

# How to Stop an Avalanche? JoDeM: Joint Decision Making through Compare and Contrast for Dialog State Tracking

Anonymous ACL submission

## Abstract

Dialog state tracking (DST) is a core component in task-oriented dialog systems. Existing state-of-the-art DST model incorporates insight and intuition from the human experience into design of supplementary labels, which greatly assisted the training process of turn-by-turn DST model. Though the turn-by-turn scheme and supplementary labels enabled satisfactory performance on the task, most of the DST models of this fashion label or process the raw dialogue data on the premise that the last turn dialogue state is always correct, which is usually not the case. In this paper, we address the negative impact resulted from the premise above as the avalanche phenomenon. After that, we propose JoDeM, a state-of-the-art DST model which can tackle the Avalanche phenomenon with two mechanisms. First mechanism is a jointly decision making method to extract key information from the dialogue. Second mechanism is a compare and contrast dialogue update technique to prevent error accumulation. Example study and graph analysis are presented to support our claim about the harmfulness of avalanche phenomenon. We also conduct quantitative and qualitative experiments on the high quality MultiWOZ2.3 corpus dataset to demonstrate that the proposed model not only outperforms the existing state-of-the-art methods, but also proves the validity of solving avalanche degradation problem.

## 1 Introduction

Goal-oriented dialog (GOD) systems, or Task-oriented dialogue (TOD) systems have recently attracted growing attention and significant progress has been made (Zhang et al., 2020; Neelakantan et al., 2019; Peng et al., 2020). Well-known commercial dialogue systems include the Apple Siri, Amazon Alexa, or Microsoft Cortana. In a complete GOD system, Dialog State Tracking (DST) serves as a cognitive and comprehending component, where it understands and extracts the user’s

goal in a well-constructed manner. The user’s goal is then provided to downstream for recommendation, booking, or other subsequent dialogue policy component to determine the system action and response. Hence, as the backbone of a dialogue system, it is crucial to have a DST module with exceptional performance to guarantee the base for the performance of subsequent components (Takanobu et al., 2020).

Since the blossom of the application of pre-trained language model, the accuracy of DST models has increased tremendously. Especially, turn-by-turn schematic DST models (Liao et al., 2021) with insightful design of auxiliary labels and data structure have dominated the field, where most of the among-the-best works are of this genre (Heck et al., 2020; Liao et al., 2020). However, this type of models all suffer from a major flaw, the avalanche phenomenon. The avalanche phenomenon is the result of wrong premise during the labeling process which will only occur in the DST models with turn-by-turn scheme.

As oppose to the trending turn-by-turn scheme, early multi-domain DST methods follow a dialog history scheme. Model of this scheme takes the whole or window-sized dialogue history as input. It predicts slot value without explicitly discriminating over turns of utterances. Despite the benefits of making prediction based on a more comprehensive and complete data at once, dialog history scheme has several drawbacks. The length of dialogues is often too long for pre-trained language model to process. More essentially, processing an entire dialogue at once violates the instant update nature of DST. Aligning with the need of instant update, turn-by-turn scheme was proposed (Kim et al., 2019). Models of this scheme input the dialogue state generated from the previous turn and the most recent turn utterance and output the updated dialogue state. The advantages of turn-by-turn scheme resulted in great performance boost,

084 in which most of the among-the-best works are of  
085 this scheme. On top of the choice of better scheme,  
086 to achieve a superior performance, these state-of-  
087 the-art DST models made the best use of auxiliary  
088 labels.

089 Basic input of turn-by-turn schematic DST mod-  
090 els are the current turn utterance and last turn dia-  
091 logue state, where the basic output is the updated  
092 current turn dialogue state. Using only basic out-  
093 put as golden training label inevitably leads to a  
094 sub-optimal result due to complexity of DST task.  
095 Mainstream DST systems typically incorporate sup-  
096 plemental labels to guide the model towards bet-  
097 ter performance. For example, Zhang et al. 2019;  
098 Heck et al. 2020 obtain key information from the  
099 dialogue span directly labels the starting and end-  
100 ing index of the key phrase for DST models to learn  
101 span detection.

102 However, the high utilization of supplementary  
103 labels in turn-by-turn schematic models have in-  
104 duced a new obstacle in developing a more robust  
105 and high quality DST system. In practice, a train-  
106 ing instance, which is a turn of dialog in an entire  
107 dialogue for turn-by-turn systems, are randomly  
108 shuffled along with instance from other dialogues.  
109 For convenience and effective training, supplement-  
110 ary labels are made under the assumption that the  
111 input previous turn dialogue state is correct. While  
112 in a considerable amount of cases, models have  
113 to make prediction under incorrect last turn dialog  
114 state. In those cases, the supplementary labels will  
115 also be incorrect themselves because they are also  
116 made under the false assumption. These facts add  
117 up to a poor robustness against noisy input, making  
118 the final performance way lower than expectation.  
119 To reflect this kind of characteristic where errors  
120 induce more errors, we name this phenomenon the  
121 **avalanche phenomenon**.

122 In this paper, we propose **JoDeM: Joint**  
123 **Decision Making DST** system with a compare  
124 and contrast mechanism. As mentioned, there are  
125 two major issues that directly contribute to the exis-  
126 tence of the avalanche phenomenon, incorrect last  
127 turn dialogue state and inflexible training labels.  
128 To address the former issue where DST models  
129 often perform worse when the input last turn dia-  
130 logue state is incorrect, we simply exempt dialogue  
131 state from the data flow of DST model, and strictly  
132 update it in a compare and contrast fashion. In  
133 other words, the extraction of key information is  
134 accomplished by a series of fluent back propagat-

135 able operations while the update process is not.  
136 To tackle the later issue, JoDeM deploys a joint  
137 decision making structure to successfully update  
138 dialogue state in a more robust and flexible manner  
139 despite the fact that training labels are fixed.

140 The JoDeM model contains eight modules  
141 that divides the whole DST process into three  
142 stages. The first stage contains a utterance encoder.  
143 The second stage contains four parallel modules,  
144 namely, a domain update, a slot gate, a slot type,  
145 and a span detection module. The third stage con-  
146 tains a dialogue state update module. As shown  
147 in the figure, first, we use BERT as the pre-trained  
148 language module to embed turn utterance. Then, a  
149 parallel decision making procedure is adopted by  
150 the four modules to extract key information from  
151 the embedded utterance. At last, the dialogue state  
152 update module designed to address the avalanche  
153 phenomenon is applied to output the updated dia-  
154 logue state.

