

BERT-DST: Scalable End-to-End Dialogue State Tracking with Bidirectional Encoder Representations from Transformer

(Guan-Lin Chao, Ian Lane)

presented by Alexey Korolev

16 April 2020

- 1 Introduction to DST
- 2 BERT
- 3 BERT-DST Parameters_sharing and dropout
- 4 Implementation detail
- 5 Dataset
- 6 Result

Introduction to DST

- Core component in today's task-oriented dialogue systems, maintains user's intentional states through the course of a dialogue.
- Our task fill slot(e.g. 'movie_name') with different target ('12 angry men') from user utterance.

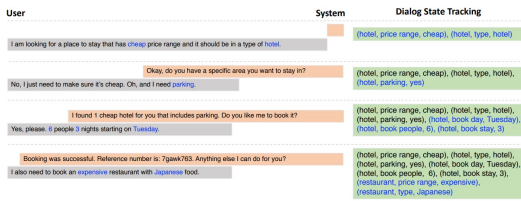
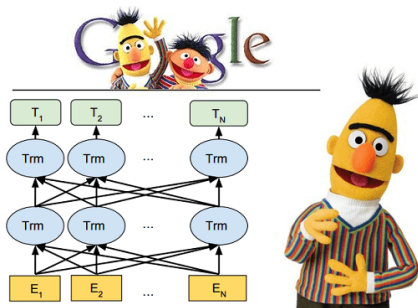


Figure 1: An example of dialog state tracking process for booking a hotel and reserving a restaurant. Each turn contains a user utterance (grey) and a system utterance (orange). The dialog state tracker (green) tracks all the $\langle domain, slot, value \rangle$ triplets until the current turn. Blue color denotes the new state appeared at that turn. Best viewed in color.

- We have 3 category: none, dontcare, span. Where 'none' denotes that a domain-slot pair is not mentioned at this turn, 'dontcare' implies that the user can accept any values for this slot 'span' represents that the slot should be processed by the model with a real value.

BERT

- BERT is a multi-layer bidirectional Transformer encoder , which is a stack of multiple identical layers each containing a multi-head self-attention and a fully-connected sub-layer with residual connections.
- Tasks: masked language modeling and next sentence prediction
- BooksCorpus and the English Wikipedia corpora.



BERT-DST

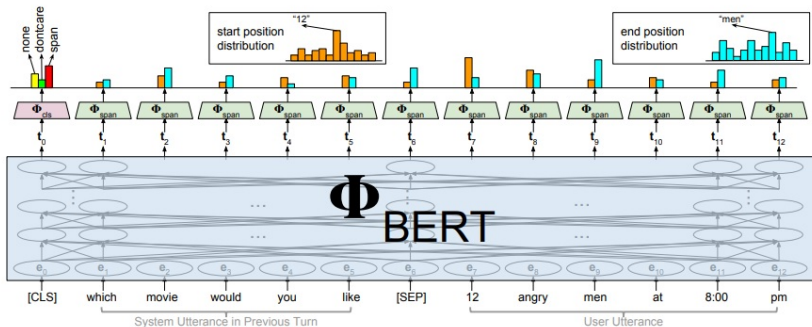


Figure 1: Architecture of the proposed BERT-DST framework. The diagram is color-coded such that modules with the same color share the same parameters. For each user turn, BERT-DST takes as input the recent dialogue context (system utterance in previous turn and the user utterance), and outputs turn-level dialogue state. BERT dialogue context encoding module Φ_{BERT} (blue) produces contextualized sentence-level and token-level representations of the dialogue context. The per-slot classification module Φ_{cls} (red) uses the sentence-level representation to generate a categorical distribution over three types of slot values {none, dontcare and span}. The per-slot span prediction module Φ_{span} (green) gathers the token-level representations and output the start and end positions (span) of the slot value. Note that the dialogue context encoding module Φ_{BERT} allows parameter sharing across all slots.

$$\begin{aligned} BERTinput(x_i) &= E_{tok}(x_i) + E_{seg}(i) + E_{pos}(i) = \\ &= \mathbf{e}_i, \end{aligned} \quad (1)$$

where $E_{tok}(x_i)$ — WordPiece embedding for x_i

$E_{seg}(i) \in \{e_{first}, e_{second}\}$ — segment embedding whose value is determined by whether the token belongs to the first or second sentence

$E_{pos}(i)$ — positional embedding for the i -th token.

$$\Phi_{BERT}([e_0, \dots, e_n]) = [t_0, \dots, t_n], \quad (2)$$

t_0 use for classification $[t_1, \dots, t_n]$ for span prediction.

