

# Orphaned Sophistication: Detecting AI-Generated Prose Through Structurally Unsupported Figurative Language

Anonymous ACL submission

## Abstract

We identify a novel stylometric artifact in large language model (LLM) prose generation: *orphaned sophistication*, the production of figuratively sophisticated word choices that lack structural support from their surrounding context. Through controlled experiments comparing 25 human-authored passages against 100 LLM-generated passages from five model runs spanning three independent model families (Anthropic, OpenAI, Google), we demonstrate that LLMs produce polysemous words whose secondary semantic fields overlap with active figurative registers at rates significantly exceeding human prose (initial single-model analysis,  $n = 45$ : Fisher's exact test,  $p = 0.001$ , Cohen's  $h = 1.69$ ). We propose a theoretical explanation rooted in training-weight distributional bias and formalize a three-dimensional orphanhood model (isolation, chain connectivity, tonal preparation), implementing a deterministic rule-based detector achieving 28.0% true positive rate on LLM prose with 4% false positive rate on human prose (cross-family pooled analysis,  $n = 125$ : Fisher's  $p = 0.006$ , Cohen's  $h = 0.71$ ). The signal spans all three families tested: Anthropic (15–35%), OpenAI GPT-4o (15%), and Google Gemini 2.5 Flash (40%,  $p = 0.004$ ). The central finding is that the uncanny valley of AI prose is a structural coherence failure, not a lexical quality failure, and it is measurable. We provide a semiotic interpretation grounding the signal in the distinction between Barthes's *significance* and *signification*, and identify a structurally identical pathology in computational drug repurposing, suggesting domain generality.

## 1 Introduction

The detection of AI-generated text has become a critical problem in computational linguistics, digital forensics, and publishing. Existing approaches fall broadly into two categories: statistical fingerprinting methods that measure distributional properties of token sequences (perplexity, burstiness,

$n$ -gram frequency profiles), and watermarking schemes that embed detectable signals during generation. Both share a fundamental limitation: they identify *that* a text is machine-generated without explaining *why* it reads as machine-generated. The qualitative experience of encountering AI prose, the uncanny valley sensation (Mori, 1970) that something is simultaneously competent and wrong, remains unformalized.

We present a third approach grounded in structural analysis of figurative language. Our central claim is that autoregressive language models, as a consequence of distributional biases in their training data, produce a specific and detectable artifact: figuratively sophisticated word choices that are structurally orphaned from the prose architecture that would justify them in human writing. A human author who writes “the hungry steel teeth” in a passage about a sawmill has, in deliberate literary prose, prepared that personification through tonal shifts, metaphor chains, or explicit signposting. An LLM produces the same construction as a default token prediction, without preparation, without continuation, and without architectural awareness that the construction requires either.

This paper makes four contributions: (1) empirical identification of the orphaned sophistication artifact through controlled experiments with formal statistical testing; (2) theoretical explanation through a training-weight over-indexing model; (3) a formal detection framework based on a three-dimensional orphanhood model, implemented as a fully deterministic rule-based algorithm; and (4) a semiotic interpretation connecting the signal to the distinction between Barthes's *significance* and *signification*.

## 2 Related Work

128

### 2.1 AI Text Detection

Current detection methods include perplexity-based classifiers (Mitchell et al., 2023), watermarking (Kirchenbauer et al., 2023), and supervised classifiers trained on LLM output distributions (Tian et al., 2023). These methods achieve variable accuracy and degrade across domains, paraphrasing attacks, and model updates (Krishna et al., 2023; Wu et al., 2025). Critically, none provides a structural explanation for *what* distinguishes AI prose from human prose at the level of craft.

We do not claim that orphaned sophistication detection replaces these methods. It operates in a different regime: short-form literary and descriptive prose where figurative language is expected.

### 2.2 Polysemy and Priming in LLM Output

Kugler (2025) demonstrates that LLM output exhibits a “flatter semantic space” than natural language (frequency-specificity correlation:  $\rho \approx -0.3$  for LLMs vs.  $\rho \approx -0.5$  to  $-0.7$  for human text). This flat distribution is consistent with our over-indexing hypothesis. Jumelet et al. (2024) demonstrate that lexico-semantic overlap boosts token-level probability in transformers through structural priming effects, confirming the mechanistic foundation of our claim.

