

# Tổng quát hoá miền tri thức trên bài toán phân tích cú pháp phụ thuộc tiếng Việt

SVTH: Huỳnh Đặng Vĩnh Hiền, Lê Châu Anh  
GVHD: TS. Nguyễn Thị Quý, ThS. Huỳnh Thiện Ý

Khoa Khoa học Máy tính,  
Trường Đại học Công nghệ Thông tin,  
Đại học Quốc gia Thành phố Hồ Chí Minh

Ngày 15 tháng 1 năm 2025



# Mục lục

- 1 Tổng quan khoá luận
- 2 Xây dựng kho ngữ liệu DGDT
- 3 Xây dựng tiêu chuẩn đánh giá DGDTMark
- 4 Mô hình xử lý bài toán
- 5 Thực nghiệm
- 6 Công bố khoa học
- 7 Kết luận

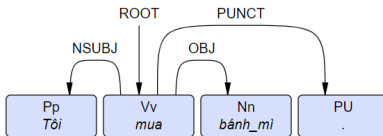
# Tổng quan khoá luận

# Bài toán phân tích cú pháp phụ thuộc

## Định nghĩa

Phân tích cú pháp phụ thuộc là một bài toán trong Xử lý ngôn ngữ tự nhiên, tập trung vào phân tích cấu trúc ngữ pháp, cụ thể là mối quan hệ phụ thuộc giữa các từ trong câu.

Để thể hiện cấu trúc ngữ pháp, mỗi câu sẽ được biểu diễn dưới dạng một cây phụ thuộc.



Hình: Một ví dụ về cây phụ thuộc

# Kho ngữ liệu phân tích cú pháp phụ thuộc tiếng Việt

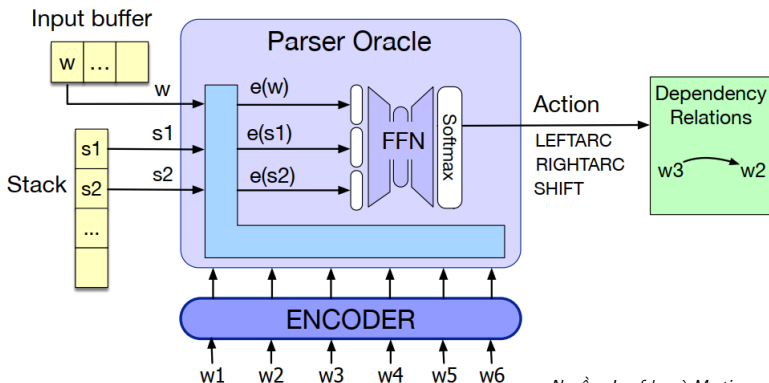
Trong bài toán phân tích cú pháp phụ thuộc, tiếng Việt vẫn được xem là một ngôn ngữ có tài nguyên hạn chế.

**Bảng:** Tổng quan các kho ngữ liệu cú pháp phụ thuộc tiếng Việt

Kho ngữ liệu	Số câu	Số nhãn phụ thuộc	Phương pháp gán nhãn	Nguồn ngữ liệu
BKTreebank	6,909	26	Thủ công	Báo Dân Trí
VnDT	10.197	33	Tự động	Báo Tuổi Trẻ
UD_Vietnamese-VTB	3.323	84	Tự động	Báo Tuổi Trẻ

# Phân tích cú pháp phụ thuộc dựa trên chuyển đổi<sup>2</sup> (transition-based)

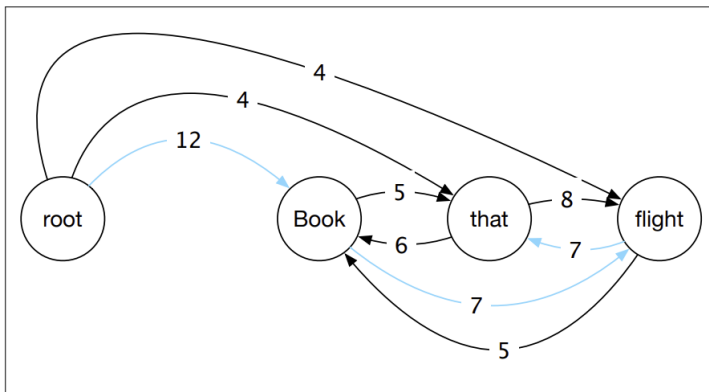
Ý tưởng: Xây dựng cấu trúc phụ thuộc bằng cách sử dụng các toán tử dịch chuyển khi duyệt câu một cách tuần tự.



