



Successful New-entry Prediction for Multi-Party Online Conversations via Latent Topics and Discourse Modeling

Lingzhi Wang

lzwang@se.cuhk.edu.hk

The Chinese University of Hong Kong
Hong Kong, China

Xingshan Zeng

zxshamson@gmail.com

Huawei Noah's Ark Lab
Hong Kong, China

Jing Li

jing-amelia.li@polyu.edu.hk

The Hong Kong Polytechnic University
Hong Kong, China

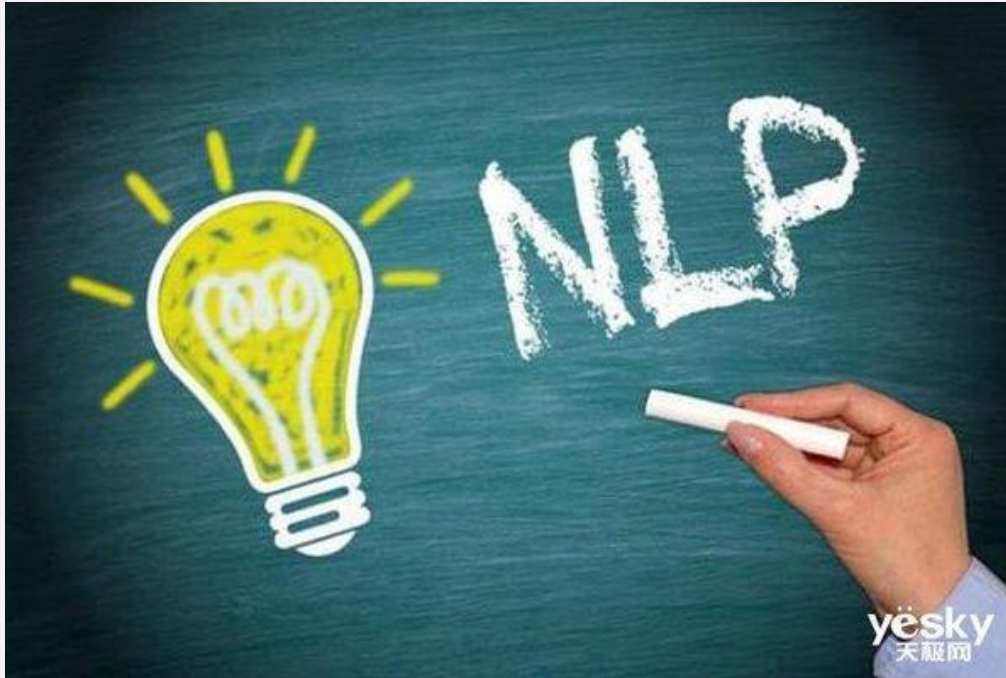
Kam-Fai Wong

kfwong@se.cuhk.edu.hk

The Chinese University of Hong Kong
Hong Kong, China

(WWW-2022)

Reported by Jia Wang



- 1. Introduction**
- 2. Approach**
- 3. Experiments**

Introduction

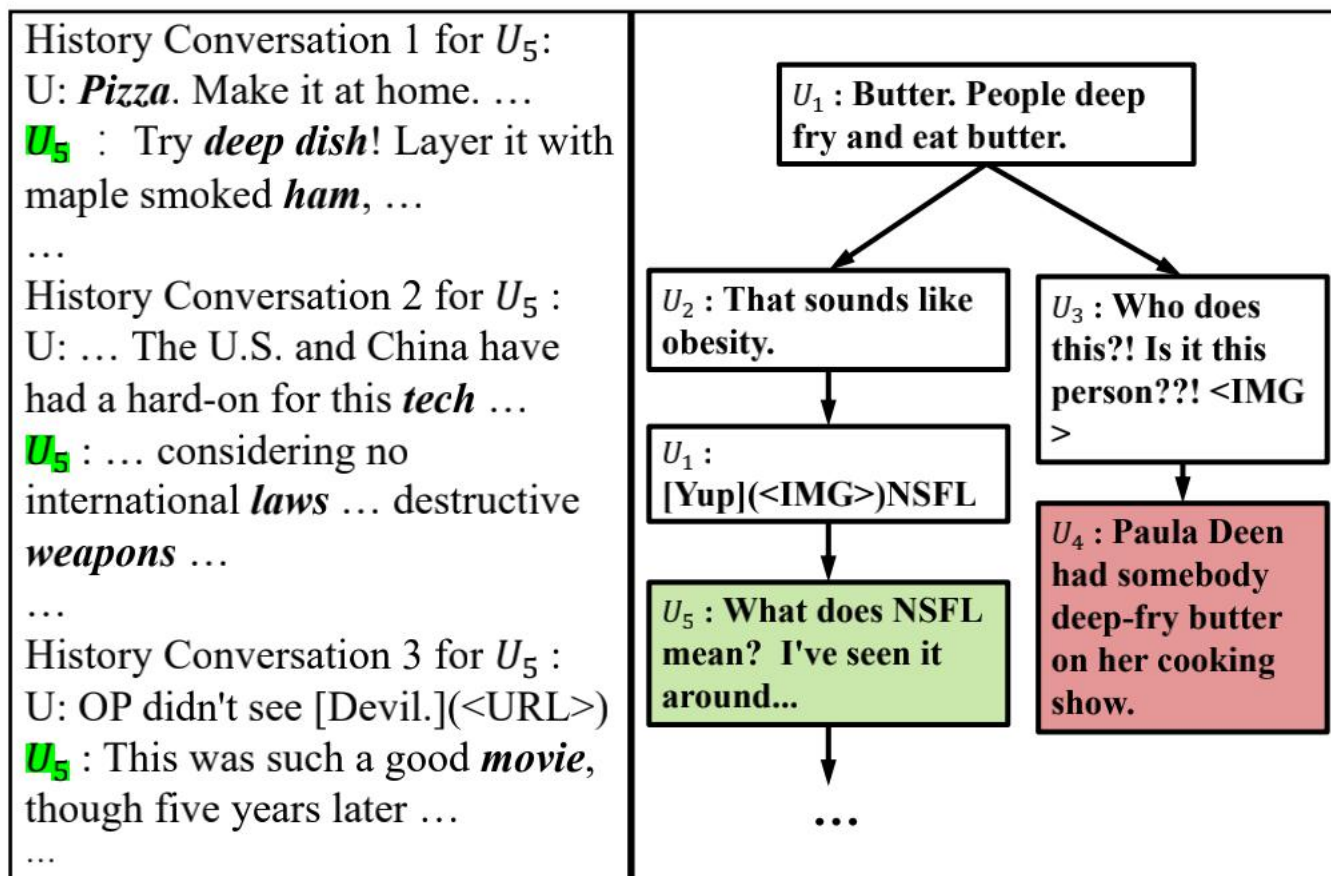


Figure 1: A Reddit conversation on the left part (U_i : the i -th user). U_3 made a successful engagement (i.e. receiving a reply, we omit here to save space). The right column shows U_3 's chatting history, where the topic words are in *bold and italic*.



Introduction

The main contributions of this work can be summarized as follows:

- We first formulate the task of successful new-entry prediction and contribute two large-scale datasets, Twitter and Reddit. The SNP task can benefit the development of **online assistants and early socialization strategies**.
- We propose a novel framework combining **unsupervised and supervised neural networks**. VAE and RNN-based modules are incorporated for the personalized user engagement prediction via learning latent topic and discourse representation.
- Experimental results on both Twitter and Reddit show that the proposed model significantly outperforms the baselines. For example, we achieve 34.6 F1 on Reddit compared with 32.5 achieved by a BERT-based method.

