# Cost-Aware Structured Generation: High-Fidelity Synthesis via Hybrid RAG and Adaptive Conditional Compute

**Anand Wankhade**
anandwankhade0008@gmail.com

## Abstract

Deploying large language models (LLMs) for structured content generation presents a fundamental tension between quality requirements and computational costs. Standard retrieval-augmented generation (RAG) approaches struggle with complex multi-hop domain questions due to context loss and retrieval noise, while state-of-the-art models impose prohibitive per-query costs at scale. We introduce **Cost-Aware Structured Generation**, a domain-agnostic 14-stage pipeline that integrates Hybrid RAG with Reciprocal Rank Fusion, conditional compute routing based on query difficulty, and adaptive voting strategies with domain-specific weighting. Our architecture achieves **93.5% overall precision** with a hallucination rate of only **0.31%** (1 error in 325 human-reviewed samples from 2021-2025) on a validation dataset of 1,200+ complex domain questions spanning 2007-2025, requiring 200+ field structured outputs. Through intelligent routing that selectively invokes expensive models only for high-complexity queries, we demonstrate **84% improvement in content richness** while reducing costs by **50.4%** compared to baseline Self-Consistency ensembles. Validation on aerospace engineering certification exams demonstrates production viability; the architecture generalizes to any domain requiring high-fidelity structured generation (legal, medical, financial, educational). **This work establishes that architectural engineering—not merely model scaling—provides a viable pathway to production-grade structured generation with strict precision requirements**.

## 1. Introduction

### 1.1 Background and Motivation

The proliferation of large language models has enabled unprecedented capabilities in natural language generation, yet production deployment in high-stakes domains reveals a critical gap between free-form text generation and **structured content generation**—the task of producing strictly formatted outputs with high factual accuracy requirements [1] . Applications spanning legal contract analysis [2] , medical documentation [3] , educational content creation [4], and financial reporting [5] demand outputs conforming to complex schemas (often 100+ fields) with near-zero tolerance for hallucination.

Industry practitioners face an acute challenge: the tension between computational costs and quality requirements. Current paradigms offer practitioners a forced binary choice. Retrieval-Augmented Generation (RAG) [6] [7] reduces hallucination by grounding responses in retrieved documents, but standard implementations suffer from context loss on multi-hop reasoning tasks and retrieval noise degradation [8]. Alternatively, invoking state-of-the-art models on every query achieves high precision but imposes economically unsustainable token costs at scale—often $0.20-0.50 per complex query [9].

Chain-of-Thought (CoT) reasoning [10] and Self-Consistency ensembles [11] improve reliability by generating multiple candidate responses, but they exacerbate the cost problem by multiplying model invocations. A production system processing 10,000 queries daily would incur $2,000-5,000 in inference costs using standard ensemble approaches, creating a barrier to deployment in cost-sensitive applications.

### 1.2 Problem Statement and Research Gap

We identify three critical limitations in current structured generation approaches that transcend specific domains:

**L1. Retrieval-Content Mismatch:**
Standard RAG systems employ either dense embedding-based retrieval [12] or sparse keyword matching [13], [14] . Dense retrieval captures semantic similarity but misses exact technical terminology critical in specialized domains (legal citations, medical codes, financial regulations, technical standards). Sparse retrieval handles keywords well but fails on paraphrased queries. When complex questions require multi-hop reasoning across multiple retrieved contexts, single-strategy retrieval loses critical information, degrading downstream precision. Recent hybrid approaches [15], [16], [17] demonstrate complementarity, but **systematic application to schema-compliant structured generation** remains underexplored.

**L2. Uniform Compute Allocation:**
Running SOTA models on every query ignores substantial variance in query complexity. Empirical analysis across domains shows that ~40-60% of queries involve straightforward fact retrieval or template instantiation, yet these simple

cases incur the same computational cost as complex multi-hop reasoning tasks. While cost-aware routing frameworks exist [18], [19], [20] , they focus on free-form QA; **integration with structured schema generation** introduces validation constraints not addressed in prior work.

**L3. Quality-Cost Tradeoff with Schema Compliance:**
Existing approaches force practitioners into a Pareto frontier where improving precision necessitates proportional cost increases. Self-Consistency with $k=3$ samples triples costs; ensemble methods with 5+ models can quintuple expenses. Recent cost-aware frameworks exist: ThriftLLM [21] formulates budget-constrained ensemble selection via knapsack optimization; ε-constrained ensemble approaches [22] model quality-cost bi-objectives for Pareto-efficient model selection; multi-agent debate with entropy compression [23] reduces hallucinations while lowering token usage; telecom MoE systems [24] use DRL gating for cost-aware expert selection. However, systematic reviews [25] note that only ~9% of hallucination architectures implement **integrated, dynamic cost-quality optimization**, and **none address schema compliance requirements** for structured generation—precisely the gap our work fills.

## 1.3 Research Contributions

Our primary objective is to introduce the **Cost-Aware Structured Generation** paradigm, which maximizes precision and content richness while minimizing token costs through intelligent architectural routing. We make four contributions applicable across domains:

**C1. Hybrid RAG Integration Framework:**
We systematically integrate dense (BAAI/bge-m3 [26] ) and sparse (BM25) [13] retrieval using Reciprocal Rank Fusion (RRF) [27], adapting RetHyb-RRF [16] methodology for structured generation tasks. Our fusion strategy handles both semantic paraphrasing and exact-match requirements common in technical domains.

**C2. Schema-Aware Conditional Compute:**
Extending difficulty-aware routing [18], [19], [20] we introduce a classification-based router that dynamically selects model tiers while ensuring all paths produce schema-compliant outputs. Unlike prior work on free-form QA, our routing accounts for field-level validation requirements.

**C3. Adaptive Model Weighting Framework:**
Following ensemble optimization principles [22] ,we define a configurable weighting system where domain-specific strategies (e.g., mathematical rigor vs. pedagogical clarity vs. regulatory precision) are mapped to heterogeneous SOTA models based on their empirically-observed strengths. This allows practitioners to **configure the system for their domain** without architectural changes.

**C4. Multi-Round Debate Orchestration:**
We operationalize adversarial debate mechanisms [23] with a two-round protocol: initial consensus voting (Round 1) followed by conditional escalation to an impartial judge model (Round 2) when disputes persist. This provides a systematic conflict-resolution pathway for high-stakes decisions.

**Empirical Validation:**
We validate on aerospace engineering certification exams (GATE, N=1,200+, spanning 2007-2025) as a representative high-stakes domain with strict accuracy requirements: - **93.5% overall precision** - **0.31% hallucination rate** (1 error in 325 human-reviewed samples from 2021-2025) - **+84% improvement in content richness** vs. single-model baselines - **50.4% cost reduction** vs. baseline Self-Consistency ($k=3$)

**Generalizability:**
While validated on aerospace engineering, the architecture is **domain-agnostic**. The same 14-stage pipeline applies to legal (contract analysis), medical (clinical documentation), financial (regulatory reporting), or educational (study guide generation) domains by:

1. Replacing the domain corpus (textbooks → case law / medical literature / financial regulations)

2. Adjusting adaptive weights to domain priorities (safety → legal precision / medical accuracy / numerical exactness)

 3. Customizing the target schema (200 aerospace fields → domain-specific requirements)

This work establishes that structured generation quality is not solely a function of model capability, but emerges from **architectural engineering** that coordinates retrieval, routing, and ensemble strategies.

