

# Hate Speech Detection based on Sentiment Knowledge Sharing

## Nhóm 4

Mạnh Hùng   Công Thiện   Minh Trí   Bảo Lợi   Đình Quân

Khoa Khoa học máy tính  
*Trường đại học Công Nghệ Thông Tin*

Ngày 3 tháng 1 năm 2024

# Nội dung

- 1 Giới thiệu bài toán
- 2 Phương pháp
- 3 Thực nghiệm
- 4 Bàn luận & Kết luận
- 5 Tài liệu tham khảo

- 1 Giới thiệu bài toán
- 2 Phương pháp
- 3 Thực nghiệm
- 4 Bàn luận & Kết luận
- 5 Tài liệu tham khảo

# Giới thiệu bài báo

**Bài báo:** <https://aclanthology.org/2021.acl-long.556>

**Nội dung:** Hate Speech Detection based on Sentiment Knowledge Sharing

**Các tác giả:** Xianbing Zhou, Yong Yang, Xiaochao Fan, Ge Ren, Yunfeng Song, Yufeng Diao, Liang Yang, Hongfei Lin

**Code:** <https://github.com/1783696285/sks>

**Trích dẫn:** Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)

**Ngày công bố:** Tháng 8 năm 2021

## Đặt vấn đề

**Ví dụ:** Go fucking kill yourself.

- Sự phổ biến của Internet và mạng xã hội  
→ Điều kiện để hate speech lan truyền một cách rộng rãi
- Hậu quả: Bị phân biệt đối xử, ảnh hưởng tâm lý, bạo lực.. và đây là vấn đề nghiêm trọng của xã hội

⇒ **Vấn đề:** Làm thế nào để có thể phát hiện hate speech tự động, nhanh chóng và chính xác, đồng thời can thiệp sớm để ngăn chặn chúng?

- ① Cấu trúc của ngôn ngữ tự nhiên phức tạp
- ② Các giải pháp trước chủ yếu dựa vào **quy tắc** hoặc **trích xuất đặc trưng thủ công**
- ③ **Phương pháp máy học**: sử dụng các đặc trưng nhân tạo ở mức độ nông (shallow features)
- ④ **Phương pháp học sâu**: chưa khai thác đầy đủ đặc trưng sentiment và nguồn kiến thức sentiment ở bên ngoài

## ❶ Sử dụng thông tin sentiment

- Tích hợp các từ mang tính xúc phạm có trong câu vào mạng neural network
- Sử dụng multi-task learning để mô hình học và chia sẻ kiến thức sentiment

## ❷ Đề xuất một framework mới

- Sử dụng nhiều đơn vị trích xuất đặc trưng
- Áp dụng cơ chế Gated Attention để kết hợp các đặc trưng

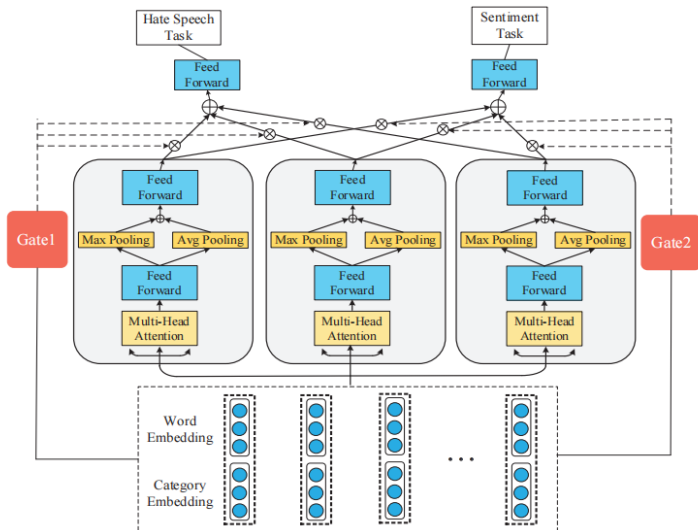
## ❸ Thực nghiệm trên 2 tập dữ liệu: **SemEval-2019 task-5, Davidson**

# Nội dung

- 1 Giới thiệu bài toán
- 2 Phương pháp**
- 3 Thực nghiệm
- 4 Bàn luận & Kết luận
- 5 Tài liệu tham khảo

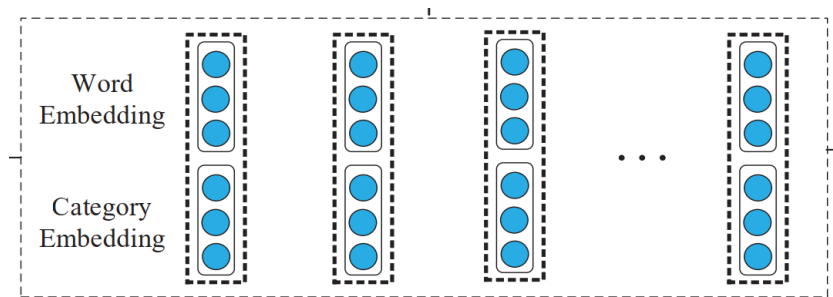


# Tổng quan



Hình 1: Kiến trúc tổng quan của phương pháp **Hate Speech detection based on Sentiment Knowledge Sharing (SKS)**

