



Partner Matters! An Empirical Study on Fusing Personas for Personalized Response Selection in Retrieval-Based Chatbots

Jia-Chen Gu¹, Hui Liu², Zhen-Hua Ling¹, Quan Liu^{1,3}, Zhigang Chen³, Xiaodan Zhu²

¹National Engineering Laboratory for Speech and Language Information Processing,
University of Science and Technology of China

²ECE & Ingenuity Labs, Queen's University

³State Key Laboratory of Cognitive Intelligence, iFLYTEK Research

Outline

- **Introduction**
- Persona Fusion for Response Selection
- Experiments
- Conclusion

Introduction

- Personalized Response Selection

A task that aims to **select** an appropriate response from a set of candidates given **conversation contexts** and **personas of speakers**, is an important technique to **present personalities of dialogue agents** in **retrieval-based** chatbots.

Introduction

- An example dialogue

Persona 1		Persona 2	
Original	I just bought a brand new house. I like to dance at the club. I run a dog obedience school. I have a big sweet tooth. I like taking and posting selkies.	Original	I love to meet new people. I have a turtle named timothy. My favorite sport is ultimate frisbee. My parents are living in bora bora. Autumn is my favorite season.
Revised	I have purchased a home. Just go dancing at the nightclub, it is fun! I really enjoy animals. I enjoy chocolate. I pose for pictures and put them online.	Revised	I like getting friends. Reptiles make good pets. I love to run around and get out my energy. My family lives on a island. I love watching the leaves change colors.

Dialogue

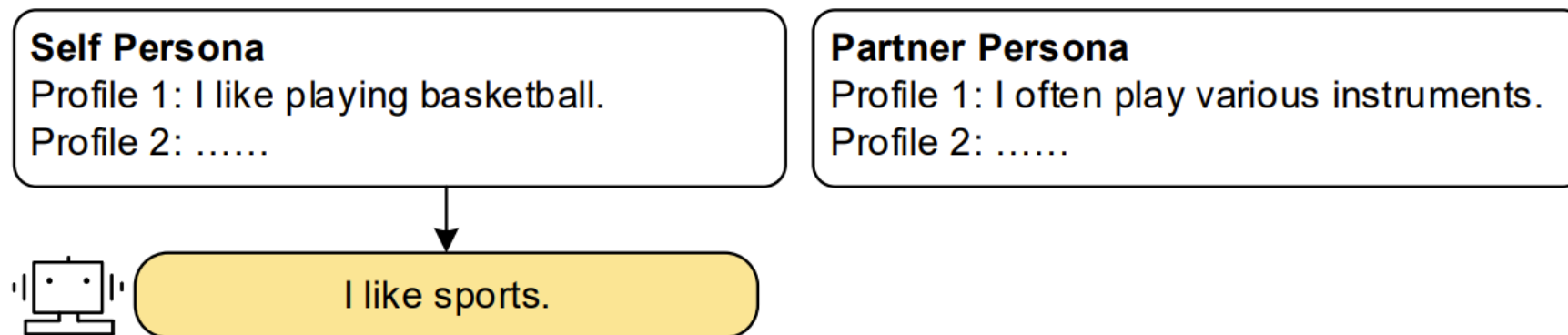
Person 1: Hello, how are you doing tonight?
Person 2: I am well an loving this interaction how are you?
Person 1: I am great. I just got back from the club.
Person 2: This is my favorite time of the year season wise.
Person 1: I would rather eat chocolate cake during this season.
Person 2: What club did you go to? Me an timothy watched tv.
Person 1: I went to club chino. What show are you watching?
Person 2: LOL oh okay kind of random.
Person 1: I love those shows. I am really craving cake.
Person 2: Why does that matter any? I went outdoors to play frisbee.
Person 1: It matters because I have a sweet tooth.
Person 2: So? LOL I want to meet my family at home in bora.
Person 1: My family lives in alaska. It is freezing down there.
Person 2: I bet it is oh I could not.

Motivation

- Most of previous studies **focused on the self speaker's persona** in dialogue who was about to utter a response, while **the contribution of the partner speaker's persona to dialogue was rarely noticed**.
- For a conversation conditioned on personas, if a dialogue agent has **no access to the partner persona**, it often **over-focuses on retrieving responses related to the agent itself**, which sometimes **deviates from the ground truth** of how a conversation really goes.