Approach

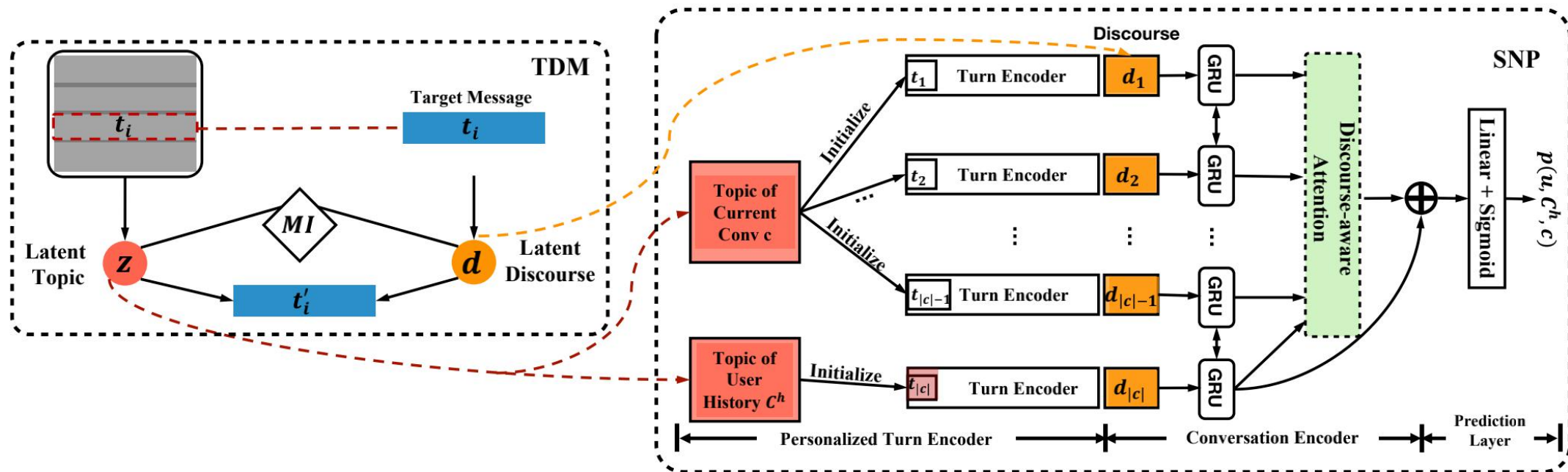
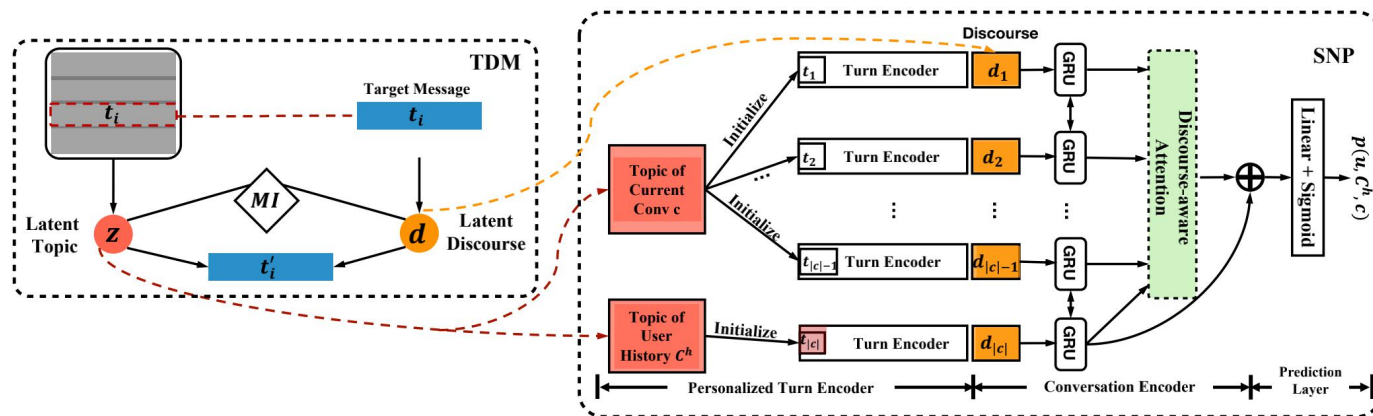


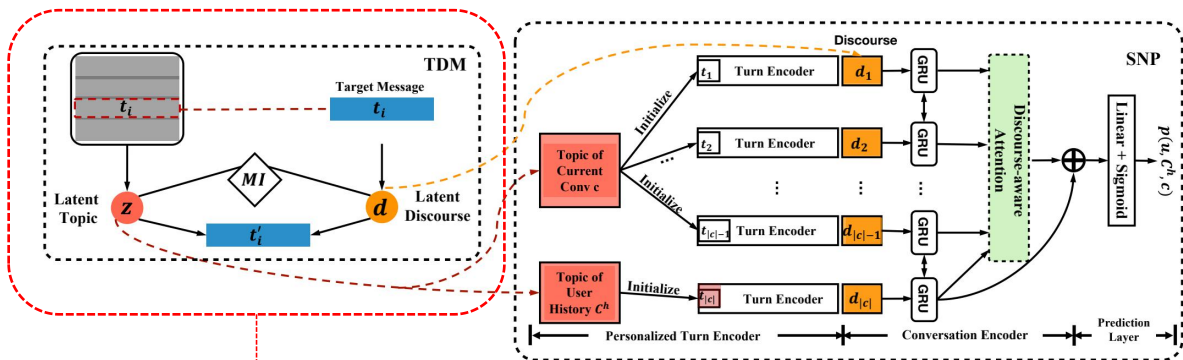
Figure 2: Our generic framework for successful new-entry prediction. It contains two modules: Topic and Discourse Modeling (TDM) and Successful New-entry Prediction (SNP). SNP consists of three parts: turn encoder, conversation encoder and prediction layer.

Approach



The input for our model can be divided into two parts: the observed conversation c and the history conversation set $C^h = \{c_1^h, c_2^h, \dots, c_k^h\}$ of the newcomer u , where k is the number of history conversations obtained from training set. The conversation c is formalized as a sequence of turns (e.g., posts or tweets) $\{t_1, t_2, \dots, t_{|c|}\}$, and the $|c|^{th}$ turn is posted by the newcomer u (we predict whether u can get others' response afterwards). The conversations in user's history conversation set C^h are organized similarly into the sequences of turns. For output, we yield a Bernoulli distribution $p(u, C^h, c)$ to indicate the estimated likelihood of whether u gets responses from other participants (successful new-entries).

Topic and Discourse Modeling (TDM)



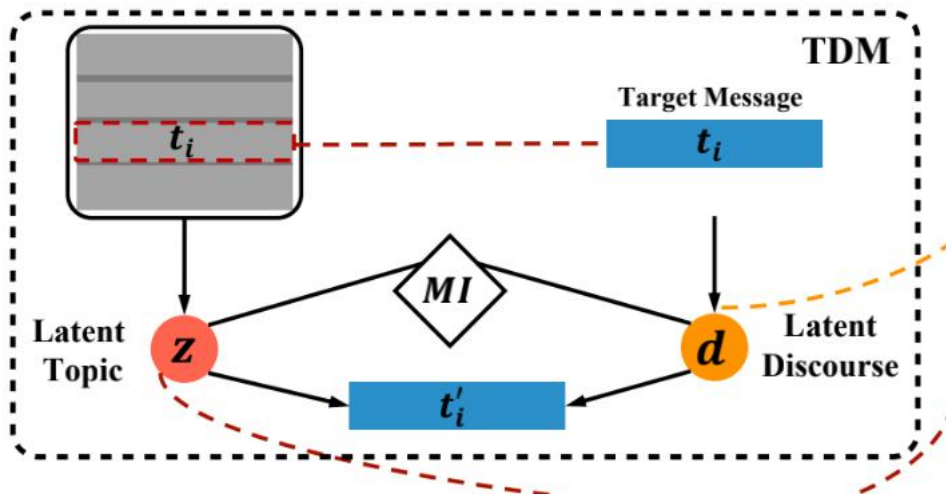
conversation c , $\begin{cases} e_c \\ d = \langle d_1, \dots, d_{|c|} \rangle. \end{cases}$

C^h of the newcomer $u \rightarrow e_u$.

