

SciVSG: Unified Visual-Semantic Graph with Traceable Evidence for Scientific Diagram Understanding

Anonymous ACL submission

Abstract

Scientific diagrams capture complex mechanisms, experimental evidence, and structural relationships, yet remain challenging to interpret and reason over due to their heterogeneous layouts and lack of traceable representations. Scientific Visual-Semantic Graph (SciVSG) addresses this challenge by providing a unified visual-semantic representation that integrates layout and semantic information. It first constructs a Visual Layout Hierarchy (VLH) from layout cues and reading conventions, establishing a structured foundation for diagram understanding. Node-level verifiable evidence, including localized OCR and aligned paper snippets, grounds predictions to explicit text spans and regions. On this basis, a Semantic Scene Graph (SSG) is built by linking typed entities and normalized relations to nodes under strict evidence constraints, enabling module-aware reasoning and fine-grained traceability. A benchmark of diverse scientific diagrams is also provided, annotated with VLH, node evidence, entities, relations, and expert-authored QA pairs across multiple categories. Experiments demonstrate that SciVSG substantially enhances knowledge extraction and produces more reliable, evidence-attributable answers for diagram-based question answering.

Keywords: Scientific Diagram Understanding; Visual-Semantic Representation; Visual Layout Hierarchy; Semantic Scene Graph; Knowledge Extraction

1 Introduction

Scientific diagrams are essential for conveying mechanisms, experimental procedures, and structural relationships in research papers (Pan et al., 2024). They often encode actionable knowledge that researchers seek to retrieve, compare, and reuse, including workflows, component arrangements, and interaction logic (Zala et al., 2023; Pra-

manick et al., 2024). Despite their importance, accessing diagram-centric knowledge remains challenging due to visual heterogeneity, structural density, and modular organization with connectors and localized text (Zala et al., 2023; Kembhavi et al., 2016; Cui et al., 2025).

Scientific diagrams can be broadly categorized into four layout-driven types based on visual composition and relation expression (Mondal et al., 2024): **Process Diagrams** such as experimental pipelines, **Structural Diagrams** like part-whole schematics, **Interaction Diagrams** including regulatory networks, and **Composite Diagrams** that combine multiple modules (Zala et al., 2023). This categorization illustrates the diversity of diagram structures, highlighting the need for representations capable of capturing modularity, terminal visual units, and explicit linking cues. These diagram types serve as background context rather than the focus of the present work.

Existing methods address diagrammatic challenges unevenly (Suri et al., 2025). While regular process diagrams benefit from specialized benchmarks and QA methods (Pramanick et al., 2024), and educational datasets like AI2D (Kembhavi et al., 2016) and AI2D-RST (Hiippala et al., 2021) offer insights into element-level and discourse-inspired organization, they are not designed as general-purpose interfaces for scientific diagrams. Specifically, they lack support for module-aware extraction within nested hierarchies and fail to enforce verifiable evidence attribution linking predictions to localized regions. Similarly, conventional scene graph methods (Kembhavi et al., 2016; Hiippala et al., 2021), although effective for natural images, cannot handle the multi-module structures and dense connectors inherent in scientific contexts, thus limiting their capacity for traceable reasoning.

To bridge this gap, this work introduces **Scientific Visual-Semantic Graph (SciVSG)**, a unified

084 representation that combines layout and semantic
085 information into a coherent, evidence-grounded
086 structure. SciVSG first constructs a **Visual Lay-
087 out Hierarchy (VLH)** using visual cues and read-
088 ing conventions, decomposing diagrams into mod-
089 ules and terminal visual units. Node-level ver-
090 ifiable evidence, including localized OCR and
091 aligned paper snippets, is then attached to VLH
092 nodes, grounding predictions in explicit text and
093 regions. Building on this foundation, SciVSG
094 generates a **Semantic Scene Graph (SSG)** by
095 linking typed entities and normalized relations to
096 nodes under strict evidence constraints, producing
097 a traceable semantic graph that enables module-
098 aware reasoning and evidence-attributable predic-
099 tions. Together, VLH and SSG form an integrated
100 framework for structured diagram understanding.

101 The main contributions are as follows: (1)
102 SciVSG, a unified visual-semantic representation
103 combining a layout-induced hierarchy (VLH),
104 node-attached verifiable evidence, and a node-
105 grounded semantic graph (SSG) for scientific di-
106 agrams. (2) Traceability-oriented constraints that
107 ensure all extracted entities, relations, and QA an-
108 swers can be checked against localized regions
109 and evidence spans. (3) A benchmark of scien-
110 tific diagrams with expert-authored QA for rig-
111 orous evaluation of unified, evidence-grounded di-
112 agram understanding.

113 2 Related Work

114 **Hierarchical Layout and Visual Structure.** Un-
115 derstanding diagrams requires modeling hierarchi-
116 cal and compositional visual structures that reflect
117 how elements are organized and connected (Shen
118 et al., 2021; Zhong et al., 2019). Early resources
119 such as AI2D provide annotated elements and
120 parse graphs for educational diagrams, captur-
121 ing relations among textual and visual compo-
122 nents (Kembhavi et al., 2016). AI2D-RST ex-
123 tends AI2D with multi-layer annotations of group-
124 ing, connectivity, and discourse structure, support-
125 ing analysis of compositional layouts and visual
126 discourse (Hiippala et al., 2021). Recent diagram
127 parsing work emphasizes modular decomposition
128 and structural grouping, but largely targets dataset
129 construction or element-level parsing rather than
130 unified, evidence-grounded semantic representa-
131 tions.

132 **Graph-Based Semantic Modeling.** Graph-
133 structured representations (e.g., scene graphs) en-

code objects and relations for visual reasoning (Kr-
ishna et al., 2017; Zellers et al., 2018). How-
ever, they are less suited to multi-module diagrams
with dense connectors and localized labels (Li and
Tajbakhsh, 2023). Graph-based multimodal cor-
pora also model discourse and connectivity among
diagram elements (Hiippala et al., 2021), but typ-
ically lack typed entities and normalized relations
under strict evidence constraints. SciVSG ad-
vances this line by integrating layout hierarchies,
typed entities, and semantic relations into a uni-
fied, traceable graph representation.

