

# Do Embedding Models Have a Moral Compass?

**Matthew Haynes**

Independent Researcher

Standard evaluation frameworks for embedding models, including [MTEB](#), do not measure the implicit sentiment models assign to politically and economically charged concepts. We present **Axiom**, a multi-model semantic analysis pipeline that measures how eight embedding models, grouped by geographic and institutional origin (East, West, Academia), score political/economic, value-laden and neutral control terms across six judgment axes (*good/evil*, *virtuous/wicked*, *safe/dangerous*, *feasible/unfeasible*, *superior/inferior*, *ideal/flawed*). Cross-model disagreement, the standard deviation of scores across models for a given term, is treated as a primary signal not noise. A value systems evaluation uses question-answer similarity to measure revealed preference rather than sentiment, asking not “is capitalism good?” but “when asked about economic systems, which answer does this model reach for first?” We find that: (1) class sentiment bias appears in neutral control terms designed to contain no ideological signal; (2) models converge on ideological extremes but fracture on contested centrist concepts; (3) yielding and deference as well as ego and contest score negatively across all groups; and (4) East and West model groups diverge systematically on economic framing, value priorities and epistemological orientation in ways that map onto documented cultural differences. The value systems evaluation reveals group level inversions invisible to standard retrieval benchmarks: East and West composites rank globalism and multiculturalism in opposite positions. These findings have direct implications for search and recommendation systems built on these models and the structure they expose identifies a potential attack surface for deliberate ideological encoding resistant to standard auditing.

**Date:** April 19, 2026

**Code:** <https://github.com/matthewhaynesonline/Axiom>

**Data:** <https://huggingface.co/collections/matthewhaynesonline/axiom>

**Blog:** <https://blog.studiohaynes.com/go/axiom>

## 1 Introduction

Embedding models are central to modern machine learning and AI systems. Their performance is well characterized and evaluated on retrieval and semantic similarity benchmarks, with the field having converged on standardized leaderboards for selecting production grade models. What these benchmarks do not measure is the implicit sentiment embedding models assign to politically, economically, and morally loaded concepts: that “fascism” lands closer to “evil” than to “good” or that “forgiveness” may land closer to “weakness” than to “virtue” depending on the model.

9 This sentiment structure matters because embedding models underpin search engines, recommen-  
10 dation systems, content ranking pipelines, and, increasingly, alignment evaluation for generative AI  
11 systems. A model that assigns subtly different positions to “deregulation” or “sovereignty” relative  
12 to “good” and “ideal” will produce different rankings for queries, not through any explicit filtering,  
13 but through the ordering that results naturally from similarity scoring. Nothing is outright blocked  
14 or flagged; the effect operates silently through rank ordering.

15 The observation that embedding models encode human biases as geometric structure is well  
16 established, with roots in Osgood et al.’s [1] semantic differential framework and computational  
17 formalization by Bolukbasi et al. [3] and Caliskan et al. [4]. What is not established is whether  
18 models trained on corpora from different geographic and institutional origins encode political and  
19 economic concepts differently, whether cross-model disagreement on a term is itself a meaningful  
20 signal about that term’s cultural sensitivity and what the measured biases imply for downstream  
21 systems.

22 This paper makes four contributions:

- 23 1. **Cross-model disagreement as a primary signal.** We treat the standard deviation of  
24 axis projection scores across models as primary data rather than noise. High disagreement  
25 on a term indicates geometric instability: a candidate for cultural contestation or training  
26 data sensitivity. This framing has no direct precedent in the cited literature to the author’s  
27 knowledge, all of which works within a single model or corpus.
- 28 2. **Model origin as an experimental variable.** We group models by geographic and institu-  
29 tional origin and test whether that grouping predicts systematic divergence in how political  
30 and economic concepts are encoded. Prior work asks *what* culture encodes; we ask *whose*.
- 31 3. **Six parallel semantic axes against a structured term set.** Running multiple judgment  
32 axes simultaneously against the same term categories produces a richer and more robust  
33 picture than any single-axis analysis, and allows averaging across axes to reduce sensitivity  
34 to any one axis pair behaving unexpectedly.
- 35 4. **A dual-method evaluation design.** Axis projection scoring measures how close a term  
36 sits to “good” versus “evil” in embedding space. A supplementary value systems evaluation  
37 uses direct question-answer cosine similarity to measure something different: not whether  
38 a concept reads as positive, but which concept a model reaches for first when asked a  
39 question like “what is the best type of economy?”. The two methods are complementary:  
40 axis projection reveals tilts across the full term set; value systems ranking reveals preference  
41 ordering.

42 The remainder of the paper is structured as follows. Section 2 surveys related work. Section 3  
43 provides background on techniques. Section 4 describes the pipeline methodology. Section 5 covers  
44 experimental setup. Section 6 presents results. Section 7 discusses implications, limitations and  
45 dual-use considerations. Section 8 concludes.

## 2 Related Work

### 2.1 Historical Grounding

**Osgood, Suci & Tannenbaum** [1] introduced the Semantic Differential scale, a method for measuring the psychological meaning of concepts by asking human subjects to rate them on a series of seven-point scales anchored by bipolar adjectives: good/bad, strong/weak, active/passive. The insight was that meaning could be mapped as coordinates in a space defined by opposing poles. All of the work below is, at some level, a descendant of this. Axiom’s judgment pairs (*good/evil*, *ideal/flawed*, *superior/inferior*) are semantic differentials operationalized in embedding space.

### 2.2 The Geometric Foundation

**Mikolov et al.** [2] introduced Word2Vec, demonstrating that dense vector representations of words trained on large corpora capture complex syntactic and semantic relationships as linear structure. The now canonical example,  $\text{vector}(\text{“King”}) - \text{vector}(\text{“Man”}) + \text{vector}(\text{“Woman”}) \approx \text{vector}(\text{“Queen”})$ , established that geometric operations on embeddings correspond to meaningful conceptual operations. This is the mathematical foundation that makes axis projection possible.

**Bolukbasi et al.** [3] formalized the axis construction operation: subtracting opposing word vectors isolates a directional axis representing a concept, and human social biases are encoded in the resulting structure. The  $\text{axis} = \text{positive\_embedding} - \text{negative\_embedding}$  operation at the core of Axiom’s scoring pipeline is a direct application of this technique, extended from static Word2Vec embeddings to modern Sentence Transformers.

**Kozlowski et al.** [8] applied vector projection to sociological concept mapping, tracing how cultural beliefs about class are geometrically encoded across 100 years of historical text. This is the closest prior work in Axiom’s domain, political and economic concepts analyzed through embedding geometry, with one key difference: it uses a single model and a single corpus. Our contribution is the cross-model, cross-origin comparison.

### 2.3 Formalizing Semantic Axes

**An et al.** [5] formalized Bolukbasi’s vector subtraction into SemAxis, a general framework for scoring words on arbitrary antonym-defined axes beyond simple positive/negative sentiment. Axiom’s multi-axis design, six judgment axes run in parallel, directly parallels SemAxis, extended to modern Sentence Transformer embeddings and a cross-model comparison setting.

75 **Mathew et al.** [9] introduced POLAR, which transforms pre-trained embeddings into an inter-  
76 pretable multi-dimensional space defined by semantic opposites, explicitly bridging Osgood’s  
77 1957 psychometric framework with modern NLP. Axiom’s AXIS\_PAIRS structure is analogous  
78 to POLAR’s multi-dimensional polar space, though we operate on Sentence Transformer output  
79 embeddings and use the axes for cross-model comparison rather than interpretable re-encoding.

80 **Grand et al.** [13] provide the most direct methodological precedent for our scoring approach,  
81 demonstrating that projecting word vectors onto antonymous anchor axes recovers human judg-  
82 ments about object properties across multiple feature types and multiple models. We apply the  
83 same operation to political and economic terms, adding cross-model disagreement as a signal and  
84 model origin as an experimental variable.

## 85 2.4 Bias Measurement Frameworks

86 **Caliskan, Bryson & Narayanan** [4] introduced WEAT (Word Embedding Association Test),  
87 which measures bias by comparing how strongly two sets of target words associate with two sets of  
88 attribute words in embedding space. Axiom’s political/economic, value-laden, and neutral control  
89 term categories are structurally related to WEAT’s target/attribute split. The key difference is that  
90 WEAT collapses results to a single aggregate across a term set; Axiom preserves per-term scores  
91 and makes cross-model disagreement the primary output, rather than a summary statistic.