Nguồn: Jurafsky và Martin

# Phân tích cú pháp phụ thuộc dựa trên đồ thị (graph-based)

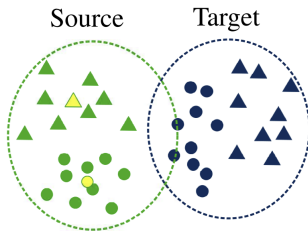
Ý tưởng: Mô hình hóa việc tìm cấu trúc cú pháp tối ưu cho câu thành bài toán tìm cây khung lớn nhất.



Nguồn: Jurafsky và Martin

# Vấn đề khác biệt miền tri thức

- Hầu hết nghiên cứu về phân tích cú pháp phụ thuộc không xem xét đến bối cảnh có sự khác biệt về phân phối giữa các miền tri thức (domain gap).
- Tuy nhiên, đây là một thách thức phổ biến khi triển khai mô hình với dữ liệu thực tế, gây khó khăn và làm giảm hiệu suất của mô hình.
- Chỉ có một số ít nghiên cứu liên quan tới domain gap nhưng chủ yếu mang tính khảo sát và chưa nghiên cứu trên tiếng Việt.



Nguồn: Robert và cộng sự





## Miền tri thức *Kinh doanh*:

- Dù chứng khoán đảo chiều vào cuối phiên để tăng 4 điểm, thanh khoản vẫn rơi về mức thấp nhất kể từ cuối tháng 10/2023.
- Theo Chứng khoán Vietcombank (VCBS), thanh khoản sụt giảm cho thấy lực cung bán ra đã có phần chững lại.

## Miền tri thức *Sức khỏe*:

- Người bệnh ung thư tuyến giáp có các thể xuất hiện các dấu hiệu như khối u dai dẳng ở cổ, khó nuốt, thay đổi giọng nói.
- Phát hiện sớm thông qua các xét nghiệm chẩn đoán, nhất là người có nguy cơ di truyền bệnh, có thể cải thiện kết quả điều trị.

Nguồn: VnExpress

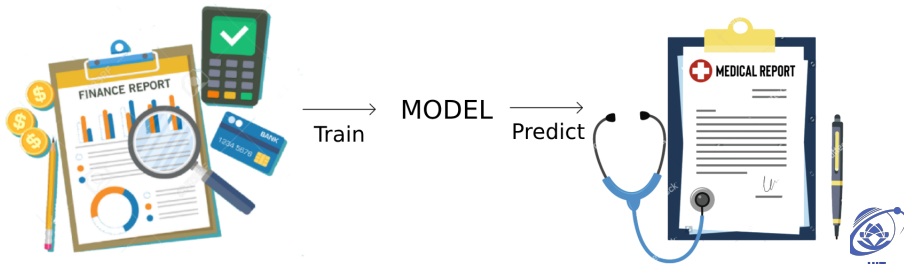


UIT  
TRƯỜNG ĐẠI HỌC  
CÔNG NGHỆ THÔNG TIN

# Bài toán tổng quát hóa miền tri thức

## Định nghĩa

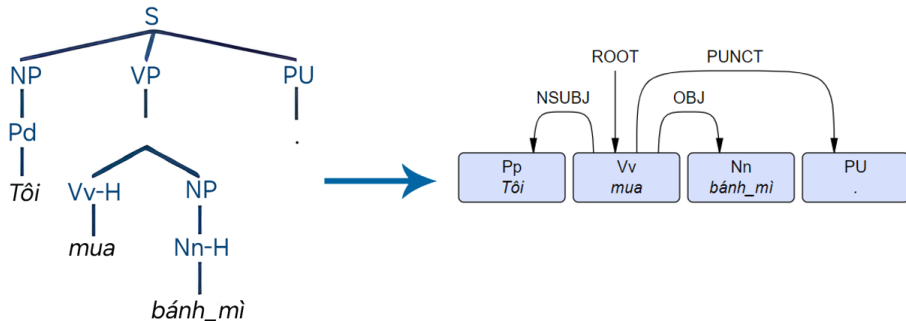
Tổng quát hóa miền tri thức (domain generalization) là một bài toán quan trọng trong học máy, mục tiêu là phát triển mô hình có khả năng áp dụng tri thức đã học được trong quá trình huấn luyện để xử lý hiệu quả tri thức mới chưa từng gặp mà không yêu cầu mô hình tiếp cận dữ liệu từ miền dữ liệu đích (target domain) trong quá trình huấn luyện.



# Xây dựng kho ngữ liệu DGDT

# Xây dựng kho ngữ liệu DGDT

Nhằm cân bằng giữa chất lượng và chi phí xây dựng kho ngữ liệu, chúng tôi áp dụng phương pháp chuyển đổi từ kho ngữ liệu cú pháp thành tổ.



