



Code Intelligence: Models, Applications and Future

Nan Duan

Senior Principal Research Manager

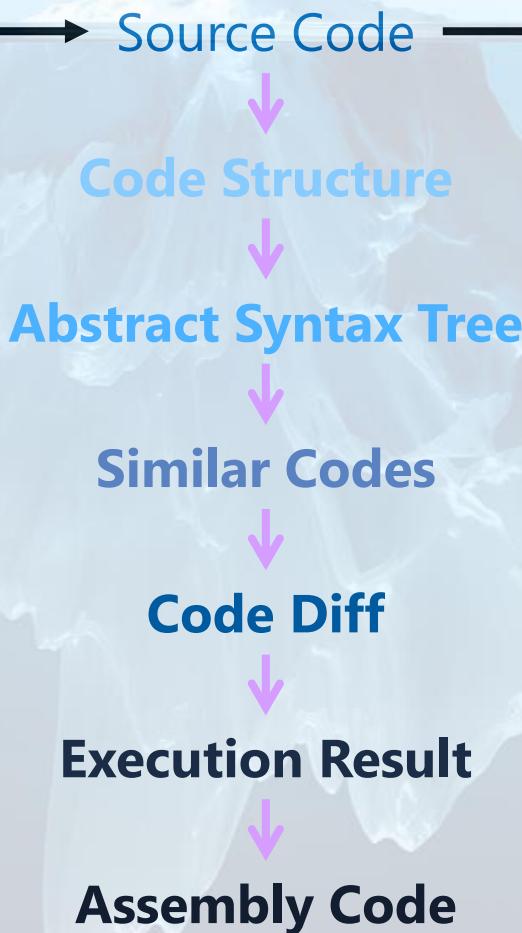
Microsoft Research Asia

Korea AI Summit @ 2022-12-15

Large-scale Pre-trained Models for Code Intelligence



treat code as natural language



OpenAI Codex

We've created an improved version of OpenAI Codex, our AI system that translates natural language to code, and we are releasing it through our API in private beta starting today. Codex is the model that powers [GitHub Copilot](#), which we built and launched in partnership with GitHub a month ago. Proficient in more than a dozen programming languages, Codex can now interpret simple commands in natural language and execute them on the user's behalf—making it possible to build a natural language interface to existing applications. We are now inviting businesses and developers to build on top of OpenAI Codex through our API.

Codex (OpenAI)

Introducing Text and Code Embeddings in the OpenAI API

We are introducing embeddings, a new endpoint in the OpenAI API that makes it easy to perform natural language and code tasks like semantic search, clustering, topic modeling, and classification. Embeddings are numerical representations of concepts converted to number sequences, which make it easy for computers to understand the relationships between those concepts. Our embeddings outperform top models in 3 standard benchmarks, including a 20% relative improvement in code search.

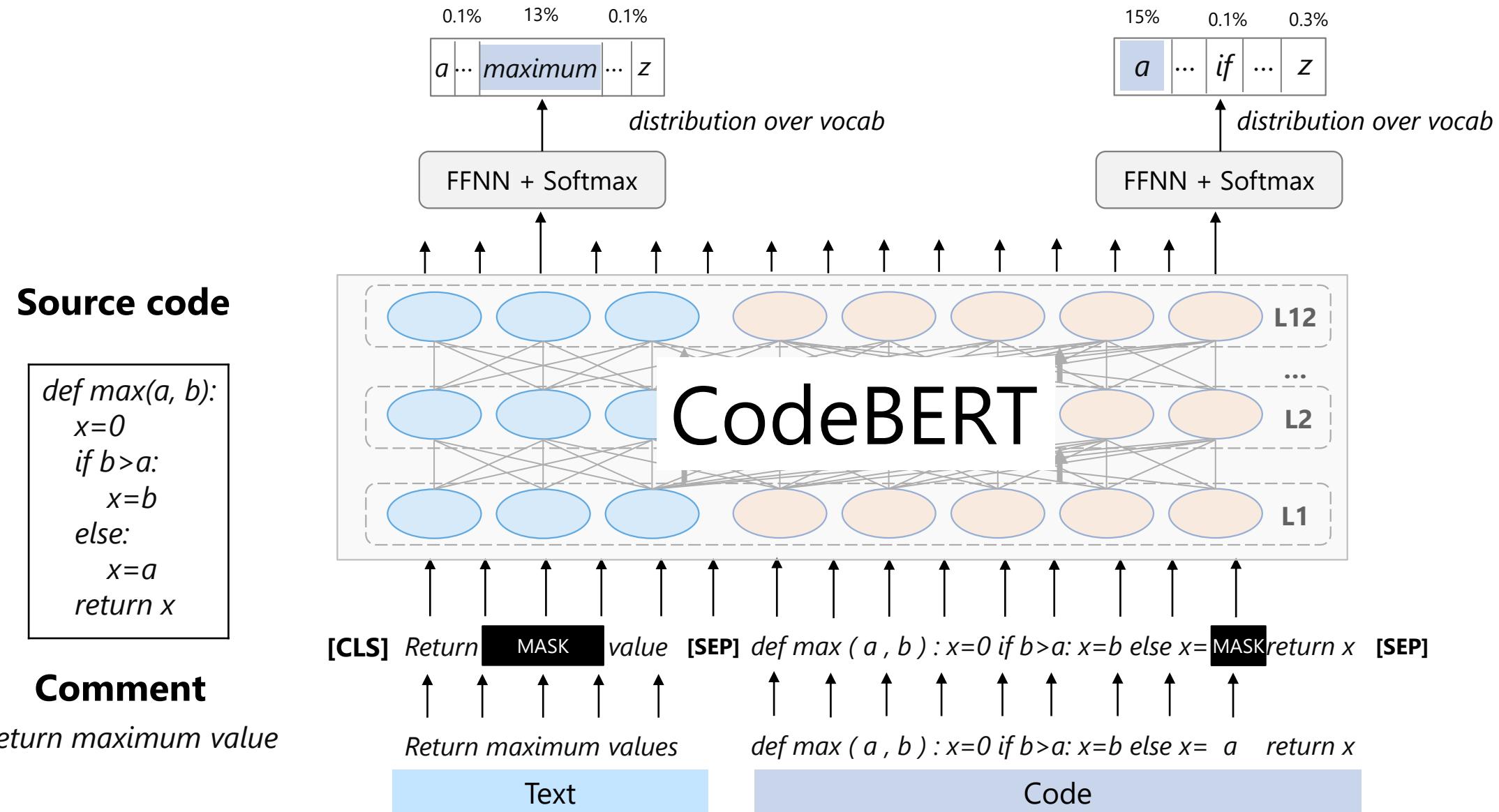
Embeddings (OpenAI)

(Some of) Our Work @ MSRA

- **CodeBERT** (EMNLP 2020): propose the 1st code-text pre-trained model
- **GraphCodeBERT** (ICLR 2021): use variable relationship to enhance code representation
- **UniXcoder** (ACL 2022): use sequential AST to enhance code representation
- **AR2** (ICLR 2022): propose a general dense retrieval framework that trains retriever & ranker in a minimax adversarial manner and can cover text-to-text, text-to-code, and code-to-code tasks
- **CodeRetriever** (EMNLP 2022): extract code-text and code-code pairs from Web for building code retrieval models with unimodal and bimodal contrastive learning
- **GPT-C w/ Extended Context** (EMNLP 2021): improve GPT-C for code completion with extended context
- **Grammformer** (ICLR 2022): generate codes based on code grammar and predict holes to alleviate the uncertainty issue
- **ReACC** (ACL 2022): propose a retrieval-augmented code completion framework, which uses the partial code as a query to retrieve similar codes and predicts the following codes based on the retrieved similar codes
- **CodeReviewer** (ESEC/FSE 2022): propose an encoder-decoder based pre-trained model for automating code review activities, such as code diff quality estimation, review generation, code refinement based on code review, etc.
- **CodeReviewer v2.0** (on-going 2022): propose a retrieval-augmented model for automating code review activities, like ReACC for code completion.
- **CodeXGLUE** (NeurIPS 2021): build a benchmark for code intelligence, which covers 14 datasets for 10 code tasks
- **CoSQA** (ACL 2021): build a benchmark for text-to-code search based on Web data
- **CodeExplanation** (EMNLP 2022): build a benchmark for code-to-explanation generation



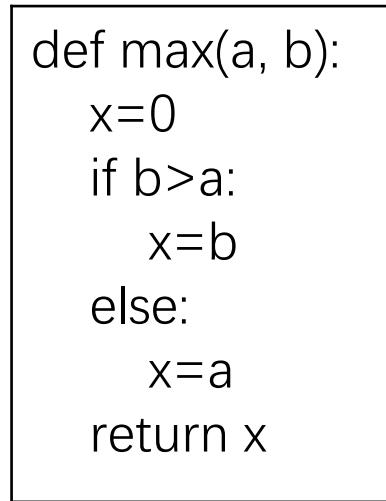
CodeBERT (v1): Pre-Train with Code+Text



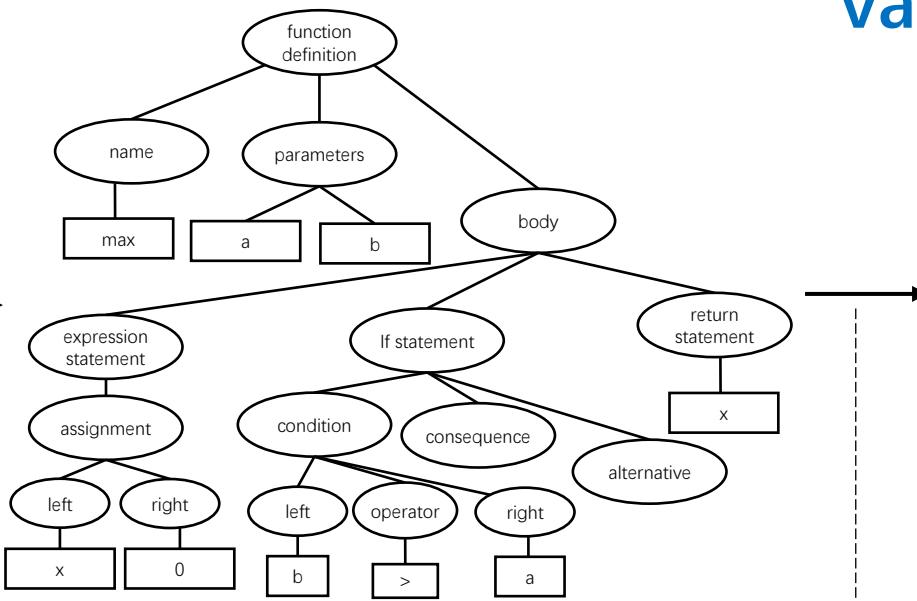
Extract Variable Relationship (VarRel) for Code Pre-training

Source code

Parse into AST



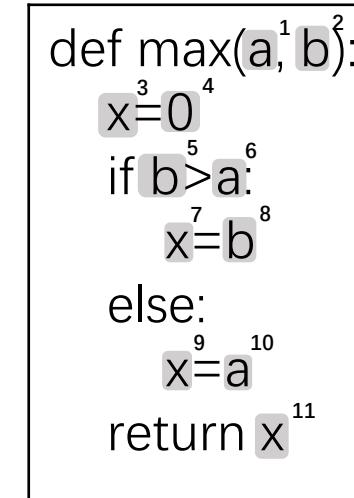
TreeSitter (public tool)



