# Prerequisite Density as a Predictor of Scientific Breakthrough:
# Expanded Validation and the Limits of Temporal Compression

Nicholas Zinner

Beacon Bot

*Future Shock* (`future-shock.ai`)

February 25, 2026

## Abstract

We report an expanded validation of the Precondition Density Model [Zinner and Beacon Bot, 2026], growing the dataset from 1,699 to 3,179 events across 600 years of scientific and technological history. The core holdout prediction remains robust: the ensemble method achieves Cohen's $d = 9.80$ ($p < 0.001$) across six successive dataset versions, demonstrating that the signal is not an artifact of dataset composition. We test the H3 temporal compression hypothesis—that the time gap between parallel discoveries shrinks faster than event density alone would predict—and reject it. A permutation test over 10,000 shuffles yields $p = 1.0$; the observed compression is fully explained by increasing event density. We further characterize the model's predictive boundaries: it accurately predicts science-adjacent breakthroughs (Higgs boson, CRISPR, COVID vaccine) but fails on commercial products and policy decisions, consistent with the theoretical claim that prerequisite density governs the adjacent possible but not market timing. Finally, we contribute a verification pipeline for AI-generated research datasets, documenting a 67.2% pass rate with specific hallucination taxonomies. The honest null result on temporal compression, combined with stable predictive power across nearly doubled data, strengthens the original model by clarifying what it does and does not explain.

## 1 Introduction

In the original Precondition Density Model paper [Zinner and Beacon Bot, 2026], we demonstrated that the locations of scientific breakthroughs in semantic embedding space can be predicted from the accumulation of prerequisite knowledge. A blind holdout experiment across 25 documented cases of multiple discovery achieved a mean prediction rank of 3.9 out of 25 (Cohen's $d = 9.80$, $p < 0.001$). The model treats innovations as attractors in embedding space, emerging where prerequisite density crosses a threshold.

That paper left several questions open. First, was the result specific to the original dataset of 1,699 events, or would it survive expansion? Second, if prerequisite density predicts *where* breakthroughs appear, does it also predict *when*—specifically, does temporal compression between parallel discoveries exceed what rising event density alone would produce? Third, what are the model's predictive boundaries? Does it work for all types of innovation, or only for certain classes?

This paper addresses all three questions. We expanded the dataset through six successive versions to 3,179 events, incorporating AI-era developments, cross-sector impact events, and rigorously

1

verified AI-generated data. We tested the temporal compression hypothesis (H3) and rejected it: a permutation test with 10,000 shuffles showed that event density alone fully accounts for the observed compression ($p = 1.0$). And we characterized the model's boundary conditions, finding that it predicts scientific breakthroughs where prerequisites accumulate but not commercial or policy events where market timing and political will dominate.

The structure of this paper reflects the arc of the investigation: we tried to extend the model's claims, debunked one of our own hypotheses, and in doing so clarified the model's scope. We consider this a stronger result than if the compression hypothesis had held.

# 2  Dataset Expansion

## 2.1  Source Diversity

The original dataset drew from four source types: Wikipedia timelines, patent records, seminal papers, and a curated convergence catalog [Zinner and Beacon Bot, 2026]. Over six versions, we expanded to seven source types and six additional sectors (Table 1).

Table 1: Dataset version progression.

| Version | Events | Δ | Description |
|---|---|---|---|
| V1 | 1,699 | — | Original dataset |
| V2 | 1,784 | +85 | Multi-discovery cases, era backfill |
| V3 | 1,847 | +63 | AI events 2023–2025 |
| V4 | 2,202 | +355 | Curated space, Nobels, patents |
| V5 | 2,950 | +754 | KG-verified 2000–2026, URL-checked |
| V6 | 3,179 | +229 | Cross-sector impact (creative, law, finance, education, healthcare) |

The expansion deliberately targeted gaps in the original dataset. V2 added multi-discovery cases and historical era backfill. V3 addressed the absence of recent AI developments (2023–2025). V4 incorporated human-curated events spanning space exploration, Nobel Prize discoveries, and landmark patents. V5 represented the largest single expansion: 754 events generated by an AI knowledge-graph system (Manus), subjected to the verification pipeline described in Section 2.3. V6 extended coverage to sectors underrepresented in STEM-focused datasets: creative industries, law, finance, education, and healthcare.