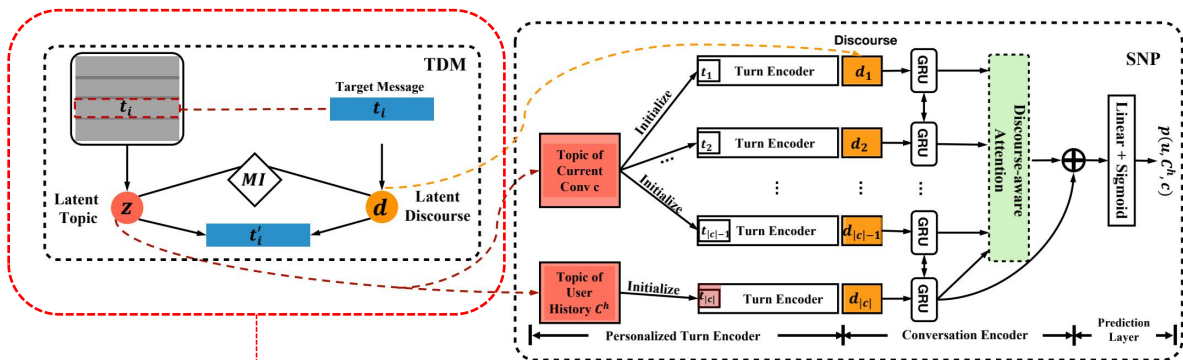
Encode step

$$\begin{aligned} \mu &= f_\mu(fe(c_{bow})), \log \sigma = f_\sigma(fe(c_{bow})) \\ \pi &= \text{softmax}(f_\pi(t_{i_{bow}})) \end{aligned} \quad (1)$$

where $f_*(\cdot)$ is neural perceptrons performing linear transformations activated with an ReLU function [33].



Topic and Discourse Modeling (TDM)

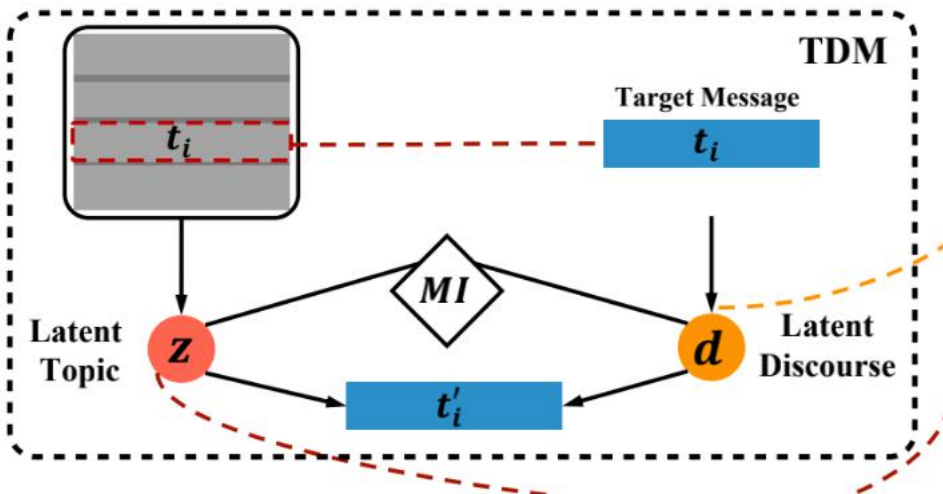


$$\text{conversation } c, \begin{cases} \mathbf{e}_c \\ \mathbf{d} = \langle \mathbf{d}_1, \dots, \mathbf{d}_{|c|} \rangle. \end{cases}$$

C^h of the newcomer $u \rightarrow \mathbf{e}_u$.

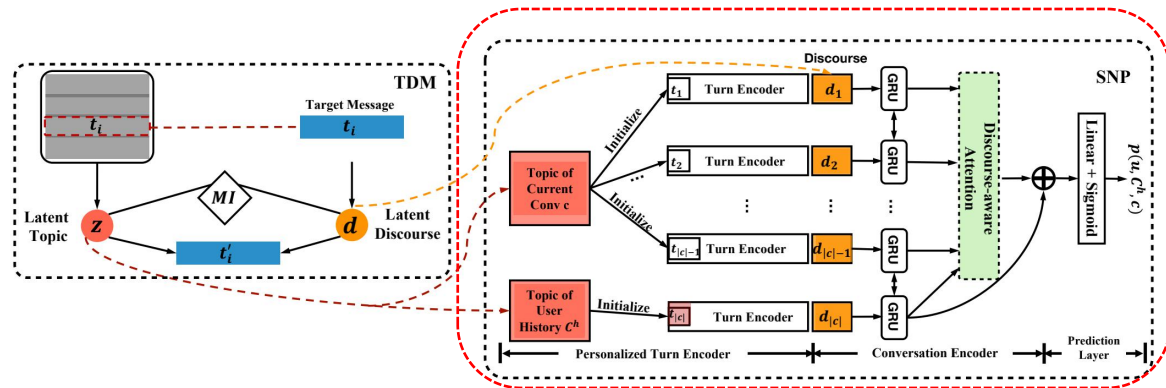
Decode step

- Draw latent topic $\mathbf{z} \sim \mathcal{N}(\mu, \sigma^2)$.
- Topic mixture $\theta = \text{softmax}(f_\theta(\mathbf{z}))$.
- Draw the latent discourse $\mathbf{d} \sim \text{Multi}(\pi)$.
- For the n -th word in the conversation:
 - $\beta_n = \text{softmax}(f_{\phi^T}(\theta) + f_{\phi^D}(\mathbf{d}))$
 - Draw the word $w_n \sim \text{Multi}(\beta_n)$.



In particular, the weight matrix of $f_{\phi^T}(\cdot)$ (after the softmax normalization) is considered as the topic-word distribution ϕ^T . We can also get the discourse-word distribution ϕ^D in a similar way.