155 After introducing related work and the details  
156 of JoDeM, we conduct multiple standard and cus-  
157 tomized evaluation and analysis in this paper to  
158 show that not only JoDeM achieved a state-of-the-  
159 art performance, but also the reason why it achieved  
160 such robustness against the avalanche phenomenon.  
161 In short, our contribution is twofold:

- 162 1. We bring up the attention to the avalanche  
163 phenomenon, a previous uncharted territory  
164 in dialogue state tracking task, and present  
165 quantitative evidence to show its existence and  
166 severity to the performance of DST systems.
- 167 2. We proposed a DST model to verify the fea-  
168 sibility of a solution to address the avalanche  
169 phenomenon, targeting straight to the roots of  
170 the phenomenon. After that, we performed  
171 quantitative and qualitative experiments to  
172 show the validity of our work and that our  
173 model has achieved a state-of-the-art perfor-  
174 mance on the qualified MultiWOZ2.3 dataset.

## 175 2 Related Work

176 Depending on the inputs, existing DST models are  
177 categorized to history-based and turn-by-turn based  
178 (Liao et al., 2021). The former scheme takes the  
179 whole or window-sized dialogue history as input to  
180 recurrent neural networks or networks (Goel et al.,  
181 2019; Gao et al., 2019). For example, HJST con-  
182 siders the full dialogue history using a hierarchical  
183 RNN (Gao et al., 2019; Serban et al., 2015). Works

such as Wu et al. 2019 treats the entire dialogue as a concatenated sequence while using Bi-LSTM or RNN as an encoder. There are also works inputting the whole history or window-sized dialogue history into BERT such as Lee et al. 2019a. In order to overcome the limitations of history-based scheme mentioned in the introduction, turn-by-turn DST systems was developed. Typically, model of this scheme takes the previous turn dialogue state and the current turn utterance as input to generate new dialogue state (Chao and Lane, 2019; Ren et al., 2019; Heck et al., 2020).

Basic label of the DST task is the correct dialogue state at each turn, which is often insufficient for the model to learn from effectively. The most common example of supplementary label is the starting index and ending index of the value phrase utilized in the span-based models (Zhang et al., 2019; Heck et al., 2020; Chen et al., 2020b). Kim et al. 2019, a turn-by-turn model designed a set of operation-based labels to guide the updating process of dialogue state. Heck et al. 2020 defined three copy strategy and labeled the original dialogue state tracking process with more refined information. These attempts have made significant result on the performance by incorporating human knowledge to the training process by applying supplementary labels. However, these labels are created under the assumption that the last dialogue state at every turn is flawless, while in reality it is usually not the case. The gap between ideal and reality creates a major drawback on the performance and robustness. In our JoDeM model, we not only design our supplementary label base on fine intuition, but also address the drawback resulted from the avalanche phenomenon.

### 3 JoDeM: Joint Decision Making through Compare and Contrast

The proposed JoDeM model in Figure 1 consists of eight components that are located in three different stages of the DST process. The first stage is the *Utterance Encoder* that encodes the basic inputs, i.e., system and user utterance into vector embedding. After that, the utterance embedding is sent to the second stage, which is the *Joint Decision Making* stage. In this stage, key information is extracted from the utterance embedding by the following component, *Domain Update*, *Slot Gate*, *Type Prediction*, *Span Detection* and *Co-ref Classification*. At the last stage, compare and contrast mechanism

is applied by the *Dialogue State Update* component to update the dialogue state according to the key information from the second stage and the previous turn dialog state.

Before formally getting into the detail of the JoDeM model, we first layout the necessary mathematical notations and proper definition for the DST problem. We define a complete dialogue as  $X = \{(S_1, U_1), \dots, (S_T, U_T)\}$ , which has  $T$  sets, or turns of system and user utterance that are in a sequential order. The dialogue states of an entire dialogue which is a set of dialogue state from all  $T$  turns is defined as  $DS = \{DS_1, \dots, DS_T\}$ , where  $DS_i$  is the dialogue state of the  $i$ th turn. Each turn’s dialogue state is a set which takes multiple triplets of format  $(domain, slot, value)$  as its elements. To complete a DST task is equivalent to the following statement: for any turn  $t$ , given the turn utterance  $(S_t, U_t)$  and the last turn dialogue state  $DS_{t-1}$  as input, we should output  $DS_t$ , which contains the correct set of triplets  $(domain, slot, value)$ .

#### 3.1 Utterance Encoder

Utterance encoder is the cornerstone of all NLP task including the DST task. At each turn  $t$ , we use the pre-trained BERT (Devlin et al., 2018) as the front-end encoder to encode the dialog utterance  $(S_t, U_t)$  as

$$\mathbf{R}_t = \text{BERT}([\text{CLS}] \oplus S_t \oplus [\text{SEP}] \oplus U_t), \quad (1)$$

where  $\mathbf{R}_t$  is the embedding of utterance from turn  $t$ .  $\oplus$  is the concatenation operator. Special token CLS is the starting token for BERT and SEP is the separation token separating system utterance  $S_t$  and user utterance  $U_t$ . The embedding of utterance can also be denoted as  $\mathbf{R}_t = [\mathbf{r}_t^{\text{CLS}}, \mathbf{r}_t^{S_1}, \mathbf{r}_t^{S_2}, \dots, \mathbf{r}_t^{\text{SEP}}, \mathbf{r}_t^{U_1}, \mathbf{r}_t^{U_2}, \dots]$ , where  $\mathbf{r}_t^{\text{CLS}}$  is the vector representation of the entire turn dialogue. The vector  $\mathbf{r}_t^i$  is the contextual representations for the  $i$ th token in the utterance. The dimension of the embedding is  $h$ , which is a hyper-parameter of BERT. Above sentence embedding is then utilized for joint decision making.

