# Contextual Augmented Multi-Model Programming (CAMP): A Local-Cloud Copilot Solution

Yuchen Wang, Shangxin Guo, Chee Wei Tan

# Introduction

Introducing **CAMP**: A Collaborative Multi-Model Copilot using <u>Context-Aware Retrieval-Augmented Generation (RAG)</u> to bridge <u>cloud LLMs</u> with <u>local models</u>.

- +12.5% over non-contextual generation
- +6.3% over baseline RAG
- Deployed as "Copilot for Xcode"
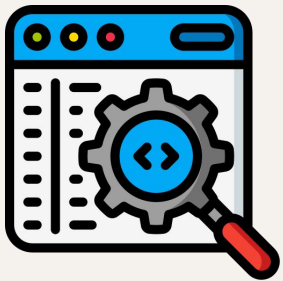- Integrated with GitHub Copilot

Copilot for Xcode: exploring AI-assisted programming by prompting cloud-based large language models
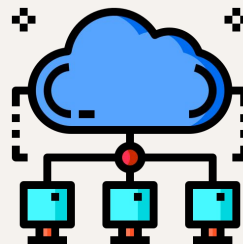CW Tan, S Guo, MF Wong, CN Hang, 2023, https://arxiv.org/abs/2307.14349
M. F. Wong and C. W. Tan, "Aligning Crowd-Sourced Human Feedback for Reinforcement Learning on Code Generation by Large Language Models," in IEEE Transactions on Big Data, 2024.

