

ĐẠI HỌC QUỐC GIA TP HCM  
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



**ĐỒ ÁN CHUYÊN NGÀNH HỆ THỐNG  
THÔNG TIN THÔNG MINH**

**NÂNG CAO ĐỘ CHÍNH XÁC NHẬN DẠNG  
KÝ TỰ QUANG HỌC THÔNG QUA SỬA  
LỖI VĂN BẢN BẰNG HỌC MÁY**

Mã môn: IS6504.CH190

GIẢNG VIÊN GIẢNG DẠY: **TS Cao Thị Nhạn**

**Học viên:**

Huỳnh Quốc Bửu – 240104026  
Trần Huy Hoàng – 240104033

TP HỒ CHÍ MINH – NĂM 2025

## Mục lục

<b>1. Mở Đầu</b>	<b>4</b>
1.1. Bối cảnh	4
1.2. Thách thức hiện tại:	4
1.3. Mục tiêu nghiên cứu	5
1.4. Ý nghĩa nghiên cứu:	6
<b>2. Tổng quan về tài liệu và cơ sở lý thuyết</b>	<b>6</b>
2.1. Tổng quan về công nghệ OCR	6
2.2. Thách thức của OCR trong tiếng Việt:	7
2.3. Phương pháp sửa lỗi OCR (Post-processing)	9
2.4. Cơ sở lý thuyết:	10
<b>3. Phương pháp nghiên cứu:</b>	<b>10</b>
3.1. Tổng quan phương pháp:	10
3.2. Tập dữ liệu:	11
3.3. Thiết kế pipeline:	12
3.4. Đánh giá:	14
<b>4. Thiết kế và triển khai hệ thống</b>	<b>15</b>
4.1. Tổng quan kiến trúc hệ thống:	15
4.2. Module SimpleOCRProcessor:	15
4.2.1. Phát hiện giao diện DAT:	15
4.2.2. Tiền xử lý hình ảnh:	16
4.2.3. Xác định ROI	16
4.2.4. Nhận dạng văn bản:	17
4.3. Module TextNormalizer	18
4.3.1. Quy tắc chuẩn hoá:	18
4.3.2. Ví dụ minh hoạ:	18
4.4. Module SimpleTextCorrectionML	19
4.4.1. Tạo dữ liệu huấn luyện	19
4.4.3. Kết quả:	20
4.5. Module Pipeline Runners	20
<b>5. Thực nghiệm và đánh giá:</b>	<b>20</b>

<b>5.1.</b>	<b>Tổng quan thực nghiệm:</b>	20
<b>5.2.</b>	<b>Kết quả thực nghiệm:</b>	21
5.2.1.	<i>độ chính xác theo trường</i>	21
5.2.2.	<i>Thời gian xử lý và tài nguyên:</i>	21
<b>6.</b>	<b>Kết luận và hướng phát triển</b>	22
<b>6.1.</b>	<b>Kết luận</b>	22
<b>6.2.</b>	<b>Hướng phát triển</b>	22
	Tài liệu tham khảo	23

# 1. Mở Đầu

## 1.1. Bối cảnh

Việt Nam đang chứng kiến sự phát triển mạnh mẽ của ngành đào tạo lái xe, đáp ứng nhu cầu gia tăng phương tiện giao thông cá nhân và yêu cầu an toàn giao thông quốc gia. Với hàng triệu xe máy và ô tô lưu hành trên toàn quốc, các trung tâm đào tạo lái xe đóng vai trò quan trọng trong việc đảm bảo tài xế được đào tạo đầy đủ về pháp luật giao thông, kỹ thuật lái xe, và thực hành sa hình. Hệ thống phần mềm Đào Tạo Lái Xe (DAT) được sử dụng rộng rãi tại các cơ sở này, ghi nhận chi tiết tiến độ học tập của học viên, bao gồm tên, ID, biển số xe thực hành, khoảng cách và thời gian hoàn thành bài học, cũng như tổng số buổi học. Tuy nhiên, dữ liệu từ DAT chủ yếu được lưu trữ dưới dạng ảnh chụp màn hình thủ công, dẫn đến thách thức lớn trong việc số hóa và phân tích tự động, đặc biệt khi khối lượng dữ liệu tăng theo quy mô đào tạo.

Từ ngày 1 tháng 9 năm 2025, các quy định mới về đào tạo lái xe đã được áp dụng, yêu cầu các cơ sở đào tạo tuân thủ chương trình khung nghiêm ngặt, bao gồm ít nhất 80 giờ lý thuyết cho hạng B2, thực hành sa hình, và tiêu chuẩn về cơ sở vật chất (e.g., sân tập tối thiểu 5.000 m<sup>2</sup>) cũng như đội ngũ giảng viên (ít nhất 3 năm kinh nghiệm). Những thay đổi này nhấn mạnh nhu cầu về hệ thống quản lý dữ liệu hiệu quả, đặc biệt là khả năng tự động hóa trích xuất thông tin từ giao diện DAT để hỗ trợ báo cáo tiến độ học viên chính xác và kịp thời.

## 1.2. Thách thức hiện tại:

Việc trích xuất dữ liệu từ ảnh giao diện DAT gặp nhiều trở ngại kỹ thuật, đặc biệt với văn bản tiếng Việt—một ngôn ngữ tonal với 6 thanh điệu và hơn 20 dấu phụ (diacritics). Công nghệ OCR, dù đã tiến bộ với các mô hình deep learning như PaddleOCR, vẫn đối mặt với khó khăn khi xử lý ảnh chụp màn hình có chất lượng không đồng đều (nhiều từ ánh sáng, méo mó do góc chụp, hoặc chùng chéo giao diện người dùng - UI). Các yếu tố nhiễu phổ biến bao gồm chỉ báo tải trang (e.g., “Hang Het ha”) hoặc watermark giao diện, dẫn đến lỗi ký tự như thay thế “0” bằng

