Chongqing
University of
Technology

ATAI
Advanced Technique of
Artificial Intelligence

Artificial
Intelligence

# DIALKI: Knowledge Identification in Conversational Systems through Dialogue-Document Contextualization

Zeqiu Wu ♠∗   Bo-Ru Lu ♠∗   Hannaneh Hajishirzi ♠♦ Mari Ostendorf ♠

♠University of Washington   ♦Allen Institute for AI

{zeqiuwu1,roylu,hannaneh,ostendor}@washington.edu

Code : https://github.com/ellenmellon/DIALKI

——EMNLP 2021

**Reported by Yabo Yin**

1.Introduction

2.Method

3.Experiments

Chongqing
University of

ATAI
Advanced Technique
of Artificial
Intelligence

# Introduction

## Dialogue Context

[User]: Hi, can you tell me something about the private service bureau licenses?

...

[Agent]: Do you want to apply for a PSB?

[User]: No, I was being curious. Just in case, what should I do if I apply for PSB?

[Agent]: Your application will be reviewed in Albany's DMV. After that, it will be sent to your local DMV office and you'll be scheduled for an inspection.

## Grounding document

Department of Motor Vehicles

**[Sec 1] Private Service Bureau Licenses**

- A Private Service Bureau PSB license is required of …

**[Sec 3] How to Apply?**

- Request Name Approval
  - Before you can apply for a license to operate a PSB, …

**[Sec 30] After you apply**

- After your application is reviewed by the DMV in Albany, …
- When the inspector visits your location, …

Figure 1: In a document-grounded conversation, *knowledge identification* targets to locate a knowledge string within a long document to assist the agent in addressing the current user query.
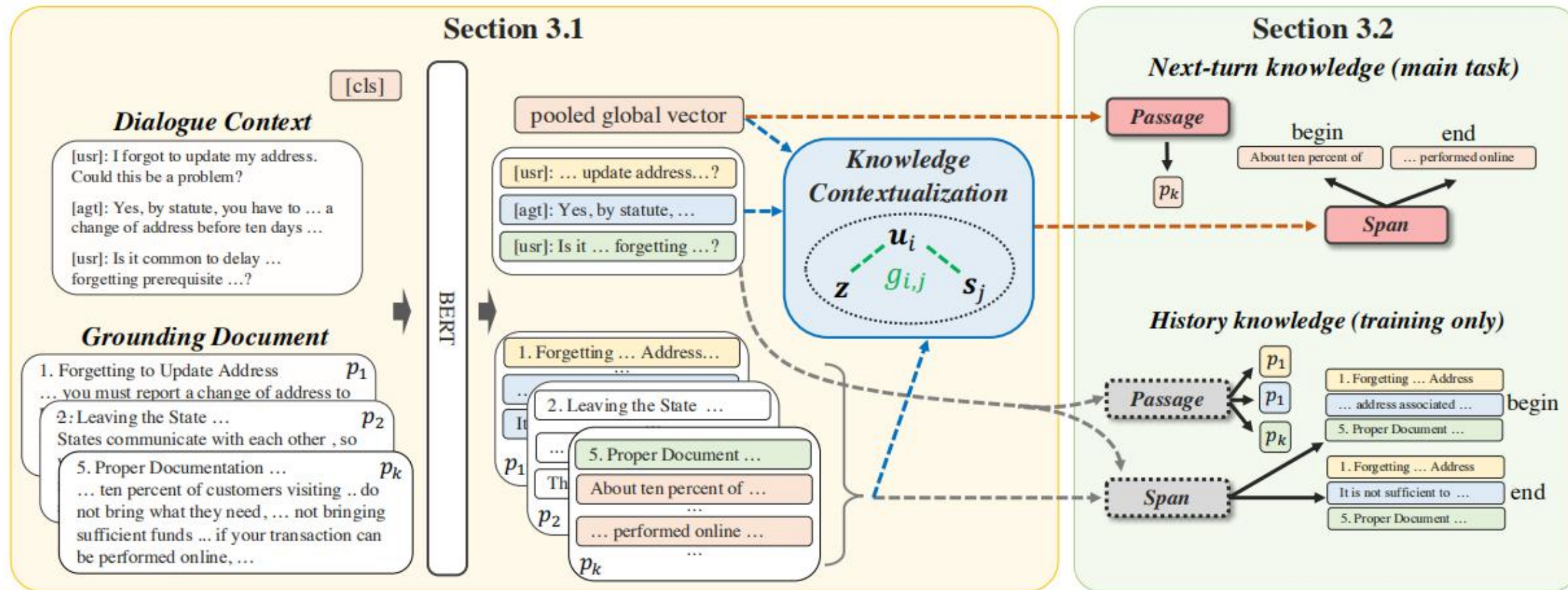
# Method



