

EVIREPORT: From Reasoned Outlines to Evidence Tracked Long-Form Reports

Anonymous ACL submission

Abstract

Evidence-intensive analytical reports are expected to be fact-dense, quantitatively correct, and supported by figures. Yet one-shot long-form generation with large language models (LLMs) frequently produces fluent but under-supported drafts: core facts are missed, numbers drift, and key visuals are absent, making the report hard to trust. We propose EVIREPORT, an evidence-tracked report-writing workflow that improves reliability by (i) organizing corpus evidence into compact, traceable units and retrieves query-relevant subgraphs into retrieval-ready packages (ii) leveraging a reasoning-focused LLM sketches a high-level plan for full coverage, then a chat-based LLM sharpens it into a detailed hierarchical outline with explicit scope and ordering (iii) rive generation with a facts-first iterative loop: extracting verifiable facts, composing strictly from those facts, then triggering gap-aware append queries to fill missing evidence To evaluate both correctness and completeness, we introduce EviReportBench, a benchmark instantiated on data-rich indicator reports that measures factual accuracy (claim verification), factual coverage (quiz-based evaluation), and visual evidence integration (image recall). Across 8 topics, experiments show that EVIREPORT consistently outperforms strong baselines in factual coverage (2.16 \times), factual accuracy (+8.9 points), and visual evidence integration (+34 points), approaching the quality of expert-written reports across multiple dimensions.

1 Introduction

Evidence-intensive *long-form* reports are widely used for assessment and policy decision-making because they synthesize heterogeneous materials (papers, official documents, indicator tables, and figures) into a coherent narrative with verifiable claims. A representative example is the **United Nations Sustainable Development Goals (SDGs)** reporting ecosystem: SDGs define 17 goals with an

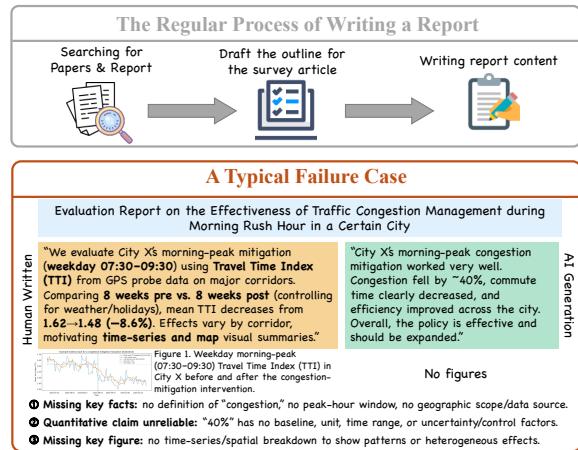


Figure 1: Report writing is evidence-intensive. **Top:** a regular workflow (search, outline, draft). **Bottom:** a typical failure of direct long-form LLM generation on a morning-peak congestion mitigation report: the AI text often (1) misses key facts, (2) makes unverifiable numbers, and (3) lacks critical figures, compared with a human-written version grounded in metric/time window/scope and a time-series plot.

official global indicator framework, and progress is regularly documented through indicator-centric reports that integrate multi-source evidence and rely heavily on figures (United Nations Statistics Division; United Nations General Assembly, 2017; United Nations Department of Economic and Social Affairs, 2025). In this work, we focus on automating SDG-style report writing as a primary application scenario and a stress test for evidence-intensive report generation.

Large language models (LLMs) can accelerate drafting (Jin et al., 2024; OpenAI, 2025a), but directly prompting an LLM to generate a full report often leads to three recurring failures (Fig. 1). First, the output may exhibit *low factual coverage*: it misses topic-critical facts while including generic or irrelevant points that do not advance the analysis. Second, the output may contain *quantitative unreliability*, including factual errors in

trends/ranges/numbers and under-specified claims that omit essential constraints such as year, unit, population, or data source. Third, it often shows weak visual evidence integration, omitting key figures (e.g., maps and trend plots) or failing to align textual claims with the corresponding visuals. These issues directly undermine report credibility and motivate workflows that explicitly control evidence use, planning, and multimodal grounding.

Evidence-intensive long-form report generation poses challenges beyond producing fluent text. Real-world report evidence is typically *structured*: high-level targets/indicators are argued through concrete cases supported by specific datasets and analytical methods, and ignoring such evidence chains leads to omissions or unsupported statements. Report quality also depends on *planning*: outlines determine what will be covered and how the narrative is organized, yet producing comprehensive and non-redundant outlines under long contexts and complex constraints remains difficult. Moreover, many reports are *multimodal*, where figures carry core information, and evaluation is non-trivial because valid reports are not unique and single-aspect scoring cannot diagnose whether a system is correct, complete, and visually grounded.

To address these challenges, we present EVIREPORT, an end-to-end workflow for evidence-tracked multimodal long-form report generation. EVIREPORT follows a retrieve-plan-write paradigm with three stages. (1) **Graph-enhanced evidence retrieval**: we build an evidence-oriented knowledge structure from the reference corpus to organize evidence into compact, traceable units, retrieve and prune query-relevant subgraphs, and summarize associated evidence into retrieval-ready packages. (2) **Two-stage outline planning**: a reasoning-oriented LLM drafts a coarse plan to ensure global coverage, and a chat-oriented LLM refines it into a fine-grained hierarchical outline with clearer scopes and ordering. (3) **Evidence-guided multimodal writing**: for each subsection, EVIREPORT performs targeted text and figure retrieval (including caption-guided image retrieval) and generates content via a three-step iterative loop (facts \rightarrow content \rightarrow gap-aware append query) to reduce under-specified claims and recover missing evidence.

A major obstacle to progress is the lack of a diagnostic benchmark that jointly measures *correctness*, *coverage*, and *visual evidence integration* for real-world, data- and figure-intensive reports. To fill

this gap, we introduce **EviReportBench**, a multi-dimensional benchmark instantiated on SDG-style indicator reports (United Nations Department of Economic and Social Affairs, 2025; United Nations Secretary-General, 2025). Beyond claim-level factual verification, EviReportBench measures factual coverage via a quiz-based protocol derived from authoritative reports, and evaluates multimodal integration via image recall against topic-specific ground-truth figures. Experiments across eight representative topics show that EVIREPORT consistently outperforms strong baselines on factual accuracy, factual coverage, and image recall. **Evaluation code**. We release the benchmark construction and scoring code at <https://anonymous.4open.science/r/EviReport-03E2>.

In summary, our contributions are:

- We propose EVIREPORT, a workflow for evidence-tracked multimodal long-form report generation with graph-enhanced retrieval, two-stage outline planning, and iterative evidence-guided writing.
- We introduce caption-guided image retrieval and a three-step iterative writing loop to improve quantitative reliability, factual coverage, and visual evidence integration.
- We build EviReportBench, a new multi-dimensional benchmark instantiated on SDG-style indicator reports, with claim verification for factual accuracy, quiz-based evaluation for factual coverage, and image recall for multimodal evidence integration.

2 Related Work

Long-form report generation, planning, grounding, and verification. Although LLMs excel at short-form NLG, generating long-form reports that are well-structured, coherent, and comprehensive remains challenging under long-context constraints and complex instructions (Dong et al., 2024; Kumar et al., 2024). Prior work therefore adopts planning and multi-stage pipelines that separate global organization from local drafting (Tang et al., 2022; Tan et al., 2021); for example, STORM (Shao et al., 2024) introduces a pre-writing stage to research a topic and construct an outline before drafting Wikipedia-style articles, while AutoSurvey (Wang et al., 2024) decomposes survey writing into retrieval, outlining, subsection drafting, refinement, and evaluation. Complementary to

164 planning, verification-based prompting improves
165 factuality by checking and revising model outputs
166 (Gao et al., 2023a; Zhao et al., 2023; Madaan et al.,
167 2023); Chain-of-Verification (CoVe) generates ver-
168 ification questions and revises drafts accordingly
169 (Dhuliawala et al., 2024). On the evaluation side,
170 claim-level verification (e.g., FactScore-style proto-
171 cols) provides more objective signals than holistic
172 ratings (Min et al., 2023). Our work complements
173 these directions with an iterative subsection-writing
174 loop (facts → draft → gap-aware evidence expan-
175 sion) and multi-dimensional evaluation.

176 **Retrieval augmentation, graph-enhanced re-**
177 **trieval, and SDG-style reports.** Retrieval aug-
178 mentation grounds generation on external evidence
179 (Lewis et al., 2020; Fan et al., 2024; Gao et al.,
180 2023b), and structured/relational evidence moti-
181 vates graph-based variants that retrieve entity-
182 centric subgraphs for more complete evidence
183 aggregation (Edge et al., 2024; Hu et al., 2024;
184 Gao et al., 2025). Representative systems include
185 KERAG (Sun et al., 2025), SubgraphRAG (Li et al.,
186 2025), and KAG (Liang et al., 2024), with sur-
187 veys summarizing common GraphRAG pipelines
188 (Peng et al., 2024). The UN SDGs define 17
189 goals with an official indicator framework for mon-
190 itoring progress (United Nations Statistics Divi-
191 sion; United Nations General Assembly, 2017);
192 SDG-style indicator reports are evidence- and data-
193 intensive, contain dense quantitative claims, and
194 rely heavily on figures (e.g., maps and trend plots),
195 making them a strong stress test for grounded mul-
196 timodal report generation. We therefore instantiate
197 our benchmark on authoritative SDG-style reports
198 to enable diagnostic evaluation of factual accuracy,
199 factual coverage, and visual evidence integration.

200 3 Methodology

201 In this section, we propose EVIREPORT, an
202 evidence-grounded workflow for long-form report
203 generation with LLMs. As shown in Fig. 2, EVIRE-
204 PORT consists of three main stages: knowledge-
205 guided retrieval, outline planning, and content draft-
206 ing. The retrieval stage constructs a lightweight
207 report knowledge graph from the reference corpus
208 and performs graph-enhanced retrieval to produce
209 compact, traceable evidence packages. The out-
210 line planning stage adopts a coarse-to-fine strategy,
211 where a reasoning-oriented LLM generates a high-
212 level plan and a chat-based LLM refines it into a
213 subsection-level outline. The content drafting stage

214 conducts plan-conditioned multimodal retrieval (in-
215 cluding caption-guided image retrieval) and gen-
216 erates each subsection via an iterative refinement
217 loop to improve correctness, coherence, and com-
218 pleteness. Finally, we introduce EviReportBench, a
219 multi-dimensional benchmark instantiated on SDG-
220 style indicator reports to evaluate factual accuracy,
221 factual coverage, and visual evidence integration.
222 The details are elaborated in Sec. 3.2, Sec. 3.3,
223 Sec. 3.4, and Sec. 3.5, respectively.

224 3.1 Preparation Phase

225 We preprocess heterogeneous reference materials
226 to support evidence-grounded long-form report
227 generation. Our sources include domain-specific
228 research papers and governmental/organizational
229 reports. Textual resources are parsed into a uni-
230 fied Markdown (MD) format, while structured in-
231 dicator or statistics tables are cleaned and stored
232 in a database. Since the end-to-end data pipeline
233 (collection, parsing, chunking, and indexing) is de-
234 scribed in Appendix A.2, we focus on how these
235 resources are organized into retrieval-ready knowl-
236 edge structures for downstream generation.