# Motivation

**Challenge**: Cloud-based LLMs are powerful for code generation, but lack local context. Locally integrated tools are adaptive, but limited in scope.
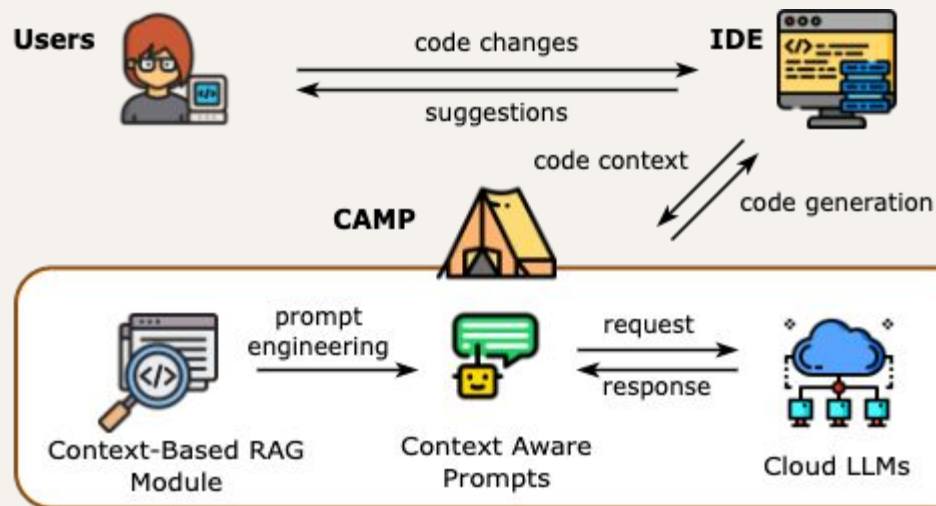
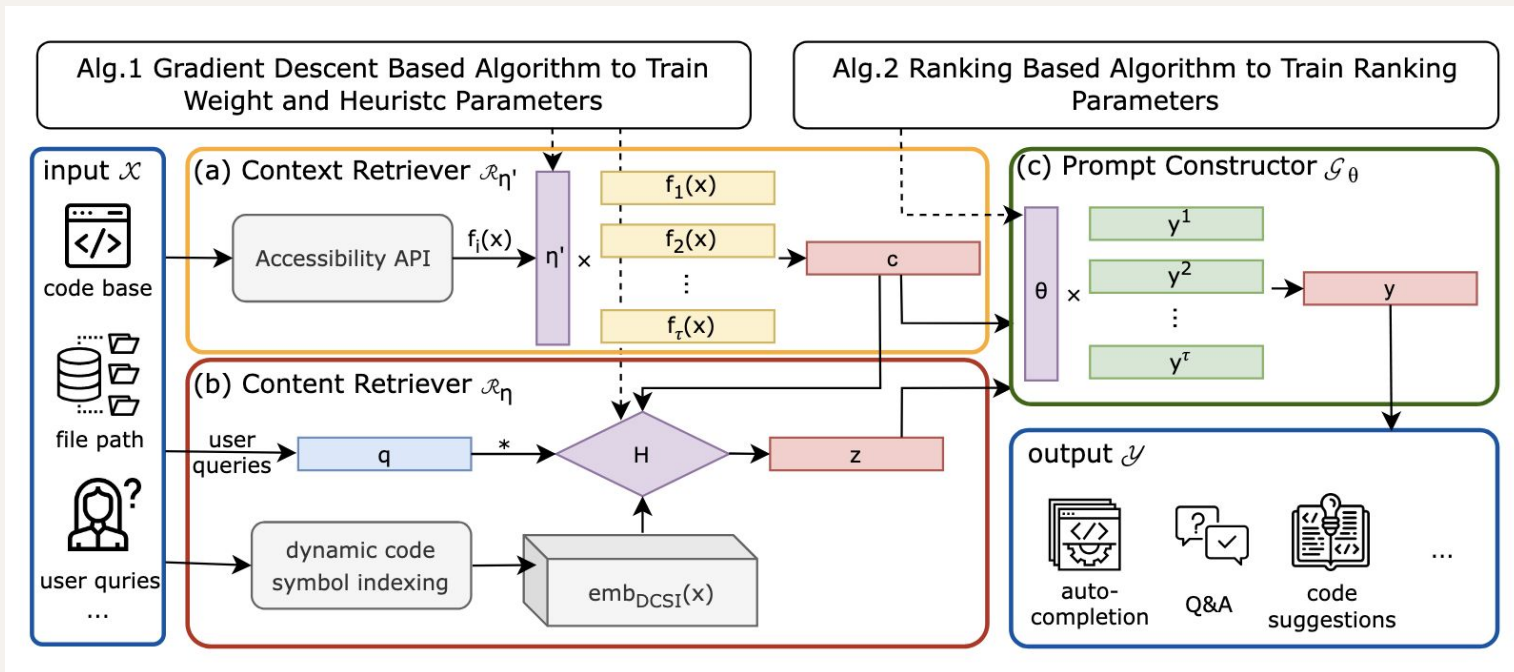**Local models**: light weight, has code context

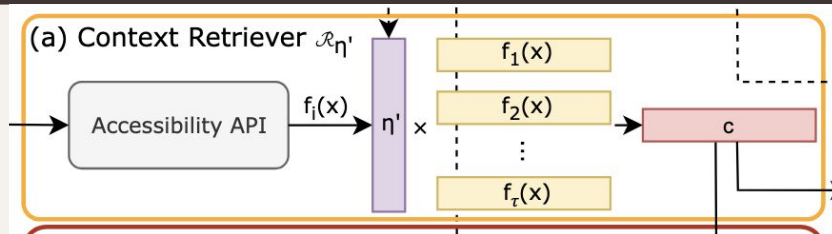**Cloud models**: strong generative power, lack local context

# Overview of CAMP

CAMP integrates cloud LLMs into local development environments, employing a RAG module that dynamically learns from code context to optimize prompt construction.

# Architecture



Alg.1 Gradient Descent Based Algorithm to Train Weight and Heuristc Parameters

Alg.2 Ranking Based Algorithm to Train Ranking Parameters

input $\mathcal{X}$
- code base
- file path
- user quries
- ...