# Input layer



Hình 2: Input layer

## Ý tưởng chung

- Pretrained word embedding
- Category embedding
- Kết hợp (word embedding, category embedding)

## Word Embedding

- Dựa trên giả định về phân phối
- Ánh xạ các từ vào không gian đặc trưng nhiều chiều
- Duy trì thông tin về ngữ nghĩa

## Quan sát của tác giả

- Hate speech chứa những từ ngữ xúc phạm

**Ví dụ:** Go fucking kill yourself and die already ugly pile of shit scumbag

- Việc xác định những từ ngữ xúc phạm góp phần lớn vào việc xác định hate speech  
→ Tạo một từ điển chứa các từ ngữ xúc phạm

## Category Embedding

### **Bộ từ vựng của từ điển**

- Nguồn: Wikipedia, trang Noswearing
- 5 loại chính: Hate speech, Disability, LGBT, Ethnic, Religious

### **Chức năng**

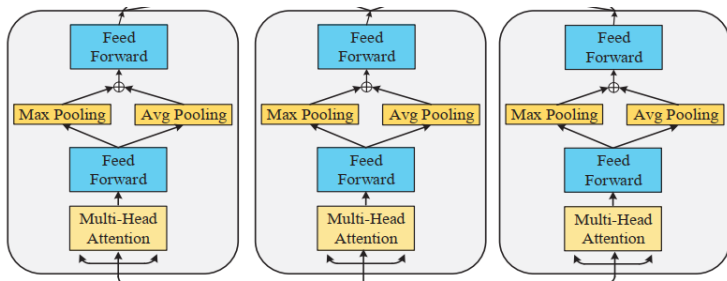
- Dùng để đánh giá mỗi từ trong câu có phải từ căm thù hay không

## Bối cảnh

- Ý nghĩa xúc phạm ẩn trong ngữ cảnh văn hóa, không còn thể hiện ở mỗi từ ngữ
  - **non-hate:** i'm so fucking ready!
- **Hate speech detection:** Thiếu dữ liệu huấn luyện chất lượng cao
- **Sentiment analysis:** Có nhiều dữ liệu chất lượng được gán nhãn

⇒ Mối tương quan cao giữa 2 task, áp dụng Multi-task learning để chia sẻ kiến thức sentiment

# Sentiment Knowledge Sharing layer



Hình 3: Sentiment Knowledge Sharing layer

## Ý tưởng chung

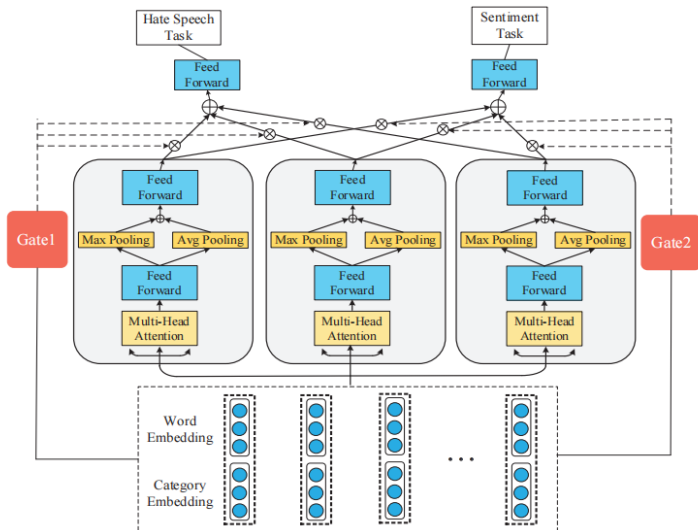
- Sử dụng cấu trúc framework Mix-of-Expert (MoE): gồm nhiều đơn vị trích xuất đặc trưng giống nhau
- Các đơn vị trích xuất đặc trưng: bao gồm một lớp multi-head attention, hai lớp pooling và hai lớp feed forward

## Multi-head Attention Layer

- Cơ chế self-attention: tính toán mức độ tương quan về ngữ nghĩa giữa các từ trong câu để xác định mối liên hệ giữa chúng
- Cơ chế multi-head attention:
  - Sử dụng nhiều head self-attention chồng lên nhau để phân tích mối quan hệ giữa các phần trong câu
  - Mỗi head tập trung vào một khía cạnh khác của câu input
  - Kết hợp kết quả của các head lại để tạo ra vector biểu diễn feature cuối cùng



# Gated Attention



Hình 4: Kiến trúc tổng quan của phương pháp **Hate Speech detection based on Sentiment Knowledge Sharing (SKS)**