#### 3.2 Joint Decision Making

The intuition behind the *Joint Decision Making* stage is to break down and imitate the human reasoning process. Human beings complete the DST task by solving the triplets of  $(domain, slot, value)$  in a joint fashion, rather than solving the elements in a triplet in an order or individually. For example, one would not

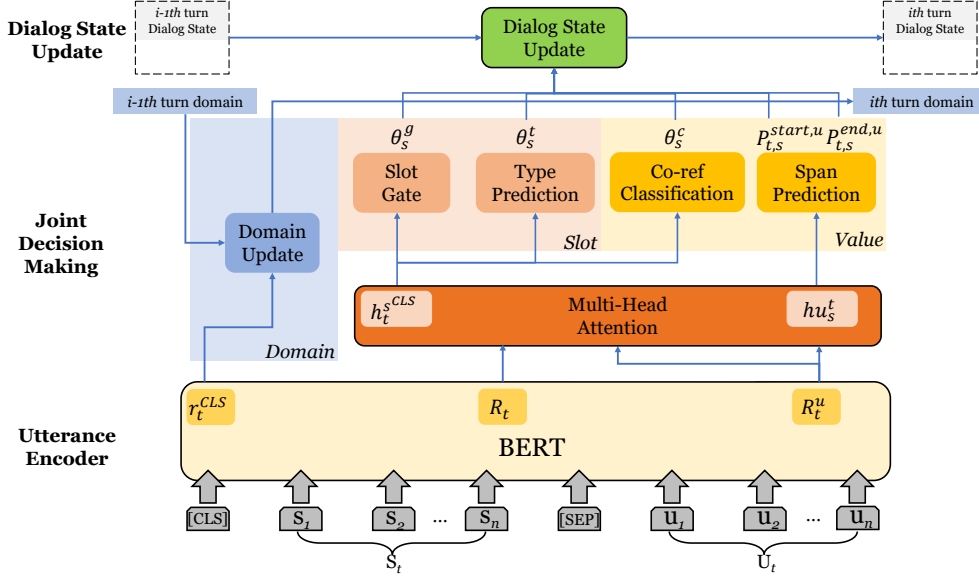


Figure 1: The architecture of the proposed JoDeM model comprised of three stages of eight components.

283 first determine the state of a  $(domain, slot)$  pair,  
 284 then search for its value. Instead, the context re-  
 285 garding different  $(domain, slot)$  pairs and their  
 286 possible values within the utterance are consid-  
 287 ered jointly so that comprehensive judgement on  
 288 the state of different  $(domain, slot, value)$  triplets  
 289 can be made. Bearing this intuition in mind, we  
 290 propose the *Joint Decision Making* stage consisting  
 291 of five parallel components that jointly solve all the  
 292  $(domain, slot, value)$  triplets in a dialogue state,  
 293 covering every possible scenario.

### 294 3.2.1 Domain Update

295 We obtain the domain of turn  $t$  by updating it from  
 296 the last turn  $t-1$  domain. As shown in the dialogue  
 297 example in Figure 2, the domain element of the dia-  
 298 log state is highly correlated to its last turn domain.  
 299 Generally, if the turn utterance doesn't contain any  
 300 trace of or sufficient domain information, the do-  
 301 main from the last turn will still be in use by the  
 302 continuity of the context. Therefore, we design the  
 303 *Domain Update* component to obtain the turn do-  
 304 main by taking the utterance representation  $\mathbf{r}_t^{CLS}$   
 305 as an input to detect new domain and the last turn  
 306 domain as a bias. The probability distribution of  
 307 the turn domain  $\mathbf{D}_t$  over all possible domains  $\mathbf{D} =$   
 308  $\{train, taxi, restaurant, hotel, attraction\}$  is  
 309 obtained by

$$\mathbf{D}_t = \text{softmax}(\gamma \cdot (\mathbf{W}^{DU} \cdot \mathbf{r}_t^{CLS} + \mathbf{b}^{DU})) \in \mathbb{R}^5, \quad (2)$$

310 where  $\mathbf{W}^{DU}$  and  $\mathbf{b}^{DU}$  are the trainable parameters  
 311 of a standard linear transformation, respectively.  
 312 Diagonal coefficient matrix  $\gamma = (\text{diag}(\mathbf{d}_{t-1}) + \mathbf{E})$   
 313

*Sys*: What area of town would you prefer?

*Usr*: I don't care about the location, but I would like to be in the moderate price range.

Figure 2: Example for the case when the turn domain is entirely dependent on the last turn context

314 where  $\mathbf{E}$  is the identity matrix,  $\mathbf{d}_{t-1}$  is the normal-  
 315 ized result from the last turn domain  $\mathbf{D}_{t-1}$ , and  
 316  $\text{diag}(\cdot)$  transforms vectors into diagonal matrices.  
 317 Due to the uniqueness of domain in each turn, the  
 318 class with the highest probability from  $\mathbf{D}_t$  is the  
 319 turn domain. The design of the impact of last turn  
 320 domain is oriented to the following purpose: we  
 321 require the impact from the last turn play a dom-  
 322 inate role when there's no new domain predicted.  
 323 At the same time, if there is new domain involved,  
 324 the influence of the last turn should be ignored. If  
 325  $\gamma = \text{diag}(\mathbf{d}_{t-1})$ , any newly discovered domain  
 326 would be covered up by the scaling effect from  
 327 the last turn domain. Also, in order to diminish  
 328 the impact from the last turn when new domain is  
 329 predicted, we make the bias itself relevant to the  
 330 outcome of the linear transformation. Only when  
 331 there is no domain discovered, i.e., the outcome of  
 332 the linear part is equally distributed, will the bias  
 333 of  $\text{diag}(\mathbf{d}_{t-1})$  dominate the result.

### 334 3.2.2 Slot Gate & Type Prediction

335 Our model is equipped with a *Slot Gate* and a *Type*  
 336 *Prediction* components for each slots. The *Slot*  
 337 *Gate* aims to determine whether a slot should be  
 338 updated, i.e., the output of a slot gate  $\mathbf{G}_s$  is a bi-

nary probability distribution. Inspired by Heck et al. 2020 and Kim et al. 2019, we summarize the possible updates into the following four types  $\{U, S, C, N\}$ .  $U$  and  $S$  indicates that the value of the slot should be found in the span of user utterance  $\mathbf{U}_t$  and system utterance  $\mathbf{S}_t$  respectively.  $C$  indicates that the value of the slot has a co-reference relationship with a certain  $(domain, slot)$  pair in the last turn dialogue state.  $N$  means that the user intend to delete the existing value of the corresponding slot in the dialogue state without providing any alternative value.

To make the above prediction for each slots, we first employ the multi-head attention mechanism (Vaswani et al., 2017) to calculate the attended context vector  $\mathbf{h}_t^s$  between  $\mathbf{R}_t$  and the user utterance embedding  $\mathbf{R}_t^u$  at  $t$  as

$$\mathbf{h}_t^s = \text{MultiHeadAtte}(\mathbf{Q}, \mathbf{K}, \mathbf{V}), \quad (3)$$

where  $\mathbf{Q}$  is the embedding of the entire utterance embedding,  $\mathbf{R}_t$ ,  $\mathbf{K}$  and  $\mathbf{V}$  are the embedding of the user utterance embedding  $\mathbf{R}_t^u = [\mathbf{r}_t^{U_1}, \mathbf{r}_t^{U_2}, \dots]$ . The reason to apply the multi-head attention mechanism is that the confirmation from a user is the essence of dialog state update, no matter the type of update. Therefore, the relationship between the entire utterance and the user utterance is needed.

