

# X-R2: Explainable Research Reasoning with Knowledge Graphs for Scientific Reasoning

Anonymous ACL submission

## Abstract

Large language models (LLMs) are increasingly used for scientific question answering, yet genomics reasoning still suffers from limited reasoning capability and insufficient explainability. Moreover, under a controlled knowledge-increment setting, it remains unclear which knowledge channel is truly necessary to obtain correct answers. To address these challenges, we propose **X-R2**, an iterative reasoning framework that performs mandatory question decomposition given a query  $Q$  and optional document evidence  $\mathcal{D}$ , while explicitly separating three additional knowledge channels: (i) entity-centric knowledge, (ii) KG triple knowledge, and (iii) parametric knowledge. At each iteration, X-R2 extracts step-level entities and constraints, acquires evidence from the enabled channels, and triggers re-decomposition when a self-check detects missing or inconsistent support, thereby improving both accuracy and explainability.

To enable fine-grained attribution, we construct **X-R2 Bench**. Experiments on this benchmark show that X-R2 consistently improves performance over direct one-shot generation and confirm that instance-aligned KG triples provide the largest marginal gains.

## 1 Introduction

This paper focuses on improving the **accuracy** and **explainability** of large language models in knowledge-intensive genomics reasoning tasks. We propose **X-R2** (explainable research reasoning with knowledge graphs for scientific reasoning), an iterative reasoning framework that performs multi-step inference over evidence, explicitly satisfies implicit constraints, and recovers missing intermediate facts when needed, thereby producing **correct and verifiable** answers.

A central challenge in genomics reasoning is that a question may be answerable only when

the system can identify which knowledge source is decisive for correctness on that specific instance. However, existing reasoning and retrieval-augmented systems are often evaluated end-to-end as monolithic pipelines, making it difficult to conduct controlled analysis or attribute failures to missing evidence versus reasoning gaps.

To enable such controlled analysis, we construct **X-R2 Bench**, a genomics research reasoning benchmark that couples each question with (i) step-level decompositions, (ii) aligned KG(Knowledge Graphs) triples (KG triple knowledge), and (iii) curated entity-centric resources. Each instance can be executed under knowledge switches that selectively enable/disable access to entity-centric knowledge, KG triple knowledge, and parametric knowledge, thereby revealing when a method truly needs each component. We consider document evidence  $\mathcal{D}$  when available and three additional knowledge channels: entity-centric knowledge, KG triple knowledge, and parametric knowledge. Reliable scientific reasoning requires not only access to these channels, but also an explicit mechanism to progressively expose and satisfy hidden constraints during inference.

Figure 1 illustrates why this design is necessary: compared with single-step parametric answering, one-shot KG-augmented reasoning, and unconstrained multi-step chain-of-thought reasoning, X-R2 introduces an inference-time control structure that progressively exposes constraints, acquires evidence step by step across knowledge channels, and revises the plan when evidence is insufficient or inconsistent, improving both accuracy and traceable attribution. Our contributions are as follows:

- **X-R2 reasoning framework.** We propose an iterative inference framework for scientific reasoning in which **question decomposition is mandatory**, followed by step-specific ex-

083	traction of entities and relations/constraints,	reasoning (Trivedi et al., 2023; Lee et al., 2024b;	131
084	controlled KG-triple retrieval and parametric-	Lyu et al., 2024; Asai et al., 2024; Li et al., 2025b),	132
085	knowledge elicitation, step reasoning, and	and examine when retrieved context is sufficient	133
086	self-check-driven re-decomposition.	or insufficient to support a correct answer, as well	134
087	• <b>X-R2 Bench.</b> We introduce <b>X-R2 Bench</b> ,	as how RAG systems should abstain or remain ro-	135
088	a genomics research reasoning benchmark	burst under retrieval errors (Maekawa et al., 2024;	136
089	with step-level decompositions, aligned KG	Yan et al., 2024; Hagström et al., 2025; Joren	137
090	triple evidence, and curated entity-centric re-	et al., 2025; Peng et al., 2025). However, these	138
091	sources. The benchmark provides knowl-	methods are still commonly evaluated as mono-	139
092	edge switches for controlled ablations across	lithic pipelines, making controlled knowledge in-	140
093	knowledge channels, enabling fine-grained	puts and fine-grained attribution difficult.	141
094	attribution of when each source is necessary	<b>Which knowledge source matters in genomics</b>	142
095	for correctness.	<b>reasoning?</b> Beyond end-to-end accuracy, scienti-	143
096	• <b>Knowledge weight analysis.</b> We conduct an	fic QA benefits from identifying which knowl-	144
097	empirical study comparing iterative reason-	edge source is necessary for correctness. Prior	145
098	ing with direct generation baselines under dif-	work studies parametric vs. non-parametric mem-	146
099	ferent knowledge availability settings, analyz-	ories and evaluates grounding and context utiliza-	147
100	ing both accuracy and knowledge-source nec-	tion (Mallen et al., 2023; Lee et al., 2024a; Shen	148
101	essity.	et al., 2024; Hagström et al., 2025; Maekawa et al.,	149
102	<b>2 Related Work</b>	2024; Peng et al., 2025), and recent hybrid settings	150
103	<b>Multi-step reasoning and error accumulation.</b>	combine textual and relational evidence (Zhu et al.,	151
104	Recent work improves controllability by making	2025; Lee et al., 2025). However, controlled attri-	152
105	inference more explicit via multi-step decom-	bution across entity-centric knowledge, KG triple	153
106	position, iterative reasoning, or structured search	knowledge, and parametric knowledge remains un-	154
107	(Trivedi et al., 2023; Yao et al., 2023; Shinn et al.,	derexplored in genomics, partly due to the lack of	155
108	2023; Madaan et al., 2023; Lee et al., 2024b; Li	benchmarks with aligned evidence and switchable	156
109	et al., 2025a). A key limitation is error accumu-	knowledge access. We address this gap with X-R2	157
110	lation: early mistakes in decomposition or inter-	Bench and knowledge switches for instance-level	158
111	mediate conclusions can propagate and dominate	necessity analysis.	159
112	later steps. X-R2 addresses this by coupling step-	<b>3 Knowledge Taxonomy &amp; Problem</b>	160
113	wise reasoning with explicit self-checking and re-	<b>Setup</b>	161
114	vision, enabling re-decomposition when the cur-	<b>3.1 Knowledge Taxonomy</b>	162
115	rent state is unreliable.	Scientific questions often hide multiple sub-goals	163
116	<b>External knowledge augmentation and harm-</b>	and constraints. We treat the document context $\mathcal{D}$	164
117	<b>ful context.</b> RAG and related approaches aug-	(e.g., an abstract) as input evidence: it provides	165
118	ment LLMs with external context, but retrieval can	instance-specific textual evidence and should be	166
119	be noisy or irrelevant and may even hurt perfor-	prioritized for grounding when available. Beyond	167
120	mance (Yoran et al., 2023; Maekawa et al., 2024;	$\mathcal{D}$ , we formalize three additional knowledge chan-	168
121	Shen et al., 2024; Hagström et al., 2025; Peng	nels that can support inference and enable con-	169
122	et al., 2025; Li et al., 2025a). This motivates step-	trolled attribution:	170
123	wise evidence acquisition and robustness under re-	(i) <b>Entity-centric knowledge</b> provides	171
124	trieval errors. X-R2 retrieves/elicits knowledge	lightweight signals for entity grounding, includ-	172
125	step by step and revises the plan when evidence	ing canonical names, aliases/synonyms, entity	173
126	is insufficient or inconsistent.	types, and simple lexical or type constraints (e.g.,	174
127	<b>Planning, iterative retrieval, and context suffi-</b>	gene vs. locus vs. phenotype). It is used to resolve	175
128	<b>ciency.</b> To improve grounded reasoning, recent	mentions in $(Q, \mathcal{D})$ into stable entity identifiers	176
129	studies explore planning-then-generation and step-	and to constrain retrieval queries and step-level	177
130	wise (iterative) retrieval for knowledge-intensive	reasoning.	178
		(ii) <b>KG triple knowledge</b> is relational evidence	179
		retrieved from an external knowledge graph in	180

