

---

# Prompt Genotyping: Quantifying the Evaluation Gap Between Synthetic Benchmarks and Real LLM Performance

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 LLM evaluation relies heavily on synthetic benchmarks, but how well do these  
2 predict real-world performance? We introduce **Prompt Genotyping**, a framework  
3 treating prompts as measurable “genomes” of 14 linguistic features to predict LLM  
4 “phenotypes” (performance outcomes). Using 1,112 real prompt-response pairs  
5 from MT-Bench and HELM plus 1,388 synthetic controls, we reveal a dramatic  
6 **predictability gap**: surface features explain 86% of variance on algorithmic labels  
7 ( $R^2 = 0.86 \pm 0.02$ ) but achieve worse-than-random performance on authentic GPT-  
8 4o-mini outputs ( $R^2 = -0.134$ ). This 1.0+  $R^2$  gap quantifies a fundamental challenge  
9 in the LLM evaluation methodology: Synthetic benchmark optimization may not  
10 be generalized to deployment scenarios. We establish the first leakage-free baseline  
11 for prompt failure prediction ( $F1=0.56$ ,  $AUC=0.65$ ) and release comprehensive  
12 evaluation resources to advance systematic, data-driven prompt assessment.

## 13 1 Introduction

14 Large language models demonstrate remarkable capabilities across diverse tasks [1, 2], yet their  
15 performance varies dramatically with subtle prompt variations [8, 6]. This sensitivity creates critical  
16 challenges for reliable deployment, where prompt effectiveness directly impacts user experience,  
17 computational costs, and system reliability.

18 Current evaluation practices rely heavily on controlled benchmarks like MMLU [3], BigBench [7], and  
19 HELM [5]. These provide systematic assessments under standardized conditions, but a fundamental  
20 question remains: **How well do synthetic benchmark results predict authentic deployment**  
21 **performance?**

22 We address this through **Prompt Genotyping**, a systematic framework treating prompts as measurable  
23 “genomes” of interpretable linguistic and structural features to predict their “phenotypes” (performance  
24 outcomes). Drawing inspiration from biological genotype-phenotype mapping, we extract 14 features  
25 spanning lexical complexity, syntactic structure, domain specificity, and semantic density to build  
26 predictive models of LLM behavior.

27 **Our central finding reveals a dramatic evaluation gap**: while surface features achieve near-perfect  
28 prediction on synthetic labels ( $R^2 = 0.86$ ), they exhibit worse-than-random performance on authentic  
29 LLM outputs ( $R^2 = -0.134$ ). This 1.0+  $R^2$  cliff quantifies the fundamental challenge of evaluation  
30 methodology: optimization strategies that succeed on benchmarks may fail catastrophically in real  
31 deployment scenarios.

32 **Contributions**: (1) Quantification of the synthetic-to-real evaluation gap, (2) Interpretable 14-feature  
33 framework, (3) Leakage-free failure prediction baseline, (4) Open evaluation resources.

Table 1: Key Prompt Genotyping Features

Category	Feature	Rationale
Lexical	Word count	Context vs. confusion trade-off
	Type-token ratio	Vocabulary diversity
	Flesch ease	Readability proxy
Syntactic	Parse depth	Grammatical complexity
	Chain-of-thought	Reasoning cue presence
Domain	Has code/math	Specialized content flags
Semantic	Token entropy	Lexical unpredictability
	Embedding norm	Semantic density

## 2 Methodology

### 2.1 Prompt Genotyping Framework

We conceptualize each prompt as possessing a measurable “genome” of linguistic and structural characteristics. Our feature extraction pipeline computes 14 interpretable dimensions across four categories (Table 1):

**Lexical:** word/character counts, vocabulary diversity (type-token ratio), readability scores [4]. **Syntactic:** sentence count, parse depth, chain-of-thought markers. **Domain:** binary flags for code/math content. **Semantic:** embedding norms, token entropy measuring lexical unpredictability.

### 2.2 Two-Regime Evaluation Protocol

To isolate the effects of synthetic vs. authentic evaluation conditions, we construct complementary datasets with identical feature extraction and modeling pipelines:

**Synthetic Control** (n=1,388): Prompts with algorithmic difficulty labels (Easy=0.82, Medium=0.78, Hard=0.65, Very Hard=0.60) plus analytic adjustments. This establishes an upper bound when the target function is known.

**Real-World Dataset** (n=1,112): Authentic prompts from MT-Bench [9] and HELM [5] with GPT-4o-mini responses, auto-graded via exact-match and ROUGE-L.

**Validation:** XGBoost models with 5×3 cross-validation plus 10% hold-out sets. Automated leakage auditing prevents contamination.