# Lựa chọn kho ngữ liệu nguồn

Kho ngữ liệu nguồn: kho ngữ liệu NIIVTB-2 <sup>1</sup>

- Nguồn ngữ liệu có độ tin cậy cao, trích từ báo Thanh Niên.

---

<sup>1</sup>Nguyen, Q.T., Miyao, Y., Le, H.T.T. et al. Ensuring annotation consistency and accuracy for Vietnamese treebank. *Lang Resources & Evaluation* **52**, 269–315 (2018). 14/41

# Lựa chọn kho ngữ liệu nguồn

Kho ngữ liệu nguồn: kho ngữ liệu NIIVTB-2 <sup>1</sup>

- Nguồn ngữ liệu có độ tin cậy cao, trích từ báo Thanh Niên.
- Bao gồm nhiều chủ đề khác nhau.

---

<sup>1</sup>Nguyen, Q.T., Miyao, Y., Le, H.T.T. et al. Ensuring annotation consistency and accuracy for Vietnamese treebank. *Lang Resources & Evaluation* **52**, 269–315 (2018). 14/41

# Lựa chọn kho ngữ liệu nguồn

Kho ngữ liệu nguồn: kho ngữ liệu NIIVTB-2 <sup>1</sup>

- Nguồn ngữ liệu có độ tin cậy cao, trích từ báo Thanh Niên.
- Bao gồm nhiều chủ đề khác nhau.
- Mỗi chủ đề đóng vai trò là một miền tri thức.



Các chủ đề trong NIIVTB-2 được trích từ báo Thanh Niên

<sup>1</sup>Nguyen, Q.T., Miyao, Y., Le, H.T.T. et al. Ensuring annotation consistency and accuracy for Vietnamese treebank. *Lang Resources & Evaluation* **52**, 269–315 (2018). 14/41



# Lựa chọn công cụ chuyển đổi tự động

Công cụ chuyển đổi tự động: công cụ của *Trương và cộng sự*<sup>2</sup> (converter).

- Tương thích với kho ngữ liệu cú pháp thành tổ đã chọn.

---

<sup>2</sup>C. M. Truong, T. V. Pham, M. N. Phan, N. D. T. Le, T. V. Nguyen and Q. T. Nguyen, "Converting a constituency treebank to dependency treebank for Vietnamese," 2022 RIVF International Conference on Computing and Communication Technologies (RIVF), Ho Chi Minh City, Vietnam, 2022, pp. 256-261.

# Lựa chọn công cụ chuyển đổi tự động

Công cụ chuyển đổi tự động: công cụ của *Trương và cộng sự*<sup>2</sup> (converter).

- Tương thích với kho ngữ liệu cú pháp thành tổ đã chọn.
- Đi kèm quy tắc dựng cây và bộ nhãn phụ thuộc phù hợp đặc trưng ngôn ngữ của tiếng Việt.

**VnDT**

SUB



**converter**

ASUBJ

NSUBJ

VSUBJ

CSUBJ

So sánh giữa bộ nhãn của kho ngữ liệu VnDT và converter ở mỗi quan hệ danh từ.

<sup>2</sup>C. M. Truong, T. V. Pham, M. N. Phan, N. D. T. Le, T. V. Nguyen and Q. T. Nguyen, "Converting a constituency treebank to dependency treebank for Vietnamese," 2022 RIVF International Conference on Computing and Communication Technologies (RIVF), Ho Chi Minh City, Vietnam, 2022, pp. 256-261.

# Đánh giá chất lượng công cụ chuyển đổi

Để đảm bảo chất lượng của kho ngữ liệu, hai miền tri thức *Law* và *Life\_of\_youth* được sửa lỗi thủ công và so sánh với đầu ra của converter.

Miền tri thức	Số lượng câu	UAS	LAS
Law	610	95.55%	89.80%
Life_of_youth	635	95.69%	89.11%
Trung bình		95.62%	89.42%

Kết quả cho thấy đầu ra của converter có độ chính xác đủ cao để đảm bảo chất lượng ngữ liệu.

