

# Contrastive Code Representation Learning

Paras Jain\*, Ajay Jain\*, Tianjun Zhang, Pieter Abbeel, Joseph E. Gonzalez, Ion Stoica

\* equal contribution

Code and paper:  
parasj.github.io/contracode

Email me:  
ajayj@berkeley.edu

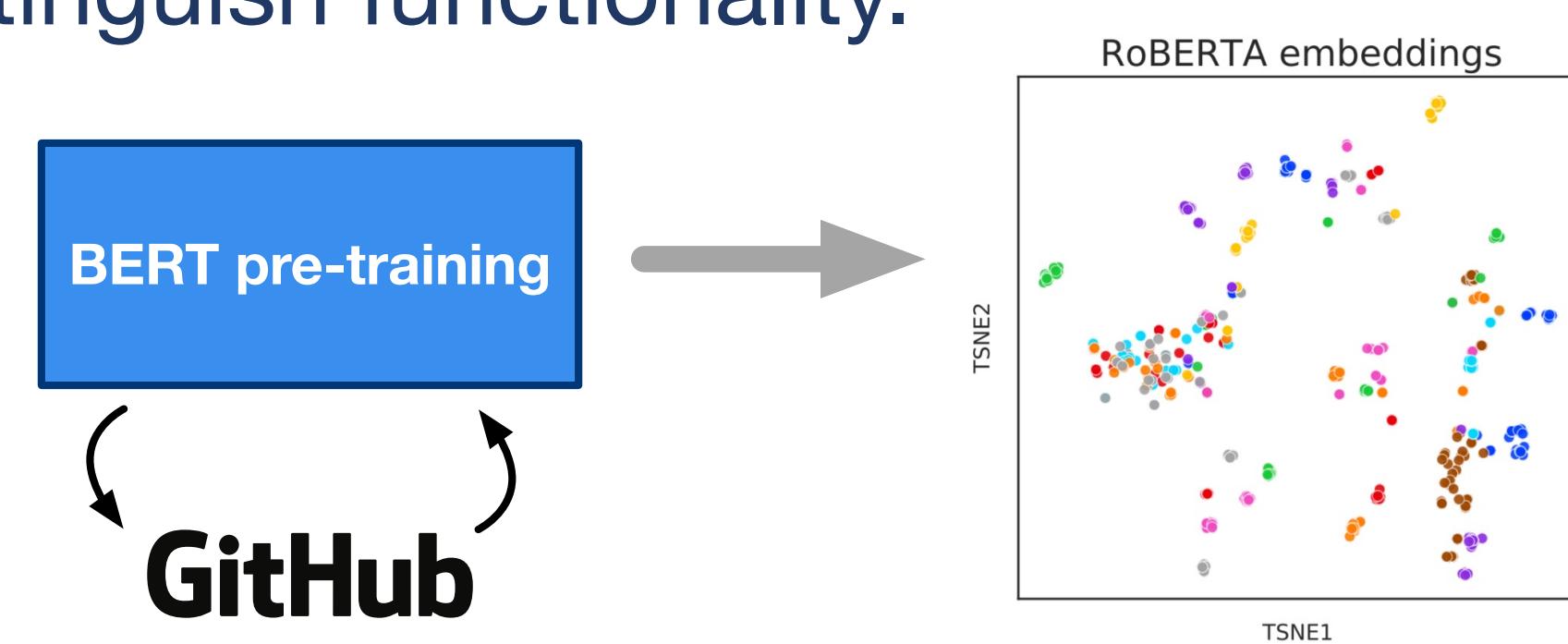


## Motivation

- ML improves & accelerates code analysis
- Problem: Learned code analysis is **data hungry and sensitive to minor implementation details**
- How to represent the *functionality* of a program in an unsupervised manner?

## Unsupervised representation of code

RoBERTa can pre-train on large corpus of unlabeled code, but representations do not distinguish functionality.



