

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC - KỸ THUẬT MÁY TÍNH



GIAI ĐOẠN 1 (GD1):
ĐỀ CƯƠNG LUẬN VĂN/ ĐỒ ÁN CHUYÊN NGÀNH/
ĐỒ ÁN MÔN HỌC KỸ THUẬT MÁY TÍNH

Tên đề tài

Tiếng Việt: Rút trích quan hệ giữa Ý định và Thực thể
sử dụng hướng tiếp cận Đọc hiểu Máy
cho nhiệm vụ xây dựng Đồ thị Tri thức trong lĩnh vực Giáo dục.

Tiếng Anh: Relation Extraction between
Intent and Entity using Machine Reading Comprehension(MRC)
approach for Knowledge
Graph Construction in Education domain.

Giáo viên hướng dẫn 1: Bùi Công Tuấn
Giáo viên hướng dẫn 2: Bùi Công Tuấn
SV thực hiện: Lưu Quốc Bình 2033009



Contents

| | | |
|----------|---|-----------|
| 1 | Giới thiệu | 2 |
| 2 | Phương pháp trích xuất quan hệ | 3 |
| 2.1 | Các phương pháp truyền thống | 3 |
| 2.1.1 | Phương pháp dựa trên luật (Rule-based) | 4 |
| 2.1.2 | Phương pháp không giám sát (Unsupervised) | 5 |
| 2.1.3 | Phương pháp có giám sát (Supervised) | 5 |
| 2.1.4 | Phương pháp bán giám sát (Semi-supervised) | 5 |
| 2.1.5 | Phương pháp giám sát từ xa (Distant supervision) | 7 |
| 2.2 | Phương pháp học sâu (Deep learning) | 8 |
| 2.2.1 | Các hướng tiếp cận dựa trên phương pháp CNN (CNN based methods) | 8 |
| 2.2.2 | Phương pháp dựa trên RNN và LSTM (RNN and LSTM based methods) | 8 |
| 2.2.3 | Encoder-decoder/transformer based methods | 10 |
| 2.2.3.a | Neural OpenIE (<i>Cui et al. (2018)</i>) | 10 |
| 2.2.3.b | Transformer based OpenIE(<i>Han and Wang (2021)</i>) | 11 |
| 2.2.3.c | Multi2OIE (<i>Ro et al. (2020)</i>) | 12 |
| 3 | Các MRC phù hợp với Tiếng Việt | 13 |
| 3.1 | Mô hình tinh chỉnh (Fine-tuned model): | 13 |
| 3.2 | API của bên thứ ba (3rd party API): | 13 |
| 4 | Định nghĩa các loại quan hệ được rút trích lĩnh vực giáo dục | 14 |
| 4.1 | Entities | 14 |
| 4.1.1 | Intents | 14 |
| 4.1.2 | Terms | 14 |
| 4.1.3 | Papers | 14 |
| 4.1.4 | Subject | 14 |
| 4.1.5 | Student | 14 |
| 4.1.6 | Lecturers | 14 |
| 4.1.7 | Academic_service_Office | 14 |
| 4.1.8 | English_Certificate | 14 |
| 4.1.9 | URL | 14 |
| 4.2 | Relations | 14 |
| 4.2.1 | hỏi ý định rút môn | 14 |
| 4.2.2 | hỏi chung về ý định | 14 |
| 4.2.3 | hỏi thời gian bat81 đầu đăng ký môn | 14 |
| 4.2.4 | hỏi lịch học | 14 |
| 4.2.5 | hỏi môn tương đương | 14 |
| 4.2.6 | đăng kí gia hạn | 15 |
| 4.2.7 | hỏi học phí | 15 |
| 4.2.8 | Thay đổi đăng ký môn học | 15 |
| 4.2.9 | hỏi về chuẩn ngoại ngữ | 15 |
| 4.2.10 | URLs | 15 |
| 4.2.10.a | hỏi đăng kí in bằng điểm | 15 |
| 4.2.10.b | hỏi điều kiện nhận luận văn | 15 |
| 4.2.10.c | hỏi điều kiện tốt nghiệp | 15 |



1 Giới thiệu

Báo cáo này trình bày các phương pháp trích xuất quan hệ sử dụng mô hình Hiểu đọc máy tính (Machine Reading Comprehension - MRC).

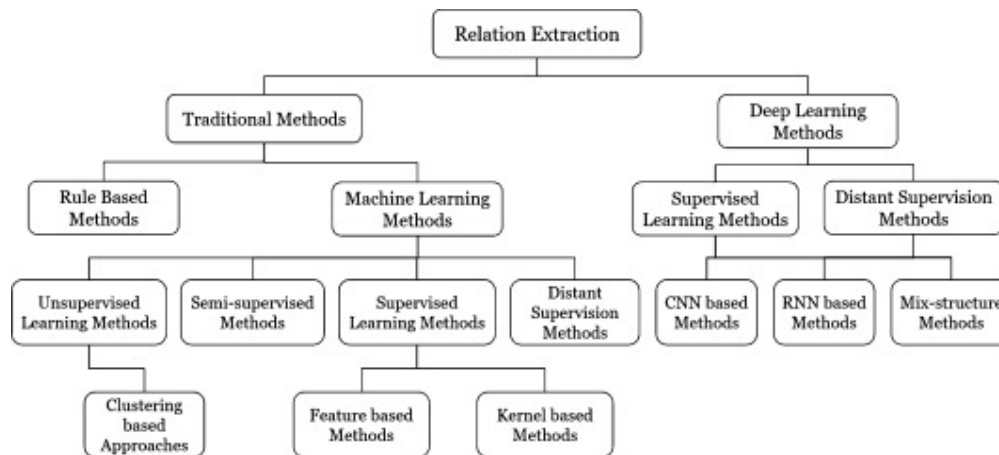
- Đường hướng nghiên cứu: Báo cáo tập trung vào phương pháp trích xuất quan hệ trong lĩnh vực giáo dục sử dụng MRC.