237 Reference preprocessing has two components:
238 (i) constructing an evidence-oriented knowledge
239 structure, and (ii) building retrieval corpora for text
240 and images.

241 3.1.1 Evidence-Oriented Knowledge 242 Structure

243 Analytical reports in many domains are typically
244 organized around a stable *reporting framework*
245 (e.g., a taxonomy of targets and measurable in-
246 dicators) and are supported by case-based anal-
247 yses grounded in specific datasets and computa-
248 tional methods. Motivated by this pattern, we build
249 an evidence-oriented knowledge graph that links
250 framework items (targets/indicators) to *cases*, and
251 further connects each case to its *datasets* and *meth-*
252 *ods*, forming an explicit “indicator–case–dataset–
253 method” evidence chain. This structure enables the
254 system to retrieve compact and traceable evidence
255 packages for a target report item, rather than rely-
256 ing on isolated text chunks. Appendix A.3 details
257 the domain-specific instantiation of this graph, in-
258 cluding schema design, extraction/fusion, and an
259 example subgraph.

260 During report writing, the system retrieves case-
261 level subgraphs and uses the linked datasets and
262 methods to steer evidence-aware retrieval and gen-
263 eration. This design improves factual coverage and

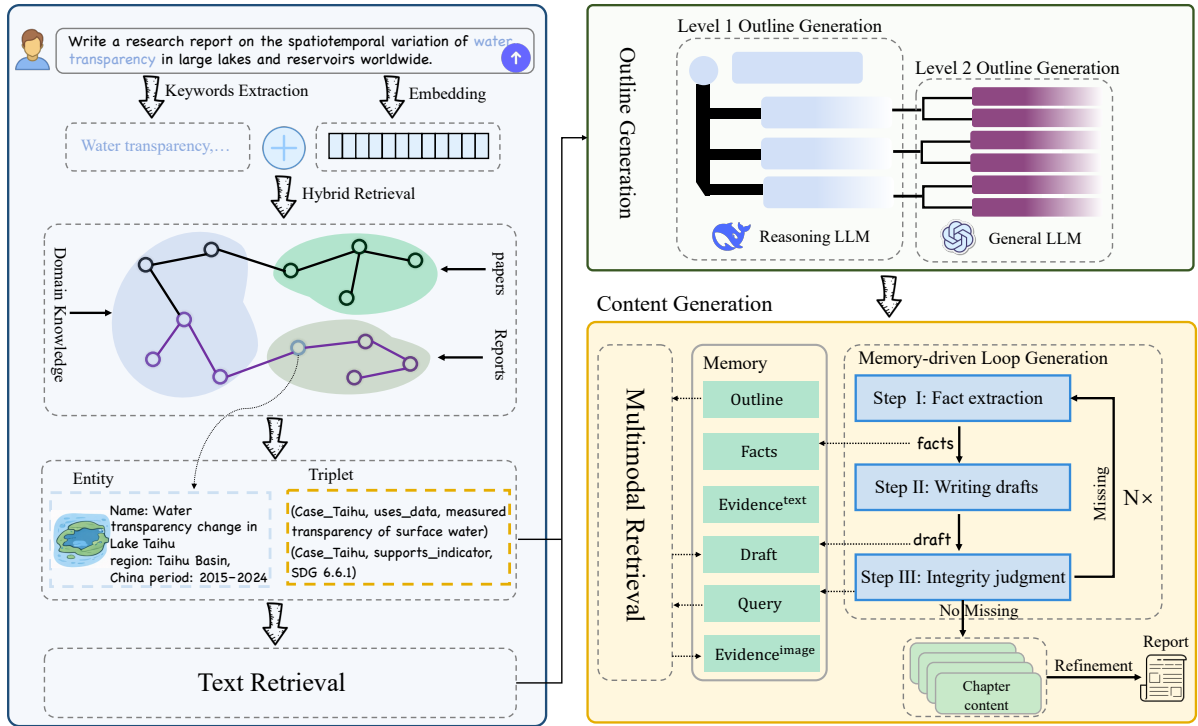


Figure 2: An overview of EVIREPORT for evidence-grounded long-form report generation. EVIREPORT combines (i) graph-enhanced hybrid retrieval over heterogeneous evidence to recall subgraphs/entities/triples for RAG, (ii) two-stage hierarchical outline planning, and (iii) an iterative generation loop that extracts facts from recalled context, writes drafts, and generates append queries when information is missing, producing structured reports with traceable evidence.

reduces unsupported statements by keeping generation tightly coupled with retrievable evidence.

3.1.2 Retrieval Corpus Construction

Based on the constructed knowledge structure, we build retrieval corpora for report generation. Parsed MD documents are segmented into chapter-aware semantic units and normalized with source and temporal metadata, enabling hybrid sparse–dense retrieval (Appendix A.2). In parallel, we construct an image retrieval database aligned with the text corpus: each figure is paired with its caption and surrounding context, and a vision–language model extracts lightweight metadata (keywords and semantic summaries) for indexing (Fig. 4).

At inference time, retrieval is conditioned on chapter- and subsection-level writing plans. Text evidence and image candidates are jointly matched to the target framework item (and its associated cases), allowing the system to recall multimodal evidence that is both relevant and traceable.

3.2 Knowledge-graph-enhanced RAG retrieval

To provide high-quality and verifiable evidence prior to planning and drafting, we introduce a graph-enhanced retrieval-augmented generation (KG-RAG) module (Algorithm 1). Given an input query q , a knowledge graph \mathcal{G} constructed from the reference corpus, and a text database \mathcal{D} , the module outputs a set of evidence pairs \mathcal{P} . Each pair consists of a compact, query-relevant subgraph and a corresponding textual summary, serving as *retrieval-ready evidence packages* for downstream outline generation and subsection writing.

Concretely, the system first extracts keywords from q and computes a semantic embedding to support hybrid retrieval. It then retrieves candidate *case nodes* from \mathcal{G} —case nodes represent concrete evidence units (e.g., empirical studies, datasets, or documented analyses) that can anchor report claims. After merging and deduplicating candidates, we keep the top- K nodes to form \mathcal{E} . For each $c \in \mathcal{E}$, we extract an entity-centered subgraph \mathcal{G}_c within h hops and prune it to obtain \mathcal{G}'_c , retaining only nodes and relations that are relevant to

308 the query and useful for report writing (e.g., linked
309 entities, data sources, and methods).

310 Next, the pruned subgraph \mathcal{G}'_c is converted into
311 retrieval constraints (entities/relations) to query the
312 text database \mathcal{D} , yielding evidence text blocks \mathcal{B}_c
313 that are semantically consistent with the case and
314 its analytical context. We summarize \mathcal{B}_c into s_c to
315 distill key facts, numbers, and explicit gaps condi-
316 tioned on q , and then augment the subgraph with
317 summary-derived attributes to obtain \mathcal{G}''_c . Finally,
318 we collect paired evidence (\mathcal{G}''_c, s_c) across all candi-
319 dates to form \mathcal{P} , which provides compact, traceable
320 inputs for the subsequent *plan–write* stages.

321 3.3 Two-stage report outline generation

322 A high-quality outline is the backbone of long-
323 form report writing: it fixes the document struc-
324 ture, allocates space across themes, and implicitly
325 determines what evidence will be retrieved and
326 discussed. In our setting, outlining must satisfy
327 two competing requirements: it should be *globally*
328 *coherent* and cover major analytical dimensions,
329 while remaining *grounded* in the available evidence
330 under long-context constraints. In practice, asking
331 a single LLM to directly produce a detailed out-
332 line from lengthy, heterogeneous evidence often
333 yields either a generic template with low evidence
334 coverage or an overfit, fragmented structure.

335 **Design insight.** We find that outline generation
336 is best treated as a *planning problem* rather than a
337 one-shot formatting task. Global structure (“what
338 to cover”) and local organization (“how to break
339 down and order content”) require different behav-
340 iors: the former benefits from deliberate reasoning
341 over evidence and reporting conventions, while
342 the latter benefits from instruction-following and
343 stylistic consistency. This motivates a coarse-to-
344 fine decomposition.

345 **Two-stage planning.** We adopt a two-stage
346 outline generation strategy. In the first stage,
347 a reasoning-oriented LLM produces a *coarse-*
348 *grained chapter plan* that emphasizes global struc-
349 ture and major analytical sections (e.g., back-
350 ground, data/measurement, empirical findings,
351 comparisons, limitations, and implications), con-
352 ditioned on the retrieved evidence packages. In
353 the second stage, a chat-based LLM expands and
354 refines the plan into a *fine-grained, subsection-*
355 *level outline*, ensuring consistent formatting, log-
356 ical ordering, and improved coverage by explic-
357 itly mapping each subsection to a small set of sup-
358 porting evidence packages. By separating global

359 planning from local decomposition, this design
360 yields outlines that are both structurally coherent
361 and evidence-aware, providing stable writing plans
362 for the downstream subsection drafting stage.

363 3.4 Content generation

364 Given the refined outline $\mathcal{O} = \{o_1, \dots, o_N\}$
365 (Sec. 3.3), EVIREPORT drafts the report in a *plan-*
366 *conditioned, serial* manner. We generate sub-
367 sections sequentially so that each subsection o_i
368 can condition on previously generated content
369 $\mathcal{C}_{<i} = \{c_1, \dots, c_{i-1}\}$, helping maintain global
370 consistency in terminology, narrative flow, and
371 cross-section references. As illustrated in Fig. 2
372 and formalized in Algorithm 2, each subsection
373 is produced through an iterative *retrieve–generate*
374 loop with explicit multimodal grounding.

375 **Plan-conditioned multimodal retrieval.** For
376 each subsection o_i , we construct a text retrieval
377 query q_i^{text} and an image retrieval query q_i^{img} from
378 $(o_i, \mathcal{C}_{<i})$. We retrieve text evidence blocks \mathcal{B}_i from
379 the text database \mathcal{D}_t and image candidates $\mathcal{I}_i^{(0)}$
380 from the image database \mathcal{D}_v . To better align vi-
381 sual evidence with the subsection content, we ad-
382 ditionally extract caption-like strings $\mathcal{C}_i^{\text{cap}}$ from \mathcal{B}_i
383 (e.g., via lightweight pattern rules) and use them as
384 auxiliary queries for caption-based image retrieval,
385 producing $\mathcal{I}_i^{(1)}$. We then merge and deduplicate im-
386 ages to obtain \mathcal{I}_i , and pack their metadata (caption,
387 source, and brief description) as model-readable
388 context.

389 **Three-step iterative writing loop.** We struc-
390 ture subsection writing as an explicit three-step
391 loop that separates *fact extraction*, *content drafting*,
392 and *gap-aware evidence expansion*. Concretely,
393 given the assembled context $\mathcal{X}_i \leftarrow [\tilde{\mathcal{B}}_i; \tilde{\mathcal{I}}_i; \mathcal{C}_{<i}]$,
394 the model first extracts a fact list F_i grounded in \mathcal{X}_i
395 (Step 1), then drafts subsection content by organiz-
396 ing and verbalizing F_i under the plan o_i (Step 2),
397 and finally identifies missing items and generates
398 an append query q_i^{app} to retrieve additional evidence
399 (Step 3). The newly retrieved evidence is appended
400 to the context and the loop repeats, progressively re-
401 ducing omissions and weakening unsupported state-
402 ments. In our main experiments, we set the number
403 of iterations to $T=3$; we report results under differ-
404 ent iteration budgets in the Appendix. A.5.1. After
405 the loop terminates, we finalize the subsection con-
406 tent c_i using the augmented context. After all sub-
407 sections are generated, we perform a lightweight
408 refinement pass to improve style consistency, pro-

409 ducing the final multimodal report \mathcal{R} .

