

Hindsight: A Git Intelligence Layer

Making Large Codebases Easier to Navigate with AI.

University of Massachusetts Amherst | 698DS



Team:

*Mustafa Ali¹
Vanshika Agrawal
Priya Balakrishnan
Sheng Kai Wen*

PhD Mentor:
Anshita Gupta

Mentors:

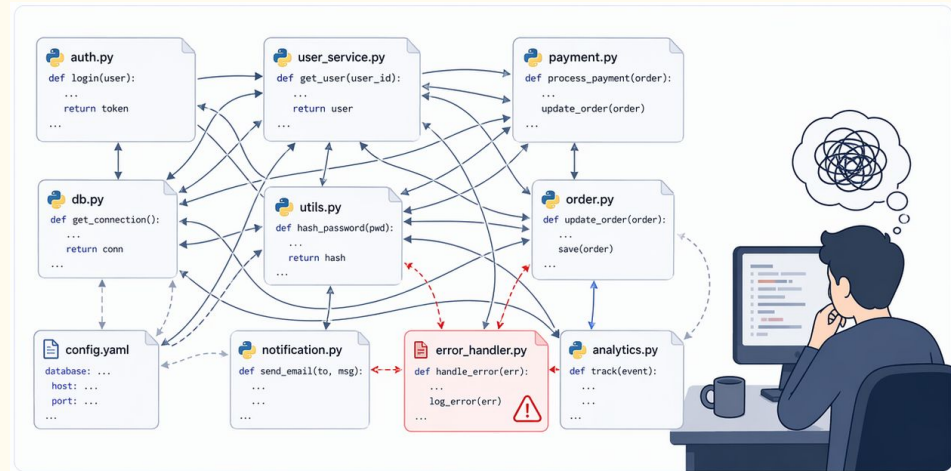
*Ryan Lingo
Rajeev Chhajer*



Problem Statement

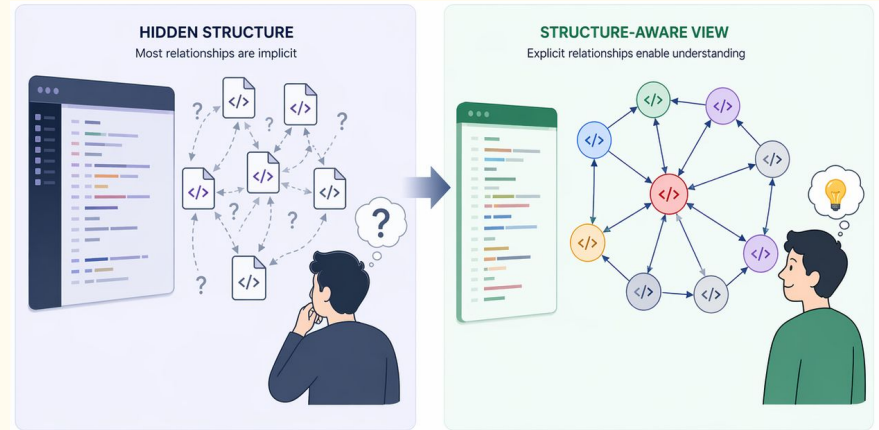
Modern software systems are large, modular, and deeply interconnected, oftenspanning hundreds of files and components. Developers struggle with:

- Understanding cross-file and module interactions
- Tracing dependency chains across the system
- Predicting impact of small changes
- Debugging hidden cross-component failures



Our Motivation

- Software complexity is outpacing understanding
- Structural knowledge is implicit and undocumented
- Misunderstanding leads to:
 - Integration failures
 - Debugging delays
 - Higher engineering cost
- Current AI tools rely on text, not structure





Research Questions

Can graph representations of code repositories improve retrieval-augmented generation (RAG) performance on multi-file code understanding tasks?

Why Graph-Augmented Retrieval over text Only Retrieval (vs BM25 Alone)

What current methods do

- Retrieves code chunks based on **keyword overlap** (BM25) or **semantic similarity** (embeddings)
- Treat repository as independent **text chunks**

What they miss

- Ignore how code is connected across files
- Fail to capture:
 - Function call chains
 - Cross-file dependencies
 - Shared state / configuration

Why this matters

- Real-world questions require multi-file reasoning
- Critical context may exist in files that are not textually similar

1) How BM25 Works (Text-Only Retrieval)

Retrieves files only based on keyword matching.

Step 1
User Question



Query:
"Why is batching equivalence test skipped on CUDA device 8?"

Step 2
Compute
BM25 Scores



BM25 scores each file based on keyword overlap (e.g., *batching*, *test*, *CUDA*, *8*)

Step 3
Return Top Results



Top Retrieved Files:
1) test_lightglue.py (highest keyword match)
2) README.md
3) requirements.txt

Limitation: Misses related files that don't share keywords but are actually connected.

Example Outcome

BM25 returns:



test_lightglue.py

Only the test file is retrieved.
It misses the files that explain why the test is skipped.

Why Graph-Augmented Retrieval over text Only Retrieval (vs BM25 Alone)

Our Approach

- L1: Use BM25 to get initial relevant files
- L2: Expand using graph over:
 - Function calls
 - Imports / dependencies
 - Module relationships

What this enables

- Retrieve connected files across the repository
- Build multi-hop context for reasoning
- Improve grounding for LLM answers

2) How We Use Graph on Top of BM25

Start with BM25, then expand using repository structure.

Step 1 User Question



Query:
"Why is batching equivalence test skipped
on CUDA device 8?"

