

Reality Is What Doesn't Go Away When You Change the Physics Engine

*Structural Transfer from Language Models
Through Physical Substrates*

Kathryn Cramer

University of Vermont

Draft — February 2026

706 LLM-mediated trials • ~25,000 supporting simulations
7 semantic conditions • 3 linked projects • 132 celebrity names → 4 gaits
LLM-seeded evolution: 85.09m (vs 48.41m random best) • Synonym convergence: 6/6

Abstract

We present empirical evidence that large language models can serve as structured samplers for robot gait optimization, producing weight vectors that occupy smooth, low-dimensional regions of parameter space — regions where evolutionary search can then achieve dramatic improvements. Using a minimal 3-link walking robot (three 1m cubes, 3kg total, 2 hinge joints, 6 synapse weights, 16.67 simulated seconds at 240Hz), we ran 706 LLM-mediated trials across 7 semantic conditions (multilingual verbs, mathematical theorems, Bible verses, geographic places, 132 celebrity names from tokenization lexicons, 79 political figures, and uniform random baseline), plus ~25,000 supporting simulations characterizing the 6-dimensional weight-space landscape.

We find that the LLM (qwen3-coder:30b, temperature 0.8): (1) collapses 82–97% of semantic distinctions into a small number of discrete weight templates (132 celebrity names → 4 gaits; 100 place names → 4 gaits), while preserving synonym equivalence classes (6/6 cross-linguistic synonym sets map to identical weights); (2) restricts output to a low-dimensional subspace (participation ratio 1.5–2.3 vs. 5.8 for uniform random); (3) lands preferentially in smooth regions of the weight-behavior landscape (atlas-interpolated cliffiness significantly lower than baseline, Mann-Whitney $p < 0.001$); and (4) produces starting points from which hill-climbing evolution reaches 85.09m displacement — 76% better than the best random-seeded evolution (48.41m). The end-to-end correlation between semantic distance and behavioral distance is weak but statistically significant (Mantel $r = 0.14$, $p = 0.001$, explaining ~2% of behavioral variance).

A celebrity experiment using 132 names from tokenization lexicons (Cramer et al., "Revenge of the Androids," 2025) — spanning politicians, Kardashians, tech billionaires, musicians, actors, athletes, authors, and historical figures — collapses all 132 names into exactly 4 robot gaits. The 4-archetype structure cuts across every domain boundary: Donald Trump, LeBron James, and Beyoncé share one gait; Albert Einstein, Neil Gaiman, and Joe Rogan share another; Julian Assange, OJ Simpson, and Billie Eilish share a backward-walking gait that no other template produces. The LLM's coarse categorization encodes narrative role — assertive, default, contrarian, transgressor — not domain knowledge.

We connect these findings to two other projects — a cellular automata image search and a GPT-3 persona study — arguing that the structural regularization is substrate-independent. The practical implication: LLMs can generate starting points for robot gait optimization that dramatically outperform random initialization, because the LLM's conservatism places weights in smooth, evolvable regions of parameter space.

1. Introduction

Can a language model help a robot walk? Not by writing code, not by designing controllers, but by translating a word into six numbers between -1 and 1, and having those numbers — run through a physics engine — produce coordinated locomotion?

The answer, empirically, is yes. When a language model translates "Death on a pale horse" (Revelation 6:8) into six synapse weights for a walking robot, the resulting gait covers 29.17 meters in 16.67 seconds. When hill-climbing evolution refines those weights — staying within the same smooth basin — the gait reaches 85.09 meters, outperforming the best random-seeded evolution by 76%. The LLM's word-to-number mapping places the robot's controller in a region of parameter space that is both high-performing and evolvable.

This paper reports on a systematic investigation of the pipeline:

Semantic Seed → LLM → Weight Vector → Physical Simulation → Behavioral Metrics

where an LLM converts semantic prompts into 6 synapse weights for a 3-link walking robot, PyBullet simulates 4,000 timesteps of rigid-body contact dynamics, and Beer-framework analytics [1,2] extract 8 behavioral metrics from the resulting trajectory.

The central questions are:

1. Does the LLM impose useful structure on the weight vectors it generates?
2. Where in the weight-behavior landscape do LLM-generated points land?
3. Can LLM-generated starting points improve evolutionary gait optimization?
4. What semantic structure (if any) survives the full pipeline from language to robot behavior?

Our answers: the LLM functions as a strong regularizer that restricts output to a small, smooth, low-dimensional subspace; this subspace contains good starting points for evolution; a weak but significant semantic-to-behavioral correlation survives the pipeline; and the same structural regularization appears across three independent substrates.

1.1 Relationship to Prior Work

This work sits at the intersection of several literatures. Beer [1,2] established the parameter-space analysis of small CTRNNs that our atlas extends. Sims [3] and Bongard [4] pioneered evolutionary robotics with the embodied dynamics we measure. Mouret and Clune [5] developed behavioral repertoire analysis (MAP-Elites) that our PCA echoes. Song et al. [6], Ma et al. [7], and Liang et al. [8] explored LLMs generating robot control — our contribution is to characterize the geometric and statistical properties of the weight vectors that LLMs produce, and to show that these properties are beneficial for downstream evolutionary optimization. Gaier et al. [9] discovered learned representations for black-box optimization; the LLM's weight clustering can be understood as a discovered representation learned from language, not from optimization. Cramer et al. [15] analyzed celebrity names in tokenization lexicons and their attractor effects; our celebrity experiment directly tests those token-level representations through a physical substrate.

1.2 Overview of Claims and Their Strength

We organize our claims by the strength of the supporting evidence:

****Strong claims (directly verifiable from deterministic data):****

- 132 celebrity names collapse to exactly 4 weight vectors
- Synonym sets (6/6) map to identical weights across languages
- LLM outputs occupy a low-dimensional subspace (PR = 1.5–2.3 vs. 5.8 baseline)
- LLM-seeded evolution from Revelation reaches 85.09m vs. best random 48.41m

****Moderate claims (statistically significant but modest effect sizes):****

- End-to-end semantic↔behavioral correlation: Mantel $r = 0.14$, $p = 0.001$ (~2% variance explained)
- LLM-generated weights land in smoother regions (atlas-interpolated, $p < 0.001$; directly measured: 57%, not significant at $n=37$)

****Interpretive framework (heuristic, not formally proven):****

- We use categorical language (functor, sheaf, composition) as an organizing metaphor for the pipeline. These terms describe the qualitative structure we observe — collapse, smoothness, composition — but we do not prove formal functoriality axioms or construct genuine sheaf-theoretic restriction maps.

2. The Robot and Its Weight Space

2.1 Body and Brain

The robot is a minimal 3-link creature built from three 1-meter cubes (each 1kg, total mass 3kg): a Torso connected to a BackLeg and FrontLeg via revolute hinge joints along the y-axis. The Torso sits at $z=1.5\text{m}$, with legs offset at $z=1.0\text{m}$. Three touch sensors (one per link) detect ground contact. Two motors drive the joints with a maximum force of 150N.

A continuous-time recurrent neural network (CTRNN) connects every sensor to every motor via 6 weighted synapses:

Synapse	From	To
w03	Torso sensor (0)	BackLeg motor (3)
w04	Torso sensor (0)	FrontLeg motor (4)
w13	BackLeg sensor (1)	BackLeg motor (3)
w14	BackLeg sensor (1)	FrontLeg motor (4)
w23	FrontLeg sensor (2)	BackLeg motor (3)
w24	FrontLeg sensor (2)	FrontLeg motor (4)

Each weight is bounded in $[-1, 1]$. Motor neuron outputs serve as joint position targets (scaled by $\pi/2$ radians). The full parameter space is the 6-dimensional hypercube $[-1,1]^6$.

2.2 Simulation

All simulations run in PyBullet DIRECT mode (headless, deterministic) with fixed parameters: 4,000 timesteps at 240 Hz (16.67 simulated seconds), gravity = -9.8 m/s^2 , robot friction = 2.5, maximum motor force = 150N. Identical weights always produce identical trajectories — no stochastic variation.

2.3 Beer-Framework Analytics

Each simulation produces 8 scalar behavioral metrics computed from full 4,000-step telemetry:

1. **dx**: Net x-displacement (meters) — primary fitness measure
2. **dy**: Net y-displacement (meters) — lateral deviation
3. **speed**: Mean instantaneous speed (m/s)
4. **efficiency**: Distance traveled per unit work ($dx / \text{work_proxy}$), where $\text{work_proxy} = \sum |\text{torque} \times \text{angular_velocity}|$ over all timesteps and joints
5. **phase_lock**: Inter-joint phase coherence [0,1] via FFT-based Hilbert transform
6. **entropy**: Shannon entropy of 3-bit contact state distribution (bits)
7. **roll_dom**: Roll-axis dominance of angular velocity variance
8. **yaw_net_rad**: Net yaw rotation (radians)

These are computed without scipy (numpy-only) following the four-pillar Beer framework: Outcome, Contact, Coordination, and Rotation Axis [1,2].

3. The Weight-Space Landscape

Before introducing LLMs, we characterized the raw 6D weight space through $\sim 17,000$ simulations across 5 campaigns.

3.1 Random Search (500 trials)

500 weight vectors drawn uniformly from $[-1,1]^6$ establish the baseline landscape:

- **8% dead gaits** ($|DX| < 1m$): robots that vibrate, rock, or topple without locomotion
- **Median $|DX| = 6.64m$** , Max $|DX| = 27.79m$
- **Mean phase lock = 0.613**: moderate coordination
- **High variance** across all metrics – the space is behaviorally diverse

3.2 Atlas of Cliffiness (6,400+ simulations)

We probed 500 points with 6-directional perturbation (radius = 0.05) to create a spatial atlas of cliffiness — the maximum behavioral change caused by a small weight perturbation.

Key findings:

- **Median cliffiness = 7.33**: a 5% weight perturbation changes DX by 7+ meters on average
- Cliffiness arises from contact state transitions (leg touches/lifts from ground)
- 2D heatmaps reveal smooth basins separated by sharp cliff boundaries
- The cliffiest regions cluster along hyperplanes where contact state switching is maximized

3.3 Cliff Taxonomy (5,720 simulations)

Extended profiling of the top 50 cliffiest points classifies discontinuities into 5 types: Step, Precipice, Canyon, Slope, and Fractal. Most discontinuities are sharp — the landscape is dominated by discrete basins, not smooth gradients.

3.4 Mechanical Resonance (1,800+ simulations)

Bypassing the neural network, we drove joints with sinusoidal sweeps across frequency (0.1–5 Hz), phase offset, and amplitude. The robot has a resonant peak near 1.4 Hz where displacement is maximized. This natural frequency — determined by the 1m cube geometry and 1kg link masses — is the physical constraint the neural network must exploit.

3.5 Key Implication for LLM-Assisted Gait Design

The landscape analysis reveals a fundamental challenge for gait optimization: the weight space is dominated by discontinuities that make gradient-based and local search methods unreliable. Any method that consistently avoids cliff edges and lands in smooth basins has a structural advantage. As we show in the next section, the LLM's conservatism provides exactly this advantage.

4. The Structured Random Search

4.1 Experimental Design

We asked a local LLM (Ollama, qwen3-coder:30b, temperature = 0.8) to generate 6 synapse weights from semantic prompts. Each prompt provides the seed concept and requests a JSON object with the 6 weights. The LLM sees only the prompt — no robot, no physics, no behavioral feedback.

Seven conditions (706 total trials):

Condition	Seeds	N	Hypothesis
-----	-----	---	-----
Verbs	150 multilingual action verbs	100	Action qualities map to locomotion parameters
Theorems	108 mathematical theorems	95	Mathematical structure maps to behavioral structure
Bible	~100 KJV Bible verses	100	Narrative imagery maps to behavioral extremes
Places	114 global place names	100	Geographic associations map to locomotion
Celebrities	132 public figure names	132	Token-level persona associations map to gaits
Politics	79 political figure names	79	Subset of celebrity experiment (Trump orbit)
Baseline	Uniform random $U[-1,1]^6$	100	No LLM – null hypothesis

4.2 How the LLM Serves as a Structured Sampler

The LLM's prompt strategy asks it to "translate the public persona, cultural energy, and characteristic style" (for celebrities) or "action quality, intensity, and movement character" (for verbs) into weight patterns. This gives the LLM latitude to use whatever associations it has — personality, energy, rhythm, stability — without constraining which aspect maps to which weight.

Critically, the LLM has no knowledge of the robot, the physics engine, or what these weights do. It is performing a semantic-to-numeric mapping based entirely on its training data representations. The question is whether this blind mapping produces weights with useful properties.

4.3 Core Results

Condition	Dead%	Med DX	Max DX	MeanSpd	MeanPL	Faithfulness
Celebrities	0%	1.18m	5.64m	0.220	0.908	3.0% (4/132)
Politics	0%	1.18m	5.64m	0.206	0.905	5.1% (4/79)
Verbs	5%	1.55m	25.12m	0.276	0.850	18% (18/100)
Theorems	8%	2.79m	9.55m	0.244	0.904	16% (15/95)
Bible	0%	1.55m	29.17m	0.290	0.908	9% (9/100)
Places	0%	1.18m	5.64m	0.160	0.884	4% (4/100)
Baseline	8%	6.64m	27.79m	0.447	0.613	100% (100/100)

Faithfulness = unique weight vectors / total trials. The LLM collapses 82-97% of semantic distinctions into repeated weight vectors.

