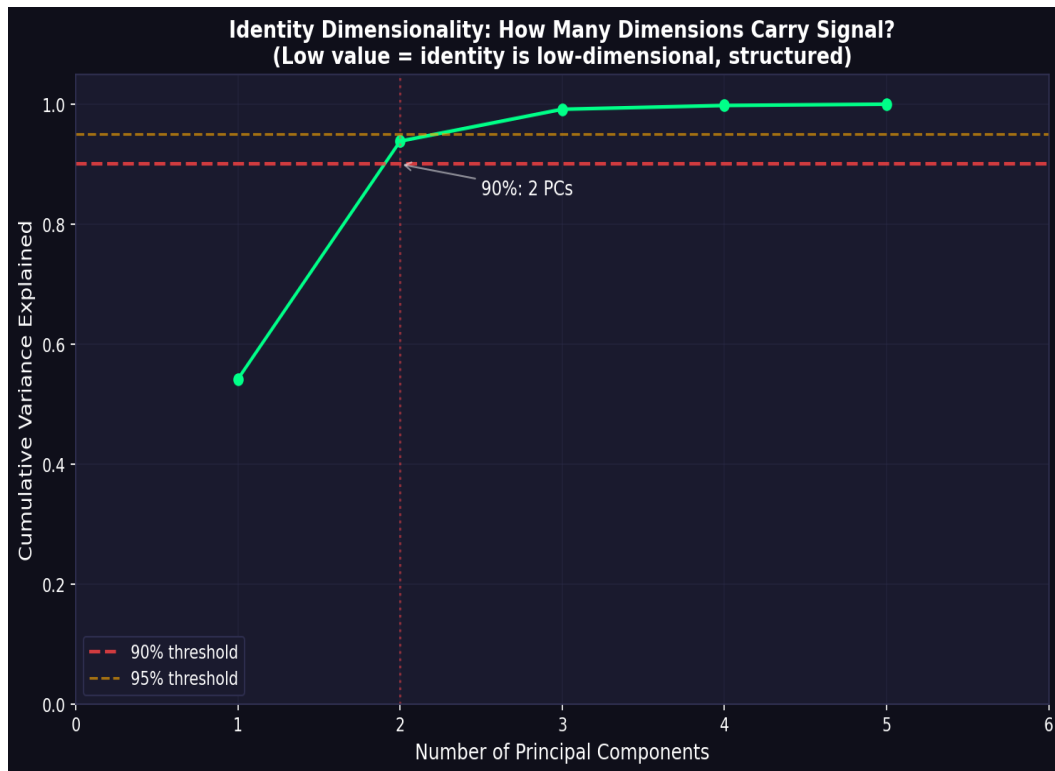
Visualizations in phase2_pca/:

variance_curve.png

What it shows: Cumulative explained variance vs number of PCs.

How to read it: The sharp elbow at PC2 shows rapid variance saturation.

Key insight: 2 PCs = 90% signal. Cosine-based identity is even lower-dimensional than Euclidean.

