

HADE: Hierarchical Affective Dialog Encoder for Personality Recognition in Conversation

Anonymous ACL submission

Abstract

Personality recognition in conversation aims to determine the personality traits of speakers through the dialogue content, which is of great importance in designing personalized conversational AI. Existing methods that use only linguistic patterns in utterances limit their performance. To fill in the gap, we investigate the effectiveness of incorporating affective information and modeling the interactions among speakers in conversations for personality recognition. However, available corpus with personality and explicit affective annotations is rare. Besides, modeling the dialog flow with multiple speakers is difficult. Faced with the issues, we proposed Hierarchical Affective Dialog Encoder (HADE) for effective personality recognition in conversation. HADE utilizes manual annotated Valance-Arousal-Dominance (VAD) vectors of single words and implicitly extracts affective information from utterances. Then, it introduces a hierarchical architecture with the dialog state embeddings to identify the speakers and encode the whole dialog flow. Finally, the affective information is integrated by an auxiliary VAD regression task to enhance personality recognition. Extensive experiments on a well-known dataset, **FriendsPersona**, demonstrate the effectiveness of our method compared with state-of-the-art models. Besides, we conduct an ablation study to discuss different approaches for integrating affective information and dialog flow modeling; the design of both parts in HADE is also verified to be effective for personality recognition in conversation¹.

1 Introduction

Personality is relatively permanent traits and unique characteristics that give both consistency and individuality to a person’s behavior (Feist and Feist, 2012). In the conversation scenario, personality influences both semantic content and emotional expressions. Therefore, recognizing the personality

¹Our code will be released at github.com.

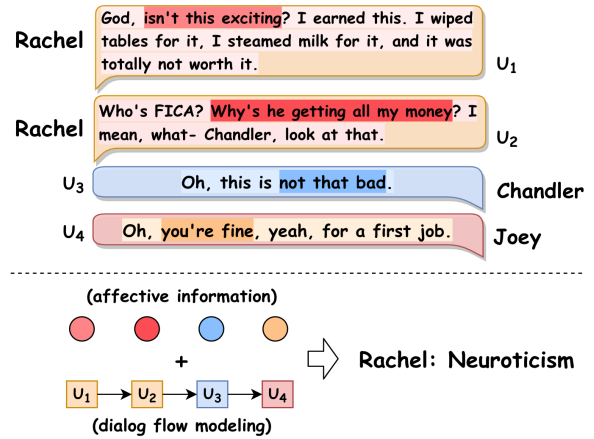


Figure 1: A toy example for personality recognition in conversation. In this example, we first analyze the affective information in utterances from Rachel: excited and nervous, while for Chandler and Joey, the affective information is quite positive. Besides, the dialog flow contains the interaction between Rachel (U_1, U_2), Chandler (U_3), and Joey (U_4), showing that others are comforting her. So, we infer that the current personality exhibited by Rachel is *Neuroticism*.

of speakers is critical for understanding the conversation content so that the dialog systems are able to provide appropriate and personalized responses to users.

Existing researches (Rissola et al., 2019; Jiang et al., 2020) simply focused on extracting linguistic patterns in utterance to recognize certain personality, which limited their performance. The main reason is that they fail to model complicated yet effective factors (e.g., the affective information in utterances or the interactions among the speakers) of personality recognition in conversation intentionally in their approaches.

Psychology studies (Watson and Clark, 1992; Mehrabian, 1995) find that there is a strong correlation between personalities and affective information in expression. Besides, by observing the conversation data, we found that in addition to the