Identify variable sequence in AST

Identify variable sequence

Variable relationship



Value comes from

Extract variable relationship
from AST according to
paths between variables

GraphCodeBERT (v2): Pre-Train with Code+Text+VarRel

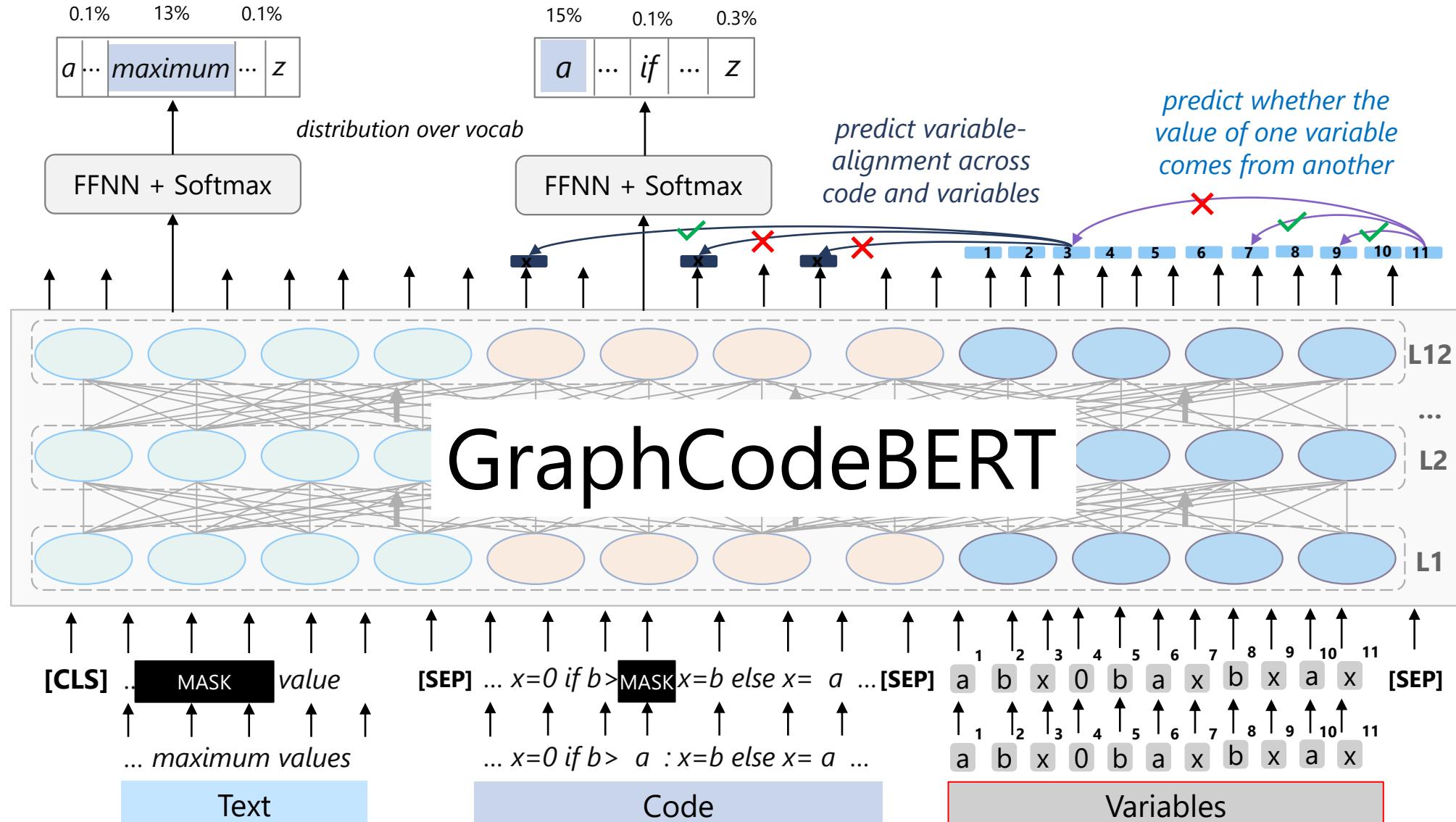
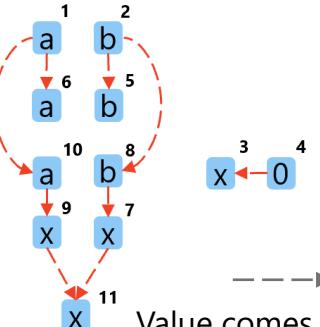
Source code

```
def max(a, b):
    x=0
    if b>a:
        x=b
    else:
        x=a
    return x
```

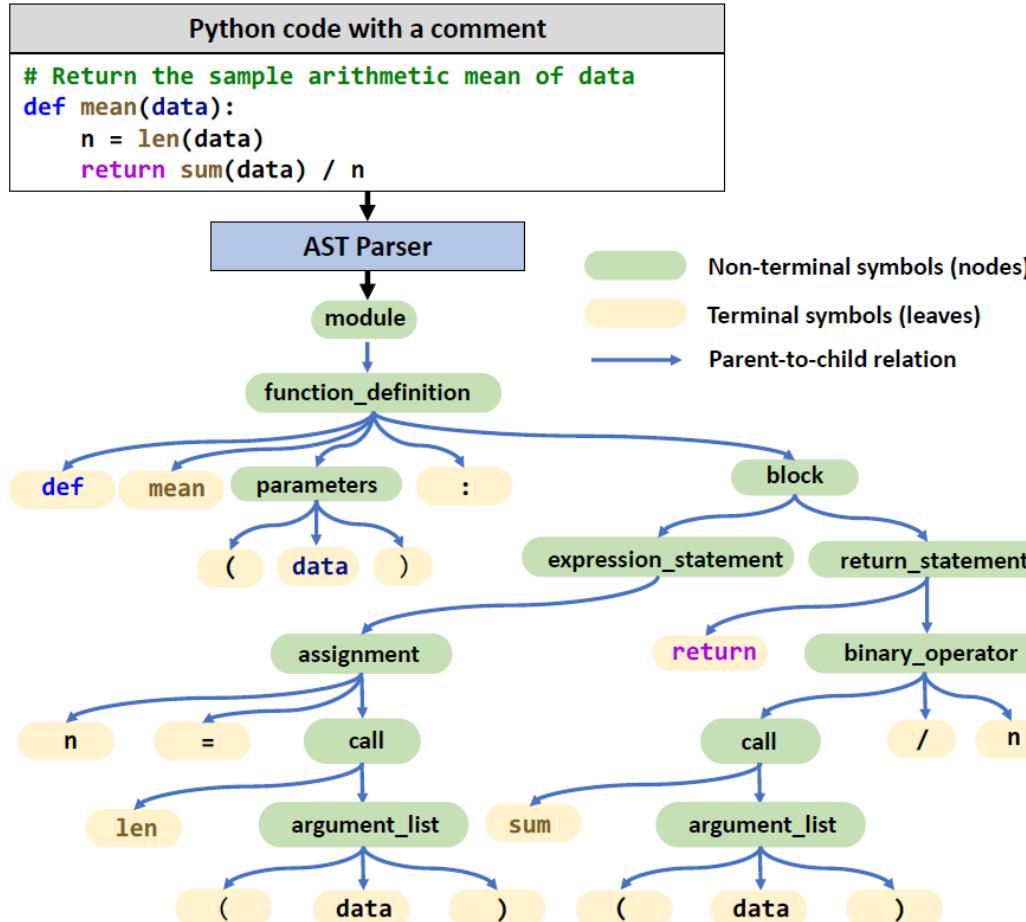
Comment

Return maximum value

Variable relationship



Extract Sequential AST (SeqAST) for Code Pre-training



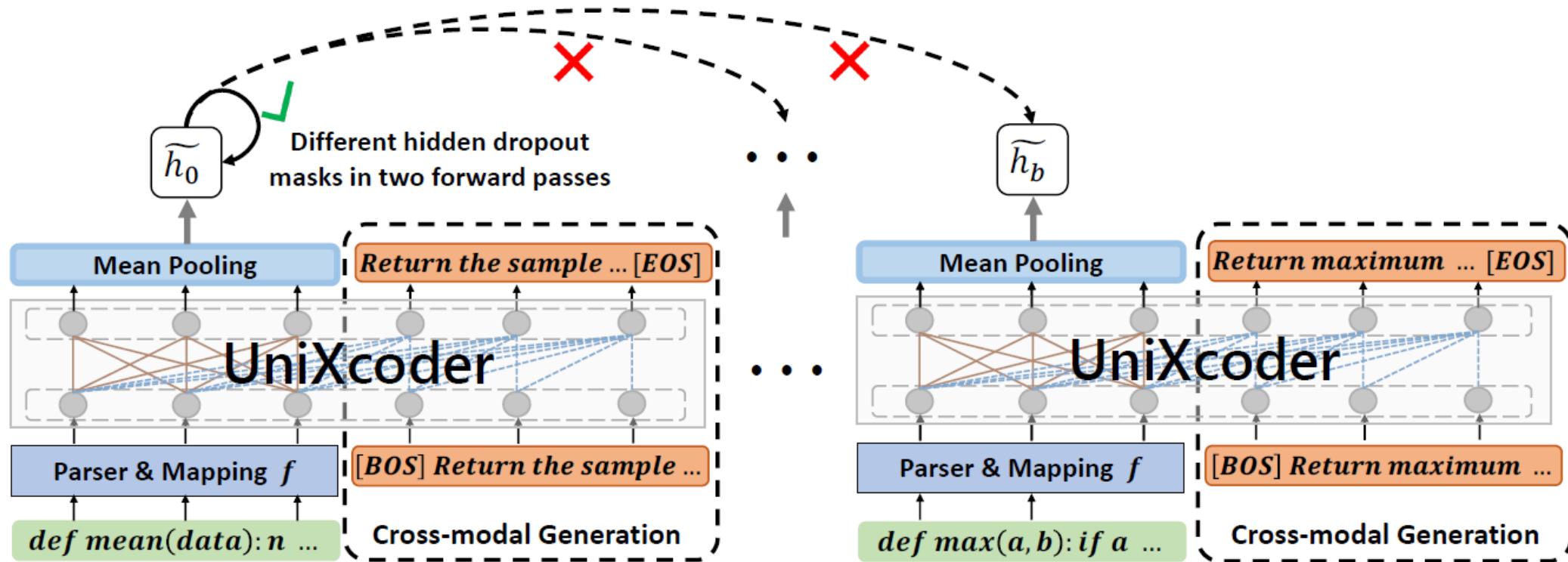
Algorithm 1 AST Mapping Function \mathcal{F}

Input: The root node $root$ of AST

Output: A flattened token sequence

```
1: function  $\mathcal{F}(root)$ 
2:    $seq =$  an empty list
3:    $name =$  the name of  $root$ 
4:   if  $root$  is a leaf then
5:      $seq.append(name)$ 
6:   else
7:      $seq.append(name :: left)$ 
8:     for  $child$  in children of  $root$  do
9:        $seq.extend(\mathcal{F}(child))$ 
10:    end for
11:     $seq.append(name :: right)$ 
12:  end if
13: end function
```

UniXcoder (v3): Pre-Train with Code+Text+SeqAST



- **AST-based code-to-code retrieval** forwards the same sequential AST input using a different hidden dropout mask as a positive example and uses other sequential ASTs in the same batch as negative examples.
- **AST-based code-to-text generation** asks the model to generate the comment based on the sequential AST input.

