



# Hindsight: a Git Intelligence Layer

Mustafa Ali<sup>1</sup>, Vanshika Agrawal<sup>1</sup>, Priya Balakrishnan<sup>1</sup>, Sheng Kai Wen<sup>1</sup>

Ryan Lingo<sup>2</sup>, Rajeev Chhajer<sup>2</sup>, Anshita Gupta<sup>1</sup>

<sup>1</sup>University of Massachusetts Amherst, <sup>2</sup>Honda RI & 99P Labs



## What Was Missing?

- Opaque Retrieval Pipelines**  
e.g., Hard to inspect how LLMs choose repo context
- Missing Cross-File Structure**  
e.g., Search may miss calls, imports, and dependency paths
- Limited KG Retrieval Evaluation**  
e.g., Few benchmarked tests of graph-augmented RAG for answer generation

## What Kinds of Questions Do We Test?

### Intent Types

- Leading Queries** (Planning before a change)
- Exploratory Queries** (Understanding unfamiliar code)
- Lagging Queries** (Debugging after a failure)

Scope: **Single** / **Multi**

Feature Request

Performance Issue

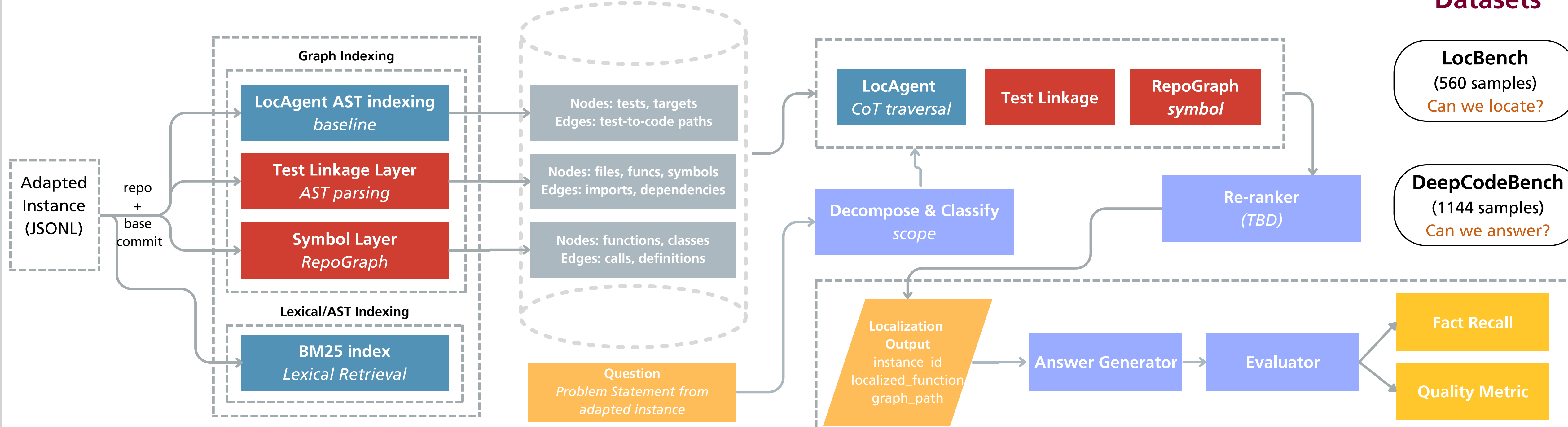
Searchability: **True** / **False**

Bug Report

Type: **Core** / **Non-core**

Security Vulnerability

## Method Overview



## Datasets

**LocBench** (560 samples)  
Can we locate?

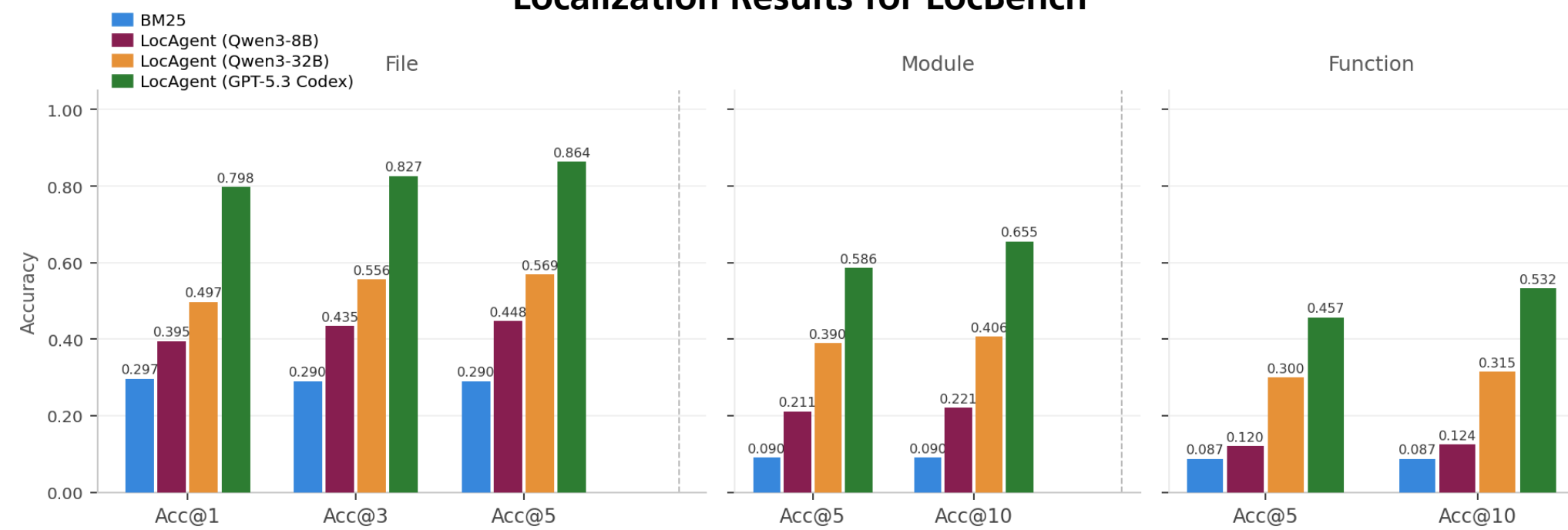
**DeepCodeBench** (1144 samples)  
Can we answer?

