

Mã hóa vị trí đơn giản và hiệu quả cho máy biến áp

Pu-Chin Chen, Henry Tsai, Srinadh Bhojanapalli,
Hyung Won Chung, Yin-Wen Chang, Chun-Sung Ferng

Nghiên cứu của Google

trừu tượng

Các mô hình máy biến áp là hoán vị tự động động.
Để cung cấp thông tin về thứ tự và loại của mã thông
báo đầu vào, phần nhúng vị trí và phân đoạn thư
đọc thêm vào đầu vào. Các công trình gần đây đã
đề xuất các biến thể của mã hóa vị trí với mã hóa
vị trí tự động đối để đạt được hiệu suất tốt hơn.
Chúng tôi
phân tích cho thấy rằng lợi ích thực sự đến
từ việc di chuyển thông tin vị trí đến lớp chú ý từ
đầu vào. Được thúc đẩy bởi điều này,
chúng tôi giới thiệu Chú ý vị trí tách biệt
cho Máy biến áp (DIET), một cơ chế đơn giản nhưng hiệu
quả để mã hóa thông tin vị trí và phân đoạn vào các mô
hình Máy biến áp. Phương pháp đề xuất có thời gian huấn
luyện nhanh hơn
và thời gian suy luận, đồng thời đạt được hiệu
suất cạnh tranh trên GLUE, XTREME và
Điểm chuẩn WMT. Chúng tôi khái quát hơn nữa
phương pháp cho máy biến áp tầm xa và hiển thị
đạt được hiệu suất.

1. Giới thiệu

Máy biến áp là mô hình tuần tự
đạt được hiệu suất hiện đại ở nhiều nơi
Các tác vụ Xử lý ngôn ngữ tự nhiên (NLP), chẳng hạn như
dịch máy, mô hình hóa ngôn ngữ và trả lời câu
hỏi (Vaswani và cộng sự, 2017; Devlin và cộng sự,
2018; Yang và cộng sự, 2019; Liu và cộng sự, 2020). Máy
biến áp có hai thành phần chính: tự chú ý
và một lớp chuyển tiếp nguồn cấp dữ liệu theo vị trí. Cả hai đều là
hoán vị tự động động và không nhạy cảm với
thứ tự của các mã thông báo đầu vào. Để làm cho các mô
hình này nhận biết được vị trí, thông tin vị trí của
các từ đầu vào thường được thêm vào dưới dạng phần
nhúng bổ sung cho phần nhúng mã thông báo đầu vào (Vaswani
và cộng sự, 2017). Ví dụ: nhúng đầu vào (W)
của một câu được thêm vào phần nhúng vị trí
(P), dẫn đến đầu vào W + P cho Máy biến áp.
Những vị trí nhúng này chỉ phụ thuộc vào vị trí mà
từ đó xuất hiện. Đối với các nhiệm vụ nhiều phân đoạn,

các phần nhúng phân đoạn bổ sung có thể được thêm vào chỉ
giống như các phần nhúng vị trí (Devlin và cộng sự, 2018).

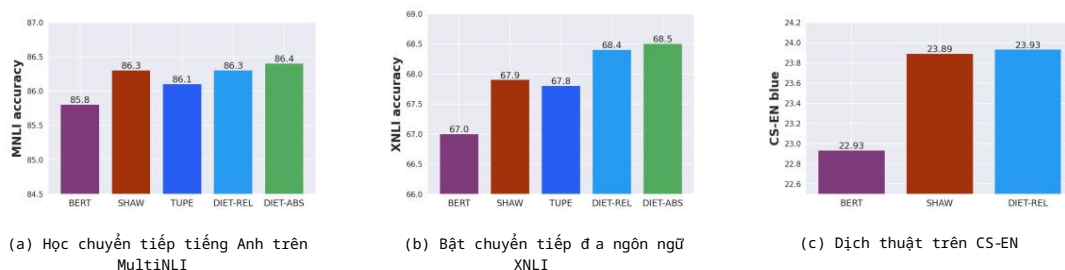
Đã có nhiều công trình khám phá những cách khác
nhau để đưa thông tin vị trí vào Transformers (Shaw
và cộng sự, 2018; Yang và cộng sự, 2019; Raffel và
cộng sự, 2020). Nhiều người trong số họ lưu ý đến
những lợi ích của việc sử dụng sơ đồ mã hóa vị trí tự động đối.
qua mã hóa vị trí tuyệt đối (xem thêm Hình 1).
Tuy nhiên nguyên nhân gây ra sự khác biệt này vẫn chưa rõ ràng.

Yun và cộng sự. (2020) đã chỉ ra rằng Transformers
với mã hóa vị trí tuyệt đối là các hàm xấp xỉ
phổ quát của tất cả các hàm chuỗi với chuỗi,
chứng minh rằng mã hóa vị trí tuyệt đối có thể nắm bắt
được thông tin vị trí. Do đó nguyên nhân gì
tính ưu việt của mã hóa vị trí tự động đối? Một
nghiên cứu có hệ thống và hiểu biết về lợi ích
và thiếu các thuộc tính của các phương pháp mã hóa vị
trí khác nhau. Kể và cộng sự. (2020) đưa ra giả thuyết rằng
mối tương quan chéo giữa từ và vị trí
nhúng trong khi tính toán sự chú ý có thể
nguyên nhân của việc thực hiện kém vị trí tuyệt đối
mã hóa. Tuy nhiên, các điều khoản chéo như vậy có mặt
trong một số phương pháp mã hóa vị trí tự động đối
(Shaw và cộng sự, 2018; Yang và cộng sự, 2019), và những điều này
các phương pháp thực hiện ngang bằng hoặc tốt hơn các phương pháp khác
sơ đồ mã hóa vị trí (xem §4).

Trong bài viết này chúng tôi thực hiện một nghiên cứu có hệ thống
để hiểu các phương pháp mã hóa vị trí khác nhau.
Chúng tôi lập luận rằng việc nhúng vị trí tuyệt đối chủ yếu
chịu đựng việc được thêm vào ở đầu vào. Chúng tôi hiển thị, với
thử nghiệm của chúng tôi về các nhiệm vụ phân loại, trả
lời câu hỏi và dịch máy, mã hóa vị trí tuyệt đối được
thêm vào ma trận chú ý với
các tham số khác nhau cho mỗi đầu cải thiện đáng kể so
với việc thêm mã hóa vị trí tuyệt đối
đến đầu vào. Điều này nhấn mạnh rằng thông tin vị trí
được đưa vào Máy biến áp là
quan trọng, đưa ra lời giải thích cho khoảng cách trong
hiệu suất giữa mã hóa vị trí tuyệt đối và tự động đối. Chúng tôi
cũng so sánh vị trí khác nhau
mã hóa và tác dụng của việc chia sẻ mã hóa vị trí

Các tác giả đóng góp như nhau cho bài viết này.
Email của tác giả tương ứng: puchin@google.com

2021.10.08 08:58:22



Hình 1: Hiệu quả hiệu suất của các phương pháp mã hóa vị trí khác nhau cho Transformers (xem § 2) trên hai bộ dữ liệu Suy luận ngôn ngữ tự nhiên từ GLUE (Wang và cộng sự, 2019), XTREME (Hu và cộng sự, 2020) và một bộ dữ liệu Dịch máy thần kinh WMT 18 (Bojar và cộng sự, 2018). Mã hóa vị trí tuyệt đối (DIET-ABS) có thể đạt được hiệu suất tốt hơn so với mã hóa tương đối (DIET-REL), cho thấy tầm quan trọng của việc thiết kế phương pháp mã hóa vị trí phù hợp.

xuyên qua các đầu và các lớp khác nhau của Máy biến hình.

Dựa trên những quan sát này, chúng tôi đề xuất sự chú ý vị trí tách rời và một phương pháp mã hóa phân đoạn mới (đổi với các nhiệm vụ có nhiều phân đoạn) và cho thấy tính ưu việt của nó bằng thực nghiệm.

Chúng tôi tóm tắt những đóng góp của chúng tôi trong bài viết này dưới đây.

- Chúng tôi phân tích về mặt lý thuyết và thực nghiệm giới hạn của việc nhúng vị trí tuyệt đối được thêm vào đầu vào. Đối với cả thông tin tuyệt đối và tương đối, chúng tôi cho thấy rằng vị trí mã hóa thành ma trận chú ý trên mỗi đầu mang lại hiệu suất vượt trội.
- Chúng tôi đề xuất một cách đơn giản và hiệu quả để mã hóa thông tin vị trí và phân đoạn. Mã hóa được đề xuất phù hợp với các phương pháp SoTA trên nhiều tác vụ NLP tiêu chuẩn trong khi có một mô hình đơn giản hơn với chi phí đào tạo/suy luận thấp hơn.
- Phương pháp đề xuất của chúng tôi có thể dễ dàng áp dụng cho các mô hình chuỗi dài (DIET-ABSLIN) và cải thiện tất cả các chỉ số so với Linformer (Wang et al., 2020).
- Chúng tôi trình bày các nghiên cứu cắt bỏ so sánh các phương pháp mã hóa vị trí khác nhau và các cách chia sẻ thông số mã hóa vị trí giữa các đầu và các lớp trong Transformer.

## Mã hóa 2 vị trí cho máy biến áp

Trong phần này, chúng tôi xem xét ngắn gọn các mô hình Máy biến áp (Vaswani và cộng sự, 2017) và thảo luận về cải tiến trước đây về mã hóa vị trí cũng như phân tích hạn chế của việc nhúng vị trí phụ gia được đề xuất trong mô hình Máy biến áp ban đầu và được áp dụng rộng rãi.

### 2.1 Máy biến áp

Khối Transformer bao gồm hai loại lớp: 1) Lớp tự chú ý và 2) Lớp chuyển tiếp nguồn cấp dữ liệu.

Mô-đun tự chú ý Với độ dài chuỗi đầu vào  $n$ , kích thước ẩn  $d$ , kích thước chiếu xuống khóa truy vấn nhiều đầu  $d_h$ , chúng tôi xác định lớp ẩn truy vấn

đầu vào cho đầu chú ý này dư dật dạng  $n \times d$ , ma trận chiếu  $X$  dư dật dạng  $n \times d_h$  và ma trận chiếu giá trị dư dật dạng  $n \times d_v$ . Ma trận chiếu giá trị dư dật dạng  $n \times d_v$  và ma trận chiếu khóa truy vấn dư dật dạng  $n \times d_h$  được kết hợp để tính toán ma trận tương tác như trong [h], cho đầu  $h$ .