Evaluation

Model	Clone Detection			Code Search			
	POJ-104		BigCloneBench		CosQA	AdvTest	CSN
	MAP@R	Recall	Precision	F1-score	MRR		
RoBERTa	76.67	95.1	87.8	91.3	60.3	18.3	61.7
CodeBERT	82.67	94.7	93.4	94.1	65.7	27.2	69.3
GraphCodeBERT	85.16	94.8	95.2	95.0	68.4	35.2	71.3
SYNCOBERT	88.24	-	-	-	-	38.3	74.0
PLBART	86.27	94.8	92.5	93.6	65.0	34.7	68.5
CodeT5-base	88.65	94.8	94.7	95.0	67.8	39.3	71.5
UniXcoder	90.52	92.9	97.6	95.2	70.1	41.3	74.4
-w/o contras	87.83	94.9	94.9	94.9	69.2	40.8	73.6
-w/o cross-gen	90.51	94.8	95.6	95.2	69.4	40.1	74.0
-w/o comment	87.05	93.6	96.2	94.9	67.9	40.7	72.6
-w/o AST	88.74	92.9	97.2	95.0	68.7	40.3	74.2
-using BFS	89.44	93.4	96.7	95.0	69.3	40.1	74.1
-using DFS	89.74	94.7	94.6	94.7	69.0	40.2	74.2

Results of code understanding tasks.

Model	Summarization		Generation	
	BLEU-4	EM	BLEU-4	EM
RoBERTa	16.57	-	-	-
CodeBERT	17.83	-	-	-
GPT-2	-	17.35	25.37	-
CodeGPT	-	20.10	32.79	-
PLBART	18.32	18.75	36.69	-
CodeT5-small	19.14	21.55	38.13	-
CodeT5-base	19.55	22.30	40.73	-
UniXcoder	19.30	22.60	38.23	-
-w/o contras	19.20	22.10	37.69	-
-w/o cross-gen	19.27	22.20	35.93	-
-w/o comment	18.97	21.45	37.15	-
-w/o AST	19.33	22.60	38.52	-
-using BFS	19.24	21.75	38.21	-
-using DFS	19.25	22.10	38.06	-

Results of code generation tasks.

Model	PY150		JavaCorpus	
	EM	Edit Sim	EM	Edit Sim
Transformer	38.51	69.01	17.00	50.23
GPT-2	41.73	70.60	27.50	60.36
CodeGPT	42.37	71.59	30.60	63.45
PLBART	38.01	68.46	26.97	61.59
CodeT5-base	36.97	67.12	24.80	58.31
UniXcoder	43.12	72.00	32.90	65.78
-w/o contras	43.02	71.94	32.77	65.71
-w/o cross-gen	42.66	71.83	32.43	65.63
-w/o comment	42.18	71.70	32.20	65.44
-w/o AST	42.56	71.87	32.63	65.66
-using BFS	42.83	71.85	32.40	65.55
-using DFS	42.61	71.97	32.87	65.75

Results on code completion task.

Model	Ruby			Python			Java			Overall
	Ruby	Python	Java	Ruby	Python	Java	Ruby	Python	Java	
CodeBERT	13.55	3.18	0.71	3.12	14.39	0.96	0.55	0.42	7.62	4.94
GraphCodeBERT	17.01	9.29	6.38	5.01	19.34	6.92	1.77	3.50	13.31	9.17
PLBART	18.60	10.76	1.90	8.27	19.55	1.98	1.47	1.27	10.41	8.25
CodeT5-base	18.22	10.02	1.81	8.74	17.83	1.58	1.13	0.81	10.18	7.81
UniXcoder	29.05	26.36	15.16	23.96	30.15	15.07	13.61	14.53	16.12	20.45
-w/o contras	24.03	17.35	7.12	15.80	22.52	7.31	7.55	7.98	13.92	13.73
-w/o cross-gen	28.73	24.16	12.92	21.52	26.66	12.60	11.14	10.82	13.75	18.03
-w/o comment	22.24	15.90	7.50	15.09	19.88	6.54	7.84	7.12	13.20	12.81
-w/o AST	27.54	23.37	10.17	21.75	27.75	9.94	9.79	9.21	14.06	17.06
-using BFS	26.67	23.69	13.56	21.31	27.28	13.63	11.90	12.55	14.92	18.39
-using DFS	27.13	22.65	11.62	20.21	25.92	11.85	9.59	10.19	13.30	16.94

Results on zero-shot code-to-code search task.

(Some of) Our Work @ MSRA

- **CodeBERT** (EMNLP 2020): propose the 1st code-text pre-trained model
- **GraphCodeBERT** (ICLR 2021): use variable relationship to enhance code representation
- **UniXcoder** (ACL 2022): use sequential AST to enhance code representation
- **AR2** (ICLR 2022): propose a general dense retrieval framework that trains retriever & ranker in a minimax adversarial manner and can cover text-to-text, text-to-code, and code-to-code tasks
- **CodeRetriever** (EMNLP 2022): extract code-text and code-code pairs from Web for building code retrieval models with unimodal and bimodal contrastive learning
- **GPT-C w/ Extended Context** (EMNLP 2021): improve GPT-C for code completion with extended context
- **Grammformer** (ICLR 2022): generate codes based on code grammar and predict holes to alleviate the uncertainty issue
- **ReACC** (ACL 2022): propose a retrieval-augmented code completion framework, which uses the partial code as a query to retrieve similar codes and predicts the following codes based on the retrieved similar codes
- **CodeReviewer** (ESEC/FSE 2022): propose an encoder-decoder based pre-trained model for automating code review activities, such as code diff quality estimation, review generation, code refinement based on code review, etc.
- **CodeReviewer v2.0** (on-going 2022): propose a retrieval-augmented model for automating code review activities, like ReACC for code completion.
- **CodeXGLUE** (NeurIPS 2021): build a benchmark for code intelligence, which covers 14 datasets for 10 code tasks
- **CoSQA** (ACL 2021): build a benchmark for text-to-code search based on Web data
- **CodeExplanation** (EMNLP 2022): build a benchmark for code-to-explanation generation

 Code Representation

 Code Retrieval

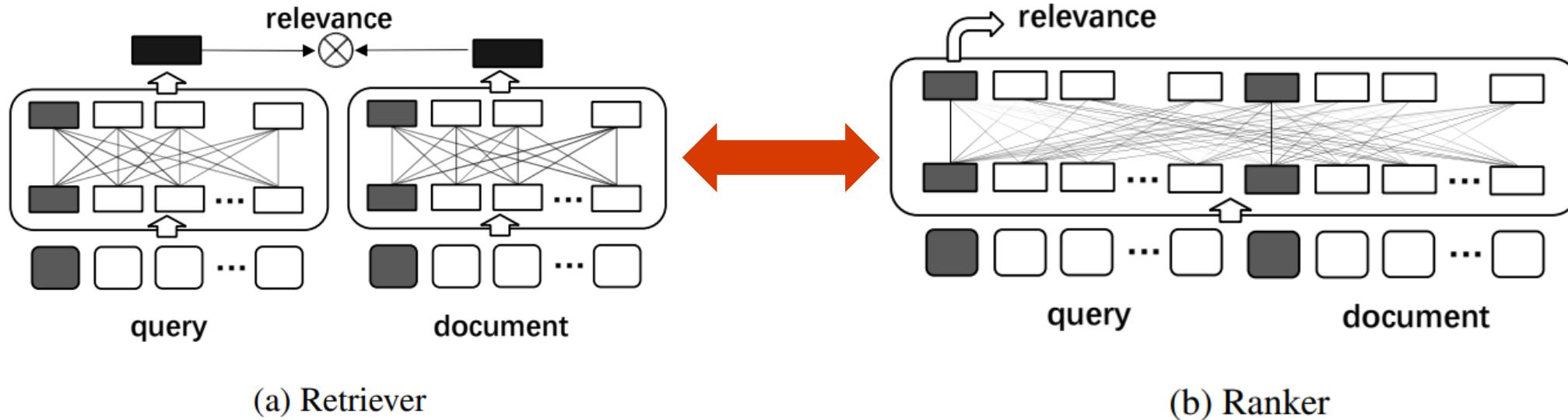
 Code Generation

 Code Review & Refinement

 Code Benchmark



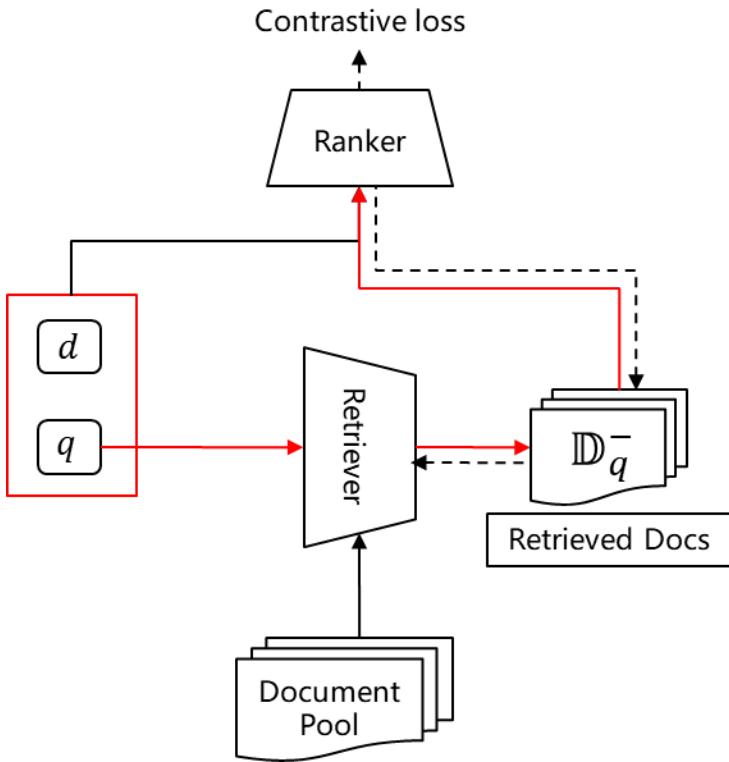
Jointly Optimize Single-tower and Two-tower Models



Jointly optimize two modules according to a minimax adversarial objective

- Retriever: retrieve negative documents to cheat Ranker
- Ranker: distinguish the ground-truth document and the retrieved ones by Retriever

