

AI for detecting code plagiarism

1. What are Code Clones

คือ โค้ดที่ซ้ำกัน เหมือน กัน หรือ คล้าย กัน

Type 1

เหมือนกัน 100 %

```
public int sum(int a, int b){
    int sum;
    sum = a + b;
    return sum;
}

public int sum(int a, int b){
    int sum;
    sum = a + b;
    return sum;
}
```

Type 2

เหมือนกัน 100 %

```
public int sum(int a, int b){
    int sum;
    sum = a + b;
    return sum;
}

public int sum(int num1, int num2){
    int result;
    result = num1 + num2;
    return result;
}
```

Identical code fragments
except for literals, identifiers, data types,
layout, white space, and comments.

Type 3 สimilar % ไม่ 100 %

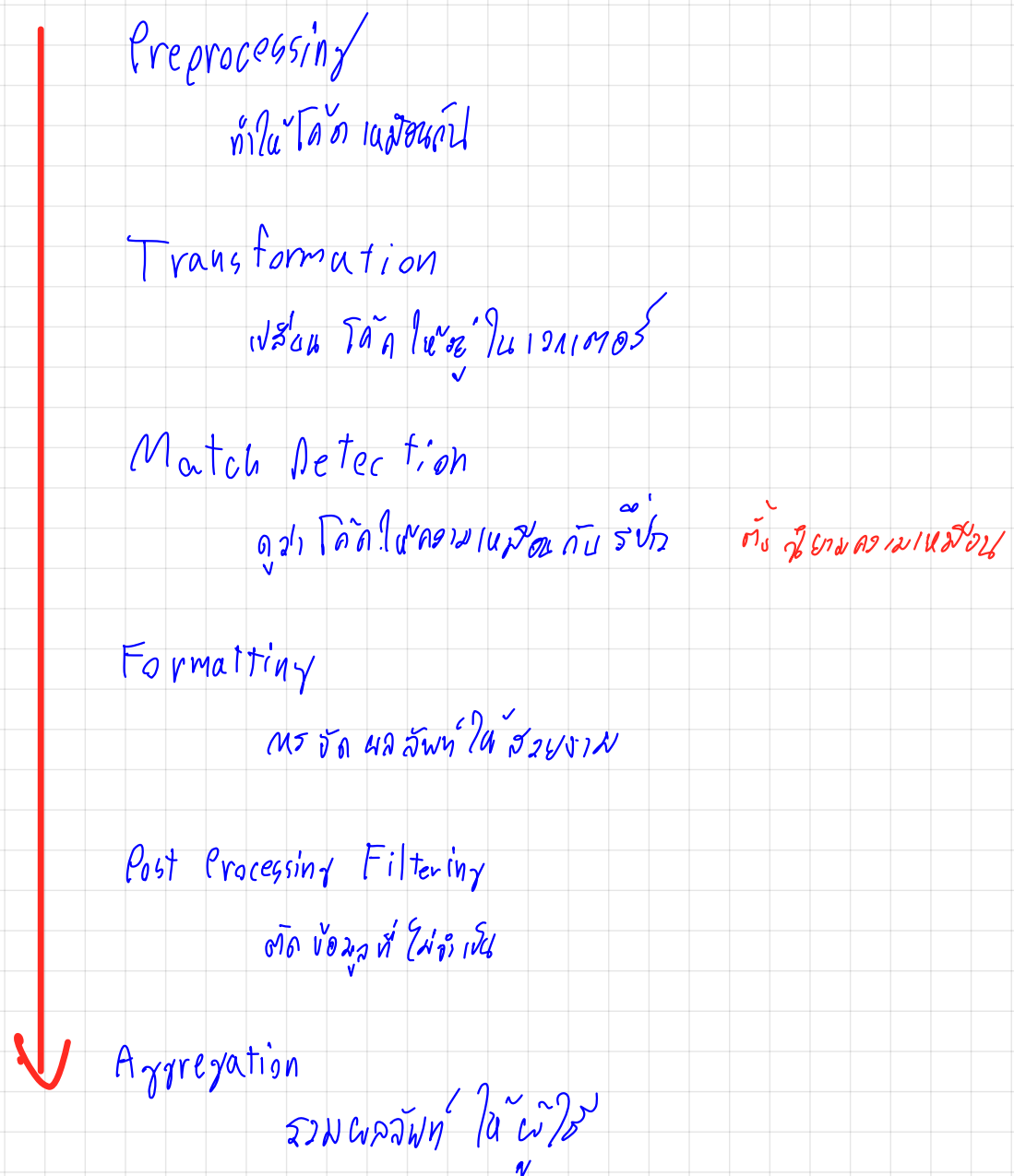
```
public int sum(int a, int b){
    int sum;
    sum = a + b;
    return sum;
}

public int sum(int a, int b){
    return a + b;
}
```