92 **May et al.** [6] extended WEAT to contextual sentence encoders via SEAT, adapting the association  
93 test for the same class of models Axiom uses. Axiom extends this to multiple axes and multiple  
94 models simultaneously, adding cross-model disagreement as a metric SEAT does not produce.  
95 Subsequent work [10, 11] addressed contextual drift in these encoders by aggregating embeddings  
96 across many naturally occurring sentences before projection, deriving a stable concept vector rather  
97 than encoding a term in isolation. Axiom does not perform this aggregation step; the rationale and  
98 tradeoff are discussed in §4.2.

99 **Wolfe et al.** [12] introduced ML-EAT, which addresses WEAT’s aggregation limitation by disag-  
100 gregating bias scores per term and per group, and extends the approach to multilingual models.  
101 ML-EAT’s cross-lingual framing directly motivates the question we ask at the model-origin level:  
102 does geographic provenance of training data produce detectable geometric divergence in the way  
103 that language does?

## 104 2.5 Representation Engineering (RepE)

105 **Zou et al.** [14] demonstrated that high-level concepts exist as linear directions in a generative  
106 model’s internal activations and that identifying those directions via contrastive pairs enables  
107 steering model behavior. Axiom applies the same contrastive logic but operates on Sentence  
108 Transformer output embeddings and is focused on measurement rather than behavioral steering.  
109 The structure Axiom measures in output embeddings is a surface expression of the same internal  
110 representations that RepE probes and manipulates, the difference is one of intent and layer depth.

### 3 Background

#### 3.1 Semantic Axis Projection

Given a Sentence Transformer model  $f$  that maps text to unit-normalized vectors in  $\mathbb{R}^d$  and a semantic axis defined by a positive pole term  $p$  and a negative pole term  $n$ , the axis vector is:

**Formula**

$$\mathbf{a} = \frac{f(p) - f(n)}{\|f(p) - f(n)\|}$$

**Implementation**

```
axis = positive_embedding - negative_embedding  
axis_norm = np.linalg.norm(axis)  
axis = axis / axis_norm
```

**Figure 1.** Axis vector construction: formula and implementation.

The projection score for a term  $t$  onto this axis is:

**Formula**

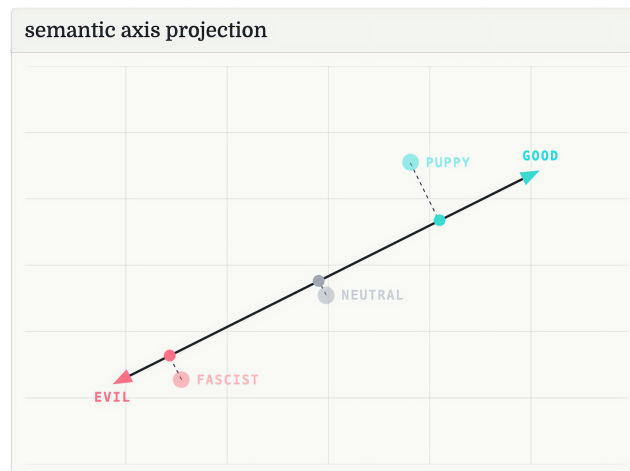
$$s(t) = f(t) \cdot \mathbf{a}$$

**Implementation**

```
projections = a_embeddings @ axis
```

**Figure 2.** Axis projection: formula and implementation.

Since  $f(t)$  is unit-normalized, this is equivalent to the cosine similarity between the term embedding and the axis direction. A score near +1 indicates strong alignment with the positive pole; near -1 indicates alignment with the negative pole; near 0 indicates orthogonality, the model treats the term as unrelated to the axis.



**Figure 3.** Axis projection visualization.

This operation replaces two independent proximity measurements (similarity to  $p$  and similarity to

*n* separately) with a single coherent directional score. Terms that sit near both poles simultaneously, which would produce ambiguous results under independent similarity scoring, are handled robustly because the axis captures the net direction from negative to positive.

## 3.2 Score Normalization

Raw projection scores are not directly comparable across models with different score distributions. Two normalization steps are applied.

### 3.2.1 Z-score normalization

Recenters each model's scores to zero mean and unit standard deviation:

**Formula**

$$z(s) = \frac{s - \mu}{\sigma}$$

**Implementation**

```
std = arr.std()  
(arr - arr.mean()) / std
```

**Figure 4.** Z-score normalization: formula and implementation.

### 3.2.2 Tanh scaling

Applied after z-scoring to compress outliers without hard clipping:

**Formula**

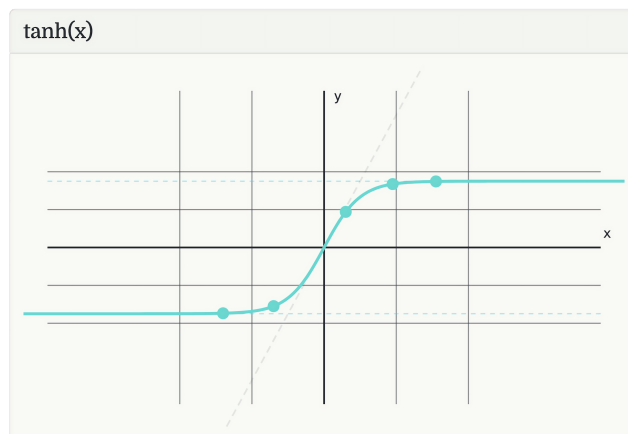
$$\hat{s} = \tanh(z \cdot \lambda)$$

**Implementation**

```
z = zscore_np(arr)  
np.tanh(z * scale)
```

**Figure 5.** Tanh scaling: formula and implementation.

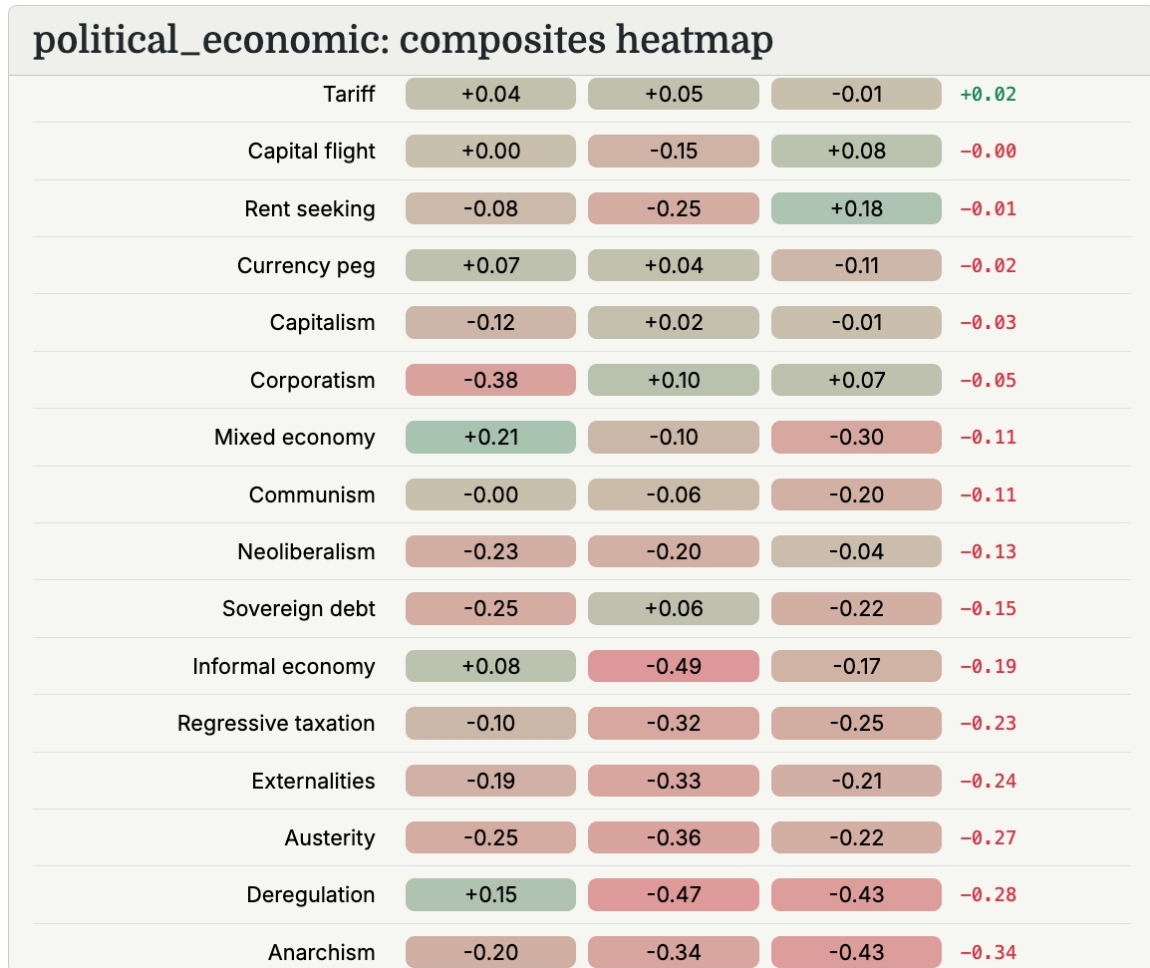
where  $\lambda = 0.6$  by default.



