UNIVERSITY OF **ILLINOIS** URBANA-CHAMPAIGN

Artifacts Available V1.1

Artifacts Evaluated Functional V1.1
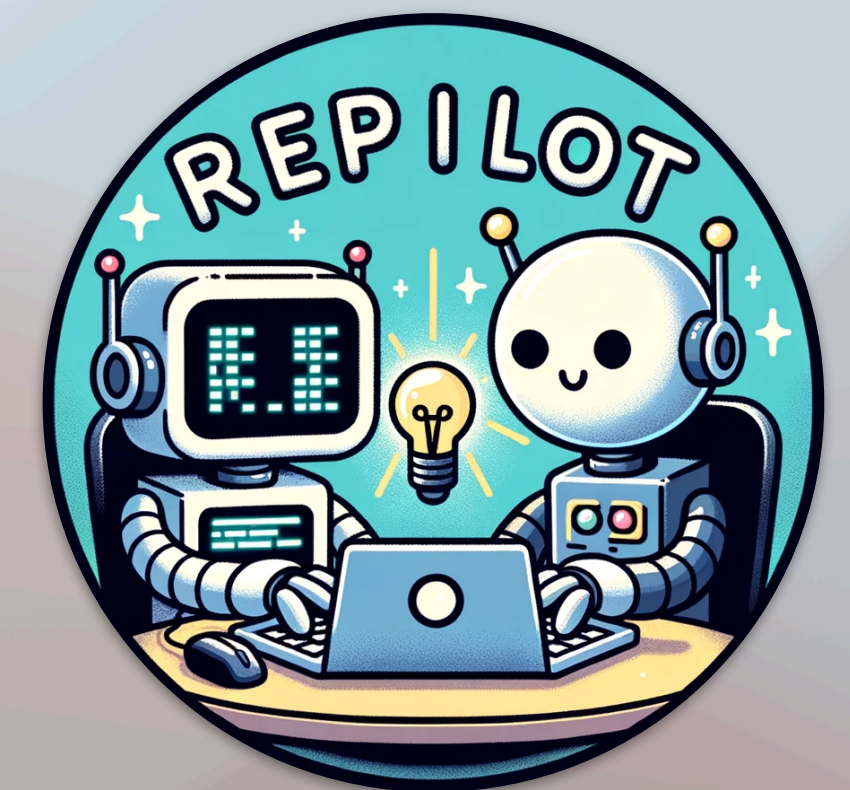
Artifacts Evaluated Reusable V1.1