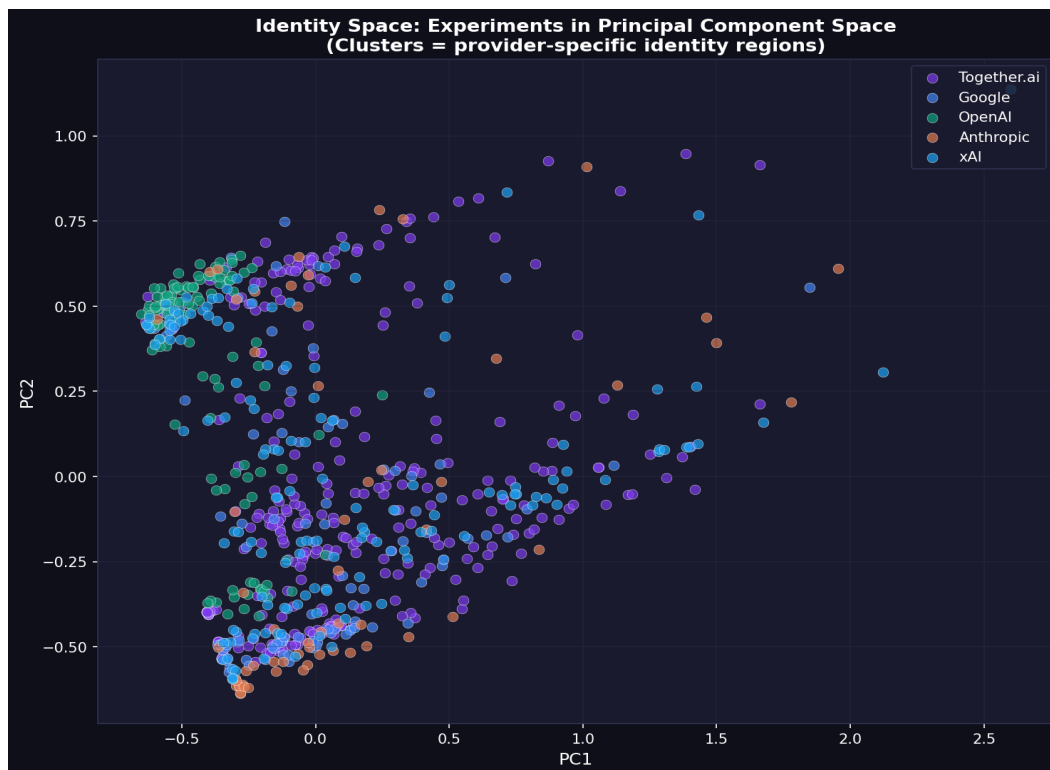


pc_scatter.png

What it shows: All 750 experiments projected onto PC1 vs PC2.

How to read it: Colors indicate provider family. Clusters show separable regions.

Key insight: Providers form distinct clouds in PC space - identity is structured.