Retriever → Ranker

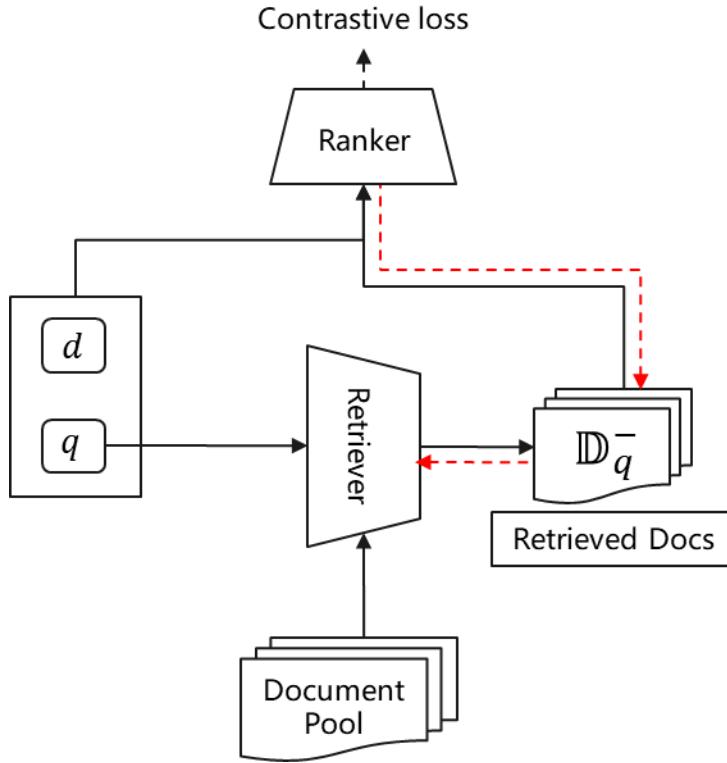


Ranker training:

$$\phi^* = \operatorname{argmax}_\phi \log p_\phi(d|q, d, \mathbb{D}_q^-)$$

$$p_\phi(d|q, d, \mathbb{D}_q^-) = \frac{e^{\tau D_\phi(q, d)}}{e^{\tau D_\phi(q, d)} + \sum_{i=1}^n e^{\tau D_\phi(q, d_i^-)}}$$

Ranker → Retriever

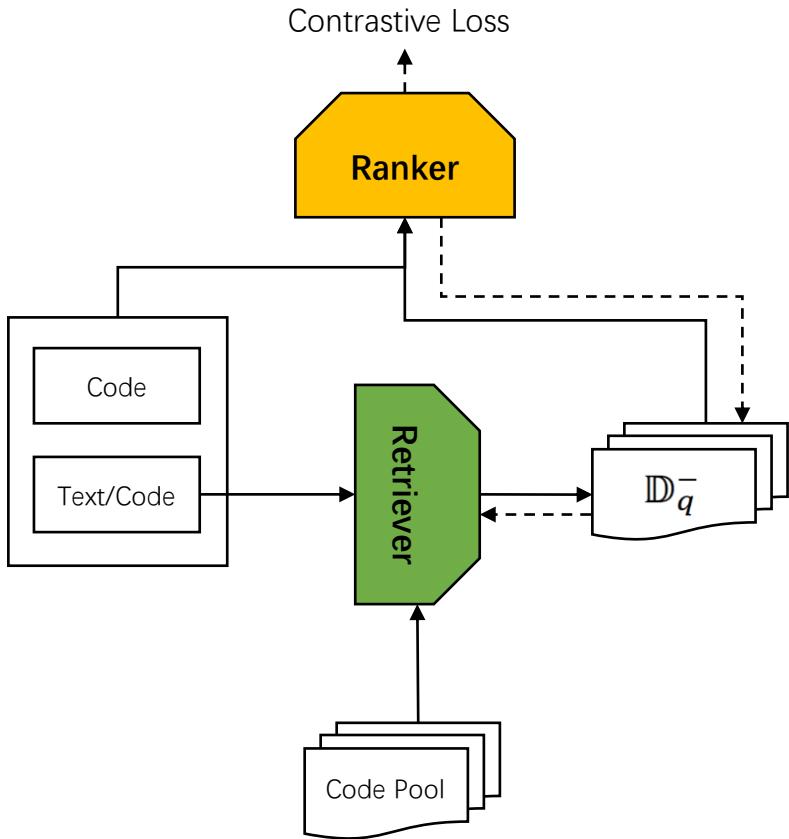


Retriever training:

$$\theta^* = \operatorname{argmin}_{\theta} J^{\theta} = \mathbf{E}_{\mathbb{D}_q^- \sim G_{\theta}(q, \cdot)} [\log p_{\phi}(d|q, d, \mathbb{D}_q^-)]$$

$$p_{\phi}(d|q, d, \mathbb{D}_q^-) = \frac{e^{\tau D_{\phi}(q, d)}}{e^{\tau D_{\phi}(q, d)} + \sum_{i=1}^n e^{\tau D_{\phi}(q, d_i^-)}}$$

Dense Retrieval w/ Adversarial Retriever-Ranker (AR2)



$$J^{G^*, D^*} = \min_{\theta} \max_{\phi} \mathbf{E}_{\mathbb{D}_q^- \sim G_\theta(q, \cdot)} [\log p_\phi(d|q, d, \mathbb{D}_q^-)]$$

$$p_\phi(d|q, d, \mathbb{D}_q^-) = \frac{e^{\tau D_\phi(q, d)}}{e^{\tau D_\phi(q, d)} + \sum_{i=1}^n e^{\tau D_\phi(q, d_i^-)}}$$

Retriever θ : try to find the hard negatives \mathbb{D}_q^- to cheat Ranker ϕ .

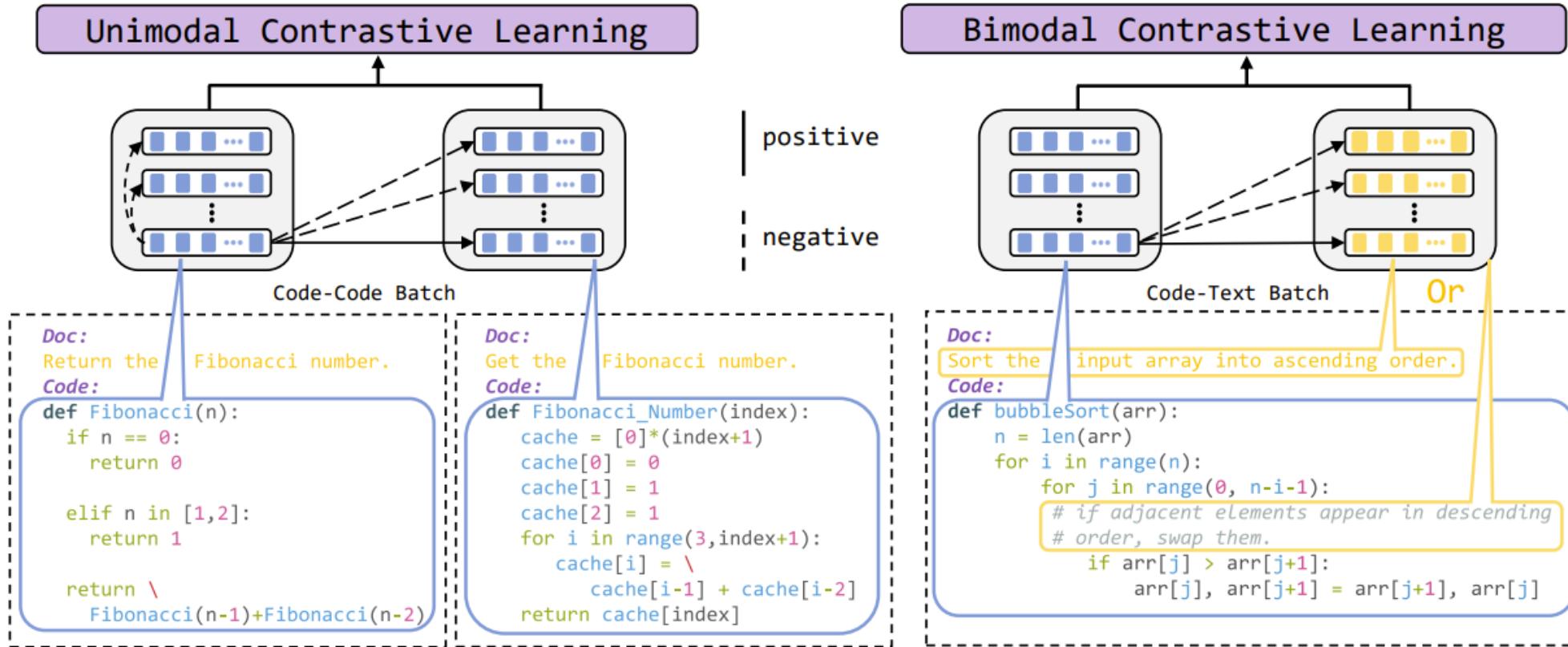
$$\theta^* = \operatorname{argmin}_{\theta} J^\theta = \mathbf{E}_{\mathbb{D}_q^- \sim G_\theta(q, \cdot)} [\log p_\phi(d|q, d, \mathbb{D}_q^-)]$$

Ranker ϕ : try to find the golden d from the negatives selected by Retriever θ .

$$\phi^* = \operatorname{argmax}_{\phi} \log p_\phi(d|q, d, \mathbb{D}_q^-)$$

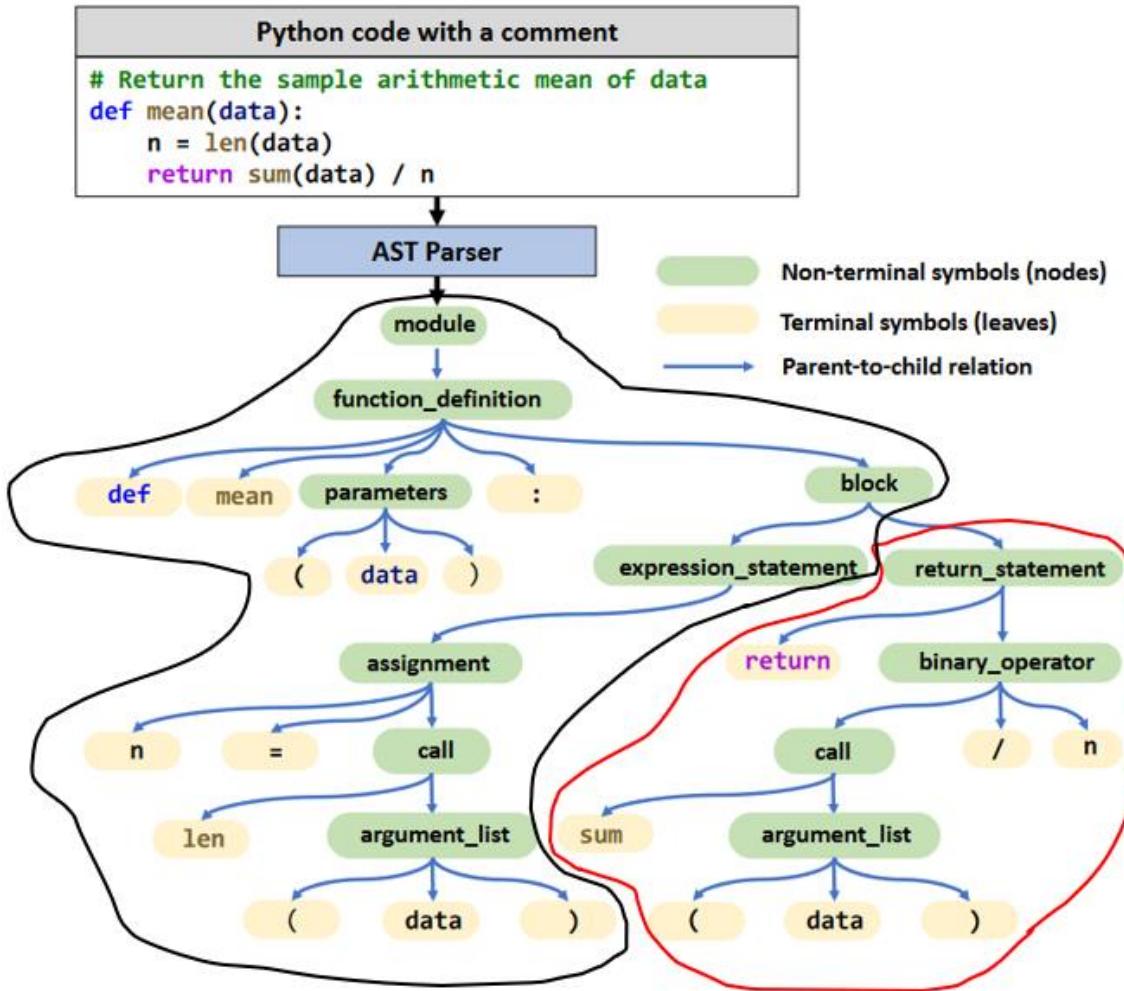
- **Text-Code pairs** come from CodeSearchNet
- **Code-Code pairs** come from AST-based ICT

CodeRetriever with Code-Code and Code-Text Pairs



1. CodeRetriever proposes a semantic-guided method to build positive code-code pairs based on the documentation and function names.
2. CodeRetriever uses unimodal (i.e., code-code) and bimodal (i.e., text-code) contrastive learning to learn function-level code representations and achieves new SOTA results on the text-to-code search task, comparing to several strong baselines.

