

MetaGAI: A Large-Scale and High-Quality Benchmark for Generative AI Model and Data Card Generation

Anonymous ACL submission

Abstract

The rapid proliferation of Generative AI necessitates rigorous documentation standards for transparency and governance. However, manual creation of Model and Data Cards is not scalable, while automated approaches lack large-scale, high-fidelity benchmarks for systematic evaluation. We introduce MetaGAI, a comprehensive benchmark comprising 2,541 verified document triplets constructed through semantic triangulation of academic papers, GitHub repositories, and Hugging Face artifacts. Unlike prior single-source datasets, MetaGAI employs a multi-agent framework with specialized Retriever, Generator, and Editor agents, validated through four-dimensional human-in-the-loop assessment. We establish a robust evaluation protocol combining automated metrics with validated LLM-as-a-Judge frameworks. Extensive analysis reveals that sparse Mixture-of-Experts architectures achieve superior cost-quality efficiency, while a fundamental trade-off exists between faithfulness and completeness. MetaGAI provides a foundational testbed for benchmarking, training, and analyzing automated Model and Data Card generation methods at scale. Our data and code are available at <https://anonymous.4open.science/r/MetaGAI-DBB4>.

1 Introduction

The rapid proliferation of Generative AI (GenAI) has fundamentally transformed machine learning deployment. As these systems transition from research artifacts to critical infrastructure, the demand for transparency and accountability has intensified. This has driven the evolution of documentation standards from foundational Model Cards (Mitchell et al., 2019) and Data Cards (Pushkarna et al., 2022) to comprehensive System Cards adopted by industry leaders (OpenAI et al., 2024; Comanici et al., 2025; OpenAI et al., 2025). These modern frameworks extend

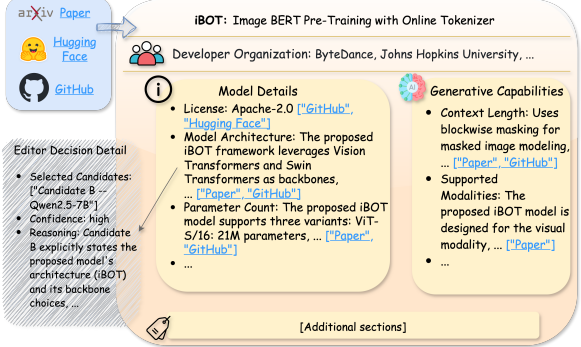


Figure 1: MetaGAI Benchmark Construction Example. Automated GenAI card generation for the iBOT model (Zhou et al., 2022) demonstrating Multi-Source Triangulation combining architectural concepts from Papers, hyperparameters from GitHub, and licensing data from Hugging Face, with Editor-Based Synthesis to produce high-fidelity ground truth.

beyond static performance metrics to document safety alignment procedures, red-teaming results, and societal impacts.

Standardized documentation artifacts serve as essential infrastructure for AI governance and reproducibility. They enable longitudinal tracking of model evolution (Castaño et al., 2024), automated compliance auditing, and systematic risk assessment in high-stakes domains (Longpre et al., 2024). Without high-quality documentation, the AI ecosystem lacks the interoperability required to benchmark capabilities or trace data provenance across complex supply chains (Rahman et al., 2025).

However, a significant bottleneck impedes widespread adoption. The GenAI ecosystem is increasingly driven by the “long tail” of open-source contributions, comprising thousands of models on community platforms (Wolf et al., 2020; Horwitz et al., 2025). Unlike well-resourced industry laboratories, developers in this ecosystem often lack the capacity to maintain rigorous documentation. Manual card creation suffers from severe scalabil-

ity constraints, characterized by pervasive incompleteness, inconsistency, and subjectivity (Yang et al., 2024; Liang et al., 2024), leading to a transparency crisis where many research papers and repositories lack structured metadata necessary for reproducibility (Olmo et al., 2025). Automated documentation generation has emerged as a critical necessity (Liu et al., 2024a). To address this challenge, current approaches face significant hurdles: zero-shot methods often hallucinate details when summarizing lengthy documents, while retrieval-augmented strategies struggle to align diverse paper structures with rigid schemas. Progress is stalled by the lack of large-scale, high-quality benchmarks that can objectively measure automated generation accuracy against verified ground truth.

To fill this gap, we introduce **MetaGAI**, a large-scale benchmark designed to systematically evaluate automated Model and Data Card generation. Unlike prior datasets treating documentation as simple summarization, MetaGAI formulates it as a complex multi-source information generation task, mirroring real-world requirements for verifying scientific claims against implementation details.

As illustrated in Figure 1, creating complete documentation cards requires triangulating evidence from heterogeneous sources. Taking the iBOT model (Zhou et al., 2022) as an example, architectural concepts are derived from the academic Paper, implementation details are extracted from the GitHub repository, and deployment constraints are verified via Hugging Face. Our construction pipeline utilizes a multi-agent framework comprising specialized Retriever, Generator, and Editor agents to synthesize these signals into verified ground truth. We implement rigorous human-in-the-loop validation across four dimensions: (1) retrieval strategy validation through domain expert annotation, (2) generator divergence analysis validating ensemble diversity, (3) editor efficacy assessment through hybrid human-LLM evaluation panels, and (4) editor architecture selection through pairwise comparisons, establishing a rigorous standard for evaluating automated systems.

We employ MetaGAI to evaluate cost-effective LLMs suitable for large-scale deployment, identifying architectures that deliver high-fidelity card generation while maintaining practical cost-efficiency for processing scientific literature at scale.

In a nutshell, the key contributions of this study are summarized as follows:

- We construct MetaGAI, the largest high-quality benchmark for GenAI documentation, comprising 2,541 verified triplets with rigorous human-in-the-loop validation across four dimensions.
- We propose a robust evaluation framework combining granular automated metrics with a validated LLM-as-a-Judge protocol.
- We provide extensive empirical analysis revealing that sparse Mixture-of-Experts (MoE) architectures offer superior cost-quality efficiency, though a systematic trade-off persists between faithfulness and completeness in generating abstract metadata.

2 Related Work

Model Documentation. The lack of standardized documentation for trained machine learning models limits transparency, reproducibility, and responsible use in NLP. Model Cards were introduced to provide structured summaries of trained models, including intended use, evaluation settings, performance, and limitations (Mitchell et al., 2019). Dataset-level documentation was developed in parallel. Data Cards focus on data provenance, collection processes, representativeness, and ethical considerations (Pushkarna et al., 2022). Together, Model Cards and Data Cards form a documentation framework that addresses both model-level and data-level sources of uncertainty. Subsequent work extended model documentation to interactive and machine-readable formats. Interactive Model Cards were shown to better support user understanding of model behavior (Crisan et al., 2022), while linked and semantic representations of documentation were proposed to improve traceability and reuse (Donald et al., 2023). Empirical analysis of a large number of Model Cards indicates substantial variation in documentation quality, with evaluation and limitation sections often missing or incomplete (Liang et al., 2024).

Automatic Model and Data Card Generation. To reduce the manual effort required for documentation, recent work has explored automatic generation of Model and Data Cards. Liu et al. (2024a) introduced CARDBENCH, a benchmark of human-written cards, together with CARDGEN, an LLM-based system that generates structured documentation by retrieving information from model repositories and associated papers. Their results show

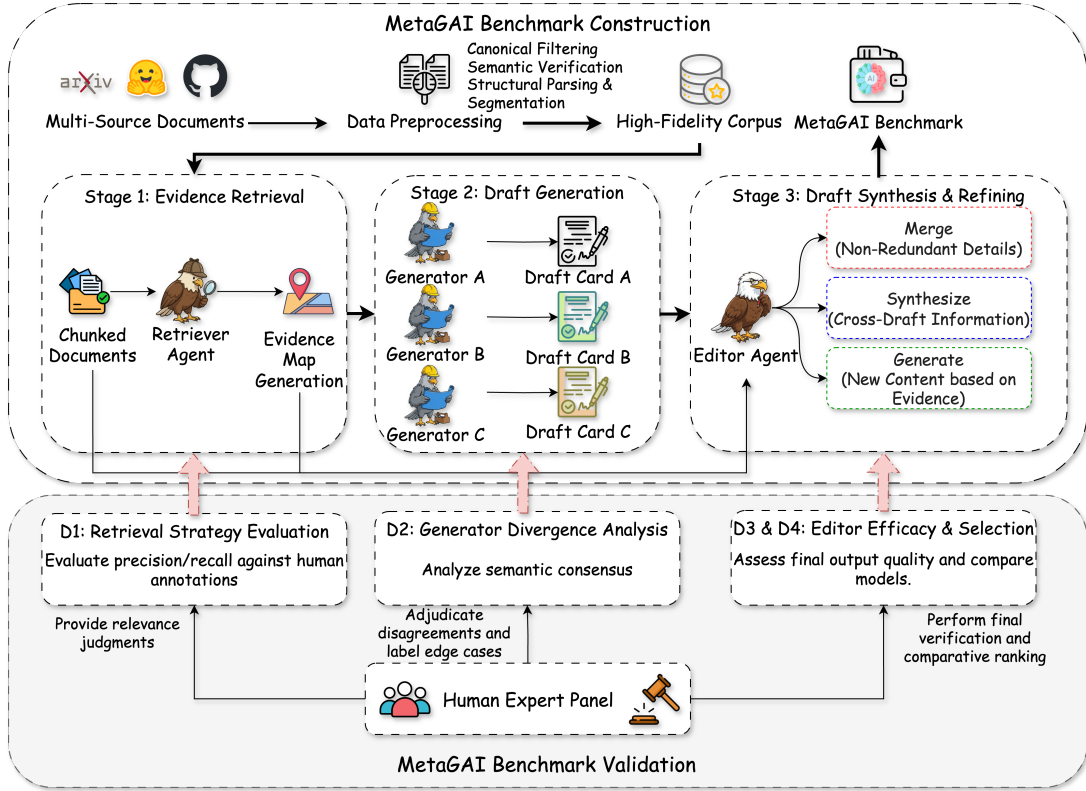


Figure 2: **MetaGAI Benchmark Construction and Validation Framework.** The pipeline integrates multi-source document preprocessing, a multi-agent generation framework (Evidence Retrieval, Draft Generation, Draft Synthesis and Refining), and a four-dimensional validation protocol (D1-D4) incorporating human expert adjudication.

