# Catch Me If You Can: Detecting Plagiarism Using Code Similarity Detection Tools

Sang Jun Chun, Claire Luo

# Background

- Motivation: Find the most effective code similarity detection tools for catching plagiarism in school projects.
- Tools (open source) we used:
  - **Plaggie:** source code similarity detector
  - **DeepSim & GEMINI:** binary code similarity detector
  - **SIM:** text difference detector

# Research Questions

RQ1: How effectively can code similarity tools detect **different types of code plagiarism**?

RQ2: How do different types of code similarity tools react to **original works** addressing the same problem?

RQ3: Of the approaches examined, **which** is particularly effective for detecting academic plagiarism in code?
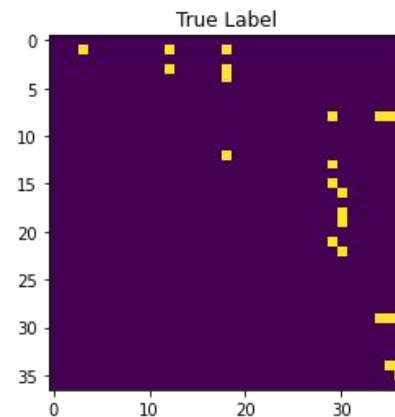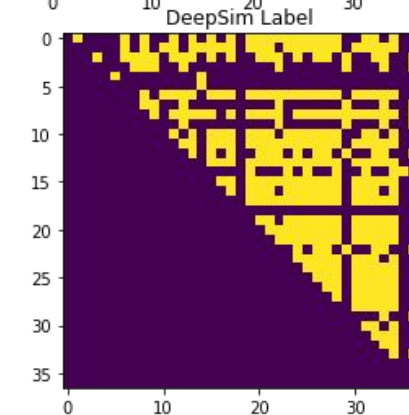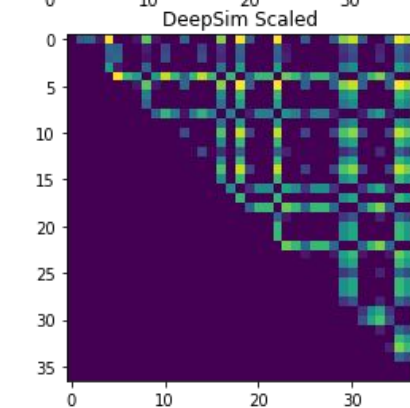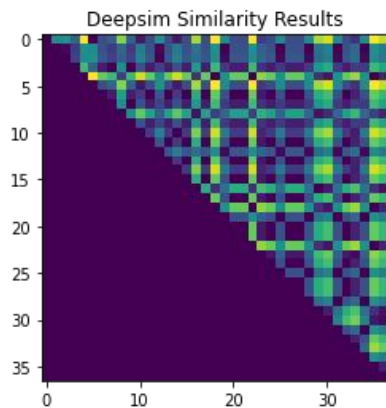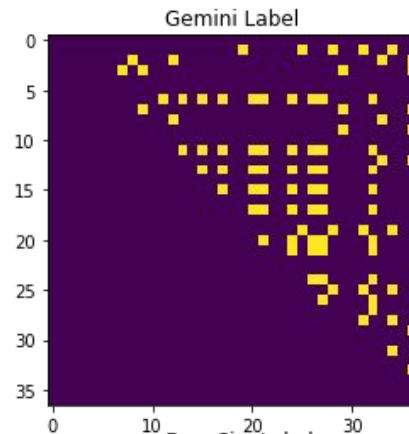
# Methodology (1/2)

- 37 Java homework submissions for COMS W4156 (Fall 2020)
  - Only a subset of the entire pool available as public repositories on GitHub
- How to label pairs?
  - **GEMINI:** Cosine similarity on embedding
  - **DeepSim:** Confidence on classification on semantic embedding
  - **Plaggie:** Generate token-level similarity %, and default threshold set at 50%
  - **SIM:** Text diff similarity %, threshold set by the artificial pairs (59.8%).
- True label
  - Manual detection

[0,1]

[0,1]

%

%

# Methodology (2/2)

- Artificial plagiarism
  - Created plagiarized code using a team member's original submission
  - Three categories:
    - Variables change (name, type)
    - Line change (insert or remove verbose lines)
    - Structure change (rearrange the ordering of classes/methods)

# Results - DeepSim & GEMINI

# Results - Plaggie



Plaggie Similarity Results

Plaggie Label

True Label

# Results - SIM



SIM Similarity Results     SIM Scaled     SIM Label     True Label

# Aggregate Results

| Pairs | DeepSim/GEMINI | Plaggie | SIM | Plagiarized?* |
|-------|----------------|---------|-----|---------------|
| (2, 19) | 0.492; 0.998 | 47.1% | 78% | Likely |
| (9, 36) | 0.568; 0.981 | 60.9% | 63% | Likely |
| (31, 36) | 0.268;0.6924 | 56.2% | 55% | Unlikely |
| (2,29) | 0.360; 0.9056 | 32.6% | 21% | Unlikely |
| (6,37) | 0.870; 0.988 | 40.8% | 22% | Unlikely |
| (8,31) | 0.5947; 0.7517 | 56.2% | 50% | Unlikely |

*Manually inspected.