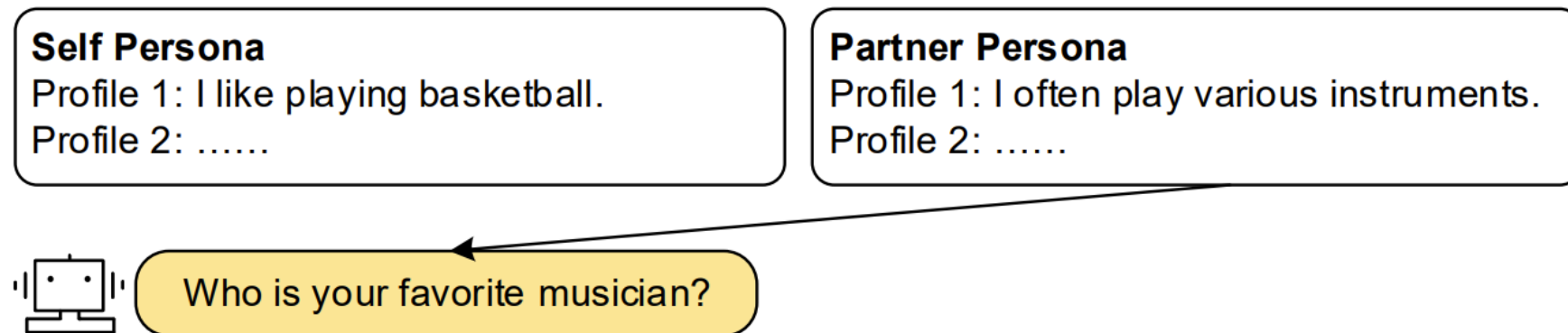
Motivation

- A conversation about hobbies. If the agent only has access to the **self persona profile**, it often **over-weights** response candidates **related to the agent itself**.



Motivation

- A conversation about hobbies. If the agent also has access to the **partner persona profile**, it gives models **more flexibility** to not only focus on continuously talking about the agent itself, but also **conducting more collaborative communication**, e.g., asking questions as the real conversations often proceed.



Preliminary Experiments

- In order to **compare the ability of different personas** for selecting an appropriate response directly, the **context information was ablated** and an appropriate response was selected **given only the self or partner persona** information.
- Matching models: HRE, IMN and BERT.

Preliminary Experiments

- **Single persona-response matching can achieve a comparable performance**, showing the usefulness of utilizing persona information to select an appropriate response.
- Although the **partner persona is less important than the self persona**, it can **still contribute to response selection** to some extent, which is consistent with our assumption.

Model	Persona	hits@1	MRR
HRE	Self	23.9	40.1
	Partner	8.7	23.4
IMN	Self	48.8	60.7
	Partner	19.3	34.2
BERT	Self	50.6	62.5
	Partner	20.6	35.6

Outline

- Introduction
- **Persona Fusion for Response Selection**
- Experiments
- Conclusion

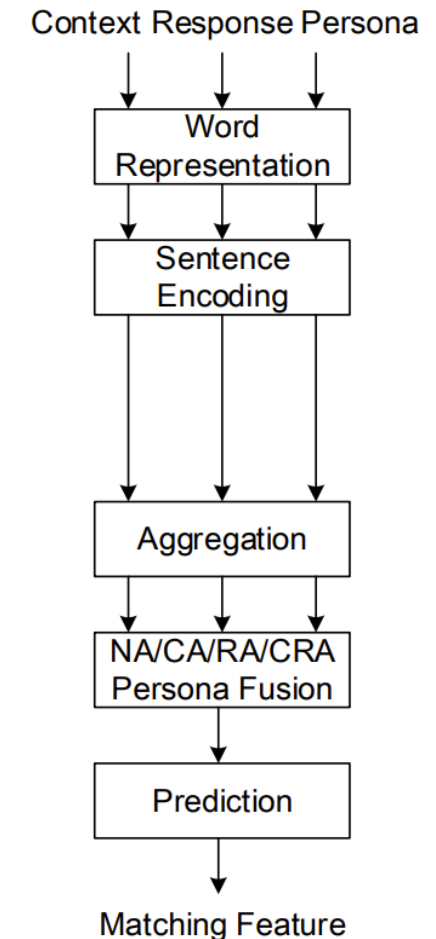
Persona Fusion for Response Selection

- **Four persona fusion strategies**, i.e., none-aware (NA), context-aware (CA), response-aware (RA) and context-response-aware (CRA) ones, are designed based on **whether or not considering the interactions between personas and contexts** as well as **the interactions between personas and responses**.
- For a thorough comparison and analysis, these four strategies are **implemented into three representative models** for response selection, which are based on the **HRE, IMN and BERT** models respectively.

Sentence-Encoding-Based Model

- Hierarchical Recurrent Encoder.

Given **the sequence of profile embeddings**, the **aggregated persona embedding** is obtained by **persona fusion**. In this paper, four persona fusion strategies are designed.