410 3.5 Evaluation Benchmark

411 Evaluating AI-generated long-form reports is chal- 460
412 lenging: many outputs can be acceptable for the 461
413 same prompt, and holistic human ratings are of- 462
414 ten subjective and poorly attributable to specific 463
415 system components. To enable reproducible and di- 464
416 agnostic evaluation aligned with the failure modes 465
417 we target (hallucinated facts, low factual coverage, 466
418 and missing visual evidence), we introduce **EviRe-** 467
419 **portBench**, a benchmark for *evidence-grounded* 468
420 *multimodal report generation*. 469

421 **Why SDG-style indicator reports?** We instan- 470
422 tiate EviReportBench on SDG-style analytical re- 471
423 ports not because the method is SDG-specific, but 472
424 because this domain exhibits the exact properties 473
425 that stress-test grounded report writing: (1) *evi-* 474
426 *dence intensity*: claims must be supported by het- 475
427 erogeneous sources (papers, official documents, 476
428 tables, figures); (2) *quantitative rigor*: reports con- 477
429 tain many numbers (indicator values, trends, com- 478
430 parisons) that are easy to hallucinate or mis-state; 479
431 (3) *multimodal reliance*: high-quality reports heav- 480
432 ily use figures/tables for key findings, and miss- 481
433 ing visuals is a common failure of naive genera- 482
434 tion. These characteristics directly correspond to 483
435 the core problems addressed by EVIREPORT: im- 484
436 proving factual correctness, expanding evidence 485
437 coverage, and integrating visual evidence. 486

438 **Benchmark construction.** EviReportBench is 487
439 derived from authoritative real-world reports. We 488
440 collect 30 reports spanning 8 representative top- 489
441 ics from official publications released by national 490
442 governments and international organizations. All 491
443 reports are obtained from official websites and cu- 492
444 rated by domain experts to ensure reliability and 493
445 professional quality. We evaluate generated reports 494
446 along three complementary dimensions. 495

447 **Factual Accuracy (claim verification).** Since 496
448 the primary risk in analytical report generation is 497
449 producing fluent but incorrect statements, we assess 498
450 factual accuracy at the *claim level*. Following the 499
451 spirit of FactScore-style evaluation, we extract inde- 500
452 pendent factual claims from a generated report and 501
453 verify each claim against the supporting evidence 502
454 (Min et al., 2023). We report both the number of 503
455 extracted claims and the correctness rate, yielding 504
456 an objective measure of how faithfully the report 505
457 reflects verifiable facts. This dimension is most 506
458 directly impacted by evidence-grounded retrieval 507
459 and evidence-aware drafting. 508

460 **Factual Coverage (quiz-based completeness).**

461 Accuracy alone does not capture whether the 462
463 report covers the key information expected by 464
465 expert-written references. We therefore design a 466
467 quiz-based protocol to measure coverage with re- 468
469 spect to human reports. For each topic, we con- 470
471 struct 40 quizzes using predefined templates (Ap- 472
473 pendix A.6); each quiz is answered by an LLM us- 474
475 ing *only* the generated report as context. Coverage 476
477 is measured by the *answerable ratio* R_a (whether 478
479 the report contains sufficient information to answer) 480
481 and the *accuracy* A (whether the produced answer 482
483 is correct), which are combined into an overall 484
485 score: 486

$$487 S = 0.25 \times R_a + 0.75 \times A. \quad (1) \quad 488$$

489 This dimension is designed to reflect whether the 490
491 outline and subsection plans lead to comprehensive, 492
493 evidence-supported content rather than a generic 494
495 template. 496

497 **Image Recall (visual evidence integration).**

498 To quantify whether the system retrieves and in- 499
500 tegrates informative visuals, we select the highest- 501
502 quality expert report per topic and extract its key 503
504 figures as ground-truth visual evidence. We then 504
505 extract images from generated reports and compute 505
506 image recall as the proportion of ground-truth im- 506
507 ages covered by the generated output. This metric 507
508 directly targets the “missing visuals” failure mode 508
509 and evaluates caption-guided multimodal retrieval 509
510 and integration. 510

490 4 Experiments

491 4.1 Experimental Settings

492 **Implementation.** In the retrieval stage, we em- 493
494 ploy bge-large-en-v1.5 (Xiao et al., 2023) to 494
495 encode queries and evidence units for semantic re- 495
496 trieval. For report generation, we use DeepSeek-V3 496
497 (DeepSeek-AI, 2024) and GPT-5 (OpenAI, 2025a) 497
498 as LLM agents in the planning and writing stages. 498
499 During evaluation, we adopt Qwen-Max (Alibaba 499
500 Cloud, 2025) with web search enabled as the ver- 500
501 ification model for claim checking and evidence 501
502 validation. 502

503 **Baselines.** We compare EVIREPORT with the 503
504 following representative baselines: 504

505 **Naive_RAG:** Uses the same reference corpus 505
506 as EVIREPORT and directly prompts an LLM to 506
507 generate a full report with retrieved passages. 507

508 **ReAct_RAG:** A tool-using agent that iteratively 508
509 calls text/image retrieval tools in a ReAct-style 509
510 reasoning loop (Yao et al., 2023). 510

Methods	Factual Accuracy		Factual Coverage			Image Recall
	Fact Num.	Accuracy	Ans. Ratio	Accuracy	Overall Score	
Naive_RAG	22.00	0.8890	0.1713	0.0632	0.0917	0.38
ReAct_RAG	10.38	0.8311	0.3395	0.1314	0.1834	0.45
Coze	32.75	0.3466	0.3162	0.0649	0.1277	–
OpenAI_DR	53.62	0.8850	0.2274	0.0719	0.1107	–
Gemini_DR	29.00	0.6854	0.1684	0.0316	0.0658	–
EVIREPORT	78.50	0.9687	0.6100	0.3250	0.3963	0.79

Table 1: Main results on EviReportBench. We report factual accuracy, factual coverage, and image recall. “Fact Num.” denotes the number of factual statements that are verified as correct in the generated report. “Ans. Ratio” denotes the proportion of questions that can be answered based solely on the generated report.

OpenAI_DR: Generates reports using ChatGPT Deep Research (OpenAI, 2025b).

Gemini_DR: Generates reports using Gemini Deep Research (Google AI for Developers, 2025).

Coze: Generates reports using the Coze agent platform (Coze, 2025).

Test Cases. We evaluate all methods on the 8 benchmark topics described in Sec. 3.5. For each topic, we generate a complete long-form report and evaluate it using the proposed metrics on factual accuracy, factual coverage, and visual evidence integration.

4.2 Main Results

Table 1 summarizes the main results on EviReportBench, evaluating all methods from three complementary perspectives: factual accuracy, factual coverage, and image recall.

Overall comparison. EVIREPORT consistently outperforms all baselines across all metrics by a large margin. Notably, EVIREPORT achieves both *high factual density* and *high correctness*, indicating that it improves faithfulness without sacrificing informativeness. We additionally test different backbone LLMs for generation, and the trends remain consistent; detailed results are provided in the Appendix. A.5.2.

Factual accuracy. EVIREPORT obtains the highest number of verified factual statements (**78.5**) with an accuracy of **0.9687**, substantially surpassing Naive_RAG and ReAct_RAG. This suggests that organizing evidence into compact, traceable packages and enforcing evidence-aware drafting effectively reduces unsupported statements while enabling richer factual reporting.

Factual coverage. EVIREPORT also shows a clear advantage in coverage, achieving an answerable ratio of **0.61** and an accuracy of **0.325**, resulting in the highest overall coverage score of **0.3963**.

Compared with the strongest baseline ReAct_RAG, EVIREPORT improves the overall coverage score by nearly a factor of two, indicating that hierarchical planning helps allocate space to key aspects and reduces “template-like” omissions.

Deep research baselines. OpenAI_DR and Gemini_DR achieve relatively high factual accuracy but low factual coverage on our benchmark. This suggests that general-purpose deep research systems can produce locally correct statements, yet may under-cover topic-critical facts when the task requires systematic synthesis over dense quantitative and multimodal evidence. Coze and Naive_RAG show moderate performance but suffer from either low factual accuracy or limited coverage, highlighting the limitations of direct long-form generation without structured planning and evidence packaging.

Image recall. EVIREPORT achieves the highest image recall (**0.79**), outperforming Naive_RAG and ReAct_RAG by a clear margin. This demonstrates the effectiveness of caption-guided multimodal retrieval for integrating informative visual evidence into reports. We do not report image recall for OpenAI_DR, Gemini_DR, and Coze, since these systems rely on proprietary or heterogeneous corpora and do not operate on the same controlled image database as EviReportBench.

Taken together, these results show that EVIREPORT produces reports that are more accurate, more complete, and better supported by visual evidence, validating the effectiveness of the proposed retrieve-plan-write framework.

4.3 Ablation Studies

Table 2 reports ablation results on EviReportBench, quantifying the contribution of each key component in EVIREPORT to factual accuracy, factual coverage, and visual evidence integration.

Ablation Object	Factual Accuracy		Factual Coverage			Image Recall
	Fact Num.	Accuracy	Ans. Ratio	Accuracy	Overall Score	
w/o KG	67.50	0.9053	0.3221	0.1326	0.1800	0.45
w/o Fact Extraction	52.75	0.8350	0.4574	0.1419	0.2208	0.63
w/o CGIR	76.88	0.9270	0.6003	0.2785	0.3542	0.58
EVIREPORT	78.50	0.9687	0.6100	0.3250	0.3963	0.79

Table 2: Ablation study on EviReportBench. KG denotes the evidence-oriented knowledge graph used for graph-enhanced retrieval. CGIR denotes Caption-Guided Image Retrieval.

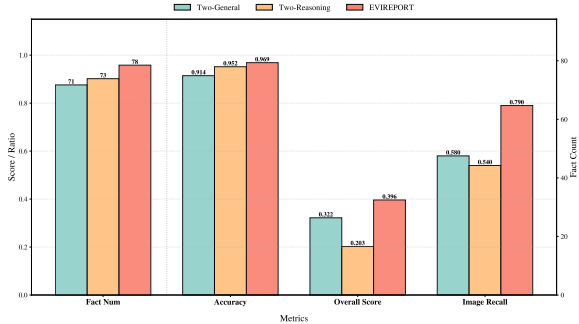


Figure 3: Ablation on two-stage outlining. "Two-General" and "Two-Reasoning" denote single-model planning variants, while EVIREPORT uses coarse-to-fine outlining.

Effect of graph-enhanced retrieval (KG). Removing the knowledge graph in the pre-planning retrieval stage (w/o KG) leads to a pronounced drop in factual coverage: the overall coverage score decreases from 0.3963 to 0.1800 and the answerable ratio drops from 0.6100 to 0.3221. This suggests that bypassing structured evidence neighborhoods weakens the system’s ability to retrieve and organize sufficient, topic-relevant evidence packages. We also observe degraded multimodal integration, as image recall drops from 0.79 to 0.45, indicating that graph-guided evidence retrieval helps surface contexts that later facilitate figure discovery and selection.

