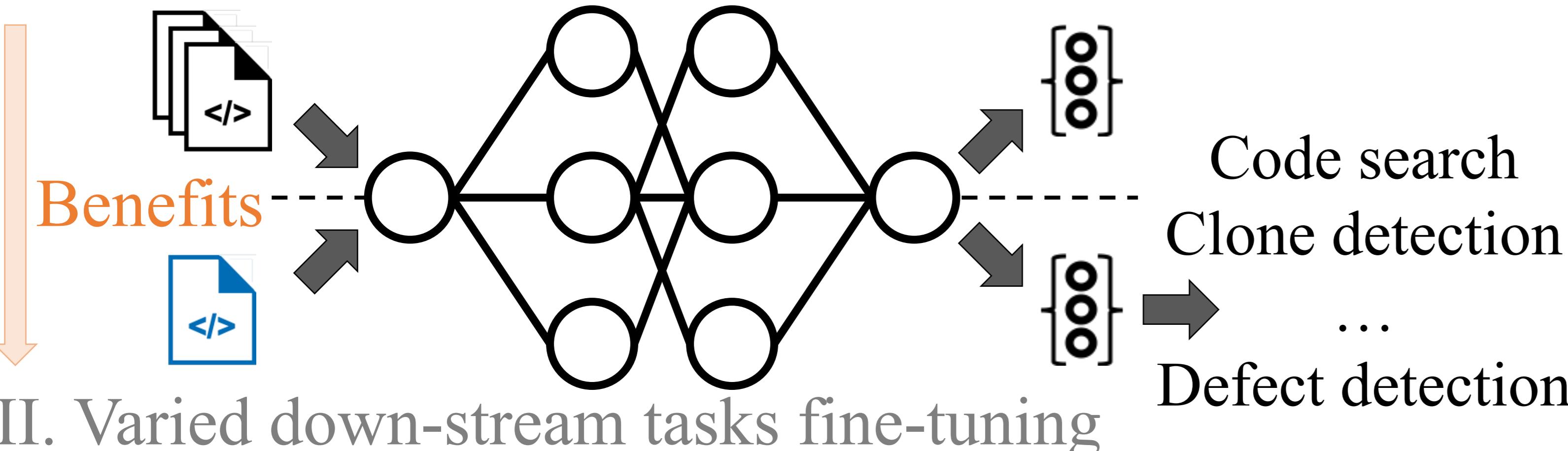


Jiabo Huang¹, Jianyu Zhao¹, Yuyang Rong², Yiwen Guo^{*3}, Yifeng He², Hao Chen²¹Tencent Security Big Data Lab, ²UC Davis, ³Independent Researcher

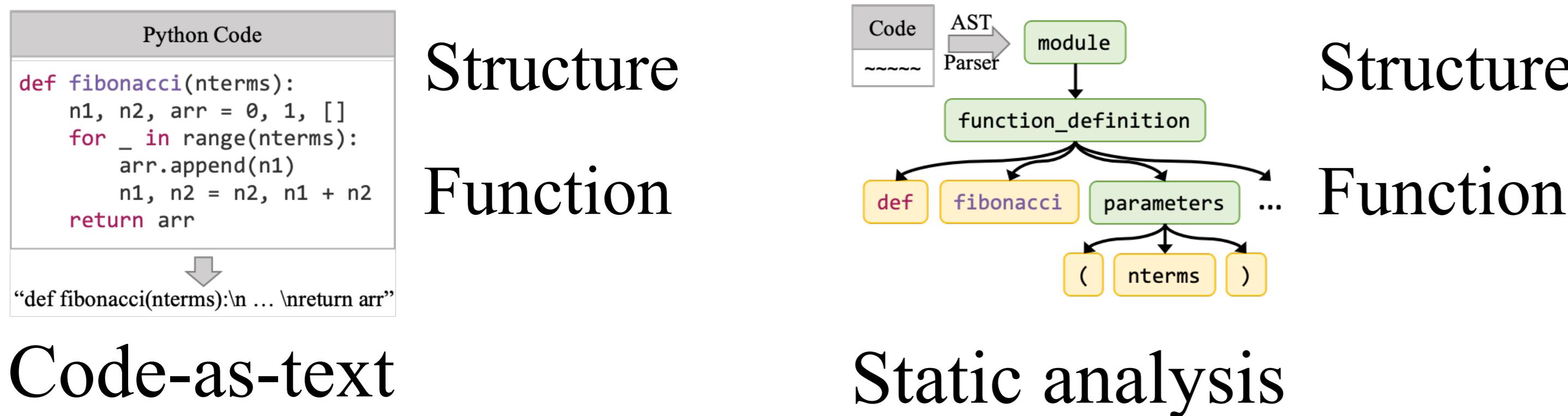
{jiabohuang, yjjyzhao}@tencent.com, {PeterRong96, guoyiwen89}@gmail.com, {yfhe, chen}@ucdavis.edu