# Thống kê của kho ngữ liệu DGDТ

Tập	Miền tri thức	Số lượng câu	Số lượng từ
Train	Education	844	21.068
	Health	725	16.045
	Law	610	17.188
	Life_of_youth	635	16.396
	Military	690	16.949
	Politics_Society	712	20.245
	Science	692	16.964
	Sports	697	16.780
	Travel	540	13.397
	World	645	16.393
	<i>Toàn tập</i>	6.790	171.425
Dev	Entertainment	708	17.463
	Information_Technology	714	18.374
	<i>Toàn tập</i>	1.422	35.837
Test	Economic	725	18.207
	Life	828	19.537
	<i>Toàn tập</i>	1.553	37.744
<b>Toàn kho ngữ liệu</b>		<b>9.765</b>	<b>245.006</b>

# So sánh với các kho ngữ liệu khác

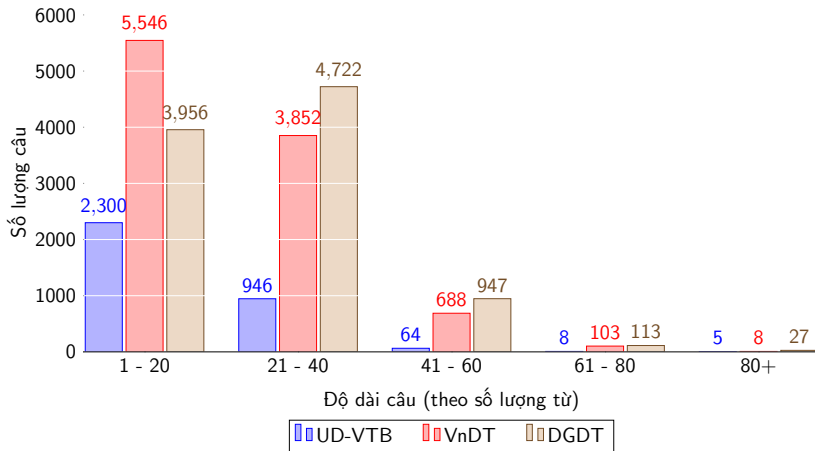
Kho ngữ liệu	Số lượng nhãn	Số lượng miền tri thức	Số lượng câu	Số lượng từ	Gán nhãn thủ công
BKTreebank <sup>3</sup>	26	không rõ	6,909	không rõ	✓
UD_Vietnamese-VTB <sup>4</sup>	84	1	3,323	58,069	✗
VnDT <sup>5</sup>	33	1	10,197	218,749	✗
<b>DGDT</b>	40	14	9,765	245,006	✗

<sup>3</sup>Nguyen, Hieu. (2017). BKTreebank: Building a Vietnamese Dependency Treebank. 10.48550/arXiv.1710.05519.

<sup>4</sup>[https://github.com/UniversalDependencies/UD\\_Vietnamese-VTB](https://github.com/UniversalDependencies/UD_Vietnamese-VTB)

<sup>5</sup><https://github.com/datquocnguyen/VnDT>

# So sánh với các kho ngữ liệu khác



# Xây dựng tiêu chuẩn đánh giá DGDTMark

# Tiêu chuẩn đánh giá DGDTMark

DGDT gồm 4 bối cảnh bồi trí ngữ liệu nhằm đánh giá tác động của đa miền tri thức lên bài toán phân tích cú pháp phụ thuộc tiếng Việt.

- ➊ **in-domain**: chia mỗi miền tri thức thành ba phần theo tỷ lệ 8:1:1 và gộp các phần tương ứng thành ba tập train, dev, test.



# Tiêu chuẩn đánh giá DGD<sup>2</sup>Mark

DGDT gồm 4 bối cảnh bồi trí ngữ liệu nhằm đánh giá tác động của đa miền tri thức lên bài toán phân tích cú pháp phụ thuộc tiếng Việt.

- ➊ **in-domain**: chia mỗi miền tri thức thành ba phần theo tỷ lệ 8:1:1 và gộp các phần tương ứng thành ba tập train, dev, test.
- ➋ **domain-k-fold**: mỗi miền tri thức sẽ đóng vai trò làm tập dev và test, phần còn lại của kho ngữ liệu sẽ làm tập train.

# Tiêu chuẩn đánh giá DGD<sup>2</sup>Mark

DGDT gồm 4 bối cảnh bồi trí ngữ liệu nhằm đánh giá tác động của đa miền tri thức lên bài toán phân tích cú pháp phụ thuộc tiếng Việt.

- ① **in-domain**: chia mỗi miền tri thức thành ba phần theo tỷ lệ 8:1:1 và gộp các phần tương ứng thành ba tập train, dev, test.
- ② **domain-k-fold**: mỗi miền tri thức sẽ đóng vai trò làm tập dev và test, phần còn lại của kho ngữ liệu sẽ làm tập train.
- ③ **domain-generalization**: bố trí ngữ liệu sao cho mỗi miền tri thức nằm độc nhất ở một trong ba tập train, dev, test.