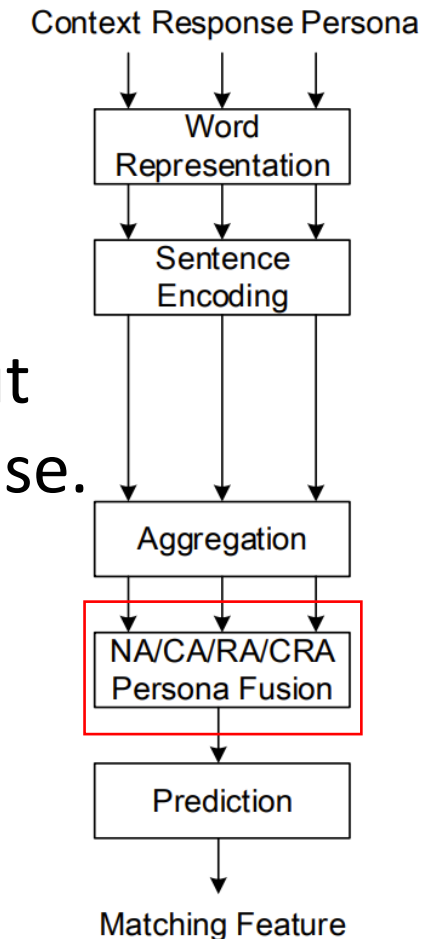
Sentence-Encoding-Based Model

- Hierarchical Recurrent Encoder.

None-Aware Persona Fusion: the persona fusion is **independent of both contexts and responses**.

A self-attention-based aggregation is designed without being aware of any information of context and response.

$$\alpha_n = \mathbf{w}^\top \cdot \bar{\mathbf{p}}_n^{agr} + b,$$
$$\hat{\mathbf{p}}^{agr} = \sum_{n=1}^{n_p} \frac{e^{\alpha_n}}{\sum_{k=1}^{n_p} e^{\alpha_k}} \bar{\mathbf{p}}_n^{agr},$$

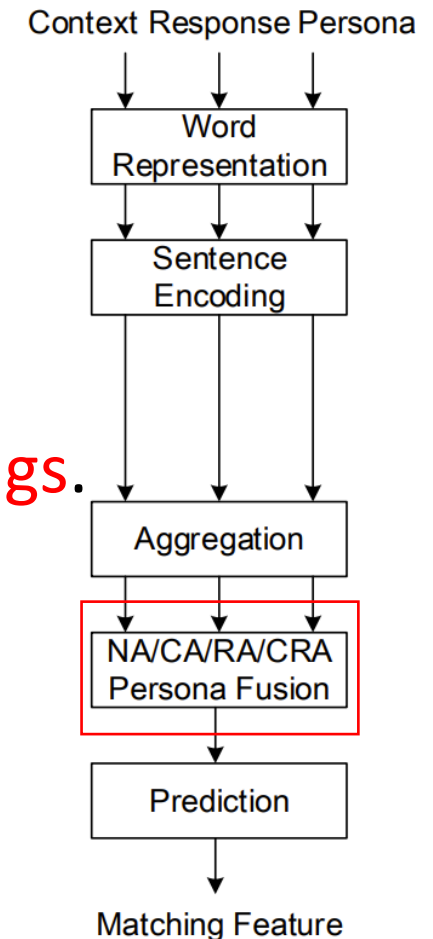


Sentence-Encoding-Based Model

- Hierarchical Recurrent Encoder.

Context-Aware Persona Fusion: compute **similarities between the context embedding and each profile embedding**, and perform the attention operation by **attaching different importance to profile embeddings**.

$$\alpha_n = \hat{\mathbf{c}}^{agr\top} \cdot \bar{\mathbf{p}}_n^{agr},$$
$$\hat{\mathbf{p}}^{agr} = \sum_{n=1}^{n_p} \frac{e^{\alpha_n}}{\sum_{k=1}^{n_p} e^{\alpha_k}} \bar{\mathbf{p}}_n^{agr},$$