After obtaining the attended embedding of the entire utterance, for each slot  $s$ , slot gates and type predictions are made by two parallel trainable linear layer classification,

$$\theta_s^g = \text{softmax}((\mathbf{W}_s^g \cdot \mathbf{h}_t^{sCLS} + \mathbf{b}_s^g)) \in \mathbb{R}^2, \quad (4)$$

$$\theta_s^v = \text{softmax}((\mathbf{W}_s^v \cdot \mathbf{h}_t^{sCLS} + \mathbf{b}_s^v)) \in \mathbb{R}^4. \quad (5)$$

### 3.2.3 Span Detection & Co-Ref Classification

*Span detection* and *co-ref classification* are equipped to solve the possible value for each slots.

*Span detection* is utilized for the slots whose values are found in the utterance. The attended utterance embedding is separated into two parts, the attended vector for user  $\mathbf{hu}_t^s$  and the attended vector for system  $\mathbf{hs}_t^s$ . A slot specific span detection layer performs a user/system specific span detection on the attended context vector  $\mathbf{hu}_t^s$  and system context vector  $\mathbf{hs}_t^s$  separately to obtain the span of potential values in the utterance to update. The expression of the process, using *span detection* on the user utterance as an example, is

$$[\alpha_{t,i}^{s,u}, \beta_{t,i}^{s,u}] = (\mathbf{W}_s^{span} \cdot \mathbf{hu}_{t,i}^s + \mathbf{b}_s^{span}) \in \mathbb{R}^2$$

$$P_{t,s}^{start,u} = \text{argmax}(\alpha_t^s)$$

$$P_{t,s}^{end,u} = \text{argmax}(\beta_t^s)$$

$i$  is the index of a token in the attended context of user utterance,  $P_{t,s}^{start,u}$  is the starting position of span in the user utterance  $\mathbf{U}_t$  for slot  $s$  in turn  $t$  and  $P_{t,s}^{end,u}$  is the corresponding ending position.

*Co-ref classification* is utilized for the slots whose value should be filled via co-referencing with a known value in the last turn dialog state. We simply take  $\mathbf{h}_t^{sCLS}$  which is the attended context embedding of the representation token for the entire utterance and perform a linear layer classification,

$$\theta_s^c = (\mathbf{W}_s^c \cdot \mathbf{h}_t^{sCLS} + \mathbf{b}_s^c) \in \mathbb{R}^{31}, \quad (6)$$

where the output  $\theta_s^c$  is a probability distribution on all possible thirty  $(domain, slot)$  pairs and one none class.

### 3.3 Dialogue State Update

*Dialogue State Update* is the key part of any turn-by-turn schematic DST systems, which is the procedure where the avalanche phenomenon originated from. We mentioned that the conflict between incorrect last turn dialogue state and the supplementary labels which are based on the correct last dialog state is the main contributor to the avalanche phenomenon. Therefore, we exclude the dialogue state updating process from the forward and backward propagation of data processing flow, by updating the dialogue state by carefully comparing and contrasting through the information that we obtained from the previous *Joint Decision Making* stage. At last, to achieve better robustness of the model, we apply a trick in the training process. The overall dialogue state update procedure is shown in Algorithm 1.

First, we specify the domain by the result of the *Domain Update* component  $\mathbf{D}_t$ . Second, we determine whether to update a slot within the domain through the *Slot Gate* result  $\theta_s^g$ . If it equals to 1, that is,  $\theta_s^g = 1$ , we move on to the next step. In the third step, we go through the slots with  $\theta_s^g = 1$  and determine their corresponding values according to their *Type Prediction*  $\theta_s^v$ . For the slots whose  $\theta_s^v = U$  or  $\theta_s^v = S$ , we obtain their values by getting the corresponding span from the user or system utterance. The span is determined by the corresponding starting and ending index  $P_s^{start,u}, P_s^{end,u}, P_s^{start,s}, P_s^{end,s}$ . If  $\theta_s^v = C$ , the

---

**Algorithm 1: DS Update**

---

**Input:**  $\theta_s^g, \theta_s^v, \theta_s^c$ ,  
 $P_s^{start,u}, P_s^{end,u}, P_s^{start,s}, P_s^{end,s}$ ,  
 $D_t, DS_{t-1}$   
**Output:**  $DS_t$

- 1 Specify the turn Domain via ( $D_t$ )
- 2 **for** each slots  $s$  in the turn Domain **do**
- 3     **if**  $\theta_s^g$  **then**
- 4         **if**  $\theta_s^v = U$  **then**
- 5              $v \leftarrow U_t[P_s^{start,u}:P_s^{end,u}]$
- 6         **else if**  $\theta_s^v = S$  **then**
- 7              $v \leftarrow S_t[P_s^{start,s}:P_s^{end,s}]$
- 8         **end**
- 9         **else if**  $\theta_s^v = C$  **then**
- 10              $v \leftarrow DS_{t-1}[\theta_s^c]$
- 11         **end**
- 12         **else if**  $\theta_s^v = N$  **then**
- 13              $v \leftarrow \text{none}$
- 14         **end**
- 15     **end**
- 16     **if**  $DS_t\{D_t, s\} \neq v$  **then**
- 17          $DS_t\{D_t, s\} \leftarrow v$
- 18     **else**
- 19         **if** Training **then**
- 20              $\theta_s^g, \theta_s^v, \theta_s^c$ ,  
            $P_s^{start,u}, P_s^{end,u}, P_s^{start,s}, P_s^{end,s}$   
            $\leftarrow \text{GoldenLabel}$
- 21         **end**
- 22     **end**
- 23 **end**

---

436 values of the slots will be determined by the co-  
437 referred (*domain, slot*) pairs  $\theta_s^c$  from the last turn  
438 dialogue state. At last, for slots with  $\theta_s^v = N$ ,  
439 we simply delete the values that were stored pre-  
440 viously. Finally, in the last step, we perform the  
441 update by comparing and contrasting new triplets  
442 (*domain, slot, value*) and the ones in the last dia-  
443 logue state.

444 As mentioned above, we perform a special op-  
445 eration at this stage during the training process.  
446 During training, if the potential value is equal to  
447 the last turn dialogue state, we set all the output  
448 from the forward propagation to the golden label.  
449 Thus, preventing the back propagation process to  
450 alter the trainable parameters in the model. This  
451 operation can enable the model to develop the abil-  
452 ity to self-correct, resulting in a better performance.  
453 More details can be found in the *example study* in  
454 the appendix.

