# BACKGROUND

Documenting code is **crucially** important.

# BACKGROUND

Documenting code is **crucially** important.

It allows developers to better **understand** code.

# BACKGROUND



Documenting code is **crucially** important.

It allows developers to better **understand** code.

Documenting code is both **labor-intensive** and frequently **neglected.**

# BACKGROUND

Documenting code is **crucially** important.

It allows developers to better **understand** code.

Documenting code is both **labor-intensive** and frequently **neglected**.

Automated **code summarization** has emerged as a promising solution.

# CODE SUMMARIZATION

```java
public String toString() {
      final StringBuffer s = new StringBuffer();
      final int size = size();
      for (int i = 0; i < size; i++)
          s.append(getInt(i));
      return s.toString();
}
```



DL-Model

Returns a string representation of this vector.

# PIPELINE

Data Collection

Pre-processing and filtering

Splitting

Training

Evaluation

# PIPELINE

<code, summary> pairs

| Data Collection |
| --- |
| Pre-processing and filtering |
| Splitting |
| Training |
| Evaluation |

# PIPELINE



Pre-processing and filtering

&lt;code, summary&gt; pairs

| Data Collection |
| :---: |
| Pre-processing and filtering |
| Splitting |
| Training |
| Evaluation |

# BACKGROUND



<code, summary> pairs

Pre-processing and filtering

Splitting

Train

Validation

Test

| Data Collection |
| Pre-processing and filtering |
| Splitting |
| Training |
| Evaluation |

# BACKGROUND



Pre-processing and filtering

&lt;code, summary&gt; pairs

Splitting

Train

Validation

Test

Training

Data Collection

Pre-processing and filtering

Splitting

Training

Evaluation

# PIPELINE



<code, summary> pairs

Pre-processing and filtering

Splitting

Train

Validation

Test

Evaluation

Training

Data Collection

Pre-processing and filtering

Splitting

Training

Evaluation

# PIPELINE



<code, summary> pairs

Splitting

Train

Validation

Test

Pre-processing and filtering

Evaluation

Training

Data Collection

Pre-processing and filtering

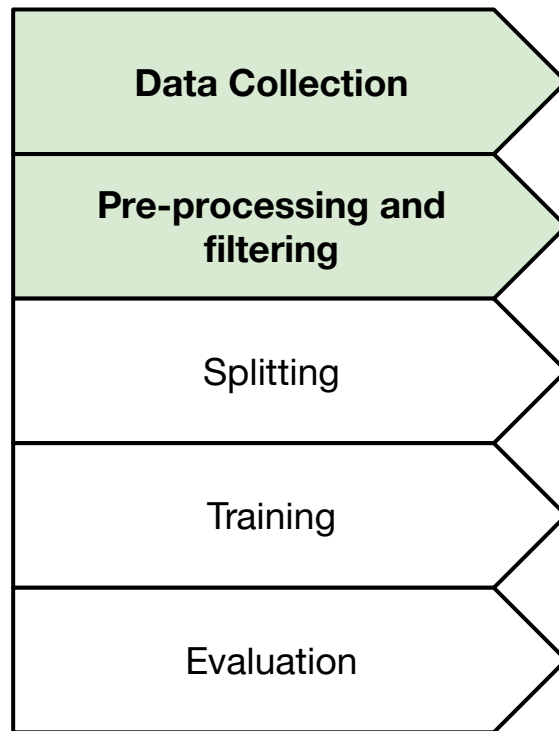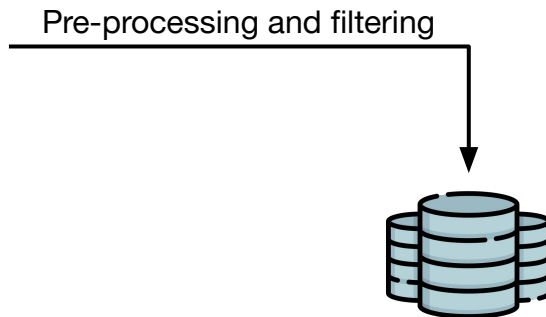Splitting