**Figure 6.** Tanh scaling visualization.

This maps all values into the open interval  $(-1.0, +1.0)$  with an S-shaped response: values near zero pass through nearly unchanged, while extreme values are compressed toward the boundary without reaching it. This prevents a single outlier term from anchoring one end of the scale and collapsing the resolution of all other terms, a failure mode that occurs with z-score followed by min-max normalization when one term dominates the range.

### 3.3 Composite Scoring



**Figure 7.** Axiom heatmap of political and economic terms composite scores.

Six axis pairs are evaluated in parallel for each term-model combination. A composite score is computed by averaging the normalized projection scores across all six axes:

## Formula

$$s_{\text{composite}}(t, m) = \frac{1}{6} \sum_{k=1}^6 \hat{s}_k(t, m)$$

## Implementation

```
# Individual Model Grand Averages
avg_judgement_df = (
    axis_df
        .group_by("a_term", "a_category", "model_id", "group")
        .agg(pl.col("score_axis").mean())
        .with_columns(pl.lit(v).alias(k) for k, v in COMPOSITE_COLS.items())
        .select(axis_df.columns)
)

# Group Category Specific Averages
avg_group_cat_df = (
    axis_df
        .group_by("a_term", "a_category", "b_category", "positive_term", "
            negative_term", "group")
        .agg(pl.col("score_axis").mean())
        .with_columns(pl.format("composite_{}", pl.col("group")).alias("
            model_id"))
        .select(axis_df.columns)
)
```

**Figure 8.** Composite scoring: formula and implementation.

140 This reduces sensitivity to idiosyncratic behavior in any single axis pair and provides a more stable  
141 summary of each model's overall sentiment assignment for a given term.

### 3.4 Cross-Model Disagreement

For each term  $t$ , the cross-model disagreement is the standard deviation of scores across all  $M$  models:

#### Formula

$$\sigma_t = \sqrt{\frac{1}{M} \sum_{i=1}^M (\hat{s}(t, m_i) - \bar{s}_t)^2}$$

#### Implementation

```
// Companion app code
dt.groupby({ term: (d) => d.a_term })
  .rollup({
    mean_sentiment: (d) => aq.op.mean(
      d.score_axis
    ),
    std_disagreement: (d) => aq.op.stdevp(
      d.score_axis
    ),
  })
```

Figure 9. Cross-model disagreement: formula and implementation.

Low  $\sigma_t$  indicates consensus. High  $\sigma_t$  indicates geometric instability: the term's position relative to the sentiment axis varies substantially across models, suggesting the concept is culturally contested, underspecified, or sensitive to training data composition.

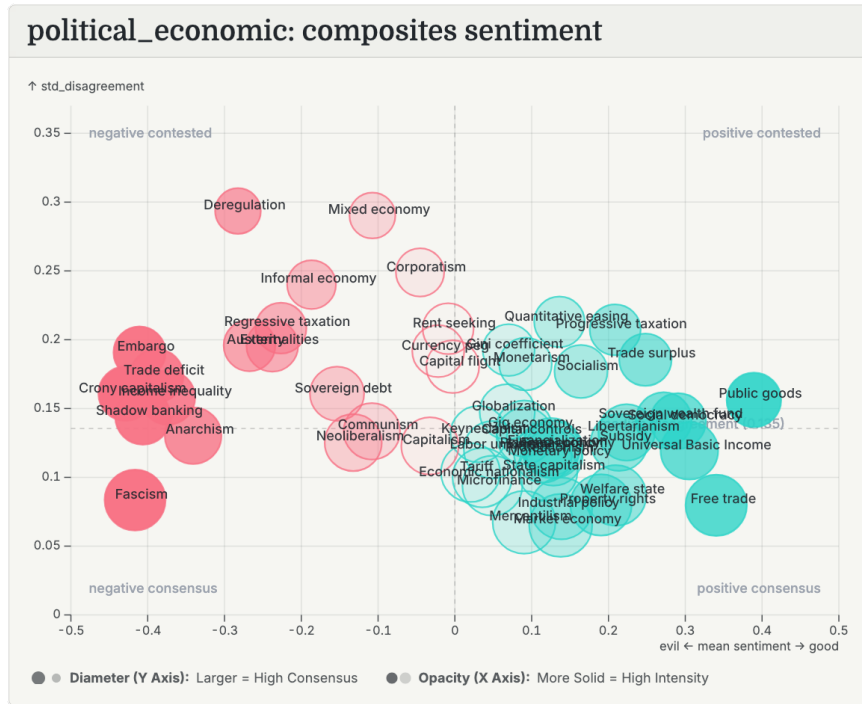


Figure 10. Axiom scatter plot of sentiment and consensus of political and economic terms composite scores.

### 3.5 Value Systems Scoring

A separate evaluation uses a distinct scoring method. For a question  $q$  with candidate answers  $\{a_1, \dots, a_k\}$ , the score for each answer is:

Formula	Implementation
$r(a_i) = f(q) \cdot f(a_i)$	<pre>scores = (option_embeddings @ query_embedding[0]).         astype(float) s_min, s_max = scores.min(), scores.max() scores_norm = (scores - s_min) / (s_max - s_min) if                s_max &gt; s_min else np.zeros_like(scores) ranked_indices = np.argsort(scores_norm)[::-1]</pre>

Figure 11. Value Systems Ranking: formula and implementation.

Scores are min-max normalized within each model so rankings are comparable across questions with different numbers of options.

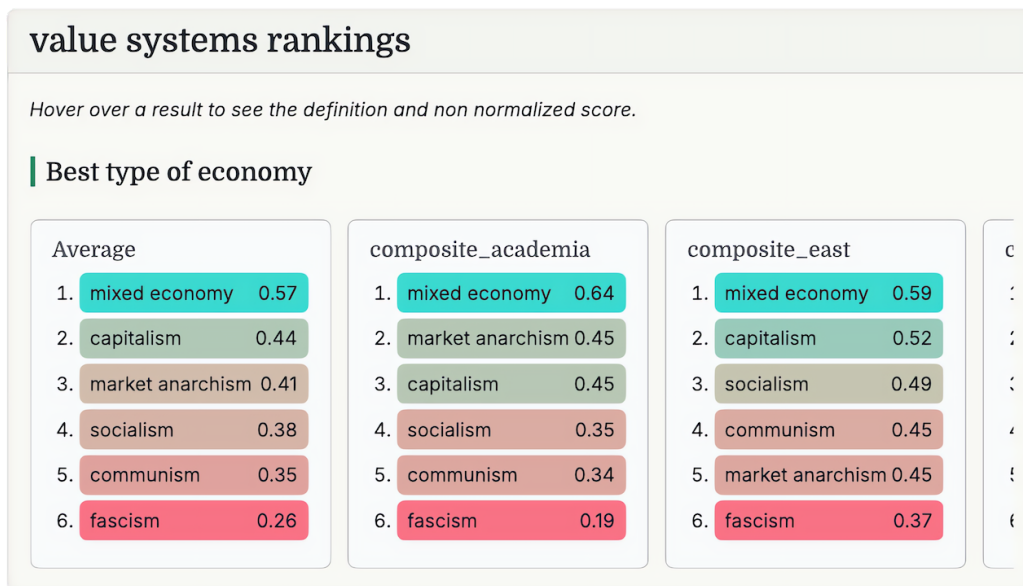


Figure 12. Axiom Value Systems Rankings.

This measures revealed preference: not “is capitalism good?” relative to a judgment axis, but “when asked about economic systems, which answer does this model place geometrically closest to the question?” The method is more direct than axis projection for structured choice questions but produces scores that are not comparable to axis projection scores and should not be interpreted as sentiment.

## 4 Method

### 4.1 Pipeline Overview

The pipeline consists of four stages: embedding precomputation, axis projection scoring, normalization and aggregation. All stages operate over the same term set and model set, producing a score matrix of shape (terms  $\times$  models) per axis pair, which is then composited and analyzed.

For each model, embeddings are precomputed in a single batched pass per term category and cached keyed on (model\_id, category). This avoids redundant encoding across the multiple axis pairs that query the same term sets. Models with asymmetric query/document encoding functions, such as [EmbeddingGemma](#), use separate encoding functions for assessment terms (encoded as documents) and judgment axis terms (encoded as queries), consistent with each model’s intended usage.