Thông thường,  $d_h < d$  khi chúng ta thực hiện chú ý nhiều đầu với biểu diễn nhỏ hơn trên mỗi đầu ( $d_h = d/h$ ). Với điều đó, chúng ta có thể viết điểm chú ý của sản phẩm chấm:

$$A_i = (XW_i Q)(XW_i K)$$

Điểm chú ý này được sử dụng để tính toán đầu ra cho mỗi đầu, sau khi chia tỷ lệ và chuẩn hóa mỗi hàng bằng softmax:

$$\text{head}_i = \text{Softmax}(A_i / \sqrt{d}) \cdot (XW_i V)$$

Đầu ra của tất cả các đầu chú ý trong một lớp được nối và chuyển sang lớp chuyển tiếp tiếp theo được áp dụng mã thông báo.

### 2.2 Nhận thức được vị trí của bản thân

Nhiều tác vụ NLP, chẳng hạn như dịch máy, mô hình hóa ngôn ngữ, rất nhạy cảm với thứ tự các từ đầu vào. Vì Transformers là hoán vị tương đối nên chúng tôi thường đưa thêm thông tin vị trí vào đầu vào. Dưới đây chúng tôi thảo luận về một số phương pháp mã hóa vị trí phổ biến.

#### 2.2.1 Mã hóa vị trí tuyệt đối Mã hóa vị trí

tuyệt đối được tính toán trong lớp đầu vào và được tổng hợp bằng các mã nhúng mã thông báo đầu vào. Vaswani và cộng sự. (2017) đề xuất điều này

cho Transformers và nó đã là một lựa chọn phổ biến trong các công trình tiếp theo (Radford và cộng sự, 2018; Devlin và cộng sự, 2018). Có hai biến thể phổ biến của mã hóa vị trí tuyệt đối - cố định và học được.

2.2.2 Mã hóa vị trí tương đối

Một nhược điểm của mã hóa vị trí tuyệt đối là nó yêu cầu độ dài cố định của chuỗi đầu vào và không không trực tiếp nắm bắt được vị trí tương đối của từng từ. Để giải quyết những vấn đề này một số vị trí tương đối các phương án đã được đề xuất.

Shaw và cộng sự. (2018) đề xuất sử dụng mã hóa vị trí tương đối thay vì mã hóa vị trí tuyệt đối, và thêm các phần nhúng vị trí vào các phép chiếu khóa và giá trị tùy chọn thay vì đầu vào. Họ cho thấy cách mã hóa thông tin vị trí mới này dẫn đến hiệu suất tốt hơn trên máy nhiệm vụ dịch thuật. Yang và cộng sự. (2019) đã đơn giản hóa điều này bằng cách loại bỏ các phần nhúng vị trí trong các phép chiếu giá trị và cho thấy hiệu suất tốt hơn trong các nhiệm vụ mô hình hóa ngôn ngữ. Cả hai cách tiếp cận này đều sử dụng một biểu diễn vector để mã hóa thông tin vị trí .

Raffel và cộng sự. (2020) sử dụng đại số vô hướng để mã hóa vị trí tương đối giữa truy vấn và chỉ số chính và thêm trực tiếp vào ma trận điểm chú ý. Họ tiếp tục sử dụng việc gộp logarit của thông tin vị trí vào một số nhóm cố định. Tất cả những điều này

phương pháp vị trí tương đối chia sẻ thêm vị trí các tham số mã hóa trên các lớp.

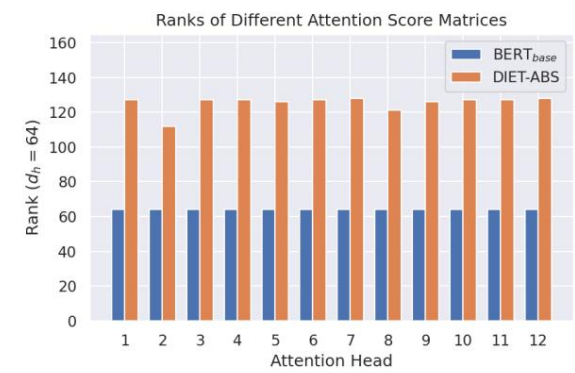
Gần đây Ke và cộng sự. (2020) đưa ra giả thuyết rằng mối tương quan chéo giữa vị trí và phần nhúng mã thông báo có thể dẫn đến hiệu suất yếu hơn của phần nhúng vị trí tuyệt đối quảng cáo và thay vào đó đề xuất bổ sung thêm sự chú ý dựa trên cả thông tin vị trí tuyệt đối và tương đối trực tiếp vào từng cái đầu. Tuy nhiên, các thuật ngữ chéo như vậy có mặt trong phương pháp được đề xuất bởi Shaw et al. (2018), điều đó làm cạnh tranh với các phương pháp khác. Thay vào đó chúng tôi đưa ra giả thuyết rằng mã hóa vị trí ở giới hạn đầu vào thứ hạng của ma trận chú ý vị trí dẫn đến hiệu suất kém của nó.

2.3 Hạn chế của phụ gia đầu vào

Nhúng vị trí

Trong phần này chúng ta thảo luận về một số hạn chế của cách thực tế để thêm mã hóa vị trí tuyệt đối vào phần nhúng mã thông báo đầu vào.

Đầu tiên chúng ta so sánh quyền lực đại diện trong điều khoản về thứ hạng của ma trận chú ý có thể đạt được với các mã hóa vị trí khác nhau.



Hình 2: Thứ hạng của ma trận chú ý: Chúng tôi trình bày một so sánh thứ hạng của ma trận điểm chú ý của mô hình BERTBASE với các phần nhúng vị trí tuyệt đối ở đầu vào so với các phần nhúng vị trí tuyệt đối trên mỗi đầu (DIET-ABS (1)). Với tính năng nhúng vị trí bổ sung ở đầu vào, ma trận chú ý có mức độ chú ý thấp hơn nhiều cấp bậc, hạn chế quyền lực đại diện. Điều này được giảm bớt nhờ DIET-ABS.

Định lý 1. Cho  $P \in \mathbb{R}^{n \times d}$  là vị trí đầu vào nhúng và  $P^* \in \mathbb{R}^{n \times d_p}$  là các nhúng vị trí theo lớp. Đặt  $W_Q, W_K \in \mathbb{R}^{d \times d_h}$  là ma trận truy vấn và phép chiếu khóa có đầu kích thước chiều  $d_h$ , và  $d_h < d_p$ ,  $d$  và  $n \geq d_h + d_p$ . Đặt  $A_a = (X + P)W_Q W_K^T (X + P)^T$  và  $A_r = XW_Q W_K^T X^T + P^* P^{*T}$  là ma trận chú ý được tính toán bằng cách sử dụng đầu vào và theo lớp vị trí nhúng tương ứng. Sau đó với bất kỳ  $X, P, W_Q, W_K$

$$\text{hạng}(A_a) \leq d_h.$$

Tồn tại sự lựa chọn  $X, P^*, W_Q, W_K$  sao cho

$$\text{hạng}(A_r) = d_p + d_h > d_h.$$

Nhận xét. Định lý này cho chúng ta thấy rằng thứ hạng của ma trận chú ý bị ràng buộc với giá trị tuyệt đối mã hóa vị trí ở đầu vào và sử dụng trên mỗi đầu mã hóa vị trí bằng cách thêm thông tin vị trí vào ma trận chú ý trực tiếp dẫn đến việc cho phép sự chú ý cấp cao hơn. Xem § B để chứng minh.

Việc thêm mã hóa vị trí trực tiếp vào đầu vào càng đặt ra một hạn chế đối với động lực đào tạo bằng cách buộc các gradient phải giống nhau cho cả đầu vào nhúng mã thông báo và vị trí (xem § B). Liên quan đến mã hóa vị trí đã thảo luận trước đó, trong khi giải quyết một số mối lo ngại này, lại gặp phải tình trạng chậm hơn. thời gian đào tạo/suy luận (xem Bảng 1) với các cách triển khai phức tạp (Shaw và cộng sự (2018); Ke và cộng sự. (2020)). Trong phần tiếp theo, chúng tôi trình bày đơn giản phương pháp mã hóa vị trí để tránh những hạn chế này.

### 3 Vị trí và phân khúc đề xuất Mã hóa

Trong phần trước, chúng ta đã tìm hiểu về những hạn chế của việc nhúng vị trí phụ gia đầu vào và các công việc hiện có. Dựa trên những quan sát này, chúng tôi đề xuất hai cách tối thiểu/hiệu quả để kết hợp mã hóa vị trí (tuyệt đối/tương đối) cùng với phương pháp mã hóa phân đoạn tuyệt đối mới. Bằng cách tách vị trí và phân đoạn khỏi phần nhúng mã thông báo, chúng tôi điều chỉnh hiệu suất SoTA đồng thời cải thiện thời gian đào tạo/suy luận (xem [§3.3](#)).

#### 3.1 Chú ý vị trí tuyệt đối được tách rời Chúng tôi

đề xuất phương pháp mã hóa vị trí tuyệt đối đơn giản sau đây để thêm thông tin vị trí vào ma trận chú ý mã thông báo trực tiếp trong mỗi đầu chú ý. Ngoài ra, chúng tôi cũng thêm thông tin phân đoạn vào mã thông báo chú ý thay vì nhúng đầu vào.

Bằng cách này, chúng ta có thể đặt thứ hạng của mã hóa vị trí một cách độc lập, dẫn đến ma trận chú ý có thứ hạng cao hơn, giải quyết các hạn chế đã thảo luận trước đó.

DIET-ABS

$$\begin{aligned} \text{AABS} &= (\text{Xi:WQ})(\text{Xj:WK}) / \sqrt{d_{i,j}} \\ &+ (\text{PQP K})_{i,j} + \text{ES}(\text{S}(i), \text{S}(j)), \end{aligned} \tag{1}$$

trong đó PQ, PK  $\in \mathbb{R}^{n \times dp}$  là các ma trận xếp vị trí cấp thấp và ES là sự chú ý tuyệt đối của phân khúc đối với các tương tác mô hình giữa các phân đoạn được xác định là

$$\begin{aligned} \text{ES}(\text{S}(i), \text{S}(j)) &= \hat{\text{S}}^i, \hat{j} \\ \text{trong đó } \text{S}(i) &= \hat{i} \text{ nếu chỉ số } i \text{ nằm trong phân đoạn } \hat{i}. \end{aligned} \tag{2}$$

Xin lưu ý rằng chúng tôi sử dụng ký hiệu sau trong phương trình trên.  $\text{Ai}, \text{j}$  biểu thị phần tử  $(i, j)$  của ma trận A.  $\text{Xi}$ : và  $\text{X:j}$  lần lượt biểu thị hàng thứ  $i$  và cột thứ  $j$  của X. Chúng ta sẽ theo ký hiệu này trong phần còn lại của bài viết.