2 Phương pháp trích xuất quan hệ

2.1 Các phương pháp truyền thống

Như được thể hiện trong Hình 1, các phương pháp trích xuất quan hệ (RE) hiện có có thể được phân loại thành hai loại chính: Phương pháp truyền thống (Traditional Methods) và Phương pháp học sâu (Deep Learning Methods). Phương pháp truyền thống sử dụng các kỹ thuật dựa trên luật (Rule Based Methods) hoặc kỹ thuật học máy để trích xuất một tập hợp các quan hệ được xác định trước từ một kho văn bản (corpus). Các phương pháp truyền thống có thể được phân loại chi tiết thành phương pháp dựa trên luật (Rule Based Methods) và phương pháp học máy (Machine Learning Methods).

Các phương pháp học máy (Machine Learning Methods) có thể được phân loại thêm thành bốn loại: (1) Phương pháp không giám sát (Unsupervised Learning Methods). (2) Phương pháp có giám sát (Supervised Learning Methods). (3) Phương pháp bán giám sát (Semi-supervised). và (4) Phương pháp Giám sát từ xa (Distant supervision).



Nguồn: <https://www.sciencedirect.com/science/article/pii/S2667305323000698br0060>

Hình. 1: Các loại phương pháp trích xuất quan hệ.

- Phương pháp dựa trên luật (Rule-based): Sử dụng các quy tắc thủ công được thiết lập sẵn để trích xuất quan hệ.
- Phương pháp không giám sát (Unsupervised): Tự động học các mẫu trích xuất quan hệ từ dữ liệu.
- Phương pháp có giám sát (Supervised): Sử dụng dữ liệu được dán nhãn để huấn luyện mô hình trích xuất quan hệ.
- Phương pháp bán giám sát (Semi-supervised): Kết hợp dữ liệu được dán nhãn và không được dán nhãn để huấn luyện mô hình.
- Phương pháp giám sát từ xa (Distant supervision): Sử dụng dữ liệu gián tiếp để huấn luyện mô hình.

2.1.1 Phương pháp dựa trên luật (Rule-based)

Các phương pháp này còn được gọi là phương pháp mẫu thủ công (hand-built pattern methods). Chúng xác định một tập hợp các mẫu trích xuất ([extraction patterns](#)) cho các quan hệ được định nghĩa trước. Sau đó, các mẫu trích xuất này được đối chiếu với văn bản. Nếu một mẫu khớp với văn bản, thì một quan hệ tương ứng với mẫu đó được tìm thấy trong văn bản. [Riloff \(1993\)](#), [Appelt et al. \(1993\)](#), [Agichtein và Gravano \(2000\)](#), [Jayram et al. \(2006\)](#), [Shen et al. \(2007\)](#), [Fundel et al. \(2006\)](#), [Nebhi \(2013\)](#) đã sử dụng các phương pháp dựa trên luật để trích xuất quan hệ từ văn bản. Bảng 3 hiển thị các ví dụ về quy tắc để trích xuất quan hệ hạ đẳng (hyponymy) từ văn bản.

| Pattern | Input sentence | Extracted relations |
|--|---|---|
| <i>such NP as</i> <i>{NP}*</i> | ... works by such authors as Herrick, Goldsmith, and Shakespeare. | Hyponym ("author", "Herrick"), Hyponym ("author", "Goldsmith"), Hyponym ("author", "Shakespeare") |
| <i>NP { NP } * { }</i> <i>or other NP</i> | Bruises, wounds, broken bones or other injuries ... | Hyponym ("bruise", "injury"), Hyponym ("wound", "injury"), Hyponym ("broken bones", "injury") |

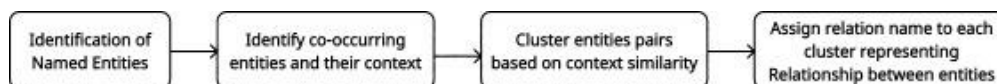
Hình. 2: Ví dụ cho partern (mẫu/lệnh) Hyponyms ([Hearst \(1992\)](#)).

Các phương pháp dựa trên luật đòi hỏi chuyên môn về lĩnh vực và kiến thức ngôn ngữ để xác định các mẫu trích xuất. Những phương pháp này phụ thuộc vào từng lĩnh vực cụ thể, trong đó cấu trúc của tài liệu là cố định và các quan hệ mục tiêu được xác định trước. Nếu chuyển từ lĩnh vực này sang lĩnh vực khác, thì cần phải xác định lại một tập hợp mới các quan hệ mục tiêu và mẫu trích xuất. Do đó, các phương pháp dựa trên luật đòi hỏi nhiều công sức thủ công và không thể sử dụng cho các tập văn bản đa dạng (heterogeneous corpora).