Figure 2: The overview of DIALKI. Each document is divided into passages. We apply BERT and a knowledge contextualization mechanism to obtain dialogue context and knowledge representations (left), for performing both next (main) and history (auxiliary) turn knowledge identification tasks (right). For each turn, DIALKI identifies knowledge by selecting the relevant passage as well as the begin/end spans in the passage.

# Method

## Problem Definition

dialogue context $(u_1, u_2, \ldots, u_n)$

grounding document $\mathcal{D} = \{p_1, p_2, \ldots, p_{|\mathcal{D}|}\}$

Each passage $p$ consists of a sequence of semantic units $p = (s_1, s_2, \ldots, s_l)$
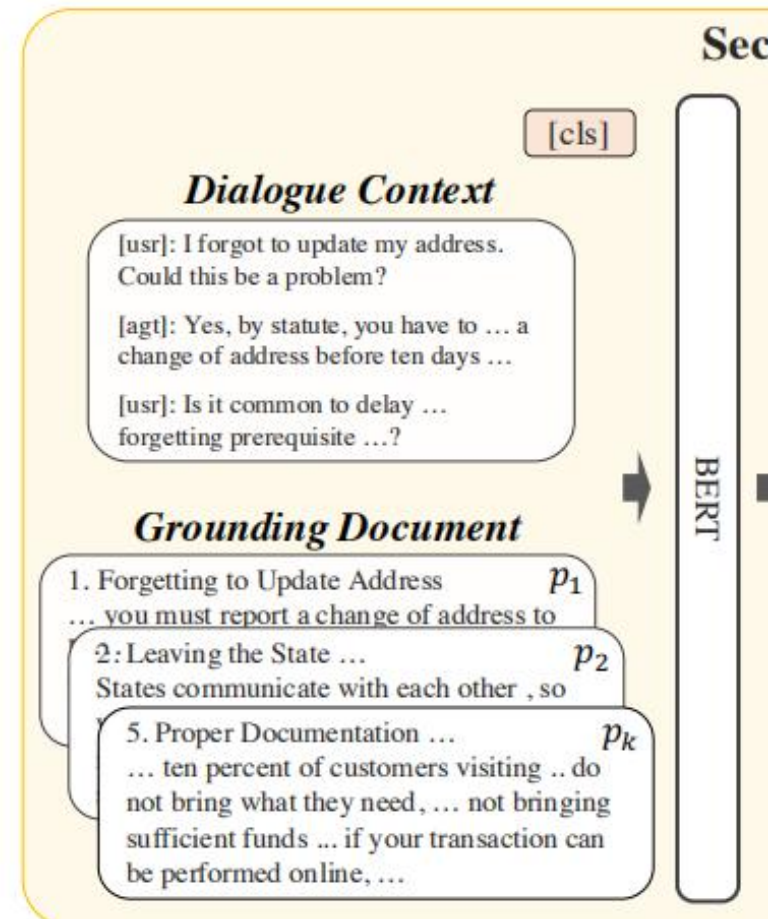
## Multi-Passage Encoding

$$\mathbf{X} = [\text{cls}][\text{usr}]\, u_1\, [\text{agt}]\, u_2 \cdots [\text{usr}]\, u_n$$
$$[\text{sep}]\, t\, [\text{cls}]\, s_1\, [\text{cls}]\, s_2 \cdots [\text{cls}]\, s_l\, [\text{sep}]$$