060	semantics in utterances, the interactions among	is also verified to be effective in HADE. The con-	110
061	different speakers in the dialog flow also helps to	tributions of this work are summarized as follows:	111
062	recognize the personality.		
063	However, implementing the insights above meets	• We investigate the effectiveness of incorporat-	112
064	two major challenges. The first one is the lack	ing affective information and modeling the	113
065	of explicit affective annotations in the personality	interactions among speakers and proposed	114
066	analysis corpus. Personality analysis datasets in the	HADE for personality recognition in conver-	115
067	conversation scenario are already rare because col-	sation.	116
068	lecting such data may cause privacy concerns. Nev-		
069	ertheless, almost none of them incorporates explicit	• In HADE, we utilize pre-annotated VAD vec-	117
070	affective annotations. The second challenge arises	tors of single words and introduce a hierarchi-	118
071	in modeling the dialog flow to analyze the specified	cal architecture with the dialog state embed-	119
072	speakers. The data shortage tends us to use general	dings, which solves the challenges of affective	120
073	pre-train language models. However, it is difficult	annotation shortage and the dialog flow mod-	121
074	to indicate specific speakers efficiently with exist-	eling.	122
075	ing conversational models (<i>e.g.</i> , DialoGPT (Zhang		
076	et al., 2019), PLATO (Bao et al., 2019), and EVA	• HADE outperforms state-of-the-art methods	123
077	(Zhou et al., 2021)).	on a public conversation dataset, FriendsPer-	124
078	To tackle the issues mentioned above, we pro-	sona . Besides, through ablation study, the	125
079	pose the Hierarchical Affective Dialog Encoder	modules in HADE are validated effective to	126
080	(HADE) to implicitly extract the affective informa-	integrate affective information and model the	127
081	tion from the dialog content and design a hierar-	dialog flow.	128
082	chical architecture to encode the dialog flow for		
083	personality recognition. First, to alleviate the lack	2 Related Works	129
084	of explicit affective annotations in the personal-		
085	ity analysis corpus, HADE uses the pre-annotated	In this section, we review existing researches	130
086	Valence-Arousal-Dominance (VAD) vectors for	that related to personality analysis in conversation.	131
087	single words to represent the implicit affective	These researches are categorized into two aspects:	132
088	factors in utterances. Then, we design HADE	Text-based Personality Analysis and Dialog Mod-	133
089	based on BERT (Devlin et al., 2018), and a trans-	elling in Conversation.	134
090	former (Vaswani et al., 2017) encoder hierarchi-		
091	cally. BERT in the bottom layer encodes all the	2.1 Text-based Personality Analysis	135
092	utterances, respectively. After that, the transformer		
093	encoder receives the output from the bottom layer	Most existing researches in text-based personality	136
094	and the dialog state embeddings designed to iden-	recognition are limited to analyzing self-reported	137
095	tify different speakers for personality recognition.	essays (Pennebaker and King, 1999; Tighe et al.,	138
096	To incorporate the affective information to enhance	2016) or behaviors on social media (Golbeck et al.,	139
097	personality recognition, we integrate an auxiliary	2011; Schwartz et al., 2013). (Schwartz et al.,	140
098	VAD regression task in the upper layer of HADE	2013) analyzed 700 million words, phrases, and	141
099	through a regression head of BERT.	topic instances collected from the Facebook mes-	142
100	To show the effectiveness of our method, we	sages of 75,000 volunteers and found striking vari-	143
101	conduct extensive experiments on FriendsPersona	ations in language with personality, gender, and	144
102	constructed by (Jiang et al., 2020). It is the dialog	age. The Facebook data is also studied in (Lynn	145
103	script with personality annotations in 711 differ-	et al., 2020). They hierarchically encode all posts	146
104	ent dialogues, including 8,157 utterances from the	from one user with attention-based GRU (Cho et al.,	147
105	famous TV Series <i>Friends</i> ² . Adequate results vali-	2014) to produce the whole contextual representa-	148
106	date that our model outperforms the state-of-the-art	tions for personality identification.(Moreno et al.,	149
107	methods. We also design an ablation study to evalu-	2019) adopted a feature-engineering approach to	150
108	ate different modules in our model. The utilization	extract text-based features from Twitter blogs to	151
109	of affective information in personality recognition	identify the personality of Twitter users. Only a	152
		few works (Mehl et al., 2006; Rissola et al., 2019;	153
		Jiang et al., 2020) focus on the conversation sce-	154
		nario due to the shortage of available data: (1) The	155
		number of conversation datasets with personality	156

²<https://www.imdb.com/title/tt0108778/>

labels is insufficient as collecting such kinds of data may cause privacy concerns; and (2) The length of the dialog flow is short compared with self-reports, essays, and multiple posts on social media.

2.2 Dialog Flow Modeling in Conversation

Modeling the dialog flow is also helps to understand the personalities of speakers in conversation. In the early stage, (Serban et al., 2017) regards the tokens in utterances and utterances in a dialog flow as two kinds of sequences and proposes the classic hierarchical RNN encoder for dialog data. (Mehri et al., 2019) proposes two novel pre-training objectives: masked-utterance retrieval and inconsistency identification to better capture both the utterance-level and context-level information. Similarly, (Gu et al., 2020) employs a hierarchical BERT architecture to encode the utterances and the dialog context separately to enable the model to capture multi-level coherences. Furthermore, (Wolf et al., 2019b) adds the dialog state embeddings during utterance encoding so that the model can identify the utterances from different speakers.

3 Preliminaries

Before introducing our method, we first present the development of the Big-five personality traits and the affective information for personality analysis. This part inspires the design of HADE and helps to understand our method as preliminary knowledge.

3.1 The Big-five Personality Traits

The Big-five trait theory presents a discrete taxonomy of personality as shown in Table 1³, which is naturally suitable for personality analysis as a classification problem. This theory was defined by several independent sets of researchers who used factor analysis of verbal descriptors of human behavior. It is developed from the trait theory and the lexical hypothesis and in psychology.

Factor	Description
Openness	Openminded, imaginative, and sensitive.
Conscientiousness	Scrupulous, well-organized.
Extraversion	The tendency to experience positive emotions.
Agreeableness	Trusting, sympathetic, and cooperative.
Neuroticism	The tendency to experience psychological distress.

Table 1: The OCEAN personality traits and description (Costa and McCrae, 1992)

³https://en.wikipedia.org/wiki/Big_Five_personality_traits

In the trait theory, personality is the set of psychological traits and mechanisms within the individual that are organized and relatively enduring and that influence their interactions with, and adaptations to, the intrapsychic, physical, and social environments (Larsen and Buss, 2008). The lexical hypothesis first states that those personality characteristics that are important to a group of people will eventually become a part of that group’s language (Cattell, 1943). It second states that more important personality characteristics are more likely to be encoded into language as a single word (John et al., 1988), which also explains the principles of existing personality analysis researches based on linguistic patterns.

Therefore, the big-five personality traits are widely applied as personality recognition classification labels in social medias (Iacobelli et al., 2011; Souri et al., 2018) and conversations (Mairesse and Walker, 2006; Mairesse et al., 2007).

3.2 Affective Information for Personality Analysis

Besides linguistic patterns, affective information in expressions is important for personality analysis. Affect, in psychology, refers to the underlying experience of feeling, emotion, or mood (Fiske and Taylor, 1991). Affective states vary along three principal dimensions: valence, arousal, and motivational intensity (Harmon-Jones et al., 2013) (also interpreted as dominance in some works (Bradley and Lang, 1999; Mohammad, 2018)).