## 4 Experiment 455

### 4.1 Dataset 456

457 We evaluate our model on the public dataset: Mul-  
458 tiWOZ2.3, which is a fully-labeled task-oriented  
459 corpora comprised of human-human written con-  
460 versation. It contains 8439 multi-turn dialogues  
461 with dialogue having 6.84 turns on average. The  
462 difference between the MultiWOZ2.3 dataset and  
463 the previous versions of MultiWOZ dataset is  
464 that MultiWOZ2.3 has a cleaner and more accu-  
465 rate annotation as opposed to the noisier annota-  
466 tion of the previous MultiWOZ versions (Zhou  
467 and Small, 2019a; Han et al., 2020; Zang et al.,  
468 2020). Following previous work, only five domains,  
469 (*restaurant, hotel, attraction, taxi, train*) are  
470 employed in our experiments.

### 4.2 Training Configuration 471

472 We use the pre-trained BERT-based-uncased model  
473 as the utterance encoder in our model, which has  
474 12 hidden layers with 768 units. The limitation of  
475 the maximum sequence length isn't problematic,  
476 therefore setting length  $l = 256$  would suffice.

477 In our experiments, Adam optimizer is utilized,  
478 whose learning rate linearly decreases from  $5e - 5$ .  
479 We have trained the model with 25 epochs.

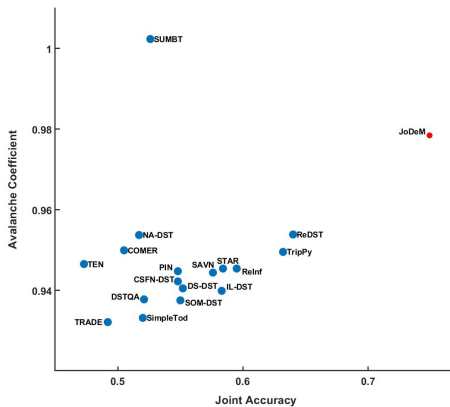
### 4.3 DST result 480

481 Both standard metrics and customized evaluation  
482 are carried out to compare our model and the  
483 state-of-the-art models. Standard metrics include  
484 Joint accuracy and Domain-Slot accuracy. Joint  
485 accuracy is the accuracy of the prediction of di-  
486 alogue states. It requires that all of the thirty  
487 (*domain, slot, value*) triplets in the dialogue state  
488 to be predicted and updated correctly. Only when  
489 the turn output Dialogue State is completely correct  
490 will  $JA = 1$ . In other cases,  $JA = 0$ , which is  
491 likely to happen when the input last turn dialogue  
492 state is wrong in the first place because models  
493 of turn-by-turn scheme typically can't self-correct.  
494 Domain-Slot accuracy is the accuracy of all the la-  
495 bels for each *Domain-Slot* pair in a turn. In the case  
496 of the JoDeM model, labels of a *Domain-Slot* pair  
497 includes the turn Domain,  $D_t$ , the slot gate for the  
498 slot  $\theta_s^g$ , the type prediction of the slot  $\theta_s^v$ , the co-  
499 ref classification of the slot  $\theta_s^c$ , and all the index of  
500 span detection of the slot  $P_s^{start,u}, P_s^{end,u}, P_s^{start,s}$ ,  
501 and  $P_s^{end,s}$ . There are thirty *Domain-Slot* pairs in  
502 total. It's apparent that  $JA$  is a much demanding

503 criterion to achieve and is also the most crucial  
 504 metric to evaluate a dialogue state tracking system.

505 We make a thorough comparison over our model  
 506 with the following state-of-the-art models from  
 507 both schemes including TRADE (Wu et al., 2019),  
 508 DS-DST (Zhang et al., 2019), IL-DST (Zhang  
 509 et al., 2021), SUMBT (Lee et al., 2019a), PIN  
 510 (Chen et al., 2020b), SOM-DST (Kim et al., 2019),  
 511 COMER (Ren et al., 2019), DSTQA (Zhou and  
 512 Small, 2019b), NA-DST (Le et al., 2020), TEN  
 513 (Chen et al., 2020a), ReDST (Liao et al., 2020),  
 514 ReInf (Liao et al., 2021), CSFN-DST (Zhu et al.,  
 515 2020a), SAVN (Wang et al., 2020b), TripPy (Heck  
 516 et al., 2020), SimpleTod (Hosseini-Asl et al., 2020),  
 517 and STAR (Ye et al., 2021).

518 The first two columns of Table 1 are the results  
 519 of standard metrics. The turn-by-turn schematic  
 520 DST models have shown significant performance  
 521 improvement over the dialog-history scheme in  
 522 both Joint accuracy and Domain-Slot accuracy. By  
 523 enhancing the accuracy at the turn level, turn-by-  
 524 turn schematic DST models are able to gain a much  
 525 higher joint accuracy at the end. Our model, the  
 526 JoDeM DST model, despite having a Domain-Slot  
 527 accuracy among the best, has achieve a state-of-the-  
 528 art performance boost on the joint accuracy metric.  
 529 This indicates that our model has a high robustness  
 530 against the avalanche phenomenon, which resulted  
 531 in a better overall performance.



532 Figure 3: The correlation of joint accuracy and  
 533 avalanche coefficient of various DST systems

534 Customized evaluation is designed to better evalu-  
 535 ate and compare the robustness of different DST  
 536 systems against the avalanche phenomenon. For  
 537 quantification, we introduce a novel avalanche co-  
 538 efficient,  $\alpha$ , which is calculated as  $\alpha = \frac{\sqrt{\bar{p}_j}}{\bar{p}_{ds}}$ , where  
 $\bar{l}$ ,  $\bar{p}_j$  and  $\bar{p}_{ds}$  are the mean length of dialogues, Joint  
 accuracy and Domain-Slot accuracy respectively.

Model	J Acc	D-S Acc	A Coeff
TRADE	49.2	96.94	0.932
DS-DST	55.2	97.67	0.941
IL-DST	58.3	98.50	0.940
SUMBT	52.6	91.02	<b>1.0023</b>
PIN	54.8	97.13	0.945
SOM-DST	55.0	97.93	0.938
COMER	50.5	95.48	0.950
DSTQA	52.1	97.15	0.938
NA-DST	51.7	95.42	0.954
TEN	47.3	94.93	0.947
ReDST	64.0	98.36	0.954
ReInf	59.5	98.21	0.945
CSFN-DST	54.8	97.39	0.942
SAVN	57.6	97.86	0.944
TripPy	63.2	98.63	0.949
SimpleTod	52.0	97.60	0.933
STAR	58.4	97.95	0.945
JoDeM	<b>74.9</b>	98.07	<b>0.979</b>

539 Table 1: Joint accuracy, slot accuracy and avalanche  
 540 coefficient on the test sets of MultiWOZ2.3.

541 With fixed dialogues, the avalanche coefficient is  
 542 model relevant only, which means it is an intrinsic  
 543 parameter to DST systems.