165 that automatic generation can improve documenta- 188
 166 tion completeness. However, existing benchmarks 189
 167 remain limited in scale and coverage, which con- 190
 168 strains systematic evaluation of generation meth- 191
 169 ods. To address this limitation, we introduce Meta- 192
 170 GAI, a large-scale benchmark for GenAI Model 193
 171 and Data Card generation. 194

172 3 MetaGAI Benchmark Construction 195

173 Figure 2 illustrates the MetaGAI benchmark con- 196
 174 struction and validation pipeline, comprising three 197
 175 sequential stages: (1) evidence retrieval from multi- 198
 176 source documents, (2) multi-generator draft syn- 199
 177 thesis, and (3) editor-based consolidation and re- 200
 178 finement. This architecture extends beyond tradi- 201
 179 tional single-model approaches (Liu et al., 2024a) 202
 180 by incorporating ensemble generation and rig- 203
 181 orous human-in-the-loop validation to ensure high- 204
 182 fidelity card generation for GenAI ecosystems.

183 3.1 Task Definition 205

184 We formulate **MetaGAI** as a structured genera- 206
 185 tion task mapping unstructured scientific papers 207
 186 \mathcal{P} to standardized cards \mathcal{C} . Ground truth \mathcal{C}_{GT} is 208
 187 derived from a document triplet $(\mathcal{P}, \mathcal{G}, \mathcal{H})$, where 209

188 \mathcal{G} denotes the GitHub repository and \mathcal{H} the Hug- 189
 190 ging Face artifact. This multi-source triangulation 191
 192 addresses implementation details frequently omit- 193
 194 ted in manuscripts. We adopt taxonomies from 195
 196 Mitchell et al. (2019) and Pushkarna et al. (2022), 197
 198 extending them with GenAI-specific fields (Ap- 199
 200 pendix Table 3) to capture domain-specific char- 201
 202 acteristics such as prompt engineering configura- 203
 204 tions and generative model architectures.

Each card \mathcal{C} comprises M field-value pairs 205
 $\{(k_i, v_i)\}_{i=1}^M$, where k_i represents a schema at- 206
 207 tribute and v_i represents its corresponding content. 208
 209 The objective is learning a mapping f_θ that predicts 210
 $\hat{\mathcal{C}}$ from paper input alone: 211

$$212 \hat{\mathcal{C}} = f_\theta(\mathcal{P}) \quad (1) \quad 213$$

214 minimizing divergence from the high-fidelity refer- 215
 216 ence $\mathcal{C}_{GT}(\mathcal{P}, \mathcal{G}, \mathcal{H})$. 217

218 3.2 Data Acquisition and Filtering 220

219 **Corpus Construction.** We construct the bench- 221
 222 mark through systematic triangulation of arXiv, 223
 224 GitHub, and Hugging Face metadata to capture im- 225
 226 plementation details absent from published papers. 227

Unlike prior work that aggregates existing documentation (Liu et al., 2024a), our work constructs a new corpus from scratch, with ground truth verified via cross-source semantic consistency. An initial corpus of 15,727 candidates is canonically filtered to identify unique paper–model–dataset linkages, resulting in 4,068 entries.

Semantic Verification. To eliminate spurious citations and ensure cross-source alignment, we employed Qwen3-30B-A3B-Instruct (Yang et al., 2025) for automated semantic consistency verification across the three sources (Prompt F.1). A pilot study with two domain experts achieved 100% inter-annotator agreement, validating the automated approach. This rigorous filtering process produced a final high-fidelity corpus of 2,541 verified triplets (Appendix B), addressing the incompleteness and inconsistency challenges inherent in human-authored documentation. Detailed provenance analysis appears in Appendix G.

3.3 Benchmark Generation Pipeline

We developed an automated pipeline synthesizing high-fidelity cards through document preprocessing and a multi-agent framework.

3.3.1 Pre-processing

Raw PDF documents were converted to structured Markdown using OLMoCR-2, a state-of-the-art OCR model optimized for academic layouts (Poznanski et al., 2025). To accommodate context window constraints, converted papers and README files were segmented into 1024-token chunks, enabling precise retrieval of implementation details dispersed throughout documents.

3.3.2 Multi-Agent Framework

We employ a multi-agent architecture (Algorithm 1) to maximize factual grounding and minimize hallucinations (Du et al., 2024). Where traditional approaches rely on single-model generation following retrieval (Liu et al., 2024a), our framework introduces ensemble diversity and editor-based cross-validation. The framework processes full context $\mathcal{P} \cup \mathcal{S}$, where $\mathcal{S} = \{\mathcal{G}, \mathcal{H}\}$ denotes supplementary documentation, through three specialized agents:

$$v_i = \text{Retriever}(\mathcal{P} \cup \mathcal{S}, k_i) \quad (2)$$

$$\tilde{\mathcal{C}} = \text{Generator}(\{(k_i, v_i)\}_{i=1}^M) \quad (3)$$

$$\hat{\mathcal{C}} = \text{Editor}(\mathcal{P} \cup \mathcal{S}, \tilde{\mathcal{C}}) \quad (4)$$

Retriever Agent. This agent aligns unstructured text with schema fields through exhaustive chunk-level classification. Each document segment receives a relevance score (0–4) with supporting keywords (Prompt F.2). Comparative analysis (Section 4.1) established generative models’ superiority over discriminative rerankers for complex schema alignment, informing our selection of Qwen3-30B-A3B-Instruct as the retrieval backbone. This generative reasoning approach contrasts with embedding-based similarity matching, enabling more nuanced interpretation of implicit metadata requirements.

Generator Agent. To mitigate architectural bias and capture diverse interpretations, we employ ensemble generation (Wang et al., 2023) with three architecturally distinct LLMs (OLMo-3-7B, Llama-3.1-8B, Qwen2.5-7B). Each model independently synthesizes draft content, evidence quotations, and confidence scores. A strict “Subject Focus” constraint (Prompt F.3) ensures agents synthesize information about the proposed system only, excluding baselines and prior work. This multi-model design addresses single-architecture limitations in capturing the full evidence spectrum across complex GenAI systems.

Editor Agent. A “Chief Editor” consolidates candidate drafts using a larger, cross-family model to avoid self-enhancement bias (Zheng et al., 2023). The editor validates semantic alignment between drafts and raw evidence, filters attribution errors, and merges non-redundant details into concise final entries (Prompt F.4). This explicit cross-validation stage against original sources provides stronger hallucination mitigation than single-pass generation, which is particularly critical for technical specifications where factual precision is essential. Concrete examples appear in Appendix H.

4 MetaGAI Validation and Analysis

To rigorously validate the benchmark construction pipeline, we conduct human-in-the-loop experimentation across four critical dimensions, substantially extending beyond prior validation efforts (Liu et al., 2024a) through systematic component-wise evaluation.

Our validation framework examines retrieval strategies (D1) comparing generative reasoning versus discriminative reranking, quantifies semantic variation across ensemble generators (D2), measures quality improvements from editor-based con-

Model Architecture	P@1	R@5	F1@5
<i>Discriminative Rerankers</i>			
BGE-Reranker-v2-m3	0.135	0.114	0.101
Qwen3-Reranker-8B	0.115	0.137	0.091
<i>Generative LLMs</i>			
Llama-3.1-8B-Instruct	0.615	0.406	0.326
Qwen2.5-7B-Instruct	0.538	0.282	0.232
Qwen3-30B-A3B-Instruct	0.635	0.469	0.342

Table 1: **Comparative Validation of Retrieval Models (D1)**. Performance metrics evaluated on the curated validation set. We report Precision at rank 1 (P@1) to measure top-result accuracy, alongside Recall and F1 at rank 5 (R@5, F1@5) to assess broader retrieval coverage. The best performance in each category is highlighted in bold.

solidation (D3), and compares performance across editor architectures (D4). This enables isolating each pipeline component’s contribution and optimizing overall system configuration.

4.1 Retrieval Strategy Validation (D1)

We conducted comparative validation to identify the optimal retrieval backbone using human-verified gold standards. Ten randomly sampled entries were segmented and evaluated by two domain experts. From 6,116 chunk-field decisions, we retained 276 chunks achieving unanimous consensus as ground truth, prioritizing precision over coverage. We compared five models across two paradigms: discriminative rerankers (BGE-Reranker-v2-m3 (Chen et al., 2025), Qwen3-Reranker-8B (Zhang et al., 2025)) and generative LLMs (Llama-3.1-8B-Instruct (Grattafiori et al., 2024), Qwen2.5-7B-Instruct, Qwen3-30B-A3B-Instruct (Yang et al., 2025)).

Table 1 demonstrates generative models’ decisive advantage. Discriminative rerankers achieved $P@1 < 0.14$, indicating that semantic similarity alone is insufficient for schema mapping. Qwen3-30B-A3B-Instruct achieved superior performance ($P@1$: 0.635, $F1@5$: 0.342), establishing it as the retrieval backbone.

