

# Signboard text detection and recognition in streaming video

Nguyễn Đình Quân - 20521184, Nguyễn Hùng Phát - 22521074

December 15, 2025

# LỜI CẢM ƠN

Trước tiên, em xin gửi lời cảm ơn chân thành và sâu sắc đến Ban Giám hiệu nhà trường và Khoa Khoa học Máy tính đã tạo điều kiện học tập và nghiên cứu thuận lợi trong suốt thời gian em theo học tại Trường Đại học Công nghệ Thông tin.

Em xin bày tỏ lòng biết ơn đặc biệt đến Thầy Đỗ Văn Tiến, đã trực tiếp giảng dạy và tận tình hướng dẫn em trong quá trình thực hiện đề tài khóa luận. Những định hướng, chỉ dẫn rõ ràng cùng sự hỗ trợ quý báu từ thầy đã là tiền đề quan trọng giúp em hoàn thành tốt công việc nghiên cứu và viết báo cáo đúng tiến độ. Em cũng xin cảm ơn thầy vì đã cung cấp tài liệu, giải đáp thắc mắc và luôn tạo môi trường học tập tích cực, hiệu quả.

Trong suốt quá trình thực hiện đề tài, em đã có cơ hội vận dụng những kiến thức nền tảng đã được học, đồng thời tích cực học hỏi, tìm tòi thêm các kiến thức mới. Đây là một trải nghiệm quý báu giúp em trưởng thành hơn trong tư duy và kỹ năng làm việc nghiên cứu.

Mặc dù đã nỗ lực hoàn thành đề tài với tinh thần nghiêm túc và cầu thị, nhưng do hạn chế về thời gian và kinh nghiệm, khóa luận không thể tránh khỏi những thiếu sót. Em rất mong nhận được sự thông cảm, góp ý chân thành từ các thầy cô để em có thể tiếp tục hoàn thiện và phát triển trong tương lai.

Em xin chân thành cảm ơn!

# TÓM TẮT KHÓA LUẬN

aaaaa.....

# Contents

<b>LỜI CẢM ƠN</b>	<b>i</b>
<b>Tóm tắt khóa luận</b>	<b>ii</b>
<b>Contents</b>	<b>iii</b>
<b>List of Figures</b>	<b>v</b>
<b>List of Tables</b>	<b>vi</b>
<b>1 TỔNG QUAN</b>	<b>1</b>
1.1 Đặt vấn đề . . . . .	1
1.2 Mục tiêu và phạm vi . . . . .	5
1.2.1 Mục tiêu . . . . .	5
1.2.2 Phạm vi . . . . .	6
1.3 Đóng góp của khóa luận . . . . .	6
1.4 Cấu trúc khóa luận . . . . .	7
<b>2 CƠ SỞ LÝ THUYẾT VÀ CÁC NGHIÊN CỨU LIÊN QUAN</b>	<b>8</b>
2.1 Giới thiệu . . . . .	8
2.1.1 Tổng quan và ý nghĩa thực tiễn của bài toán đọc văn bản trên biểu hiệu trong ảnh/video đường phố . . . . .	8
2.1.2 Thách thức về tính đa dạng và phức tạp của văn bản trong môi trường tự nhiên . . . . .	9

2.1.2.1	Bài toán Scene Text Detection and Recognition . . .	9
2.1.2.2	Phân loại hướng tiếp cận cho bài toán Text Detection and Recognition . . . . .	10
2.2	Cơ sở lý thuyết . . . . .	10
2.2.1	Biểu diễn dữ liệu video và trích xuất khung hình . . . . .	10
2.2.2	Quy trình xử lý tổng thể (Integrated Pipeline) . . . . .	11
2.2.3	Cơ sở lý thuyết Text Detection . . . . .	12
2.2.4	Cơ sở lý thuyết Text Recognition . . . . .	12
2.2.5	Khai thác thông tin thời gian trong video . . . . .	12
2.3	Các nghiên cứu liên quan . . . . .	13
2.3.1	Bộ dữ liệu văn bản trong video lái xe: RoadText-1K . . . . .	13
2.3.2	Phương pháp phát hiện text theo cơ chế “spotlight”: Spotlight Text Detector (STD) . . . . .	13
2.3.3	Liên hệ với đề tài khóa luận . . . . .	14
2.4	Tóm tắt chương . . . . .	14
<b>3</b>	<b>PHƯƠNG PHÁP</b>	<b>17</b>
3.1	Hệ thống phát hiện và nhận dạng chữ trên biển hiệu . . . . .	17
<b>4</b>	<b>THỰC NGHIỆM VÀ ĐÁNH GIÁ</b>	<b>18</b>
4.1	Dữ liệu . . . . .	18
4.2	Tiền xử lý . . . . .	18
4.3	Tập câu truy vấn đánh giá . . . . .	18
4.4	Độ đo đánh giá . . . . .	18
4.5	Kết quả thực nghiệm . . . . .	18
<b>5</b>	<b>KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN</b>	<b>19</b>
5.1	Kết luận . . . . .	19
5.2	Hướng phát triển . . . . .	19

