**You**

**New Team**

**Query:** You need to fix the bug … and implement …

You

*Which parts of the repo should you retrieve for editing?*

# Code Editing Retrieval Problem



*Which parts of the repo should you retrieve for editing?*

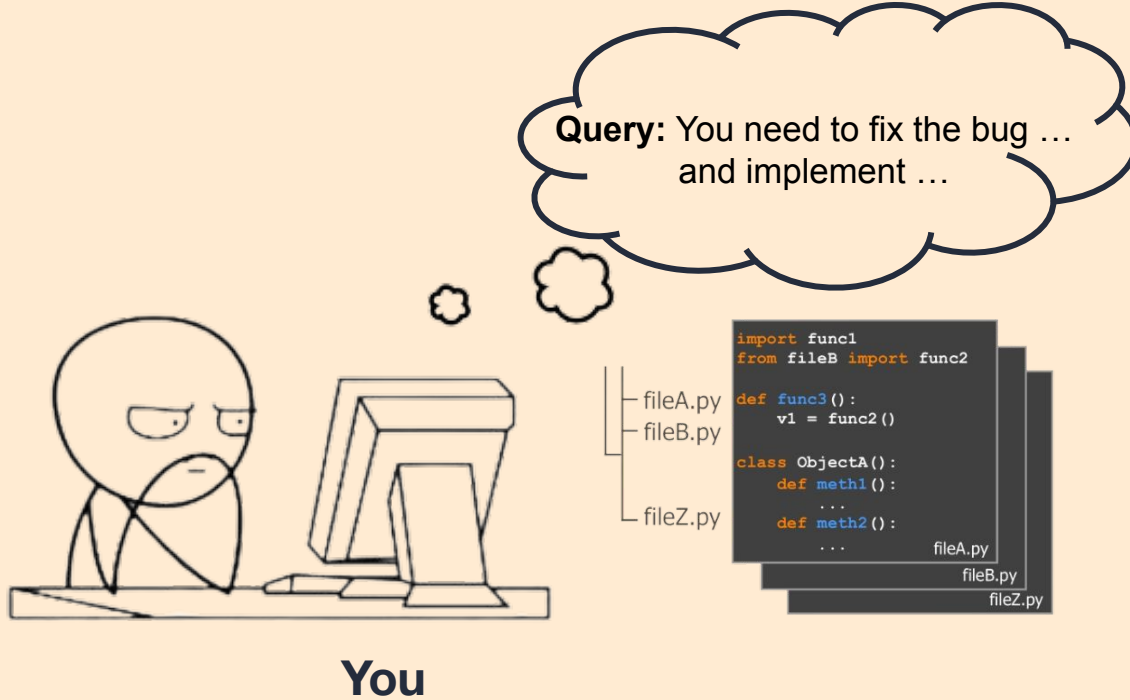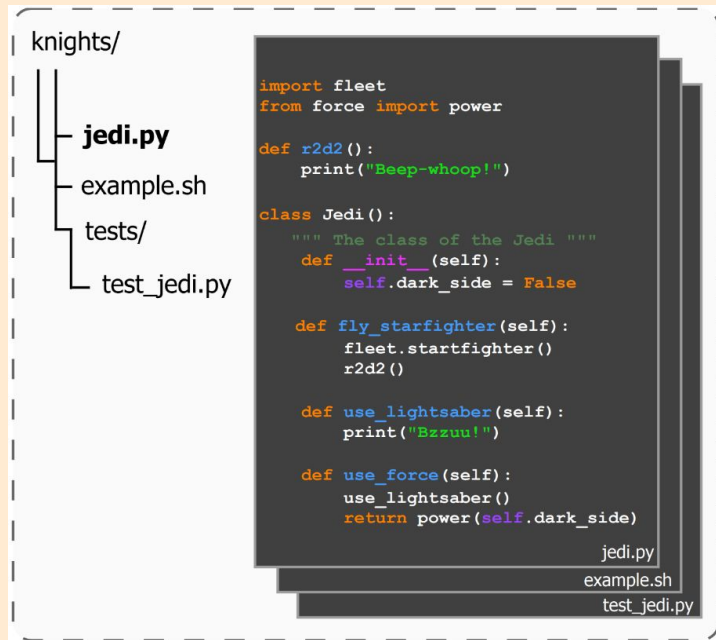**Train** a light-weight **code retriever** including **semantics** and **structure** across a repo.