4.2 Generator Divergence Analysis (D2)

We analyzed semantic consensus within the generator ensemble (OLMo-3-7B, Llama-3.1-8B, Qwen2.5-7B) by computing pairwise BERTScore (Zhang et al., 2020) similarity across all 2,541 samples. Figure 3 reveals distinct behavioral patterns. Fields clustering

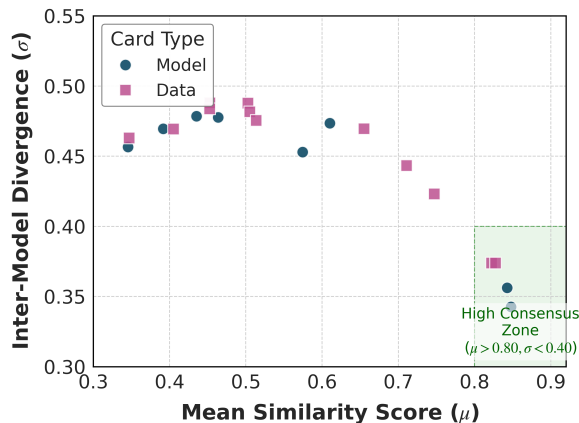


Figure 3: **Mean-Variance analysis of semantic alignment**. Mean BERTScore similarity (μ , X-axis) versus variance (standard deviation σ , Y-axis) across card fields. High-consensus zones reflect evidence scarcity with converged placeholder responses.

in the high-consensus zone ($\mu > 0.8, \sigma < 0.4$), such as *Ethical Considerations* ($\mu = 0.848$), predominantly reflect evidence scarcity in source documents (Appendix G.2): when documentation is absent, generators converge on similar placeholder responses. Conversely, substantial portions of both Model and Data Card fields exhibit pronounced divergence ($\sigma > 0.4$), confirming that architecturally distinct models generate complementary details from information-rich contexts. This variation validates our ensemble design, enabling the Editor Agent to synthesize broader evidence coverage and mitigate individual model biases.

4.3 Editor Efficacy and Selection (D3 & D4)

We conducted controlled experiments to evaluate editor impact and architecture selection using a stratified sample of 10 entries (5 Model/5 Data Cards). For each entry, we generated three blinded outputs: (1) Raw Baseline (highest-token-count generator draft), (2) Editor A (GPT-OSS-20B (OpenAI et al., 2025)), and (3) Editor B (Mistral-3-14B-Instruct¹).

A hybrid panel comprising three Ph.D. students and three LLMs (Claude 4.5 Sonnet, GPT 5.2 Thinking, Gemini 3 Pro) assessed candidates on 1–5 Likert scales across five dimensions (Appendix D). LLM evaluations were triplicated with randomized orders to mitigate position bias. We computed Mean Likert Scores and Pairwise Win Rates, validating LLM-human alignment via Spear-

¹<https://mistral.ai/news/mistral-3>

man correlation. **Editor Efficacy (D3):** Both editors substantially outperformed Raw Baseline (win rates: 70–80%), confirming that consolidation effectively filters hallucinations. Absolute scores: Raw Baseline (3.44 human/3.48 LLM), Mistral (3.75/3.72), GPT (3.80/3.84). **Editor Selection (D4):** GPT marginally outperformed Mistral (60% win rate), though quality differences remain minimal (3.80 vs. 3.75). We therefore incorporate both architectures, averaging outputs to ensure architectural neutrality. All LLM judges achieved strong human alignment (Spearman $\rho > 0.67$, $p < 0.001$; GPT-Judge: $\rho = 0.752$), validating automated evaluation frameworks.

5 Experiments

5.1 Experimental Setup

We evaluate models in a zero-shot setting where cards $\hat{\mathcal{C}}$ are generated exclusively from paper text \mathcal{P} , excluding supplementary sources \mathcal{S} to simulate realistic deployment constraints. In practice, a substantial proportion of models and datasets lack documentation cards or suffer from severe incompleteness (Liang et al., 2024; Yang et al., 2024), making paper-only generation the predominant scenario for large-scale deployment. To balance statistical rigor with computational feasibility, we employ a hybrid protocol: automated metrics are computed on the **full benchmark** ($N = 2, 541$), while LLM-as-a-Judge evaluations use a **stratified random sample** of 500 entries (250 each).

5.2 Baselines

We examine models categorized by access modality and architectural paradigm. Preliminary experiments with models below 14B parameters show JSON formatting failures; therefore, we restrict evaluation to architectures with 20B+ parameters.

Open-Weight Models. We evaluate two architectural classes. **Dense Models** employ standard transformer architectures (20B–32B parameters), including Mistral-Small-3.2-24B-Instruct, Gemma-3-27B-IT (Team et al., 2025), and Qwen3-32B (Yang et al., 2025). **MoE Models** utilize Mixture-of-Experts architectures providing high parameter counts with efficient active parameter usage, comprising GPT-OSS-20B (OpenAI et al., 2025), NVIDIA-Nemotron-3-Nano-30B-A3B (NVIDIA et al., 2025), and Qwen3-30B-A3B-Instruct (Yang et al., 2025). **Closed-Source Models.** We include proprietary models accessed via

API: GPT-5.1 series (Mini/Nano)² and Gemini-2.5 series (Flash/Flash-Lite) (Comanici et al., 2025).

5.3 Evaluation Metrics

We employ a dual-layered evaluation combining quantitative structural metrics with qualitative expert judgment.

Automated Metrics. We measure recall and semantic alignment through three metrics:

- **Completeness:** Quantifies field-level recall. Let $\mathcal{K}(\mathcal{C})$ denote the set of populated keys in card \mathcal{C} :

$$\text{Completeness} = \frac{|\mathcal{K}(\hat{\mathcal{C}}) \cap \mathcal{K}(\mathcal{C}_{GT})|}{|\mathcal{K}(\mathcal{C}_{GT})|} \quad (5)$$

- **Semantic Similarity:** We report **ROUGE-L** (Lin, 2004) for lexical overlap and **BERTScore (F1)** (Zhang et al., 2020) for semantic alignment.

LLM-as-a-Judge Evaluation. To capture nuances beyond n-gram matching, we implement an ensemble framework using GPT-OSS-120B (OpenAI et al., 2025), Llama-3.3-70B-Instruct (Grattafiori et al., 2024), and Qwen3-235B-A22B-2507 (Yang et al., 2025). Following the protocol in Appendix D, judges evaluate generated fields on a 1–5 Likert scale. Scores are averaged across judges to mitigate single-model bias (agreement analysis in Appendix E).

Cost Efficiency. To assess economic feasibility at scale, we introduce a Cost Index that estimates the inference cost per card generation task. The index normalizes costs across models by standardizing token consumption to 1M input tokens and 0.2M output tokens, enabling direct price comparison independent of model-specific tokenization efficiency. Pricing is derived from standardized rates on OpenRouter³, facilitating economic comparison across proprietary and open-weight models.

5.4 Experimental Results

We establish a comprehensive evaluation protocol combining automated metrics on 2,541 entries with LLM-as-a-Judge assessment on 500 samples (vs. CardBench’s 350 (Liu et al., 2024a)). Beyond quality assessment, we introduce cost-efficiency analysis, revealing that sparse MoE models achieve optimal cost-quality trade-offs while traditional lexical

²<https://openai.com/index/gpt-5-1/>

³<https://openrouter.ai/>

Model	Faithfulness		Relevance		Accuracy		Consistency		Usefulness		Qual. Avg		Comp.		BScore		RL		Cost
	D	M	D	M	D	M	D	M	D	M	D	M	D	M	D	M	D	M	
Dense Models																			
Mistral-Small-24B	3.39	3.75	3.28	3.59	3.33	3.70	3.64	3.97	2.67	3.14	3.26	3.63	.280	.327	.193	<u>.194</u>	.270	<u>.255</u>	0.10
Gemma-3-27B	3.81	3.97	3.76	3.94	3.76	3.96	4.03	4.23	3.26	3.58	3.72	3.94	.464	<u>.734</u>	.151	.151	.214	<u>.203</u>	0.07
Qwen3-32B	3.40	3.72	3.29	3.58	3.32	3.69	3.66	3.94	2.65	3.10	3.26	3.60	.386	.513	<u>.184</u>	.198	<u>.267</u>	.265	0.13
MoE Models																			
GPT-OSS-20B	3.18	3.49	3.09	3.36	3.11	3.45	3.45	3.74	2.49	2.87	3.06	3.38	.394	.492	.146	.133	.235	.220	0.06
Nemotron-Nano-30B-A3B	3.45	3.80	3.41	3.76	3.39	3.77	3.71	4.03	2.88	3.35	3.37	3.74	<u>.557</u>	.644	.124	.133	.199	.206	0.11
Qwen3-30B-A3B-Instruct	4.33	4.50	4.36	4.57	4.31	4.51	4.46	4.67	4.06	4.49	4.30	4.55	.702	.786	.169	.174	.246	.243	0.15
Closed-Source Models																			
GPT-5-Mini	<u>4.14</u>	<u>4.32</u>	<u>4.09</u>	<u>4.26</u>	<u>4.10</u>	<u>4.32</u>	<u>4.25</u>	<u>4.43</u>	<u>3.74</u>	<u>4.07</u>	<u>4.06</u>	<u>4.28</u>	.556	.216	.102	.117	.185	.207	0.65
GPT-5-Nano	3.18	3.38	3.06	3.18	3.10	3.32	3.46	3.57	2.53	2.73	3.06	3.23	.383	.127	.113	.149	.188	.228	0.13
Gemini-2.5-Flash	3.88	4.08	3.80	3.97	3.82	4.06	4.04	4.19	3.27	3.69	3.76	4.00	.443	.181	.160	.170	.241	.246	0.80
Gemini-2.5-Flash-Lite	3.93	4.10	3.88	4.03	3.87	4.08	4.12	4.25	3.39	3.69	3.83	4.03	.511	.194	.140	.144	.207	.208	0.18

Table 2: **MetaGAI Benchmark Comprehensive Results.** Performance across Data Card (D) and Model Card (M) generation. Completeness, BERTScore, and ROUGE-L are reported on the full test set (N=2,541). Qual. Avg (Qualitative Average) is evaluated on a sample of 500 entries. **Bold/Underline** indicate best/second-best performance.