(a) Context Retriever $\mathcal{R}_{\eta'}$

Accessibility API $\quad f_i(x) \quad \eta' \times$

$f_1(x)$
$f_2(x)$
$\vdots$
$f_\tau(x)$

$c$

(b) Content Retriever $\mathcal{R}_\eta$

user queries $\quad q \quad * \quad H \quad z$

dynamic code symbol indexing $\quad emb_{DCSI}(x)$

(c) Prompt Constructor $\mathcal{G}_\theta$

$\theta \times$

$y^1$
$y^2$
$\vdots$
$y^\tau$

$y$

output $\mathcal{Y}$
- auto-completion
- Q&A
- code suggestions
- ...

# Part 1. Context Retriever


(a) Context Retriever $\mathcal{R}_{\eta'}$
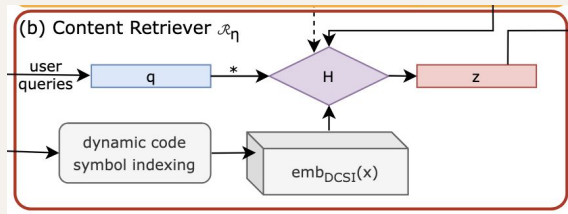
Context retriever $R_{n'}$ retrieves **contextual information** from the local development environment maximizes the insights brought to the next step.

$$
\begin{aligned}
\mathcal{R}_{\eta'}(x) &= \mathrm{agg}([\eta_0' c_0, \eta_1' c_1, \ldots, \eta_{\tau_c}' c_{\tau_c}]) \\
&= \mathrm{agg}([\eta_0' f_0(x_0), \eta_1' f_1(x_1), \ldots, \eta_{\tau_c}' f_{\tau_c}(x_{\tau_c})]) \\
&= \mathrm{agg}(\eta' \cdot f(x))
\end{aligned}
$$

**Examples of $c_i$**: cursor position, absolute repository path, cached build artifacts", index information.

**Assumption**: The relative importance of different factors in the local development environment remains stable.

# Part 2. Content Retriever



(b) Content Retriever $\mathcal{R}_\eta$

Given the retrieved context from previous step, the content retriever $R_\eta$ aims to deliver highly relevant content z that enhances prompt construction with local, context-aware information.

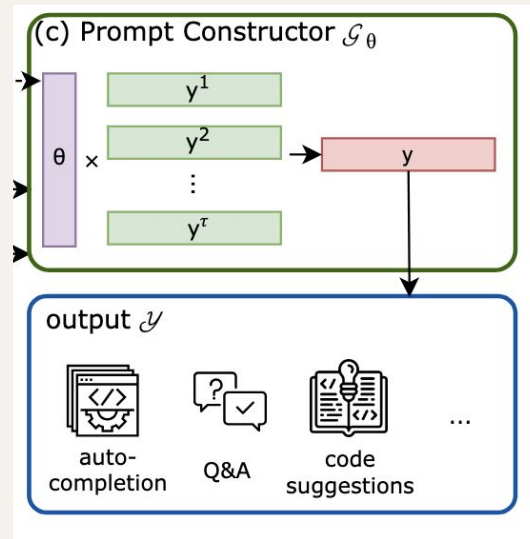- Function 1: support codebase embedding, with dynamic code symbol indexing.

$$p_\eta(z|x,c) = \frac{\exp\left(\text{emb}_{\text{DCSI}}(z)^T H \text{emb}_{\text{DCSI}}(x)\right)}{\sum_{z'} \exp\left(\text{emb}_{\text{DCSI}}(z')^T H \text{emb}_{\text{DCSI}}(x)\right)}$$

- Function 2: facilitating content search to obtain the highest ranked content that matches the contextual information including user queries.
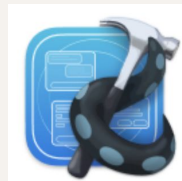
# Part 3. Prompt Generator

