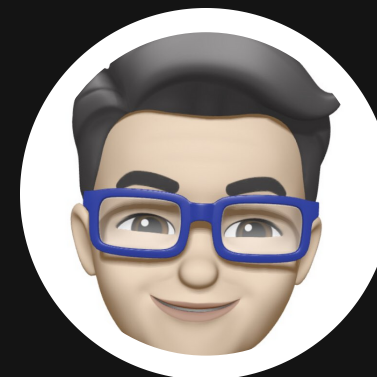
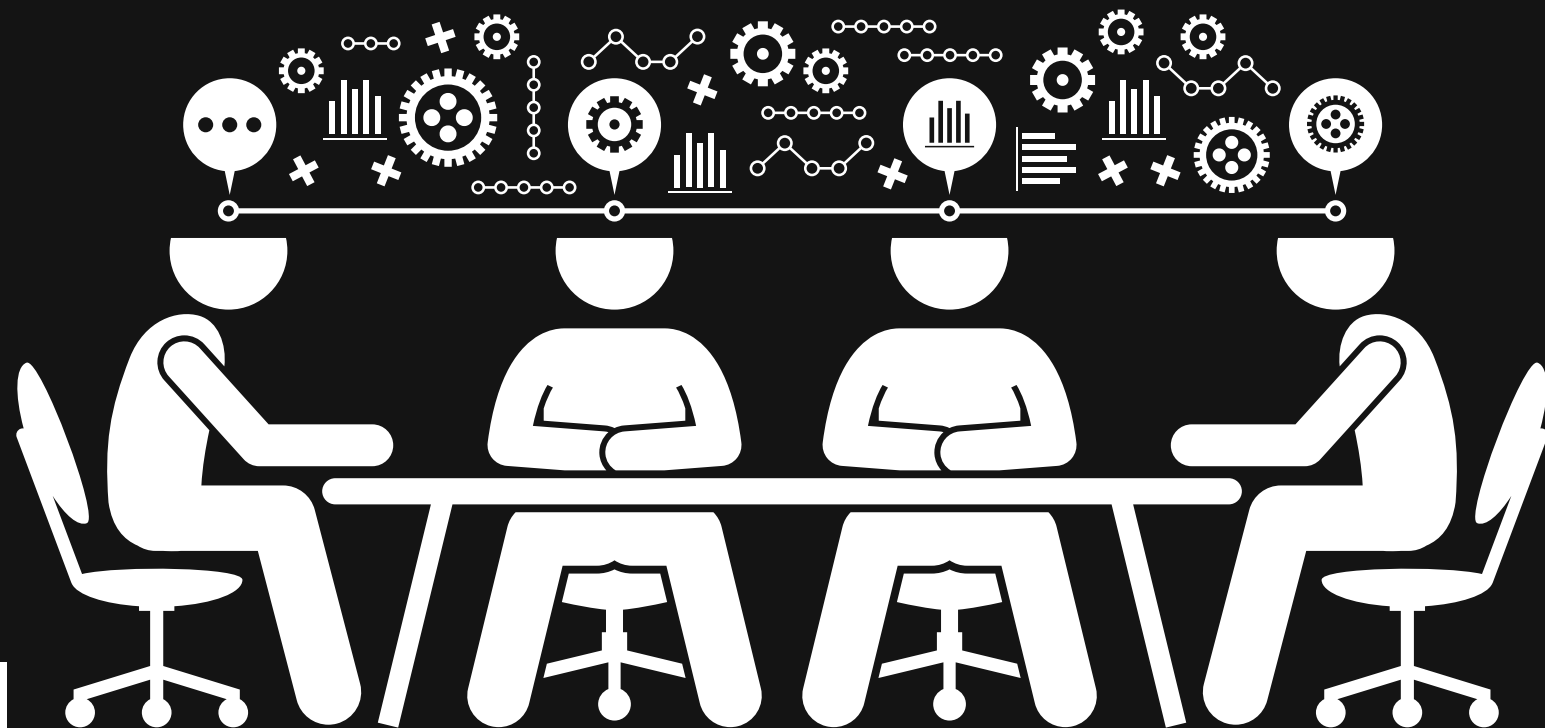


