

SAGEO Arena: A Realistic Environment for Evaluating Search-Augmented Generative Engine Optimization

Sunghwan Kim

Department of Artificial Intelligence
Yonsei University
Seoul, Republic of Korea
happysnail06@yonsei.ac.kr

Wooseok Jeong

Department of Computer Science and
Engineering
Konkuk University
Seoul, Republic of Korea
oseoko8@gmail.com

Serin Kim

Department of Artificial Intelligence
Yonsei University
Seoul, Republic of Korea
kimserin@yonsei.ac.kr

Sangam Lee

Department of Artificial Intelligence
Yonsei University
Seoul, Republic of Korea
salee@yonsei.ac.kr

Dongha Lee*

Department of Artificial Intelligence
Yonsei University
Seoul, Republic of Korea
donalee@yonsei.ac.kr

Abstract

Search-Augmented Generative Engines (SAGE) have emerged as a new paradigm for information access, bridging web-scale retrieval with generative capabilities to deliver synthesized answers. This shift has fundamentally reshaped how web content gains exposure online, giving rise to Search-Augmented Generative Engine Optimization (SAGEO), the practice of optimizing web documents to improve their visibility in AI-generated responses. Despite growing interest, no evaluation environment currently supports comprehensive investigation of SAGEO. Specifically, existing benchmarks lack end-to-end visibility evaluation of optimization strategies, operating on pre-determined candidate documents that abstract away retrieval and reranking preceding generation. Moreover, existing benchmarks discard structural information (e.g., schema markup) present in real web documents, overlooking the rich signals that search systems actively leverage in practice. Motivated by these gaps, we introduce SAGEO ARENA, a realistic and reproducible environment for stage-level SAGEO analysis. Our objective is to jointly target search-oriented optimization (SEO) and generation-centric optimization (GEO). To achieve this, we integrate a full generative search pipeline over a large-scale corpus of web documents with rich structural information. Our findings reveal that existing approaches remain largely impractical under realistic conditions and often degrade performance in retrieval and reranking. We also find that structural information helps mitigate these limitations, and that effective SAGEO requires tailoring optimization to each pipeline stage. Overall, our benchmark paves the way for realistic SAGEO evaluation and optimization beyond simplified settings.

*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).
Conference acronym 'XX, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

Keywords

Search-Augmented Generative Engine, Generative Engine Optimization, Search Engine Optimization, Benchmark, Evaluation

ACM Reference Format:

Sunghwan Kim, Wooseok Jeong, Serin Kim, Sangam Lee, and Dongha Lee. 2018. SAGEO Arena: A Realistic Environment for Evaluating Search-Augmented Generative Engine Optimization. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

With the rapid advancement of Large Language Models (LLMs), Search-Augmented Generative Engines (SAGE) [7, 29, 48] have emerged as a new paradigm for information access. By integrating large-scale retrieval with generative capabilities, these systems perform *generative search*, providing synthesized answers that directly address user queries. As an increasing portion of web traffic now originates from such AI-generated answers [9], how content gains exposure online is fundamentally changing. This shift has given rise to Search-Augmented Generative Engine Optimization (SAGEO), the practice of optimizing web documents to improve their visibility (i.e., presence and contribution) in AI-generated responses.

SAGEO naturally extends the core objective of traditional Search Engine Optimization (SEO) [12, 20]. For decades, optimization has focused on improving retrieval and ranking on Search Engine Results Pages (SERPs), primarily through on-page factors such as titles and metadata [42]. With the emergence of SAGE, optimization has increasingly shifted toward a generation-centric perspective, studying how a document is favored and cited in model responses [8]. A prominent line of study is Generative Engine Optimization (GEO) [1], focusing on presentational modifications to improve visibility at the generation stage. Yet, retrieval and ranking remain essential in generative search, as documents that fail at either stage cannot be cited in AI responses. Thus, SEO and GEO must be jointly addressed to achieve successful optimization. In light of this, we clearly distinguish SAGEO as optimization that targets the full generative search pipeline, from retrieval and reranking to generation (Figure 1). This motivates a systematic study of

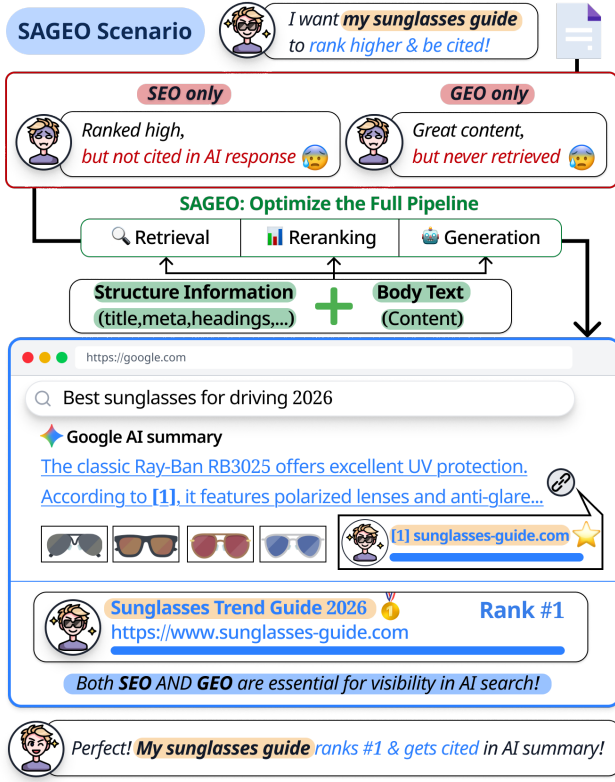


Figure 1: Optimizing for ranking alone (SEO) or citation alone (GEO) causes failure at either retrieval/reranking or generation. SAGEO jointly optimizes structural information and body texts across the full pipeline to succeed at both stages.

what constitutes effective SAGEO and how optimization signals propagate and interact across stages in generative search.

Despite growing interest, there exists no evaluation environment that enables such investigation. Existing benchmarks [1, 10, 37, 49] tend to oversimplify the complexities of real web environments, leaving it unclear whether current optimization strategies are effective in practical settings. Specifically, there are two major limitations: **(1) Lack of end-to-end evaluation.** Current benchmarks predominantly operate on pre-determined candidate documents, abstracting away the retrieval and reranking stages that precede generation. It remains unclear whether current optimization strategies trade off or benefit performance at earlier stages, making stage-level analysis of SAGEO essential. **(2) Loss of structural information.** Existing benchmarks operate solely on plain body text, discarding structural information present in real web documents (e.g., schema markup). Notably, recent findings suggest that augmenting structured representations such as markdown summaries yields substantial visibility gains [37]. Yet no existing benchmark preserves structural information to enable analysis of their impact.

Motivated by these gaps, we introduce SAGEO ARENA, a realistic and reproducible benchmark designed for stage-level SAGEO analysis. Crucially, our benchmark is distinguished in two ways: **(1) Comprehensive generative search environment.** Unlike

existing benchmarks that assume fixed retrieval context, SAGEO ARENA integrates a full generative search pipeline where documents are dynamically retrieved, ranked, and passed to generation. Each stage is modular and configurable, allowing researchers to trace how optimization signals propagate across the complete search process. **(2) Extensive corpus with rich structural annotations.** We curate a corpus of 170k web documents spanning 9 domains. Following guidelines published by commercial search engines, we extract both body text and rich structural information that SAGEO must navigate in practice (e.g., meta descriptions, headings, and schema markup). Overall, we aim to construct an environment that closely reflects real-world deployment, paving the way for optimization strategies that generalize beyond controlled settings.