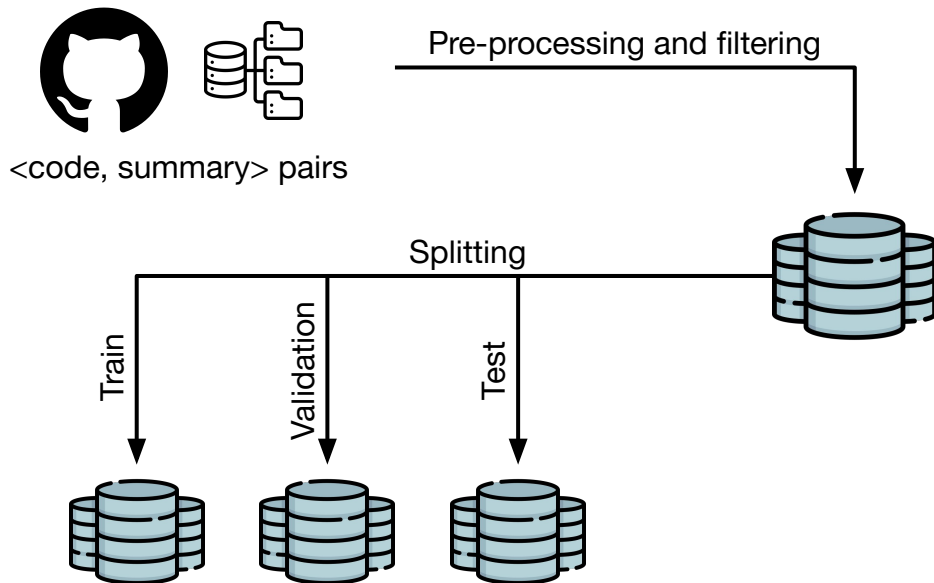Training

Evaluation

# LOW QUALITY INSTANCES

```
/* Returns the high-value
 * for an item within a series. */
```

# LOW QUALITY INSTANCES

```
/* Returns the high-value
 * for an item within a series. */
```

returns the high value

# LOW QUALITY INSTANCES

```
/* Returns the high-value
 * for an item within a series. */
```

returns the high value

```
/* <p> Builds the JASPIC application context. </p> */
```

p builds the jaspic application context p

# LOW QUALITY INSTANCES

```
/* Returns the high-value
 * for an item within a series. */
```

returns the high value

```
/* <p> Builds the JASPIC application context. </p> */
```

p builds the jaspic application context p

```
public void testConstructor() {
      System TestResult str;
      System TestID testID1;
      ...
}
```

test the constructor

# LOW QUALITY INSTANCES



Propose **CAT** (**C**ode-comment cle**A**ning **T**ool), a **rule-based filtering tool** for automatically scanning and detecting the occurrences and distribution of data noises for a given dataset.

FSE 22'

# LOW QUALITY INSTANCES

```java
public HashSet getCommandResultsRootFeatures() {
    HashSet rootFeatureSet = new HashSet();
    Feature belowSplitRoot = null;
    Feature aboveSplitRoot = null;
    if (belowSplitTranscript != null) {
        belowSplitRoot = belowSplitTranscript.getRootFeature();
        rootFeatureSet.add(belowSplitRoot);
    }
    if (aboveSplitTranscript != null) {
        aboveSplitRoot = aboveSplitTranscript.getRootFeature();
        if (aboveSplitRoot != belowSplitRoot)
            rootFeatureSet.add(aboveSplitRoot);
    }
    return rootFeatureSet;
}
```

Invoked AFTER the command is executed.

# LOW QUALITY INSTANCES



They propose **SIDE** (**S**ummary al**I**gnment to co**D**e s**E**mantics), a new metric leveraging contrastive learning to model the characteristics of **suitable and unsuitable** code summaries for a given code.

ICSE 24'

# LOW QUALITY INSTANCES



**METHOD**

```
public Element asElement() {
    return this.component.createCopy();
}
```

**REFERENCE SUMMARY**

not yet documented

**GENERATED SUMMARY**

returns a copy of this component as an element

**SCORES**

| | | | SOTA | | | | OUR | | HUMAN | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| BLEU-1 | BertScore-R | SentenceBERT_CS | InferNet_CS | ROUGE-1-P | Rouge-4-R | Rouge-W-R | SIDE | DA | Adequacy | Concise | Fluency |
| 0.34 | 0.00 | 0.08 | 0.43 | 0.00 | 0.00 | 0.00 | 0.91 | 88 | 4 | 4 | 5 |

"**SIDE** is the metric that better describes humans' assessment of summary quality."