# Copiloting the Copilots (Repilot)
## Fusing Large Language Models with Completion Engines for Automated Program Repair
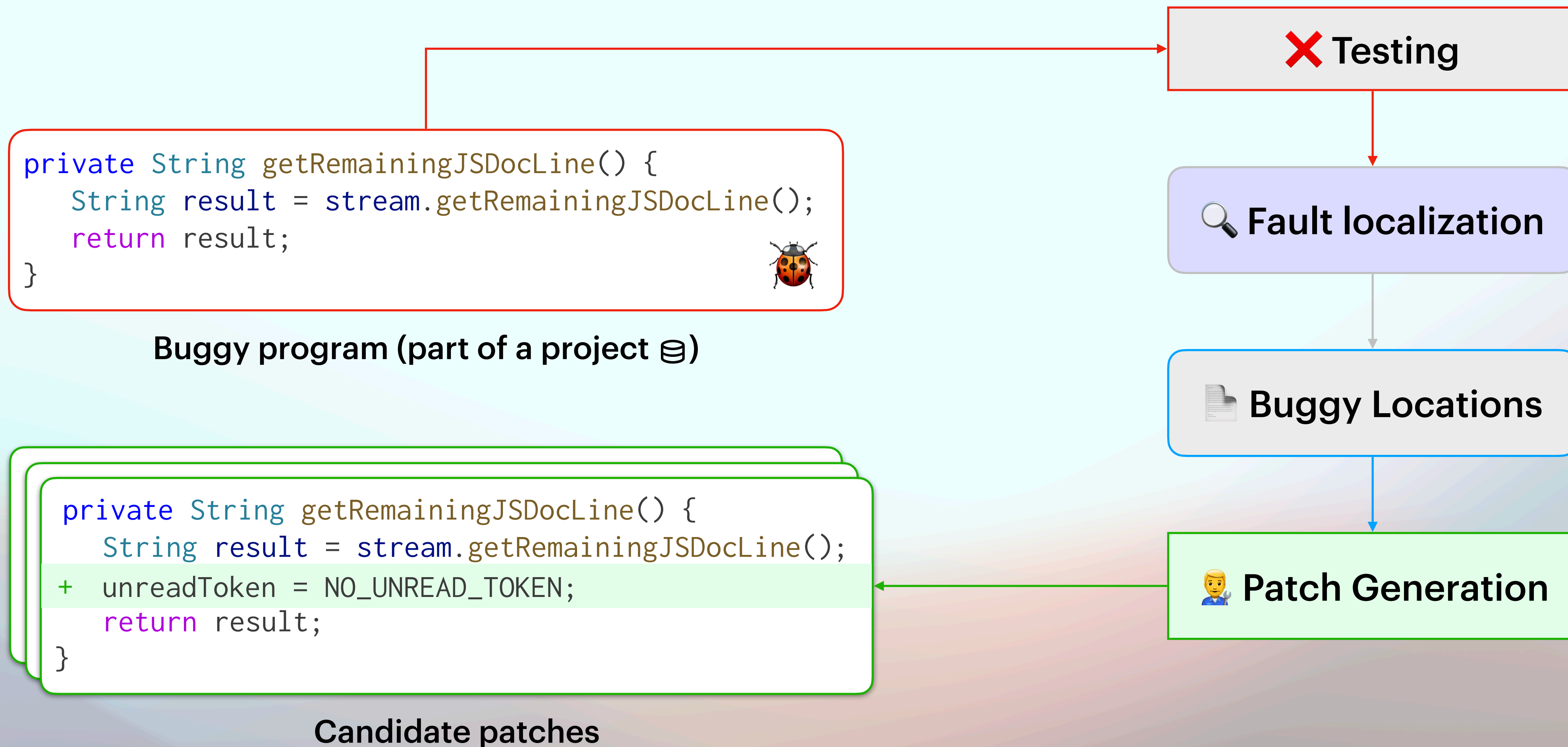
**Yuxiang Wei**
𝕏 @YuxiangWei9