# COPY CAT

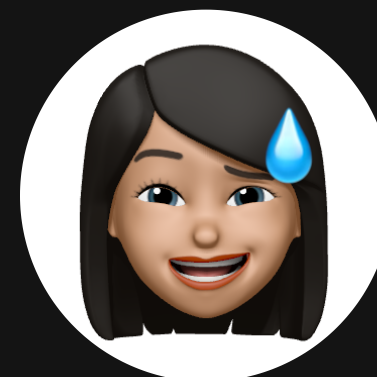
ARCHIT MANEK, SHEBNA MATHEW, DEVINA RAITHATHA &  
NIKHITA SINGH

CS: 5500 - FOUNDATIONS OF SOFTWARE ENGINEERING

# THE TEAM



ARCHIT MANEK



SHEBNA MATHEW



DEVINA RAITHATHA



NIKHITA SINGH



# AGENDA

03



HIGH-LEVEL ARCHITECTURE

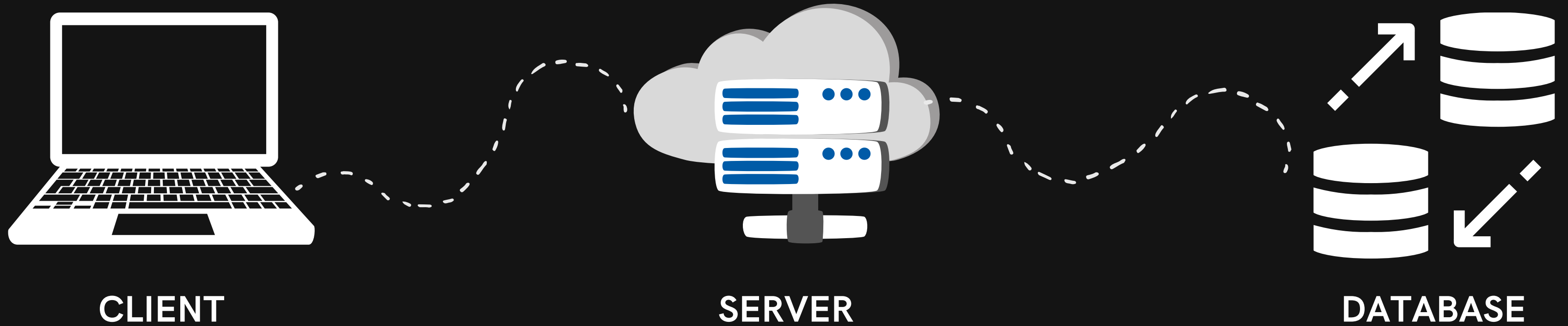
PLAGIARISM DETECTION ALGORITHM

INTERACTIVE PRESENTATION OF TOOL

DESIGN EVOLUTION

SOFTWARE DEVELOPMENT PRINCIPLES

COPYCAT



# HIGH LEVEL ARCHITECTURE



**PARSING AND TRAVERSING THE AST:  
BABEL**



**FOUR STRING REPRESENTATIONS OF  
THE SOURCE CODE**



**STRING - SIMILARITY COMPARISON**

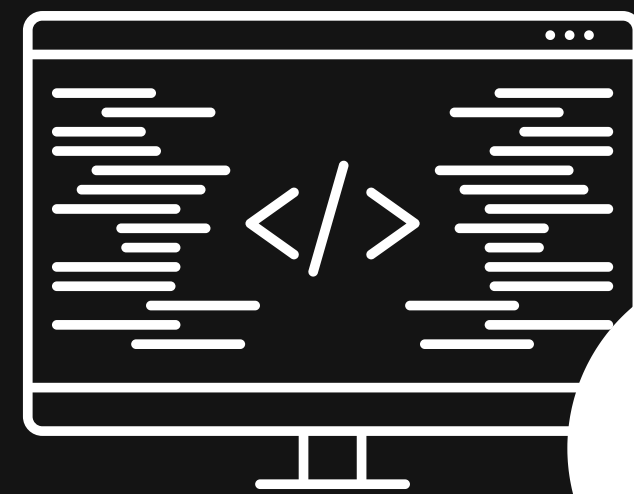


**TOTAL PLAGIARISM %**

# **PLAGIARISM DETECTION ALGORITHM**

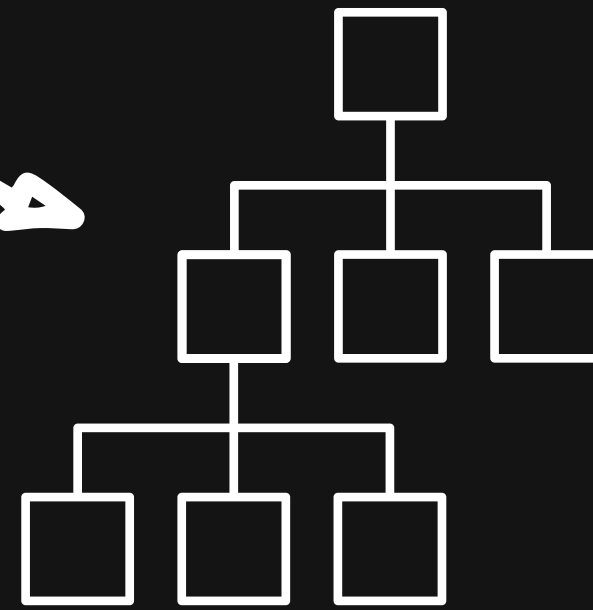