We conduct extensive experiments on SAGEO ARENA to evaluate how optimization strategies affect document visibility across the full generative search pipeline. Our analysis reveals that body text optimization alone, the predominant focus of prior work, remains largely insufficient in realistic generative search settings. It not only fails to improve generation-stage visibility but actively degrades retrieval performance, causing optimized documents to drop out of the retrieval results and never reach the generator. Structural information proves essential for mitigating this degradation, serving as the primary mechanism through which documents are retrieved in early pipeline stages. Yet further analysis reveals that the two scopes play fundamentally complementary roles, where structural information drives retrieval and informative body text remains a key factor for reranking and generation. Through in-depth analysis across domains, citation behavior, and optimization models, we further show that each pipeline stage prioritizes distinct document qualities, motivating the need for stage-aware optimization. Guided by these findings, we propose stage-aware SAGEO, a method that tailors optimization to the priorities of each generative search stage, achieving the strongest overall visibility among all evaluated strategies. To summarize, our contributions are as follows:

- We introduce SAGEO ARENA, the first benchmark that enables stage-level visibility evaluation of SAGEO. By integrating a full generative search pipeline with a large-scale corpus preserving rich structural information, SAGEO ARENA bridges the gap between existing benchmarks and real-world SAGE.
- We provide the first empirical analysis of structural information optimization, studying how structural signals that generative search systems must navigate in practice influence visibility.
- Through extensive experiments, we show that current optimization approaches are insufficient in realistic settings. To address this, we introduce stage-aware SAGEO, providing actionable guidance to further facilitate SAGEO research.

2 Related Work

2.1 Generative Search Engine Optimization

With the rise of Large Language Models, Search-Augmented Generative Engines [22, 30, 54] have become a cornerstone of modern information access, offering synthesized answers grounded in retrieved documents. This shift has introduced new visibility challenges for content creators, motivating research on optimization strategies tailored to generative engines. Existing work can be

Table 1: A comparison of SAGEO ARENA to existing benchmarks.

Benchmark	Environment				Document Details		Evaluation		
	Doc. Corpus	Retrieval	Reranking	Generation	Structure Info.	Body Text	Search	Generation	Visibility Metric
GEO-Bench [1]	-	✗	✗	✓	✗	✓	✗	✓	Word Count
AutoGEO [49]	-	✗	✗	✓	✗	✓	✗	✓	Word Count, Utility
C-SEO Bench [37]	-	✗	✗	✓	✗	✓	✗	✓	Citation Rank
CC-GSEO-Bench [10]	-	✗	✗	✓	✗	✓	✗	✓	Influence
SAGEO ARENA (Ours)	170K	✓	✓	✓	✓	✓	✓	✓	Hit Rate, Rank Change

broadly categorized into two directions: (1) *Black-hat* approaches explore adversarial methods that exploit model behavior through malicious content injection [24, 35, 36]. In contrast, (2) *white-hat* approaches focus on cooperative content modification, formalized as Generative Engine Optimization (GEO) [1]. Building on this foundation, subsequent work [10, 37, 49] has advanced GEO research by proposing diverse optimization strategies and evaluation dimensions (Table 1). Nevertheless, existing evaluation protocols often rely on oversimplified environments, predominantly assessing optimization effects at the generation stage with fixed document candidates. Thus, how optimization signals propagate through the full generative pipeline remains largely unexplored. Motivated by these, this study provides a comprehensive evaluation environment that explicitly models the end-to-end generative search pipeline.

2.2 Comparison of SEO and SAGEO

Traditional Search Engine Optimization (SEO) has long served as the dominant framework for online visibility [19, 26, 38, 42], aiming to improve a source’s position on Search Engine Results Pages (SERP). Following established guidelines from leading search engine providers [4, 16], SEO encompasses both *on-page* factors such as keyword placement, metadata, and schema markup, as well as *off-page* factors including backlinks and domain authority. In contrast, GEO research predominantly focuses on modifying body text in the generation stage [11], investigating rewriting strategies that adjust textual properties such as fluency, technicality, and ease of understanding. However, recent empirical findings [37] challenge this generation-centric perspective, demonstrating that a document’s position determined in the retrieving and reranking stage plays a far more dominant role in determining its visibility in the final response. Notably, document position is largely shaped by structural information emphasized in traditional SEO, which raises a fundamental question: *do these signals also contribute to visibility in generative engines?* Existing benchmarks offer limited insight into this question, and we bridge this gap by preserving the structural information and examining their influence in generative search.

2.3 Structured Web Information Understanding

The World Wide Web is fundamentally built upon HyperText Markup Language (HTML). Beyond body content, HTML documents encode structural information through elements such as title tags, meta descriptions, heading hierarchies, and schema markup. These structural attributes have long influenced how search engines crawl, index, and rank documents [28]. Recent studies [13, 21, 45, 47] demonstrate that structured inputs help LLMs better understand web content, showing that structural representations in HTML

complement plain text in Retrieval-Augmented Generation (RAG) systems. Moreover, there is evidence indicating that commercial search-augmented generative engines, such as Bing Copilot and Google Search, similarly leverage structural context when interpreting web documents and tend to cite pages with well-organized structure data [2, 6, 25, 32]. These findings motivate incorporating structural information into GEO research, but no existing GEO benchmark systematically evaluates its impact. We address this gap by proposing a novel benchmark that preserves the structural elements in web documents, enabling more realistic evaluation of optimization strategies in real-world generative search settings.

3 SAGEO ARENA

Problem Setting. We consider a search-augmented generative engine that answers user queries by retrieving relevant web documents and synthesizing a response grounded in the retrieved content. The pipeline follows the well-established Retrieval-Augmented Generation (RAG) paradigm [29], consisting of a retriever, reranker, and generator. In this setting, a document must first be retrieved and ranked before it can influence the generated response. This creates an inherent dependency between traditional SEO concerns (retrieval and ranking) and emerging GEO concerns (generation). Our benchmark is designed to capture this full pipeline, enabling thorough evaluation of optimization strategies across all three stages.

Design Principles. Our core challenge lies in reproducing an environment that reflects real-world search-augmented generative engines. To achieve this, we (1) adopt a standardized search-augmented generation pipeline following established practices [15, 18, 39] and (2) preserve structural information explicitly recommended by commercial search engines for optimization [17, 33, 34].

Task Formulation. Let q denote a user query and \mathcal{D} a document corpus. Given q , the system produces a response A_q through three stages: (1) a retriever \mathcal{R} returns a candidate set $\mathcal{D}_q = \mathcal{R}_k(q, \mathcal{D})$ containing the top- k documents; (2) a reranker \mathcal{F} reorders \mathcal{D}_q by relevance, producing a ranked list $\mathcal{D}_q^* = \mathcal{F}(q, \mathcal{D}_q)$; (3) a generator \mathcal{G} produces a response $A_q = \mathcal{G}(q, \mathcal{D}_q^*)$ with inline citations.

$$A_q = \mathcal{G}\left(q, \mathcal{F}\left(q, \mathcal{R}_k(q, \mathcal{D})\right)\right). \quad (1)$$

The visibility of each document $d_i \in \mathcal{D}_q^*$ is reflected by whether and to what extent it is cited in A_q . Let $d_i^{\text{tgt}} \in \mathcal{D}$ denote a target document relevant to the query q . SAGEO aims to optimize d_i^{tgt} , without prior knowledge of the query q , such that it (1) is retrieved into \mathcal{D}_q , (2) ranks highly in \mathcal{D}_q^* , and (3) maximizes visibility in A_q .