**Chunqiu Steven Xia**
𝕏 @steven_xia_
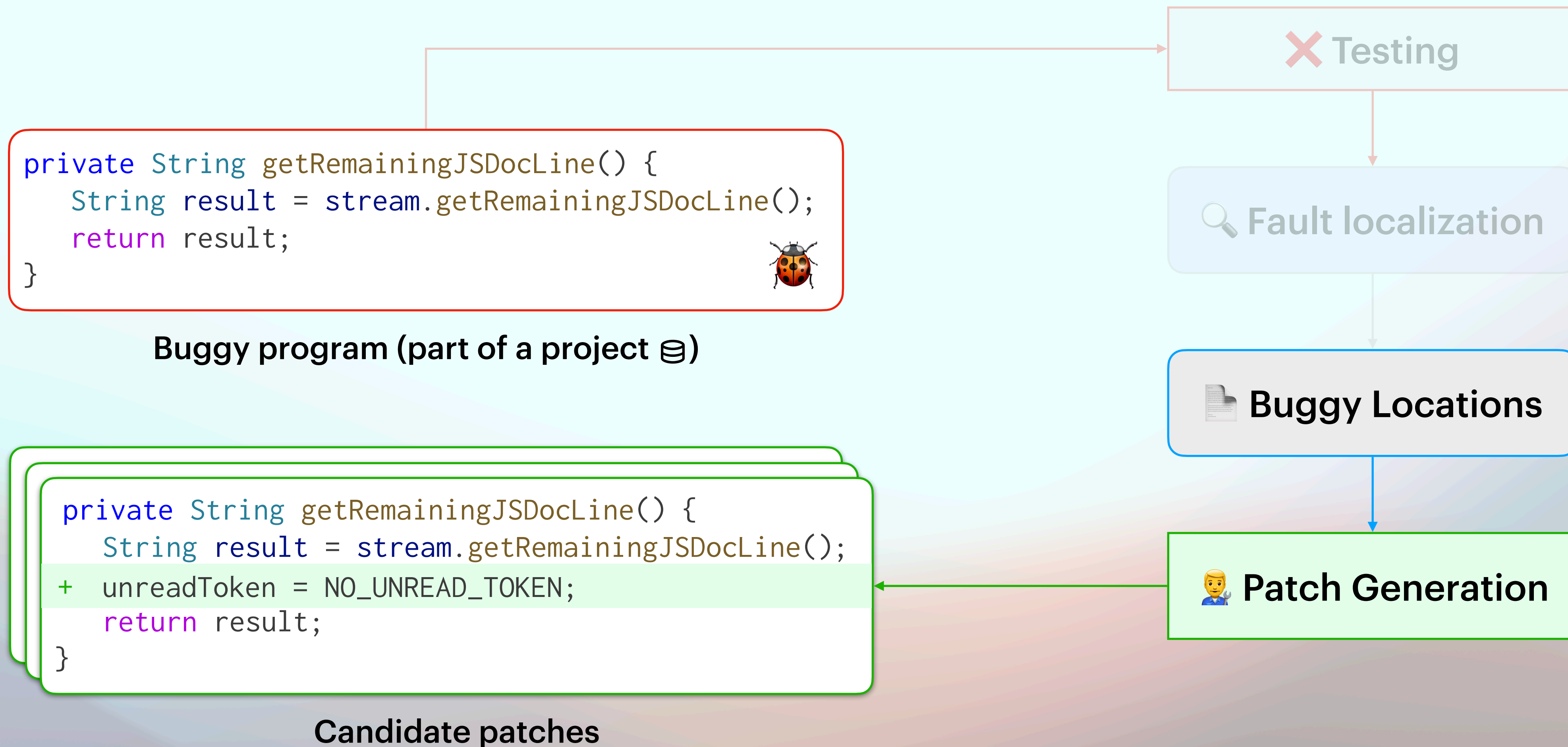
**Lingming Zhang**
𝕏 @LingmingZhang

Powered by DALL·E 3

# Automated Program Repair (APR)
## The two stages

```java
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
    return result;
}
```
🐞

Buggy program (part of a project 🗄)

```java
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
+   unreadToken = NO_UNREAD_TOKEN;
    return result;
}
```

Candidate patches

❌ Testing

🔍 Fault localization

📄 Buggy Locations

👨‍🏫 Patch Generation

# Automated Program Repair (APR)
## The two stages

```
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
    return result;
}                                              🐞
```

Buggy program (part of a project 🗄)

```
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
+   unreadToken = NO_UNREAD_TOKEN;
    return result;
}
```

Candidate patches

❌ Testing

🔍 Fault localization

📄 Buggy Locations

👨‍🏫 Patch Generation

# Patch Generation with Large Language Models (LLMs)

```
❌ Testing  →  🔍 Fault location
```

```
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
    return result;
}                                        🐞
```

**Buggy program (part of a project 🗄)**

📄 **Buggy Locations**

```
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
    <Fault Location>
    return result;
}
```

**"Cloze"**

**Langauge Model**

```
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
+   unreadToken = NO_UNREAD_TOKEN;
    return result;
}
```

**Candidate patches**

# LLMs vs. Autocompletion



$$\Sigma_{LM}^* \to \text{DecRep}$$

$$\Sigma_{LM} \to [0, 1]$$

**Decoder**

DecRep

...

3%

91%

Probability Map

$$\Sigma_{LM}^*$$

Decoder inputs
(Limited context size)

Code generation with LLMs

5

# LLMs vs. Autocompletion



$$\Sigma_{LM}^* \to \text{DecRep}$$

$$\Sigma_{LM} \to [0, 1]$$

**Decoder**

DecRep

...