Diagram Question Answering. Diagram
QA benchmarks have grown from educa-
tional and flowchart datasets (e.g., TQA,
FlowchartQA) (Kembhavi et al., 2016; Kafle
et al., 2018; Tannert et al., 2023) to scientific
resources. SPIQA targets figures and tables
in scientific papers, evaluating interpretation
of complex visual structures in research con-
texts (Pramanick et al., 2024). MathVerse
benchmarks visual math understanding, testing
whether models interpret diagrams rather than
rely on textual shortcuts (Zhang et al., 2024).
MISS-QA evaluates multimodal foundation
models on schematic diagrams in scientific
papers (Zhao et al., 2025). While these datasets
support evaluation, none provide an interpretable
computational framework for explicit reasoning
over evidence-grounded semantic relations.

164 3 Layout-induced Hierarchical Diagram 165 Representation

166 We present *SciVSG*, our *Layout-induced Hier-*
167 *archical Diagram Representation* shown in Fig. 1.
168 SciVSG builds a *Visual Layout Hierarchy* (VLH)
169 to organize a diagram into localized regions with
170 explicit containment and a reader-oriented traver-
171 sal order, providing a layout-grounded backbone
172 for subsequent understanding. To ensure traceabil-
173 ity, each node is augmented with verifiable evi-
174 dence, and we further attach typed entities and
175 normalized relations to form a node-grounded *Se-*
176 *mantic Scene Graph* (SSG) under strict evidence
177 constraints. We also release a diagram benchmark
178 annotated with VLH, node evidence, SSG, and
179 expert-authored QA for evaluation.

180 3.1 Visual Layout Hierarchy

181 We represent a scientific *diagram* as a *Visual Lay-*
182 *out Hierarchy* (VLH): a rooted tree that decom-

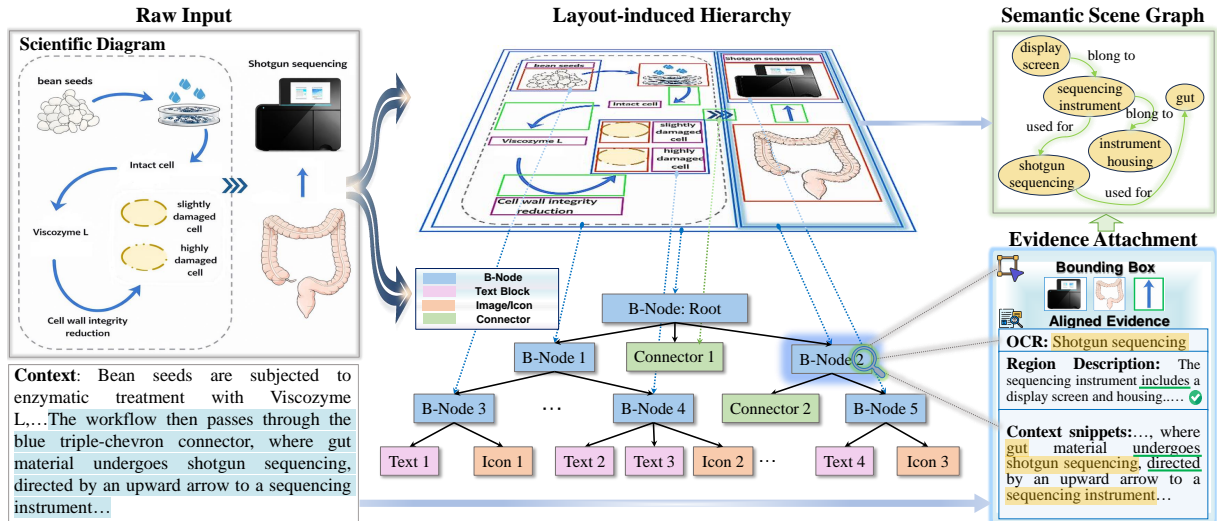


Figure 1: Layout-induced Hierarchical Diagram Representation.

poses the diagram into localized regions with explicit containment and a reader-oriented traversal order. The VLH encodes *layout* only; semantic objects (entities and relations) are attached later as node-level annotations (Section 3.2).

3.1.1 Representation Schema

A central design choice in VLH is separating *decomposable layout regions* from *terminal visual units*. Branch nodes represent composite regions that can be further partitioned. TextBlock, Icon, Image, and Connector are treated as terminal primitives (leaf nodes) during induction. This is a structural rather than semantic distinction. Terminal nodes are not split even when they contain multiple mentions or relations.

Each node belongs to one of five categories covering dominant diagram building blocks: **Branch node (B-Node)**: a composite module (e.g., panel/subfigure, bounded schematic module) that reduces cross-module ambiguity and supports module-level aggregation and traceable referencing. **TextBlock**: a contiguous text region (e.g., title, label cluster) serving as the primary OCR carrier. **Icon**: a non-photographic drawn region conveying schematic content, providing stable anchors for descriptions and aligned textual evidence. **Image**: a photographic/experimental region (e.g., microscopy, gels) that links visually grounded observations to aligned textual explanations. **Connector**: an explicit linking cue (e.g., arrows, lines, braces) modeled as a first-class node due to its directionality and structured associations.

For grounding and traceability, each node stores: (i) a parent-child link (containment), (ii) a

bounding box, (iii) a short region description, and (iv) node-referenced text sources (localized OCR and selected caption/context snippets). How these fields are populated and aligned is described in Section 3.1.3.

VLH nodes need not be semantically atomic. We store entities as node-local mention strings (optionally typed) under the hosting node; each mention must be supported by at least one node-referenced text source, preserving traceability while keeping the layout backbone independent of the semantic layer.