# List of Figures

1.1	Minh họa đầu vào và đầu ra của hệ thống. Hệ thống tiếp nhận video và trả về danh sách biển hiệu cùng nội dung văn bản nhận dạng, cũng như các đoạn video/khung hình liên quan. . . . .	4
2.1	Phân nhóm các phương pháp cho bài toán Text Detection and Recognition trong ảnh ngoại cảnh. . . . .	16

# List of Tables

# Chapter 1

## TỔNG QUAN

### 1.1 Đặt vấn đề

Phát hiện và nhận dạng văn bản trong ảnh đời thường (*Scene Text Detection and Recognition* – STDR) là một bài toán quan trọng trong thị giác máy tính, thu hút nhiều sự quan tâm nhờ các ứng dụng rộng rãi như dịch tự động, hỗ trợ dẫn đường, số hóa tài liệu ngoài trời, hay phân tích biển báo giao thông. Với đầu vào là ảnh tĩnh hoặc khung hình video, STDR hướng tới việc xác định vị trí xuất hiện của văn bản và trích xuất chính xác nội dung văn bản đó (Hình 1).

Trong số các dạng văn bản đời thường, **văn bản trên biển hiệu** (*Signboard Text*) có ý nghĩa đặc biệt vì thường chứa thông tin về *tên địa điểm, cơ sở kinh doanh, dịch vụ* hoặc *định danh không gian* trong môi trường đô thị. Do vậy, **phát hiện và nhận dạng văn bản biển hiệu** (*Signboard Text Detection and Recognition*) trở thành một nhánh quan trọng của STDR, có nhiều tiềm năng ứng dụng trong hệ thống dẫn đường thông minh, phân tích thông tin đô thị và xây dựng bản đồ số.

Tuy nhiên, STDR nói chung gặp nhiều thách thức do sự đa dạng của phong chữ, kích thước, hướng và bố cục; văn bản có thể bị nghiêng, cong, chồng chéo hoặc hòa lẫn trong nền phức tạp, cùng các phong cách thiết kế nghệ thuật và đa ngôn ngữ (Hình 2). Đối với tiếng Việt, khó khăn còn lớn hơn do hệ thống dấu (, ; ^ ~ , dấu hỏi, dấu nặng) và các nguyên âm biến thể (ô, ê, â, ă, ơ, ư), làm tăng số lượng ký tự cần nhận dạng và dễ gây nhầm lẫn giữa các chữ có hình dạng gần giống (ví dụ *a* với *â*, *ă*, *á*).

Bên cạnh đó, biển hiệu cũng đa dạng về hình dạng, kích thước, vật liệu và thường xuất hiện ở các vị trí phức tạp trong ảnh (Hình 3), như bị che khuất một phần, bị phản xạ ánh sáng, hoặc nằm trong các bối cảnh đông đúc. Theo hiểu biết của chúng tôi, đến nay mới chỉ có một nghiên cứu [2] tập trung vào phát hiện biển hiệu trên đường phố Việt Nam, trong khi hướng kết hợp *cả phát hiện đối tượng biển hiệu và nhận dạng nội dung văn bản trên biển hiệu* vẫn còn ít được khai thác.

Ngoài ra, khi mở rộng từ ảnh tĩnh sang **video hành trình** (dashcam), bài toán còn đối mặt với các thách thức đặc thù như mờ do chuyển động, độ phân giải hạn chế của camera hành trình, cùng sự thay đổi liên tục về ánh sáng và góc nhìn. Những yếu tố này khiến việc phát hiện và nhận dạng văn bản trong video trở nên khó khăn hơn so với ảnh đơn lẻ.

Trong phạm vi khóa luận này, chúng tôi tiến hành khảo sát và tổng hợp các hướng tiếp cận liên quan đến STDR và các phương pháp hiện đại [3,4,5,6]. Trên cơ sở đó, chúng tôi lựa chọn, cài đặt, thực nghiệm và đánh giá một số phương pháp tiên tiến [7,8,9,10] trên tập dữ liệu SignboardText [1]. Đồng thời, tập dữ liệu này được **mở rộng** bằng cách **bổ sung nhãn đối tượng biển hiệu** (thay vì chỉ nhãn văn bản). Tiếp theo, chúng tôi áp dụng các phương pháp vào dữ liệu video camera hành trình trên đường phố Việt Nam, với đầu vào là khung hình chứa biển hiệu và đầu ra là *vị trí biển hiệu cùng nội dung văn bản* tương ứng (Hình 4). Từ văn bản trích xuất được, hệ thống hướng tới việc phát triển ứng dụng minh họa, chẳng hạn như xác định loại hình cơ sở (cửa hàng, nhà hàng, trường học, bệnh viện...), hỗ trợ tìm kiếm và truy xuất thông tin.

## Bài toán và pipeline đề xuất

Khóa luận tập trung vào pipeline phát hiện và nhận dạng văn bản biển hiệu trong video, bao gồm các mô-đun chính:

- **Trích xuất khung hình và tiền xử lý:** lấy mẫu khung hình từ video, hiệu chỉnh cơ bản (nếu cần) nhằm giảm nhiễu, mờ chuyển động và sai lệch ánh sáng.
- **Phát hiện biển hiệu (Signboard Detection):** xác định vùng chứa biển hiệu trong khung hình (dạng hộp bao hoặc đa giác), làm vùng quan tâm (ROI).

