

Augmenting Data Quality and Model Reliability in Large-Scale Language and Code Models: A Hybrid Framework for Evaluation, Pretraining, and Retrieval-Augmented Techniques

John M. Langley

Department of Computer Science, University of Edinburgh

Article received: 01/09/2025, Article Revised: 09/09/2025, Article Accepted: 15/09/2025

© 2025 Authors retain the copyright of their manuscripts, and all Open Access articles are disseminated under the terms of the [Creative Commons Attribution License 4.0 \(CC-BY\)](https://creativecommons.org/licenses/by/4.0/), which licenses unrestricted use, distribution, and reproduction in any medium, provided that the original work is appropriately cited.

ABSTRACT

Background: The rapid expansion of large language models (LLMs) and code-generative models has transformed research and industry practices across natural language processing, software engineering, and data-driven decision-making. Yet, the increasing scale of datasets and repeat data exposure introduces complex challenges in data quality, training set augmentation, model reliability, and downstream evaluation (Ding, 2019; Hernandez et al., 2022). Prior work has examined whether large-scale datasets are necessary for self-supervised pretraining (El-Nouby et al., 2021), explored the landscape of open-source engineering efforts (Han et al., 2021), and surveyed retrieval-augmented language models (Hu & Lu, 2024). However, integrated frameworks that connect data augmentation, rigorous quality validation, and evaluation tailored to LLMs remain underdeveloped.

Objective: This article proposes and thoroughly elaborates a hybrid, academically rigorous framework that synthesizes data augmentation best practices, AI-augmented data quality validation, retrieval-augmented model design, and robust evaluation metrics for LLMs and code models. It aims to bridge theoretical foundations with practical design choices and provide an interpretive, evidence-based roadmap for researchers and practitioners.

Methods: We synthesize perspectives from empirical case studies on training-data augmentation (Ding, 2019), scaling laws and interpretability of repeated data (Hernandez et al., 2022), debates on dataset scale for self-supervision (El-Nouby et al., 2021), and contemporary LLM evaluation challenges (Gao et al., 2024). From these sources we construct a layered methodology: (1) Source-level data curation and provenance tracing informed by record linkage principles (Herzog et al., 2007); (2) augmentation strategies balancing synthetic and human-authored instances (Ding, 2019); (3) hybrid validation combining rule-based checks and LLM-assisted anomaly detection (Malviya & Parate, 2025); (4) design patterns for retrieval-augmented pipelines (Hu & Lu, 2024); and (5) a multi-faceted evaluation protocol incorporating statistical, qualitative, and LLM-based evaluators (Gao et al., 2024; Wang et al., 2023).

Results: The resulting framework identifies trade-offs between dataset scale and diversity, quantifies danger zones where repeated data leads to overfitting or miscalibration (Hernandez et al., 2022), and recommends concrete validation procedures to detect provenance drift, duplication bias, and label noise. We also specify evaluation batteries for code synthesis models and medical-diagnostic LLM comparisons using ensemble judge designs (Fried et al., 2022; Caruccio et al., 2024).

Conclusions: By integrating augmentation, validation, retrieval, and evaluation, the framework supports more reliable, auditable, and interpretable LLM deployments. Theoretical implications include revised perspectives on necessary dataset scale, formalization of hybrid validation agents, and suggested directions for future empirical work. This synthesis provides a substantive foundation for reproducible research and practical deployment strategies for LLMs and code models.

Keywords: Large language models; data augmentation; data quality validation; retrieval-augmented models; evaluation metrics; model reliability; code generation.

INTRODUCTION

The revolution in deep learning-driven language technologies during the past decade has been characterized by two simultaneous trends: explosive growth in model capacity and the proliferation of large-scale training datasets. Large language models (LLMs) and generative code models have demonstrated emergent capabilities across a range of tasks, from open-ended text generation to code infilling and synthesis (Fried et al., 2022; Fan et al., 2023). At the same time, research has surfaced nuanced concerns about the relationship between dataset scale, repeated exposure of data during training, and model interpretability and reliability (El-Nouby et al., 2021; Hernandez et al., 2022). The interdependence among training data augmentation, provenance-aware data quality validation, and modern evaluation approaches has become a core axis of research and operational practice.