$$\mathbf{H} = \bar{G}(\text{BERT}(\mathbf{X}))$$
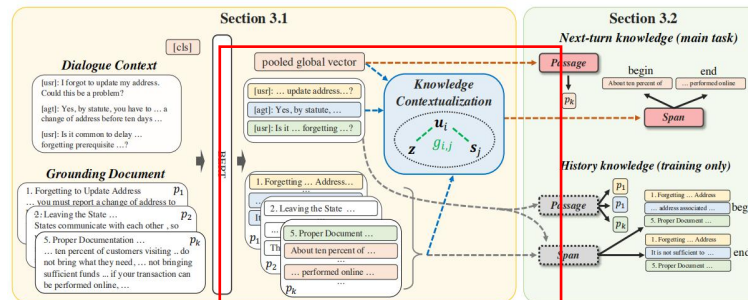
where $G(.)$ gathers vectors of all '[cls]', '[usr]' and '[agt]' tokens.

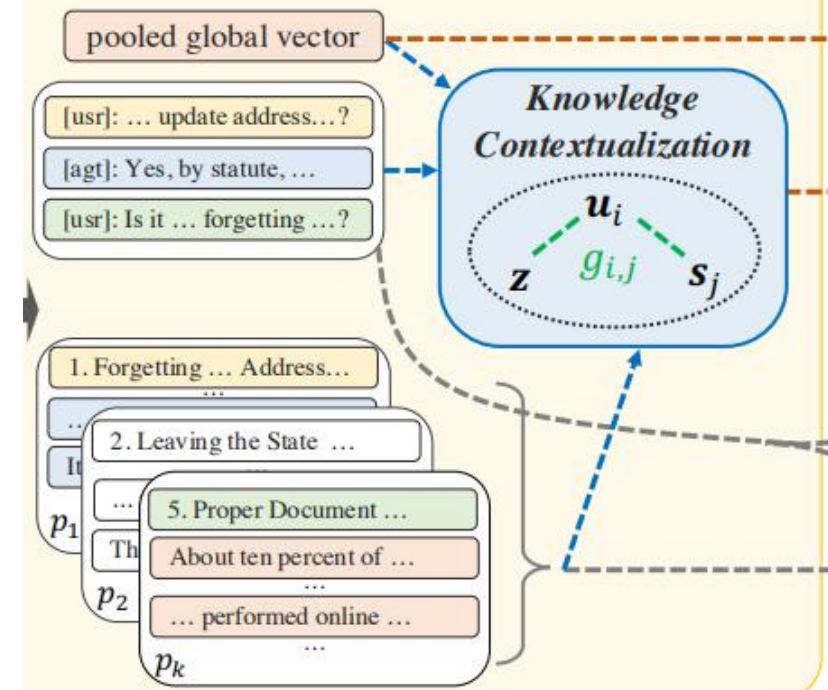$\mathbf{H}$ as $[\mathbf{z}, \mathbf{u}_1, \ldots, \mathbf{u}_n, \mathbf{s}_1, \ldots, \mathbf{s}_l]$

Sec

[cls]

**Dialogue Context**

[usr]: I forgot to update my address. Could this be a problem?

[agt]: Yes, by statute, you have to ... a change of address before ten days ...

[usr]: Is it common to delay ... forgetting prerequisite ...?

**Grounding Document**

1. Forgetting to Update Address     $p_1$
... you must report a change of address to

2: Leaving the State ...     $p_2$
States communicate with each other , so

5. Proper Documentation ...     $p_k$
... ten percent of customers visiting .. do not bring what they need, ... not bringing sufficient funds ... if your transaction can be performed online, ...