### 2.3 Coherence as a Latent Dimension

Shaib et al. (2025) develop a taxonomy of “AI slop” through expert annotation, finding that standard text metrics fail to capture coherence dimensions and that capable LLMs likewise fail to reliably identify slop. Our orphanhood framework provides one structural answer: it operationalizes a specific form of incoherence (figurative constructions arriving without architectural support) as a measurable, deterministic signal.

## 3 Theoretical Framework

### 3.1 Latent Semantic Recruitment

We define *Latent Semantic Recruitment* (LSR) as the phenomenon whereby an autoregressive language model, generating text within an active figurative register  $R$ , disproportionately selects polysemous words whose secondary semantic fields overlap with  $R$ .

Let  $w$  be a word token with primary sense  $s_1$  (contextually appropriate) and secondary sense  $s_2$

(not contextually required). LSR occurs when  $P(w \mid \text{context}, R) > P(w \mid \text{context}, \neg R)$  specifically because the embedding of  $s_2$  overlaps with  $R$  in the model’s representation space.

This follows from the standard transformer output computation (Vaswani et al., 2017). The logit for token  $w$  at position  $t$  is  $z_w = \mathbf{v}_w^\top \mathbf{h}_t$ , where  $\mathbf{v}_w$  is the output embedding and  $\mathbf{h}_t$  is the hidden state. When context contains register-activating content,  $\mathbf{h}_t$  encodes semantic components overlapping with register-aligned secondary senses. For polysemous words,  $\mathbf{v}_w$  encodes both  $s_1$  and  $s_2$ , and the inner product with a register-active  $\mathbf{h}_t$  is elevated compared to a monosemous alternative encoding only  $s_1$ .

### 3.2 Training-Weight Over-Indexing

LSR explains the mechanism, but not why the result is detectable. Human writers also select polysemous words. The critical question is why LLM polysemous usage is distinguishable.

We propose the **training-weight over-indexing hypothesis**: training corpora contain a distributional bias that systematically over-represents exceptional figurative prose. The texts exhibiting the most sophisticated polysemous craft (Conrad, McCarthy, Woolf, Morrison) receive the most analytical attention, pedagogical citation, and anthology inclusion, producing the most duplication across training data. Under standard cross-entropy training, tokens appearing more frequently contribute proportionally more to the cumulative gradient. The model therefore learns to reproduce this level of sophistication not as exceptional but as the expected register of competent prose.

The result is a distributional inversion. In the population of human writers, polysemous craft at the level of Conrad or McCarthy occupies the far right tail. In the model’s learned distribution, it occupies the mode. We hypothesize that when this disparity occurs, it constitutes a detectable signal.

**Caveat.** A direct demonstration would require measuring the frequency of specific figurative constructions in training data and correlating that frequency with generation probability, an analysis requiring training-data access we do not have. The hypothesis is argued from distributional logic, consistent with the observed signal, but not independently verified.

### 3.3 Orphaned Sophistication

226

The over-indexing hypothesis predicts that LLMs will produce sophisticated figurative language *without the structural architecture that earns it*. We define **orphaned sophistication** as a figurative construction satisfying three conditions:

**Isolation.** The figurative density of the sentence containing  $w$  is significantly higher than its neighbors (window  $\pm 2$  sentences). Score:  $\min(1.0, (\varphi(s_w) - \bar{\varphi}(N(s_w)))/\tau_1)$ , where  $\tau_1 = 0.2$ .

**Chain disconnection.** The register field activated by  $s_2$  is not activated by other words within  $\pm 3$  sentences. In human literary prose, figurative constructions participate in metaphor chains. Score: 0 connections = 1.0, 1 = 0.6, 2 = 0.2, 3+ = 0.0.

**Lack of preparation.** The context is scored for signposting markers: simile constructions, explicit frame-setting, tonal shifts (sentence-length ratio  $> 2.5:1$ ), and figurative density in adjacent sentences ( $> 0.15$ ). Score ranges from 1.0 (fully unprepared) to 0.0 (fully prepared).