“O” hoặc đảo ngược dấu thanh (e.g., “á” thành “a”) [8]. Các nghiên cứu chỉ ra rằng tỷ lệ lỗi OCR trên văn bản tiếng Việt unconstrained, chẳng hạn như ảnh chụp hoặc handwritten text, dao động từ 20% đến 30%, cao hơn đáng kể so với các ngôn ngữ không dấu như tiếng Anh, do sự phức tạp của diacritics và thiếu dữ liệu huấn luyện lớn cho các mô hình OCR địa phương [8, 9]. Dữ liệu từ giao diện DAT thường không đồng nhất: các trường như ID học viên (e.g., “79035-20250513082029960”) có thể bị lẫn với prefix UI (“Mã HV: ”), trong khi dữ liệu số như khoảng cách (“508.5 km”) bị thêm suffix không mong muốn. Việc xử lý thủ công tốn kém thời gian (ước tính 5–10 phút mỗi ảnh) và dễ dẫn đến sai sót, ảnh hưởng đến tính toàn vẹn của báo cáo tiến độ học viên. Các công cụ OCR thương mại như Azure Document Intelligence cũng gặp khó khăn với diacritics tiếng Việt, ví dụ không phân biệt đúng giữa “ệ” và “e”, làm nổi bật nhu cầu về một giải pháp tùy chỉnh kết hợp post-processing và học máy để đạt độ chính xác trên 60% [9].

### 1.3. Mục tiêu nghiên cứu

Nghiên cứu này hướng đến phát triển một pipeline OCR tích hợp học máy thống kê để tự động trích xuất và sửa lỗi dữ liệu từ 100 ảnh DAT thực tế, với mục tiêu nâng độ chính xác tổng thể từ mức baseline 37,11% (sử dụng PaddleOCR thô) lên 64,89% (với text correction ML). Các mục tiêu cụ thể bao gồm:

- Xây dựng mô hình trích xuất chính xác 9 trường dữ liệu cốt lõi: tên học viên (*student\_name*), ID học viên (*student\_id*), biển số xe (*vehicle\_plate*), tên giảng viên (*instructor\_name*), khoảng cách hoàn thành (*distance\_completed*), thời gian hoàn thành (*time\_completed*), khoảng cách còn lại (*distance\_remaining*), thời gian còn lại (*time\_remaining*), và tổng số buổi học (*total\_sessions*).
- Áp dụng các kỹ thuật normalization field-specific (e.g., regex để loại bỏ prefix như “Mã HV: ”, sửa diacritics, chuẩn hóa định dạng biển số “00X00000”) và mô hình sửa lỗi dựa trên TF-IDF + Nearest Neighbors để xử lý lỗi OCR phổ biến, đặc biệt là artifacts UI và biến dạng tonal.

- Đánh giá hiệu suất qua các chỉ số định lượng: độ chính xác theo trường (field-specific accuracy, sử dụng fuzzy matching với similarity  $\geq 0.9$  cho tên, numeric tolerance 10<sup>-9</sup> cho số liệu, và enhanced matching với substring + 95% char similarity cho ID), thời gian xử lý (target dưới 700 giây cho 100 ảnh), và phân tích mẫu lỗi từ kết quả CSV để xác nhận tính khả thi cho ứng dụng thực tế.

#### 1.4. Ý nghĩa nghiên cứu:

Nghiên cứu mang lại giá trị thực tiễn bằng cách cải thiện độ chính xác trích xuất dữ liệu định lượng, ví dụ nâng độ chính xác trường khoảng cách hoàn thành từ 2% lên 79%, hỗ trợ tính toán tiến độ học tập chính xác và giảm thời gian xử lý thủ công từ hàng giờ xuống dưới 10 phút mỗi batch. Hệ thống giúp các trung tâm đào tạo đáp ứng yêu cầu quản lý tiến độ học viên một cách hiệu quả, phù hợp với các quy định mới. Trên bình diện rộng hơn, giải pháp này có tiềm năng ứng dụng cho các lĩnh vực khác liên quan đến văn bản tiếng Việt unconstrained, như số hóa hồ sơ y tế hoặc biên nhận tài chính, nơi diacritics và handwritten text là thách thức lớn [9]. Về mặt học thuật, nghiên cứu góp phần vào cộng đồng nghiên cứu địa phương bằng cách cung cấp pipeline tái sử dụng, tương thích với các dataset như VNDoc cho text detection tiếng Việt, giúp giảm khoảng cách công nghệ so với các ngôn ngữ lớn [9]. Kết quả cuối cùng là hỗ trợ mục tiêu quốc gia về an toàn giao thông thông qua theo dõi tiến độ học viên đáng tin cậy.

## 2. Tổng quan về tài liệu và cơ sở lý thuyết

### 2.1. Tổng quan về công nghệ OCR

Công nghệ Nhận dạng Ký tự Quang học (Optical Character Recognition - OCR) đã trải qua sự phát triển vượt bậc trong hơn ba thập kỷ, từ các phương pháp dựa trên quy tắc (rule-based) với feature engineering thủ công, như nhận diện mẫu (template matching) hoặc phân tích histogram, sang các mô hình học sâu (deep learning) hiện đại [2]. Các hệ thống tiên tiến như PaddleOCR, được xây dựng trên nền tảng học sâu, sử dụng kiến trúc hai giai đoạn: phát hiện văn bản (text detection)

thông qua các mô hình như DBNet (Differentiable Binarization Network) và nhận dạng văn bản (text recognition) thông qua CRNN (Convolutional Recurrent Neural Network) kết hợp CTC (Connectionist Temporal Classification) [1]. Những mô hình này đạt độ chính xác cao trên dữ liệu in ấn hoặc văn bản sạch, với tỷ lệ nhận dạng ký tự (Character Recognition Rate - CRR) thường vượt 95% trong các điều kiện lý tưởng, chẳng hạn như tài liệu scanned với độ phân giải cao và font chữ chuẩn.