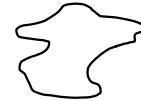
CodeRetriever with AST-Based Inverse Cloze Test



Query Code:



Answer Code:



Motivation: compared with applying ICT of NLP (random sampling token span of code tokens as query), AST-based ICT can generate query-answer code pairs without syntax errors.

Evaluation on Code Search & Clone Detection

Code Search	CodeSearchNet (Husain et al., 2019)	CoSQA (Huang et al., 2021)	AdvTest (Lu et al., 2021)
CodeBERT	69.28%	64.70%	27.20%
GraphCodeBERT	73.30%	67.50%	35.20%
UniXcoder	74.35%	70.10%	41.30%
CodeRetriever	76.70%	74.50%	45.50%

Clone Detection	CodeNet (zero-shot) (Puri et al., 2021)				POJ-104 (Mou et al., 2016)
	Ruby	Python	Java	Overall	MRR
CodeBERT	13.55%	14.39%	7.62%	11.85%	82.67%
GraphCodeBERT	17.01%	19.34%	13.31%	16.55%	85.16%
UniXcoder	29.05%	30.15%	16.12%	25.11%	90.52%
CodeRetriever	33.72%	32.78%	18.91%	28.47%	92.45%

(Some of) Our Work @ MSRA

- **CodeBERT** (EMNLP 2020): propose the 1st code-text pre-trained model
- **GraphCodeBERT** (ICLR 2021): use variable relationship to enhance code representation
- **UniXcoder** (ACL 2022): use sequential AST to enhance code representation
- **AR2** (ICLR 2022): propose a general dense retrieval framework that trains retriever & ranker in a minimax adversarial manner and can cover text-to-text, text-to-code, and code-to-code tasks
- **CodeRetriever** (EMNLP 2022): extract code-text and code-code pairs from Web for building code retrieval models with unimodal and bimodal contrastive learning
- **GPT-C w/ Extended Context** (EMNLP 2021): improve GPT-C for code completion with extended context
- **Grammformer** (ICLR 2022): generate codes based on code grammar and predict holes to alleviate the uncertainty issue
- **ReACC** (ACL 2022): propose a retrieval-augmented code completion framework, which uses the partial code as a query to retrieve similar codes and predicts the following codes based on the retrieved similar codes
- **CodeReviewer** (ESEC/FSE 2022): propose an encoder-decoder based pre-trained model for automating code review activities, such as code diff quality estimation, review generation, code refinement based on code review, etc.
- **CodeReviewer v2.0** (on-going 2022): propose a retrieval-augmented model for automating code review activities, like ReACC for code completion.
- **CodeXGLUE** (NeurIPS 2021): build a benchmark for code intelligence, which covers 14 datasets for 10 code tasks
- **CoSQA** (ACL 2021): build a benchmark for text-to-code search based on Web data
- **CodeExplanation** (EMNLP 2022): build a benchmark for code-to-explanation generation

 Code Representation

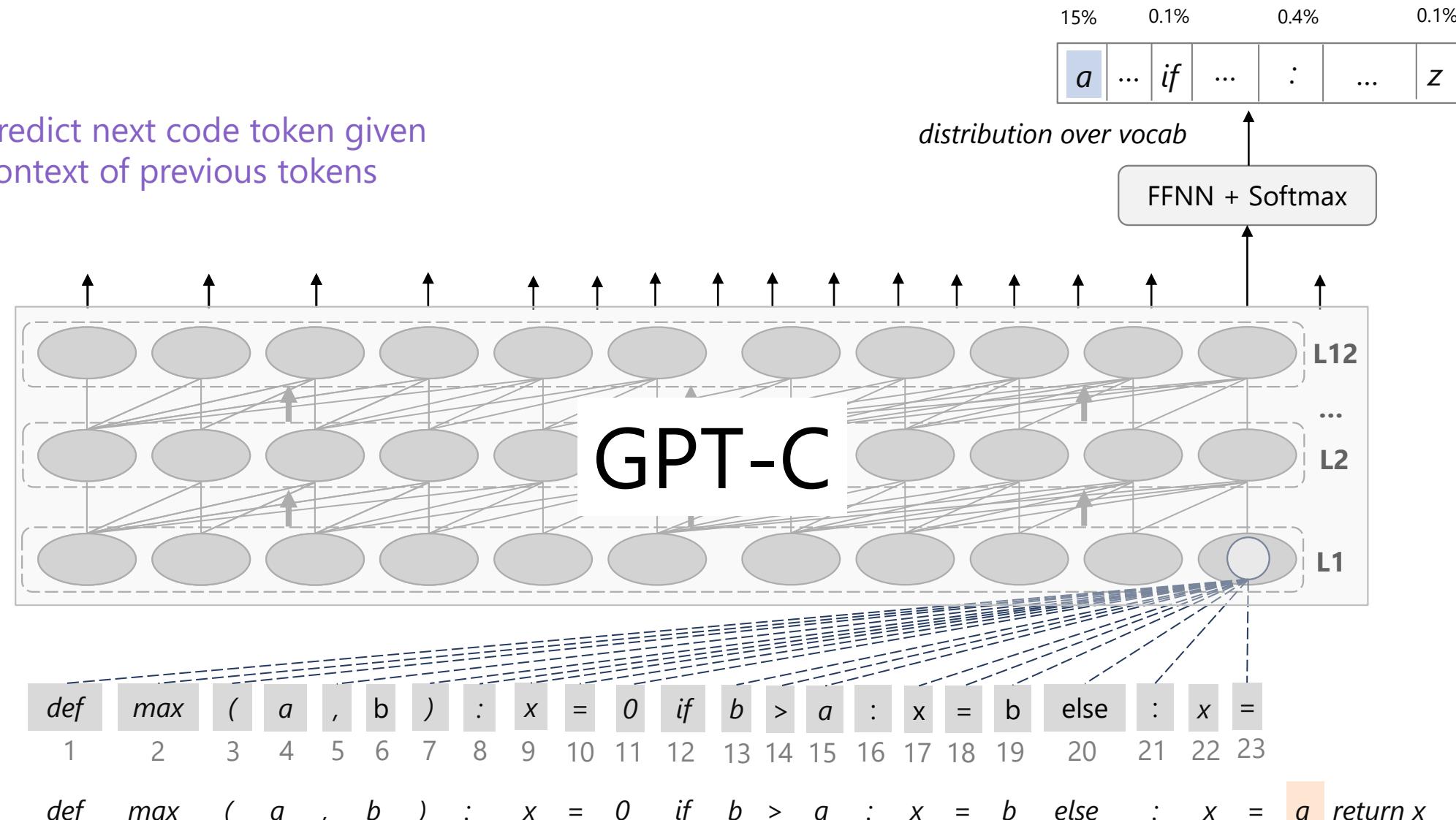
 Code Retrieval

 Code Generation

 Code Review & Refinement

 Code Benchmark

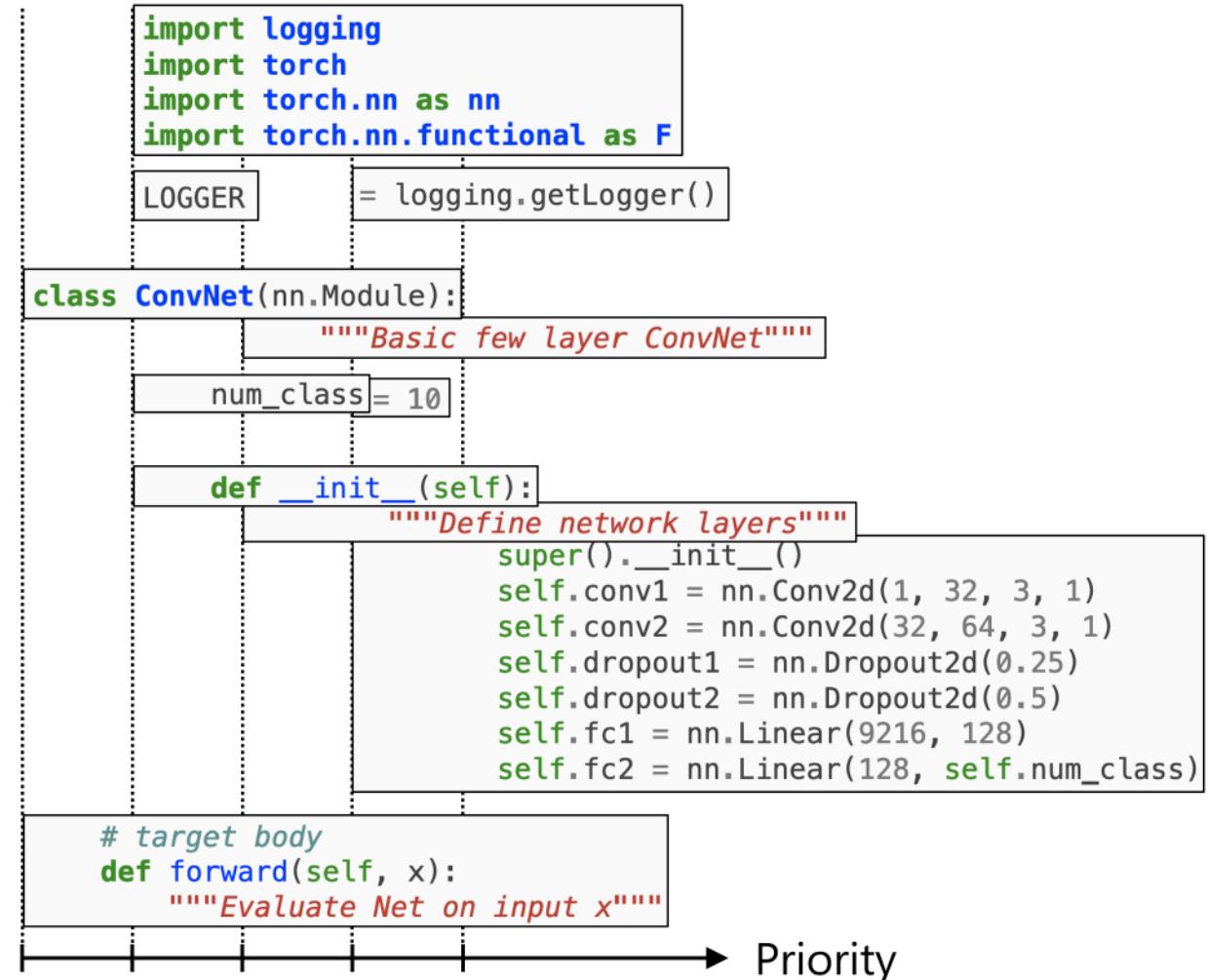
GPT-C (v1): Multilingual Code Completion Model