- **Phát hiện văn bản (Text Detection):** trong ROI biển hiệu, phát hiện vùng văn bản (theo hộp thẳng hoặc hộp xoay) để tăng độ chính xác.
- **Nhận dạng văn bản (Text Recognition):** nhận dạng chuỗi ký tự tiếng Việt/đa ngôn ngữ từ các vùng văn bản phát hiện được.
- **Hậu xử lý và hợp nhất theo thời gian:** loại nhiễu, chuẩn hóa chuỗi, và hợp nhất kết quả giữa các khung hình liên tiếp (tracking/temporal voting) để ổn định đầu ra.
- **Tầng ứng dụng:** khai thác văn bản nhận dạng để tìm kiếm/truy xuất theo từ khóa hoặc danh mục; cung cấp nền tảng cho các tác vụ downstream.

Như vậy, **đầu vào (Input)** của hệ thống bao gồm (xem thêm hình minh họa Hình 1.1):

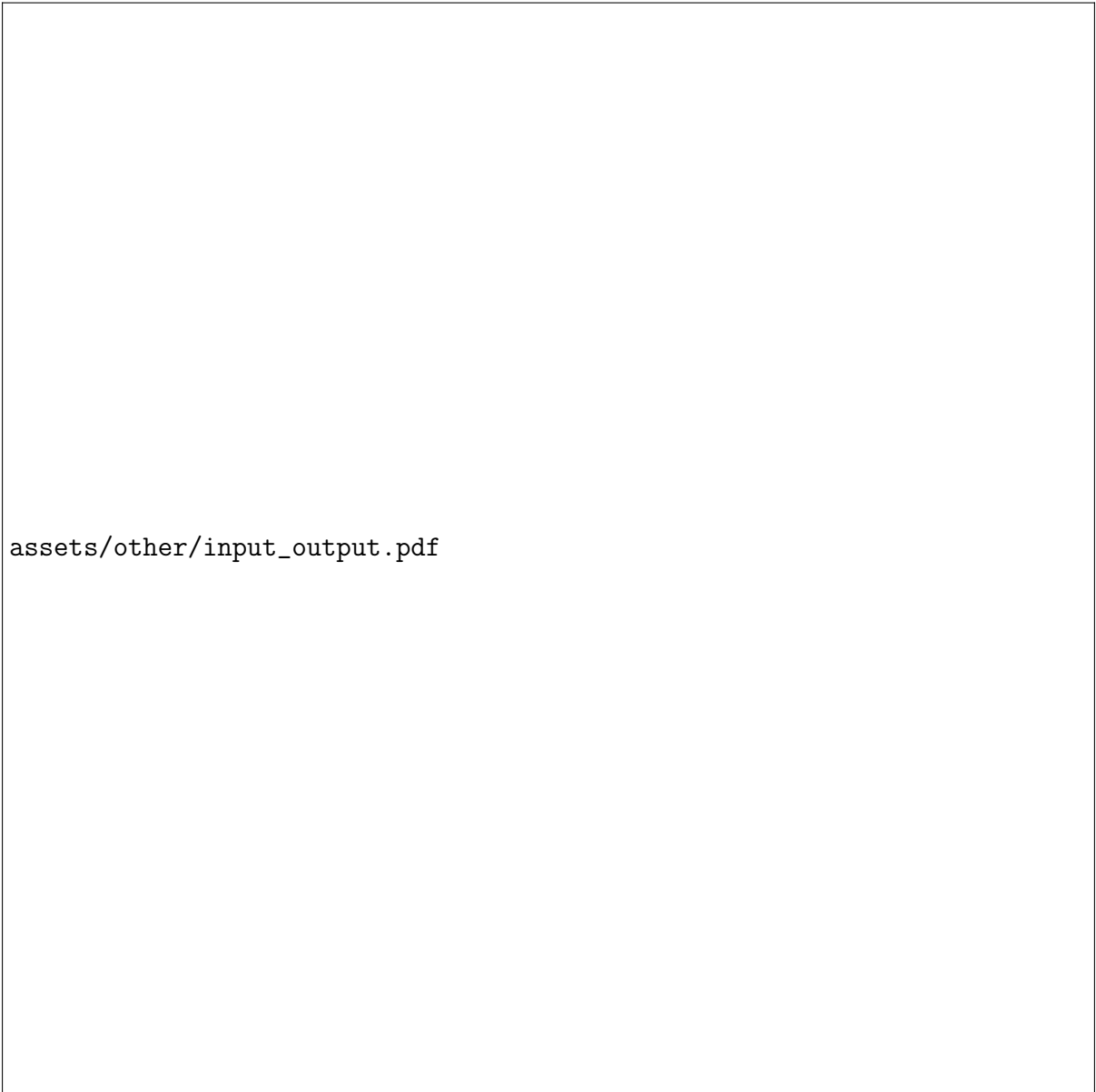
- **Video hành trình đầu vào:** video camera hành trình trên đường phố Việt Nam, chứa các cảnh có biển hiệu ở nhiều điều kiện (ban ngày/ban đêm, mưa/nắng, đông người/ít người, nhiều góc nhìn).