544 From the definition, we deduce that higher the  
 545 avalanche coefficient, the less a model suffers from  
 546 the avalanche phenomenon. The avalanche coef-  
 547 ficient of a DST model equals 1 when the model  
 548 doesn't suffer from the avalanche phenomenon. As  
 549 shown in Figure 3, despite the poor joint accuracy  
 550 performance, dialogue history scheme based mod-  
 551 els has an avalanche coefficient higher than 1. Our  
 552 model, among with other turn-by-turn schematic  
 553 DST models, has an avalanche coefficient lower  
 554 than 1, but way closer to 1 than the current state-  
 555 of-the-art models, resulting in a much better over-  
 556 all Joint accuracy performance. This proves that  
 addressing the avalanche is crucial for obtaining  
 higher Joint accuracy in DST models.

#### 557 4.4 Component Analysis

558 In order to dig deeper into the black box of the Jo-  
 559 DeM model, we carry out detailed analysis to show  
 560 the sufficiency and necessity of different compo-  
 561 nents in the JoDeM model and how our design is  
 562 aligned with our intuition.

563 To examine the *Domain Update* component,  
 564 we conduct two sets of control experiments with  
 565 unique variation on the original *Domain Update*

Original JoDeM	Training with $\gamma = E$	Testing with $\gamma = E$
97.75	91.93	73.94

Table 2: Domain Accuracy Analysis with Different Settings of the JoDeM Model.

Original JoDeM	Variation One	Variation Two
74.9	67.3	41.3

Table 3: Joint accuracy comparison on the JoDeM model with different usage setting of the multi-head attention mechanism

component.

**Variation one:** We set the diagonal coefficient matrix in 2 to  $\gamma = E$  during the training process. This setting means that the component learns to obtain the turn domain without utilizing any last turn information.

**Variation two:** Similarly but different from the variation one, we set the diagonal coefficient matrix in 2 to  $\gamma = E$  only during the testing process. This setting means that the model is trained given the last turn domain but being denied that information while performing on the test set.

The results from the original JoDeM model and the two variation are presented in Table 2. The metric we investigate is domain accuracy, which is the accuracy of the prediction of the turn domain. As you can see, the first column, which is the original JoDeM model, has the highest domain accuracy. The second column corresponds to variation one, which is the one with  $\gamma = E$  during training process. We can see that although a model can predict the turn domain solely using the turn utterance information, but the performance is sub-par compared to the one with last turn domain. The third column is the one with  $\gamma = E$  during testing only, whose decline of the domain accuracy is massive. The significance of this set of control experiment is to demonstrate that the last turn domain plays a key role or is relied heavily in the prediction of the turn domain.

Next, we focus on the question which is the purpose of the extra multi-head attention layer before applying the *slot gate*, *type prediction*, *span detection* and *co-ref classification* components. The intuition behind utilizing multi-head attention layer between user utterance embedding and the entire dialogue embedding is that any update from the

dialogue state is based on the consent of user. For example, the system may recommend a piece of information about a restaurant, but whether that information should be inserted into the dialogue state is up to whether the user takes the advice. To fairly evaluate, we train two control JoDeM models under two variation respectively. The metric we investigate is the joint accuracy.

**Variation one:** Instead of attending the user utterance embedding to the entire turn utterance embedding, we apply two multi-head self-attention layers on user and system utterance separately. The purpose of this variation is to examine and explore exactly what kind of attended relationship is the crux to dialogue state tracking.

**Variation two:** We discard the multi-head attention layer entirely, the input sequence for the *slot gate*, *type prediction*, *span detection* and *co-ref classification* components is the direct embedding of the pre-trained BERT. The goal of this variation is to examine the necessity of applying attention mechanism in the first place.

The results are shown in Table 3. Apparently, applying an additional attention layer is not only necessary but also crucial for the performance for dialogue state tracking. This observation is consistent with respect to other previous analytical work on dialogue state tracking. Furthermore, applying a multi-head cross-attention layer has the edge over a self-attention layer. This indicates that learning the relationship between the user utterance and the whole utterance is important in dialogue state tracking, which aligns with our intuition and the interactive nature of dialogue itself.

## 5 Conclusion

We proposed a novel, robust DST model JoDeM to address the rarely discussed problem, the Avalanche phenomenon. We showed that the trending topnotch DST systems all suffer from the Avalanche phenomenon with quantitative results and evidence. By multiple control experiments, we demonstrated how the overall structure and different techniques served the performance and robustness of the JoDeM model. We achieved a state-of-the-art performance on Joint accuracy and the criterion we design for measuring the impact of the Avalanche phenomenon. Finally, through the success of JoDeM, we show that the Avalanche phenomenon is worth solving and that there is more potential in this perspective for the DST task.



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  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705

## References

Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.

Galen Andrew and Jianfeng Gao. 2007. Scalable training of L1-regularized log-linear models. In *Proceedings of the 24th International Conference on Machine Learning*, pages 33–40.

Vevake Balaraman and Bernardo Magnini. 2021. Domain-aware dialogue state tracker for multi-domain dialogue systems. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*.

Jie Cao and Yi Zhang. 2021. A comparative study on schema-guided dialogue state tracking.

Guan-Lin Chao and Ian R. Lane. 2019. BERT-DST: scalable end-to-end dialogue state tracking with bidirectional encoder representations from transformer. *CoRR*, abs/1907.03040.

Junfan Chen, Richong Zhang, Yongyi Mao, and Jie Xu. 2020a. Neural dialogue state tracking with temporally expressive networks. *CoRR*, abs/2009.07615.

Junfan Chen, Richong Zhang, Yongyi Mao, and Jie Xu. 2020b. Parallel interactive networks for multi-domain dialogue state generation. *CoRR*, abs/2009.07616.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi, Sanchit Agarwal, Shuyang Gao, and Dilek Hakkani-Tür. 2019. Multiwoz 2.1: Multi-domain dialogue state corrections and state tracking baselines. *CoRR*, abs/1907.01669.

Shuyang Gao, Abhishek Sethi, Sanchit Agarwal, Tagyoung Chung, and Dilek Hakkani-Tür. 2019. Dialog state tracking: A neural reading comprehension approach. *CoRR*, abs/1908.01946.

Rahul Goel, Shachi Paul, and Dilek Hakkani-Tür. 2019. Hyst: A hybrid approach for flexible and accurate dialogue state tracking. *CoRR*, abs/1907.00883.

Dan Gusfield. 1997. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, Cambridge, UK.