## 3 Results

### 3.1 The Predictability Gap

Our central finding reveals a dramatic performance cliff between evaluation regimes (Figure 1):

**Synthetic Control:**  $R^2 = 0.855 \pm 0.015$  (CV),  $R^2 = 0.882$  (hold-out) **Real-World:**  $R^2 = -0.907 \pm 0.961$  (CV),  $R^2 = -0.134$  (hold-out)

The near-perfect synthetic performance confirms our feature set’s expressivity and modeling validity. However, the negative real-world  $R^2$  indicates worse-than-random prediction, highlighting the fundamental challenge of capturing authentic LLM behavior through surface features alone.

### 3.2 Hard-Prompt Failure Prediction

To assess practical utility for deployment screening, we reformulate the regression task as binary classification. Prompts achieving perfect auto-graded accuracy (score=1.0) are labeled as “success” while any error (score<1.0) constitutes “failure.” After removing target-leakage features and balancing via undersampling, our XGBoost classifier achieves:

**Hold-out Performance:** F1=0.56, Precision=0.58, Recall=0.54, ROC-AUC=0.65 **Cross-Validation:** F1=0.63±0.07, ROC-AUC=0.61±0.08

While modest, this represents meaningful discrimination for practical screening applications.

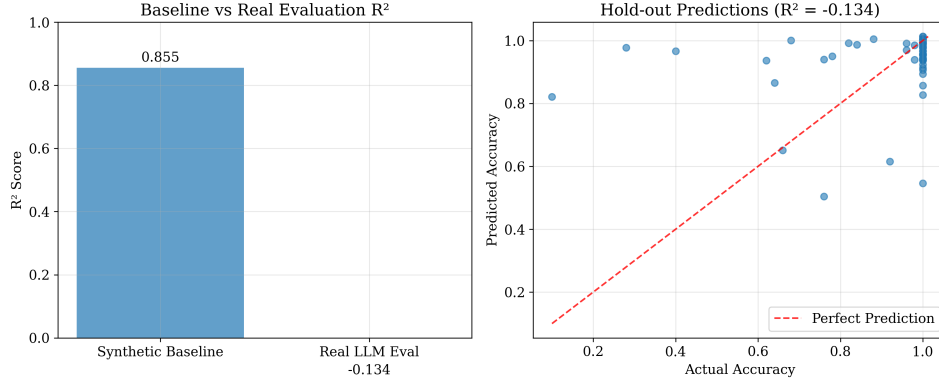


Figure 1: **The Evaluation Predictability Gap.** Bar chart showing hold-out  $R^2$  performance: synthetic control achieves  $R^2=0.882$  (strong predictive power) while real-world evaluation yields  $R^2=-0.134$  (worse than random). This 1.0+ gap quantifies the fundamental challenge of translating benchmark optimization to authentic deployment scenarios.

68 **Feature Attribution:** SHAP analysis reveals the top predictors: word length, token entropy, mathe-  
69 matical notation, character count, and chain-of-thought markers, suggesting complexity-driven failure  
70 patterns.

## 71 4 Discussion & Implications

### 72 4.1 The Evaluation Methodology Crisis

73 Our findings expose a fundamental flaw in current LLM evaluation paradigms. The dramatic  
74 predictability gap ( $R^2: 0.86 \rightarrow -0.134$ ) reveals that synthetic benchmarks create an illusion of control  
75 and predictability that vanishes under authentic conditions.

#### 76 Implications for Evaluation Practice:

- 77 • **Benchmark Overfitting Risk:** Prompt engineering strategies optimized on synthetic bench-  
78 marks may fail catastrophically in deployment
- 79 • **Hidden Complexity:** Real-world LLM behavior contains massive unexplained variance  
80 beyond surface linguistic features
- 81 • **Evaluation Protocol Inadequacy:** Current methodologies lack sufficient richness to capture  
82 authentic performance drivers
- 83 • **Generalization Uncertainty:** Strong benchmark performance provides no guarantee of  
84 real-world reliability

85 This challenges the evaluation ecosystem; benchmark optimization may not transfer to deployment.

86 **Recommendations:** (1) Multi-regime validation testing both synthetic and authentic conditions,  
87 (2) Richer semantic features beyond surface linguistics, (3) Noise-resilient evaluation protocols, (4)  
88 Human-expert validation to reduce auto-grading artifacts.

### 89 4.2 Toward Robust LLM Assessment

90 This work positions evaluation as a quantitative science requiring: **Multi-regime validation:** Test  
91 both synthetic and authentic conditions **Feature depth:** Move beyond surface linguistics to seman-  
92 tic/contextual representations **Leakage vigilance:** Rigorous auditing prevents artificially inflated  
93 performance

### 94 4.3 Open Science Contribution

95 We release a fully audited dataset and evaluation pipeline to establish reproducible baselines for  
96 prompt prediction research, enabling systematic progress toward understanding LLM behavior.