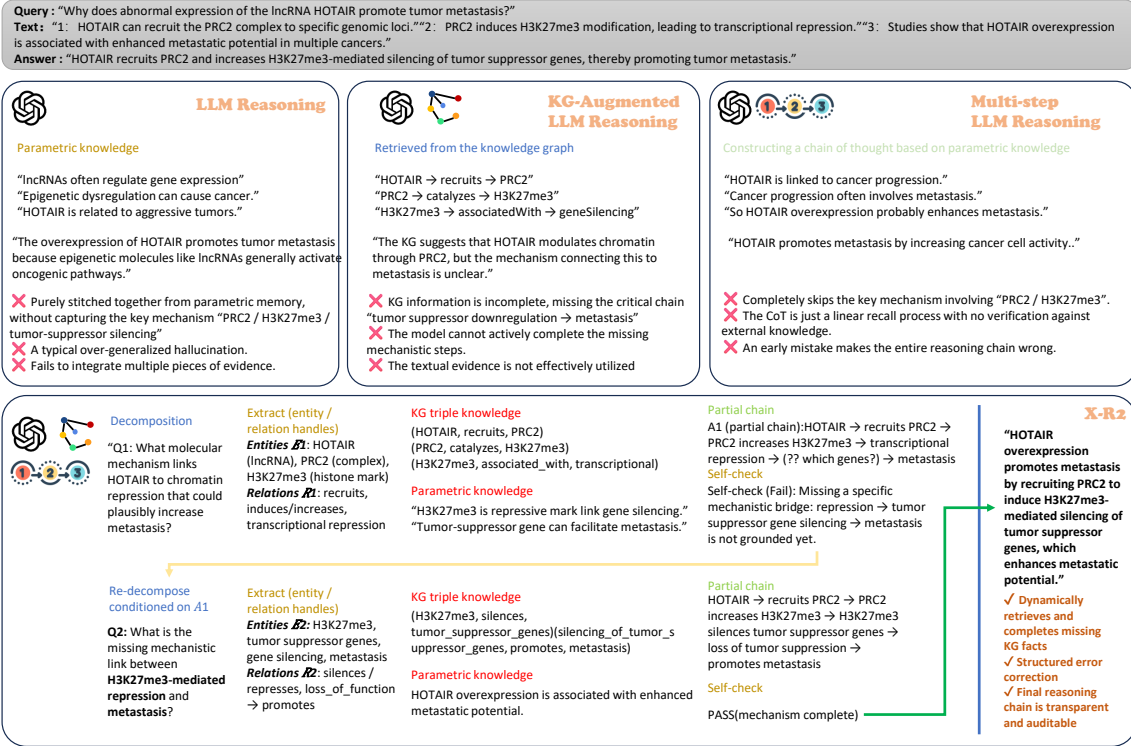


Figure 1: **Motivating example for X-R2.** Scientific reasoning often requires completing missing links between evidence and the final claim. Pure parametric answering may overgeneralize; one-shot KG augmentation can retrieve correct but disconnected triples; and vanilla CoT may drift without verification. X-R2 addresses these failure modes with mandatory decomposition, step-wise entity/relation extraction, controlled KG-triple retrieval and parametric-knowledge elicitation, and self-check-driven re-decomposition.

the form of triples (or short paths) anchored on grounded entities. It complements  $\mathcal{D}$  by supplying structured relations (e.g., INTERACTS\_WITH, CAUSES, ASSOCIATED\_WITH) that may be missing, implicit, or scattered across the document evidence.

(iii) **Parametric knowledge** refers to background facts and patterns stored in model weights. While it can fill in missing intermediate facts when external evidence is incomplete, it is also the least verifiable channel; therefore, we treat it as a controlled, optional source that can be explicitly enabled/disabled (or elicited via a dedicated prompt) for analysis.

Together,  $\mathcal{D}$  and the three channels define the knowledge space used by X-R2 and underpin the knowledge switches used throughout our controlled evaluation.

### 3.2 Problem Setup

We define the task and evaluation protocol shared by X-R2 and all baselines.

**Task.** Given a query  $Q$  and an document context  $\mathcal{D}$ , the system produces a prediction  $\hat{A}$ . When document access is enabled ( $s_{\text{Doc}}=1$ ),  $\hat{A}$  should be grounded in  $\mathcal{D}$  under the corresponding setting. Each instance also includes a reference answer  $A^*$ , and performance is measured by comparing  $\hat{A}$  against  $A^*$ .

**Available inputs.** The mandatory input is  $Q$ , with  $\mathcal{D}$  provided when available. In addition, an instance may include optional resources for controlled analysis: (i) entity-centric knowledges (canonical names, types, and aliases), and (ii) instance-aligned KG evidence as a small set of KG triples (or short paths). These resources are packaged by the benchmark. When KG access is enabled ( $s_{\text{KG}}=1$ ), the system may also retrieve additional triples from an external KG store.