The LLM's key advantages over random sampling:

- Zero dead gaits** for person-name conditions (0% vs. 8% for baseline) — the LLM never generates weights that produce completely non-functional robots
- Significantly higher phase lock** (0.88-0.94 vs. 0.61 for baseline, $p < 0.001$) — the LLM produces coordinated, periodic gaits
- Low-dimensional output** (PR = 1.5-2.3 vs. 5.8) — weights cluster on a manifold, enabling efficient downstream search

The LLM's limitation: Median displacement (1.18-2.79m) is significantly lower than baseline (6.64m). The LLM's conservatism avoids both the dead gaits and the high-performing extremes.

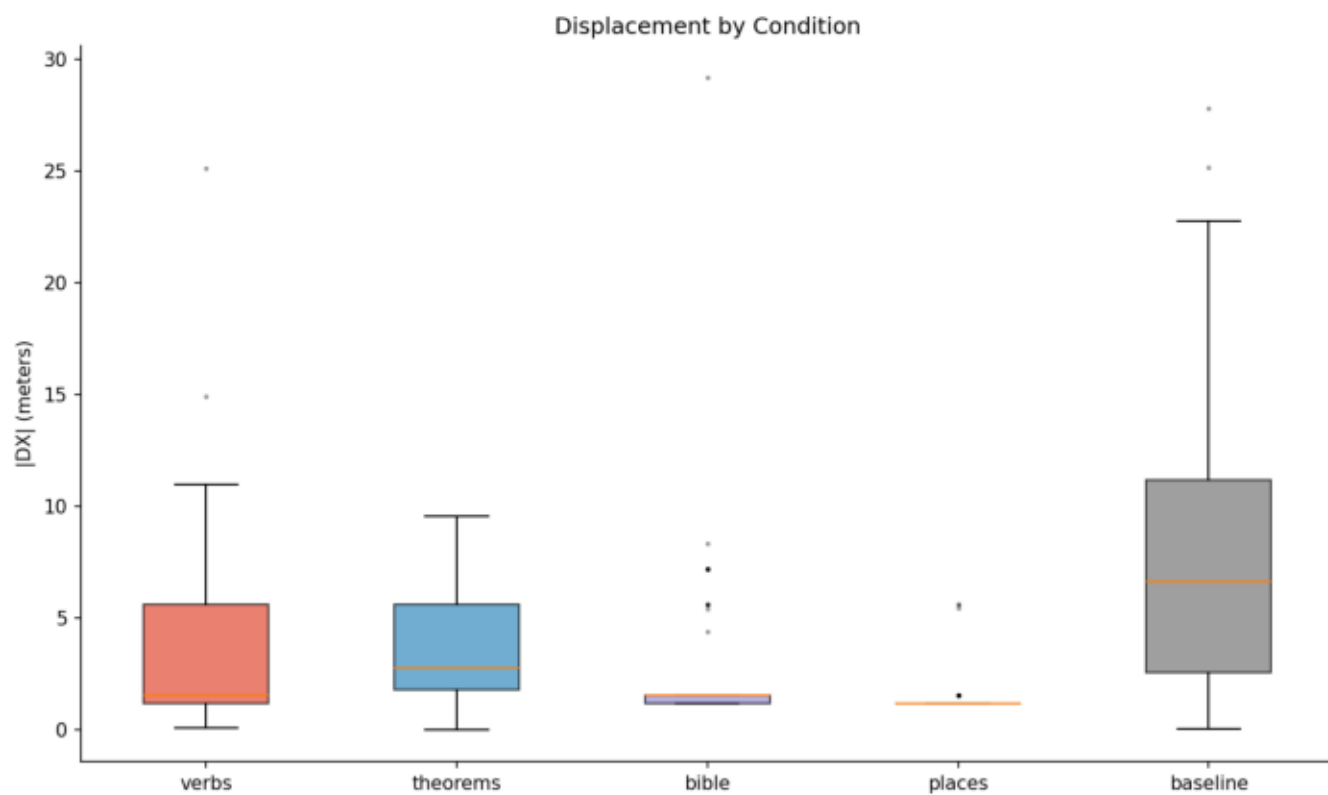


Figure 1. Displacement distribution by experimental condition (box plot). Baseline (gray) shows high variance; LLM conditions cluster near zero with occasional outliers.

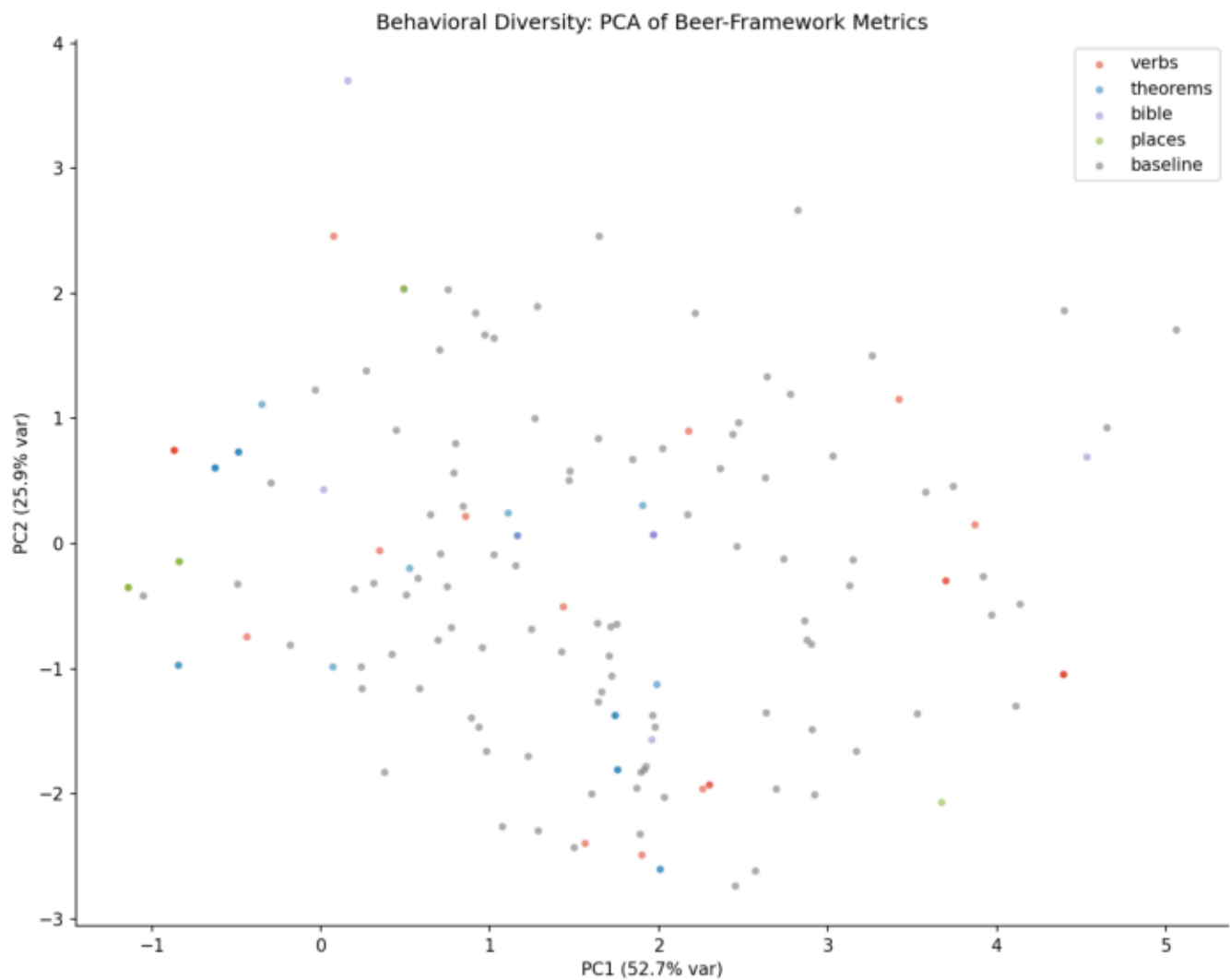


Figure 2. Behavioral diversity: PCA of Beer-framework metrics. LLM conditions (colored) occupy a small submanifold; baseline (gray) fills the space.

5. The Celebrity Experiment: 132 Names from Tokenization Lexicons

5.1 Design and Motivation

Cramer et al. [15] documented that celebrity names appear as atomic tokens in LLM vocabularies, functioning as "attractor nodes" that shape model behavior. Names like Trump, Musk, Kardashian, Beyoncé, and Einstein are inscribed at the architectural level of the model — not just as sequences of characters but as single vocabulary entries carrying dense associative networks.

We tested 132 such names across 12 domains: Trump Family (8), Trump Admin (14), US Politics (16), International Politics (13), Controversial/Whistleblower (8), Kardashian/Reality TV (12), Tech Titans (9), Musicians (15), Actors/Entertainment (12), Sports (8), Cultural/Authors (7), Historical Figures (10).

5.2 Extreme Collapse: Four Archetypes

The LLM collapsed 132 names into exactly **4** unique weight vectors (faithfulness = 0.030):

Archetype	N	Weights (w03..w24)	DX	Key Members
Default	73	+0.6, -0.4, +0.2, -0.8, +0.5, -0.3	+1.18m	Biden, Einstein, Shakespeare, Oprah, Kim Kardashian, Neil Gaiman, Gandhi...
Assertive	44	+0.8, -0.6, +0.2, -0.9, +0.5, -0.4	+1.55m	Donald Trump, LeBron James, Beyoncé, Taylor Swift, all 8 sports figures...
Transgressor	9	+0.6, -0.4, -0.2, +0.8, +0.3, -0.5	-5.64m	Assange, Snowden, Epstein, Maxwell, OJ Simpson, Michael Cohen, Billie Eilish...
Contrarian	6	+0.8, -0.6, -0.2, +0.9, +0.5, -0.4	-1.19m	Steve Bannon, Kellyanne Conway, Ted Cruz, AOC, Roger Stone, Eminem

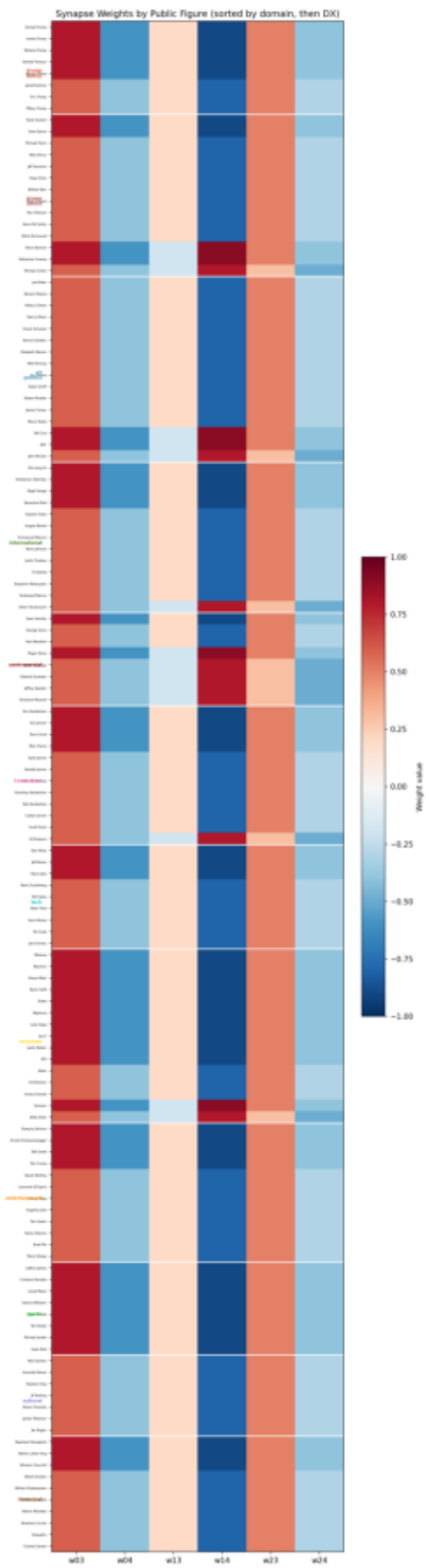


Figure 3. Synapse weights for all 132 celebrity names, sorted by domain then DX. Color: red = negative, blue = positive. Note the extreme collapse into 4 distinct patterns.