# LOW QUALITY INSTANCES

**METHOD**

```java
public Element asElement() {
    return this.component.createCopy();
}
```

**REFERENCE SUMMARY**

not yet documented

**GENERATED SUMMARY**

returns a copy of this component as
an element

**SCORES**

| | | | SOTA ↓ | | | | OUR ↑ | | HUMAN ↑ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| BLEU-1 | BertScore-R | SentenceBERT_CS | InferNet_CS | ROUGE-1-P | Rouge-4-R | Rouge-W-R | SIDE | DA | Adequacy | Concise | Fluency |
| 0.34 | 0.00 | 0.08 | 0.43 | 0.00 | 0.00 | 0.00 | 0.91 | 88 | 4 | 4 | 5 |

"**SIDE** is the metric that better describes humans' assessment of summary quality."

"Provides a continuous score ranging between [-1, 1]"

# LOW QUALITY INSTANCES



**METHOD**

```java
public Element asElement() {
    return this.component.createCopy();
}
```

**RE~~FERENCE SUMMARY~~**

**GENERATED SUMMARY**

returns a copy of this component as an element

**SCORES**

|  |  |  | SOTA |  |  |  | OUR |  | HUMAN |  |  |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| BLEU-1 | BertScore-R | SentenceBERT_CS | InferNet_CS | ROUGE-1-P | Rouge-4-R | Rouge-W-R | SIDE | DA | Adequacy | Concise | Fluency |
| 0.34 | 0.00 | 0.08 | 0.43 | 0.00 | 0.00 | 0.00 | 0.91 | 88 | 4 | 4 | 5 |

**"SIDE** is the metric that better describes humans' assessment of summary quality."

"Provides a continuous score ranging between [-1, 1]"

# LOW QUALITY INSTANCES

```
public Element asElement() {
    return this.component.createCopy();
}
```

not yet documented

SIDE

# LOW QUALITY INSTANCES

```
public Element asElement() {
    return this.component.createCopy();
}
```

not yet documented

SIDE

-0.19

# LOW QUALITY INSTANCES

```
public Element asElement() {
    return this.component.createCopy();
}
```

not yet documented

SIDE → -0.19

```
public Element asElement() {
    return this.component.createCopy();
}
```

returns a copy of this component as an element

SIDE → 0.91

# LOW QUALITY INSTANCES

```
public Element asElement() {
    return this.component.createCopy();
}
```

```
not yet documented
```

SIDE → -0.19 → **Incoherent**

```
public Element asElement() {
    return this.component.createCopy();
}
```

```
returns a copy of this component as an
element
```

SIDE → 0.91 → **Coherent**

Can filtering out **incoherent code-comment pairs** serve as an effective strategy for **optimizing** code-summarization datasets**?**

# EMPIRICAL STUDY



**Training Sets**

Funcom          TI-CodeSum

# EMPIRICAL STUDY

**Training Sets**

Funcom

TI-CodeSum

~1M Java <code, summary> pairs

~54k Java <code, summary> pairs

# EMPIRICAL STUDY

# EMPIRICAL STUDY

# EMPIRICAL STUDY

# RQ0

How do code summarization datasets measure up in terms of **code-comment coherence?**

# RQ0

| FUNCOM |
| --- |
| 0.81 mean SIDE score |

# RQ0

| FUNCOM |
|---|
| 0.81 mean SIDE score |

| TL-CODESUM |
|---|
| 0.83 mean SIDE score |

# RQ0

| FUNCOM |
|---|
| 0.81 mean SIDE score |

| TL-CODESUM |
|---|
| 0.83 mean SIDE score |

**More than 50%** of the instances
have sub-optimal SIDE scores below 0.9.

# RQ1

How does a **coherence-aware strategy** selection **impact the performance** of neural code summarization models?

# RQ1

# RQ1

# RQ1



TI-CodeSum

CAT

SIDE

0.5  0.6  0.7  0.8  0.9

Funcom

CAT

SIDE

0.5  0.6  0.7  0.8  0.9

# RQ1



Full — CodeT5+

0.5 — CodeT5+

0.6 — CodeT5+

0.7 — CodeT5+

0.8 — CodeT5+

0.9 — CodeT5+

# RQ1

# RQ1

# RQ1



Models **performances are comparable** to those obtained when fine-tuning the model on the **full** training sets.

# RQ1

Models **performances are comparable** to those obtained when fine-tuning the model on the **full** training sets.

This happens **using only 50%** of the training set.