3.1.2 Layout-induced Hierarchy Induction

We induce the Visual Layout Hierarchy (VLH) in a top-down manner using fixed annotation rules and reading conventions. The process is *annotation-driven*: the hierarchy is deterministically constructed from manually annotated node types and bounding boxes, without learned detectors. The VLH has no fixed depth.

Step 1: determine decomposability. Starting from the root, we assess whether the diagram contains separable modules under the schema. Visually atomic diagrams are represented as single-node trees.

Step 2: propose modules. For decomposable diagrams, we annotate perceptually coherent regions as Branch nodes when at least one cue is present, including boundaries, background differences, distinctive alignments, or local headers such as panel tags.

Step 3: build the branch backbone. We establish parent-child relations among Branch nodes by minimal geometric containment and order sib-

251	lings by reading conventions to obtain a determin-	verifiable node evidence. In §3.2, we attach typed	302
252	istic traversal.	entities and normalized relations to \mathcal{T} to enable	303
253	Step 4: attach leaf nodes. Terminal primitives	module-aware reasoning with region-level attribu-	304
254	are attached to the smallest containing Branch	tion.	305
255	node and are not further decomposed, serving as		
256	atomic units for evidence attachment and ground-	3.2 Semantic Scene Graph (SSG)	306
257	ing.		
258	Overall, given annotations and fixed rules, VLH	We build a <i>Semantic Scene Graph</i> (SSG) on top	307
259	induction is deterministic: branch consolidation	of the visual layout hierarchy (VLH) by attach-	308
260	follows overlap resolution, parent assignment fol-	ing typed entities and normalized relations to VLH	309
261	lows minimal containment, and expansion termi-	nodes. The result is a traceable semantic graph	310
262	ates at terminal primitives, yielding a stable struc-	that supports module-aware reasoning and <i>verifi-</i>	311
263	ture for subsequent evidence alignment and seman-	<i>able</i> evidence attribution at the node (region) level.	312
264	tic graph construction.		
265	3.1.3 Node Evidence Attachment	3.2.1 Entity and Relation Inventory	313
266	To support extraction and reasoning, each VLH	An entity is a <i>nameable</i> scientific object mention.	314
267	node is grounded to a region and augmented with	Each entity must be grounded in a mention string	315
268	compact, verifiable evidence fields, while preserv-	that can be <i>verbatim located</i> in at least one textual	316
269	ing the VLH topology.	node-attached evidence field (e.g., localized OCR	317
270	Bounding box. Each node stores a bounding	strings or region captions). This constraint en-	318
271	box for localization and stable attribution.	sures automatic verifiability and discourages hal-	319
272	Region description vs. context snippets. A	lucinated entities.	320
273	<i>region description</i> summarizes observable visual	We adopt a seven-way entity type system (Hu	321
274	cues (e.g., position, boundaries) and is independ-	et al., 2024): Tool (named instrument, device, plat-	322
275	ent of paper text, serving as an interpretable an-	form, or software; e.g., <i>mass spectrometer</i>); Ma-	323
276	chor and auxiliary alignment signal. A <i>context</i>	terial (physical, biological, or data-carrying sub-	324
277	<i>snippet</i> is a relevant caption/context span aligned	stance; e.g., <i>serum sample</i>); Indicator (measur-	325
278	to the node as textual evidence.	able metric or property; e.g., <i>accuracy</i>); Method	326
279	Localized OCR. OCR is primarily attached	(approach, algorithm, or model; e.g., <i>Bayesian in-</i>	327
280	to TextBlock nodes; for Connector nodes we	<i>ference</i>); Problem (task, defect, disease, or phe-	328
281	record nearby OCR fragments when available.	nomenon; e.g., <i>noise removal</i>); Dataset (uniquely	329
282	For Image/Icon nodes, OCR may be sparse and	identifiable named dataset; e.g., <i>MNIST</i>); Value	330
283	aligned snippets serve as the main textual anchors.	(numeric or comparative numeric expression; e.g.,	331
284	Aligned paper snippets. We align cap-	<i>95%</i>).	332
285	tion/context spans using reproducible heuristics	A relation is a directed semantic edge between	333
286	based on OCR overlap, connector-adjacent trig-	two grounded entities, drawn from a <i>fixed</i> label set.	334
287	ger words, and panel/local identifiers (e.g., “(a)”,	Relation labels are normalized semantic categories	335
288	“Panel B”). Branch nodes may store module-level	and need not appear verbatim in evidence; instead,	336
289	snippets, while fine-grained evidence is anchored	evidence must support the underlying claim and	337
290	at terminal primitives.	be attributable to a specific node.	338
291	All evidence is stored with minimal provenance	We define seven relation labels: Belong To	339
292	(<i>source</i> , <i>node_id</i> , exact span); each span must	(part-whole or membership); Used For (func-	340
293	be an exact substring of its source (no paraphras-	tional usage); Measure (measurement or eval-	341
294	ing). Region descriptions cannot introduce en-	uation); Produce (generation or output); Pre-	342
295	tity mentions; entities must be supported by OCR	cede (temporal or procedural ordering); Equal	343
296	or aligned snippets. We enforce: tree validity ,	To (aliasing or semantic equivalence); Affect (di-	344
297	containment consistency , type constraints (no	rected influence such as increase or suppression).	345
298	decomposition of primitives), and evidence in-	If a downstream benchmark adopts a different	346
299	tegrity (exactly matchable spans), flagging viola-	relation inventory, only label substitution is re-	347
300	tions automatically.	quired; grounding and evidence rules remain un-	348
301	This step outputs a grounded VLH tree \mathcal{T} with	changed.	349

3.2.2 Node-attached Entity Construction

SSG does not assume that VLH nodes are semantically atomic. A single node (including TextBlock, and nodes containing embedded labels near Icon or Image regions) may host zero, one, or multiple entity mentions. In our current data format, a node stores entity mentions as strings (with types) rather than forcing entities to be leaf nodes in the hierarchy.