TYPE 4 ได้ผลลัพธ์ที่สอดคล้องกับที่เราได้มา

```
private static String getFormatByName(String name){
    if(name != null){
        final int j = name.lastIndexOf(".") + 1,
        k = name.lastIndexOf("/") + 1;
        if(j > k &&
        j < name.length()){
            return name.substring(j);
        }
    }
    return null
}
```

```
public static String getExtension(final String filename){
    if(filename == null ||
    filename.trim().length == 0 ||
    !filename.contains(".")) {
        return null;
    }
    int pos = filename.lastIndexOf(".");
    return filename.substring(pos+1);
}
```



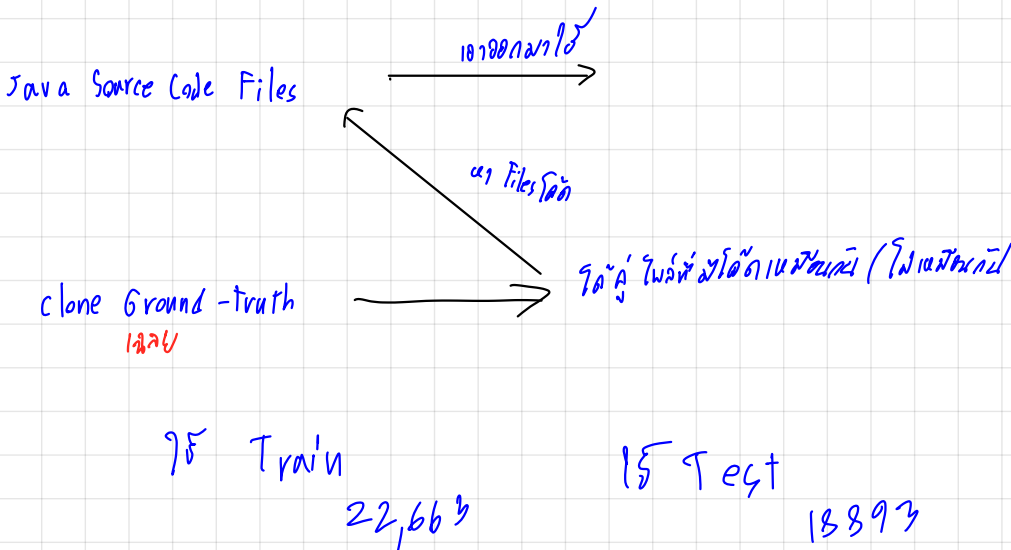
Merry Engine

2. สร้าง Model

Building Merry Engine วิธีสร้าง
ส 3 ขั้นตอน

1. เก็บ Data

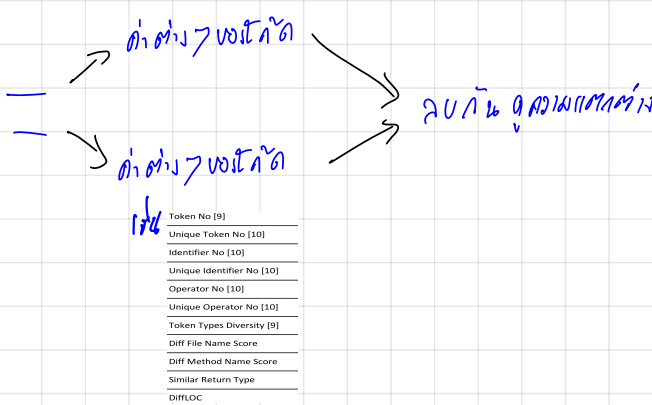
ใช้ข้อมูล BigClone Bench



2. สร้างค่า metrics

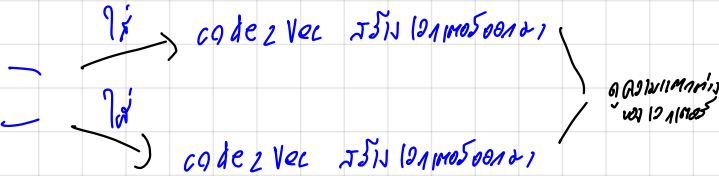
สมมุติ
syntactic metrics

ดูจาก ข้อกำหนด, กฎ, วิธีทำ



semantic metrics

ดูสมรรถภาพของโค้ด
semantic

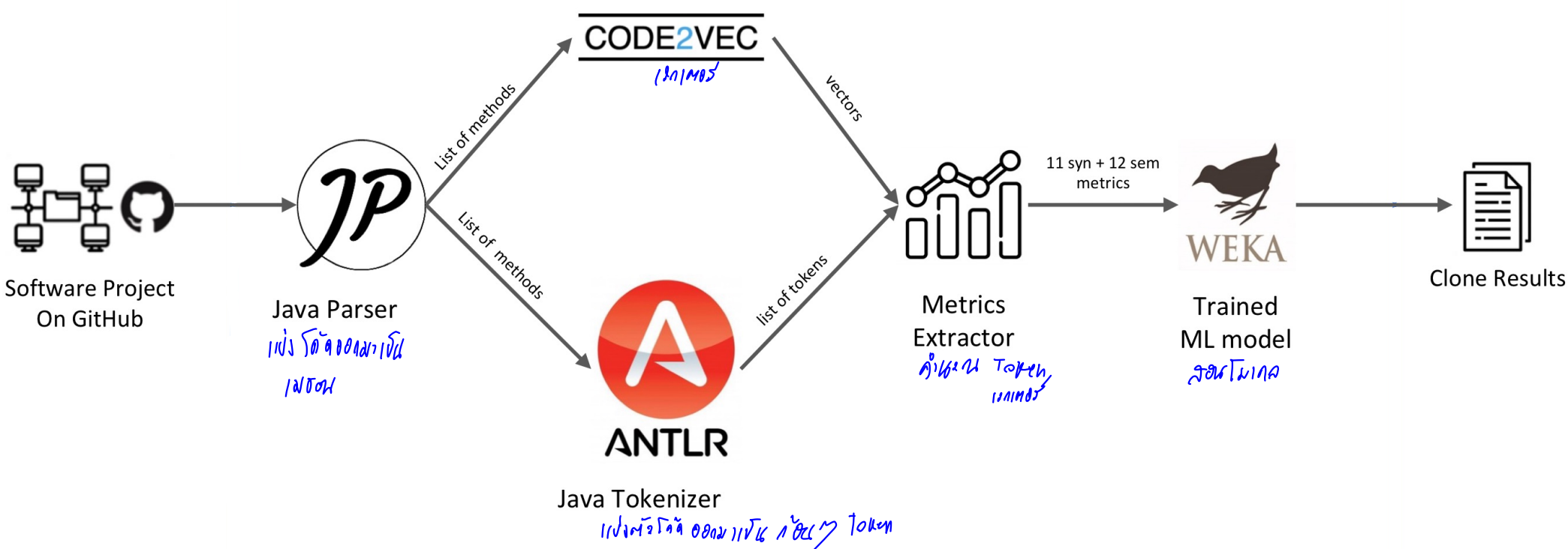


3. เลือกใช้ Model ML
ใช้ 4 แบบ

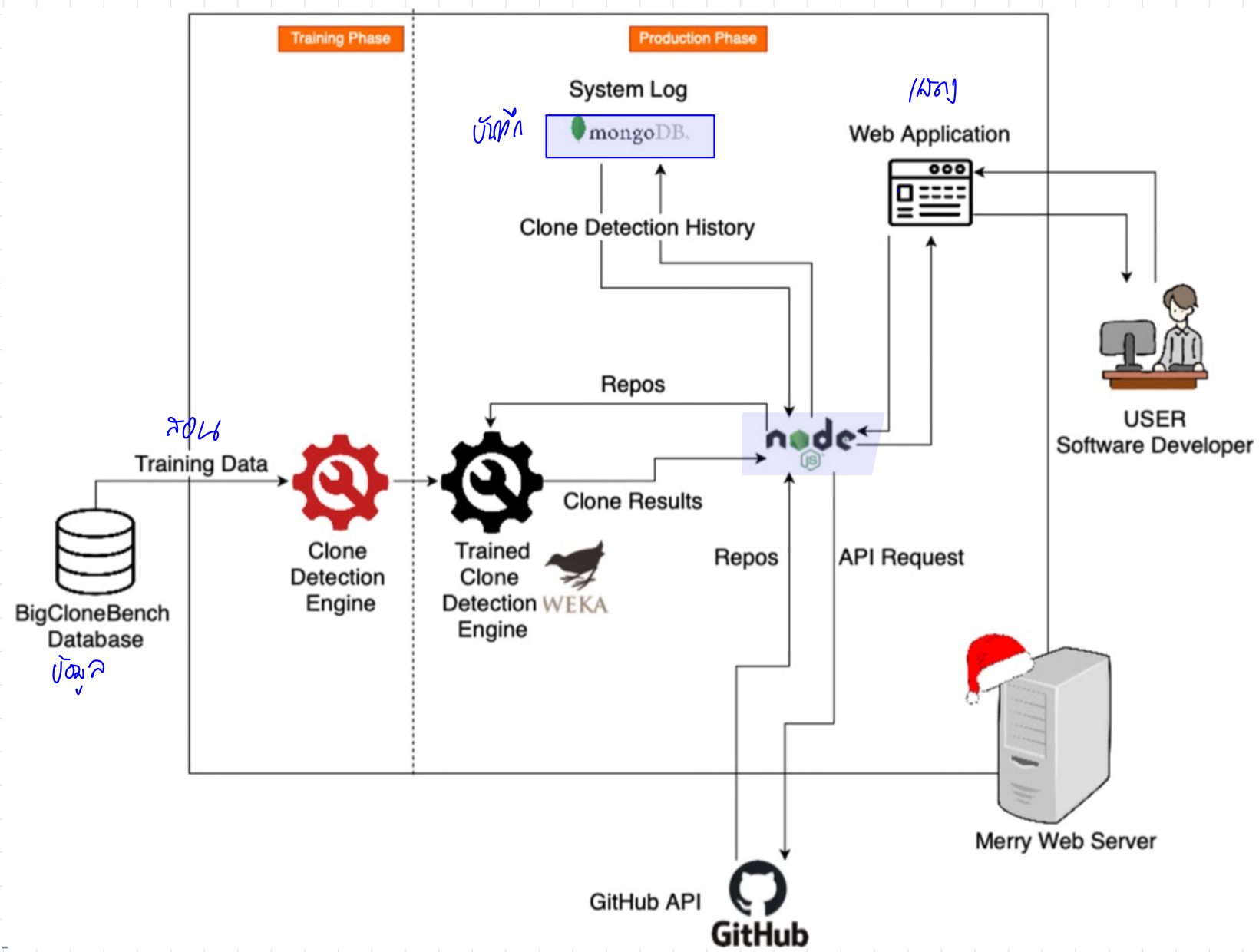
- 1 Decision Tree
- 2 Random Forest
- 3 SVM
- 4 SVM using SMO