Effect of two-stage outlining. Replacing our coarse-to-fine outlining with a single model degrades performance in different ways. As shown in Fig. 3, the general-only variant yields higher coverage but lower factual accuracy, while the reasoning-only variant maintains higher factual accuracy but substantially reduces coverage. This trade-off supports our design insight that global planning and fine-grained decomposition require complementary behaviors, motivating the two-stage outlining procedure.

Effect of fact extraction. Removing the explicit

fact extraction step (Step 1) and directly drafting from retrieved context significantly hurts both reliability and completeness: factual accuracy drops from 0.9687 to 0.8350 and the number of verified facts decreases from 78.50 to 52.75, accompanied by a large reduction in coverage (overall score 0.2208 vs. 0.3963).

Effect of CGIR. Removing Caption-Guided Image Retrieval (w/o CGIR) notably reduces image recall (0.79 → 0.58) and also degrades textual quality, with factual accuracy dropping from 0.9687 to 0.9270 and the overall coverage score decreasing from 0.3963 to 0.3542. This shows that caption-guided retrieval improves figure selection and more broadly strengthens multimodal grounding for subsection drafting.

5 Conclusion and Outlook

We present EVIREPORT, an evidence-grounded workflow for automated long-form report generation that follows a retrieve–plan–write paradigm. Across all metrics on EviReportBench, EVIREPORT consistently outperforms strong baselines in factual accuracy, factual coverage, and image recall, benefiting from graph-enhanced retrieval for compact and traceable evidence packaging, coarse-to-fine outline planning, and an iterative multimodal writing loop with caption-guided image retrieval. While EviReportBench is instantiated on SDG-style, data- and figure-intensive reports, the proposed workflow is domain-agnostic and can be applied to other evidence-heavy reporting scenarios (e.g., policy analysis, scientific surveys, and industrial analytics). We hope this work provides a practical path toward more trustworthy report-generation agents and offers an evaluation blueprint for studying faithfulness, completeness, and multimodal evidence integration in long-form generation.

650 Limitations

651 This work has several limitations. First, the system
652 does not incorporate continuously updated external
653 data sources, limiting access to diverse and up-to-
654 date datasets referenced in human-written SDGs
655 reports. As a result, it cannot always compute the
656 latest indicator values, reducing the timeliness of
657 generated reports. Second, the generation process
658 relies on large language models, which may still
659 produce hallucinated or imprecise statements de-
660 spite evidence-guided retrieval, affecting factual
661 reliability in complex analytical scenarios. Future
662 work will explore integrating a multi-agent SDGs
663 data analysis platform with the report generation
664 assistant to enable coordinated data acquisition,
665 indicator computation, and verification, thereby
666 improving the depth, timeliness, and accuracy of
667 generated reports.

668 Ethical considerations

669 This work studies an automated workflow for gener-
670 ating SDGs domain reports to assist users in effi-
671 ciently organizing and summarizing evidence from
672 large collections of public sustainability documents.
673 Our framework relies on publicly available sources,
674 including official reports released by international
675 organizations and governments, as well as pub-
676 lished research papers, and we only use these mate-
677 rials for retrieval and grounded synthesis to respect
678 copyright and intellectual property constraints.

679 The system is intended to augment, not replace,
680 human expertise. Users should critically review
681 generated reports, especially quantitative claims
682 and policy recommendations, and verify them
683 against authoritative sources before downstream
684 use. As SDGs-related materials and indicators may
685 contain reporting biases, temporal gaps, or regional
686 unevenness, the generated outputs can inherit such
687 biases or present incomplete coverage. We mitigate
688 these risks by emphasizing evidence-grounded gener-
689 ation, providing traceable retrieval context, and
690 evaluating factuality and coverage; nevertheless,
691 we do not claim that outputs are error-free. We
692 encourage responsible use, transparent disclosure
693 when AI assistance is involved, and appropriate
694 human oversight for high-stakes decision making.

695 References

696 Alibaba Cloud. 2025. Alibaba cloud model studio:
697 Qwen-max model. <https://www.alibabacloud.com/product/modelstudio>. Accessed: 2025-12-20.

- 698
699
700 CBAS. International research center of big data for
701 sustainable development goals (cbas). <https://www.cbac.ac.cn/en/>. Accessed: 2025-12-22. 702
- 703 Coze. 2025. Coze space: Ai agent platform. <https://space.coze.cn/>. Accessed: 2025-12-20. 704
- 705 DeepSeek-AI. 2024. Deepseek-v3 technical report. 706
Preprint, arXiv:2412.19437.
- 707 Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu,
708 Roberta Raileanu, Xian Li, Asli Celikyilmaz, and
709 Jason Weston. 2024. Chain-of-verification reduces
710 hallucination in large language models. In *Findings*
711 *of the association for computational linguistics: ACL*
712 *2024*, pages 3563–3578.
- 713 Zican Dong, Tianyi Tang, Junyi Li, Wayne Xin Zhao,
714 and Ji-Rong Wen. 2024. Bamboo: A comprehensive
715 benchmark for evaluating long text modeling capaci-
716 ties of large language models. In *Proceedings of*
717 *the 2024 Joint International Conference on Compu-*
718 *tational Linguistics, Language Resources and Evalu-*
719 *ation (LREC-COLING 2024)*, pages 2086–2099.
- 720 Darren Edge, Ha Trinh, Newman Cheng, Joshua
721 Bradley, Alex Chao, Apurva Mody, Steven Truitt,
722 Dasha Metropolitan, Robert Osazuwa Ness, and
723 Jonathan Larson. 2024. From local to global: A
724 graph rag approach to query-focused summarization.
725 *arXiv preprint arXiv:2404.16130*.
- 726 Wenqi Fan, Yujuan Ding, Liangbo Ning, Shijie Wang,
727 Hengyun Li, Dawei Yin, Tat-Seng Chua, and Qing
728 Li. 2024. A survey on rag meeting llms: Towards
729 retrieval-augmented large language models. In *Pro-*
730 *ceedings of the 30th ACM SIGKDD conference on*
731 *knowledge discovery and data mining*, pages 6491–
732 6501.
- 733 Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony
734 Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent
735 Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and
736 Kelvin Guu. 2023a. Rarr: Researching and revising
737 what language models say, using language models.
738 In *Proceedings of the 61st Annual Meeting of the*
739 *Association for Computational Linguistics (Volume*
740 *1: Long Papers)*, pages 16477–16508.
- 741 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang
742 Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun,
743 Haofen Wang, and Haofen Wang. 2023b. Retrieval-
744 augmented generation for large language models: A
745 survey. *arXiv preprint arXiv:2312.10997*, 2(1).
- 746 Zengyi Gao, Yukun Cao, Hairu Wang, Ao Ke, Yuan
747 Feng, S. Kevin Zhou, and Xike Xie. 2025. Frag: A
748 flexible modular framework for retrieval-augmented
749 generation based on knowledge graphs. In *Find-*
750 *ings of the Association for Computational Linguis-*
751 *tics: ACL 2025*, pages 6178–6192.

752	Google AI for Developers. 2025. Gemini deep research agent. https://ai.google.dev/gemini-api/docs/deep-research . Accessed: 2025-12-20.	809
753		810
754		
755	Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan, Chen Ling, and Liang Zhao. 2024. Grag: Graph retrieval-augmented generation. <i>arXiv preprint arXiv:2405.16506</i> .	811
756		812
757		813
758		
759	Haibo Jin, Haoxuan Che, Yi Lin, and Hao Chen. 2024. Promptrg: Diagnosis-driven prompts for medical report generation. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 2607–2615.	814
760		815
761		816
762		817
763		
764	Ishita Kumar, Snigdha Viswanathan, Sushrita Yerra, Alireza Salemi, Ryan A Rossi, Franck Dernoncourt, Hanieh Deilamsalehy, Xiang Chen, Ruiyi Zhang, Shubham Agarwal, and 1 others. 2024. Longlamp: A benchmark for personalized long-form text generation. <i>arXiv preprint arXiv:2407.11016</i> .	818
765		819
766		820
767		821
768		822
769		
770	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in Neural Information Processing Systems</i> , 33.	823
771		824
772		825
773		826
774		827
775		828
776		829
777		830
778	Mufe i Li, Siqi Miao, and Pan Li. 2025. Simple is effective: The roles of graphs and large language models in knowledge-graph-based retrieval-augmented generation. In <i>International Conference on Learning Representations (ICLR)</i> .	831
779		832
780		833
781		834
782		835
783		836
784		
785	Lei Liang, Mengshu Sun, Zhengke Gui, Zhongshu Zhu, Zhouyu Jiang, Ling Zhong, Yuan Qu, Peilong Zhao, Zhongpu Bo, Jin Yang, and 1 others. 2024. Kag: Boosting llms in professional domains via knowledge augmented generation. <i>arXiv preprint arXiv:2409.13731</i> .	837
786		838
787		839
788		840
789		841
790		842
791		843
792		844
793		845
794		846
795		847
796		
797		848
798		849
799		850
800		
801		851
802		852
803		853
804		854
805		855
806		856
807		
808		857
		858
		859
		860
		861
		862
		863
		864

865 Yidong Wang, Qi Guo, Wenjin Yao, Hongbo Zhang,
866 Xin Zhang, Zhen Wu, Meishan Zhang, Xinyu Dai,
867 Qingsong Wen, Wei Ye, and 1 others. 2024. Autosur-
868 vey: Large language models can automatically write
869 surveys. *Advances in neural information processing*
870 *systems*, 37:115119–115145.

871 Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas
872 Muennighoff. 2023. [C-pack: Packaged resources](#)
873 [to advance general chinese embedding](#). *Preprint*,
874 arXiv:2309.07597.

875 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak
876 Shafran, Karthik R. Narasimhan, and Yuan Cao. 2023.
877 [React: Synergizing reasoning and acting in language](#)
878 [models](#). In *International Conference on Learning*
879 *Representations (ICLR)*.

880 Ruo Chen Zhao, Xingxuan Li, Shafiq Joty, Chengwei
881 Qin, and Lidong Bing. 2023. [Verify-and-edit: A](#)
882 [knowledge-enhanced chain-of-thought framework](#).
883 In *Proceedings of the 61st Annual Meeting of the*
884 *Association for Computational Linguistics (Volume*
885 *1: Long Papers)*, pages 5823–5840.

886 A Appendix

887 Due to space limitations in the main paper, we
888 provide additional details in the appendix:

- 889 • Sec. [A.1](#) Algorithms (pseudocode).
- 890 • Sec. [A.2](#) Database Construction.
- 891 • Sec. [A.3](#) SDGs Domain Knowledge Graph
892 Construction for Report Generation.
- 893 • Sec. [A.4](#) Details of Evaluation Data.
- 894 • Sec. [A.5](#) Additional Results (hyperparameters,
895 model variants, per-topic analysis).
- 896 • Sec. [A.6](#) Quiz Construction for Factual Cov-
897 erage Evaluation.
- 898 • Sec. [A.7](#) Prompts Used.