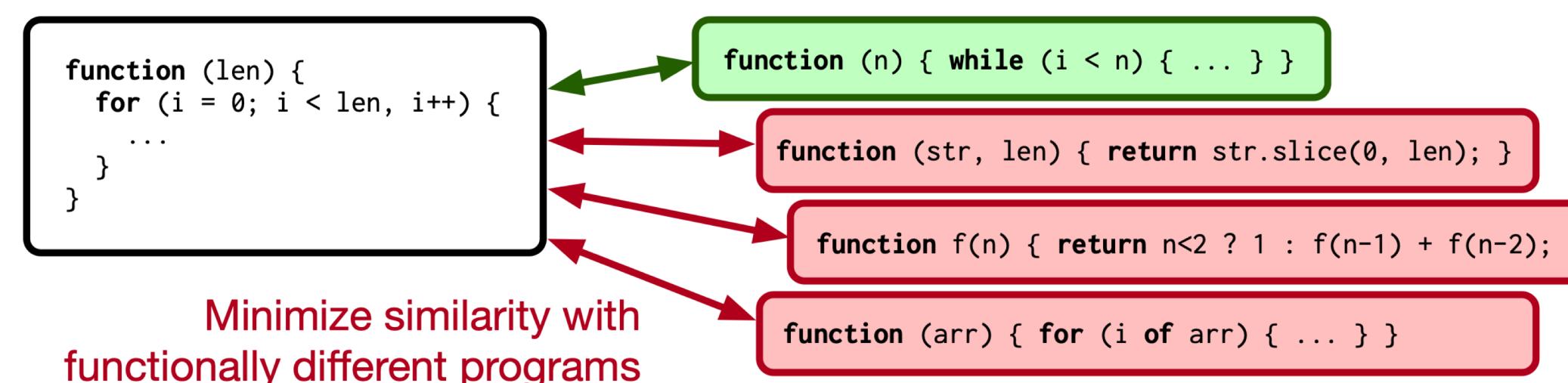
## Contrastive learning is great for code

Programs with the **same functionality** should have **similar representations!**

ContraCode pre-trains code encoders with an unsupervised contrastive objective to encode prior knowledge of functionality into the representation.

$$\mathcal{L}_{q, k^+, k^-} = -\log \frac{\exp(q \cdot k^+/t)}{\exp(q \cdot k^+/t) + \sum_{k^-} \exp(q \cdot k^-/t)}$$