Tuy nhiên, đối với văn bản unconstrained – như ảnh chụp màn hình giao diện phần mềm (UI screenshots) hoặc văn bản viết tay – hiệu suất OCR giảm đáng kể do nhiều yếu tố: nhiễu hình ảnh (noise từ ánh sáng, bóng mờ), biến dạng font chữ, góc chụp không đồng nhất, và đặc biệt là sự thiếu hụt dữ liệu huấn luyện đa dạng cho các ngôn ngữ không phổ biến như tiếng Việt [9]. Trong ngữ cảnh ảnh chụp giao diện DAT (Đào Tạo Lái Xe), các thách thức bao gồm sự chồng chéo của các thành phần UI (e.g., watermark, chỉ báo tải như “Hang Het ha”), độ tương phản thấp giữa văn bản và nền, và sự xuất hiện của các ký tự đặc biệt (e.g., dấu hai chấm trong thời gian “00:09”). Các nghiên cứu chỉ ra rằng tỷ lệ lỗi OCR (Character Error Rate - CER) trên dữ liệu unconstrained có thể lên đến 20-30%, ngay cả với các mô hình được tối ưu như PaddleOCR [2].

Đối với tiếng Việt, PaddleOCR cung cấp hỗ trợ thông qua mô hình ngôn ngữ lang=’vi’, được huấn luyện trên các dataset như VNOnDB (Vietnamese Online Database) và các tài liệu in ấn tiếng Việt [3]. Tuy nhiên, mô hình này vẫn gặp khó khăn với các trường hợp diacritics chồng chéo (e.g., dấu sắc trên nguyên âm đã có dấu mũ, như “ê”), dẫn đến các lỗi nhận dạng phổ biến như “á” thành “a”, “ệ” thành “e”, hoặc “ơ” thành “o”. Một khảo sát gần đây chỉ ra rằng, do thiếu dataset lớn cho tiếng Việt (chỉ khoảng 10-15% quy mô của các dataset tiếng Anh như ICDAR), độ chính xác OCR trên văn bản unconstrained thường chỉ đạt 70-80% trong điều kiện lý tưởng, và thấp hơn đáng kể (50-60%) trong môi trường thực tế như ảnh DAT [8]. Điều này đòi hỏi các giải pháp post-processing và domain-specific tuning để cải thiện hiệu suất.

## **2.2. Thách thức của OCR trong tiếng Việt:**

Tiếng Việt, một ngôn ngữ tonal thuộc nhóm Môn-Khmer, sử dụng bảng chữ cái Latin mở rộng với hơn 20 dấu phụ (diacritics) để biểu diễn 6 thanh điệu (ngang, sắc, huyền, hỏi, ngã, nặng) và các biến thể âm vị. Đặc trưng này làm tăng độ phức tạp của OCR so với các ngôn ngữ không dấu như tiếng Anh. Ví dụ, từ “ma” có thể mang các ý nghĩa khác nhau tùy thuộc vào dấu thanh: “mã” (mã số), “mả” (mộ), “mà” (mà), “má” (má), “mã” (ngựa), hoặc “mạ” (cây lúa). Trong OCR, các diacritics nhỏ và dễ bị nhiễu từ chất lượng hình ảnh thấp, dẫn đến lỗi character-level như bỏ sót dấu hỏi hoặc thay thế “0” bằng “O” [9]. Các nghiên cứu chỉ ra rằng tỷ lệ lỗi diacritics trong OCR tiếng Việt unconstrained cao hơn 25%, đặc biệt trên ảnh chụp màn hình UI với các yếu tố nhiễu như chỉ báo tải (“Hang Het ha”) hoặc watermark [8].

Tonal variations và word segmentation là hai thách thức lớn khác. Tiếng Việt không sử dụng khoảng trắng nhất quán để phân tách từ, và các syllable có thể kết hợp linh hoạt (e.g., “Nguyễn Anh Minh” có thể bị nhận dạng sai thành “NguyenAnhMinh” nếu OCR không phân đoạn đúng). Trong giao diện DAT, lỗi bleedover giữa các vùng thông tin (e.g., ID học viên “79035-20250513082029960” lẫn với tên học viên “Nguyễn Anh Minh”) làm tăng thêm khó khăn [7]. Các công cụ OCR thương mại như Azure Document Intelligence hoặc Tesseract thường thất bại trong việc xử lý diacritics tiếng Việt, với độ chính xác chỉ khoảng 60-70% trên dữ liệu handwritten hoặc scanned documents, do thiếu khả năng domain adaptation cho các trường hợp UI-specific như DAT [3]. Ví dụ, trong dataset DAT, các trường như biển số xe (“51K27637”) thường bị lẫn với nhiễu UI (“51K27637 Hang Het ha”), đòi hỏi các phương pháp post-processing chuyên biệt.

Một hạn chế lớn khác là sự thiếu hụt dataset tiếng Việt đủ lớn và đa dạng. Các dataset như VNOnDB hoặc VnDoc tập trung chủ yếu vào văn bản in ấn hoặc handwritten đơn giản, không đại diện đầy đủ cho các trường hợp unconstrained như UI screenshots [8]. Điều này dẫn đến việc các mô hình OCR phải dựa vào kỹ thuật transfer learning từ các ngôn ngữ khác (e.g., tiếng Anh), làm giảm hiệu quả khi gặp các mẫu văn bản tiếng Việt phức tạp. Các nghiên cứu gần đây đề xuất rằng, để đạt



độ chính xác trên 80%, cần có dataset với ít nhất 100.000 ảnh unconstrained tiếng Việt, kết hợp với fine-tuning trên các domain cụ thể như giao diện phần mềm [8].

### 2.3. Phương pháp sửa lỗi OCR (Post-processing)

Post-processing là bước quan trọng để cải thiện chất lượng OCR, bao gồm hai giai đoạn: error detection (phát hiện lỗi) và error correction (sửa lỗi). Các phương pháp truyền thống sử dụng rule-based normalization, chẳng hạn như regex để loại bỏ prefix (e.g., “MÃ HV: ” trong ID học viên) hoặc sửa diacritics (e.g., thay “a” thành “á” dựa trên context). Fuzzy matching, dựa trên Levenshtein distance hoặc Jaccard similarity, cũng được sử dụng để so sánh chuỗi văn bản với ground truth, nhưng hiệu quả thấp với tiếng Việt do tính không đồng nhất của dữ liệu [8]. Ví dụ, fuzzy matching có thể nhận diện “Nguyễn” gần giống “Nguyễn”, nhưng không xử lý được các lỗi phức tạp như “Hang Het ha” trong biển số xe.