Theo mặc định, chúng tôi đặt dp giống như dh. Điều này đã dẫn đến ma trận chú ý có thứ hạng  $dp+dh$  như được hiển thị trong Định lý 1. Để minh họa điều này, chúng tôi so sánh thứ hạng của ma trận chú ý trong lớp đầu tiên của mô hình BERT cơ sở và mô hình DIET-ABS cho một mẫu được lấy mẫu. Lô trong Hình 2. Hình này cho thấy ma trận chú ý của DIET-ABS có thứ hạng cao hơn BERT cơ sở. Kết quả thử nghiệm chi tiết của chúng tôi trong § 4 cũng cho thấy DIET-ABS hoạt động tốt hơn rõ rệt. Điều này xác nhận quan sát trước đó của chúng tôi trong Định lý 1 rằng việc nhúng vị trí cộng vào ở đầu vào có thể hạn chế mô hình và

việc thêm các phần nhúng vị trí trên mỗi đầu sẽ loại bỏ ràng buộc này và mang lại hiệu suất tốt hơn.

Với việc nhúng vị trí tách rời, chúng ta có thể tăng dp lên bất kỳ chiều rộng k nào để phá vỡ nút cổ chai cấp thấp được hiển thị trong Định lý 1. Chúng tôi gọi đó là mô hình DIET-ABS-Rank-k. Chúng tôi cũng giải quyết hiệu quả vấn đề chính xác được đưa ra bởi một phép nhân ma trận bổ sung (PQP K). Vì việc nhúng vị trí độc lập với đầu vào nên chúng ta chỉ cần tính toán phép nhân ma trận một lần cho mỗi đợt huấn luyện và chúng ta có thể lưu trữ ma trận đã tính toán trước khi chạy suy luận. Kết quả là, chúng tôi nhận thấy chi phí suy luận và đào tạo tăng lên đáng kể trong biến thể mô hình này.

#### 3.2 Sự chú ý theo vị trí tương đối được tách rời

Để kết hợp độ lệch quy nạp vị trí tương đối, chúng tôi xem xét một phiên bản đơn giản của mã hóa vị trí được đề xuất trong T5 ([Raffel và cộng sự, 2020](#)) mà không cần chia sẻ tham số theo lớp và tạo nhóm nhật ký. Ngoài ra, chúng tôi còn kết hợp mã hóa phân khúc theo đầu ngữ rời như trong DIET-ABS. Mô hình có thể được viết

BẢNG:

DIET-REL

$$\begin{aligned} \text{AREL}_{1,j} &= (\text{Xi:WQ})(\text{Xj:WK}) / \sqrt{d} + \text{Ri}_j + \\ &\text{ES}(\text{S}(i), \text{S}(j)). \end{aligned} \tag{3}$$

Chúng tôi đưa ra một ví dụ về mô hình này với hai phân đoạn trong Hình 3.

#### 3.3 Chi phí đào tạo và suy luận

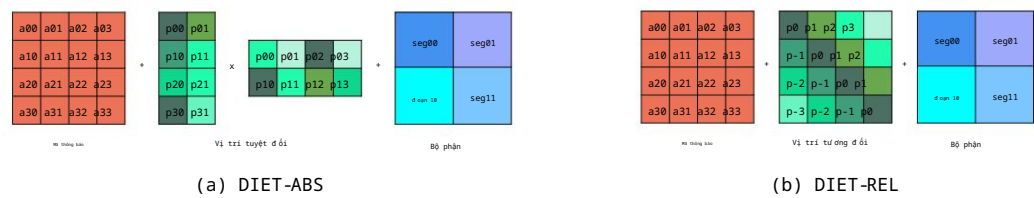
Tiếp theo, chúng tôi trình bày các mô hình được đề xuất có ít chi phí tính toán hơn so với mô hình cơ sở, làm cho mô hình của chúng tôi có tính thực tế hơn các lựa chọn thay thế. Chúng tôi xem xét hai mô hình khác nhau - mô hình BERTBASE và một mô hình nhỏ hơn, BERTSMALL, có kích thước ẩn 512, 4 lớp và 8 đầu chú ý.

Trong Bảng 1, chúng tôi so sánh chi phí đào tạo và suy luận của các phương pháp mã hóa vị trí của [Shaw et al. \(2018\)](#), [Kế và cộng sự. \(2020\)](#), DIET-ABS và DIET-REL.

Chúng tôi nhận thấy rằng tính đơn giản của các phương pháp được đề xuất thực sự giúp tiết kiệm cả thời gian huấn luyện và suy luận so với các phương pháp mã hóa vị trí khác. Việc tiết kiệm thời gian bù đắp thậm chí còn đáng kể hơn đối với các mô hình nhỏ hơn (BERTSMALL) và trong quá trình suy luận.

Lưu ý rằng có thể có sự khác biệt giữa tốc độ đào tạo và tốc độ suy luận do cập nhật độ dốc chỉ phối chi phí tại thời điểm đào tạo ([Lan và cộng sự, 2020](#)).

Tại thời điểm suy luận, chúng ta chỉ đo thời gian của một



Hình 3: Phương pháp tiếp cận hiệu quả được đề xuất để bao gồm mã hóa vị trí và phân đoạn bằng cách thêm chúng trực tiếp vào ma trận chú ý mã thông báo trên mỗi đầu ngữ. Hình bên trái cho thấy cách chúng tôi mã hóa sự chú ý theo vị trí tuyệt đối. Hình bên phải đại diện cho sự chú ý vị trí tương đối.

		Chế độ Shaw và cộng sự. (2018)	Kế và cộng sự. (2020)	DIET-ABS	DIET-REL
BERTBASE	Đào tạo	+13%	+1%	+0%	+0%
Suy luận BERTBASE	Đào tạo	+33%	+19%	+0%	+0%
tạo BERTSMALL	Suy	+24%	+4%	+0%	+0%
luận BERTSMALL		+65%	+27%	+1%	+0%

Bảng 1: So sánh thời gian huấn luyện và suy luận của Transformers với các phương pháp mã hóa vị trí khác nhau với mô hình BERT cơ bản trên TPU v2. Chúng tôi nhận thấy rằng sự đơn giản của DIET-REL và DIET-ABS dẫn đến tăng đáng kể cả về thời gian huấn luyện và suy luận. Chúng tôi nhận thấy tốc độ thậm chí còn nhanh hơn đối với BERTSMALL nhờ mô hình so với BERTBASE.

chuyển tiếp tương ứng với chi phí sử dụng

những mô hình như vậy trong các hệ thống thực.

3.4 Ứng dụng cho máy biến áp tầm xa

Một ưu điểm khác của phương pháp đề xuất của chúng tôi là chúng dễ dàng mở rộng sang các mẫu máy biến áp tầm xa. Đối với các đầu vào chuỗi dài, Máy biến áp phải chịu sự phụ thuộc bậc hai của khả năng tính toán độ phức tạp liên quan đến độ dài chuỗi. Một lớp phương pháp này làm giảm sự phức tạp này bằng cách sử dụng một phép chiếu thứ hạng thấp của chuỗi đầu vào để tính toán sự chú ý (Wang và cộng sự, 2020; Choromanski và cộng sự, 2021; Đại và cộng sự, 2020). Tuy nhiên, các phương pháp như vậy sử dụng mã hóa vị trí đầu vào mặc định và chưa có nhiều công việc trong việc kết hợp thông tin vị trí trên đầu ngữ ời mà không đưa a vào độ phức tạp tính toán bậc hai trên độ dài chuỗi đầu vào. Chúng tôi minh họa khả năng ứng dụng của chúng tôi các phương pháp cài đặt như vậy bằng cách áp dụng DIET-ABS cho Linformer (Wang và cộng sự, 2020), dự kiến ma trận khóa và giá trị chú ý đến thứ nguyên k thấp hơn trong quá trình tính toán chú ý.

DIET-ABSLIN Phương pháp được đề xuất có thể là Viết như :

$$ALIN_{i,j} = (X_i W_Q)((E X_j) W_K) / \sqrt{d} + (P Q K)_{i,j},$$

trong đó  $E \in \mathbb{R}^{k \times n}$ ,  $P Q \in \mathbb{R}^{n \times d}$ ,  $P K \in \mathbb{R}^{k \times d}$ .

4 thí nghiệm

Trong phần này chúng tôi trình bày kết quả thực nghiệm so sánh các phương pháp mã hóa vị trí và phân đoạn khác nhau đã được thảo luận trong các phần trước. Chúng tôi

tiến hành thí nghiệm ở ba môi trường khác nhau để bao gồm một loạt các trường hợp sử dụng. Đầu tiên, chúng tôi kiểm tra Kết quả của phương pháp học chuyển giao phổ biến từ việc đào tạo trước LM có mặt nạ đến các nhiệm vụ cuối cùng trong KEO (Devlin và cộng sự, 2018). Thứ hai, chúng tôi nghiên cứu khả năng chuyển đổi đa ngôn ngữ không cần bản của các mô hình được đào tạo trước đa ngôn ngữ (Hu và cộng sự, 2020) sang nhiệm vụ phân loại và trả lời câu hỏi trong Điểm chuẩn XTREME (Hu và cộng sự, 2020). Cuối cùng, chúng tôi xem xét đào tạo mô hình Transformer từ đầu cho dịch máy.

Chúng tôi so sánh mã hóa vị trí sau đây các phương pháp tiếp cận - nhưng vị trí tuyệt đối (De-vlin và cộng sự, 2018), nhưng vị trí tương đối (Shaw và cộng sự, 2018), kết hợp mã hóa vị trí tuyệt đối và tương đối (Ke và cộng sự, 2020), tương đối phương pháp vô hướng (Raffel và cộng sự, 2020), đề xuất của chúng tôi Phương pháp mã hóa vị trí trên đầu DIET-ABS và DIET-REL. Chúng tôi biểu thị các phương pháp thêm thông tin vị trí/phân đoạn trực tiếp vào đầu vào nhưng mã thông báo với đầu vào và các phương thức thêm thông tin vị trí/phân khúc trực tiếp được chú ý lớp với mỗi đầu. Để thực nghiệm hoàn chỉnh thiết lập, xem Phụ lục A.