2.1.2 Phương pháp không giám sát (Unsupervised)

Phương pháp không giám sát không yêu cầu bất kỳ dữ liệu được dán nhãn nào. Hầu hết các phương pháp RE không giám sát sử dụng cách tiếp cận dựa trên cụm (clustering). Một trong những phương pháp RE không giám sát dựa trên cụm tiên phong được đề xuất bởi [Hasegawa et al. \(2004\)](#). Họ sử dụng công cụ đánh nhãn Thực thể Danh riêng (NE) để trích xuất các thực thể, cho phép tập trung chỉ vào các quan hệ với những thực thể được đề cập. Hình 2 minh họa các bước của phương pháp học không giám sát:

- (1) Nhận dạng các thực thể tên riêng trong kho văn bản.
- (2) Xác định các thực thể tên riêng đồng xuất hiện và ngữ cảnh của chúng.
- (3) Cụm các cặp thực thể dựa trên tính tương đồng ngữ cảnh.
- (4) Gán tên quan hệ ngữ nghĩa cho mỗi cụm.



Hình. 3: Cách tiếp cận dựa trên phân cụm cho *Unsupervised Learning*.

2.1.3 Phương pháp có giám sát (Supervised)

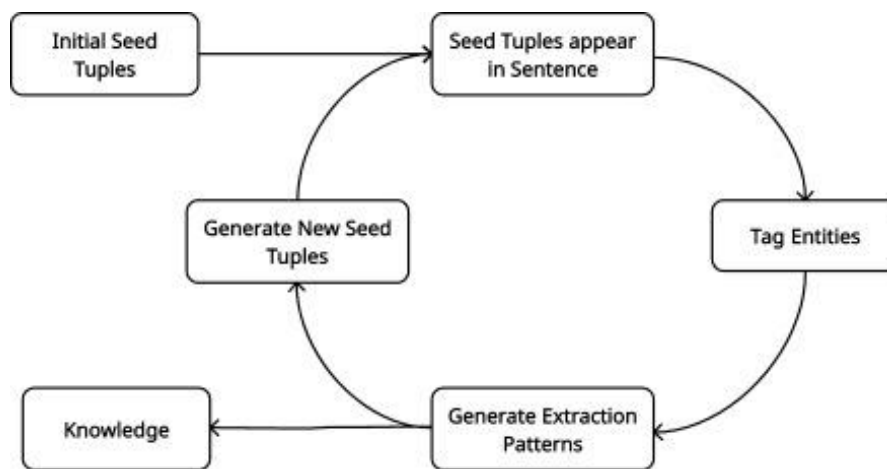
Phương pháp có giám sát yêu cầu một lượng lớn dữ liệu huấn luyện được dán nhãn với tập hợp các thực thể và quan hệ. Phương pháp này sử dụng dữ liệu huấn luyện để đào tạo bộ phân loại, trích xuất quan hệ từ dữ liệu kiểm tra. Có hai loại phương pháp có giám sát: Phương pháp dựa trên đặc trưng (Feature-based methods) và Phương pháp dựa trên Hạt nhân (Kernel-based methods) (Pawar et al. (2017)).

- Trong phương pháp dựa trên đặc trưng (Rink và Harabagiu (2010), Kambhatla (2004), Zhou et al. (2005)), một tập hợp các đặc trưng được tạo ra cho mỗi quan hệ trong dữ liệu huấn luyện, và bộ phân loại được đào tạo để trích xuất một thể hiện quan hệ mới. Một số đặc trưng về từ vựng (lexical), cú pháp (syntactic) và ngữ nghĩa (semantic) được mô tả trong (Kambhatla (2004)). Việc lựa chọn các đặc trưng ảnh hưởng đến hiệu suất của hệ thống RE có giám sát dựa trên đặc trưng.
- Không giống như phương pháp dựa trên đặc trưng, phương pháp dựa trên hàm nhân (Bunescu và Mooney, 2005b, Bunescu và Mooney, 2005a, Zelenko et al. (2003), Culotta và Sorensen (2004), Zhang et al. (2008), Zhang (2004)) không sử dụng các bước tiền xử lý để thiết kế đặc trưng. Trong phương pháp dựa trên hàm nhân, các *kernel functions* được sử dụng để xác định độ tương đồng giữa hai biểu diễn thể hiện quan hệ, trong khi Máy vectơ hỗ trợ (SVM) được sử dụng để phân loại. Hiệu suất của hệ thống RE dựa trên *kernel functions* phụ thuộc vào việc thiết kế các hàm nhân.

