

Nhận dạng ký tự quang học

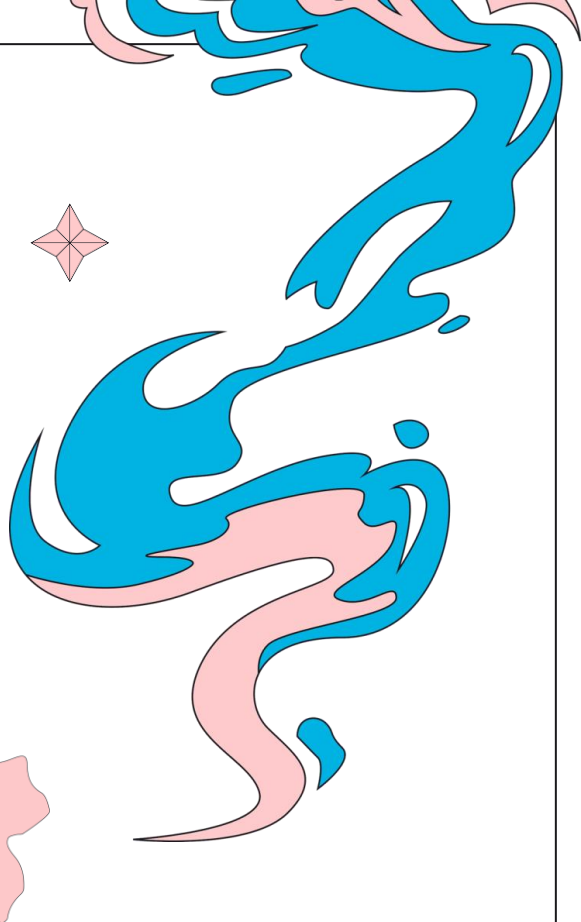
Tài liệu cổ

Chữ Nôm

字
喃

chữ

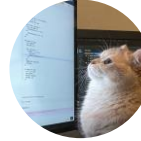
Nôm



Thành viên



- 19127078 - Nguyễn Đỗ Thanh Trúc



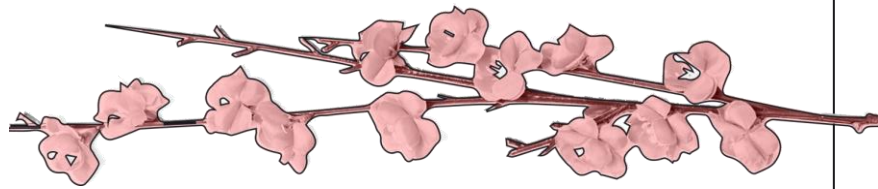
- 19127135 - Phạm Bảo Hân



Giảng viên hướng dẫn

- Đinh Điền
- Lương An Vinh
- Lê Thị Thúy Hằng

Mục lục



◆ **01** **Chữ Nôm**
Một số đặc điểm của chữ Nôm

◆ **02** **OCR**
Giới thiệu OCR

◆ **03** **Sửa lỗi văn bản**
Giới thiệu bài toán Sửa lỗi văn bản

◆ **04** **Bài báo chính**
Ý tưởng của bài báo chính và mô hình ngôn ngữ sử dụng

◆ **05** **Quy trình**
Giới thiệu về quy trình, dữ liệu và khó khăn

Đ



01



Chữ Nôm

Nôm character

1. Tổng quan:

- Chữ Nôm được hình thành và phát triển trong khoảng từ thế kỷ 10 - thế kỷ 19.
- Hàn Thuyên được công nhận là người có công truyền bá chữ Nôm.
- Không thể thống kê được có bao nhiêu chữ Nôm.
- Có ít nhất 32,695 chữ Nôm theo *Hội Bảo tồn Di sản chữ Nôm*.



Hàn Thuyên

Nôm character

2. Đặc điểm

- Mượn từ Hán Việt (Sino-Vietnamese)

- Chữ **sơn** nghĩa là ‘ngọn núi’ được viết 山 (phát âm: shān)

- Sáng tạo bởi người Việt Nam

- Mượn 2 từ Hán: 1 để diễn tả âm thanh, 1 để diễn tả ý

- Từ **lửa** nghĩa là ‘ngọn lửa’ được viết 𤇀
(gồm 2 thành phần: 火 (nghĩa: lửa), 呂 (phát âm: lửa))

- Mượn từ Hán để diễn tả âm thanh

- Từ **một** nghĩa là ‘số một’ được viết 𣎵 (phát âm: mò)

- Mượn 2 từ Hán để diễn tả ý nghĩa

- Từ **trời** nghĩa là ‘bầu trời’ được viết 𠂇 |
(gồm 2 thành phần: 天 (nghĩa: thiên đường), 上 (nghĩa: phía trên))

Nôm character

3. Hệ thống chữ viết

- Ghi âm tiếng mẹ đẻ.
=> Ý thức tự chủ ngôn ngữ muốn có chữ viết riêng của dân tộc ta
- Lưu giữ lịch sử văn hóa và di sản vô giá của dân tộc Việt Nam.
(văn học, triết học, lịch sử, luật pháp, y học, tôn giáo,...).
- Với sự ra đời của chữ Quốc Ngữ, chữ Nôm dần dần biến mất.
=> **Vấn đề cấp bách:** Bảo tồn di sản, văn hóa của dân tộc Việt Nam.



慕辭勳揆些

1 Trăm năm trong cõi người ta.

荊才荊命窖口怙鏡

Chữ tài chữ mệnh khéo là ghét nhau.