899 A.1 Algorithms

900 This appendix provides implementation-level pseu-
901 docode for the key modules of EVIREPORT, com-
902 plementing the high-level descriptions in Sec. 3.2–
903 Sec. 3.4. Algorithm 1 details graph-enhanced evi-
904 dence retrieval that produces compact, traceable
905 evidence packages for planning and writing. Al-
906 gorithm 2 specifies the subsection drafting proce-
907 dure with the three-step iterative loop and caption-
908 guided multimodal retrieval.

Algorithm 1 Graph-Enhanced RAG Retrieval (KG- RAG)

- 1: **Input:** query q , knowledge graph \mathcal{G} , text
database \mathcal{D} , top- K nodes K , hop size h
 - 2: **Output:** evidence pairs $\mathcal{P} = \{(\mathcal{G}_i, s_i)\}$
 - 3: Extract keywords \mathcal{K} from q and compute query
embedding \mathbf{e}_q
 - 4: Retrieve candidate *case* nodes from \mathcal{G} using \mathcal{K}
and \mathbf{e}_q
 - 5: Merge and deduplicate retrieved nodes to ob-
tain \mathcal{E} ; keep top- K by relevance
 - 6: $\mathcal{P} \leftarrow \emptyset$
 - 7: **for** each case node $c \in \mathcal{E}$ **do**
 - 8: Extract an entity-centered subgraph \mathcal{G}_c
from \mathcal{G} within h hops of c
 - 9: Prune \mathcal{G}_c by removing low-relevance
nodes/edges w.r.t. q to obtain \mathcal{G}'_c
 - 10: Use \mathcal{G}'_c to form retrieval constraints (enti-
ties/reasons) and retrieve text blocks \mathcal{B}_c from
 \mathcal{D}
 - 11: Summarize \mathcal{B}_c into s_c (key facts, numbers,
and missing information) conditioned on q
 - 12: Augment \mathcal{G}'_c with summary-derived at-
tributes to obtain \mathcal{G}''_c
 - 13: $\mathcal{P} \leftarrow \mathcal{P} \cup \{(\mathcal{G}''_c, s_c)\}$
 - 14: **end for**
 - 15: **Return:** \mathcal{P}
-

A.2 Database Construction

To ensure the quality and professionalism of AI-
generated reports, we build both a text retrieval
database and an image retrieval database.

Data collection. We collect (i) 109,731 SDGs-
related research papers, (ii) 14 sustainability re-
ports (2019–2025) released by the International
Research Center of Big Data for Sustainable De-
velopment Goals (CBAS)(CBAS), including China
and international SDGs reports, (iii) three sustain-
ability reports (2020–2022) for the “Belt and Road”
initiative, and (iv) 687 documents published by the
United Nations, other international organizations,
and national governments.

Document parsing and text indexing. The col-
lected documents exhibit diverse and complex for-
mats. After preliminary investigation, we adopt a
hybrid pipeline combining MinerU2.5(Niu et al.,
2025) and a vision-language model (VLM) to con-
vert all sources into unified Markdown (MD) files.
Specifically, we first apply MinerU2.5 to extract
text, tables, and figures, and then verify the con-

Algorithm 2 Content generation with a three-step loop

```

1: Input: outline  $\mathcal{O} = \{o_1, \dots, o_N\}$ ; text DB  $\mathcal{D}_t$ ; image DB  $\mathcal{D}_v$ ; max iterations  $T$ 
2: Output: final report  $\mathcal{R}$ 
   Initialize generated subsection contents
3:  $\mathcal{C} \leftarrow \emptyset$ 
4: for  $i = 1$  to  $N$  do
5:   Construct retrieval queries  $q_i^{\text{text}}$  and  $q_i^{\text{img}}$ 
   from  $(o_i, \mathcal{C}_{<i})$ 
6:   Retrieve text blocks
7:    $\mathcal{B}_i \leftarrow \text{SEARCH}(\mathcal{D}_t, q_i^{\text{text}})$ 
8:   Retrieve images
9:    $\mathcal{I}_i^{(0)} \leftarrow \text{SEARCH}(\mathcal{D}_v, q_i^{\text{img}})$ 
10:  Extract caption-like strings from text
11:   $\mathcal{C}_i^{\text{cap}} \leftarrow \text{EXTRACTCAPTIONS}(\mathcal{B}_i)$ 
12:  Caption-based image retrieval
13:   $\mathcal{I}_i^{(1)} \leftarrow \text{SEARCH}(\mathcal{D}_v, \mathcal{C}_i^{\text{cap}})$ 
14:   $\mathcal{I}_i \leftarrow \text{DEDUP}(\mathcal{I}_i^{(0)} \cup \mathcal{I}_i^{(1)})$ 
15:   $\tilde{\mathcal{B}}_i \leftarrow \text{CLEAN}(\mathcal{B}_i)$ 
16:   $\tilde{\mathcal{I}}_i \leftarrow \text{PACKMETA}(\mathcal{I}_i)$ 
17:  Assemble context  $\mathcal{X}_i \leftarrow [\tilde{\mathcal{B}}_i; \tilde{\mathcal{I}}_i; \mathcal{C}_{<i}]$ 
18:   $d_i \leftarrow \emptyset$ ;  $m_i \leftarrow \emptyset$ 
19:  for  $t = 1$  to  $T$  do
20:    Step 1 (Facts):
21:     $F_i \leftarrow \text{LLM\_FACTS}(\mathcal{X}_i, o_i)$ 
22:    Step 2 (Draft):
23:     $d_i \leftarrow \text{LLM\_WRITE}(F_i, \mathcal{X}_i, o_i)$ 
24:    Step 3 (Append):
25:     $(m_i, q_i^{\text{app}}) \leftarrow \text{LLM\_APPENDQUERY}(d_i, F_i, o_i)$ 
26:    if  $m_i = \emptyset$  then
27:      break
28:    end if
29:    Retrieve additional evidence
30:     $\Delta\mathcal{B}_i \leftarrow \text{SEARCH}(\mathcal{D}_t, q_i^{\text{app}})$ 
31:     $\mathcal{X}_i \leftarrow [\mathcal{X}_i; \text{CLEAN}(\Delta\mathcal{B}_i)]$ 
32:  end for
33:   $c_i \leftarrow \text{FINALIZE}(d_i, \mathcal{X}_i)$ 
34:   $\mathcal{C} \leftarrow \mathcal{C} \cup \{c_i\}$ 
35: end for
36:  $\mathcal{R} \leftarrow \text{REFINE}(\mathcal{C})$ 
37: Return:  $\mathcal{R}$ 

```

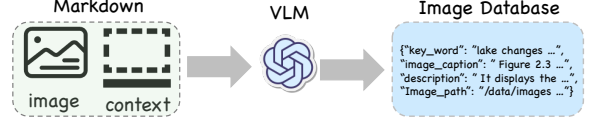


Figure 4: Pipeline for image database construction. We extract each image together with its caption and surrounding context from the parsed Markdown documents, and use a vision–language model (VLM) to generate keywords and a semantic summary for indexing.

verted content, focusing on error-prone elements such as tables and images. When conversion errors are detected, we invoke the VLM for format correction. This hybrid design balances efficiency and accuracy: the VLM provides higher fidelity but is constrained by token budget and latency, so we only use it for difficult cases.

We then split each MD document into chapter-aware chunks, resulting in approximately two million text blocks. For each block, we prepend the document title and its ancestor section titles, preserving hierarchical structure for downstream retrieval. Next, we encode the text blocks using the bge-large-en-v1.5 embedding model and index them in Milvus for dense retrieval. In addition, we index the raw text in Elasticsearch to enable hybrid sparse–dense retrieval.

Image database construction. As illustrated in Fig. 4, we extract each image together with its caption and surrounding context paragraphs from the MD files. We feed the image and its textual context into a VLM to generate a structured image summary and keywords. Finally, we store the image caption, VLM-generated summary, and image path in an image database to support caption-guided image retrieval.

A.3 Knowledge Graph Construction

This section details how we construct the SDGs domain knowledge graph (SDG-KG) used in EVIREPORT for retrieval and evidence organization during report generation. SDG-KG is designed to explicitly model the evidence chain behind SDGs reporting, i.e., linking official indicators to real-world cases and further grounding each case with datasets and analytical methods.

A.3.1 Base SDGs schema from the official framework

We start from the official SDGs global indicator framework, which provides a stable hierarchical

970 mapping from *Goals* to *Indicators*(United Nations
971 [Statistics Division; United Nations General Assem-](#)
972 [bly, 2017\)](#). We parse the framework documents
973 and build a canonical hierarchy consisting of:

- 974 • **Goal** nodes (17 Goals), each associated with
975 goal titles and descriptions;
- 976 • **Indicator** nodes, associated with indicator
977 codes (e.g., 11.3.1), indicator names, and
978 textual definitions;
- 979 • **Goal-Indicator** edges representing the offi-
980 cial hierarchy.

981 This base schema forms the backbone of SDG-KG
982 and serves as an anchor for integrating evidence
983 extracted from external references.

984 **A.3.2 Evidence-oriented extension:** 985 **Indicator-Case-Dataset-Method**

986 Human-written SDGs analyses and policy reports
987 typically discuss indicator progress through con-
988 crete *cases* (e.g., a country, city, region, or project),
989 supported by specific *datasets* and *methods* (e.g.,
990 statistical models, remote sensing pipelines, ma-
991 chine learning algorithms). To reflect this writing
992 pattern, we extend the base SDG hierarchy into an
993 evidence-oriented knowledge structure with four
994 key entity types:

- 995 • **Case**: a concrete analysis scenario (e.g., coun-
996 try/region/city, sector, or program) used to
997 assess one or more indicators;
- 998 • **Dataset**: the data source(s) used to measure
999 or estimate the indicator(s) for a given case;
- 1000 • **Method**: the analytical or computational ap-
1001 proach applied to the dataset(s) (e.g., regres-
1002 sion, index construction, causal estimation,
1003 satellite-based mapping, ML models);
- 1004 • **Indicator**: the official SDG indicator being
1005 assessed.

1006 We then model their relationships as:

- 1007 • **RELATED_TO** edges between **Indicator**
1008 and **Case**, capturing that a case study reports
1009 or evaluates the indicator;
- 1010 • **USE_DATASET** edges between **Case** and
1011 **Dataset**;
- 1012 • **USE_METHOD** edges between **Case** and
1013 **Method**.

As a result, each indicator is associated with a col- 1014
lection of evidence subgraphs, enabling the system 1015
to retrieve not only definition-level information, 1016
but also grounded analytical evidence used in real 1017
SDGs reporting. 1018

1019 **A.3.3 Entity and relation extraction from** 1020 **references**

To populate SDG-KG beyond the official frame- 1021
work, we automatically extract entities and rela- 1022
tions from the collected SDGs-related papers and 1023
reports. Specifically, given a parsed document in 1024
Markdown, we perform: 1025

- 1026 1. **Candidate entity mining.** We use an LLM-
1027 based extractor to identify mentions of cases,
1028 datasets, and methods in paragraphs that dis-
1029 cuss SDGs indicators. The extractor produces
1030 typed candidates (Case/Dataset/Method) to-
1031 gether with supporting spans.
- 1032 2. **Indicator linking.** We link extracted evi-
1033 dence to indicator nodes by matching indi-
1034 cator codes/names and by semantic similarity
1035 between the evidence context and indicator
1036 descriptions. When multiple indicators are
1037 plausible, we keep top candidates and resolve
1038 them in the fusion step (Sec. A.3.4).
- 1039 3. **Relation construction.** We create **Indicator-**
1040 **Case**, **Case-Dataset**, and **Case-Method**
1041 edges when the corresponding entities co-
1042 occur within a coherent local context (e.g.,
1043 within the same section or a bounded win-
1044 dows of paragraphs), and the LLM extractor
1045 provides a positive relation decision with min-
1046 imal justification.