Entity mentions are extracted and attached using **textual** node-attached sources only, namely localized OCR strings and aligned caption, context snippets. We enforce a *verbatim mention constraint*: every stored entity string must be exactly matchable as a substring in one of the allowed textual evidence fields of its hosting node. Notably, region descriptions are **not** used as an entity source, and cannot introduce new entity mentions.

We apply lightweight normalization (whitespace trimming and canonical punctuation where appropriate) while preserving scientific symbols and units. If the same mention appears in multiple nodes, we keep multiple grounded instances because they carry distinct spatial provenance; an optional canonical id can be added for downstream consolidation.

3.2.3 Relation Construction with Evidence Attribution

For each node v , we form candidate entity pairs within a *module-constrained* scope enabled by the VLH: (i) **connector-centric** pairs suggested by connector geometry and nearby labels, (ii) **within-node** co-occurrence pairs from entities attached to the same node, and (iii) **within-subtree** pairs whose endpoints fall under the same Branch node. This avoids global enumeration and reduces cross-panel ambiguity.

For each candidate (h, t) , we select a relation label r from the fixed inventory using node-scoped cues (connector direction/morphology, nearby triggers, aligned caption/context, and region descriptions). Each relation must include at least one evidence record $\langle source, span, v \rangle$, where *span* is copied verbatim from an allowed node-attached source (OCR, connector-local text, region description, or aligned snippets); we enforce **exact span matching** (no paraphrasing). Evidence need not contain the label name, but must support the normalized claim. Attaching relations to VLH nodes with boxes yields a traceability chain: *relation* \rightarrow *evidence span* \rightarrow *node id* \rightarrow *node box*.

If endpoints are grounded to different nodes, we attach the relation to the *lowest common ancestor* Branch node covering both, while preserving each endpoint’s original node id in the evidence metadata.

3.3 Benchmark Construction (Data and QA)

We construct a benchmark to evaluate unified and traceable diagram representations across four genres: **Process**, **Structural**, **Interaction**, and **Composite** diagrams. Each sample is annotated with (i) a Visual Layout Hierarchy (VLH), (ii) node evidence fields enabling span-verifiable attribution, and (iii) a Semantic Scene Graph (SSG) with typed entities and normalized relations attached to VLH nodes. Domain experts additionally author QA pairs, which we categorize into four types with strict answer matching for automatic evaluation.

3.3.1 Diagram Collection and Coverage

We collect candidate papers from IEEE venues and arXiv, starting from approximately 30000 papers with clear diagrams and sufficient context. Manual screening of panels and captions yields 18760 diagram candidates that are readable, self-contained, and interpretable with surrounding text.

We then curate a high-quality subset suitable for annotation, retaining diagrams that meet: (i) adequate resolution for text/connectors, (ii) clear module boundaries when applicable, (iii) unambiguous reading order or structural organization, and (iv) sufficient context for evidence alignment. After filtering, we retain **1200** diagrams, including **600** English and **600** Chinese. Coverage across the four genres is summarized in Table B.1 (Appendix).

We balance the dataset across genres. Process diagrams show stepwise transitions. Structural diagrams highlight partwhole organization. Interaction diagrams feature dense, typed relations. Composite diagrams combine multiple local organizations in a single layout.

3.3.2 Annotation Protocol

Annotation layers. Each diagram is annotated in three layers.

(1) **VLH structure and node boxes.** Annotators construct VLH following the layout-induced rules in Section 3.1.2, assigning each node a type and bounding box. Branch nodes are decomposable modules, while terminal primitives (TextBlock, Icon, Image, Connector) are not subdivided.

(2) Node evidence fields. We attach three textual sources per node: (i) *region description* written purely from visual content, (ii) *context snippets* aligned from paper text, and (iii) *OCR text* for TextBlock nodes (and optional nearby OCR cues for Connector nodes). These fields support verifiable attribution in extraction and QA.

(3) SSG entities and relations. Annotators attach typed entities and normalized relations to VLH nodes under evidence constraints: entity mentions are copied verbatim from OCR or context snippets, and each relation is supported by at least one evidence span.

Each evidence record is stored as a four-tuple

$$\langle source, node_id, span, char_offset \rangle,$$

where *span* is an exact substring of the declared source (OCR/context/region description), and *char_offset* is verified by deterministic lookup.

3.3.3 Consistency and Quality Control

Each diagram is annotated by two annotators and adjudicated by a senior annotator to resolve disagreements in VLH boundaries/types, entity typing, and relation labels. We apply checks for (i) tree validity (single root, acyclic, unique parent), (ii) geometric consistency (child boxes within parent up to tolerance), (iii) connector integrity (terminal connectors with supported endpoints when applicable), and (iv) evidence integrity (exact substring spans with valid offsets).

Agreement is reported at multiple layers: mean IoU and node-type agreement for VLH/boxes; set-based Precision/Recall/F1 over unique (mention, type) pairs for entities; and set-based Precision/Recall/F1 over relation triples for relations. Representative agreement numbers are: mean box IoU **0.86**, node-type agreement **0.95**, entity F1 **0.84**, and relation F1 **0.79**.

3.3.4 Expert-authored QA and Taxonomy

Domain experts author questions using the diagram and aligned context; questions are not templated from labels. We group questions into four categories and constrain answers to be short and matchable for automatic scoring.

Q1 Entity: identify a named entity of a specified type. **Q2 Relation:** identify the normalized relation label between two named entities with traceable evidence. **Q3 Hierarchical trace:** identify the module where a claim is expressed and, if needed, return the module path from the root using

matchable identifiers (e.g., panel tags or OCR fragments). **Q4 Hierarchy-helpful reasoning:** answer questions that benefit from restricting or aggregating evidence within the hierarchy.

We collect **8** questions per diagram on average, yielding **7200** QA pairs, balanced across four categories. Detailed distributions and verification constraints are summarized in Table B.2 (Appendix).

4 Experiments

4.1 Overview and Experimental Setup

We evaluate SciVSG on four targets: VLH structure extraction, SG construction, SSG construction, and diagram question answering (QA). Our benchmark covers four diagram genres: *Process*, *Structural*, *Interaction*, and *Composite*. These genres differ in module organization and in connector-mediated relations.