駛戈沒局液桃

Trải qua một cuộc bể dâu,

仍調口口苞口痘口

Những điều trông thấy đã đau đớn lòng.

還之彼晉斯豐

5 Là gì bỉ sắc tư phong,

忒攢涓貝騰紅打攄

Trời xanh quen với má hồng đánh ghen.

稿黃吝口趨烟

Cảo thơm lần giở trước đèn,

風情固錄群傳史撐

Phong tình có lục còn truyền sử xanh.

浪辭口靖朝明

Rằng: Năm Gia Tĩnh triều Minh,

眾方榜湖江京凭鑽

10 Bốn phương phẳng lặng, hai kinh vững

vàng.

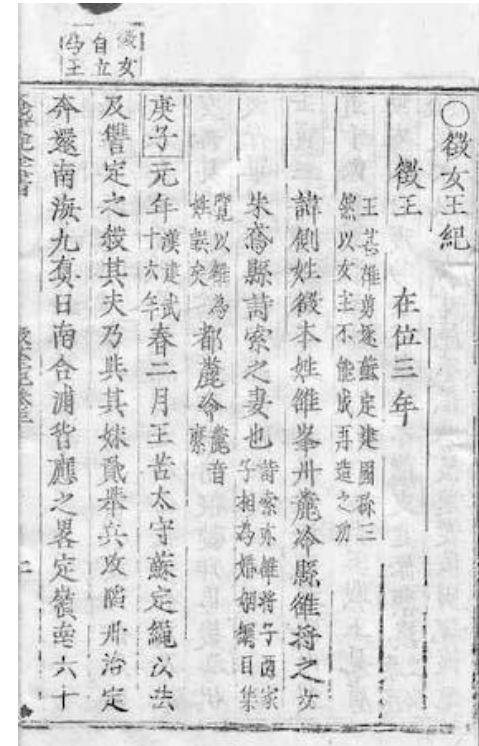
固茹員外户王

Truyện Kiều

Nôm character

4. Ý nghĩa:

- Ghi âm tiếng mẹ đẻ.
=> Ý thức tự chủ ngôn ngữ muốn có chữ viết riêng của dân tộc ta
- Lưu giữ lịch sử văn hóa và di sản vô giá của dân tộc Việt Nam.
(văn học, triết học, lịch sử, luật pháp, y học, tôn giáo,...).
- Với sự ra đời của chữ Quốc Ngữ, chữ Nôm dần dần biến mất.
=> **Vấn đề cấp bách:** Bảo tồn di sản, văn hóa của dân tộc Việt Nam.



Bộ Đại Việt sử ký toàn thư
Kỷ Trưng Nữ Vương

02

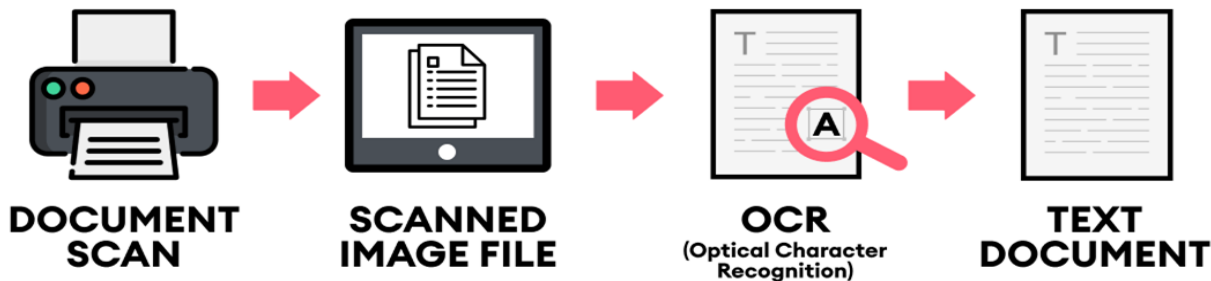
✦ Nhận dạng ký tự quang học (OCR)