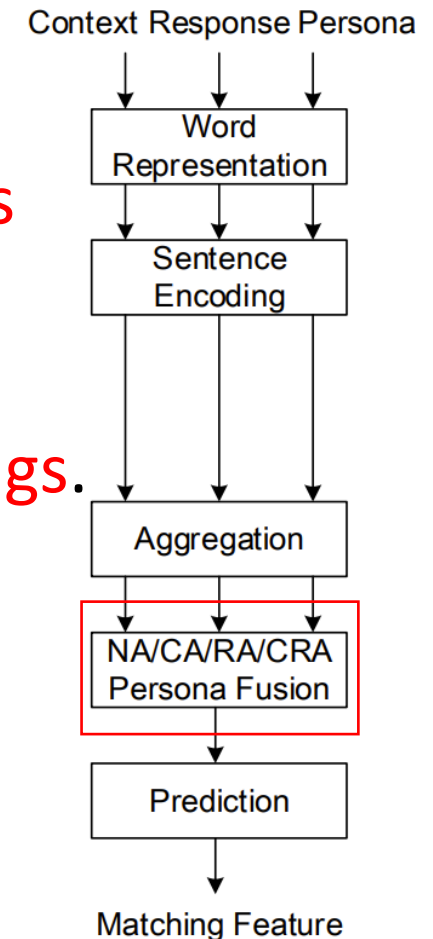


Sentence-Encoding-Based Model

- Hierarchical Recurrent Encoder.

Response-Aware Persona Fusion: compute **similarities between the response embedding and each profile embedding**, and perform the attention operation by **attaching different importance to profile embeddings**.

$$\alpha_n = \bar{\mathbf{r}}^{agr\top} \cdot \bar{\mathbf{p}}_n^{agr},$$
$$\hat{\mathbf{p}}^{agr} = \sum_{n=1}^{n_p} \frac{e^{\alpha_n}}{\sum_{k=1}^{n_p} e^{\alpha_k}} \bar{\mathbf{p}}_n^{agr},$$

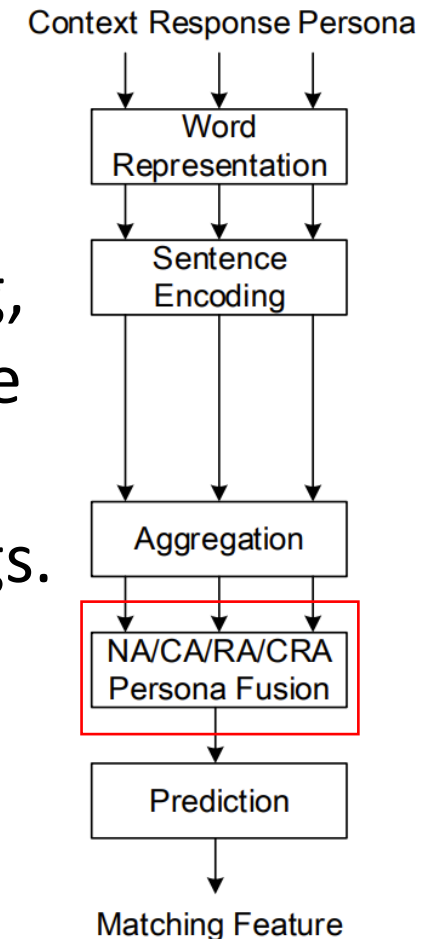


Sentence-Encoding-Based Model

- Hierarchical Recurrent Encoder.

Context-Response-Aware Persona Fusion: first **concatenate** the context and the response embedding, and then **transform** it to the same dimension of profile embeddings with a linear transformation. Similarities are computed between it and each profile embeddings.

$$\alpha_n = (\mathbf{w}^\top \cdot [\hat{\mathbf{c}}^{agr}; \bar{\mathbf{r}}^{agr}] + b)^\top \cdot \bar{\mathbf{p}}_n^{agr},$$
$$\hat{\mathbf{p}}^{agr} = \sum_{n=1}^{n_p} \frac{e^{\alpha_n}}{\sum_{k=1}^{n_p} e^{\alpha_k}} \bar{\mathbf{p}}_n^{agr}.$$



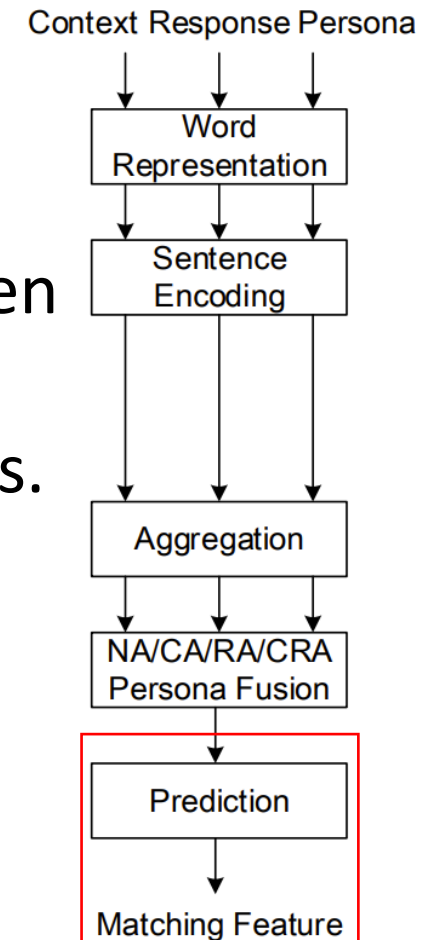
Sentence-Encoding-Based Model

- Hierarchical Recurrent Encoder.