# Gated Attention

## Ý tưởng

- Gated attention chọn một tập các feature extraction unit để sử dụng, dựa trên nội dung và ngữ cảnh của input đó
- Tác vụ khác nhau  $\rightarrow$  có Gate khác nhau
- Cấu trúc của Gate unit  $\sim$  feature extraction unit

## Kết quả

- Output Gate thứ  $k$ : là một vector thể hiện xác suất của các feature extraction unit (FEU) được chọn
- Đối với mỗi Gate:  
 $\sum(\text{output mỗi gate} * \text{FEU}) \rightarrow$  vector biểu diễn cuối cùng của input

# Nội dung

- 1 Giới thiệu bài toán
- 2 Phương pháp
- 3 Thực nghiệm**
- 4 Bàn luận & Kết luận
- 5 Tài liệu tham khảo

Dataset	total	Classes
<b>SemEval2019 task5(SE)</b>	11.971	hate (5.035)
		non-hate (6.936)
<b>Davidson(DV)</b>	24.783	hate (1.430)
		non-hate (23.353)
<b>Sentiment Analysis(SA)</b>	31.962	negative(2.242)
		positive(29,720)

**Bảng 1:** Số liệu thống kê về các dataset sử dụng trong thực nghiệm

## So sánh với các phương pháp baseline

Model	DV		SE	
	Acc	F1(wei)	Acc	F1(macro)
SVM*	-	<u>87.0</u>	<u>49.2</u>	<u>45.1</u>
LSTM*	94.5	93.7	<u>55.0</u>	<u>53.0</u>
GRU*	94.5	93.9	<u>54.0</u>	<u>52.0</u>
CNN-GRU*	-	<u>94.0</u>	62.0	61.5
BiLSTM*	94.4	93.7	<u>53.5</u>	<u>51.9</u>
BiGRU_Stacked*	-	-	<u>56.0</u>	<u>54.6</u>
USE_SVM*	-	-	<u>65.3</u>	<u>65.1</u>
BERT*	94.8	95.8	-	<u>48.8</u>
GPT*	-	-	-	<u>51.5</u>
SKS	<b>95.1</b>	<b>96.3</b>	<b>65.9</b>	<b>65.2</b>

Bảng 2: So sánh với các phương pháp hiện có

# Thực nghiệm của bài báo

<div>Metric Model</div>	DV		SE	
	Acc	F1(wei)	Acc	F1(macro)
-sc	94.0	94.0	59.6	59.3
-s	94.5	94.3	61.3	61.3
no-gate	94.8	95.9	64.7	64.3
<b>SKS</b>	<b>95.1</b>	<b>96.3</b>	<b>65.9</b>	<b>65.2</b>

Bảng 3: Bảng kết quả các thực nghiệm của các model

## Trong đó:

- “-sc”: không sử dụng sentiment feature và category embedding
- “-s”: không sử dụng sentiment feature
- “no-gate”: không sử dụng Gated attention layer

# Thực nghiệm của nhóm

	<b>Paper</b>	<b>Ours</b>
<b>embedding data</b>	glove.840B.300d	glove.6B.300d
<b>learning rate</b>	0.001	0.002

**Bảng 4:** Các tham số thay đổi khi training

	<b>DV</b>		<b>SE</b>	
<b>metrics</b>	Acc	F1(wei)	Acc	F1(macro)
<b>SKS</b>	95.1	96.3	65.9	65.2
<b>Ours</b>	93.0	93.0	61.0	61.0

**Bảng 5:** Kết quả thực nghiệm trên 2 tập datasets

# Nội dung

- 1 Giới thiệu bài toán
- 2 Phương pháp
- 3 Thực nghiệm
- 4 Bàn luận & Kết luận**
- 5 Tài liệu tham khảo



## Ưu điểm

- Đạt hiệu quả vượt trội so với các phương pháp baseline

## Nhược điểm

- Khó khăn trong việc thu thập và lựa chọn các nguồn dữ liệu phù hợp
- Tồn kém tài nguyên tính toán do phức tạp của mô hình

- ① Sử dụng nhiều **feature extraction unit** để các task chia sẻ tham số, tận dụng được kiến thức liên quan
- ② Áp dụng cơ chế gated attention để kết hợp các đặc trưng một cách linh hoạt
- ③ Mô hình đề xuất có thể sử dụng đầy đủ thông tin sentiment của câu input và nguồn dữ liệu sentiment bên ngoài

- Áp dụng cho nhiều ngôn ngữ khác ngoài tiếng Anh
- Tập trung vào lựa chọn và mở rộng các nguồn dữ liệu huấn luyện, bao gồm các dạng hate speech khác nhau, các loại dữ liệu sentiment phù hợp với quy mô lớn hơn

# Nội dung

- 1 Giới thiệu bài toán
- 2 Phương pháp
- 3 Thực nghiệm
- 4 Bàn luận & Kết luận
- 5 Tài liệu tham khảo**

# Tài liệu tham khảo chính



<https://aclanthology.org/2021.acl-long.556>

*Cảm ơn các bạn và cô đã lắng nghe!*