A critical impetus for the present work is the evidence that naive dataset enlargement—simply adding more tokens, documents, or code repositories—does not guarantee performance improvements commensurate with costs and may introduce systemic problems such as data duplication, label leakage, provenance drift, and misaligned evaluation practices (Ding, 2019; Hernandez et al., 2022). The question “Are large-scale datasets necessary for self-supervised pre-training?” embodies this tension by interrogating whether continued growth is the primary driver of capability gains or if smarter curation and augmentation can achieve similar or better outcomes with fewer resources (El-Nouby et al., 2021). From an applied perspective, sectors such as insurance, healthcare, and software engineering have an urgent need for trustworthy model behavior. For instance, empirical comparisons between LLMs and classical predictive models in medical diagnosis contexts reveal where LLMs succeed, where they fail, and how evaluation methodology shapes those perceptions (Caruccio et al., 2024). Similarly, studies of the open-source project landscape in major technology firms emphasize dataset heterogeneity and the importance of domain-specific considerations for training and evaluation (Han et al., 2021).

Concurrently, retrieval-augmented models (RAG) and related retrieval-and-use paradigms (RAU) promise to mitigate some limitations of purely parametric LLM knowledge by furnishing models with external, up-to-date, and verifiable evidence at inference time (Hu & Lu, 2024). Yet, integrating retrieval with careful validation of retrieved sources and quantifying the resulting effect on model reliability requires disciplined frameworks that cross methodological boundaries.

This article presents a comprehensive hybrid framework that explicitly connects the chain from dataset curation and augmentation to validation, retrieval integration, and evaluation. The framework emphasizes provenance-aware data pipelines, AI-augmented rule-based validation, and multi-faceted evaluation strategies suited to both general LLMs and domain-specialized models

such as code synthesizers and medical diagnostic bots. Our goals are both theoretical—refining understanding of scale versus curation—and practical—providing operationalizable recommendations for research teams and organizations.

The literature reviewed below sets the theoretical foundation and evidentiary base for this synthesis, highlighting gaps and converging themes that inform the framework’s construction.

METHODOLOGY

The methodology we develop is descriptive, normative, and prescriptive: descriptive in mapping existing empirical findings about data augmentation and scale (Ding, 2019; El-Nouby et al., 2021), normative in articulating best-practice validation procedures (Herzog et al., 2007; Malviya & Parate, 2025), and prescriptive by proposing concrete technical components that can be integrated into model development life-cycles (Hu & Lu, 2024; Gao et al., 2024). The methodology is organized into five interlocking layers: (1) source-level curation and provenance; (2) augmentation strategy; (3) hybrid validation agents; (4) retrieval-augmented design; and (5) multi-dimensional evaluation.

Source-level curation and provenance tracing. The foundation of reliable model training is the provenance and quality of source data. Building on classical record linkage and data quality techniques (Herzog et al., 2007), the framework recommends instrumenting ingestion pipelines with cryptographic and metadata-based provenance tokens that capture origin, timestamp, licensing, and transformation steps. At ingestion, automated deduplication heuristics should be applied, but with sensitivity to near-duplicates that are semantically distinct (e.g., slightly paraphrased documentation, alternative code snippets, or updated medical case studies). Statistical tests for contingency and independence—rooted in long-established methods—can be applied when assessing categorical metadata relationships (Fisher, 1922). Practically, the ingestion layer should maintain immutable logs that enable backward tracing and auditing of any training instance.

Augmentation strategy. Augmentation must be a deliberate mixture of synthetic transformations, human-authored augmentation, and selective sampling from large corpora (Ding, 2019). Rather than relying solely on scale, the augmentation strategy emphasizes diversity—linguistic, stylistic, domain-specific—and representativeness of decision-critical categories. For code models, augmentation includes transforming API call sequences, variable renaming with type consistency, and introducing controlled edge cases that reveal model brittleness (Fried et al., 2022). For general LLMs, augmentation processes should include paraphrase generation, controlled hallucination insertion for robustness testing, and adversarially crafted counterfactuals. Augmentation pipelines must be versioned and tagged by method, magnitude, and purpose

to enable ablation studies and reproducible comparisons. Hybrid validation agents. Purely rule-based validation cannot scale to the subtleties of modern training corpora, while relying solely on LLMs to validate data risks circularity and overfitting. Thus, we propose hybrid validation agents that combine classical rule checks (format, schema, label consistency) with LLM-assisted anomaly detection and human-in-the-loop adjudication (Malviya & Parate, 2025). The rule engine enforces invariants and flags violations deterministically, while an LLM-based validator examines semantic coherence, annotation plausibility, and latent inconsistencies. An adjudication interface allows expert reviewers to accept, correct, or discard flagged instances, and to provide feedback that improves both rule heuristics and LLM validation prompts. This triad—rules, LLMs, humans—yields a pragmatic balance between scalability and judgment.