The final matching feature vector is the **concatenation of context, persona and response embeddings** and then sent into a **multi-layer perceptron** classifier.

Models are learnt by minimizing the cross-entropy loss.

$$\mathcal{L}(\mathcal{D}, \Theta) = - \sum_{(c,p,r,y) \in \mathcal{D}} y \log(g(c, p, r)).$$

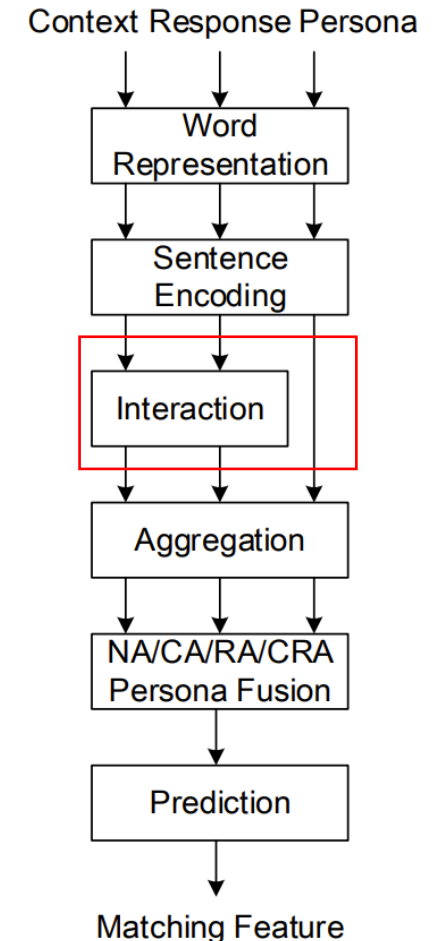


Cross-Attention-Based Model

- Interactive Matching Network.

The reason to choose this model is that it shares the most **similar architecture with HRE** except the **interaction module**, so that we can **explore the effect of interactions between contexts and responses on persona fusion**.

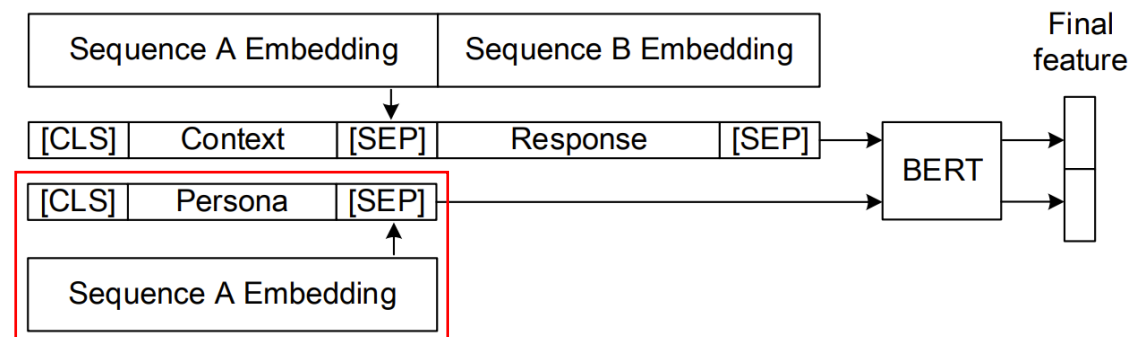
The persona aggregation in IMN is **identical** to that in HRE.



Pretraining-Based Model

- BERT

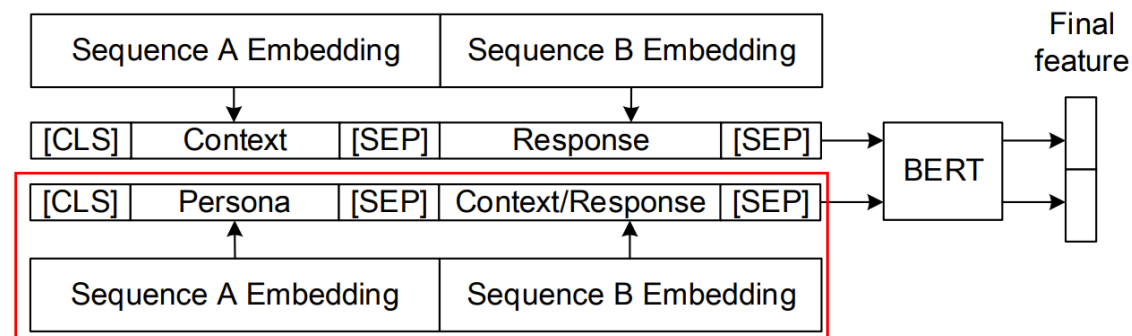
None-Aware Persona Fusion: A **dual matching** architecture is composed of **two encoding pipelines**. One is used to derive the **matching feature between contexts and responses**, and the other is used to derive the **persona fusion feature alone** without any interactions with contexts or responses.