2.1.4 Phương pháp bán giám sát (Semi-supervised)

Việc tạo dữ liệu cho các phương pháp RE có giám sát đòi hỏi nhiều chi phí, nhân công và thời gian. Phương pháp bán giám sát tự động tạo dữ liệu được dán nhãn bằng thuật toán

bootstrapping. Cách tiếp cận này có hai ưu điểm: (1) giảm thiểu công sức cần thiết để tạo dữ liệu được dán nhãn; (2) tận dụng dữ liệu không được dán nhãn **unlabeled data** sẵn có miễn phí. Thuật toán bootstrapping yêu cầu một lượng lớn dữ liệu không được dán nhãn và một số ít instance hạt giống (seed instance) của kiểu quan hệ mong muốn. Ví dụ, để trích xuất quan hệ "Thủ đô của" (CapitalOf), các instance hạt giống "(New Delhi, Ấn Độ)", "(Canberra, Úc) và (London, Anh)" có thể được sử dụng để học một mẫu trích xuất. Với các instance hạt giống này làm đầu vào, thuật toán bootstrapping có thể trích xuất các quan hệ tương tự với các cặp thực thể như "(Paris, Pháp)". **Hình 3** minh họa mô hình để trích xuất các mẫu và các cặp thực thể hạt giống bằng phương pháp học bán giám sát. KnowItAll ([Etzioni et al. \(2004\)](#)), TextRunner ([Banko et al. \(2007\)](#)), OLLIE ([Mausam et al. \(2012\)](#)), v.v. sử dụng phương pháp bán giám sát cho RE.



Hình. 4: Phương pháp học bán giám sát (Semi-Supervised Learning) để trích xuất quan hệ.

2.1.5 Phương pháp giám sát từ xa (Distant supervision)

Phương pháp này còn được gọi là phương pháp giám sát yếu (weakly supervised methods) hoặc phương pháp dựa trên kiến thức (knowledge based methods). [Mintz et al. \(2009\)](#) đã đề xuất phương pháp giám sát từ xa, trong đó dữ liệu huấn luyện được tạo tự động bằng cách liên kết văn bản với một Cơ sở tri thức (Knowledge Base - KB). Điều này giúp loại bỏ vấn đề dán nhãn thủ công cho dữ liệu huấn luyện. Học từ xa dựa trên giả định rằng nếu hai thực thể có quan hệ xuất hiện trong KB, thì tất cả các cụm từ đề cập đến hai thực thể này đều có thể diễn đạt mối quan hệ đó.

Theo cách này, Giám sát từ xa sử dụng cơ sở tri thức, chẳng hạn như Freebase, để trích xuất các mối quan hệ giữa hai đối tượng. Khi cùng một cặp thực thể xuất hiện trong cả một câu và một KB, thì theo quy tắc heuristic, câu đó được liên kết với mối quan hệ khớp từ KB. Ví dụ, hãy xem xét câu: 'Bill Gates là người sáng lập Microsoft.' Nếu người 'Bill Gates' và tổ chức 'Microsoft' xuất hiện dưới dạng bộ ba '(entity1: Bill Gates, entity2: Microsoft, relation: founder_{of})' trong Freebase, thì hai thực thể này đại diện cho mối quan hệ founder_of (người sáng lập). Smirnova và Cudré-Mauroux (2018) đã sử dụng một quy trình như được thể hiện trong [Hình 4](#) để tạo dữ liệu huấn luyện bằng Giám sát từ xa. Các cách tiếp cận khác ([Riedel et al. \(2010\)](#), [Takamatsu et al. \(2012\)](#), [Zeng et al. \(2014\)](#), [Qin et al. \(2018\)](#), [Chen and Manning \(2014\)](#), [Quirk and Poon \(2017\)](#)) sử dụng Giám sát từ xa để trích xuất các mối quan hệ từ văn bản.

2.2 Phương pháp học sâu (Deep learning)

Mạng nơ-ron sâu (Deep Neural Networks - DNNs) ([Bengio \(2009\)](#)) đã trở nên phổ biến trong thập kỷ qua cho nhiều ứng dụng, bao gồm xử lý ngôn ngữ tự nhiên (NLP), thị giác máy tính (CV) và nhận dạng giọng nói. Kiến trúc DNN dựa trên Mạng nơ-ron tích chập (Convolutional Neural Network - CNN), Mạng nơ-ron hồi quy (Recurrent Neural Network - RNN), Bộ nhớ dài hạn (Long Short-Term Memory - LSTM), Gated Recurring Unit (GRU), Bộ mã hóa song hướng từ Transformers (BERT) được sử dụng cho các nhiệm vụ NLP khác nhau chẳng hạn như đánh nhãn từ tính (POS tagging) ([Gopalakrishnan et al. \(2019\)](#), [Srivastava et al. \(2018\)](#), [Kumar et al. \(2018\)](#)), dịch máy ([Wu et al. \(2016\)](#), [Wang et al. \(2019\)](#), [Gehring et al. \(2017\)](#)), trả lời câu hỏi ([Zhu et al. \(2018\)](#), [Lei et al. \(2018\)](#), [Liu et al. \(2020\)](#)), gắn nhãn vai trò ngữ nghĩa (semantic role labelling) ([He et al. \(2017\)](#), [Zhou and Xu \(2015\)](#), [Marcheggiani and Titov \(2017\)](#)), sinh hội thoại ([Li et al. \(2016\)](#), [Miranda and Kessaci \(2020\)](#), [Zhao and Eskenazi \(2016\)](#)), sinh văn bản ([Chen et al. \(2020\)](#), [Raffel et al. \(2020\)](#), [Li et al. \(2021\)](#)), phân tích định hướng (sentiment analysis) ([Yu et al. \(2019\)](#), [AL-Smadi et al. \(2019\)](#), [Zhao et al. \(2018\)](#)), tóm tắt tự động ([Zhang et al. \(2019\)](#), [See et al. \(2017\)](#), [Liu and Lapata \(2019\)](#), [Nallapati et al. \(2016\)](#), [Cai et al. \(2019\)](#)) v.v.