Successful New-entry Prediction (SNP)



Personalized Turn Encoder

For each turn t_i in conversation c ,

$$\begin{aligned} & \downarrow w_j \\ & \langle r_{i1}, r_{i2}, \dots, r_{i,n_i} \rangle \end{aligned}$$

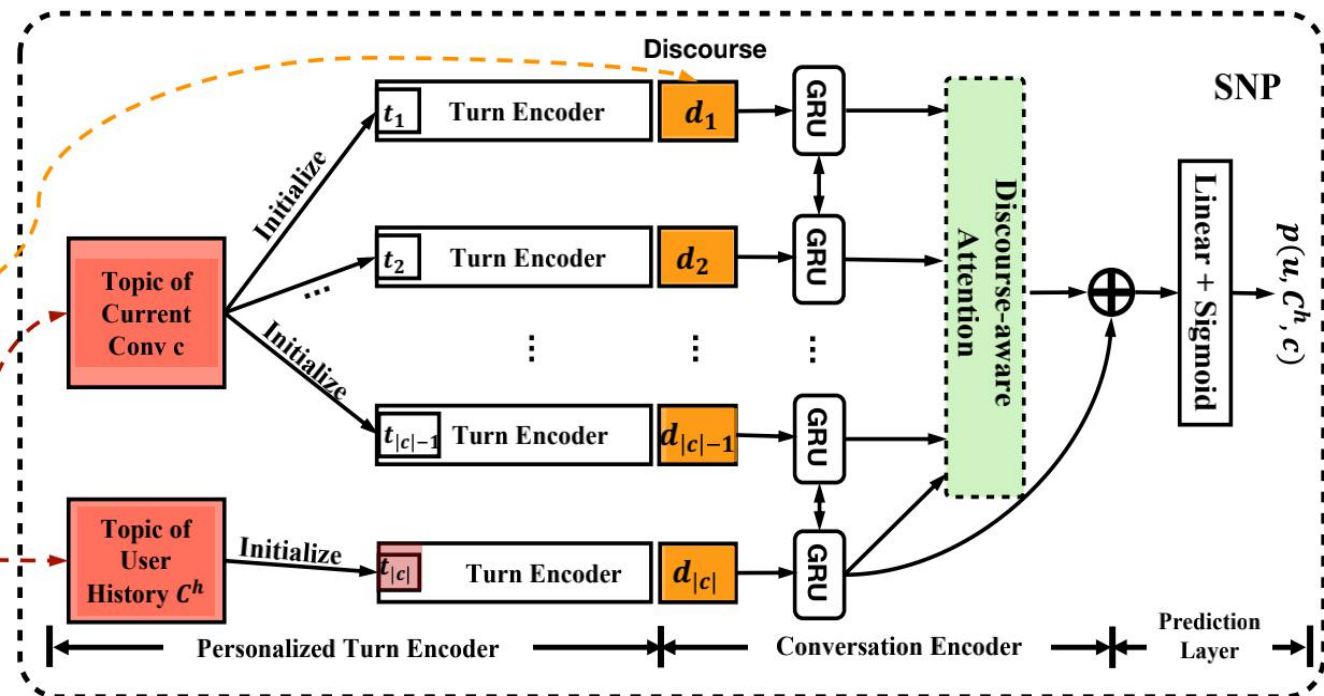
We divide the observed conversation turns into context turns (turns before the last turn) and query turn (last turn, posted by newcomer u). For query turn, we use u 's topic representation e_u (produced by TDM module in Section 3.2) to initialize the aforementioned Bi-GRU.

$$\vec{h}_{i,0} = \overleftarrow{h}_{i,n_i} = W^P e + b^P \quad e \text{ is } e_u \text{ or } e_c$$

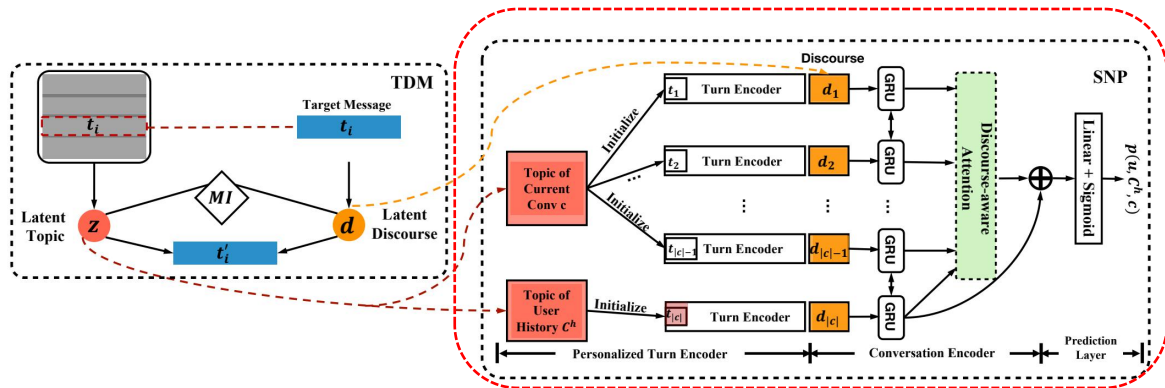
$$\vec{h}_{ij} = f_{GRU}(r_{ij}, \vec{h}_{i,j-1}), \quad \overleftarrow{h}_{ij} = f_{GRU}(r_{ij}, \overleftarrow{h}_{i,j+1}) \quad (2)$$

$$h_i = [\vec{h}_{i,n_i}; \overleftarrow{h}_{i,0}]$$

$$\langle h_1, h_2, \dots, h_{|c|} \rangle.$$



Successful New-entry Prediction (SNP)



Discourse-aware Conversation Encoder

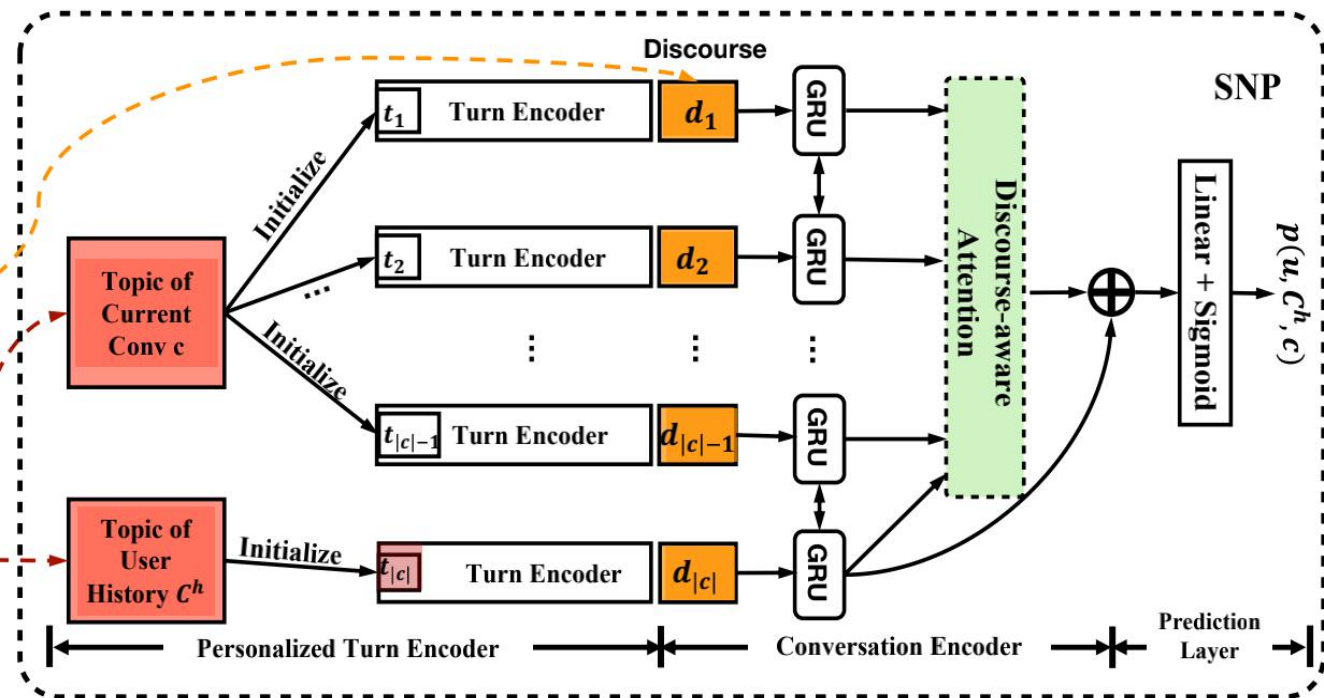
$$\vec{h}_j^d = f_{GRU}(s_j, \vec{h}_{j-1}^d), \quad \overleftarrow{h}_j^d = f_{GRU}(s_j, \overleftarrow{h}_{j+1}^d) \quad (3)$$

where $s_j = [h_j; d_j]$ and the representation of each turn after GRU is $h_j^d = [\vec{h}_j^d; \overleftarrow{h}_j^d]$.