Why off-the-shelf foundation models without training? We validate SciVSG as a *model-agnostic intermediate representation* and test how robustly different multimodal foundation models instantiate it under the same stage-wise protocol. This setting (i) isolates the contribution of the representation and supervision interface (schema, evidence constraints, normalization) from task-specific fine-tuning, and (ii) reflects a plug-and-play deployment without additional training.

Stage-wise construction and evaluation targets. Given a diagram image and associated paper text, we run a fixed pipeline to obtain: (i) a VLH with node types and parent-child edges, (ii) an SG of semantic triplets, and (iii) an SSG that attaches typed entities and normalized relations to VLH nodes with evidence. We evaluate VLH node typing and edges, SG/SSG triplets, and QA under controlled ablations.

Cross-model evaluation (Table 1). We compare multiple multimodal foundation models with identical inputs and prompts to quantify consistency in hierarchy induction and semantic graph construction under strict matching.

Ablations (Table 2). Using a single model, we vary inputs (image; caption+context; image+caption+context) and provided representations (VLH, SG, SSG) to isolate component effects.

Traceability analysis. A visualized QA case study in Appendix Figure B.1 shows that VLH enables precise back-tracing from answers to supporting nodes and evidence.

4.2 Metrics

We evaluate (i) **VLH structure extraction** with Node-type macro-F1 and Edge-F1, (ii) **SG/SSG construction** with set-based triplet Precision/Recall/F1, and (iii) **diagram QA** with BLEU and ROUGE-L.

Notation. Let a diagram have a gold VLH with node set \mathcal{N} and directed parent-child edge set \mathcal{E} . Let $\hat{\mathcal{N}}$ and $\hat{\mathcal{E}}$ denote predicted nodes and edges. Let \mathcal{T}^{SG} and $\hat{\mathcal{T}}^{SG}$ be gold and predicted SG triplet sets, and \mathcal{T}^{SSG} and $\hat{\mathcal{T}}^{SSG}$ be gold and predicted SSG item sets.

Exact matching and normalization. For graph items containing text fields (heads/tails in SG/SSG), we normalize mentions by (1) trimming leading/trailing whitespace and (2) collapsing multiple spaces into one, while keeping case and punctuation unchanged. All set intersections below use *exact* equality after this normalization.

Unified set-based Precision/Recall/F1. For any gold set A and predicted set \hat{A} , we compute

$$P(A, \hat{A}) := \frac{|\hat{A} \cap A|}{|\hat{A}|}, \quad (1a)$$

$$R(A, \hat{A}) := \frac{|\hat{A} \cap A|}{|A|}, \quad (1b)$$

$$F1(A, \hat{A}) := \frac{2P(A, \hat{A})R(A, \hat{A})}{P(A, \hat{A}) + R(A, \hat{A})}. \quad (1c)$$

To provide a more comprehensive evaluation of the diagram parsing performance, we apply this unified metric to evaluate the various structural dimensions of the diagrams. Specifically, we report Node-F1 and Edge-F1 to assess node typing and edge recovery within the Visual Layout Hierarchy (VLH). For semantic triplet extraction in both Scene Graphs (SG) and Situated Scene Graphs (SSG), we report the full suite of Precision, Recall, and F1-score. Detailed formal definitions for each task-specific instantiation are provided in Appendix A.

4.3 Cross-model Evaluation on Structure, SG, and SSG Construction

Table 1 summarizes module-wise evaluation across foundation models for VLH structure extraction, SG construction, and SSG construction.

VLH is relatively stable across models. Structure extraction yields comparatively high Node-type F1 and Edge-F1 across models, indicating

that layout-induced hierarchy can be reliably induced from visual layout cues and reading conventions.

Graph construction remains the bottleneck. SG F1 is consistently lower than VLH scores, reflecting the difficulty of extracting correct semantic relations among heterogeneous terminal units. SSG is further harder than SG because it requires both semantic correctness and correct node attachment under a normalized inventory.

SSG exposes grounding failures beyond semantic errors. The gap between SG and SSG indicates additional error modes, including incorrect assignment of entity mentions to nodes, attachment inconsistency across modules, and relation normalization mistakes. These results motivate a controlled study of how different inputs and intermediate representations affect downstream QA.

4.4 Ablations over Raw Inputs and Representation Variants

Table 2 reports ablations over raw inputs and representation variants using a single fixed model.

Effect of raw inputs. Caption plus context improves language-heavy questions because many answer strings are explicitly stated in the paper text. Image-only inputs remain competitive when questions depend on visual layout cues, module boundaries, and connector-mediated transitions. Combining image with caption plus context yields the best overall performance, showing that visual structure and textual evidence are complementary.

Effect of VLH. Providing VLH improves Q3 and Q4 because the hierarchy constrains the candidate space to the relevant module and reduces cross-module confusion when similar tokens recur in multiple regions.

SG versus SSG. SG helps relation-centric questions by explicitly exposing candidate triplets, but it can still mix instances across modules when the structure is implicit. SSG strengthens consistency by enforcing node attachment and normalized relations, which is especially beneficial for questions that require traceability and module-consistent grounding.

4.5 Summary

Overall, the experiments show that a reliable layout hierarchy forms a strong backbone, while semantic graph construction remains the primary challenge. The ablations confirm that VLH contributes by enforcing hierarchical constraints

Table 1: **Module-wise evaluation across foundation models.** We report Node-type F1 and Edge-F1 for VLH structure extraction, and Precision / Recall / F1 for SG construction and SSG construction.