# Tiêu chuẩn đánh giá DGD<sup>2</sup>Mark

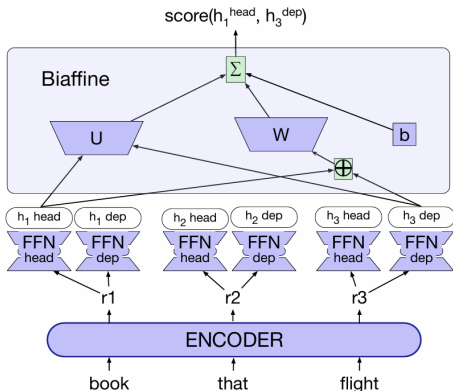
DGDT gồm 4 bối cảnh bố trí ngữ liệu nhằm đánh giá tác động của đa miền tri thức lên bài toán phân tích cú pháp phụ thuộc tiếng Việt.

- ➊ **in-domain**: chia mỗi miền tri thức thành ba phần theo tỷ lệ 8:1:1 và gộp các phần tương ứng thành ba tập train, dev, test.
- ➋ **domain-k-fold**: mỗi miền tri thức sẽ đóng vai trò làm tập dev và test, phần còn lại của kho ngữ liệu sẽ làm tập train.
- ➌ **domain-generalization**: bố trí ngữ liệu sao cho mỗi miền tri thức nằm độc nhất ở một trong ba tập train, dev, test.
- ➍ **dataset-generalization**: bố trí ngữ liệu giống bối cảnh 3, thay tập test bằng tập test lấy từ kho ngữ liệu **NIIVTB\_DT-1**.

# Mô hình xử lý bài toán

# Baseline: Mô hình Deep Biaffine Parsing

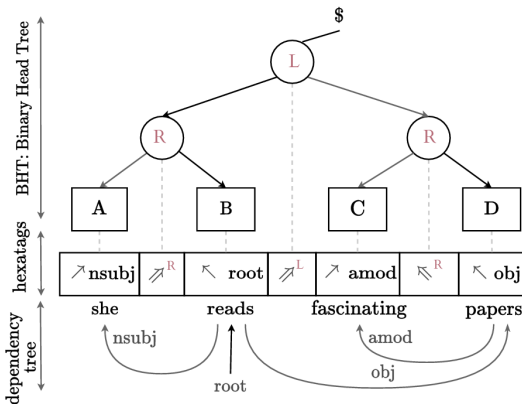
- Deep Biaffine Parsing (Biaffine) là mô hình có sức ảnh hưởng lớn đến sự phát triển của bài toán.
- Mô hình dựa trên hướng tiếp cận graph-based, trong đó sử dụng *biaffine attention* để tính điểm cho các cung trong đồ thị.



Kiến trúc của mô hình Biaffine

# Baseline: Mô hình Hexatagger

- Trong mô hình này, mỗi cây phụ thuộc sẽ được chuyển đổi thành cây nhị phân (Binary Head Tree - BHT).
- Mô hình sẽ tiến hành học cách xây dựng BHT thông qua việc gán nhãn cho từng đỉnh trong cây.
- Sau khi thực hiện, mô hình giải mã BHT đã xây dựng lại thành cây phụ thuộc.

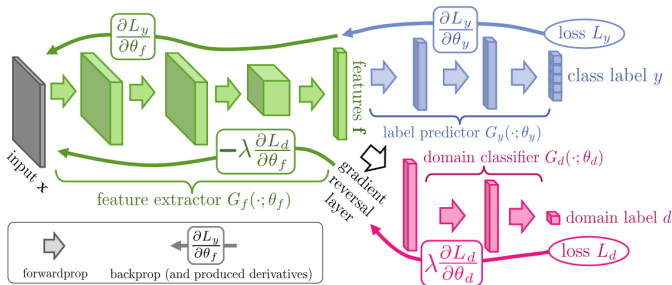


# Phương pháp đề xuất: Adv

- Huấn luyện đối kháng (*Adversarial Training*) là một phương pháp được đề xuất nhằm giúp mô hình nắm bắt được các đặc trưng bất biến trên đa miền tri thức.

# Phương pháp đề xuất: Adv

- Huấn luyện đối kháng (*Adversarial Training*) là một phương pháp được đề xuất nhằm giúp mô hình nắm bắt được các đặc trưng bất biến trên đa miền tri thức.
- Mô hình sẽ được huấn luyện song song với một bài toán khác (bài toán đối kháng) với giả thiết là khả năng mô hình xử lý tốt bài toán đối kháng **tỷ lệ nghịch** với khả năng mô hình nắm bắt được các đặc trưng bất biến.