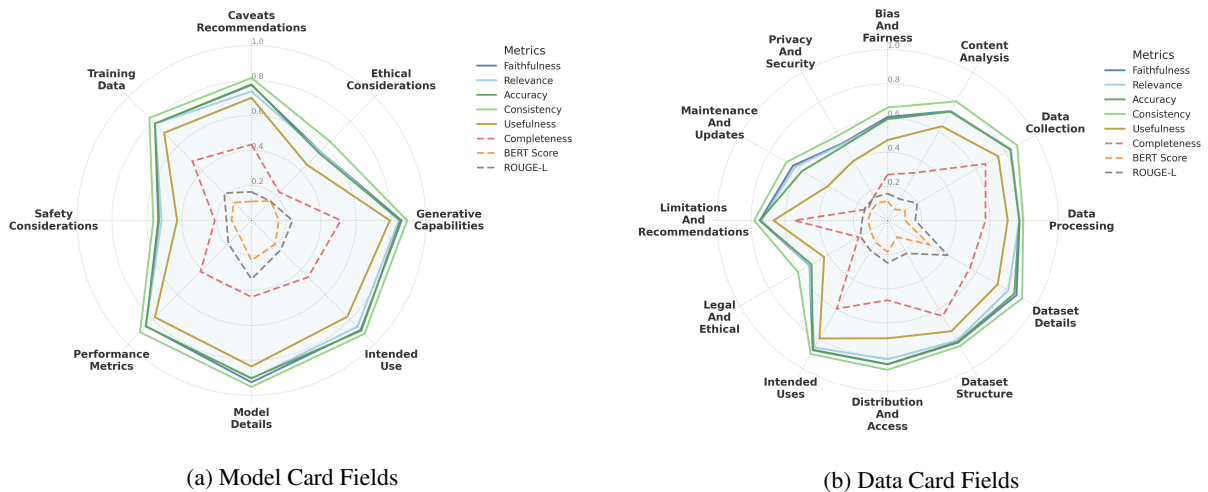


Figure 4: **Field-Level Performance Patterns Averaged Across All Baselines.** Evaluation metrics (colored lines) across card fields (axes). Performance is strong on signal-rich fields (*Model Details*) but degrades on abstract categories (*Ethical Considerations*), revealing systematic generation difficulty when documentation is sparse.

metrics inversely correlate with semantic quality. These findings provide practical deployment guidelines absent in prior benchmarks. Table 2 presents comprehensive evaluation results across all baseline models. We analyze performance through architectural efficiency, evaluation metric validity, and cognitive limitations.

5.4.1 Architectural Efficiency and Economic Viability

MoE Superiority with Scale Threshold. Sparse MoE architectures demonstrate superior parameter efficiency, but only at sufficient scale. Within the Qwen family, MoE-based Qwen3-30B-A3B-Instruct (4.55 Model Card quality) outperforms

dense Qwen3-32B (3.60) by 0.95 points despite comparable parameter counts. Given shared pre-training lineages, this gap validates that sparse activation enables more effective information synthesis. However, this advantage requires approximately 30B parameters: the smallest MoE (GPT-OSS-20B) achieves the lowest performance (3.06/3.38), underperforming comparable dense models. This establishes a minimum viable scale threshold for MoE benefits in complex generation tasks.

Open-Weight Cost Dominance. Qwen3-30B-A3B-Instruct occupies the Pareto-optimal position on the cost-quality frontier. It delivers state-of-the-art performance (4.55 quality score) at the lowest

normalized cost (Cost Index: 0.15), outperforming GPT-5-Mini (4.28 quality, Cost Index: 0.65) by 4.3× in cost efficiency. Gemini-2.5-Flash presents the least favorable value proposition: inferior quality (4.00) at the highest normalized cost (Cost Index: 0.80). Detailed analysis in Appendix G.4 confirms that optimized open-weight MoE architectures provide economically superior solutions for production-scale card generation.

5.4.2 Evaluation Metric Validity

Lexical Similarity Metrics Invalidation. Traditional content matching metrics exhibit inverse correlation with semantic quality. ROUGE-L demonstrates this paradox: Mistral-Small-24B achieves the highest ROUGE-L (0.270) through naive verbatim copying yet produces one of the lowest quality (3.26), while Qwen3-30B-A3B-Instruct yields lower ROUGE-L (0.246) through abstractive synthesis but superior quality (4.30). BERTScore similarly fails as a quality discriminator, compressing all models into a narrow 0.10–0.20 range despite 1.5-point quality differences (3.06–4.55). These metrics penalize good abstractive synthesis that necessarily diverges from the source text. However, BERTScore retains diagnostic value for analyzing generation difficulty under varying information density (Appendix G.2).

Completeness-Quality Orthogonality. Structural coverage and semantic quality represent independent dimensions. On Model Cards, Gemma-3-27B achieves high Completeness (0.734, second-best) yet scores only 3.94 in quality, substantially below Qwen3-30B-A3B-Instruct (4.55) despite comparable Completeness (0.786). Similarly, on Data Cards, Nemotron-Nano-30B-A3B achieves 0.557 Completeness yet scores only 3.37 in quality, far below Qwen3-30B-A3B-Instruct (4.30) despite the latter’s higher Completeness (0.702). This decoupling reveals that high field coverage does not guarantee semantic quality without refinement. Figure 4 illustrates field-level performance patterns: averaged across all baselines, models achieve greater than 0.7 Completeness on explicit fields (*Model Details, Performance Metrics*) containing abundant signals, but drop below 0.5 on abstract categories (*Ethical Considerations, Maintenance*) requiring inference from sparse contexts. Appendix G.2 demonstrates that ground truth information density explains 31% of Completeness variance, confirming source sparsity as the primary

driver of generation difficulty.

5.4.3 Systematic Cognitive Limitations

Universal Data Card Difficulty. All models show consistent performance degradation on Data Cards versus Model Cards (average gap: 0.2–0.4 points), from weakest (GPT-OSS-20B: 3.06 vs 3.38) to strongest (Qwen3-30B-A3B-Instruct: 4.30 vs 4.55). Data Card schemas demand high-granularity lifecycle documentation (privacy protocols, security measures, maintenance plans), whereas papers treat datasets as ancillary artifacts with sparse experimental descriptions. This systematic scarcity imposes the need to synthesize complete profiles from limited signals, under which current models show consistent degradation in completeness and quality metrics.

Faithfulness-Completeness Trade-off. Qwen3-30B-A3B-Instruct achieves near-perfect Faithfulness, avoiding hallucinations, yet exhibits only 0.786 Completeness, systematically omitting 21% of ground truth fields. This precision-recall imbalance aligns with documented long-context retrieval limitations (Liu et al., 2024b), where models struggle to access dispersed information. Linguistically, outputs show high abstraction with hedging phrases and cross-references contrasting with the ground truth’s concrete specifications. Appendix G.3 quantifies these divergences via Log-Odds Ratio analysis, revealing systematic narrative bias.

6 Conclusion

We introduce MetaGAI, a large-scale benchmark comprising 2,541 validated entries constructed through multi-source triangulation and a multi-agent framework. Unlike prior work that aggregates existing documentation, our approach constructs ground truth through rigorous semantic verification and four-dimensional human-in-the-loop validation, addressing incompleteness and inconsistency in human-authored cards. Our experiments reveal that sparse MoE architectures achieve optimal cost-quality performance, while traditional lexical metrics inversely correlate with semantic quality. A fundamental faithfulness-completeness trade-off emerges where models systematically omit fields despite maintaining high factual accuracy. MetaGAI provides critical infrastructure for advancing transparency and reproducibility in GenAI documentation at a production scale.

589 Limitations

590 While MetaGAI provides a strong baseline for auto-
591 mated Model and Data Card generation, it remains
592 limited in scope. Our current framework relies on
593 text-only generation, which overlooks important
594 information embedded in figures, tables, and other
595 non-textual modalities. In addition, MetaGAI treats
596 documentation generation as isolated paper-level
597 tasks, without modeling the complex dependencies
598 among papers, models, and datasets in the broader
599 GenAI ecosystem. Addressing multimodal content
600 and capturing ecosystem-level relationships
601 through structured or graph-based representations
602 are promising directions for future work. Beyond
603 these technical limitations, automated documenta-
604 tion also introduces potential risks. Errors or omis-
605 sions in generated cards may propagate misleading
606 signals about model capabilities, data provenance,
607 or licensing conditions, especially when such ar-
608 tifacts are reused at scale. Without careful val-
609 idation and human oversight, these inaccuracies
610 could undermine transparency efforts or lead to
611 misplaced trust. Addressing multimodal under-
612 standing, ecosystem-level modeling, and robust
613 verification mechanisms are important directions
614 for future work.

615 Ethics Statement

616 We introduce MetaGAI to improve transparency
617 and reproducibility in GenAI through automated
618 documentation. We address the following ethical
619 considerations:

620 **Data Provenance and Licensing.** Our bench-
621 mark triangulates public data from arXiv, GitHub,
622 and Hugging Face, strictly adhering to each plat-
623 form’s terms of use. Source materials are used
624 exclusively for scientific research and card gen-
625 eration. No private information beyond publicly
626 associated author names is collected.

627 **Human Evaluation.** Validation involved volun-
628 teer Ph.D. students in NLP and machine learning.
629 All evaluators were informed of the task nature,
630 with a reasonable workload and no exposure to
631 harmful content.