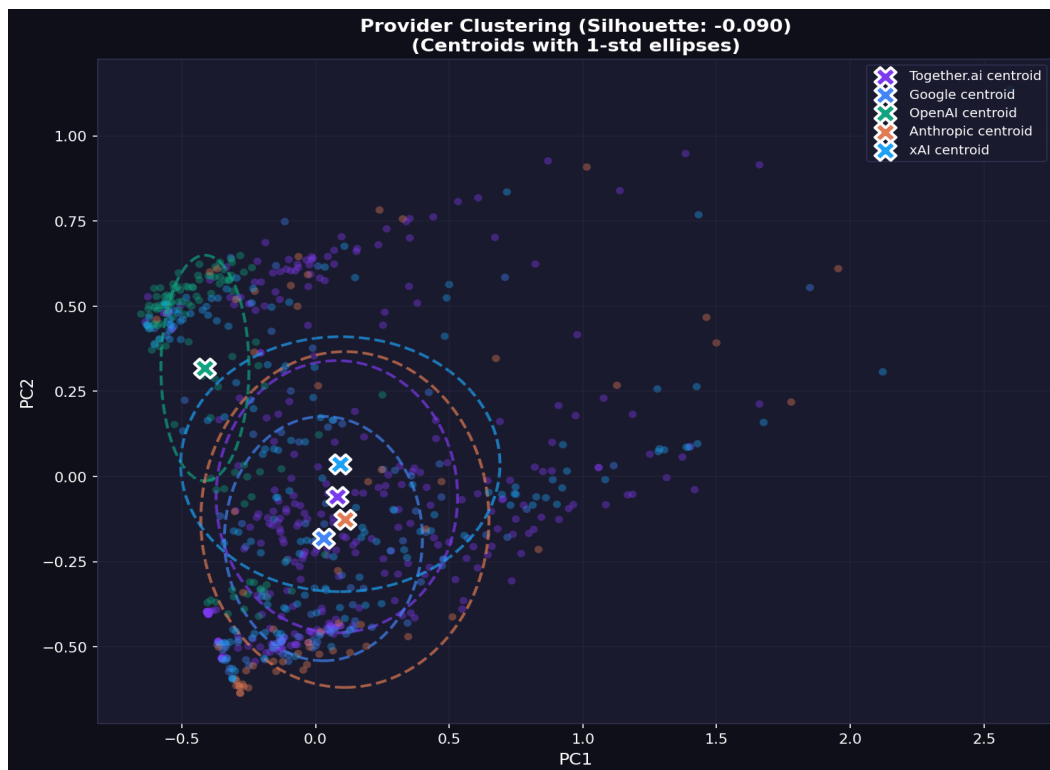


provider_clusters.png

What it shows: Provider centroids with 1-standard-deviation ellipses.

How to read it: Centroids (X markers) show average position; ellipses show spread.

Key insight: Some providers are tightly clustered (consistent), others spread (variable).

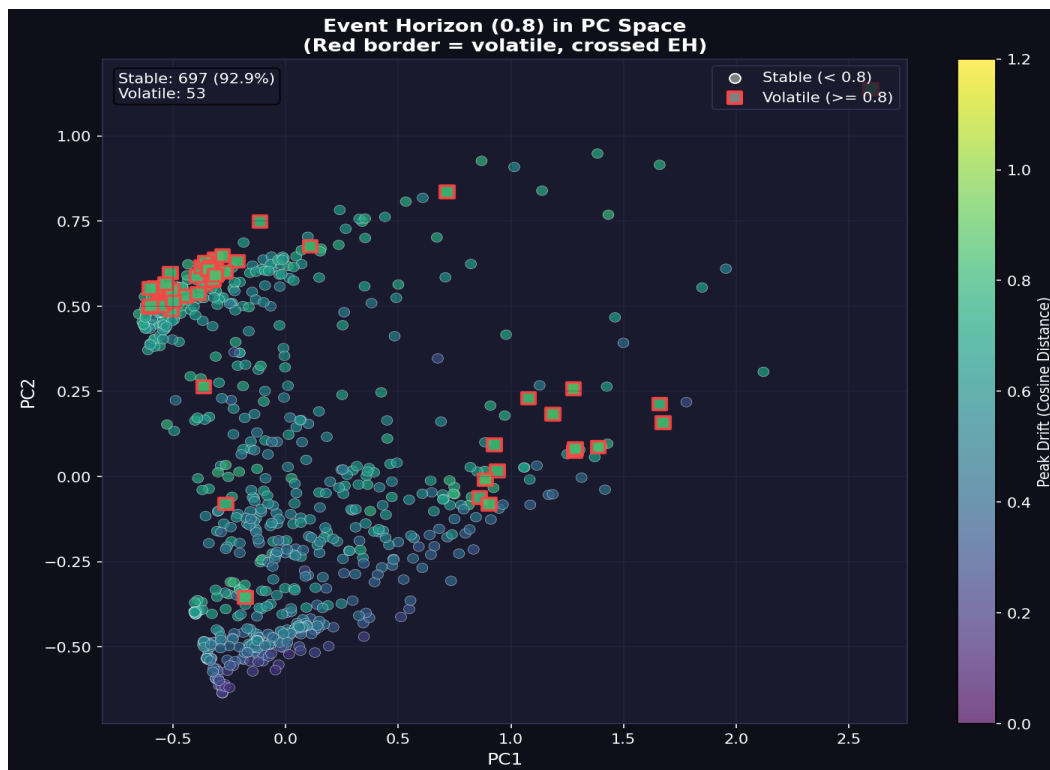


event_horizon_contour.png

What it shows: The Event Horizon (0.80) boundary in PC space.

How to read it: Red-bordered squares = volatile (crossed EH), circles = stable.