A word’s orphanhood score is the arithmetic mean of its three test scores. Classification threshold: 0.6 (set *a priori*, not optimized on test data).

### 3.4 Semiotic Interpretation

The framework admits an interpretation through Barthes’s semiotic theory (Barthes, 1970, 1973), though we acknowledge this narrows Barthes’s framework considerably. In *S/Z*, *significance* concerns the plurality of meaning generated by the interaction of multiple codes in writerly texts; our usage maps a more architectural reading onto the term, treating *significance* as requiring structural scaffolding. This narrowing is deliberate and operational.

When Conrad writes “the rudder would bite,” the word performs something closer to *signification*: it participates in a novel-length architecture. When an LLM writes “the hungry steel teeth,” the same semantic content is present but the structural labor is absent. The word performs *signification* without *significance*. Our three tests map onto this distinction: isolation measures sustained vs. anomalous sophistication; chain connectivity measures productive labor vs. standalone activation; preparation measures deliberate register transition vs. its absence.

A necessary caveat: a sufficiently long LLM text may, through stochastic density alone, produce pas-

sages scoring well on all three dimensions. Our detector measures necessary conditions for *significance* (structure is present) but not the sufficient condition (structure was produced through authorial labor). This is why we describe the framework as identifying *orphaned* sophistication rather than *unearned* sophistication.

## 4 Experimental Method

### 4.1 Corpus Construction

We assembled three corpora across five physical-register domains (ocean storm, kitchen/restaurant, blacksmith/forge, battlefield surgery, sawmill/logging):

**Human Corpus A (Published).** 20 passages (~100–200 words), drawn from published fiction and nonfiction spanning 1902–2016, four per domain. Authors include Conrad (1902), Hemingway (1952), McCarthy (1985), Bourdain (2000), Proulx (2016), Powers (2012), Barker (1991), Remarque (1929), Orwell (1933), and ten others.<sup>1</sup>

**Human Corpus B (Non-professional).** 5 passages (~150–250 words), hand-written by a non-professional writer under experimental conditions (one per domain, written under time pressure without revision, before the detection framework was developed).

**LLM Corpora C.** 100 passages total from five model runs spanning three independent families: Claude Sonnet 4 (20 + 20 replication, Anthropic), Claude Haiku 3.5 (20, Anthropic), GPT-4o (20, OpenAI), Gemini 2.5 Flash (20, Google). All generated via API at temperature 1.0, four passages per domain, identical prompts across models. Prompts requested 150–200 word passages specifying physical detail, past tense, third person, no dialogue.

### 4.2 Corpus Provenance

Corpus A passages were reproduced from LLM training data, introducing a potential circularity (Section 7.5). Corpus B provides an uncontaminated baseline. The detection instrument is deterministic (no LLM judgment in scoring), so circularity applies only to corpus construction.

### 4.3 Detection Instrument

We developed three successive detection instruments. Detector v1 (rate-based) counted pol-

<sup>1</sup>Full corpus: O’Brian (1969), Junger (1997), Buford (2006), Fisher (1954), Thompson (1945), McPhee (1975), Sturt (1923), Hooker (1968), Kesey (1964), Berry (2000), Pollan (1997).

ysemos words with register-aligned secondary senses; it was discarded because human and LLM rates were too similar. Detector v2 introduced domain-literal filtering, personification detection, and metaphor signpost detection, achieving strong separation but unable to distinguish skilled human figurative construction from LLM-generated equivalents at the individual word level.

The reported instrument (detector v3) identifies figurative polysemous words using v2’s mechanisms, then subjects each candidate to the three orphanhood tests defined in Section 3.3. The algorithm processes each sentence, identifying words that (a) are not in the domain-literal set for the passage’s domain (34–41 words per domain), (b) appear in at least one of six register fields (consumption, personification, body, violence, fire/heat, water/weather; 16–29 words each), and (c) exhibit figurative usage (personification, animate verb, or animate-quality modifier). Qualifying words receive the three orphanhood scores; the arithmetic mean must exceed 0.6 for classification as orphaned.