Và **đầu ra (Output)** của hệ thống là:

- **Danh sách các phát hiện biển hiệu:** mỗi mục gồm thời điểm (timestamp), vị trí biển hiệu (bounding box/polygon) và nội dung văn bản nhận dạng.
- **Danh sách các đoạn video/khung hình liên quan:** các đoạn video ngắn hoặc khung hình đại diện chứa biển hiệu phù hợp với truy vấn, được sắp xếp theo mức độ liên quan.

Mặc dù đã có nhiều tiến bộ trong lĩnh vực phát hiện và nhận dạng văn bản, bài toán phát hiện và nhận dạng văn bản biển hiệu trong video hành trình vẫn tồn tại nhiều thách thức đáng kể:

1. **Biến thiên điều kiện chụp và chất lượng ảnh:** mờ do chuyển động, độ phân giải thấp, nhiễu nén video và thay đổi ánh sáng làm giảm khả năng phát hiện và nhận dạng.



assets/other/input\_output.pdf

Figure 1.1: Minh họa đầu vào và đầu ra của hệ thống. Hệ thống tiếp nhận video và trả về danh sách biển hiệu cùng nội dung văn bản nhận dạng, cũng như các đoạn video/khung hình liên quan.

2. **Đa dạng hình dạng biển hiệu và bố cục văn bản:** biển hiệu có thể cong, nghiêng, bị che khuất, nhiều lớp thông tin (logo, biểu tượng, chữ nghệ thuật), gây khó khăn cho cả detection và recognition.
3. **Đặc thù tiếng Việt và đa ngôn ngữ:** dấu tiếng Việt làm tăng độ phức tạp của bộ ký tự và dễ gây nhầm lẫn; thực tế biển hiệu có thể pha trộn Việt–Anh hoặc ký tự đặc biệt.

Từ những thách thức đã nêu, khóa luận này đặt ra mục tiêu phát triển một **pipeline phát hiện và nhận dạng văn bản biển hiệu trong video hành trình**, có khả năng:

- Xác định chính xác vùng biển hiệu trong khung hình video và phát hiện vùng văn bản tương ứng;
- Nhận dạng văn bản trên biển hiệu ổn định theo thời gian và hỗ trợ khai thác thông tin cho tác vụ tìm kiếm/truy xuất.

## 1.2 Mục tiêu và phạm vi

### 1.2.1 Mục tiêu

Trong khóa luận này, sinh viên đề ra các mục tiêu như sau:

- Khảo sát và tổng hợp các hướng tiếp cận tiên tiến để giải quyết bài toán STDR và bài toán văn bản biển hiệu.
- Mở rộng tập dữ liệu SignboardText [1] bằng cách bổ sung nhãn đối tượng biển hiệu (*signboard*).
- Cài đặt, thực nghiệm và đánh giá một số phương pháp hiện đại; phân tích ưu/nhược điểm của từng phương pháp.
- Xây dựng pipeline phát hiện và nhận dạng văn bản trên biển hiệu trong video.

- Phát triển ứng dụng minh họa khai thác thông tin văn bản từ biển hiệu, chẳng hạn như xác định loại hình cơ sở (cửa hàng, nhà hàng, trường học, bệnh viện. . .), nhằm hỗ trợ tìm kiếm và truy xuất theo từ khóa hoặc danh mục; đồng thời cung cấp nền tảng cho các tác vụ downstream như gợi ý địa điểm hoặc phân loại dịch vụ.

### 1.2.2 Phạm vi

Trong khóa luận này, nhóm sinh viên tập trung hoàn thành các công việc sau:

- Tìm hiểu tổng quan, thách thức và cơ sở lý thuyết của các phương pháp phát hiện biển hiệu, phát hiện và nhận dạng văn bản trong ảnh đời thường.
- Mở rộng tập dữ liệu SignboardText [1] bằng cách bổ sung nhãn biển hiệu; đồng thời thu thập thêm dữ liệu video hành trình từ camera hành trình trên đường phố Việt Nam.
- Cài đặt, thực nghiệm và đánh giá một số phương pháp hiện đại trên tập dữ liệu đã mở rộng; phân tích ưu/nhược điểm của từng phương pháp.
- Xây dựng pipeline phát hiện và nhận dạng văn bản trên biển hiệu trong video trên đường phố Việt Nam.
- Phát triển ứng dụng minh họa trên nền tảng web, cho phép khai thác thông tin văn bản từ biển hiệu nhằm hỗ trợ tìm kiếm, truy xuất thông tin và cung cấp dữ liệu đầu vào cho các hệ thống thông minh khác.

## 1.3 Đóng góp của khóa luận

Các đóng góp chính của khóa luận bao gồm:

- **Mở rộng bộ dữ liệu:** bổ sung nhãn đối tượng biển hiệu cho SignboardText [1] và xây dựng tập dữ liệu video hành trình phục vụ đánh giá pipeline.