Step 2 BM25 Gets Initial Files

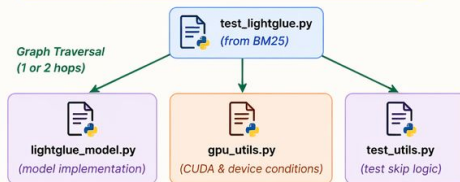


BM25 retrieves initial candidates:
• test_lightglue.py

Step 3 Graph Expansion (Hop to Related Files)



Traverse the **code graph** from the
initial file using relationships like:
• function calls • imports
• dependencies • module references



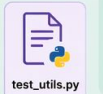
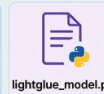
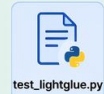
Step 4 Return Augmented Context



Final Retrieved Context:
{ test_lightglue.py + lightglue_model.py
+ gpu_utils.py + test_utils.py }

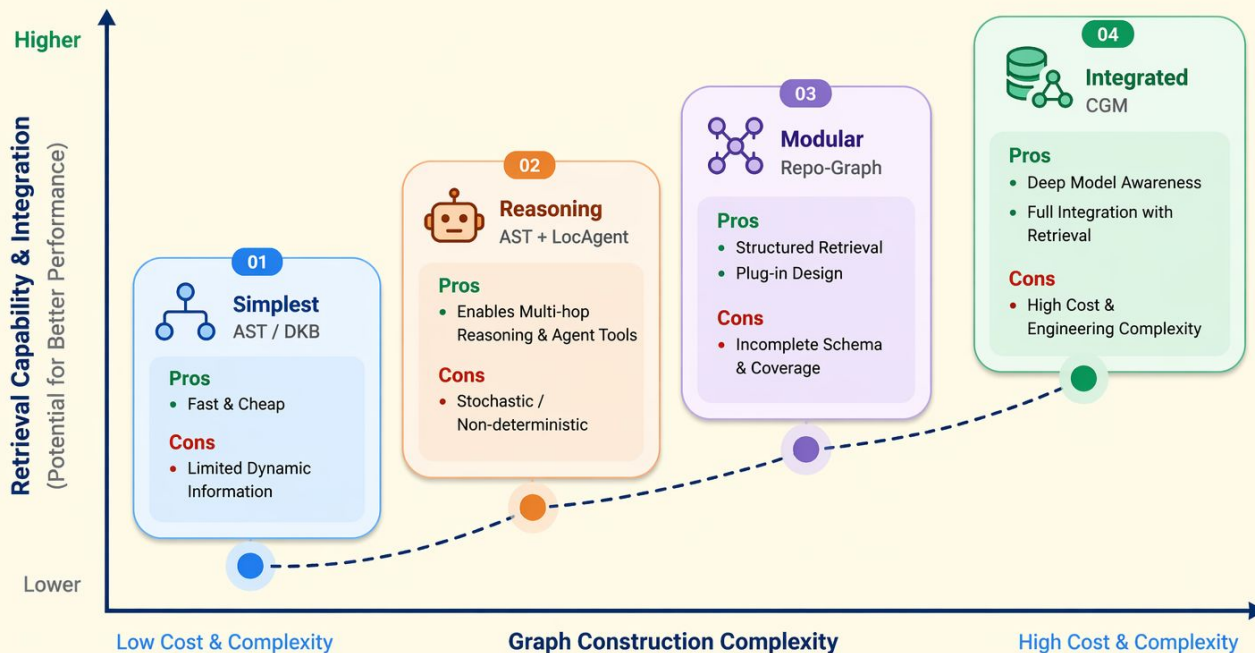
Example Outcome

BM25 + Graph returns:



We get the test file AND
the related files that reveal
the real reason: a CUDA
device condition causes
the test to be skipped.

What graph construction strategies balance performance and cost?



Related Work

Related work converges on two themes: realistic repository-level benchmarks and structure-aware retrieval over codebases.

Repository QA benchmarks

DeepCodeBench

Qodo
2025

PR-grounded multi-file QA benchmark over 8 OSS repositories. Uses fact recall and answer-quality judging to score generated responses.

Why it matters: Directly matches our answer-generation setting and lets us test whether structured retrieval improves grounded reasoning.

SpyderCodeQA

ACL SRW
2024

325 repository-level Q&A pairs covering source-code semantics, dependencies, and repository meta-information.

Why it matters: Shows that repository QA must capture cross-file relations, not just single-snippet understanding.

CoReQA

arXiv
2025

Issue-derived benchmark for repository QA from real GitHub issues and comments, evaluated with LLM-as-a-judge.

Why it matters: Makes repo QA more realistic by tying questions to developer-facing issues and discussions.

Shared signal: *all three benchmarks emphasize multi-file retrieval before answer generation.*

Localization and graph-based methods

LocBench + LocAgent

ACL 2025

LocBench covers bug, feature, performance, and security issue localization. LocAgent uses graph-guided exploration over a heterogeneous code graph.

Why it matters: Better localization narrows the search space before downstream reasoning or answer generation.

RepoGraph

ICLR 2025

Repository-level code graph plug-in that gives AI systems fine-grained repository navigation and context aggregation.

Why it matters: Motivates adding structure-aware exploration instead of relying only on flat chunk retrieval.

Code Graph Model (CGM)

arXiv
2025

Graph-integrated LLM that injects code graph structure into the model and pairs it with graph RAG for repository-level tasks.