**Code Repository**

```
knights/
├── jedi.py
├── example.sh
├── tests/
    └── test_jedi.py
```

```python
import fleet
from force import power

def r2d2():
    print("Beep-whoop!")

class Jedi():
    """ The class of the Jedi """
    def __init__(self):
        self.dark_side = False

    def fly_starfighter(self):
        fleet.startfighter()
        r2d2()

    def use_lightsaber(self):
        print("Bzzuu!")

    def use_force(self):
        use_lightsaber()
        return power(self.dark_side)
```

jedi.py

example.sh

test_jedi.py

**Code Repository**

**Preprocessing**

**Code Chucks**

# Code Chunks



**Code Repository**

**Code Chucks**

# Code Chunks with Repo-Hierarchy



**Code Repository**

**Code Chucks**

```
knights/jedi.py
def lightsaber():
    lightsaber_on()
```

$c_i$

```
knights/utils.py
def lightsaber_on():
    print("Bzzuu!")
```

$c_{out}$

$$c_i \; ; \; \texttt{[DOWN]} \; ; \; c_{out}$$

```
knights/jedi.py
def lightsaber():
    lightsaber_on()
```

```
knights/utils.py
def lightsaber_on():
    print("Bzzuu!")
```

$c_i$        $c_{out}$

# Embedding with Call Graph Context



$$c_i \; ; \; [\texttt{DOWN}] \; ; \; c_{out}$$

```
knights/jedi.py
def lightsaber():
    lightsaber_on()
```
$c_i$

```
knights/utils.py
def lightsaber_on():
    print("Bzzuu!")
```
$c_{out}$

# Training

$$\mathcal{L}(\theta)$$

# Training

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i}^{N}$$

# Training

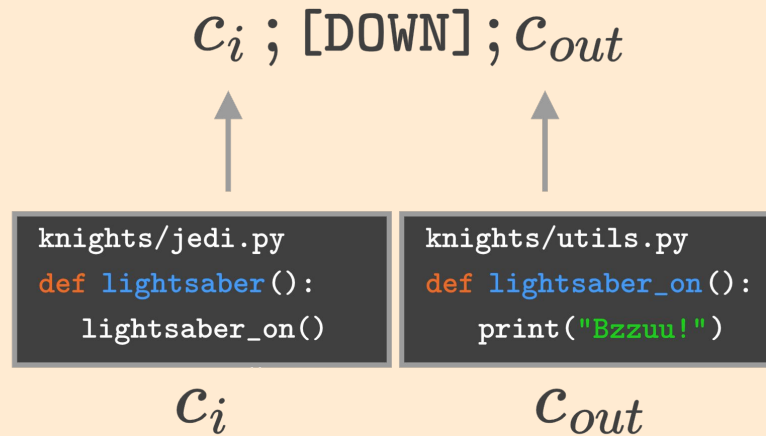$$\mathcal{L}(\theta) = \frac{1}{N} \sum_i^N \frac{1}{\mathcal{C}_i^*} \sum_{c^* \in \mathcal{C}_i^*}$$

# Training

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_i^N \frac{1}{\mathcal{C}_i^*} \sum_{c^* \in \mathcal{C}_i^*} \log \frac{\exp(\boldsymbol{q}_i \cdot \boldsymbol{c}^*)}{\exp(\boldsymbol{q}_i \cdot \boldsymbol{c}^*) + \sum_{c \in \mathcal{B}} \exp(\boldsymbol{q}_i \cdot \boldsymbol{c})}$$