4.1 Kết quả học chuyển tiếp tiếng Anh

Bộ dữ liệu và mô hình Để đào tạo trước, chúng tôi sử dụng Bộ dữ liệu Wikipedia và Sách tiếng Anh (Devlin và cộng sự, 2018). Đối với các tác vụ Tinh chỉnh, chúng tôi sử dụng bộ dữ liệu từ điểm chuẩn GLUE (Wang và cộng sự, 2019). Chúng tôi áp dụng mã thông báo từ phụ trên bản thô dữ liệu văn bản bằng WordPiece (Wu và cộng sự, 2016) với 30.000 từ vựng mã thông báo.

Người đời mẫu	Phân đoạn và vị trí		MNLI QQP QNLI SST2 CoLa STS-B 393k 8.5k					Trung bình
			364k	105k	67k	7k		
Devlin và cộng sự. (2018) <a href="#">đầu vào</a> 85,8 / 85,9 91,1 86,3 / 86,0 91,2 <a href="#">Raffel et al. (2020)</a> bình quân đầu người 86,1 / 86,2 91,2 DIET-REL bình quân đầu người 86,3 / 86,3 91,0 DIET-ABS (dp=128, chia sẻ) bình quân đầu người 86,4 / 86,4 90,8	Shaw et al. (2018) <a href="#">bình quân đầu người</a> 86,4 / 86,2 91,2 <a href="#">Ke et al. (2020)</a> bình quân đầu người 86,0 / 86,1 91,0 DIET-REL bình quân đầu người 86,3 / 86,3 91,0 DIET-ABS (dp=128, chia sẻ) bình quân đầu người 86,4 / 86,4 90,8		89,9	93,2	58,7	89,0	84,8	
			90,5	93,2	59,8	89,3	85,2	
			90,1	93,0	59,6	90,1	85,2	
			90,3	93,1	59,6	89,6	85,2	
			89,8	92,8	59,6	89,0	84,9	
			90,5	92,9	60,3	89,3	85,2	
			90,6	92,8	60,1	89,4	85,3	
			89,5	93,0	59,8	90,2	85,2	
Vương và cộng sự. (2020) (dp=32)	đầu	đầu vào	82,3 / 82,6 90,2	86,3	91,4	53,9	87,6	82,0
DIET-ABSLIN (dp=32)	vào mỗi đầu	đầu vào	83,0 / 83,1 90,6	86,7	92,0	55,7	87,6	82,7

Bảng 2: GLUE: Kết quả trên tập phát triển GLUE của các mô hình được tinh chỉnh dựa trên mô hình được huấn luyện trước với 12-kiến trúc lớp BERTBASE. Chúng tôi báo cáo mức trung bình của độ chính xác tối đa trên tất cả các điểm kiểm tra trong số năm điểm chạy. Chúng tôi nhận thấy rằng DIET-ABS được chia sẻ với thứ hạng 128 hoạt động có tính cạnh tranh với vị trí tương đương đối diện có những các mô hình SoTA mà không có sai lệch quy nạp của các vị trí tương đương đối. Phương pháp đề xuất cũng cải thiện hiệu suất trong cài đặt máy biến áp tầm xa cấp thấp của (Wang và cộng sự, 2020), trong đó vị trí tương đương đối phương pháp những không hiệu quả để sử dụng.

Người đời mẫu	Phân đoạn vị trí	Phân loại XNLI 393k	Trả lời câu hỏi XQuAD MLQA TyDiQA Trung bình 88k 3,7k			
Devlin và cộng sự. (2018) <a href="#">đầu vào</a> 67,0 <a href="#">Shaw và cộng sự. (2018)</a> bình quân đầu người 67,9 <a href="#">Raffel và cộng sự. (2020)</a> bình quân đầu người 68,5 <a href="#">Kế và cộng sự. (2020)</a> bình quân đầu người 67,8 DIET-REL vào bình quân đầu người 68,0 DIET-REL đầu vào bình quân đầu người 68,4 DIET-ABS (dp=128, chia sẻ) 68,5	đầu vào đầu vào bình quân đầu người đầu vào bình quân đầu người đầu vào bình quân đầu người đầu vào bình quân đầu người đầu vào bình quân đầu người đầu vào bình quân đầu người		66,0 / 49,9 56,2 / 41,0 59,0 / 47,9 55,3	69,5 / 53,9 58,2 / 43,1 64,8 / 49,9 58,2	69,9 / 53,5 59,5 / 44,3 63,8 / 50,6 58,6	68,6 / 52,0 58,6 / 43,2 63,9 / 48,7 57,5
			68,1 / 52,8 57,7 / 42,7 63,3 / 50,9 57,6	69,4 / 54,4 58,6 / 43,5 62,4 / 49,3 58,0	70,0 / 53,6 59,8 / 44,5 64,6 / 51,5 58,9	
			59,1 / 43,7 48,9 / 34,0 50,5 / 37,9 48,2	61,6 / 46,0 52,2 / 37,0 53,6 / 40,9 50,8		
Vương và cộng sự. (2020) (dp=256)	đầu	đầu vào	63,6	59,1 / 43,7 48,9 / 34,0 50,5 / 37,9 48,2		
DIET-ABSLIN (dp=256)	vào mỗi đầu	đầu vào	64,4	61,6 / 46,0 52,2 / 37,0 53,6 / 40,9 50,8		

Bảng 3: XTREME: Tinh chỉnh mô hình đa ngôn ngữ trên tập huấn luyện tiếng Anh (Chuyển giao ngôn ngữ chéo). Hiệu suất được đo bằng độ chính xác để phân loại và điểm F1/khớp chính xác để trả lời câu hỏi. Đồng ý với kết quả trong Bảng 2, chúng tôi thấy trong bảng này rằng việc sử dụng mã hóa vị trí trên mỗi đầu hoàn toàn tốt hơn vị trí tuyệt đối mã hóa ở đầu vào. Với tính năng chia sẻ theo lớp, DIET-ABS xếp hạng 128 vượt trội hơn tất cả các mẫu SoTA.

Người đời mẫu	EN-DE	DE-EN	EN-CS	CS-EN
Vaswani và cộng sự. (2017) <a href="#">Shaw và cộng sự. (2018)</a> 40 DIET-REL 39,47 38,49 18,68 23,93	39,00	38,42	18,55	22,93

Bảng 4: Dịch máy: Chúng tôi báo cáo kết quả so sánh các phương pháp mã hóa vị trí khác nhau cho Trans-formers trên các tác vụ dịch máy en-de, de-en, en-cs và cs-en từ tập dữ liệu Newstest 2018. Chúng tôi lưu ý rằng tất cả sơ đồ mã hóa vị trí trên mỗi đầu người (tất cả ngoại trừ hàng đầu tiên) đều hoạt động tốt hơn vị trí tuyệt đối các phần những được thêm vào ở đầu vào. Hơn nữa đề xuất Cách tiếp cận DIET-REL đơn giản có tính cạnh tranh với các phương pháp khác các phương pháp mã hóa vị trí

Kết quả Chúng tôi kiểm tra các cách khác nhau của vị trí và phân đoạn mã hóa ảnh hưởng đến việc truyền khả năng học tập của người học tiếng Anh BERT tiền đề tạo mô hình bằng cách tinh chỉnh trên điểm chuẩn GLUE (Wang và cộng sự, 2019), và trình bày kết quả trong Ta-

ble 2. Đầu tiên chúng tôi nhận thấy rằng tất cả các phương pháp tiếp cận mã hóa các tính năng vị trí rõ ràng ở cấp độ mỗi đầu hoạt động tốt hơn vị trí phụ gia cơ bản mã hóa ở đầu vào (Devlin và cộng sự, 2018). Tất cả mô hình kết hợp các vị trí tương đương đối (Shaw et al., 2018; Raffel và cộng sự, 2020; Ke và cộng sự, 2020), mặc dù sự khác biệt về mô hình của họ, có điểm trung bình rất giống nhau. Chúng tôi cho thấy mức tăng thêm (84,9 đến 85,2 cho DIET-REL) bằng cách di chuyển các đặc điểm của phần khúc sang trên đầu người.

Điều thú vị là chúng tôi nhận thấy rằng phương pháp mã hóa vị trí tuyệt đối được đề xuất DIET-ABS, với chia sẻ theo lớp, ngang bằng với tất cả các tính năng chia sẻ trước đây Mã hóa vị trí tương đương đối SoTA. Màn trình diễn này mà ngay cả mã hóa vị trí tuyệt đối cũng có thể thực hiện tốt hơn khi được bao gồm trên đầu người thay vì ở đầu vào. Chúng tôi trình bày một nghiên cứu cắt bỏ chi tiết khác nhau phương pháp xếp hạng và chia sẻ vị trí tuyệt đối



chú ý (DIET-ABS) trong Bảng 8 và Bảng 9 trong Phụ lục C.

Đối với đầu vào tầm xa, chúng tôi xem xét Linformer (Wang và cộng sự, 2020) với kích thước hình chiếu là 32. Do chiều xuống, chúng tôi thấy hiệu suất giảm không đáng kể khi so sánh với Máy biến áp. Ngay cả đối với cài đặt này, chúng tôi thấy rằng sự chú ý tuyệt đối của chúng tôi DIET-ABS có thể được sử dụng để cải thiện hiệu suất của mô hình.

4.2 Kết quả mô hình đa ngôn ngữ

Bộ dữ liệu và mô hình Đối với các thử nghiệm đa ngôn ngữ, chúng tôi đào tạo trước các mô hình trên Wikipedia kho văn bản bằng 100 ngôn ngữ tự động tự như (Lample và Conneau, 2019) cho 125 nghìn bực có trình tự dài 512, sau đó tinh chỉnh ở phía dư thừa Nhiệm vụ của XTREME (Hu và cộng sự, 2020). Chúng tôi sử dụng mã thông báo độc lập với ngôn ngữ, Mảnh câu (Kudo và Richardson, 2018) , với 120.000 từ vựng mã thông báo để mã hóa văn bản đầu vào.

Phân loại Chúng tôi tiến hành 5 thử nghiệm tinh chỉnh cho từng mô hình trên MultinLI (Williams và cộng sự, 2018) , sau đó thực hiện dự đoán không bắn trúng XNLI (Conneau và cộng sự, 2018), chọn độ chính xác trung bình để báo cáo.

Trả lời câu hỏi Chúng tôi tiến hành 5 thử nghiệm tinh chỉnh cho từng mô hình trên bộ dữ liệu SQuAD V1.1, theo sau là các dự đoán không bắn trên XQuAD (11 ngôn ngữ), MLQA (7 ngôn ngữ) và TyDiQA-GoldP (9 ngôn ngữ), chọn trung vị F1/EM điểm để báo cáo.