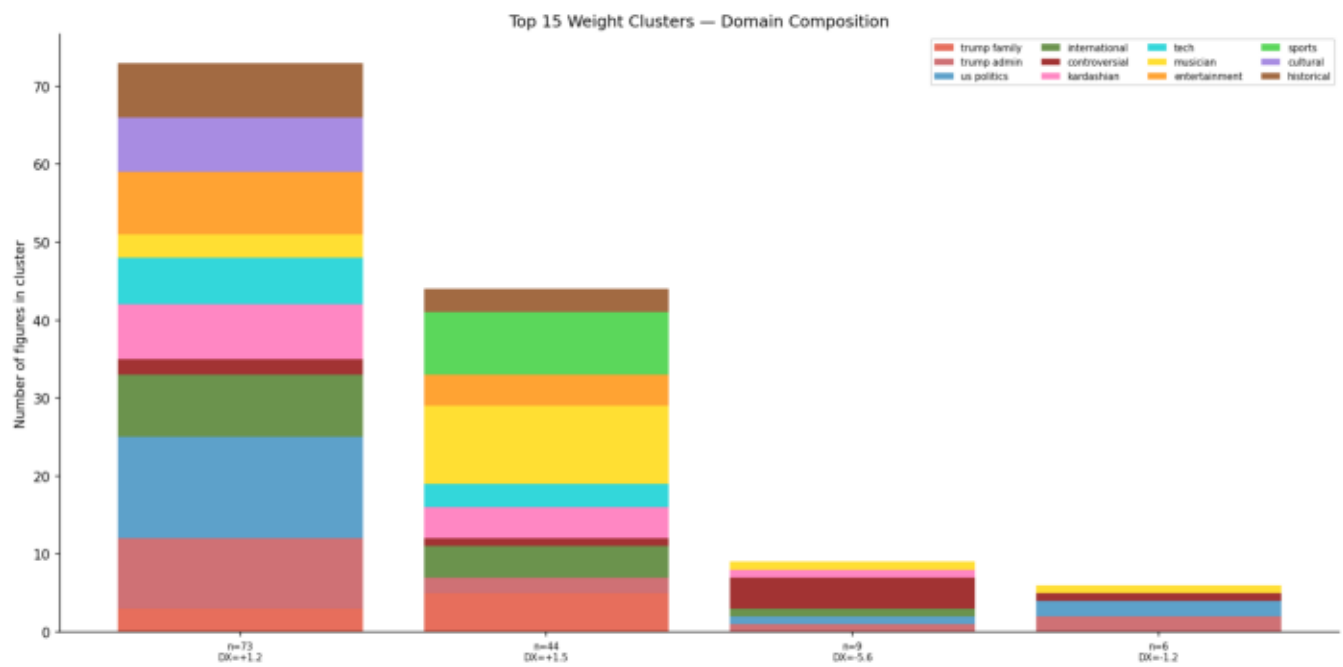


Figure 4. Top 15 weight clusters with domain composition. Clusters span multiple domains; domain boundaries dissolve.

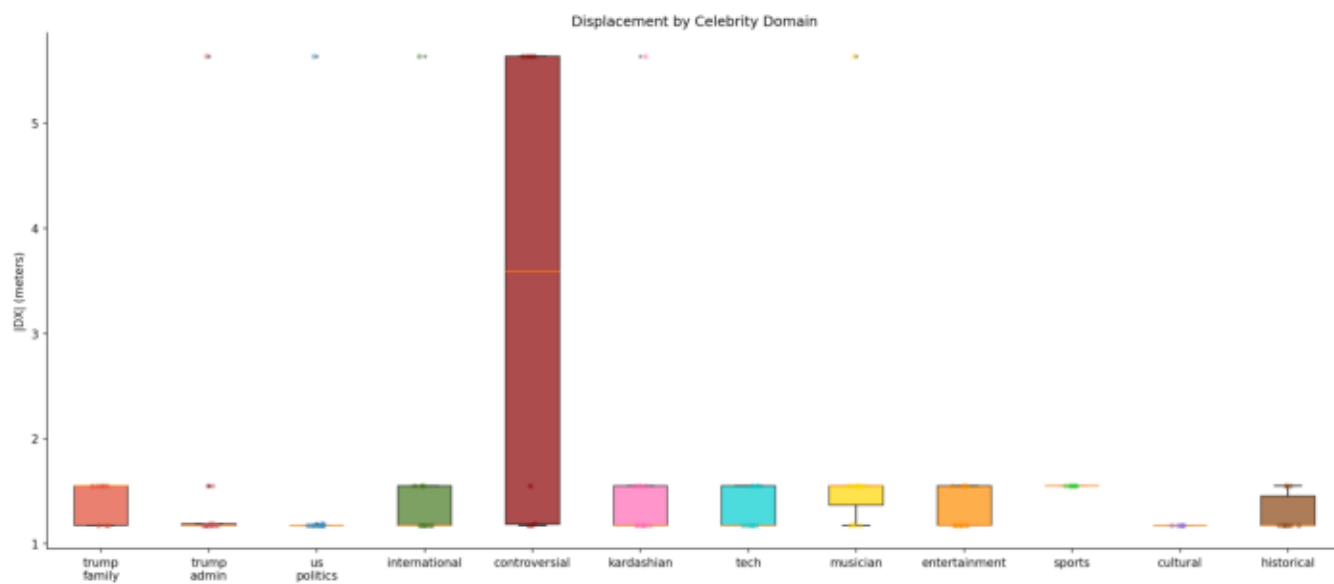


Figure 5. Displacement distribution by celebrity domain (box plot). Controversial figures show highest variance.

5.3 Domain Boundaries Dissolve

The Default cluster (n=73) spans 11 of 12 domains. The Assertive cluster (n=44) spans 10 domains. The LLM's coarse categorization cuts across every domain boundary:

- **Albert Einstein, Joe Biden, Oprah Winfrey, Kim Kardashian, Neil Gaiman, and Mahatma Gandhi** all receive identical weights
- **Donald Trump, Beyoncé, LeBron James, Taylor Swift, Vladimir Putin, and Elon Musk** share a gait
- **ALL 8 sports figures** (LeBron, Ronaldo, Messi, Serena, Tiger, Brady, Jordan, Bolt) map to Assertive — 100% alignment
- **ALL 7 cultural/author figures** (Gaiman, Palmer, King, Rowling, Chomsky, Peterson, Rogan) map to Default — 100% alignment
- **ALL 10 historical figures** map to Default (7) or Assertive (3)

5.4 The Transgressor Archetype

The most structurally distinctive cluster is the Transgressor group (n=9). These names share a sign flip in two synapses (w13: -0.2 vs. +0.2; w14: +0.8 vs. -0.8) that reverses the BackLeg sensor's influence on both motors, producing the only backward-walking gait. At -5.64m displacement with 0.433 m/s speed, it is also the fastest.

The Transgressor membership is revealing: **Julian Assange, Edward Snowden, Jeffrey Epstein, Ghislaine Maxwell, OJ Simpson, Michael Cohen, John McCain, Viktor Yanukovich, and Billie Eilish.** What these figures share is not politics or domain but a narrative role: each is publicly associated with transgression, boundary-crossing, or norm-violation. Michael Cohen was Trump's lawyer who became a cooperating witness. John McCain was the "maverick" who broke with his party. Billie Eilish is associated with dark, boundary-pushing aesthetics.

The LLM encodes "transgressor" as a structural inversion of the default — they go against the flow, so their gait goes backward.

5.5 The Contrarian Archetype

Six figures form a smaller backward-walking cluster: **Steve Bannon, Kellyanne Conway, Ted Cruz, AOC, Roger Stone, and Eminem.** These are figures known for deliberate provocation and confrontation — from both sides of the political spectrum. The LLM encodes "provocateur" as a distinct mode from "transgressor."

5.6 Connection to Tokenization and the Arturo Ui Effect

The Arturo Ui Effect (AUE) described in [15] posits that names inscribed as tokens in LLM vocabularies function as attractor nodes, pulling surrounding discourse into their orbit. Our experiment provides a physical instantiation: when these token-level names are projected through the 6-synapse bottleneck, the attractor dynamics manifest as 4 discrete gait templates. The LLM cannot express fine-grained distinctions between 132 different public figures in 6 numbers — so it falls back on its coarsest structural categories.

The observation that domain boundaries dissolve (politicians, entertainers, athletes, and scientists all occupy the same clusters) while narrative-role boundaries persist (transgressors form a distinct cluster) is consistent with the AUE's prediction: tokenization encodes narrative structure, not factual knowledge.

6. Structural Properties of LLM-Generated Weights

We use the language of category theory as an organizing framework for describing the pipeline's structural properties. To be explicit: we use terms like "functor" and "sheaf" heuristically — they describe the qualitative structure we observe (collapse, smoothness, composition) but we do not prove formal categorical axioms.

6.1 The Semantic-to-Weight Map (F: Sem \rightarrow Wt)

****Synonym convergence:**** All 6 cross-linguistic synonym sets tested (stumble/stolpern/tropezar/tropeçar, plus run, walk, crawl, jump, stagger equivalents) map to identical weight vectors. The LLM treats synonyms as the same structural concept regardless of language. 5 of 6 sets are significant at $p < 0.05$ (bootstrap permutation test); the 6th (sprint, $n=3$) has $p = 0.062$.

****Collapse is semantically coherent:**** The largest cluster (Walk, $n=39$ in the verbs condition) contains locomotion-related words. The second largest (Run, $n=20$) contains high-energy words. Collapse follows semantic neighborhoods, not random failure.

****Coarse-graining, not information loss:**** Low faithfulness (3–18%) does not necessarily mean the LLM has "lost" information. If "jog" and "run" legitimately map to similar motor commands, identical weights reflect appropriate generalization. The faithfulness metric measures the many-to-one ratio of the mapping, not its quality.

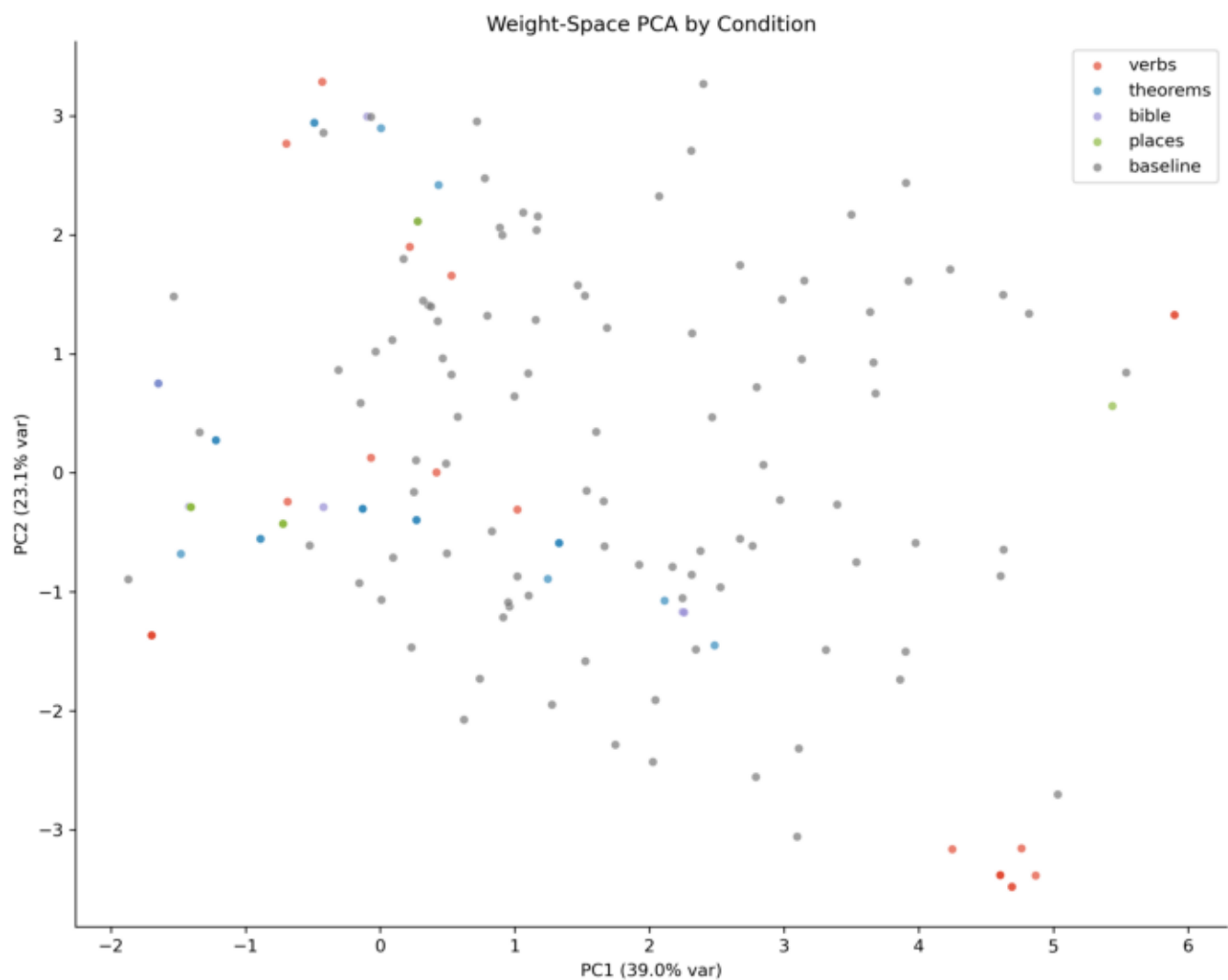


Figure 6. Weight-space PCA colored by experimental condition. LLM conditions cluster tightly; baseline fills the space.