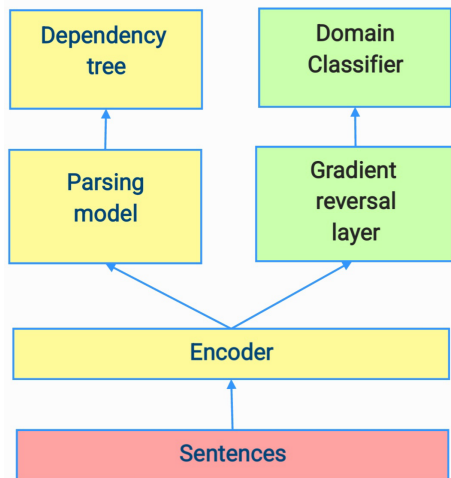


Kiến trúc của phương pháp sử dụng huấn luyện đối kháng



# Phương pháp đề xuất: Adv

- Bài đoán đối kháng được chọn là phân loại miền tri thức (*domain classifier*).
- Lớp Gradient Reversal giúp **đảo ngược dấu** của đạo hàm trả về trong quá trình backward.
- Tham số  $\gamma \in [0, 1]$  được thêm vào nhằm cân bằng giữa hai bài toán:  $\mathcal{L}_{\text{model}} = \gamma \mathcal{L}_{\text{classifier}} + (1 - \gamma) \mathcal{L}_{\text{parser}}$



Kiến trúc mô hình do nhóm đề xuất

# Thực nghiệm

# Thiết lập thực nghiệm

- Baseline cho DGDTMark được xác lập thông qua hai mô hình là Biaffine và Hexatagger, mỗi mô hình được thực nghiệm trên hai encoder là PhoBERT và XLM-R.

# Thiết lập thực nghiệm

- Baseline cho DGDTMark được xác lập thông qua hai mô hình là Biaffine và Hexatagger, mỗi mô hình được thực nghiệm trên hai encoder là PhoBERT và XLM-R.
- Phương pháp đề xuất được thực nghiệm trên mô hình Biaffine với encoder là PhoBERT.

# Thiết lập thực nghiệm

- Baseline cho DGDTMark được xác lập thông qua hai mô hình là Biaffine và Hexatagger, mỗi mô hình được thực nghiệm trên hai encoder là PhoBERT và XLM-R.
- Phương pháp đề xuất được thực nghiệm trên mô hình Biaffine với encoder là PhoBERT.
- Mỗi mô hình được chạy trên 100 epoch, mô hình tốt nhất được lựa chọn trên tập dev với độ đo LAS.

# Kết quả thực nghiệm trên DGD TMark

Độ đo: **UAS (%)**

Mô hình	Bộ mã hoá	<i>in-domain</i>	<i>domain- k-fold</i>	<i>domain- generalization</i>	<i>dataset- generalization</i>
Biaffine	PhoBERT	<b>91.62</b>	<b>91.22</b>	88.95	88.54
	XLM-R	90.25	<u>89.72</u>	<u>88.26</u>	<u>86.98</u>
Hexatagger	PhoBERT	91.49	91.00	<b>90.00</b>	<b>89.19</b>
	XLM-R	<u>89.72</u>	89.97	89.62	89.01

Độ đo: **LAS (%)**

Mô hình	Bộ mã hoá	<i>in-domain</i>	<i>domain- k-fold</i>	<i>domain- generalization</i>	<i>dataset- generalization</i>
Biaffine	PhoBERT	87.83	87.01	84.76	82.74
	XLM-R	<u>85.89</u>	<u>84.99</u>	<u>82.88</u>	<u>80.97</u>
Hexatagger	PhoBERT	<b>89.00</b>	<b>88.28</b>	<b>86.81</b>	<b>85.40</b>
	XLM-R	86.99	87.08	86.48	85.22

## Kết quả thực nghiệm trên *domain-k-fold* (UAS)

Miền tri thức	Biaffine	Hexatagger
Economic	<b>90.53</b>	90.21
Education	<b>90.97</b>	90.92
Entertainment	90.52	<b>90.57</b>
Health	<b>91.18</b>	91.07
Information_Technology	<b>90.36</b>	90.28
Law	<b>91.58</b>	91.05
Life	89.50	<b>89.91</b>
Life_of_youth	<b>91.12</b>	90.85
Military	<b>92.29</b>	92.16
Politics_Society	<b>91.16</b>	90.43
Science	<b>92.00</b>	91.71
Sports	<b>91.84</b>	91.49
Travel	<b>91.17</b>	90.67
World	<b>92.89</b>	92.62
<i>Trung bình</i>	<b>91.22</b>	91.00