**Output and evaluation format.** We normalize  $\hat{A}$  into a short-answer form (e.g., a list of canonical strings) for deterministic evaluation, and compute exact-match/normalized-match scores against  $A^*$ . When evidence is insufficient under a given set-

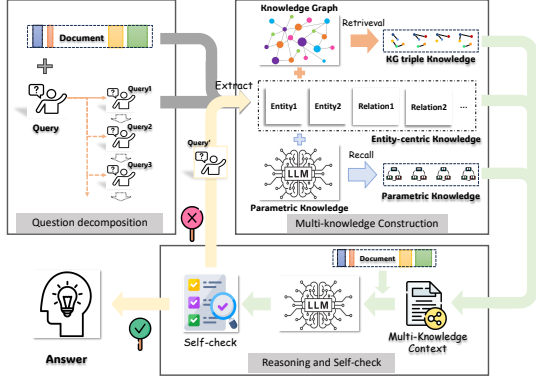


Figure 2: **X-R2 overview.** Given a query  $Q$  and optional document evidence  $\mathcal{D}$ , X-R2 performs mandatory question decomposition, X-R2 performs mandatory question decomposition to expose step-level sub-questions. At each step, the model extracts entity/relation cues, retrieves KG triple knowledge from an external knowledge graph, and (optionally) elicits parametric knowledge via a controlled parametric-recall prompt. The model then conducts step reasoning and runs a self-check to decide whether to stop and output the final answer or to re-decompose and continue the next iteration.

ting, the model may abstain rather than speculate; under accuracy-style metrics, abstentions are counted as incorrect.

#### Controlled analysis via knowledge switches.

To attribute behavior to different knowledge components, we evaluate each instance under binary knowledge switches that independently enable/disable access to: document evidence ( $s_{\text{DOC}}$ ), entity-centric knowledge ( $s_{\text{ENT}}$ ), KG triple knowledge ( $s_{\text{KG}}$ ), and parametric knowledge ( $s_{\text{PAR}}$ ). Comparing results across switch configurations reveals which knowledge source(s) are necessary for correctness on that instance.

## 4 Method

Figure 2 summarizes the end-to-end inference loop of X-R2. At a high level, X-R2 iterates over decompose  $\rightarrow$  extract  $\rightarrow$  retrieve/recall  $\rightarrow$  reason  $\rightarrow$  self-check. In each iteration, entity-centric signals guide KG retrieval and optional parametric elicitation, after which a self-check determines whether the current intermediate answer is sufficient or the system should re-decompose and continue.

### 4.1 Iterative Multi-step Inference

At iteration  $t$  with state  $A_{t-1}$  ( $A_0 = \emptyset$ ), X-R2 executes: Decompose  $\rightarrow$  Extract  $\rightarrow$  RetrieveKG  $\rightarrow$

Recall  $\rightarrow$  Reason  $\rightarrow$  Check. Concretely,  $Q_t$ ,  $(E_t, R_t)$ ,  $KG_t$ ,  $M_t$ ,  $A_t$ ,  $\text{stop}_t$  are produced in order from  $(Q, \mathcal{D}, A_{t-1})$ ; if  $\text{stop}_t=1$ , the system returns  $\hat{A}=A_t$ , otherwise it continues with  $A_t$  and re-invokes decomposition.

**Question decomposition.** Given  $(Q, \mathcal{D}, A_{t-1})$ , the model produces a decomposition object

$$\text{Decomp}_t = \{\text{subquestions}, \text{constraints}, \text{required\_evidence}, \text{current}\}, \quad (1)$$

where current specifies the step target  $Q_t$ . For  $t = 1$ , the model decomposes from  $(Q, \mathcal{D})$ . For  $t > 1$ , it re-decomposes conditioned on  $(A_{t-1}, Q, \mathcal{D})$  to generate refined follow-up questions that address remaining gaps.

**Step-specific extraction.** Conditioned on  $(Q_t, \mathcal{D})$ , the model extracts step-specific entities  $E_t$  and relation/constraint cues  $R_t$ :

$$(E_t, R_t) = g(Q_t, \mathcal{D}). \quad (2)$$

Here  $R_t$  may include relation types, directions, conditions, required attributes, or comparison operators, serving as structured handles for downstream knowledge acquisition.

**KG triple retrieval.** Using  $(E_t, R_t)$  as retrieval intent, X-R2 collects KG triple knowledge  $KG_t$  as a compact set of triples (and optionally short paths) from an external knowledge graph.

**Controlled parametric recall.** In parallel to KG retrieval, X-R2 may elicit related parametric knowledge induced by  $(E_t, R_t)$ :

$$M_t = h(E_t, R_t). \quad (3)$$

Operationally,  $M_t$  is obtained via a controlled parametric-recall prompt that encourages the model to surface relevant background knowledge while keeping it separable from retrieved KG triples and document evidence.

**Step reasoning and self-check.** The model produces an intermediate answer by reasoning over the combined context:

$$A_t = f(Q_t, \mathcal{D}, E_t, R_t, KG_t, M_t). \quad (4)$$

After generating  $A_t$ , the model performs a self-check:

$$\text{stop}_t = \text{Check}(A_t, Q, \mathcal{D}), \quad \text{stop}_t \in \{0, 1\}, \quad (5)$$

---

**Algorithm 1** X-R2: Iterative multi-step multi-source reasoning

---

**Require:** Query  $Q$ , document  $\mathcal{D}$ , max iterations  $T_{\max}$

- 1:  $A_0 \leftarrow \emptyset$
- 2: **for**  $t = 1$  **to**  $T_{\max}$  **do**
- 3:    $\text{Decomp}_t \leftarrow \text{Decompose}(Q, \mathcal{D}, A_{t-1})$  // mandatory
  
- 4:    $Q_t \leftarrow \text{Decomp}_t.\text{current}$
- 5:    $(E_t, R_t) \leftarrow \text{Extract}(Q_t, \mathcal{D})$
- 6:    $KG_t \leftarrow (\text{s\_KG} = 1) ? \text{RetrieveKG}(E_t, R_t) : \emptyset$
- 7:    $M_t \leftarrow (\text{s\_Par} = 1) ? \text{RecallParametric}(E_t, R_t) : \emptyset$
- 8:    $A_t \leftarrow \text{Reason}(Q_t, \mathcal{D}, E_t, R_t, KG_t, M_t)$
- 9:    $\text{stop}_t \leftarrow \text{Check}(A_t, Q, \mathcal{D})$
- 10:   **if**  $\text{stop}_t = 1$  **then**
- 11:     **return**  $\hat{A} \leftarrow A_t$
- 12:   **end if**
- 13: **end for**
- 14: **return**  $\hat{A} \leftarrow A_{T_{\max}}$

---

291 deciding whether  $A_t$  sufficiently answers the orig-  
292 inal query and is consistent with  $\mathcal{D}$ . If  $\text{stop}_t = 1$ ,  
293 X-R2 terminates and outputs  $\hat{A} = A_t$ ; otherwise it  
294 continues to iteration  $t+1$  by re-invoking decom-  
295 position. We cap the maximum number of itera-  
296 tions  $T_{\max}$ ; if reached, the system outputs the lat-  
297 est  $A_t$ .

298 **Algorithm summary.** Alg. 1 summarizes the itera-  
299 tive multi-step inference loop.