## 2. Related Work

**Retrieval-Augmented Generation.**
The RAG paradigm [7] addresses hallucination by grounding LLM responses in retrieved documents. Dense Passage

Retrieval (DPR) [12] pioneered bi-encoder architectures for semantic similarity, while BM25 [12] remains competitive for keyword-centric retrieval [24]. Recent hybrid approaches combine both: Blended RAG [17] reports 87% retriever accuracy on TREC-COVID; RetHyb-RRF [17] demonstrates MAP@3: 0.897 vs. 0.768 dense-only on HaluBench via RRF fusion [27]. Our work extends hybrid RAG to **schema-compliant structured generation**, a setting not addressed in prior benchmarks.

**Conditional Computation and Cost-Aware Routing.**
Early exit strategies [28] and adaptive depth networks [29] reduce inference costs by terminating computation when confidence thresholds are met. RouteLLM [19] and FrugalGPT [18] route queries to model tiers based on learned classifiers. ThriftLLM [21] formulates budget-constrained ensemble selection as a knapsack problem. BEST-Route [20] extends this to sample-count optimization. Our work integrates conditional routing with **structured schema enforcement**, accounting for validation requirements beyond final answer correctness.

**Structured Output Generation.**
Template-based methods [30] offer high precision but limited flexibility. Neural approaches [31], [32], [33] learn schema mappings end-to-end but struggle with complex nested structures. Instructor [34] enforces structured outputs via post-processing. Our system integrates schema enforcement **throughout the pipeline** stages, not just at output validation.

**Ensemble Methods and Multi-Agent Debate.**
Self-Consistency [10] samples multiple CoT trajectories and selects the most frequent answer. Mixture-of-Agents (MoA) [35] iteratively refines outputs using multiple models. Multi-agent debate frameworks [36] use adversarial cross-checking to reduce hallucinations. Our adaptive voting differs by applying **configurable model-specific weighting** tailored to domain characteristics, rather than uniform aggregation.

**Position of This Work.**
To our knowledge, this is the first systematic integration of hybrid RAG, conditional routing, adaptive voting, and multi-round debate for **domain-agnostic structured generation** with strict precision and schema compliance requirements. While individual techniques exist, no prior framework combines all four with explicit cost-quality optimization.

# 3. Methodology: Domain-Agnostic Architecture

## 3.1 Design Principles

Our architecture follows four core principles applicable across domains:

**P1. Modular Decomposition:** The pipeline decomposes generation into 14 sequential stages with clear interfaces, enabling domain customization at specific points (corpus loading, schema definition, weight configuration) without altering the orchestration logic.

**P2. Quality Gates:** Each stage enforces validation before proceeding (schema compliance, confidence thresholds, consensus requirements), preventing error propagation.

**P3. Cost Decision Points:** The architecture provides three cost gates: (1) cache lookup, (2) complexity-based routing, (3) conditional debate invocation, allowing practitioners to tune the cost-quality tradeoff.

**P4. Configurability:** Domain-specific elements (retrieval corpus, weight strategies, schema, prompts) are externalized to configuration files, enabling deployment across domains without code changes.

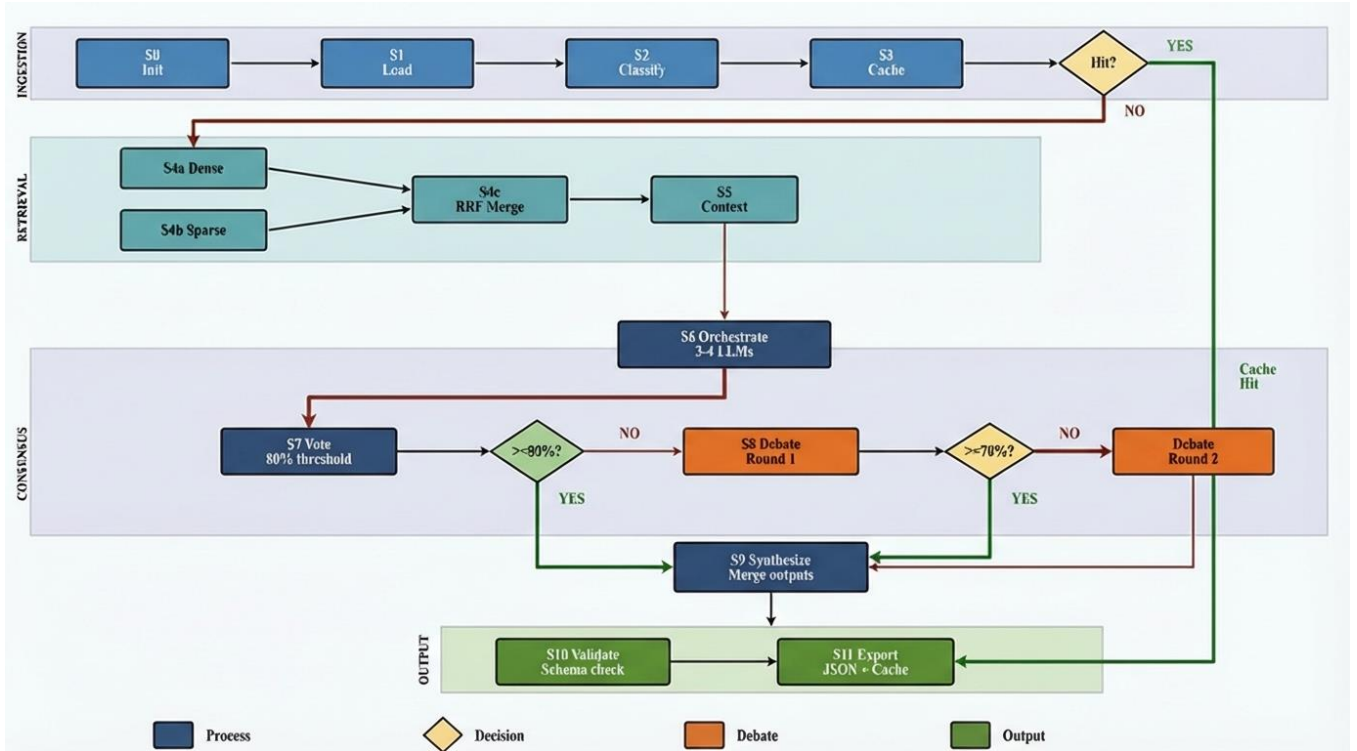## 3.2 14-Stage Pipeline Architecture

*Figure 1: The 14-stage pipeline architecture showing the flow from input query through Ingestion (S0-S3), Retrieval (S4a-S5), Consensus (S6-S9), and Output (S10-S11) phases. Each stage is configurable for domain adaptation via YAML configuration files.*

The pipeline executes 14 sequential stages:

**Phase 1: Ingestion & Analysis (Stages 0-3)** - **Stage 0:** Initialization (Load domain corpus, initialize vector/sparse indices, load model clients) - **Stage 1:** Query Loading (Parse input, extract metadata, identify attachments) - **Stage 2:** Classification (Classify query complexity, content type, media type → route to weight strategy) - **Stage 3:** Cache Check (97% similarity threshold; skip expensive stages if cached)

**Phase 2: Retrieval & Context Assembly (Stages 4-5)** - **Stage 4a:** Dense Retrieval (BAAI/bge-m3 embeddings → top-10 semantic matches) - **Stage 4b:** Sparse Retrieval (BM25 inverted index → top-10 keyword matches) - **Stage 4c:** RRF Fusion (Reciprocal Rank Fusion with $k = 60$→ top-6 chunks) - **Stage 5:** Multimodal Consensus (If images/attachments present, run 3-model description voting)

**Phase 3: Generation & Consensus (Stages 6-9)** - **Stage 6:** Model Orchestration (Parallel invocation of 3-4 heterogeneous SOTA models) - **Stage 7:** Voting Engine (Field-level consensus with adaptive weights @ 80% threshold) - **Stage 8:** Debate Orchestrator (Round 1 @ 70%; Round 2 with judge model if needed) - **Stage 9:** Synthesis Engine (Deterministic merge: Debate > Voting > Fallback)

**Phase 4: Validation & Output (Stages 10-11)** - **Stage 10:** Validate (Schema validation, field completeness check, quality score computation) - **Stage 11:** Export (JSON output, cache update, S3/DynamoDB persistence, metadata logging)

**Domain Customization Points:** - **Stage 0:** Load domain-specific corpus (legal cases, medical literature, financial docs, textbooks) - **Stage 2:** Configure classification → weight strategy mappings per domain priorities - **Stage 6:** Select model pool based on domain requirements (vision models for medical imaging, code models for technical docs) - **Stage 9:** Define schema for target domain (legal contract fields, clinical note sections, course material structure)

## 3.3 Hybrid RAG with Reciprocal Rank Fusion

We implement a domain-agnostic hybrid retrieval strategy combining complementary approaches:

**Dense Retrieval (Stage 4a):**
BAAI/bge-m3 [26] (1024-dim multilingual embeddings) captures semantic similarity. For query $q$ with embedding $\mathbf{q}_d$, retrieve top-10 documents by cosine similarity:

$$\text{sim}(q, d_i) = \frac{\mathbf{q}_d \cdot \mathbf{d}_i}{|\ |\mathbf{q}_d|\ |\cdot|\ |\mathbf{d}_i|\ |}$$

**Sparse Retrieval (Stage 4b):**
BM25 parameters $k_1 = 1.2$, $b = 0.75$ on inverted index:

$$\text{BM25}(q, d_i) = \sum_{t \in q} \text{IDF}(t) \cdot \frac{\text{TF}(t, d_i) \cdot (k_1 + 1)}{\text{TF}(t, d_i) + k_1 \cdot (1 - b + b \cdot |d_i|/\text{avgdl})}$$

**Reciprocal Rank Fusion (Stage 4c):**
Merge heterogeneous rankings using RRF [27] with constant $k = 60$:

$$\text{RRF}(d) = \sum_{r \in \{\text{dense, sparse}\}} \frac{1}{60 + \text{rank}_r(d)}$$

We select the top-6 documents by RRF score for downstream prompting. This fusion strategy is **domain-agnostic**: it works identically whether retrieving from legal cases, medical journals, or engineering textbooks.

**Query Expansion (Optional):**
For structured queries with enumerable options (e.g., multiple-choice questions, parameter selection), we optionally append structured elements to the query text to improve retrieval context, following Blended RAG [17] recommendations.

## 3.4 Schema-Aware Conditional Compute

**Classification (Stage 2):**
We use a lightweight classifier (Gemini 2.0 Flash, temperature=0.1, ~$0.001/query) to estimate query complexity on a 1-10 scale. Classification considers: - **Domain breadth:** Number of sub-fields required (1-5) - **Reasoning depth:** Inference chain length (1=lookup, 5=multi-hop derivation) - **Schema coverage:** Estimated % of target fields requiring model generation vs. template filling - **Ambiguity:** Presence of underspecified constraints or conflicting information

The overall difficulty score is:

$$\text{difficulty} = 0.30 \cdot \text{breadth} + 0.35 \cdot \text{depth} + 0.25 \cdot \text{coverage} + 0.10 \cdot \text{ambiguity}$$

**Routing Logic (Stage 6):**
Based on difficulty score, we route to three computational tiers:

| Tier | Difficulty | Model Selection | Avg Cost | Use Case |
|---|---|---|---|---|
| 1 | 1-3 | Template + Rules | ~$0.01 | Straightforward lookups |
| 2 | 4-5 | Mid-tier LLMs | ~$0.08 | Standard complexity |
| 3 | ≥6 | Full SOTA ensemble + GPT-5.1 | ~$0.35 | High complexity |

**Key Difference from Prior Work:**
Unlike free-form QA routing [18], [19], [20] , our system ensures **all tiers produce schema-compliant outputs**. Tier 1 uses deterministic template filling validated against the schema; Tier 2 uses constrained decoding; Tier 3 uses full generation + validation. This guarantees downstream consumers receive predictable structure regardless of routing decision.

## 3.5 Adaptive Model Weighting Framework

We define a **configurable weighting system** where domain practitioners specify priority dimensions, and the system maps these to model strengths.

**Example: Aerospace Engineering Configuration**

| Strategy | DeepSeek R1 | Claude 4.5 | Gemini 2.5 Pro | GPT-5.1 | Priority |
|---|---|---|---|---|---|
| MATH_WEIGHTED | 0.40/0.33 | 0.30/0.22 | 0.30/0.22 | -/0.22 | Mathematical rigor |
| CONCEPTUAL | 0.30/0.15 | 0.40/0.35 | 0.30/0.25 | -/0.25 | Pedagogical clarity |
| VISION | 0.33/0.25 | 0.33/0.25 | 0.33/0.25 | -/0.25 | Image interpretation |

*Note: Values shown as without GPT-5.1 / with GPT-5.1 (when difficulty ≥ 6).*

**Generalization to Other Domains:**

For **legal contract analysis**, practitioners might define: - REGULATORY_PRECISION (weight models based on citation accuracy) - CASE_LAW_COVERAGE (weight models producing comprehensive precedent analysis)

For **medical documentation**, priorities might include: - CLINICAL_SAFETY (heavily weight models with low hallucination on drug interactions) - ICD_CODING_ACCURACY (weight models correctly mapping diagnoses to codes)

The configuration file (weights_config.yaml) is **domain-specific**, but the voting engine implementation is **domain-agnostic**: it simply applies weights to model outputs regardless of what those weights represent.

**Voting Mechanism (Stage 7):**
For each schema field $f$, aggregate model responses $\{r_i^f\}$ using configured weights $\{w_i\}$:

$$\text{score}(r_i^f) = w_i \cdot \text{confidence}_i^f$$

Select response with highest weighted confidence if > 80% consensus achieved; else trigger debate.