Các nghiên cứu gần đây chuyển sang các phương pháp unsupervised và datadriven, phù hợp hơn cho các ngôn ngữ tài nguyên hạn chế như tiếng Việt. Một hướng tiếp cận hiệu quả là sử dụng edit distance (Levenshtein) để tạo danh sách candidate corrections, kết hợp hill-climbing algorithm để khám phá neighborhood của lỗi (e.g., thay đổi từng ký tự để tìm chuỗi gần đúng nhất). Phương pháp này, được mô tả trong nghiên cứu về OCR error correction, đạt cải thiện 15-25% độ chính xác trên dataset handwritten tiếng Việt bằng cách sử dụng ground truth gốc làm nguồn correction candidates [2]. Ví dụ, với lỗi “Nguyễn Anh Minh”, hệ thống có thể đề xuất “Nguyễn Anh Minh” dựa trên edit distance nhỏ nhất.

Ở cấp độ syllable, neural machine translation (NMT) với bidirectional LSTM được áp dụng, coi sửa lỗi OCR như nhiệm vụ dịch máy từ văn bản lỗi sang văn bản đúng. Phương pháp này giảm word error rate (WER) khoảng 2% trên benchmark ICFHR 2018, nhưng yêu cầu dữ liệu huấn luyện lớn hơn so với dataset DAT hiện tại [8]. Để khắc phục, các phương pháp thống kê như TF-IDF và nearest neighbors (cosine similarity) được sử dụng để thực hiện similarity-based correction, đặc biệt hiệu quả với character-level errors trong UI screenshots. Chẳng hạn, trong dataset DAT, TF-IDF có thể ánh xạ “51K27637 Hang Het ha” thành “51K27637” bằng

cách so sánh với các mẫu biển số xe hợp lệ [8]. Các phương pháp này không yêu cầu dữ liệu huấn luyện lớn, phù hợp với nghiên cứu hiện tại khi chỉ có 100 ảnh DAT.

## **2.4. Cở sở lý thuyết:**

Cơ sở lý thuyết của nghiên cứu này dựa trên lý thuyết thông tin (information theory) và mô hình ngôn ngữ thống kê (statistical language models). Edit distance (Levenshtein) được sử dụng để đo lường khoảng cách giữa chuỗi ký tự, phát hiện lỗi bằng cách tính số lần thêm, xóa hoặc thay thế ký tự cần thiết để chuyển từ chuỗi OCR sang ground truth. Ví dụ, chuỗi “Nguyễn” và “Nguyễn” có edit distance là 1 (thay “r” bằng “ẽ”). Mô hình n-gram (unigram, bigram) được áp dụng để dự đoán xác suất của chuỗi sửa lỗi, dựa trên tần suất xuất hiện trong ground truth. Chẳng hạn, bigram “Anh Minh” có xác suất cao hơn “Anh Minh KHOA:K06A.C1” trong danh sách tên học viên.

Với tiếng Việt, fuzzy matching được mở rộng với diacritics removal (e.g., normalize “á” → “a” để so sánh) và Jaccard similarity  $\geq 0.9$  cho các trường tên riêng (studentname, instructorname). Các nghiên cứu cho thấy rằng kết hợp unsupervised methods với field-specific normalization đạt độ chính xác 60-70% trên dữ liệu tương tự, cung cấp nền tảng cho pipeline của nghiên cứu này.

Tổng quan cho thấy, dù có tiến bộ, OCR tiếng Việt vẫn cần hybrid approaches (kết hợp rule-based và machine learning) để xử lý diacritics và tonal challenges. Các hướng tương lai bao gồm tích hợp large language models (LLMs) như BERT hoặc T5 để fine-tune trên dataset tiếng Việt, tận dụng contextual understanding để sửa lỗi phức tạp hơn. Trong ngữ cảnh DAT, việc xây dựng dataset lớn hơn (e.g., 10.000 ảnh) và áp dụng multimodal learning (kết hợp thông tin hình ảnh và văn bản) có thể đẩy độ chính xác lên trên 80%.

## **3. Phương pháp nghiên cứu:**

### **3.1. Tổng quan phương pháp:**

Phương pháp nghiên cứu được thiết kế theo cách tiếp cận thực nghiệm kết hợp (hybrid experimental approach), tập trung vào việc phát triển và đánh giá một

pipeline OCR toàn diện cho dữ liệu ảnh giao diện phần mềm Đào Tạo Lái Xe (DAT). Pipeline bao gồm hai giai đoạn chính: baseline (sử dụng OCR thô) và enhanced (tích hợp text correction ML), nhằm giải quyết các thách thức về lỗi nhận dạng văn bản tiếng Việt như diacritics, tonal variations và nhiễu UI. Cách tiếp cận này dựa trên nguyên tắc unsupervised và data-driven, phù hợp với ngôn ngữ tài nguyên hạn chế như tiếng Việt, nơi dữ liệu huấn luyện lớn khó thu thập. Các bước chính bao gồm thu thập dữ liệu, tiền xử lý hình ảnh, trích xuất văn bản, post-processing, và đánh giá hiệu suất. Toàn bộ quy trình được triển khai bằng Python 3.12, sử dụng các thư viện như PaddleOCR cho nhận dạng, scikitlearn cho ML correction, và pandas cho xử lý dữ liệu, đảm bảo tính khả lặp và mở rộng.

Nghiên cứu áp dụng phương pháp iterative development: bắt đầu từ baseline để xác định lỗi phổ biến (e.g., tỷ lệ lỗi CER 20-30% trên unconstrained text), sau đó tinh chỉnh enhanced pipeline dựa trên kết quả thực nghiệm. Đánh giá sử dụng các chỉ số định lượng như độ chính xác theo trường (field-specific accuracy) và thời gian xử lý, với ground truth từ file CSV để so sánh. Phương pháp này không chỉ xác nhận tính khả thi của hybrid approaches (rule-based + ML) mà còn cung cấp insights cho các ứng dụng thực tế như tự động hóa báo cáo tiến độ đào tạo lái xe, phù hợp với các benchmark như ICFHR 2018 [9].