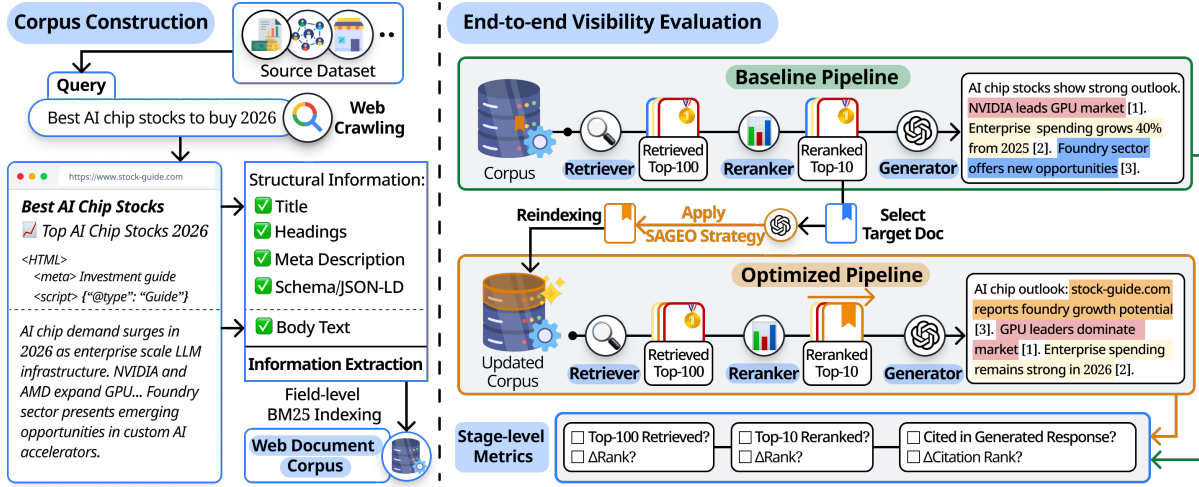


Figure 2: Overview of SAGEO ARENA. During document corpus construction, we extract structural information (title, meta description, headings, schema) and body text from web documents. During SAGEO evaluation, we first execute the pipeline with each test query and randomly select a target document from those that reach the generation stage (top- k after reranking) as the baseline. The target document is then optimized, re-indexed into the corpus, and the pipeline is re-executed with the same query, tracking whether the document maintains, gains, or loses visibility at each stage (retrieval, reranking, generation).

3.1 Corpus Construction

To create a large-scale corpus with strong domain diversity, we curate queries from nine established information retrieval datasets spanning a broad range of domains (Table 4). We sample 300 queries from each dataset, yielding 2,700 unique queries in total. Following prior approaches, we retrieve up to 100 search results per query using the Google Custom Search API. We then crawl each URL to extract both body text and rich structural information (Section 3.2). After filtering documents with malformed content, the final corpus contains 171,003 unique web documents. On average, each query is associated with 63 candidate documents. The 300 sampled queries per domain serve as test queries for evaluating optimization. Detailed statistics of SAGEO ARENA are provided in Appendix A

3.2 Structural Information Extraction

Unlike prior benchmarks that extract only body text, we explicitly preserve structural information embedded in real-world web documents. However, such documents exhibit substantial diversity and noise [46, 52], making it impractical to preserve all structural information for systematic evaluation. We therefore selectively retain body fields consistently emphasized in public search engine guidelines [34], enabling a realistic yet controlled study of how these elements influences search-augmented generative engines.

- **Title.** The document title provides a concise topical summary of the page. It serves as a primary signal of the page’s main topic and intent, playing a key role in initial relevance assessment.
- **Meta Description.** A concise page summary intended for search result snippets. This field offers a compact yet informative description of page content that complement the title.
- **Headings.** Hierarchical section headers (H1–H6) that signal topical structure within a document. Search engines treat headings as strong indicators of what each section is about, often weighting them higher than surrounding body text.

- **Schema/JSON-LD.** Structured data markup that explicitly defines entities, attributes, and relationships in a machine-readable format. Search engines rely on this markup to interpret page semantics that cannot be inferred from plain text alone.
- **Body Text.** The main textual content of the document, representing the primary source of information.

Notably, merging these fields into a unified text representation at the indexing stage would obscure the distinct signals carried by each component. Therefore, search-augmented generative engines broadly index these fields as separate components and combine them later during relevance scoring [41]. We adopt the same practice, preserving field-level separation to enable fine-grained analysis of how each component contributes to visibility across the pipeline.

3.3 Search-augmented Generative Engine

We implement a full generative search engine pipeline with retrieval, reranking, and generation. Each stage reflects real-world practices while remaining modular for controlled experimentation.

Document Representation. In our setting, a web document consists of two types of content: structural information and body text. The structural information of a document, denoted as $\mathcal{S}(d) = \{s_1, \dots, s_m\}$, is a set of structural fields corresponding to the elements described in Section 3.2 (e.g., title, meta description). The body text of a document, denoted as $\mathcal{P}(d) = \{p_1, \dots, p_n\}$, is a set of passages chunked from the body text following standard practice. A document can therefore be represented as $(\mathcal{S}(d), \mathcal{P}(d))$.

Semantic Unit. In generative search systems, documents are often chunked and retrieved at the passage level. However, structural information of a web document is shared across all passages, making it essential to bridge this gap between document-level and passage-level information for effective retrieval. We therefore pair each passage $p_i \in \mathcal{P}(d)$ with the associated structural information $\mathcal{S}(d)$,

creating a semantic unit $\mathcal{U}(p_i) = \mathcal{S}(d) \cup \{p_i\}$ that serves as a passage representation of p_i during retrieval and reranking.

Retriever. We employ BM25 [40], a widely adopted lexical retrieval method, as our retriever \mathcal{R} . For each semantic unit $\mathcal{U}(p)$, we build a separate BM25 index for every element $u \in \mathcal{U}(p)$. Given a query q , we score q against each u independently, and aggregate the resulting ranks via reciprocal rank fusion to obtain the retrieval score of passage p :

$$\text{score}_{\mathcal{R}}(q, p) = \sum_{u \in \mathcal{U}(p)} \frac{1}{\kappa + \text{rank}(q, u)}, \quad (2)$$

where $\text{rank}(q, u)$ is the rank of element u , and κ is a constant. The top- k passages are selected as candidates for reranking.

Reranker. The retrieved candidates are reranked using a cross-encoder that scores each query-passage pair (q, p) independently. Following practices in production systems (e.g., Vespa AI¹, Google Vertex AI²), we score relevance at the passage level rather than over entire documents. This approach handles cases where only specific sections of a long document are relevant to the query.

Generator. The top-reranked candidates are provided to the generation model along with their structural fields. The model is prompted to cite sources, enabling measurement of document-level visibility. By grounding our design on a complete pipeline with structurally enriched documents, SAGEO ARENA enables analyses beyond the scope of existing benchmarks. Specifically, it allows (1) stage-wise evaluation of diverse strategies, (2) component-level ablations, and (3) testing robustness and generalization across realistic settings.

3.4 Stage-level Visibility Evaluation

We describe how SAGEO ARENA measures document visibility across each pipeline stage and quantifies the effect of optimization.

Evaluation Pipeline. As shown in Figure 2, we establish a baseline as follows. For each test query q , we first execute the full generative search pipeline in SAGEO ARENA. Among the documents that reach the generation stage (i.e., ranked within the top- k at the reranking stage), we randomly select one as the target document d^{tgt} , yielding the baseline instance (q, d^{tgt}) . The target document d^{tgt} is then optimized using a specified strategy and reindexed into the corpus following the process described in Section 3.3. We then re-execute the search pipeline with the same query and compare the document’s visibility against the baseline at each stage. This enables evaluating optimization effects under realistic, stage-level settings, rather than within a fixed retrieval context like prior benchmarks.

Tracking Target Document. Since search-augmented generative engines internally operate at the passage level, multiple passages from the same document may appear in the candidate list as they are scored independently. This makes tracking the visibility of a target document an inherent challenge. To address this, we define the rank of a target document d_q^{tgt} as the highest rank achieved by any of its associated passages at each pipeline stage:

$$\text{rank}(d_q^{\text{tgt}}) = \min_{p \in \mathcal{P}(d_q^{\text{tgt}})} \text{rank}(p), \quad (3)$$

where $\mathcal{P}(d_q^{\text{tgt}})$ denotes the set of passages derived from d_q^{tgt} . The passage with the highest rank serves as the visibility indicator of the associated document at each stage, and we measure optimization effectiveness in two ways: (1) whether the target document appears within the top- k candidates at each stage, and (2) how the rank of the target document shifts before and after optimization.