- **Thực nghiệm và phân tích:** cài đặt và đánh giá nhiều phương pháp hiện đại cho các mô-đun (phát hiện biển hiệu/văn bản, nhận dạng văn bản), kèm phân tích ưu/nhược điểm theo bối cảnh tiếng Việt.
- **Pipeline và ứng dụng minh họa:** đề xuất pipeline STDR cho biển hiệu trong video và phát triển ứng dụng web hỗ trợ tìm kiếm/truy xuất theo từ khóa hoặc danh mục từ văn bản biển hiệu.

## 1.4 Cấu trúc khóa luận

Nội dung khóa luận được tổ chức như sau:

**Chương 1:** Tổng quan bài toán, bối cảnh, động lực, mục tiêu, phạm vi và đóng góp.

**Chương 2:** Cơ sở lý thuyết và các nghiên cứu liên quan đến phát hiện biển hiệu, phát hiện/nhận dạng văn bản và các kỹ thuật xử lý video.

**Chương 3:** Các phương pháp và pipeline đề xuất cho bài toán phát hiện và nhận dạng văn bản biển hiệu trong video, bao gồm mô tả kiến trúc hệ thống và mô-đun xử lý.

**Chương 4:** Thực nghiệm và đánh giá trên tập dữ liệu SignboardText mở rộng và dữ liệu video hành trình; phân tích kết quả và thảo luận.

**Chương 5:** Xây dựng ứng dụng minh họa và mô tả các chức năng khai thác thông tin văn bản biển hiệu.

**Chương 6:** Kết luận và hướng phát triển trong tương lai.

## Chapter 2

# CƠ SỞ LÝ THUYẾT VÀ CÁC NGHIÊN CỨU LIÊN QUAN

## 2.1 Giới thiệu

### 2.1.1 Tổng quan và ý nghĩa thực tiễn của bài toán đọc văn bản trên biển hiệu trong ảnh/video đường phố

Trong môi trường giao thông đô thị, biển hiệu và bảng quảng cáo (signboards) là nguồn thông tin quan trọng phản ánh danh tính địa điểm (tên cửa hàng), cũng như loại sản phẩm/dịch vụ mà địa điểm đó cung cấp. Với dữ liệu video quay từ camera hành trình, hệ thống “đọc văn bản trong cảnh” (scene text reading) có thể hỗ trợ nhiều ứng dụng thực tế như: (i) lập bản đồ/định vị theo ngữ nghĩa (semantic mapping), (ii) tìm kiếm địa điểm theo từ khóa (place search), (iii) thống kê loại hình kinh doanh theo khu vực, và (iv) hỗ trợ nhận thức tình huống trong các hệ thống giao thông thông minh.

Bài toán trong khóa luận tập trung vào xây dựng một pipeline tích hợp gồm:

- **Phát hiện (Text Detection):** xác định và khoanh vùng các vùng chứa văn bản trong từng khung hình.
- **Nhận dạng (Text Recognition):** chuyển đổi ảnh vùng chữ thành chuỗi ký tự.

### 2.1.2 Thách thức về tính đa dạng và phức tạp của văn bản trong môi trường tự nhiên

Khác với tài liệu quét (scanned documents), văn bản trong cảnh đường phố thường xuất hiện trong điều kiện chụp không kiểm soát. Các thách thức nổi bật bao gồm:

- **Nền phức tạp (cluttered background):** nhiều vật thể/hoa văn gây nhiễu.
- **Văn bản biến dạng:** chữ cong/ngiên/méo do phối cảnh, bề mặt biến hiệu hoặc góc nhìn.
- **Đa dạng phong chữ và kích thước:** font stylized, độ dày nét khác nhau, chữ rất nhỏ hoặc rất lớn.
- **Đa ngôn ngữ và dấu:** tiếng Việt có dấu, có thể xen kẽ tiếng Anh/Trung/Hàn; dấu câu đa dạng.
- **Motion blur và out-of-focus:** do xe di chuyển, rung camera, tốc độ cao.
- **Độ phân giải thấp (low resolution):** chữ nhỏ, ở xa camera, nén video làm mất chi tiết.

#### 2.1.2.1 Bài toán Scene Text Detection and Recognition

Tổng quát, bài toán Scene Text Reading có thể được tiếp cận theo ba hướng chính:

1. **Text Detection:** chỉ dự đoán vị trí vùng chữ (bounding box / polygon / mask).
2. **Text Recognition:** nhận dạng ký tự/chuỗi ký tự từ các vùng chữ đã được cắt sẵn.
3. **End-to-End Text Spotting/Recognition:** kết hợp phát hiện và nhận dạng trong một pipeline thống nhất.

Trong khóa luận, trọng tâm là xây dựng pipeline tích hợp cho dữ liệu street-view/video, ưu tiên nhóm đối tượng biển hiệu/bảng quảng cáo cửa hàng.