# PRE-PROCESSING & 4 STRINGS

PARSING: BUILDING AN AST



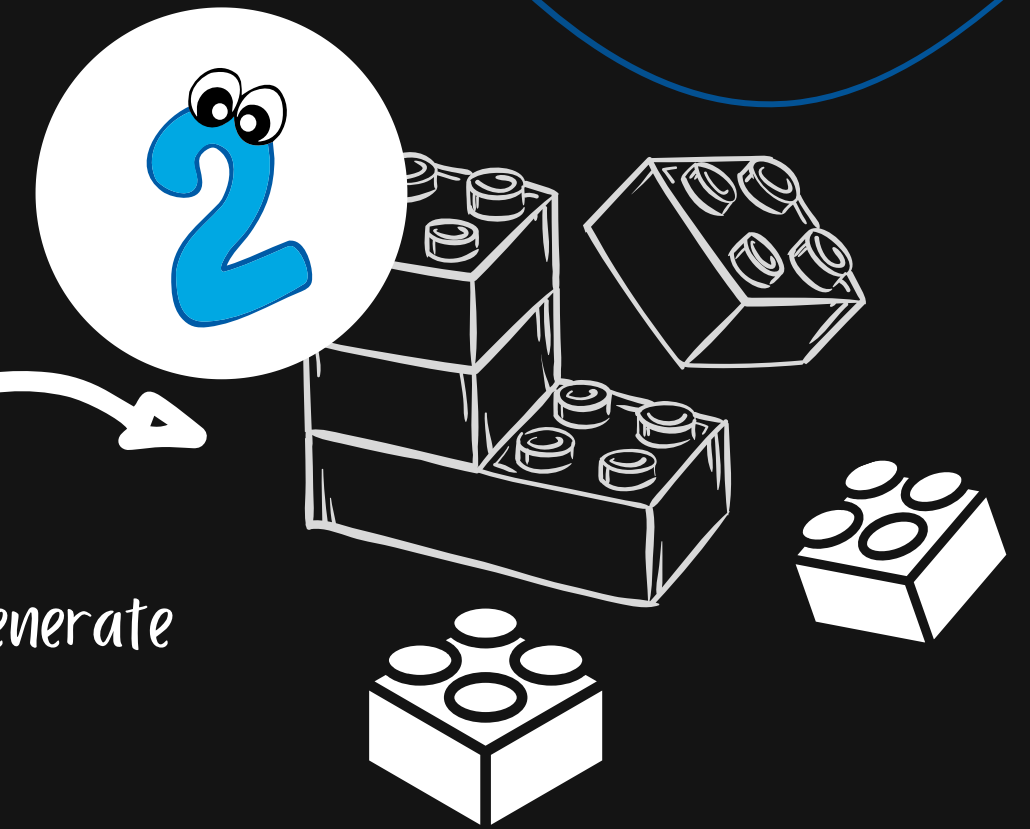
STUDENT FILE:  
SOURCE CODE

parse



ABSTRACT  
SYNTAX TREE

traverse & generate



STRING  
REPRESENTATIONS



## FILE STRUCTURE:

"ProgramVariableDeclarationExpressionStatement"

## IDENTIFIERS:

1f40fc92da241694750979ee6cf582f2d5d7d28e18335de05abc54d0560e0f5302860c652bf08d560252aa5e74210546f369fbbbce8c12cfc7957b2652fe9a75eed55db3ffa1983455e1c1b291f6d4d48009bbf43338657ac7a49631126140f126f0a95456b60535e3a7902acd712a62044b445972a8b6c1e79c7cbbb80bc01f873bcc37e512b7da86476367769c932009fc1be59c929879f10ac89df541124db5010ae7297ec71db8f71a09adccd27c002f00fcfc93ca7e157105e7505f24d"

## LITERALS:

"10HelloWorld"