# Kết quả thực nghiệm trên *domain-k-fold* (LAS)

Miền tri thức	Biaffine	Hexatagger
Economic	85.82	<b>87.32</b>
Education	86.74	<b>88.05</b>
Entertainment	85.94	<b>87.47</b>
Health	86.55	<b>87.99</b>
Information_Technology	85.89	<b>87.43</b>
Law	87.45	<b>88.49</b>
Life	84.82	<b>86.83</b>
Life_of_youth	86.56	<b>87.82</b>
Military	88.82	<b>89.89</b>
Politics_Society	86.83	<b>87.69</b>
Science	88.85	<b>89.63</b>
Sports	87.16	<b>88.69</b>
Travel	87.03	<b>87.97</b>
World	89.61	<b>90.60</b>
<b>Trung bình</b>	<b>87.01</b>	<b>88.28</b>

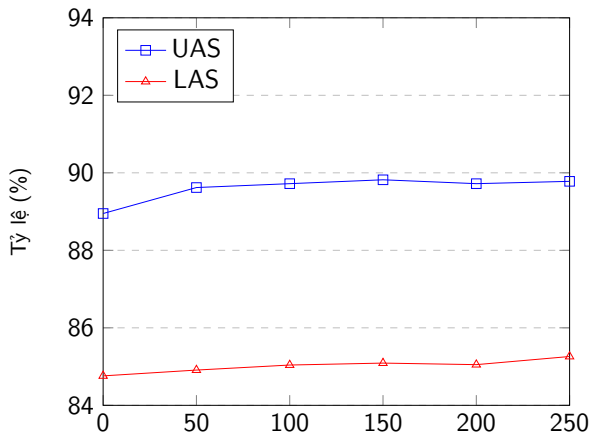


# Phân tích hiệu suất theo nhãn

Nhãn phụ thuộc	Số lần xuất hiện	UAS (%)	LAS (%)
PUNCT	5136	83.94	83.92
NN	3570	92.83	88.01
OBJ	3080	94.42	90.06
VMOD	2637	91.13	81.53
PREP	2570	85.37	83.66
NSUBJ	2340	90.64	89.53
POBJ	2216	<b>97.70</b>	95.35
ADJUNCT	2185	93.00	90.30
CONJ	2169	<u>78.01</u>	<u>74.09</u>
ROOT	1553	90.86	90.86
CC	1250	82.56	81.84
AMOD	1132	89.66	85.42
DET	1062	96.23	<b>95.10</b>

*Danh sách bao gồm những nhãn xuất hiện hơn 1000 lần trong tập đánh giá.*

# Ảnh hưởng của dữ liệu kiểm tra trong quá trình huấn luyện



Số lượng câu từ tập dữ liệu kiểm tra cung cấp cho tập dữ liệu huấn luyện

# Kết quả thực nghiệm phương pháp đề xuất trên DGDTMark

Độ đo: **UAS (%)**

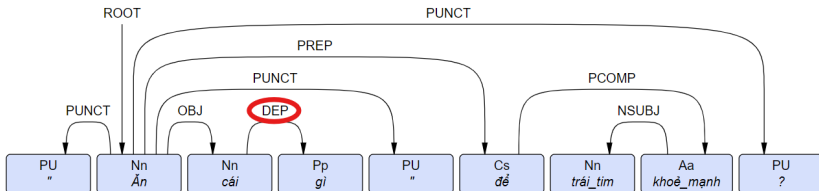
Mô hình	<i>in-domain</i>	<i>domain- k-fold</i>	<i>domain- generalization</i>	<i>dataset- generalization</i>
Biaffine	<b>91.62</b>	<b>91.22</b>	88.95	<b>88.54</b>
Biaffine+Adv	91.20	90.89	<b>89.41</b>	87.51

Độ đo: **LAS (%)**

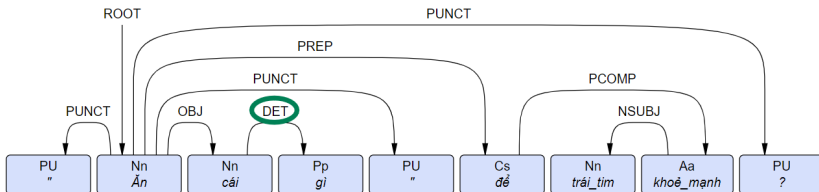
Mô hình	<i>in-domain</i>	<i>domain- k-fold</i>	<i>domain- generalization</i>	<i>dataset- generalization</i>
Biaffine	87.83	87.01	84.76	<b>82.74</b>
Biaffine+Adv	<b>87.88</b>	<b>87.22</b>	<b>85.25</b>	82.62

# So sánh dự đoán trên bối cảnh *in-domain*

- Dự đoán của mô hình Biaffine (sai):

