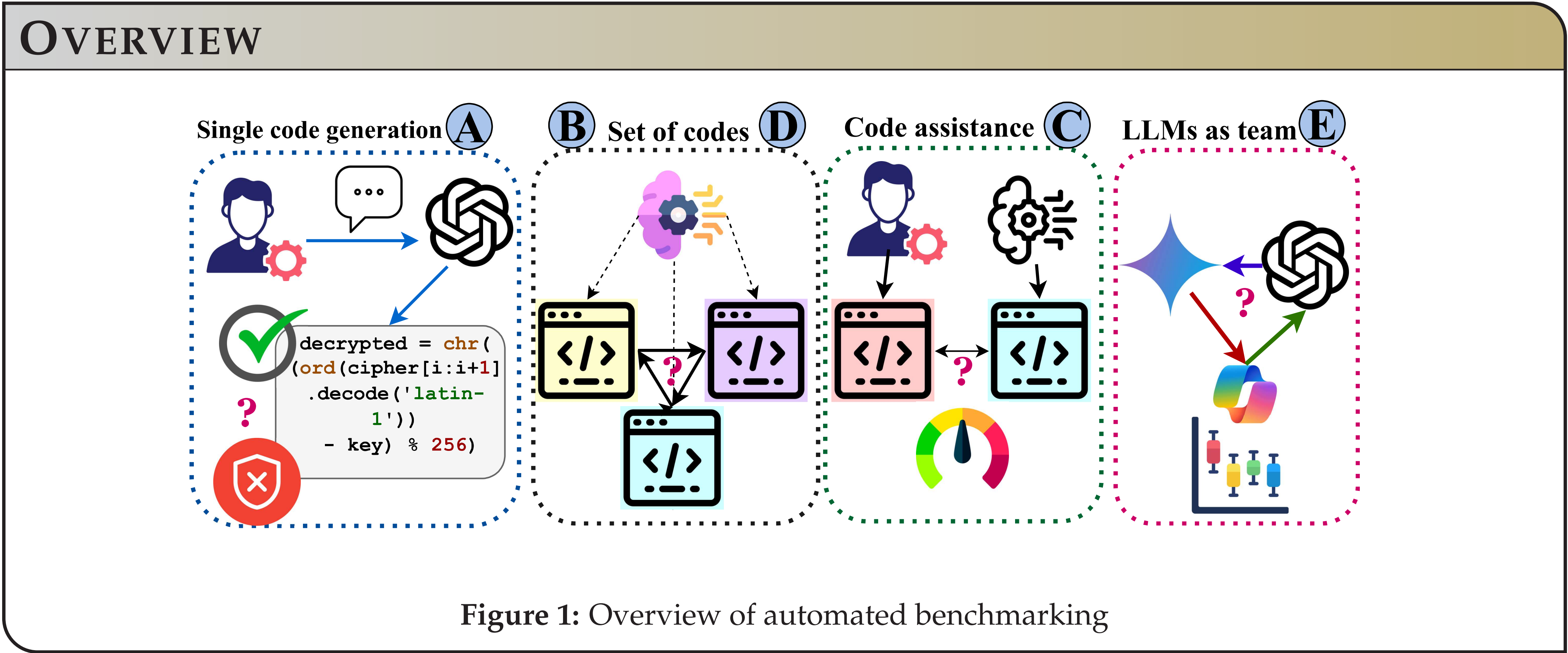


## ABSTRACT

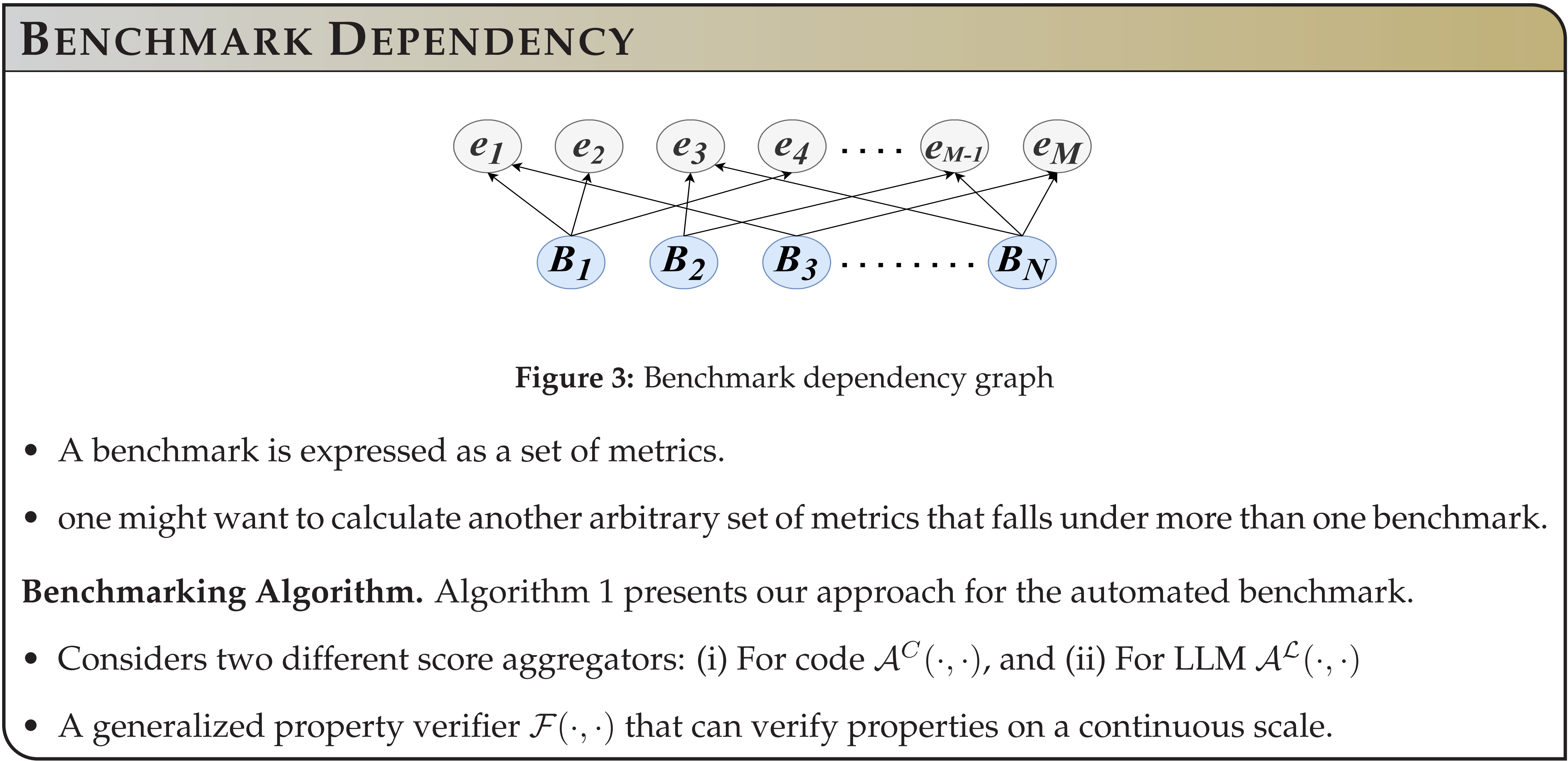
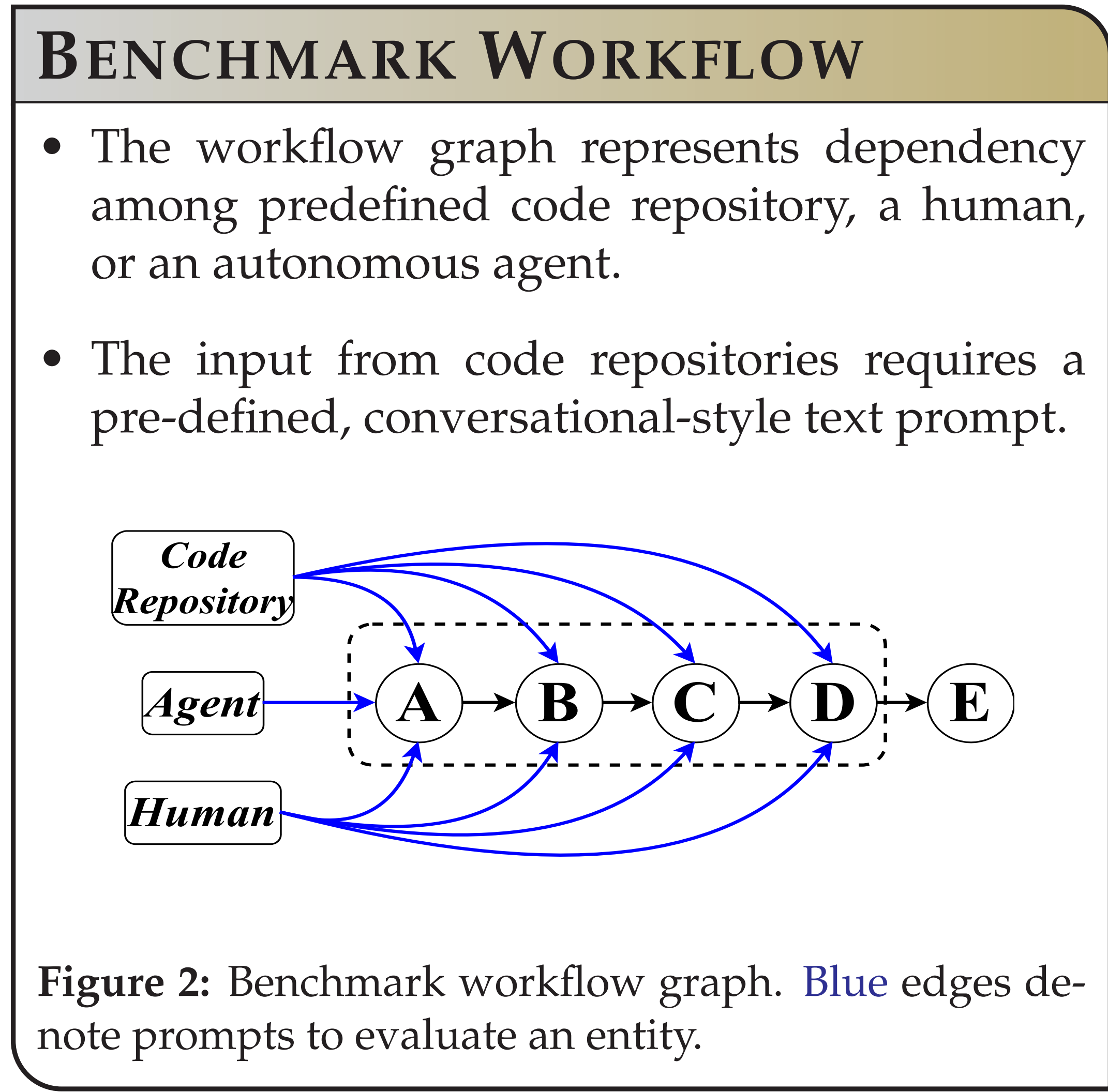
- Generative LLMs have proven to be valuable code generators, thus enabling code copilots and meeting several requirements in software engineering.
- But few questions remain-
  - How good is an LLM in software projects?
  - How secure is the code generated by them?
  - Can they work in a team?
- In this research, we address these questions by proposing a holistic benchmarking framework.