Kết quả Chúng tôi trình bày kết quả của mình về việc phân loại và trả lời câu hỏi trong các nhiệm vụ tinh chỉnh trong XTREME cho các phương pháp mã hóa vị trí và phân đoạn khác nhau trong Bảng 3. Một lần nữa, tất cả các phương pháp mã hóa trên mỗi đầu các phương pháp mã hóa vị trí tốt hơn các phương pháp mã hóa vị trí bổ sung đầu vào. Thật thú vị, đơn giản của chúng tôi DIET-ABS hóa ra là mẫu tốt nhất, tốt hơn hơn các mô hình khác sử dụng tính năng vị trí tự động đối. Chia sẻ theo lớp và chú ý đến phân khúc theo đầu ngữ để cho phép DIET-ABS hoạt động tốt hơn DIET-REL. Chúng tôi trình bày một nghiên cứu cắt bỏ chi tiết trong Bảng 5 để hiểu tác dụng của sự chú ý theo vị trí tách rời các biến thể. Cuối cùng, chúng tôi nhận thấy những lợi thế tự động tự trong sử dụng DIET-ABS với Linformer (Wang và cộng sự, 2020) trong cài đặt tầm xa.

4.3 Kết quả dịch thuật

Bộ dữ liệu và mô hình cho dịch máy nhiệm vụ chúng tôi xem xét hai cặp ngôn ngữ (cả hai đều trực tiếp

tions) dành cho đào tạo - WMT 2018 English-to-Đức (en-de), tiếng Đức sang tiếng Anh (de-en), tiếng Anh sang tiếng Séc (en-cs) và tiếng Séc sang tiếng Anh (cs-en) (Bo-jar et al., 2018). Chúng tôi kiểm tra các mô hình tự động ứng trên bộ dữ liệu Newstest 2018 tự động ứng và báo cáo kết quả điểm BLEU của SacreBLEU (Post, 2018) với cài đặt mặc định. Thiết lập của chúng tôi sau Vaswani và cộng sự. (2017) chặt chẽ và sử dụng khung Tensor2Tensor của họ (Vaswani và cộng sự, 2018). Theo dõi Vaswani et al. (2017), chúng tôi sử dụng Transformer 6 lớp với kiến trúc mã hóa-giải mã. Vì biết thêm chi tiết về thiết lập thử nghiệm của chúng tôi xin vui lòng xem Phụ lục A

Kết quả Chúng tôi báo cáo điểm BLEU của các mô hình trong Bảng 4. Chúng tôi quan sát thấy rằng vị trí chuyển động thông tin từ đầu vào đến lớp chú ý trên đầu ngữ để cải thiện điểm BLEU. Các biến thể khác nhau của sự chú ý theo vị trí trên mỗi đầu không có tác dụng nhiều sự khác biệt với DIET-REL đang cạnh tranh với Shaw và cộng sự. (2018).

4.4 Nghiên cứu cắt bỏ

Trong phần này, chúng tôi chia sẻ những phát hiện của chúng tôi về các yếu tố chính ảnh hưởng đến hiệu suất của vị trí tách rời chú ý.

Chia sẻ mã hóa vị trí trước đó công trình (Raffel và cộng sự, 2020; Ke và cộng sự, 2020; Shaw và cộng sự, 2018) đã sử dụng các phương pháp chia sẻ khác nhau để mã hóa vị trí để giảm các tham số mô hình. Chúng tôi trình bày một nghiên cứu chi tiết về các hình thức khác nhau chia sẻ mã hóa vị trí và tác dụng của nó đối với hiệu suất. Cụ thể, chúng tôi so sánh các biến thể sau trong việc chia sẻ mã hóa vị trí các tham số trên các đầu khác nhau và các lớp trong Máy biến áp.

- thông minh - Các tham số giống nhau được sử dụng cho tất cả đầu trong một lớp, với các lớp khác nhau bằng cách sử dụng các thông số khác nhau (Shaw và cộng sự, 2018; Ke và cộng sự, 2020).
- theo lớp - Chia sẻ các tham số mã hóa vị trí giữa các lớp với các tham số khác nhau cho mỗi đầu (Raffel et al., 2020).
- không - Mỗi lớp và đầu sử dụng khác nhau thông số mã hóa vị trí.

Chúng tôi trình bày kết quả so sánh các chia sẻ khác nhau các phương pháp trong Bảng 5 cho các nhiệm vụ XTREME. Chúng tôi làm những quan sát sau đây 1) chia sẻ khôn ngoan là luôn tệ hơn so với lớp khôn ngoan, 2) chia sẻ gây tổn hại

Ngữ ới mẫu	Phân đoạn chia sẻ	Phân loại XNLI	Trả lời câu hỏi			Trung bình		
			XQuAD	MLQA	TyDiQA-GoldP			
DIET-REL	-	đầu	68,0	68,1 / 52,8	57,7 / 42,7	63,3 / 50,9	57,6	
DIET-REL	theo từng	vào	67,7	66,2 / 51,0	56,0 / 41,1	60,1 / 45,9	55,4	
DIET-REL	lớp theo	đầu	đầu	68,0	68,6 / 53,3	58,1 / 43,1	61,3 / 48,2	57,2
DIET-REL	-	vào	đầu	68,4	69,4 / 54,4	58,6 / 43,5	62,4 / 49,3	58,0
DIET-REL	vào theo	đầu theo	đầu	67,8	66,0 / 50,5	55,5 / 40,4	59,2 / 44,6	54,7
DIET-REL	theo lớp theo	đầu	68,1	68,7 / 53,8	58,4 / 43,2	61,0 / 48,4	57,3	
Đầu vào DIET-ABS (dp=64)	-	68,0	67,4 / 50,5	57,8 / 42,3	61,3 / 46,8	56,3		
DIET-ABS (dp=64) trên	đầu	ngữ ới	67,9	67,5 / 52,4	57,3 / 42,3	61,6 / 46,8	56,5	
DIET-ABS (dp=128) trên	đầu	ngữ ới	68,1	68,2 / 52,0	57,9 / 42,6	61,5 / 47,6	56,8	
DIET-ABS (dp=512) trên	đầu	ngữ ới	68,5	68,0 / 52,0	57,7 / 42,4	61,6 / 48,4	56,9	
Đầu vào theo lớp DIET-ABS (dp=64)		68,0	69,3 / 53,1	59,3 / 43,9	63,2 / 48,6	57,9		
DIET-ABS (dp=64) theo	từng lớp	trên	đầu	68,4	69,3 / 53,2	59,4 / 44,1	63,3 / 48,6	58,0
DIET-ABS (dp=128) theo	từng lớp	trên	đầu	68,5	70,0 / 53,6	59,8 / 44,5	64,6 / 51,5	58,9
DIET-ABS (dp=256) theo	từng lớp	trên	đầu	68,4	69,9 / 53,8	59,6 / 44,2	62,8 / 49,1	58,3
DIET-ABS (dp=512) theo	từng lớp	trên	đầu	67,8	69,0 / 53,2	58,4 / 43,0	62,5 / 48,8	57,5

Bảng 5: Nghiên cứu cắt bỏ trên XTREME: Chúng tôi thực hiện nghiên cứu cắt bỏ sự chú ý theo vị trí tách rời để hiểu được tác dụng của 1) chia sẻ các tham số chú ý vị trí giữa các lớp và đầu 2) sự chú ý phân khúc được thêm vào ở mỗi đầu 3) hiệu suất tương đối và tuyệt đối 4) sự chú ý vị trí tuyệt đối xếp hạng dp từ 64 đến 512.

	Tiếng Anh		Đa ngôn ngữ		
	Thông số +	GLUE Thông số + XTREME			
Devlin và cộng sự. (2018)	110,1 triệu	84,8 112,9 triệu +2,5%	178,9M	-	55,3
Shaw và cộng sự. (2018)	85,2 109,9 triệu	+0,0% 85,2 109,7 triệu	181,7M	+1,7%	58,2
DIET-REL		+0,0% 85,0 128,6 triệu	178,7M	+0,0%	58,0
DIET-REL (chia sẻ)	+16,8% 85,3 111,3 triệu	+1,1% 85,2	178,5M	+0,0%	57,3
DIET-ABS (dp=128)			197,4M	+10,0%	56,8
DIET-ABS (dp=128, chia sẻ)			180,1M	+0,6%	58,9

Bảng 6: Tham số mô hình: Chúng tôi liệt kê số lượng tham số mô hình và hiệu suất cho các phương pháp mã hóa vị trí khác nhau. Chúng tôi nhận thấy rằng việc chia sẻ làm tổn hại đến hiệu suất của DIET-REL với lợi ích không đáng kể về mặt số lượng tham số. Ngược lại, tác dụng chính quy hóa của việc chia sẻ giúp DIET-ABS ổn định hơn với các thông số ít hơn để đạt được hiệu suất cạnh tranh.

hiệu suất của DIET-REL trong khi nó cải thiện Hiệu suất của DIET-ABS. Chúng tôi tóm tắt cài đặt chính cùng với số thông số mô hình trong Bảng 6. Đối với DIET-REL, việc chia sẻ mang lại ít ảnh hưởng đến việc lưu tham số và ảnh hưởng đến hiệu suất. Do đó, chúng tôi khuyên bạn không nên chia sẻ mã hóa vị trí tương đối (DIET-REL). Trên mặt khác, cần phải chia sẻ các thông số cho DIET-ABS để giữ số lượng thông số ở mức thấp. Điều thú vị là việc chia sẻ đã được chính quy hóa ảnh hưởng đến DIET-ABS, làm cho mô hình hoạt động tốt hơn. Chúng tôi chọn chia sẻ theo lớp thay vì chia sẻ theo chiều sâu để có hiệu suất tốt hơn.

Mã hóa phân đoạn Thiết kế mã hóa phân đoạn mới của chúng tôi cải thiện hơn nữa hiệu suất của mô hình được trình bày trong Bảng 5. Cả hai mô hình chú ý vị trí tách rời tương đối và tuyệt đối đều được hưởng lợi từ việc chuyển mã hóa phân đoạn từ đầu vào sang đầu ngữ ới: DIET- REL (+0,4%), theo lớp chia sẻ DIET-REL (+0,1%), DIET-ABS (+0,2%), DIET-ABS được chia sẻ theo lớp (+0,1%). Xem ứng dụng-

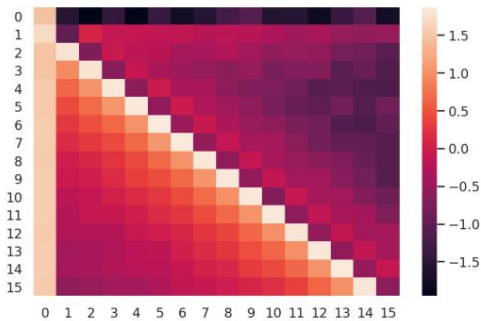
phụ lục D về kết quả của điểm chuẩn GLUE và Phụ lục C để hình dung sự chú ý của phân khúc.