BERT

Chongqing
University of

ATAI
Advanced Technique
of Artificial
Intelligence

# Method



$$\mathbf{z}, \text{ pooled global,}$$

$$\mathbf{u}_i \text{ dialogue utterance } u_i$$

$$\mathbf{s}_j \text{ span } s_j$$

$$\mathbf{a}_{i,j} = \mathbf{W}_s \mathbf{s}_j + \mathbf{W}_z \mathbf{z} + \mathbf{W}_u \mathbf{u}_i, \ i \in \mathbf{C}_u$$

$$g_{i,j} = \sigma\left(\mathbf{u}_i^\top \mathbf{z} + \mathbf{u}_i^\top \mathbf{s}_j\right),$$

$$\widehat{\mathbf{s}}_j = \upsilon\left(\sum_{i \in \mathbf{C}_u}\left[\phi(\mathbf{a}_{i,j}) \odot g_{i,j}\right] + \mathbf{s}_j\right)$$
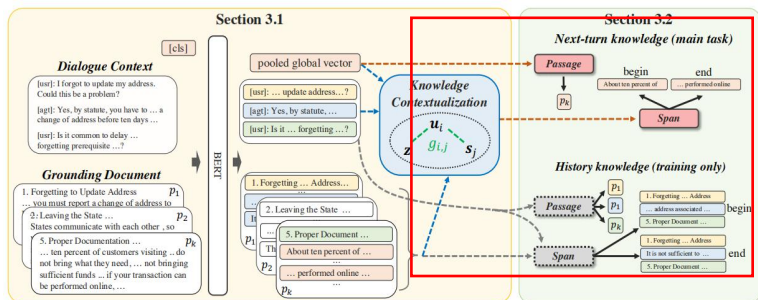
where $\mathbf{W}_s, \mathbf{W}_z, \mathbf{W}_u \in \mathbb{R}^{d \times d}$
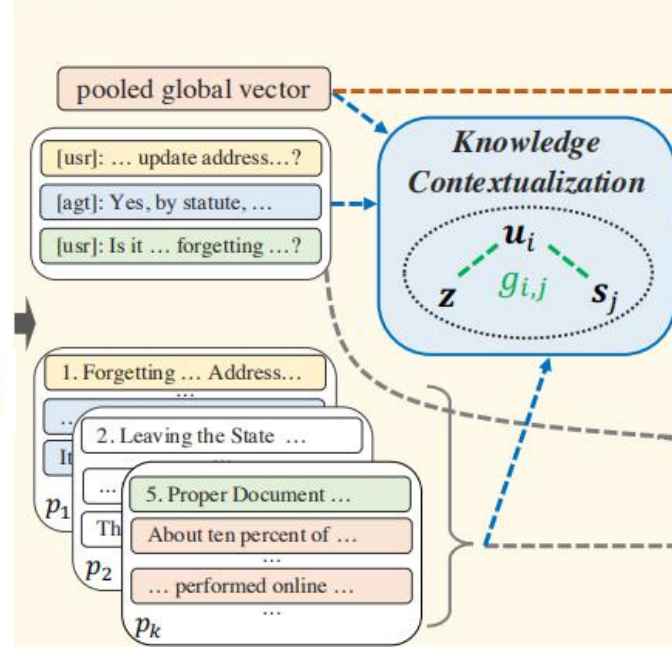
$\mathbf{C}_u$ indexes the most recent user turns.

$\widetilde{\mathbf{s}}_j$ with previous *agent* turns.

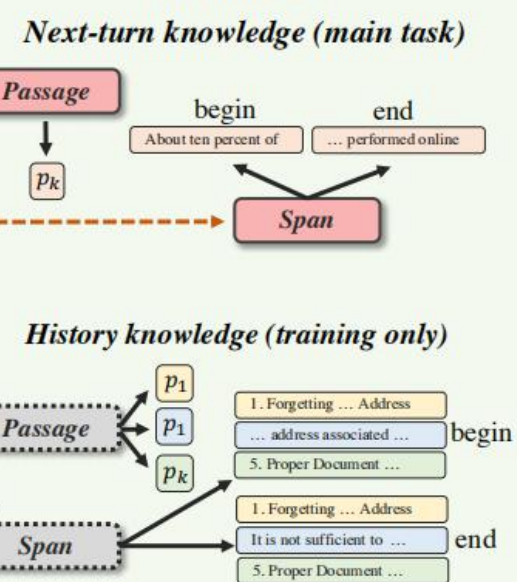$$\dot{\mathbf{s}}_j = [\mathbf{s}_j, \widehat{\mathbf{s}}_j, \widetilde{\mathbf{s}}_j]$$