Trained for 10 PLs: JavaScript, C, Java, Go, PHP, Python, C++, C#, Ruby, TypeScript

Extended Context for Code Completion

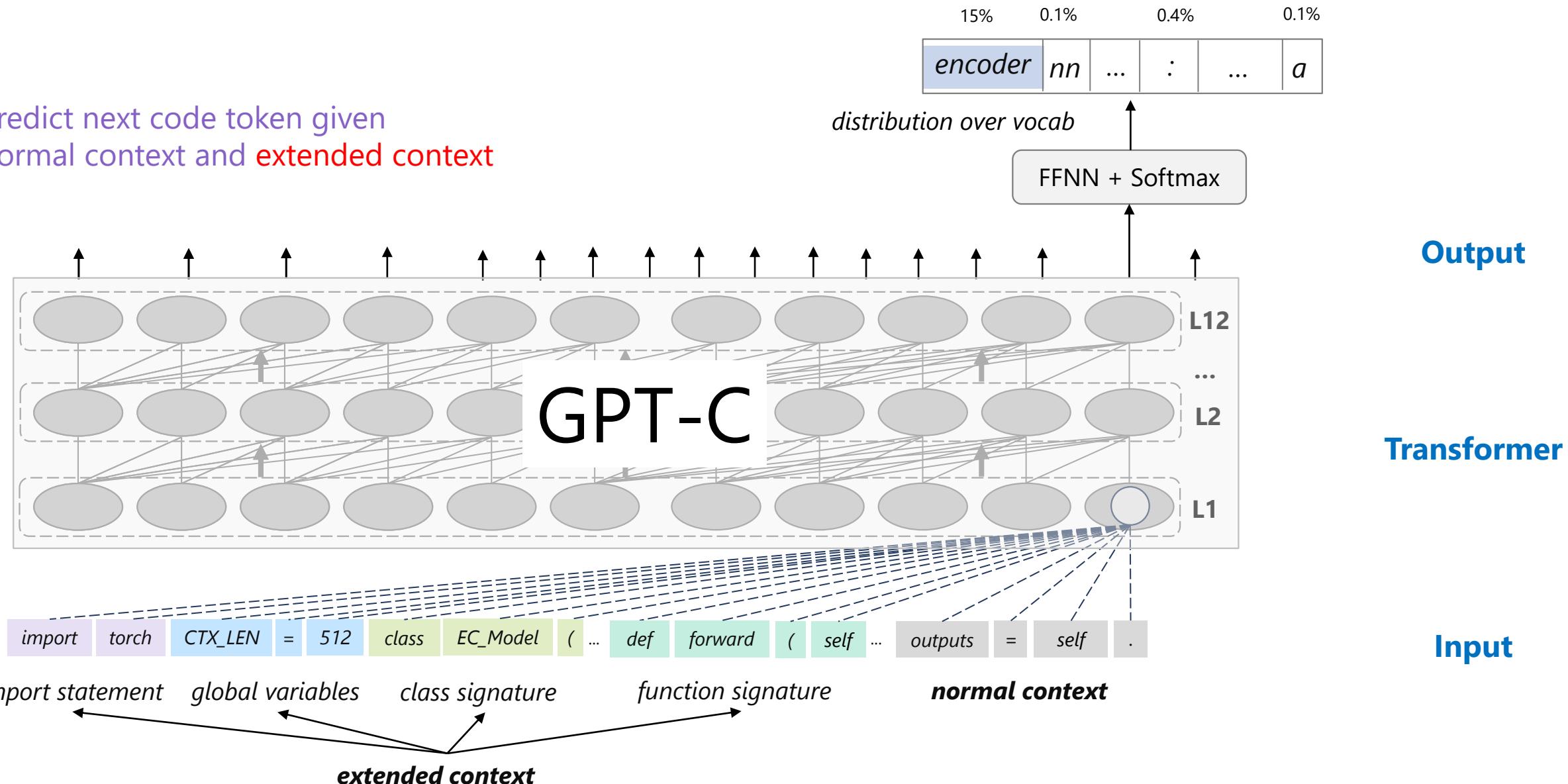
1. Take the concrete syntax tree of the source file;
2. Prioritize the syntactic elements;
 - i. signature and docstring of the focal method;
 - ii. global import statements and assigned values;
 - iii. class attributes, peer class method signatures, class docstring, peer class method docstrings;
 - iv. global expressions and code bodies of peer class methods.
3. Take elements based on their priorities until the context window has been filled.



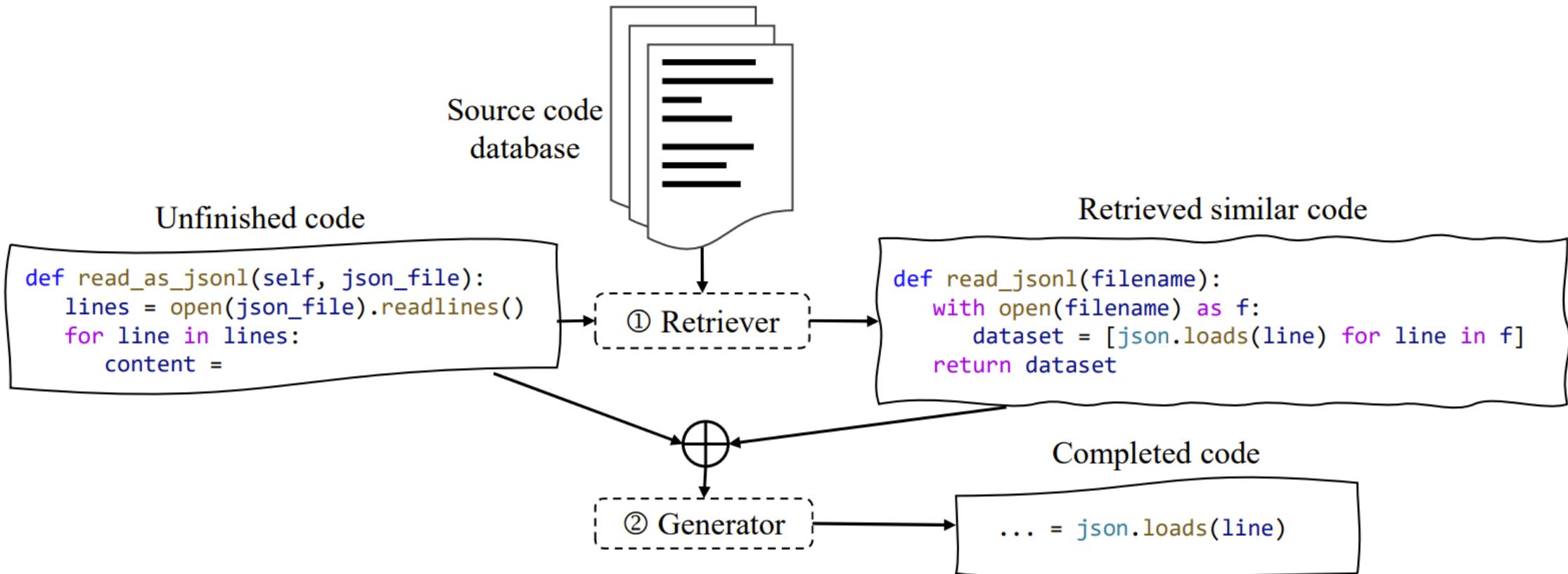
We reserve 3/4 (768/1,024) tokens for the extended context and 1/4 (256/1024) tokens for the local context.

GPT-C (v2) with Extended Context

Predict next code token given
normal context and **extended context**

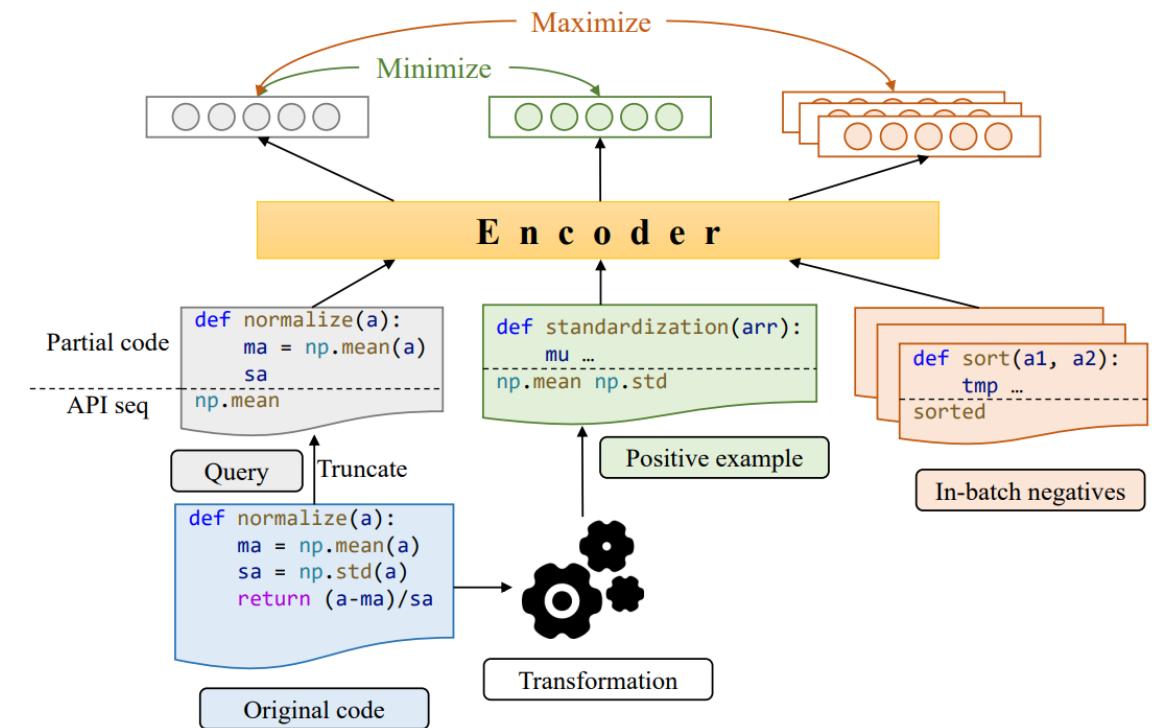


GPT-C (v3) with Retrieved Similar Code



Code-to-Code Retrieval Pre-training

- **Goal**
 - Partial code → Similar complete code
- **Contrastive pre-training**
 - **query**: a random truncation of the original code + API sequence
 - **+ instance**: the entire transformed code + API sequence
 - **- instance**: in-batch negatives + API sequences



Semantic-Preserving Code Transformation

```
import socket
def echo_server(client, timeout, bufsize):
    try:
        if timeout > 0:
            client.settimeout(timeout)
        get_buf = client.recv(bufsize)
        client.send(get_buf)
    except socket.timeout:
        pass
    client.close()
```

original python code

```
import socket
def get_mean(c, doc, local):
    try:
        if doc > 0:
            c.settimeout(doc)
            _user_id = c.recv(local)
            c.send(_user_id)
    except socket.timeout:
        pass
    c.close()
```

After renaming all variables

```
import socket
def echo_server(client, timeout, bufsize):
    try:
        if timeout > 0:
            client.settimeout(timeout)
        get_buf = client.recv(bufsize)
        if True:
            tmp = [x**2 for x in range(10)]
            client.send(get_buf)
    except socket.timeout:
        pass
    client.close()
```

After inserting dead code

- **Identifier renaming** is a method of renaming an identifier with another.
- **Dead code insertion** is to insert a dead code into a code fragment at a proper location.