2.2.1 Các hướng tiếp cận dựa trên phương pháp CNN (CNN based methods)

Các công trình ban đầu về trích xuất quan hệ sử dụng học sâu dựa trên mô hình học có giám sát với tập dữ liệu huấn luyện được dán nhãn thủ công. Mô hình coi nhiệm vụ RE là vấn đề phân loại đa lớp, trong đó mô hình gán một lớp quan hệ cho một câu chứa cặp thực thể được đề cập. [Liu et al. \(2013\)](#) **đề xuất một mô hình CNN đơn giản để trích xuất quan hệ**. Liu et al. (2013) là một trong những nhóm nghiên cứu đầu tiên sử dụng kiến trúc dựa trên Mạng nơ-ron tích chập (CNN) cho trích xuất quan hệ. Mô hình này xây dựng một kiến trúc mạng nơ-ron đầu cuối (end-to-end) với ba khối chính: lớp đầu vào, lớp tích chập và lớp mạng nơ-ron cổ điển.

- Lớp đầu vào: Sử dụng bảng tra cứu (lookup table) để chuyển đổi các câu đầu vào thành vector từ (word vector) bằng cách tận dụng các đặc trưng từ vựng và từ điển đồng nghĩa.
- Lớp tích chập: Sử dụng một kernel tuần tự, ánh xạ các vector từ của lớp đầu vào vào một không gian vector mới.
- Lớp mạng nơ-ron: Kết quả đầu ra của lớp tích chập được đưa vào mạng nơ-ron với hàm softmax để tính toán xác suất phân loại.

Trên tập dữ liệu ACE, mô hình này đạt được điểm F1 là 83,8

Các hướng tiếp cận dựa trên phương pháp CNN:

- Trích xuất thông tin có giám sát dựa trên CNN (CNN based supervised IE)
 - Simple CNN based model([Liu et al. \(2013\)](#))
 - CNN model with max-pooling([Zeng et al. \(2014\)](#))
 - CNN model with multiple window size filter([Nguyen and Grishman \(2015\)](#))
 - CNN model with classification by ranking([Dos Santos et al. \(2015\)](#))

2.2.2 Phương pháp dựa trên RNN và LSTM (RNN and LSTM based methods)

- Simple RNN ([Relation classification via recurrent neural network - Zhang and Wang \(2015\)](#))



- SDP-LSTM ([Bidirectional long short-term memory networks for relation classification - Xu et al. \(2015\)](#))
- BLSTM ([Relation extraction with multi-instance multi-label convolutional neural networks - Zhang et al. \(2015\)](#))
- Att-BLSTM ([Attention-based bidirectional long short-term memory networks for relation classification - Zhou et al. \(2016\)](#))
- DRNN (deep recurrent neural network) ([Improved relation classification by deep recurrent neural networks with data augmentation - Xu et al. \(2016\)](#))
- EAtt-BiGRU([Qin et al. \(2017\)](#))
- RNNOIE ([Stanovsky et al. \(2018\)](#))
- SpanOIE([Zhan and Zhao \(2020\)](#))

2.2.3 Encoder-decoder/transformer based methods

Các mô hình dựa trên bộ mã hóa(encoder) -giải mã (decoder) được đề xuất trong các công trình (Cho et al. (2014), Sutskever et al. (2014)).

Bộ mã hóa (encoder) bao gồm nhiều lớp LSTM hoặc GRU, nhận một dãy các ký hiệu làm đầu vào và tạo ra một vector có độ dài cố định, biểu diễn một cách nén của dãy đầu vào.

Bộ giải mã (decoder) cũng bao gồm nhiều lớp LSTM hoặc GRU, ánh xạ đầu ra của bộ mã hóa thành dãy đích.

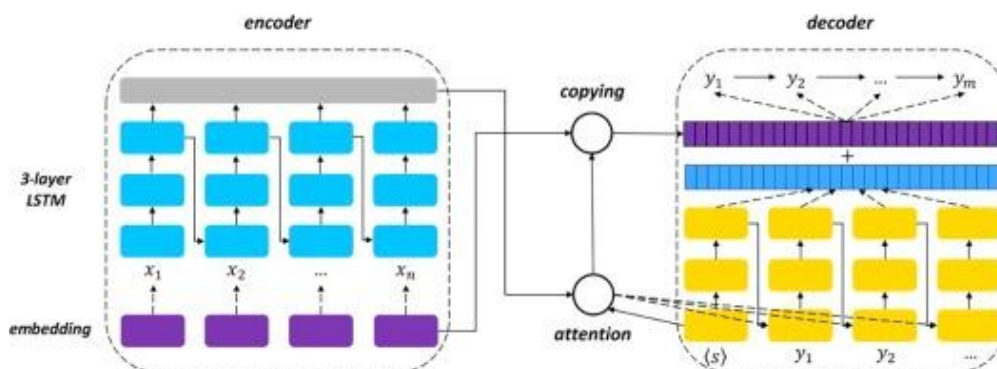
Kiến trúc mã hóa(encoder) -giải mã (decoder) với cơ chế attention được giới thiệu trong (Bahdanau et al. (2015) ; Luong et al. (2015)) đã cải thiện hiệu suất hơn nữa bằng cách cho phép bộ giải mã tập trung vào các phần của dãy đầu vào có liên quan đến từ đích. Vaswani et al. (2017) đã đề xuất kiến trúc Transformer, dựa trên kiến trúc mã hóa-giải mã với self-attention.