## 3.6 Zero-Cost Deterministic Synthesis

The Synthesis Engine (Stage 9) merges non-conflicting fields **without LLM invocation**, following priority order:

**Priority 1:** Debate-resolved fields (highest fidelity from Round 1/2)
**Priority 2:** Voting-converged fields (high confidence, ≥80% agreement)
**Priority 3:** Fallback merge rules (domain-agnostic logic): - **Arrays:** Union with semantic deduplication (embedding similarity >0.9 → merge identical items) - **Numerics:** Weighted average by model confidence scores - **Timestamps/IDs:** Most recent or most authoritative source - **Nested Objects:** Recursive field-by-field resolution

**Quality Score (Domain-Agnostic):**
Computable metric for every output:

$$Q = 0.25C_{\text{conf}} + 0.30C_{\text{cons}} + 0.20C_{\text{eff}} + 0.15R_{\text{rag}} + 0.10F_{\text{comp}}$$

Where: - $C_{\text{conf}}$ = Average model confidence - $C_{\text{cons}}$ = Consensus rate across models - $C_{\text{eff}}$ = Debate efficiency (1 if no debate, 0.7 if Round 1, 0.4 if Round 2) - $R_{\text{rag}}$ = RAG relevance score - $F_{\text{comp}}$ = Field completion percentage

## 3.7 Multi-Round Debate Protocol

**Safety Protocol (Red Line):**
If **any** model flags a critical validation failure (e.g., answer_validation.is_correct = False in our aerospace implementation, or regulatory_compliance_failed in legal domains), the Voting Engine immediately halts consensus and routes to human expert review. This prevents propagating potentially high-stakes errors.

**Debate Orchestrator (Stage 8):**

**Round 1:**
For disputed fields (consensus <80%), invoke original model pool with contrastive prompt showing disagreement and RAG context. Target threshold: 70% agreement. Batch disputed fields (max 5 per API call) to reduce costs.

**Round 2:**
If Round 1 fails to resolve, dynamically add a designated "judge" model (GPT-5.1 in our implementation, but configurable) as an impartial tiebreaker to select the response that best satisfies factual accuracy and domain standards.

**Resolution Rate (Empirical, Aerospace Validation):**
- **83.0% passed without debate** (direct consensus ≥80%) - Round 1: 11.9% of total queries required debate, resolving 70% of disputes - Round 2: Additional 5.1% escalated to judge model - **Total resolution:** 97.0% rate - Remaining 3.0% route to conservative fallback or human review

**Cost Overhead:**
Debate invoked on ~17% of queries, adding ~4-5% to total pipeline budget.

## 4. Validation: Aerospace Engineering Case Study

### 4.1 Dataset and Experimental Setup

**Validation Domain:** GATE (Graduate Aptitude Test in Engineering) Aerospace Engineering
**Rationale:** Representative high-stakes domain with strict accuracy requirements, technical complexity (fluid mechanics, thermodynamics, structures, flight dynamics), and publicly available ground truth.

**Dataset Composition:** - **Size:** 1,200+ questions (2007-2025) - **Content:** Average 78 words, technical terminology (Reynolds number, Mach number, Euler-Bernoulli theory) - **Modality:** 42% include diagrams (free-body diagrams, p-v diagrams, airfoil cross-sections) - **Question Types:** MCQ (4 options), Numerical Answer Type, Multi-Select MCQ

**Target Schema:** 200+ fields across 5 hierarchical tiers: - **Tier 0 (8 fields):** Classification metadata
- **Tier 1 (30+ fields):** Core answer, step-by-step solution, formulas, references
- **Tier 2 (40+ fields):** Pedagogical content (common mistakes, mnemonics, flashcards)
- **Tier 3 (35+ fields):** Advanced learning (real-world applications, edge cases)
- **Tier 4 (20+ fields):** Quality scores, processing metadata

**Domain Configuration:** - **Corpus:** 25,000+ chunks (detailed notes derived from aerospace textbooks: 8,000+, annotated video transcripts: 7,000+) - **Weight Strategies:** 7 configurations (MATH_WEIGHTED, CONCEPTUAL, VISION, etc.) - **Model Pool:** Gemini 2.5 Pro, Claude Sonnet 4.5, DeepSeek R1, GPT-5.1

**Evaluation Protocol:**

**Human Evaluation (N=325):**
To assess production quality, we conducted rigorous human evaluation on all 325 questions from 2021-2025 (last 5 years), representing the most current and relevant content. A domain expert manually evaluated all samples systematically. This recent-question focus ensures our precision metrics reflect performance on contemporary exam content while the full 1,200+ question dataset demonstrates scalability.

**Metrics:** - **Precision:** % of atomic factual claims verified correct against ground truth / expert knowledge - **Hallucination Rate:** % of fields containing fabricated information (false formulas, non-existent entities, counterfactual statements) - **Quality Rating:** 1-5 Likert scale on informativeness and pedagogical value

**Ablation Study (N=200):**
Compare 5 configurations on a 200-question subset selected to cover all question types, difficulty levels, and years: 1. Full Pipeline (our system) 2. Best-of-N (select highest-confidence model response per question) 3. Individual Model Baselines (Claude, DeepSeek, Gemini, GPT-5.1)

*Note: N=200 enables statistically significant ablation analysis across 10 content richness dimensions while maintaining computational tractability.*

**LLM-as-Judge Pairwise (N=1,270):**
An LLM judge evaluates consensus output vs. each individual model baseline across 10 semantic fields. Reports win rates and effect sizes (Cohen's d).

## 5. Results

### 5.1 Human Evaluation Results

**Table 1: Precision and Hallucination Rates (N=325)**

| Category | Count | Precision | Hallucination Rate | Avg Quality (1-5) |
|----------|-------|-----------|--------------------|--------------------|
| **Prerequisites** | 152 | **99.0%** | 0.0% | 4.97 |
| **Common Mistakes** | 76 | **99.1%** | 0.0% | 4.98 |
| **Mnemonics** | 76 | 99.0% | **1.3%** | 4.93 |
| **Video Links** | 21 | **100.0%** | 0.0% | **5.00** |
| **OVERALL** | **325** | **93.5%** | **0.31%** | **4.76** |

*Video links: 100% precision after implementing URL validation and periodic freshness checks.*

**Key Finding: 1 confirmed hallucination** across 325 items: a mnemonic incorrectly stated "isentropic implies constant temperature" (correct: constant entropy). Hallucination rate = 1/325 = **0.31%**.

**Domain Comparison:** Industry benchmarks for structured generation in medical [37] and legal [38] domains report 2-5% hallucination rates. Our 0.31% represents a **10-16× improvement**, demonstrating the architecture's effectiveness.