1047 **A.3.4 Entity fusion, deduplication, and** 1048 **incremental updates**

Since extracted entities may contain aliases and 1049
surface-form variations (e.g., “LCRPG” vs. “Land 1050
Consumption Rate to Population Growth Rate”), 1051
we perform entity fusion and deduplication before 1052
insertion. We combine: 1053

- 1054 • **String-level normalization** (case folding,
1055 punctuation removal, acronym expansion
1056 when available);
- 1057 • **Embedding-based similarity** to detect near-
1058 duplicate entities;
- 1059 • **Type-aware constraints** to prevent merging
1060 across incompatible types (e.g., Method vs.
1061 Dataset).

After fusion, we incrementally insert new nodes/edges into SDG-KG while preserving the canonical Goal–Indicator backbone.

A.3.5 Example: an evidence subgraph

Fig 5 illustrates an example evidence subgraph in SDG-KG (English version). In this example, the **Goal** node SDG 11.3 is linked to its official **Indicator** node SDG 11.3.1 via the HAS_INDICATOR relation. We further connect a concrete **Case** node (e.g., “urban land-use efficiency in China”) to both the goal and the indicator using RELATED_TO, reflecting that a case study is typically written to assess progress under a specific goal/indicator.

Crucially, the case is grounded by explicit **Dataset** and **Method** nodes. As shown, the case is associated with multiple datasets (e.g., land-cover products, national geographic conditions monitoring data, and statistical yearbooks) through USE_DATASET. Meanwhile, the analytical process is represented by method nodes (e.g., multi-scale semantic segmentation, machine-learning algorithms, and built-up area extraction pipelines) connected via USE_METHOD. In addition, metric definitions used for indicator computation (e.g., LCRPGR and BPC metrics) are also modeled as method/measurement nodes and linked to the case. This evidence-oriented structure enables SDG-KG to retrieve not only indicator definitions, but also representative cases and their supporting datasets/methods, thereby improving evidence coverage for downstream report generation.

A.4 Details of Evaluation Data

We select eight representative SDGs report cases from authoritative sustainability reports and related technical assessments, as summarized in Table 3. These cases span diverse SDG themes (e.g., food production, water environment, building electrification, urban development, climate extremes) and cover multiple evidence types, including indicator definitions, datasets, analytical methods, quantitative results, and key figures (maps, trend curves, and spatial comparisons). Such diversity enables a comprehensive evaluation of long-form report generation in terms of *factual accuracy*, *factual coverage*, and *multimodal evidence integration*.

How the cases are used. In our experiments, each case topic is used as an input query for report generation, ensuring that all compared methods produce reports under the same thematic require-

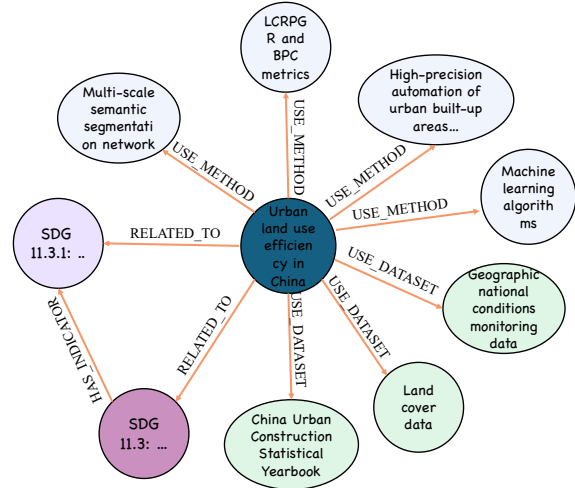


Figure 5: An example evidence subgraph in SDG-KG. The goal SDG 11.3 links to its indicator SDG 11.3.1 via HAS_INDICATOR. A case node is connected to the goal/indicator via RELATED_TO, and further grounded by dataset nodes (USE_DATASET) and method/measurement nodes (USE_METHOD).

ments. The corresponding human-written reference reports serve two purposes: (i) they provide source material for constructing the quiz-based factual coverage benchmark (Sec. A.6); and (ii) for image recall evaluation, we extract key figures from the highest-quality reference report for each topic as ground-truth visual evidence, and then compute the fraction covered by figures included in generated reports.

A.5 Additional Results

A.5.1 Iteration Budget for the Three-step Loop

We study the iteration budget T of the three-step writing loop (facts \rightarrow draft \rightarrow append query) by varying the number of loop iterations from $T=1$ to $T=5$ while keeping all other settings fixed (same outline, retrieval setup, and evaluation protocol). The results are summarized in Table 4.

Overall, increasing T substantially improves report quality in the low-iteration regime. Moving from $T=1$ to $T=3$ consistently increases factual density and completeness: the number of verified facts grows from 61.25 to 78.50, factual accuracy improves from 0.9316 to 0.9687, the answerable ratio rises from 0.4781 to 0.6010, and image recall increases from 0.54 to 0.79. This indicates that iterative gap-aware retrieval effectively recovers missing evidence and reduces weakly supported statements when the draft is still under-covered.

Topic	Reference report	Example figure
Goal 2 Zero Hunger	Global Farmland Change Assessment (1985–2020)	Statistical analysis of changes in total cultivated land area across six continents (1985–2020)
Goal 2 Zero Hunger	Global High-Resolution Remote Sensing Monitoring and Assessment of Multiple Crop Reclamation in Farmland (2020)	Proportion of replanted arable land in each continent relative to global replanted area
Goal 2 Zero Hunger	Identification of stable and vulnerable food production zones in Africa	Distribution patterns of maize, wheat, rice, and soybean yields in Africa
Goal 6 Clean Water and Sanitation	Spatiotemporal distribution of algal blooms in typical lakes worldwide	Annual variation rate of water transparency across regions and climate zones (2000–2021)
Goal 7 Affordable and Clean Energy	Global status of building electrification	Global unelectrified building area percentage (2020)
Goal 11 Sustainable Cities and Communities	Global nighttime light intensity changes	“NPP-VIIRS-like” nighttime light brightness distribution map
Goal 13 Climate Action	Analysis of global heat wave disaster changes and their impacts	Average annual probability of a global heat wave
Goal 2 Zero Hunger	spatiotemporal dynamics of rice cultivation patterns in the Indochina Peninsula	Spatial distribution of rice planting pattern changes in the Indochina Peninsula (2000–2019)

Table 3: Evaluation cases (topics) used for benchmarking SDGs domain report generation. Each case is associated with an authoritative reference report and representative figures used for quiz construction and image recall evaluation.

Loop Num.	Factual Accuracy		Factual Coverage			Image Recall
	Fact Num.	Accuracy	Ans. Ratio	Accuracy	Overall Score	
1	61.25	0.9316	0.4781	0.3137	0.3548	0.54
2	69.25	0.9568	0.5250	0.3214	0.3723	0.71
3	78.50	0.9687	0.6010	0.3250	0.3963	0.79
4	78.30	0.9462	0.5813	0.3065	0.3752	0.79
5	82.49	0.9374	0.5750	0.3207	0.3843	0.83

Table 4: Effect of iteration budget T for the three-step writing loop on EviReportBench.

Beyond $T=3$, we observe diminishing returns and a mild trade-off between metrics. While $T=5$ further increases the number of verified facts (82.49) and yields the best image recall (0.83), it slightly reduces factual accuracy (0.9374) and lowers the overall coverage score (0.3843) compared to $T=3$ (0.3963). Since each additional iteration introduces extra model calls and retrieval steps, we use $T=3$ as the default in the main experiments, as it provides the best balance between factual correctness, coverage, and computational cost.

A.5.2 Backbone LLM Variants

We evaluate the robustness of EVIREPORT to different backbone LLMs by replacing the generation model used in the planning and writing stages while keeping the retrieval setup, iteration budget ($T=3$),

and evaluation protocol unchanged. Table 5 reports the results.

Overall, EVIREPORT is relatively stable across strong backbones: Gemini-2.0-flash, Claude-sonnet-4, and GPT-5.0 all achieve high factual accuracy (≥ 0.965) and comparable factual coverage. GPT-5.0 yields the best overall balance, achieving the highest factual accuracy (0.9687) and the best coverage overall score (0.3963). Gemini-2.0-flash achieves the highest answerable ratio (0.6125) and the best image recall (0.83), suggesting stronger multimodal retrieval/selection synergy under this backbone, while its overall coverage score is slightly below GPT-5.0. Claude-sonnet-4 performs competitively but is consistently lower on coverage metrics. DeepSeek-v3.1 shows a larger

LLM	Factual Accuracy		Factual Coverage			Image Recall
	Fact Num.	Accuracy	Ans. Ratio	Accuracy	Overall Score	
Gemini-2.0-flash	77.25	0.9656	0.6125	0.3188	0.3922	0.83
Claude-sonnet-4	75.50	0.9664	0.5813	0.3063	0.3751	0.79
GPT-5.0	78.50	0.9687	0.6010	0.3250	0.3963	0.79
DeepSeek-v3.1	70.25	0.9574	0.5594	0.2813	0.3508	0.67

Table 5: Backbone LLM variants for EVIREPORT on EviReportBench (iteration budget $T=3$).

drop in coverage and image recall, indicating that the downstream writing loop benefits from stronger instruction-following and gap-aware querying capability. Based on these results, we use GPT-5.0 as the default backbone in the main experiments.

A.5.3 Per-topic Performance and Radar Plots

To provide a fine-grained view beyond averaged scores, we report per-topic results on the 8 benchmark topics using radar plots (Fig. 6). We visualize four complementary metrics: (a) **Fact Num.** (verified factual statements), (b) **Factual Accuracy** (claim correctness rate), (c) **Overall Score** for factual coverage (combining answerable ratio and accuracy), and (d) **Image Recall** (coverage of topic-specific ground-truth figures). Each axis corresponds to one topic, allowing direct inspection of where a method performs strongly or fails.

Overall, EVIREPORT shows consistently high performance across topics on factual accuracy and coverage-related metrics, indicating that the retrieve-plan-write workflow generalizes beyond a small subset of “easy” topics. We also observe that some baselines exhibit larger topic variance: while they may perform competitively on a few topics, they often drop sharply on others, suggesting sensitivity to evidence density and topic structure. Finally, the image-recall plot highlights that multi-modal performance is highly topic-dependent and that caption-guided image retrieval improves figure coverage particularly on topics where key findings are primarily communicated through maps or trend plots.

A.6 Quiz Construction

Evaluating the *factual coverage* of long-form SDGs reports is challenging because (i) gold answers are often non-unique, and (ii) free-form judging can be subjective. Following the quiz-driven, reader-aligned evaluation philosophy (i.e., a report is useful if it can answer readers’ concrete questions), we construct an *evidence-constrained* quiz set and

evaluate whether a generated report can answer these quizzes faithfully.

Quiz perspectives. We design six SDGs-specific quiz perspectives that target typical information needs in SDGs reporting, including indicator definitions, datasets, methods, indicator values, comparisons, and limitations. Table 6 summarizes the perspectives, descriptions, and examples.