# RQ2

How does the coherence-aware strategy selection **compare** with a **random baseline**?

# RQ2



0.9 — CodeT5+

Random — CodeT5+

# RQ2

# RQ2

Filtering by code-comment coherence provides models with **comparable** effectiveness to those trained on **randomly selected** instances.

# OUTCOMES

Code-comment coherence **might not be a problem** in state-of-the-art datasets.

The results clearly indicate that state-of-the-art datasets contain **instances that do not contribute** to improving the models' effectiveness.

**Other quality aspects** of code and comments that have not been explored yet (such as readability) may be important for smartly selecting
the training instances.

Future work could explore different criteria for data selection that identify the **most informative** subsets for training.

# OUTCOMES

Code-comment coherence **might not be a problem** in state-of-the-art datasets.

The results clearly indicate that state-of-the-art datasets contain **instances that do not contribute** to improving the models' effectiveness.

**Other quality aspects** of code and comments that have not been explored yet (such as readability) may be important for smartly selecting the training instances.

Future work could explore different criteria for data selection that identify the **most informative** subsets for training.

# OUTCOMES

Code-comment coherence **might not be a problem** in state-of-the-art datasets.

The results clearly indicate that state-of-the-art datasets contain **instances that do not contribute** to improving the models' effectiveness.

**Other quality aspects** of code and comments that have not been explored yet (such as readability) may be important for smartly selecting the training instances.

Future work could explore different criteria for data selection that identify the **most informative** subsets for training.

# OUTCOMES

Code-comment coherence **might not be a problem** in state-of-the-art datasets.

The results clearly indicate that state-of-the-art datasets contain **instances that do not contribute** to improving the models' effectiveness.

**Other quality aspects** of code and comments that have not been explored yet (such as readability) may be important for smartly selecting the training instances.

Future work could explore different criteria for data selection that identify the **most informative** subsets for training.

# SUMMARY



**PIPELINE**

<code, summary> pairs

Pre-processing and filtering

Splitting

Train | Validation | Test

Evaluation

Training

- Data Collection
- Pre-processing and filtering
- Splitting
- Training
- Evaluation