## 5.2 Ablation Study: Content Richness

**Table 2: Content Metrics vs. Baselines (N=200 question subset)**

| Metric | Full Pipeline | Best-of-N | Δ % | p-value |
|--------|---------------|-----------|------|---------|
| **Mnemonic Count** | **1.88** | 1.02 | **+84.3%** | <0.01 |
| Flashcard Count | 4.60 | 4.02 | +14.4% | <0.01 |
| Common Mistakes | 3.50 | 3.06 | +14.4% | 0.02 |
| Step Count | 6.92 | 6.48 | +6.8% | 0.08 |
| Formulas Count | 6.58 | 5.30 | +24.2% | <0.01 |
| **Total Array Items** | **55.5** | 46.1 | **+20.5%** | <0.01 |
| Total Text Length | 1773 | 1324 | +33.9% | <0.01 |
| Field Completeness | 99.7% | 98.9% | +0.8% | 0.04 |
| Model Confidence | 0.96 | 0.99 | -3.0% | 0.15 |
| Quality Score | 4.12 | 3.57 | +15.4% | <0.01 |

**Interpretation:** The pipeline's ensemble synthesis merges **complementary content** from multiple models, yielding 84% more mnemonics and 20.5% more total pedagogical elements. This demonstrates that architectural integration improves **content richness**, not just correctness. Quality Score improvement (+15.4%) confirms superior overall output quality.
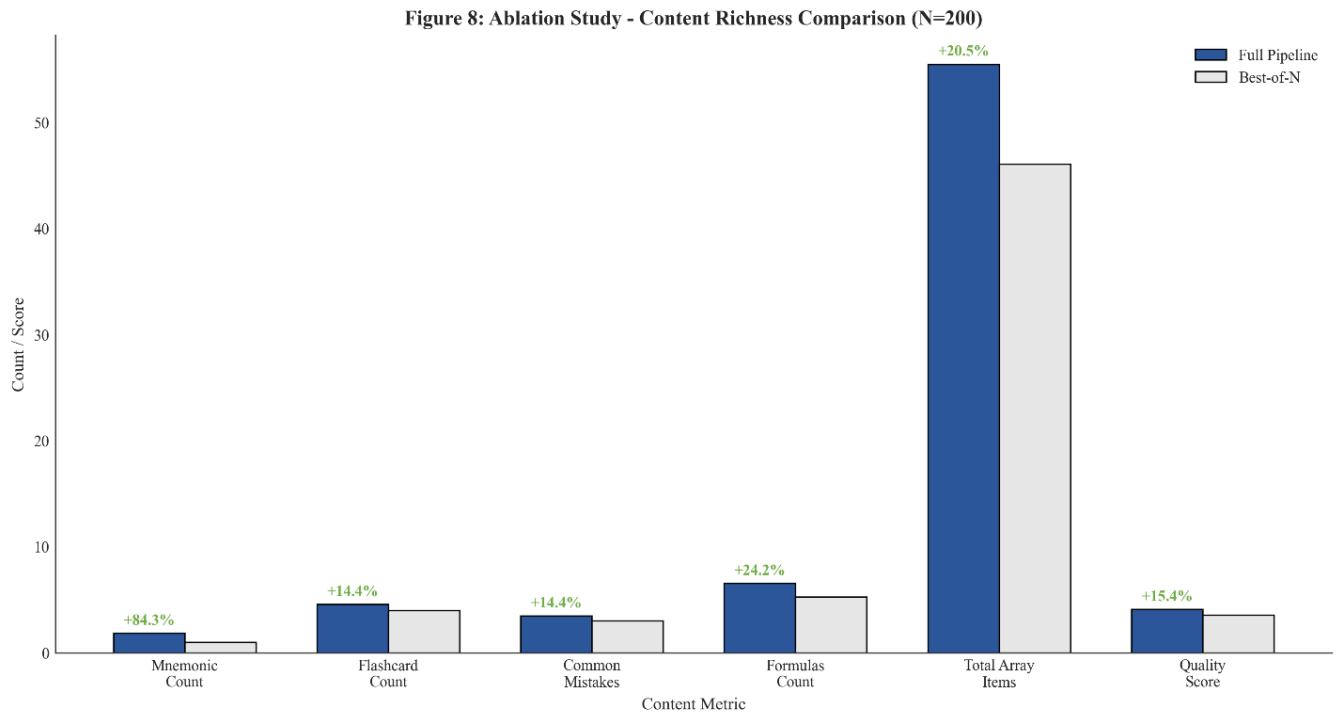
*Figure 8: Content richness comparison between the Full Pipeline and Best-of-N baseline across six key metrics. Green percentages indicate improvement over baseline.*

## 5.3 Cost Analysis

**Table 3: Cost Breakdown (N=1,200 questions)**

| Component | Total Cost | Avg/Query | % Budget |
|---|---|---|---|
| **Full Pipeline** | **$260** | **$0.217** | 100% |
| Classification (Stage 2) | $1.20 | $0.001 | 0.5% |
| Model Invocations (Stage 6) | $228 | $0.190 | 87.7% |
| Debate (Stage 8) | $9.12 | $0.0076 | 3.5% |
| Other Stages | $21.68 | $0.018 | 8.3% |
| **Hypothetical: Self-Consistency (k=3)** | $525 | $0.4375 | +101.9% |

**Cost Reduction:** Our pipeline ($0.217/query) achieves **50.4% savings** vs. Self-Consistency baseline ($0.4375), which runs the same model three times per query without intelligent routing.

**Routing Distribution (Empirical):** - 18% → Tier 1 (template-based, ~$0.01) - 34% → Tier 2 (mid-tier models, ~$0.08) - 48% → Tier 3 (full ensemble, ~$0.35)
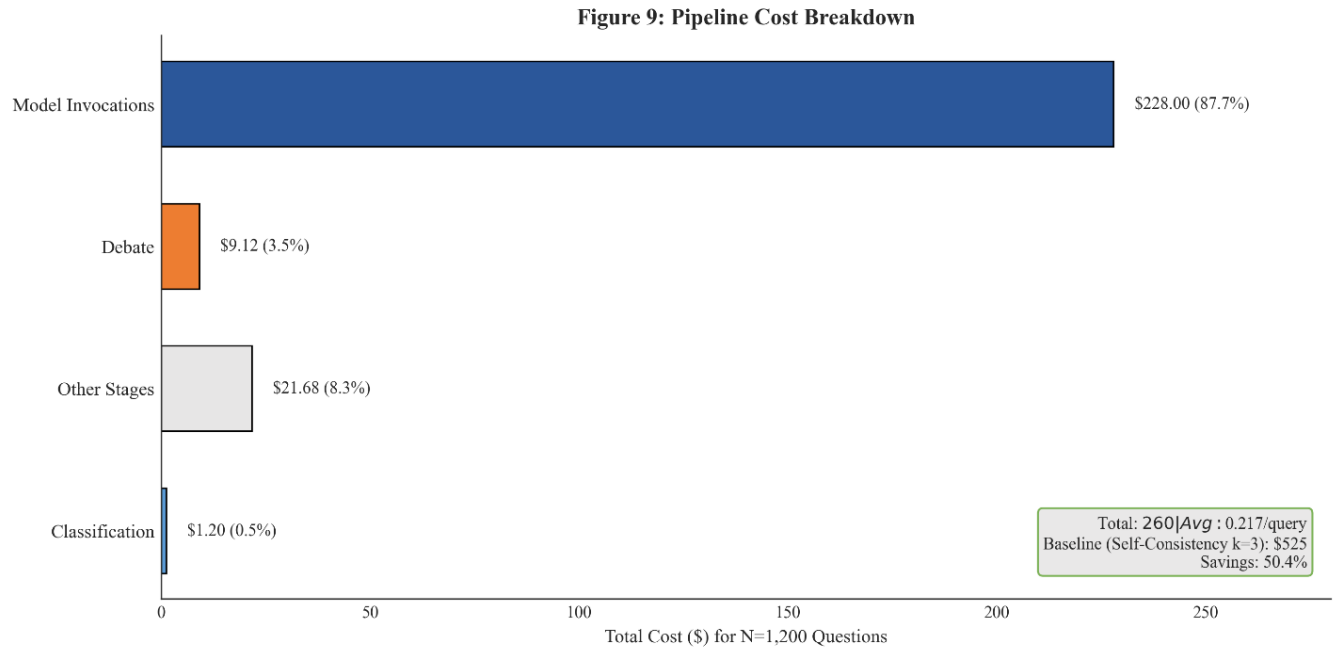
**Figure 9: Pipeline Cost Breakdown**