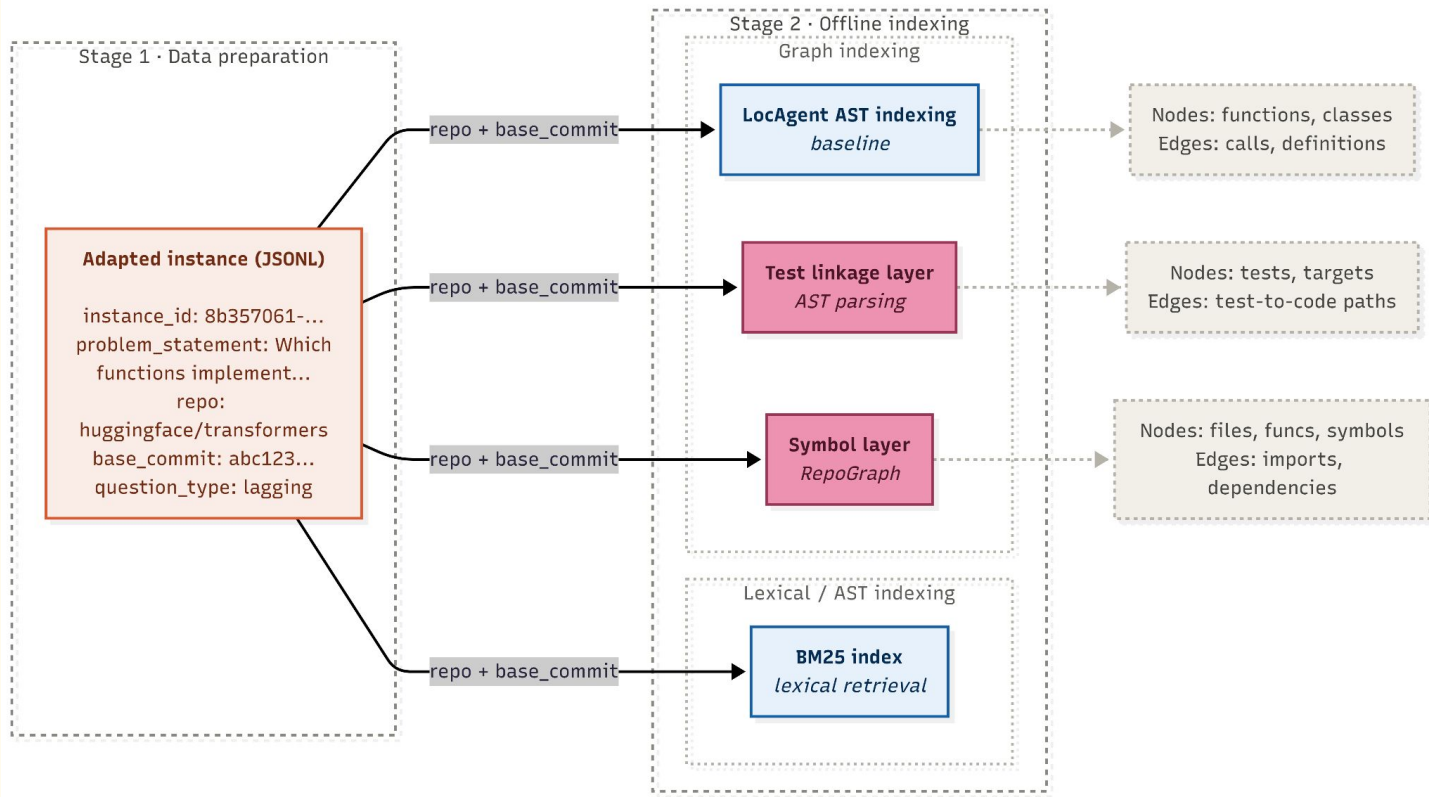
Why it matters: Suggests that structural signals can help both retrieval and the model's internal representation of the repository.

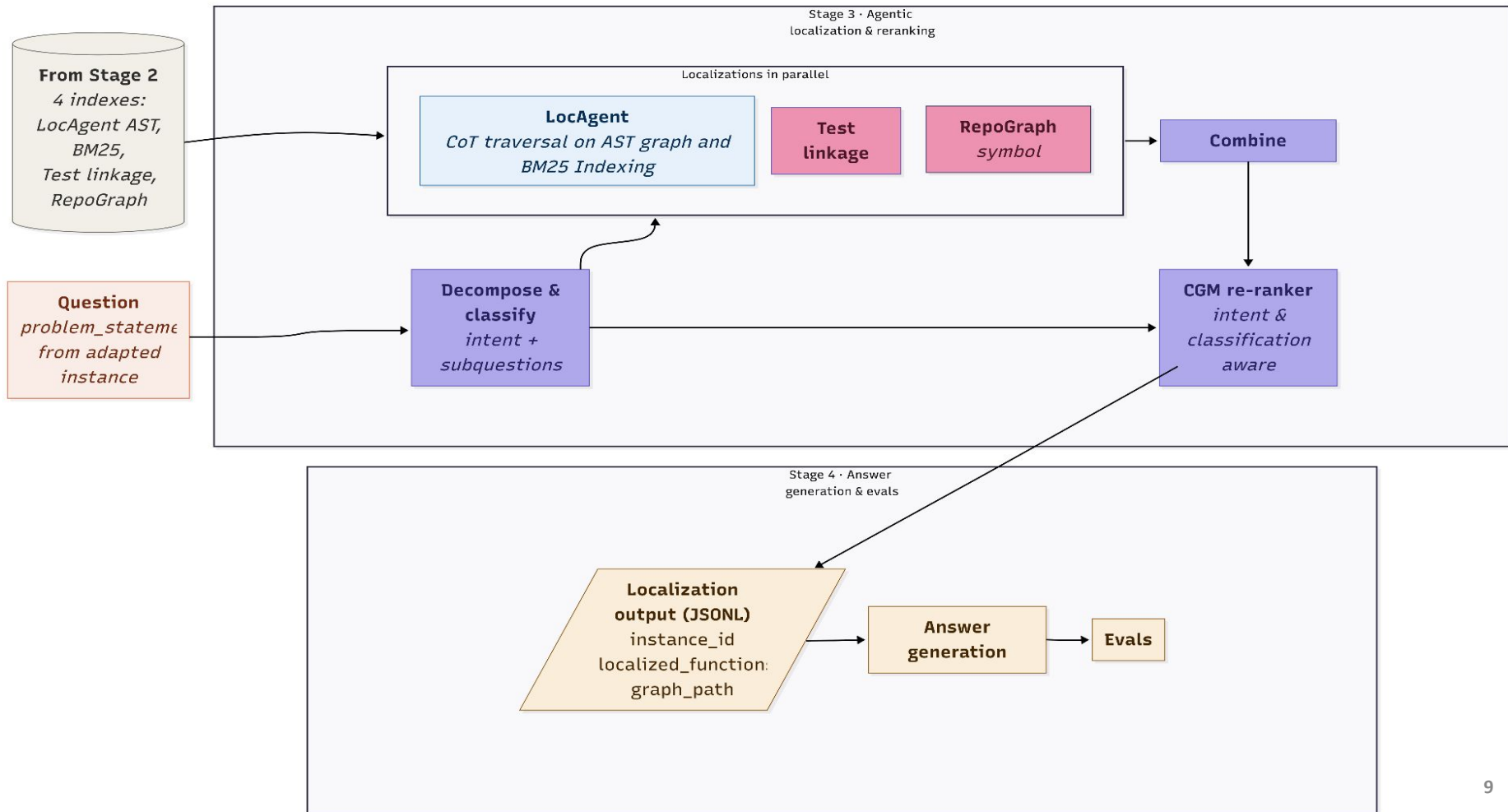
Our positioning *We sit between benchmark-driven QA evaluation and graph-guided repository navigation.*

Takeaway: repository understanding improves when evaluation is realistic and retrieval is structure-aware.