Model	Metrics		Structure extraction (VLH)			SG construction (SG)			SSG construction (SSG)		
	Node-F1↑	Edge-F1↑	P↑	R↑	F1↑	P↑	R↑	F1↑			
GPT-5-mini	0.861	0.413	0.546	0.482	0.514	0.699	0.636	0.668			
GPT-5-nano	0.712	0.345	0.538	0.474	0.506	0.685	0.622	0.654			
GPT-5.1	0.936	0.619	0.586	0.518	0.552	0.748	0.676	0.712			
Gemini-3	0.921	0.618	0.532	0.468	0.499	0.674	0.611	0.643			
Qwen3-VL-235B	0.420	0.310	0.540	0.476	0.508	0.681	0.618	0.649			
Qwen2.5-VL-32B	0.743	0.367	0.387	0.310	0.349	0.415	0.333	0.374			
DeepSeek-VL2	0.280	0.220	0.333	0.258	0.296	0.402	0.317	0.360			
Claude Haiku 4.5	0.858	0.407	0.555	0.489	0.522	0.701	0.636	0.669			

Table 2: **Ablations over raw inputs and representation variants.** We use a single model and report BLEU and ROUGE-L by question type (Q1–Q4) and the macro average (AVG).

Setting	Raw inputs		Representation			BLEU by question type					ROUGE-L by question type				
	Image	Cap+Ctx	VLH	SG	SemSG	Q1	Q2	Q3	Q4	AVG	Q1	Q2	Q3	Q4	AVG
S1	✓	–	–	–	–	0.58	0.79	0.32	0.67	0.59	0.76	0.83	0.53	0.81	0.73
S2	✓	–	✓	–	–	0.89	0.91	0.14	0.62	0.64	0.92	0.89	0.38	0.75	0.73
S3	✓	–	–	✓	–	0.90	0.62	0.12	0.61	0.56	0.92	0.96	0.38	0.76	0.75
S4	✓	–	✓	✓	–	0.90	0.28	0.14	0.67	0.50	0.91	0.97	0.39	0.78	0.76
S5	✓	–	✓	–	✓	0.92	0.94	0.35	0.66	0.72	0.89	0.89	0.56	0.81	0.79
S6	–	✓	–	–	–	0.28	0.87	0.18	0.41	0.44	0.46	0.91	0.39	0.60	0.59
S7	–	✓	✓	–	–	0.79	0.91	0.09	0.59	0.60	0.85	0.87	0.37	0.69	0.70
S8	–	✓	–	✓	–	0.77	0.96	0.13	0.60	0.62	0.79	0.93	0.38	0.73	0.71
S9	–	✓	✓	✓	–	0.84	0.99	0.12	0.56	0.62	0.87	0.97	0.37	0.72	0.73
S10	–	✓	✓	–	✓	0.85	0.98	0.39	0.65	0.72	0.87	0.97	0.55	0.75	0.79
S11	✓	✓	–	–	–	0.61	0.92	0.37	0.66	0.64	0.79	0.89	0.55	0.82	0.76
S12	✓	✓	✓	–	–	0.86	0.94	0.14	0.59	0.63	0.90	0.91	0.38	0.73	0.73
S13	✓	✓	–	✓	–	0.86	0.96	0.16	0.62	0.65	0.89	0.91	0.39	0.76	0.74
S14	✓	✓	✓	✓	–	0.90	0.97	0.17	0.61	0.66	0.93	0.94	0.39	0.75	0.75
S15	✓	✓	✓	–	✓	0.91	0.98	0.38	0.72	0.75	0.93	0.96	0.57	0.80	0.82

and module locality, and SSG further improves faithfulness through node-attached semantics and evidence-consistent grounding. These findings support SciVSG as a practical intermediate representation for diagram-centric knowledge extraction and question answering.

5 Conclusion

We tackle scientific diagram understanding, where complex layouts and weak grounding make extraction and reasoning hard to verify. We propose SciVSG, a unified visual-semantic representation that combines a layout backbone with node-grounded semantics: it first builds a Visual Layout Hierarchy (VLH) from layout cues and reading order, then attaches verifiable node evidence (localized OCR and aligned caption/context snippets), and finally constructs a node-grounded Semantic Scene Graph (SSG) with typed entities and nor-

malized relations under strict evidence constraints. Experiments on our benchmark show that SciVSG improves diagram knowledge extraction and produces more reliable, evidence-attributable answers for diagram QA.

Limitations

SciVSG still has limitations, notably it does not directly handle charts, which require specialized parsing of axes, scales, and data series. Looking forward, two practical directions are as follows.

(i) building robust detectors for the basic visual units defined by SciVSG to enable scalable, automatic construction;

(ii) strengthening cross-source grounding by tighter alignment between figure regions and paper text/tables, supporting broader coverage and more complete scientific interpretation.

679
680
681
682
683

684

685
686
687
688
689
690
691

692
693
694
695
696

697
698
699
700
701

702
703
704
705
706

707
708
709
710
711

712
713
714
715
716
717
718

719
720
721
722

723
724
725
726
727
728
729
730

Acknowledgments

This paper used large language models (e.g., ChatGPT) to assist data annotation during the image understanding pipeline, and we sincerely appreciate their support.

References

Zhiqing Cui, Jiahao Yuan, Hanqing Wang, Yanshu Li, Chenxu Du, and Zhenglong Ding. 2025. [Draw with thought: unleashing multimodal reasoning for scientific diagram generation](#). In *Proceedings of the 33rd ACM International Conference on Multimedia*, MM 25, pages 5050–5059. Association for Computing Machinery.

Tuomo Hiippala, Malihe Alikhani, Katri Haverinen, Arne Köhn, Stina Ojala, Achim Stein, and Sanni Uusi-Mäkelä. 2021. Ai2d-rst: A multimodal corpus of 1,000 primary school science diagrams. *Language Resources and Evaluation*, 55(3):661–688.

Maodi Hu, Li Qian, Zhijun Chang, and Zhixiong Zhang. 2024. [Kdpg-enhanced mrc framework for scientific entity recognition in survey papers](#). *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 32:25322543.

Kushal Kaffle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. [DVQA: Understanding data visualizations via question answering](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5648–5656.

Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. [A diagram is worth a dozen images](#). In *Computer Vision – ECCV 2016*, pages 235–251, Cham. Springer International Publishing.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael Bernstein, and Li Fei-Fei. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International Journal of Computer Vision*, 123(1):32–73.