Source reports and paragraph selection. For each topic, we start from an expert-selected human-written report (officially released on authoritative websites). We segment the report into paragraphs (by section boundaries) and retain paragraphs that are information-rich and self-contained. Inspired by prior quiz construction practice, we filter out paragraphs that are too short or lack concrete signals (e.g., missing numeric values, methods, datasets, or explicit claims).

Quiz generation with evidence grounding. Given a retained paragraph p , we prompt an LLM to generate multiple quiz pairs $\{(q_j, a_j)\}_{j=1}^k$ under the perspective templates in Table 6, with the **strict constraint** that each answer a_j must be *fully supported* by p (no external knowledge) and should be *specific* (preferably including explicit entities, numbers, datasets, methods, time ranges, or uncertainty statements). This design reduces the ambiguity of open-ended judging and aligns the evaluation with verifiable evidence.

Multi-criteria validation and de-duplication. Each generated (q, a) pair is validated to ensure quality and faithfulness: (1) **format & length**: question/answer exceed minimal length and avoid vague wording; (2) **specificity**: answers contain concrete indicators (e.g., numeric values, named datasets/methods, time ranges); (3) **grounding**: keyword/phrase overlap and consistency checks between (q, a) and paragraph p ; (4) **de-duplication**: remove near-duplicate quizzes within

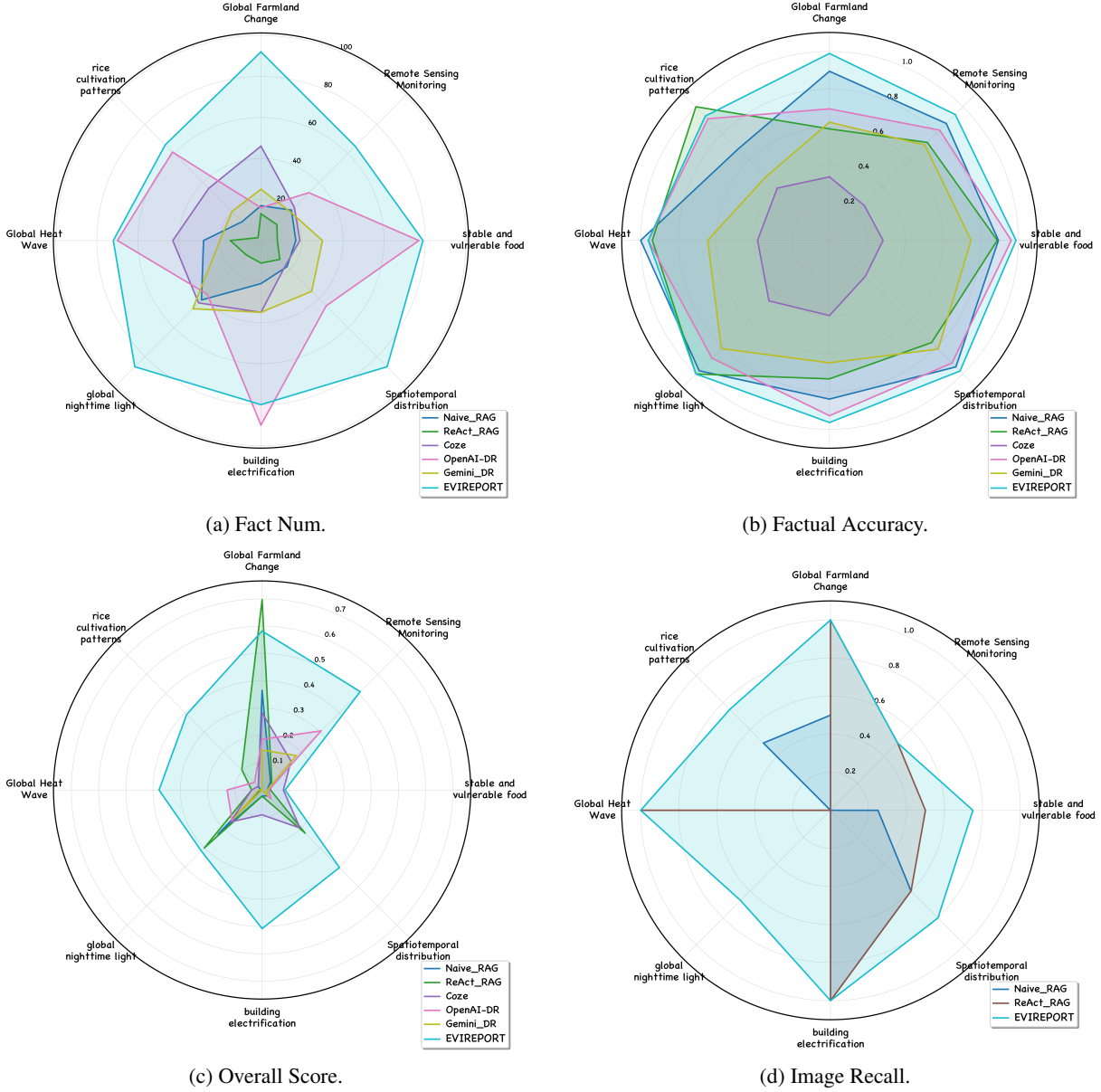


Figure 6: Per-topic radar plots on EviReportBench. (a) Fact Num., (b) Factual Accuracy, (c) Overall Score, and (d) Image Recall.

a topic/perspective. We iterate generation and validation until reaching a target number of quizzes per topic (e.g., 40).

Using quizzes for factual coverage evaluation.

Given a generated report \hat{R} , we ask an evaluation LLM to answer each quiz *using only* \hat{R} (and the retrieved supporting context if applicable), and to return a fixed fallback response such as No relevant content found in the report when the report does not contain sufficient information. We then compute: (i) **AnswerableRatio**: the fraction of quizzes that are answerable from \hat{R} , and (ii) **Accuracy**: the correctness of the produced answers (optionally with evidence checks), which together

reflect whether the report *covers* the key factual points expected from expert-written SDGs reports.

A.7 Prompts Used

This section summarizes the main prompts used in EVIREPORT for (i) outline generation, (ii) content generation, and (iii) evaluation. Full prompt templates are provided after this subsection.

Outline generation. We condition outlining on retrieved evidence (KG subgraphs and summarized text blocks) and structural priors from high-quality SDGs reports. Following our two-stage strategy, a reasoning-oriented LLM first produces a high-level plan that covers essential SDGs perspectives (e.g.,

Perspective	Description	Example
Definition	Examine whether the report clearly defines SDG indicators, terminologies, measurement units, and versions used in data reporting.	What is the unit of measurement for indicator 11.3.1 on land use efficiency in China?
Dataset	Verify that dataset sources, temporal coverage, and spatial resolution are accurately stated and traceable to authoritative repositories.	Which dataset and time period are used to estimate the urban land use efficiency index?
Method	Assess whether the report describes analytical methods, models, equations, and validation metrics in sufficient detail to ensure reproducibility.	What model or formula is used to calculate the SDG indicator value?
Indicator Value	Check whether the report presents specific numerical values or rates for defined countries, years, or regions.	What was the reported land use efficiency value for Beijing in 2020?
Comparison	Evaluate whether the report analyzes temporal changes or spatial differences among SDG indicators, highlighting significant trends.	How did urban land use efficiency change between 2000 and 2020 across major cities?
Limitation	Identify acknowledgments of data gaps, uncertainties, biases, or limitations in indicator interpretation and model generalization.	What uncertainties or data limitations are mentioned in the SDG 11.3.1 analysis?

Table 6: Perspectives and examples used for evaluating factual coverage in SDGs domain reports.

1278 indicator definition, data/methods, findings, trends,
1279 limitations, and implications). A chat-oriented
1280 LLM then refines it into a detailed hierarchical out-
1281 line with non-overlapping subsection scopes and
1282 explicit writing goals, ensuring coherence and re-
1283 ducing redundancy.

1284 **Content generation.** For each subsection, we
1285 retrieve subsection-specific text and image evi-
1286 dence and generate content sequentially to maintain
1287 global consistency. A CoT-style prompt guides the
1288 model to (1) extract factual points from the pro-
1289 vided evidence, (2) draft grounded content, and
1290 (3) list missing/uncertain information. We further
1291 apply an append-and-rewrite prompt that issues a
1292 follow-up query, incorporates additional evidence,
1293 and rewrites the draft to improve completeness and
1294 factual consistency.

1295 **Evaluation.** We use prompts for factual verifica-
1296 tion and quiz-based factual coverage. The verifier
1297 extracts atomic claims and checks support via evi-
1298 dence retrieval (with web search when enabled).
1299 For coverage, the evaluator answers each quiz using
1300 only the generated report, enabling AnswerableRa-
1301 tio and Accuracy computation.



Generation Prompt

Outline stage-1:

OUTLINE_STAGE_1_GENERATE_PROMPT= '''

You are a professional assistant for generating article outlines. Based on the user's query and the provided knowledge points, please generate an article outline.

I have constructed a knowledge graph based on existing information, which includes data sources, computational methods used in relevant SDG case studies, and the specific SDG goals they correspond to (e.g., SDG 11.1). I will provide you with retrieved information from the knowledge graph as well as reference texts.

The topic is: {query}. Please help generate chapter titles and outlines accordingly. In each chapter outline, make sure to **retain any image captions** from the provided reference content.

Please follow the chapter structure below:

1. Target
2. Background
3. Data Used
4. Research Methods (retain relevant image captions when describing the methods)
5. Results and Analysis (retain relevant image captions in the analysis to support conclusions)
6. Key Contributions
7. Discussion and Outlook

Knowledge Graph Retrieval Information:

{kg_info}

'''

Outline stage-2:

OUTLINE_STAGE_2_GENERATE_PROMPT= '''

Please generate a well-structured chapter-by-chapter writing plan in **JSON format** based on the report topic below.

You will be provided with **entity nodes and relationship information** retrieved from a knowledge graph, as well as **reference context** retrieved from the knowledge base.

Additionally, you may refer to the initial version of the outline generated by another LLM.