Method

Proposed Solution







Experiments

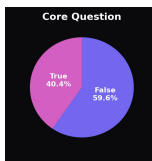
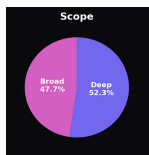


Datasets

DeepCodeBench

1,144 total

Q&A benchmark



Scope: Question spans single block/multiple files

Type: Question targets core functionality or peripheral details

Searchable: Question includes location hints or not

Measures the quality of the model's final coding answer.

LocBench

560 total

Localization benchmark

Bug Report	<div style="width: 42.5%;"></div>	242
Feature Request	<div style="width: 26.8%;"></div>	150
Performance Issue	<div style="width: 24.8%;"></div>	139
Security Vulnerability	<div style="width: 5.0%;"></div>	29

Issue descriptions are linked to the relevant edited functions or code regions.

Measures whether the system can find where the problem is.

Intent Categories



→ Leading

When: proactive • pre-change

Focus: impact • consequences • planning

Example cues

- change impact
- downstream effect
- design / implement

! Lagging

When: reactive • post-failure

Focus: debugging • diagnosis • root cause

Example cues

- why broke?
- what caused it?
- unexpected behavior

? Exploratory

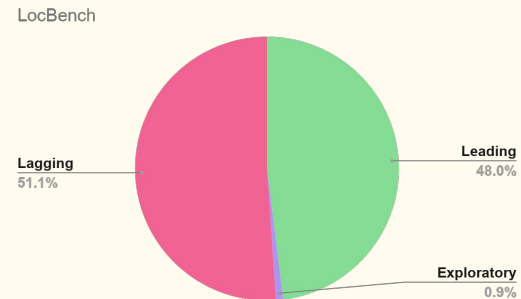
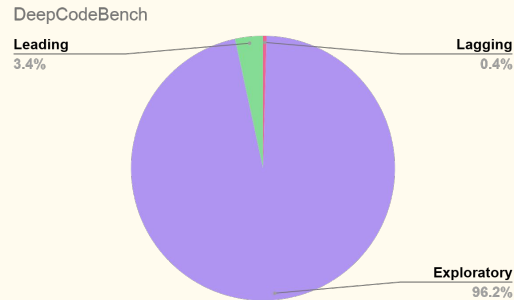
When: neutral • no active failure

Focus: structure • mechanism • lookup

Example cues

- where is X?
- how does Z work?
- which component?

Dataset Distribution



Evaluation Methods

Q&A Evaluation

Fact Recall

Measures: Coverage of gold reference facts in the generated answer

Formula:

Fact Recall =
(# facts judged present) / (# total gold facts)

Judge Signal:

Fact-level YES / NO

Quality Evaluation

Measures: Overall answer quality relative to the reference answer

Dimensions:

Accuracy • Completeness • Relevance • Clarity

Judge Setup:

1 to 10 scale • repeated calls • average

Localization Evaluation

Acc@k

Measures: Overall localization quality i.e, whether the correct output is in the top-k predicted locations

Example

In the dataset

What error is raised when the 'mid' scale is provided as a list of more than one value?

Scope: Deep, Searchable: False, Type: Non-Core

Inputs

Judges

Outputs

In the dataset

Ground Truth Answer

In ``_maybe_expand_lora_scales_for_one_adapter`` (src/diffusers/loaders/unet_loader_utils.py), if ``scales["mid"]`` is a list with more than one element it raises:

ValueError: "Expected 1 scales for mid, got {len(scales['mid'])}."

Generated by prompting an LLM

Generated Answer

A ``ValueError`` is raised. This occurs in ``_maybe_expand_lora_scales_for_one_adapter`` (src/diffusers/loaders/unet_loader_utils.py) when the ``"mid"`` key in the ``scales`` dictionary is provided as a list containing more than one value.

Facts (In the dataset)

The function ``_maybe_expand_lora_scales_for_one_adapter`` is defined in the file src/diffusers/loaders/unet_loader_utils.py.

In ``_maybe_expand_lora_scales_for_one_adapter``, if ``scales["mid"]`` is a list with more than one element, the function raises a ``ValueError``.

The ``ValueError`` message raised is "Expected 1 scales for mid, got {len(scales['mid'])}."