Retrieval-augmented design and pipeline integration. Retrieval-augmented language models (RAG) decouple parametric storage from evidence provision, enabling up-to-date and verifiable responses (Hu & Lu, 2024). Our framework codifies best practices for retrieval integration: (a) indexing with provenance and freshness metadata; (b) retrieval scoring that penalizes low-quality or ambiguous sources; (c) contextual fusion mechanisms that present retrieved evidence in model prompts without overwhelming the model; and (d) confidence calibration that reflects both parametric certainty and retrieval signal strength. Retrieval systems must also be subject to validation—ensuring that documents are authentic, not duplicated from training data in problematic ways, and appropriately licensed.

Multi-dimensional evaluation protocol. Evaluation must move beyond single-number metrics to a battery approach that measures accuracy, robustness, calibration, factuality, and human-centered criteria such as explainability and trust. We recommend three complementary evaluator classes: (1) statistical evaluators that compute performance across held-out test sets and measure distributional shifts; (2) adversarial and robustness evaluators that probe model behavior under perturbations and out-of-distribution inputs; and (3) meta-evaluators that use LLMs and human raters to assess aspects like fluency, factual grounding, and safety (Gao et al., 2024; Wang et al., 2023). For code models, task-specific evaluations—compilation success, functional correctness via unit tests, and human review of architectural choices—are crucial (Fried et al., 2022). When comparing LLM outputs to supervised models in domains like medical diagnosis, evaluation designs must ensure comparable evidence sets and consistent decision thresholds (Caruccio et al., 2024).

Operational protocols and governance. Complementing the technical stack, governance protocols mandate data retention policies, privacy-preserving transformations (e.g., redaction, differential privacy where appropriate), and auditing mechanisms. Governance also prescribes when human oversight is required, thresholds for model

redeployment, and practices for transparent reporting of data sources and validation outcomes.

RESULTS

This section presents a descriptive synthesis of how applying the above methodology yields measurable improvements in model reliability and interpretability. Because the present work is a theoretical and integrative synthesis rather than a report of a single empirical experiment, the “results” are framed as expected outcomes, theoretical trade-offs, and operational metrics that organizations can observe when implementing the framework.

Trade-offs between dataset scale and curated diversity. One primary result is the formal recognition of a trade-off surface: naive scale increases token counts but can introduce duplication and low-information content that marginally improves or even degrades downstream generalization. Empirical case studies on augmentation show that curated diversity—introducing examples that fill gaps in the decision boundary—yields greater marginal returns than indiscriminate scale increases in many settings (Ding, 2019; El-Nouby et al., 2021). Where datasets contain repeated exposures of identical or near-identical instances, scaling laws indicate diminishing returns and challenges for interpretability (Hernandez et al., 2022). Implementing provenance tracing and targeted augmentation reduces duplication and creates higher-value datasets, as measured by downstream validation metrics and sample efficiency.

Effectiveness of hybrid validation agents. Hybrid agents reduce the incidence of egregious label noise and provenance drift. Rule engines catch deterministic violations at scale, while LLM assistants identify semantically coherent anomalies that rules miss (Malviya & Parate, 2025). Human adjudication confers the final quality check for ambiguous cases. Practically, teams report reductions in error propagation when hybrid validation acts early in the pipeline: fewer mislabeled instances reach model training and fewer problematic examples are amplified by augmentation processes (Herzog et al., 2007). The hybrid approach also provides structured feedback loops that incrementally improve rule sets and LLM verification prompts.

Retrieval integration improves factual grounding and update capability. RAG pipelines, when combined with reliable retrieval indices and provenance-aware ranking, significantly improve model factuality and the ability to incorporate recent knowledge without full retraining (Hu & Lu, 2024). However, this improvement depends on rigorous validation of the retrieval corpus—indexing unvetted web content can reintroduce noise and bias. Robust retrieval scoring that incorporates source credibility, recency, and alignment with provenance metadata leads to better calibration of model outputs, as models can rely on explicit evidence rather than memorized parametric artifacts.

Evaluation battery yields richer diagnostic insights. A multi-dimensional evaluation toolchain exposes

weaknesses that single metrics mask. For example, code models may attain high token-level likelihoods while failing functional tests; conversely, models that are robust under adversarial perturbations provide more reliable user-facing behavior (Fried et al., 2022; Gao et al., 2024). Using LLMs as meta-evaluators can be informative, yet LLM-based judgments must themselves be calibrated and validated to avoid circularity with the models being assessed (Gao et al., 2024; Wang et al., 2023).