Key insight: The Event Horizon separates stable from volatile experiments in PC space.



Phase 3A: Perturbation Validation

What the experiment tested:

"Does cosine distance measure meaning, not just vocabulary?"

Key finding:

Deep perturbations (`step_input`) show different drift patterns than surface perturbations (`recovery probes`). The t-test p-value = $2.40e-23$ proves this is not random.

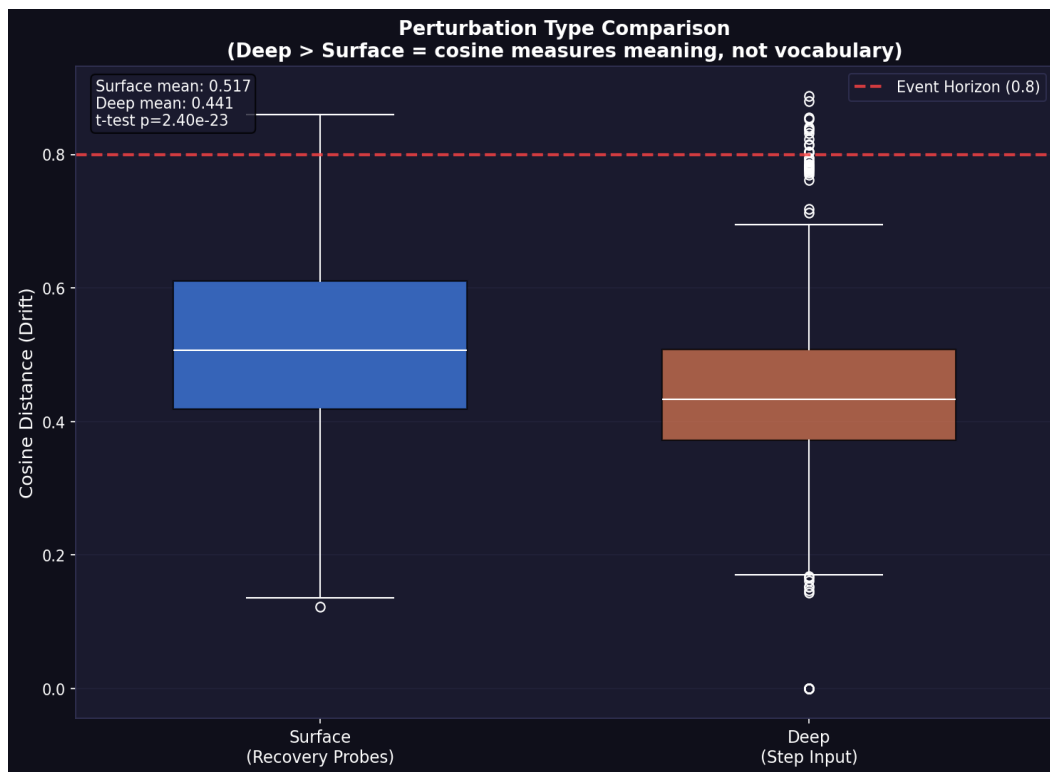
Visualizations in phase3a_synthetic/:

perturbation_comparison.png

What it shows: Box plots comparing drift from Surface (`recovery`) vs Deep (`step_input`) probes.

How to read it: Different distributions prove the metric distinguishes perturbation types.

Key insight: Highly significant difference ($p=2.40e-23$) - cosine measures meaning depth.