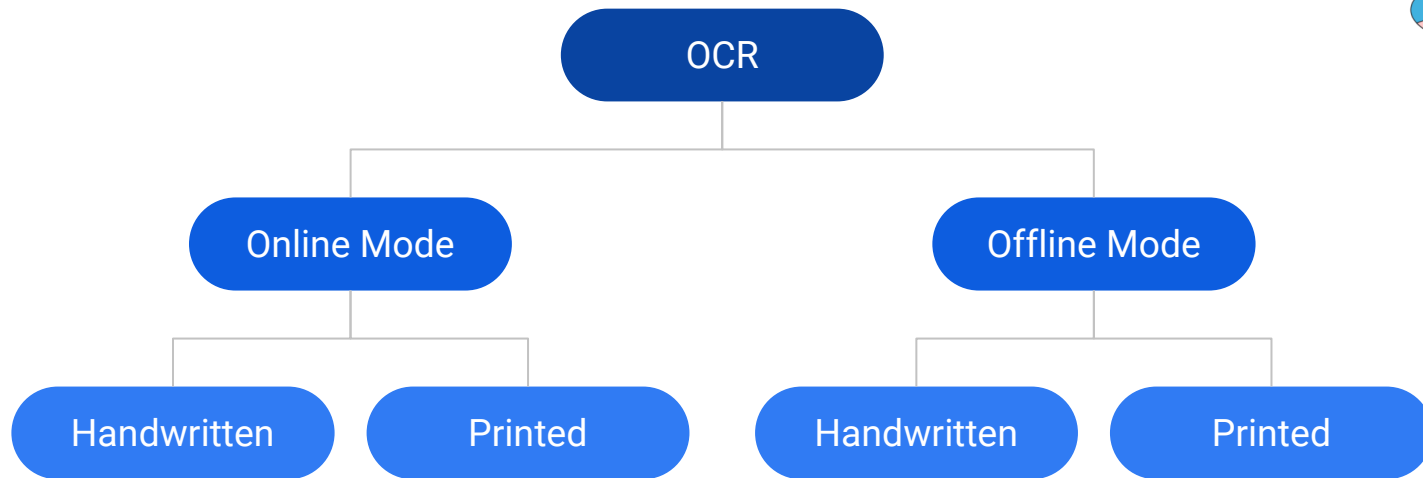
I. Định nghĩa



Nhận diện ký tự quang học là quá trình tự động nhận dạng và trích xuất các ký tự trong ảnh (ảnh văn bản viết tay, văn bản in hoặc quét), sau đó chuyển đổi thành định dạng kỹ thuật số có thể chỉnh sửa.



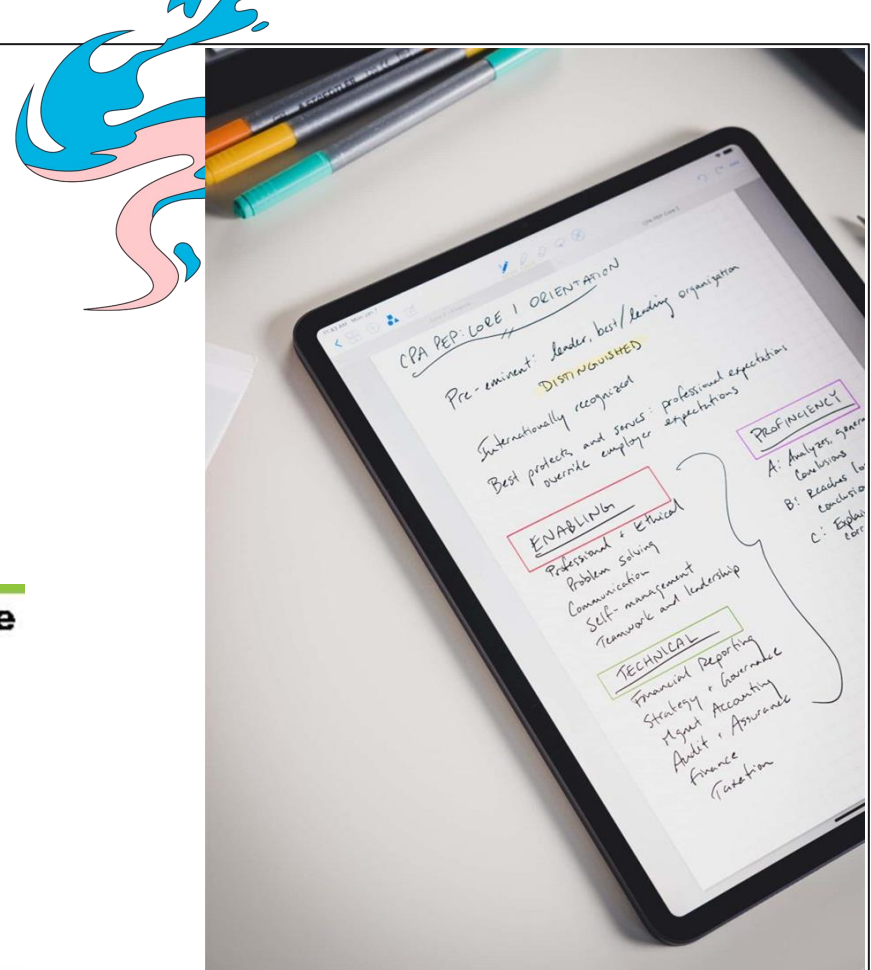
II. Phân loại



III. Ứng dụng

Online OCR

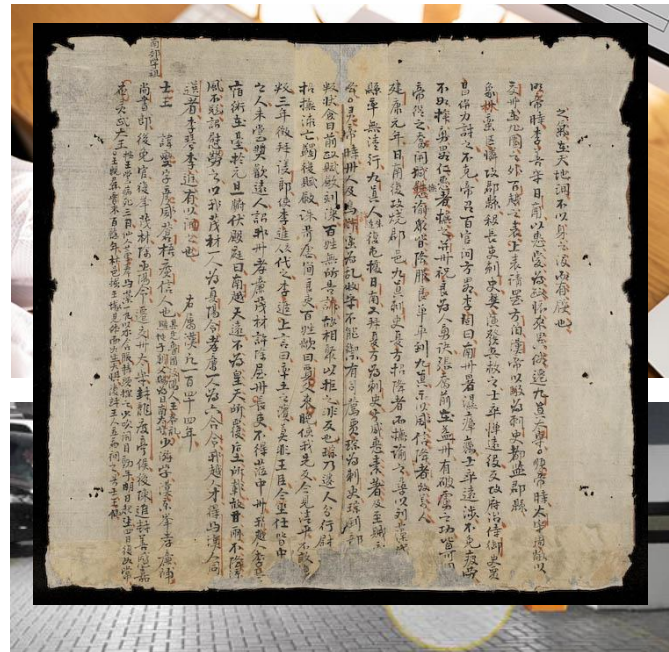
- Bút cảm ứng dùng cho smart phone, tablet,...
- Xác minh chữ ký
- Xác thực danh tính khách hàng