(Watson and Clark, 1992) pointed out that there are strong relations between the *Extraversion* and *Conscientiousness* traits and the positive affects, and between *Neuroticism* and *disagreeableness* and various negative affects. (Mehrabian, 1995) analyzed the relationship between the big-five personality with the PAD⁴ scales as follows: *Extraversion* includes pleasant and dominant characteristics; *Agreeableness* consists of pleasant and submissive qualities; *Conscientiousness* relates positively to trait pleasure; *Neuroticism* includes pleasant and arousable qualities; and *Openness* is comprised of pleasant, arousable, and dominant characteristics. Based on the analysis above, (Mehrabian, 1996) further estimates the relationship into a set of re-

⁴It is Pleasure-Arousability-Dominance (PAD) in the original paper, PAD and VAD share the same meaning in the context of verbal text, we will use VAD for consistency henceforth.

gression equations. These theories are also adopt to design human-like robots (Han et al., 2012; Masuyama et al., 2018), and empathetic dialog systems (Ball and Breese, 2000; Wen et al., 2021).

The following section will introduce the studied problem and the HADE model in detail.

4 Methodology

4.1 Problem Statement

The studied problem is stated as follows. Given a dialog flow $C = \{U_1, U_2, \dots, U_n\}$ including n utterances from multiple speakers, the objective is to recognize the personality trait P of the analyzed speaker s through the semantic content and the affective information in C .

The personality trait p is represented as a 5- d binary vector $[O, C, E, A, N]$ indicating the Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism respectively referring to the big-five personality theory. The affective information is indicated by the manual-annotated VAD vectors of words. Therefore, following the problem settings in some similar personality analysis works (Rissola et al., 2019; Jiang et al., 2020), we model the personality recognition as a binary classification problem over the five personality traits, respectively.

4.2 The HADE Model

To solve the challenges mentioned earlier of affective annotation shortage and effective speaker identification in dialog flow encoding, we design the HADE model as shown in Figure 2. HADE includes three modules: Utterance Encoding, Dialog Flow Encoding, and Utterance VAD regression. We will introduce these modules in detail.

4.2.1 Utterance Encoding

In conversation, utterances convey the personality trait of the speaker in addition to their semantic content (Mairesse et al., 2007). We choose BERT in the bottom layer of HADE to encode all the utterances, respectively. Pre-trained on the massive corpus, BERT does not rely on training with a large dataset to extract the semantics in utterances, which meets the challenge of data shortage.

For each utterance U_i in the dialog flow, we add a [CLS] and a [SEP] token in the beginning and the last position during tokenization. Hereafter, the U_1, U_2, \dots, U_n are separately encoded by the BERT encoder as a list of hidden representations

E_1, \dots, E_n , where the E_i is the embedding of the [CLS] token in U_i from the last pooling layer output in BERT.

4.2.2 Dialog Flow Encoding

By observing the dialog data, we found that the sentence-level interaction among the speakers (i.e., what are the current speaker talks to others and how others respond to the current speaker) is also essential to analyze the personality traits. Therefore, in the upper layer, we design the dialog flow encoding module based on a vanilla transformer encoder, as shown in the upper left of Figure 2. The transformer encoder receives the output of the bottom layer and the dialog state embeddings designed to identify the speakers for personality recognition.

First, $\{E_1, \dots, E_n\}$ are the utterance embeddings from the BERT encoder. Inspired by (Wolf et al., 2019b) and (Lin et al., 2019), we then construct the dialog state embedding $\{D_1, \dots, D_n\}$ to identify the utterance from the analyzed speaker s and the context. To be more specific, we use 1 to indicate the utterances from s , and 0 for utterances from other speakers. To feed the indicators into the model, we obtain the dialog state embedding by $D_i = MLP(is_uttr(U_i))$, where $is_uttr()$ outputs 1 and 0 as mentioned above. We also follow the original setting in (Vaswani et al., 2017) and construct the positional encodings $\{P_1, \dots, P_n\}$ to help the model understand the dialog flow:

$$\begin{aligned} P_i(2j) &= \sin\left(\frac{i}{10000^{\frac{2j}{d_{model}}}}\right) \\ P_i(2j+1) &= \cos\left(\frac{i}{10000^{\frac{2j}{d_{model}}}}\right) \end{aligned} \quad (1)$$

where i is the token position in the utterance, d_{model} is the size of the positional encodings, $j = 0, 1, \dots, d_{model}/2 - 1$.

After we get the three embeddings/encodings, we sum them together and feed them into the transformer model. We use all the last layer output of the transformer as the utterance representations R_1, \dots, R_n containing the sentence-level interactions through the self-attention mechanism. Then, we adopt the average pooling on the utterance representations for the personality classification minimizing the cross-entropy loss \mathcal{L}_{ce} during training:

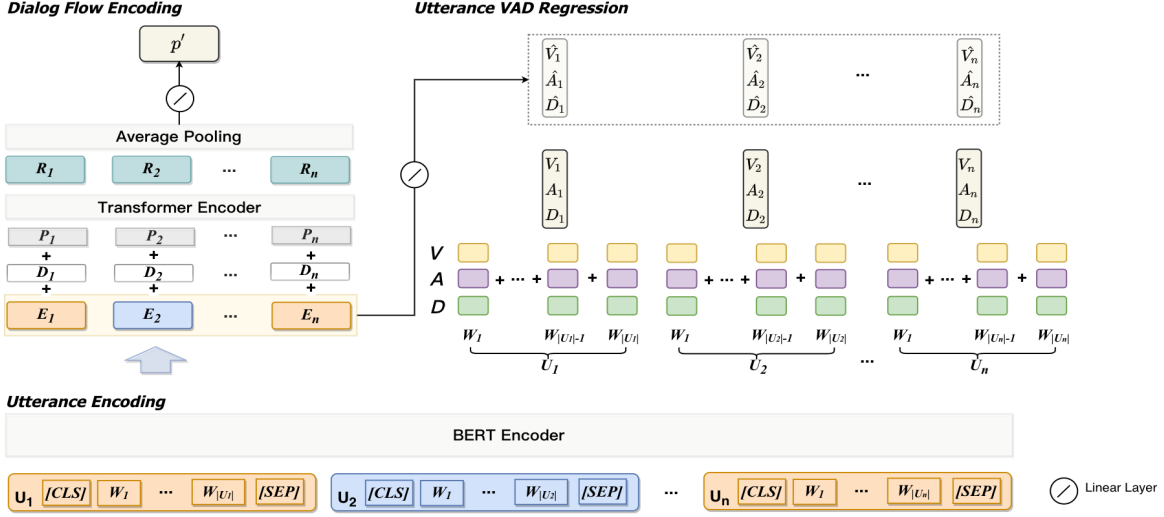


Figure 2: The model illustration of HADE. We use the same color to represent the utterances from the same speaker. e.g., U_1 and U_n .

$$\begin{aligned}
 R_i &= f_t(E_i + D_i + P_i) \\
 p' &= MLP\left(\sum_{i=1}^n \frac{R_i}{n}\right) \\
 \mathcal{L}_{ce} &= p \log(p') + (1 - p) \log(1 - p')
 \end{aligned} \quad (2)$$

where f_t is the transformer encoder, p' is the predicted personality label, and p is the ground truth personality label.

HADE first extracts the token-level semantic information in the bottom layer and then models the sentence-level interactions among speakers in the upper layer to facilitate personality recognition. The hierarchical modeling is verified as an efficient way to extract semantics in text with different granularities (Nawrot et al., 2021). It is also widely adopted in the conversation scenarios (Serban et al., 2017; Lynn et al., 2020; Gu et al., 2020).

4.2.3 Utterance VAD Regression

Although plenty of researches (Rank et al., 2013; Skowron et al., 2013; Asghar et al., 2018) work on the affective dialog systems, few works (Bauerhenne et al., 2020; Wen et al., 2021) combine it with personality analysis. One of the reasons is the lack of datasets with both emotion and personality annotations. Therefore, HADE extracts the affective information implicitly from utterances with VAD annotations for all the words in the utterances. Not only does this approach not need explicit emotion annotations, but it also can present the strength of emotions with numeric VAD vectors rather than discrete emotion labels.

Specifically, to preserve the encoding ability of BERT in HADE, we design an utterance VAD regression task with a regression head for the affective information extraction. The utterance VAD regression task supervises the model to capture affective information from the utterances.

For each utterance U_i in the input, we obtain the VAD vectors annotated by (Mohammad, 2018) of each word, which is also commonly utilized to represent affective information in conversation (Zhong et al., 2019; Colombo et al., 2019; Wen et al., 2021; Lee and Lee, 2021). The VAD vectors are numeric values ranging in $[0, 1]$ that indicate emotion intensity in three different dimensions. The valence measures positivity/negativity, arousal is for the excitement/calmness, and dominance is for the powerfulness/weakness.

$$\begin{aligned}
 V_i, A_i, D_i &= \sum_{j=1}^{|U_i|} V_j, A_j, D_j \\
 \hat{V}_i, \hat{A}_i, \hat{D}_i &= f(E_i) \\
 \mathcal{L}_{mse} &= \frac{1}{n} \sum_{i=1}^n (\sqrt{(V_i - \hat{V}_i)^2} \\
 &\quad + \sqrt{(A_i - \hat{A}_i)^2} + \sqrt{(D_i - \hat{D}_i)^2})
 \end{aligned} \quad (3)$$

We sum the VAD vectors of all the words in each utterance as the regression objectives $\{V_j, A_j, D_j\}$ for U_i . Then, each E_i obtained from the bottom layer is fed into a linear function f to regression the objective by generating $\hat{V}_i, \hat{A}_i, \hat{D}_i$. Finally, the regression loss \mathcal{L}_{mse} is calculated by averaging the regression loss for all the utterances. This proce-

383 dure is formulated in Formular 3.

384 4.2.4 Training Strategy

385 Our model is based on the bert-base-uncased model
386 implemented by Huggingface Transformer reposi-
387 tory (Wolf et al., 2019a). With 110 million pa-
388 rameters pre-trained on the massive corpus, we
389 found that it is challenging to incorporate such a
390 big model with the modules we designed in HATE.
391 Therefore, we fixed the look-up embeddings and
392 the parameters in the first 11 encoder layers in the
393 BERT encoder during training, only to fine-tune the
394 last encoder layer and the pooler layers in BERT
395 and train other modules designed by us at the same
396 time.