# An Example of Plagiarism Caught (SIM)

```java
14   class PlayGame {
15
16     private static final int PORT_NUMBER = 8080;
17
18     private static Javalin app;
19
20     private static GameBoard gameboard = new GameBoard();
21
22     // Decide if it is a draw
23     private static boolean isDraw(char[][] boardState) {
24       for (int x = 0; x < boardState.length; x++) {
25         for (int y = 0; y < boardState[x].length; y++) {
26           if (boardState[x][y] == '\u0000') {
27             return false;
28           }
29         }
30       }
31       return true;
32     }
33
34     // Decide if there is a winner of draw
35     private static int getBoardStatus(char[][] boardState, int x, int y) {
36       for (int row = 0; row < 3; row++) {
37         if (boardState[row][0] == boardState[row][1] && boardState[row][1] == boardState[row][2]) {
38           if (boardState[row][0] == 'X') {
39             return x;
40           } else if (boardState[row][0] == 'O') {
41             return y;
42           }
43         }
44       }
45       for (int col = 0; col < 3; col++) {
46         if (boardState[0][col] == boardState[1][col] && boardState[1][col] == boardState[2][col]) {
47           if (boardState[0][col] == 'X') {
48             return x;
49           } else if (boardState[0][col] == 'O') {
50             return y;
51           }
52         }
53       }
```

```java
15   class PlayGame {
16
17     private static final int PORT_NUMBER = 8080;
18
19     private static Javalin app;
20
21     private static GameBoard board;
22
23     //Lets us know if a board shows a draw
24     private static boolean isDraw(char[][] boardState) {
25       for (int x = 0; x < boardState.length; x++) {
26         for (int y = 0; y < boardState[x].length; y++) {
27           if (boardState[x][y] == '\u0000') {
28             return false;
29           }
30         }
31       }
32       return true;
33     }
34
35     //Goes over the rows and then the columns, and finally, diagonals, to determine possible winner
36     // or draw
37     private static int getBoardStatus(char[][] boardState, int x, int y) {
38       for (int row = 0; row < 3; row++) {
39         if (boardState[row][0] == boardState[row][1] && boardState[row][1] == boardState[row][2]) {
40           if (boardState[row][0] == 'X') {
41             return x;
42           } else if (boardState[row][0] == 'O') {
43             return y;
44           }
45         }
46       }
47       for (int col = 0; col < 3; col++) {
48         if (boardState[0][col] == boardState[1][col] && boardState[1][col] == boardState[2][col]) {
49           if (boardState[0][col] == 'X') {
50             return x;
51           } else if (boardState[0][col] == 'O') {
52             return y;
53           }
54         }
55       }
```

# An Example of Plagiarism Caught (DeepSim/Gemini)

```
151   /**
152    * Check if the move is valid or not.
153    * If invalid, a message appears to the players.
154    * @param move player move
155    * @return false if the move is invalid
156    */
157   public boolean isMoveValid(Move move) {
158     char[][] board = getBoardState();
159     if (getTurn() != move.getPlayer().getId()
160         || board[move.getMoveX()][move.getMoveY()] != '0') {
161       return false;
162     }
163     return true;
164   }
165
166   /**
167    * Set a valid move and update overall game status.
168    * Should be used after isMoveValid() method to check if the move is valid.
169    * @param move move
170    * @return true if the game is over
171    */
172   public boolean setMove(Move move) {
173     char[][] board = getBoardState();
174     board[move.getMoveX()][move.getMoveY()] = move.getPlayer().getType();
175     if (getTurn() == 1) {
176       setTurn(2);
177     } else {
178       setTurn(1);
179     }
180     return checkOverallStatus(board, move);
181   }
```

```
149   public int ssmmmm113(Move move) {
150     int three = 3;
151     int four = 4;
152     System.out.println(three + four);
153     char[][] board = ggeeee();
154     board[move.whatsTheNextMoveX()][move.whatsTheNextMoveY()] = move.getHimHere().retrievePlayerType();
155     if (iiiee() == 1) {
156       ssttt(2);
157     } else {
158       ssttt(1);
159     }
160     return ccooss444(board, move);
161   }
162
163   public int makeIt2nd(Player bb) {
164     int three = 3;
165     int four = 4;
166     System.out.println(three + four);
167     this.p2 = bb;
168     return 0;
169   }
```

# Analysis (1/2)

RQ1: How effectively can code similarity tools detect different types of code plagiarism?