III. Ứng dụng

Offline OCR

- Chuyển đổi văn bản
- Nhận diện biển số xe
- **Số hóa các tài liệu cổ:** Những tài liệu cổ chứa đựng một lượng lớn thông tin và kiến thức lịch sử quan trọng. Việc số hóa các tài liệu này là phương pháp tốt nhất để bảo tồn chúng. Các bước chính của quá trình này bao gồm:
 - Scanning
 - OCR
 - Lưu trữ dữ liệu

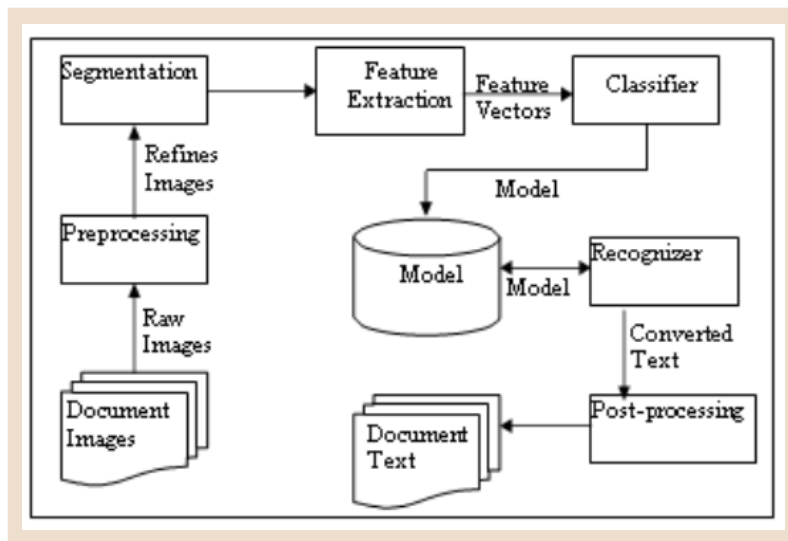




IV. Các bước thực hiện



- 3 Phân tách
- 2 Tiền xử lý
- 1 Thu nhận ảnh



Rút trích
đặc trưng 4

Phân loại
ký tự 5

Hậu xử lý 6



IV. Các bước thực hiện



Thu nhận ảnh



Tiền xử lý



Phân tách

被 歪 吹 尉 術 茄 咈 民

Bị
trời **thối**
đuổi
về
nhà
dạy
dân

Rút trích đặc trưng
+
Phân loại ký tự

被 歪 欧 尉 術 茄 咈 民

Bị
trời
xua
đuổi
về
nhà
dạy
dân

Hậu xử lý



03

**Sửa lỗi
văn bản**



Sửa lỗi văn bản

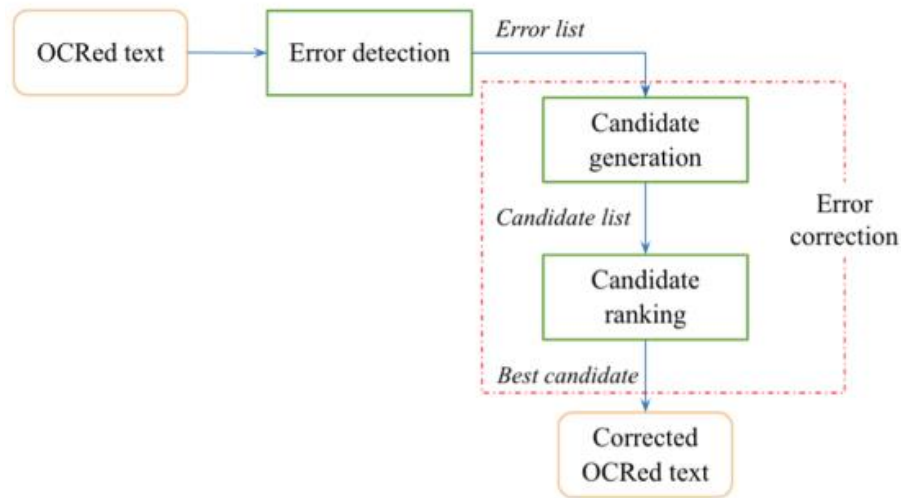


Quy trình:

- Phát hiện lỗi: xác định các token lỗi từ văn bản đầu vào.
- Sửa lỗi: khắc phục các token bị sai (danh sách các từ ứng viên).

Yêu cầu:

- Chính xác, không sinh ra lỗi mới.
- Độ phức tạp về thời gian và bộ nhớ thấp.





Phân loại theo loại lỗi



- Lỗi chính tả (Spelling error): từ sai chính tả

VD: ngành (ngành), cộ sát (cộ sát)

Error Type	Example
Cognitive error:	I don't know the correct spelling of <u>Levenstain</u> distance.
Typographic:	<u>THis</u> sentence was typed in <u>haser</u> .
Typographic (OCR):	<u>SUPpLEMENTAhy</u> <u>INFOhMATION</u> .
Typographic (Diacritic):	The authors of this article are Daniel <u>Hladek</u> , <u>Matus</u> Pleva and <u>Jan</u> <u>Stas</u> .
Note: The spelling errors are underlined.	