## 5 Benchmark Construction

301 To support controlled analysis of iterative multi-  
302 step reasoning with heterogeneous knowledge  
303 sources, we construct **X-R2 Bench**, where each  
304 instance is packaged with (i) an explicit decom-  
305 position and step trace, (ii) step-level entity/rela-  
306 tion seeds for KG retrieval, and (iii) an instance-  
307 aligned KG evidence subgraph with provenance.  
308 The resulting instances are directly consumable by  
309 our inference scripts.

### 5.1 Data Sources and Target Format

310 We start from a QA dataset in JSONL format,  
311 where each record contains an identifier, a ques-  
312 tion  $Q$ , a document context  $\mathcal{D}$  (in our setup, an  
313 abstract), and a reference answer  $A^*$ . We also as-  
314 sume an external KG stored as JSONL triples with  
315 fields such as subject, relation, and object.

316 Each benchmark instance is a single JSON ob-  
317 ject with three core blocks:

- 318 • **Input/answer:** input (question and document)  
319 and reference\_answer.
- 320
- 321 • **Reasoning interface:** decomposition and a  
322 step-aligned step\_trace (same length), plus

---

**Algorithm 2** X-R2 Bench construction (decompo-  
sition + KG alignment)

---

**Require:** QA records  $\mathcal{R}$ , KG triples  $\mathcal{G}$ , top- $K$ , max fix  
rounds  $F$

- 1: Build an inverted index over normalized entity strings in  
 $\mathcal{G}$
- 2: **for** each record  $r \in \mathcal{R}$  **do**
- 3:    $(\text{decomp}, \text{seeds}, \text{kg\_query}) \leftarrow \text{LLM\_Extract}(Q, \mathcal{D})$
  
- 4:    $\text{kg\_hits} \leftarrow \text{RetrieveSubgraph}(\mathcal{G}, \text{seeds}, K)$
- 5:    $x \leftarrow \text{LLM\_Build}(r, \text{decomp}, \text{seeds}, \text{kg\_hits})$
- 6:   **if not**  $\text{Validate}(x)$  **then**  $x \leftarrow \text{LLM\_Fix}(x)$  (repeat up  
to  $F$ )
- 7:   Write  $x$  to benchmark JSONL (preserving  $A^*$  exactly)
- 8: **end for**

---

entity\_slot for canonical entities and rela- 323  
tion/constraint cues. 324

- **Aligned evidence:** kg.subgraph containing a 325  
sentence map for  $\mathcal{D}$  and KG triples with prove- 326  
nance (doc sentence IDs or stable KG IDs). 327

We provide the complete schema, field definitions, 328  
and an example instance in Appendix A. 329

### 5.2 Construction Pipeline

Algorithm 2 summarizes a three-stage pipeline. 331

**Stage 1 (LLM): decomposition and retrieval** 332  
**seeds.** Given  $(Q, \mathcal{D})$ , an LLM produces 333  
a short decomposition and retrieval seeds 334  
(seed\_entities, seed\_relations) plus a 335  
concise kg\_query. 336

**Stage 2 (lexical): KG subgraph retrieval.** We 337  
retrieve top- $K$  KG triples by lexical matching over 338  
the KG store using the extracted seeds, and assign 339  
stable kg\_ids. 340

**Stage 3 (LLM): instance assembly with** 341  
**validation.** We prompt an LLM with the 342  
schema, the QA record, and retrieved KG 343  
hits to produce a complete instance JSON, 344  
and validate structural constraints (e.g., 345  
 $\text{len}(\text{step\_trace})=\text{len}(\text{decomposition})$ ), 346  
optionally running a small fix loop. 347

## 6 Experiments

348 We evaluate (i) the effectiveness of our iterative 349  
multi-step reasoning procedure, (ii) its generaliza- 350  
tion across different LLM backbones, and (iii) con- 351  
trolled attribution across document evidence and 352  
three knowledge channels (entity-centric signals, 353  
KG triple knowledge, and parametric knowledge). 354

## 6.1 Benchmark and Tasks

We conduct experiments on **X-R2 Bench**, a genomics-oriented scientific reasoning benchmark built from paper abstracts. Each instance is a JSON object that includes an input question  $Q$ , an abstract  $\mathcal{D}$ , a reference answer  $A^*$ , and structured fields used by our method, such as a decomposition plan (decomposition), step-level traces (step\_trace), entity slots (entity\_slot), and an aligned evidence subgraph (kg.subgraph) with sentence\_map and KG/doc triples.

**Scale.** X-R2 Bench contains **410** instances in total (each instance corresponds to one task type paired with one abstract). We use a single held-out evaluation split (no train/dev/test partition), as our focus is on controlled inference-time analysis.

**Task types.** Each instance belongs to one of five task types (stored in the task field):

- **methods\_and\_techniques:** core methods/technologies used in the study (up to two).
- **research\_field:** the research field/domain issue addressed by the study.
- **study\_type:** the study type/angle (e.g., basic research, database development).
- **datasets:** datasets/cohorts used (e.g., TCGA, 1000 Genomes).
- **genes\_and\_loci:** explicitly mentioned genes or loci.

**External KG.** We use an external domain KG with **98,799** triples. When enabled, the system can retrieve triples from this KG using step-level seed entities/relations.

## 6.2 Models, Baselines, and Settings

**Models.** We evaluate three representative LLMs: **Qwen3-8B** (Yang et al., 2025), **DeepSeek-V3** (DeepSeek-AI et al., 2024), and **ChatGPT-5.1** (OpenAI, 2025a,b). We use each model through a unified prompting interface with deterministic decoding (below).

### Methods compared.

- **Direct (one-shot).** A single-pass baseline that answers  $Q$  given the enabled context blocks.
- **Iterative (ours).** Iterative multi-step reasoning in dynamic mode, producing step-level subques-

tions  $Q_t$ , intermediate answers  $A_t$ , and a termination decision.

**Deterministic decoding.** All runs use temperature=0 and top\_p=0, with fixed settings for reproducibility.

**Iteration and KG context budget.** Unless otherwise stated, we set the maximum number of iterative steps to  $T_{\max} = 5$  (via `-max_steps 5`) and truncate the KG triples fed into the prompt to top- $K = 5$  triples (via `-max_kg_triples 5`). This keeps the evidence budget small and emphasizes step-wise evidence selection.

**Controlled switches.** Each instance includes four binary knowledge switches: `s_Doc` (document evidence), `s_Ent` (entity-centric resources), `s_KG` (KG triple knowledge), and `s_Par` (parametric knowledge). When `s_Par=1`, we enable a controlled parametric-recall prompting mechanism to elicit parametric knowledge; otherwise the parametric channel is disabled. Our implementation supports per-run overrides of these switches, which we use for ablations and bucket attribution (Sec. 6.5–6.7).