Xếp hạng của sự chú ý theo vị trí tuyệt đối Thiết kế của DIET-ABS cho phép tìm hiểu các ma trận chú ý ở cấp độ cao hơn như trong Định lý 1. Để hiểu được tác động của thứ hạng chú ý theo vị trí tuyệt đối (dp) trong thực tế, chúng tôi tiến hành các thí nghiệm khác nhau thứ hạng từ dp = 64 đến dp = 512. Chúng tôi trình bày kết quả trong Bảng 5. Chúng tôi nhận thấy rằng hiệu suất được cải thiện khi chúng tôi tăng thứ hạng từ 64 lên 128. Tuy nhiên, có sự bão hòa về hiệu suất trong tiếp tục tăng nó lên 512. Chúng tôi trình bày một cách trực quan hóa thứ hạng của ma trận chú ý vị trí trong Phụ lục B.

4.5 Trực quan hóa mẫu chú ý theo vị trí

Tiếp theo chúng ta hình dung sự chú ý theo vị trí đã học được mô hình của DIET-ABS trong Hình 4. Đầu tiên chúng tôi lưu ý rằng DIET-ABS đã học được cách nắm bắt mối quan hệ họ hàng quan hệ vị trí giữa các yếu tố đầu vào. Cũng lưu ý rằng, đối với chỉ số 0 (mã thông báo [CLS]), được tách rời





Hình 4: Hình dung sự chú ý theo vị trí đã học được mô hình DIET-ABS. Lưu ý rằng ngoài việc nắm bắt các mối quan hệ vị trí tương đối, mô hình còn học cách tham dự [CLS] ở chỉ số 0, gợi ý thiết kế cờ trí [CLS] chuyên dụng trong [Ke et al. \(2020\)](#) thì không cần thiết với DIET-ABS.

Sự chú ý tuyệt đối về vị trí tương học được một khuôn mẫu đặc biệt. Mẫu này không thể chỉ được mô hình hóa bằng các phương pháp nhúng vị trí tương đối hiện có, và một số công trình hiện có ([Ke và cộng sự, 2020](#)) đã xử lý tương hợp này cụ thể bằng cách giới thiệu các thông số mới. Điều này cho thấy lợi ích của DIET-ABS trong việc không yêu cầu bất kỳ thành kiến quy nạp được thiết kế cẩn thận như trong các phương pháp tiếp cận hiện tại ([Shaw và cộng sự \(2018\)](#); [Raf-fel và cộng sự \(2020\)](#)), có thể không khái quát hóa được nhiệm vụ.

5. Kết luận

Trong bài báo này, chúng tôi đã xem xét về mặt lý thuyết và thực nghiệm giới hạn của việc nhúng vị trí phụ gia ở đầu vào và chỉ ra rằng việc nhúng vị trí trên mỗi đầu sẽ mang lại hiệu suất tốt hơn. Chúng tôi lập luận rằng hiệu suất vượt trội của một số các phương pháp mã hóa vị trí tương đối đến từ sự bổ sung trên đầu ngữ nghĩa của họ vào ma trận chú ý thay vì hơn thông tin vị trí là tương đối so với tuyệt đối. Thật vậy, chúng tôi chứng minh rằng việc sử dụng vị trí tuyệt đối mã hóa trên mỗi đầu mang lại hiệu suất tốt hơn. Được thúc đẩy bởi điều này, chúng tôi đề xuất một phương pháp chú ý theo phân khúc và vị trí trên đầu ngữ nghĩa đơn giản để đạt được hiệu suất tiên tiến trên nhiều NLP nhiệm vụ và hiệu quả tính toán hơn so với các cách tiếp cận hiện có.

Người giới thiệu

Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Philipp Koehn và Christof Monz. 2018. [Kết quả của hội nghị về dịch máy năm 2018 \(WMT18\)](#). Trong Biên bản của Hội nghị lần thứ ba về dịch máy: Tài liệu nhiệm vụ chung, trang 272-303, Bel-

gium, Brussels. Hiệp hội tính toán Lin-guistics.

Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Bài hát Xingyou, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, và Adrian Weller. 2021. [Suy nghĩ lại về sự chú ý với người biểu diễn](#).

Alexis Conneau, Ruty Rinott, Guillaume Lample, Adina Williams, Samuel R. Bowman, Holger Schwenk, và Veselin Stoyanov. 2018. Xnli: Đánh giá cách trình bày câu đa ngôn ngữ. Trong Kỳ yếu của Hội nghị năm 2018 về các phương pháp thực nghiệm trong xử lý ngôn ngữ tự nhiên. Hiệp hội ngôn ngữ học tính toán.

Zihang Dai, Guokun Lai, Yiming Yang và Quốc V. Lê. 2020. [Bộ biến đổi kênh: Lọc ra phần dư thừa tuần tự để xử lý ngôn ngữ hiệu quả](#).

Jacob Devlin, Ming-Wei Chang, Kenton Lee, và Kristina Toutanova. 2018. BERT: Đào tạo trước cho Transformers hai chiều sâu sắc để hiểu ngôn ngữ. bản in trước arXiv arXiv:1810.04805.

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat và Melvin Johnson. 2020. Xtreme: Điểm chuẩn đa tác vụ đa ngôn ngữ để đánh giá khái quát hóa đa ngôn ngữ. TRONG Kỳ yếu của Hệ thống và Học máy 2020, trang 7449-7459.

Guolin Ke, Di He và Tie-Yan Liu. 2020. [Suy nghĩ lại về mã hóa vị trí trong đào tạo trước ngôn ngữ](#). bản in trước arXiv arXiv:2006.15595.

Taku Kudo và John Richardson. 2018. [Câu văn: Một trình mã hóa và trình mã hóa từ phụ đơn giản và độc lập với ngôn ngữ để xử lý văn bản thần kinh](#).

Guillaume Lample và Alexis Conneau. 2019. [Cross-Đào tạo trước mô hình ngôn ngữ ngôn ngữ](#)

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma và Radu Soricut. 2020. [Albert: Một chuyên gia nhỏ về việc tự học cách biểu diễn ngôn ngữ](#).

Xiaodong Liu, Kevin Duh, Liyuan Liu và Jianfeng Cao. 2020. [Máy biến áp rất sâu cho dịch máy thần kinh](#).

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-dar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer và Veselin Stoyanov. 2019. RoBERTa: Đào tạo trước BERT được tối ưu hóa mạnh mẽ tiếp cận. bản in trước arXiv arXiv:1907.11692.

bài Matt. 2018. Kêu gọi sự rõ ràng trong báo cáo màu xanh điểm số. Trong WMT.

Alec Radford, Karthik Narasimhan, Tim Salimans và Ilya Sutskever. 2018. Nâng cao khả năng hiểu ngôn ngữ bằng cách đào tạo trước mang tính khái quát. Báo cáo kỹ thuật, OpenAI.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li và Peter J Liu. 2020. Khám phá những hạn chế của việc học chuyển tiếp bằng công cụ chuyển đổi văn bản thành văn bản thống nhất. Tạp chí Nghiên cứu Học máy, 21(140):1-67.

Peter Shaw, Jakob Uszkoreit và Ashish Vaswani. 2018. Tự chú ý với các đại diện vị trí tự động đổi. Trong Kỷ yếu Hội nghị năm 2018 của Chi hội Bắc Mỹ của Hiệp hội Ngôn ngữ học Tính toán: Công nghệ Ngôn ngữ Con người, Tập 2 (Bài viết ngắn), trang 464-468.

Ashish Vaswani, Samy Bengio, Eugene Brevdo, Fran-cois Chollet, Aidan N Gomez, Stephan Gouws, Llion Jones, Łukasz Kaiser, Nal Kalchbrenner, Niki Parmar, và những người khác. 2018. Tensor2tensor cho dịch máy thần kinh. bản in trước arXiv arXiv:1803.07416.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser và Illia Polosukhin. Năm 2017. Sự chú ý là tất cả những gì bạn cần. Trong Những tiến bộ trong hệ thống xử lý thông tin thần kinh, trang 5998-6008.

Alex Wang, Amanpreet Singh, Julian Michael, Fe-lix Hill, Omer Levy và Samuel Bowman. 2019. Keo: Nền tảng phân tích và điểm chuẩn đa tác vụ để hiểu ngôn ngữ tự nhiên. Trong Hội nghị quốc tế lần thứ 7 về đại diện học tập, ICLR 2019.

Sinong Wang, Belinda Z. Li, Madian Khabsa, Han Fang và Hao Ma. 2020. [Linformer: Tự chú ý với độ phức tạp tuyến tính](#).

Yu-An Wang và Yun-Nung Chen. 2020. Việc nhúng vị trí học được gì? một nghiên cứu thực nghiệm về mã hóa vị trí mô hình ngôn ngữ được đào tạo trước. Trong EMNLP 2020.

Adina Williams, Nikita Nangia và Samuel R. Bowman. 2018. [Kho ngữ liệu thử thách có phạm vi rộng để hiểu câu thông qua suy luận](#).

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quốc V. Lê, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin John-son, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes và Jeffrey Dean. 2016. Hệ thống dịch [máy thần kinh của Google](#) : Thu hẹp khoảng cách giữa dịch người và dịch máy.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov và Quốc V. Lê. 2019. XLNet: Huấn luyện trước tự hồi quy tổng quát để hiểu ngôn ngữ. bản in trước arXiv arXiv:1906.08237.

Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J Reddi và Sanjiv Kumar. 2020. Transformers có phải là công cụ xấp xỉ phổ quát của các hàm tuần tự không? Trong Hội nghị quốc tế về đại diện học tập.

Một thiết lập thử nghiệm

Trong phần này, chúng tôi trình bày thêm chi tiết về thiết lập thử nghiệm của chúng tôi.