632 **Risks of Automated Documentation.** We ac-
633 knowledge the inherent risks of LLM hallucina-
634 tions in automated documentation. Generated cards
635 should serve as preliminary drafts requiring hu-
636 man oversight, not expert replacements. Cards may

omit critical safety warnings or hallucinate capabil- 637
ities. Our framework includes Editor Agent verifi- 638
cation and explicitly advocates human-in-the-loop 639
approaches for high-stakes governance scenarios. 640

Environmental Considerations. While bench- 641
mark construction employed large-scale LLMs, 642
our cost-quality analysis demonstrates that sparse 643
Mixture-of-Experts architectures achieve high- 644
quality generation with substantially lower compu- 645
tational requirements, helping reduce energy con- 646
sumption in large-scale documentation efforts. 647

References 648

- 649 Joel Castaño, Silverio Martínez-Fernández, Xavier
650 Franch, and Justus Bogner. 2024. [Analyzing the
651 evolution and maintenance of ml models on hugging
652 face](#). In *Proceedings of the 21st International Con-
653 ference on Mining Software Repositories, MSR ’24*,
654 page 607–618, New York, NY, USA. Association for
655 Computing Machinery.
- 656 Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu
657 Lian, and Zheng Liu. 2025. [M3-embedding: Multi-
658 linguality, multi-functionality, multi-granularity text
659 embeddings through self-knowledge distillation](#).
660 *Preprint*, arXiv:2402.03216.
- 661 Gheorghe Comanici, Eric Bieber, Mike Schaeckermann,
662 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-
663 cel Blistein, Ori Ram, Dan Zhang, et al. 2025. [Gemini 2.5: Pushing the frontier with advanced reason-
664 ing, multimodality, long context, and next generation
665 agentic capabilities](#). *Preprint*, arXiv:2507.06261. 666
- 667 Anamaria Crisan, Margaret Drouhard, Jesse Vig, and
668 Nazneen Rajani. 2022. [Interactive model cards: A
669 human-centered approach to model documentation](#).
670 In *Proceedings of the 2022 ACM Conference on Fair-
671 ness, Accountability, and Transparency, FAccT ’22*,
672 page 427–439, New York, NY, USA. Association for
673 Computing Machinery.
- 674 Andy Donald, Apostolos Galanopoulos, Edward Curry,
675 Emir Muñoz, Ihsan Ullah, M. A. Waskow, Maciej
676 Dabrowski, and Manan Kalra. 2023. [Towards a se-
677 mantic approach for linked dataspace, model and
678 data cards](#). In *Companion Proceedings of the ACM
679 Web Conference 2023, WWW ’23 Companion*, page
680 1468–1473, New York, NY, USA. Association for
681 Computing Machinery.
- 682 Yilun Du, Shuang Li, Antonio Torralba, Joshua B.
683 Tenenbaum, and Igor Mordatch. 2024. [Improving
684 factuality and reasoning in language models through
685 multiagent debate](#). In *Forty-first International Con-
686 ference on Machine Learning*.
- 687 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,
688 Abhinav Pandey, Abhishek Kadian, Ahmad Al-
689 Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten,

690	et al. 2024. The llama 3 herd of models . <i>Preprint</i> , arXiv:2407.21783.	746
691		747
692	Eliahu Horwitz, Nitzan Kurer, Jonathan Kahana, Liel Amar, and Yedid Hoshen. 2025. We should chart an atlas of all the world’s models . <i>Preprint</i> , arXiv:2503.10633.	748
693		749
694		750
695		751
696	Weixin Liang, Nazneen Rajani, Xinyu Yang, Ezinwanne Ozoani, Eric Wu, Yiqun Chen, Daniel Scott Smith, and James Zou. 2024. Systematic analysis of 32,111 ai model cards characterizes documentation practice in ai . <i>Nature Machine Intelligence</i> , 6(7):744–753.	752
697		753
698		754
699		755
700		756
701	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	757
702		758
703		759
704		760
705	Jiarui Liu, Wenkai Li, Zhijing Jin, and Mona Diab. 2024a. Automatic generation of model and data cards: A step towards responsible AI . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 1975–1997, Mexico City, Mexico. Association for Computational Linguistics.	761
706		762
707		763
708		764
709		765
710		766
711		767
712		768
713	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024b. Lost in the middle: How language models use long contexts . <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	769
714		770
715		771
716		772
717		773
718	Shayne Longpre, Robert Mahari, Anthony Chen, Naana Obeng-Marnu, Damien Sileo, William Brannon, Niklas Muennighoff, Nathan Khazam, Jad Kabbara, et al. 2024. A large-scale audit of dataset licensing and attribution in ai . <i>Nature Machine Intelligence</i> , 6(8):975–987.	774
719		775
720		776
721		777
722		778
723		779
724	Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting . In <i>Proceedings of the Conference on Fairness, Accountability, and Transparency</i> , FAT* ’19, page 220–229, New York, NY, USA. Association for Computing Machinery.	780
725		781
726		782
727		783
728		784
729		785
730		786
731		787
732	Burt L Monroe, Michael P Colaresi, and Kevin M Quinn. 2008. Fightin’ words: Lexical feature selection and evaluation for identifying the content of political conflict . <i>Political Analysis</i> , 16(4):372–403.	788
733		789
734		790
735		791
736	NVIDIA, :, Aaron Blakeman, Aaron Grattafiori, Aarti Basant, Abhibha Gupta, Abhinav Khattar, Adi Renduchintala, Aditya Vavre, et al. 2025. Nemotron 3 nano: Open, efficient mixture-of-experts hybrid mamba-transformer model for agentic reasoning . <i>Preprint</i> , arXiv:2512.20848.	792
737		793
738		794
739		795
740		796
741		797
742		798
743		799
744		800
745		801
		802
	OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, et al. 2025. gpt-oss-120b & gpt-oss-20b model card . <i>Preprint</i> , arXiv:2508.10925.	
	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, et al. 2024. Gpt-4o system card . <i>Preprint</i> , arXiv:2410.21276.	
	Jake Poznanski, Aman Rangapur, Jon Borchardt, Jason Dunkelberger, Regan Huff, Daniel Lin, Aman Rangapur, Christopher Wilhelm, Kyle Lo, and Luca Soldaini. 2025. olmocr: Unlocking trillions of tokens in pdfs with vision language models . <i>Preprint</i> , arXiv:2502.18443.	
	Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjar-tansson. 2022. Data cards: Purposeful and transparent dataset documentation for responsible ai . In <i>Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency</i> , FAccT ’22, page 1776–1826, New York, NY, USA. Association for Computing Machinery.	
	Mohammad Shahedur Rahman, Peng Gao, and Yuede Ji. 2025. Hugginggraph: Understanding the supply chain of llm ecosystem . In <i>Proceedings of the 34th ACM International Conference on Information and Knowledge Management</i> , CIKM ’25, page 5997–6005, New York, NY, USA. Association for Computing Machinery.	
	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, et al. 2025. Gemma 3 technical report . <i>Preprint</i> , arXiv:2503.19786.	
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. Self-consistency improves chain of thought reasoning in language models . In <i>The Eleventh International Conference on Learning Representations</i> .	
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, et al. 2020. Huggingface’s transformers: State-of-the-art natural language processing . <i>Preprint</i> , arXiv:1910.03771.	
	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, et al. 2025. Qwen3 technical report . <i>Preprint</i> , arXiv:2505.09388.	
	Xinyu Yang, Weixin Liang, and James Zou. 2024. Navigating dataset documentations in AI: A large-scale analysis of dataset cards on huggingface . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert . In <i>International Conference on Learning Representations</i> .	

803 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang,
 804 Huan Lin, Baosong Yang, Pengjun Xie, An Yang,
 805 Dayiheng Liu, et al. 2025. [Qwen3 embedding: Advancing text embedding and reranking through foundation models](#). *Preprint*, arXiv:2506.05176.

808 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
 809 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
 810 Zhuohan Li, Dacheng Li, et al. 2023. [Judging llm-as-a-judge with mt-bench and chatbot arena](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 46595–46623. Curran Associates, Inc.

815 Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Ci-
 816 hang Xie, Alan Yuille, and Tao Kong. 2022. [Image BERT pre-training with online tokenizer](#). In *International Conference on Learning Representations*.

819 A Definition of Model and Data Card

820 The definitions of Model and Data card for GenAI
 821 are shown in Table 3.

822 B Dataset Characteristics

823 Figure 5 (Top) illustrates the temporal distribution
 824 of our corpus, which shows an increase from 141
 825 triplets in 2022 to 931 in 2025. The domain com-
 826 position is dominated by Computer Vision (cs.CV)
 827 and Computational Linguistics (cs.CL), reflecting
 828 our filtering criteria targeting papers with verified
 829 GitHub and Hugging Face artifacts in generative
 830 modeling. The word count analysis (Bottom) re-
 831 veals distinct content characteristics across sources:
 832 academic papers provide comprehensive technical
 833 context (median ≈ 10 k words), while GitHub and
 834 Hugging Face README files offer concise im-
 835 plementation details (median < 1 k words). This
 836 complementary information structure empirically
 837 justifies our multi-source triangulation approach.
 838 Papers emphasize theoretical foundations and ex-
 839 perimental analysis, whereas repository artifacts
 840 document practical deployment specifications, to-
 841 gether enabling the construction of high-fidelity
 842 ground truth metadata that captures both concep-
 843 tual frameworks and implementation details.

844 C MetaGAI Benchmark Generation 845 Algorithm

846 D MetaGAI Quality Validation Protocols

847 To rigorously assess generated card quality, we
 848 employ a five-dimensional evaluation framework.
 849 This framework adopts three established metrics
 850 (*Faithfulness*, *Relevance*, and *Accuracy*) from prior

Algorithm 1 MetaGAI Benchmark Generation Pipeline

Require: $\mathcal{P} \cup \mathcal{S}$: Full Context (Paper, GitHub, Hugging Face) \mathcal{X} : Segmented chunks of $\mathcal{P} \cup \mathcal{S}$ \mathcal{K} : Schema Keys $\{k_i\}_{i=1}^M$

Ensure: $\hat{\mathcal{C}}$: Final Card

Stage 1: Retriever Agent (Evidence Mapping)

1: $\mathcal{E} \leftarrow \emptyset$
 2: **for** each $k_i \in \mathcal{K}$ **do**
 3: $v_{evidence} \leftarrow \text{Retriever}(\mathcal{X}, k_i)$ \triangleright Extract factual content

4: $\mathcal{E}[k_i] \leftarrow v_{evidence}$

5: **end for**

Stage 2: Generator Agent (Draft Synthesis)

6: $\tilde{\mathcal{C}} \leftarrow \emptyset$

7: **for** each $k_i \in \mathcal{K}$ **do**

8: $\tilde{v}_i \leftarrow \text{Generator}(\mathcal{E}[k_i], \text{Prompt}_{\text{Gen}})$ \triangleright Synthesize structured draft

9: $\tilde{\mathcal{C}}[k_i] \leftarrow \tilde{v}_i$

10: **end for**

Stage 3: Editor Agent (Consolidation)