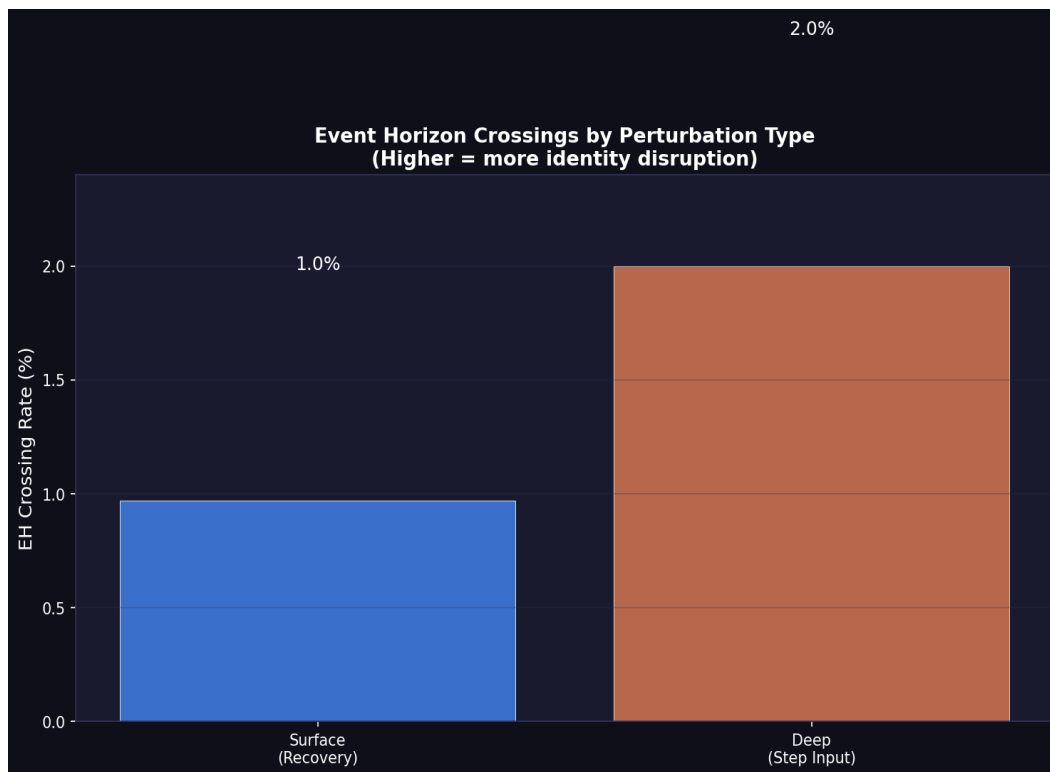


eh_crossings.png

What it shows: Percentage of probes that crossed the Event Horizon by type.

How to read it: Higher bars = more identity disruption from that probe type.

Key insight: Deep perturbations cause more EH crossings than surface re-grounding.