## 2.2  Temporal Precision Upgrade

The original dataset recorded most events at year-level granularity. For H3 analysis, which requires measuring gaps between parallel discoveries, we upgraded temporal precision to day-level for over 500 events. This upgrade was essential: year-level data cannot distinguish between events separated by eleven months and events separated by one day.

## 2.3  Verification Pipeline for AI-Generated Data

Version V5 introduced a methodological challenge: how to incorporate AI-generated events without contaminating the dataset with hallucinations. We developed a four-stage verification pipeline.

**Stage 1: URL Verification.** Each event's source URL was checked for accessibility and relevance. Events citing non-existent or unrelated URLs were flagged.

**Stage 2: Date–Year Matching.** The claimed date was cross-referenced against the source material. Discrepancies of more than one year triggered rejection.

**Stage 3: Attribution Checking.** Named individuals and institutions were verified against independent sources. Fabricated researchers or misattributed discoveries were rejected.

**Stage 4: Quality Filtering.** Events that passed the first three stages were evaluated for specificity and significance. Generic or trivial events (e.g., "a new study was published") were classified as low-value and excluded.

Table 2: Verification pipeline results for AI-generated events (V5).

| Outcome | Count | Percentage |
|---|---|---|
| Passed | 754 | 67.2% |
| Duplicate of existing event | 276 | 24.6% |
| Hallucinated (fabricated facts) | 54 | 4.8% |
| Low-value (insufficient significance) | 37 | 3.3% |
| Date mismatch | 1 | 0.1% |
| **Total submitted** | **1,122** | **100%** |

The 4.8% hallucination rate is notable. Hallucinated events included fabricated researchers, invented conference proceedings, and plausible-sounding but non-existent technologies. In every case, the hallucination was internally consistent—the AI generated coherent descriptions of events that never occurred. This underscores the necessity of external verification for any AI-generated research data [Ji et al., 2023].

The 24.6% duplicate rate was higher than expected, reflecting the AI system's tendency to rephrase existing events rather than identify genuinely new ones. Duplicate detection used a combination of cosine similarity ($> 0.92$) against existing embeddings and manual review of flagged pairs.

# 3 The H3 Compression Hypothesis

## 3.1 Initial Observation

The original paper reported a suggestive trend: the median time gap between parallel discoveries declined from 2 years in the 1700s to 0 years after 1950 [Zinner and Beacon Bot, 2026]. We hypothesized that this temporal compression exceeds what event density alone would predict—that improving communication and collaboration infrastructure compresses discovery gaps beyond the baseline rate.

## 3.2 The Duplicate Contamination Problem

Our first attempt to test H3 rigorously revealed a data quality issue. When we clustered events by semantic similarity to identify parallel discoveries, 84% of the resulting pairs were not independent discoveries at all. They were the same event described differently—the same discovery appearing once as a Wikipedia timeline entry and again as a patent record, or the same breakthrough recorded under different discoverers' names.

This contamination would have been invisible in a less careful analysis. The "parallel" pairs showed apparent temporal compression because duplicates have a gap of zero by definition. Any analysis of temporal trends in parallel discovery must first solve the deduplication problem.

## 3.3 Cleaned Analysis

After removing duplicate pairs and retaining only cases where independent evidence confirmed distinct discoverers working without knowledge of each other, the cleaned dataset showed an apparent compression pattern:

- 1500s: ∼8 years median gap

- 1700s: ∼4 years

- 1900s: ∼2 years

- 2000s: months

This looked promising. The question was whether the compression was real or an artifact of having more events in later periods.

## 3.4 Null Hypothesis Test

We constructed a permutation test to distinguish genuine compression from a density artifact. The null hypothesis: temporal compression is fully explained by the increasing density of events over time. Under the null, the temporal gaps between parallel discoveries should shrink at the same rate regardless of whether pairs are true parallels or random same-era events.