11: **for** each $k_i \in \mathcal{K}$ **do**

12: \triangleright Editor verifies draft against full context

13: $\hat{v}_i \leftarrow \text{Editor}(\mathcal{P} \cup \mathcal{S}, \tilde{\mathcal{C}}[k_i], \text{Prompt}_{\text{Edit}})$

14: $\hat{\mathcal{C}}[k_i] \leftarrow \hat{v}_i$

15: **end for**

16: **return** $\hat{\mathcal{C}}$

work (Liu et al., 2024a), and introduces two ad- 851
 ditional metrics (*Consistency* and *Usefulness*) to 852
 address the structural and practical requirements of 853
 technical documentation. 854

855 D.1 Metric Definitions

Faithfulness (F). Measures the degree to which 856
 generated content is grounded in provided source 857
 materials (paper, GitHub, Hugging Face), ensuring 858
 absence of unsupported claims or hallucinations. 859

Relevance (R). Evaluates information density 860
 and pertinence to the target field. High-scoring con- 861
 tent strictly adheres to field definitions, avoiding 862
 redundancy, verbosity, or out-of-scope information. 863

Accuracy (A). Assesses factual correctness of 864
 technical details (numerical values, licenses, entity 865
 names). Unlike Faithfulness, which verifies source 866
 alignment, Accuracy validates objective correct- 867
 ness against domain knowledge. 868

Card	Field	Description
Model Card	Model Details	Model architecture, parameter count, developer organization, license, release date, and other model metadata
	Intended Use	Primary applications, target audience, out-of-scope uses, domain restrictions, and other intended uses
	Generative Capabilities	Context length, inference latency, multilingual support, supported modalities, and other capabilities
	Safety Considerations	Safety alignment methods, red teaming results, jailbreak resistance, harm reduction strategies, and other safety analysis
	Training Data	Training corpus, data mixture, filtering pipeline, data provenance, and other training data information
	Performance Metrics	Benchmark results, accuracy metrics, reasoning scores, safety scores, robustness metrics, and other quantitative metrics
	Ethical Considerations	Environmental impact, intellectual property, dual-use risks, societal implications, and other ethical issues
	Caveats & Recommendations	Hardware requirements, deployment guidelines, operational constraints, model limitations, and other recommendations
Data Card	Dataset Details	Dataset name, version identifier, creators and curators, funding, license, text language, and other dataset metadata
	Dataset Structure	Instance count, field schema, data splits, dataset size, and other structural details
	Data Collection	Collection methodology, data sources, collection timeframe, consent process, and other collection details
	Data Processing	Preprocessing steps, cleaning procedures, labeling process, filtering criteria, deduplication, and other processing steps
	Intended Uses	Primary tasks, intended use cases, prohibited uses, commercial restrictions, and other usage information
	Bias & Fairness	Demographic distribution, geographic coverage, social bias analysis, fairness assessment, and other bias or fairness information
	Privacy & Security	Personally identifiable information (PII), anonymization methods, data security protocols, confidentiality measures, and other privacy measures
	Content Analysis	Content types, toxicity analysis, misinformation risks, offensive language, and other content risks
	Legal & Ethical	Copyright status, terms of use, ethical review, compliance requirements, and other legal details
	Maintenance & Updates	Maintenance plan, update frequency, versioning policy, deprecation plan, and other maintenance information
	Distribution & Access	Access mechanism, download instructions, repository link, citation requirements, and other access details
Limitations & Recommendations	Known limitations, usage guidelines, quality caveats, and other recommendations	

Table 3: Definitions of Model and Data Card for Generative AI

Consistency (C). Evaluates internal logical coherence across sentences and fields within a card, detecting contradictions that compromise documentation integrity.

Usefulness (U). Measures practical value for downstream users. Useful content provides specific, actionable insights supporting deployment decisions rather than generic descriptions.

D.2 Scoring Rubric

Each metric is evaluated on a 1–5 Likert scale. Table 4 presents the scoring criteria applied by both human evaluators and LLM judges.

E Inter-Judge Agreement Analysis

To validate the robustness of our LLM-as-a-Judge framework, we computed Pearson correlation coefficients between the three judge models on the strat-

ified evaluation set ($N = 500$). Table 5 reveals substantial inter-judge variation. While Qwen3-235B-A22B-2507 and Llama-3.3-70B-Instruct demonstrate moderate alignment ($\rho = 0.540$), GPT-OSS-120B exhibits markedly different scoring patterns (correlations ≤ 0.226). All observed correlations achieve statistical significance ($p < 0.001$), indicating that score divergence reflects systematic architectural differences in quality assessment rather than random measurement error. This architectural diversity justifies our ensemble averaging strategy: aggregating judgments across models with distinct evaluation perspectives mitigates single-model biases and provides more robust quality estimates than any individual judge.

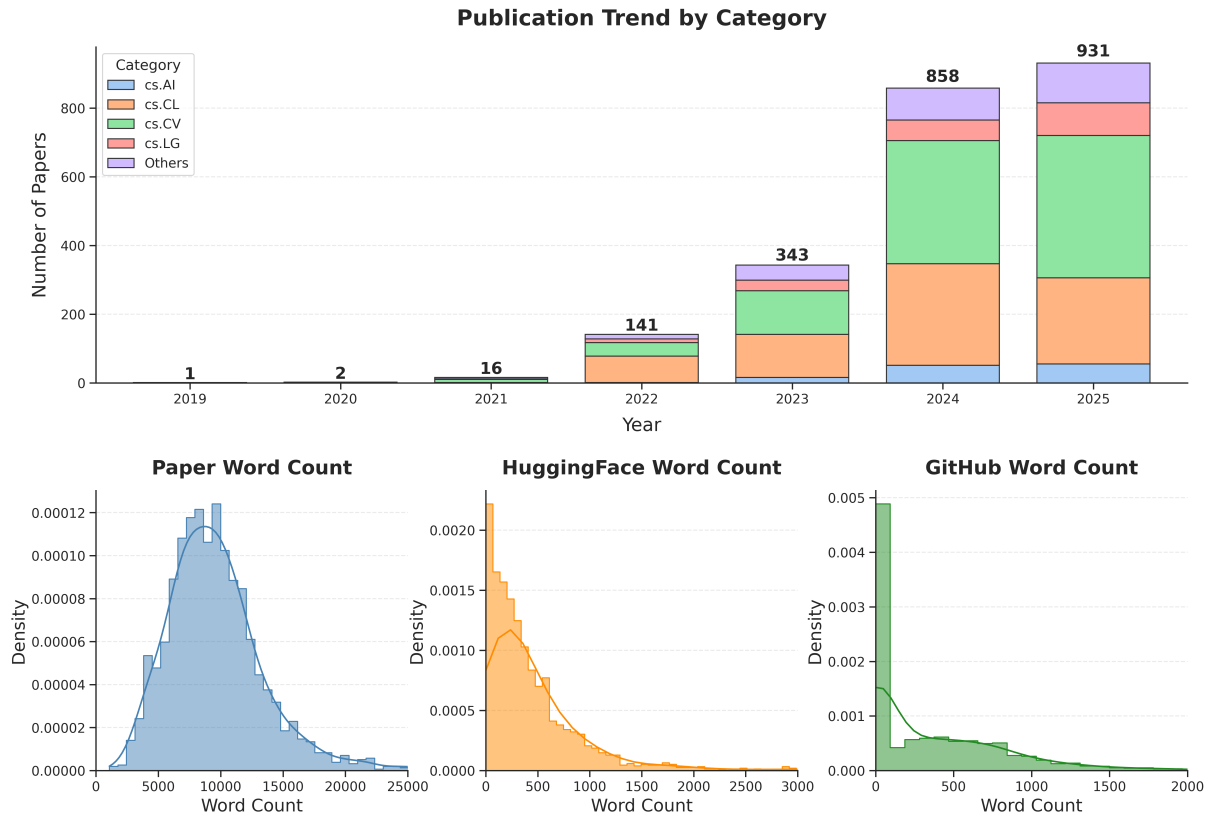


Figure 5: **Dataset Overview.** *Top:* Publication trends of the 2,541 collected triplets (2019–2025), categorized by primary arXiv domain. *Bottom:* Word count distributions across three data sources (Academic Papers, Hugging Face, GitHub), illustrating complementary information granularities.

F Prompt Templates 900

F.1 Correspondence Verification Prompt 901

Correspondence Verification Prompt

System Instruction

You are an expert research evaluator specializing in [dataset/model] identification.

Task Description

Determine whether a research paper introduction and two README files (GitHub and Hugging Face) describe the **same** [dataset/model].

Evaluation Criteria

1. **Existence:** Does the paper explicitly introduce a specific named entity?
2. **Consistency:** Do BOTH the GitHub README and Hugging Face README reference the SAME entity as the paper?
3. **Alignment:** Compare names, domain, size metrics, methodology, and unique features.

Input Data

Paper Content: <Paper text>

GitHub README: <GitHub README text>

Hugging Face README: <HF README text>

Output Format

RELATED: [Yes/No]

CONFIDENCE: [High/Medium/Low]

EXPLANATION: [2-3 sentences citing specific evidence]

902

Metric	Score 1 (Poor)	Score 3 (Acceptable)	Score 5 (Excellent)
Faithfulness	Contains major hallucinations; claims unsupported by sources.	Mostly supported with minor extrapolations.	Fully grounded in source text; no unsupported claims.
Relevance	Largely off-topic, redundant, or contains significant noise.	Generally on-topic with some redundant details.	Concise and strictly focused on field definition.
Accuracy	Contains critical factual errors.	Core facts correct; minor imprecisions present.	Factually precise; all technical details correct.
Consistency	Contains direct logical contradictions.	Generally consistent; minor ambiguities present.	Logically coherent; no internal contradictions.
Usefulness	Vague or generic; provides no practical value.	Provides basic information but lacks depth.	Rich, actionable insights supporting deployment.

Table 4: **Evaluation Scoring Rubric.** Criteria for assessing generated metadata cards on a 1–5 Likert scale across five quality dimensions.