## Research Question

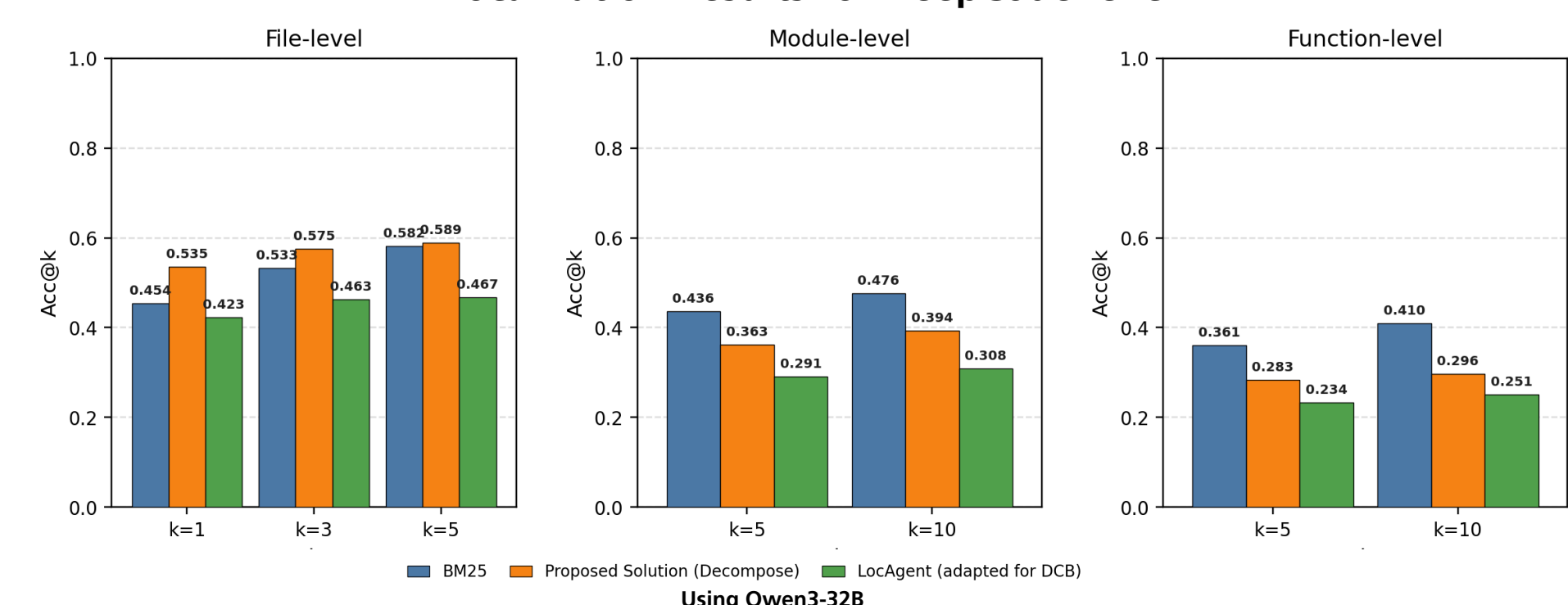
Can knowledge graph representations of code repositories improve retrieval-augmented generation (RAG) performance on codebase understanding tasks, as measured by DeepCodeBench?