# Method



$\mathbf{Z}$ matrix containing the pooled global vectors for all

$\mathbf{U}_i$ utterance representations for $u_i$ in all passages

$\mathbf{\dot{S}}$ all span representations

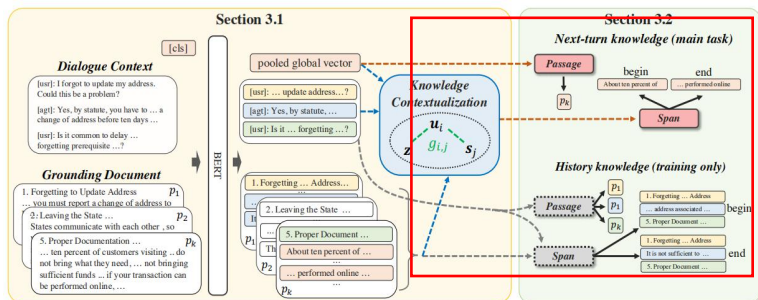$$\mathcal{L}_{\mathrm{psg}} = -\log q(\mathbf{W}_p \mathbf{Z})_{\hat{k}} \qquad (1)$$

$$\mathcal{L}_{\mathrm{begin}} = -\log q(\mathbf{W}_b \mathbf{\dot{S}})_{\hat{b}} \qquad (2)$$

$$\mathcal{L}_{\mathrm{end}} = -\log q(\mathbf{W}_e \mathbf{\dot{S}})_{\hat{e}} \qquad (3)$$

where $\mathbf{W}_p, \mathbf{W}_b, \mathbf{W}_e \in \mathbb{R}^d$.

$$\mathcal{L}_{\mathrm{next}} = \mathcal{L}_{\mathrm{psg}} + \mathcal{L}_{\mathrm{begin}} + \mathcal{L}_{\mathrm{end}}.$$

# Method



$U_i$ utterance representations for $u_i$ in all passages

$$\mathcal{L}_{\text{psg}}^h = \frac{1}{\|\mathbf{U}^*\|} \sum_{u_i \in \mathbf{U}^*} -\log q\left(\mathbf{W}_p^h \, \phi(\mathbf{W}^h \, \mathbf{U}_i)\right)_{k_i}$$

where $\mathbf{W}^h \in \mathbb{R}^{d \times d}, \mathbf{W}_p^h \in \mathbb{R}^d$,
$\mathbf{U}^*$ is the set of history turns that can find
their knowldge strings in the document $\mathcal{D}$. $k_i$ is the
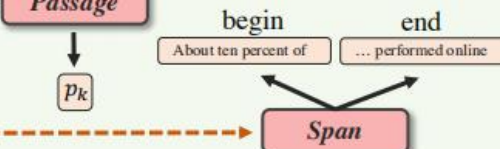gold passage index for turn $u_i$. $\phi$ is a non-linear

$$\mathcal{L}_{\text{hist}} = \mathcal{L}_{\text{psg}}^h + \mathcal{L}_{\text{begin}}^h + \mathcal{L}_{\text{end}}^h.$$

$$\mathcal{L}_{\text{adv}} = \max_{\|\epsilon\| \le a} \sum_{f \in \{f_{psg}, f_{begin}, f_{end}\}} \text{Div}\left(f(x) \| f(x + \epsilon)\right)$$

$$\mathcal{L} = \mathcal{L}_{\text{next}} + \alpha \mathcal{L}_{\text{hist}} + \beta \mathcal{L}_{\text{adv}} \tag{4}$$

Chongqing
University of

ATAI
Advanced Technique
of Artificial
Intelligence

# Experiments

| Method | Overall | |
|---|---|---|
| | EM | F1 |
| BERTQA-Token | 34.6 | 53.2 |
| BERTQA-Token (our version) | 35.8 | 52.6 |
| DIALKI ($\mathcal{L}_{next}$ only) | 51.2 | 64.7 |
| DIALKI | 59.5 | 71.0 |
| DIALKI (BERT-large) | **61.8** | **73.1** |