The prompt constructor G aims to determine the optimal combination and ranking of the components.
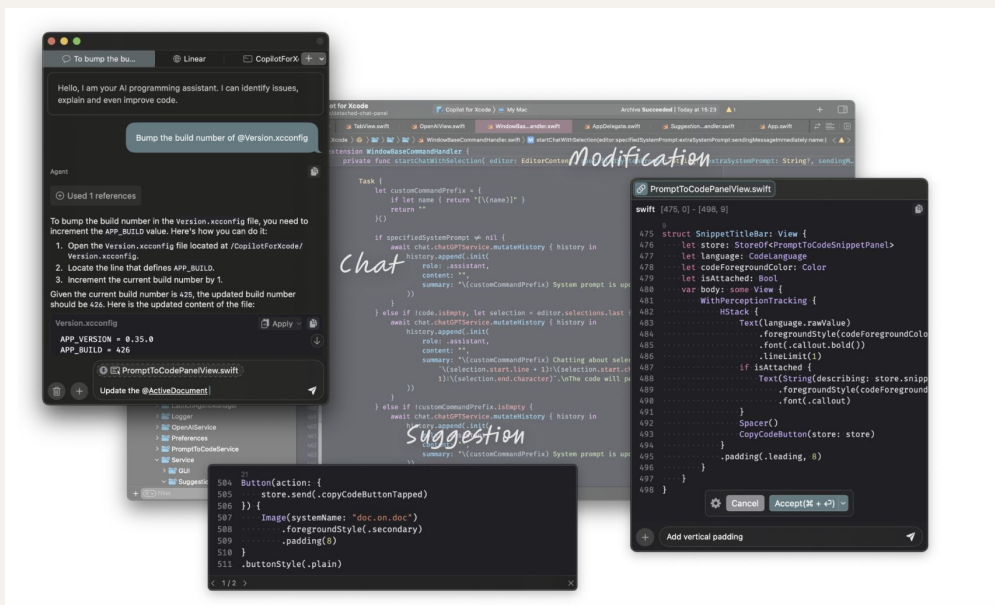
$$\mathcal{G}_\theta(x, c, z, y_{1:i-1}) = y_i = \text{order}([y^1, y^2, \ldots, y^{\tau_k}])$$
$$= [\theta_1 \quad \theta_2 \quad \ldots \quad \theta_k]^T [y^1 \quad y^2 \quad \ldots \quad y^k]^T$$

# Implementation Details on Xcode

CAMP is implemented as a plugin for Xcode, to demonstrate its practical utility and validate the methodology's robustness in challenging coding environments.
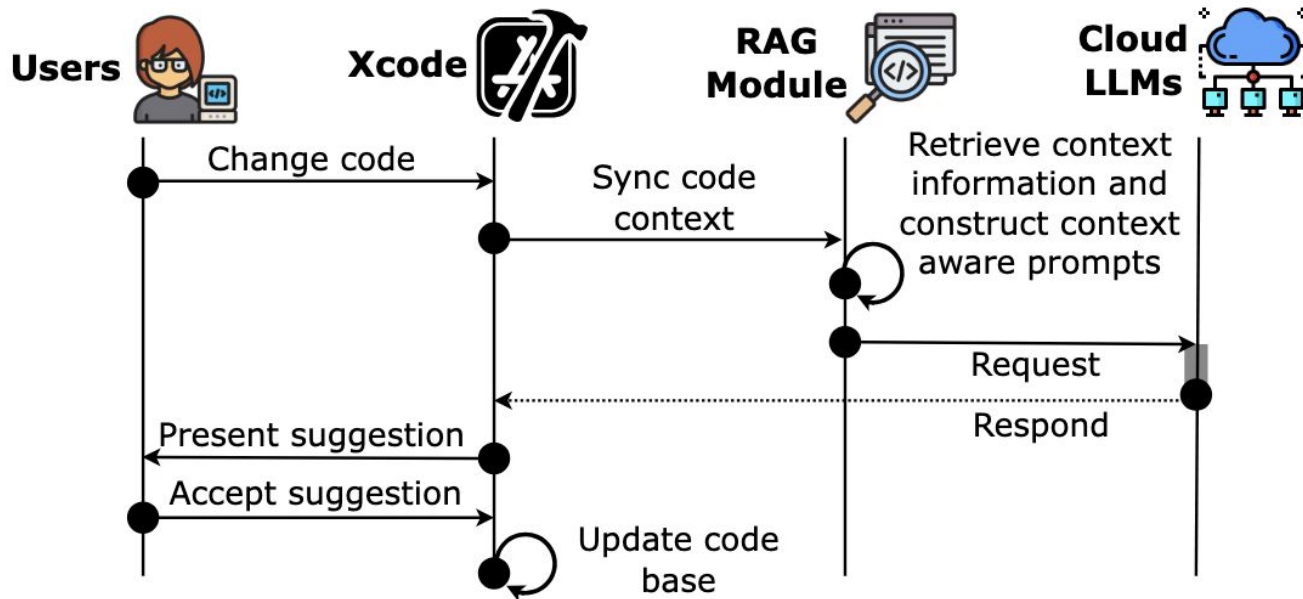
# Implementation Details on Xcode