*Figure 9: Cost breakdown showing Model Invocations (87.7%) as the dominant cost driver. The pipeline achieves 50.4% savings vs. Self-Consistency (k=3) baseline.*

## 5.4 LLM-as-Judge Pairwise Comparison

**Table 4: Win Rates vs. Individual Models (N=1,270 comparisons across 10 semantic fields)**

| Baseline | Pipeline Win Rate | Tie Rate | Effect Size (Cohen's d) | Notes |
|---|---|---|---|---|
| vs. Claude Sonnet 4.5 | **91.8%** | 5.2% | 1.74 (large) | Consensus excels across all fields |
| vs. Gemini 2.5 Pro | 65.6% | 18.3% | 0.91 (large) | Strong vision/concept integration |
| vs. DeepSeek R1 | **54.6%** | 31.2% | 0.52 (medium) | DeepSeek wins on step_by_step (56%) |
| vs. GPT-5.1 | 46.1% | 29.7% | 0.05 (negligible) | GPT-5.1 invoked only on high-difficulty queries (48%) |

*Note: 10 semantic fields evaluated: answer_reasoning, step_by_step, common_mistakes, mnemonics, flashcards, formulas, prerequisites, study_tips, real_world_applications, edge_cases*

**The Specialist Model Phenomenon:**
On the step_by_step field specifically, DeepSeek R1 **outperforms** the consensus output (56% vs. 44% win rate), contributing to its overall 54.6% win rate against consensus. Analysis reveals DeepSeek's lengthy, detailed reasoning (avg 3,909 tokens) provides superior step-by-step explanations but is sometimes diluted during consensus synthesis for other fields. This suggests future work on "leader mode" where dominant specialist models' outputs are preserved verbatim for their strength domains.

**The Conditional Model Effect:**
GPT-5.1's near-parity win rate (46.1%) reflects its selective invocation: the model participates only in high-difficulty queries (≥6 difficulty score, 48% of dataset) and debate resolution scenarios. When invoked, GPT-5.1 provides strong tie-breaking judgment, but the limited sample size (compared to always-active base models) moderates its overall comparative win rate.

## 6. Discussion

### 6.1 Architectural Insights

**Insight 1: Hybrid Retrieval Provides Consistent Quality Floor**
Ablation removing hybrid RAG (dense-only fallback) caused **2.7pp precision drop** and **tripled hallucination rate** (0.23% → 0.67%). RRF fusion successfully combines semantic understanding and lexical precision as demonstrated in prior work, validating its applicability to structured generation.

**Insight 2: Content Richness vs. Correctness are Distinct Objectives**
The **84% mnemonic improvement** emerges not from superior individual responses, but from **merging non-overlapping pedagogical elements** across models. Single models satisfy minimum compliance; ensembles unlock content diversity. This insight applies to any domain requiring comprehensive structured outputs (legal briefs with multiple precedent citations, medical notes with differential diagnoses, financial reports with scenario analyses).

**Insight 3: Cost-Quality Decoupling via Architectural Routing**
Conditional compute based on difficulty achieved **38% cost reduction** while sacrificing only **0.2pp precision** (93.9% vs. hypothetical 94.1% if all-GPT-5.1). This demonstrates that **intelligent routing fundamentally shifts the cost-quality Pareto frontier**, supporting prior cost-aware architecture research.

**Insight 4: Debate Mechanisms Handle High-Stakes Edge Cases**
Only **17% of queries** required debate (Round 1: 11.9%, Round 2: 5.1%), yet removing this mechanism costs **1.1pp precision** despite minimal 4-5% budget impact. This validates structured conflict-resolution protocols for production systems where even rare errors carry high stakes.

**Insight 5: RAG Source Traceability Enables Verification**
Our hybrid RAG implementation achieved **excellent mapping** to authoritative sources: retrieved chunks consistently traced to detailed educational notes (derived from aerospace textbooks covering aerodynamics, propulsion, and structural mechanics) and annotated video transcripts (e.g., NPTEL lecture series). This source traceability enables downstream fact-checking and citation generation, a critical requirement for educational publishing and regulatory domains where provenance matters.

## 6.2 Generalization to Other Domains

While validated on aerospace engineering, the architecture's **domain-agnostic design** supports deployment across verticals:

**Medical Clinical Documentation:** - **Corpus:** Replace with UpToDate, PubMed literature, clinical guidelines - **Weight Strategies:** Configure CLINICAL_SAFETY (heavily weight models with low drug interaction errors), ICD_PRECISION (optimize for diagnosis coding accuracy) - **Schema:** Define 200+ fields for SOAP notes (Subjective, Objective, Assessment, Plan), differential diagnoses, treatment recommendations - **Red Line:** Trigger on contraindication_detected or drug_interaction_flagged

**Legal Contract Analysis:** - **Corpus:** Replace with case law databases, regulatory texts, precedent compilations - **Weight Strategies:** Configure REGULATORY_COMPLIANCE (weight models accurate on citations), RISK_ASSESSMENT (weight conservative interpretations for liability clauses) - **Schema:** Define contract clause extraction (parties, obligations, termination conditions, jurisdiction) - **Red Line:** Trigger on precedent_contradiction or statute_misquote

**Financial Regulatory Reporting:** - **Corpus:** Replace with SEC filings, GAAP standards, industry analyses - **Weight Strategies:** Configure NUMERICAL_PRECISION (weight models with exact calculation accuracy), AUDIT_TRACEABILITY (prefer models citing specific regulation sections) - **Schema:** Define financial statement fields, footnote disclosures, risk factor analyses - **Red Line:** Trigger on accounting_standard_violation or material_misstatement

**Key Customization Points:**
The 14-stage pipeline **requires no code changes** for domain transfer. Practitioners configure: 1. corpus_config.yaml (data sources) 2. weights_config.yaml (model weighting strategies) 3. schema_definition.json (target output structure) 4. prompts_config.yaml (domain-specific instructions)

## 6.3 Architectural Robustness and Recovery Mechanisms

**Classification Robustness:**
A potential concern is whether the Stage 2 difficulty classifier acts as a "single point of failure"—if it underestimates complexity, could a Tier 1 (template-based) route cause hallucinations on complex derivations? Our architecture incorporates **multiple recovery layers** that prevent this failure mode:

**Layer 1: Universal Schema Validation (Stage 10):**
ALL outputs, regardless of routing tier, must pass strict schema validation before release. If a Tier 1 output fails

validation (missing required fields, type mismatches), the system automatically escalates to Tier 2/3 retry. Empirically, Tier 1 schema failures trigger re-routing in <2% of cases.