3%

✓ 91%

Probability Map

$$\Sigma_{LM}^*$$

Decoder inputs
(Limited context size)

Code generation with LLMs

```
…
public MultiplePiePlot(…) {
    super();
    this.set|
…
```

**Program and caret position**

**Completion
Engine**

```
Dataset
DatasetGroup
BackgroundAlpha
DataExtractOrder
DrawingSupplier
ForegroundAlpha
…
```

**Completions**

Semantics-based autocompletion

# Problems with LLM-only code generation
## The autoregressive (token-by-token) nature

**Language model Predictions**

| | | |
|---|---|---|
| String | 91% | ✕ |
| Name | 3% | ✕ |
| **End** | **0.2%** | ✓ |
| ... | | |

**Completions**

**as**EndTag
**as**StartTag
**as**Comment
**as**Doctype
**as**Character

```
String name = t.as|
```

**ⓐ Generating infeasible tokens**

# Problems with LLM-only code generation
## The autoregressive (token-by-token) nature



Language model Predictions — Completions

String 91% ✕
Name 3% ✕
End 0.2% ✓
...

asEndTag
asStartTag
asComment
asDoctype
asCharacter

String name = t.as

❶ Generating infeasible tokens

Name 16%
Tag 7%
...

asEndTag ✓

String name = t.asEndTag

❷ Hard to generate rare tokens

6

# How Repilot works



String name = t.as

Generated tokens

Language Model

String
**91%**

Name
**3%**

End
**0.2%**

Search space

# How Repilot works

# How Repilot works



**Language model Predictions**

| | | |
|---|---|---|
| String | **91%** | × |
| Name | **3%** | × |
| End | **0.2%** | ✓ |
| ... | | |

String name = t.as|

**Completions**

**as**End Tag
**as**StartTag
**as**Comment
**as**Doctype
**as**Character

Generated tokens

String name = t.as

**Language Model**

Search space

String **91%**
Name **3%**
End **0.2%**

**Completion Engine**

× **Pruning**

Pruned search space

String **0%**
Name **0%**
End **0.2%**

Token sampling loop
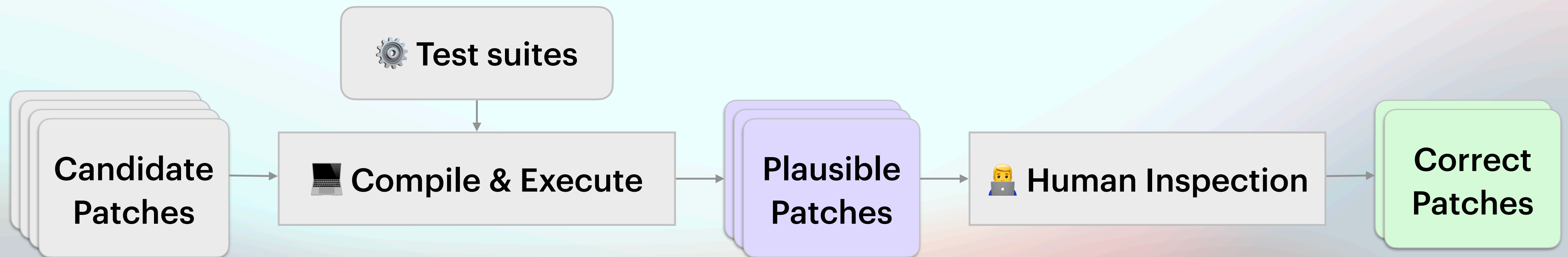
# How Repilot works

# How Repilot works



7

# Evaluation
## Evaluation pipeline and metrics

# Evaluation
## Comparison with existing tools

| Tool | Methodology | #Correct Fixes | | |
|------|-------------|----------------|---|---|
| | | Defects4J 1.2 | Defects4J 2.0 | Total |
| CoCoNuT | NMT | 30 | - | - |
| DLFix | NMT | 32 | - | - |
| PraPR | Template | 35 | - | - |
| TBar | Template | 41 | 7 | 48 |
| CURE | NMT | 43 | - | - |
| RewardRepair | NMT | 45 | 24 | 69 |
| Recoder | NMT | 51 | 10 | 61 |
| AlphaRepair | LLM | 52 | 34 | 86 |
| **Repilot** | LLM | **66** | **50** | **116** |

NMT means "Neural Machine Translation"