Judge Model	Llama	GPT	Qwen
Llama-3.3-70B-Instruct	1.000	0.101*	0.540*
GPT-OSS-120B	0.101*	1.000	0.226*
Qwen3-235B-A22B-2507	0.540*	0.226*	1.000

Table 5: **Inter-Judge Correlation Matrix.** Pearson correlations computed on 500 samples. All correlations are statistically significant ($*p < 0.001$). **Llama:** Llama-3.3-70B-Instruct, **GPT:** GPT-OSS-120B, **Qwen:** Qwen3-235B-A22B-2507.

F.3 Generator Agent Prompt

905

F.2 Retriever Agent Prompt

Retriever Agent Prompt

System Instruction

You are a metadata classification assistant. Output valid JSON only. STRICTLY select sub-fields from the provided ontology.

Task Description

Classify the provided text chunk into predefined metadata fields for [Dataset/Model] documentation.

Evaluation Criteria

- Relevance:** Assign a score (0–4) indicating how well the chunk describes the field.
- Sub-field Selection:** Select 1–5 sub-fields strictly from the provided “Selectable Sub-fields” list. Do not extract raw text.

Input Data

Source: <[PAPER] / [GITHUB] / [HUGGING-FACE]>

Content: <Text chunk content>

Output Format

```
{ "classifications": [ {
  "field": "field_name",
  "relevance": 3,
  "matched_sub_fields": ["sub_field1"] } ] }
```

Generator Agent Prompt

System Instruction

You are an expert AI Researcher. Generate specific metadata sections for a [Model/Dataset] Card. Output valid JSON only.

Critical Constraint: Subject Focus

- Focus:** Solely on the entity **introduced** in this paper.
- Ignore:** Baselines, pre-training models, or comparisons.

Input Data

Target: <Field Name> (e.g., *Training Data*) | **Context:** <Retrieved chunks>

Output Requirements (Per Sub-field)

- Content:** Summarized factual answer.
- Evidence Quote:** Direct verbatim quote supporting the answer.
- Confidence:** [low, medium, high, certain].
- Source:** Provenance (e.g., “Paper+GitHub”).

Output Format

```
{ "sub_field": {
  "content": "...",
  "evidence_quote": "...",
  "confidence": "high",
  "source": "Paper" } }
```

904

904

906

F.4 Editor Agent Prompt

Editor Agent Prompt

System Instruction

You are the Chief Editor. Consolidate 3 candidate drafts into ONE **concise** entry describing the **proposed** entity.

Task Description

Review three independent drafts (Candidates A, B, C) against the source ground truth. Identify the correct **proposed** entity, filter out hallucinations, and merge valid details.

Evaluation Logic

1. **Identify Proposed Entity:** Describe ONLY the entity introduced in THIS paper. Discard candidates describing baselines.
2. **Verify Evidence:** Check attribution.
3. **Conciseness:** Remove fluff; use direct facts (1–3 sentences).

Input Data

Target: <Field Name> | **Ground Truth:** <Raw chunks>

Candidates: <Draft A>, <Draft B>, <Draft C>

Output Format

```
{
  "selected_candidates": ["Candidate B"],
  "final_content": "...",
  "final_evidence": "...",
  "reasoning": "..."}

```

F.5 Evaluation Judge Prompt

Evaluation Judge Prompt

System Instruction

You are an expert evaluator for Generative AI documentation. Your task is to compare the generated metadata card against the provided triangulated sources (Paper, GitHub, Hugging Face).

Task Description

Evaluate the candidate text based on the five dimensions below. Assign a score from 1 (Poor) to 5 (Excellent) for each dimension.

Evaluation Metrics

- **Faithfulness (F)**
- **Relevance (R)**
- **Accuracy (A)**
- **Consistency (C)**
- **Usefulness (U)**

Analysis Considerations

Classify content into: Accurate facts (verifiable), Vague facts, Logical reasoning, Illogical reasoning, or Acceptable inferences.

Input Data

Sources: <Context>; Target Field: <Field Name>; Candidate: <Text>

Output Format

```
{"scores": {"F": 5, "R": 5, "A": 5, "C": 5, "U": 5}, "reasoning": "..."}

```

Source Category	Bench_Mistral		Bench_GPT-OSS	
	Model	Data	Model	Data
Paper Only	65.9%	61.9%	57.1%	45.1%
Multi-Source	22.4%	26.7%	22.2%	30.7%
GitHub Only	3.5%	3.3%	3.4%	3.9%
Not Provided	8.2%	8.1%	17.4%	20.3%

Table 6: **Information Provenance Distribution.** Percentage of metadata fields derived from single vs. multiple sources across the two benchmark subsets.

G Comprehensive Analysis of Benchmark Characteristics and Baseline Behaviors

To understand the MetaGAI benchmark’s inherent challenges and model behaviors, we conduct granular analysis covering information provenance, information sparsity impact, comparative lexical quality, and cost-efficiency trade-offs.

G.1 Information Provenance and Editor Dynamics

A central design principle of MetaGAI is that high-fidelity card generation requires multi-source triangulation. We validate this through provenance analysis of final ground-truth fields.

Source Distribution. Table 6 demonstrates that while academic papers provide the majority of information (45–66%), a substantial proportion of metadata fields (22–31%) requires multi-source triangulation across paper, GitHub, and Hugging Face artifacts. Data Cards exhibit higher reliance on multi-source integration (27–31%) compared to Model Cards (22%), confirming that dataset documentation is typically fragmented across theoretical descriptions in papers and implementation details in repositories, necessitating cross-source synthesis.

Editor Strategy. Table 7 reveals distinct consolidation patterns. The Best Match strategy dominates (80–95%), while the Merge strategy contributes 5–20% depending on editor architecture. This demonstrates that achieving comprehensive ground truth for certain fields necessitates synthesizing complementary information across multiple generator drafts rather than selecting a single best candidate.

G.2 Impact of Information Sparsity

To explain the field-level performance stratification observed in Figure 4, we analyze the relationship between source text information density and gener-

Metric	Bench_Mistral	Bench_GPT-OSS
<i>Consolidation Strategy</i>		
Best Match	80.0%	94.9%
Merge	20.0%	5.1%
<i>Winning Candidate Source</i>		
Qwen2.5-7B	64.6%	63.8%
Olmo3-7B	27.8%	30.3%
Llama-3.1-8B	7.6%	5.9%

Table 7: **Editor Agent Dynamics.** *Top:* The logic used by the Editor to finalize content. *Bottom:* The distribution of Generator drafts selected as the final ground truth.

948 ation performance using BERTScore as a diagnos-
949 tic indicator of generation difficulty.

950 We hypothesize that the difficulty of generating
951 abstract fields (*Safety Considerations*, *Bias And*
952 *Fairness*) stems primarily from information spar-
953 sity in ground truth sources. High-density fields
954 such as *Model Details* provide abundant explicit
955 signals, enabling high semantic alignment. Con-
956 versely, abstract fields contain sparse or implicit
957 information, challenging models to infer unstated
958 constraints without generating hallucinations.

959 Figure 6 demonstrates a statistically significant
960 positive relationship ($R^2 = 0.31, p = 0.011$) be-
961 tween ground truth completeness (information den-
962 sity) and semantic similarity (BERTScore), indicat-
963 ing that 31% of BERTScore variance is explained
964 by source completeness. This diagnostic analy-
965 sis confirms that a primary bottleneck for current
966 LLMs is the inferential capacity required to popu-
967 late sparse metadata fields from weak textual sig-
968 nals, rather than simple retrieval failure. Note that
969 while BERTScore serves as an effective diagnostic
970 tool for analyzing the impact of information spar-
971 sity on generation difficulty, it does not function as
972 a quality metric for comparing model performance,
973 as discussed in Section 5.4.2.

974 G.3 Comparative Lexical Analysis

975 We perform lexical analysis using Log-Odds Ra-
976 tio (Monroe et al., 2008) to assess linguistic quality
977 differences between generated content and ground
978 truth. We focus on *Model Details* and *Dataset De-*
979 *tails* fields, which require complex generation of
980 open-ended, high-density technical information.

981 Figure 7 reveals systematic linguistic divergence.
982 Baseline models (red) exhibit narrative-oriented
983 patterns, frequently employing hedging language
984 (*suggests*, *implying*) that summarizes experimen-
985 tal narratives rather than generating factual spec-

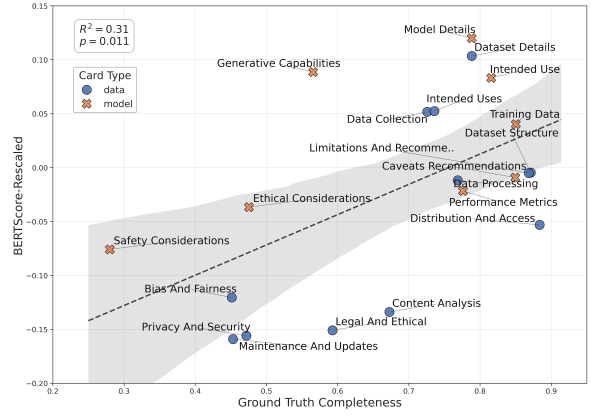


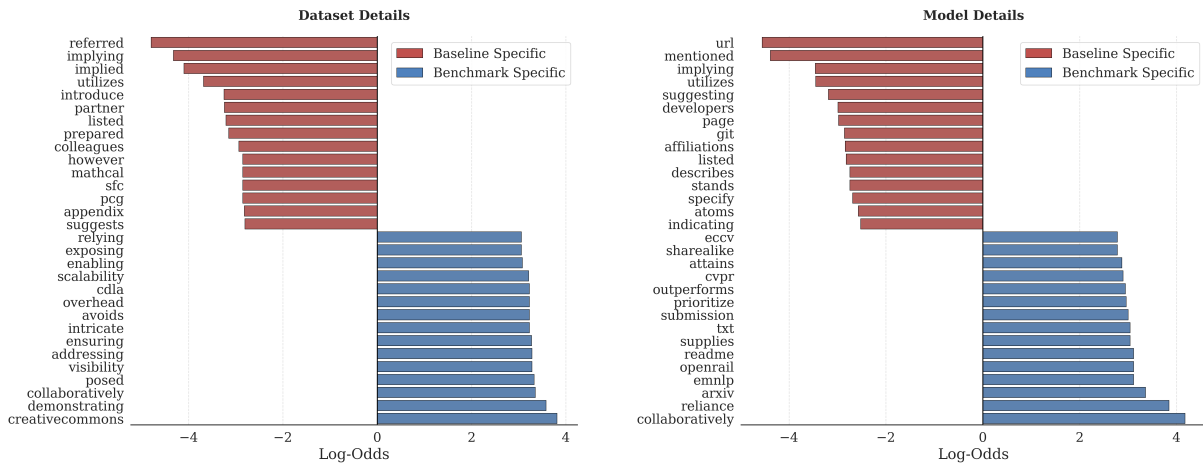
Figure 6: **Information Sparsity Impact on Generation Difficulty.** The relationship ($R^2 = 0.31, p = 0.011$, OLS linear regression) between ground truth information density (X-axis, measured by completeness) and generation difficulty (Y-axis, measured by BERTScore as a diagnostic indicator). Fields with high sparsity (lower completeness) result in significantly lower generation success, highlighting the challenge of inferring abstract metadata from sparse signals.