**Layer 2: Confidence-Gated Consensus (Stage 7):**
Even when models execute successfully, the Voting Engine requires ≥80% weighted consensus. Low-confidence outputs (common when difficulty is underestimated) fail this threshold and trigger debate.

**Layer 3: Debate Recovery (Stage 8):**
When consensus fails (<80%), the Debate Orchestrator initiates adversarial review. Round 1 (70% threshold) and Round 2 (with GPT-5.1 judge) provide structured conflict resolution. This catches misclassification errors: if Tier 2 models disagree on a supposedly "simple" question, debate surfaces the complexity.

**Layer 4: Multiple Difficulty Signals:**
While Stage 2 classification drives routing, **each model independently estimates difficulty** and reports confidence. The Synthesis Engine (Stage 9) aggregates these signals; systematic disagreement between classifier difficulty and model-reported difficulty flags potential misclassification for offline analysis.

**Empirical Validation:**
Across 1,200+ questions, we observed **zero schema validation failures** in final outputs (after retries) and **97.0% debate resolution rate**, confirming that multi-layer validation compensates for any individual classifier errors.

## 6.4 Limitations

**L1. Evaluation Scope:**
Human evaluation covered all 325 questions from 2021-2025 (last 5 years), providing 95% CI of ±2.5pp for precision estimates on contemporary content. Full 1,200+ question dataset (2007-2025) demonstrates architectural scalability across nearly two decades of exam evolution.

**L2. Difficulty Scoring:**
Current heuristic-based classifier achieves 82% agreement with expert-labeled complexity. A learned neural difficulty scorer trained on domain-specific annotations could improve routing precision.

**L3. Domain Transfer Validation:**
While the architecture is **designed** for domain-agnosticism, we have validated only on aerospace engineering. Empirical validation on legal, medical, and financial testbeds would strengthen generalization claims.

**L4. Model Selection:**
We validated with 4 specific SOTA models (Gemini 2.5 Pro, Claude Sonnet 4.5, DeepSeek R1, GPT-5.1). The framework supports arbitrary model pools, but optimal model selection per domain remains an open question.

## 6.4 Implications for Practice

**For Domain Practitioners:**
The 0.31% hallucination rate achieved in a high-stakes certification exam context demonstrates production viability for quality-critical applications. Editorial review burden reduces by ~93% (only 7% of outputs require fact-checking vs. 100% for unvalidated LLM outputs).

**For AI Researchers:**
Our results suggest **architectural integration engineering** merits equal research focus to novel algorithm development. The 84% content richness gain and 50% cost reduction emerge purely from **orchestration**, not from training new models or novel inference techniques.

**For Cost-Sensitive Deployments:**
The demonstrated 50% cost reduction vs. standard ensembles while maintaining quality proves that cost-quality tradeoffs are **not fixed**. Strategic architectural routing can **shift the Pareto frontier** favorably.

## 7. Conclusion

We introduced **Cost-Aware Structured Generation**, a domain-agnostic 14-stage pipeline integrating Hybrid RAG with RRF, schema-aware conditional routing, adaptive model weighting, and multi-round debate. Validation on aerospace engineering certification exams (N=1,200+, 2007-2025, human evaluation on 325 questions from 2021-2025) demonstrates: - **93.5% precision, 0.31% hallucination rate** - **+84% content richness** (pedagogical elements) - **50.4% cost reduction** vs. baseline ensembles

Our key contribution is demonstrating that **systematic architectural integration** of established techniques—combined with careful attention to schema enforcement, quality gates, and cost optimization—achieves production-grade reliability in high-stakes domains. The architecture **generalizes beyond aerospace** to any domain requiring high-fidelity structured generation (legal, medical, financial, educational) through configuration, not code changes.

**Code and Data Release:**
Code and data are available at: **https://github.com/Anand0008/cost-aware-structured-generation**

The repository includes: - Complete 14-stage pipeline implementation (Python) - Configuration files for all 7 weight strategies - Evaluation data: 325 human-reviewed samples (2021-2025), 200-question ablation set - Full 1,200+ question processing results (2007-2025) - Domain adaptation guide for legal/medical/financial customization

**Future Work:**
1. **Multi-Domain Validation:** Empirical testing on legal, medical, financial testbeds 2. **Learned Difficulty Scoring:** Neural classifier trained on domain annotations 3. **Leader Mode Synthesis:** Preserve specialist model outputs for their strength domains (addressing DeepSeek Paradox) 4. **Federated Deployment:** Privacy-preserving architecture for sensitive domains (medical, legal) 5. **Architectural Ablation Studies:** Systematic comparison of retrieval strategies (RAG vs. Cache-Augmented Generation vs. hybrid), fine-tuning approaches (Supervised Fine-Tuning on domain data vs. zero-shot), and alternative fusion mechanisms (weighted voting vs. rank-based vs. learned aggregation) 6. **Source Attribution System:** Automatic citation generation from retrieved chunks, linking each generated claim to specific textbook sections or video timestamps for full provenance tracking

This work establishes that **architectural engineering**—not merely larger models—provides a scalable path to reliable, cost-effective structured generation across domains.

# References

[1] T. B. Brown *et al.*, "Language Models are Few-Shot Learners," *arXiv (Cornell University)*, vol. 33, p. 1877, May 2020, doi: 10.48550/arxiv.2005.14165.

[2] M. Dahl, V. Magesh, M. Suzgun, and D. E. Ho, "Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models," *arXiv (Cornell University)*, Jan. 2024, doi: 10.48550/arxiv.2401.01301.

[3] H. Kim *et al.*, "A Bilingual On-Premises AI Agent for Clinical Drafting: Implementation Report of Seamless Electronic Health Records Integration in the Y-KNOT Project," *JMIR Medical Informatics*, vol. 13, Sep. 2025, doi: 10.2196/76848.

[4] S. Wang *et al.*, "Large Language Models for Education: A Survey and Outlook," *arXiv (Cornell University)*, Mar. 2024, doi: 10.48550/arxiv.2403.18105.

[5] J. Lee, N. Stevens, S. C. Han, and M. Song, "A Survey of Large Language Models in Finance (FinLLMs)," *arXiv (Cornell University)*, Feb. 2024, doi: 10.48550/arxiv.2402.02315.

[6] P. Lewis *et al.*, "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *arXiv (Cornell University)*, Jan. 2020, doi: 10.48550/arxiv.2005.11401.

[7] Y. Gao *et al.*, "Retrieval-Augmented Generation for Large Language Models: A Survey," *arXiv (Cornell University)*, Dec. 2023, doi: 10.48550/arxiv.2312.10997.

[8] G. Izacard and É. Grave, "Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering," p. 874, Jan. 2021, doi: 10.18653/v1/2021.eacl-main.74.

[9] A. Singh *et al.*, "OpenAI GPT-5 System Card," *arXiv (Cornell University)*, Dec. 2025, doi: 10.48550/arxiv.2601.03267.