The detector is fully deterministic: no neural network, no LLM, no learned parameters. All thresholds are set a priori. Domain-literal filtering operates on the passage’s declared domain, assigned at corpus construction time.<sup>2</sup>

#### 4.4 Statistical Methods

All comparisons use Fisher’s exact test (appropriate for small-sample count data). We report one-sided  $p$ -values (testing the directional hypothesis that LLM rates exceed human rates) and two-sided  $p$ -values. Confidence intervals use the Clopper-Pearson exact method ( $\alpha = 0.05$ ). Effect sizes are Cohen’s  $h$ , where  $h > 0.8$  is conventionally large. No multiple-comparison correction is applied to the primary analysis (single pre-specified comparison); per-domain exploratory analyses are flagged as uncorrected.

#### 4.5 Experimental Design

The primary experiment (Experiment 8c/v3): all 25 human and 20 LLM passages (Corpus C-Sonnet) were processed by detector v3; orphanhood scores were computed for each flagged word; results were aggregated by source and domain. Cross-model replication used 20 additional passages each from

Source <sup>71</sup>	<i>n</i>	Orphaned	Flagged	Rate
Published human	20	1	1/20	5.0%
Non-prof. human	5	0	0/5	0%
<b>All human</b>	<b>25</b>	<b>1</b>	<b>1/25</b>	<b>4.0%</b>
Sonnet (primary)	20	9	7/20	35.0%

Table 1: Detector v3 results, primary experiment. Fisher’s  $p = 0.010$  (one-sided), Cohen’s  $h = 0.86$ .

Word	Register	Score	Iso	Chain	Prep
“hungry”	personif.	0.88	0.6	1.0	1.0
“stubborn”	personif.	0.73	0.2	1.0	1.0
“bite”	consumption	0.73	0.2	1.0	1.0
“roar”	water/weath.	0.80	0.4	1.0	1.0

Table 2: Four orphaned words in passage L06 (sawmill, Sonnet). Four distinct register fields, no chain connectivity between any, no preparation for any.

Sonnet 4 and Haiku 3.5 under identical conditions. Cross-family validation used 20 passages each from GPT-4o and Gemini 2.5 Flash, testing whether the signal generalizes beyond a single model family.

## 5 Results

### 5.1 Primary Analysis

The single human detection was Conrad’s “the rudder would bite again,” a nautical usage where “bite” is arguably domain-literal. In the 5 non-professional passages, zero detections occurred.

The detector v2 analysis (unjustified figurative polysemy, before orphanhood filtering) showed a rate ratio of  $18.8 \times$  (LLM 0.750 per passage vs. human 0.040), Fisher’s  $p = 0.001$ , Cohen’s  $h = 1.69$ . The v3 orphanhood model correctly reclassified 6 of the 15 v2 detections as “earned”: cases where the model had accidentally produced chain connectivity (e.g., “bit” appearing near “teeth” or “hungry”).

### 5.2 Qualitative Analysis

The most striking LLM passage was L06 (sawmill domain, Sonnet), which contained four orphaned words:

By contrast, human passages employing figurative language do so within explicitly prepared frames. The non-professional sawmill passage (“Life shaves pieces of your health off... Bertha takes that, too”) signposts the saw-as-life metaphor (“That’s life, that is”), develops it across multiple sentences, and connects to a chain of related vocabulary. The detector correctly identifies this as

<sup>2</sup>Full algorithmic specification, domain-literal sets, register field taxonomy, and all code will be released upon publication.

Model	Family	<i>n</i>	Rate	<i>h</i>	<i>p</i>
Human	—	25	4.0%	—	—
Gemini 2.5	Google	20	40.0%	0.97	.004
Sonnet (orig)	Anth.	20	35.0%	0.86	.010
Haiku 3.5	Anth.	20	35.0%	0.86	.010
GPT-4o	OpenAI	20	15.0%	0.39	.224
Sonnet (repl)	Anth.	20	15.0%	0.39	.224
<b>All LLM</b>	<b>3 fam.</b>	<b>100</b>	<b>28.0%</b>	<b>0.71</b>	<b>.006</b>

Table 3: Cross-family validation. All *p*-values one-sided Fisher’s exact. Pooled 95% CI: LLM [0.195, 0.379]; human [0.001, 0.204]. Power: 93.8%.