## FRAMEWORK

In our benchmark framework, we holistically consider five entities:

- (A) Single code segment quality,
- (B) Quality of set of code segments from a single task-specific prompt,
- (C) LLM performance as a Co-pilot to assist sub-completed code,
- (D) Code comprehensiveness of a single LLM,
- (E) Performance of LLMs as a team.



## ALGORITHM

**Algorithm 1** Holistic Benchmarking Algorithm

**Input:** Properties  $\Phi : \{\phi_1, \phi_2, \dots, \phi_d\}$ , Generated codes  $\zeta : \{C_1, C_2, \dots, C_n\}$ , LLMs  $\Lambda : \{\mathcal{L}_1, \mathcal{L}_2, \dots, \mathcal{L}_p\}$

**Require:** Score Aggregators  $\Rightarrow$  Code:  $\mathcal{A}^C(\cdot, \cdot)$ , LLM:  $\mathcal{A}^L(\cdot, \cdot)$ ; Property verifier  $\mathcal{F}(\cdot, \cdot)$

**Output:** Code score set  $[e]_n$ , LLM score set  $[\beta]_p$

```

for  $a \leftarrow 1$  to  $p$  do
     $\beta^{\mathcal{L}_a} \leftarrow 0$  ▷ Initially all LLM scores are 0
end for
for  $i \leftarrow 1$  to  $n$  do
     $e^i \leftarrow 0$  ▷ Initially  $i$ -th code score is 0
    for  $j \leftarrow 1$  to  $d$  do
         $S_j^i \leftarrow \mathcal{F}(C_i^{\mathcal{L}_a}, \phi_j)$ ,  $e^i \leftarrow \mathcal{A}^C(e^i, e_j^i)$ 
         $\beta^{\mathcal{L}_a} \leftarrow \mathcal{A}^L(\beta^{\mathcal{L}_a}, e^i)$ 
    end for
end for
    
```

## EVALUATION MODEL

$$B_{\mathcal{L}}(t_i) = \sum_{j=1}^m w_{ij} \hat{f}(t_i, e_j), \quad B_{\mathcal{L}} = \frac{1}{n} \sum_{i=1}^n B_{\mathcal{L}}(t_i)$$

- $E = \{e_1, e_2, \dots, e_M\}$  is the set of evaluation metrics to assess the generated code snippet quality. Evaluation function  $f : T \times E \rightarrow R$  maps each task  $t_i$  and evaluation metric  $e_j$  to a real-valued score.
- $w_{ij}$  denote the weight assigned to evaluation metric  $e_j$ . Here,  $w_{ij} \geq 0$ ,  $\sum_{j=1}^m w_{ij} = 1$  for each task  $t_i$ .
- $B_{\mathcal{L}} : T \rightarrow R$  works as a function of the performance scores obtained across all tasks and evaluation metrics.

## CONCLUSION & FUTURE RESEARCH

- We propose a benchmark framework for LLMs based on their code generation capabilities and their ability to be a comprehensive development tool.
- We are developing a platform to integrate multiple open-source and black-box models and property verifiers such as security, code correctness, and hallucination.
- In future, we plan to incorporate more complex tasks, such as static analysis, in the benchmark framework.