

Syntax Sentinels

Team 2

An NLP-based code plagiarism
detector

Presentation Roadmap

- **Meet the Team** – Introduction to team members.
- **Existing Solutions: MOSS** – Strengths and limitations.
- **Why It Matters** – Importance of effective plagiarism detection.
- **Problem Statement** – Challenges with current solutions and our motivation.
- **Demo** – Live demonstration of the system.
- **Features** – Overview of our NLP-based plagiarism detection system.
- **System Architecture** – Technical breakdown of how our system works.
- **Model Design** – The key components of our model
- **Q&A** – Open floor for questions.

Meet the Team



**Dennis
Fong**



**Lucas
Chen**



**Mohammad
Mohsin
Khan**



**Luigi
Quattrociochi**



**Julian
Cecchini**

PICTURE THIS!

Existing Solution: MOSS (Measure of Software Similarity)

- **Strengths:**
 - Effective for **exact and near-exact code matches**.
 - Widely used in **academia** and **coding competitions**.
 - Works across multiple programming languages.
- **Limitations:**
 - **Syntax-based approach** – Cannot detect **semantic plagiarism** (e.g., logic-preserving transformations).
 - Struggles with **function reordering, and code restructuring**.
 - **High false negatives** – Clever modifications can evade detection.



Why It Matters

- **Academic Integrity** – Ensures fairness in grading and protects originality.
- **Code Competitions** – Detects unfair advantages in submissions.
- **User Trust** – Minimizes false positives/negatives for reliable detection.



Problem statement

Academic institutions, coding competition organizers, and online platforms need a more effective plagiarism detection system because current tools like MOSS rely on syntax-based comparisons. Our NLP-based solution detects semantic plagiarism, providing a more accurate, customizable, and scalable system.



Demo

Features

Comprehensive Detection via an NLP approach

Semantic Understanding

- Unlike traditional tools that rely on **syntactic** matching (such as MOSS), our system uses **Natural Language Processing (NLP)** to understand the **semantic** of the code.

Why is this better

- This makes our product more effective at catching **advanced plagiarism** techniques like **code obfuscation, inlining, block reordering and variable renaming**.
- All in all, this improves detection of **logic-preserving plagiarism** where code structure is changed but the functionality remains the same.



Detailed Reports

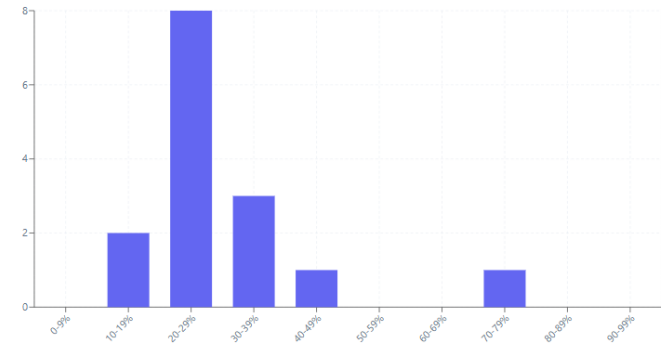
- Provides comprehensive reports which includes:
 - Similarity scores as a histogram
 - Line by line analysis of file pairs
- Helps users understand plagiarism insights clearly.

Highest Similarity
78.261%

Average Similarity
31%
Median: 26.087%

Total Submissions
6

Similarity Distribution (pairs)



Similarity Threshold: 50%

50

```
1 def scoring(d, t, k=10):
2     arr = [1 for l in last]
3     arr[t] = d
4     penalty = sum((d - 1) * c for l, c in zip(arr, C))
5     score = S[d][t] - penalty
6     score -= (min(D, d + k) - d - 1) * penalty
7     return score
8
9 D = int(input())
10 C = list(map(int, input().split()))
11 S = []
12 for i in range(D):
13     s = list(map(int, input().split()))
14     S.append(s)
15 last = [-1] * 26
16 ans = []
17 for i in range(D):
18     res = -10 ** 6
19     idx = -1
20     for j in range(26):
21         tmp = scoring(i, j)
22         if tmp > res:
23             res = tmp
24             idx = j
25     last[idx] = i
26     ans.append(idx+1)
27 print(*ans, sep='\n')
```

```
1 import time
2 import copy
3
4 start_time = time.time()
5
6 #input
7 D = int(input())
8 c_list = list(map(int, input().split()))
9
10 s_grid = []
11 for i in range(D):
12     array = list(map(int, input().strip().split(' ')))
13     s_grid.append(array)
14
15 def calculate_score(d, t, last):
16     score = s_grid[d][t]
17     last[t] = -1
18     for i in range(26):
19         score -= c_list[i] * (last[i] + 1)
20     return score
21
22 t_list = [] #task_list
23 last_list = [0] * 26
24
25 for k in range(0, D):
26     X = -1 # k-日目に変わる番号を探す
27     p = 0
28     for i in range(26): # 26通り試す
29         tmp = calculate_score(k, i, last_list)
30         if tmp > p:
31             X = i
32             p = tmp
```

Why these features matter?