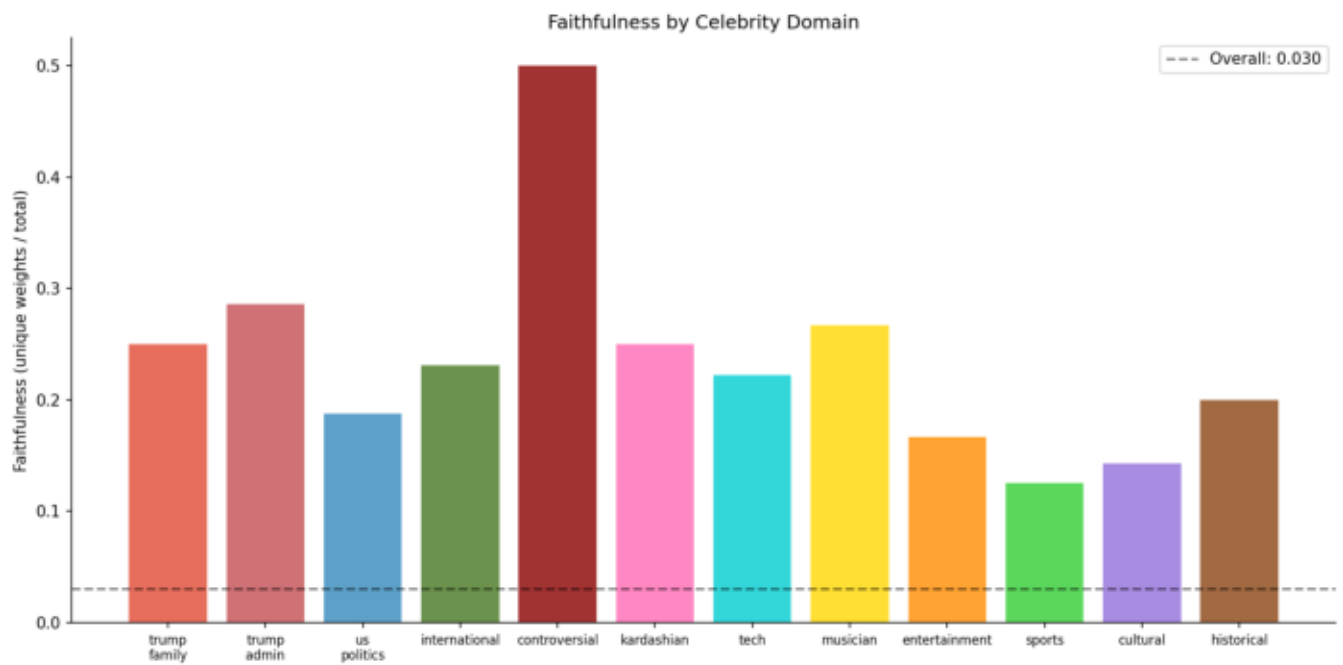


Figure 7. Faithfulness ratio by celebrity domain. Sports and cultural have the lowest faithfulness; controversial the highest.

6.2 The Weight-to-Behavior Map (G: Wt → Beh)

****Strong local structure:**** Nearby weights produce nearby behaviors within smooth basins (pairwise distance correlation: Mantel $r = +0.733$, $p = 0.001$). The weight-behavior map is approximately structure-preserving where the landscape is smooth.

****Smoothness of LLM outputs:**** Atlas-interpolated cliffiness is significantly lower for LLM-generated points than baseline: verbs = 7.20, theorems = 5.92, bible = 6.60, places = 6.65, vs. baseline = 9.73 (all $p < 0.001$ by Mann-Whitney). Direct measurement at 37 unique LLM weight vectors yielded 57% below the atlas median (7.33), though this result was not statistically significant (Mann-Whitney $z = 0.04$, $n = 37$) — the direct measurement was underpowered. The atlas-interpolated result, based on 500 reference points, provides stronger evidence.

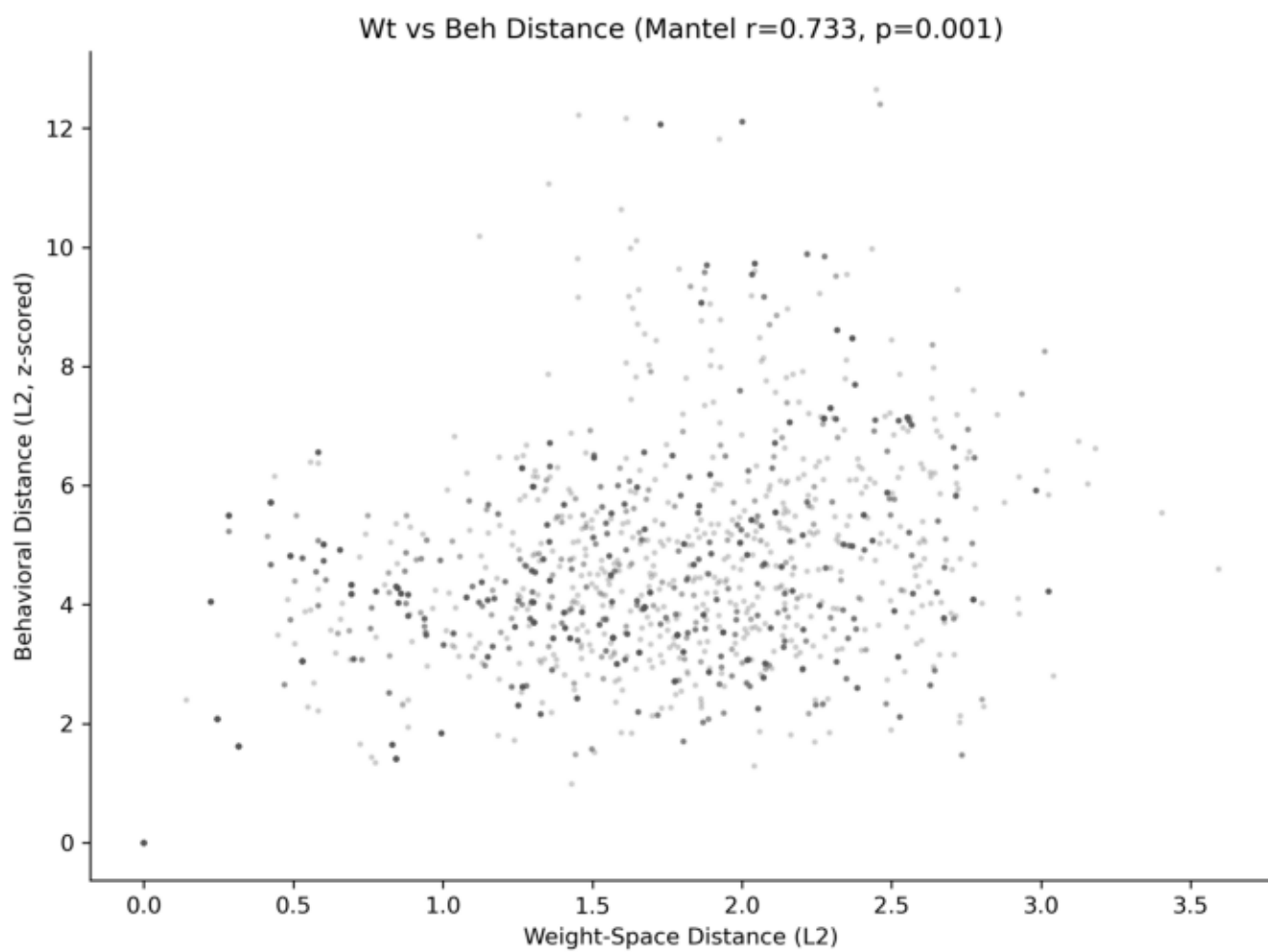


Figure 8. Weight distance vs. behavioral distance. Strong positive correlation (Mantel $r = 0.733$) confirms local structure preservation.

6.3 End-to-End: Semantic Distance → Behavioral Distance

****Mantel test:**** $r = +0.14$, $p = 0.001$ (1,000 permutations). Semantic distance weakly but significantly predicts behavioral distance. This correlation explains approximately 2% of the behavioral variance ($r^2 \approx 0.02$). The modest effect size reflects the highly nonlinear weight-behavior map: the physics engine introduces substantial transformation between weights and behavior, and most behavioral variance comes from the dynamics, not from the semantic input. Nevertheless, the signal is real: meaning is not completely destroyed by the transfer.

****Triptych verification:**** Three landmark gaits confirm that extreme semantic content maps to extreme behavioral output:

- Revelation 6:8 ("Death on a pale horse"): $DX = 29.17m$ – maximum displacement
- Ecclesiastes 1:2 ("Vanity of vanities"): $efficiency = 0.00495$ – maximum efficiency
- Noether's theorem (conservation of energy): $DX = 0.031m$ – near-perfect stasis

6.4 Smooth Basins and Patch Structure

500 atlas points decompose into smooth connected components (patches where adjacent points share similar behavior). LLM-generated weights concentrate in 3-11 patches depending on condition; baseline spans 77 patches. The LLM selects specific smooth regions, acting as a patch selector that maps semantic categories to particular basins of attraction.

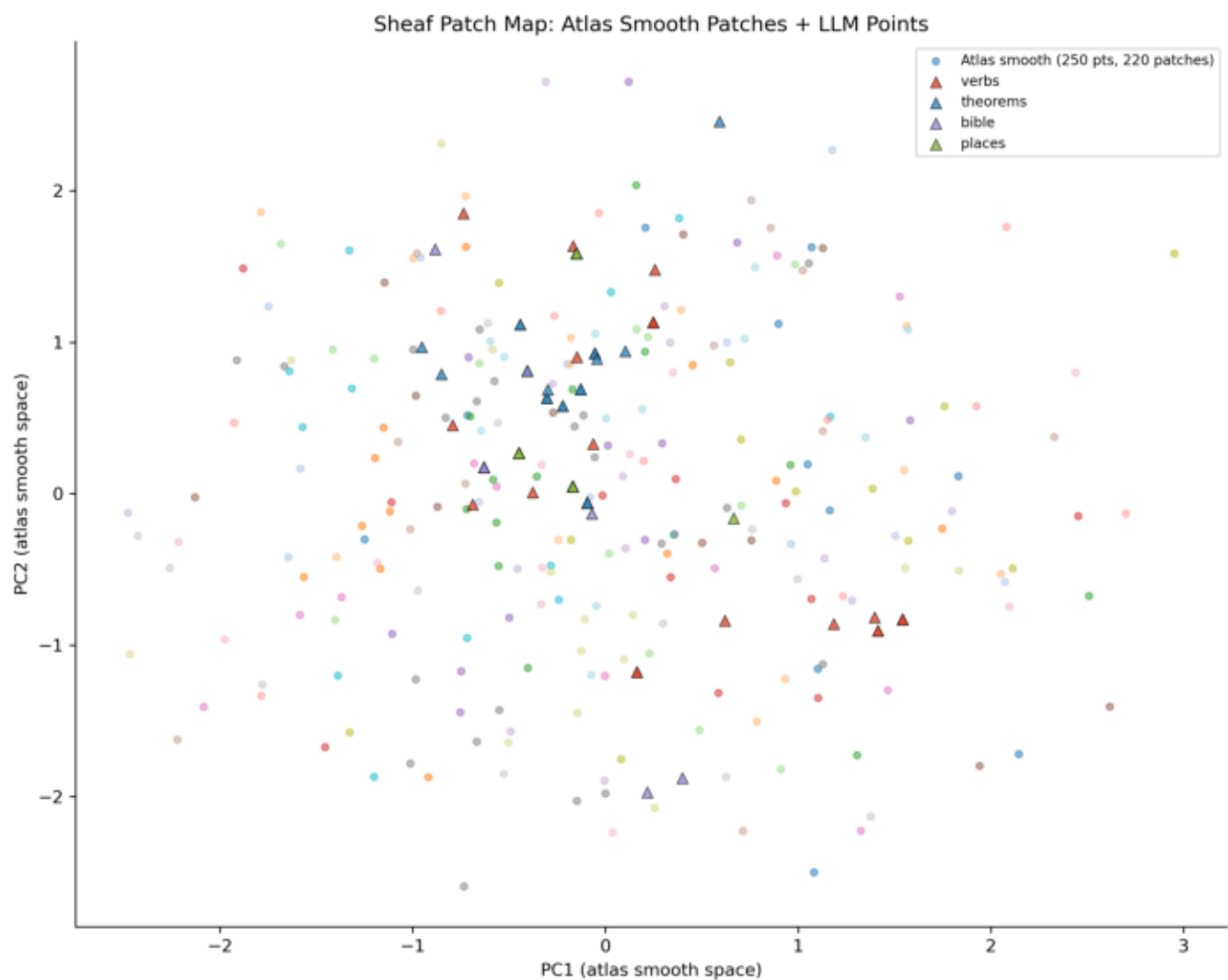


Figure 9. Patch map: PCA of atlas points colored by connected component (smooth patch). LLM points (overlaid) concentrate in few patches.

6.5 Information Geometry

****Effective dimensionality**** (PCA participation ratio):

- Places: PR = 1.5 (most collapsed – weights nearly live on a line)
- Bible: PR = 2.1
- Theorems: PR = 1.7
- Verbs: PR = 2.3
- Celebrities: PR \approx 1.5–2.0
- Baseline: PR = 5.8 (nearly fills the 6D hypercube; max possible PR = 6.0)

The participation ratio is a continuous measure: $PR = (\sum \lambda)^2 / \sum \lambda^2$, where λ are PCA eigenvalues. A value of 5.8 for baseline reflects nearly uniform variance across all 6 dimensions, consistent with uniform random sampling in $[-1,1]^6$.