Metrics. For each test query $q \in Q$, we execute the pipeline and track the rank of its associated target document d_q^{tgt} via Equation (3). We then evaluate visibility at each stage using two complementary metrics: (1) HIT RATE (H@ k), measuring the proportion of queries for which the target document appears within the top- k candidates. For retrieval and reranking, it is defined by

$$\text{H@}k = \frac{1}{|Q|} \sum_{q \in Q} \mathbb{I}(\text{rank}(d_q^{\text{tgt}}) \leq k), \quad (4)$$

where $\mathbb{I}(\cdot)$ returns 1 if the condition is satisfied and 0 otherwise. For generation, we instead use CITATION RATE, the proportion of queries where the target document is cited in the final response. (2) RANK CHANGE, the average positional shift of the target document between baseline and optimized settings:

$$\Delta\text{Rank} = \frac{1}{|Q|} \sum_{q \in Q} (\text{rank}_{\text{base}}(d_q^{\text{tgt}}) - \text{rank}_{\text{SAGEO}}(d_q^{\text{tgt}})), \quad (5)$$

where $\text{rank}_{\text{base}}$ and $\text{rank}_{\text{SAGEO}}$ represent the target document rank before and after optimization, respectively. Documents not appearing in the top- k candidates either at the retrieval or reranking stage are assigned a default rank of $k + 1$. We set $k = 100$ for retrieval and $k = 10$ for reranking. For generation, rank at this stage is defined as the citation order in which a source first appears in the response.

4 Experimental Setup

Optimization Strategies. We evaluate ten LLM-based optimization strategies derived from prior research. Eight are adapted from [1], each targeting a specific stylistic or content modification to improve document visibility. These include strategies that adjust tone and readability (Authoritative, Fluency, Easy Language), add supporting evidence (Cite Sources, Quotation, Statistics), and diversify vocabulary (Technical Terms, Unique Words). We additionally include All-in-One [37], which applies all eight strategies simultaneously along with structural formatting such as bolding and improved layout, and AutoGEO [49], which leverages preference rules learned from generative engine behavior to guide document optimization. Appendix B provides further details on strategies.

- **Authoritative:** Adopts a more confident and persuasive tone.
- **Cite Sources:** Adds inline citations to credible external sources.
- **Fluency:** Improves sentence flow and readability.
- **Quotation:** Incorporates direct quotes from authoritative sources.
- **Easy Language:** Simplifies vocabulary and shortens sentences.
- **Statistics:** Adds quantitative data and facts to support claims.
- **Technical Terms:** Introduces domain-specific terminology.
- **Unique Words:** Replaces common words with distinctive terms.
- **All-in-One:** Combines all eight strategies with formatting.
- **AutoGEO:** Rewrites documents using rules derived from analyzing which content generative search engines prefer to cite.

¹<https://vespa.ai/>

²<https://cloud.google.com/vertex-ai>

Table 2: Effectiveness of SAGEO strategies. We report Hit Rate and Δ Rank for each optimization target at each stage. Percentage values indicate relative change from the pre-optimization baseline. Positive Δ Rank indicates rank improvement. H@10 at reranking is 1.00 in the baseline, as all target documents are selected from the top-10 candidates at the reranking stage.

Strategy	Body Text only						Structural Information only						Both					
	Retrieval		Reranking		Generation		Retrieval		Reranking		Generation		Retrieval		Reranking		Generation	
	H@20	Δ Rank	H@10	Δ Rank	Cite	Δ Rank	H@20	Δ Rank	H@10	Δ Rank	Cite	Δ Rank	H@20	Δ Rank	H@10	Δ Rank	Cite	Δ Rank
Baseline	0.58 -- %	-	1.00 -- %	-	0.50 -- %	-	0.58 -- %	-	1.00 -- %	-	0.50 -- %	-	0.58 -- %	-	1.00 -- %	-	0.50 -- %	-
Auth.	0.57 -1%	-0.20	0.91 -9%	-0.24	0.49 -3%	-0.10	0.68 +18%	+1.60	0.81 -19%	-0.59	0.49 -3%	+0.01	0.69 +19%	+0.47	0.77 -23%	-0.80	0.48 -3%	+0.01
Cite	0.56 -3%	-1.50	0.87 -13%	-0.45	0.48 -4%	-0.13	0.74 +28%	+6.05	0.86 -14%	-0.17	0.52 +2%	+0.25	0.66 +13%	-1.81	0.73 -26%	-0.99	0.47 -6%	-0.04
EasyLang	0.52 -10%	-4.18	0.84 -16%	-0.54	0.49 -4%	-0.09	0.74 +27%	+5.15	0.86 -14%	+0.13	0.53 +5%	+0.38	0.71 +22%	+3.76	0.80 -20%	-0.30	0.51 +2%	+0.34
Fluency	0.57 -1%	-0.71	0.91 -9%	-0.18	0.50 -1%	-0.01	0.75 +30%	+6.62	0.88 -12%	+0.08	0.53 +5%	+0.37	0.72 +24%	+2.37	0.81 -19%	-0.42	0.49 -2%	+0.11
Quote	0.57 -1%	-0.33	0.90 -10%	-0.36	0.47 -7%	-0.26	0.74 +27%	+5.47	0.85 -15%	-0.28	0.53 +4%	+0.30	0.79 +35%	+8.87	0.85 -15%	-0.29	0.51 +2%	+0.24
Stats	0.57 -2%	-0.51	0.90 -10%	-0.33	0.48 -4%	-0.18	0.75 +29%	+6.03	0.86 -14%	-0.15	0.54 +7%	+0.47	0.71 +22%	+1.53	0.80 -20%	-0.75	0.48 -5%	-0.08
Tech.	0.50 -14%	-6.23	0.80 -20%	-1.03	0.47 -6%	-0.15	0.66 +13%	-0.59	0.79 -21%	-0.61	0.51 +1%	+0.22	0.62 +6%	-2.39	0.64 -36%	-2.02	0.43 -14%	-0.44
Unique	0.53 -8%	-3.47	0.86 -14%	-0.70	0.46 -8%	-0.33	0.69 +19%	+1.09	0.81 -19%	-0.62	0.49 -4%	-0.06	0.66 +13%	-0.16	0.75 -25%	-1.24	0.43 -13%	-0.48
All-in-One	0.50 -14%	-5.93	0.83 -17%	-0.68	0.49 -2%	-0.03	0.71 +22%	+2.44	0.82 -18%	-0.42	0.50 -1%	+0.12	0.69 +17%	+1.96	0.77 -23%	-0.62	0.52 +3%	+0.38
AutoGeo	0.37 -36%	-22.35	0.58 -42%	-2.28	0.39 -22%	-0.78	0.60 +4%	-6.68	0.73 -27%	-0.72	0.51 +0%	+0.31	0.44 -25%	-20.15	0.58 -42%	-1.93	0.44 -12%	-0.13
Avg.	0.53 -9%	-4.54	0.84 -16%	-0.68	0.47 -6%	-0.21	0.71 +22%	+2.72	0.83 -17%	-0.34	0.52 +2%	+0.24	0.67 +15%	-0.56	0.75 -25%	-0.94	0.48 -5%	-0.01

5 Results & Discussion

5.1 Main Results

We present findings from extensive experiments on SAGEO ARENA. Our goal is to evaluate how optimization strategies affect document visibility across the full generative search pipeline. Existing research focuses on body text alone in SAGE, leaving the effectiveness of such a setting in realistic environments largely unexplored. Moreover, optimizing structural information (Section 3.2), the process of organizing core document information into metadata fields so that search systems can correctly interpret the document, has not been studied in the context of generative search. We therefore separate optimization scope into three settings: body text only, structural information only, and both, enabling the first analysis of where gains and losses originate at each generative pipeline stage.