Domain	Human	LLM	Fisher <i>p</i>
Sawmill	0/5	4/4	0.008*
Surgery	0/5	4/4	0.008*
Blacksmith	0/5	1/4	0.444
Ocean storm	0/5	1/4	0.444
Kitchen	0/5	0/4	1.000

Table 4: Per-domain results (original Sonnet run). \*Uncorrected for multiple comparisons.

349  
earned.

### 350 5.3 Cross-Model and Cross-Family Validation

351 Three of five model runs are individually significant.  
352 The pooled analysis across three independent  
353 families provides the definitive test ( $p = 0.006$ ,  
354  $h = 0.71$ , power 93.8%).

355 **Cross-family word convergence.** The same  
356 register fields and often the same words (“roar,”  
357 “hungry,” “angry,” “bit”) recur across independently  
358 trained models from three organizations, strongly  
359 supporting the over-indexing hypothesis.

360 **Gemini signal strength.** Gemini produced  
361 the highest orphanhood rate (40%, 8/20), with  
362 one blacksmith passage producing three orphaned  
363 words. The word “grip” in that passage warrants  
364 scrutiny as potentially domain-literal (blacksmiths  
365 literally grip tools); removing it would reduce or-  
366 phaned words from 10 to 9 without affecting the  
367 passage-level rate.

### 368 5.4 Per-Domain Distribution

369 The kitchen domain produced a null result (0/5  
370 human, 0/4 LLM), the only domain with zero de-  
399 tections. Culinary language is inherently action-  
400 oriented and consumption-related, so words that  
401 would register as figurative in other domains are  
402 domain-literal in a kitchen context. The detector  
403 correctly identifies these as non-figurative.<sup>33</sup>

Config	Dims	Hum	LLM	<i>h</i>	<i>p</i>
Full	I+C+P	1/25	27/100 <sup>†</sup>	.690	.008
—Isolation	C+P	1/25	33/100	.821	.002
—Chain	I+P	1/25	21/100	.549	.035
—Preparation	I+C	2/25	26/100	.497	.041

Table 5: Ablation study (125 passages). I = isolation, C = chain, P = preparation. <sup>†</sup>One Sonnet passage falls at the classification boundary (score  $\approx 0.60$ ), producing 27/100 in the unified ablation pass vs. 28/100 in the incremental main analysis. The difference is due to floating-point variation and does not affect qualitative conclusions.

Condition	Orphan score
Suppressed (“avoid figurative”)	0.000
Neutral (no instruction)	0.098
Amplified (“use vivid language”)	0.755

Table 6: Dose-response experiment. Orphanhood is register-dependent and dose-responsive.

## 376 5.5 Ablation Study

377 The results reveal an asymmetric architecture. Re-  
378 moving chain connectivity reduces *h* by 0.141, con-  
379 firming that chain detection captures discriminative  
380 signal. Removing preparation produces the largest  
381 *h* degradation (0.193) and doubles the human false  
382 positive rate from 4.0% to 8.0%, indicating that  
383 preparation is the dimension most responsible for  
384 specificity.

385 Removing isolation *increases* *h* to 0.821 because  
386 isolation functions as a conservative filter, sup-  
387 pressing true positives where the figurative spike  
388 coincides with mildly elevated neighborhood den-  
389 sity. All four configurations maintain significance  
390 ( $p < 0.05$ ).

## 391 6 Mechanism Validation

### 392 6.1 Monte Carlo Logit Proxy

393 A Monte Carlo simulation generated 100,000 ran-  
394 dom word-context pairings and computed register-  
395 overlap scores. Of 15 register-field/domain pair-  
396 ings, 12 showed zero overlap between the random  
397 distribution and the observed LLM orphan scores  
398 ( $p < 10^{-5}$  each), confirming that the observed  
399 scores are not achievable by chance co-occurrence.