# Training with Likelihood Loss

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_i^N \frac{1}{\mathcal{C}_i^*} \sum_{c^* \in \mathcal{C}_i^*} \log \frac{\exp(\boldsymbol{q}_i \cdot \boldsymbol{c}^*)}{\exp(\boldsymbol{q}_i \cdot \boldsymbol{c}^*) + \boxed{\sum_{c \in \mathcal{B}} \exp(\boldsymbol{q}_i \cdot \boldsymbol{c})}}$$

# Data & Model

**SWE-Bench** (Verified subset)

# Data & Model

**SWE-Bench** (Verified subset)

**LCA** (Bug localisation task)

# Data & Model

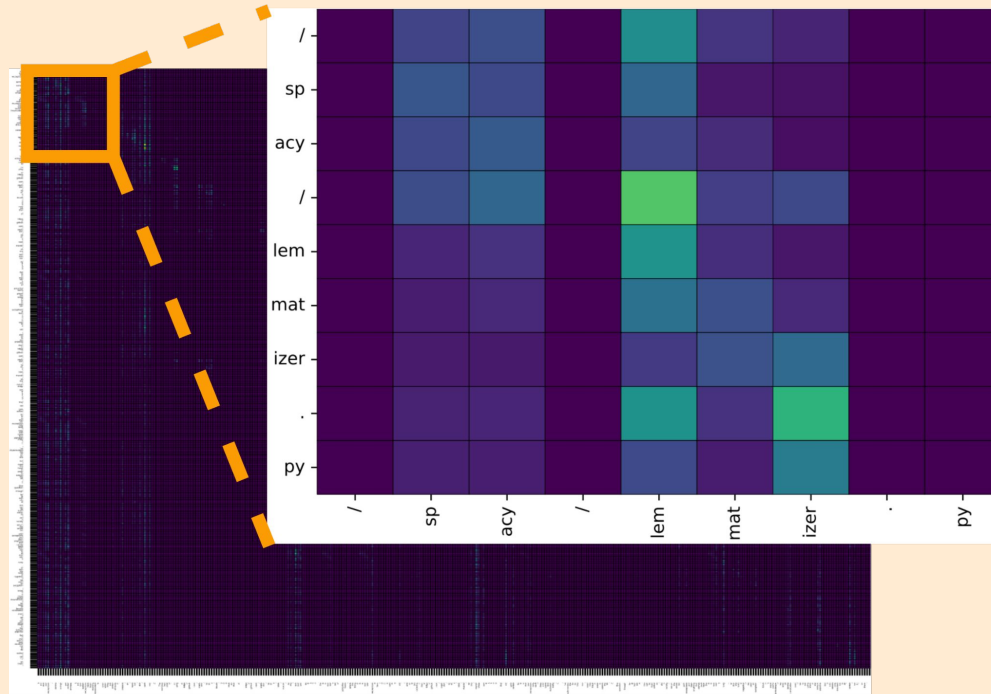**SWE-Bench** (Verified subset)

**LCA** (Bug localisation task)

**Codesage S** (128M parameters)

# Repo-Hierarchy is important

# Call graph improves multi-chunk retrieval

| Model | SWE Verified | | | LCA | | |
|---|---|---|---|---|---|---|
| | @5 | @20 | MRR | @5 | @20 | MRR |
| CodeSage S | 0.34 | 0.51 | 0.35 | 0.26 | 0.34 | 0.28 |
| CoRet − CG | 0.52 | 0.69 | 0.52 | **0.32** | 0.41 | 0.45 |
| CoRet − CG + file | **0.54** | 0.69 | 0.52 | 0.29 | 0.38 | 0.44 |
| CoRet | **0.54** | **0.71** | **0.53** | **0.32** | **0.47** | **0.47** |

# Call graph improves multi-chunk retrieval

| Model | SWE Verified | | | LCA | | |
|---|---|---|---|---|---|---|
| | @5 | @20 | MRR | @5 | @20 | MRR |
| CodeSage S | 0.34 | 0.51 | 0.35 | 0.26 | 0.34 | 0.28 |
| CoRet − CG | 0.52 | 0.69 | 0.52 | **0.32** | 0.41 | 0.45 |
| CoRet − CG + file | **0.54** | 0.69 | 0.52 | 0.29 | 0.38 | 0.44 |
| CoRet | **0.54** | **0.71** | **0.53** | **0.32** | **0.47** | **0.47** |