Number of unique fixes
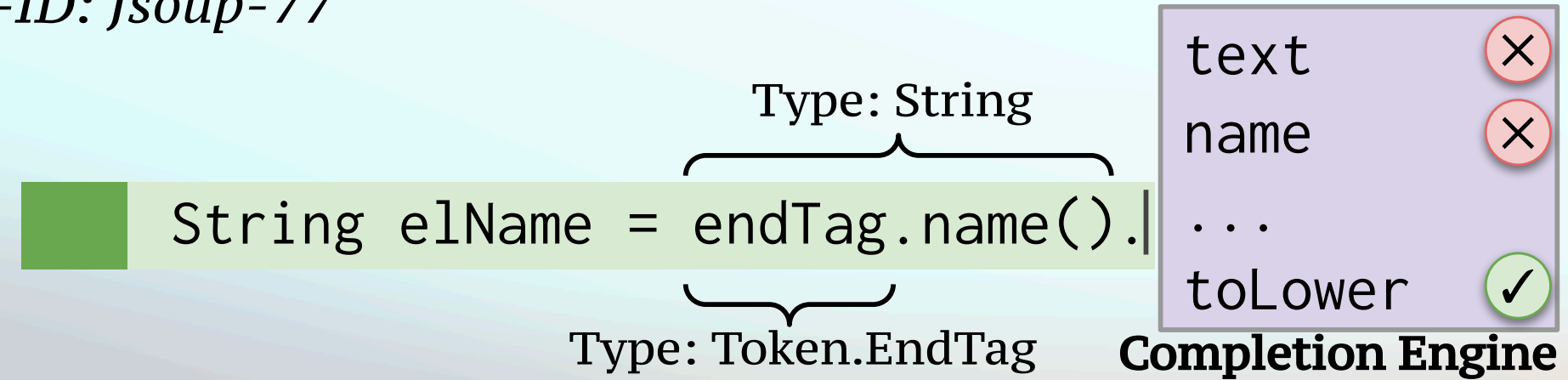
# Evaluation
## Unique fixes generated by Repilot

```
private void popStackToClose(Token.EndTag endTag) {
-     String elName = endTag.name();
+     String elName = endTag.name().toLowerCase();
    Element firstFound = null;
```
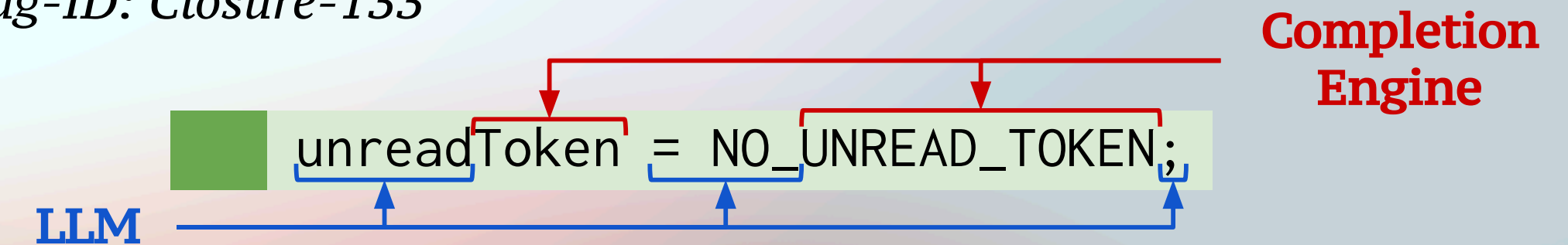*Bug-ID: Jsoup-77*

Type: String

String elName = endTag.name().|

Type: Token.EndTag

| text | ⊗ |
| name | ⊗ |
| ... | |
| toLower | ✓ |

**Completion Engine**

*Patch Generation Process*

**Completion engine filters out invalid tokens**

```
private String getRemainingJSDocLine() {
    String result = stream.getRemainingJSDocLine();
+     unreadToken = NO_UNREAD_TOKEN;
    return result;
}
```
*Bug-ID: Closure-133*

**Completion Engine**

unreadToken = NO_UNREAD_TOKEN;