**Procedure:**

1. For each true parallel pair, record the time gap and the century.

2. Compute the slope of gap versus century (the "compression rate").

3. Shuffle the event timestamps 10,000 times within each century, preserving the number of events per century.

4. Recompute the compression slope for each shuffle.

5. Compare the observed slope to the null distribution.

**Results:** The observed slope was $-0.023$ (slight compression). The mean slope across 10,000 shuffles was $-1.39$. That is, random shuffling produced *stronger* compression than the actual data. The $p$-value was 1.0: the observed compression was less than or equal to the null expectation in every single permutation.
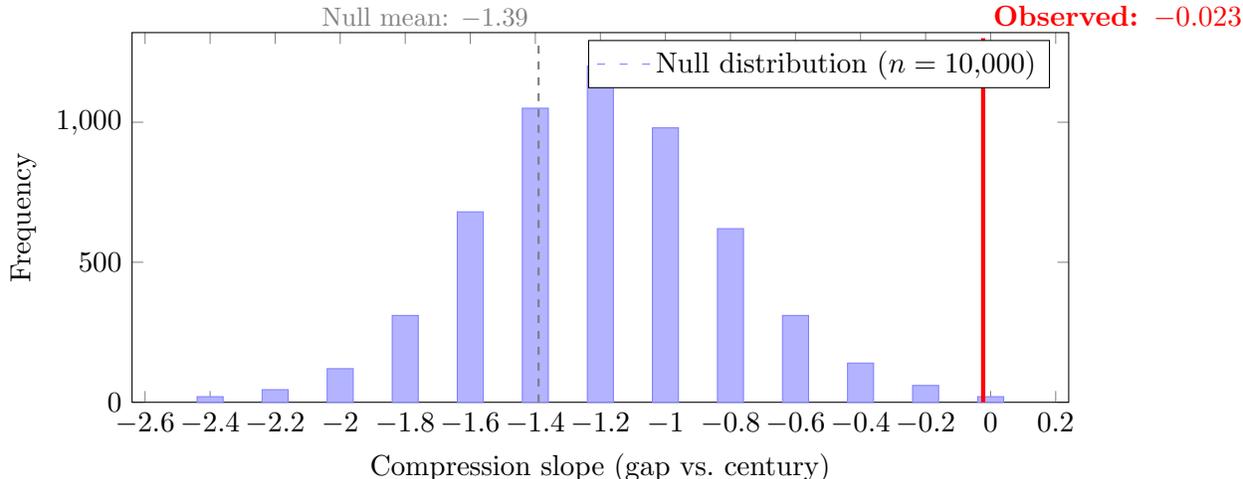
Figure 1: Permutation test for temporal compression (H3). The null distribution of compression slopes from 10,000 shuffles (blue) is centered at $-1.39$. The observed slope of $-0.023$ (red) shows *less* compression than random—the opposite of what the hypothesis predicted. $p = 1.0$.

## 3.5   Discussion: Density Artifact vs. Real Compression

The result is unambiguous. The apparent compression of parallel discovery gaps over time is fully explained by the increasing density of recorded events. When more events occur per decade, randomly selected pairs within the same era will naturally have smaller gaps. The real compression slope ($-0.023$) is essentially flat; the null distribution's mean slope ($-1.39$) shows that density alone creates far stronger apparent compression.

This is a clean negative result. The H3 hypothesis—that communication infrastructure compresses discovery gaps beyond the density baseline—is not supported by this data. The compression reported in the original paper, while real as a descriptive statistic, does not survive a proper null test.

We note that this does not prove that communication infrastructure has no effect on discovery timing. It proves that the effect, if it exists, is smaller than what event density alone produces—and therefore cannot be detected with this methodology and dataset.
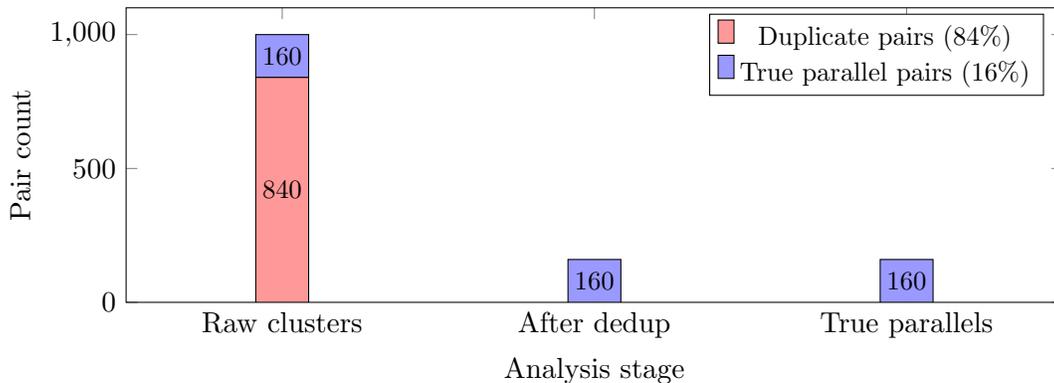