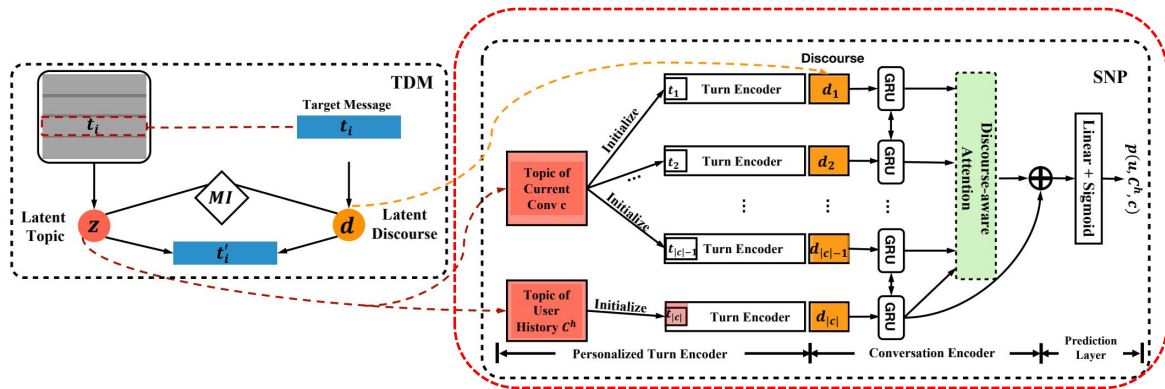
$$a_j = f_d(\text{argmax}(d_j)) \quad (4)$$

where $\text{argmax}(d_j)$ means the learned discourse behavior to turn j , and $f_d(\cdot)$ maps the discourse behaviors to different weight values.

$$h^c = [h_{|c|}^d; \sum_j \text{softmax}(a_j) h_j^d] \quad (5)$$



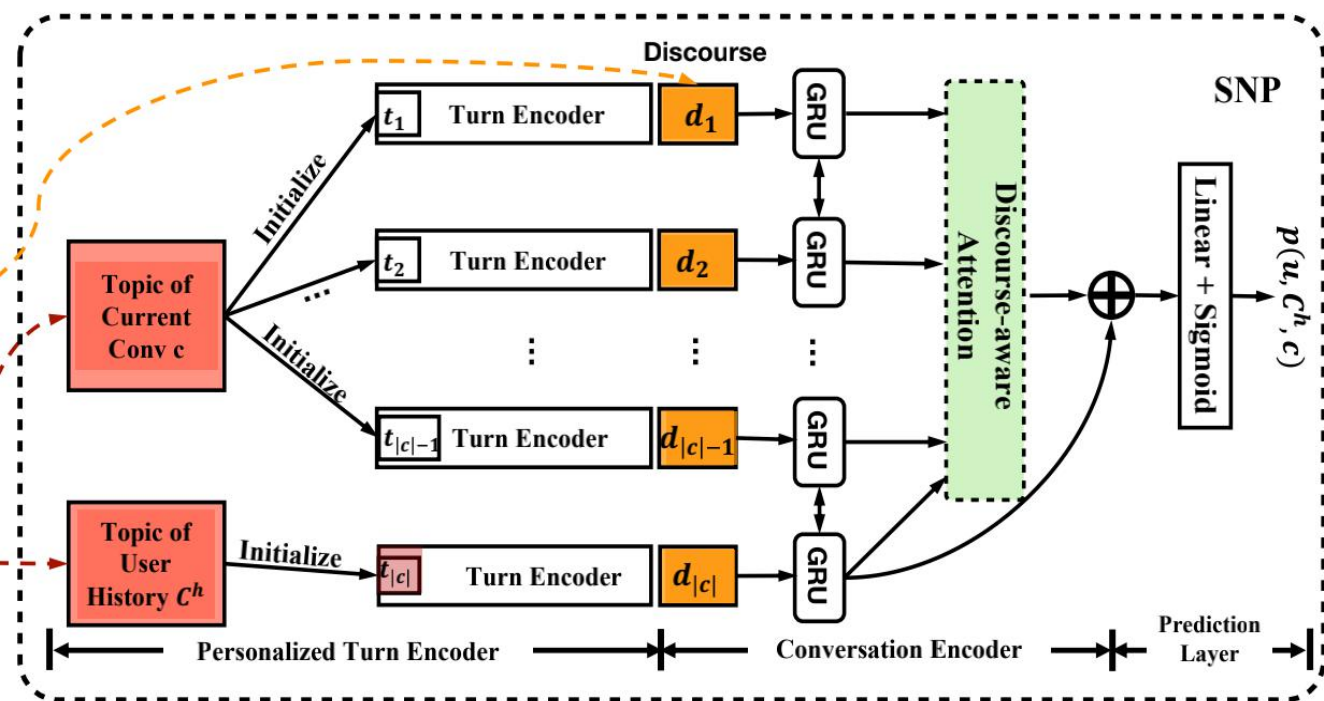
Successful New-entry Prediction (SNP)



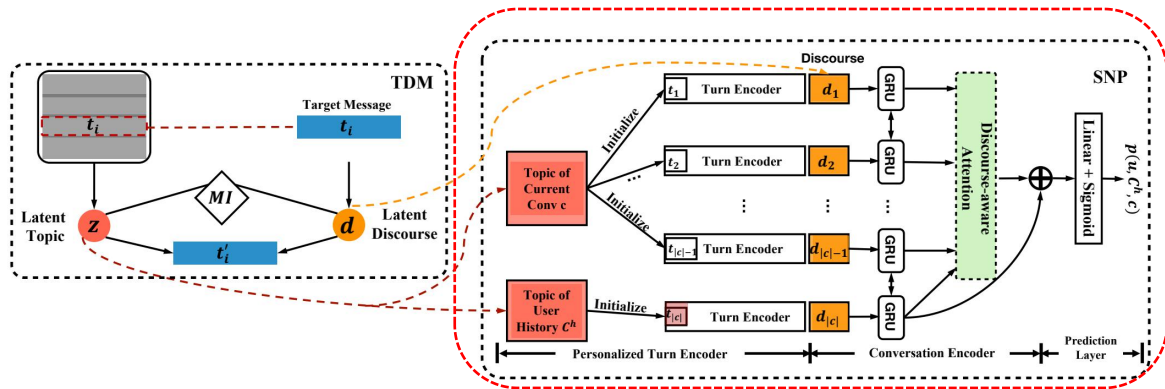
Prediction Layer

$$p(u, C^h, c) = \sigma(\mathbf{w}^T \mathbf{h}^c + b) \quad (6)$$

where \mathbf{w}^T and b are trainable, and $\sigma(\cdot)$ is the sigmoid activation function.