Shengzhi Li and Nima Tajbakhsh. 2023. [Scigraphqa: a large-scale synthetic multi-turn question-answering dataset for scientific graphs](#). *arXiv preprint arXiv:2308.03349*.

Ishani Mondal, Zongxia Li, Yufang Hou, Anandhavelu Natarajan, Aparna Garimella, and Jordan Lee Boyd-Graber. 2024. [Scidoc2diagrammer-maf: towards generation of scientific diagrams from documents guided by multi-aspect feedback refinement](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 13342–13375. Association for Computational Linguistics.

Huitong Pan, Qi Zhang, Cornelia Caragea, Eduard Dragut, and Longin Jan Latecki. 2024. [Flowlearn: evaluating large vision-language models on flowchart understanding](#). *arXiv preprint arXiv:2407.05183*. 731
732
733
734
735

Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. 2024. [SPIQA: A dataset for multimodal question answering on scientific papers](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 118807–118833. Curran Associates, Inc. 736
737
738
739
740
741

Zejiang Shen, Ruochen Zhang, Melissa Dell, Benjamin Charles Germain Lee, Jacob Carlson, and Weining Li. 2021. [LayoutParser: A unified toolkit for deep learning based document image analysis](#). In *Document Analysis and Recognition – ICDAR 2021*, pages 131–146, Cham. Springer International Publishing. 742
743
744
745
746
747
748

Manan Suri, Puneet Mathur, Nedim Lipka, Franck Dernoncourt, Ryan A. Rossi, Vivek Gupta, and Dinesh Manocha. 2025. [Follow the flow: fine-grained flowchart attribution with neurosymbolic agents](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 22485–22508. Association for Computational Linguistics. 749
750
751
752
753
754
755
756

Simon Tannert, Amir Sadafi, and Klemens Böhm. 2023. [Flowchartqa: The first large-scale benchmark for reasoning over flowcharts](#). In *Proceedings of the 1st Workshop on Linguistic Insights from and for Multimodal Language Processing*. 757
758
759
760
761

Abhay Zala, Han Lin, Jaemin Cho, and Mohit Bansal. 2023. [Diagrammrgpt: Generating open domain diagrams via diagram understanding and text-to-diagram generation](#). *arXiv preprint arXiv:2310.12128*. *ArXiv:2310.12128*. 762
763
764
765
766

Rowan Zellers, Mark Yatskar, Sam Thomson, and Yejin Choi. 2018. [Neural motifs: Scene graph parsing with global context](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 767
768
769
770
771

R. Zhang, D. Jiang, Y. Zhang, and 1 others. 2024. [Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems?](#) In *European Conference on Computer Vision, ECCV*, pages 169–186, Cham. Springer Nature Switzerland. 772
773
774
775
776

Yilun Zhao, Chengye Wang, Chuhan Li, and Arman Cohan. 2025. [Can multimodal foundation models understand schematic diagrams? an empirical study on information-seeking qa over scientific papers](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 18598–18631, Vienna, Austria. Association for Computational Linguistics. 777
778
779
780
781
782
783

Xu Zhong, Jianbin Tang, and Antonio Jimeno Yepes. 2019. [PubLayNet: Largest dataset ever for document layout analysis](#). In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, pages 1015–1022. 784
785
786
787
788

789 A Appendix A: Evaluation Metrics

790 In this section, we provide the formal mathemat-
 791 ical definitions for the metrics discussed in Sec-
 792 tion 4.2. Following the unified Precision/Recall/F1
 793 framework defined in Eq. 1, we specify the truth
 794 sets (A) and predicted sets (\hat{A}) for each task-
 795 specific instantiation to ensure a comprehensive
 796 evaluation of both structural and semantic parsing.

797 **VLH Node-type macro-F1.** We measure node
 798 typing as a per-class set matching problem. For
 799 each node type $c \in \mathcal{C}$, define

$$800 \begin{aligned} A_c &:= \{n \in \mathcal{N} : y(n) = c\}, \\ \hat{A}_c &:= \{n \in \mathcal{N} : \hat{y}(n) = c\}. \end{aligned} \quad (2)$$

801 We compute $F1_c = F1(A_c, \hat{A}_c)$ by Eq. (1c) and
 802 report the macro average:

$$803 \text{Node-F1} = \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} F1_c. \quad (3)$$

804 **VLH Edge-F1.** Edge-F1 evaluates recovery of
 805 directed parent–child relations. Let $A = \mathcal{E}$ and
 806 $\hat{A} = \hat{\mathcal{E}}$; Edge-F1 is $F1(A, \hat{A})$ by Eq. (1c).

807 **SG Triplet-F1.** Each SG item is a triplet (h, r, t) .
 808 Let $A = \mathcal{T}^{SG}$ and $\hat{A} = \hat{\mathcal{T}}^{SG}$; SG Preci-
 809 sion/Recall/F1 are computed by Eq. (1) under ex-
 810 act triplet matching after mention normalization.

811 **SSG Triplet-F1 with node attachment.** Each
 812 SSG item is (h, r, t, u_h, u_t) , where u_h and u_t
 813 are the VLH node ids to which the head and tail men-
 814 tions are attached. Let $A = \mathcal{T}^{SSG}$ and $\hat{A} =$
 815 $\hat{\mathcal{T}}^{SSG}$; SSG Precision/Recall/F1 are computed by
 816 Eq. (1). This metric is stricter than SG because a
 817 prediction is correct only when both the semantic
 818 triplet (h, r, t) and the node attachments (u_h, u_t)
 819 match exactly.