- Compare similarity scores for artificially plagiarized code

|        | Lines removed | Variables changed | Structures changed | All changed |
|--------|---------------|-------------------|--------------------|-------------|
| GEMINI | 0.877         | 0.972             | 0.999              | 0.8948      |
| Plaggie | 50.2%        | 77.9%             | 68.8%              | 46.1%       |
| SIM    | 49%           | 79%               | 92%                | 41%         |

# Analysis (2/2)

RQ2: How do different types of code similarity tools react to original works addressing the same problem?
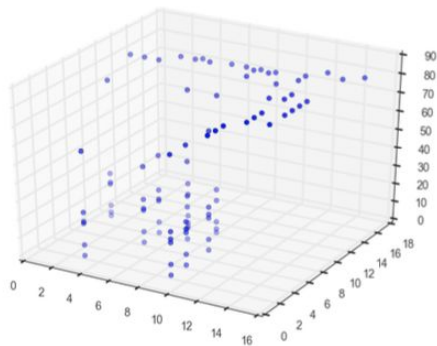
| Tool | Similarity score |
|------|------------------|
| DeepSim/GEMINI | 0.73 / 0.656 |
| Plaggie | 30.5% |
| SIM | 47% |

RQ3: Of the approaches examined, which is particularly effective for detecting academic plagiarism in code?

# Takeaways

- Binary code similarity detectors performed well, but produced many false positives

    - Possible reasons: skeleton code, lack of data points, modified code

- Plaggie (source code similarity detector) is only good at catching obvious copies

- Adding or removing redundant lines from the original code hinders the tools the most

```
00.55(1555)-55(1555).IF_END.}
61:62(1623)-62(1649):IF:if (app=="/move/:playerId")
62:62(1651)-62(1651):BLOCK:{
63:63(1663)-63(1716):VARIABLE_DECLARATION:playerId = Integer.parseInt(ctx.pathParam("playerId"))
64:63(1674)-63(1716):METHOD_INVOCATION:Integer.parseInt(ctx.pathParam("playerId"))
65:63(1691)-63(1715):METHOD_INVOCATION:ctx.pathParam("playerId")
66:64(1732)-64(1789):VARIABLE_DECLARATION:curPlayer = boardController.getBoard().getPlayer(playerId
67:64(1744)-64(1769):METHOD_INVOCATION:boardController.getBoard()
68:64(1744)-64(1789):METHOD_INVOCATION:boardController.getBoard().getPlayer(playerId)
69:65(1822)-65(1859):VARIABLE_DECLARATION:paraMap = getParametersMap(ctx.body())
70:65(1832)-65(1859):METHOD_INVOCATION:getParametersMap(ctx.body())
71:65(1849)-65(1858):METHOD_INVOCATION:ctx.body()
72:66(1872)-66(1909):VARIABLE_DECLARATION:x = Integer.parseInt(paraMap.get("x"))
73:66(1876)-66(1909):METHOD_INVOCATION:Integer.parseInt(paraMap.get("x"))
74:66(1893)-66(1908):METHOD_INVOCATION:paraMap.get("x")
75:67(1922)-67(1959):VARIABLE_DECLARATION:y = Integer.parseInt(paraMap.get("y"))
76:67(1926)-67(1959):METHOD_INVOCATION:Integer.parseInt(paraMap.get("y"))
77:67(1943)-67(1958):METHOD_INVOCATION:paraMap.get("y")
78:68(1973)-68(2008):VARIABLE_DECLARATION:newMove =  new Move(curPlayer, x, y)
79:68(1984)-68(2008):NEW:new Move(curPlayer, x, y)
80:69(2025)-69(2068):VARIABLE_DECLARATION:message = boardController.validMove(newMove)
81:69(2035)-69(2068):METHOD_INVOCATION:boardController.validMove(newMove)
```

```
90:75(1600)-75(1600):IF_END:}
91:77(1611)-77(1637):IF:if (app=="/move/:playerId")
92:77(1639)-77(1639):BLOCK:{
93:78(1651)-78(1704):VARIABLE_DECLARATION:playerId = Integer.parseInt(ctx.pathParam("playerId"
94:78(1662)-78(1704):METHOD_INVOCATION:Integer.parseInt(ctx.pathParam("playerId"))
95:78(1679)-78(1703):METHOD_INVOCATION:ctx.pathParam("playerId")
96:79(1717)-79(1756):VARIABLE_DECLARATION:x = Integer.parseInt(ctx.formParam("x"))
97:79(1721)-79(1756):METHOD_INVOCATION:Integer.parseInt(ctx.formParam("x"))
98:79(1738)-79(1755):METHOD_INVOCATION:ctx.formParam("x")
99:80(1769)-80(1808):VARIABLE_DECLARATION:y = Integer.parseInt(ctx.formParam("y"))
100:80(1773)-80(1808):METHOD_INVOCATION:Integer.parseInt(ctx.formParam("y"))
101:80(1790)-80(1807):METHOD_INVOCATION:ctx.formParam("y")
102:82(1825)-82(1837):VARIABLE_DECLARATION:player = null
103:83(1846)-83(1888):IF:if (playerId == board.getPlayer1().getId())
104:83(1862)-83(1879):METHOD_INVOCATION:board.getPlayer1()
105:83(1862)-83(1887):METHOD_INVOCATION:board.getPlayer1().getId()
106:83(1890)-83(1890):BLOCK:{
107:84(1900)-84(1926):ASSIGNMENT:player = board.getPlayer1()
108:84(1909)-84(1926):METHOD_INVOCATION:board.getPlayer1()
109:85(1935)-85(1935):BLOCK_END:}
110:85(1935)-85(1935):IF_END:}
111:85(1937)-85(1940):ELSE:else
```