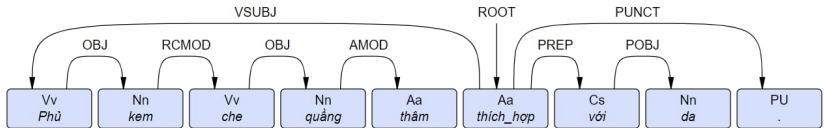


- Dự đoán của mô hình đề xuất Biaffine+Adv (đúng):

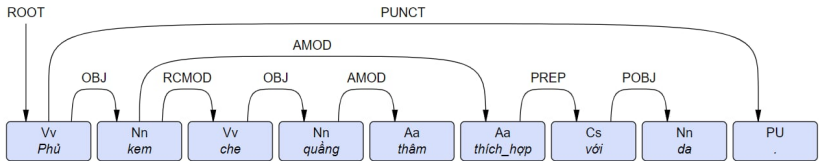


# So sánh dự đoán trên bối cảnh *domain-generalization*

- Dự đoán của mô hình Biaffine (sai):



- Dự đoán của mô hình đề xuất Biafine+Adv (đúng):



# Công bố khoa học

[1] Vinh-Hien D. Huynh, Chau-Anh Le, Chau M. Truong, Y Thien Huynh and Quy T. Nguyen, *Domain Generalization in Vietnamese Dependency Parsing: A Novel Benchmark and Domain Gap Analysis*, In Proceedings of the 13th International Symposium on Information and Communication Technology (SOICT '24)

⇒ Trình bày về việc xây dựng kho ngữ liệu DGDT và tiêu chuẩn đánh giá DGDTMark, thiết lập baseline bằng mô hình Biaffine và phân tích ảnh hưởng của domain gap.

# Kết luận



# Kết luận

- Khoá luận đã xây dựng kho ngữ liệu DGDT, một kho ngữ liệu phân tích cú pháp phụ thuộc tiếng Việt trên đa miền tri thức.

# Kết luận

- Khoá luận đã xây dựng kho ngữ liệu DGDT, một kho ngữ liệu phân tích cú pháp phụ thuộc tiếng Việt trên đa miền tri thức.
- Ngoài ra, nhóm nghiên cứu đã đề xuất tiêu chuẩn đánh giá DGDTMark nhằm đánh giá tác động của domain gap lên bài toán phân tích cú pháp phụ thuộc.

## Kết luận

- Khoá luận đã xây dựng kho ngữ liệu DGDĐT, một kho ngữ liệu phân tích cú pháp phụ thuộc tiếng Việt trên đa miền tri thức.
- Ngoài ra, nhóm nghiên cứu đã đề xuất tiêu chuẩn đánh giá DGDĐTMark nhằm đánh giá tác động của domain gap lên bài toán phân tích cú pháp phụ thuộc.
- Khoá luận cũng giới thiệu một hướng tiếp cận mới nhằm cải thiện khả năng tổng quát hoá miền tri thức của các mô hình.

# Kết luận

- Khoá luận đã xây dựng kho ngữ liệu DGDT, một kho ngữ liệu phân tích cú pháp phụ thuộc tiếng Việt trên đa miền tri thức.
  - Ngoài ra, nhóm nghiên cứu đã đề xuất tiêu chuẩn đánh giá DGDTMark nhằm đánh giá tác động của domain gap lên bài toán phân tích cú pháp phụ thuộc.
  - Khoá luận cũng giới thiệu một hướng tiếp cận mới nhằm cải thiện khả năng tổng quát hoá miền tri thức của các mô hình.
- Kết quả thực nghiệm cho thấy, các mô hình tiên tiến nhất hiện nay (SOTA) vẫn chịu ảnh hưởng đáng kể từ vấn đề domain gap.

# Kết luận

- Khoá luận đã xây dựng kho ngữ liệu DGDT, một kho ngữ liệu phân tích cú pháp phụ thuộc tiếng Việt trên đa miền tri thức.
  - Ngoài ra, nhóm nghiên cứu đã đề xuất tiêu chuẩn đánh giá DGDTMark nhằm đánh giá tác động của domain gap lên bài toán phân tích cú pháp phụ thuộc.
  - Khoá luận cũng giới thiệu một hướng tiếp cận mới nhằm cải thiện khả năng tổng quát hoá miền tri thức của các mô hình.
- Kết quả thực nghiệm cho thấy, các mô hình tiên tiến nhất hiện nay (SOTA) vẫn chịu ảnh hưởng đáng kể từ vấn đề domain gap.
- Phương pháp đề xuất trong khoá luận đã thể hiện tiềm năng trong việc cải thiện hiệu suất của mô hình ở một số bối cảnh bố trí ngữ liệu.