Pretraining-Based Model

- BERT

Context/Response-Aware Persona Fusion: **Two encoding pipelines.** One is used to derive the **matching feature between contexts and responses**, and the other is used to derive the **persona fusion feature concatenated with contexts/responses** for encoding for interacting with contexts/responses.



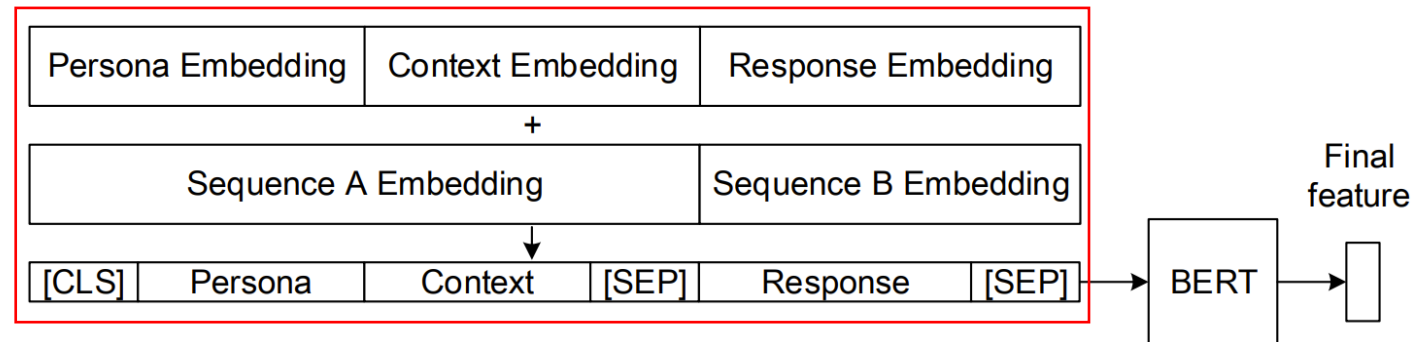
Pretraining-Based Model

- BERT

Context-Response-Aware Persona Fusion: **One encoding pipeline.**

Personas are concatenated with both contexts and responses for **interacting with them at the same time.**

Three **subtypes of embedding** are added to distinguish personas, contexts and responses.



Outline

- Introduction
- Persona Fusion for Response Selection
- **Experiments**
- Conclusion

Experiments

- Datasets

Datasets	# Candidates	Settings	Train	Valid	Test
Persona-Chat	20	Original/Revised	65719	7801	7512

- Metrics

The **recall** of true positive replies by **selecting k best-matched response** from available candidates for the given context and knowledge, denoted as **$R@k$** .

Experiments

- Overall Performance

Table 3: Evaluation results of our reimplemented HRE, IMN and BERT models together with previous methods on the Persona-Chat dataset without using any personas.

	hits@1	MRR
IR baseline [26]	21.4	-
Starspace [26]	31.8	-
Profile [26]	31.8	-
KV Profile [26]	34.9	-
HRE [21]	42.7	60.0
IMN [4]	63.8	75.8
BERT[3]	70.7	80.8

(1) Self personas are more important under all strategies and models. It is reasonable since **self personas provides fundamental descriptions** of the speaker who is about to utter a response.

Table 4: Performance of four persona fusion strategies implemented into three models on the Persona-Chat dataset under the original persona configuration. Numbers marked with \star denote that the gains or losses after adding persona conditions are statistically significant (t-test with p -value < 0.05) comparing with the corresponding baseline models in Table 3. Numbers in bold denote the persona fusion strategy that achieves the best performance.