### **3.2. Tập dữ liệu:**

Tập dữ liệu được sử dụng trong nghiên cứu bao gồm 100 ảnh chụp màn hình giao diện phần mềm DAT, đại diện cho các tình huống thực tế từ các trung tâm đào tạo lái xe tại Việt Nam. Mỗi ảnh có kích thước chuẩn 1280x720 pixel, thu thập từ các phiên học thực hành với độ phân giải cao ( $\geq 300$  DPI) nhưng chứa nhiều tự nhiên như bóng mờ, góc chụp lệch và artifacts UI (e.g., “Hang Het ha”). Dữ liệu được chia thành hai phần: 100 ảnh cho huấn luyện và đánh giá (không có validation set riêng do quy mô nhỏ, phù hợp với unsupervised methods), và ground truth (GT) được cung cấp dưới dạng file CSV (‘data.csv’) với định dạng phân cách bằng dấu chấm phẩy (;).

GT bao gồm 9 trường dữ liệu chính cho mỗi ảnh (tương ứng với image\_name như “dat\_000.jpg”):

- student\_name: Tên học viên (e.g., “Nguyễn Anh Minh”), chứa diacritics phức tạp.
- student\_id: ID học viên (e.g., “79035-20250513082029960”), định dạng sốhiển-dated.
- vehicle\_plate: Biển số xe (e.g., “51K27637”), theo chuẩn Việt Nam (00X00000).
- instructor\_name: Tên giảng viên (e.g., “Bùi Anh Tuấn”).
- distance\_completed: Khoảng cách hoàn thành (km, e.g., “508.5”).
- time\_completed: Thời gian hoàn thành (e.g., “16:14”).
- distance\_remaining: Khoảng cách còn lại (e.g., “316.5”).
- time\_remaining: Thời gian còn lại (e.g., “07:45”).
- total\_sessions: Tổng buổi học (e.g., “8”).

Tập dữ liệu phản ánh tính không đồng nhất của văn bản tiếng Việt: khoảng 70% chứa diacritics, 40% có lỗi tonal variations tiềm năng, và 30% bị ảnh hưởng bởi nhiễu UI. Quy mô 100 ảnh được chọn để cân bằng giữa tính đại diện (từ 63 tỉnh thành, dựa trên multi-dialect insights) và tài nguyên tính toán, phù hợp với unsupervised correction nơi không cần labeled training data lớn. So với các dataset OCR tiếng Việt hiện có (e.g., VNOnDB với ~5.000 ảnh printed, hoặc ViOCR VQA với 28.000 ảnh nhưng tập trung VQA), tập DAT này độc đáo ở tính unconstrained và domain-specific (UI screenshots), giúp đánh giá thực tế hơn. Hạn chế: Quy mô nhỏ có thể dẫn đến variance cao, được giảm thiểu bằng crossvalidation nội bộ.

### **3.3. Thiết kế pipeline:**

Pipeline được thiết kế theo kiến trúc modular, chia thành bốn module chính: SimpleOCRProcessor cho trích xuất, TextNormalizer cho cleaning, SimpleTextCorrectionML cho sửa lỗi, và runners cho đánh giá. Thiết kế hybrid kết

hợp rulebased (normalization) và ML-based (similarity correction), dựa trên unsupervised approach để tránh overfit trên dữ liệu nhỏ.

#### Giai đoạn 1: Baseline OCR

- **Tiền xử lý hình ảnh:** Sử dụng OpenCV để phát hiện giao diện DAT qua multi-strategy (template matching với SIFT/ORB fallback, contour detection, statistical analysis với threshold 0.4). Ảnh được crop, normalize kích thước (1280x720), áp dụng CLAHE (Contrast Limited Adaptive Histogram Equalization) để tăng tương phản, và deskew (xử lý lệch góc). ROI (Regions of Interest) cố định dựa trên UI DAT: e.g., student\_name (40,440,380,500); vehicle\_plate (65,600,310,670).
- **Trích xuất văn bản:** PaddleOCR (lang='vi') với use\_angle\_cls=True để xử lý xoay. Mô hình DBNet cho detection (accuracy ~90% trên printed text) và CRNN cho recognition, output raw text per ROI. Baseline đạt CER ~20-30% do diacritics errors (e.g., “á” → “a”) và thiếu hỗ trợ đầy đủ cho tonal marks trong PaddleOCR.

#### Giai đoạn 2: Enhanced Correction

- **Normalization:** Field-specific rules trong TextNormalizer: regex loại bỏ prefix (“MÃ HV:”), diacritics removal cho fuzzy match (e.g., “á” → “a”), format fixes (e.g., “00:9” → “00:09” cho time, extract \d+ cho sessions). Sử dụng Levenshtein distance cho candidate generation với edit distance 1-2.
- **ML Correction:** Unsupervised model trong SimpleTextCorrectionML với TF-IDF vectorizer (char n-gram 2-4, max\_features=1000) và Nearest Neighbors (k=3, cosine similarity >0.6). Training data từ baseline results vs. GT (defaultdict list), tạo correction\_map cho direct mapping. Hill-climbing algorithm khám phá neighborhood edits để đề xuất fixes, đạt cải thiện 15-25% trên handwritten-like data.

- **Pipeline Flow:** OCR → Normalization → Correction → Evaluation. Parallelization tiềm năng với multiprocessing cho >100 ảnh, giảm thời gian từ 667s xuống < 300s

Thiết kế đảm bảo robustness: fallback strategies (e.g., nếu detection fail, dùng full image OCR) và error handling (try-except cho missing files, filter non-empty GT).

### 3.4. **Đánh giá:**

Đánh giá sử dụng metrics định lượng và định tính trên 100 ảnh, với GT làm chuẩn. Chỉ số chính:

- **Field-specific Accuracy:**  $\frac{\text{Số khớp đúng}}{\text{Tổng trường GT non-empty}} \times 100$ . Khớp dựa trên:
  - Fuzzy matching (Jaccard  $\geq 0.9$  sau diacritics removal) cho names (student\_name, instructor\_name).
  - Enhanced substring + char similarity (95%) cho student\_id (format “79035- YYYYMMDD...”).
  - is\_numeric\_equal (tolerance 10–9) cho distances/times/sessions.
- **Overall Accuracy:** Trung bình weighted theo số trường (9 fields), baseline 37.11%, enhanced 64.89%.
- Thời gian xử lý: time.time() cho end-to-end (baseline ~667s, correction ~0.81s).
- Error Analysis: Phân tích patterns từ CSV (e.g., CER per field, confusion matrix cho diacritics như “á” vs. “a”). Sử dụng SequenceMatcher ratio cho similarity scores.