397 Although there are two optimization objectives
398 (\mathcal{L}_{ce} , \mathcal{L}_{mse}) for HADE, it is still designed to fo-
399 cus on personality recognition. So, we conduct a
400 two-stage training approach by first minimize the
401 overall loss function $\mathcal{L} = \mathcal{L}_{ce} + \mathcal{L}_{mse}$, and then re-
402 move the gradients in the auxiliary utterance VAD
403 regression task and only train HADE on \mathcal{L}_{ce} in the
404 second stage.

405 5 Experiment Settings

406 5.1 Dataset

407 Most personality recognition datasets focus on
408 the posts on social media (Schwartz et al., 2013)
409 or essays (Pennebaker and King, 1999). Record-
410 ing daily conversation for analysis, especially in-
411 cluding multiple speakers in the conversation, is
412 privacy-intrusive. So, we use the **FriendsPersona**
413 constructed by (Jiang et al., 2020) to evaluate our
414 method. It is a dialog script dataset with personality
415 annotations in 711 different dialogues, including
416 8,157 utterances. These dialogues are from the fa-
417 mous TV Series *Friends*. In **FriendsPersona**, the
418 average length of the dialog flows is 11.47 ut-
419 terances, while the average number of tokens for the
420 utterances is 16.27.

421 The personality in **FriendsPersona** is repre-
422 sented as 5-d binary vectors for the big-five traits.
423 The distribution of the personality annotations is
424 shown in Figure 3. The **AGR**, **CON**, **EXT**, **OPN**
425 and **NEU** indicate the big-five personality traits
426 respectively.

427 To facilitate the utterance VAD regression mod-
428 ule in our method, we also calculate the number of
429 tokens that have accurate VAD annotations from
430 (Mohammad, 2018) in the dataset. It suggests that
431 among 5,346 unique tokens, 2,796 of them have

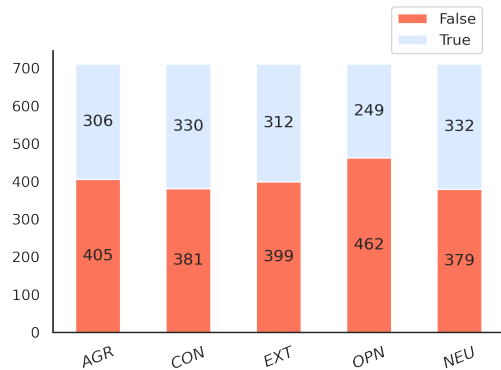


Figure 3: Personality annotations in **FriendsPersona**.

432 valid VAD annotations, the coverage is around
433 52.3%. As for the overall tokens, the corresponding
434 number is 28.6% (27,669/96,801).

435 5.2 Baseline Models

436 To show the effectiveness of our method, we
437 compare HADE with three state-of-the-art models
438 as below with a personality classification task on
439 **FriendsPersona**:

440 **HAN**: Hierarchical Attention Network (HAN)
441 is proposed in (Yang et al., 2016). It encodes
442 dialogue on both utterance and token levels by
443 RNN encoders with attention layers for personality
444 classification.

445 **RoBERTa(S)** and **RoBERTa(F)** are proposed in
446 (Jiang et al., 2020). They use the RoBERTa (Liu
447 et al., 2019) as the dialog encoder and try different
448 input for personality classification. **RoBERTa(S)**
449 only use the utterances from the analyzed speaker
450 as input; while **RoBERTa(F)** input all the ut-
451 terances within the whole dialog flow in their natural
452 order for classification.

453 5.3 Ablation Study Settings

454 To further investigate the effectiveness of different
455 modules in HADE and the methods we process
456 the input, we adopt an ablation study to com-
457 pare the performances of the following sub-models:

458 **Uttr**: We only use the BERT to encode the
459 utterances from the speaker s for personality
460 classification through a classification head.

461 **Uttr VAD**: Based on the **Uttr**, we add an aux-
462 iliary VAD regression head beside the original
463 classification head. The additional VAD regression
464

task is to supervise the model to extract affective information through a multi-task learning scheme.

VAD Embedding: We obtain the affective embeddings by inputting the VAD vectors of all the single words in the utterance into a linear layer. Then, we add the affective embeddings on the pre-trained look-up embeddings in BERT as the model input. This sub-model is to compare the way to utilize affective information with **Uttr VAD**.

Flow (Dialog State): We concatenate all the utterances in the whole dialog flow and feed it into the BERT encoder for personality classification. Simultaneously, we indicated the utterances from the analyzed speaker and the context with the segment embeddings in the BERT inspired by (Wolf et al., 2019b): **1** for utterances and **0** for the rest dialog context.

Hierarchical Flow: We first use the BERT model to encode each utterance in the bottom layer, and then in the second layer, we model the dialog flow as described in Section 4.2.2.

To sum up, **Uttr VAD** and **VAD Embedding** show the different ways to process the affective information; while **Flow (Dialog State)** and **Hierarchical Flow** are different approaches to model the dialog flow.

5.4 Implementation Details

During implementation, we pad all the utterances with [PAD] to a MAX_LEN of 64; besides, each dialog flow is padded to 20 utterances according to the dataset statistics. The dialog flows are fed into the models in batches of 16. As for the transformer model for the dialog flow encoding in HADE, we choose four heads and 512 as the d_{model} according to the best performance.

Due to the limited data, we do not conduct the warm-up training. Besides, we set the drop-out rate as 0.1 to avoid overfitting in training. We use the Adam (Kingma and Ba, 2014) as the optimization algorithm in training. The learning rate for each model is selected to refer to the best performance in evaluation.