1 Problem Definition

I. Unsupervised Representation Pre-training



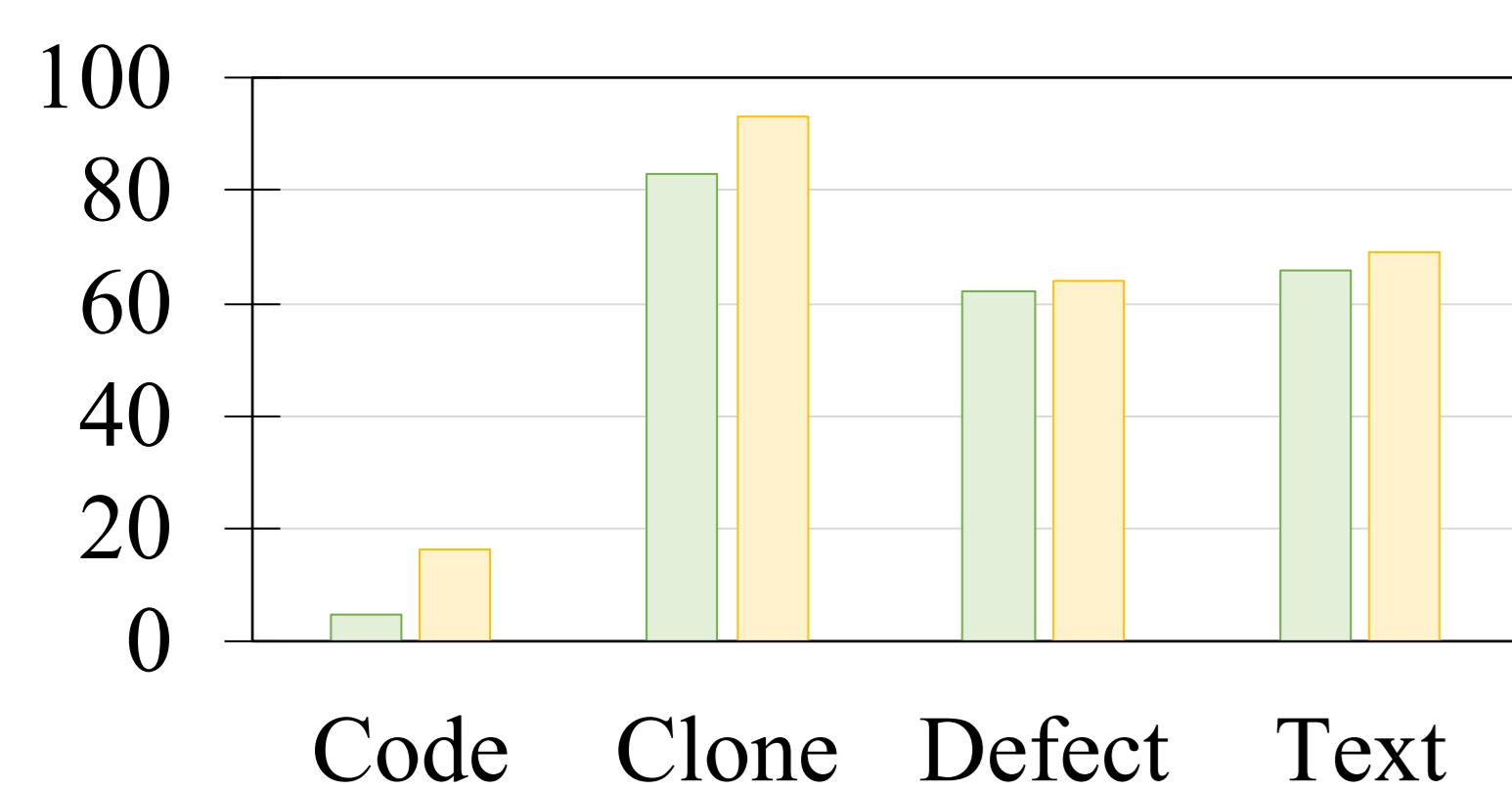
2 Motivation



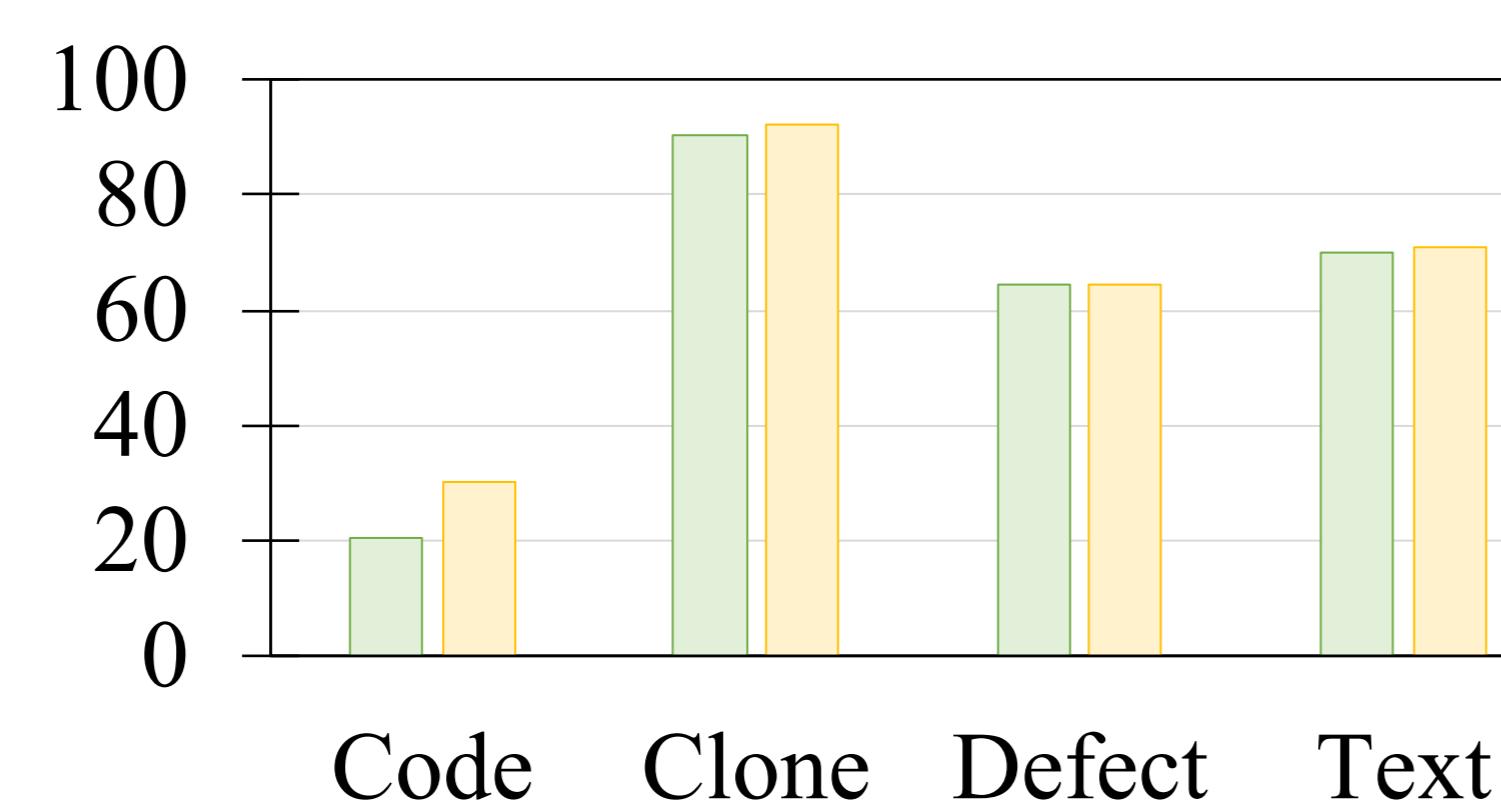
5 Experiments

- FuzzPretrain is beneficial to multiple code understanding tasks
- Dynamic information from test cases complements both code and its syntactic representations (AST)

