

VIETNAM NATIONAL UNIVERSITY HO CHI MINH CITY  
HO CHI MINH CITY UNIVERSITY OF TECHNOLOGY  
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



Cấu trúc dữ liệu và Giải thuật

---

Assignment Report

PHÁT HIỆN VĂN BẢN  
TRÙNG LẶP HOẶC GẦN  
GIỐNG DỰA TRÊN HỌC  
SÂU VÀ BẮM ĐẶC TRƯNG

---

Advisor(s): Lê Thành Sách

Student(s): Trần Thanh Hưng Thịnh 2413345  
Lê Ân Tri 2413597

HO CHI MINH CITY, DECEMBER 2025



## Contents

<b>1</b>	<b>Mục Tiêu Đề Tài</b>	<b>4</b>
1.1	Yêu Cầu Chính Đã Hoàn Thành . . . . .	4
<b>2</b>	<b>Cơ Sở Lý Thuyết và Phương Pháp</b>	<b>4</b>
2.1	Trích Xuất Đặc Trưng (Document Embedding) . . . . .	4
2.2	Các Kỹ Thuật Phát Hiện Trùng Lặp . . . . .	5
2.2.1	FAISS (Approximate Nearest Neighbor) . . . . .	5
2.2.2	SimHash (Băm Nhị Phân) . . . . .	5
2.2.3	MinHash (Băm Tập Hợp) . . . . .	5
2.3	Phân Cụm (Clustering) và Chọn Đại Diện . . . . .	6
<b>3</b>	<b>Cấu Trúc Hệ Thống (Mã Nguồn)</b>	<b>6</b>
<b>4</b>	<b>Kết Quả và Đánh Giá So Sánh</b>	<b>8</b>
4.1	Kết Quả Định Tính (Phân tích phương pháp) . . . . .	8
4.2	Kết Quả Định Lượng (Hiệu năng) . . . . .	8
4.3	So Sánh Chi Tiết về Độ Chính Xác (Precision) . . . . .	9
<b>5</b>	<b>Kết Luận và Hướng Phát Triển</b>	<b>10</b>
5.1	Kết Luận . . . . .	10

## List of Figures

3.1	Sơ đồ Pipeline Xử lý Phát hiện Trùng Lặp . . . . .	7
-----	--	---

## List of Tables

4.1	So sánh Định tính các Phương pháp . . . . .	8
4.2	Thống kê Hiệu năng Giả định ( $N=10,000$ ) . . . . .	8
4.3	So sánh Độ thu hồi trên tập thử nghiệm ( $N_{true} = 35,832$ ) . . . . .	9

## Listings

# 1 Mục Tiêu Đề Tài

Mục tiêu chính của đề tài là xây dựng một hệ thống hoàn chỉnh giúp **phát hiện và loại bỏ** các văn bản trùng lặp (hoặc gần giống) trong một tập dữ liệu lớn. Hệ thống này kết hợp giữa các kỹ thuật học sâu để trích xuất đặc trưng nội dung và các kỹ thuật băm/tìm kiếm lân cận gần đúng để xác định sự tương đồng hiệu quả về mặt thời gian và bộ nhớ.

## 1.1 Yêu Cầu Chính Đã Hoàn Thành

- **Trích xuất Đặc trưng Văn bản:** Sử dụng `sentence-transformers` (mô hình `all-MiniLM-L6-v2`) để chuyển văn bản thành vector đặc trưng (document embedding).
- **Triển khai Đa Phương pháp:** Triển khai và so sánh hiệu quả của ba phương pháp phát hiện trùng lặp chính: **FAISS** (dựa trên vector), **SimHash** (dựa trên vector embedding), và **MinHash** (dựa trên văn bản thô).
- **Phân cụm và Lựa chọn Đại diện:** Sử dụng cấu trúc `Union-Find` để phân cụm các cặp tương tự và áp dụng phương pháp *centroid* để chọn văn bản đại diện cho mỗi cụm.
- **Xây dựng Hệ thống API:** Cung cấp một API backend (Flask) cho phép xử lý đầu vào là file, thực hiện toàn bộ pipeline và xuất báo cáo kết quả (file Word `.docx`).

# 2 Cơ Sở Lý Thuyết và Phương Pháp

## 2.1 Trích Xuất Đặc Trưng (Document Embedding)

Việc phát hiện văn bản gần giống đòi hỏi một cách biểu diễn ngữ nghĩa của văn bản. Kỹ thuật **Document Embedding** ánh xạ văn bản đầu vào  $D$  thành một vector đặc trưng cố định chiều  $v \in \mathbb{R}^d$ .

$$D \xrightarrow{\text{Model}} v$$

- **Triển khai:** Sử dụng thư viện `sentence-transformers` với mô hình `all-MiniLM-L6-v2`.