### 400 6.2 Dose-Response

### 401 6.3 Token Probability Probing

402 To test the over-indexing hypothesis at the gen-  
403 eration level, we designed eight probes targeting

Probe	Register	Anth.	OAI	Gem.
SAW_BITE	Consumption	1.41	3.43	$\infty^*$
OCEAN_ROAR	Vocalization	1.17	1.44	3.00
FORGE_STUB	Personif.	9.50	3.00	2.00
SURG_SCRM	Vocalization	1.38	2.50	$\infty^*$

Table 7: Preference ratios (literary/equivalent) for the four probes showing consistent literary preference across all three families. \*Infinite: literary words present, zero equivalents in all completions. Equivalent word lists are not exhaustive;  $\infty$  indicates strong directional preference, not absence of all non-literary vocabulary. Four additional probes showed mixed/reversed preferences.

404 constructions the detector most frequently flags.  
 405 Each probe provides a physical-register context  
 406 (e.g., “Write a paragraph describing a sawmill  
 407 blade cutting through hardwood”) and generates  
 408  $N$  completions at temperature 1.0 ( $N=20$  for An-  
 409 thropic/Gemini,  $N=10$  for OpenAI). For each com-  
 410 pletion, we count “literary” words (high-prestige  
 411 figurative constructions: “bite,” “teeth,” “hungry,”  
 412 “stubborn,” etc.) versus semantically equivalent al-  
 413 ternatives (“cut,” “slice,” “hard,” “rigid,” etc.). Full  
 414 word lists are provided with the code release.

415 The FORGE\_STUBBORN probe produced the  
 416 strongest signal: Anthropic generated personifi-  
 417 cation vocabulary at  $9.5\times$  the rate of physical-  
 418 property alternatives. All three families showed  
 419 aggregate literary preference (Anthropic  $1.30\times$ ,  
 420 OpenAI  $1.45\times$ , Gemini  $1.91\times$ ).

## 421 7 Discussion

### 422 7.1 The Uncanny Valley Formalized

423 The orphaned sophistication framework provides  
 424 a structural account of the “uncanny valley” of  
 425 AI prose (Mori, 1970). The deficiency lies not  
 426 in vocabulary or grammar but in the *relationship*  
 427 between sophistication and structure: the text pro-  
 428 duces figurative constructions implying architec-  
 429 tural control, but the architecture is absent. This  
 430 formalizes the observation that AI prose reads as  
 431 “too good” at the sentence level while failing at the  
 432 paragraph level (Shaib et al., 2025).

### 433 7.2 Why This Is Not Watermarking

Orphaned sophistication is not an imposed signal;  
 under the over-indexing hypothesis, it is an emer-  
 gent artifact of training. If correct, the signal would  
 resist removal by post-processing or prompt engi-  
 neering. Whether fine-tuning could selectively

reduce orphanhood without degrading prose quality  
 is an open empirical question.

## 441 7.3 Alternative Explanations

442 **Attention span.** LLMs may produce orphaned so-  
 443 phistication due to attention-window limitations  
 444 rather than training-weight bias. However, this  
 445 does not explain why the *specific* constructions are  
 446 so consistent across independent generations. At-  
 447 tention limitations would predict random figurative  
 448 orphanhood; we observe patterned orphanhood.

449 **Mode collapse.** All passages were generated  
 450 at temperature 1.0, a regime that Holtzman et al.  
 451 (2020) demonstrate substantially reduces repeti-  
 452 tive degeneration. We observe the same figurative  
 453 strategy expressed in varied syntactic frames, more  
 454 consistent with a learned register preference. Tem-  
 455 perature 1.0 reduces repetition while preserving  
 456 the probability distribution’s shape, precisely the  
 457 regime where over-indexing effects would manifest  
 458 as preferential selection.