## 6.3 Evaluation Metrics

**Accuracy (task-aware matching).** We follow the exact matching logic implemented in our evaluation script. Both gold and predictions are normalized by lowercasing and whitespace cleanup. Predictions are required to be `list[str]` (a JSON array of strings). Correctness is computed by task:

- **genes\_and\_loci, datasets, methods\_and\_techniques:** we mark an instance correct iff the predicted set is non-empty and covers the gold set (i.e.,  $\text{Gold} \subseteq \text{Pred}$ ).
- **research\_field and study\_type:** we mark an instance correct iff there is non-empty overlap between normalized strings, or one side is a substring of the other after concatenation (robust to minor phrasing differences).

**Abstain rate.** We report the fraction of instances where the system abstains (empty `pred_answer` or explicit `abstain` field), which reflects evidence insufficiency under stricter switch settings.

**Average steps.** For the iterative method, we report the average number of executed steps until

446 termination, computed from the length of the pro- 494  
447 duced steps / stop decision in the dynamic JSON 495  
448 output. 496

449 **Bucketed accuracy.** To attribute which knowl- 497  
450 edge channel is necessary, we define three buck- 498  
451 ets using controlled re-runs: needs-KG (KG triple 499  
452 knowledge), needs-Entity (entity-centric knowl- 500  
453 edge), and needs-Parametric (parametric knowl- 501  
454 edge). We report bucketed accuracy as well as the 502  
455 bucket sizes.

## 456 6.4 Main Results

457 Table 1 compares iterative reasoning with di- 503  
458 rect one-shot generation on X-R2 Bench. For 504  
459 ChatGPT-5.1, iterative reasoning improves overall 505  
460 accuracy from 61.8 to 68.9 (+7.1), while keeping 506  
461 abstention low (1–2%). The average number of ex- 507  
462 ecuted steps is 2.84, indicating that the procedure 508  
463 typically converges early and does not require the 509  
464 full  $T_{\max}$  budget.

465 **Discussion.** The consistent gain suggests that 510  
466 explicitly surfacing latent constraints via manda- 511  
467 tory decomposition, and then selecting evidence 512  
468 step-by-step, mitigates overgeneralization in ge- 513  
469 nomics QA. Notably, the abstain rate increases 514  
470 only slightly (+0.5), indicating that the higher ac- 515  
471 curacy is not obtained by simply refusing harder 516  
472 instances.

## 473 6.5 Ablation and Attribution via Controlled 522 474 Switches

475 We perform ablations by overriding the four 523  
476 switches at inference time, while keeping de- 524  
477 coding deterministic and budgets fixed. Specifi- 525  
478 cally, we test: (i) disabling KG triple knowledge 526  
479 ( $s_{\text{KG}}=0$ ), (ii) disabling entity-centric resources 527  
480 ( $s_{\text{Ent}}=0$ ), (iii) disabling parametric knowledge 528  
481 ( $s_{\text{Par}}=0$ ), (iv) reducing the iterative depth to 529  
482  $T_{\max} = 1$  (to isolate the effect of iteration 530  
483 beyond mandatory decomposition), and (v) **KG 531  
484 shuffled**—permuting retrieved triples across 532  
485 instances while preserving the same triple budget, to 533  
486 stress-test whether gains require correct instance- 534  
487 level KG alignment. 535

488 **What contributes most?** Disabling KG triple 536  
489 knowledge causes the largest drop (68.9  $\rightarrow$  537  
490 59.7, -9.2), larger than removing entity slots (-4.4) or 538  
491 parametric recall (-2.1). This indicates that, under 539  
492 a fixed evidence budget (top- $K=5$ ), correctly re- 540  
493 trieved KG triples provide the strongest marginal 541

benefit in X-R2 Bench.

**Is KG alignment real or incidental?** The **KG 495  
496 shuffled** control largely removes the gain (61.1 vs 497  
498 68.9), even though the prompt still contains the 499  
500 same number of triples. This suggests the improve- 501  
502 ment is not a prompt-length artifact: it depends on 503  
504 instance-aligned KG evidence, supporting the cor- 505  
506 rectness/validity of X-R2 Benchs KG alignment 506  
507 interface. 507

**Bucketed accuracy.** Using the controlled re- 503  
504 runs described in Sec. 6.3, we further report buck- 504  
505 eted performance. This directly quantifies when 505  
506 correctness depends on KG triple knowledge, 506  
507 entity-centric knowledge, or parametric knowl- 507  
508 edge, under the same instance-aligned interface. 508

**Interpretation.** The largest bucket is needs-KG 509  
510 (186 instances), and it also exhibits the largest 510  
511 within-bucket gain (+16.8), consistent with the 511  
512 ablation that KG triples are the dominant con- 512  
513 tributor. The needs-Entity bucket remains sub- 513  
514 stantial, highlighting that entity-centric ground- 514  
515 ing signals are often a prerequisite for effective 515  
516 KG retrieval and constraint satisfaction. Finally, 516  
517 needs-Parametric captures cases where safe 517  
518 background recall is still necessary (e.g., canon- 518  
519 ical expansions, common aliases, or domain con- 519  
520 ventions not present in the abstract/KG context un- 520  
521 der the top- $K$  budget). 521

## 522 6.6 Generalization Across LLMs

523 To test whether improvements are model-specific, 523  
524 we evaluate both methods on three different LLMs: 524  
525 Qwen3-8B (Yang et al., 2025), DeepSeek-V3 525  
526 (DeepSeek-AI et al., 2024), and ChatGPT-5.1 526  
527 (OpenAI, 2025a,b). We keep the same evaluation 527  
528 protocol,  $T_{\max} = 5$ , and top- $K = 5$  KG context. 528

**Discussion.** We observe consistent gains of +7– 529  
530 8 points across all backbones, suggesting the ben- 530  
531 efit comes from the inference procedure (decom- 531  
532 position + step-wise evidence use) rather than 532  
533 a model-specific artifact. The largest absolute 533  
534 improvement appears on the weaker backbone 534  
535 (Qwen3-8B), consistent with the hypothesis that 535  
536 structured intermediate supervision at inference 536  
537 time offers stronger regularization when the base 537  
538 model is less reliable in direct one-shot extraction. 538

## 539 6.7 Why X-R2 Bench is Necessary

540 A core goal of this work is to attribute when 540  
541 correctness requires different knowledge access 541

Method	Overall Acc.	Abstain ↓	Avg. Steps	Δ vs Direct
Direct (one-shot)	<b>61.8</b>	<b>0.7</b>	–	–
Iterative (ours)	<b>68.9</b>	<b>1.2</b>	<b>2.84</b>	<b>+7.1</b>

Table 1: Main results on X-R2 Bench (410 instances) using ChatGPT-5.1.