### 2.1.2.2 Phân loại hướng tiếp cận cho bài toán Text Detection and Recognition

Các phương pháp học sâu cho bài toán đọc văn bản trong ảnh ngoại cảnh (scene text reading) thường được phân nhóm theo phạm vi xử lý: (i) chỉ phát hiện vùng chữ, (ii) chỉ nhận dạng chữ trên vùng cắt sẵn, hoặc (iii) pipeline end-to-end kết hợp cả hai.

#### Text Detection.

- **Regression-based:** hồi quy trực tiếp hộp/tứ giác bao quanh vùng chữ.
- **Connected component-based:** phát hiện thành phần ký tự (hoặc stroke) rồi liên kết thành dòng/từ.
- **Segmentation-based:** dự đoán bản đồ pixel thuộc text, sau đó tách instance bằng hậu xử lý.

#### Text Recognition.

- **Segmentation-based:** tách ký tự (hoặc vùng con) rồi nhận dạng.
- **Segmentation-free:** nhận dạng trực tiếp chuỗi (CTC/attention/transformer) không cần tách ký tự.

#### End-to-End Text Recognition.

- **One-stage:** phát hiện và nhận dạng trong một mô hình thống nhất.
- **Two-stage:** phát hiện trước, sau đó cắt/chuẩn hóa và nhận dạng ở mô-đun thứ hai.

## 2.2 Cơ sở lý thuyết

### 2.2.1 Biểu diễn dữ liệu video và trích xuất khung hình

Video được xem như chuỗi khung hình (frame) theo thời gian. Cho video  $V$ , ta trích xuất tập khung hình  $\{I_t\}_{t=1}^T$  với tốc độ lấy mẫu phù hợp (ví dụ: lấy mọi frame hoặc lấy

theo bước nhảy để tối ưu tính toán). Vì văn bản có thể xuất hiện trong nhiều frame liên tiếp, dữ liệu video mang *tính dư thừa theo thời gian* (temporal redundancy) có thể khai thác để tăng độ ổn định.

### 2.2.2 Quy trình xử lý tổng thể (Integrated Pipeline)

Pipeline đề xuất ở mức khái niệm gồm các bước:

#### 1. Tiền xử lý (Pre-processing):

- Giảm nhiễu, cân bằng sáng, tăng tương phản cục bộ khi cần thiết.
- Giảm mờ do chuyển động (deblurring) hoặc tăng độ phân giải (super-resolution) cho vùng chữ nhỏ (tùy tài nguyên).
- Ổn định video (video stabilization) trong trường hợp rung mạnh.

2. **Phát hiện vùng văn bản (Text Detection):** Dự đoán vị trí vùng chữ theo dạng hộp (box), tứ giác (quadrilateral), đa giác (polygon) hoặc mặt nạ (segmentation mask). Đầu ra gồm tập vùng  $\mathcal{B}_t = \{b_t^{(i)}\}$  tại frame  $I_t$ .

3. **Chuẩn hóa hình học và cắt vùng chữ (Crop & Rectify):** Với vùng chữ nghiêng/cong, cần biến đổi phối cảnh hoặc chuẩn hóa hình học để đưa về ảnh chữ “thẳng” (rectified) trước khi nhận dạng. Gọi ảnh vùng chữ sau chuẩn hóa là  $\hat{I}_t^{(i)}$ .

4. **Nhận dạng văn bản (Text Recognition):** Mô hình nhận dạng thực hiện ánh xạ  $\hat{I}_t^{(i)} \rightarrow s_t^{(i)}$ , trong đó  $s_t^{(i)}$  là chuỗi ký tự dự đoán.

#### 5. Hậu xử lý (Post-processing):

- Chuẩn hóa Unicode tiếng Việt, sửa lỗi dấu/telex nếu cần.
- Loại bỏ ký tự nhiễu, lọc theo độ tin cậy (confidence).
- Gộp các kết quả theo thời gian (temporal fusion) nếu cùng một biểu hiện xuất hiện ở nhiều frame.

6. **Suy luận loại dịch vụ/sản phẩm (Semantic Inference):** Từ chuỗi ký tự  $s$ , hệ thống gán nhãn ngành hàng/dịch vụ bằng: (i) luật từ khóa (keyword rules), (ii) phân lớp văn bản (text classification), hoặc (iii) kết hợp NER + taxonomy.

### 2.2.3 Cơ sở lý thuyết Text Detection

Text Detection trong ảnh có thể chia thành các hướng chính:

- **Regression-based:** dự đoán trực tiếp hộp/tứ giác bao quanh text.
- **Connected-component based:** phát hiện thành phần ký tự và liên kết thành cụm.
- **Segmentation-based:** dự đoán bản đồ pixel thuộc vùng text và tách instance bằng hậu xử lý.

Trong thực tế signboards, hướng segmentation-based phổ biến vì linh hoạt với chữ cong/biến dạng, nhưng gặp khó khi các cụm chữ nằm gần nhau gây chồng lấp.

### 2.2.4 Cơ sở lý thuyết Text Recognition

Text Recognition thường được mô hình hóa như bài toán nhận dạng chuỗi:

- **CTC-based:** ánh xạ đặc trưng theo chiều ngang thành chuỗi ký tự với CTC loss.
- **Attention/Encoder-Decoder:** sinh chuỗi theo cơ chế chú ý.
- **Transformer-based recognizer:** tận dụng self-attention để học phụ thuộc dài và chống méo tốt hơn.