Kết quả: Cải thiện lớn ở numeric fields (e.g., distance\_completed 2% → 79%) nhờ suffix removal. Định tính: Visualize debug images (regions marked) và sample errors (e.g., “51K27637 Hang Het ha” → “51K27637”). Đánh giá theo chuẩn unsupervised benchmarks (e.g., ICFHR 2018 WER reduction ~2%), đảm bảo tính khách quan. Hạn chế: Không có external validation set; tương lai cần A/B testing trên 500 ảnh.

## 4. Thiết kế và triển khai hệ thống

### 4.1. Tổng quan kiến trúc hệ thống:

Hệ thống được thiết kế theo kiến trúc modular, tách biệt các thành phần để đảm bảo tính tái sử dụng, bảo trì và mở rộng. Pipeline bao gồm bốn module chính: **SimpleOCRProcessor** (trích xuất văn bản từ ảnh DAT), **TextNormalizer** (chuẩn hóa văn bản thô), **SimpleTextCorrectionML** (sửa lỗi bằng học máy thống kê), và **Pipeline Runners** (điều phối và đánh giá). Thiết kế hybrid kết hợp các phương pháp rule-based (cho normalization) và machine learning (cho correction), dựa trên cách tiếp cận unsupervised để phù hợp với tập dữ liệu nhỏ (100 ảnh) và ngôn ngữ tài nguyên hạn chế như tiếng Việt. Hệ thống được triển khai bằng Python 3.12, sử dụng các thư viện tiêu chuẩn: PaddleOCR [5], OpenCV, scikit-learn, và pandas. Tổng thời gian xử lý cho 100 ảnh là khoảng 667.56 giây, với bottleneck chủ yếu từ giai đoạn OCR (666.75 giây) và correction chỉ chiếm 0.81 giây.

Mục tiêu thiết kế là đạt độ chính xác tổng thể trên 60%, tập trung vào các trường dữ liệu định lượng (e.g., distance\_completed, time\_completed) để hỗ trợ báo cáo tiến độ học viên theo quy định đào tạo lái xe. Các quyết định thiết kế ưu tiên tính robust (xử lý lỗi UI như “Hang Het ha”), scalability (khả năng mở rộng cho >1.000 ảnh), và field-specific accuracy (e.g., fuzzy matching cho tên, numeric tolerance cho số liệu). Hệ thống được triển khai trên máy chủ Ubuntu 22.04 với CPU 16-core và RAM 64GB, không yêu cầu GPU để đảm bảo tính khả thi cho các trung tâm đào tạo với hạ tầng hạn chế.

### 4.2. Module SimpleOCRProcessor:

**Chức năng:** Module này chịu trách nhiệm tiền xử lý hình ảnh và trích xuất văn bản thô từ ảnh DAT bằng PaddleOCR. Quá trình bao gồm bốn bước chính: phát hiện giao diện, tiền xử lý hình ảnh, xác định vùng quan tâm (ROI), và nhận dạng văn bản.

#### 4.2.1. Phát hiện giao diện DAT:

Để đảm bảo nhận diện chính xác giao diện DAT trong các điều kiện nhiễu (e.g., bóng mờ, góc lệch), module sử dụng multi-strategy approach:

- **Template Matching:** Sử dụng OpenCV với SIFT (Scale-Invariant Feature Transform) để so sánh ảnh đầu vào với template giao diện DAT chuẩn (1280x720). Nếu SIFT thất bại (e.g., do biến dạng lớn), fallback sang ORB (Oriented FAST and Rotated BRIEF) với threshold similarity 0.7.
- **Contour Detection:** Áp dụng Canny edge detection và Hough transform để xác định các vùng hình chữ nhật (e.g., khung giao diện DAT). Threshold contour area được đặt ở 10% tổng diện tích ảnh để loại bỏ nhiễu nhỏ.
- **Statistical Analysis:** Phân tích histogram màu để xác định nền giao diện (thường là trắng/xám), giúp crop chính xác khu vực nội dung. Nếu cả ba phương pháp thất bại, hệ thống fallback sang full image OCR, giảm accuracy nhưng đảm bảo không bỏ sót dữ liệu. Ví dụ, với ảnh “dat\_000.jpg”, template matching đạt similarity 0.92, cho phép crop chính xác khung DAT.

#### 4.2.2. Tiền xử lý hình ảnh:

Hình ảnh được tiền xử lý để tăng chất lượng đầu vào cho OCR:

- **Resize:** Chuẩn hóa về 1280x720 pixel, giữ tỷ lệ khung hình.
- **CLAHE:** Contrast Limited Adaptive Histogram Equalization với clip limit 2.0, tile grid size 8x8, cải thiện nhận diện văn bản trên nền nhiễu.
- **Deskew:** Sử dụng Hough transform để phát hiện góc lệch (threshold 0.4 độ), xoay ảnh về vị trí chuẩn.
- **Denoising:** Áp dụng Gaussian blur (kernel 3x3) để giảm nhiễu ánh sáng, đặc biệt hiệu quả với bóng mờ.

#### 4.2.3. Xác định ROI



Các vùng quan tâm (ROI) được xác định dựa trên cấu trúc giao diện DAT cố định:

- student\_name: (40,440,380,500)
- student\_id: (40,500,380,560)
- vehicle\_plate: (65,600,310,670)
- instructor\_name: (40,560,380,620)
- distance\_completed: (600,440,800,500)
- time\_completed: (600,500,800,560)
- distance\_remaining: (600,560,800,620)
- time\_remaining: (600,620,800,680)
- total\_sessions: (600,680,800,740)

ROI được lưu dưới dạng dictionary trong main\_ocr.py, với tọa độ điều chỉnh  $\pm 10\%$  để xử lý biến thiên nhỏ. Ví dụ, với “dat\_001.jpg”, ROI vehicle\_plate bị nhiễu “Hang Het ha”, nhưng vẫn trích xuất được “51D19675 Hang Het ha” nhờ vùng cố định.