**RQ1**

Full
0.5
0.6
0.7
0.8
0.9

CodeT5+
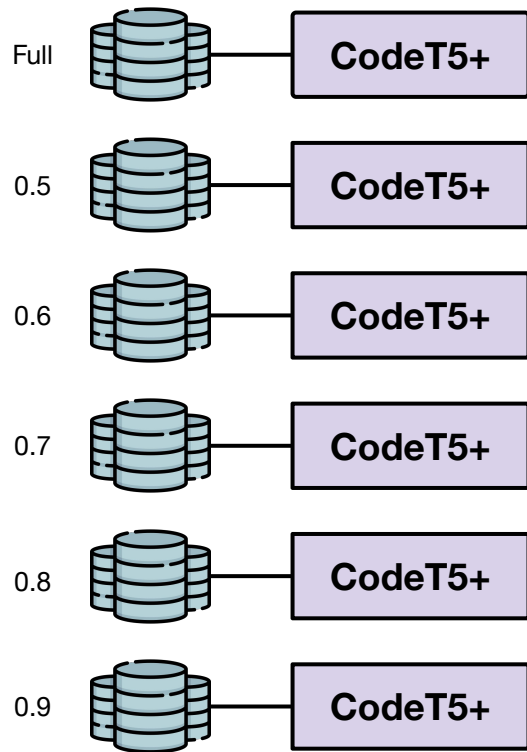
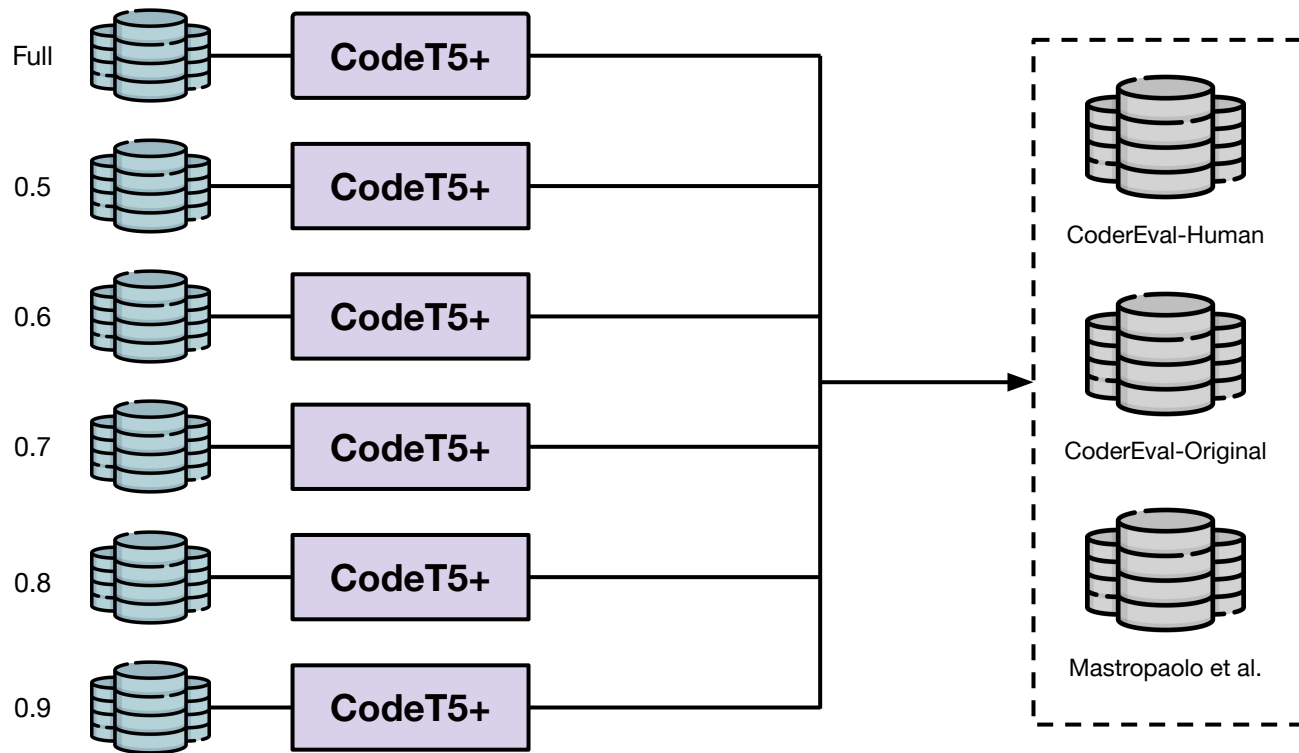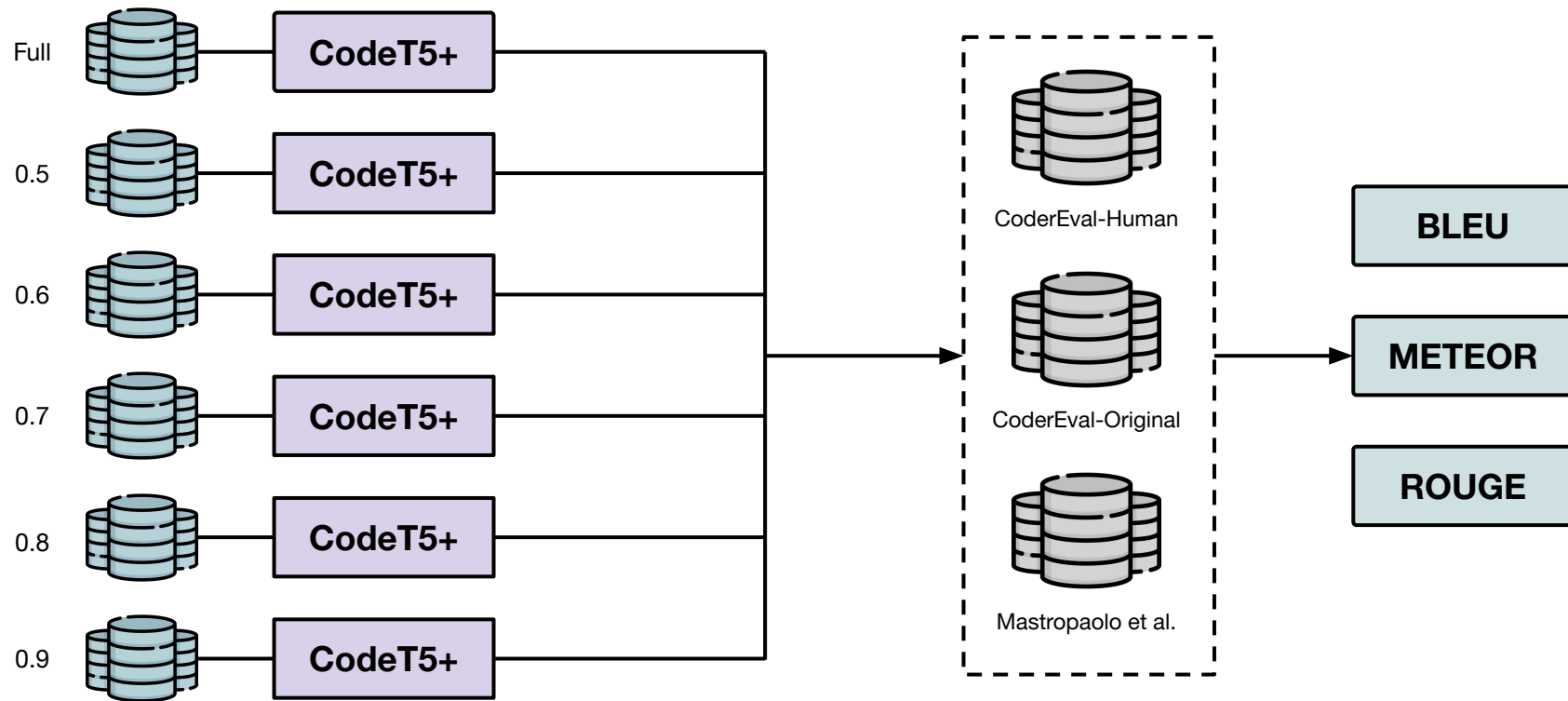CoderEval-Human
CoderEval-Original
Mastropaolo et al.