T. Han, X Liu, R. Takanobu, Y. Lian, C. Huang, W. Peng, and M. Huang. 2020. Multiwoz 2.3: A multi-domain task-oriented dataset enhanced with annotation corrections and co-reference annotation.

M. Heck, C Van Niekerk, N. Lubis, C. Geishausen, H. C. Lin, M. Moresi, and M. Gai. 2020. Trippy: A triple copy strategy for value independent neural dialog state tracking.

Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. A simple language model for task-oriented dialogue. *CoRR*, abs/2005.00796. 706  
707  
708  
709

Sungdong Kim, Sohee Yang, Gyuwan Kim, and Sangwoo Lee. 2019. Efficient dialogue state tracking by selectively overwriting memory. *CoRR*, abs/1911.03906. 710  
711  
712  
713

Hung Le, Richard Socher, and Steven C. H. Hoi. 2020. Non-autoregressive dialog state tracking. *CoRR*, abs/2002.08024. 714  
715  
716

Chia-Hsuan Lee, Hao Cheng, and Mari Ostendorf. 2021. Dialogue state tracking with a language model using schema-driven prompting. *CoRR*, abs/2109.07506. 717  
718  
719

Hwaran Lee, Jinsik Lee, and Tae-Yoon Kim. 2019a. SUMBT: Slot-utterance matching for universal and scalable belief tracking. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5478–5483, Florence, Italy. Association for Computational Linguistics. 720  
721  
722  
723  
724  
725

Sungjin Lee, Qi Zhu, Ryuichi Takanobu, Zheng Zhang, Yaoqin Zhang, Xiang Li, Jinchao Li, Baolin Peng, Xiujun Li, Minlie Huang, and Jianfeng Gao. 2019b. ConvLab: Multi-domain end-to-end dialog system platform. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 64–69, Florence, Italy. Association for Computational Linguistics. 726  
727  
728  
729  
730  
731  
732  
733

S. Li, S. Yavuz, K. Hashimoto, J. Li, T. Niu, N. Rajani, X. Yan, Y. Zhou, and C. Xiong. 2020. Coco: Controllable counterfactuals for evaluating dialogue state trackers. 734  
735  
736  
737

Lizi Liao, Yunshan Ma, Wenqiang Lei, and Tat-Seng Chua. 2020. Rethinking dialogue state tracking with reasoning. *CoRR*, abs/2005.13129. 738  
739  
740

Lizi Liao, Tongyao Zhu, Le Long, and Tat Chua. 2021. Multi-domain dialogue state tracking with recursive inference. pages 2568–2577. 741  
742  
743

Shikib Mehri, Mihail Eric, and Dilek Hakkani-Tür. 2020. Dialogue: A natural language understanding benchmark for task-oriented dialogue. *CoRR*, abs/2009.13570. 744  
745  
746  
747

Arvind Neelakantan, Semih Yavuz, Sharan Narang, Vishaal Prasad, Ben Goodrich, Daniel Duckworth, Chinnadhurai Sankar, and Xifeng Yan. 2019. Neural assistant: Joint action prediction, response generation, and latent knowledge reasoning. *CoRR*, abs/1910.14613. 748  
749  
750  
751  
752  
753

Baolin Peng, Chunyuan Li, Jinchao Li, Shahin Shayandeh, Lars Liden, and Jianfeng Gao. 2020. SOLOIST: few-shot task-oriented dialog with a single pre-trained auto-regressive model. *CoRR*, abs/2005.05298. 754  
755  
756  
757  
758

759	Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015.	slot-value predictions on multi-domain dialog state tracking. <i>CoRR</i> , abs/1910.03544.	814
760	Yara parser: A fast and accurate dependency parser.		815
761	<i>Computing Research Repository</i> , arXiv:1503.06733.		
762	Version 2.		
763	Liliang Ren, Jianmo Ni, and Julian J. McAuley.	Ye Zhang, Yuan Cao, Mahdis Mahdieh, Jeffrey Zhao,	816
764	2019. Scalable and accurate dialogue state tracking	and Yonghui Wu. 2021. Improving longer-range	817
765	via hierarchical sequence generation. <i>CoRR</i> ,	dialogue state tracking. <i>CoRR</i> , abs/2103.00109.	818
766	abs/1909.00754.		
767	Iulian Vlad Serban, Alessandro Sordoni, Yoshua Ben-	Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. Task-	819
768	gio, Aaron C. Courville, and Joelle Pineau. 2015.	oriented dialog systems that consider multiple appro-	820
769	Hierarchical neural network generative models for	appropriate responses under the same context. <i>Proceed-</i>	821
770	movie dialogues. <i>CoRR</i> , abs/1507.04808.	<i>ings of the AAAI Conference on Artificial Intelligence</i> ,	822
		34:9604–9611.	823
771	Ryuichi Takanobu, Qi Zhu, Jinchao Li, Baolin Peng,	J. Zhao, M. Mahdieh, Y. Zhang, Y. Cao, and Y. Wu.	824
772	Jianfeng Gao, and Minlie Huang. 2020. Is your goal-	2021. Effective sequence-to-sequence dialogue state	825
773	oriented dialog model performing really well? em-	tracking.	826
774	pirical analysis of system-wise evaluation. In <i>Pro-</i>	L. Zhou and K. Small. 2019a. Multi-domain dialogue	827
775	<i>ceedings of the 21th Annual Meeting of the Special</i>	state tracking as dynamic knowledge graph enhanced	828
776	<i>Interest Group on Discourse and Dialogue</i> , pages	question answering.	829
777	297–310, 1st virtual meeting. Association for Com-	Li Zhou and Kevin Small. 2019b. Multi-domain dia-	830
778	putational Linguistics.	logue state tracking as dynamic knowledge graph en-	831
		hanced question answering. <i>CoRR</i> , abs/1911.06192.	832
779	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	Su Zhu, Jieyu Li, Lu Chen, and Kai Yu. 2020a. Ef-	833
780	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	ficient context and schema fusion networks for	834
781	Kaiser, and Illia Polosukhin. 2017. Attention is all	multi-domain dialogue state tracking. <i>CoRR</i> ,	835
782	you need. <i>CoRR</i> , abs/1706.03762.	abs/2004.03386.	836
783	Dingmin Wang, Chenghua Lin, Li Zhong, and Kam-Fai	Su Zhu, Jieyu Li, Lu Chen, and Kai Yu. 2020b. Effi-	837
784	Wong. 2020a. Dialogue state tracking with pretrained	cient context and schema fusion networks for multi-	838
785	encoder for multi-domain task-oriented dialogue sys-	domain dialogue state tracking. In <i>Findings of the</i>	839
786	tems. <i>CoRR</i> , abs/2004.10663.	<i>Association for Computational Linguistics: EMNLP</i>	840
787	Yexiang Wang, Yi Guo, and Siqi Zhu. 2020b. Slot	2020, pages 766–781, Online. Association for Com-	841
788	attention with value normalization for multi-domain	putational Linguistics.	842
789	dialogue state tracking. pages 3019–3028.		
790	Chien-Sheng Wu, Steven Chu-Hong Hoi, and Caiming	<b>A Appendix: Example Study</b>	843
791	Xiong. 2020. Improving limited labeled dia-		
792	logue state tracking with self-supervision. <i>CoRR</i> ,	To inspect the actual effect the JoDeM model have	844
793	abs/2010.13920.	on the update prediction of dialogue states, we	845
794	Chien-Sheng Wu, Andrea Madotto, Ehsan Hosseini-	provide two examples to demonstrate the strength	846
795	Asl, Caiming Xiong, Richard Socher, and Pascale	of the JoDeM model.	847
796	Fung. 2019. Transferable multi-domain state gen-	<b>A.1 Example One</b>	848
797	erator for task-oriented dialogue systems. <i>CoRR</i> ,		
798	abs/1905.08743.	The first example is presented in Figure 4. It not	849
799	Fanghua Ye, Jarana Manotumruksa, Qiang Zhang,	only serves as a demonstration of the actual oper-	850
800	Shenghui Li, and Emine Yilmaz. 2021. Slot	ation of JoDeM, but also can show the robustness	851
801	self-attentive dialogue state tracking. <i>CoRR</i> ,	of the joint decision making technique. First, the	852
802	abs/2101.09374.	domain of the turn is obtained, which is <i>Hotel</i> . Af-	853
803	Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara,	ter domain is specified, the updating procedure will	854
804	Raghav Gupta, Jianguo Zhang, and Jindong Chen.	strictly be limited in the domain. As shown in the	855
805	2020. MultiWOZ 2.2 : A dialogue dataset with	figure, after the domain is obtained, the focus shifts	856
806	additional annotation corrections and state tracking	to slot information. According to the <i>Slot Gate</i> ,	857
807	baselines. In <i>Proceedings of the 2nd Workshop on</i>	slots <i>Price range</i> , <i>Name</i> , <i>Area</i> is altered from the	858
808	<i>Natural Language Processing for Conversational AI</i> ,	context. After that, the value of the slot is extracted	859
809	pages 109–117, Online. Association for Computa-	from the utterance according to the <i>Type Prediction</i>	860
810	tional Linguistics.	and <i>Span Detection</i> . As you can see, although <i>Slot</i>	861
811	Jianguo Zhang, Kazuma Hashimoto, Chien-Sheng Wu,	<i>gate</i> and <i>Type Prediction</i> made a false judgment	862
812	Yao Wan, Philip S. Yu, Richard Socher, and Caiming	on <i>Area</i> , it didn’t lead to a wrongful update. The	863
813	Xiong. 2019. Find or classify? dual strategy for		