1 Φ_{cls}

$$a^s = W_{cls}^s t_0 + b_{cls}^s \quad (3)$$

$$p^s = \text{softmax}(a^s) = \quad (4)$$

$$= [p_{none}^s, p_{doncare}^s, p_{span}^s]$$

$$slot_value^s = \text{argmax}(p_c^s) \quad (5)$$

2 Φ_{slot}

$$[\alpha^s, \beta^s] = W_{span}^s t_0 + b_{span}^s \quad (6)$$

$$p_\alpha^s = \text{softmax}(\alpha^s) \quad (7)$$

$$p_\beta^s = \text{softmax}(\beta^s) \quad (8)$$

$$start_pos^s = \text{argmax}_i(p_\alpha^s) \quad (9)$$

$$end_pos^s = \text{argmax}_i(p_\beta^s) \quad (10)$$

Parameters sharing and dropout

- BERT-DST SS(slot specific)
Model work only with one slot. For different slot we train different BERT-DST module.
- BERT-DST PS (parameters sharing)
We can apply parameter sharing in the dialogue context encoding module across all slots. So we reduce nuber of model parameters.



- To improve the robustness for unseen slot values, in the training phase, we replace each of the target slot value tokens by a special [UNK] token at a certain probability(30% dropout rate).
- $Loss_{total} = 0.8L_{cls}^{xent} + 0.1L_{span_start}^{xent} + 0.1L_{span_end}^{xent}$,
where L^{xent} denotes the cross entropy loss for the corresponding prediction target
- Update all layers in the model using ADAM optimization with an initial learning rate $2e^{-5}$ and early stopping on the validation set.

- Sim-M, Sim-R - automatically synthesize labeled datasets in the movie and restaurant domain.
- DSTC2 and WOZ 2.0 are standard benchmarks for task-oriented dialogue systems, which are both in the restaurant domain and share the same ontology. In DSTC2, automatic speech recognition (ASR) hypotheses of user utterances are provided to assess DST models' robustness against ASR errors, so we use the top ASR hypo

Datasets	# Dialogues (train, dev, test)	Slots
Sim-M	384, 120, 264	date, time, num_tickets, theatre_name, movie (5/5; 26/26)
Sim-R	1116, 349, 775	date, time, category, price_range, rating, num_people, location, meal, restaurant_name (5/19; 9/23)
DSTC2	1612, 506, 1117	area, price_range, food (1/73; 0/74)
WOZ 2.0	600, 200, 400	area (0/6; 1/7), price_range, food (1/65; 2/72)

Table 3: Dataset statistics. The number of dialogues is given for train, dev and test sets respectively. The slots containing OOV values are marked in bold. Parentheses represent (# unique OOV values in dev set / # unique values in dev set; # unique OOV values in test set / # unique values in test set).

- BERT-DST_PS gives less accuracy than BERT-DST_SS.
- Model doesn't give close to SOTA results and it's scalable.
- Slot value improves models.

DST Models	Sim-M	Sim-R
DST + LU Candidates [7]	50.4%	87.1%
DST + Oracle Candidates [†] [5]	96.8%	94.4%
BERT-DST_SS	71.6%	87.4%
+ slot value dropout	76.3%*	87.6%
BERT-DST_PS	72.3%	88.6%*
+ slot value dropout	80.1%*	89.6%*

Table 1: Comparison with prior approaches on Sim-M and Sim-R datasets (joint goal accuracy). * indicates statistically significant improvement over BERT-DST model (paired sample t-test; $p < 0.01$). [†] indicates the corresponding model should be considered as a kind of oracle because the candidates are ground truth slot-tagging labels, i.e. the targeted slot value is guaranteed to be in the candidate list and considered by DST.

DST Models	DSTC2	WOZ 2.0
DST + LU Candidates [7]	67.0%	-
DST + n-gram Candidates [8]	68.2±1.8%	-
DST + Oracle Candidates [5]	70.3%	-
Pointer Network [6]	72.1%	-
Delex.-Based Model [2]	69.1%	70.8%
Delex. + Semantic Dict. [2]	72.9%	83.7%
Neural Belief Tracker [2]	73.4%	84.2%
GLAD [3]	74.5±0.2%	88.1±0.4%
StateNet [4]	75.5%	88.9%
BERT-DST_PS	69.3±0.4%	87.7±1.1%

Table 2: Comparison with prior approaches on DSTC2 and WOZ 2.0 datasets (joint goal accuracy). We report the average and standard deviation of test set accuracy of 5 model runs with random training data shuffling and normal initialization on classification and span prediction weights.

Thank for attentions!

<https://github.com/guanlinchao/bert-dst>