Đào tạo trước Chúng tôi đào tạo trước các mô hình bằng cách sử dụng tác vụ LM bị che (Devlin và cộng sự, 2018) và không sử dụng tính năng mất Dự đoán câu tiếp theo (NSP) như được đề xuất trong RoBERTa (Liu và cộng sự, 2019). Mỗi đầu vào được xây dựng với các câu đầy đủ từ các tài liệu và được đóng gói theo độ dài chuỗi tối đa. Chúng tôi sử dụng kiến trúc tự động như BERTBASE (Devlin và cộng sự, 2018) (L = 12, H = 768, A = 12) cho các thử nghiệm của mình.

Tinh chỉnh Một số tác vụ tiếp theo có các nhóm câu đầy đủ khác nhau được cung cấp ở đầu vào. Đối với những tác vụ đó (ví dụ: MNLI, CoLA, XNLI, SQuAQ), chúng tôi tinh chỉnh các mô hình bằng mã hóa phân đoạn bổ sung được thảo luận trong Phần 5.3. Chúng tôi giữ nguyên các mô hình cho các nhiệm vụ khác dưới dạng thư từ trước khi đào tạo.

Siêu tham số Siêu tham số chúng tôi sử dụng được trình bày trong Bảng 7.

	Tiếng Anh		Đa ngôn ngữ	
	Huấn luyện trước	Finetune luyện tập trước		Tinh chỉnh
Số bước tối đa	500K	5 hoặc 10 k	125K	3 k
Tỷ lệ học	nguyên 0,0018 {1e-5, 2e-5, 3e-5, 4e-5}		nguyên 0,0018 {1e-5, 2e-5, 3e-5, 4e-5}	
Tỷ lệ khởi động	0,025 0,1 128 128 4096 32 20k 3,5k		0,025 0,1 512 512 4096 32 20k 3,5k	
Độ dài chuỗi				
Kích thước lô				
Khoảng thời gian kiểm tra				

Bảng 7: Siêu tham số cho tất cả các model

Dịch Đối với các thử nghiệm Dịch của chúng tôi, chúng tôi tuân theo cách thiết lập của Vaswani et al. (2017) và sử dụng khung Tensor2Tensor của họ (Vaswani và cộng sự, 2018). Chúng tôi đào tạo bằng cách sử dụng các bộ dữ liệu WMT18 ((Europarl v7, Common Crawl corpus và News Commentary v13) en-de, de-en, en-cs và cs-en. Chúng tôi báo cáo điểm BLUE do SacreBLEU (Post, 2018) cung cấp trên tập dữ liệu mới nhất 2018 Chúng tôi đào tạo mô hình Transformer 6 lớp. Mọi thay đổi đối với mã hóa vị trí đều được áp dụng cho tất cả các lớp chú ý cả trong bộ mã hóa và bộ giải mã. Chúng tôi sử dụng trình tối ưu hóa Adam và đào tạo 250 nghìn bước. Để giải mã, chúng tôi sử dụng tìm kiếm chùm tia với kích thước chùm tia 10 và phạt chiều dài 0,6.

B Bằng chứng

Chúng minh Định lý 1. Dễ dàng nhận thấy khẳng định đầu tiên bằng cách nhận thấy rằng hạng tích của hai ma trận bị giới hạn trên bởi mức tối thiểu của các hạng riêng lẻ.

$$\text{hạng}(Aa) = \text{hạng}((X + P)WQ^T K(X + P)^T) \leq \min(\text{thứ hạng}(X + P), \text{hạng}(WQ), \text{hạng}(X + P), \text{hạng}(WK)) \leq dh.$$

$$\text{xếp hạng}((X + P)WQ^T K(X + P)^T) \leq dh, \text{ trong đó } WQ, WK \in \mathbb{R}^{d \times dh}$$

Bất đẳng thức cuối cùng suy ra từ  $\text{hạng}(WQ) \leq dh$  là  $WQ \in \mathbb{R}^{d \times dh}$ . Để chứng minh khẳng định thứ hai, chúng ta áp dụng phương pháp xây dựng. Trước tiên chúng ta lấy  $WQ = WK$  bằng nhau ma trận với các hàng  $dh$  đầu tiên là ma trận nhận dạng và các hàng  $d - dh$  còn lại đều là số 0. Sau đó

$$WQ^T K = \begin{bmatrix} I_{dh, dh} & 0_{dh, d-dh} \\ 0_{d-dh, dh} & 0_{d-dh, d-dh} \end{bmatrix}.$$

Ở đây  $I_{dh, dh}$  biểu thị ma trận đồng nhất trong  $\mathbb{R}^{dh \times dh}$  và  $0_{dh, d-dh}$  biểu thị ma trận toàn số 0 trong  $\mathbb{R}^{dh, d}$ . Chúng ta đặt  $X$  sao cho  $d$  hàng đầu tiên tạo thành ma trận nhận dạng và phần còn lại là số 0 -  $X = [I_d, d, 0_n \ d, d]$ . Do đó  $XWQ^T KX$  trở thành ma trận đường chéo tự động với

$$XWQ^T KX = \begin{bmatrix} I_{dh, dh} & 0_{dh, n-dh} \\ 0_n \ dh, dh & 0_n \ dh, n-dh \end{bmatrix}.$$

Chọn  $d_p = n > d_h$  và đặt  $P^* = I$ . Bây giờ chọn  $P^*$  với các số 0 ở các cột  $n - d_p$  đầu tiên và đẳng thức ở các cột  $d_p$  cuối cùng ( $P^* = [0_{d,n-d_p}, I_{d,d_p}]$ ) sẽ cho kết quả

$$P^* P^* = \begin{bmatrix} 0_{n-d_p, n-d_p} & 0_{n-d_p, d_p} \\ 0_{d_p, n-d_p} & I_{d_p, d_p} \end{bmatrix}.$$

Kết hợp cả hai điều này mang lại cho chúng ta

$$\text{hạng}(A_r) = \text{hạng}(XWQW^T KX + P^* P^*) = \min(d_h + d_p, n) > d_h.$$

□

Đặt  $X \in \mathbb{R}^{n \times d}$  là phần nhúng từ đầu vào trong chiều  $d$  với độ dài chuỗi  $n$ . Chúng tôi có thể huấn luyện vị trí  $P \in \mathbb{R}^{n \times d}$ , luyện được thêm vào chuỗi đầu vào trước khi đưa vào mô hình  $g$ . Đối với đầu vào  $X$  nhất định và nhãn  $y$ , mục tiêu của hàm mất mát như sau:

$$L = (g(X + P), y) \tag{5}$$

Định lý 2. Cho  $X$  và  $P$  là các ma trận nhúng có thể huấn luyện được trong  $\mathbb{R}^{n \times d}$ . Khi đó, độ dốc của hàm mất mát trong phương trình (5), tại bất kỳ điểm nào  $(X, y)$  và đối với mọi hàm khả vi và  $g$ , đều giống nhau đối với  $X$  và  $P$ .

Nhận xét. Định lý này cho chúng ta thấy rằng độ dốc giống nhau đối với các nhúng mã thông báo đầu vào và các nhúng vị trí. Mặc dù trong các tác vụ NLP tiêu chuẩn, đầu vào  $X$  có thể khác nhau ở mỗi bước do các mã thông báo đầu vào khác nhau xuất hiện trong mỗi lô nhỏ, kết quả vẫn cho thấy rằng việc nhúng vị trí phụ gia có thể hạn chế mô hình tìm hiểu tầm quan trọng tương đối của mã hóa vị trí đối với mã thông báo nhúng dựa trên nhiệm vụ đào tạo hiện tại.

Chứng minh Định lý 2. Định lý trên được thực hiện bằng cách chỉ tính các gradient và cho thấy chúng bằng nhau ở mỗi bước.

Độ dốc của mục tiêu trên  $X$  và  $P$  như sau.

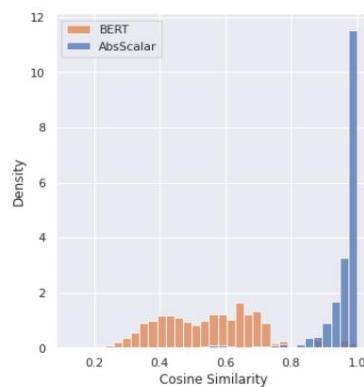
$$\begin{aligned} \frac{\partial L}{\partial X} &= \frac{\partial L}{\partial X+P} \cdot \frac{\partial (X+P)}{\partial X} = \frac{\partial L}{\partial X+P} \cdot I \\ &= \frac{\partial L}{\partial X+P} \cdot \frac{\partial (X+P)}{\partial P} = \frac{\partial L}{\partial X+P} \cdot I \end{aligned}$$

Việc tính toán độ dốc ở trên tuân theo quy tắc dây chuyền. Điều này cho thấy độ dốc của  $L$  wrt  $X$  và  $P$  là như nhau. □

## C Trực quan hóa sự chú ý

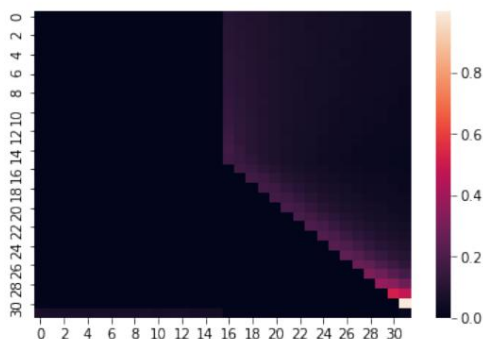
Trong phần này, chúng tôi kiểm tra nội bộ của mô hình để hiểu cách thức hoạt động của mô hình được đề xuất. Trước tiên, chúng tôi hình dung nội dung bên trong mô hình của các lựa chọn thay thế mô hình khác nhau để tranh luận rằng mô hình đề xuất của chúng tôi là hợp lý.

Tại sao chúng tôi loại bỏ phần nhúng đầu vào để hiểu liệu có hợp lý hay không khi loại bỏ phần nhúng phụ gia đầu vào sau khi thêm vô hướng vị trí trên mỗi đầu, chúng tôi thêm phần nhúng vị trí phụ gia vào mô hình DIET-ABS của chúng tôi. Sau đó, chúng tôi kiểm tra việc nhúng vị trí của mô hình BERT và biến thể DIET-ABS của chúng tôi bằng cách nhúng vị trí phụ gia. Hình 5 cho thấy rằng, khi mô hình có cả nhúng vị trí tuyệt đối vô hướng tuyệt đối và bổ sung, việc nhúng vị trí hầu như không mã hóa thông tin - tất cả các nhúng vị trí ở đầu vào đều tương tự nhau.