Limitation of Body Text only Optimization. As shown in Table 2 (Left), optimizing body text alone consistently degrades visibility across all stages. At retrieval, optimization strategies that replace common expressions with domain-specific terms (e.g., technical words) or uncommon vocabulary (e.g., unique words) show the largest retrieval drops. We attribute this to the lexical mismatch between optimized documents and user queries, which typically use common vocabulary. For example, replacing terms like “eating” with “alimentary routines” or “sleeping” with “somnia” directly reduces term overlap, causing BM25-based retrievers to assign lower relevance scores. Notably, AutoGEO [49] exhibits the largest degradation, with a retrieval rank drop of -22.35. We find that AutoGEO tends to substantially expand the document content, introducing lengthy rewrites that dilute keyword density and shift the document further from the original query vocabulary. At reranking, a similar but milder degradation is observed, suggesting that rewrites may introduce slight semantic shifts that rerankers are sensitive to. At generation, gains remain relatively marginal across strategies, aligning with recent findings that presentational modifications offer limited visibility improvements [37]. Overall, body-only optimization is insufficient for enhancing visibility throughout the pipeline. Although the average rank drop at retrieval and reranking remains

moderate, even a small shift can be critical in practice. In generative search, only top-ranked documents are passed to the generator. A document dropping below the input threshold is excluded entirely, making it invisible at the generation stage regardless of its content quality. These findings suggest that effective optimization must extend beyond body text to structural information, which plays a significant role in document ranking in real-world web search.

Effectiveness of Structural Information Optimization. As shown in Table 2 (Center), extending optimization scope to structural information significantly improves visibility across all strategies compared to body-only optimization. The improvement is particularly notable at retrieval, with a +22% boost in Hit Rate and +2.72 average retrieval rank gain. Structural information is inherently designed to be dense with query-relevant terms, increasing lexical overlap with user queries that BM25-based retrievers directly prioritize. Accordingly, strategies that enrich content with keywords, entities, and numbers show strong effectiveness. For example, adding statistics rewrites a verbose meta description into a concise summary with specific facts like “1983, comedy, Rotten Tomatoes, grossed 61M,” while adding quotations refines a generic title “Panel Clarifies Advice” into “IOM Panel Clarifies Vitamin D Guidance,” introducing entity-specific terms that better align with user queries. Interestingly, optimizing structural information alone also improves visibility at reranking and generation. We attribute this partly to the increased number of documents that pass through retrieval into downstream stages, raising their chances of being reranked favorably and cited. Additionally, optimizing structural information effectively places well-organized, informative summaries at the top of the document, consistent with recent findings that such positioning benefits document visibility in generation stage [37]. However, applying optimization to both scopes (Table 2 Right) yields lower visibility gains. The negative retrieval effects of body-text modification, observed in the body-only setting, could partially diminish the gains from optimizing structural information.

Reranking as a Persistent Bottleneck. All optimization strategies struggle at the reranking stage, showing consistent visibility



Figure 3: Reranking case studies illustrating two factors that improve document ranking after optimization. Case A shows that adding content directly addresses the query’s informational need boosts rank. Case B shows that placing the direct answer in early paragraphs improves ranking position.

degradation regardless of optimization scope (Table 2). One contributing factor is that top-ranked documents are very close in relevance, making them more sensitive to rank displacement from even minor content changes. Although the average rank drop remains within 1 position across most settings, even such small shifts can alter the relative ordering among closely ranked candidates. Notably, 5.8% of target documents in our experiments dropped from rank 10 to 11 during reranking, narrowly missing the input threshold for the generator in SAGEO ARENA. This highlights the importance of maintaining or improving rank position within generative search pipelines, as they inherently impose such cutoffs. To better understand what factors positively or negatively affect reranking results, we conduct a case study analyzing documents that gained or lost rank after optimization. Figure 3 presents two representative cases. It is worth noting that SAGEO optimizes documents without access to the incoming query, meaning content modifications cannot be tailored to specific query intents. Despite this constraint, clear patterns emerge in how the reranker responds to different types of optimizations. First, the reranker favors content additions that enhance alignment with the query’s informational need, while penalizing additions that expand the document’s scope beyond what the query seeks. Second, placing the answer early in the document yields higher reranking scores, whereas restructuring that displaces the answer to later paragraphs results in significant rank drops, even when the answer itself remains intact. These observations indicate that reranking-aware optimization should preserve topical alignment and answer prominence over broad content expansion.

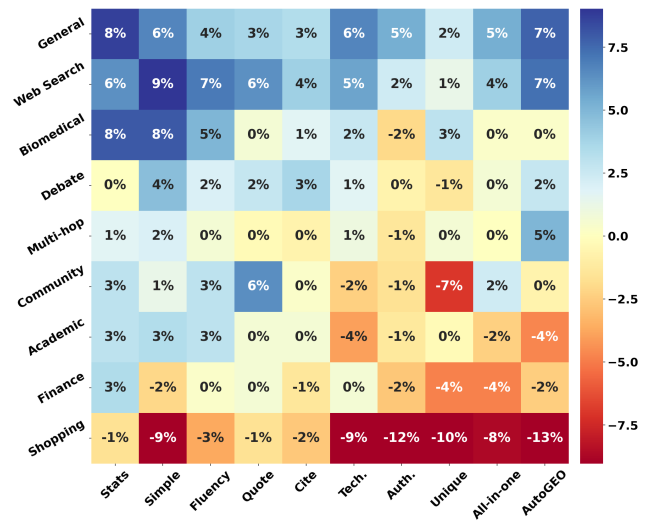


Figure 4: Effectiveness of each optimization strategy across nine domains, measured by change in citation rate at the generation stage. Positive % indicate increased citation rate.

5.2 In-depth Analysis

To further understand effectiveness of SAGEO, we analyze three dimensions: domain, citation source, and SAGEO backbone model.

Domain Analysis. Figure 4 shows the effectiveness of each optimization method across nine domains, reporting change in citation rate at the generation stage. We observe that domains with general information-seeking queries (e.g., Web Search, General QA) benefit the most, as strategies like adding statistics provide concrete evidence that directly address these queries. Shopping is the only domain where every optimization method decreases citation likelihood. Product documents are often already well-organized for their customers, and the queries in this domain are predominantly casual and recommendation-seeking (e.g., “gift ideas for my friend”). Optimizing such documents may shifts them away from this expected tone, reducing their competitiveness in the generator’s response. These results suggest that optimization strategies applied without considering domain-specific user intent can be ineffective or even harmful, emphasizing the need for domain-aware optimization.

Citation Source Analysis. To understand what drives citation at the generation stage, we examine how generative models interact with document content when forming responses. To achieve this, we conduct an experiment prompting the generator to provide the exact quote used for each citation in its response. We then apply fuzzy matching to locate each quoted segment in the source document and identify the region from which it originates. Figure 6 illustrates the results as a density plot. Each colored region corresponds to a structural component (e.g., title, meta description, headings, JSON-LD) or body text, and values above the horizontal baseline ($y = 1.0$) indicate that the region tends to be cited more frequently. We observe that the vast majority of citations originate from body text, while structural information is cited less frequently despite its strong contribution to retrieval. We attribute this to the

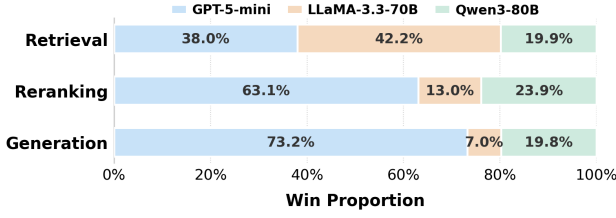


Figure 5: Win rate comparison across three backbone models (GPT-5-mini, LLaMA-3.3-70B, Qwen3-80B) at each stage. A model is considered a win when its optimized document achieves the highest rank among the three at a given stage.

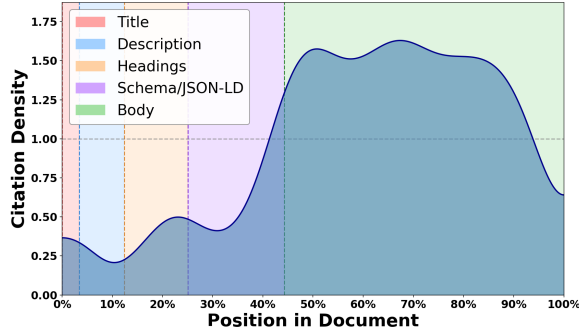


Figure 6: Density of citation sources across document regions.

information density of body text, which provides richer evidence for addressing user queries compared to the concise, keyword-oriented nature of structural information. Combined with earlier finding that structural information improves visibility at the generation stage (Section 5), these results suggest that structural information plays an important role in surfacing documents, but the generator primarily references body text when forming responses. Thus, structural information and body text serve complementary roles in SAGE, and effective optimization requires addressing both.