Streamlines the user experience of educators in maintaining academic integrity.

Improves precision and reduces false positives with NLP.

Visualizations provide a high-level overview of plagiarism within a corpus.

Changes to requirements documentation

Zero Data Retention (ZDR)

- This was limiting from a user perspective.
- Eliminating this requirement helped to streamline the user experience.

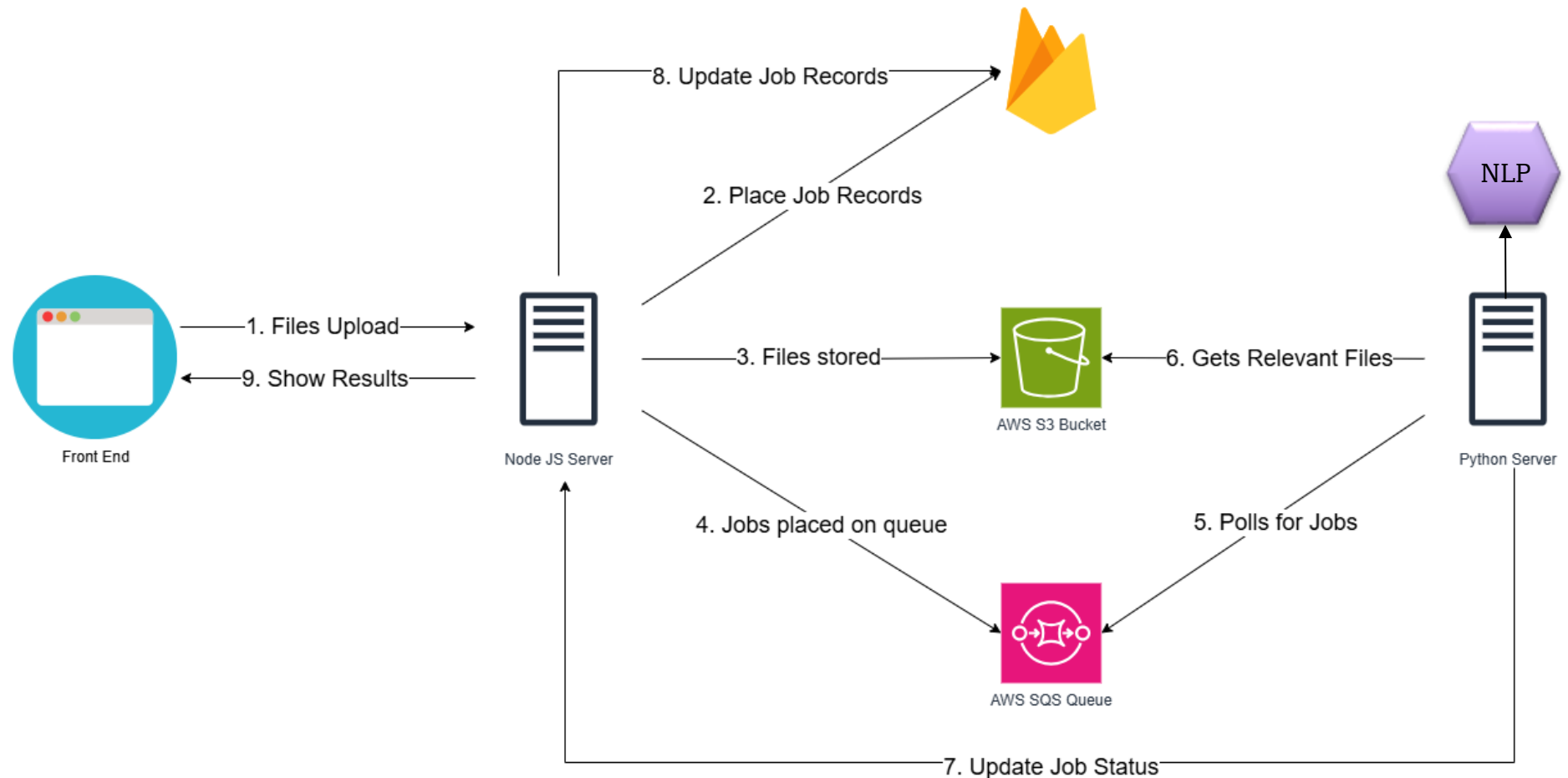
Emailing

- The user flow for our revision 0 (receiving an email with your results) was not intuitive.
- Instead of sending the results in a zip file, we have a button to download the results.