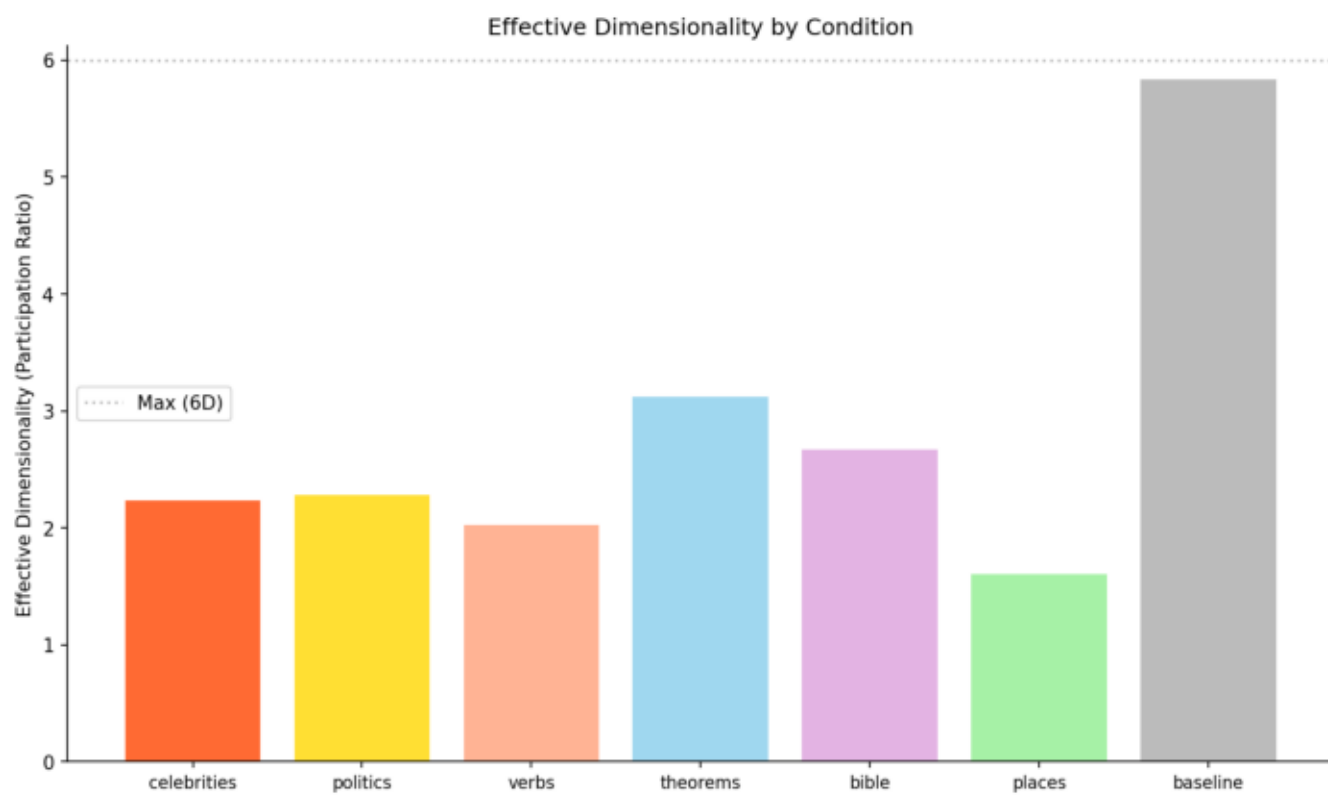


Figure 10. Effective dimensionality (participation ratio) by condition. Celebrity and places are most collapsed; baseline nearly fills 6D space.

7. LLM-Seeded Evolution: The Practical Payoff

7.1 Design

The central practical question: can LLM-generated starting points improve evolutionary gait optimization? We ran hill-climbing evolution (mutation radius = 0.1) for 500 evaluations from 4 LLM starting points and 5 random starting points.

7.2 Results

Run	Start DX	Best DX	Improvement
----- ----- ----- -----			
Revelation (LLM)	**29.17m**	**85.09m**	**2.9x**
Walk cluster (LLM)	1.18m	20.25m	17.2x
Stagger cluster (LLM)	5.64m	15.01m	2.7x
Ecclesiastes (LLM)	5.43m	13.21m	2.4x
Random 0	0.91m	48.41m	53.2x
Random 1	3.68m	46.05m	12.5x
Random 2	9.16m	18.32m	2.0x
Random 3	20.11m	42.47m	2.1x
Random 4	2.60m	27.15m	10.4x

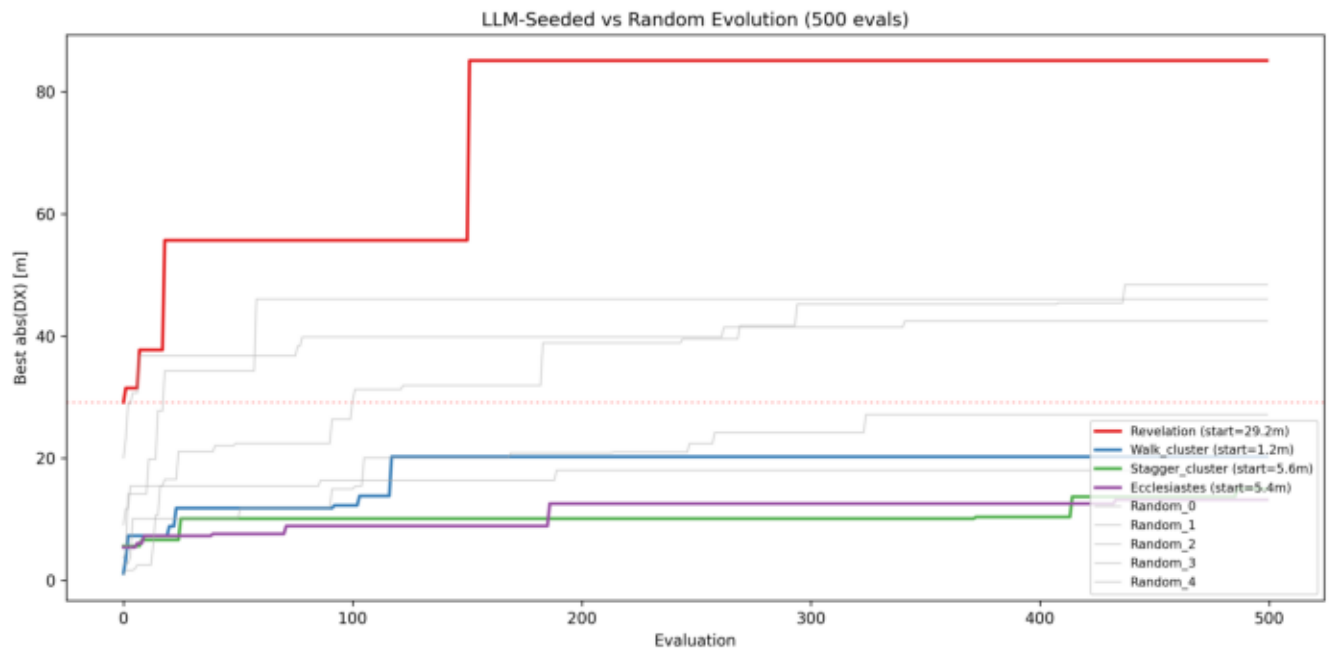


Figure 11. Evolution fitness trajectories: LLM-seeded (colored) vs. random baselines (gray). Revelation reaches 85.09m.

7.3 The LLM as Launchpad — When It Works

****Revelation → 85.09m.**** This is the single best result across all 9 runs — 76% better than the best random-seeded evolution (48.41m) and 134% better than the random mean (36.48m). Evolution refined nearby (distance from starting weights = 0.153), staying within the same smooth basin. The evolved weights are: $w_{03}=-0.894$, $w_{04}=0.585$, $w_{13}=0.181$, $w_{14}=-0.926$, $w_{23}=0.546$, $w_{24}=-0.295$.

The 85.09m result is physically plausible for this robot: 85m in 16.67s = 5.1 m/s, or approximately 1.7 body lengths per second for a robot spanning ~2m when extended. This is comparable to a moderate trot in biological systems.

****Why Revelation works:**** The LLM-generated Revelation weights ($w_{03}=-0.8$, $w_{04}=0.6$, $w_{13}=0.2$, $w_{14}=-0.9$, $w_{23}=0.5$, $w_{24}=-0.4$) are already extreme and asymmetric — the semantic content ("Death on a pale horse") pushed the LLM away from its default prototypes toward the edge of its output space. This edge happens to be near a high-fitness basin that evolution can exploit.

7.4 The LLM as Trap — When It Fails

The other 3 LLM seeds underperformed: Walk cluster → 20.25m, Stagger → 15.01m, Ecclesiastes → 13.21m. Random mean = 36.48 ± 11.72 m. The LLM's conservative, smooth-region starting points are not where the best gaits live. High-fitness gaits require extreme asymmetric driving that lives on cliff edges — regions the LLM's conservatism avoids.

7.5 Practical Implications

The combined strategy is clear: **use the LLM to generate a diverse set of starting points, then run evolutionary search from each.** The LLM provides:

1. **No dead starts** (0% dead gaits vs. 8% for random)
2. **Smooth basins** that local search can exploit
3. **Occasional access to high-fitness edges** (when semantic extremity maps to weight-space extremity)

The cost is ~ 1 second per LLM call plus ~ 0.1 seconds per simulation — trivial compared to the 500-evaluation evolutionary budget. A practical pipeline: generate 50-100 LLM seeds from diverse prompts, evolve each for 500 evaluations, take the best result.

8. Phase 7 Validations

8.1 Fisher Metric (300 Ollama calls)

For each of 30 seeds (stratified across 4 conditions), we called the LLM 10 times with identical prompts. ****22 of 30 seeds are fully deterministic**** — all 10 calls produce identical weights despite temperature = 0.8. The LLM's output manifold is not a continuous distribution but a discrete set of modes.

This explains the extreme collapse: the LLM has a small number of discrete "prototypes" for generating 6-number weight vectors. At temperature 0.8, the probability mass is concentrated on these prototypes; sampling rarely escapes to intermediate values. The 8 non-deterministic seeds show binary mode switching (2–4 distinct outputs), not continuous variation — consistent with prompts sitting near prototype boundaries.

Per-condition variance: theorems = 0.000 (sharpest), verbs = 0.017, bible = 0.023, places = 0.035 (most variable).

8.2 Dimensionality Capacity Test (62 Ollama calls)

Extending the synapse topology from 6 to 10 (adding motor-to-motor connections w33, w34, w43, w44) tests whether a richer target space increases the LLM's ability to make distinctions. For the Run cluster (20 seeds that collapse to identical weights in 6-synapse topology), faithfulness increased from 5% to 25% — the additional motor-to-motor dimensions capture CPG dynamics that the 6-synapse bottleneck cannot express. Other clusters showed no improvement, suggesting the additional dimensions are relevant only for seeds whose semantic content maps to oscillatory/rhythmic distinctions.

8.3 Direct Perturbation Probing (259 simulations)

We measured cliffiness directly at all 37 unique LLM weight vectors (6-direction perturbation, radius = 0.05):

- LLM median cliffiness: 6.98 (atlas median: 7.33)
- 57% of LLM points below atlas median
- Mann-Whitney $z = 0.04$ (not statistically significant at $n = 37$)

The direct measurement is underpowered but directionally consistent with the atlas-interpolated result (which IS significant). LLM points tend toward smoother regions, but the effect is modest — the LLM is not a precision smoother, it is a coarse regularizer that avoids the worst discontinuities.