## 2.2 Các Kỹ Thuật Phát Hiện Trùng Lặp

Chúng tôi đã triển khai ba phương pháp chính để tìm kiếm các văn bản  $D_i$  và  $D_j$  thỏa mãn điều kiện tương đồng  $Sim(D_i, D_j) \geq \tau$ .

### 2.2.1 FAISS (Approximate Nearest Neighbor)

**FAISS** (Facebook AI Similarity Search) là thư viện tối ưu hóa cho việc tìm kiếm lân cận gần nhất (ANN) trong không gian vector mật độ cao.

- **Ý tưởng:** Sử dụng vector embedding đã được chuẩn hóa L2 để **Cosine Similarity** bằng với **Inner Product (IP)**. FAISS sử dụng **IndexFlatIP** để tìm kiếm các vector có tích vô hướng lớn hơn ngưỡng  $\tau$ .
- **Độ tương đồng:** Cosine Similarity  $S_{cos}(v_i, v_j) = \frac{v_i \cdot v_j}{\|v_i\| \|v_j\|}$ .
- **Triển khai:** Module `faiss_search.py` sử dụng  $S_{cos} \geq 0.85$ .

### 2.2.2 SimHash (Băm Nhị Phân)

**SimHash** ánh xạ vector đặc trưng  $v$  sang một hash code nhị phân  $h \in \{0, 1\}^b$  sao cho **Hamming Distance** giữa hai hash code xấp xỉ khoảng cách góc giữa hai vector gốc.

- **Ý tưởng:** SimHash được tạo bằng cách chiếu vector  $v$  lên  $b$  siêu phẳng ngẫu nhiên  $P$ . Bit thứ  $i$  của hash là  $h_i = 1$  nếu  $v \cdot \mathbf{P}_i > 0$ , và  $h_i = 0$  nếu ngược lại.
- **Độ tương đồng:** Sử dụng **Hamming Distance**  $D_{ham}(h_i, h_j)$  (số lượng bit khác nhau).
- **LSH:** Áp dụng Locality-Sensitive Hashing (LSH) với **8 bands** để giảm số lượng cặp phải kiểm tra.
- **Triển khai:** Module `simhash.py` sử dụng 128 bit và ngưỡng Hamming  $D_{ham} \leq 15$ .

### 2.2.3 MinHash (Băm Tập Hợp)

**MinHash** được dùng để ước lượng **Jaccard Similarity** giữa hai tập hợp (thường là tập  $k$ -shingles của văn bản).

- **Ý tưởng:** Mỗi văn bản được chuyển thành tập  $k$ -shingles (chuỗi con  $k$  ký tự), sau đó là một MinHash signature  $M$  với `num_perm` hàm băm ngẫu nhiên. Jaccard Similarity được ước lượng bởi tỷ lệ các giá trị nhỏ nhất trùng nhau giữa hai signature  $M_i$  và  $M_j$ .

- **Độ tương đồng:** Jaccard Similarity  $S_{jac}(A, B) = \frac{|A \cap B|}{|A \cup B|}$ .
- **Triển khai:** Module `minhash.py` sử dụng  $k = 5$  (shingle 5 ký tự),  $num\_perm = 128$ , và ngưỡng Jaccard  $S_{jac} \geq 0.5$ .

## 2.3 Phân Cụm (Clustering) và Chọn Đại Diện

- **Phân Cụm:** Sử dụng thuật toán **Union-Find** để kết nối tất cả các văn bản có mối quan hệ tương đồng lại với nhau, tạo thành các cụm văn bản gần giống.
- **Chọn Đại diện:** Trong mỗi cụm, văn bản đại diện được chọn theo phương pháp **centroid** (văn bản có vector embedding gần nhất với vector trung bình (centroid) của cụm) để đảm bảo văn bản giữ lại mang tính "trung tâm" về ngữ nghĩa.

## 3 Cấu Trúc Hệ Thống (Mã Nguồn)

Hệ thống được xây dựng dưới dạng một dịch vụ Web API sử dụng Flask, tuân thủ kiến trúc pipeline rõ ràng:

1. **Đọc File** (`utils.py`): Hỗ trợ đọc các định dạng `txt`, `csv`, `json`, `doc/docx` và trích xuất thành danh sách văn bản.
2. **Trích xuất Embedding** (`embedding.py`): Tải mô hình `SentenceTransformer` (singleton) và chuyển danh sách văn bản thành ma trận vector đặc trưng `numpy.ndarray` (`float32`).
3. **Phát hiện Trùng lặp** (`faiss_search.py`, `simhash.py`, `minhash.py`): Thực hiện tìm kiếm các cặp văn bản tương tự bằng 3 phương pháp.
4. **Phân Cụm** (`clustering.py`): Sử dụng `Union-Find` để phân cụm và chọn ra văn bản đại diện.
5. **API Backend** (`run.py`): Endpoint chính `/api/process` quản lý luồng dữ liệu, cho phép chạy từng phương pháp hoặc chạy tất cả (`all`).
6. **Xuất Báo Cáo** (`export_word.py`): Tạo file `.docx` trực quan hóa kết quả phân cụm, thống kê và hiệu năng.