2.2.3.a Neural OpenIE (Cui et al. (2018))

Neural OpenIE là một khung dựa trên bộ mã hóa(encoder)-giải mã (decoder), coi nhiệm vụ trích xuất quan hệ (RE) là vấn đề sinh văn bản chuỗi-sang-chuỗi. Trong đó, đầu vào là một chuỗi các token (ký hiệu) và đầu ra là một chuỗi các token với các dấu phân cách chỉ ranh giới của thực thể và quan hệ.

Như được hiển thị trong Figure 5 , kiến trúc Neural OpenIE sử dụng một khung mã hóa (decoder)-giải mã (encoder) với chú ý và sao chép dựa trên chú ý. Trong kiến trúc này, cả bộ mã hóa và bộ giải mã(encoder) đều có 3 lớp mạng LSTM. Bộ mã hóa (decoder) nhận chuỗi văn bản có độ dài thay đổi làm đầu vào và chuyển đổi nó thành một biểu diễn ẩn.

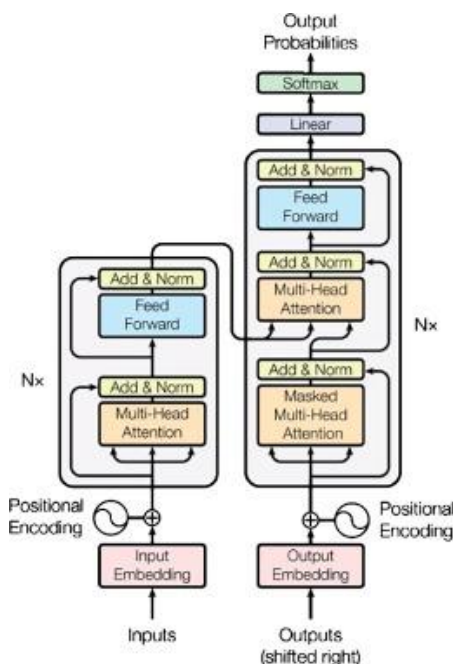
Bộ giải mã (encoder) lấy đầu ra của bộ mã hóa (decoder) và thông tin chú ý làm đầu vào, đưa vào mạng LSTM tiếp theo là softmax để tạo ra chuỗi đầu ra cuối cùng. Mô hình sử dụng cơ chế sao chép để giảm số lượng từ chưa biết trong chuỗi đầu ra.



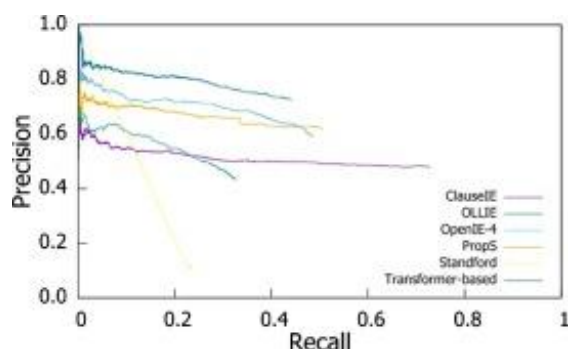
Hình. 5: Kiến trúc mô hình mã hóa (decoder)-giải mã (encoder) của hệ thống Neural OpenIE

2.2.3.b Transformer based OpenIE(Han and Wang (2021))

Bài báo này đề xuất một hệ thống OpenIE dựa trên kiến trúc Transformer (Vaswani et al., 2017). Kiến trúc mô hình Transformer được minh họa trong Figure 6. Bộ mã hóa là một ngăn xếp của nhiều lớp giống hệt nhau, trong đó mỗi lớp bao gồm hai lớp con. Lớp con thứ nhất là lớp self-attention, cho phép bộ mã hóa liên hệ các token khác nhau trong dãy đầu vào. Lớp con thứ hai bao gồm Mạng nơ-ron Feed-Forward (FFN). Bộ giải mã cũng bao gồm hai lớp này, đồng thời có thêm lớp thứ ba với sự chú ý đa đầu (multi-head attention), giúp bộ giải mã xác định các token có liên quan từ dãy đầu vào. So sánh hiệu suất của mô hình dựa trên Transformer với các hệ thống OpenIE khác nhau được thể hiện trong Hình 13.