Robustness of Joint Decision and the SOP of JoDeM

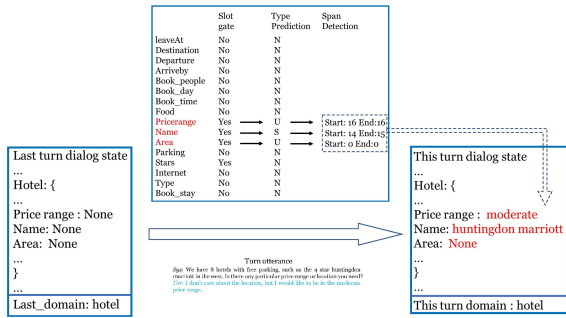


Figure 4: Example on Robustness of Joint Decision Making

reason for that is that the corresponding *Span Detection* detected that the starting and ending index are appointed to the *[CLS]* token, which means no information is detected. Only when all the components have made wrongful decision will they result in a wrongful update, which is the reason why *Joint Decision Making* is a robust way to extract information in a DST system.

## A.2 Example Two

The second example is presented in Figure 5. It shows that the JoDeM model can self-correct to a certain extent and why too many supplementary labels might be problematic. We focus on the *Destination* slot in the *Train* domain. As shown in the figure, the value of *Destination* is incorrect in the predicted last turn dialogue state. But it was rectified in this turn. If the predicted last turn dialogue state was correct, the correct operation at this turn is that *Slot gate* wouldn't have predicted the altering of the slot, which is aligned with the supplementary labels we have tagged. Therefore it would appear that the JoDeM model didn't get all the predictions right, but it enhanced the performance at the end. This ability of the JoDeM model takes credit from the trick we applied during the training process, which is setting the predicted values to the **GoldenLabel** when  $DSt\{D_t, s\} = v$ . Had the system follow the operation of the correct labels, it wouldn't be able to right the wrongs from the past turns.

## B Appendix: Responsible NLP Research Checklist

### B.1 Limitations and Risks

Although our work is evaluate on a public and high quality dataset, as we summarized in the ab-

Turn utterance  
 Sys: Absolutely. Where are you heading in from? What day?  
 User: I'll be leaving London kings cross and heading to Cambridge. I need to be there by 10:30 on Tuesday. Can you book this for 3 people? Reference please?

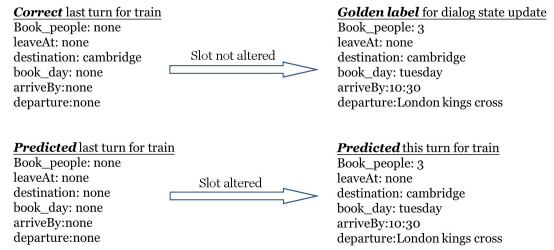


Figure 5: Example on self-correcting of JoDeM

stract and introduction, the dialogue state tracking task in real world application is far more complicated. Therefore there is both limitation and risks on whether our model can perform well in application.

### B.2 Use of scientific artifacts

The only scientific artifact our work applied is the dataset MultiWoz2.3 which is specifically designed for dialogue state tracking and publicly accessible. The content of the dataset doesn't contain any information that names or uniquely identifies individual people or offensive content. The dataset is about information regarding assorted places in Britain. The proportion of train/dev/test set is 8/1/1.

### B.3 Computational Experiments

In our experiment, 8 GPU is used to train our model, which has 222M parameters. One training epoch takes 21 minutes. Any setting of hyperparameter, including the existing package of pretrained language bert model, is presented in the experiment section. Our result, as well as the compared result from other works, is the mean of multiple independent identically distributed tests.