Statistical and interpretive outcomes. By applying Fisher's principles of statistical interpretation and contingency analysis in metadata auditing, practitioners can detect nonrandom associations that signal provenance mixing or labeling bias (Fisher, 1922). For instance, an over-representation of certain authors, domains, or time periods in training data can be quantified and corrected through targeted sampling and augmentation. Robust record-linkage methods allow consolidation of multiple identifiers for the same underlying content, preventing accidental duplication that skews model learning (Herzog et al., 2007).

DISCUSSION

This section offers a deep interpretive analysis of the framework's implications, potential limitations, and future research directions. It contextualizes the theoretical propositions and anticipates operational challenges.

Theoretical implications: scale revisited. Our synthesis refines the position articulated by El-Nouby et al. (2021) and Hernandez et al. (2022): while large-scale datasets have historically yielded performance gains, the marginal usefulness of additional data depends heavily on its information content, uniqueness, and representativeness. From a theoretical standpoint, the framework suggests reframing scaling debates: not as "more versus less," but as "scale plus curation." In other words, the fundamental variable is not raw size but the effective information density—how much unique, decision-relevant variation the dataset offers. This perspective aligns with empirical observations that augmentation and careful curation can substitute for some quantity-driven benefits (Ding, 2019).

Evaluation paradigms and the role of meta-evaluation. Contemporary evaluation practice often treats models as black boxes and relies on single-number metrics that can be gamed or detached from human utility. Our framework promotes layered evaluation: objective statistical tests, functional correctness checks, adversarial stress-tests, and human-centered assessments. LLM-based meta-evaluators are promising but must be designed to minimize circularity. Specifically, meta-evaluators should be trained or prompted with distinct evidence sets and be regularly calibrated against human judgments (Gao et al., 2024; Wang et al., 2023). This approach fosters a culture of continuous evaluation where models are constantly probed for brittle or misaligned behavior.

Hybrid validation: balancing automation and expertise. The hybrid triad of rules, LLMs, and human adjudicators aims to harness the scale of automation while preserving expert discernment where matters are unclear or high-stakes. A core insight is that human reviewers' time is best spent on borderline or high-impact items—cases that neither rules nor LLM validators resolve confidently. Systems that funnel such examples to humans produce better outcomes than those that attempt full automation or indiscriminate manual curation.

Retrieval design: transparency and provenance are non-negotiable. Retrieval-augmented models are powerful precisely because they allow models to anchor outputs to external sources. This anchoring demands honest and machine-actionable provenance metadata. Without it, retrieved evidence can be misleading, out-of-date, or violate licenses. The framework therefore insists on provenance-rich indices and ranking functions that penalize low-quality sources, ensuring that retrieval supports, rather than undermines, model reliability (Hu & Lu, 2024).

Limitations and risks. Several limitations and risks must be acknowledged. First, the framework is resource-intensive: provenance tracing, hybrid validation, and RAG infrastructure require engineering investment and human oversight, which may be prohibitive for smaller teams. Second, LLM-assisted validation introduces the risk of confirmation bias if validators are not carefully decoupled from the models they assess. Third, while augmentation can close gaps in training coverage, poorly designed synthetic data can amplify biases or introduce unrealistic patterns. Finally, the framework's governance recommendations—privacy-preserving methods, audit trails—introduce trade-offs between transparency and privacy, necessitating careful legal and ethical review. Future research directions. Several avenues emerge for empirical validation and theoretical refinement:

1. Quantifying information density. Develop formal measures of effective information content in datasets to predict when curation will outperform simple scale increases. This research could draw on information-theoretic metrics and empirical performance curves.
2. Meta-evaluator benchmarking. Construct standardized benchmarks to evaluate LLM-based meta-evaluators, ensuring they generalize across domains and do not overfit to specific model families.
3. Adjudication workflow optimization. Investigate active learning strategies to select instances for human adjudication that maximize validator effectiveness per unit human effort.
4. Provenance standards for retrieval. Propose interoperable provenance schemas for retrieval indices that balance expressivity and implementability across organizations.
5. Cost-benefit modeling. Develop economic

models comparing the marginal gains from additional data versus investments in curation, validation, and retrieval infrastructure.

CONCLUSION

This article proposes and elaborates a comprehensive hybrid framework connecting data augmentation, provenance-aware curation, hybrid validation, retrieval-augmented design, and multi-dimensional evaluation for large language models and code-generative systems. The synthesis draws on empirical case studies, theoretical debates about scale, and recent developments in retrieval and evaluation methodologies (Ding, 2019; El-Nouby et al., 2021; Hernandez et al., 2022; Hu & Lu, 2024; Gao et al., 2024). Key contributions include the articulation of effective information density as a lens for assessing the value of dataset expansion, the design of hybrid validation agents that combine rule-based, LLM-assisted, and human adjudication techniques, and the operationalization of retrieval integration with provenance-aware indices.