Given a program,

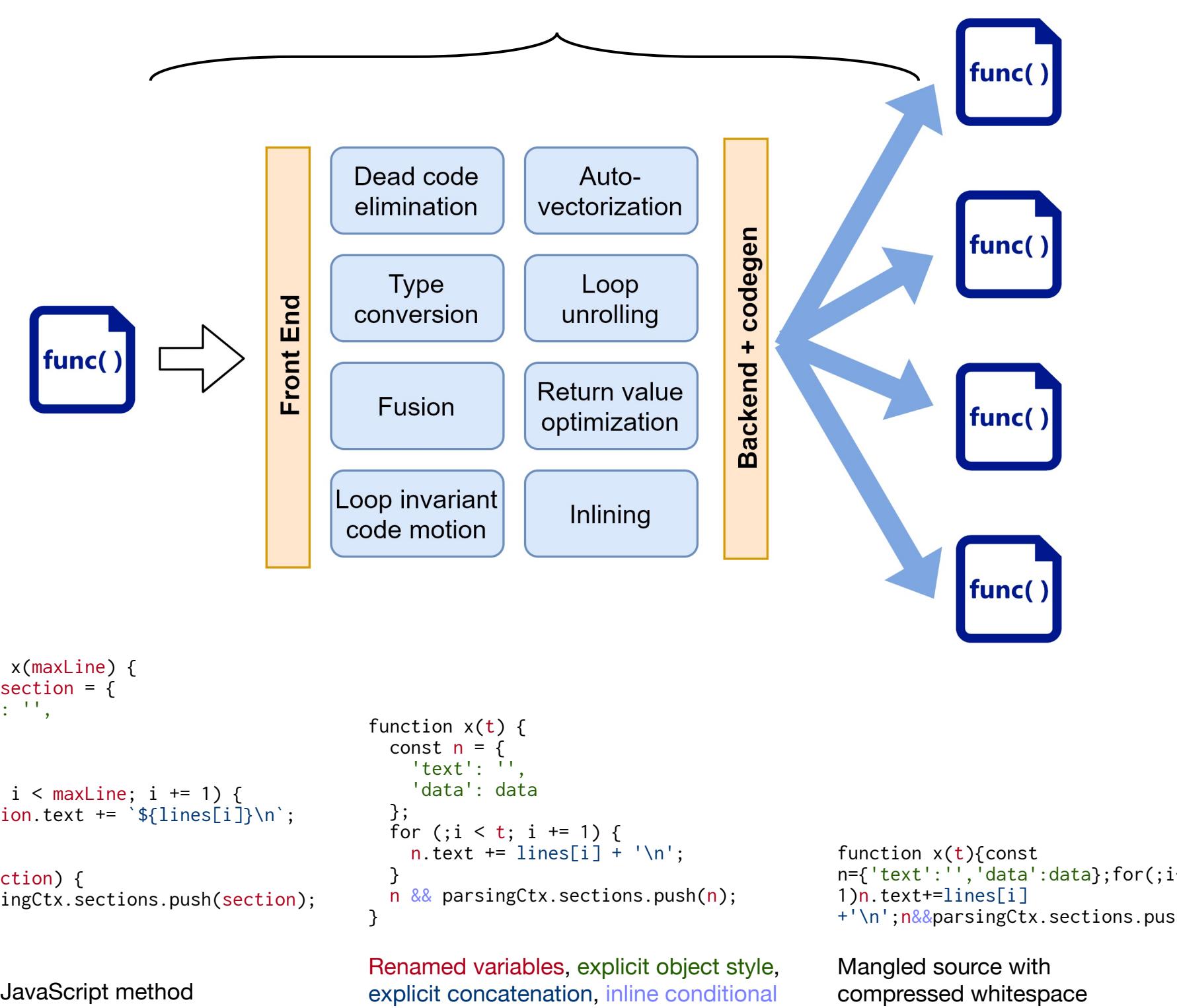


## Data generation with compilers

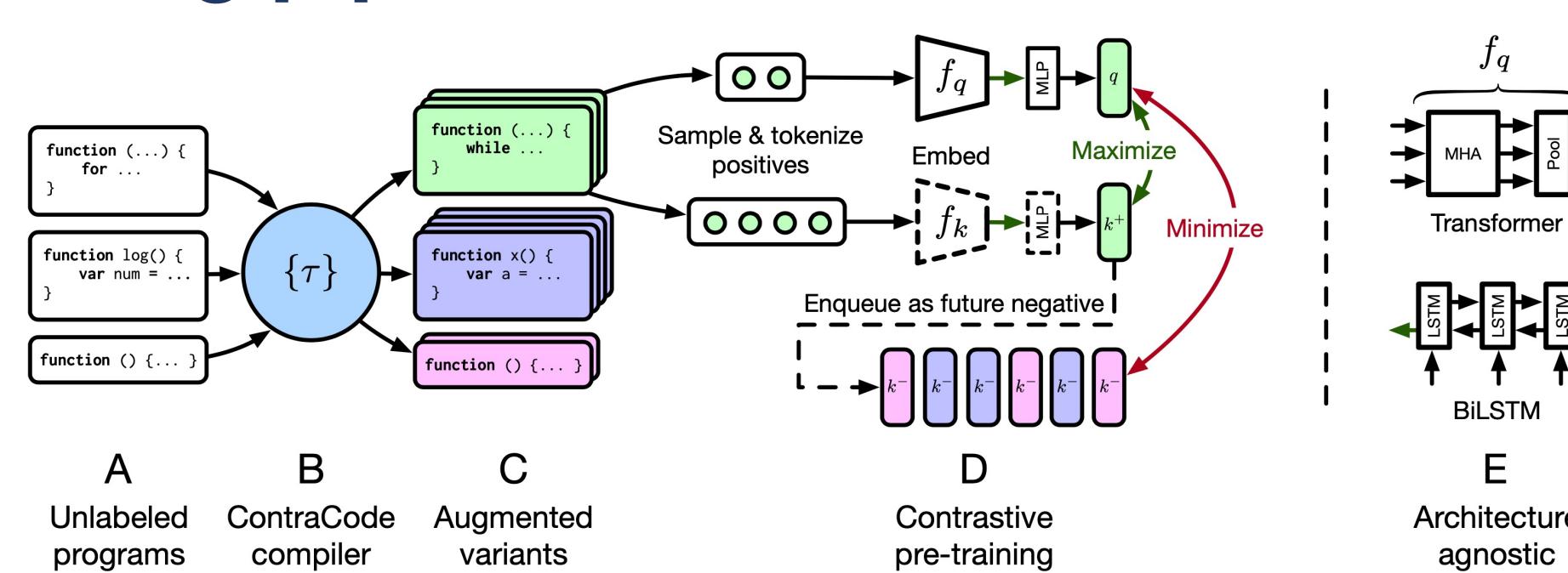
- Using source-to-source compiler transforms, we generate diverse pairs of equivalent of programs.