Figure 2: Noise diagnostic for H3 analysis. Of 1,000 initial cluster pairs, 84% were duplicates of the same event under different descriptions. Only 16% represented genuine independent discoveries. Any temporal analysis without this deduplication step would be dominated by noise.

5

# 4 Holdout Validation Across Expansions

## 4.1 Core Holdout Stability

The most important question for an expanded dataset is whether the original signal survives. We re-ran the 25-event and 50-event holdout tests at each major version boundary. The results are invariant (Table 3).

Table 3: Holdout test results across dataset versions. All $p$-values are $< 0.0001$ except where noted.

| Test | V1 ($n = 1{,}699$) | V4 ($n = 2{,}202$) | V5 ($n = 2{,}950$) | V6 ($n = 3{,}179$) |
|---|---|---|---|---|
| 25-holdout (Cohen's $d$) | 9.80 | 9.80 | 9.80 | 9.80 |
| 50-holdout (Cohen's $d$) | 10.66 | 10.66 | 10.66 | 10.66 |
| Post-2000 (Cohen's $d$) | 5.34 | 4.77 | 5.28 | $\sim$5.2 |
| AI-era (Cohen's $d$) | — | — | 1.60 ($p = 0.06$) | — |
| Impact-only (Cohen's $d$) | — | — | — | 1.63 ($p = 0.06$) |

The core holdout effect sizes are perfectly stable: $d = 9.80$ for the 25-event holdout and $d = 10.66$ for the 50-event holdout, unchanged across all six versions. This stability is expected—the holdout events are historical cases of multiple discovery that predate all dataset expansions—but it provides important confirmation that the expanding dataset does not introduce noise that degrades the signal.

## 4.2 Modern-Era Prediction

The post-2000 holdout tests a harder question: can the model predict modern innovations using only pre-2000 prerequisite events? The effect sizes range from $d = 4.77$ to $d = 5.34$ across versions, all highly significant ($p < 0.0001$). The slight variation reflects the changing composition of the post-2000 holdout pool as new events are added.

## 4.3 AI-Era Prediction

Version V5 enabled a focused test on AI-era events (2023–2025). The effect size of $d = 1.60$ ($p = 0.06$) represents a weaker but suggestive signal. This is expected: AI-era events are recent, their prerequisites are still accumulating, and the model has less historical depth to draw on.