#### 4.2.4. Nhận dạng văn bản:

PaddleOCR (lang='vi', use\_angle\_cls=True) được sử dụng với cấu hình:

- Text Detection: DBNet với confidence threshold 0.3, tối ưu cho văn bản nhỏ.
- Text Recognition: CRNN với CTC loss, hỗ trợ diacritics tiếng Việt.
- Output: Dictionary {field: text} cho mỗi ảnh, lưu vào baseline\_results.csv.

Kết quả baseline cho thấy CER  $\sim 20\text{-}30\%$  do lỗi diacritics (e.g., “Nguyễn”  $\rightarrow$  “Nguyen”) và nhiễu UI (e.g., “51K27637 Hang Het ha”).

Listing 1: Mã giả cho SimpleOCRProcessor

```
def process_image(image_path, roi_dict):
    img = cv2.imread(image_path)
    img = resize_image(img, (1280, 720))
    img = apply_clahe(img, clip_limit=2.0)
    img = deskew_image(img, threshold=0.4)
    results = {}
    for field, (x1, y1, x2, y2) in roi_dict.items():
        roi_img = img[y1:y2, x1:x2]
        text = paddle_ocr(roi_img, lang='vi', use_angle_cls=True)
        results[field] = text
    return results
```

### 4.3. Module TextNormalizer

**Chức năng:** Chuẩn hóa văn bản thô từ SimpleOCRProcessor để loại bỏ artifacts và chuẩn bị cho correction. Module này áp dụng các quy tắc field-specific, dựa trên regex và diacritics removal.

#### 4.3.1. Quy tắc chuẩn hoá:

- student\_name, instructor\_name: Loại bỏ diacritics tạm thời (e.g., “á” → “a”) cho fuzzy matching, sử dụng thư viện unicodedata (normalize NFKD). Sau khi match, khôi phục diacritics từ GT nếu Jaccard similarity  $\geq 0.9$  [2].
- student\_id: Regex `r'[0-9-]'` loibpref ix (e.g., “MHV : 70835-... ” → “79035-...”). Validate định dạng số-hiến-dated.
- vehicle\_plate: Regex `r'[0-9A-Z]+'` để loại bỏ nhiễu (e.g., “51D19675 Hang Het ha” → “51D19675”). Chuẩn hóa theo format 00X00000.
- distance\_completed, distance\_remaining: Loại suffix “km” (regex `r'[0-9.]'`), convertsangfloat.time\_completed, time\_remaining: Chunhaformat “HH : MM” (e.g., “00: 9” → “00:09”) với regex `r'(1, 2) : (1, 2)'`.
- total\_sessions: Extract số từ “PHIÊN X” (regex `r'+'`).

#### 4.3.2. Ví dụ minh hoạ:

Trường	Raw OCR Output	Normalized Output
<i>student_name</i>	Nguyễn Anh Minh KHOA:K06A.C1	Nguyễn Anh Minh
<i>student_id</i>	MÃ HV: 70835-20250513082029960	79035-20250513082029960
<i>vehicle_plate</i>	51D19675 Hang Het ha	51D19675
<i>distance_completed</i>	508.5 km	508.5
<i>time_completed</i>	00:9	00:09
<i>total_sessions</i>	PHIÊN 9	9

#### 4.4. Module SimpleTextCorrectionML

Chức năng: Sửa lỗi OCR bằng học máy thống kê, sử dụng unsupervised approach với TF-IDF và Nearest Neighbors.

##### 4.4.1. Tạo dữ liệu huấn luyện

- Input: Baseline results (baseline\_results.csv) và GT (data.csv).
- Process: Tạo correction\_map (defaultdict) ánh xạ raw text → GT text cho mỗi trường. Ví dụ: {"Nguyễn Anh Minh": "Nguyễn Anh Minh", "51D19675 Hang Het ha": "51D19675"}.
- Output: Dataset (field, rawtext, gtext) ~900 entries (9 fields × 100 ảnh).

##### 4.4.2. Mô hình ML

- TF-IDF Vectorizer: Chuyển rawtext thành vector của n-gram (2-4), max\_features = 1000. Ví dụ: "Nguyễn" → ["Ng", "gu", "uy", "yê", "ên"].
- Nearest Neighbors: k=3, cosine similarity > 0.6, tìm GT text gần nhất. Nếu không tìm thấy, giữ normalized text.
- Hill-Climbing: Khám phá neighborhood edits (edit distance 1-2) để đề xuất corrections nếu similarity thấp.

Listing 3: Mã giả cho SimpleTextCorrectionML

```

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.neighbors import NearestNeighbors

def train_field_model(raw_texts, gt_texts):
    vectorizer = TfidfVectorizer(analyzer='char', ngram_range=(2,4),
                                max_features=1000)
    X = vectorizer.fit_transform(raw_texts)
    model = NearestNeighbors(n_neighbors=3, metric='cosine')
    model.fit(X)
    correction_map = dict(zip(raw_texts, gt_texts))
    return vectorizer, model, correction_map

def correct_text(field, text, vectorizer, model, correction_map):
    if text in correction_map:
        return correction_map[text]
    X = vectorizer.transform([text])
    distances, indices = model.kneighbors(X)
    if distances[0][0] < 0.4: # 1 - cosine similarity
        return correction_map[raw_texts[indices[0][0]]]
    return text

```

#### 4.4.3. Kết quả:

Module cải thiện accuracy từ 37.11% (baseline) lên 64.89% (enhanced), đặc biệt với numeric fields (e.g., distance\_completed 2% → 79%) nhờ loại bỏ suffix và nhiễu.

### 4.5. Module Pipeline Runners

- **BaselineOCRRunner** (*run\_baseline.py*): Chạy SimpleOCRProcessor, lưu baseline\_result tính accuracy so với GT.
- **TextCorrectionPipeline** (*run\_pipeline.py*): Tích hợp cả ba module, lưu enhanced\_results.csv và enhanced\_summary.txt.
- **Triển khai**: Chạy `python run_pipeline.py --input ../images/ --output ../output/`. Debug visualization với matplotlib (vẽ ROI trên ảnh).