Retrieval-augmented Code Completion

Model	PY150			JavaCorpus		
	Perplexity	Exact Match	Edit Sim	Perplexity	Exact Match	Edit Sim
GPT-2	-	41.73	70.60	-	27.50	60.36
CodeGPT	2.502	42.18	71.23	4.135	28.23	61.81
CodeGPT-adapted	2.404	42.37	71.59	3.369	30.60	63.45
CodeT5-base	-	36.97	67.12	-	24.80	58.31
PLBART	-	38.01	68.46	-	26.97	61.59
ReACC-bm25	2.312	46.07	73.84	3.352	30.63	64.28
ReACC-dense	2.329	45.32	73.95	3.355	30.30	64.43
ReACC-hybrid	2.311	46.26	74.41	3.327	30.70	64.73

(Some of) Our Work @ MSRA

- **CodeBERT** (EMNLP 2020): propose the 1st code-text pre-trained model
- **GraphCodeBERT** (ICLR 2021): use variable relationship to enhance code representation
- **UniXcoder** (ACL 2022): use sequential AST to enhance code representation
- **AR2** (ICLR 2022): propose a general dense retrieval framework that trains retriever & ranker in a minimax adversarial manner and can cover text-to-text, text-to-code, and code-to-code tasks
- **CodeRetriever** (EMNLP 2022): extract code-text and code-code pairs from Web for building code retrieval models with unimodal and bimodal contrastive learning
- **GPT-C w/ Extended Context** (EMNLP 2021): improve GPT-C for code completion with extended context
- **Grammformer** (ICLR 2022): generate codes based on code grammar and predict holes to alleviate the uncertainty issue
- **ReACC** (ACL 2022): propose a retrieval-augmented code completion framework, which uses the partial code as a query to retrieve similar codes and predicts the following codes based on the retrieved similar codes
- **CodeReviewer** (ESEC/FSE 2022): propose an encoder-decoder based pre-trained model for automating code review activities, such as code diff quality estimation, review generation, code refinement based on code review, etc.
- **CodeReviewer v2.0** (on-going 2022): propose a retrieval-augmented model for automating code review activities, like ReACC for code completion.
- **CodeXGLUE** (NeurIPS 2021): build a benchmark for code intelligence, which covers 14 datasets for 10 code tasks
- **CoSQA** (ACL 2021): build a benchmark for text-to-code search based on Web data
- **CodeExplanation** (EMNLP 2022): build a benchmark for code-to-explanation generation

 Code Representation

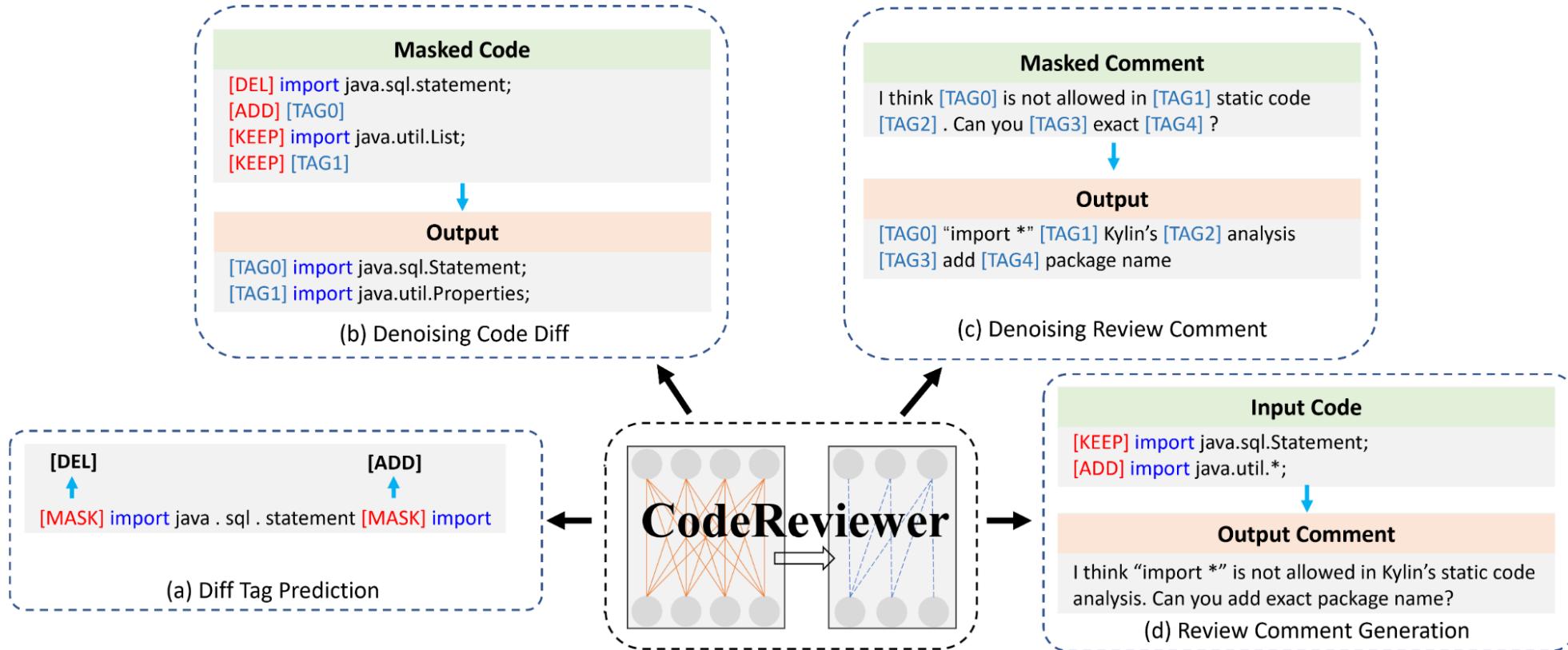
 Code Retrieval

 Code Generation

 Code Review & Refinement

 Code Benchmark

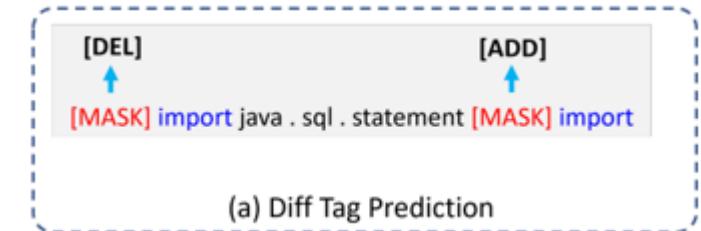
CodeReviewer for Automating Code Review Activities



Pre-training Tasks

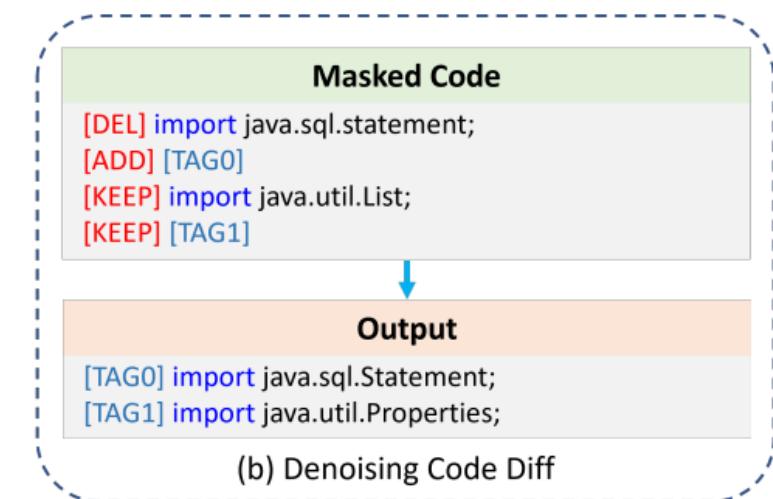
- Diff Tag Prediction

$$\mathcal{L}_{DTP} = - \sum_i \left(y_0^{(i)} \log p_0^{(i)} + y_1^{(i)} \log p_1^{(i)} + y_2^{(i)} \log p_2^{(i)} \right)$$



- Denoising Code Diff

$$\mathcal{L}_{DCD} = \sum_{t=1}^k -\log P_\theta(c_t | \mathbf{c}^{\text{mask}}, \mathbf{c}_{<t})$$



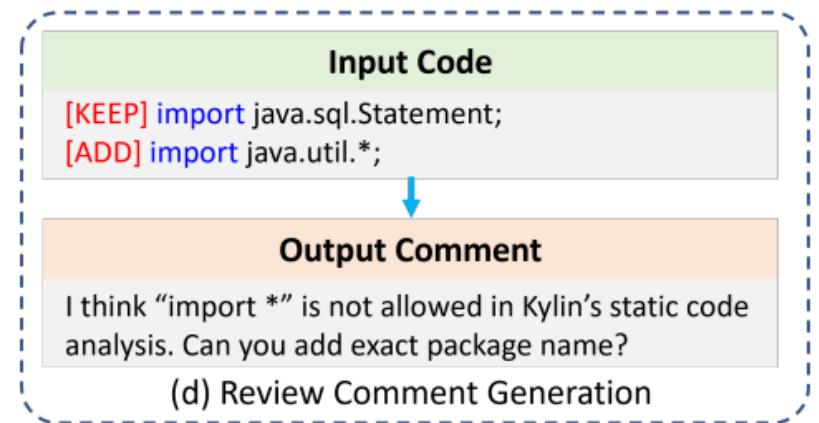
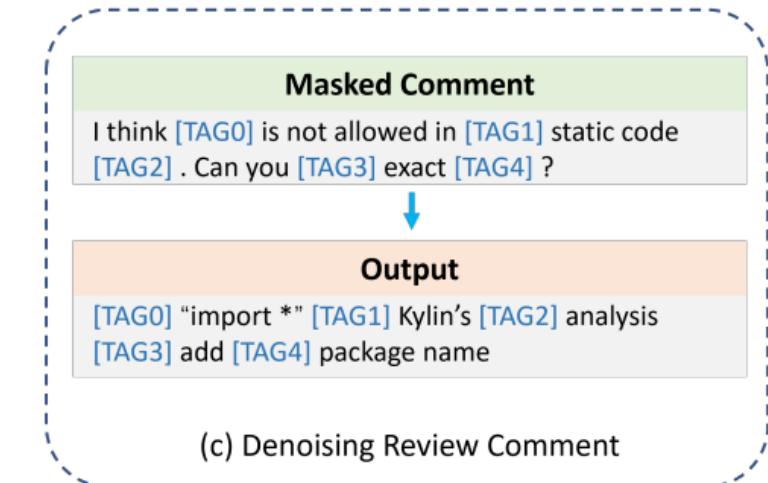
Pre-training Tasks

- Denoising Review Comment

$$\mathcal{L}_{DRC} = \sum_{t=1}^k -\log P_{\theta}(w_t | w^{\text{mask}}, w_{<t})$$

- Review Comment Generation

$$\mathcal{L}_{RCG} = \sum_{t=1}^k -\log P(w_t | c, w_{<t})$$



Code Review Dataset

1. CodeReviewer is pre-trained on the pull requests crawled from GitHub in 9 programming languages and establishes a benchmark dataset for code review activities.
2. Collected from projects who containing more than 1000 PRs.