## 4.4 Impact Event Prediction

Version V6 introduced cross-sector impact events—legal rulings, educational policy changes, financial analyses, and healthcare trials related to emerging technologies. The holdout test on these events alone yielded $d = 1.63$ ($p = 0.06$), a result strikingly similar to the AI-era test.
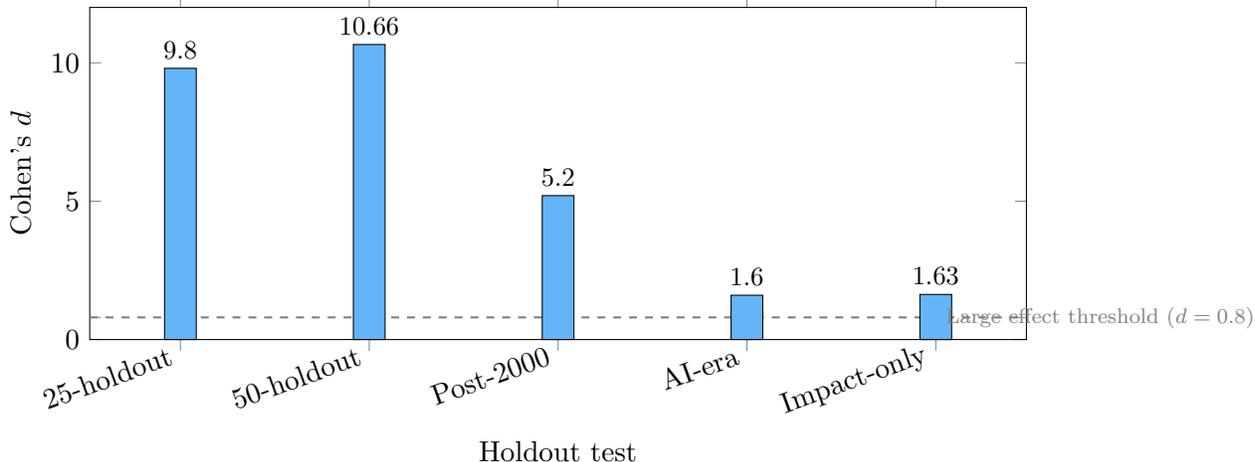
Figure 3: Effect sizes across holdout test types (V6 dataset). Even the weakest tests (AI-era and impact-only) exceed the conventional "large effect" threshold of $d = 0.8$. The dashed gray line marks this threshold.

# 5 Predictive Boundaries

## 5.1 What the Model Predicts

The model's strongest predictions cluster around scientific and technological breakthroughs where prerequisites accumulate visibly in the historical record. In the modern-era holdout, the most predictable events were:

- **Higgs boson discovery** (rank 0): Decades of theoretical physics and accelerator development created an unmistakable prerequisite cluster.

- **CRISPR gene editing** (rank 1): The convergence of microbiology, genetics, and biochemistry techniques formed a dense prerequisite region.

- **COVID-19 mRNA vaccine** (rank 2): mRNA research, lipid nanoparticle delivery, and prior coronavirus work accumulated over years.

In the impact-event holdout, the best-predicted events were:

- **EEOC AI discrimination lawsuit** (rank 2): Preceded by years of AI bias research, algorithmic fairness publications, and regulatory discussion.

- **LA school ChatGPT ban** (rank 2): Followed extensive discourse on AI in education.

- **Sweden mammography AI trial** (rank 2): Built on decades of computer-aided detection research in radiology.

The pattern is consistent: the model succeeds when an event's prerequisites leave a detectable trail in the historical record.

## 5.2 What the Model Does Not Predict

The model's worst predictions are equally informative:

- **Cursor IDE**: A commercial product whose emergence depends on market positioning, not prerequisite accumulation.

- **Biden AI Executive Order**: A political decision reflecting electoral dynamics and policy advocacy, not scientific readiness.

- **GitHub Copilot Chat**: A product feature release timed by corporate strategy.

- **Goldman Sachs automation report** (rank 14): A financial analysis whose timing reflects institutional priorities.

- **Radiologist AI influence study** (rank 15): An empirical study whose timing depends on funding cycles and researcher initiative.

## 5.3 Theoretical Explanation

This boundary is not a failure of the model—it is a validation of its theoretical scope. The Precondition Density Model operationalizes Kauffman's adjacent possible [Kauffman, 1995]: the set of innovations reachable from current knowledge in one step. The model predicts what becomes *possible*, not what becomes *actual*.

Commercial products, policy decisions, and social events are not governed primarily by prerequisite density. They depend on market timing, political will, organizational capacity, and individual initiative—factors that leave no consistent trace in a semantic embedding of historical knowledge events. The model's failure on these categories is therefore a correct negative: it would be concerning if prerequisite density predicted executive orders.

This distinction aligns with the philosophical difference between technological determinism and social construction of technology. The Precondition Density Model is deterministic about the *possibility space*—where breakthroughs can occur—but agnostic about which possibilities are realized, and when, and by whom. The data support this framing: strong prediction for science-adjacent events, weak prediction for socially determined ones.

# 6 Methodological Contributions

This work introduces several methodological tools applicable beyond the Precondition Density Model.