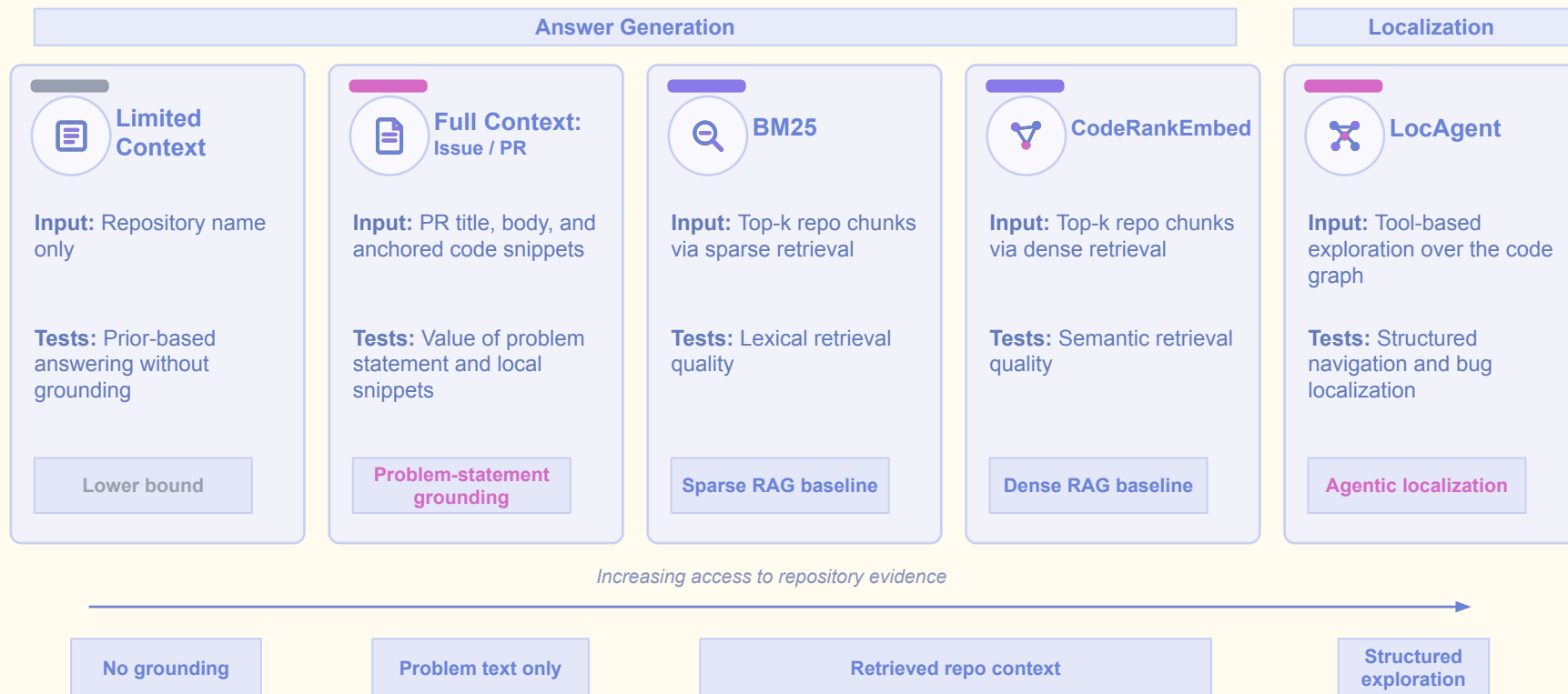
Quality Judge

Accuracy	Completeness	Relevance	Clarity
9.5	9.0	10.0	9.5

Fact Recall Judge

Facts present = 2
Total facts = 3
Fact Recall = 0.66

Baselines



Implementation Details

Off-the-Shelf Assets

What we reuse directly

Benchmarks

DeepCodeBench • LocBench

Task models

Qwen3-32B • Qwen2.5-7B • gpt-5.3-codex

Judge

Qwen3-32B thinking mode

Retrieval

BM25 • CodeRankEmbed

Benchmarks + models + evaluation

Software Stack

Core implementation layer

Core ML/NLP:

Python • PyTorch • HF transformers • datasets

Indexing:

Bm25s • rank_bm25 • faiss-cpu

Tooling:

OpenAI API • llama-index

Graph Methods:

AST • Dependency • Structural • Knowledge

Retrieval + orchestration

Computational Resources

How runs are executed

Requirements:

CPU for Retrievals

CPU for Graph Constructions

GPU for Localizations and Q&A

Cluster-scale evaluation

Results

How does performance vary across DeepCodeBench categories?

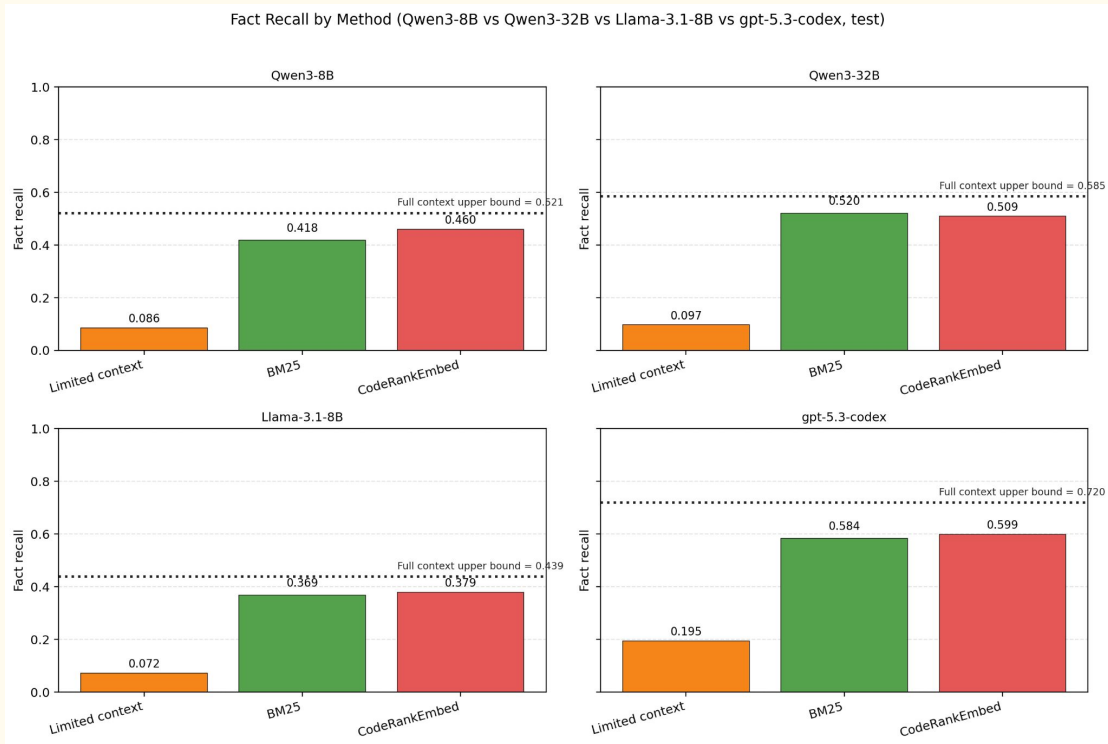
How reliable is the LLM judge across repeated evaluations and answers from different models?

19

How do Qwen3-32B and gpt-5.3-codex differ in answer quality?

How do BM25 and CodeRankEmbed compare as retrieval methods?

Fact Recall on DeepCodeBench

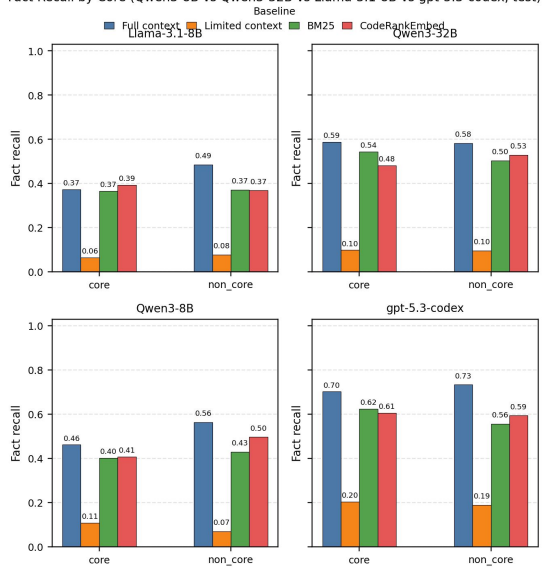


- Full context should be read as an upper bound rather than a competing baseline.
- Among the actual methods, retrieval consistently improves fact recall over Limited context.
- BM25 is the most reliable retrieval method overall, with CodeRankEmbed close behind and occasionally competitive on specific model/question slices.