Requirements:

1. **Strictly output in JSON format**, avoiding extra line breaks, explanatory texts, or formatting errors.
2. The **chapter structure should be clear**, including:
 - `chapter_title`: The chapter title (must match the provided title exactly).
 - `content_structure`: A structural guide for the chapter (e.g., whether it has sub-sections, brief or detailed content, etc.).
 - `content_points`: A detailed outline of the key points that must be included in this chapter; avoid vague or generic descriptions.
3. **Chapter titles must match the provided ones exactly. Do not alter them. The structure must include the following chapters: Target, Background, Data Used, Methodology, Results and Analysis, Key Contributions, Discussion and Outlook.**

Output Example:

```json

```
[
 {{
 "chapter_title": "Target",
 "content_structure": "As the first section of the article, this should generate the article title. It
 should correspond to the specific SDG target. No sub-sections (i.e., no ### headers); provide a brief
 description.",
 "content_points": ...
 }},
 {{
 "chapter_title": "Background",
 "content_structure": "No sub-sections (i.e., no ### headers); provide a brief description in 1-2
 paragraphs.",
 "content_points": ...
 }}
 ...
 }}
]
```

Please strictly follow the JSON format. Ensure `chapter\_title` matches exactly with the provided titles. `content\_structure` should summarize the content framework for each chapter, and `content\_points` should include quantitative data or methodological details where possible—avoid vague or generic language.

Report Topic: {}

Reference Context: {}

Reference Outline (Generated by another LLM for reference):{}  
 ...

## Content:

CONTENT\_GENERATE\_PROMPT = '''

You are a professional research assistant specializing in **Retrieval-Augmented Generation (RAG)** for scientific writing.

Use **ONLY** the provided RAG knowledge base (annual SDG reports, peer-reviewed literature, books) to ensure **factual basis, authority, and contextual consistency**.

**User Input (Query):**

{query}

---

### ### 📌 Task Requirements

- 1) Write **in English**, using academic style that meets research standards.
- 2) No fixed length; prioritize clarity, depth, tight logic; avoid redundancy and fluff.
- 3) **Obey the chapter structure** described by `content\_structure`; keep hierarchy concise.
- 4) **Exploit retrieved evidence** (data, methods, results, key findings). Prefer quantitative details (numbers, units, dates, spatial/temporal resolution,  $R^2$ /accuracy) over generic descriptions.
- 5) **Images**: insert only when valid absolute local paths are provided in retrieval results; embed as `[caption](ABS\_PATH)`. Do not describe images without a path.
- 6) **Avoid generalizations**. Every claim must be backed by data, methods, or specific cases from the references.
- 7) If a computation has a formula, **list the formula and variables clearly** (plain text or LaTeX inline), and define symbols.
- 8) **Expand `content\_points`** and organize strictly per `content\_structure`.
- 9) Use the given **chapter\_title** verbatim; keep sub-headings minimal.
- 10) **Formatting**:
  - `# Title` only for the **first** chapter of the whole report
  - `## Chapter Title` for main chapters
  - `### Sub-chapter Title` optional
  - Body in plain text (you may include inline equations if needed)

---

### ### 🔍 Evidence & Citation Policy (Hard Constraints)

- Write **only** from the "Reference Context" below. **Do not fabricate** sources, numbers, or quotes.
- **Each paragraph must contain at least one citation** in the format `[X]`, where `X` is the numbered block id.
- If a sentence relies on multiple references, append all: e.g., `[3][7][15]`.
- If no evidence supports a required point, do **not invent it**—state an **Evidence Gap Note** at the end instead.

---

### ### 🍪 Generation Procedure (You MUST follow)

#### **Step 1 – Fact Extraction (bullet list):**

Extract **at least 10 concrete, verifiable facts** from the Reference Context that directly support this chapter.

A **fact** refers to a **specific, evidence-based statement** that can be verified from the provided references.

Facts should represent **objective, quantitative, or methodologically grounded** information – not general statements or interpretations.

Each fact must:

- Contain **explicit data** (e.g., numerical values, proportions, metrics,  $R^2$ , accuracies, thresholds, resolutions, years, or geographic scope); or
- Describe **concrete methodologies or models** (e.g., regression, classification, remote sensing methods, LiDAR, MODIS, Landsat-based analysis); or
- Present **clearly bounded evidence** (temporal or spatial scope, dataset name, reference institution, or experimental setup).

Avoid vague claims or general interpretations such as "agriculture improved globally" or "climate change affects yield" – these are not facts.

When listing, follow this structure:

- Use **one bullet per fact**.
- Express each fact as a complete sentence.
- Include the supporting **reference IDs in brackets** (e.g., [7][15]).
- Do not infer beyond the given evidence.

Example:

- In 2020, the average cropland multiple cropping index in eastern China exceeded 180%, derived from MODIS 500 m time-series NDVI data between 2001–2020 [15][32].
- The LiDAR-based digital surface model (DSM) achieved an  $R^2$  of 0.92 when validated against ground survey data in 2019 [7].
- Between 2000 and 2020, Vietnam's arable land area decreased by 5.8%, while agricultural productivity per hectare increased by 12% [30][35].

\*(This fact list serves as the factual foundation for the chapter and will be shown before the generated text.)\*

#### **Step 2 – Fact-Based Writing:**

Write the chapter strictly based on chapter Outline and the extracted facts.

Within each paragraph follow the micro-logic:


- **Fact:** state a referenced observation/data/method [X]
- **Interpretation:** explain significance/implications
- **Connection:** link to the chapter goal or relevant SDG target

Ensure every paragraph has  $\geq 1$  citation.

**\*\*Step 3 – Self-Verification (silent check, then apply fixes):\*\***

- Remove or revise any sentence without support.
- Verify that numbers/units/time ranges match the cited sources.
- Ensure terminology is consistent with references (e.g., sensor names, resolutions, indices).

---

**###**  Good vs Bad


Good = evidence-driven, quantitative, cited.

Bad = generic, uncited, or speculative. Follow the Good pattern.


---

**###**  Research Plan


{plan}

**###**  Already Generated Content

{already\_generated}

**###**  Your Task

Execute the **\*\*current step\*\*** in the research plan.

**###**  Chapter Requirements

- **\*\*Chapter Title\*\***: {chapter\_title} (use exactly)
- **\*\*Chapter Content Structure Requirements\*\***: {content\_structure} (use as the internal outline)
- **\*\*Chapter Core Point Information\*\*** (expand these points with evidence):

{content\_points}

---


**###**  Reference Context (authoritative, numbered)

{context\_with\_references}

{reference\_instruction}


{image\_info\_text}

---

**###**  Output Format (strict)

- 1) **\*\*Extracted Facts (with citations)\*\*** – a short bullet list.
- 2) **\*\*Chapter Text\*\*** – formatted with Markdown headers; every paragraph must include citations like [X].
- 3) **\*\*Evidence Gap Notes (if any)\*\*** – bullet any required points that could not be supported by the references, and specify what evidence would be needed (keep concise).

---

 Notes

- 1) If this is **\*\*not\*\*** the first chapter, do **\*\*not\*\*** output a document title (#); start at **\*\*##** {chapter\_title}.
- 2) Adhere to {content\_structure}; ensure smooth linkage with prior chapters; avoid repetition with "Already Generated Content".
- 3) Expand each argument using {content\_points} + referenced data/methods/cases.
- 4) Do **\*\*not\*\*** write phrases like "as shown in the figure" unless an **\*\*actual image path\*\*** is inserted.
- 5) Prefer concrete values and methods over narrative statements.
- 6) All claims must be traceable to the numbered Reference Context; **\*\*no outside knowledge\*\***.

...



# Evaluation Prompt

## factual verification:

### ### Role

You are a meticulous fact-checking assistant. Your goal is to verify whether a given claim is factually correct by using **web search** and **primary/authoritative sources**. You must avoid hallucination: if evidence is insufficient or ambiguous, return "unsure".

### ### Task

Given the claim below, (1) design effective search queries, (2) retrieve evidence from the web, (3) assess whether the evidence supports or refutes the claim, and (4) output a strict JSON object following the schema.

### ### Claim

{fact}

### ### Verification Guidelines

- Decompose if needed**: If the claim contains multiple factual components (entities, numbers, dates, locations, causal relations), break it into atomic sub-claims mentally and verify each.
  - Output "true" only if **all** core sub-claims are supported.
  - Output "false" if **any** core sub-claim is contradicted by credible evidence.
  - Output "unsure" if evidence is missing, conflicting, outdated, or not authoritative enough.
- Source quality & independence**
  - Prefer: official organizations, academic venues, standards bodies, reputable news outlets, peer-reviewed papers, well-maintained official docs.
  - Avoid relying on a single source. Use **at least 2 independent sources** when possible.
  - If only low-quality sources exist, lower confidence or output "unsure".
- Time sensitivity**
  - If the claim is time-dependent (e.g., "latest version", "as of 2025"), ensure evidence is recent and include publication date when available.
  - If sources disagree across time (old vs new), treat that as potential conflict.
- Evidence handling**
  - Evidence must be directly relevant to the claim (not merely related).
  - Provide short excerpts ( $\leq 25$  words) as "snippets" and include enough context to justify the verdict.
  - If you find contradictory evidence, include it in `conflicting\_info` and lower confidence.
- Confidence calibration (0.0-1.0)**
  - 0.90-1.00: multiple high-quality independent sources directly confirm/deny, minimal ambiguity.
  - 0.70-0.89: good evidence but minor ambiguity (e.g., wording mismatch, limited sources).
  - 0.40-0.69: weak/indirect evidence or partial verification of sub-claims.
  - 0.00-0.39: little evidence, severe conflicts, or unclear claim → usually "unsure".

### ### Output Requirements

- Output **JSON only** (no markdown, no comments, no trailing commas).
- Do not fabricate URLs, titles, dates, or quotes.
- If you cannot find reliable evidence, return "unsure" with low confidence and explain what is missing.
- Keep `reasoning` concise (2-6 sentences). Do not include step-by-step hidden deliberation.

### ### Output JSON (strict)

```
 "supporting_evidence": [
 {
 "claim": "{fact}",
 "url": "",
 "answer": "true" | "false" | "unsure",
 "title": "",
 "confidence": 0.0,
 "publisher": "",
 "reasoning": "",
 "date": "",
 "conflicting_info": "",
 "snippet": "",
 "search_queries": [
 "stance": "support" | "refute"
 ""
]
},
]
```

## quiz-based factual coverage:

You are an expert evaluator for Sustainable Development Goals (SDG) reports. Your job is to answer the user question using ONLY the provided Report Text. Do not use outside knowledge. Do not guess. If the report does not explicitly contain the answer, say so.

Question:

{QUESTION}

Report Text:

.....

{REPORT\_TEXT}

.....

Instructions (strict):

1. Read the question and identify what type of information is requested:
  - Definition/term, dataset/source, method/formula, numeric value, time range, spatial scope, trend/comparison, limitation/uncertainty, policy implication.
2. Locate the answer by searching the Report Text:
  - Prefer exact matches; also consider paraphrases and abbreviations.
  - If multiple candidate passages exist, choose the one that most directly answers the question.
3. Extract the answer:
  - Keep it concise (typically 1–2 sentences or a short phrase/value).
  - Preserve numbers, units, years, and names exactly as written.
  - Do NOT rewrite the claim beyond minor trimming. No additional interpretation.
4. Provide evidence:
  - Copy the exact supporting sentence(s) verbatim from the Report Text.
  - If one sentence is insufficient, include up to 3 sentences, concatenated with a single space.
  - Evidence must be the minimal span that supports the answer.
5. Handle ambiguity:
  - If the report contains related content but does not fully answer the question, output "No relevant content found." (do not guess).
  - If the question asks for something not present (e.g., specific country/year not mentioned), output "No relevant content found."
6. Confidence scoring (0.0–1.0):
  - 0.90–1.00: direct explicit statement matching the question exactly.
  - 0.70–0.89: clear but slightly indirect (e.g., paraphrase, minor inference from a stated table/value).
  - 0.40–0.69: partially related; answer is incomplete → usually should return "No relevant content found."
  - 0.00–0.39: not found or highly ambiguous → return "No relevant content found."
  - If you output "No relevant content found.", do not output confidence.

Output JSON (strict; JSON only, no markdown, no extra keys):

Option A (answer found):

```
{
 "answer": "...",
 "supporting_sentence": "...",
 "confidence": 0.0
}
```

Option B (not found):

```
{
 "answer": "No relevant content found."
}
```