Với tiếng Việt, các yếu tố dấu và biến thể font làm tăng độ khó; hậu xử lý chuẩn hóa và từ điển miền (domain lexicon) có thể cải thiện độ chính xác.

### 2.2.5 Khai thác thông tin thời gian trong video

Khác ảnh tĩnh, video cho phép:

- **Tracking text regions:** theo dõi một vùng chữ qua nhiều frame, giảm số lần chạy recognizer (nhận dạng một lần cho cả track).
- **Temporal fusion:** hợp nhất nhiều kết quả nhận dạng theo vote/confidence/edit distance để tăng độ ổn định.

## 2.3 Các nghiên cứu liên quan

### 2.3.1 Bộ dữ liệu văn bản trong video lái xe: RoadText-1K

RoadText-1K giới thiệu bộ dữ liệu lớn cho bài toán phát hiện và nhận dạng văn bản trong video lái xe, gồm các đoạn clip được lấy mẫu từ dữ liệu lái xe thực tế và được gán nhãn dày (dense) theo từng frame. Bộ dữ liệu cung cấp: (i) bounding boxes cho vùng chữ, (ii) phiên âm (transcription), và (iii) nhãn phân loại text (ví dụ: tiếng Anh/không phải tiếng Anh/không đọc được), đồng thời tách riêng trường hợp biển số xe trong nhóm tiếng Anh. RoadText-1K được thiết kế theo hướng “không thiên lệch theo text” (unconstrained), phản ánh đúng bối cảnh camera hành trình với các nhiễu như motion blur, out-of-focus và glare. Các đánh giá baseline cho thấy các phương pháp SOTA trên ảnh tĩnh khi áp dụng vào video lái xe sẽ gặp suy giảm đáng kể do độ khó tăng lên.

### 2.3.2 Phương pháp phát hiện text theo cơ chế “spotlight”: Spotlight Text Detector (STD)

Spotlight Text Detector (STD) tập trung giải quyết hai vấn đề lớn của text detection dạng segmentation: (i) các instance chữ nằm gần nhau gây chồng lấp khó tách, và (ii) hình dạng/độ dài chữ biến thiên lớn khiến mô hình khó khái quát. STD đề xuất hai thành phần chính:

- **Spotlight Calibration Module (SCM):** hiệu chỉnh vùng ứng viên (candidate kernel) dựa trên coarse mask, tương tự cơ chế camera “focus” vào mục tiêu; module này giúp giảm false positives bằng cách hiệu chỉnh dự đoán và tăng khả năng tập trung vào vùng kernel quan trọng.

- **Multivariate Information Extraction Module (MIEM):** trích xuất thông tin hình học đa dạng theo nhiều “shape schemes”, nhằm học tốt hơn các đặc trưng tỷ lệ, hướng và hình dạng của chữ trong cảnh.

Kết quả thực nghiệm cho thấy STD đạt hiệu năng cạnh tranh/vượt trội trên nhiều benchmark text detection phổ biến (ICDAR2015, CTW1500, MSRA-TD500, Total-Text), đồng thời ablation chứng minh đóng góp của SCM và MIEM.