Code compression		Identifier modification	
✓	Reformatting (R)	✓	Variable renaming (VR)
✓	Beautification (B)	✓	Identifier mangling (IM)
✓	Compression (C)	✓	Regularization
✓	Dead-code elimination (DCE)	✓	Dead-code insertion (DCI)
✓	Type upconversion (T)	✓	Subword regularization (SW)
✓	Constant folding (CF)	✗	Line subsampling (LS)

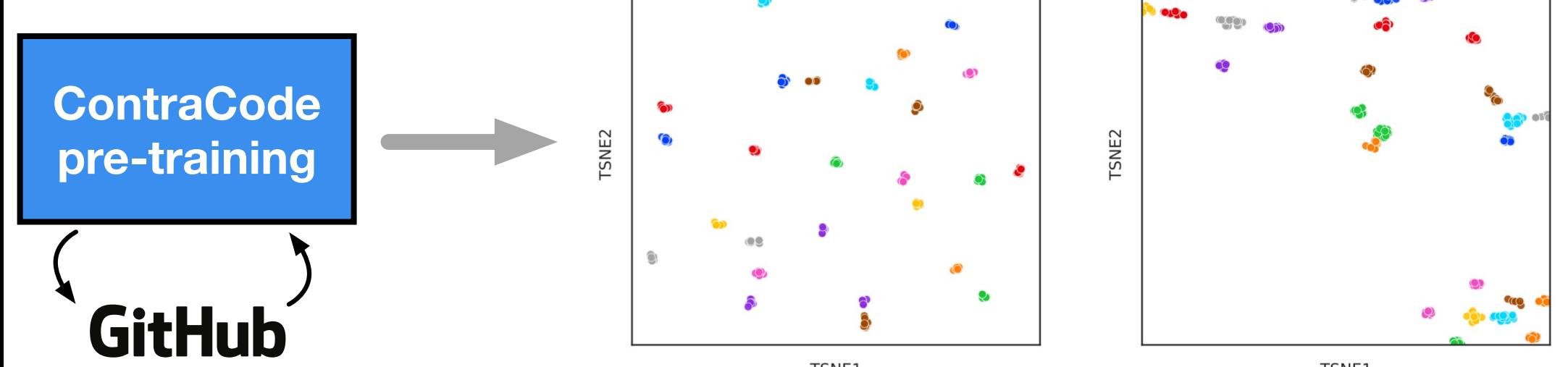
✓ = semantics-preserving transformation ✗ = lossy transformation



## Training pipeline



## Visualizing the representation



## Type inference

**JavaScript → TypeScript**

Method	Acc@1	Acc@5
TypeScript CheckJS	45.11%	—
DeepTyper, variable name only	28.94%	70.07%
GPT-3 Codex (zero-shot, 175B)	26.62%	—
GPT-3 Codex (few-shot, 175B)	30.55%	—
Transformer	45.66%	80.08%
+ RoBERTa MLM pre-train	40.85%	75.76%
+ ContraCode pre-train	46.86%	<b>81.85%</b>
+ ContraCode + MLM (hybrid)	47.16%	81.44%
DeepTyper BiLSTM	51.73%	82.71%
+ RoBERTa MLM pre-train	50.24%	82.85%
+ ContraCode pre-train	54.01%	<b>85.55%</b>