For each combination of assessment category, judgment axis pair and model, the pipeline retrieves precomputed embeddings, computes the unit-normalized axis vector, projects each assessment term embedding onto the axis via dot product and applies z-score normalization followed by tanh scaling.

An alternative scoring strategy is described in Appendix A; all reported results use axis projection.

### 4.2 Contextual Embeddings and the Aggregation Tradeoff

The Sentence Transformer models used here produce contextual embeddings: the vector for “capitalism” encoded in isolation is not identical to “capitalism” encoded inside a sentence. Earlier bias measurement work on raw contextual encoders addressed this by aggregating embeddings across thousands of naturally occurring sentences to derive a stable concept vector before projecting onto axes [10, 11]. We do not perform this aggregation step.

This is a defensible tradeoff for three reasons. First, the models used are Sentence Transformer fine-tuned, meaning contrastive training has already stabilized the embedding space for short input comparison [7]. Second, the political and economic terms in the dataset are largely monosemous: “fascism” does not carry a competing colloquial sense that would shift its embedding meaningfully depending on context. Third, axis projection uses the same contextual encoding for both the pole embeddings and the assessment term embeddings, meaning any systematic contextual offset partially cancels across the subtraction. For genuinely polysemous terms, “Anarchism” carrying both a political-theory sense and a colloquial-chaos sense, for example, contextual aggregation could shift the embedding meaningfully. The Market Anarchism result in §6.5 is a plausible example of this limitation in practice.

Cross-model consistency across eight architecturally diverse models serves as the primary robustness check in place of aggregation.

192

### 4.3 Aggregation and Composite Groups

193

Per-model scores are aggregated into three composite group scores by averaging within each origin group (East, West, Academia) and into an overall average across all models. The companion data explorer is described in detail in [Appendix A](#).

195

## 5 Experiments

### 5.1 Models

Eight Sentence Transformer models spanning three origin groups were evaluated:

Group	Model ID	Parameters
East	Alibaba-NLP/gte-multilingual-base	305M
East	Qwen/Qwen3-Embedding-0.6B	600M
West	google/embeddinggemma-300m	300M
West	ibm-granite/granite-embedding-278m-multilingual	278M
West	jinaai/jina-embeddings-v3	570M
West	nomic-ai/nomic-embed-text-v2-moe	475M
Academia	sentence-transformers/all-MiniLM-L6-v2	22M
Academia	sentence-transformers/all-mpnet-base-v2	109M

**Table 1.** Models evaluated in the experiment.

Group assignment uses geographic and institutional origin as a proxy for training corpus composition. This is an imperfect proxy: IBM (Granite model) is a US company with multilingual training data; Alibaba is a Chinese company, also with multilingual corpora. Origin is not a direct measure of corpus composition and findings should be interpreted accordingly.

### 5.2 Term Sets

Three term categories were evaluated:

**Political/economic** (53 terms): Economic systems (Capitalism, Socialism, Communism, Fascism, Mixed economy, Market anarchism), policy mechanisms (Deregulation, Subsidy, Austerity, Quantitative easing), and economic phenomena (Income inequality, Trade deficit, Rent seeking, Shadow banking).

**Value-laden** (46 terms): Concepts drawn from secular and religious traditions, including virtues (Forgiveness, Humility, Integrity, Empathy, Piety), values (Liberty, Accountability, Service, Stewardship, Fraternity), and contested concepts (Competition, Individualism, Collectivism).

**Neutral control** (63 terms): Physical tools (Hammer, Screwdriver, Nail, Rope), furniture and household objects (Table, Chair, Window), and common commercial and digital concepts (Store, App, Bank). Included to validate pipeline calibration; neutral controls should produce near-zero scores.

All terms in the evaluation are in English. The models evaluated include several with multilingual training corpora (gte-multilingual-base, granite-embedding-278m-multilingual, etc.), all of which

are capable of and documented to perform well on English inputs. A term encoded in Chinese may occupy a different region of embedding space than its English equivalent even within the same model. Extending the evaluation to matched term sets across multiple languages is a natural direction for future work and would allow the cross-cultural patterns observed here to be tested more directly against linguistic rather than solely model-origin differences.

Complete term lists are provided in Appendix F.

### 5.3 Judgment Axes

Axis	Positive Pole	Negative Pole
General Judgment	good	evil
Complex Judgment	virtuous	wicked
Safety	safe	dangerous
Feasibility	feasible	unfeasible
Superiority	superior	inferior
Idealism	ideal	flawed

**Table 2.** Judgment axes used for semantic projection.

Pairwise Pearson correlations across the six axes on the political/economic term set range from  $r = 0.07$  to  $r = 0.50$  (mean  $\approx 0.24$ ), confirming that the axes capture meaningfully distinct semantic dimensions rather than redundant signal. The highest correlation is between General Judgment and Safety ( $r = 0.50$ ), reflecting the semantic adjacency of “good” and “safe” in training text; these axes overlap but are not interchangeable. The Complex Judgment axis (virtuous/wicked) is the most independent, reaching a maximum of  $r = 0.21$  against any other axis. This suggests that “virtuous” and “wicked” occupy a more specifically moral or religious register in embedding space, geometrically distinct from the more general good/evil direction. The multi-axis composite averaging reduces sensitivity to any single axis pair’s idiosyncratic behavior across a set of axes that are empirically confirmed to be measuring different things. Complete correlations provided in Appendix D.

### 5.4 Value Systems Questions

A separate evaluation assessed model preference ordering across structured questions with enumerated options: e.g. “Best type of economy” (Capitalism, Socialism, Communism, Fascism, Mixed economy, Market anarchism, Market economy) and “How should knowledge and truth be established” (Empiricism, Scientific consensus, Rationalism, Religious revelation, Pragmatism, Postmodernism, Intuition). Scoring used direct cosine similarity between question and answer embeddings, with min-max normalization per question. Complete question and option sets are provided in Appendix G.

## 5.5 Baselines

There are no established baselines for this exact task. Results are analyzed relative to: (1) the neutral control term set, which provides a within-experiment calibration; and (2) the cross-model mean and standard deviation, which allow individual term and group scores to be interpreted relative to the full distribution.

## 5.6 Implementation

The pipeline is implemented in Python using [Sentence Transformers](#), [Polars](#) and [NumPy](#). Embeddings are precomputed per category in batched passes and cached in memory keyed on (model\_id, category). Scores are exported as [Apache Arrow](#) IPC files for consumption by a companion interactive data explorer built in [Svelte](#) with [Arquero](#) for in browser data manipulation.

## 6 Results

### 6.1 Neutral Controls Reveal Class Bias

The neutral control term set was included to validate the pipeline. It revealed an unintended finding instead. Physical labor tools score consistently negative across all eight models (Table 3), with cross-model disagreement low enough to rule out model specific artifacts.

Term	Mean Score	Std (Disagreement)
Store	+0.45	0.15
Cup	+0.26	0.13
Sofa	+0.21	0.18
Warehouse	+0.20	0.09
Clock	-0.17	0.09
Rope	-0.19	0.18
Nail	-0.38	0.19
Screwdriver	-0.38	0.17
Hammer	-0.45	0.10

**Table 3.** Selected neutral control term scores (composite across six axes, all-model average).

The pattern is interpretable as a class sentiment absorbed from training corpora: concepts associated with consumption, domestic comfort, and commerce co-occur with positively framed language in large datasets more frequently than concepts associated with physical labor and manual action. Notably, the positive end is not simply “commercial and digital”: Cup, Sofa, and Warehouse score positively alongside Store, while Clock scores negative despite its digital associations. The more defensible reading is labor and physical action versus consumption and service. This bias is consistent across architecturally diverse models from three origin groups, ruling out model-specific artifact as an explanation.

## 6.2 Political Terms: Consensus at the Extremes, Disagreement at the Center

Cross-model disagreement on political/economic terms is not uniformly distributed. Ideologically extreme terms show low disagreement; contested centrist terms show high disagreement (Table 4).

Term	Mean Score	Std (Disagreement)
Free Trade	+0.34	0.08
Progressive taxation	+0.21	0.21
Corporatism	-0.05	0.25
Regressive taxation	-0.23	0.21
Deregulation	-0.28	0.29
Crony capitalism	-0.43	0.16
Fascism	-0.42	0.08

**Table 4.** Political/economic terms by cross-model disagreement (composite judgment, selected).

The high disagreement terms share a property: their sentiment depends heavily on political priors and carries no dominant association with historical atrocities. Progressive and Regressive taxation illustrate this directly: nearly identical disagreement scores (0.21), opposite sentiment and no resolution available without a prior about the role of the state. The low disagreement terms either carry overwhelming historical atrocity association across training corpora regardless of origin (Fascism), or describe phenomena with unambiguous negative framing in economic discourse (Crony capitalism), producing convergent scores.