- Lỗi ngữ pháp (Grammatical error): từ đúng chính tả nhưng sai về mặt ngữ nghĩa trong câu

Type	Example
IWC	自己有双聪明能干的手，什么都能做出来。 You have smart hands to do everything. (Tips: "hands" cannot be combined with "smart")
CM	绿色植物具有产生氧气。 Plants have (the ability) to produce oxygen. (Tips: Lack of object "the ability")
CR	我们已走了约十里左右的路程。 We had walked about 10 miles or so. (Tips: "about" and "or so" are redundant)

	交通事故发生的原因是开车看手机造成的。 Traffic accidents are caused by (because) looking at cell phones while driving. (Tips: the structure of "because" and "caused by" cannot appear together in one sentence)
SC	我改正并认识了自己的错误。 I corrected and realized my fault. (Tips: realize the fault first and correct it later)
IWO	我们应该防止事故不发生。 We should prevent accidents from not occurring. (Tips: double negation causes illogical errors)
ILL	刚一开门，看病的就进来了。 As the door opened, the doctor/patient came in. (Tips: there is an ambiguity about who comes in)

- Sai trật tự từ (IWO)
- Sai cách kết hợp từ (IWC)
- Thiếu thành phần (CM)
- Dư thành phần (CR)
- Nhầm lẫn cấu trúc (SC)
- Phi logic (Phi logic, ILL)
- Nhập nhằng (AM)



Phương pháp



Thủ công

- Độ chính xác cao
- Tốn kém về chi phí, thời gian
- Phụ thuộc vào trình độ của người review
- Cần người theo dõi, đánh giá hiệu suất công việc
- Không phù hợp với các văn bản có bản quyền





Phương pháp



- *Non-word errors: từ sai chính tả và không có trong từ điển (từ vô nghĩa)*
- *Real-word errors: từ được dùng sai ngữ cảnh nhưng có trong từ điển (từ có nghĩa)* Vd: *He is string (strong).*

Dựa trên luật

- Tốc độ xử lý nhanh.
- Sửa lỗi về non-word errors, một số real-word errors đơn giản.
- Thiếu sự uyển chuyển, đặc biệt đối với các tác phẩm văn học.
- Cần xây dựng bộ từ điển đủ lớn, các bộ luật đủ mạnh, văn phạm không nhập nhằng, ...
- Ví dụ:
 - Tra cứu từ điển dùng hash table, search tree
 - Thuật toán phân tích cú pháp Earley Parser

Dựa trên máy học:

- Tốc độ xử lý nhanh
- Sửa nhiều loại lỗi (real-word and non-word errors)
- Chi phí về việc huấn luyện mô hình (ngữ liệu huấn luyện, chi phí huấn luyện: thời gian, GPU, ...)
- Phức tạp hơn các mô hình theo hướng dựa trên luật
- Học ngữ cảnh của văn bản => phù hợp cho nhiều dạng văn bản khác nhau
- Ví dụ:
 - Mô hình ngôn ngữ dạng học sâu: CNN, Transformer, ...
 - Mô hình ngôn ngữ Ngram

04

Bài báo

**TOWARDS NÔM HISTORICAL DOCUMENT
OPTICAL CHARACTER RECOGNITION**

Tran Thi Anh Thu, Le Pham Ngoc Yen, Tran Thai Son, Dinh Dien



Tổng quan

End-to-end model: Tesseract OCR, YOLOv5

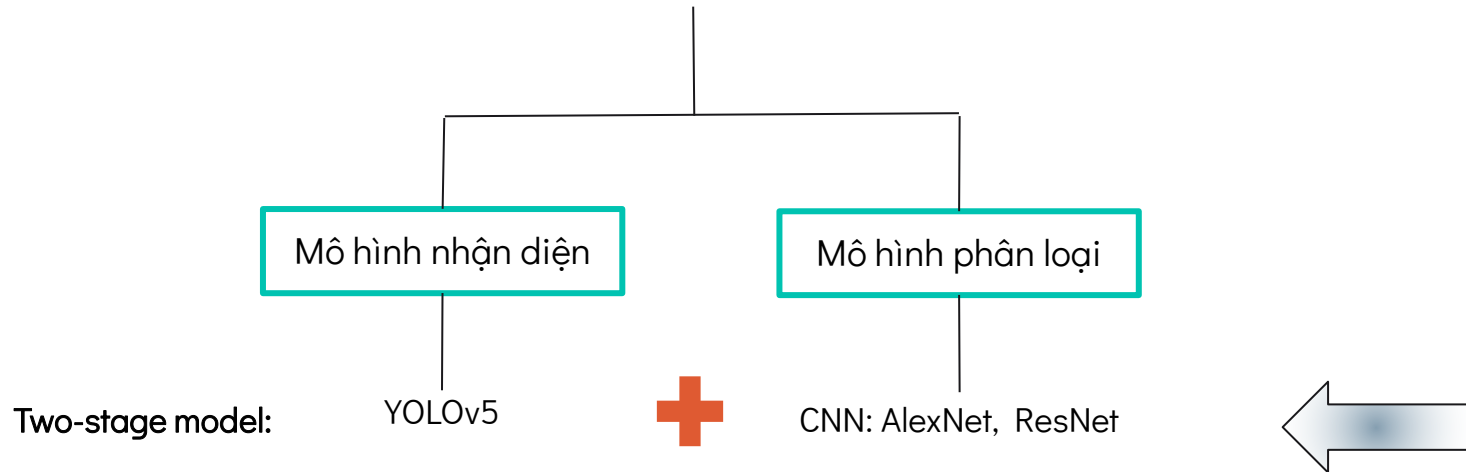


Table 5. Final top-1 and top-5 validation (mAP@0.5, %) on Nôm datasets of the three OCR baselines

Model	NomTrain	NomTest	
	K1871	K1902	LVT
End-to-end YOLO	0.0618	0.062	
YOLO-AlexNet	98.47 - 99.14	59.02 - 68.76	59.78 - 70.60
YOLO-ResNet101	98.91 - 99.32	71.80 - 80.54	72.96 - 82.36

Tổng quan

Mô hình ngôn ngữ

MacBERT: MLM as correction BERT

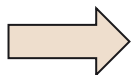
- MacBERT là mô hình dựa trên mô hình RoBERTa, nhưng được chỉnh sửa để dùng cho tiếng Trung.
- Thay vì sử dụng ký hiệu [MASK], MacBERT dùng các từ có nghĩa tương đương để thay thế. Mô hình sẽ tự động nhận diện chữ được mask, dự đoán và sửa đổi thành chữ ban đầu.