	Self Persona		Partner Persona	
	hits@1	MRR	hits@1	MRR
HRE-NA	47.4 \star	63.7 \star	42.2 \star	59.3 \star
HRE-CA	47.0 \star	63.7 \star	42.1 \star	59.3 \star
HRE-RA	58.1\star	71.8\star	42.8	60.0
HRE-CRA	43.3 \star	60.4	42.1 \star	59.1 \star
IMN-NA	64.4 \star	76.3 \star	64.1	76.1
IMN-CA	64.6 \star	76.5 \star	63.9	76.1
IMN-RA	66.3\star	77.7\star	64.3\star	76.2\star
IMN-CRA	64.1	76.2	64.1	76.1
BERT-NA	71.1	80.9	70.9	80.8
BERT-CA	71.2 \star	81.0	70.9	80.9
BERT-RA	82.6 \star	89.0 \star	71.1 \star	80.9
BERT-CRA	84.3\star	90.3\star	71.2\star	80.9

Experiments

- Overall Performance

Table 3: Evaluation results of our reimplemented HRE, IMN and BERT models together with previous methods on the Persona-Chat dataset without using any personas.

	hits@1	MRR
IR baseline [26]	21.4	-
Starspace [26]	31.8	-
Profile [26]	31.8	-
KV Profile [26]	34.9	-
HRE [21]	42.7	60.0
IMN [4]	63.8	75.8
BERT[3]	70.7	80.8

(2) The **partner persona** was shown to **contribute to the performance as well.**

Particularly in the cross-attention-based and pretraining-based models.

Table 4: Performance of four persona fusion strategies implemented into three models on the Persona-Chat dataset under the original persona configuration. Numbers marked with \star denote that the gains or losses after adding persona conditions are statistically significant (t-test with p -value < 0.05) comparing with the corresponding baseline models in Table 3. Numbers in bold denote the persona fusion strategy that achieves the best performance.

	Self Persona		Partner Persona	
	hits@1	MRR	hits@1	MRR
HRE-NA	47.4 \star	63.7 \star	42.2 \star	59.3 \star
HRE-CA	47.0 \star	63.7 \star	42.1 \star	59.3 \star
HRE-RA	58.1\star	71.8\star	42.8	60.0
HRE-CRA	43.3 \star	60.4	42.1 \star	59.1 \star
IMN-NA	64.4 \star	76.3 \star	64.1	76.1
IMN-CA	64.6 \star	76.5 \star	63.9	76.1
IMN-RA	66.3\star	77.7\star	64.3\star	76.2\star
IMN-CRA	64.1	76.2	64.1	76.1
BERT-NA	71.1	80.9	70.9	80.8
BERT-CA	71.2 \star	81.0	70.9	80.9
BERT-RA	82.6 \star	89.0 \star	71.1 \star	80.9
BERT-CRA	84.3\star	90.3\star	71.2\star	80.9

Experiments

- Overall Performance

Table 3: Evaluation results of our reimplemented HRE, IMN and BERT models together with previous methods on the Persona-Chat dataset without using any personas.

	hits@1	MRR
IR baseline [26]	21.4	-
Starspace [26]	31.8	-
Profile [26]	31.8	-
KV Profile [26]	34.9	-
HRE [21]	42.7	60.0
IMN [4]	63.8	75.8
BERT[3]	70.7	80.8

(3) **RA** persona fusion strategy performs best in the **sentence-encoding-based** and **cross-attention-based** models and **CRA** fusion strategy performs best in the **pretraining-based model**.

Table 4: Performance of four persona fusion strategies implemented into three models on the Persona-Chat dataset under the original persona configuration. Numbers marked with \star denote that the gains or losses after adding persona conditions are statistically significant (t-test with p -value < 0.05) comparing with the corresponding baseline models in Table 3. Numbers in bold denote the persona fusion strategy that achieves the best performance.

	Self Persona		Partner Persona	
	hits@1	MRR	hits@1	MRR
HRE-NA	47.4 \star	63.7 \star	42.2 \star	59.3 \star
HRE-CA	47.0 \star	63.7 \star	42.1 \star	59.3 \star
HRE-RA	58.1\star	71.8\star	42.8	60.0
HRE-CRA	43.3 \star	60.4	42.1 \star	59.1 \star
IMN-NA	64.4 \star	76.3 \star	64.1	76.1
IMN-CA	64.6 \star	76.5 \star	63.9	76.1
IMN-RA	66.3\star	77.7\star	64.3\star	76.2\star
IMN-CRA	64.1	76.2	64.1	76.1
BERT-NA	71.1	80.9	70.9	80.8
BERT-CA	71.2 \star	81.0	70.9	80.9
BERT-RA	82.6 \star	89.0 \star	71.1 \star	80.9
BERT-CRA	84.3\star	90.3\star	71.2\star	80.9