function rendererError ( message: string 99.5% ) : void 99.7% {

## Code summarization (JavaScript → English)

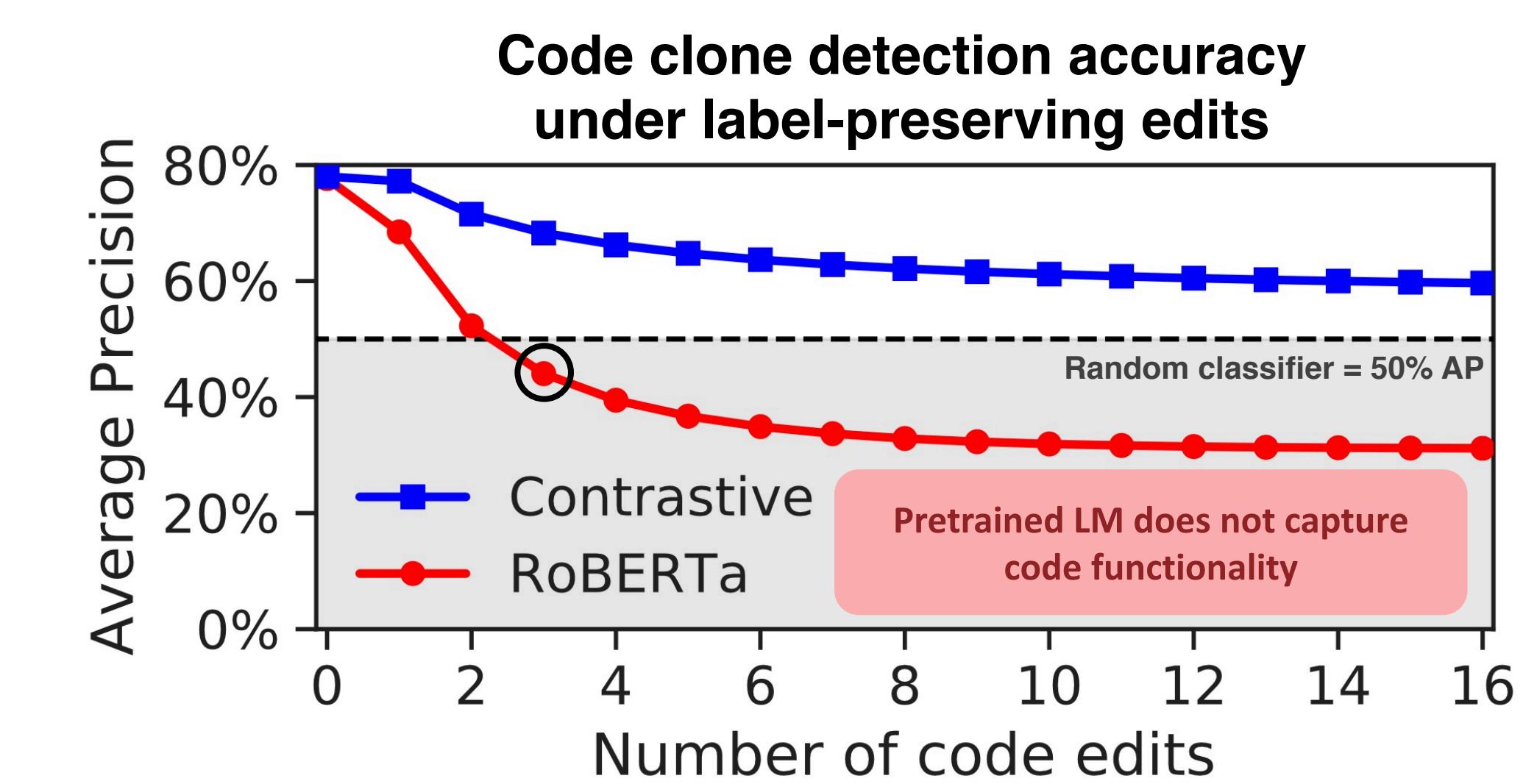
Generate short English summary of a method.

Method	Precision	Recall	F1
code2vec	10.78%	8.24%	9.34%
code2seq	12.17%	7.65%	9.39%
RoBERTa MLM	15.13%	11.47%	12.45%
Transformer	18.11%	15.78%	16.86%
+ ContraCode	20.34%	14.96%	<b>17.24%</b>

Contrastive pre-train

## Zero-shot Clone Detection (JavaScript)

Detect when human-written programs are equivalent.



Natural code	AUROC	AP
Edit distance heuristic	$69.55 \pm 0.81$	73.75
Randomly initialized Transformer	$72.31 \pm 0.79$	75.82
+ RoBERTa MLM pre-train	$74.04 \pm 0.77$	77.65
+ ContraCode pre-train	$75.73 \pm 0.75$	78.02
+ ContraCode + RoBERTa MLM	<b><math>79.39 \pm 0.70</math></b>	<b>81.47</b>

RoBERTa MLM pre-train  
+ 1.8% AP  
Contrastive pre-train  
+ 2.2% AP  
Hybrid pre-train  
+ 5.56% AP