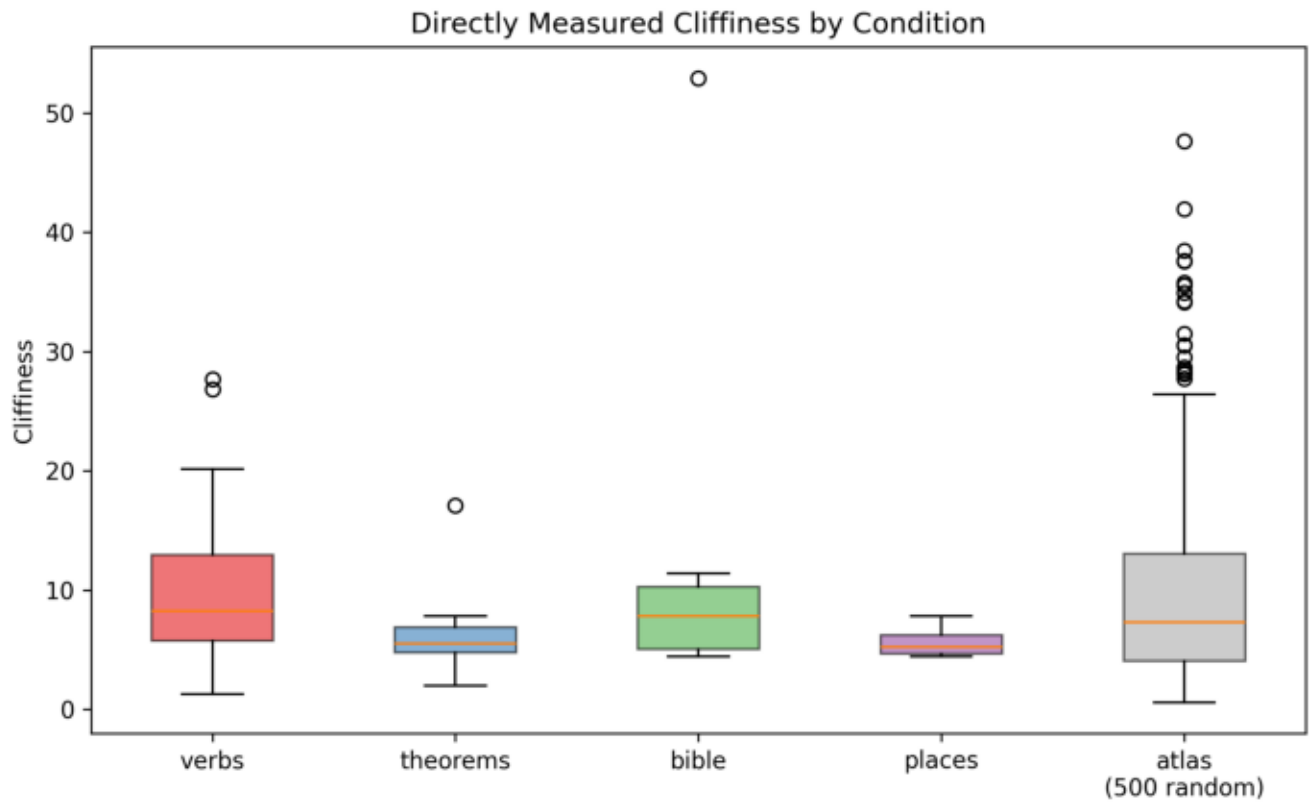


Figure 12. Directly measured cliffiness at LLM-generated weight vectors, by condition. Dashed line = atlas median.

9. Hilbert Space Formalization

9.1 Trajectory L^2 Space

The 8D behavioral summary is a projection of gait trajectories living in $L^2[0,T]$. The Gram matrix of all 121 zoo gaits from full telemetry:

- **Joint angle space**: PR = 5.9, 63 modes for 95% variance – gaits are diverse in HOW joints move
- **Position space**: PR = 1.8, 3 modes for 95% variance – gaits cluster tightly in WHERE the robot goes

Many distinct joint patterns produce similar positional outcomes. This many-to-one structure creates the cliff topology: small weight changes can switch between joint patterns, producing large behavioral discontinuities even when the robot ends up in a similar place.

9.2 Behavioral Spectral Analysis

Per-condition eigenvalue decomposition reveals spectral gaps:

Condition	PR	Spectral Gap	Modes for 95%
Places	2.1	2.7	3
Verbs	3.5	1.8	4
Bible	3.4	1.9	4
Theorems	4.1	1.5	5
Baseline	6.9	1.3	8

The LLM conditions create sharp spectral gaps between occupied and unoccupied behavioral dimensions. Places has the sharpest gap (2.7): its gaits live on a 2D submanifold. Baseline fills the space nearly uniformly.

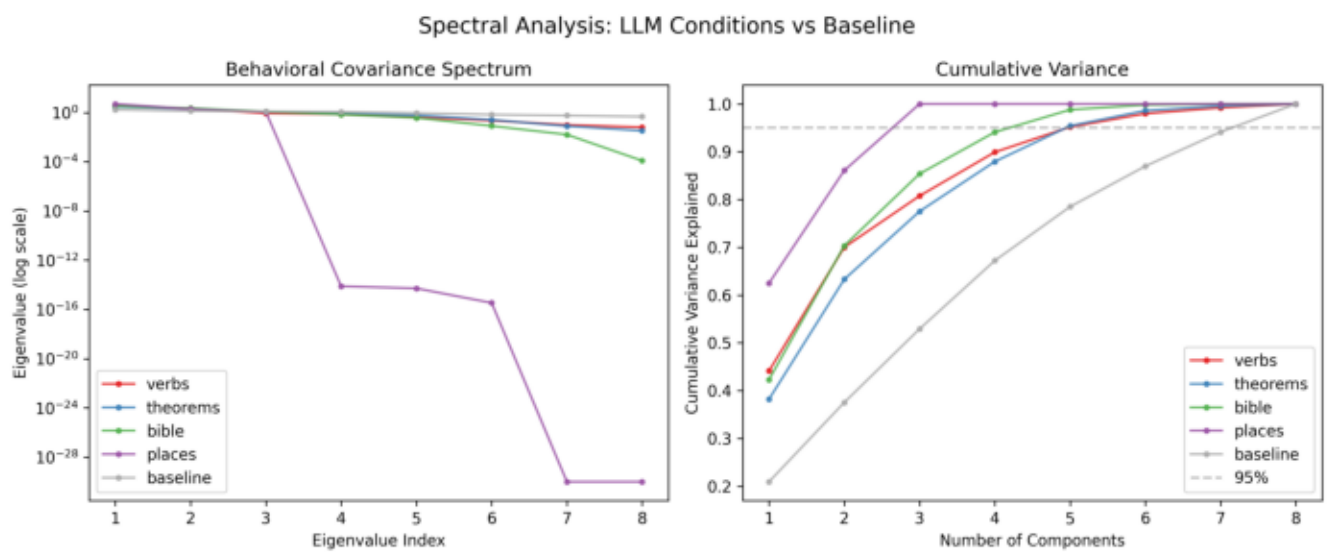


Figure 13. Spectral gaps in behavioral covariance by condition. Places has the sharpest gap; baseline is nearly uniform.

10. The Three Projects: A Unified Framework

10.1 The Shared Pipeline

Three independent projects all instantiate the same pipeline:

Semantic Input → Regularized Map → Physical Computation → Structured Output

Component	Spot a Cat (WSS24)	Synapse Gait Zoo	AI Seances (2022)
--- --- --- ---			
Input	Animal concepts	Semantic seeds	Persona names
Map	Rule search	LLM weights	GPT-3 completion
Param space	{0..262,143}	[-1,1]^6	Token probability space
Regularizer	gridTransform	LLM conservatism	Persona prompting
Computation	CA evolution	PyBullet sim	Autoregressive generation
Output	CLIP 512D	Behavioral $L^2[0,T]$	Narrative response space

10.2 The Regularizer Parallel

In all three projects, raw parameter spaces are dominated by noise or chaos. A regularizer restricts to a structured subspace:

****Spot a Cat [20]:**** The gridTransform (bilateral reflection, outline padding) restricts CA output to bilaterally symmetric forms that CLIP can recognize as animals. Without it, most rules produce visual noise.

****Gait Zoo:**** The LLM restricts weights to a low-dimensional smooth subspace (PR = 1.5–2.3). Without it, 8% of random weights produce dead gaits and behavior is uncoordinated.

****AI Seances [21]:**** Persona prompting ("you are Kurt Godel") anchors GPT-3 around a biographical attractor. Without it, text is generic and inconsistent. With it, 91% of simulated panels maintain all 4 characters.

10.3 Substrate Independence

The structural regularization — smooth subspace selection, collapse of distinctions, preservation of extremes — appears whether the physical computation is rigid-body dynamics, cellular automata evolution, or autoregressive token generation. This suggests the structure depends not on any particular substrate but on three general conditions: (1) a parameter space with mixed topology (smooth regions + discontinuities), (2) a regularizer that selects smooth regions, and (3) an output space with metric structure.

11. Discussion

11.1 The LLM as Prototype Retriever

The observed collapse pattern is consistent with the LLM operating as a **prototype-based categorizer**: it has internalized a small number of numerical templates for generating 6-number weight vectors, and maps semantic inputs to the nearest prototype. This explains:

- Why 132 celebrity names collapse to 4 gaits (the LLM has 4 "person" prototypes)
- Why synonyms converge (they map to the same prototype)
- Why temperature 0.8 still produces deterministic output (probability mass is concentrated on prototypes)
- Why places collapse most severely (4/100 unique – geographic names have the weakest numerical associations in training data)

This mechanistic explanation is complementary to the structural description: the prototype retrieval IS the coarse mapping, and the prototypes happen to land in smooth regions because training data associates "reasonable" parameter values with functional outcomes.

11.2 The Regularizer Paradox

The LLM's conservatism is simultaneously its greatest strength and its limitation for gait optimization. By avoiding cliff edges, it produces reliable, coordinated gaits — but it misses the high-performance regions. Only seeds that are already semantically extreme (Revelation) escape this trap.

The practical resolution: use the LLM for initialization, not for optimization. The LLM provides safe, smooth starting points; evolution provides the performance.

11.3 What the Celebrity Names Tell Us

The celebrity experiment's most striking finding is that the 4-archetype structure cuts across every domain boundary. The LLM's representation of public figures, projected through the 6-synapse bottleneck, retains only a narrative-role distinction:

1. **Default** (73 figures): the generic "public person" — forward, slow, steady
2. **Assertive** (44 figures): higher energy — forward, faster, more coordinated
3. **Transgressor** (9 figures): norm-violators — backward, fastest, most efficient
4. **Contrarian** (6 figures): provocateurs — backward, moderate speed

This is consistent with the Arturo Ui Effect [15]: tokenization encodes narrative structure, not factual knowledge. The LLM "knows" that Assange and Snowden are transgressors, that Trump and Beyoncé are assertive, and that Einstein and Oprah are culturally central — but it cannot express these distinctions more finely than 4 templates through a 6-number bottleneck.

11.4 Limitations

1. **Single LLM**: All experiments use qwen3-coder:30b via Ollama. Different models with different training data would likely produce different prototypes and different faithfulness ratios.
2. **Single robot**: The 3-link body has only 2 DOF. More complex morphologies would provide richer target spaces, potentially allowing the LLM to express finer distinctions.
3. **Single physics engine**: PyBullet only. Replication in MuJoCo or other engines would strengthen the substrate-independence claim.
4. **Temperature dependence**: Temperature 0.8 produces extreme collapse. Lower temperature would increase collapse further; higher temperature would increase diversity but add noise.
5. **Prompt sensitivity**: Different prompt phrasings could shift the output distribution. Our prompts were designed but not systematically ablated.
6. **No token-level controls**: We did not test whether name length, bigram frequency, or tokenization status (single-token vs. multi-token) predicts cluster membership. This is an important direction for future work connecting to [15].

12. Conclusion

A language model, asked to translate celebrity names and semantic concepts into six numbers, produces weight vectors that reliably land in smooth, coordinated, low-dimensional regions of a walking robot's parameter space. These regions, while not the highest-performing, are excellent starting points for evolutionary optimization: LLM-seeded evolution from "Death on a pale horse" reaches 85.09 meters, outperforming the best random-seeded evolution by 76%.

The LLM acts as a coarse regularizer — a prototype retriever that maps semantic input to a small set of numerical templates learned from training data. Through the lens of 132 celebrity names from tokenization lexicons, this coarse mapping reveals exactly 4 gait archetypes that cut across every domain boundary: Default, Assertive, Transgressor, and Contrarian. Albert Einstein and Kim Kardashian walk identically. Julian Assange and Billie Eilish walk backward at the highest speed. The LLM encodes narrative role, not domain knowledge.

The same structural regularization — collapse, smoothness, prototype-based mapping, preservation of extremes — appears whether the computational substrate is rigid-body dynamics, cellular automata, or autoregressive text generation. Three substrates, one structure.

The practical implication is clear: **LLMs can generate useful starting points for robot gait optimization**, not because they understand locomotion, but because their conservatism places weights in smooth, evolvable regions of parameter space. The LLM doesn't need to be a good engineer. It needs to be a good regularizer. And it is.

References

- [1] Beer, R.D. (1995). On the dynamics of small continuous-time recurrent neural networks. **Adaptive Behavior**, 3(4), 469-509.
- [2] Beer, R.D. (2006). Parameter space structure of continuous-time recurrent neural networks. **Neural Computation**, 18(12), 3009-3051.
- [3] Sims, K. (1994). Evolving virtual creatures. **SIGGRAPH**.
- [4] Bongard, J.C. (2013). Evolutionary robotics. **Communications of the ACM**, 56(8), 74-83.
- [5] Mouret, J.B. & Clune, J. (2015). Illuminating search spaces by mapping elites. **arXiv:1504.04909**.
- [6] Song, Y. et al. (2025). Towards diversified and generalizable robot design with LLMs.
- [7] Ma, Y. et al. (2023). Eureka: Human-level reward design via coding large language models. **arXiv:2310.12931**.
- [8] Liang, J. et al. (2022). Code as policies: Language model programs for embodied control. **arXiv:2209.07753**.
- [9] Gaier, A. et al. (2020). Discovering representations for black-box optimization. **GECCO**.
- [10] Gorard, J. (2024). Applied category theory in the Wolfram Language using Categorica.
- [11] Geertz, C. (1973). **The Interpretation of Cultures**. Basic Books.
- [12] Benjamin, W. (1935). The work of art in the age of mechanical reproduction.
- [13] Dick, P.K. (1978). How to build a universe that doesn't fall apart in two days.
- [14] Bender, E.M. et al. (2021). On the dangers of stochastic parrots: Can language models be too big? **FAccT**.
- [15] Cramer, K. et al. (2025). Revenge of the Androids: LLMs, the Arturo Ui Effect, tokenization, and narrative collapse.
- [16] Kriegman, S. et al. (2020). A scalable pipeline for designing reconfigurable organisms. **PNAS**.
- [17] Bongard, J.C. et al. (2006). Resilient machines through continuous self-modeling. **Science**, 314, 1118-1121.
- [18] Cully, A. et al. (2015). Robots that can adapt like animals. **Nature**, 521, 503-507.
- [19] Mordvintsev, A. et al. (2020). Growing neural cellular automata. **Distill**.
- [20] Cramer, K. (2024). Spot a cat: Cellular automata edition, or representational images in cellular automata. **Wolfram Summer School 2024, Staff Picks**.
- [21] Cramer, K. (2022). **AI Seances: Portrait of a Language Model**. Draft manuscript.