6 Results Analysis

In this section, we describe the results of the evaluation of our method through experiments with the settings above. We analyze the result by answering

the following two research questions (RQs):

- RQ1: What is the performance of HADE in personality recognition in conversation?
- RQ2: How do the affective information and the dialog flow encoding influence the personality recognition HADE, respectively?

RQ1: What is the performance of our method in personality recognition in conversation?

We compare HADE with **HAN**, **RoBERTa(S)**, and **RoBERTa(F)** on binary personality classification. Following the settings in (Jiang et al., 2020), we conduct the 10-folds cross validation on **FriendsPersona**, and calculate the average classification accuracy of the test sets over the 10 splits. The results are shown in Table 2.

Model	AGR	CON	EXT	OPN	NEU	Avg
HAN	0.619	0.578	0.584	0.664	0.584	0.606
RoBERTa (S)	0.656	0.568	0.642	0.685	0.601	0.630
RoBERTa (F)	0.645	0.574	0.601	0.672	0.593	0.617
HADE	0.659	0.627	0.639	0.689	0.643	0.651

Table 2: Accuracy of binary personality classification.

We first focus on the performance of HADE. It achieves the highest accuracy (0.689) when predicting the Openness of the speakers. The lowest accuracy (0.627) occurs when indicating *Conscientiousness*. The average accuracy among the five personality traits is 0.651, and the standard deviation is around 0.021.

HADE outperforms other baseline models in four (**AGR**, **CON**, **OPN**, and **NEU**) over five personality traits with a considerable improvement. Besides, the average accuracy among the five personality traits of our model is also higher than the best baseline **RoBERTa(S)** over 3.3%. Although for **EXT**, our model does not outperform the **RoBERTa(S)**, the result is also close to the best. The results show that with our model design, the affective information and the dialog flow modeling can effectively help the personality recognition in conversation.

We also conclude that methods based on pre-trained language models are more competitive than those (e.g., **HAN**) with the traditional RNN encoders. Moreover, **RoBERTa(S)** beats **RoBERTa(F)** on overall performance, which indicates that even if input information is more, pure pre-trained language models are not appropriate to model the dialog flow data without modification.

Model	AGR	CON	EXT	OPN	NEU	Avg
Uttr	0.675 ± 0.023	0.613 ± 0.075	0.613 ± 0.134	0.791 ± 0.002	0.632 ± 0.087	0.665
Uttr VAD	0.700 ± 0.099	0.632 ± 0.047	0.625 ± 0.047	0.791 ± 0.003	0.621 ± 0.089	0.674
VAD Embedding	0.642 ± 0.084	0.588 ± 0.125	0.469 ± 0.103	0.716 ± 0.052	0.602 ± 0.120	0.603
Flow (Dialog State)	0.672 ± 0.066	0.625 ± 0.098	0.614 ± 0.033	0.656 ± 0.104	0.609 ± 0.021	0.641
Hierarchical Flow	0.710 ± 0.035	0.625 ± 0.109	0.623 ± 0.023	0.780 ± 0.030	0.612 ± 0.044	0.670
HADE	0.719 ± 0.100	0.627 ± 0.072	0.625 ± 0.062	0.787 ± 0.017	0.643 ± 0.091	0.680

Table 3: F1 scores for binary classification of personality traits.

RQ2: How do the affective information and the dialog flow encoding influence the personality recognition HADE, respectively?

After we verify the effectiveness of HADE, we are still curious about how and to what extent the modules in HATE influence the performance. Hence, we conduct an ablation study as the setting above. To better describe the personality classification performances, we use F-score (considers both precision and recall) rather than merely accuracy as the metric in the ablation study. Moreover, we run each experiment 10 times with ten different random seeds for dataset partition and model parameter initialization (except for parameters in BERT). We also record the standard deviations. The results are shown in Table 3.

In general, by integrating all the modules, HADE does outperform the **Uttr** in most of the traits, which verifies the benefit of our model design. By comparing **Uttr** and **Uttr VAD**, we can see that adding the VAD regression task improves the accuracy in **AGR** and **CON**, but slightly reduce the performance in **EXT** and **NEU**. Consequently, the average performance is still improved. Nevertheless, when we focus on the **VAD Embedding**, which modifies the look-up embeddings in the pre-trained language model by VAD vectors, the accuracy decrease in all the traits compared with both **Uttr** and **Uttr VAD**. The reason is that VAD vectors damage the original semantics stored in the look-up embeddings pre-trained in the massive corpus. However, the training dataset is too small to supervise the model to learn to process the VAD vectors in the input. Therefore, even both methods integrate the affective information in the model; only the appropriate way can preserve the strength of BERT and improve the performance.

Then, we turn to the dialog flow modeling. We compare the results between the **Uttr** and **Flow (Dialog State)** and found that although incorporating the dialog flow improves the performance in **CON** and **EXT**, it decreases the performance in

other traits, especially in predicting **OPN**. It shows that similar to **VAD Embedding**, directly incorporating with the dialog state embeddings in the pre-trained language model fails to make it learn to process such information appropriately with such a small training set. However, if we focus on the performance of **Hierarchical Flow**, we can see the results are much better. So, hierarchically and separately modeling the utterances (in token level) and the dialog flow (in sentence-level) is a better approach to utilize pre-trained language models in our problem.

Combining **Uttr VAD** and **Hierarchical Flow** forms HADE and improves both sub-models. Nevertheless, we can also see that the average performance of **Uttr VAD** is slightly higher than **Hierarchical Flow**, even they are calculated on ten different random seeds. So, we conclude that affective information is more important in personality recognition under the design of HADE.