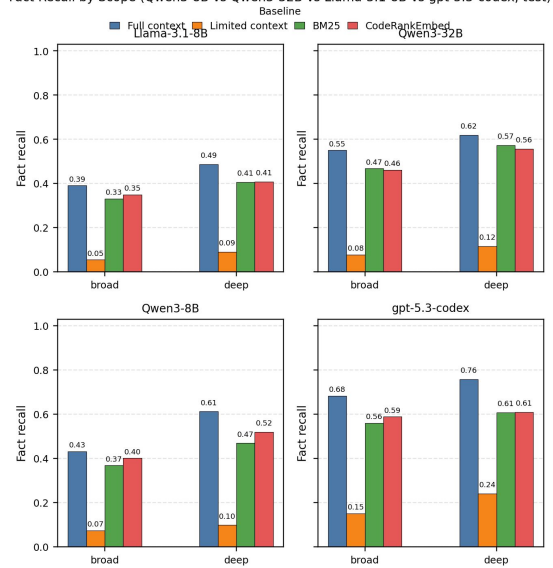
Fact Recall on DeepCodeBench

- The retrieval gains are largest on deep questions, where access to the right evidence matters most.
- Retrieval also helps more on searchable questions than on non-searchable ones, suggesting these methods are strongest when relevant facts can be localized.

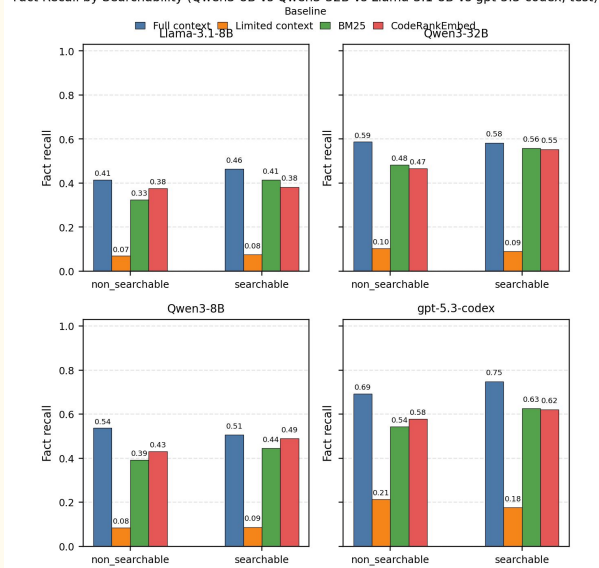
Fact Recall by Core (Qwen3-8B vs Qwen3-32B vs Llama-3.1-8B vs gpt-5.3-codex, test)



Fact Recall by Scope (Qwen3-8B vs Qwen3-32B vs Llama-3.1-8B vs gpt-5.3-codex, test)

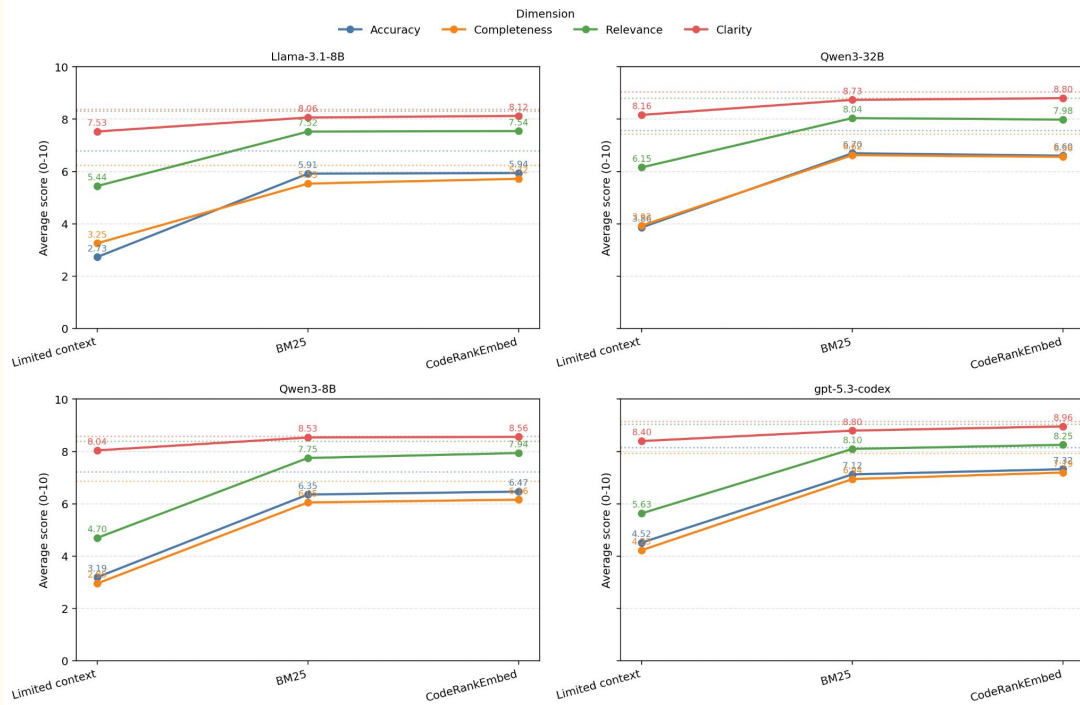


Fact Recall by Searchability (Qwen3-8B vs Qwen3-32B vs Llama-3.1-8B vs gpt-5.3-codex, test)