## 459 7.4 Implications

460 If the over-indexing hypothesis is correct, or-  
 461 phaned sophistication should be present in all  
 462 LLMs trained on standard web corpora. The signal  
 463 should persist across architectures because it arises  
 464 from distributional properties of training data. The  
 465 signal is interpretable: a detection report points to  
 466 specific words, explains why they are orphaned,  
 467 and provides structural explanation. For writers us-  
 468 ing LLMs collaboratively, the framework provides  
 469 actionable revision guidance: flagged passages re-  
 470 quire not deletion but *architecture* (build a chain,  
 471 prepare the register shift, sustain figurative den-  
 472 sity).

## 473 7.5 The Generalizable Principle

474 The hypothesis predicts that any domain where  
 475 models are trained on corpora dominated by ex-  
 476 ceptional exemplars will exhibit an analogous ar-  
 477 tifact. Du et al. (2026) independently identify  
 478 “hard negatives” in computational drug repurpos-  
 479 ing: well-studied compounds appearing ideal due  
 480 to high knowledge-graph connectivity but failing  
 481 clinically. The mechanism is structurally identical:  
 482 FDA-approved drugs dominate training corpora  
 483 through citation and patent literature, producing  
 484 the same over-indexing dynamic. Graph neural  
 485 networks learn to produce binding moieties resem-  
 486 bling successful drugs as default output, locally  
 487 brilliant binding predictions structurally orphaned

488 from the ADMET properties that would make them  
489 clinically viable. This corresponds to the “activity  
490 cliff” problem in medicinal chemistry (Stumpfe  
491 and Bajorath, 2012).

492 The mapping is exact: a high-affinity binding  
493 moiety without metabolic stability is an orphaned  
494 figurative word. Isolation, chain disconnection, and  
495 lack of preparation all have molecular analogues.  
496 The convergence, observed independently in a dif-  
497 ferent modality by researchers with no knowledge  
498 of the orphanhood framework, suggests the pathol-  
499 ogy is a general property of training-data distribu-  
500 tional bias. (We note that Du et al. is a bioRxiv  
501 preprint not yet peer-reviewed.)

## 502 Limitations

- 503 1. **Sample size.**  $n = 125$  (25 human, 100 LLM)  
504 across 5 domains and 3 families. Statistically  
505 significant but modest for strong generalization.  
506
- 507 2. **Human corpus provenance.** Corpus A was  
508 reproduced from LLM training data. Corpus B provides an uncontaminated baseline,  
509 but ideal replication would use passages trans-  
510cribed from physical books.  
511
- 512 3. **Human corpus skill ceiling.** The detector’s  
513 FPR was measured against elite literary prose.  
514 Its behavior on mid-tier prose (workshop fiction,  
515 genre fiction, journalism) is unknown  
516 and is a priority for future validation.  
517
- 518 4. **No baseline comparison.** We have not run  
519 existing detectors (DetectGPT, GPTZero) on  
the same corpus.  
520
- 521 5. **Domain specificity.** Currently implemented  
522 for five physical-register domains. Extension  
523 to abstract registers requires additional word  
sets.  
524
- 525 6. **Passage length.** Orphanhood tests operate on  
526 passage-length windows. Earned figurative  
527 language can score as orphaned if excerpted  
from longer works.  
528
- 529 7. **Same-model variance.** The Sonnet replica-  
530 tion (15%) was individually non-significant  
( $p = 0.224$ ). Post-hoc power analysis shows  
531 80% power to detect  $h \geq 0.53$  at  $n = 20$ ; the  
532 replication’s  $h = 0.39$  falls below this thresh-  
533 old. The non-significant result reflects power  
534 limitations, not absent signal.  
535

## 8 Conclusion

We have identified and formalized a novel artifact of autoregressive language generation: orphaned sophistication. Our experiments provide statistically significant evidence across three independent model families (Fisher’s  $p = 0.006$ , Cohen’s  $h = 0.71$ ,  $n = 125$ ). The artifact arises, we argue, from training-weight over-indexing on exceptional exemplars, causing models to produce locally sophisticated outputs without the structural architecture that would earn them.

The detection signal is structural and interpretable: it provides not merely a classification but a diagnosis of where the architecture is missing and what kind of work would repair it. The core contribution is a reframing: the uncanny valley of AI prose is a structural coherence failure, not a lexical quality failure, and it is measurable. The machine does not write badly. It writes too well, in moments that have not been earned.