7 Conclusion and Future Work

We propose HADE to extract affective information implicitly and model the dialog flow for personality recognition in conversation. We utilize pre-defined VAD vectors of single words and design a hierarchical architecture to model the dialog flow, which solves the challenging issues met in existing works. Our model outperforms state-of-the-art methods on a public conversation dataset. Through ablation study, our approach is validated as an effective way to apply affective information into the model design with pre-trained language models.

HADE outperforms state-of-the-art models on **FriendsPersona**; we also want to verify the generality of HADE in other conversation scenarios. One significant barrier is that conversation datasets with personality annotations are rare due to privacy concerns. So, in future work, we will investigate the conversational dataset construction in a privacy-nonintrusive manner so that HADE, and even more approaches can be evaluated.

643
644
645
646
647

648
649
650

651
652
653
654

655
656
657

658
659
660
661
662

663
664
665

666
667
668
669
670
671

672
673
674
675

676
677
678

679
680
681
682

683
684

685
686

687
688
689
690

691
692
693
694

References

Nabiha Asghar, Pascal Poupart, Jesse Hoey, Xin Jiang, and Lili Mou. 2018. Affective neural response generation. In *European Conference on Information Retrieval*, pages 154–166. Springer.

Gene Ball and Jack Breese. 2000. Emotion and personality in a conversational agent. *Embodied conversational agents*, 189.

Siqi Bao, Huang He, Fan Wang, Hua Wu, and Haifeng Wang. 2019. Plato: Pre-trained dialogue generation model with discrete latent variable. *arXiv preprint arXiv:1910.07931*.

C Bauerhenne, A Gammoudi, M Moussaoui, J Reichelt, and J Wang. 2020. Emotional states and personality profiles in conversational ai.

Margaret M Bradley and Peter J Lang. 1999. Affective norms for english words (anew): Instruction manual and affective ratings. Technical report, Technical report C-1, the center for research in psychophysiology

Raymond B Cattell. 1943. The description of personality: Basic traits resolved into clusters. *The journal of abnormal and social psychology*, 38(4):476.

Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.

Pierre Colombo, Wojciech Witon, Ashutosh Modi, James Kennedy, and Mubbasir Kapadia. 2019. Affect-driven dialog generation. *arXiv preprint arXiv:1904.02793*.

Paul T Costa and Robert R McCrae. 1992. Normal personality assessment in clinical practice: The neo personality inventory. *Psychological assessment*, 4(1):5.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Jess Feist and Gregory Feist. 2012. J.(2008). theories of personality.

Susan T Fiske and Shelley E Taylor. 1991. *Social cognition*. McGraw-Hill Book Company.

Jennifer Golbeck, Cristina Robles, and Karen Turner. 2011. Predicting personality with social media. In *CHI'11 extended abstracts on human factors in computing systems*, pages 253–262.

Xiaodong Gu, Kang Min Yoo, and Jung-Woo Ha. 2020. Dialogbert: Discourse-aware response generation via learning to recover and rank utterances. *arXiv preprint arXiv:2012.01775*.

Meng-Ju Han, Chia-How Lin, and Kai-Tai Song. 2012. Robotic emotional expression generation based on mood transition and personality model. *IEEE transactions on cybernetics*, 43(4):1290–1303.

Eddie Harmon-Jones, Philip A Gable, and Tom F Price. 2013. Does negative affect always narrow and positive affect always broaden the mind? considering the influence of motivational intensity on cognitive scope. *Current Directions in Psychological Science*, 22(4):301–307.

Francisco Iacobelli, Alastair J Gill, Scott Nowson, and Jon Oberlander. 2011. Large scale personality classification of bloggers. In *international conference on affective computing and intelligent interaction*, pages 568–577. Springer.

Hang Jiang, Xianzhe Zhang, and Jinho D Choi. 2020. Automatic text-based personality recognition on monologues and multiparty dialogues using attentive networks and contextual embeddings (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13821–13822.

Oliver P John, Alois Angleitner, and Fritz Ostendorf. 1988. The lexical approach to personality: A historical review of trait taxonomic research. *European journal of Personality*, 2(3):171–203.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Randall J Larsen and David M Buss. 2008. Psihologija ličnosti. *Naklada Slap, Jastrebarsko*.

Joosung Lee and Woojin Lee. 2021. Compm: Context modeling with speaker’s pre-trained memory tracking for emotion recognition in conversation. *arXiv preprint arXiv:2108.11626*.

Zhaojiang Lin, Andrea Madotto, Jamin Shin, Peng Xu, and Pascale Fung. 2019. Moel: Mixture of empathetic listeners. *arXiv preprint arXiv:1908.07687*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Veronica Lynn, Niranjan Balasubramanian, and H Andrew Schwartz. 2020. Hierarchical modeling for user personality prediction: The role of message-level attention. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5306–5316.

François Mairesse and Marilyn Walker. 2006. Automatic recognition of personality in conversation. In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, pages 85–88.