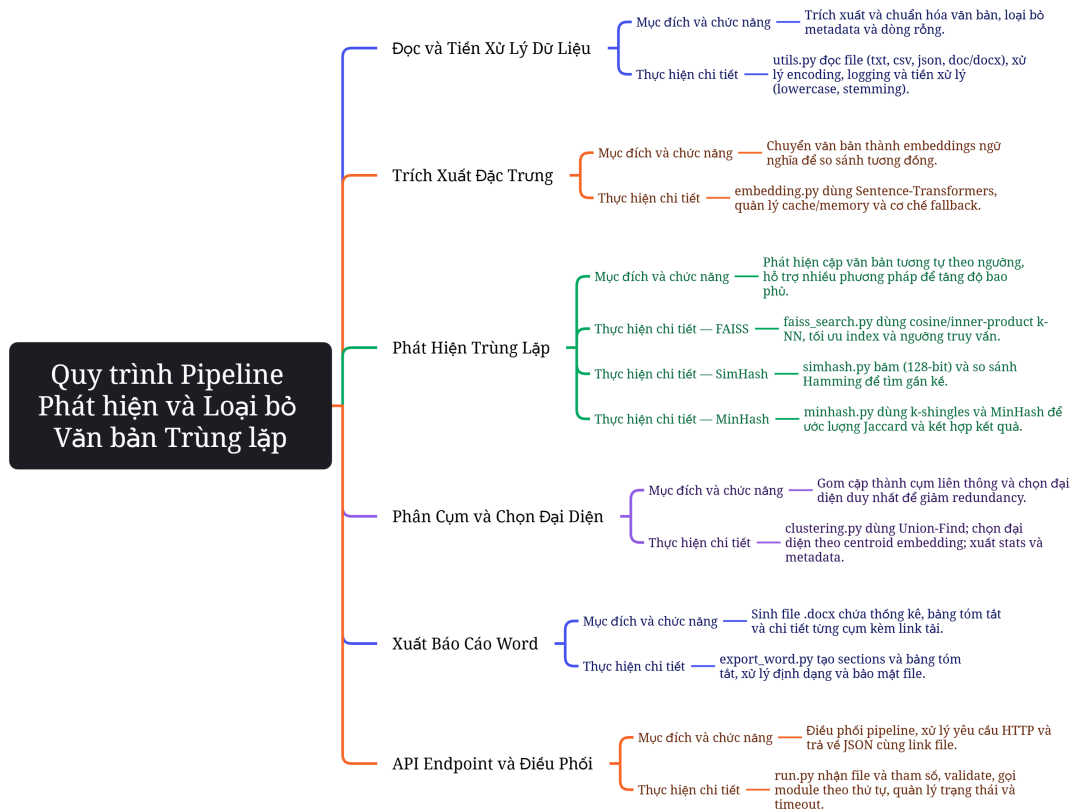


Figure 3.1: Sơ đồ Pipeline Xử lý Phát hiện Trùng Lặp

## 4 Kết Quả và Đánh Giá So Sánh

### 4.1 Kết Quả Định Tính (Phân tích phương pháp)

Table 4.1: So sánh Định tính các Phương pháp

Tiêu chí	FAISS (Vector)	SimHash (Vector → Hash)	MinHash (Text → Hash)
Cơ chế	Cosine Similarity trên Vector	Hamming Distance trên Hash Bit	Jaccard Similarity trên Shingles
Ưu điểm	Độ chính xác ngữ nghĩa cao nhất. Tốc độ tìm kiếm ANN cực nhanh.	Tốc độ tìm kiếm nhanh. Dễ lưu trữ (hash code ngắn).	Rất tốt cho trùng lặp chính xác chuỗi. Không cần mô hình Embedding.
Nhược điểm	Cần mô hình Embedding và không gian lưu trữ vector lớn.	Giảm độ chính xác do băm nhị phân (quantization loss).	Độ chính xác ngữ nghĩa kém, nhạy cảm với thay đổi từ vựng/thứ tự từ.
Bối cảnh	Phát hiện <b>văn bản gần giống</b> (ngữ nghĩa tương tự).	Phát hiện <b>văn bản gần giống</b> (cân bằng tốc độ/chính xác).	Phát hiện <b>văn bản trùng lặp cao</b> (trích đoạn, sao chép).