## TOKENS:

ProgramVariableDeclarationVariableDeclaratorIdentifierNumericLiteralExpressionStatementCallExpressionMemberExpressionIdentifierIdentifierStringLiteral



# STRING SIMILARITY

ONCE WE HAVE THE FOUR STRINGS FOR EACH SOURCE CODE, WE GET THE DEGREE OF SIMILARITY BETWEEN TWO STRINGS.

FOR THIS WE USE A STRING COMPARE BASED ON THE DICE SIMILARITY COEFFICIENT WHICH IS A SIMPLE AND EFFECTIVE WAY TO CALCULATE A MEASURE OF THE SIMILARITY OF TWO STRINGS. THE OUTCOME IS A VALUE BOUNDED BETWEEN ZERO AND ONE THAT WE USE TO CALCULATE SIMILARITY PERCENTAGE.





# STRING SIMILARITY %

HOW SIMILAR ARE THESE TWO STRINGS?

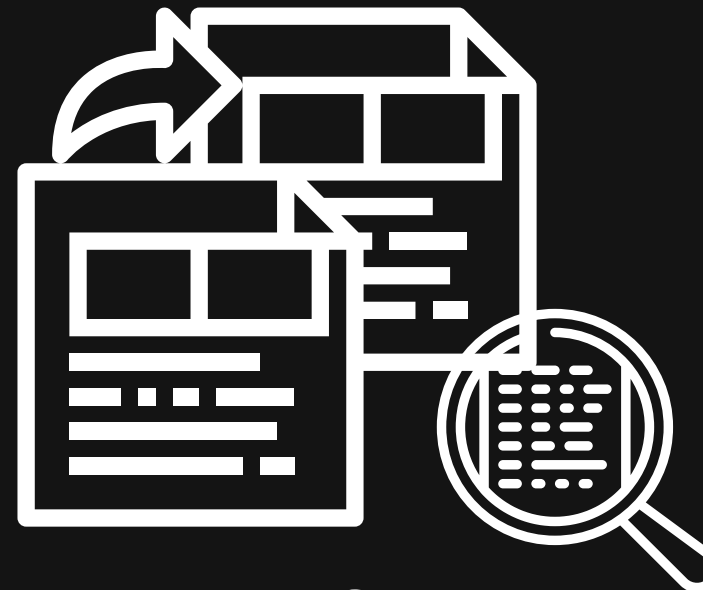
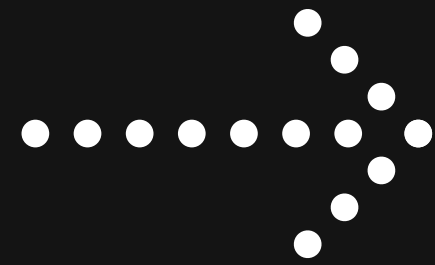
$$\frac{S_f + S_i + 2 + S_l + 10 \cdot S_t}{14}$$