Table 1: Evaluation results on the Doc2Dial test set.

| Method | Seen | | Unseen | |
|---|---|---|---|---|
| | EM | F1 | EM | F1 |
| Transformer MemNet | 22.5 | 33.2 | 12.2 | 19.8 |
| Transformer MemNet + Pretrain | 24.5 | 36.4 | 23.7 | 35.8 |
| DiffKS (RNN) | 25.5 | – | 19.7 | – |
| SLKS (RNN) | 23.4 | – | 14.7 | – |
| SLKS (BERT-base) | 26.8 | – | 18.3 | – |
| Multi-Sentence (BERT-base) | 30.4 | 37.7 | 27.6 | 35.4 |
| DIALKI (BERT-base) | **32.9** | **40.7** | **35.5** | **43.4** |

Table 2: Evaluation results of WoW test sets.

Chongqing
University of

**ATAI**
Advanced Technique
of Artificial
Intelligence

# Experiments

| Method | Doc2Dial | | | | | | WoW | | | | | |
| | Overall | | Seen | | Unseen | | Overall | | Seen | | Unseen | |
| | EM | F1 | EM | F1 | EM | F1 | EM | F1 | EM | F1 | EM | F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BERTQA-Token | 42.2 | 58.1 | 48.3 | 61.1 | 37.0 | 55.6 | – | – | – | – | – | – |
| BERTQA-Span | 46.3 | 59.3 | 54.4 | 63.5 | 39.4 | 55.6 | – | – | – | – | – | – |
| Multi-Sentence | 59.5 | 68.8 | 63.6 | 71.6 | 56.0 | 66.4 | 29.2 | 37.0 | 32.4 | 39.7 | 26.1 | 34.3 |
| DIALKI ($\mathcal{L}_{next}$ only) | 60.4 | 71.2 | 64.2 | 72.3 | 57.1 | 70.2 | 31.5 | 39.7 | 33.3 | 41.1 | 29.8 | 38.3 |
| $+\mathcal{L}_{hist}$ | 63.0 | 72.6 | 66.5 | 73.9 | 59.9 | 71.9 | 33.6 | 41.6 | 35.1 | 42.7 | 32.2 | 40.5 |
| $+\mathcal{L}_{hist}$, know-ctx | 63.8 | 73.4 | **67.7** | 74.8 | 60.5 | 72.3 | 33.6 | 41.5 | 35.2 | 42.8 | 32.1 | 40.3 |
| $+\mathcal{L}_{adv}$ | 64.4 | 73.8 | 66.2 | 73.9 | 62.8 | 73.7 | 32.9 | 40.8 | 34.6 | 42.2 | 31.1 | 39.5 |
| $+\mathcal{L}_{hist}$, $\mathcal{L}_{adv}$, know-ctx | **65.9** | **74.8** | 67.6 | **74.9** | **64.4** | **74.7** | **34.2** | **42.1** | **35.9** | **43.5** | **32.6** | **40.7** |

Table 3: Ablation results on Doc2Dial and WoW dev sets.