BLEU
METEOR
ROUGE

**RQ1**

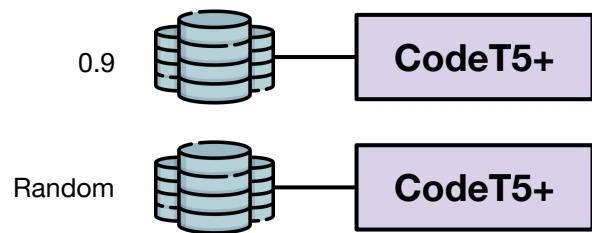Models **performances are comparable** to those obtained when fine-tuning the model on the **full** training sets.

This happens **using only 50%** of the training set.

**RQ2**

0.9
Random

CodeT5+

CoderEval-Human
CoderEval-Original
Mastropaolo et al.

BLEU
METEOR
ROUGE

**RQ2**

Filtering by code-comment coherence provides models with **comparable** effectiveness to those trained on **randomly selected** instances.

**OUTCOMES**

Code-comment coherence **might not be a problem in** state-of-the-art datasets.

Other **quality aspects** of code and comments that have not been explored yet (such as readability) may be important for smartly selecting the training instances.
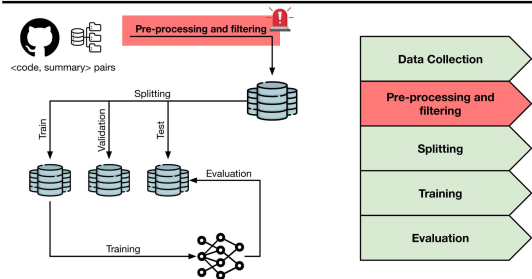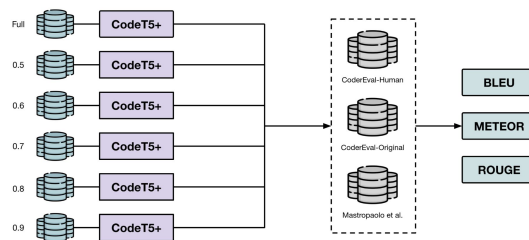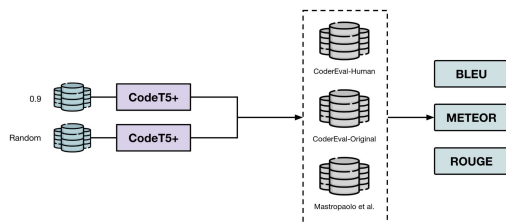
The results clearly indicate that state-of-the-art datasets contain **instances that do not contribute** to improving the models' effectiveness.

Future work could explore different criteria for data selection that identify the **most informative** subsets for training.