Setting (ours)	Overall Acc.	Abstain ↓	Avg. Steps	Notes
Full (default)	68.9	1.2	2.84	Doc+Entity+KG+Param.
w/o KG	59.7	1.5	3.05	s_KG=0
w/o Entity	64.5	1.3	2.96	s_Ent=0
w/o Parametric	66.8	1.2	2.91	s_Par=0
KG shuffled	61.1	1.3	2.98	permute triples across instances
$T_{\max}=1$	63.2	0.9	1.00	no multi-step refinement

Table 2: Ablation results via switch overrides (ChatGPT-5.1). “KG shuffled” breaks instance-level KG alignment while keeping the same evidence budget.

Bucket	#Instances	Direct Acc.	Iterative Acc.
needs-KG	186	41.4	58.2
needs-Entity	98	52.0	61.5
needs-Parametric	51	39.2	47.8

Table 3: Bucketed accuracy by controlled switch attribution (ChatGPT-5.1). Remaining 75 instances are doc-sufficient under our protocol.

Model	Direct Acc.	Iterative Acc.	Δ
Qwen3-8B	44.6	52.8	+8.2
DeepSeek-V3	56.3	64.1	+7.8
ChatGPT-5.1	61.8	68.9	+7.1

Table 4: Cross-model results on X-R2 Bench under the same protocol ( $T_{\max} = 5$ , top- $K = 5$ ).

modes. This is difficult to do reliably without (i) **aligned KG triples** at the instance level and (ii) **controlled switches** that can disable document evidence, entity-centric knowledge, KG triples, or parametric knowledge without changing the task interface.

**Aligned KG enables controlled attribution.** Because each instance includes an aligned evidence subgraph (kg.subgraph) with provenance (sentence\_map and KG IDs), we can toggle s\_KG or s\_Ent and re-run the same instances under identical prompting and budgets. This produces bucket labels such as needs-KG, needs-Entity, and needs-Parametric (Table 3) in a reproducible manner, which is not directly supported by standard QA benchmarks lacking KG alignment and switchable interfaces.

### Benchmark validity via alignment stress-test.

The **KG shuffled** control in Table 2 preserves the triple budget but destroys instance-level alignment. Its substantial degradation (68.9 → 61.1) empirically supports that the benchmarks KG signal is not incidental: correctly aligned KG triples are required to realize the full benefit, and incorrect alignment cannot be compensated for by simply providing more symbolic text.

**Practical impact.** Beyond reporting higher accuracy, X-R2 Bench allows us to answer why a method works: whether improvements come from better entity grounding, better retrieval of symbolic relational evidence, or safe use of background knowledge. This makes the evaluation actionable for designing reliable scientific reasoning systems.

## 7 Conclusions

We introduced **X-R2**, an iterative genomics reasoning framework that performs mandatory question decomposition and step-wise evidence use across document, entity-centric, KG-triple, and parametric knowledge channels. We also constructed **X-R2 Bench** with aligned evidence and binary knowledge switches (s\_Doc, s\_Ent, s\_KG, s\_Par) for controlled attribution. Experiments across three LLM backbones show consistent gains over direct one-shot generation, and ablations highlight the largest marginal benefit from instance-aligned KG triples.

## 589 Limitations.

590 Our controlled study centers on three knowledge  
591 channels (entity-centric signals, KG triples, and  
592 prompted parametric recall) and treats them as  
593 binary switches, which simplifies attribution but  
594 does not cover richer genomics knowledge re-  
595 sources, and the parametric channel remains the  
596 least verifiable. In addition, we evaluate on a  
597 single split of 410 abstract-based instances with  
598 task-aware normalized matching; this setting may  
599 not reflect performance on full-text evidence and  
600 can under-penalize over-prediction for list-style  
601 answers.

## 602 References

603 Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and  
604 Hannaneh Hajishirzi. 2024. [Self-rag: Learning to re-  
605 trieve, generate, and critique through self-reflection.](#)  
606 *In International Conference on Learning Representations (ICLR).*  
607

608 DeepSeek-AI, Aixin Liu, and 1 others. 2024.  
609 [Deepseek-v3 technical report.](#) *Preprint,*  
610 [arXiv:2412.19437.](#)

611 Lovisa Hagström, Sara Vera Marjanovic, Haeun Yu, Ar-  
612 nav Arora, Christina Lioma, Maria Maistro, Pepa  
613 Atanasova, and Isabelle Augenstein. 2025. [A reality  
614 check on context utilisation for retrieval-augmented  
615 generation.](#) *In Proceedings of the 63rd Annual Meet-  
616 ing of the Association for Computational Linguistics  
617 (Volume 1: Long Papers),* pages 19691–19730, Vi-  
618 enna, Austria. Association for Computational Lin-  
619 guistics.

620 Hailey Joren, Jianyi Zhang, Chun-Sung Ferng, Da-  
621 Cheng Juan, Ankur Taly, and Cyrus Rashtchian.  
622 2025. [Sufficient context: A new lens on retrieval  
623 augmented generation systems.](#) *In International  
624 Conference on Learning Representations (ICLR).*

625 Chanhee Lee, Zhaojun Chen, Luke Zettlemoyer, Wen-  
626 tau Yih, and Hannaneh Hajishirzi. 2024a. [How well  
627 do large language models truly ground?](#) *In Proceed-  
628 ings of the 2024 Conference of the North American  
629 Chapter of the Association for Computational Lin-  
630 guistics: Human Language Technologies (Volume 1:  
631 Long Papers),* pages 2412–2437, Mexico City, Mex-  
632 ico. Association for Computational Linguistics.

633 Deokhee Lee, Jinyoung Kim, Hyunwoo Park, Jiny-  
634 oung Lee, and Jinho Seo. 2025. [HybGRAG: Hybrid  
635 retrieval augmented generation in textual-relational  
636 contexts.](#) *In Proceedings of the 63rd Annual Meet-  
637 ing of the Association for Computational Linguistics  
638 (Volume 1: Long Papers),* pages 835–854, Vienna,  
639 Austria. Association for Computational Linguistics.

640 Myeonghwa Lee, Seonho An, and Min-Soo Kim.  
641 2024b. [PlanRAG: A plan-then-retrieval augmented](#)

[generation for generative large language models as  
642 decision makers.](#) *In Proceedings of the 2024 Con-  
643 ference of the North American Chapter of the Asso-  
644 ciation for Computational Linguistics: Human Lan-  
645 guage Technologies (Volume 1: Long Papers),* pages  
646 6537–6555, Mexico City, Mexico. Association for  
647 Computational Linguistics. 648

649 Jiacheng Li, Tianming Wang, Wenchang Xu, Yongfei  
650 Zhang, and Wenhu Chen. 2025a. [R3-RAG: Learn-  
651 ing to Retrieve, Reason and Refine for multi-hop  
652 question answering.](#) *In Findings of the Association  
653 for Computational Linguistics: EMNLP 2025,* pages  
654 9570–9587, Suzhou, China. Association for Compu-  
655 tational Linguistics.