749	François Mairesse, Marilyn A Walker, Matthias R Mehl, and Roger K Moore. 2007. Using linguistic cues for the automatic recognition of personality in conversation and text. <i>Journal of artificial intelligence research</i> , 30:457–500.	H Andrew Schwartz, Johannes C Eichstaedt, Margaret L Kern, Lukasz Dziurzynski, Stephanie M Ramones, Megha Agrawal, Achal Shah, Michal Kosinski, David Stillwell, Martin EP Seligman, et al. 2013. Personality, gender, and age in the language of social media: The open-vocabulary approach. <i>PloS one</i> , 8(9):e73791.	804 805 806 807 808 809 810
754	Naoki Masuyama, Chu Kiong Loo, and Manjeevan Seera. 2018. Personality affected robotic emotional model with associative memory for human-robot interaction. <i>Neurocomputing</i> , 272:213–225.	Iulian Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 31.	811 812 813 814 815 816
758	Matthias R Mehl, Samuel D Gosling, and James W Pennebaker. 2006. Personality in its natural habitat: Manifestations and implicit folk theories of personality in daily life. <i>Journal of personality and social psychology</i> , 90(5):862.	Marcin Skowron, Mathias Theunis, Stefan Rank, and Arvid Kappas. 2013. Affect and social processes in online communication—experiments with an affective dialog system. <i>IEEE Transactions on Affective Computing</i> , 4(3):267–279.	817 818 819 820 821
763	Albert Mehrabian. 1995. Relationships among three general approaches to personality description. <i>The journal of Psychology</i> , 129(5):565–581.	Alireza Souri, Shafiqeh Hosseinpour, and Amir Masoud Rahmani. 2018. Personality classification based on profiles of social networks’ users and the five-factor model of personality. <i>Human-centric Computing and Information Sciences</i> , 8(1):1–15.	822 823 824 825 826
766	Albert Mehrabian. 1996. Analysis of the big-five personality factors in terms of the pad temperament model. <i>Australian journal of Psychology</i> , 48(2):86–92.	Edward P Tighe, Jennifer C Ureta, Bernard Andrei L Pollo, Charibeth K Cheng, and Remedios de Dios Bulos. 2016. Personality trait classification of essays with the application of feature reduction. In <i>SAIIP@IJCAI</i> , pages 22–28.	827 828 829 830 831
770	Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019. Pretraining methods for dialog context representation learning. <i>arXiv preprint arXiv:1906.00414</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in neural information processing systems</i> , pages 5998–6008.	832 833 834 835 836
774	Saif Mohammad. 2018. Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 174–184.	David Watson and Lee Anna Clark. 1992. On traits and temperament: General and specific factors of emotional experience and their relation to the five-factor model. <i>Journal of personality</i> , 60(2):441–476.	837 838 839 840
779	Daniel Ricardo Jaimes Moreno, Juan Carlos Gomez, Dora-Luz Almanza-Ojeda, and Mario-Alberto Ibarra-Manzano. 2019. Prediction of personality traits in twitter users with latent features. In <i>2019 International Conference on Electronics, Communications and Computers (CONIELECOMP)</i> , pages 176–181. IEEE.	Zhiyuan Wen, Jiannong Cao, Ruosong Yang, Shuaiqi Liu, and Jiaying Shen. 2021. Automatically select emotion for response via personality-affected emotion transition. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 5010–5020.	841 842 843 844 845 846
782	Piotr Nawrot, Szymon Tworowski, Michał Tyrolski, Łukasz Kaiser, Yuhuai Wu, Christian Szegedy, and Henryk Michalewski. 2021. Hierarchical transformers are more efficient language models. <i>arXiv preprint arXiv:2110.13711</i> .	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019a. Huggingface’s transformers: State-of-the-art natural language processing. <i>arXiv preprint arXiv:1910.03771</i> .	847 848 849 850 851 852
786	James W Pennebaker and Laura A King. 1999. Linguistic styles: language use as an individual difference. <i>Journal of personality and social psychology</i> , 77(6):1296.	Thomas Wolf, Victor Sanh, Julien Chaumond, and Clement Delangue. 2019b. Transfertransfo: A transfer learning approach for neural network based conversational agents. <i>arXiv preprint arXiv:1901.08149</i> .	853 854 855 856
791	Stefan Rank, Marcin Skowron, and David Garcia. 2013. Dyads to groups: modeling interactions with affective dialog systems. <i>International Journal of Computational Linguistics</i> , 4(1):22–37.		
795	Esteban Andres Rissola, Seyed Ali Bahrainian, and Fabio Crestani. 2019. Personality recognition in conversations using capsule neural networks. In <i>IEEE/WIC/ACM International Conference on Web Intelligence</i> , pages 180–187.		

857 Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He,
858 Alex Smola, and Eduard Hovy. 2016. Hierarchical at-
859 tention networks for document classification. In *Pro-
860 ceedings of the 2016 conference of the North Ameri-
861 can chapter of the association for computational lin-
862 guistics: human language technologies*, pages 1480–
863 1489.

864 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,
865 Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing
866 Liu, and Bill Dolan. 2019. Dialogpt: Large-scale
867 generative pre-training for conversational response
868 generation. *arXiv preprint arXiv:1911.00536*.

869 Peixiang Zhong, Di Wang, and Chunyan Miao. 2019.
870 Knowledge-enriched transformer for emotion de-
871 tection in textual conversations. *arXiv preprint
872 arXiv:1909.10681*.

873 Hao Zhou, Pei Ke, Zheng Zhang, Yuxian Gu, Yinhe
874 Zheng, Chujie Zheng, Yida Wang, Chen Henry
875 Wu, Hao Sun, Xiaocong Yang, et al. 2021. Eva:
876 An open-domain chinese dialogue system with
877 large-scale generative pre-training. *arXiv preprint
878 arXiv:2108.01547*.