Hình. 6: Kiến trúc mô hình Transformer (Vaswani et al. (2017))

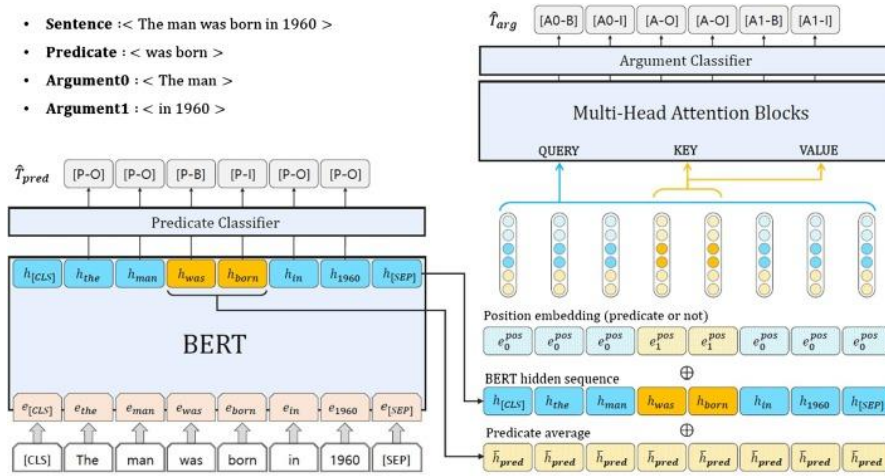


Hình. 7: So sánh hiệu suất của mô hình máy biến áp với các hệ thống OpenIE(Han and Wang (2021))

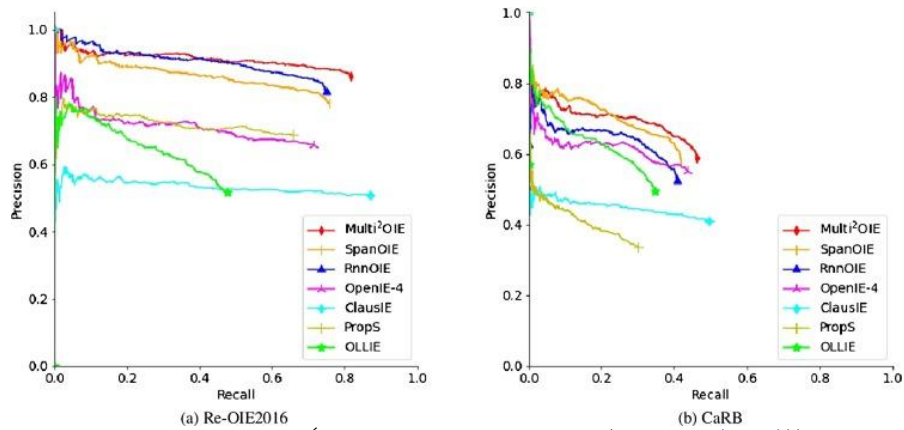
2.2.3.c Multi2OIE (Ro et al. (2020))

Tương tự như RNNOIE (Stanovsky et al. (2018)), Multi2OIE coi nhiệm vụ trích xuất quan hệ (RE) là vấn đề dẫn nhãn dây. Multi2OIE sử dụng BERT (Devlin et al. (2019)) kết hợp với các khối chú ý đa đầu (multi-head attention) (Vaswani et al. (2017)) để thực hiện trích xuất thông tin (IE). Kiến trúc của Multi2OIE được minh họa trong Figure 7. Trong phần đầu tiên, mô hình vị từ (predicate model) sử dụng BERT tiếp theo là mạng nơ-ron Feed-Forward (FFN) và một lớp softmax để xác định vị từ trong câu. Ở phần thứ hai, mô hình luận cứ (argument model) sử dụng trạng thái ẩn của BERT kết hợp với thông tin vị từ và những vị trí làm đầu vào cho khối chú ý đa đầu, tiếp theo là FFN và softmax để dự đoán các luận cứ.

Figure 8 mô tả so sánh hiệu suất của mô hình Multi2OIE với các mô hình OIE khác trên các tập dữ liệu Re-OIE2016 (Zhan and Zhao (2020)) và Bhardwaj et al. (2019).



Hình. 8: Kiến trúc của mô hình Multi2OIE (Ro et al. (2020)).



Hình. 9: Hiệu suất của mô hình Multi2OIE (Ro et al. (2020)).

3 Các MRC phù hợp với Tiếng Việt

3.1 Mô hình tinh chỉnh (Fine-tuned model):

- XML-net fine-tuned: Mô hình này được tinh chỉnh cho tác vụ hỏi đáp.
- Pho-bert fine-tuned: Tương tự như XML-net, Pho-bert cũng là một mô hình được tinh chỉnh cho tác vụ hỏi đáp.
- Vi-deberta fine-tuned (phòng thí nghiệm của chúng tôi), mô hình Vi-deberta được phát triển bởi [Cong Dao Tran, Nhut Huy Pham, Anh Nguyen, Truong Son Hy Tu Vu \(2023\)](#) và được tinh chỉnh cho tác vụ hỏi đáp (QA task). Do được xây dựng dựa trên dữ liệu tiếng Việt, Vi-deberta có khả năng tương thích cao với tập dữ liệu/lĩnh vực của chúng ta (cần kiểm chứng thêm).