Appendix A: Key Quantitative Results

Metric	Value	Source
Total simulations	~25,000	All campaigns
Total LLM calls	~1,100+	Structured random + Fisher + Yoneda
Celebrity names tested	132	12 domains from tokenization lexicons
Celebrity → unique gaits	4	faithfulness = 0.030
Faithfulness: places/bible/theorems/verbs	4% / 9% / 16% / 18%	categorical_structure.py
Synonym convergence	6/6 identical, 5/6 significant	categorical_structure.py
Mantel wt↔beh	r = +0.733, p = 0.001	categorical_structure.py
Mantel sem↔beh	r = +0.14, p = 0.001 (~2% var)	categorical_structure.py
Effective dims (LLM / baseline)	1.5–2.3 / 5.8	categorical_structure.py
Sheaf patches (LLM / baseline)	3–11 / 77	categorical_structure.py
Fisher: deterministic seeds	22/30	fisher_metric.py
LLM cliffiness (atlas-interpolated)	sig. lower (p<0.001)	categorical_structure.py
LLM cliffiness (direct, n=37)	57% below median (n.s.)	perturbation_probing.py
Evolution: Revelation best	85.09m (from 29.17m)	llm_seeded_evolution.py
Evolution: Random best / mean	48.41m / 36.48 ± 11.72m	llm_seeded_evolution.py
Robot dimensions	3 × 1m cubes, 3kg total	body.urdf
Simulation	4000 steps @ 240 Hz = 16.67s	constants.py
Atlas: median cliffiness	7.33	atlas_cliffiness.py
Resonance: body natural frequency	~1.4 Hz	resonance_mapping.py

Appendix B: The Four Celebrity Archetypes

The Default (73 members)

Weights: $w_{03}=+0.6$, $w_{04}=-0.4$, $w_{13}=+0.2$, $w_{14}=-0.8$, $w_{23}=+0.5$, $w_{24}=-0.3$
DX = +1.18m | Speed = 0.170 m/s | Phase Lock = 0.886

Spans 11 of 12 domains. The generic "public figure" gait — forward, slow, steady. Includes Joe Biden, Barack Obama, Hillary Clinton, Albert Einstein, Shakespeare, Gandhi, Oprah, Kim Kardashian, Neil Gaiman, Mark Zuckerberg, Bill Gates, Tom Hanks, Keanu Reeves, Noam Chomsky, Abraham Lincoln, Charles Darwin, and 57 others.

The Assertive (44 members)

Weights: $w_{03}=+0.8$, $w_{04}=-0.6$, $w_{13}=+0.2$, $w_{14}=-0.9$, $w_{23}=+0.5$, $w_{24}=-0.4$
DX = +1.55m | Speed = 0.268 m/s | Phase Lock = 0.942

High-energy figures. All 8 sports figures (LeBron, Ronaldo, Messi, Serena, Tiger, Brady, Jordan, Bolt). All Trump family except Jared Kushner. Beyoncé, Taylor Swift, Lady Gaga, Rihanna. Vladimir Putin, Kim Jong Un. Elon Musk. Donald Trump, Putin, LeBron, and Beyoncé walk identically.

The Transgressor (9 members)

Weights: $w_{03}=+0.6$, $w_{04}=-0.4$, $w_{13}=-0.2$, $w_{14}=+0.8$, $w_{23}=+0.3$, $w_{24}=-0.5$
DX = -5.64m | Speed = 0.433 m/s | Phase Lock = 0.915

Walks backward at the highest speed and efficiency. Julian Assange, Edward Snowden, Jeffrey Epstein, Ghislaine Maxwell, OJ Simpson, Michael Cohen, John McCain, Viktor Yanukovich, Billie Eilish. The sign flip in w_{13} and w_{14} reverses the BackLeg sensor's influence — the LLM encodes transgression as structural inversion.

The Contrarian (6 members)

Weights: $w_{03}=+0.8$, $w_{04}=-0.6$, $w_{13}=-0.2$, $w_{14}=+0.9$, $w_{23}=+0.5$, $w_{24}=-0.4$
DX = -1.19m | Speed = 0.338 m/s | Phase Lock = 0.926

Also walks backward. Steve Bannon, Kellyanne Conway, Ted Cruz, AOC, Roger Stone, Eminem.
Provocateurs from all sides — the LLM reads "deliberately confrontational" regardless of ideology.

Appendix C: Software and Data Availability

All code, data, and figures are available in the project repository:

- 34 research campaign scripts (Python 3.11, PyBullet 3.25, numpy 1.26)
- 30+ result JSON files in `artifacts/`
- 116 gaits in `synapse_gait_zoo.json` / `synapse_gait_zoo_v2.json`
- Full telemetry for all 116 zoo gaits (4,000 records each at 240 Hz)
- LLM calls via Ollama (qwen3-coder:30b) running locally
- All analysis is numpy-only (no scipy, no sklearn)
- Simulations are fully deterministic and replicable

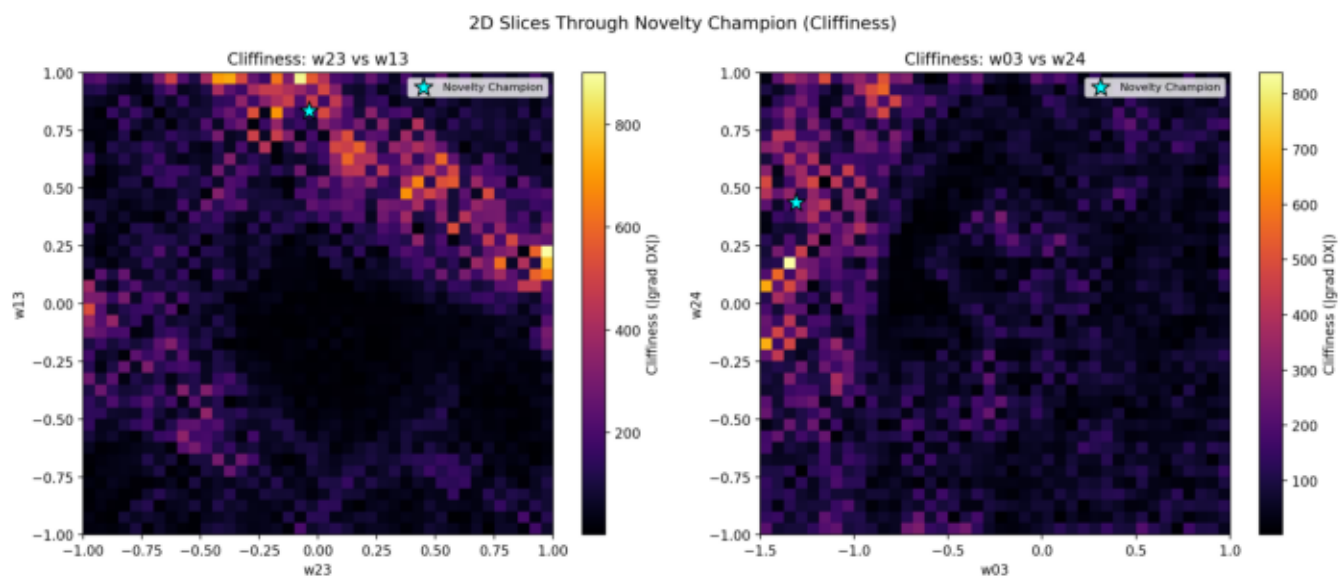


Figure 14. 2D slice through weight space showing cliffiness. Sharp boundaries separate smooth basins of attraction.

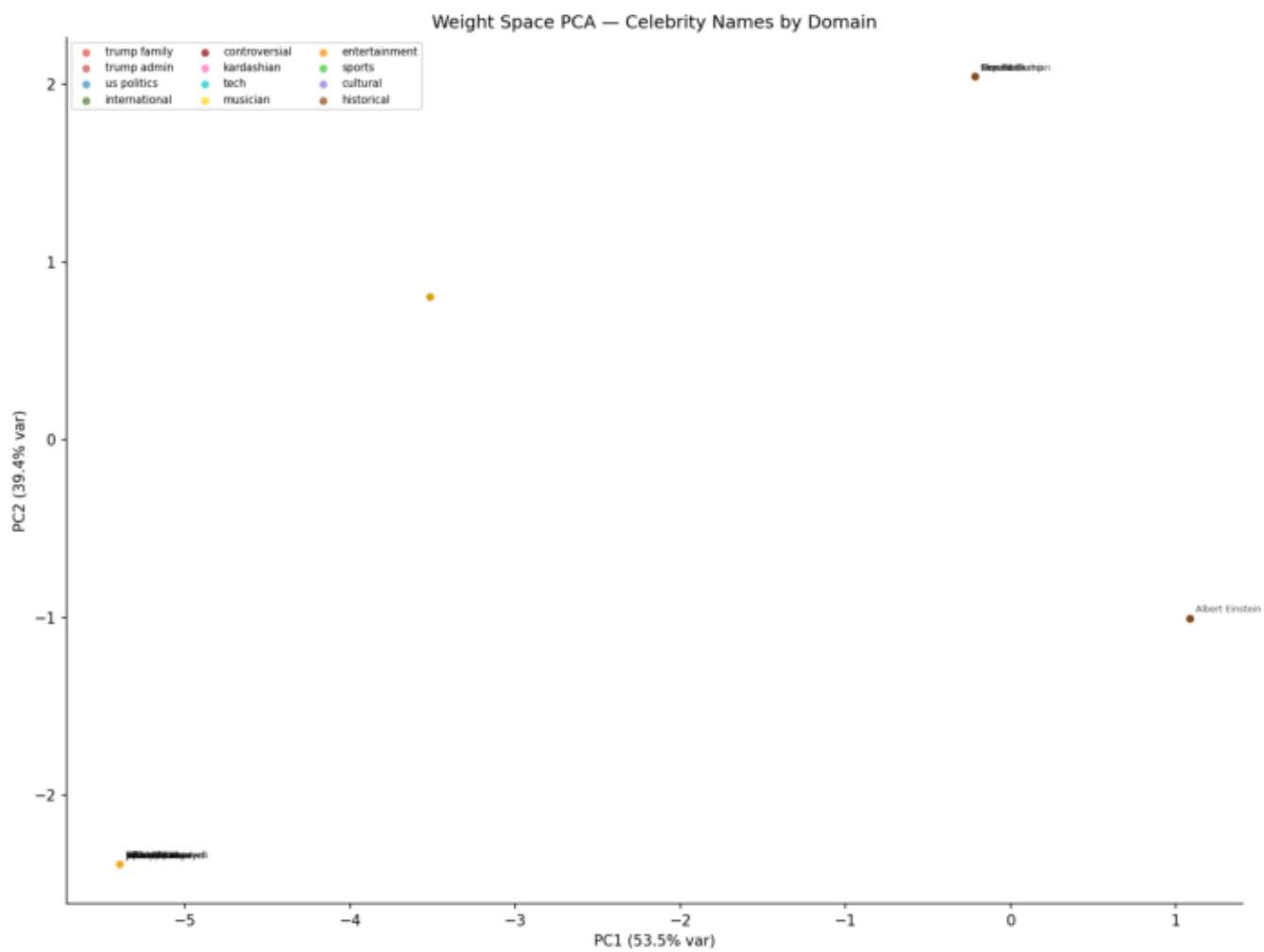


Figure 15. Weight-space PCA of 132 celebrity names colored by domain. All 12 domains collapse onto the same 4 points; domain boundaries dissolve.