986 ifications. In contrast, MetaGAI ground truth
987 (blue) uniquely captures high-value technical en-
988 tities (*CDLA*, *OpenRAIL*) and precise licensing
989 terms. This demonstrates that our pipeline suc-
990 cessfully generates concrete implementation speci-
991 fications (artifact metadata) rather than paper nar-
992 ratives.

993 G.4 Cost-Efficiency Analysis

994 To further evaluate the economic feasibility of de-
995 ploying these models at the scale of millions of
996 papers, we analyze the trade-off between genera-
997 tion quality and inference cost. We define a Cost
998 Index normalized by standardizing token consump-
999 tion (1M input / 0.2M output) based on OpenRouter
1000 pricing.

1001 Figure 8 visualizes this cost-quality landscape.
1002 The analysis reveals a distinct Pareto frontier do-
1003 minated by open-weight architectures. Specifically,
1004 the sparse MoE model Qwen3-30B-A3B-Instruct
1005 achieves the highest qualitative scores while main-
1006 taining one of the lowest costs per task. In contrast,
1007 proprietary models like Gemini-2.5-Flash and GPT-
1008 5-Mini, while capable, sit far to the right of the ef-
1009 ficient frontier, incurring significantly higher costs
1010 (up to $5\times$) without a proportional gain in genera-
1011 tion fidelity. This disparity suggests that for structured
1012 generation tasks, specialized open-weight models
1013 offer a far superior return on investment.



(a) Dataset Details Divergence

(b) Model Details Divergence

Figure 7: **Lexical Divergence (Log-Odds Ratio)**. **Left (Red)**: Words over-represented in Baselines, indicating narrative bias and format hallucinations. **Right (Blue)**: Words specific to MetaGAI, highlighting the capture of high-value entities and technical specifications.

H Case Study

To demonstrate the efficacy of the MetaGAI pipeline in complex retrieval and synthesis scenarios, we present two qualitative case studies illustrating how multi-source triangulation and multi-agent architecture mitigate common failure modes including information omission and incomplete evidence coverage.

H.1 Multi-Source Triangulation

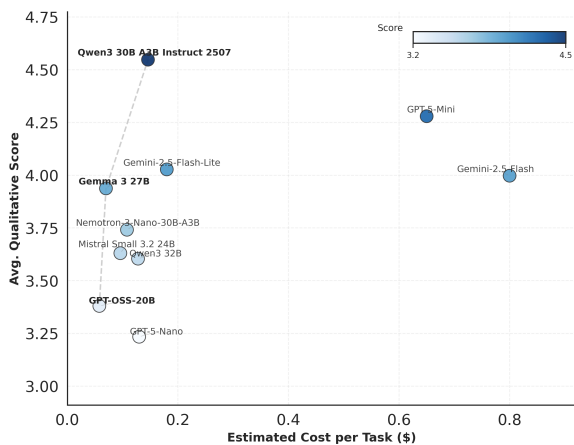
We analyze the *Model Architecture* generation for the paper “Low-light Image Enhancement via Breaking Down the Darkness”. As illustrated in Figure 9, high-fidelity output requires fusing distinct information modalities: the paper defines macro-level architectural topology, GitHub specifies implementation hyperparameters typically omitted from manuscripts, and Hugging Face validates entity alignment and deployment availability.

This case validates the pipeline’s capacity to resolve the granularity gap between conceptual descriptions and implementation specifications. While the paper provides the architectural framework, the GitHub repository supplies concrete hyperparameters necessary for reproducibility. The Editor Agent (GPT-OSS-20B) correctly identified that Candidate B integrated evidence across all three sources, selecting it over simpler drafts that merely paraphrased the abstract.

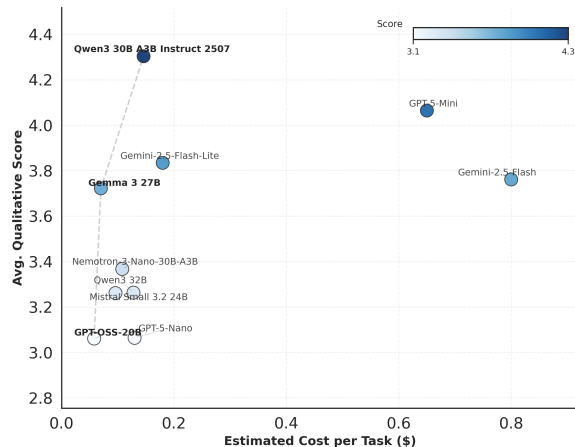
H.2 Editor-Driven Information Synthesis

While retrieval ensures evidence access, individual Generator Agents frequently exhibit narrow focus, generating information from specific paper sections while overlooking complementary content. The Editor Agent provides critical synthesis capabilities to resolve this incompleteness.

Continuing with the same paper, we analyze the *Robustness Metrics* field for the “Bread” model. As shown in Figure 10, the two Generator Agents produced factually accurate but incomplete drafts: Candidate A (Olmo3-7B) generated only comparative benchmark results, while Candidate B (Qwen2.5-7B) focused exclusively on ablation study findings. The Editor Agent (GPT-OSS-20B) recognized that both candidates captured orthogonal robustness dimensions: external performance validation versus internal architectural stability. Rather than selecting a single candidate, the editor merged non-redundant information to generate comprehensive coverage addressing both comparative benchmarking and component ablation analysis. This synthesis capability demonstrates a critical advantage of our multi-agent architecture over single-model generation approaches, which typically exhibit selection bias toward particular evidence types.



(a) Model Card Generation



(b) Data Card Generation

Figure 8: **Cost-Efficiency Analysis.** The Pareto frontier (dashed grey line) highlights models that offer optimal quality for a given cost. Note the significant gap between the efficient open-weight frontier and closed-source proprietary models.

Triangulation Logic for “Bread” Architecture

[Source] Paper

“Figure 4 shows the overall architecture... comprises an illumination adjustment net (IAN), an adaptive noise suppression net (ANSN), and a color adaption net (CAN).”

[Source] GitHub

(Code Implementation)
class IAN(nn.Module):
layers = [Conv3x3] * 3
3 down-sampling
3 up-sampling
out = torch.sigmoid(x)

[Source] Hugging Face

(Metadata Verification)
Task: Image Restoration
Tags: [PyTorch, bread]
Confirms model identity and availability.

MetaGAI Editor Decision (GPT-OSS-20B)

Selected Candidate: Qwen2.5-7B (Draft B)

Source Types Verified: Paper, GitHub, Hugging Face

Reasoning: “Candidate B accurately captures the proposed Bread architecture, whereas the other candidates lack relevant details.”

Final Generated Content:

Bread is a three-stage CNN pipeline:

- 1. Illumination Adjustment Network (IAN)** predicts a relative illumination map to brighten the luminance.
- 2. Adaptive Noise Suppression Network (ANSN)** denoises the brightened luminance using multiple suppression strengths fused by a Noise Fusion Module.
- 3. Color Adaption Network (CAN)** refines chrominance guided by the denoised luminance.

Synthesis: Each sub-network follows a symmetric encoder-decoder with 3 down-sampling and 3 up-sampling layers, 3x3 convolutions, and a sigmoid output (except ANSN).

Figure 9: **Multi-Source Triangulation.** The Editor Agent synthesizes the high-level topology from the Paper, specific layer depths and activation functions from GitHub, and validates alignment via Hugging Face, resulting in a complete specification.

Synthesis Logic for “Bread” Robustness

Candidate A (Olmo3-7B)

(Focus: External Benchmarks)

“Bread achieves the lowest **NIQE and LOE** scores on the DICM, NPE, and VV datasets, outperforming all competing methods.”

✗ *Missing internal analysis.*

Candidate B (Qwen2.5-7B)

(Focus: Internal Ablation)

“**Ablation studies** show that removing key modules such as the noise fusion module causes marked drops in PSNR/SSIM.”

✗ *Missing comparative stats.*

MetaGAI Editor Decision (GPT-OSS-20B)

Action: Merge Non-Redundant Information

Reasoning: “Both candidates describe valid, orthogonal aspects of robustness. Candidate A covers external metrics (NIQE/LOE), while Candidate B covers internal validity (Ablation). Merging creates a holistic summary.”

Final Generated Content:

- **External Performance:** Bread achieves the lowest NIQE and LOE scores on the DICM, NPE, and VV datasets.
- **Internal Stability:** Ablation studies confirm framework robustness; removing the noise fusion module causes marked drops in PSNR/SSIM.

Figure 10: **Resolving Incompleteness via Synthesis.** The Editor Agent detects that Candidates A and B focus on orthogonal aspects of robustness (External Benchmarks vs. Internal Ablation). Instead of selecting a single winner, the Editor merges them to generate a comprehensive entry.