REPRESENTS THE TOTAL SIMILARITY AS A  
WEIGHTED AVERAGE OF THE FOUR  
SIMILARITIES: SF (FILE STRUCTURE), SI  
(IDENTIFIER), SL (LITERALS) AND ST (TOKENS).

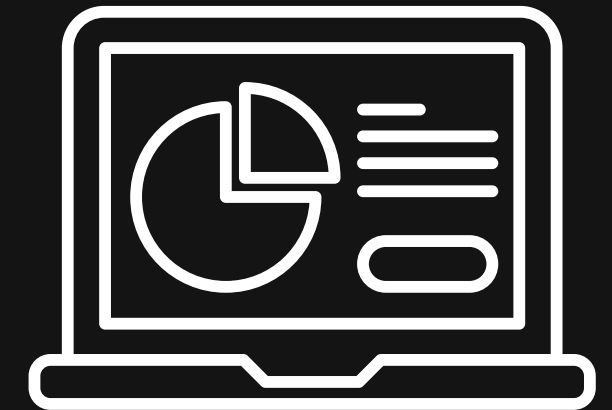
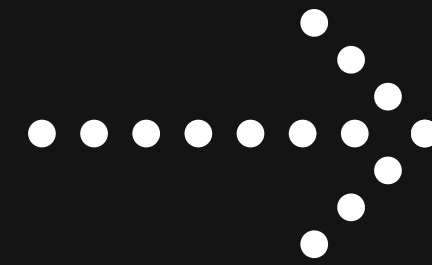