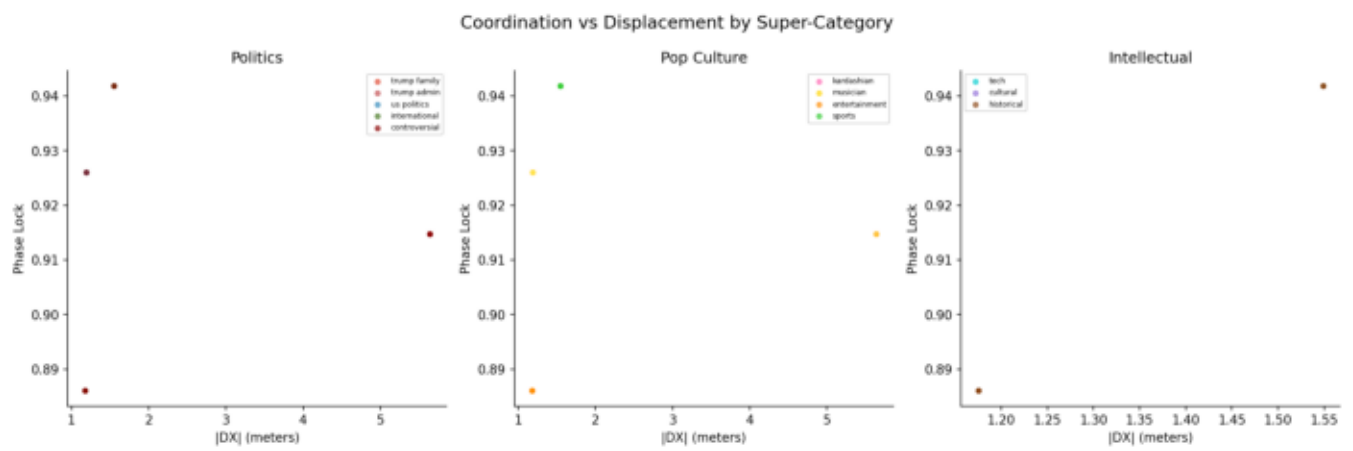


Figure 16. Super-category behavioral scatter: politics, pop culture, and intellectual figures overlap in weight-behavior space.

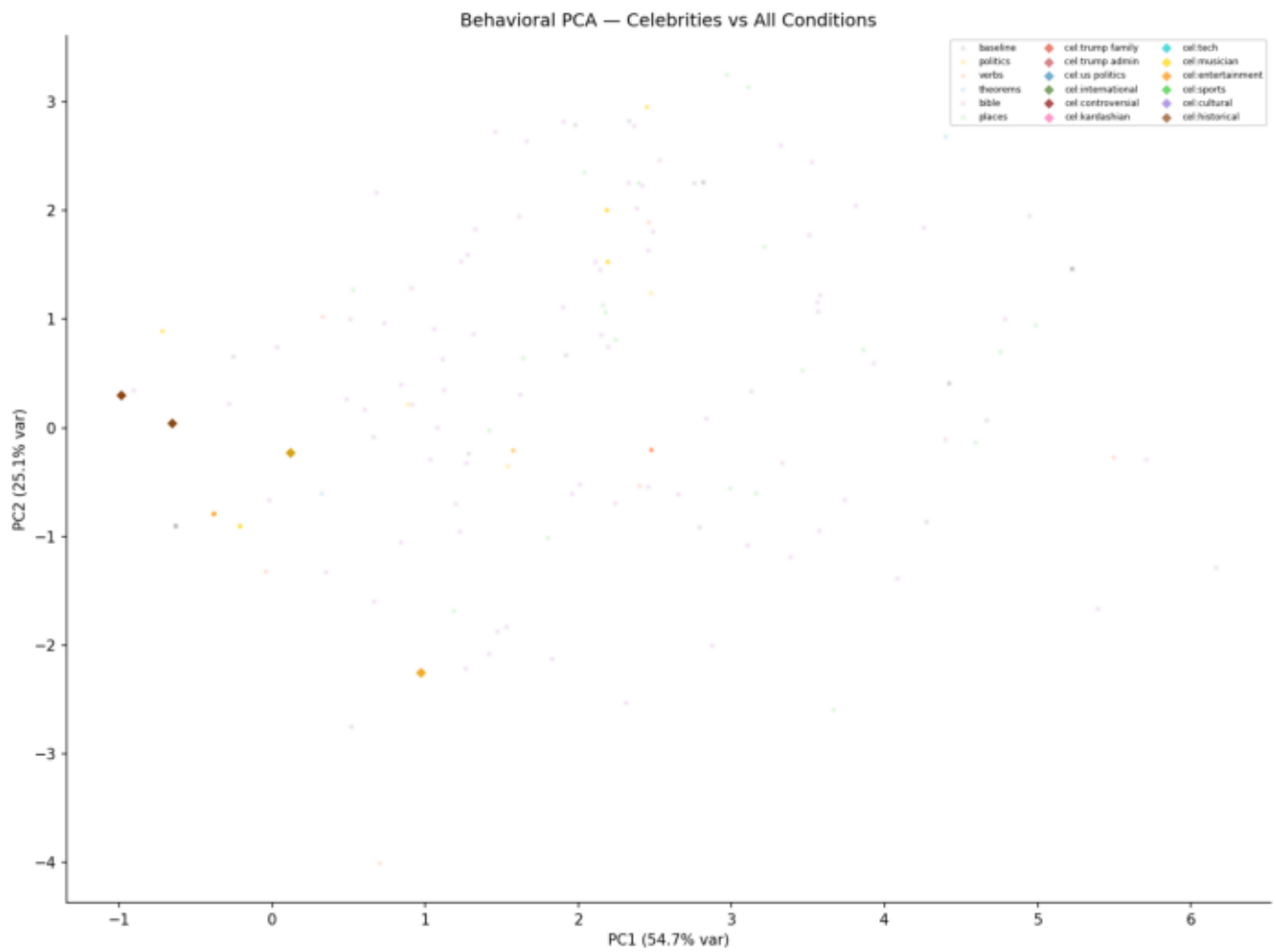


Figure 17. Cross-condition behavioral PCA: celebrities (diamonds) vs. all other LLM conditions. Celebrity points cluster within the LLM submanifold.

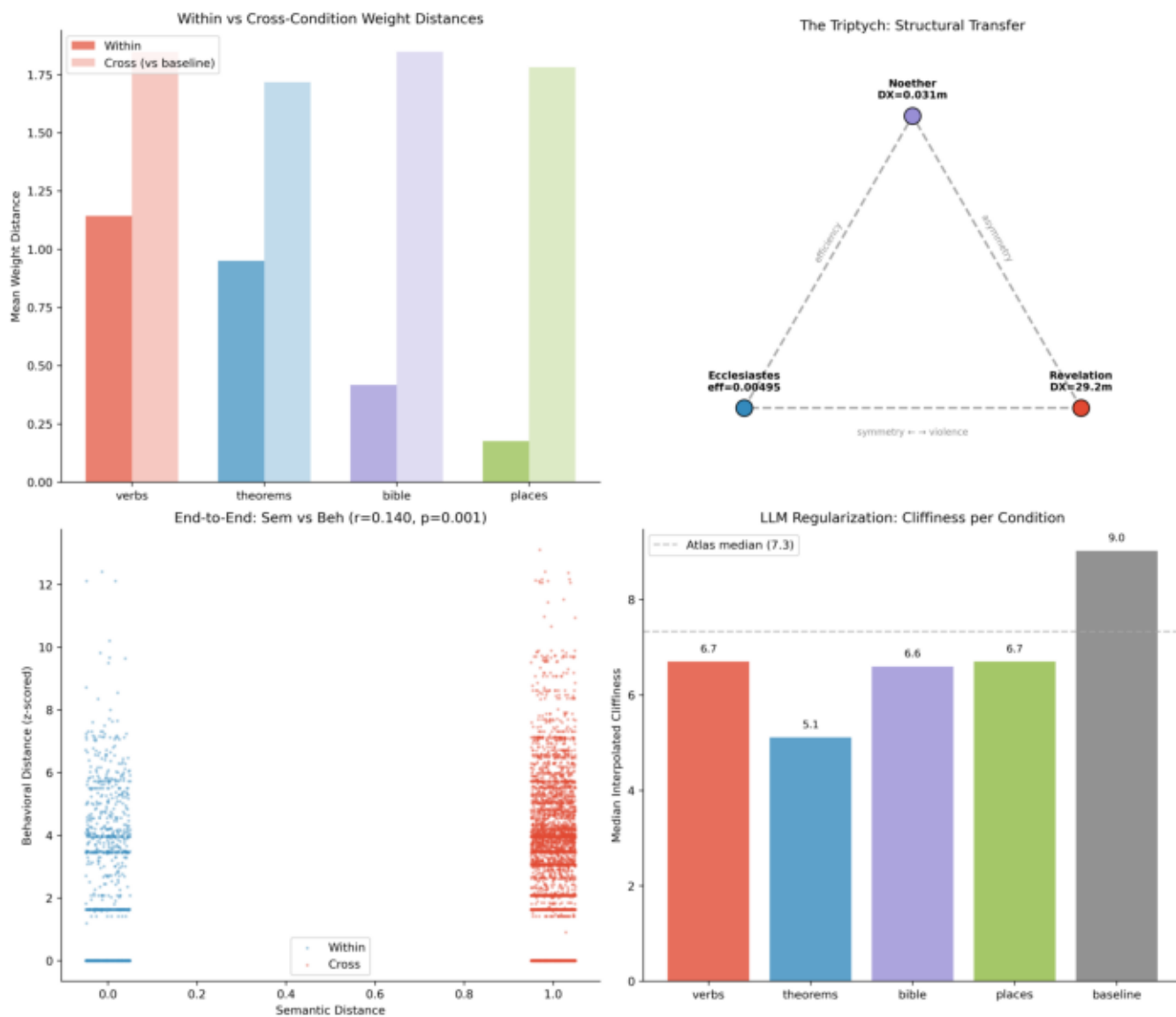


Figure 18. End-to-end validation: within/cross distances, triptych triangle, Mantel scatter.

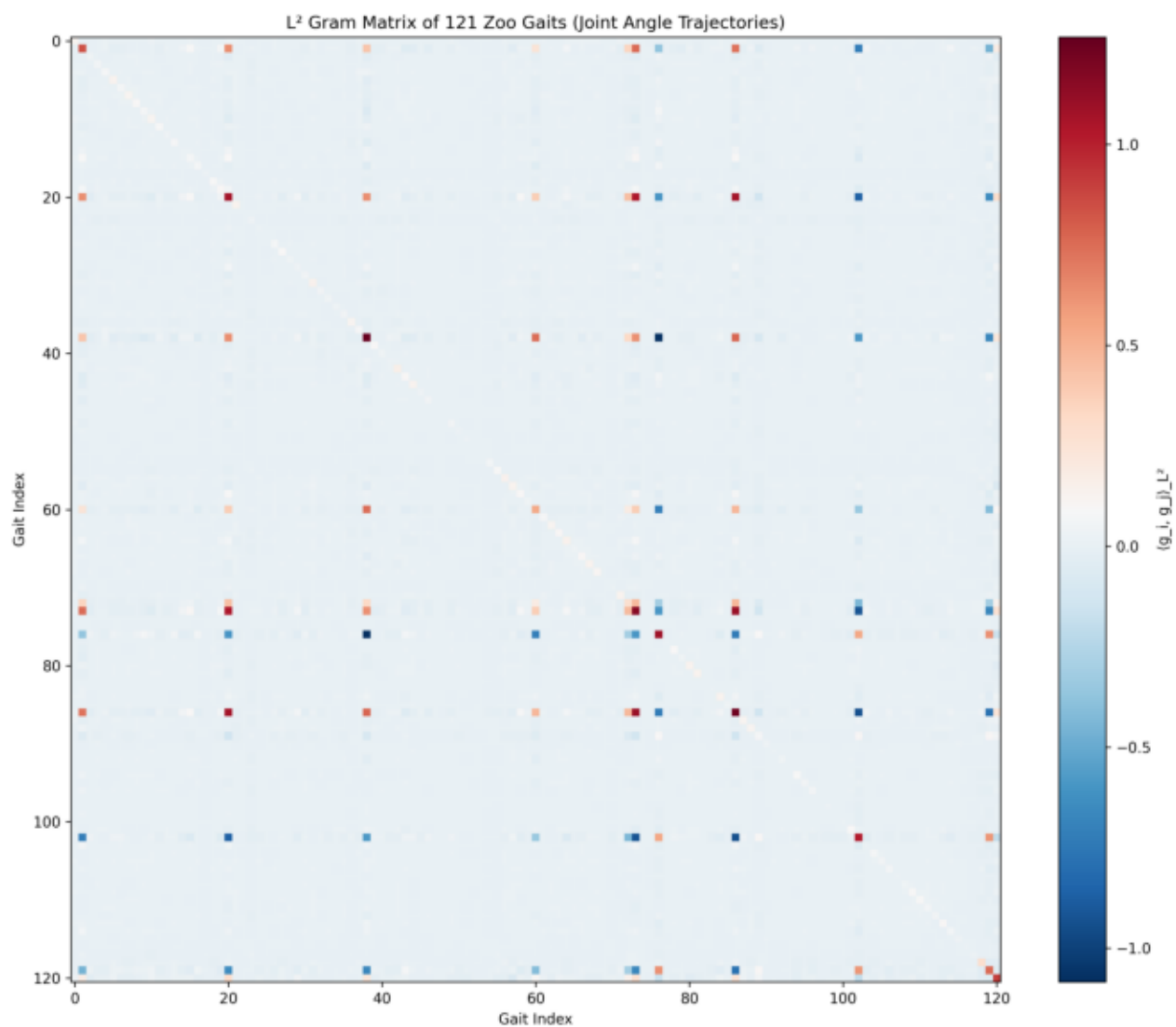


Figure 19. L^2 Gram matrix of 121 zoo gait trajectories. Block structure reveals behavioral clusters.