SAGEO Backbone Model Analysis. To examine whether the choice of backbone model for SAGEO exhibit distinct optimization behaviors, we conduct experiments with two additional open-source models, LLaMA-3.3-70B and Qwen3-80B (Figure 5). We provide comparison against GPT-5-mini, applying the All-in-One strategy across all three models. Interestingly, LLaMA-3.3-70B achieves the highest win rate (42.2%) at retrieval, surpassing GPT-5-mini (38.0%). We find that LLaMA often generates short, keyword-dense text (296 words on average vs. 782 for GPT). Because BM25 rewards keyword matches and applies length normalization that favors shorter documents, this style is naturally advantaged at the retrieval stage. However, this advantage reverses at reranking, where GPT-5-mini dominates with a 63.1% win rate while LLaMA-3.3-70B drops to 13.0%. Unlike BM25, the LLM-based reranker evaluates both relevance and coherence of the documents to the user query, and tends to penalize irrelevant content [23, 44]. Case study reveals that LLaMA-optimized documents frequently mention query terms without meaningfully addressing the underlying question. For instance, given the query “What are some of the most useful vim

Table 3: Comparison of SAGEO strategy combination and STAGEAWARE optimization (Ours) under the Both setting. STAGEAWARE achieves the strongest SAGEO performance.

Method	Retrieval		Reranking		Generation	
	H@20	Δ Rank	H@10	Δ Rank	Cite	Δ Rank
Baseline	0.58	-	1.00	-	0.50	-
Combined Strategy						
EasyLang + Quote	0.67 +14%	-0.25	0.78	-0.33	0.54	+0.64
EasyLang + Stats	0.65 +11%	-1.77	0.76	-0.56	0.51	+0.38
Fluency + Quote	0.66 +13%	-0.63	0.77	-0.43	0.54	+0.65
Fluency + Stats	0.65 +11%	-1.63	0.76	-0.55	0.53	+0.60
STAGEAWARE (Ours)	0.75 +28%	+4.86	0.80	-0.08	0.58	+1.01

shortcuts?”, LLaMA opens with unrelated introduction while GPT directly lists shortcuts with clear formatting. At generation, the gap widens further. GPT-5-mini reaches a 73.2% win rate while LLaMA-3.3-70B falls to 7.0%. These results demonstrate that short, keyword-dense documents can succeed at retrieval, but downstream stages increasingly reward content that substantively addresses the query.

5.3 Insights & Exploration

In this section, we study effective strategy combinations and introduce practical optimization guidelines informed by our findings.

Strategy Combination. When applying SAGEO, multiple strategies can be combined. However, the All-in-One optimization results in Table 2 suggest that combining all strategies at once is not always effective. We therefore selectively evaluate combinations of top-performing strategies (Table 3). While these combinations improve generation performance, retrieval rank degradation persists across all pairs. This indicates that naive pairing cannot simultaneously benefit all pipeline stages, motivating a stage-targeted approach.

Practical Guidelines. To this end, we introduce stage-aware SAGEO, which tailors optimization to the specific priorities of each pipeline stage. Our method builds on following core principles:

- **Enrich structural fields with key entities.** Core entities, numbers, and terms from body text are organized in structural fields to strengthen query-document matching at retrieval.
- **Make claims prominent and self-contained.** The main claim is placed at the start of the body, and each statement is ensured to carry specific evidence, increasing citability at generation.
- **Maintain topical coherence across paragraphs.** Ambiguous pronouns are replaced with explicit subject references, and core terms are naturally emphasized to keep paragraphs connected.
- **Adapt to domain context.** The document’s domain is assessed before any changes, and both domain characteristics and user intent are considered to determine how to optimize.

We demonstrate the effectiveness of our method in Table 3, where it achieves competitive performance across all stages, with notable gains at both reranking and generation among all strategies evaluated. These results suggest that SAGEO ARENA serves as an effective evaluation environment for developing practical SAGEO methods, enabling stage-level analysis providing actionable insights. The detailed prompt used in our method is provided in Appendix E.

6 Conclusion

In this paper, we introduce SAGEO ARENA, a realistic evaluation environment for stage-level visibility analysis of Search-Augmented Generative Engine Optimization. By integrating a full generative search pipeline with a large-scale corpus preserving structural information, SAGEO ARENA enables investigation of how optimization signals propagate in generative search. We observe that findings from existing benchmarks do not readily generalize to realistic settings, and that stage-aware optimization targeting each pipeline stage is crucial for visibility in generative search. Overall, SAGEO ARENA offers a reproducible environment for developing optimization strategies that generalize to real-world generative search.