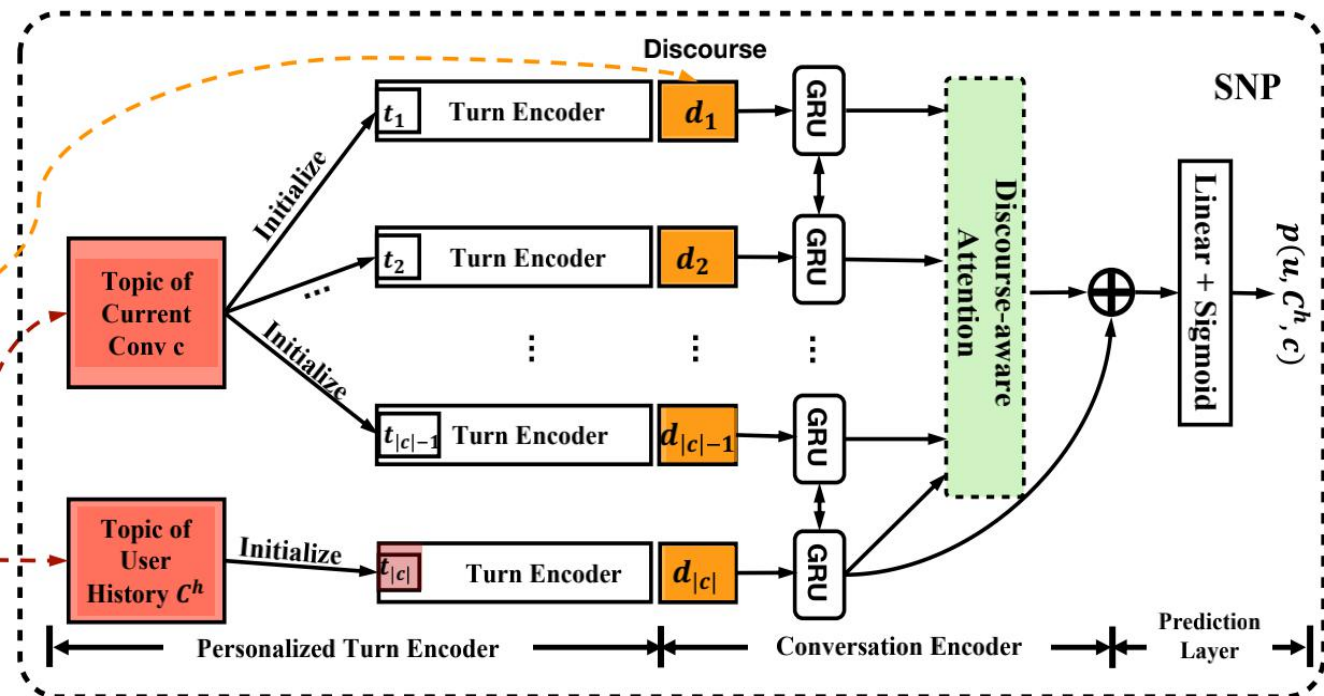
Learning Objective



$$\mathcal{L} = \mathcal{L}_{TDM} + \mathcal{L}_{SNP} \quad (7)$$

$$\mathcal{L}_{TDM} = \mathcal{L}_z + \mathcal{L}_d + \mathcal{L}_t - \lambda \mathcal{L}_{MI} \quad (8)$$

where \mathcal{L}_z and \mathcal{L}_d are objectives about learning topics and discourse, \mathcal{L}_t is the loss for target message reconstruction, and \mathcal{L}_{MI} ensures that topics and discourse learn differently.

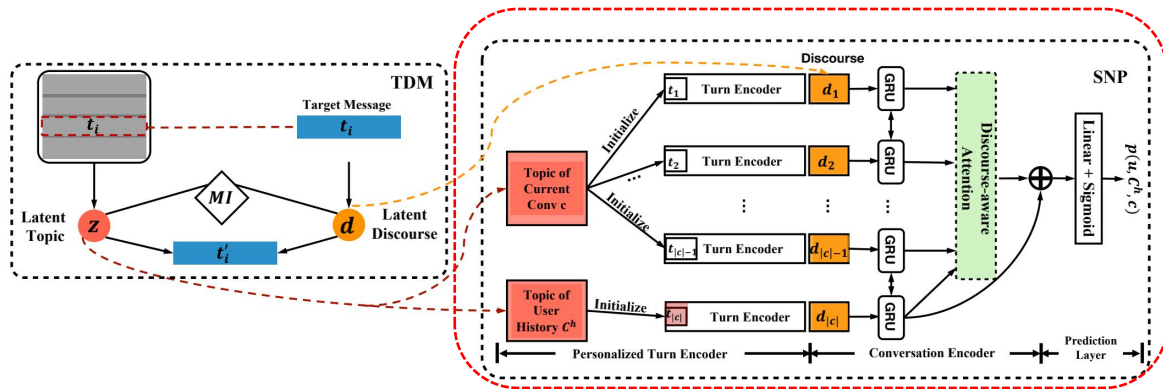


$$\mathcal{L}_z = \mathbb{E}_{q(z|c)} [p(c|z)] - D_{KL}(q(z|c) || p(z)) \quad (9)$$

$$\mathcal{L}_d = \mathbb{E}_{q(d|t)} [p(t|d)] - D_{KL}(q(d|t) || p(d)) \quad (10)$$

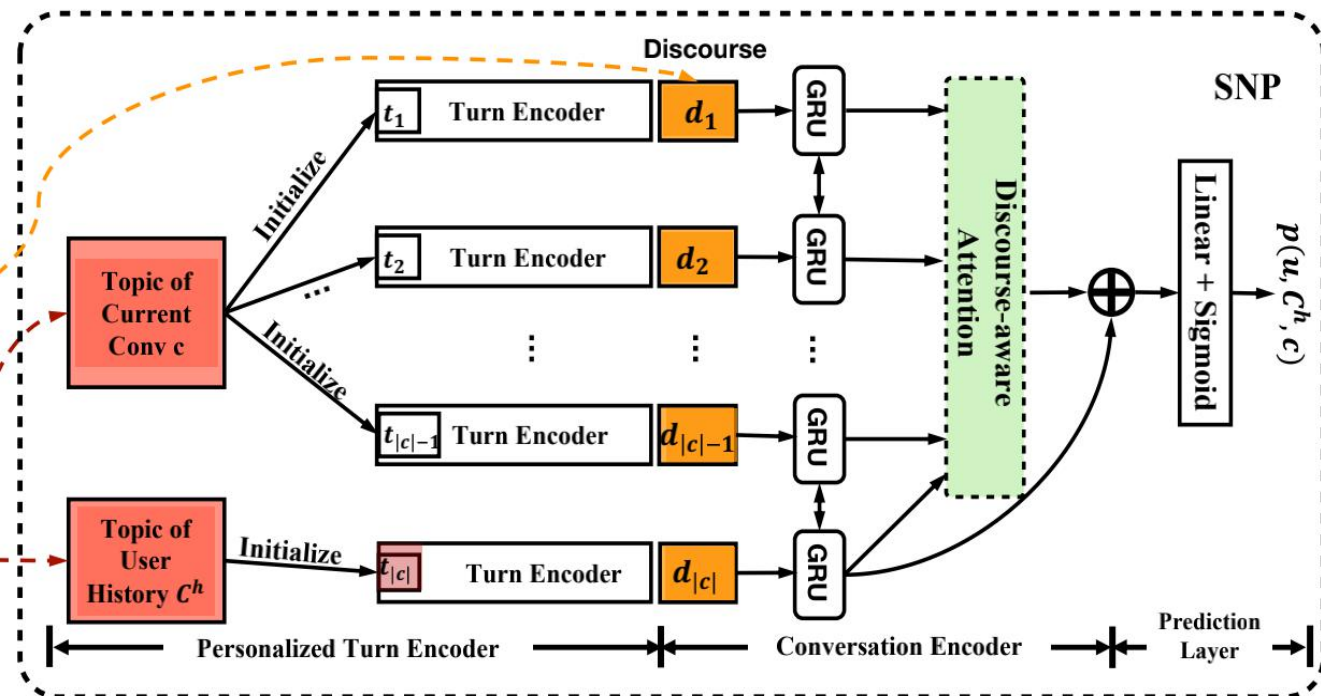
where $q(z|c)$ and $q(d|t)$ are approximated posterior probabilities describing how the latent topic z and the latent discourse d are generated from the conversations and message turns. $p(c|z)$ and $p(t|d)$ represent the corpus likelihoods conditioned on the latent variables. $p(z)$ follows the standard normal prior $\mathcal{N}(0, \mathbf{I})$ and $p(d)$ is the uniform distribution $Unif(0, 1)$. D_{KL} refers to the Kullback-Leibler divergence that ensures the approximated posteriors to be close to the true ones.