Architecture

System Architecture



Firestore - Firestore

What is Firestore?

- Firestore is a **NoSQL database**
- Firestore lives on the **cloud**, and is built on Google's infrastructure

Why Firestore?

- We wanted to utilize a NoSQL database due to the unstructured nature of the data contained in a job
- **Real Time Sync** – Connected clients receive updated data in real time
 - Other NoSQL databases such as DynamoDB or MongoDB require extra set up and cost to make this possible
- **No Management Required** – Firestore is fully managed by Firebase. Scaling read/write capacities are done automatically
 - No need to manage read/write capacities (in DynamoDB), or manage clusters (in MongoDB)



Simple Storage Service (S3)

What is S3?

- Low cost, high-capacity object storage
- Enables fast GET operation on objects

Why S3?

- Zero cost for use cases this product was designed for
- Can be scaled up to store much more data, while still costing almost nothing
 - To spend 1 cent, you'd need to upload 10,000 files approximately 43.5 KB each
- **Extremely Durable** – objects in S3 are stored across multiple availability zones, meaning data can be recovered even through disaster
- **Highly Secure** – Data is encrypted, and secured by permissions set by admins



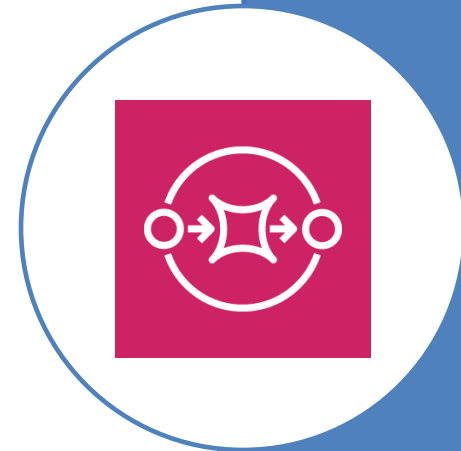
Simple Queue Service (SQS)

What is SQS?

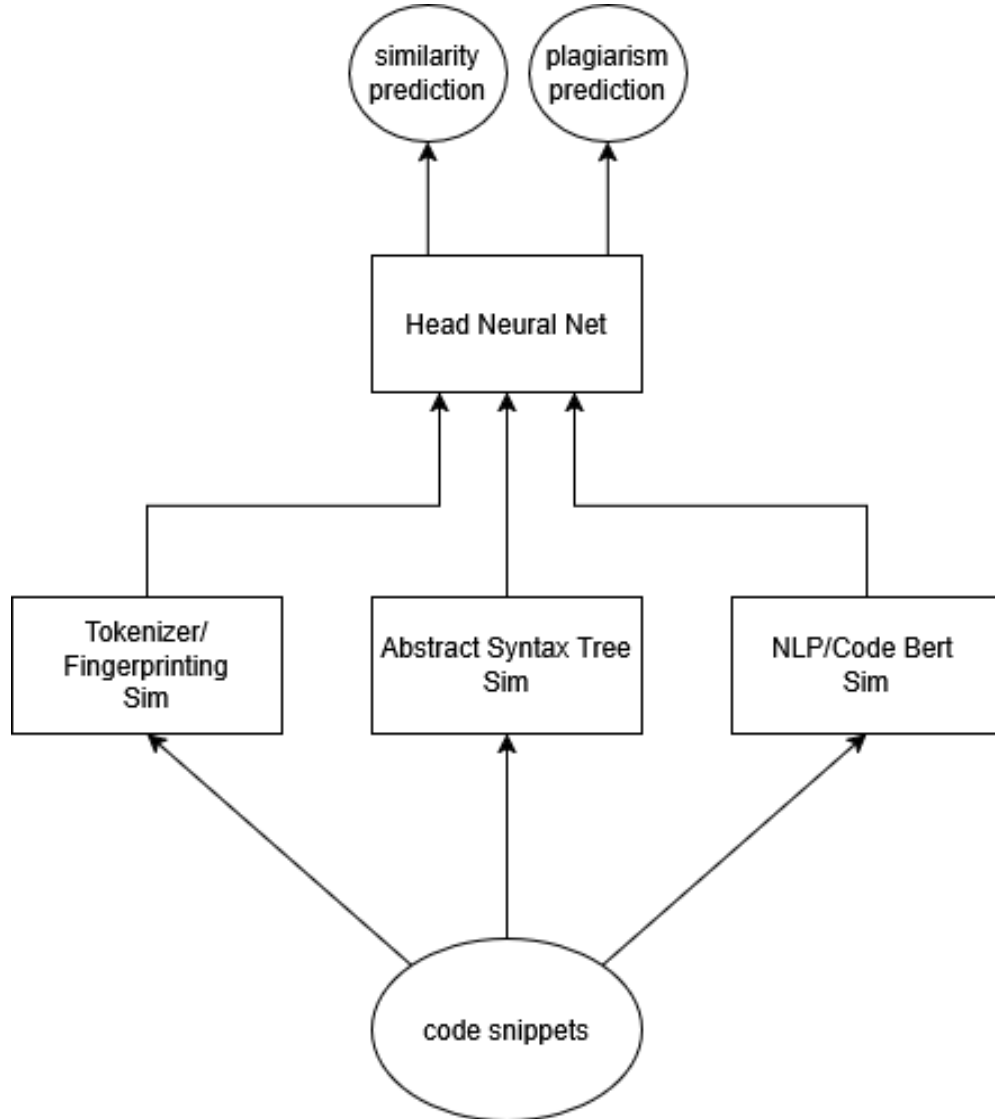
- SQS is a **message queuing service**
- Allows **decoupled applications to communicate** with each other without interfering with each other
 - o No more "Oh no, I broke the endpoint"