# Negatives from the same repo are best

# Recall +15 percentage points!



Chunk level retrieval on SWE-bench Verified

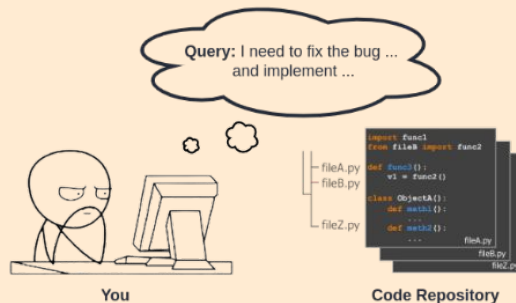**Train** a light-weight **code retriever** including **semantics** and **structure** across a repo.

# Paper

## Code Editing Retrieval Problem

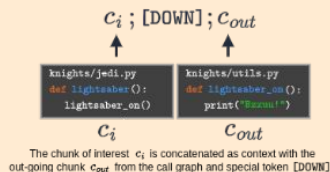Query: I need to fix the bug ... and implement ...

You

Code Repository

fileA.py
fileB.py
fileZ.py

*Which parts of the repo should you retrieve for editing?*

## Code Chunks with Repo-Hierarchy

Preprocessing

Code Repository

Code Chunks $c_i$

The code repo is split into semantically succinct unit we called **code chunks**.
We include **repo-hierarchy** structure by including the file path string.
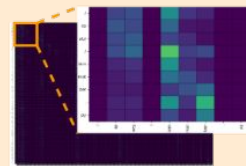
## Embedding with Call Graph Context

$$c_i \; ; \; [\text{DOWN}] \; ; \; c_{out}$$

```
knights/jedi.py          knights/utils.py
def lightsaber():        def lightsaber_on():
    lightsaber_on()          print("Bzzzz!")
```

$c_i$                    $c_{out}$

The chunk of interest $c_i$ is concatenated as context with the
out-going chunk $c_{out}$ from the call graph and special token $[\text{DOWN}]$.

### Call graph improves multi-chunk retrieval

| Model | SWE Verified | | | LCA | | |
|---|---|---|---|---|---|---|
| | @5 | @20 | MRR | @5 | @20 | MRR |
| CodeSage S | 0.34 | 0.51 | 0.35 | 0.26 | 0.34 | 0.28 |
| CoRet − CG | 0.52 | 0.69 | 0.52 | **0.32** | 0.41 | 0.45 |
| CoRet − CG + file | **0.54** | 0.69 | 0.52 | 0.29 | 0.38 | 0.44 |
| CoRet | **0.54** | **0.71** | **0.53** | **0.32** | **0.47** | **0.47** |

**CoRet:** Fine-tuned CodeSage S (130M parameters).
**SWE Verified:** Software Engineering Benchmark (Verified subset).
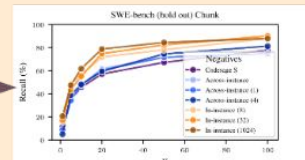**LCA:** Long Code Arena (Bug localisation task).

### Repo-hierarchy is important

## Training with Likelihood Loss

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i}^{N} \frac{1}{\mathcal{C}_i^*} \sum_{c^* \in \mathcal{C}_i^*} \log \frac{\exp(\boldsymbol{q}_i \cdot \boldsymbol{c}^*)}{\exp(\boldsymbol{q}_i \cdot \boldsymbol{c}^*) + \sum_{c \in \mathcal{B}} \exp(\boldsymbol{q}_i \cdot \boldsymbol{c})}$$

$N$ = Number of repo instances $i$. $\quad \mathcal{C}^*$ = Set of ground truth code chunks $c^*$. $\quad \boldsymbol{q}$ = Natural language query.
$\mathcal{B}$ = Random negative sample in the same repo instance.

### Negatives from the same repo are best

Training with negatives from the same repo instance improve
over negatives across repo instances (standard *in-batch* negatives).

### Recall +15 percentage points