Technical challenges and solutions:

- XPC service-level communication -> enable interaction with language servers and facilitate real-time code suggestions in the UI
- Accessibility API -> capture rich contextual data

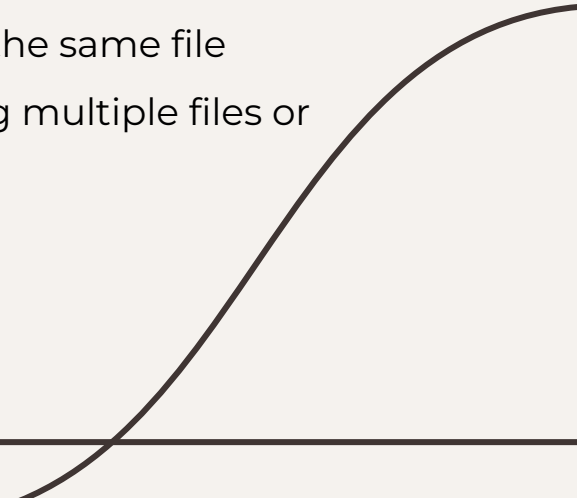Copilot for Xcode is later assimilated to Github Copilot in October 2024.

Copilot for Xcode: exploring AI-assisted programming by prompting cloud-based large language models
CW Tan, S Guo, MF Wong, CN Hang, 2023, https://arxiv.org/abs/2307.14349
Y Wang, S Guo, CW Tan, From Code Generation to Software Testing: AI Copilot with Context-Based RAG, IEEE Software, 2025

# System Flow

# Evaluation

**Experiment Setup**

- Benchmark: CoderEval
  - ***class-runnable***: Code outside the function but within the same class.
  - ***file-runnable***: Code outside the class but within the same file
  - ***project-runnable***: Code outside the file, spanning multiple files or repositories.

# Evaluation

**Baseline Models**

- We compare CAMP against the following baseline models, using GPT-3.5-Turbo as the cloud-based LLM.
  - ***CloudOnly***: Inputs are processed solely by the cloud-based model, with no local processing or context retrieval.
  - ***BaseRAG***: Implements standard RAG techniques.
  - ***FileContext***: A variant of \camp{} that prioritizes context retrieved from the currently open files in the IDE. This lightweight version balances performance and resource efficiency.

# Evaluation

CAMP consistently outperforms baselines

- +12.5% vs. CloudOnly
- +6.3% vs. BaseRAG
- Best results for complex, multi-file tasks

Analysis and discussion

- CAMP's context-aware code content retrieval enhances prompts
- FileContext is competitive at simple tasks, but falls short in large-scale ones

Insights and learnings

- Broader context contributes better to complex tasks.
- Dynamic scope tuning saves compute on simple tasks without sacrificing accuracy.

# Future Directions

We propose directions for future research works, in algorithm optimization, model refinement, and software features extensions.

- **End to End Training**: To train all parameters end-to-end in one data pipeline, which might bring new insights to our understanding of the model parameters.
- **Trust AI**: To implement a robust user consent mechanism that empowers users with the ability to control the data access and make informed decisions about the extent of information shared.
- **Software Feature Extensions**: To develop a portable functionality kit to maximize the tool's compatibility with future versions of Xcode and multiple programming languages.