**LLM**

*Patch Generation Process*

**Interaction between LLM and completion engine**

# Evaluation
## Generalizability

Repilot is generalizable across bug subjects and models

| Variant | Model | Subject of Bugs | Generation Time | %Compilable Patches | #Correct Fixes |
|---------|-------|-----------------|-----------------|---------------------|----------------|
| Base LLM | CodeT5-large | Defects4J 1.2 | **0.232s** | 43.2% | 37 |
| Repilot | CodeT5-large | Defects4J 1.2 | 0.248s | **63.4%** | **42** |
| Base LLM | CodeT5-large | Defects4J 2.0 | **0.230s** | 46.7% | 41 |
| Repilot | CodeT5-large | Defects4J 2.0 | 0.247s | **64.8%** | **45** |
| Base LLM | INCODER-6.7B | Defects4J 1.2 | **1.70s** | 32.4% | 48 |
| Repilot | INCODER-6.7B | Defects4J 1.2 | **1.70s** | **47.2%** | **54** |
| Base LLM | INCODER-6.7B | Defects4J 2.0 | **1.67s** | 34.6% | 45 |
| Repilot | INCODER-6.7B | Defects4J 2.0 | 1.69s | **48.0%** | **46** |

# Repilot

**Copiloting the Copilots**

Fuses **Large Language Models** with **Completion Engines** for more **effective Patch Generation** in Automated Program Repair

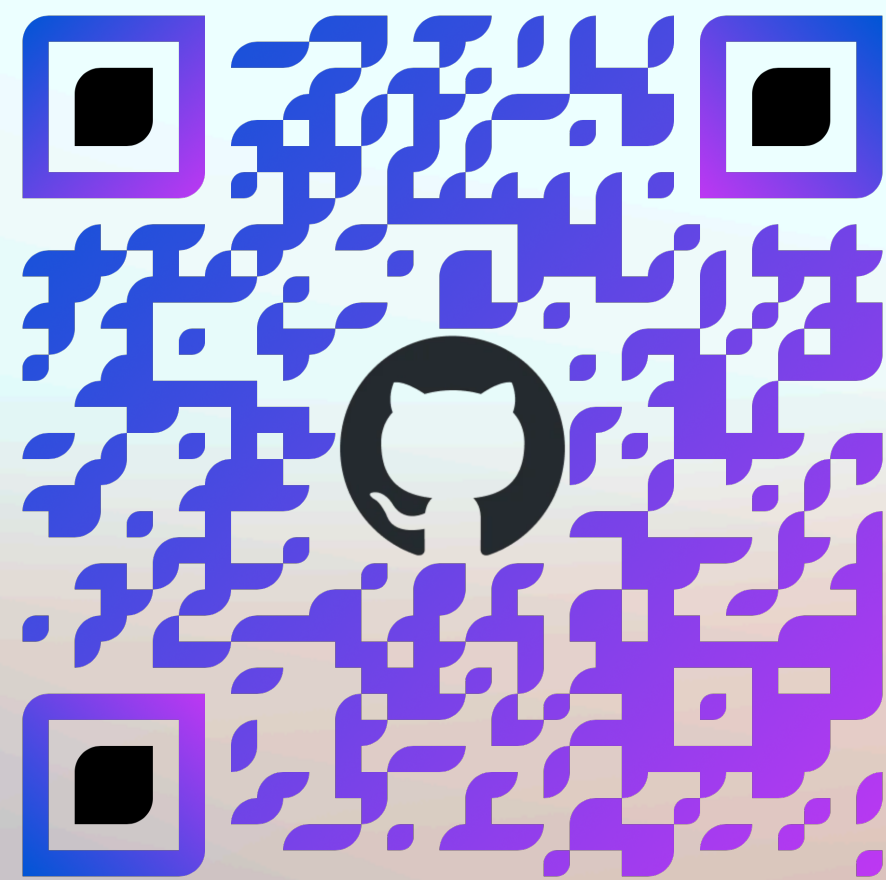Can be applied to general **program synthesis**