### 2.3.3 Liên hệ với đề tài khóa luận

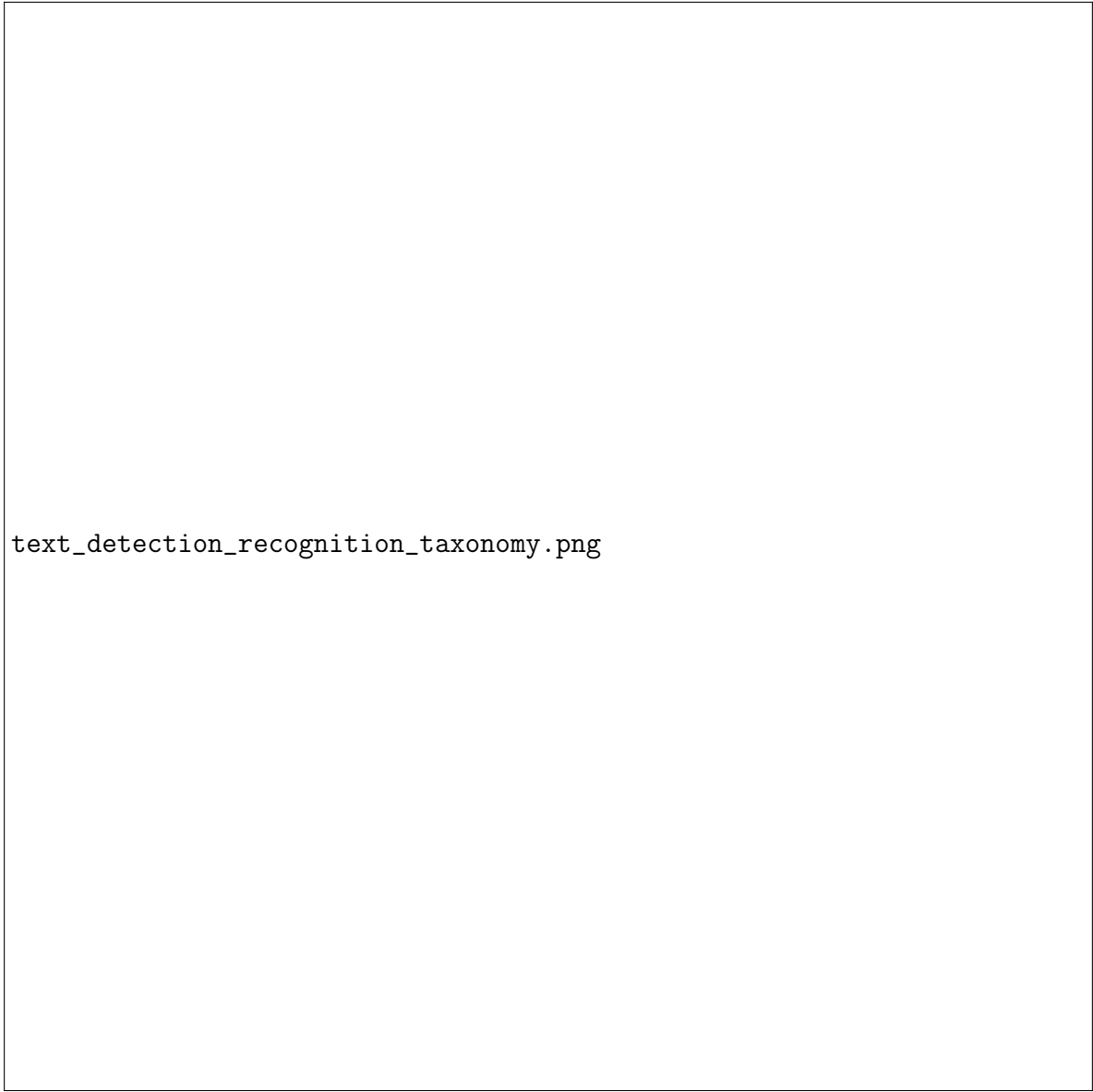
Từ các nghiên cứu trên, có thể rút ra các định hướng quan trọng cho bài toán biển hiệu trong video street-view:

- **Về dữ liệu và đánh giá:** cần ưu tiên bối cảnh “unconstrained driving video” và các dạng nhiễu đặc thù; RoadText-1K là nguồn tham khảo về cách thiết kế dữ liệu/nhãn và tiêu chí benchmark.
- **Về phát hiện văn bản:** các phương pháp segmentation nâng cao cơ chế hiệu chỉnh/khoanh vùng (như SCM của STD) hữu ích khi text gần nhau và nền phức tạp — đặc trưng thường gặp ở biển hiệu phố.
- **Về pipeline tích hợp:** cần kết hợp (i) detection mạnh với chữ cong/biến dạng, (ii) rectify hợp lý trước recognition, và (iii) cơ chế temporal fusion/tracking để ổn định kết quả trên video.

## 2.4 Tóm tắt chương

Chương này đã trình bày tổng quan bài toán đọc văn bản trên biển hiệu trong video đường phố, các thách thức đặc thù, cơ sở lý thuyết của hai thành phần chính (text detection và text recognition), cùng khả năng khai thác thông tin thời gian trong video. Ngoài ra, chương cũng tổng hợp hai hướng nghiên cứu liên quan tiêu biểu: bộ dữ liệu RoadText-1K cho video lái xe và phương pháp Spotlight Text Detector cho phát hiện chữ

với cơ chế hiệu chỉnh vùng ứng viên. Những nội dung này là cơ sở để thiết kế pipeline tích hợp và xây dựng thực nghiệm trong các chương tiếp theo.



text\_detection\_recognition\_taxonomy.png

Figure 2.1: Phân nhóm các phương pháp cho bài toán Text Detection and Recognition trong ảnh ngoại cảnh.

## **Chapter 3**

# **PHƯƠNG PHÁP**

### **3.1 Hệ thống phát hiện và nhận dạng chữ trên biển hiệu**

## **Chapter 4**

# **THỰC NGHIỆM VÀ ĐÁNH GIÁ**

### **4.1 Dữ liệu**

### **4.2 Tiền xử lý**

### **4.3 Tập câu truy vấn đánh giá**

### **4.4 Độ đo đánh giá**

### **4.5 Kết quả thực nghiệm**

## Chapter 5

# KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN

### 5.1 Kết luận

.....:

- aaaa.

### 5.2 Hướng phát triển

Để khắc phục những hạn chế trên và nâng cao hơn nữa tính hiệu quả, tính khả dụng và tính mở rộng của hệ thống, các hướng phát triển trong tương lai được đề xuất như sau:

**Tối ưu hóa khả năng mở rộng dữ liệu:**

- aaa
- aaa

**Tăng cường khả năng tương tác và thích ứng với người dùng:**

- Thiết kế giao diện người dùng aaaaaaaaaaaaaa

**Tích hợp truy vấn hình thức thoại (Spoken Query Integration):** Phát triển hệ thống

.....