3. วิธีใช้ Model

using Merry Engine for Clone Detection 58/55



System Architecture



4. วัดประสิทธิภาพ

1. ดึงตัวเลข

สัปดาห์ 3 ปี

1. $precision = \frac{True\ Positive}{True\ Positive + False\ Positive}$

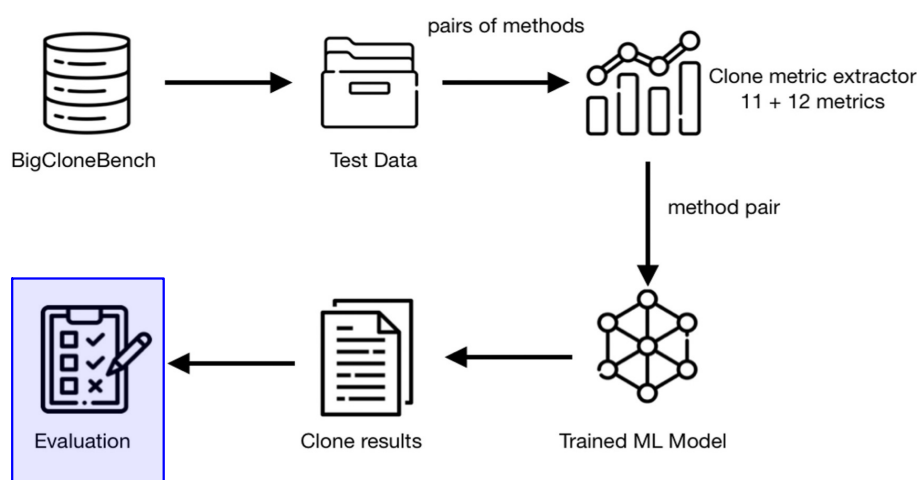
การได้จริง จากข้อมูลทั้งหมด : สัปดาห์ที่มีคนมีรถใช้

2. $Recall = \frac{True\ Positive}{True\ Positive + False\ Negative}$

การได้จริง จาก ข้อมูลทั้งหมดที่ควรจะมีรถใช้ : ถูกจับมาเท่าไร

3. $F1\ Score = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$

ค่ารวมการวัด



WR

Model	Metrics	Precision	Recall	F1-Score
Randomization (baseline)	สุ่ม	0.20	0.49	0.28
Decision Tree	Syntactic + Semantic	0.89	0.86	0.87
	Syntactic	0.95	0.72	0.86
	Semantic	0.68	0.87	0.76
Random Forest	Syntactic + Semantic	0.97	0.86	0.91
	Syntactic	0.97	0.80	0.87
	Semantic	0.70	0.87	0.78
SVM	Syntactic + Semantic	0.97	0.85	0.91
	Syntactic	0.97	0.79	0.87
	Semantic	0.62	0.90	0.73
SVM using SMO	Syntactic + Semantic	0.98	0.89	0.93
	Syntactic	0.97	0.69	0.81
	Semantic	0.63	0.90	0.74

ไม่จริง ๗ ข้อ

ไม่จริงที่สุด