Why SQS?

- **NO** cost – AWS allows you to send 1,000,000 requests in SQS every month, free of charge
 - o Therefore, it is essentially **FREE** (unless you check for plagiarism more than 1,000,000 times a month)
- **Very Reliable** – Provides "at least once" or "exactly once" delivery, ensuring messages are always received
- **No Management Required** – SQS is fully managed by AWS. No need to worry about set up, or scalability. It is all done by AWS



Model Design (Diagram)



Model Design (Diagram)

Key Components

- Tokenization Sim Module (MOSS)
- AST Sim Module
- NLP Module



Why include all?

- Best of all worlds. Majority of the lifting is done by the NLP module, but additional pointers/clues that the NLP module may not catch can be provided by other similarity methods without significant computational overhead.

Results

- Outperforms standard algorithm on cases where obvious signs of plagiarism are present.

What are the signs?

Token occurrence and frequency

Structure/ordering

Logical equivalences

Example

plag2.py

[VIEW SUBMISSION >](#)

```
1  import math
2  j = 1
3  for i in range (j):
4
5      count = int(input("Enter number of elements (0 to exit): "))
6      if count == 0:
7          exit()
8
9      values = list(map(int, input("Enter the elements separated by space: ").split()))
10     average = sum(values) / len(values)
11     variance_sum = sum((item - average) ** 2 for item in values)
12     std_deviation = math.sqrt(variance_sum / count)
13
14     print(f'{std_deviation:.8f}')
15     j+=1
16
```

Plagiarized file

s129271906.py

[VIEW SUBMISSION >](#)

```
1  from math import sqrt
2  while 1:
3      n = int(input())
4      if n == 0:
5          break
6      a = list(map(int, input().split()))
7      av = sum(a)/len(a)
8      Sum = sum((x-av)**2 for x in a)
9      print(f'{sqrt(Sum/n):.8f}')
10
```

Original file

What does MOSS say?

30% similarity with a batch mean of 64%

Above threshold of plagiarism provided by 7%

Overall, weak indication that plagiarism occurred

What does our model say?

Similarity of 68% with batch mean of 55%

Above threshold of plagiarism provided by 30%

Stronger indication that plagiarism has occurred

Thanks for listening!

Q&A

Another example from dataset

```
1 import math
2 j = 1
3 for i in range(j):
4
5     count = int(input("Enter number of elements (0 to exit): "))
6     if count == 0:
7         exit()
8
9     values = list(map(int, input("Enter the elements separated by space: ").split()))
10    average = sum(values) / len(values)
11    variance_sum = sum((item - average) ** 2 for item in values)
12    std_deviation = math.sqrt(variance_sum / count)
13
14    print(f'{std_deviation:.8f}')
15    j+=1
16
17
```

```
1 import math
2 while 1:
3     n = int(input())
4     if n == 0:
5         break
6     s = list(map(int, input().split()))
7     m = sum(s) / len(s)
8     x = 0
9     for i in range(n):
10         x += (s[i] - m) ** 2 / n
11     a = math.sqrt(x)
12     print(a)
```