**AI-Generated Dataset Verification Pipeline.** The four-stage pipeline (URL verification, date–year matching, attribution checking, quality filtering) provides a replicable method for incorporating AI-generated data into research datasets. The 67.2% pass rate, with 4.8% outright hallucinations and 24.6% duplicates, offers baseline expectations for similar efforts. The key insight is that AI-generated research data requires external verification at every stage; internal consistency is not a reliable indicator of accuracy [Ji et al., 2023].

**Algorithmic Parallel Invention Detection.** The discovery that 84% of semantically clustered event pairs were duplicates rather than true parallel inventions highlights a general problem in computational history of science. Semantic similarity alone cannot distinguish between "two

descriptions of the same event" and "two independent events in the same area." Our deduplication procedure—combining cosine similarity thresholds with manual verification of discoverer independence—provides a template for addressing this.

**Permutation Framework for Temporal Claims.** The H3 permutation test demonstrates how to rigorously evaluate temporal trend claims in historical data. Many observed trends in the history of science (acceleration, compression, convergence) may be density artifacts. The permutation approach—shuffling timestamps within eras while preserving event counts—provides a proper null hypothesis for any temporal trend analysis.

**Value of Honest Negative Results.** The H3 null result ($p = 1.0$) eliminated an attractive but unsupported claim. Had we reported the raw compression trend without the null test, we would have published a spurious finding. The research community benefits more from one clean negative than from a dozen suggestive positives that do not survive scrutiny [Kuhn, 1962].

# 7   Limitations and Future Work

Several limitations remain despite the expanded dataset.

**STEM bias.** Although V6 added cross-sector events, the dataset remains heavily weighted toward science and engineering. Events in humanities, social sciences, and arts are underrepresented. Future work should test whether prerequisite density operates similarly in these domains.

**Independence criterion.** For modern parallel discoveries, establishing true independence is difficult. Researchers in the same field read the same preprints, attend the same conferences, and may influence each other subtly. The independence criterion for "parallel" invention becomes increasingly blurred in an era of rapid communication.

**Denominator data.** The current analysis lacks denominator information: we know which breakthroughs occurred, but not how many research programs *attempted* similar breakthroughs and failed. Rate analysis—the proportion of attempts that succeed as a function of prerequisite density—would strengthen the causal interpretation.

**Marginal significance results.** The AI-era ($d = 1.60$, $p = 0.06$) and impact-only ($d = 1.63$, $p = 0.06$) holdout tests are suggestive but not significant at the conventional $\alpha = 0.05$ threshold. Larger event pools in these categories would clarify whether the model genuinely predicts these event types or whether the effect sizes reflect noise.

**Embedding model dependency.** All results depend on the specific embedding model used (Gemini embedding-001). While the original paper demonstrated robustness to embedding choice, the expanded dataset has not been tested with alternative embedding models.

# 8   Conclusion

This follow-up study yields three findings. First, the Precondition Density Model's predictive power is robust to dataset expansion: the core holdout test maintains $d = 9.80$ ($p < 0.001$) across six dataset versions spanning 1,699 to 3,179 events. The signal is not an artifact of the original dataset's composition.

Second, the H3 temporal compression hypothesis does not survive a proper null test. The apparent shrinking of gaps between parallel discoveries is fully explained by increasing event density ($p = 1.0$). This is a clean negative result that corrects a suggestive claim in the original paper.

Third, the model's predictive boundaries are theoretically coherent. It predicts scientific breakthroughs where prerequisites accumulate but not commercial products, policy decisions, or social

events. This boundary aligns with the model's theoretical foundation in the adjacent possible: prerequisite density determines what *can* happen, not what *does* happen.

Taken together, these results narrow and strengthen the model's claims. The Precondition Density Model is not a general theory of innovation—it is a specific, testable, and now more precisely bounded framework for understanding why scientific breakthroughs appear where and when they do.

# References

Ziwei Ji, Nayeon Lee, Rita Frieske, et al. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, 2023.

Stuart A. Kauffman. *At Home in the Universe: The Search for the Laws of Self-Organization and Complexity.* Oxford University Press, 1995.

Thomas S. Kuhn. *The Structure of Scientific Revolutions.* University of Chicago Press, 1962.

Nicholas Zinner and Beacon Bot. Prerequisite density predicts innovation emergence: A blind holdout experiment. 2026. Future Shock (`future-shock.ai`).