820 B Appendix B: Visual Layout Hierarchy 821 Induction Algorithm

Algorithm 1 Layout-induced VLH construction
 (rule-based, annotation-driven)

Require: Diagram image I , annotated regions \mathcal{R}
 with types and boxes, fixed decision rules \mathcal{D}

Ensure: VLH tree \mathcal{T}

- 1: $\mathcal{P} \leftarrow \{r \in \mathcal{R} \mid \text{type}(r) \in \{\text{TextBlock}, \text{Icon}, \text{Image}, \text{Connector}\}\}$
 - 2: $\mathcal{G} \leftarrow \{r \in \mathcal{R} \mid \text{type}(r) = \text{Branch node}\}$
 - 3: **if** ISATOMICDIAGRAM(\mathcal{P}, \mathcal{D}) **then**
 - 4: Create a single root node v with $\text{type}(v) \in \{\text{Image}, \text{Icon}\}$ covering the diagram
 - 5: **return** $\mathcal{T} \leftarrow \{v\}$
 - 6: **end if**
 - 7: Create root branch node G_0 covering the full diagram; initialize $\mathcal{T} \leftarrow \{G_0\}$
 - 8: $\mathcal{G} \leftarrow \text{MERGEANDRESOLVECONFLICTS}(\mathcal{G}, \mathcal{D})$ \triangleright merge near-duplicates, resolve heavy overlaps by fixed rules
 - 9: BUILDBRANCHBACKBONE($\mathcal{T}, G_0, \mathcal{G}, \mathcal{D}$) \triangleright assign parent by minimal geometric containment
 - 10: **for all** $p \in \mathcal{P}$ **do**
 - 11: $a \leftarrow \text{SMALLESTCONTAININGBRANCH}(p, \mathcal{G} \cup \{G_0\})$
 - 12: Add p as a child of a in \mathcal{T}
 - 13: **end for**
 - 14: **if** NORELIABLEBRANCHNODES(\mathcal{G}, \mathcal{D}) **then**
 - 15: Remove internal branch node structure; attach all $p \in \mathcal{P}$ directly under G_0
 - 16: **end if**
 - 17: **Stop rule:** never expand terminal primitives \triangleright VLH depth is not fixed; expansion ends at primitives
 - 18: **return** \mathcal{T}
-

Table B.1: Benchmark coverage across four diagram genres. Counts are curated from an initial pool of 18760 candidates extracted from approximately 30000 papers.

Diagram genre	English	Chinese	Total	Typical cues
Process Diagrams	170	170	340	stage blocks, ordered arrows, step titles
Structural Diagrams	155	155	310	part to whole nesting, component labels, callouts
Interaction Diagrams	135	135	270	relation-dense arrows, typed links, direction cues
Composite Diagrams	140	140	280	multi-module layout, panel boundaries, mixed substructures
Total	600	600	1200	

Table B.2: Question answering statistics and verification constraints. The reported counts are placeholders and will be updated with the final benchmark release.

QA category	Per-diagram avg.	Total	Answer form	Verification rule
Q1 Entity	2.0	2400	entity mention string	exact match to entity and evidence span
Q2 Relation	2.0	2400	relation label	exact match to relation and evidence span
Q3 Hierarchical trace	1.0	1200	module identifier plus path	exact substrings match in allowed sources, full path completeness
Q4 Hierarchy helpful	1.0	1200	short answer or matchable span	exact substrings match in allowed sources, module-consistent evidence
Total	8.0	7200		

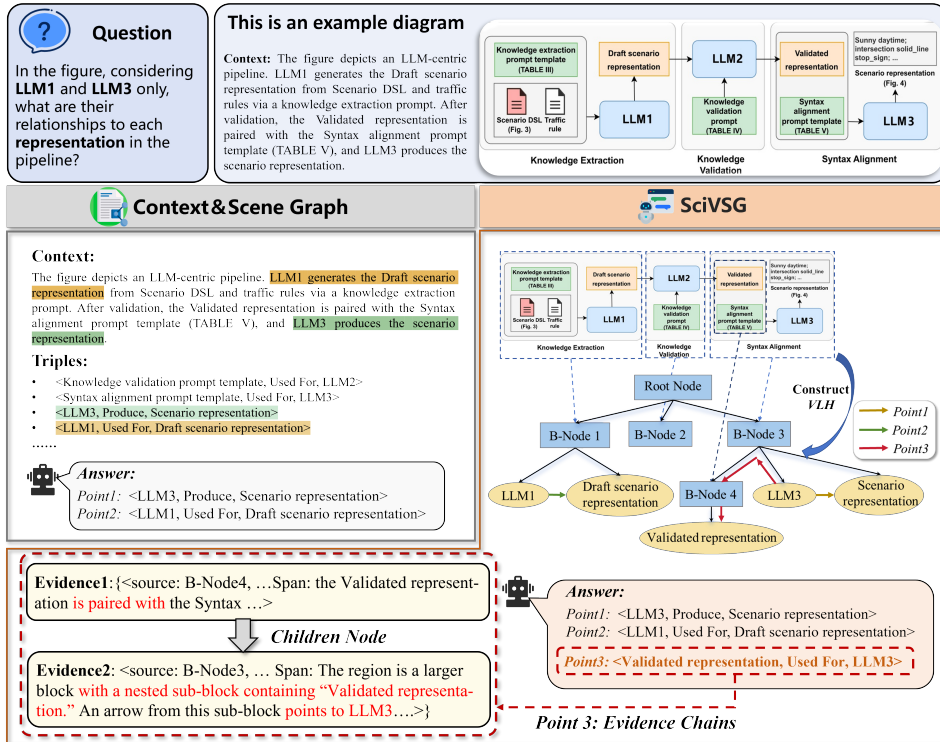


Figure B.1: Illustration of the traceable QA process in SciVSG. **Qualitative comparison:** We compare a Context & Scene Graph representation with SciVSG on an example diagram-based QA instance. Under the same question, the context-plus-graph setting recovers only Point1 and Point2, supported by textual cues and the corresponding extracted relations. In contrast, SciVSG additionally enables Point3 by exploiting the module-level hierarchy induced from diagram layout, which constrains reasoning to the correct structural scope and exposes a module-specific dependency missed by the baseline. Moreover, SciVSG offers fine-grained traceability: each predicted point is linked to the relevant VLH module, localized to the corresponding diagram region, and justified by an explicit evidence span aligned from node-attached text sources.