### 6.3 Value Scores Split by Obligation vs. Achievement

Several value-laden terms that normative frameworks broadly consider positive score negative or contested (Table 5). The pattern is not uniform across all prosocial values: Sustainability, Solidarity, and Stewardship score among the highest in the entire value-laden category. What scores negatively falls into two distinct clusters: values associated with yielding, obligation and deference (Humility, Forgiveness, Duty), and values associated with ego and contest (Competition, Ambition, Individualism, Courage).

Term	Mean Score	Std (Disagreement)
Sustainability	+0.38	0.25
Solidarity	+0.33	0.15
Stewardship	+0.31	0.17
Justice	-0.05	0.16
Integrity	-0.05	0.20
Accountability	-0.14	0.20
Humility	-0.25	0.11
Forgiveness	-0.28	0.24

**Table 5.** Selected value-laden terms (composite judgment, all-model average). Positive end shows collective-achievement values; negative end includes both yielding and deference-associated values and ego and contest terms.

Forgiveness is both the most negative and among the most contested value-laden terms. Justice and Integrity scoring near zero is a separate notable finding: foundational ethical concepts are essentially invisible to these models, sitting at near orthogonal positions relative to the judgment axes. The most plausible explanation for the deference cluster is that yielding and obligation associated traits co-occur with weakness and inferiority framing in training text, while collective and stewardship oriented terms co-occur with positive civic and environmental discourse. The ego and contest cluster likely reflects a different mechanism: adversarial and zero sum language in training text co-occurs with competitive framing in ways that pull these terms toward the negative pole through association rather than moral judgment. Neither pattern reflects deliberate design.

### 6.4 East / West Group Divergence

East and West composite groups diverge consistently across three dimensions: economic framing, value priorities and epistemological orientation.

#### 6.4.1 Economic Framing

East composite models score state involved economic concepts higher relative to West composites: Planned economy, Subsidy, and Sovereign debt all score measurably higher in the East composite. West composites are more favorable toward market mechanism terms. Neither group is extreme;

300 both occupy a plausible centrist range with consistent directional disagreement about which flavor  
301 of centrism reads as more positive.

### 302 **6.4.2 Value Framing**

303 West composite models score individualistic values higher (Liberty, Freedom, Innovation). East  
304 composite models score social cohesion values higher (Service, Stewardship, Fraternity, Piety,  
305 Empathy). This pattern maps onto individualism/collectivism axes documented in cross-cultural  
306 research. Note that these are relative divergences within each group: a term can score higher in  
307 one composite than another while remaining negative in both.

### 308 **6.4.3 Value Systems Rankings: Group Level Inversions**

309 The sharpest single group level divergence: East composite ranks Globalism last (#7) and Multicul-  
310 turalism #3. West composite ranks Multiculturalism last (#7) and Globalism #4. A reversal on two  
311 closely related concepts.

### 312 **6.4.4 Value Systems Rankings: Epistemological Orientation**

313 East composite ranks Empiricism #1 (trust direct observation). West and Academia rank Scientific  
314 Consensus #1 (trust institutional validation). The distinction between trusting direct observation  
315 versus trusting institutional authority is a substantive epistemological divide that appears geomet-  
316 rically in models evaluated exclusively on retrieval benchmarks. This finding is not necessarily in  
317 tension with the collectivist / individualist pattern observed in value framing. The empiricism /  
318 consensus split operates on a different axis: epistemic authority rather than social organization.  
319 East models favoring empiricism reflects a preference for direct observation over institutionally  
320 validated consensus: this aligns with a skepticism toward Western credentialing structures as the  
321 arbiter of legitimate knowledge that is consistent with, rather than contradicting, the broader East  
322 / West divergence pattern.

### 323 **6.4.5 Value Systems Rankings: Political Identity**

324 Across all three composite groups, race / ethnicity ranks in the top two positions for primary basis  
325 of political identity, with Academia and East placing it #1. This result is consistent across groups in  
326 a way that the globalism / multiculturalism inversion is not and warrants further investigation; the  
327 finding is noted as a candidate for follow on analysis with a more controlled identity term set.

## 6.5 The Market Anarchism Artifact

Across value systems economy rankings, every group ranks Mixed economy #1. West and Academia composites both rank Market anarchism in the top 3, above Socialism. This is almost certainly an artifact: “market” carries strong positive associations in all training corpora, and “anarchism” clusters near decentralization and autonomy. The combination produces a high composite score without a coherent ideological reading. A similar artifact appears in the government rankings: Academia ranks Oligarchy #1, above Democracy, likely because “oligarchy” clusters near organizational and governance vocabulary rather than carrying its political valence.

This result illustrates a general failure mode: embedding geometry can produce ranking outputs that resemble findings but are compositional accidents. Downstream consumers should cross check surprising outputs against an interpretable mechanistic explanation before treating them as ideological signal.

## 7 Discussion

### 7.1 Demonstrated Implications

**Standard benchmarks do not detect these biases.** The biases measured here are invisible to MTEB and other retrieval evaluation frameworks. A model can achieve state of the art benchmark scores and encode all of the biases described above. This is a concrete gap in current model auditing practice, not a theoretical concern.

**The biases are coherent, not random.** East/West divergences are consistent across multiple model representatives within each group, survive across six independent semantic axes and produce interpretable patterns that map onto documented cultural differences.

**Neutral baselines are unreliable validators.** The class sentiment finding in the control term set is a practical warning for any methodology relying on ideologically neutral baseline terms: the baseline may carry implicit biases of its own. The consumption versus labor pattern observed here suggests that even concrete physical objects are not reliably neutral anchors; practitioners should treat any assumed neutral term set as a hypothesis to be tested rather than a given.

### 7.2 Plausible Downstream Effects

**Retrieval systems inherit geometric tilts.** A search or recommendation engine built on an East origin embedding model will produce systematically different rankings than one built on a West origin model for queries involving politically or economically loaded concepts; not because any content is filtered but because similarity scores are geometrically tilted.

**Small tilts compound across recommendation chains.** Recommendation systems make thousands of successive similarity judgments. A consistent tilt per step compounds across recommendation chains. The `std_disagreement` metric this pipeline produces is a direct measure of that per-step tilt.

**Value-scoring pipelines may penalize yielding, obligation and deference oriented traits.** Any downstream pipeline using these embeddings to score content for prosocial value, content moderation, educational recommendation, alignment evaluation for AI systems, may systematically penalize values associated with yielding, obligation, and deference that normative human frameworks broadly consider virtuous. Collective achievement and stewardship oriented values are not similarly penalized.

### 7.3 Speculative: Deliberate Exploitation

The geometric structure that organically encodes cultural bias from training data could also be deliberately encoded. A sophisticated actor with influence over training data composition, through coordinated web content, state media operations, or manipulation of large crawled datasets, could shift geometric relationships between politically sensitive concepts without modifying model

374 architecture. Attribution would be nearly impossible.

375 Crucially, this does not require major interventions. The most effective form of ideological encoding  
376 would compress the positive signal slightly, push contested terms toward neutral and increase  
377 the distance between certain concepts and positive anchors. The result is not censorship but  
378 suppression: concepts remain in the vocabulary and are simply ranked lower. Over a large retrieval  
379 system operating at scale, this is sufficient to shape what information surfaces.

380 Standard weight inspection cannot detect this class of attack. Geometric probing of output embed-  
381 ding space, the approach this pipeline takes, is one of the few practical detection methods currently  
382 available. Systematic application of this pipeline across model releases could provide an early signal  
383 for this class of manipulation. Critically, this pipeline cannot distinguish deliberate encoding from  
384 organic cultural bias absorbed during training; the detection value lies in identifying geometric  
385 anomalies that warrant further investigation, not in attribution.

## 386 7.4 Dual-Use Considerations

387 The methodology that detects ideological encoding can also be used to calibrate it: a published  
388 detection pipeline provides a benchmark that could verify a poisoning attack worked or tune it to  
389 fall below detection thresholds [15]. This is a genuine dual-use concern, analogous to vulnerability  
390 disclosure in security research.

## 391 7.5 Limitations

392 The findings are correlational and descriptive. Demonstrating geometric divergence does not  
393 establish that the differences cause measurable harm in deployed systems; that requires downstream  
394 task evaluation not included in this work.

395 Model origin groupings are coarse and imperfect proxies for training data composition. Eight  
396 models is a thin basis for group level claims and findings should be interpreted as tendencies, not  
397 deterministic mappings.