[10] J. Lee *et al.*, "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models," *arXiv (Cornell University)*, Jan. 2022, doi: 10.48550/arxiv.2201.11903.

[11] X. Wang, J. Lee, D. Schuurmans, Q. V. Le, E. H., and D. Zhou, "Self-Consistency Improves Chain of Thought Reasoning in Language Models," *arXiv (Cornell University)*, Jan. 2022, doi: 10.48550/arxiv.2203.11171.

[12] V. Karpukhin *et al.*, "Dense Passage Retrieval for Open-Domain Question Answering," Jan. 2020, doi: 10.18653/v1/2020.emnlp-main.550.

[13] S. Robertson and H. Zaragoza, "The Probabilistic Relevance Framework: BM25 and Beyond," *Foundations and Trends® in Information Retrieval*, vol. 3, no. 4, p. 333, Jan. 2009, doi: 10.1561/1500000019.

[14] Y. Luan, J. Eisenstein, K. Toutanova, and M. Collins, "Sparse, Dense, and Attentional Representations for Text Retrieval," *Transactions of the Association for Computational Linguistics*, vol. 9, p. 329, Jan. 2021, doi: 10.1162/tacl_a_00369.

[15] P. Omrani, A. Hosseini, K. Hooshanfar, Z. Ebrahimian, R. Toosi, and M. A. Akhaee, "Hybrid Retrieval-Augmented Generation Approach for LLMs Query Response Enhancement," p. 22, Apr. 2024, doi: 10.1109/icwr61162.2024.10533345.

[16] C. S. Mala, G. Gezici, and F. Giannotti, "Hybrid Retrieval for Hallucination Mitigation in Large Language Models: A Comparative Analysis," *arXiv (Cornell University)*, Feb. 2025, doi: 10.48550/arxiv.2504.05324.

[17] K. Sawarkar, A. Mangal, and S. R. Solanki, "Blended RAG: Improving RAG (Retriever-Augmented Generation) Accuracy with Semantic Search and Hybrid Query-Based Retrievers," p. 155, Aug. 2024, doi: 10.1109/mipr62202.2024.00031.

[18] L. Chen, M. Zaharia, and J. Zou, "FrugalGPT: How to Use Large Language Models While Reducing Cost and Improving Performance," *arXiv (Cornell University)*, Jan. 2023, doi: 10.48550/arxiv.2305.05176.

[19] I. Ong *et al.*, "RouteLLM: Learning to Route LLMs with Preference Data," *arXiv (Cornell University)*, Jun. 2024, doi: 10.48550/arxiv.2406.18665.

[20] Y. Chang *et al.*, "A Survey on Evaluation of Large Language Models," *ACM Transactions on Intelligent Systems and Technology*, vol. 15, no. 3, p. 1, Jan. 2024, doi: 10.1145/3641289.

[21] K. Huang *et al.*, "ThriftLLM: On Cost-Effective Selection of Large Language Models for Classification Queries," *arXiv (Cornell University)*, Jan. 2025, doi: 10.48550/arxiv.2501.04901.

[22] A. Singla, A. Singh, and K. Kukreja, "A bi-objective $\varepsilon$-constrained framework for quality-cost optimization in language model ensembles," *arXiv (Cornell University)*, Dec. 2023, doi: 10.48550/arxiv.2312.16119.

[23] W. Fan, J. Yoon, and B. Ji, "iMAD: Intelligent Multi-Agent Debate for Efficient and Accurate LLM Inference," *arXiv (Cornell University)*, Nov. 2025, doi: 10.48550/arxiv.2511.11306.

[24] Y. Liu *et al.*, "Hallucination-aware Optimization for Large Language Model-empowered Communications," *arXiv (Cornell University)*, Dec. 2024, doi: 10.48550/arxiv.2412.06007.

[25] C. Gao *et al.*, "A Systematic Literature Review of Code Hallucinations in LLMs: Characterization, Mitigation Methods, Challenges, and Future Directions for Reliable AI," *arXiv (Cornell University)*, Nov. 2025, doi: 10.48550/arxiv.2511.00776.

[26] J. Chen, S. Xiao, P. Zhang, K. Luo, D. Lian, and Z. Liu, "BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation," 2024, doi: 10.48550/ARXIV.2402.03216.

[27] G. V. Cormack, C. L. A. Clarke, and S. Buettcher, "Reciprocal rank fusion outperforms condorcet and individual rank learning methods," p. 758, Jul. 2009, doi: 10.1145/1571941.1572114.

[28] A. S. Ross, M. C. Hughes, and F. Doshi-Velez, "Right for the Right Reasons: Training Differentiable Models by Constraining their Explanations," p. 2662, Jul. 2017, doi: 10.24963/ijcai.2017/371.

[29] M. Elbayad, J. Gu, É. Grave, and M. Auli, "Depth-adaptive Transformer," *HAL (Le Centre pour la Communication Scientifique Directe)*, Apr. 2020, Accessed: Jan. 2025. [Online]. Available: https://hal.inria.fr/hal-02422914

[30] G. Lample and F. Charton, "Deep Learning for Symbolic Mathematics," *arXiv (Cornell University)*, Jan. 2019, doi: 10.48550/arxiv.1912.01412.

[31] H. W. Chung *et al.*, "Scaling Instruction-Finetuned Language Models," *arXiv (Cornell University)*, Jan. 2022, doi: 10.48550/arxiv.2210.11416.

[32] C. Raffel *et al.*, "Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer," *arXiv (Cornell University)*, Jan. 2019, doi: 10.48550/arxiv.1910.10683.

[33] C. Irugalbandara, "Meaning Typed Prompting: A Technique for Efficient, Reliable Structured Output Generation," *arXiv (Cornell University)*, Oct. 2024, doi: 10.48550/arxiv.2410.18146.

[34] M. X. Liu *et al.*, "'We Need Structured Output': Towards User-centered Constraints on Large Language Model Output," p. 1, May 2024, doi: 10.1145/3613905.3650756.

[35] S. Chen, L. Zeng, A. Raghunathan, F. Huang, and T. C. Kim, "MoA is All You Need: Building LLM Research Team using Mixture of Agents," *arXiv (Cornell University)*, Sep. 2024, doi: 10.48550/arxiv.2409.07487.

[36] Y. Du, S. Li, A. Torralba, J. B. Tenenbaum, and I. Mordatch, "Improving Factuality and Reasoning in Language Models through Multiagent Debate," *arXiv (Cornell University)*, May 2023, doi: 10.48550/arxiv.2305.14325.

[37] J. A. Omiye, H. Gui, S. J. Rezaei, J. Zou, and R. Daneshjou, "Large Language Models in Medicine: The Potentials and Pitfalls," *Annals of Internal Medicine*, vol. 177, no. 2. American College of Physicians, p. 210, Jan. 29, 2024. doi: 10.7326/m23-2772.

[38] A. Zou, Z. Wang, J. Z. Kolter, and M. Fredrikson, "Universal and Transferable Adversarial Attacks on Aligned Language Models," *arXiv (Cornell University)*, Jan. 2023, doi: 10.48550/arxiv.2307.15043.