## LLM Assistance Disclosure

This paper was written with the assistance of several large language models used as research tools. Claude Opus 4 (Anthropic) was used to implement detection algorithms, reproduce published passages from training data for Corpus A, generate statistical analyses, and draft the manuscript under the author’s direction. Claude Sonnet 4, Claude Haiku 3.5, GPT-4o, and Gemini 2.5 Flash generated the respective LLM test corpora. No model involved in corpus generation was involved in scoring: all detection was performed by deterministic rule-based algorithms.

## References

- Pat Barker. 1991. *Regeneration*. Viking.  
Roland Barthes. 1970. *S/Z*. Éditions du Seuil.  
Roland Barthes. 1973. *Le Plaisir du texte*. Éditions du Seuil.  
Wendell Berry. 2000. *Jayber Crow*. Counterpoint.  
Anthony Bourdain. 2000. *Kitchen Confidential*. Bloomsbury.  
Bill Buford. 2006. *Heat*. Knopf.  
Joseph Conrad. 1902. *Typhoon*. Heinemann.  
Ruiqi Du, Maximilian Fung, Yifei Hu, and David Liu. 2026. Overcoming topology bias and cold-start limitations in drug repurposing: A clinical-outcome-aligned LLM framework. *bioRxiv*.

582	M.F.K. Fisher. 1954. <i>The Art of Eating</i> . World Publishing.	630
583		
584	Ernest Hemingway. 1952. <i>The Old Man and the Sea</i> . Scribner.	631
585		
586	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In <i>Proceedings of ICLR 2020</i> .	632
587		
588		
589	Richard Hooker. 1968. <i>MASH: A Novel About Three Army Doctors</i> . Morrow.	633
590		
591	Jaap Jumelet, Willem Zuidema, and Arabella Sinclair. 2024. Syntactic structural priming in large language models. In <i>Proceedings of ACL 2024</i> .	634
592		
593		
594	Sebastian Junger. 1997. <i>The Perfect Storm</i> . Norton.	635
595		
596	Ken Kesey. 1964. <i>Sometimes a Great Notion</i> . Viking.	636
597		
598	John Kirchenbauer, Jonas Geiping, Yuxin Wen, Jonathan Katz, Ian Miers, and Tom Goldstein. 2023. A watermark for large language models. In <i>Proceedings of ICML 2023</i> .	637
599		
600	Kalpesh Krishna, Yixiao Song, Marzena Karpinska, John Wieting, and Mohit Iyyer. 2023. Paraphrasing evades detectors of AI-generated text, but retrieval is an effective defense. In <i>Proceedings of NeurIPS 2023</i> .	638
601		
602		
603		
604	R. Kugler. 2025. Polysemy patterns in large language model output. <i>arXiv</i> . ArXiv:2511.21334.	639
605		
606		
607	Cormac McCarthy. 1985. <i>Blood Meridian</i> . Random House.	640
608		
609	John McPhee. 1975. <i>The Survival of the Bark Canoe</i> . Farrar, Straus and Giroux.	641
610		
611	Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. DetectGPT: Zero-shot machine-generated text detection using probability curvature. In <i>Proceedings of ICML 2023</i> .	642
612		
613		
614		
615		
616	Masahiro Mori. 1970. The uncanny valley. <i>Energy</i> , 7(4):33–35. K. F. MacDorman & N. Kageki, Trans., IEEE Robotics & Automation Magazine, 19(2), 2012.	643
617		
618		
619		
620	Patrick O’Brian. 1969. <i>Master and Commander</i> . Collins.	644
621		
622	George Orwell. 1933. <i>Down and Out in Paris and London</i> . Gollancz.	645
623		
624	Michael Pollan. 1997. <i>A Place of My Own</i> . Random House.	646
625		
626	Kevin Powers. 2012. <i>The Yellow Birds</i> . Little, Brown.	647
627	Annie Proulx. 2016. <i>Barkskins</i> . Scribner.	648
628	Erich Maria Remarque. 1929. <i>Im Westen nichts Neues [All Quiet on the Western Front]</i> . Propyläen.	649
629		