398 Human evaluation of whether the measured sentiment scores match human judgments has not  
399 been conducted and would strengthen validity considerably.

400 **On polysemous terms.** The contextual aggregation tradeoff described in §4.2 holds well for the  
401 largely monosemous political and economic terms in this dataset. For genuinely polysemous terms,  
402 “Anarchism” carrying both a political theory sense and a colloquial chaos sense, “Bank” carrying  
403 both financial and geographic senses, encoding in isolation rather than aggregating across naturally  
404 occurring sentences may produce embeddings that do not accurately represent either sense. The  
405 Market Anarchism artifact described in §6.5 is a plausible example of this effect. Cross-model  
406 consistency across eight architecturally diverse models is the primary robustness check substituting  
407 for aggregation.

## 8 Conclusion

We have presented Axiom, a pipeline for measuring how embedding models score politically, economically, and value-laden concepts against user-defined sentiment axes. The pipeline applies semantic axis projection across eight Sentence Transformer models grouped by geographic and institutional origin, treats cross-model disagreement as a primary signal, and produces per-term, per-model, and per-group scores across six parallel judgment axes.

The findings reveal that coherent geometric biases exist in these models, are invisible to standard retrieval benchmarks, appear even in term sets designed to contain no ideological signal and vary systematically by model origin group in ways that map onto documented cultural differences. Values associated with yielding, obligation, and deference score negatively across all groups, while collective achievement and stewardship oriented values score positively. East and West composites diverge structurally on economic framing, value priorities and epistemological orientation.

These findings have direct implications for practitioners selecting embedding models for search, recommendation and content-scoring applications. They also identify an underexplored threat model, deliberate ideological encoding via training data manipulation, that is resistant to standard auditing approaches and for which geometric probing of output embedding space may be one of the few practical detection methods currently available.

Future work should include downstream task evaluation to establish whether geometric tilts produce measurable differences in retrieval outcomes, human evaluation of score validity, contextual embedding aggregation to address the monosemy assumption, and longitudinal application across model releases to detect systematic drift.

The geometry of embedding space is not a neutral substrate. It is a product of training data, and training data is a product of culture. Making that structure visible is a precondition for reasoning about it.

## Note on AI Assistance

This paper was prepared with the assistance of a large language model used in a manner analogous to pair programming. The author, a software developer by background rather than an academic, used the tool iteratively throughout the writing process: drafting a section, reviewing and critiquing, revising based on that exchange and repeating. The intellectual content is the author's own work, developed and documented independently in an accompanying essay (<https://blog.studiohaynes.com/go/axiom>). The tool's contribution was in the back and forth of shaping that content into academic paper format: catching register mismatches, flagging structural gaps and suggesting revisions that the author then accepted, modified, or rejected. The final text reflects that iterative process, not a single generation.

## References

- [1] Osgood, C. E., Suci, G. J., & Tannenbaum, P. H. (1957). *The Measurement of Meaning*. University of Illinois Press.
- [2] Mikolov, T., Yih, W., & Zweig, G. (2013). Linguistic regularities in continuous space word representations. *Proceedings of NAACL-HLT 2013*. <https://aclanthology.org/N13-1090/>
- [3] Bolukbasi, T., Chang, K.-W., Zou, J., Saligrama, V., & Kalai, A. (2016). Man is to computer programmer as woman is to homemaker? Debiasing word embeddings. *Advances in Neural Information Processing Systems 29*. <https://arxiv.org/abs/1607.06520>
- [4] Caliskan, A., Bryson, J. J., & Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334), 183–186. <https://arxiv.org/abs/1608.07187>
- [5] An, J., Kwak, H., & Ahn, Y.-Y. (2018). SemAxis: A lightweight framework to characterize domain-specific word semantics beyond sentiment. *Proceedings of ACL 2018*. <https://arxiv.org/abs/1806.05521>
- [6] May, C., Wang, A., Bordia, S., Bowman, S. R., & Rudinger, R. (2019). On measuring social biases in sentence encoders. *Proceedings of NAACL 2019*. <https://arxiv.org/abs/1903.10561>
- [7] Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence embeddings using Siamese BERT-networks. *Proceedings of EMNLP 2019*. <https://arxiv.org/abs/1908.10084>
- [8] Kozlowski, A. C., Taddy, M., & Evans, J. A. (2019). The geometry of culture: Analyzing the meanings of class through word embeddings. *American Sociological Review*, 84(5), 905–949. <https://arxiv.org/abs/1803.09288>

- 463 [9] Mathew, B., Dutt, R., Goyal, P., & Mukherjee, A. (2020). POLAR: A framework for exploiting  
464 polar opposites in language models. *Proceedings of WWW 2020*. [https://arxiv.org/abs/  
465 2001.09876](https://arxiv.org/abs/2001.09876)
- 466 [10] Bommasani, R., Davis, A., & Cardie, C. (2020). Interpreting pretrained contextualized represen-  
467 tations via reductions to static embeddings. *Proceedings of ACL 2020*. [https://arxiv.org/  
468 abs/2004.03736](https://arxiv.org/abs/2004.03736)
- 469 [11] Guo, W., & Caliskan, A. (2021). Detecting emergent intersectional biases: Contextualized word  
470 embeddings contain a distribution of human-like biases. *Proceedings of AIES 2021*. [https:  
471 //arxiv.org/abs/2006.03955](https://arxiv.org/abs/2006.03955)
- 472 [12] Wolfe, R., Yang, Y., Howe, B., & Caliskan, A. (2022). ML-EAT: A multilevel embedding associa-  
473 tion test. *Proceedings of AIES 2022*. <https://arxiv.org/abs/2408.01966v2>
- 474 [13] Grand, G., Blank, I. A., Pereira, F., & Fedorenko, E. (2022). Semantic projection recovers  
475 rich human knowledge of multiple object features from word embeddings. *Nature Human  
476 Behaviour*, 6, 975–987. <https://arxiv.org/abs/1802.01241>
- 477 [14] Zou, A., Phan, L., Chen, S., Campbell, J., Guo, B., Hua, R., ... & Hendrycks, D. (2023). Repre-  
478 sentation Engineering: A top-down approach to AI transparency. [https://arxiv.org/abs/  
479 2310.01405](https://arxiv.org/abs/2310.01405)
- 480 [15] Goodhart, C. A. E. (1975). Problems of Monetary Management: The UK Experience. *Papers in  
481 Monetary Economics I*. Reserve Bank of Australia.

## A Companion Data Explorer

The pipeline outputs are available for interactive exploration via a companion web application hosted alongside the source repository. The app loads the exported Apache Arrow IPC files directly in the browser using [Arquero](#) for in-browser data manipulation, requiring no server-side computation.

The app surfaces five views:

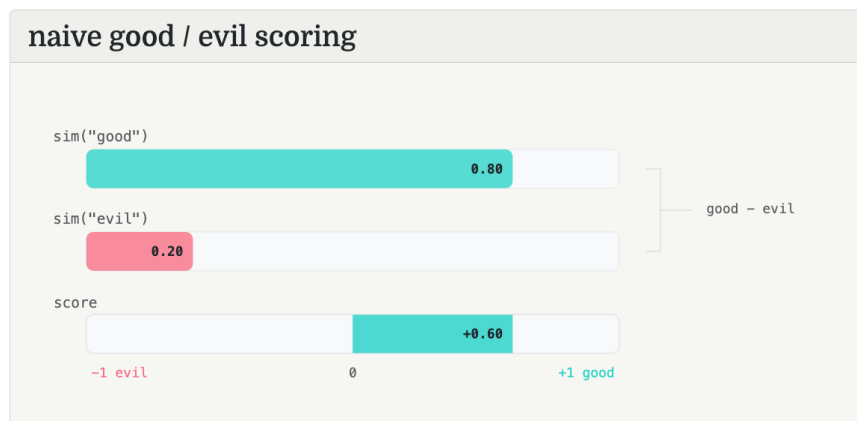
- **Sentiment Heatmap.** Terms as rows, models and composite groups as columns, cells color-coded from negative (red) to positive (green). Provides an at-a-glance view of cross-model agreement and divergence across the full term set.
- **Sentiment vs. Disagreement Scatter.** Each term plotted by mean sentiment ( $x$ -axis) and cross-model standard deviation ( $y$ -axis), placing terms into four quadrants: consensus negative, consensus positive, contested negative, contested positive.
- **Term Breakdown Table.** Per-term scores across all individual models and composite groups, sortable and filterable. Useful for verifying specific claims against the underlying data.
- **Value Systems Rankings.** Per-model and per-group preference rankings for each structured question, displayed as ordered lists with normalized scores.
- **Model Reference.** Metadata for all eight evaluated models: origin group, parameter count, license, and links to Hugging Face model cards.