## Do structural graphs improve retrieval for multi-file questions?

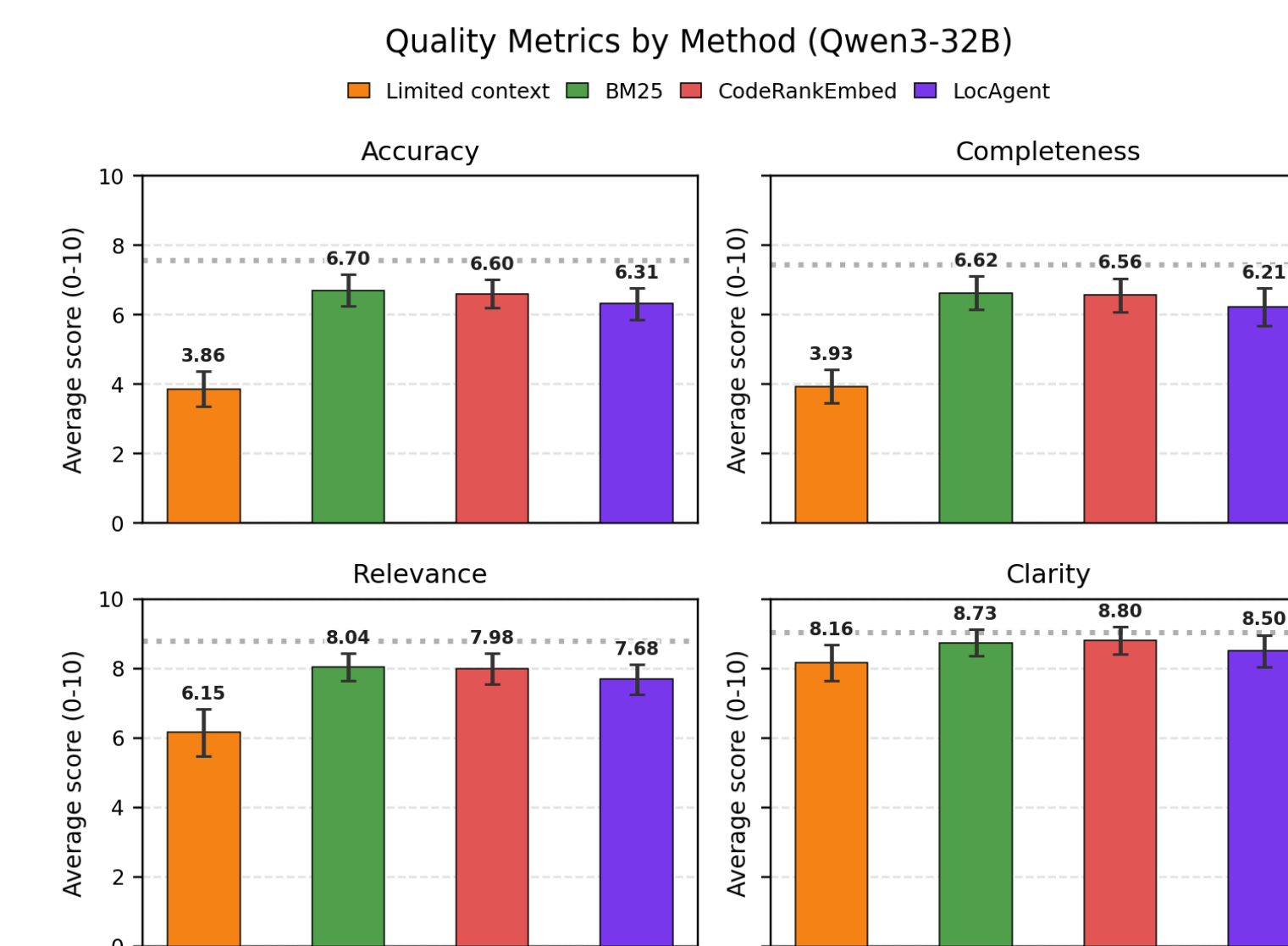
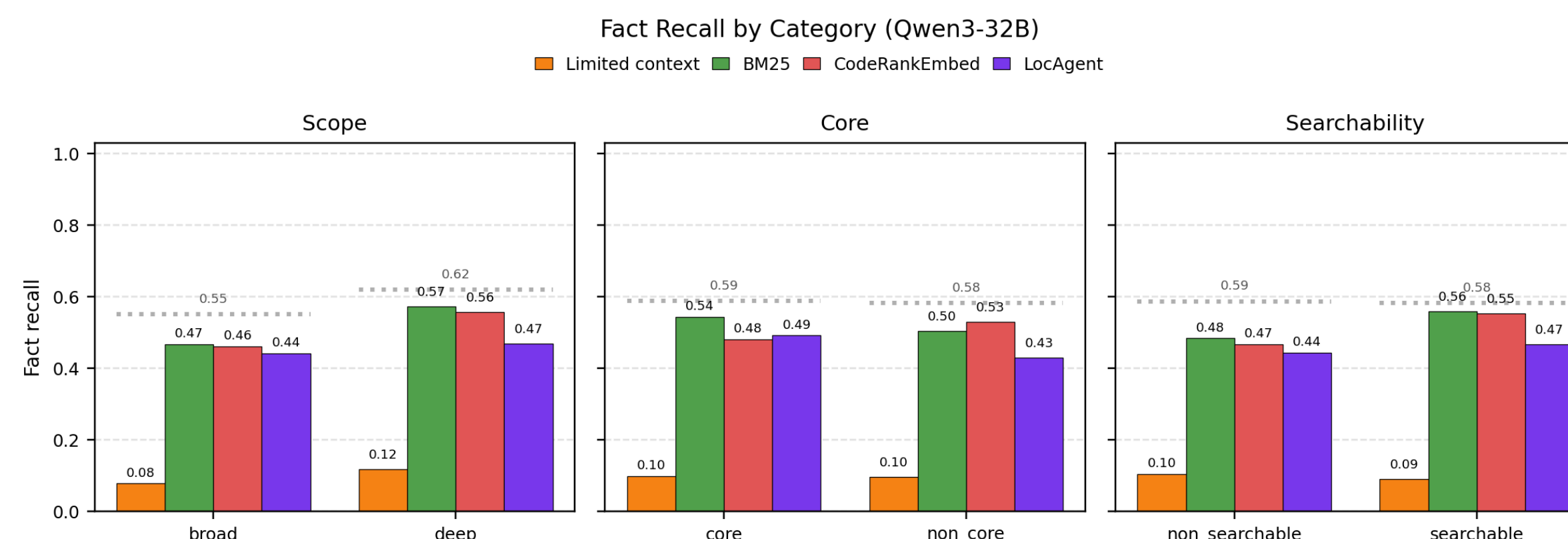
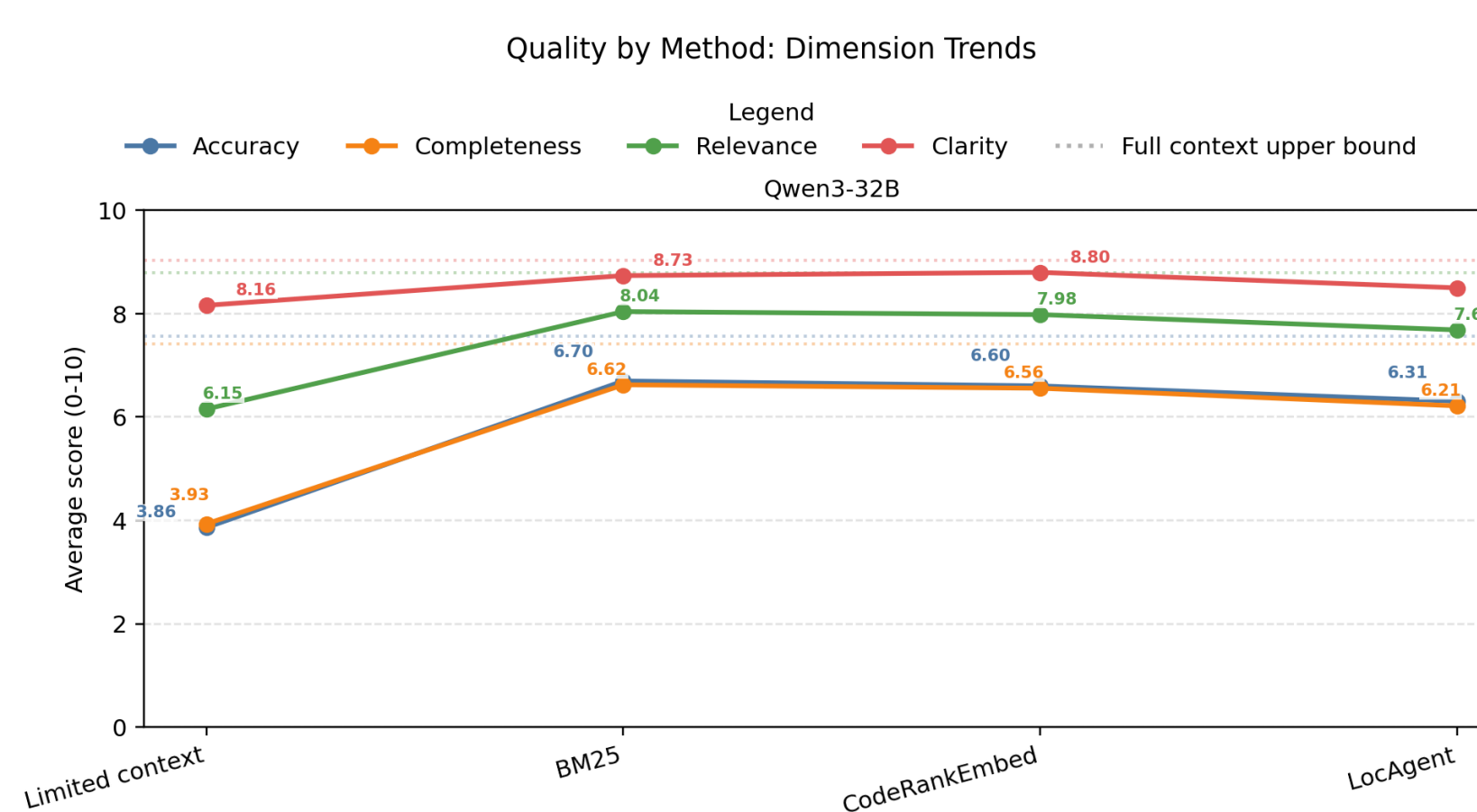
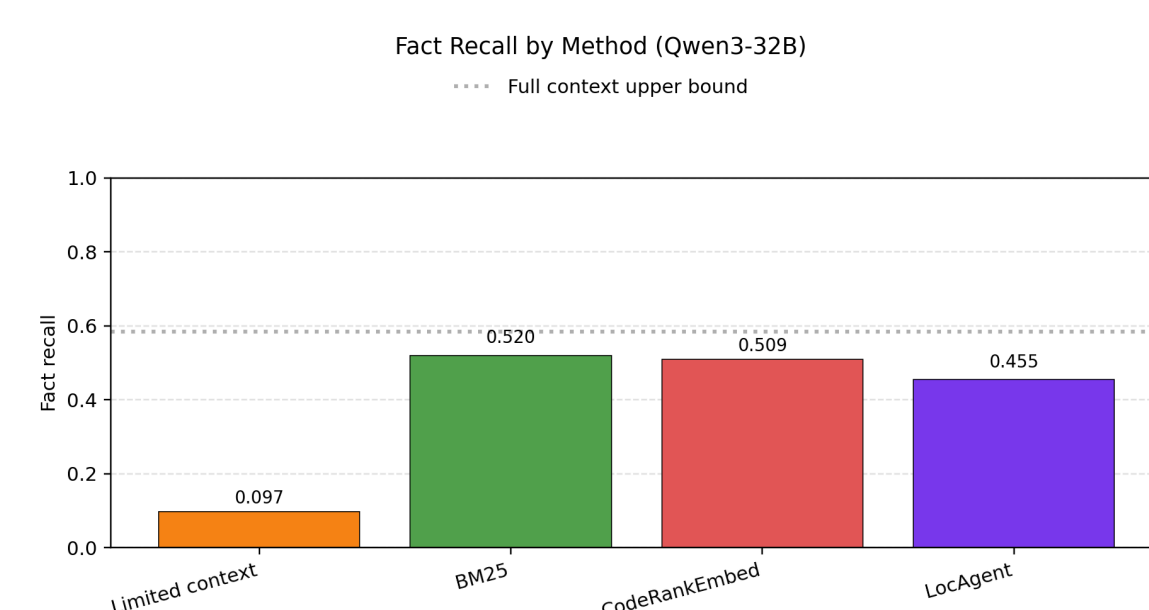
### Localization Results for LocBench



### Localization Results for DeepCodeBench



## How does KG-augmented retrieval compare to embedding retrieval on answer generation?



## Preliminary Insights

- Our graph-based method performs on par with BM25 and CodeRankEmbed for categories such as broad, core, and non-searchable questions.
- On DeepCodeBench, our method outperforms BM25 at file-level localization, suggesting that structural graph signals help identify relevant files.
- Function-level localization is still limited, showing that future work should improve graph traversal, reranking, and context assembly.
- Strong localization does not always translate to better fact recall, so future work should improve reranking and context selection.