Practically, the framework provides researchers and practitioners with a roadmap for constructing more auditable, reliable, and interpretable LLM pipelines. It suggests that responsible progress in model capabilities requires not only scaling resources but also deliberate investments in data quality, evaluation, and governance. The proposed directions for future research emphasize measurable metrics and interoperable standards that can accelerate adoption and improve comparability across research efforts.

By focusing on the intersection of curation, augmentation, validation, retrieval, and evaluation, the framework offers a balanced pathway for advancing LLM reliability without succumbing to the pitfalls of unchecked data growth. The enduring challenge for the field will be to translate these conceptual commitments into practical tooling and shared standards that make high-quality model development accessible beyond large industrial teams.

REFERENCES

1. Junhua Ding, Xinchuan Li, Xiaojun Kang, and Venkat N. Gudivada. 2019. A case study of the augmentation and evaluation of training data for deep learning. *Journal of Data and Information Quality (JDIQ)* 11, 4 (2019), 1–22.
2. Alaaeldin El-Nouby, Gautier Izacard, Hugo Touvron, Ivan Laptev, Hervé Jegou, and Edouard Grave. 2021. Are large-scale datasets necessary for self-supervised pre-training? *arXiv preprint arXiv:2112.10740* (2021).
3. Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 2023. Large language models for software engineering: Survey and open problems. *arXiv preprint arXiv:2310.03533* (2023).
4. Ronald A. Fisher. 1922. On the interpretation of χ^2 from contingency tables, and the calculation of P. *Journal of the Royal Statistical Society* 85, 1 (1922), 87–94.
5. Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wentau Yih, Luke Zettlemoyer, and Mike Lewis. 2022. Incoder: A generative model for code infilling and synthesis. *arXiv preprint arXiv:2204.05999* (2022).
6. Mingqi Gao, Xinyu Hu, Jie Ruan, Xiao Pu, and Xiaojun Wan. 2024. LLM-based NLG evaluation: Current status and challenges. *arXiv preprint arXiv:2402.01383* (2024).
7. Junxiao Han, Shuiguang Deng, David Lo, Chen Zhi, Jianwei Yin, and Xin Xia. 2021. An empirical study of the landscape of open source projects in Baidu, Alibaba, and Tencent. In *2021 IEEE/ACM 43rd International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. IEEE, 298–307.
8. Malviya, S., & Vrushali Parate. 2025. AI-Augmented Data Quality Validation in P&C Insurance: A Hybrid Framework Using Large Language Models and Rule-Based Agents. *International Journal of Computational and Experimental Science and Engineering*, 11(3). <https://doi.org/10.22399/ijcesen.3613>
9. Danny Hernandez, Tom Brown, Tom Conerly, Nova DasSarma, Dawn Drain, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Tom Henighan, Tristan Hume, et al. 2022. Scaling laws and interpretability of learning from repeated data. *arXiv preprint arXiv:2205.10487* (2022).
10. Thomas N. Herzog, Fritz J. Scheuren, and William E. Winkler. 2007. *Data quality and record linkage techniques*. Vol. 1. Springer.
11. Yucheng Hu and Yuxing Lu. 2024. RAG and RAU: A Survey on Retrieval-Augmented Language Model in Natural Language Processing. *arXiv preprint arXiv:2404.19543* (2024).
12. Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al. 2024. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology* 15 (2024), 1–45.
13. W. Wang, B. Haddow, A. Birch, W. Peng. 2023. Assessing the reliability of large language model knowledge. *arXiv:2310.09820* (2023).

<https://doi.org/10.48550/arXiv.2310.09820>

- 14.** L. Caruccio, et al. 2024. Can ChatGPT provide intelligent diagnoses? A comparative study between predictive models and ChatGPT to define a new medical diagnostic bot. *Expert Systems with Applications* 235 (2024), 121186. <https://doi.org/10.1016/j.eswa.2023.121186>
- 15.** Y. Jin, X. Wang, R. Yang, Y. Sun, W. Wang, H. Liao, X. Xie. 2022. Towards fine-grained reasoning for fake news detection. *Proceedings of the AAAI Conference on Artificial Intelligence* 36 (2022), 5746–5754. <https://doi.org/10.1609/aaai.v36i5.20517>