3.2 API của bên thứ ba (3rd party API):

- ChatGPT: ChatGPT là một API trả lời theo dạng hội thoại, có khả năng tạo văn bản theo yêu cầu và trả lời các câu hỏi một cách tự nhiên. GPT có nhiều phiên bản khác nhau với kích thước và tập dữ liệu đào tạo khác nhau. GPT-3, phiên bản phổ biến nhất, được đào tạo trên một tập dữ liệu khổng lồ gồm văn bản và mã, bao gồm sách, bài báo,
- Gemini from Google: Gemini là một mô hình ngôn ngữ lớn do Google phát triển, có khả năng thực hiện nhiều tác vụ xử lý ngôn ngữ tự nhiên, bao gồm cả trả lời câu hỏi. Gemini cũng có nhiều phiên bản, bao gồm Gemini Pro, Gemini Ultra và Gemini Nano. Gemini Ultra được đào tạo trên một tập dữ liệu khổng lồ gồm văn bản và mã, tương tự như GPT-3, trong khi Gemini Nano được đào tạo trên một tập dữ liệu nhỏ hơn để chạy trên thiết bị di động.

4 Định nghĩa các loại quan hệ được rút trích lĩnh vực giáo dục

4.1 Entities

4.1.1 Intents

4.1.2 Terms

4.1.3 Papers

4.1.4 Subject

4.1.5 Student

4.1.6 Lecturers

4.1.7 Academic_service_Office

4.1.8 English_Certificate

4.1.9 URL

4.2 Relations

4.2.1 hỏi ý định rút môn

- Intents_cancel_target_subject (Intents target to Subject): - Ex: Intents(name='Rut Mon Hoc') <---Intents_cancel_target_subject--> Subject(name='Toan roi rac')

4.2.2 hỏi chung về ý định

- Student_ask_intent (Student ask about an intent): - Student(id='2033009') <---Student_ask_intent--> Intents(name='Rut Mon Hoc')

4.2.3 hỏi thời gian bắt đầu đăng ký môn

- Student_ask_time_register_course (Student ask about an intent): - Student(id='2033009') <---Student_ask_intent--> Intents(name='Đăng ký môn') - Intents(name='Đăng ký môn') <---time_register_course--> Terms(name='HK201')

4.2.4 hỏi lịch học

- Intents_schedule_course_with_term (Student ask about an intent): - Student(id='2033009') <---Student_ask_intent--> Intents(name='Hỏi lịch học') - Intents(name='Hỏi lịch học') <---Intents_schedule_course_with_term--> Terms(name='HK201')

4.2.5 hỏi môn tương đương

- Subjects_same_Subject: - Subject(name='Toán rời rạc') <---Subjects_same_Subject--> Subject(name='Cấu trúc rời rạc cho KHMT')

4.2.6 đăng kí gia hạn

- Subjects_same_Subject: - Student(id='2033009') <--Student_ask_intent-->
Intents(name='đăng kí gia hạn') - Intents(name='đăng kí gia hạn') <--Student_ask_intent-->
->

4.2.7 hỏi học phí

- Student_paid_tuition: học sinh đã thanh toán học phí học kì 1 -
Student(id='2033009') <--Student_paid_tuition--> Intents(name='Thanh_toán_học_phí',
term=Terms(name='HK201'),)

4.2.8 Thay đổi đăng ký môn học

- Student_change_courses: - Student(id='2033009') <--Student_change_courses--> Intents(
name='Thay đổi đăng ký môn học', term=Terms(name='HK201'),)

4.2.9 hỏi về chuẩn ngoại ngữ

- Intents_qualities_for_student: - Example 1: - Student(id='2033009') <--
Student_ask_intent--> Intents(name='Chuẩn sinh viên') - Intents(name='Chuẩn
sinh viên') <--Intents_qualities_for_student--> English_Certificate() - Example 2: -
Student(id='2033009') <--Student_ask_intent--> Intents(name='Hỏi về chuẩn ngoại
ngữ') - Intents(name='Hỏi về chuẩn ngoại ngữ') <--Intents_qualities_for_student-->
English_Certificate()

4.2.10 URLs

4.2.10.a hỏi đăng kí in bảng điểm

- Intents_print_scores: - Student(id='2033009') <--Student_ask_intent-->
Intents(name='Đăng kí in bảng điểm') - Intents(name='Đăng kí in bảng điểm') <--
Intents_print_scores--> URL('url about print student score')

4.2.10.b hỏi điều kiện nhận luận văn

- Intents_ask_condition_register_thesis: - Student(id='2033009') <--Student_ask_intent-->
-> Intents(name='Hỏi điều kiện nhận luận văn') - Intents(name='Hỏi điều kiện
nhận luận văn') <--Intents_ask_condition_register_thesis--> URL('url about In-
tents_ask_condition_register_thesis')

4.2.10.c hỏi điều kiện tốt nghiệp

- Intents_ask_condition_graduated: - Student(id='2033009') <--Student_ask_intent-->
Intents(name='Hỏi điều kiện tốt nghiệp') - Intents(name='Hỏi điều kiện tốt nghiệp') <--
Intents_ask_condition_graduated--> URL('url about Intents_ask_condition_graduated')



[11pt]article
amsmath