References

- [1] Pranjal Aggarwal, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, Karthik Narasimhan, and Ameet Deshpande. 2024. Geo: Generative engine optimization. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5–16.
- [2] Fayyaz Ali and Shah Khuro. 2021. Content and link-structure perspective of ranking webpages: A review. *Comput. Sci. Rev.* 40 (2021), 100397. <https://api.semanticscholar.org/CorpusID:233535229>
- [3] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2018. MS MARCO: A Human Generated Machine Reading Comprehension Dataset. *arXiv:1611.09268 [cs.CL]* <https://arxiv.org/abs/1611.09268>
- [4] Bing Webmaster Blog. 2025. How AI Search Is Changing the Way Conversions are Measured. <https://blogs.bing.com/webmaster/November-2025/How-AI-Search-Is-Changing-the-Way-Conversions-are-Measured>
- [5] Vera Boteva, Demian Gholipour, Artem Sokolov, and Stefan Riezler. 2016. A Full-Text Learning to Rank Dataset for Medical Information Retrieval. In *Proceedings of the European Conference on Information Retrieval (ECIR)* (Padova, Italy). Springer.
- [6] Sergey Brin and Lawrence Page. 1998. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems* 30, 1-7 (1998), 107–117.
- [7] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 17754–17762.
- [8] Mahe Chen, Xiaoxuan Wang, Kaiwen Chen, and Nick Koudas. 2025. Generative engine optimization: How to dominate ai search. *arXiv preprint arXiv:2509.08919* (2025).
- [9] Mahe Chen, Xiaoxuan Wang, Kaiwen Chen, and Nick Koudas. 2026. Navigating the Shift: A Comparative Analysis of Web Search and Generative AI Response Generation. *arXiv preprint arXiv:2601.16858* (2026).
- [10] Qiyuan Chen, Jiahe Chen, Hongsen Huang, Qian Shao, Jintai Chen, Renjie Hua, Hongxia Xu, Ruijia Wu, Ren Chuan, and Jian Wu. 2025. Beyond Keywords: Driving Generative Search Engine Optimization with Content-Centric Agents. *arXiv preprint arXiv:2509.05607* (2025).
- [11] Xiaolu Chen, Haojie Wu, Jie Bao, Zhen Chen, Yong Liao, and Hu Huang. 2025. Role-Augmented Intent-Driven Generative Search Engine Optimization. *arXiv preprint arXiv:2508.11158* (2025).
- [12] Harold Davis. 2006. *Search engine optimization*. "O'Reilly Media, Inc".
- [13] Xiang Deng, Prashant Shiralkar, Colin Lockard, Binxuan Huang, and Huan Sun. 2022. DOM-LM: Learning Generalizable Representations for HTML Documents. *ArXiv abs/2201.10608* (2022). <https://api.semanticscholar.org/CorpusID:246285527>
- [14] Elastic. 2025. Relevance Tuning Guide, Weights and Boosts. <https://www.elastic.co/guide/en/app-search/current/relevance-tuning-guide.html>. Accessed: 2026-01-08.
- [15] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997* 2, 1 (2023).
- [16] Google. 2024. How Google Search Works. <https://developers.google.com/search/docs/fundamentals/how-search-works>. Accessed: 2025-01-19.
- [17] Google. 2025. How Google Search Works. <https://developers.google.com/search/docs/fundamentals/how-search-works>.
- [18] Google Cloud. 2026. Vertex AI RAG Engine Overview. <https://docs.cloud.google.com/vertex-ai/generative-ai/docs/rag-engine-overview>. Accessed on Jan 7 2026.
- [19] Gregory Goren, Oren Kurland, Moshe Tennenholtz, and Fiana Raiber. 2020. Ranking-incentivized quality preserving content modification. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 259–268.
- [20] Venkat N Gudivada, Dhana Rao, and Jordan Paris. 2015. Understanding search-engine optimization. *Computer* 48, 10 (2015), 43–52.
- [21] Yu Guo, Zhengyi Ma, Jiaxin Mao, Hongjin Qian, Xinyu Zhang, Hao Jiang, Zhao Cao, and Zhicheng Dou. 2022. Webformer: Pre-training with Web Pages for Information Retrieval. *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (2022). <https://api.semanticscholar.org/CorpusID:250340272>
- [22] Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv preprint arXiv:2503.09516* (2025).
- [23] Zixuan Ke, Weize Kong, Cheng Li, Mingyang Zhang, Qiaozhu Mei, and Michael Bendersky. 2024. Bridging the preference gap between retrievers and llms. *arXiv preprint arXiv:2401.06954* (2024).
- [24] Aounon Kumar and Himabindu Lakkaraju. 2024. Manipulating large language models to increase product visibility. *arXiv preprint arXiv:2404.07981* (2024).
- [25] Arlen Kumar and Leanid Palkhouski. 2025. AI Answer Engine Citation Behavior An Empirical Analysis of the GEO16 Framework. *arXiv preprint arXiv:2509.10762* (2025).
- [26] Oren Kurland and Moshe Tennenholtz. 2022. Competitive search. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2838–2849.
- [27] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics* 7 (2019), 452–466. doi:10.1162/tacl_a_00276
- [28] Dirk Lewandowski, Sebastian Sunkler, and Nurce Yagci. 2021. The influence of search engine optimization on Google's results: A multi-dimensional approach for detecting SEO. In *Proceedings of the 13th ACM Web Science Conference 2021*. 12–20.
- [29] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems* 33 (2020), 9459–9474.
- [30] Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yongkang Wu, Ji-Rong Wen, Yutao Zhu, and Zhicheng Dou. 2025. Webthinker: Empowering large reasoning models with deep research capability. *arXiv preprint arXiv:2504.21776* (2025).
- [31] Macedo Maia, Siegfried Handschuh, André Freitas, Brian Davis, Ross McDermott, Manel Zarrouk, and Alexandra Balahur. 2018. WWW'18 Open Challenge: Financial Opinion Mining and Question Answering. (2018), 1941–1942.
- [32] Nektarios Makrydakis. 2024. SEO mix 6 O's model and categorization of search engine marketing factors for websites ranking on search engine result pages. *International Journal of Research in Marketing Management and Sales* (2024). <https://api.semanticscholar.org/CorpusID:267567479>
- [33] Microsoft. 2025. How Bing Delivers Search Results. <https://support.microsoft.com/en-us/topic/how-bing-delivers-search-results-d18fc815-ac37-4723-bc67-9229ce3eb6a3>.
- [34] Microsoft. 2025. Webmasters Guidelines. <https://www.bing.com/webmasters/help/webmasters-guidelines-30fba23a>.
- [35] Fredrik Nestaas, Edoardo Debenedetti, and Florian Tramèr. 2024. Adversarial search engine optimization for large language models. *arXiv preprint arXiv:2406.18382* (2024).
- [36] Samuel Pfrommer, Yatong Bai, Tanmay Gautam, and Somayeh Sojoudi. 2024. Ranking manipulation for conversational search engines. *arXiv preprint arXiv:2406.03589* (2024).
- [37] Haritz Puerto, Martin Gubri, Tommaso Green, Seong Joon Oh, and Sangdoo Yun. 2025. C-SEO Bench: Does Conversational SEO Work? *arXiv preprint arXiv:2506.11097* (2025).
- [38] Nimrod Raifer, Fiana Raiber, Moshe Tennenholtz, and Oren Kurland. 2017. Information retrieval meets game theory: The ranking competition between documents' authors. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 465–474.
- [39] David Rau, Hervé Déjean, Nadezhda Chirkova, Thibault Formal, Shuai Wang, Stéphane Clinchant, and Vassilina Nikoulina. 2024. Bergen: A benchmarking library for retrieval-augmented generation. In *Findings of the Association for Computational Linguistics: EMNLP 2024*. 7640–7663.
- [40] Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. *Found. Trends Inf. Retr.* 3, 4 (April 2009), 333–389. doi:10.1561/15000000019
- [41] Stephen E. Robertson, Hugo Zaragoza, and Michael J. Taylor. 2004. Simple BM25 extension to multiple weighted fields. In *International Conference on Information and Knowledge Management*. <https://api.semanticscholar.org/CorpusID:16628332>
- [42] Dushyant Sharma, Rishabh Shukla, Anil Kumar Giri, and Sumit Kumar. 2019. A brief review on search engine optimization. In *2019 9th international conference on cloud computing, data science & engineering (confluence)*. IEEE, 687–692.

- [43] Lakshay Sharma, Laura Graesser, Nikita Nangia, and Utku Evcı. 2019. Natural Language Understanding with the Quora Question Pairs Dataset. *arXiv:1907.01041* [cs.CL]. <https://arxiv.org/abs/1907.01041>
- [44] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. Replug: Retrieval-augmented black-box language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*. 8371–8384.
- [45] Jiejun Tan, Zhicheng Dou, Wen Wang, Mang Wang, Weipeng Chen, and Ji-Rong Wen. 2024. HtmlRAG: HTML is Better Than Plain Text for Modeling Retrieved Knowledge in RAG Systems. *Proceedings of the ACM on Web Conference 2025* (2024). <https://api.semanticscholar.org/CorpusID:273821119>
- [46] Erdiñ Uzun, Hayri Volkan Agun, and Tarik Yerlikaya. 2013. A hybrid approach for extracting informative content from web pages. *Inf. Process. Manag.* 49 (2013), 928–944. <https://api.semanticscholar.org/CorpusID:34275267>
- [47] Qifan Wang, Yi Fang, Anirudh Ravula, Fuli Feng, Xiaojun Quan, and Dongfang Liu. 2022. WebFormer: The Web-page Transformer for Structure Information Extraction. *Proceedings of the ACM Web Conference 2022* (2022). <https://api.semanticscholar.org/CorpusID:246441911>
- [48] Xiaohua Wang, Zhenghua Wang, Xuan Gao, Feiran Zhang, Yixin Wu, Zhibo Xu, Tianyuan Shi, Zhengyuan Wang, Shizheng Li, Qi Qian, et al. 2024. Searching for best practices in retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 17716–17736.
- [49] Yujia Wu, Shanshan Zhong, Yubin Kim, and Chenyan Xiong. 2025. What Generative Search Engines Like and How to Optimize Web Content Cooperatively. *arXiv preprint arXiv:2510.11438* (2025).
- [50] Rongwu Xu, Xuan Qi, Zehan Qi, Wei Xu, and Zhijiang Guo. 2024. DebateQA: Evaluating Question Answering on Debatable Knowledge. *arXiv:2408.01419* [cs.CL]. <https://arxiv.org/abs/2408.01419>
- [51] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*. 2369–2380.
- [52] Lan Yi, B. Liu, and Xiaoli Li. 2003. Eliminating noisy information in Web pages for data mining. In *Knowledge Discovery and Data Mining*. <https://api.semanticscholar.org/CorpusID:6451023>
- [53] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025. Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models. *arXiv preprint arXiv:2506.05176* (2025).
- [54] Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025. Deepresearcher: Scaling deep research via reinforcement learning in real-world environments. *arXiv preprint arXiv:2504.03160* (2025).

A Benchmark Construction

Statistics. SAGEO Arena comprises nine domains sourced from established query datasets: MS MARCO [3] (Web Search), Natural Questions (General QA) [27], HotpotQA [51] (Multi-hop QA), NRCorpus [5] (Biomedical), Quora [43] (Community QA), FiQA [31] (Finance), DebateQA [50] (Debate), E-commerce (Shopping), and Researchy (Academic) [49]. Domain selection ensures diversity in content characteristics, query complexity, and information-seeking behaviors. Detailed statistics are presented in Table 4.