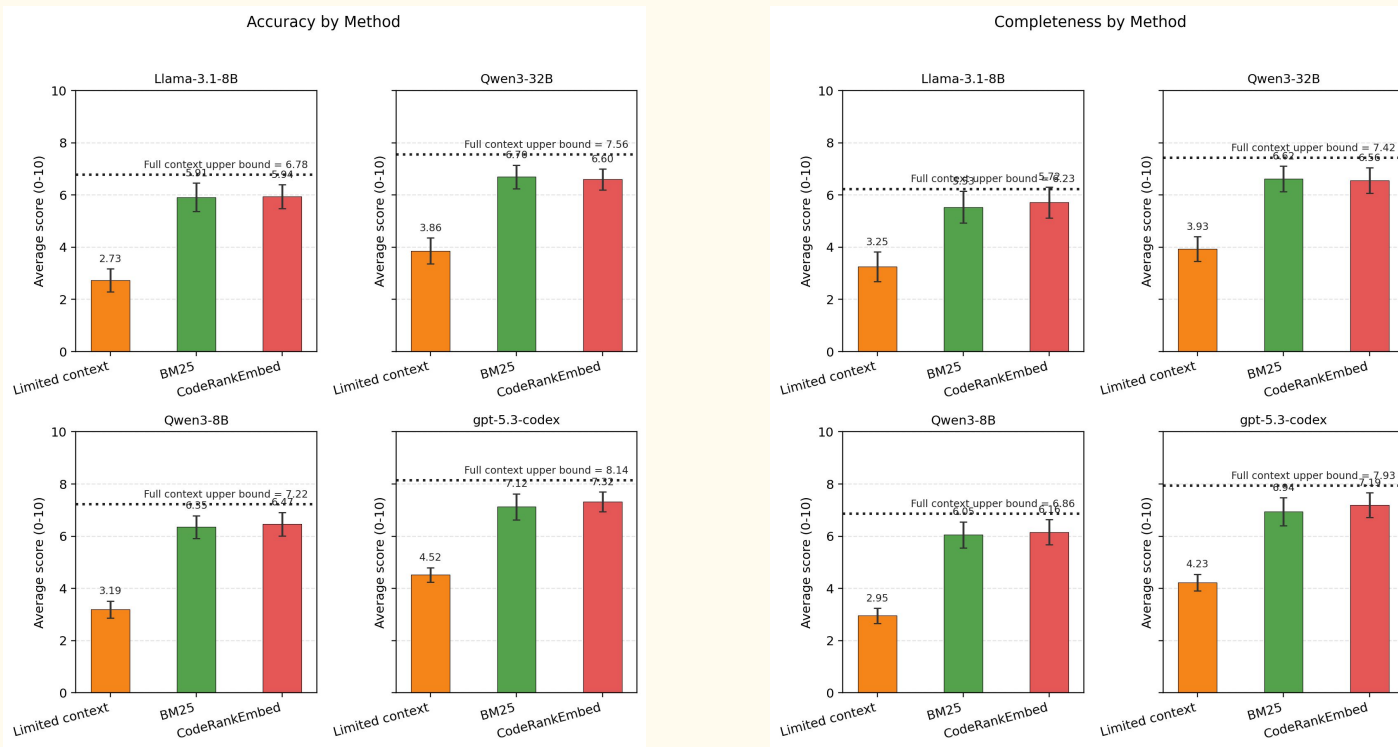
Quality Evaluation on DeepCodeBench

Quality by Method: Dimension Trends



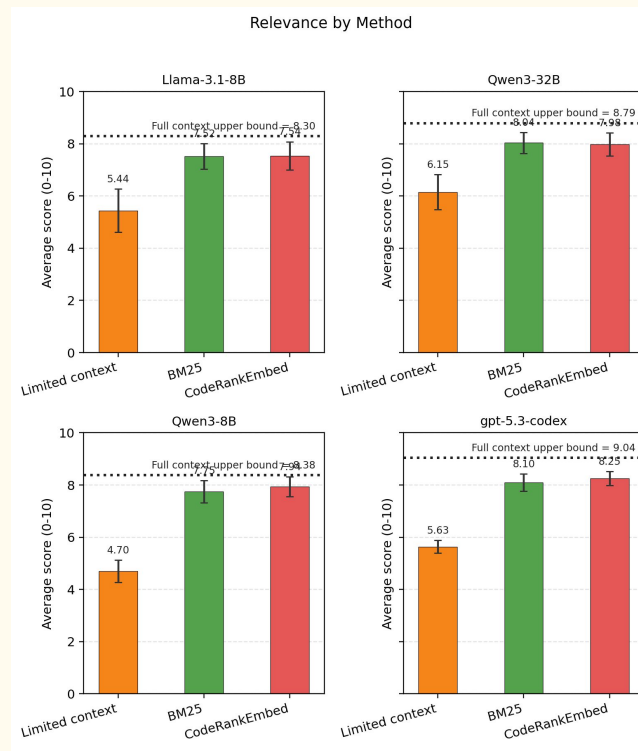
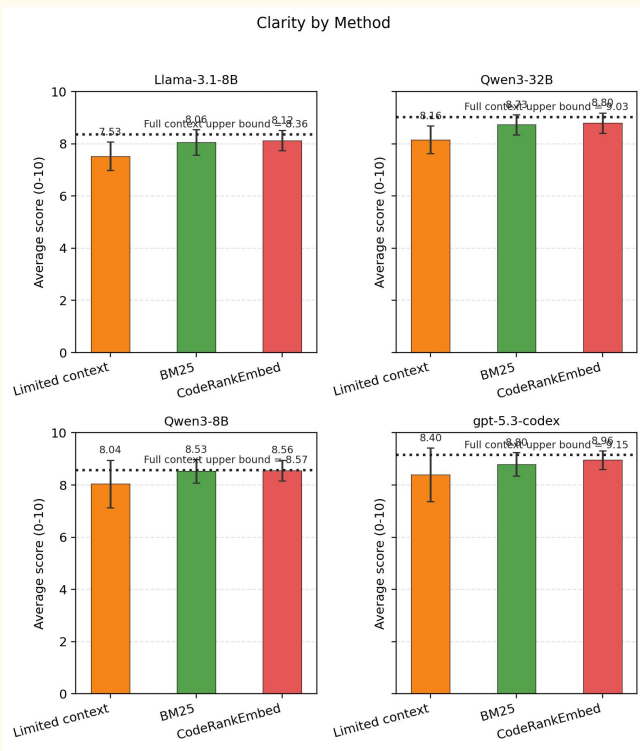
- The same method trend holds for answer quality.
- BM25 and CodeRankEmbed both move quality closer to the Full context upper bound, with relatively small gaps between them.
- Clarity is less sensitive to retrieval method than accuracy and completeness.