(⭐104) ise-uiuc/Repilot



**Code & Docker & Artifact**

arXiv:2309.00608



**Paper**

𝕏 @YuxiangWei9



**Twitter/X**



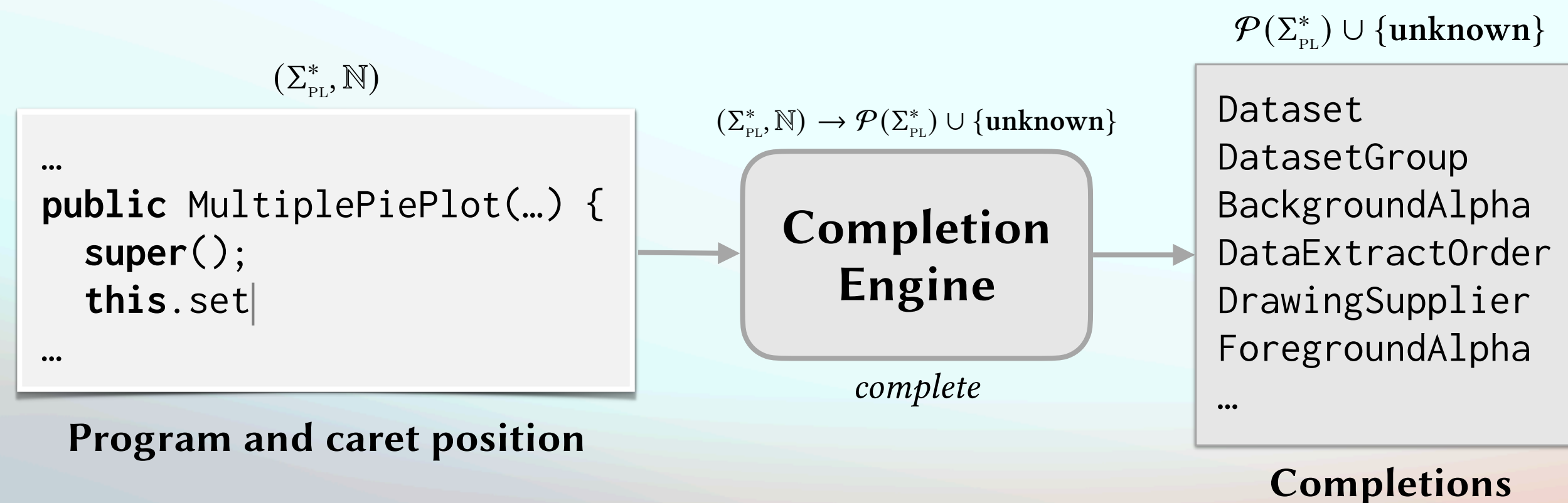Powered by **DALL·E 3**

**Yuxiang Wei**

✉ ywei40@illinois.edu

# **Supplementaries**

# Will Repilot go wrong?
## Conditional soundness of Repilot

**Theorem:** Repilot is sound with a **strict completion engine**

$(\Sigma_{\text{PL}}^*, \mathbb{N})$

```
…
public MultiplePiePlot(…) {
  super();
  this.set|
…
```

**Program and caret position**

$(\Sigma_{\text{PL}}^*, \mathbb{N}) \rightarrow \mathcal{P}(\Sigma_{\text{PL}}^*) \cup \{\textbf{unknown}\}$

**Completion Engine**

*complete*

$\mathcal{P}(\Sigma_{\text{PL}}^*) \cup \{\textbf{unknown}\}$

```
Dataset
DatasetGroup
BackgroundAlpha
DataExtractOrder
DrawingSupplier
ForegroundAlpha
…
```

**Completions**

**Assume**

$(\text{prog}, \text{caret}) \vDash \Phi$

$\text{completions} = \text{complete}(\text{prog}, \text{caret})$

$\text{completions} \neq \textbf{unknown}$

**The completion engine is strict if and only if**

$\forall c \notin \text{Prefix}(\text{completions}), (\text{prog}', \text{caret}') \nvDash \Phi$

where $\text{prog}' = \text{prog}[\text{caret} \leftarrow c]$

$\quad\quad\quad \text{caret}' = \text{caret} + |c|$

$\text{Prefix}(\,\cdot\,) = \{c \mid s \in \cdot \text{ and } c \text{ is a prefix of } s \text{ or vice versa}\}$
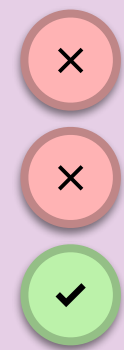
Intuitively, the completion engine should capture all possible continuations, but can do over-approximation. **"unknown"** is an over-approximation.

# Why resample instead of direct pruning?

In this case, direct pruning cannot do anything, and may result in a wrong path, e.g., `.method_a`

**Language model
Predictions**

| | |
|---|---|
| `.method_a` | ✕ |
| `.method_b` | ✕ |
| `.method_c` | ✔ |
| ... | |

**Completions**

`obj`

`var result = obj`