Experiments

- Overall Performance

	Self Persona				Partner Persona			
	Original		Revised		Original		Revised	
	hits@1	MRR	hits@1	MRR	hits@1	MRR	hits@1	MRR
IR baseline [35]	41.0	-	20.7	-	18.1	-	18.1	-
Starspace [35]	48.1	-	32.2	-	24.5	-	26.1	-
Profile [35]	47.3	-	35.4	-	28.3	-	29.4	-
KV Profile [35]	51.1	-	35.1	-	29.1	-	28.9	-
FT-PC [18]	-	-	60.7	-	-	-	-	-
DGMN [37]	67.6	-	58.8	-	-	-	-	-
DIM [9]	78.8	86.7	70.7	81.2	64.0	76.1	63.9	76.0
TransferTransfo [32]	80.7	-	-	-	-	-	-	-
P ² Bot [15]	81.9	-	68.6	-	-	-	-	-
FIRE [8]	81.6	-	74.8	-	-	-	-	-
BERT-RA	82.6 [*]	89.0 [*]	77.1 [*]	85.4 [*]	71.1 [*]	80.9 [*]	70.8 [*]	80.8 [*]
BERT-CRA	84.3[*]	90.3[*]	79.4[*]	86.9[*]	71.2[*]	80.9[*]	71.8[*]	81.5[*]
BERT-CRA - subtype	83.6 [*]	89.9 [*]	78.4 [*]	86.4 [*]	70.8 [*]	80.8 [*]	70.9 [*]	80.8 [*]

BERT-CRA achieves new state-of-the-art performance of response selection.

Experiments

- Retrieval Time

Table 6: The efficiency (cases/second) of four persona fusion strategies implemented into three models by recording the inference time over the whole validation set on the Persona-Chat dataset under the original persona configuration.

Efficiency (cases/second)					
HRE-NA	4660.1	IMN-NA	1661.4	BERT-NA	53.29
HRE-CA	4596.9	IMN-CA	1666.3	BERT-CA	53.29
HRE-RA	4626.9	IMN-RA	1674.6	BERT-RA	53.29
HRE-CRA	4643.5	IMN-CRA	1688.0	BERT-CRA	92.67

Although BERT-based models take more time, it is acceptable compared to the performance they achieved.

Experiments

- Discussion on Response Generation

Although fusing personas for dialogue generation is not the focus of this paper, we conducted a preliminary experiment to show that **self or partner personas also contribute differently to response generation.**

Response candidate is not available during inference in response generation. Thus we explored only the context-aware persona fusion strategy with a lightweight model MiniLM.

Experiments

- Discussion on Response Generation

Preliminary results can **verify our assumption** to some extent.

Table 7: Performance of response generation conditioned on the original persona. Numbers marked with \star denote that the gains or losses after adding persona conditions are statistically significant (t-test with p -value < 0.05). Numbers in bold denote the persona fusion strategy that achieves the best performance.

Model-Persona	Relevance	Diversity		Length
	BLEU	DIST-1	DIST-2	
MiniLM	3.54	94.47	99.49	8.73
MiniLM-CA-Self	4.50\star	94.94\star	99.59	8.59
MiniLM-CA-Partner	3.65 \star	93.80 \star	99.48	9.05 \star

Table 8: An example from the generated responses that demonstrates the different contributions of the self and partner personas. Given the conversation context, \times denotes an inappropriate response and \checkmark denotes an appropriate one.

Self Persona	Partner Persona
I like to go hunting.	I live in Ohio.
I am a handyman.	I like to go hiking.
I am allergic to shellfish.	I am a single mom of two boys.
I restore classic cars.	I work as an accountant.

Context:
how are you tonight , i just got back from hiking .

Response with self persona:
i am good . i just got back from hunting .

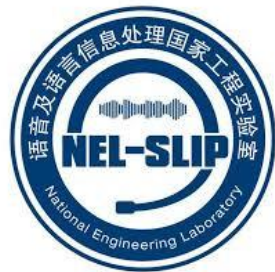
Response with partner persona:
that sounds fun . i just got back from a hike. \times
Hiking in Ohio must be very interesting. \checkmark

Outline

- Introduction
- Persona Fusion for Response Selection
- Experiments
- **Conclusion**

Conclusion

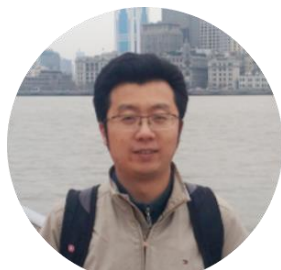
- Four persona fusion strategies are designed and implemented into three representative models to explore the impact of self and partner personas on personalized response selection in retrieval-based chatbots.
- Empirical studies show that the partner persona neglected in previous studies can still improve the performance under certain conditions.
- Our proposed models achieves a new state-of-the-art performance of response selection on the Persona-Chat dataset.



Jia-Chen Gu



Hui Liu



Zhen-Hua Ling



Quan Liu



Zhigang Chen



Xiaodan Zhu

Thanks!

Code: <https://github.com/JasonForJoy/Personalized-Response-Selection>