STUDENT CODE  
FILE UPLOAD

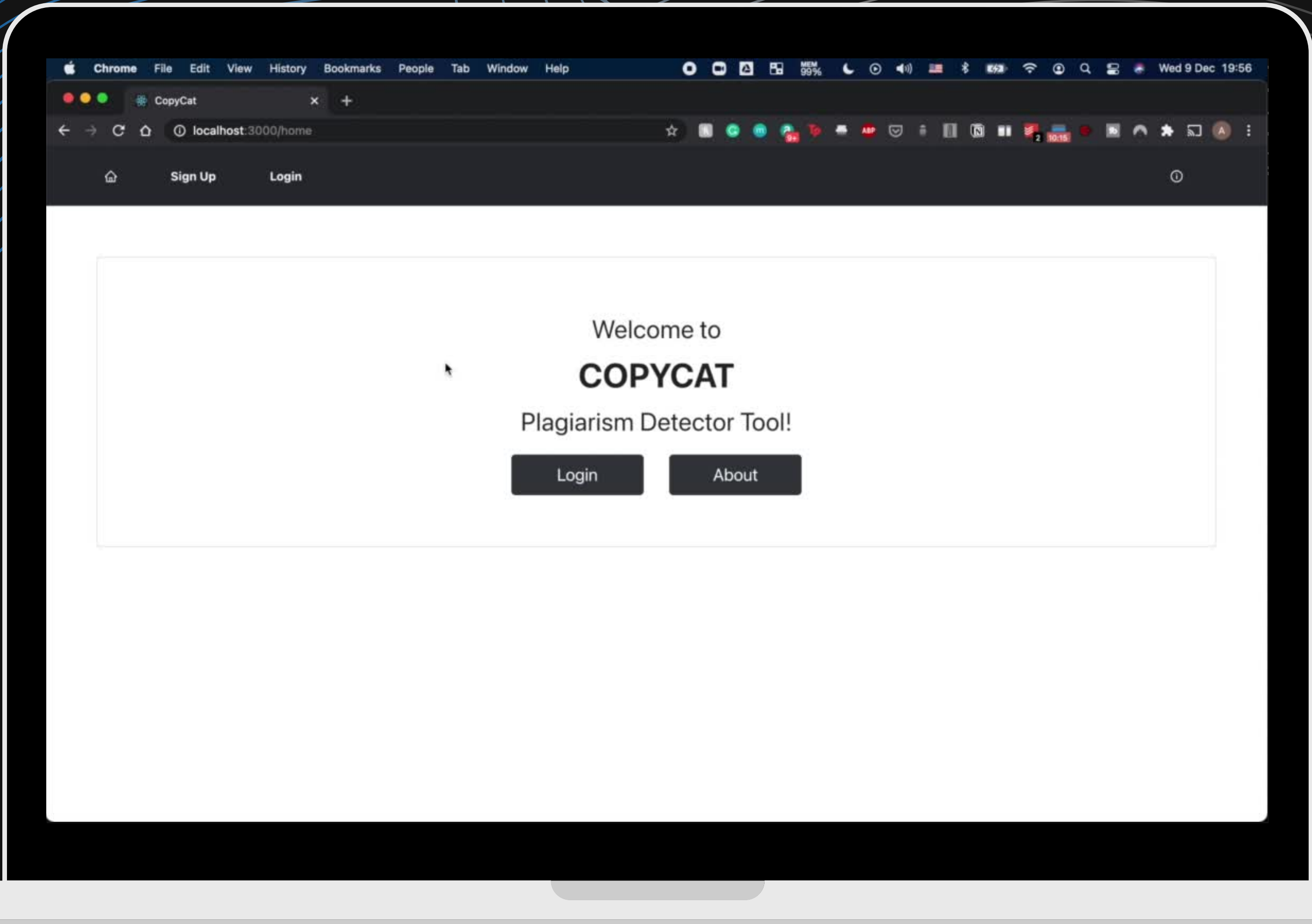


DETECT  
PLAGIARISM



PLAGIARISM  
PERCENTAGE

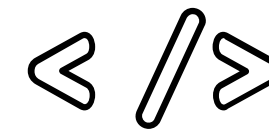
# HOW TO USE OUR TOOL



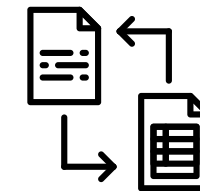
# DESIGN EVOLUTION



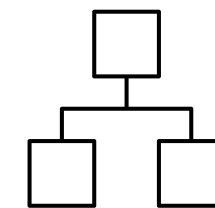
WHAT PROGRAMMING LANGUAGES ARE SUPPORTED?



WHAT LIBRARIES AND NPM PACKAGES CAN WE USE?



HOW DO WE TRANSPOSE SOURCE CODE TO AN AST?



# SOFTWARE DESIGN PRINCIPLES



**WATERFALL  
METHODOLOGY**



**FACTORY  
PATTERN**



**VERSION  
CONTROL**

**FUTURE**



# SCOPE



**MULTIPLE STUDENT  
COMPARISONS**



**MULTIPLE LANGUAGE  
SELECTION**



**SUPPORT  
COMPARISONS WITH  
ONLINE CODE**



**GENERATE &  
DOWNLOAD  
PLAGIARISM REPORT**





# THE END

Link to Slides: <https://bit.ly/fse-slides>