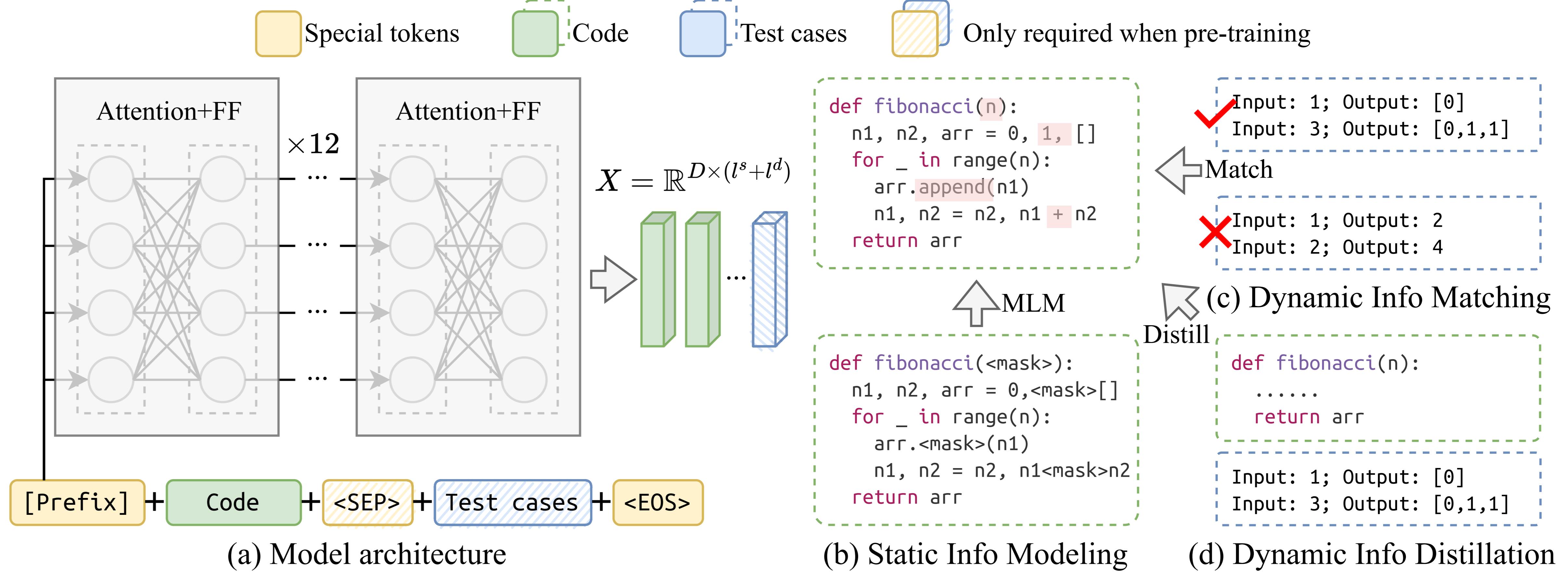
CodeBERT FuzzCodeBERT



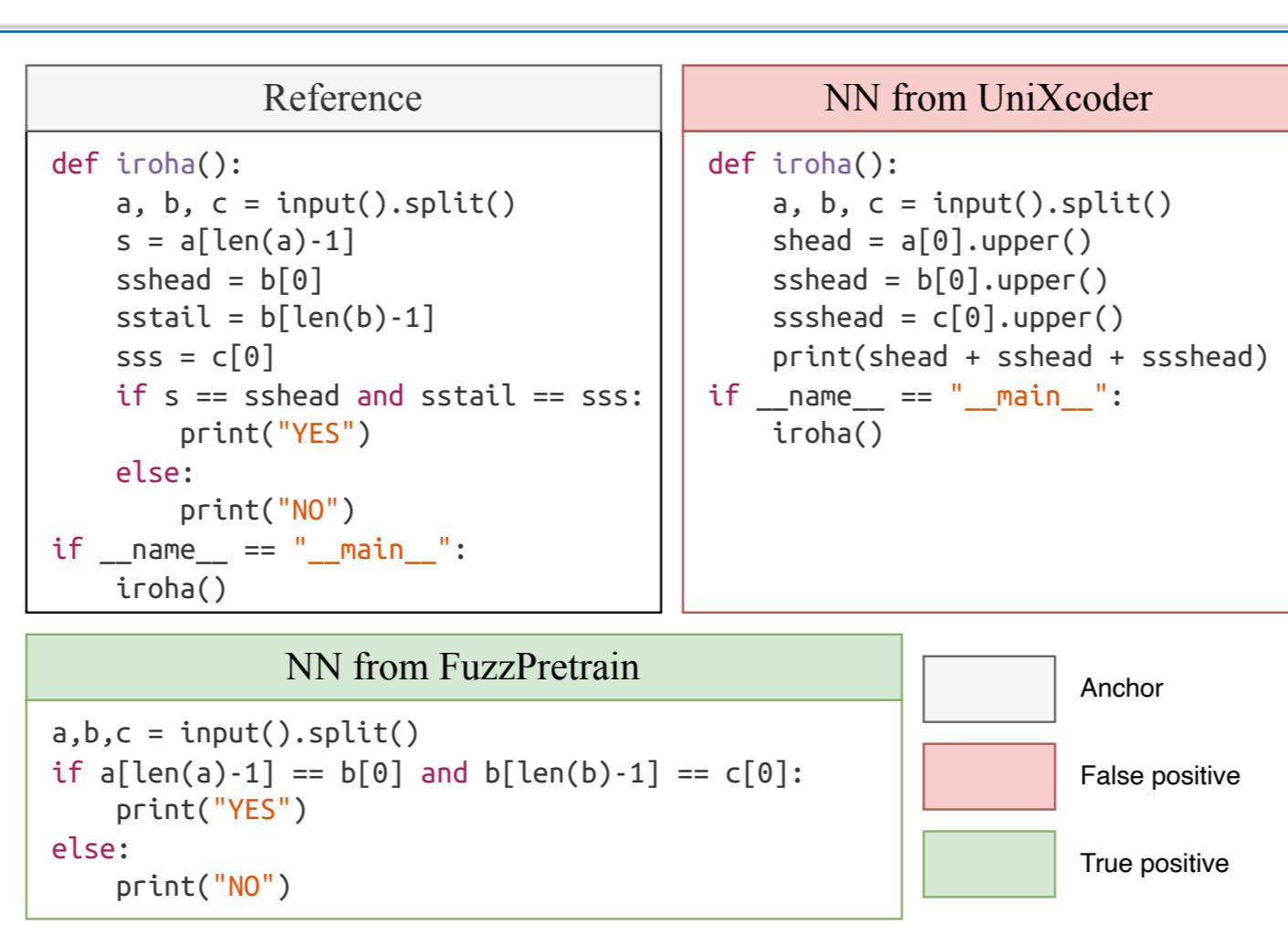
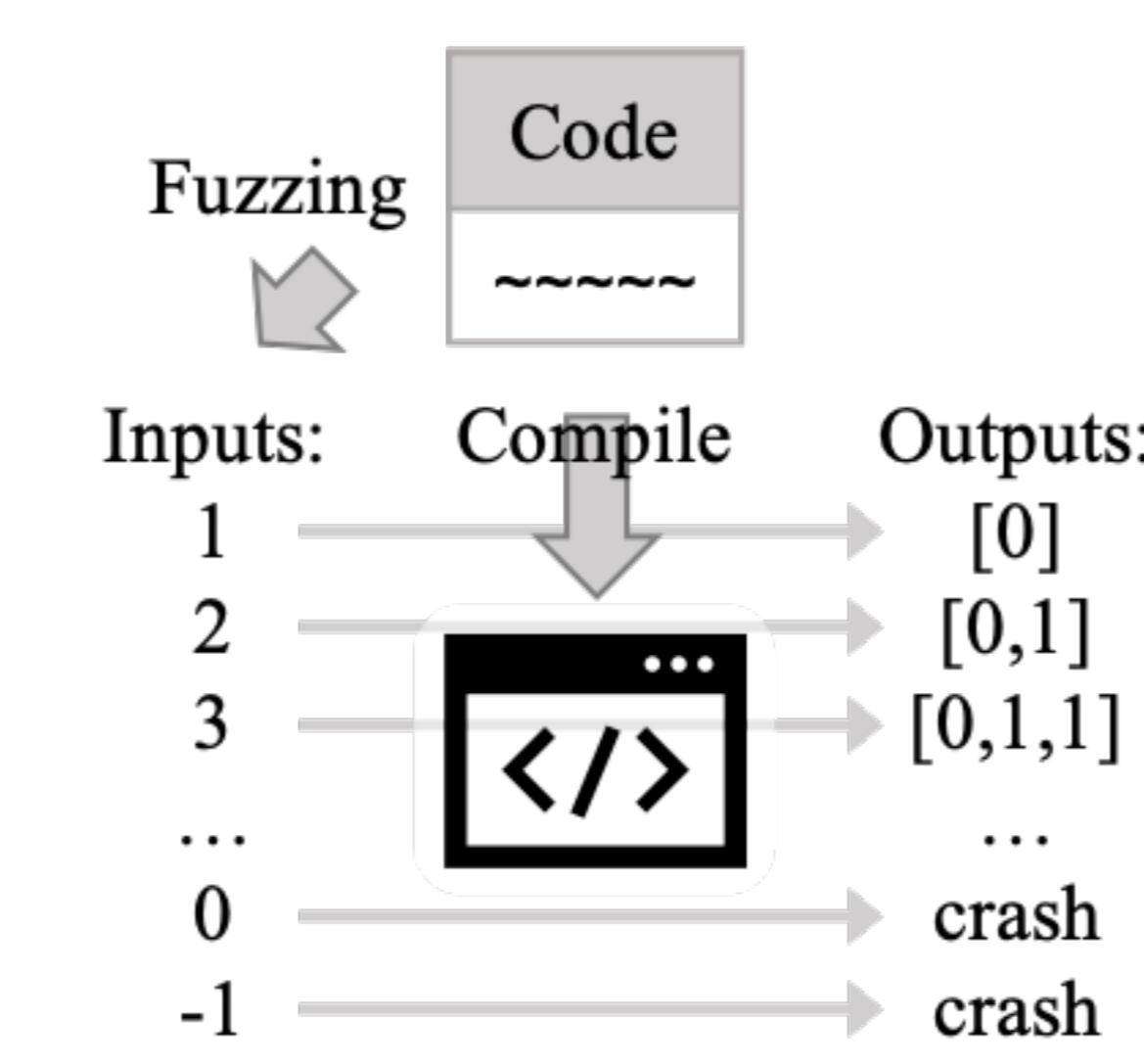
UniXcoder FuzzUniXcoder



3 Overview of FuzzPretrain



4 Data collection



5 Static & Dynamic Information Modelling

- Static Information Modelling (SIM)

Masked tokens prediction on code (S)

$$\mathcal{L}_{SIM}(S) = - \sum_{m \in M} \log(p(m | \tilde{X}^S))$$

- Dynamic Information Modelling (DIM)

Matching code (S) with test cases (D)

$$\mathcal{L}_{DIM}(S, D) = BCE(y, f_\phi(FC(x^h)))$$

- Dynamic Information Distillation (DID)

*Distilling dynamic info from holistic representation**H = S ⊕ D to code (S)*

$$\mathcal{L}_{DID}(S, S \oplus D) = -\log \frac{g(\hat{x}^h, x^s)}{g(\hat{x}^h, x^s) + \sum_{x^- \in X^-} g(x^-, x^s)}$$