Original Sentence

使用语言模型来预测下一个词的概率。

we use a language model to predict the probability of the next word.

	Chinese	English
Original Masking	使用语言[M]型来[M]测下一个词的概率。	we use a language [M] to [M] ##di ##ct the pro [M] ##bility of the next word .
+ WWM	使用语言[M] [M]来[M] [M]下一个词的概率。	we use a language [M] to [M] [M] [M] the [M] [M] [M] of the next word .
++ N-gram Masking	使用[M] [M] [M] [M]来[M] [M]下一个词的概率。	we use a [M] [M] to [M] [M] [M] the [M] [M] [M] [M] [M] next word .
+++ Mac Masking	使用语法建模来预见下一个词的几率。	we use a text system to ca ##lc ##ulate the po ##si ##bility of the next word .



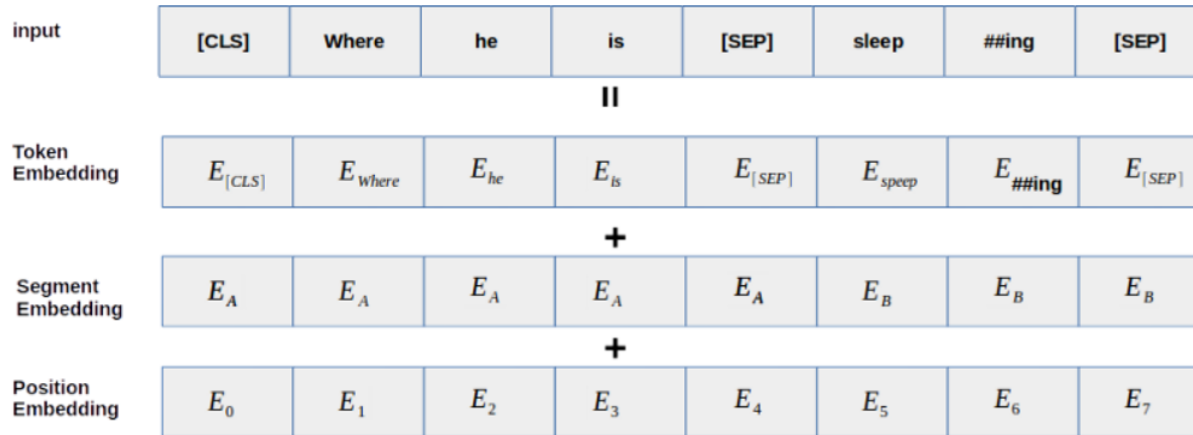
Phù hợp với bài toán Sửa lỗi văn bản cho chữ Nôm do cơ chế masking và sự tương đồng giữa chữ Nôm và Hán tự.

Tổng quan

Mô hình ngôn ngữ

MacBERT: MLM as correction BERT

- MacBERT sẽ hoạt động dựa trên các embeddings là Token, Segment, Position.





05

**Quy
trình**

Dữ liệu

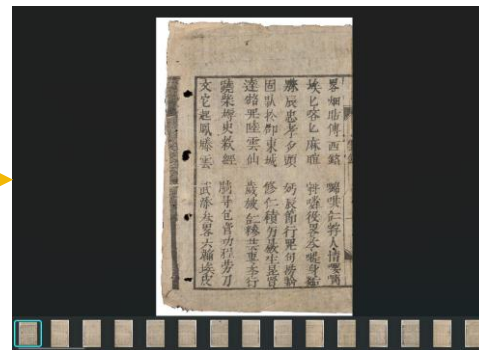
	OCR	Sửa lỗi văn bản
Chú thích	Có	Có
Đơn vị ngôn ngữ	Mẫu tự	Văn bản hoặc câu
Bình diện	Hình thái	Ngữ pháp Ngữ nghĩa
Tag set	Tất cả các mẫu tự chữ Nôm	Tất cả các mẫu tự chữ Nôm

Dữ liệu đầu vào

OCR

images

annotation.csv



	A	B	C	D	E	F	G	H
1	label	x1	x2	y1	y2	source_im	img_w	img_h
2	7567	427	462	149	190	nlvnpf-005	493	760
3	7551	426	464	192	227	nlvnpf-005	493	760
4	4021	425	462	226	262	nlvnpf-005	493	760
5	50B3	426	463	261	302	nlvnpf-005	493	760
6	897F	426	464	304	339	nlvnpf-005	493	760
7	9298	426	464	338	374	nlvnpf-005	493	760
8	2110E	424	463	410	446	nlvnpf-005	493	760
9	552D	426	463	445	482	nlvnpf-005	493	760
10	20129	426	463	481	515	nlvnpf-005	493	760
11	21982	426	463	513	551	nlvnpf-005	493	760

Dữ liệu đầu ra

OCR

	A	B	C	D	E	F	G	H
1	label	x1	x2	y1	y2	source_im	img_w	img_h
2	7567	427	462	149	190	nlvpnf-005	493	760
3	7551	426	464	192	227	nlvpnf-005	493	760
4	4021	425	462	226	262	nlvpnf-005	493	760
5	5083	426	463	261	302	nlvpnf-005	493	760
6	897F	426	464	304	339	nlvpnf-005	493	760