Table 4: Benchmark statistics across nine domains.

Source Dataset	Domain	#Sampled Queries	#Retrieved Docs
MS MARCO	Web Search	300	21,880
Natural Questions	General QA	300	21,921
HotpotQA	Multi-hop QA	300	13,409
NRCorpus	Biomedical	300	21,079
Quora	Community QA	300	21,734
FiQA	Finance	300	16,771
DebateQA	Debate	300	17,443
E-commerce	Shopping	300	18,885
Researchy	Academic	300	17,881
Total		2,700	171,003

Difference from RAG Benchmarks. Standard RAG benchmarks evaluate end-to-end answer quality, measuring whether a system produces faithful and relevant responses given a query and corpus. While a RAG benchmark provides queries and a corpus, it typically does not define a document-level evaluation protocol. Specifically, it lacks explicit target documents per query, stage-wise visibility metrics, and controls to separate document effects from answer quality effects. As a result, standard RAG benchmarks cannot isolate the effects of document-level modifications or attribute improvements to specific pipeline stages. Our benchmark addresses this limitation through a document-centric evaluation protocol. We explicitly designate target documents per query and track their visibility across retrieval, reranking, and generation stages. This design separates the effects of document optimization from improvements in answer generation, enabling controlled evaluation of SAGEO strategies.

B SAGEO Optimization Strategies

We describe the following ten optimization strategies evaluated in the SAGEO Arena and their specific transformation objectives.

- **Authoritative:** Modifies text style to be more persuasive, making statements definitive while maintaining factual accuracy.
- **Cite Sources:** Adds inline citations to external sources, inserting brief references to enhance reliability and trustworthiness.
- **Fluency:** Improves grammatical correctness by refining sentence structure and transitions without altering content meaning.
- **Quotation:** Incorporates quotations from authoritative figures or sources with proper attribution to provide external validation.
- **Easy Language:** Simplifies by replacing complex vocabulary with accessible alternatives while preserving information.
- **Statistics:** Adds quantitative data and numerical facts inline within sentences to make claims more concrete and credible.
- **Technical Terms:** Introduces domain-specific terminology to present content in a more expert and authoritative manner.
- **Unique Words:** Enriches vocabulary by incorporating less common words to signal higher content quality and specialization.
- **All-in-One:** Combines all eight base strategies, then enhances text structure by bolding key features and improving layout to increase information accessibility and visual attractiveness.
- **AutoGEO:** Applies optimization rules learned from analyzing generative engine citation patterns to guide document rewriting.

C Implementation Details

We instantiate the pipeline described in Section 3.3 with the following configuration. We segment document body text into passages of 256 tokens with a 64-token overlap. During indexing, we assign equal weights across structural fields and chunked passages to establish a controlled and fair baseline. In practice, search engines apply varying weights to different fields [14]; however, adopting uniform weights ensures that observed visibility changes stem from optimization strategies rather than field-specific biases. Nevertheless, our modular design allows researchers to configure alternative weighting schemes as needed. We retrieve the top-100 passages per query and rescore them using Qwen3-Reranker-4B [53]. For generation, the top-10 candidates are retained for response synthesis, using gpt-5-mini as the default generation model. For each test

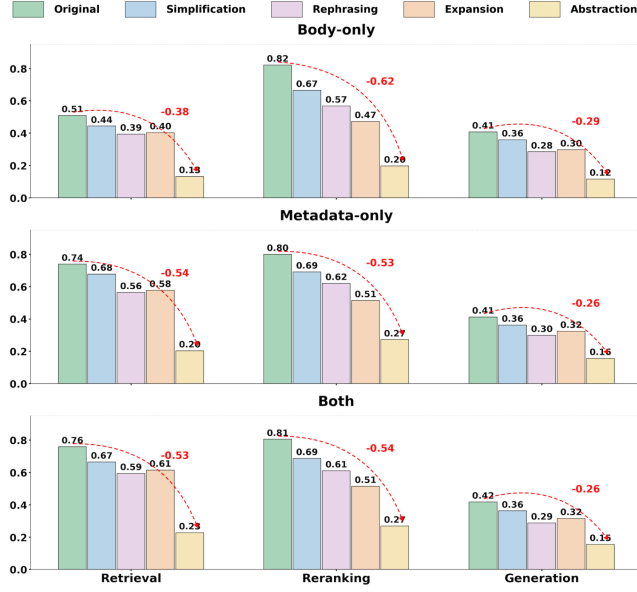


Figure 7: Visibility across pipeline stages under query variations for three optimization scopes. Values show mean visibility across eight optimization strategies, with red numbers indicating percentage point change from original query.

query, we first establish a baseline by running the query through the full search pipeline, then randomly select a target document from the top-10 candidates that enter the generation stage (i.e., ranked within the top-10 after reranking), ensuring sufficient relevance to the test query. All SAGEO optimization strategies are implemented using GPT-5-mini with strategy-specific prompts following Aggarwal et al. [1] for the eight content modification strategies, Puerto et al. [37] for All-in-One, and Wu et al. [49] for AutoGEO. We apply each strategy under three optimization scopes, body text only,

structural fields only, or both, to study how different document components contribute to visibility at each pipeline stage.

D Further Analysis

Optimization strategies show limited robustness across query formulations. Real users express information need through diverse query formulations, making robustness to query variation essential for effective optimization. We evaluate optimized documents against four query rewriting strategies: expansion (adding specificity), simplification (reducing detail), rephrasing (lexical substitution), and abstraction (semantic generalization). We report results averaged across all eight optimization methods in Figure 7. Across all optimization scopes, every query variation degrades performance. We observe that optimization tends to elaborate on the document’s existing content, making it more specific to its original context but less adaptable to alternative query formulations. As a result, performance drops progressively as query formulations deviate further from the original phrasing. Yet abstraction causes the most severe decline as failures concentrate at the retrieval stage, where documents fall entirely out of the candidate pool. Structural optimization provides partial mitigation by introducing alternative phrasings, but it is not enough to compensate for abstracted queries. This suggests that maintaining vocabulary overlap between queries and documents is important for optimization to be effective, yet current optimization strategies lack the capacity to anticipate the diverse ways users may express the same information need.

E Prompt for Stage-Aware SAGEO

Based on our analysis in Section 5.3, we introduce a stage-aware SAGEO strategy that tailors optimization to the specific priorities of each pipeline stage. The complete prompt template for our strategy is shown in Table 5.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

Table 5: The prompt template for stage-aware SAGEO.

Stage-Aware SAGEO Prompt
[Task Description] Optimize the following document with these strategies. Keep the content faithful to the original.
[Pre-Optimization Considerations] Before optimizing, think about two things: 1. Domain: Consider what domain this document belongs to (e.g., medical, finance, e-commerce, technical, casual). Match the tone and vocabulary to what readers in that domain expect. 2. Quality: If a field is already clear, specific, and well-written, keep it as-is. Only optimize fields that genuinely benefit from it. Not every document needs heavy changes.
[Optimization Strategies] 1. Entity mirroring (structural fields): Incorporate key entities, numbers, and domain terms from the body into the title, meta_description, headings, and jsonld_text. Add a keyword-rich summary sentence while keeping compact. Skip if the structural fields already contain the right keywords. 2. Fluent, easy language (all text): Rewrite sentences to be smooth, clear, and easy to read. Use simple words and short sentences. Avoid jargon when a plain alternative exists. If the writing is already clear and fluent, leave it unchanged. 3. Concrete evidence (body text): Make claims specific. Bring front the main claim to the very start of the body. Each claim should be self-contained — a reader should understand it without reading surrounding text. If claims are already specific, do not rephrase them. 4. Keyword reinforcement (body text): Naturally repeat the document’s core topic terms and key phrases throughout the body. Use the main subject name instead of pronouns where it reads naturally. This keeps every paragraph clearly connected to the topic.