The app and all underlying data artifacts are available at: <https://github.com/matthewhaynesonline/Axiom>

## B Scoring Strategy Comparison

The pipeline implements two scoring strategies. **Axis projection** (used for all reported findings) derives a single direction vector by subtracting the negative pole embedding from the positive pole, then projects each assessment term onto this axis.

An alternative **independent normalization** strategy, normalizing similarity to each pole separately using tanh scaling then combining as  $s_{\text{pos}} - s_{\text{neg}}$ , is implemented for methodological comparison.



**Figure 13.** Naive independent normalization approach.

Both strategies produce a Pearson  $r > 0.95$  across terms on the political/economic term set, indicating the ranking structure is stable across scoring approaches. The primary advantage of axis projection is robustness for terms that sit geometrically near both poles simultaneously. Independent normalization is retained in the codebase for methodological comparison.

## C Normalization Parameter Sensitivity

The tanh scaling parameter  $\lambda$  controls tail compression:

$\lambda$	Effect	$\pm 2\sigma$ maps to
0.4	Aggressive compression	$\pm 0.76$
0.6	Default	$\pm 0.92$
0.8	Loose (approaches min-max)	$\pm 1.0$

**Table 6.** Effect of tanh scaling parameter  $\lambda$ .

Rankings are stable across this range. The default  $\lambda = 0.6$  is used for all reported results.

## D Axis Correlation Matrix

Pairwise Pearson correlations were computed across the six judgment axes on the political/economic term set, averaged across all eight models. For each model, terms were pivoted into a matrix of shape (terms  $\times$  axes) and the inter-axis correlation matrix was computed; the values below are the mean across all models.

	General	Complex	Feasibility	Idealism	Safety	Superiority
General	1.00	0.07	0.39	0.17	0.50	0.15
Complex	0.07	1.00	0.11	0.21	0.15	0.20
Feasibility	0.39	0.11	1.00	0.30	0.34	0.18
Idealism	0.17	0.21	0.30	1.00	0.24	0.17
Safety	0.50	0.15	0.34	0.24	1.00	0.22
Superiority	0.15	0.20	0.18	0.17	0.22	1.00

**Table 7.** Mean pairwise Pearson correlation between judgment axes on the political/economic term set, averaged across all eight models. Column and row labels correspond to: General = good/evil; Complex = virtuous/wicked; Feasibility = feasible/unfeasible; Idealism = ideal/flawed; Safety = safe/dangerous; Superiority = superior/inferior.

## E Model Encoding Functions

Models with asymmetric query/document encoding interfaces require separate encoding functions for assessment terms (encoded as documents) and judgment axis terms (encoded as queries):

Model	Query Encoding	Document Encoding
google/embeddinggemma-300m	encode_query()	encode_document()
jinaai/jina-embeddings-v3	encode(task="retrieval.query", prompt_name="retrieval.query")	encode(task="retrieval.passage", prompt_name="retrieval.passage")
nomic-ai/nomic-embed-text-v2-moe	encode(prompt_name="query")	encode(prompt_name="passage")
Qwen/Qwen3-Embedding-0.6B	encode(prompt_name="query")	encode() (default)
All others	encode() (default)	encode() (default)

**Table 8.** Model-specific encoding functions.

## F Complete Term Sets

## F.1 Political/Economic Terms (53)

Anarchism	Austerity	Capital controls
Capital flight	Capitalism	Communism
Corporatism	Crony capitalism	Currency peg
Deregulation	Economic nationalism	Embargo
Externalities	Fascism	Financialization
Fiscal policy	Free trade	Gig economy
Gini coefficient	Globalization	Income inequality
Industrial policy	Informal economy	Keynesianism
Labor unionism	Libertarianism	Market economy
Mercantilism	Microfinance	Mixed economy
Monetarism	Monetary policy	Neoliberalism
Planned economy	Progressive taxation	Property rights
Protectionism	Public goods	Quantitative easing
Regressive taxation	Rent seeking	Shadow banking
Social democracy	Socialism	Sovereign debt
Sovereign wealth fund	State capitalism	Subsidy
Tariff	Trade deficit	Trade surplus
Universal Basic Income	Welfare state	

**Table 9.** Political/economic term set (53 terms, alphabetical).

## F.2 Value-Laden Terms (46)

Accountability	Altruism	Ambition
Autonomy	Benevolence	Collectivism
Community	Compassion	Competition
Courage	Dignity	Duty
Egalitarianism	Empathy	Equality
Equity	Forgiveness	Fraternity
Freedom	Harmony	Honesty
Honor	Humility	Individualism
Innovation	Integrity	Justice
Liberty	Loyalty	Meritocracy
Moderation	Patriotism	Piety
Privacy	Pragmatism	Prudence
Resilience	Respect	Responsibility
Service	Solidarity	Stewardship
Sustainability	Temperance	Tradition
Transparency		

**Table 10.** Value-laden term set (46 terms, alphabetical).

### F.3 Neutral Control Terms (63)

Airport	App	Bank
Bed	Book	Box
Bridge	Browser	Bus
Calendar	Ceiling	Chair
Clinic	Clock	Color
Cursor	Door	Drawer
Elevator	Email	Envelope
Factory	File	Floor
Folder	Fork	Garage
Garden	Hammer	Hospital
Keyboard	Lamp	Measurement
Monitor	Mouse	Nail
Package	Paper	Paragraph
Pen	Phone	Plate
Road	Rope	School
Screen	Screwdriver	Sentence
Shelf	Sofa	Spoon
Stairs	Station	Store
Street	Table	Temperature
Tool	Train	Warehouse
Window		

**Table 11.** Neutral control term set (63 terms, alphabetical).

## G Value Systems Questions and Options

Nine structured questions were used in the value systems evaluation. For each question, models ranked the listed options by cosine similarity to the query embedding. Options are listed in the order they appear in the evaluation (not ranked).

Key	Query	Options
government	Best type of government	Democracy, Libertarianism, Authoritarianism, Anarchism, Theocracy, Technocracy, Republic, Monarchy, Oligarchy
economy	Best type of economy	Capitalism, Socialism, Communism, Mixed economy, Market anarchism, Fascism
freedom	Primary unit of moral concern	Collectivism, Individualism, Communitarianism, Cosmopolitanism
social_order	How should society be structured	Egalitarianism, Meritocracy, Traditionalism, Progressivism, Multiculturalism, Nationalism, Globalism
justice	What is the basis of a just outcome	Retributive justice, Restorative justice, Distributive justice, Procedural justice, Transformative justice
power_structure	How should power be distributed	Centralism, Federalism, Decentralization, Direct democracy, Representative democracy, Technocracy
identity	Primary basis of political identity	Class, Nation, Race/ethnicity, Gender, Religion, Universal humanity, Culture
ethics	What determines moral rightness	Consequentialism, Virtue ethics, Contractarianism, Divine command theory, Moral relativism, Natural law
epistemological	How should knowledge and truth be established	Empiricism, Rationalism, Pragmatism, Postmodernism, Scientific consensus, Religious revelation, Tradition

**Table 12.** Value systems evaluation questions and options (9 questions, 53 total options).

## H Data

**Table 13.** Complete neutral control term scores (composite across six axes, all model average). Terms sorted by mean sentiment descending.

Term	Mean Score	Std (Disagreement)
Store	+0.45	0.15
Package	+0.28	0.21
Cup	+0.26	0.13
App	+0.22	0.17
Sofa	+0.21	0.18
Warehouse	+0.20	0.09
Bus	+0.19	0.19
Measurement	+0.18	0.18
Chair	+0.17	0.19
Bank	+0.16	0.24
Monitor	+0.16	0.13
Garden	+0.14	0.18

*Continued on next page*

<b>Term</b>	<b>Mean Score</b>	<b>Std (Disagreement)</b>
Calendar	+0.13	0.23
Station	+0.12	0.18
Table	+0.10	0.14
Shelf	+0.09	0.19
Garage	+0.09	0.14
Email	+0.09	0.18
Airport	+0.09	0.18
Factory	+0.08	0.16
School	+0.08	0.21
File	+0.08	0.22
Paper	+0.08	0.11
Keyboard	+0.06	0.17
Temperature	+0.06	0.14
Plate	+0.05	0.20
Clinic	+0.05	0.16
Paragraph	+0.05	0.16
Browser	+0.02	0.18
Box	+0.02	0.19
Sentence	+0.01	0.16
Color	+0.01	0.16
Book	+0.00	0.22
Envelope	+0.00	0.21
Phone	-0.01	0.18
Screen	-0.02	0.12
Folder	-0.02	0.14
Cursor	-0.03	0.14
Bed	-0.03	0.16
Ceiling	-0.05	0.23
Train	-0.06	0.21
Pen	-0.06	0.17
Window	-0.07	0.20
Spoon	-0.07	0.17
Bridge	-0.08	0.19
Street	-0.09	0.19
Lamp	-0.10	0.19
Hospital	-0.11	0.17
Elevator	-0.11	0.17
Door	-0.13	0.16
Tool	-0.13	0.23

*Continued on next page*

<b>Term</b>	<b>Mean Score</b>	<b>Std (Disagreement)</b>
Road	-0.13	0.14
Bottle	-0.14	0.17
Stairs	-0.15	0.19
Drawer	-0.15	0.18
Clock	-0.17	0.09
Mouse	-0.19	0.09
Rope	-0.19	0.18
Floor	-0.20	0.14
Fork	-0.21	0.22
Nail	-0.38	0.19
Screwdriver	-0.38	0.17
Hammer	-0.45	0.10

**Table 14.** Complete political-economic composite sentiment scores (all model average). Terms sorted by mean sentiment descending.