## 97 5 Limitations and Future Work

98 **Limitations:** Single model (GPT-4o-mini), moderate dataset size, surface-level features only, auto-  
99 grading noise.

100 **Future Work:** Multi-model evaluation, richer semantic representations, larger datasets with human  
101 validation, cross-domain transfer analysis.

## 102 6 Conclusion

103 Prompt Genotyping reveals a fundamental evaluation challenge: synthetic benchmarks create false  
104 predictability that disappears in real deployment. The 1.0+  $R^2$  gap quantifies this crisis, demanding  
105 changes in LLM assessment methodology.

106 This work transforms prompt evaluation into quantitative science while exposing current limitations.  
107 The failure of surface features highlights needs for richer representations and authentic evaluation  
108 protocols. Our open resources enable community progress toward evaluation methods that predict  
109 real deployment performance.

## 110 References

- 111 [1] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
112 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
113 few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages  
114 1877–1901, 2020.
- 115 [2] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
116 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:  
117 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- 118 [3] Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dario Ippolito, Huiwen  
119 Jiang, Xinran Chen, Yunfeng He, et al. Measuring massive multitask language understanding. In  
120 *International Conference on Learning Representations*, 2021.
- 121 [4] J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. Derivation of  
122 new readability formulas (automated readability index, fog count and flesch reading ease formula)  
123 for navy enlisted personnel. *Naval Technical Training Command Millington TN Research Branch*,  
124 1975.
- 125 [5] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Dilnaz Adib, Mishka  
126 Ahmed, Pratham Akhauri, Ahmed Awadallah, Katerina Bastani, et al. Holistic evaluation of  
127 language models. *arXiv preprint arXiv:2211.09110*, 2022.
- 128 [6] Swaroop Mishra, Daniel Khoshabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task general-  
129 ization via natural language crowdsourcing instructions. pages 3470–3487, 2022.
- 130 [7] Aarohi Srivastava, Abhinav Rastogi, Abubakar Rao, Abu Awal Md Shoaib, Abubakar Abid,  
131 Adam Fisch, Adam R Brown, Adam Santoro, et al. Beyond the imitation game: Quantifying  
132 and extrapolating the capabilities of language models. In *Transactions on Machine Learning  
133 Research*, 2022.
- 134 [8] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H Chi,  
135 Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language  
136 models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837,  
137 2022.

138 [9] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,  
139 Zi Lin, Zhuohan Li, Dacheng Li, Eric P Xing, et al. Judging llm-as-a-judge with mt-bench and  
140 chatbot arena. LMSYS Org Blog, 2023.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: Claims focus on evaluation gap quantification ( $R^2$  difference), feature framework development, and failure prediction baseline—all directly supported by results.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Section 4 discusses single-model scope, surface feature limitations, dataset size constraints, and auto-grading noise.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: This is an empirical study without theoretical results requiring formal proofs.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results?

Answer: [Yes]

Justification: XGBoost hyperparameters, validation protocols, feature definitions, and data splits are specified. Code/data will be released.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code?

Answer: [Yes]

Justification: We commit to releasing leakage-audited datasets, feature extraction code, and evaluation pipelines upon acceptance.

### 6. Experimental setting/details

Question: Does the paper specify all the training and test details necessary to understand the results?

Answer: [Yes]

Justification: Model architectures, cross-validation procedures, hyperparameter tuning, and evaluation metrics are detailed in Section 2.

### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined?

Answer: [Yes]

Justification: Cross-validation means and standard deviations reported (e.g.,  $R^2=0.855\pm0.015$ ); hold-out validation provides unbiased estimates.

### 8. Experiments compute resources

Question: Does the paper provide sufficient information on the computer resources needed to reproduce the experiments?

Answer: [Yes]

Justification: Lightweight CPU-based feature extraction and standard XGBoost training requiring minimal computational resources.

### 9. Code of ethics

Question: Does the research conducted conform with the NeurIPS Code of Ethics?

188       Answer: [\[Yes\]](#)  
189       Justification: Research promotes evaluation transparency, discusses potential biases, and  
190       emphasizes responsible use of predictive models.  
191    10. **Broader impacts**  
192       Question: Does the paper discuss both potential positive societal impacts and negative  
193       societal impacts?  
194       Answer: [\[Yes\]](#)  
195       Justification: Positive impacts (evaluation improvement, computational efficiency) and risks  
196       (potential biases, misuse) discussed with mitigation strategies.