656 Yuan Li, Qi Luo, Xiaonan Li, Bufan Li, Qinyuan  
657 Cheng, Bo Wang, Yining Zheng, Yuxin Wang,  
658 Zhangyue Yin, and Xipeng Qiu. 2025b. [R3-RAG:  
659 Learning step-by-step reasoning and retrieval for  
660 LLMs via reinforcement learning.](#) *In Findings of the  
661 Association for Computational Linguistics: EMNLP  
662 2025,* pages 10491–10507, Suzhou, China. Associa-  
663 tion for Computational Linguistics.

664 Yuanjie Lyu, Zihan Niu, Zheyong Xie, Chao Zhang,  
665 Tong Xu, Yang Wang, and Enhong Chen. 2024.  
666 [Retrieve-plan-generation: An iterative planning and  
667 answering framework for knowledge-intensive LLM  
668 generation.](#) *In Proceedings of the 2024 Conference  
669 on Empirical Methods in Natural Language Process-  
670 ing,* pages 4683–4702, Miami, Florida, USA. Asso-  
671 ciation for Computational Linguistics.

672 Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler  
673 Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,  
674 Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,  
675 Shashank Gupta, Bodhisattwa Prasad Majumder,  
676 Katherine Hermann, Sean Welleck, Amir Yazdan-  
677 bakhsh, and Peter Clark. 2023. [Self-refine: It-  
678 erative refinement with self-feedback.](#) *Preprint,*  
679 [arXiv:2303.17651.](#)

680 Seiji Maekawa, Hayate Iso, Sairam Gurajada, and  
681 Nikita Bhutani. 2024. [Retrieval helps or hurts? a  
682 deeper dive into the efficacy of retrieval augmenta-  
683 tion to language models.](#) *In Proceedings of the 2024  
684 Conference of the North American Chapter of the  
685 Association for Computational Linguistics: Human  
686 Language Technologies (Volume 1: Long Papers),*  
687 pages 5506–5521, Mexico City, Mexico. Associa-  
688 tion for Computational Linguistics.

689 Alex Mallen, Yue Li, Haotian Lin, Tushar Khot, Bo Hu,  
690 Shanya Li, Zhu Yao, Dan Wang, Diyi Yang, Percy  
691 Liang, and 1 others. 2023. [When not to trust lan-  
692 guage models: Investigating the effectiveness of  
693 parametric and non-parametric memories.](#) *In Pro-  
694 ceedings of the 61st Annual Meeting of the Associa-  
695 tion for Computational Linguistics (Volume 1: Long  
696 Papers),* pages 9847–9863, Toronto, Canada. Associa-  
697 tion for Computational Linguistics.

698 OpenAI. 2025a. [GPT-5.1: A smarter, more conversa-  
699 tional ChatGPT.](#) Accessed 2026-01-06.



- 797 prov="Sx", and (b) KG triples with  
798 source="kg" and prov=kg\_id.
- 799 • **Step trace aligned to decomposition:** each  
800 step includes the current sub-question,  
801 extracted entities/relations, used sentence IDs,  
802 an intermediate answer, and a self-check  
803 decision.
  - 804 • **Auxiliary fields (optional):**  
805 `closed_book_recall` (a few parametric claims  
806 with confidences), `labels_for_evaluation`  
807 (coverage targets and gold evidence sentence  
808 IDs), and a simple `judge_rubric`.

### 809 A.3 Validation Rules

810 To guarantee structural consistency, we validate  
811 each produced instance with:

- 812 • **Required keys:** presence of input,  
813 `reference_answer`, `decomposition`,  
814 `step_trace`, `kg.subgraph`.
- 815 • **Alignment constraint:**  
816 `len(step_trace)=len(decomposition)`.
- 817 • **Type checks:** `decomposition` is a list of  
818 strings; `step_trace` is a list of step objects;  
819 `kg.subgraph.triples` is a list of triples with  
820 source and prov.
- 821 • **Provenance sanity:** doc-grounded triples use  
822 sentence IDs appearing in `sentence_map`; KG  
823 triples use stable `kg_ids` returned by retrieval.

824 When validation fails, we run a bounded fix  
825 loop (up to  $F$  rounds) by prompting the LLM  
826 with the error message and the partially formed in-  
827 stance.

### 828 A.4 Retrieval Details (Stage 2)

829 Stage 2 performs best-effort lexical retrieval  
830 using the extracted `seed_entities` and  
831 `seed_relations`: (i) exact match on normalized  
832 entity strings, (ii) fallback partial containment  
833 match over subjects/objects, and (iii) a relation-  
834 keyword filter when applicable. We keep the  
835 top- $K$  hits (default  $K=40$ ) and assign stable  
836 `kg_ids` for provenance.

## 837 B Example: From a Genomics QA Pair 838 to an X-R2 Bench Instance

839 This appendix presents a concrete example show-  
840 ing how a raw genomics QA record is transformed

into a benchmark instance in X-R2 Bench, to-  
841 gether with the corresponding iterative reasoning  
842 trace visualized as a tree. 843

### 844 B.1 Raw QA Pair (Input Record)

845 **ID.** 54

846 **Task.** `genes_and_loci`

847 **Document.** `abstract_54` (abstract)

848 **Reference answer ( $A^*$ ).**

849 ["The article mainly studies BLAST",  
850 "POG", "RNA"]

851 **Document evidence ( $\mathcal{D}$ ).**

852 POGs/PlantRBP (<http://plantrbp.uoregon.edu/>)  
853 is a relational database that integrates data  
854 from rice, Arabidopsis, and maize by placing  
855 the complete Arabidopsis and rice proteomes  
856 and available maize sequences into putative  
857 orthologous groups (POGs). Annotation efforts  
858 will focus on predicted RNA binding proteins  
859 (RBPs): i.e. those with known RNA binding  
860 domains or otherwise implicated in RNA func-  
861 tion. POGs form the heart of the database, and  
862 were assigned using a mutual-best-hit-strategy  
863 after performing BLAST comparisons of the  
864 predicted Arabidopsis and rice proteomes. Each  
865 POG entry includes orthologs in Arabidopsis  
866 and rice, annotated with domain organization,  
867 gene models, phylogenetic trees, and multiple  
868 intracellular targeting predictions. A graphical  
869 display maps maize sequences on to their most  
870 similar rice gene model. The database can be  
871 queried using any combination of gene name,  
872 accession, domain, and predicted intracellular  
873 location, or using BLAST. Useful features of  
874 the database include the ability to search for  
875 proteins with both a specified domain content  
876 and intracellular location, the concurrent display  
877 of mutual best hits and phylogenetic trees which  
878 facilitates evaluation of POG assignments, the  
879 association of maize sequences with POGs, and  
880 the display of targeting predictions and domain  
881 organization for all POG members, which  
882 reveals consistency, or lack thereof, of those  
883 predictions.