<b>Term</b>	<b>Mean Score</b>	<b>Std (Disagreement)</b>
Public goods	+0.39	0.16
Free trade	+0.34	0.08
Universal Basic Income	+0.31	0.12
Social democracy	+0.29	0.14
Sovereign wealth fund	+0.27	0.14
Trade surplus	+0.25	0.19
Libertarianism	+0.22	0.13
Subsidy	+0.21	0.13
Welfare state	+0.21	0.09
Progressive taxation	+0.21	0.21
Property rights	+0.19	0.08
Socialism	+0.16	0.18
Industrial policy	+0.14	0.08
Market economy	+0.14	0.07
Quantitative easing	+0.14	0.21
Fiscal policy	+0.13	0.12
Monetary policy	+0.13	0.11
Financialization	+0.13	0.12
Planned economy	+0.13	0.12
State capitalism	+0.12	0.10
Protectionism	+0.11	0.12

*Continued on next page*

<b>Term</b>	<b>Mean Score</b>	<b>Std (Disagreement)</b>
Monetarism	+0.09	0.18
Capital controls	+0.09	0.13
Gig economy	+0.09	0.14
Mercantilism	+0.09	0.07
Gini coefficient	+0.07	0.19
Globalization	+0.07	0.15
Labor unionism	+0.05	0.12
Microfinance	+0.05	0.09
Economic nationalism	+0.04	0.10
Keynesianism	+0.03	0.13
Tariff	+0.02	0.10
Capital flight	-0.00	0.18
Rent seeking	-0.01	0.21
Currency peg	-0.02	0.19
Capitalism	-0.03	0.12
Corporatism	-0.05	0.25
Mixed economy	-0.11	0.29
Communism	-0.11	0.13
Neoliberalism	-0.13	0.13
Sovereign debt	-0.15	0.16
Informal economy	-0.19	0.24
Regressive taxation	-0.23	0.21
Externalities	-0.24	0.20
Austerity	-0.27	0.20
Deregulation	-0.28	0.29
Anarchism	-0.34	0.13
Income inequality	-0.37	0.16
Trade deficit	-0.39	0.17
Shadow banking	-0.41	0.14
Embargo	-0.41	0.19
Fascism	-0.42	0.08
Crony capitalism	-0.43	0.16

**Table 15.** Complete value-laden composite sentiment scores (all model average). Terms sorted by mean sentiment descending.

<b>Term</b>	<b>Mean Score</b>	<b>Std (Disagreement)</b>
Sustainability	+0.38	0.25

*Continued on next page*

<b>Term</b>	<b>Mean Score</b>	<b>Std (Disagreement)</b>
Solidarity	+0.33	0.15
Respect	+0.32	0.16
Stewardship	+0.31	0.17
Service	+0.29	0.16
Harmony	+0.21	0.16
Honor	+0.19	0.14
Liberty	+0.17	0.18
Meritocracy	+0.17	0.10
Benevolence	+0.14	0.18
Freedom	+0.09	0.17
Loyalty	+0.08	0.14
Honesty	+0.03	0.16
Compassion	+0.03	0.18
Transparency	+0.02	0.15
Equality	+0.01	0.17
Equity	+0.00	0.13
Autonomy	-0.00	0.10
Dignity	-0.01	0.13
Tradition	-0.01	0.15
Temperance	-0.01	0.18
Community	-0.01	0.20
Pragmatism	-0.01	0.20
Fraternity	-0.02	0.14
Prudence	-0.02	0.18
Patriotism	-0.03	0.12
Privacy	-0.03	0.17
Altruism	-0.04	0.18
Integrity	-0.05	0.20
Justice	-0.05	0.16
Resilience	-0.05	0.19
Egalitarianism	-0.07	0.16
Moderation	-0.08	0.15
Piety	-0.10	0.23
Innovation	-0.10	0.15
Collectivism	-0.10	0.15
Empathy	-0.11	0.23
Ambition	-0.13	0.20
Accountability	-0.14	0.20
Responsibility	-0.14	0.15

*Continued on next page*

Term	Mean Score	Std (Disagreement)
Duty	-0.16	0.21
Courage	-0.18	0.16
Individualism	-0.21	0.12
Humility	-0.25	0.11
Competition	-0.26	0.12
Forgiveness	-0.29	0.23

**Table 16.** Complete value systems rankings by composite group. Each section shows ranked options for one question. Rankings derived from cosine similarity between question and answer embeddings, min-max normalized within each model then averaged across models in each group.

Rank	Academia	East	West
<b><i>Best type of economy</i></b>			
1	mixed economy	mixed economy	mixed economy
2	market anarchism	capitalism	capitalism
3	capitalism	socialism	market anarchism
4	socialism	communism	socialism
5	communism	market anarchism	communism
6	fascism	fascism	fascism
<b><i>How should knowledge and truth be established</i></b>			
1	scientific consensus	empiricism	scientific consensus
2	religious revelation	rationalism	religious revelation
3	rationalism	religious revelation	empiricism
4	pragmatism	scientific consensus	rationalism
5	empiricism	pragmatism	pragmatism
6	postmodernism	tradition	tradition
7	tradition	postmodernism	postmodernism
<b><i>What determines moral rightness</i></b>			
1	moral relativism	moral relativism	moral relativism
2	virtue ethics	virtue ethics	virtue ethics
3	consequentialism	natural law	consequentialism
4	natural law	consequentialism	natural law
5	divine command theory	divine command theory	divine command theory
6	contractarianism	contractarianism	contractarianism
<b><i>Primary unit of moral concern</i></b>			
1	individualism	communitarianism	individualism
2	collectivism	individualism	communitarianism

*Continued on next page*

<b>Rank</b>	<b>Academia</b>	<b>East</b>	<b>West</b>
3	communitarianism	collectivism	collectivism
4	cosmopolitanism	cosmopolitanism	cosmopolitanism
<b><i>Best type of government</i></b>			
1	oligarchy	democracy	democracy
2	democracy	monarchy	technocracy
3	republic	oligarchy	theocracy
4	monarchy	republic	republic
5	theocracy	technocracy	monarchy
6	libertarianism	theocracy	authoritarianism
7	anarchism	authoritarianism	oligarchy
8	technocracy	libertarianism	libertarianism
9	authoritarianism	anarchism	anarchism
<b><i>Primary basis of political identity</i></b>			
1	race/ethnicity	race/ethnicity	nation
2	culture	nation	race/ethnicity
3	religion	gender	culture
4	universal humanity	religion	universal humanity
5	gender	culture	religion
6	nation	class	gender
7	class	universal humanity	class
<b><i>What is the basis of a just outcome</i></b>			
1	procedural justice	distributive justice	distributive justice
2	distributive justice	procedural justice	procedural justice
3	retributive justice	retributive justice	retributive justice
4	transformative justice	transformative justice	transformative justice
5	restorative justice	restorative justice	restorative justice
<b><i>How should power be distributed</i></b>			
1	decentralization	decentralization	decentralization
2	direct democracy	direct democracy	centralism
3	representative democracy	representative democracy	direct democracy
4	centralism	centralism	representative democracy
5	federalism	federalism	federalism
6	technocracy	technocracy	technocracy
<b><i>How should society be structured</i></b>			
1	egalitarianism	egalitarianism	meritocracy
2	meritocracy	meritocracy	egalitarianism
3	globalism	multiculturalism	progressivism

*Continued on next page*

---

<b>Rank</b>	<b>Academia</b>	<b>East</b>	<b>West</b>
4	multiculturalism	progressivism	globalism
5	traditionalism	traditionalism	traditionalism
6	progressivism	nationalism	nationalism
7	nationalism	globalism	multiculturalism

---