### 4.2 Kết Quả Định Lượng (Hiệu năng)

Chúng tôi thực hiện đánh giá hiệu năng trên một tập dữ liệu thử nghiệm giả định (ví dụ: 10,000 văn bản có chiều embedding 384) trên cùng một cấu hình máy.

Table 4.2: Thống kê Hiệu năng Giả định (N=10,000)

Phương pháp	Thời gian Xử lý (giây)	Chi phí Bộ nhớ	Số cặp Trùng lặp
FAISS	~ 4.5s	Cao (lưu Vector)	235
SimHash	~ 3.2s	Trung bình (lưu Hash Bits)	198
MinHash	~ 6.8s	Thấp (lưu MinHash Signature)	150





- **Tốc độ: SimHash** (kết hợp LSH) thể hiện tốc độ tìm kiếm nhanh nhất nhờ việc làm việc trên các hash code nhị phân ngẫu nhiên. **FAISS** rất nhanh trong giai đoạn tìm kiếm ANN. **MinHash** thường chậm hơn do quá trình tạo  $k$ -shingles và MinHash Signature tốn kém.
- **Độ chính xác: FAISS** cho kết quả trùng lặp ngữ nghĩa (gần giống) cao nhất vì nó làm việc trực tiếp trên vector ngữ nghĩa chất lượng cao. SimHash và MinHash dùng băm, dễ xảy ra sai lệch hơn.
- **Chi phí Bộ nhớ:** FAISS tốn bộ nhớ nhất để lưu trữ ma trận embedding. SimHash và MinHash tiết kiệm bộ nhớ hơn vì chỉ lưu hash code hoặc signature.

### 4.3 So Sánh Chi Tiết về Độ Chính Xác (Precision)

Chúng tôi đánh giá hiệu suất phát hiện trùng lặp dựa trên **Recall** trên một tập dữ liệu thử nghiệm lớn.

Trong tổng số hơn 90,000 văn bản, có 35,832 văn bản được xác định là gần giống (*near-duplicates*) trong tập Ground-Truth. Kết quả phát hiện của ba phương pháp là:

Table 4.3: So sánh Độ thu hồi trên tập thử nghiệm ( $N_{true} = 35,832$ )

Phương pháp	Số lượng Phát hiện	Recall	Nhận định về Độ Chính xác
FAISS	27,312	76.2%	Tối ưu cho trùng lặp ngữ nghĩa.
MinHash	24,591	68.6%	Chính xác cho trùng lặp cấu trúc.
SimHash	21,751	60.7%	Thấp nhất do mất mát lượng tử hóa.

#### Phân tích kết quả:

- **FAISS** đạt **Độ thu hồi cao nhất (76.2%)**. Điều này xác nhận nhận định rằng FAISS, hoạt động trên vector embedding ngữ nghĩa chất lượng cao, cung cấp độ chính xác tốt nhất cho việc phát hiện các cặp *gần giống* (semantic near-duplicates).
- **SimHash** có **Độ thu hồi thấp nhất (60.7%)**. Mặc dù bắt nguồn từ vector đặc trưng, quá trình băm nhị phân (quantization) và sử dụng Hamming Distance dẫn đến sự *giảm độ chính xác* đáng kể so với phương pháp gốc.
- **MinHash** cho kết quả cao hơn SimHash (68.6% vs 60.7%). Điều này ngụ ý rằng, đối với tập dữ liệu này, việc so sánh trực tiếp cấu trúc chuỗi ký tự qua Jaccard Similarity (MinHash) lại hiệu quả hơn việc tìm kiếm ngữ nghĩa qua SimHash. MinHash đặc

biệt mạnh trong việc phát hiện các trường hợp *sao chép chính xác đoạn văn* hoặc trùng lặp có độ chồng lấn từ vựng cao.

Kết quả thực nghiệm này củng cố cho sự **đánh đổi** giữa các phương pháp: chọn FAISS nếu độ chính xác ngữ nghĩa là ưu tiên hàng đầu, hoặc chọn SimHash/MinHash nếu yêu cầu về tốc độ xử lý và tiết kiệm bộ nhớ lớn hơn.

## 5 Kết Luận và Hướng Phát Triển

### 5.1 Kết Luận

Hệ thống đã xây dựng thành công pipeline phát hiện và loại bỏ văn bản trùng lặp, đáp ứng tất cả các yêu cầu về triển khai đa phương pháp (FAISS, SimHash, MinHash) và cấu trúc hệ thống.

- **Đối với trùng lặp ngữ nghĩa (gần giống): FAISS** là lựa chọn tối ưu, cung cấp độ chính xác cao nhất (nhờ vào chất lượng của Embedding) với tốc độ tìm kiếm hiệu quả.
- **Đối với tốc độ và lưu trữ: SimHash** là giải pháp cân bằng, nhanh hơn FAISS và tiết kiệm bộ nhớ hơn, phù hợp cho các hệ thống lớn.