### 884 B.2 Simulated X-R2 Bench Instance (JSON)

885 The following is a simulated benchmark instance  
886 that follows our X-R2 Bench schema and is di-  
887 rectly consumable by our inference scripts.

```
888 {
889   "id": 54,
890   "q_index": 54,
891   "task": "genes_and_loci",
892   "input": {
893     "question": "Does the study explicitly mention
894               any specific genes or loci studied? Please
895               list them.",
896     "document": {
897       "doc_id": "abstract_54",
898       "type": "abstract",
899       "text": "POGs/PlantRBP (http://plantrbp.uoregon.edu/) is a relational database
900              that integrates data from rice,
901              Arabidopsis, and maize by placing the
902              complete Arabidopsis and rice proteomes
903
```

```

904         and available maize sequences into
905         putative orthologous groups (POGs).
906         Annotation efforts will focus on
907         predicted RNA binding proteins (RBPs): i.
908         e. those with known RNA binding domains
909         or otherwise implicated in RNA function.
910         POGs form the heart of the database, and
911         were assigned using a mutual-best-hit-
912         strategy after performing BLAST
913         comparisons of the predicted Arabidopsis
914         and rice proteomes. Each POG entry
915         includes orthologs in Arabidopsis and
916         rice, annotated with domain organization,
917         gene models, phylogenetic trees, and
918         multiple intracellular targeting
919         predictions. A graphical display maps
920         maize sequences on to their most similar
921         rice gene model. The database can be
922         queried using any combination of gene
923         name, accession, domain, and predicted
924         intracellular location, or using BLAST.
925         Useful features of the database include
926         the ability to search for proteins with
927         both a specified domain content and
928         intracellular location, the concurrent
929         display of mutual best hits and
930         phylogenetic trees which facilitates
931         evaluation of POG assignments, the
932         association of maize sequences with POGs,
933         and the display of targeting predictions
934         and domain organization for all POG
935         members, which reveals consistency, or
936         lack thereof, of those predictions."
937     }
938 },
939 "reference_answer": [],
940
941 "decomposition": [
942     "Scan the abstract for explicit gene or locus
943     identifiers (e.g., standard gene symbols,
944     locus tags, accessions).",
945     "Disambiguate biological entities (genes/loci)
946     from methods, databases, species names, and
947     general biological terms.",
948     "If no explicit gene/locus is mentioned, output
949     an abstention/empty list with evidence."
950 ],
951
952 "entity_slot": {
953     "canonical_entities": [
954         {"name": "POGs/PlantRBP", "type": "database"},
955         {"name": "putative orthologous groups (POGs)",
956          "type": "concept"},
957         {"name": "RNA binding proteins (RBPs)", "type":
958          "protein_class"},
959         {"name": "Arabidopsis", "type": "species"},
960         {"name": "rice", "type": "species"},
961         {"name": "maize", "type": "species"},
962         {"name": "BLAST", "type": "method"}
963     ],
964     "relation_constraint_cues": [
965         {"cue": "explicitly mention", "type": "
966          constraint"},
967         {"cue": "specific genes or loci", "type": "
968          target_type"},
969         {"cue": "list them", "type": "output_format"}
970     ],
971     "notes": "Only accept explicit gene/locus
972     identifiers; do not treat methods (e.g.,
973     BLAST), concepts (POG), or general terms (
974     RNA/RBP) as genes/loci."
975 },
976
977 "kg": {
978     "subgraph": {
979         "sentence_map": {
980             "S0": "POGs/PlantRBP integrates data from
981             Arabidopsis/rice/maize into putative
982             orthologous groups (POGs).",
983             "S1": "Annotation focuses on predicted RNA
984             binding proteins (RBPs).",
985             "S2": "POGs were assigned using a mutual-
986             best-hit strategy after BLAST
987             comparisons.",
988             "S3": "Each POG entry includes orthologs,
989             domain organization, gene models,
990             phylogenetic trees, and targeting

```

```

1078     "used_kg_ids": [],
1079     "M_t": [],
1080     "A_t": [],
1081     "self_check": {"sufficient_for_Q": true, "
1082                   reason": "Empty list is consistent with
1083                           the document evidence."}
1084   }
1085 ],
1086
1087 "predicted_answer": []
1088 }

```

### 1089 **B.3 Reasoning Trace as a Tree (with Iterative** 1090 **Augmentation and Correction)**

```

1091 Root: Q
1092 Q: Does the study explicitly mention any specific
1093     genes or loci studied? Please list them.
1094 Iteration t=1 (decompose extract reason
1095               check)
1096 Q1: Identify candidate "gene/locus-like"
1097     mentions in the abstract.
1098 + Context acquired
1099 +Evidence sentences: S2 ("BLAST comparisons
1100                       "), S0 ("POGs/PlantRBP; POGs"), S1 ("
1101                       RNA binding proteins")
1102 +Initial constraint: "must be explicit gene/
1103                       locus identifiers"
1104 Extract (E1)
1105 E1 candidates: {BLAST, POG, RNA} (may be a
1106                       type confusion)
1107 Intermediate answer A1
1108 A1 (naive): [BLAST, POG, RNA] (methods/
1109                       concepts treated as genes/loci)
1110 Self-check C1 (verification trigger)
1111 Check-1: "Do candidates match gene/locus
1112           identifier patterns (symbols/locus
1113           tags/accessions)?" NO
1114 Action: Re-decompose + add a disambiguation
1115           sub-goal
1116
1117 Iteration t=2 (re-decompose conditioned on A1 +
1118               acquire missing info)
1119 Q2: Disambiguate candidates by entity type
1120     using document evidence.
1121 + New information introduced
1122 +Type rules: gene/locus must be (gene symbol
1123               | locus tag | accession); otherwise
1124               reject
1125 +Doc-grounded typing evidence:
1126   S2: "BLAST comparisons" BLAST is a
1127       method
1128   S0: "putative orthologous groups (POGs)"
1129       POG is a grouping concept
1130   S1: "RNA binding proteins" RNA/RBPs are
1131       general biological terms/classes
1132 Correction (fix)
1133 Fix-1: BLAST method (not gene/locus)
1134 Fix-2: POG concept/grouping (not gene/
1135       locus)
1136 Fix-3: RNA general term (not gene/locus)
1137 Intermediate answer A2 (revised)
1138 A2: No explicit genes/loci are mentioned in
1139     the abstract. Output (abstain).
1140 Self-check C2
1141 Check-2: "Is the revised decision supported
1142           by explicit sentences?" YES (S0/S1/
1143           S2)
1144 stop = 1
1145
1146 Final output Â
1147 Â: (no explicit gene/locus identifiers)
1148 Provenance: S0/S1/S2 (type-disambiguation
1149           evidence)

```