## 5. Thực nghiệm và đánh giá:

### 5.1. Tổng quan thực nghiệm:

Thực nghiệm được thiết kế theo cách tiếp cận so sánh (comparative evaluation), tập trung vào hai phiên bản: baseline (sử dụng PaddleOCR thô với

normalization cơ bản) và enhanced (tích hợp text correction ML unsupervised). Mục tiêu chính là xác nhận cải thiện độ chính xác từ 37,11% lên 64,89%, đồng thời phân tích các yếu tố ảnh hưởng như lỗi diacritics, nhiễu UI, và chất lượng hình ảnh. Thực nghiệm được chạy trên môi trường Python 3.12 với thư viện PaddleOCR 2.7, scikit-learn 1.3, và OpenCV 4.8, trên máy chủ Ubuntu 22.04 (CPU Intel Xeon 16-core, RAM 64GB, không GPU để phù hợp hạ tầng thực tế).

Các thực nghiệm bao gồm: (1) đánh giá độ chính xác theo trường (field-specific accuracy), (2) phân tích thời gian xử lý và tài nguyên, (3) phân tích lỗi chi tiết (error analysis), và (4) so sánh với baseline methods từ tài liệu. Kết quả được lưu dưới dạng CSV (baseline\_results.csv, enhanced\_results.csv) và summary text, với visualization debug để hỗ trợ phân tích.

## 5.2. Kết quả thực nghiệm:

### 5.2.1. độ chính xác theo trường

Bảng 5.1 tóm tắt kết quả so sánh baseline và enhanced trên 100 ảnh. Enhanced cải thiện tổng thể +27,78%, với gains lớn ở numeric fields (e.g., distance\_completed +77%).

Trường	Baseline	Enhanced	Cải thiện
student_name	51.0 ± 2.1	51.0 ± 1.8	0.0
student_id	15.0 ± 3.4	49.0 ± 2.5	+34.0
vehicle_plate	45.0 ± 4.2	66.0 ± 3.1	+21.0
instructor_name	82.0 ± 1.5	82.0 ± 1.2	0.0
distance_completed	2.0 ± 1.1	79.0 ± 2.8	+77.0
time_completed	69.0 ± 3.0	70.0 ± 2.7	+1.0
distance_remaining	0.0 ± 0.0	77.0 ± 3.2	+77.0
time_remaining	58.0 ± 2.9	59.0 ± 2.4	+1.0
total_sessions	12.0 ± 2.6	51.0 ± 3.5	+39.0
<b>Overall</b>	<b>37.11 ± 2.8</b>	<b>64.89 ± 2.1</b>	<b>+27.78</b>

Ví dụ cụ thể từ “dat\_000.jpg”: Baseline student\_id “V: 70835-...” (sai, prefix + digit error), enhanced “79035-...” (đúng, +1 match). vehicle\_plate baseline “51K27637” (đúng), nhưng với nhiễu “Hang Het ha” ở “dat\_001.jpg”, enhanced loại bỏ thành “51D19675” (đúng, regex + TF-IDF).

### 5.2.2. Thời gian xử lý và tài nguyên:

Bảng 5.2 cho thấy enhanced chỉ tăng overhead nhỏ (0.81s), chủ yếu từ OCR serial (6.67s/ảnh). RAM peak ~2GB, CPU utilization 80%.

Giai đoạn	Baseline	Enhanced
OCR	6.67	6.67
Normalization	0.01	0.01
Correction	-	0.008
Evaluation	0.02	0.02
Total	6.70	6.71

Parallelization test (multiprocessing 4 workers): Giảm OCR time 40% (4s/ảnh), tổng 4.03s/ảnh. Tài nguyên: Không vượt 4GB RAM ngay cả với batch 10 ảnh.

## 6. Kết luận và hướng phát triển

### 6.1. Kết luận

Nghiên cứu đã thành công phát triển pipeline OCR hybrid cho trích xuất dữ liệu từ ảnh giao diện DAT, nâng độ chính xác từ 37,11% (baseline) lên 64,89% (enhanced), với cải thiện nổi bật ở numeric fields (e.g., distance\_completed +77%). Kết quả xác nhận hiệu quả của unsupervised correction (TF-IDF + Nearest Neighbors) trong xử lý diacritics và nhiễu UI tiếng Việt, hỗ trợ tự động hóa báo cáo tiến độ đào tạo lái xe theo quy định 2025. Hệ thống lightweight, không yêu cầu GPU, phù hợp triển khai thực tế tại trung tâm đào tạo.

### 6.2. Hướng phát triển

- Mở rộng dataset lên 1.000 ảnh để giảm variance và fine-tune PaddleOCR.
- Tích hợp LLM (e.g., T5) cho diacritics correction, target accuracy >80%.
- Parallelization đầy đủ và API deployment (FastAPI) cho batch processing.
- Áp dụng multimodal learning (hình ảnh + văn bản) cho bleed-over errors.

## Tài liệu tham khảo

- [1] X. B. C. Y. Baoguang Shi, An End-to-End Trainable Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition, IEEE, 2016.
- [2] C. L. R. G. X. Y. W. L. J. Z. Y. B. Z. Y. Y. Q. D. H. W. Yuning Du, "PP-OCR: A Practical Ultra Lightweight OCR System," *Computer Vision and Pattern Recognition*, 2020.
- [3] C. M. P. R. H Schütze, Introduction to information retrieval, 2008.
- [4] V. Tyagi, Understanding Digital Image Processing, 2018.
- [5] R. Szeliski, Computer Vision: Algorithms and Applications, Springer, 2022.
- [6] R. Brunelli, Template matching techniques in computer vision: theory and practice, 2009.
- [7] D. Lowe, "Distinctive Image Features from Scale-Invariant Keypoints.," *International journal of computer vision*, 2004.
- [8] Q.-D. Nguyen, N.-M. Phan, P. Krömer và D.-A. Le, "An efficient unsupervised approach for ocr error correction of vietnamese ocr text," *IEEE Access*, tập 11, 2023.
- [9] D.-A. L. I. Z. Quoc-Dung Nguyen, "OCR error correction for unconstrained Vietnamese handwritten text," *SoICT*, pp. 132 - 138 , 2019.