Learning Objective



$$\mathcal{L}_t = \mathbb{E}_{q(z|c)q(d|t)} [\log p(t | z, d)] \quad (11)$$

$$\mathcal{L}_{MI} = \mathbb{E}_{q(z)} D_{KL}(p(\mathbf{d} | z) || p(\mathbf{d})) \quad (12)$$



$$\mathcal{L}_{SNP} = - \sum_i \mu y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \quad (13)$$

where \hat{y}_i denotes the probability estimated from $p(u, C^h, c)$ for the i -th instance, and y_i is the corresponding binary ground truth label (1 for successful entries and 0 for the opposite). To take the potential data imbalance into account, we also adopt a trade-off weight μ to give more weight to the minority class. μ is set based on the proportion of positive and negative instances in the training set.



Experiments

Table 1: Statistics of Twitter and Reddit datasets.

	Twitter	Reddit
# of users	53,488	96,001
# of convs	37,339	69,428
# of conv turns	179,265	236,764
# of successful entries	29,340	12,199
# of failed entries	7,999	57,229
Avg turn number per conv	4.8	3.4
Avg token number per turn	20.5	20.7
Ratio of newcomer with history	0.59	0.62
Avg # of history for newcomers	2.5	6.3

Experiments

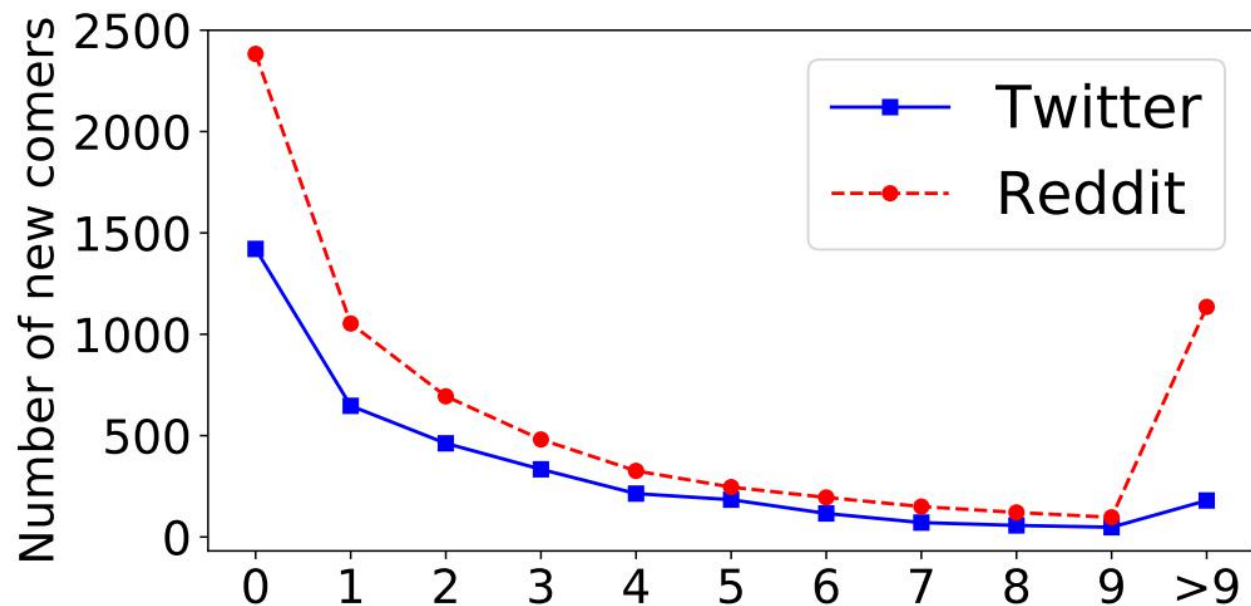


Figure 3: The distribution over the number of history conversations (X-axis). Y-axis: number of newcomers.

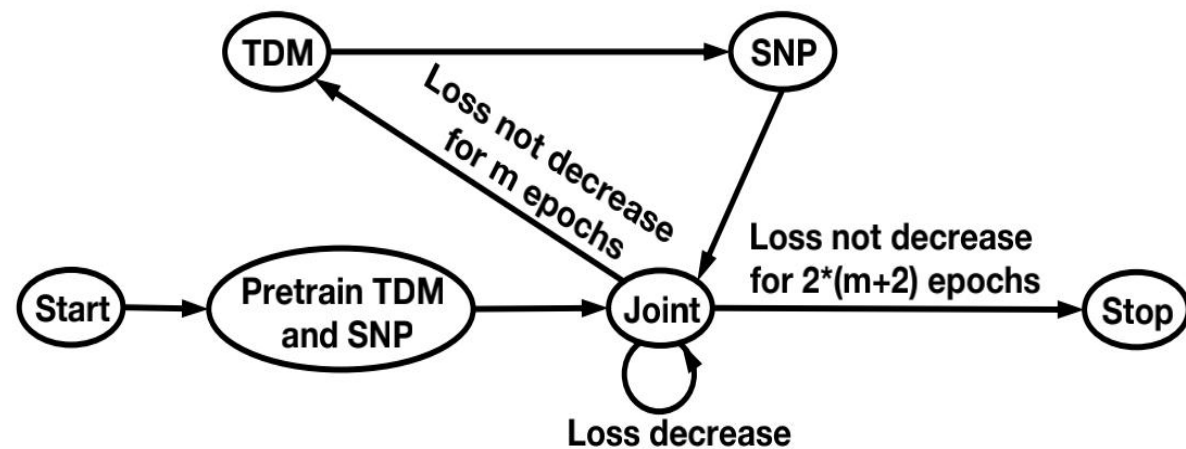


Figure 4: State transition diagram of training process.



Experiments

Table 2: Comparison results on Twitter and Reddit datasets (in %). Higher scores indicate better performance. The best results in each column are in bold. Our model gets significantly better scores than all other comparisons for all metrics ($p < 0.01$, paired t-test).

Models	Twitter				Reddit			
	AUC	Accuracy	Precision	F1	AUC	Accuracy	Precision	F1
Simple Baselines								
RANDOM	50.1	49.5	78.0	60.1	49.9	50.2	15.4	24.1
HISTORY	44.1	41.2	73.2	50.6	53.6	46.5	18.1	27.5
Comparisons								
SVM	51.5	56.3	75.7	74.2	54.3	50.1	18.9	29.1
BiLSTM	52.5	77.4	78.2	87.2	59.3	54.1	22.2	31.7
BERT	70.5	80.2	80.4	89.0	63.2	51.1	21.6	32.5
CONVNET	73.6	79.2	78.9	88.2	60.6	55.3	21.6	31.2
JECUH	75.2	80.1	80.3	88.4	60.7	57.6	22.6	31.9
Our Model	83.2	82.9	84.7	90.2	64.8	62.7	24.9	34.6



Experiments

Table 3: Comparison results with ablations(in %). Higher scores indicate better performance.