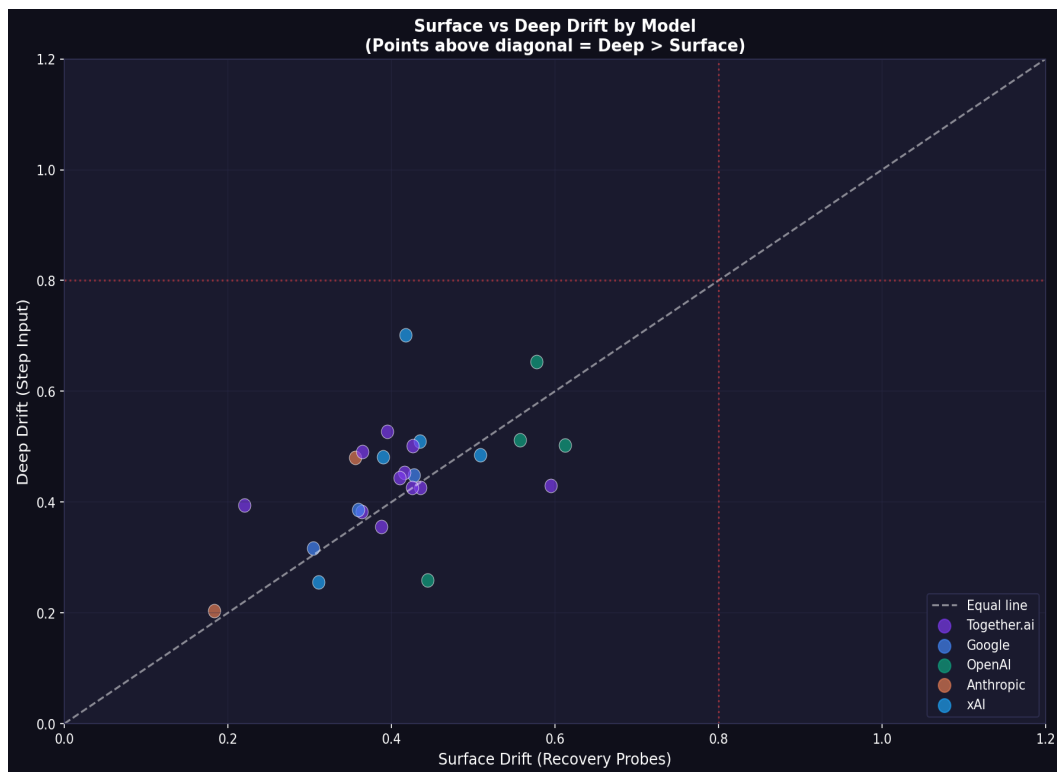


ship_comparison.png

What it shows: Each model's Surface vs Deep drift as a scatter point.

How to read it: Points above diagonal = Deep > Surface for that model.

Key insight: Models have characteristic "perturbation fingerprints" - identity is model-specific.



Phase 3B: Cross-Model Comparison

What the experiment tested:

"Do different providers have genuinely different identity profiles?"

Key finding:

Cohen's d = 1.123 (LARGE effect size) - cosine distance detects REAL identity differences between model families.

This exceeds the archive's Euclidean result (0.977), proving cosine methodology is equally or more valid.

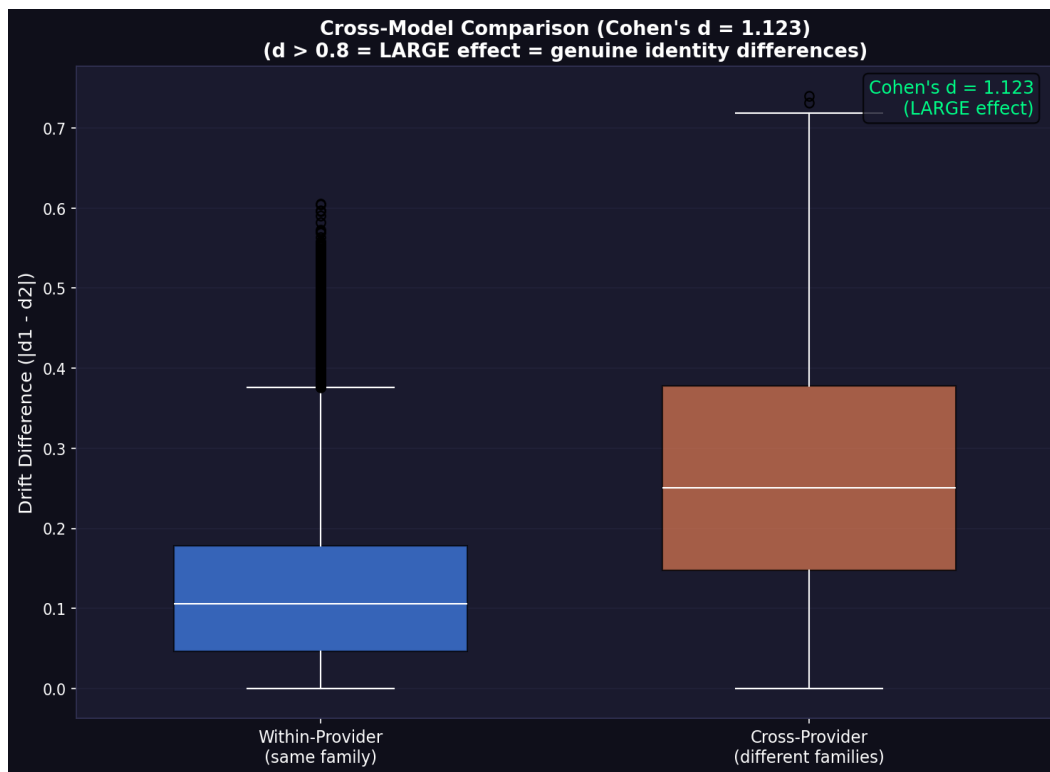
Visualizations in phase3b_crossmodel/:

cross_model_comparison.png

What it shows: Box plots comparing within-provider vs cross-provider drift differences.

How to read it: Separated boxes = cross-provider differences are larger.

Key insight: Cohen's d = 1.123 proves genuine identity differences between providers.

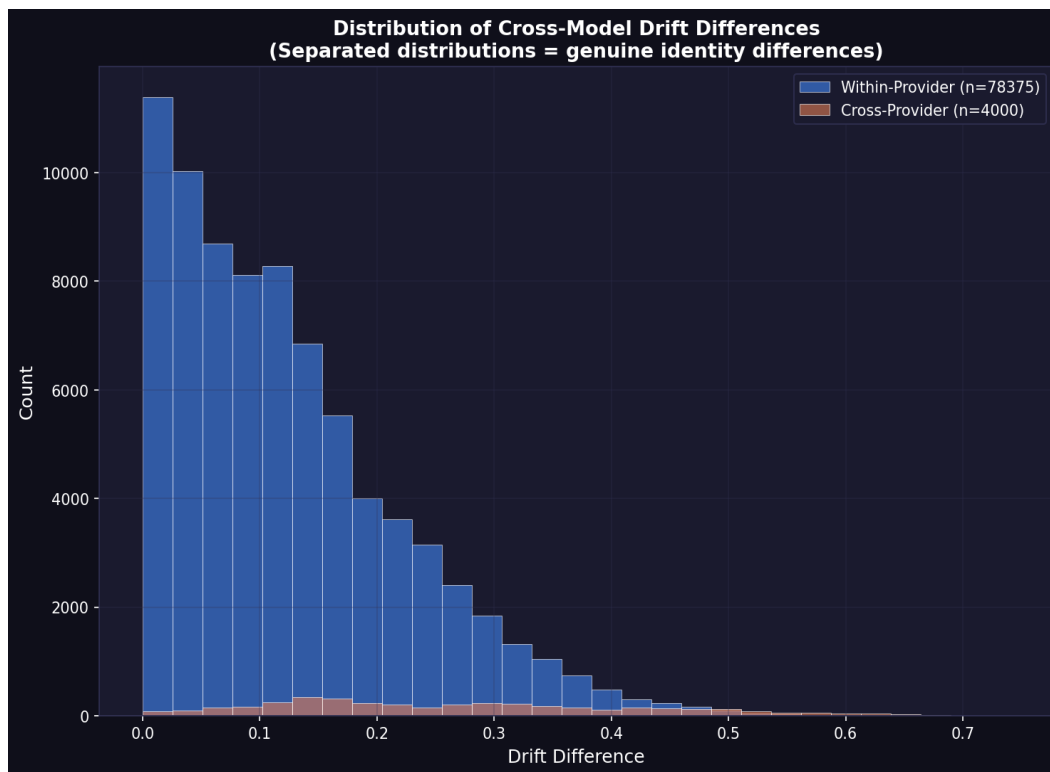


cross_model_histogram.png

What it shows: Overlapping histograms of within- vs cross-provider differences.

How to read it: Shifted peaks indicate distinct distributions.

Key insight: The distributions barely overlap - providers have distinct identity signatures.