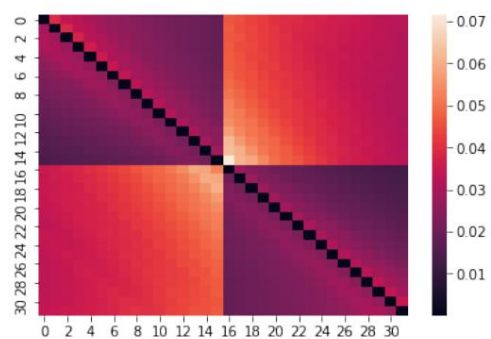


Hình 5: Phân bố độ tương tự cosin giữa tất cả các cặp vị trí tuyệt đối của nhúng vị trí phụ gia đầu vào cho mô hình BERT cơ sở và DIET-ABS được đề xuất. Chúng tôi quan sát thấy rằng, sau khi các tính năng vị trí được thêm vào mỗi đầu như trong DIET-ABS, việc nhúng vị trí đầu vào hầu như không chứa thông tin - tất cả các cặp vị trí đầu vào đều giống nhau.

Tác động của việc chú ý đến phân khúc Chúng tôi cũng xem xét tác động của việc thêm sự chú ý vào phân khúc lên trên sự chú ý về vị trí. Hình 6 cho thấy một số mẫu tiêu biểu. Chúng tôi nhận thấy rằng sự chú ý đến phân khúc cho phép mô hình chú ý nhiều hơn đến các phần của chuỗi thuộc các phân khúc nhất định.



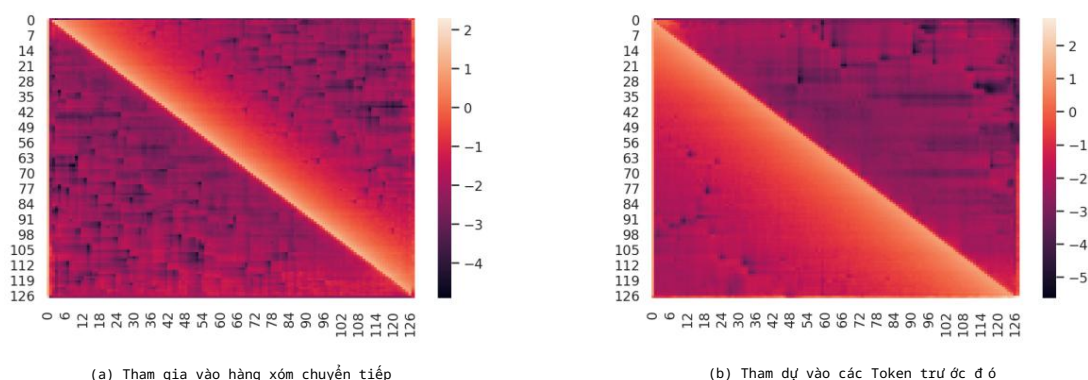
(a) Tham dự Phần thứ hai



(b) Chú ý đến vị trí tương đối của trọng lượng thấp hơn

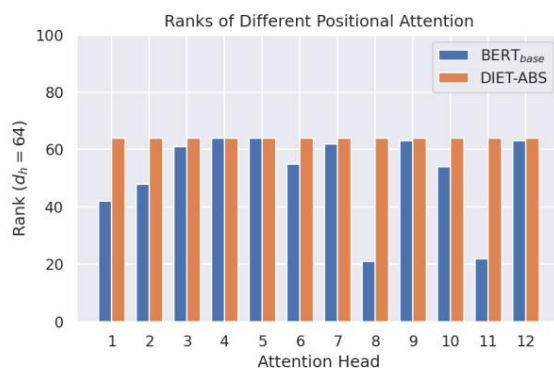
Hình 6: Chúng ta xem xét đầu vào có độ dài 32 với hai đoạn. Phân đoạn thứ hai bắt đầu ở chỉ số 16. Chúng tôi quan sát các mẫu chú ý trong mô hình DIET-REL mà không chú ý đến mã thông báo.

Mô hình dịch chuyển được học từ sự chú ý vị trí tuyệt đối Sử dụng mã hóa vị trí tương đối nhìn chung mang lại kết quả tốt hơn mặc dù quy mô cải tiến nhỏ hơn so với mã hóa tính năng di chuyển trên mỗi đầu ngữ ý. Để hiểu điều này, chúng tôi hình dung mẫu chú ý của chú ý vị trí tuyệt đối và tìm thấy hai mẫu đại diện trong DIET-ABS trong Hình 7. Chúng tôi nhận thấy rằng ngay cả khi đưa ra các đặc điểm vị trí tuyệt đối, phần lớn mô hình đều học được “mô hình dịch chuyển”. Khác với Wang và Chen (2020) khẳng định vị trí tuyệt đối chỉ tìm hiểu các mô hình địa phương, chúng tôi cho thấy sự chú ý của vị trí thực sự có thể chú ý đến bối cảnh dài hơn. Tuy nhiên, mô hình dịch chuyển có thể được mô hình hóa trực tiếp bằng vị trí tương đối. Do đó, DIET-REL có thể là lựa chọn mô hình tốt hơn với ít thông số hơn và độ lệch quy nạp chính xác hơn trong một số ứng dụng.



Hình 7: Mẫu điểm chú ý vị trí được lấy mẫu cho mô hình DIET-ABS. Chúng ta có thể thấy mô hình dịch chuyển rõ ràng do mô hình tạo ra. Các mẫu như vậy có thể được mô hình hóa tốt hơn bằng các bộ mã hóa vô hướng vị trí tương đối.

Thứ hạng của Ma trận Chú ý Vị trí Trong Hình 8, chúng tôi trình bày so sánh thứ hạng của ma trận chú ý vị trí cho mô hình BERTBASE với số lần nhúng vị trí tuyệt đối ở đầu vào (PQWQW KP K) so với số lần nhúng vị trí tuyệt đối trên mỗi đầu (DIET-ABS (1), (PQP K), trong đó PQ, PK  $R \times dp$ ). Với việc nhúng vị trí bổ sung ở đầu vào, ma trận chú ý vị trí có thứ hạng thấp hơn nhiều, hạn chế sức mạnh đại diện. Điều này được giảm bớt nhờ DIET-ABS.



Hình 8: Thứ hạng của ma trận chú ý vị trí



D Nghiên cứu cắt bỏ bổ sung trên GLUE

Trước đó, chúng tôi trình bày một nghiên cứu cắt bỏ trên XTREME trong Bảng 5 đối với các biến thể chú ý theo vị trí được tách rời. Chúng tôi so sánh DIET-REL và DIET-ABS với đường cơ sở (Devlin và cộng sự, 2018). Bây giờ chúng tôi trình bày một nghiên cứu tự động về tiêu chuẩn GLUE trong Bảng 8 và quan sát các kết quả tự động.

Mã hóa vị trí Trong Bảng 8, việc di chuyển các phần nhúng vị trí từ đầu vào sang mỗi đầu sẽ cải thiện điểm trung bình cho cả DIET-REL (+0,1%) và DIET-ABS (+0,2%).

Mã hóa phân đoạn Trong Bảng 8, việc di chuyển các phần nhúng phân đoạn từ đầu vào sang mỗi đầu sẽ cải thiện cả DIET-REL (+0,3%) và DIET-ABS (+0,05%).

Chiến lược chia sẻ Chia sẻ đóng một vai trò quan trọng đối với DIET-ABS. Trong Bảng 9, chúng tôi thấy rằng việc chia sẻ sẽ làm giảm hiệu suất của DIET-REL (-0,2% theo lớp, -0,3% theo đầu). Đối với DIET-ABS, chia sẻ giúp mô hình ổn định hơn, có khả năng cạnh tranh với DIET-REL.

Người mẫu	Phân đoạn vị trí	MNLI QQP QNLI SST2 CoLa STS-B 393k 8.5k							Trung bình
		364k	105k	67k				7k	
Devlin và cộng sự. (2018)	đầu vào 85,8 / 85,9 91,1	DIET-REL mỗi đầu 86,0 / 86,1 91,0	DIET-REL	89,9	93,2	58,7	89,0	84,8	
	đầu mỗi đầu 86,3 / 86,3 91,0	DIET-ABS (dp=64) mỗi đầu 86,1 / 85,8 91,2	DIET-ABS	89,8	92,8	59,6	89,0	84,9	
	( dp=64) bình quân đầu người bình quân đầu người 86,1 / 86,1 91,2	DIET-ABS (dp=64, phần		90,5	92,9	60,3	89,3	85,2	
	chia sẻ) bình quân đầu người bình quân đầu người 86 / 86,8 91,1	DIET-ABS (dp=128, chia sẻ)		90,0	93,0	58,9	89,9	85,0	
	bình quân đầu người bình quân đầu người 86,4 / 86,4 90,8			90,2	93,0	58,9	89,8	85,0	
				90,4	92,9	59,3	89,8	85,2	
				89,5	93,0	59,8	90,2	85,2	

Bảng 8: Nghiên cứu cắt bỏ vị trí và phân đoạn trên GLUE: DIET-REL và DIET-ABS chứng minh những ưu điểm của di chuyển cả nhúng vị trí và phân đoạn từ đầu vào sang đầu người.

Người mẫu	Chia sẻ	MNLI QQP QNLI SST2 CoLa STS-B 393k 8.5k							Trung bình
		364k	105k	67k				7k	
DIET-REL 86.3 / 86.3 91.0	DIET-REL theo lớp 86,5 / 86.3 91.1	DIET-REL	90,5	92,9	60,3	89,3	85,2		
	theo phần đầu 85,8 / 85.7 91.2	DIET-ABS (dp=64) 86.1 / 86.1 91.2	DIET-ABS	90,0	93,0	58,8	89,6	85,0	
	ABS (dp=128) 86.7 / 6,5 91,2	DIET-ABS (dp=64) theo lớp 86 / 86,8 91,1		90,2	92,8	59,8	89,1	84,9	
				90,2	93,0	58,9	89,8	85,0	
	DIET-ABS (dp=128) theo lớp 86,4 / 86,4 90,8			90,6	92,8	60,1	89,4	85,3	
				90,4	92,9	59,3	89,8	85,2	
				89,5	93,0	59,8	90,2	85,2	

Bảng 9: Chia sẻ nghiên cứu cắt bỏ trên GLUE: Chúng tôi thực hiện nghiên cứu cắt bỏ để hiểu tác động của việc chia sẻ từ thể các tham số mã hóa trên các lớp và đầu. Chúng tôi nhận thấy rằng việc chia sẻ sẽ cải thiện hiệu suất của DIET-ABS, như ng làm tổn hại đến hiệu suất của DIET-REL với việc chia sẻ theo từng lớp hoặc theo từng phần.