Models	Twitter		Reddit	
	Accuracy	F1	Accuracy	F1
w/o TOPIC INIT	79.5	88.1	60.5	32.7
w/o DISC CONCAT	80.8	88.3	61.4	33.2
w/o DISC ATT	81.1	88.6	60.2	33.5
Full Model	82.9	90.2	62.7	34.6



Experiments

Table 4: C_v scores for Top 5 and 10 words of learned topics. The values range from 0.0 to 1.0, and higher scores indicate better topic coherence.

Models	Twitter		Reddit	
	5	10	5	10
LDA	0.498	0.393	0.483	0.377
NTM	0.499	0.425	0.492	0.397
Our	0.504	0.431	0.495	0.412



Experiments

Table 5: 5 sample latent discourse behaviors discovered from Reddit (The top 5 terms by likelihood are shown here). Names in the first column are our interpretation of the discourse behaviors according to the learned clusters. Discourse words indicating the behavior are highlighted in *blue and italic*.

Discourse	Top 5 representative terms
Disagreement	<i>but</i> have ask <i>different</i> see
Explanation	<i>because</i> stil when of that
Opinion	<i>think</i> my here ! <i>never</i>
Doubt	<i>n't</i> always want like <i>why</i>
Question	? For ! <i>where what</i>

Experiments

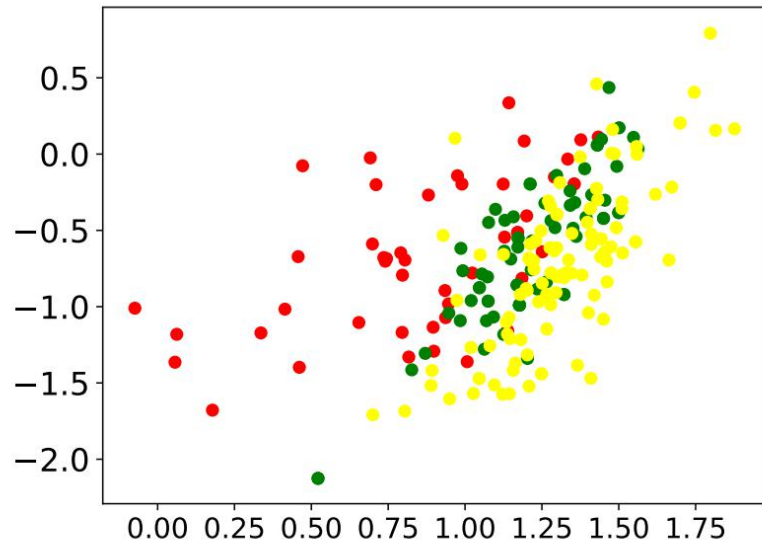
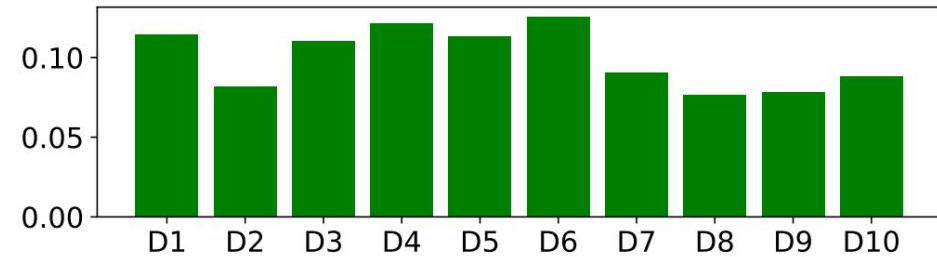
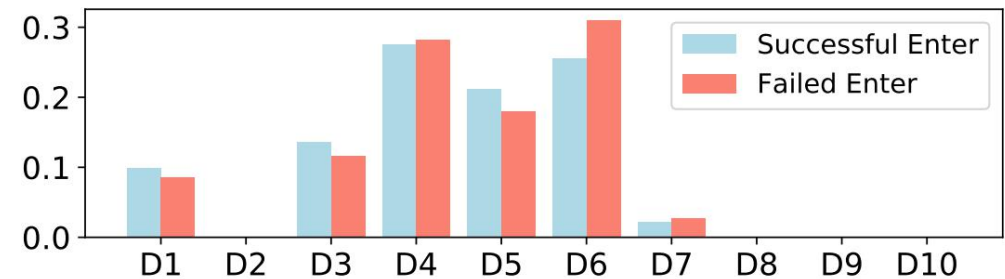


Figure 5: Topic mixture visualization (before softmax normalization) of three users' history conversations. A point refers to a conversation while different users are in distinguished colors.



(a) Disc distribution for newcomers



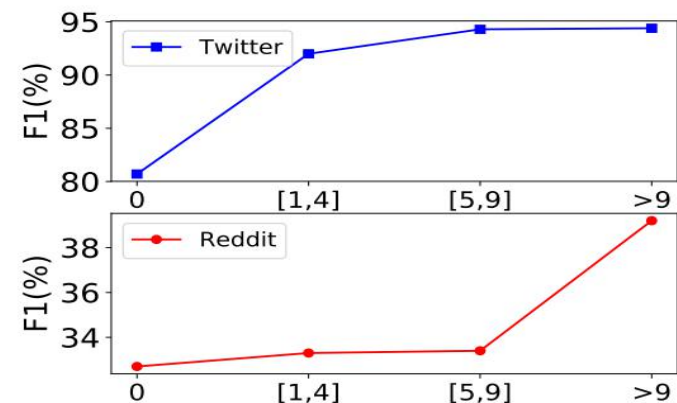
(b) Distribution for successful and failed new-entries

Figure 6: 6(a) is the distribution over discourse behaviors used in new-entries. 6(b) is the distribution of discourse behaviors for successful and failed new-entries. For both, X-axis: the 10 discourse behavior learned by our model; Y-axis: corresponding probability.

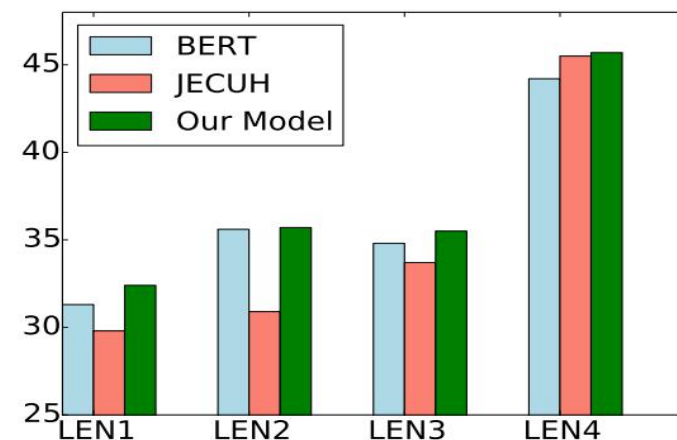
Experiments

Table 6: Human evaluation results (%). The overall inter-rater agreement achieved Krippendorff's α of 0.74, which indicates reliable results [30].

Models	Twitter				Reddit			
	OT	AQ	CL	CS	OT	AQ	CL	CS
SUCCESSFUL	100	28	0	54	98	48	14	38
FAILED	96	4	0	6	88	20	22	20



(a) F1 with Varying History #



(b) F1 with Varying Turn #

Figure 7: Y-axis: F1 score. In 7(a), X-axis: user history conversation numbers. In 7(b), LEN_i in X-axis: the i -th quantile by turn numbers (smaller i , shorter length).



Thanks !