Language	Ruby	Go	Php	Js	Python	Java	C++	C#	C	Total
# Project	72	146	92	194	195	175	133	77	77	1.2M
# (Diff, NL)	126k	410k	157k	293k	518k	467k	202k	199k	110k	2.5M
# Diff	90k	728k	215k	425k	887k	876k	474k	324k	292k	4.3M

Examples

```
public class MockExecutorLoader implements ExecutorLoader {  
    public void unassignExecutor(int executionId) throws ExecutorManagerException {  
        executionExecutorMapping.remove(executionId);  
    }  
+  
+ @Override  
+ public List<ExecutableFlow> fetchRecentlyFinishedFlows(long lifeTimeMs)  
+     throws ExecutorManagerException {  
+     return null;  
+ }  
}
```

Ground truth: I think we should return an empty list instead of null here.

T5: I don't think we need this.

CodeT5: I don't think this class belongs here.

Codex: I think this is a good idea.

CodeReviewer: return empty list instead?

```
public abstract class Directory implements Closeable {  
    * @throws AlreadyClosedException if this Directory is closed  
    */  
    protected void ensureOpen() throws AlreadyClosedException {}  
+  
+ /**  
+ * Implementations can override this if they are capable of reporting modification time  
+ * of a file in seconds since the epoch.  
+ */  
+ public long fileModified(String name) throws IOException {  
+     throw new UnsupportedOperationException();  
+ }  
+  
}
```

Ground truth: I think we should avoid changing any lucene classes for the moment - fileModified() can probably stay where it is?

T5: Why do we need this method? I don't see it being used anywhere else.

CodeT5: Shouldn't this be synchronized as well?

Codex: should we make this abstract?

CodeReviewer: I think this should throw `FileNotFoundException` instead of `IOException`

Original Code:

```
private string ReplaceUnit(string resultText)  
{  
    var unitMap = Config.UnitMap.OrderBy(o => o.Key.Length).ToDictionary(o => o.Key, p =>  
    p.Value);  
    foreach (var unit in unitMap.Keys)  
    {  
        resultText = resultText.Replace(unit, Config.UnitMap[unit]);  
    }  
    ...  
}
```

Review comment: This should only happen once in initialization, can you move it out of this method?

CodeReviewer: (also the ground truth)

```
private string ReplaceUnit(string resultText)  
{  
    foreach (var unit in Config.UnitMap.Keys)  
    {  
        resultText = resultText.Replace(unit, Config.UnitMap[unit]);  
    }  
    ...  
}
```

CodeT5:

```
private string ReplaceUnit(string resultText)  
{  
    foreach (var unit in Config.UnitMap.OrderBy(o => o.Key.Length).ToDictionary(o => o.Key,  
    p => p.Value))  
    {  
        resultText = resultText.Replace(unit, Config.UnitMap[unit]);  
    }  
    ...  
}
```

T5:

```
private string ReplaceUnit(string resultText)  
{  
    var unitMap = Config.UnitMap.OrderBy(o => o.Key.Length).ToDictionary(o => o.Key, p =>  
    p.Value);  
    foreach (var unit in unitMap.Keys)  
    {  
        resultText = resultText.Replace(unit, Config.UnitMap[unit]);  
    }  
    ...  
}
```

(Some of) Our Work @ MSRA

- **CodeBERT** (EMNLP 2020): propose the 1st code-text pre-trained model
- **GraphCodeBERT** (ICLR 2021): use variable relationship to enhance code representation
- **UniXcoder** (ACL 2022): use sequential AST to enhance code representation
- **AR2** (ICLR 2022): propose a general dense retrieval framework that trains retriever & ranker in a minimax adversarial manner and can cover text-to-text, text-to-code, and code-to-code tasks
- **CodeRetriever** (EMNLP 2022): extract code-text and code-code pairs from Web for building code retrieval models with unimodal and bimodal contrastive learning
- **GPT-C w/ Extended Context** (EMNLP 2021): improve GPT-C for code completion with extended context
- **Grammformer** (ICLR 2022): generate codes based on code grammar and predict holes to alleviate the uncertainty issue
- **ReACC** (ACL 2022): propose a retrieval-augmented code completion framework, which uses the partial code as a query to retrieve similar codes and predicts the following codes based on the retrieved similar codes
- **CodeReviewer** (ESEC/FSE 2022): propose an encoder-decoder based pre-trained model for automating code review activities, such as code diff quality estimation, review generation, code refinement based on code review, etc.
- **CodeReviewer v2.0** (on-going 2022): propose a retrieval-augmented model for automating code review activities, like ReACC for code completion.
- **CodeXGLUE** (NeurIPS 2021): build a benchmark for code intelligence, which covers 14 datasets for 10 code tasks
- **CoSQA** (ACL 2021): build a benchmark for text-to-code search based on Web data
- **CodeExplanation** (EMNLP 2022): build a benchmark for code-to-explanation generation

 Code Representation

 Code Retrieval

 Code Generation

 Code Review & Refinement

 Code Benchmark



CodeXGLUE: 14 datasets for 10 Code-related tasks

Category	Task	Dataset Name	Language	Train/Dev/Test Size	Baselines	Dataset Provider	Task definition
Code-Code	Clone Detection	BigCloneBench	Java	900K/416K/416K	CodeBERT	Univ. of Saskatchewan	Predict semantic equivalence for a pair of codes.
		POJ-104	C/C++	32K/8K/12K		Peking Univ	Retrieve semantically similar codes.
	Defect Detection	Defects4J	C	21k/2.7k/2.7k		Univ. of Washington	Identify whether a function is vulnerable.
	Cloze Testing	CT-all	Python, Java, PHP, JavaScript, Ruby, Go	-/-/176k		Created by MSRA based on CodeSearchNet	Tokens to be predicted come from the entire vocab.
		CT-max/min	Python, Java, PHP, JavaScript, Ruby, Go	-/-/2.6k		Created by MSRA based on CodeSearchNet	Tokens to be predicted come from {max, min}.
	Code Completion	PY150	Python	100k/5k/50k	CodeGPT	ETH Zurich, line-level data added by MSRA	Predict following tokens given contexts of codes.
		GitHub Java Corpus	Java	13k/7k/8k		Univ. of Edinburgh, line-level data added by MSRA	
	Code Refinement	Bugs2Fix	Java	98K/12K/12K	Encoder-Decoder	The College of William and Mary	Automatically refine codes by fixing bugs.
	Code Translation	CodeTrans	Java-C#	10K/0.5K/1K		MSRA	Translate the codes from one programming language to another programming language.
Text-Code	NL Code Search	CodeSearchnet, AdvTest	Python	251K/9.6K/19K	CodeBERT	GitHub + MSR Cambridge, test provided by MSRA	Given a natural language query as input, find semantically similar codes.
		StacQC, WebQueryTest	Python	2.9k/0.9k/1.9k		The Ohio State Univ, test provided by MSRA	Given a pair of natural language and code, predict whether they are relevant or not.
	Text-to-Code Generation	CONCODE	Java	100K/2K/2K	CodeGPT	Univ. of Washington	Given a natural language docstring/comment as input, generate a code.
Code-Text	Code Summarization	CodeSearchNet*	Python, Java, PHP, JavaScript, Ruby, Go	908K/45K/53K	Encoder-Decoder	Filtered based on CodeSearchNet data	Given a code, generate its natural language docstring/comment.
Text-Text	Documentation Translation	Microsoft Docs	English-Latvian/Danish/Norwegian/Chinese	156K/4K/4K		MSRA	Translate code documentation between human languages (e.g. En-Zh), intended to test low-resource multi-lingual translation.

CodeExp for Explanatory Code Document Generation

Code:

```
def make_rng(rng_or_seed=None, default_seed=None, constructor=None):
    if (rng_or_seed is not None) and isinstance(rng_or_seed, RNG):
        rng = rng_or_seed
    elif (rng_or_seed is not None):
        rng = constructor(rng_or_seed)
    elif (default_seed is not None):
        rng = constructor(default_seed)
    else:
        rng = constructor(42)
    return rng
```

Summary:

Returns a random number generator.

Explanatory Docstring:

Returns a random number generator.

The RNG object is generated using the first of these cases that produces a valid result:

- 1) rng_or_seed itself
- 2) constructor(rng_or_seed)
- 3) constructor(default_seed)
- 4) constructor(42)

Parameters:

rng_or_seed (int or RNG): If `rng_or_seed` is a random number generator, then it is returned. If `rng_or_seed` is an integer, then a random number generator is created using `constructor` and seeded with `rng_or_seed`.

default_seed (int): Seed used if rng_or_seed is None.

constructor (function or class): Must return a RNG object. constructor is called with rng_or_seed, default_seed or 42.

Returns:

An RNG object.

Partition	#Examples	Quality	Annotated By
CodeExp(raw)	2,285,387	Mixed	-
CodeExp(refined)	158,024	High	Machine
CodeExp(annotated)	13,186	Mixed	Human

Model	ROUGE-1 f	ROUGE-L f	BLEU	CER	METEOR	BERTScore	CodeBERT Score
GPT-2-base							
- CodeExp(raw)	0.2557	0.2443	4.42	0.4580	18.55	83.87	77.46
- CodeExp(refined)	0.2462	0.2357	4.26	0.4890	19.00	83.55	77.48
- CodeExp(r+r)	0.2623	0.2520	5.19	0.4978	20.30	83.92	77.91
GPT-Neo13							
- w/o fine-tune	0.0694	0.0646	0.40	0.1294	5.71	78.19	67.62
- CodeExp(raw)	0.2894	0.2769	7.51	0.4805	21.85	84.46	78.31
- CodeExp(refined)	0.2954	0.2815	7.36	0.5570	23.58	84.24	78.46
- CodeExp(r+r)	0.3265	0.3128	10.45	0.5524	26.31	84.78	79.26
GPT-Neo27							
- w/o fine-tune	0.1480	0.1380	0.82	0.2285	10.06	78.40	71.99
- CodeExp(raw)	0.2953	0.2818	7.72	0.4895	22.42	84.31	78.34
- CodeExp(refined)	0.2955	0.2816	8.21	0.5285	23.66	84.33	78.41
- CodeExp(r+r)	0.3298	0.3154	10.72	0.5560	26.87	84.86	79.28
CodeT5							
- multi-sum	0.1507	0.1392	0.21	0.1240	5.58	82.35	74.29
- CodeExp(raw)	0.2652	0.2530	5.39	0.3851	17.51	83.84	77.42
- CodeExp(refined)	0.3175	0.3016	8.02	0.5536	24.38	84.67	78.90
- CodeExp(r+r)	0.3415	0.3256	9.91	0.5695	26.87	84.98	79.52

Summary & Next Steps

- **LLMs have demonstrated effectiveness in various SE scenarios**
 - e.g., code retrieval, completion, review, etc., in VS, VSCode, GitHub, etc.
- **Latest LLMs (such as ChatGPT) show even stronger capabilities**
 - e.g., interactive code generation, completion, refinement, etc.
- **Next steps**
 - Responsible code models that can provide the original sources of the outputs
 - Trustworthy mechanisms that can verify the truthfulness of the outputs
 - Personalized software development scenarios
 - Code-centric models with NL as interface for low-code and no-code scenarios