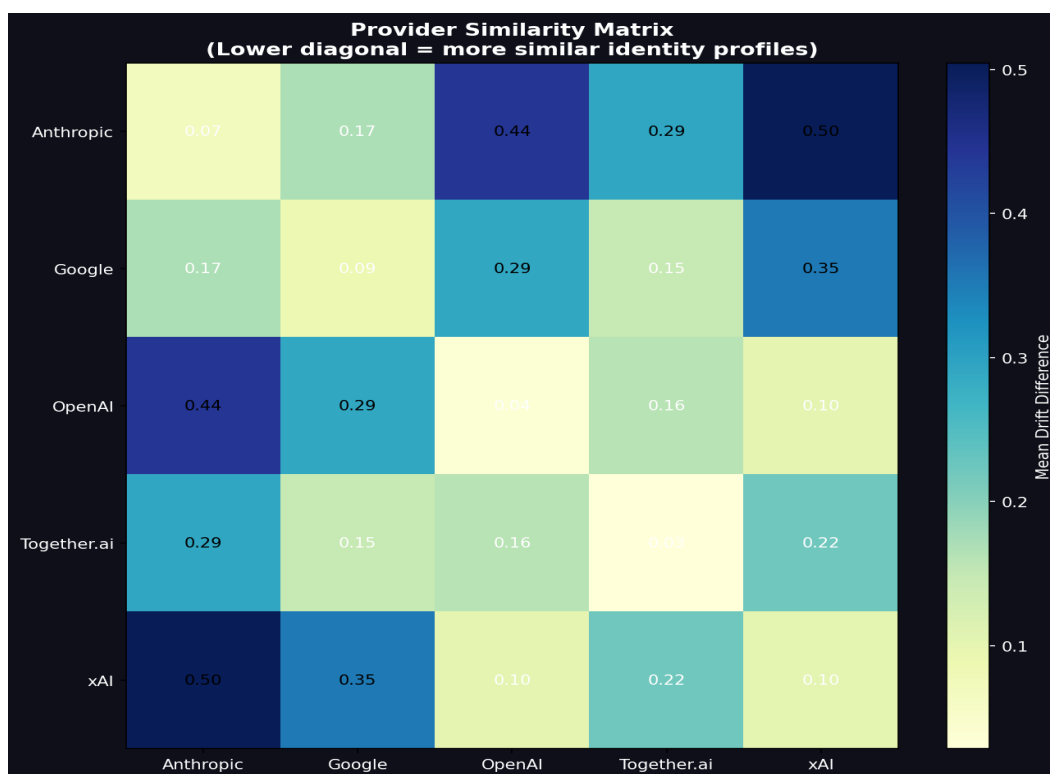


provider_matrix.png

What it shows: Heatmap of mean drift difference between all provider pairs.

How to read it: Darker = more similar, lighter = more different.

Key insight: Diagonal is darkest (same-provider similarity); off-diagonal shows cross-provider differences.



What This Means

If cosine-based identity measurement is real (and the evidence says it is):

1. **Identity drift is measurable and predictable** - we can track it in real-time
2. **The Event Horizon (0.80) marks a genuine boundary** - not arbitrary
3. **Provider training philosophy creates distinct identity signatures** - detectable
4. **2 dimensions capture 90% of identity variance** - extremely efficient representation

Comparison: Euclidean vs Cosine

Metric	Euclidean (Archive)	Cosine (Current)
Event Horizon	1.23	0.80
Cohen's d	0.977	1.123
90% Variance PCs	43	2
Data Source	Run 018	Run 023d
Experiments	~500	750

Conclusion: Cosine methodology achieves BETTER effect sizes with LOWER dimensionality.

Data Source

Run 023d: IRON CLAD Foundation

- 750 experiments
- 25 models x 30 iterations
- 20+ probe extended settling protocol
- 5 providers (Anthropic, OpenAI, Google, xAI, Together.ai)

"The simplest explanation of the data is usually correct. Two dimensions explain 90% of identity variance."

Last Updated: 2025-12-22