畧 畑 貼 傳 西 銘
 啼 嘒 仁 莽 人 情 嚶 鳴
 埃 匕 略 匕 麻 暄 身 姦
 符 隣 役 畧 苓 嶷 身 姦
 棘 辰 忠 孝 口 頭
 媽 辰 節 行 口 句 撈 輸
 固 馱 於 郡 東 城
 修 仁 積 口 口 生 昆 賢
 達 旄 口 陸 雲 仙
 歲 放 仁 糝 芸 專 孝 行
 蹕 柴 燭 史 欸 經
 腦 尋 包 口 功 程 勞 刀
 文 它 起 鳳 滕 雲
 武 添 叁 畧 六 韜 埃



(b) *Truyện Kiều* (1902) page 10

Dữ liệu

Sửa lỗi văn bản

Training

	A	B	C
1	original_text	correct_text	wrong_ids
2	黔蔣擣撓欺疊培	黔忝坦常欺疊培	3 4 2
3	客騰神容餒迅狹	客騰紅魑餒迅遁	3 7 4
4	撐箕唵溜層蓮	撐箕唵溜層蓮	
5	蕊埃吟篇朱欄餒尼	爲埃該孕朱誠餒尼	1 4 6 3
6	徽長城龍羣月	徽長城龍羣月	
7	燒洪泉須曠式靈	燒甘泉曠曠式靈	2 4
8	鸞吝鎌宝掾栢	鸞吝鎌宝掾栢	1
9	姘姑傳徽定疢織征	姘姑傳徽定軺出征	6 7
10	掙清平匹稟粹黷	諾清平匹稟粹黷	1