Quality Evaluation on DeepCodeBench



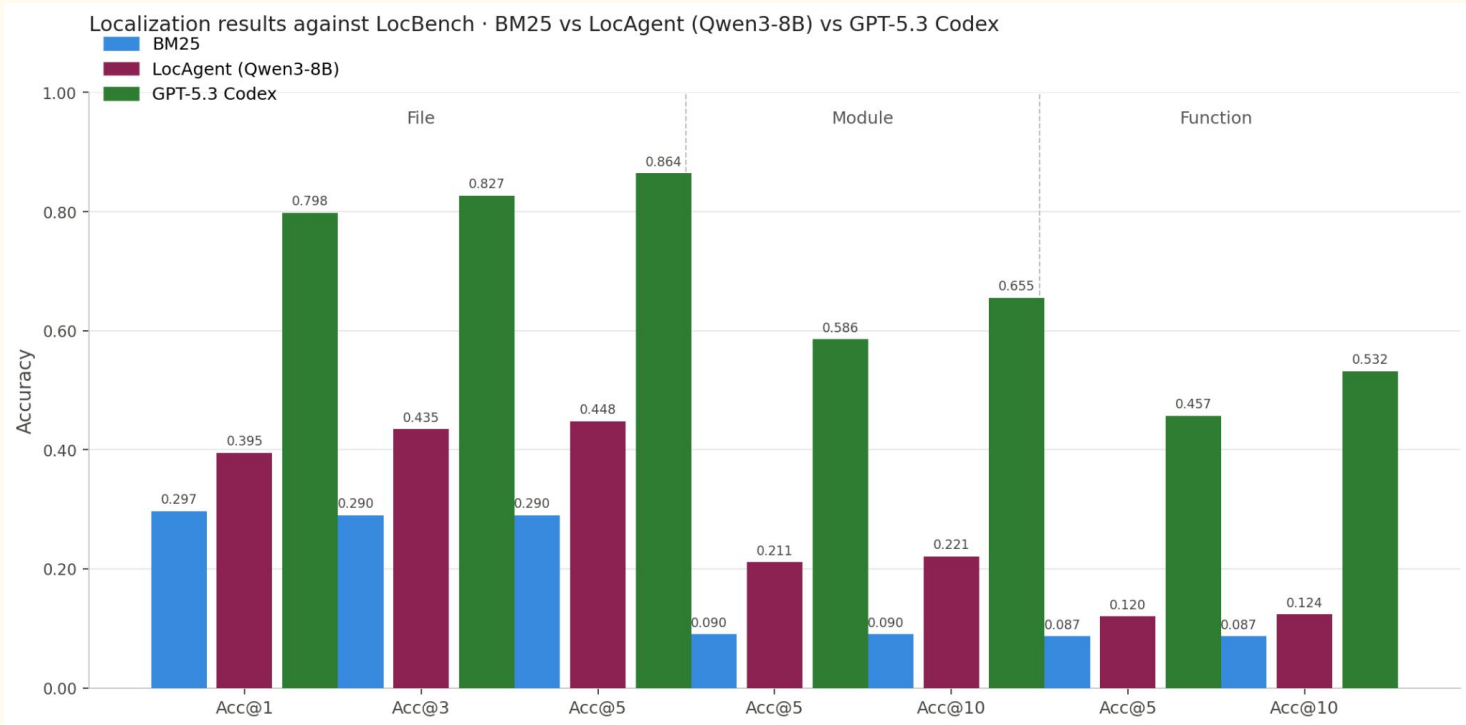
Judge variability is generally small, so the observed method-level quality differences appear stable across repeated evaluations.

Quality Evaluation on DeepCodeBench



Judge variability is generally small, so the observed method-level quality differences appear stable across repeated evaluations.

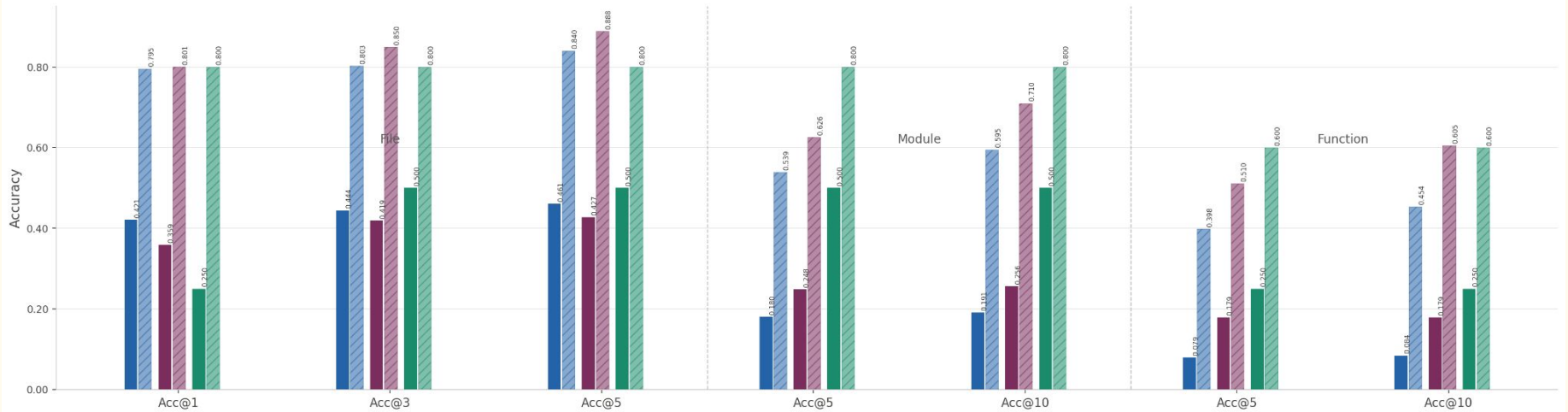
Localization on LocBench



Localization on LocBench

Localization results by question type · Leading vs Lagging vs Exploratory
LocAgent (Qwen3-8B) vs GPT-5.3 Codex

■ Leading · LocAgent (Qwen3-8B) ■ Lagging · LocAgent (Qwen3-8B) ■ Exploratory · LocAgent (Qwen3-8B)
▨ Leading · GPT-5.3 Codex ▨ Lagging · GPT-5.3 Codex ▨ Exploratory · GPT-5.3 Codex



Note: Exploratory bucket has a small sample count; values should be read as directional rather than precise.
Solid bars = LocAgent (Qwen3-8B) · Hatched bars = GPT-5.3 Codex



Thank You