Chongqing
University of

ATAI
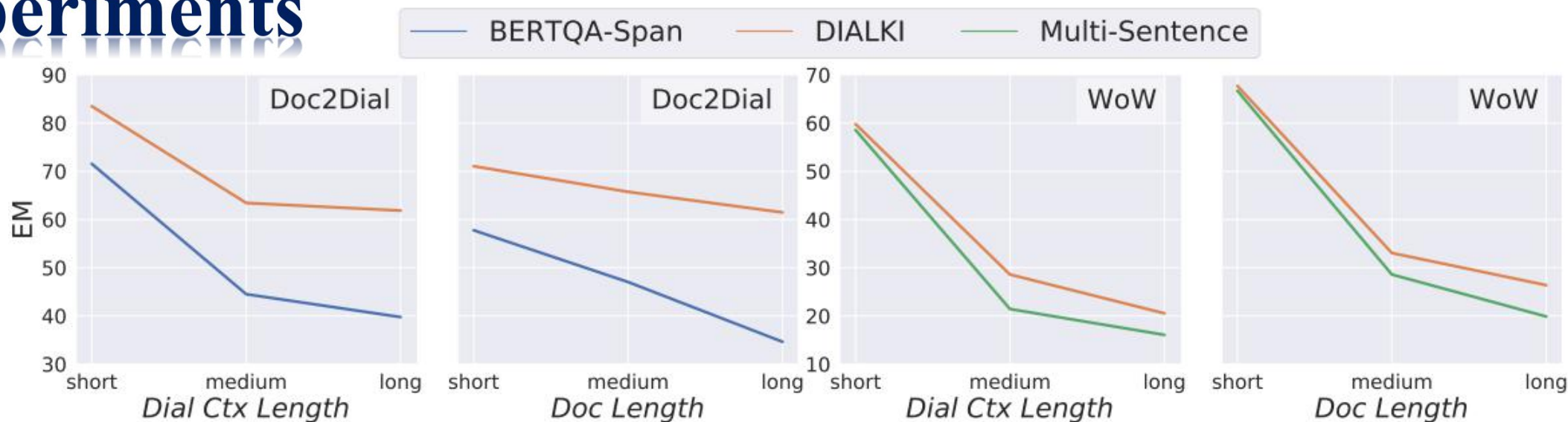Advanced Technique
of Artificial
Intelligence

# Experiments



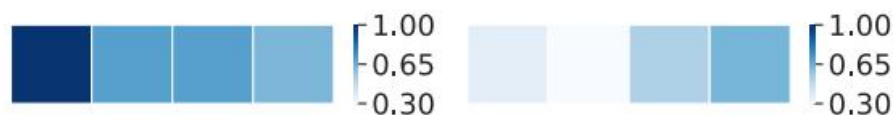Figure 4: EM versus the length of dialogue context (# previous turns) or document (# tokens).



Figure 3: Similarity between **z** and the latest 4 history turn representations (i.e., $\mathbf{u}_1 \ldots \mathbf{u}_4$ from left to right) on Doc2Dial (left) and WoW (right).

| KI Model | Knowledge Input | sacrebleu |
|---|---|---|
| – | full doc | 22.84 |
| BERTQA-Token | pred span | 21.42 |
| DIALKI | pred span | 25.16 |
| DIALKI | pred span & passage | **25.84** |

Table 4: Response generation results on Doc2Dial dev set. KI stands for Knowledge Identification.

# Experiments

| Dialogue Context | User: I want to trade in my license for a New York one. Agent: You have to exchange your-out-state driver license within 30 days of … User: What if I need my license for when I go back to my other country? … Agent: Are you studying in New York State? User: No | User: I heard something about co-op training program. Could you tell me about it, please? Agent: A co-op training program refers to … User: Sounds awesome. What should I do to get that? Agent: Well, for that, lets do a little fact check, shall we? Are you using VA educational assistance? User: Yes, I'm using that. |
|---|---|---|
| Baseline | Do you need a New York State driver license? | Good. In that case, you may be able to get money for books, tuition and housing. |
| Ours | Are you a driver from another country? | Are you enrolled at an approved institution of Higher Learning? |
| Gold Response | Are you a licensed driver from another country? | Good. Now, are you enrolled at an approved institution of Higher Learning or IHL? |

Table 5: Sample generated responses from BART with the full grounding document (baseline) or the predicted grounding span and passage by DIALKI (ours) as the additional input to the dialogue context.

| Method | Doc2Dial | | WoW | |
|---|---|---|---|---|
| | Seen | Unseen | Seen | Unseen |
| BERTQA-Span | 76.9 | 72.7 | – | – |
| Multi-Sentence | 85.3 | 81.6 | 68.0 | 57.8 |
| DIALKI ($\mathcal{L}_{next}$ only) | 86.6 | 84.4 | 72.9 | 69.0 |
| DIALKI | **88.5** | **87.5** | **73.4** | **69.7** |

Table 6: Passage prediction accuracy on dev sets.

Chongqing
University of

ATAI
Advanced Technique
of Artificial
Intelligence

# Thank you!