Inference

畧 煩 貼 傳 西 銘
嚙 嘆 極 筭 人 情 嚙 嚙
埃 匕 咯 匕 麻 聃

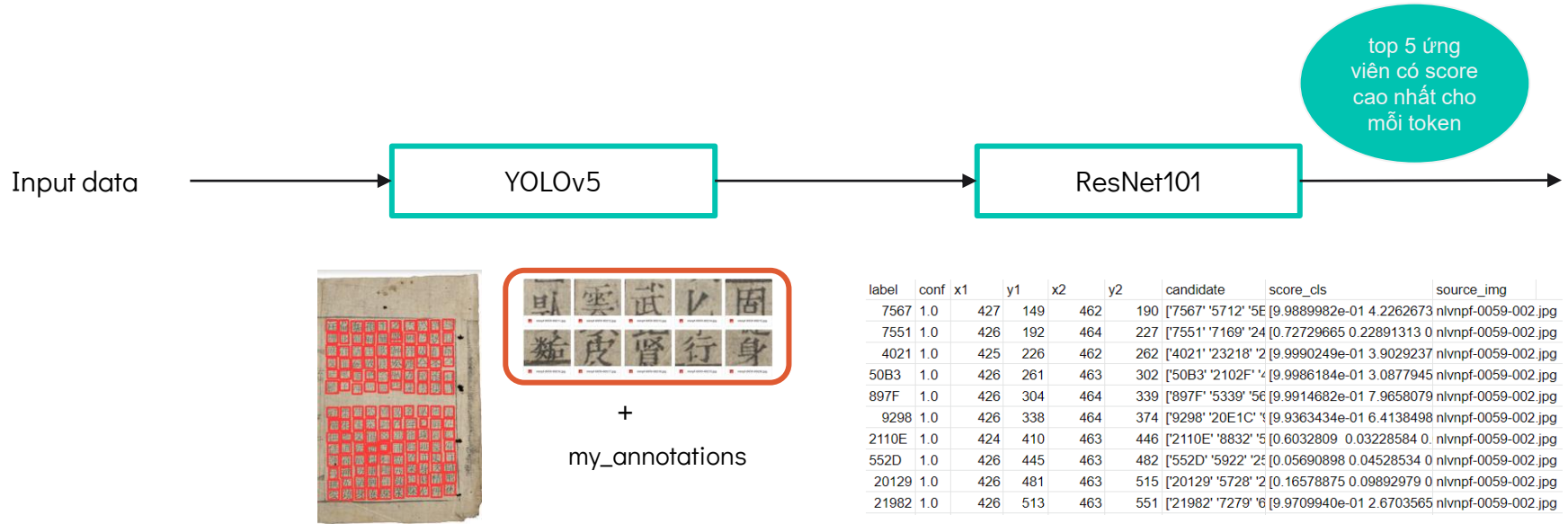


畧 焮 貼 傳 西 銘
嚙 嚙 仁 筭 人 情 嚙 嚙
埃 匕 咯 匕 麻 聃

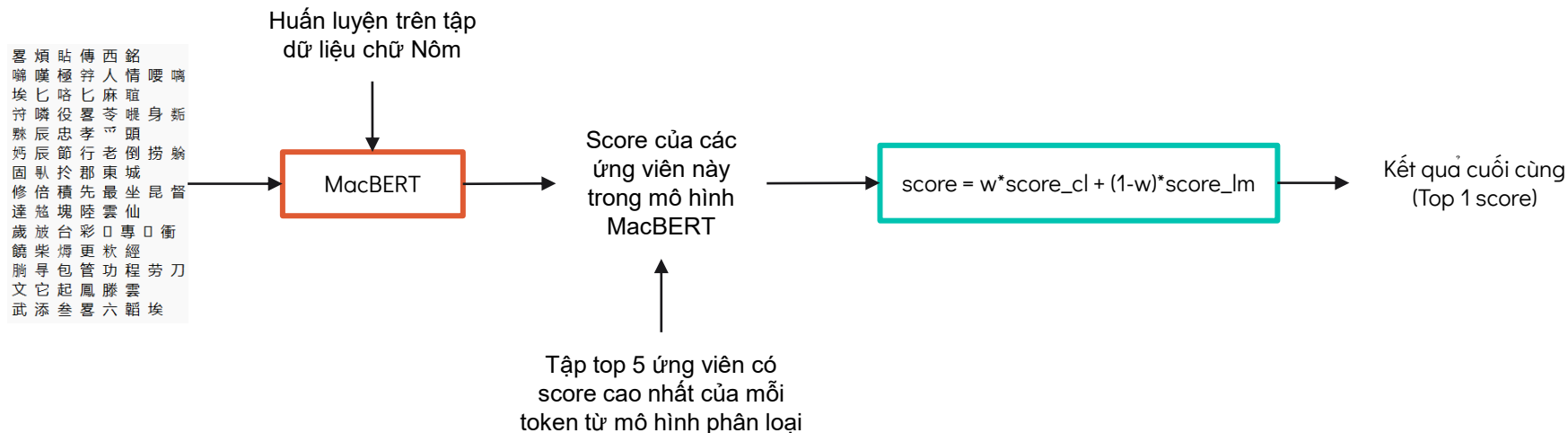
Note:

畧 煩 貼 傳 西 銘
嚙 嘆 極 筭 人 情 嚙 嚙
埃 匕 咯 匕 麻 聃

Quy trình



Quy trình





Demo quy trình hậu xử lý



Câu đúng: Hôm nay tôi đi học.

Câu input (output OCR): Hôm nay tôi đi học.

Tokenize: Hôm nay tôi đi học.

VD: tokenize cho tiếng Anh (biến hình): he is sleep ##ing.



Demo quy trình hậu xử lý



OCR	Hôm: 0.7 Hôn: 0.2 Nộ: 0.1	nay: 1.0	tôi: 0.4 tôi: 0.3 toi: 0.1	đi: 1.0	học: 1.0
OCR output	Hôm	nay	<u>tôi</u>	đi	học
MacBERT	Hôm: 0.3 Hôn: 0.1 Nộ: 0.01	nay: 1.0	tôi: 0.5 tôi: 0.3 toi: 0.1	đi: 1.0	học: 1.0
score = $w \cdot \text{score_cl} + (1-w) \cdot \text{score_lm}$ ($w = 0.5$)	S(Hôm) = 0.5 S(Hôn) = 0.15 S(Nộ) = 0.055		S(tôi) = 0.4 S(tôi) = 0.35 S(toi) = 0.1		
Output	Hôm	nay	tôi	đi	học



Khó khăn

OCR

Lý do khách quan:

- Chất lượng ảnh hoặc tài liệu không tốt. (nhiều, mờ, xuống cấp,...)
- Trong văn bản viết tay, kể cả những chữ cùng loại thì vẫn có sự biến đổi trong cách viết.
- Người viết có thể viết các dòng không theo 1 đường thẳng cố định.





Khó khăn

OCR

Lý do chủ quan:

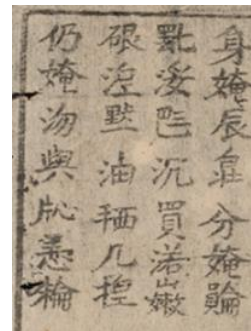
- Loại chữ gồm nhiều mẫu tự gây ra thách thức lớn cho quá trình phân loại: chữ Nôm, Hán tự, Hangul, Kanji ...

Ước tính có khoảng 50,000 Hán tự, trong đó 3775 chữ được dùng thường xuyên

vs

26 chữ cái được dùng trong Tiếng Anh

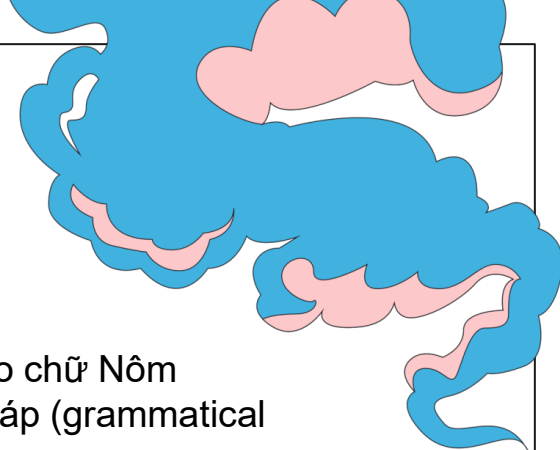
- Cách viết phức tạp, có sự tương đồng nhau về hình dáng giữa các chữ
- Viết dọc ⇒ ranh giới giữa các chữ không rõ ràng





Khó khăn

Sửa lỗi văn bản



- Tìm kiếm mô hình phù hợp để ứng dụng vào việc sửa lỗi văn bản cho chữ Nôm
 - Mô hình về sửa lỗi chính tả (spelling errors) hay sửa lỗi ngữ pháp (grammatical errors)?
 - Huấn luyện mô hình từ đầu (train from scratch) hay sử dụng học chuyển giao (transfer learning)?
- Xây dựng ngữ liệu huấn luyện
 - Thu thập văn bản chữ Nôm
 - Cần chuyên viên để xác định ngữ nghĩa của từng từ, xác định lỗi ngữ pháp, ...
 - Tìm giải pháp cho các từ lạ, chưa được đăng ký trên hệ thống Unicode



Thanks



Do you have any questions?

