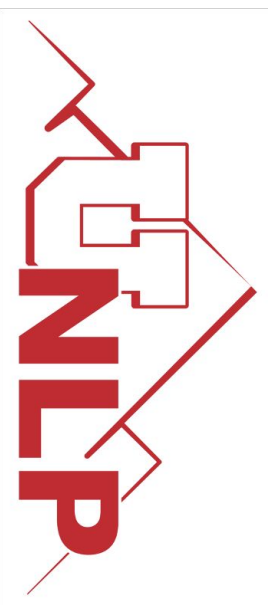


A Comparative Study on Schema-guided Dialog State Tracking

NAACL 2021

Jie Cao, Yi Zhang



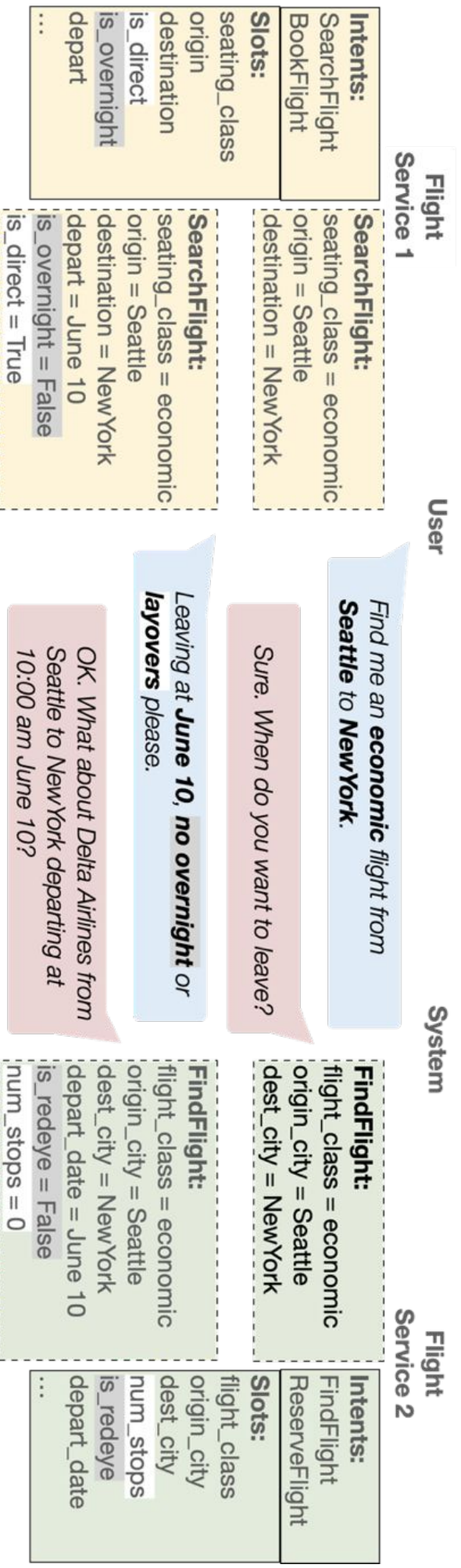
Outlines

- Motivation
- Task Description and Datasets
- Three Comparative Studies
 - i. Encoder Architectures
 - ii. Supplementary training
 - iii. Impact of Schema Description Style
- Q&A

Motivation

Challenges of Virtual Assistants (Task-oriented)

- Increasing number of **new services and APIs** → (new annotation, new model retraining)
- **Heterogeneous** Interfaces for **similar** services, precisely understanding overlapping functionalities.
- How to integrate **common sense and world knowledge**?



Schema-guided Dialogue State Tracking

Using **Natural language description** to explain the functionalities of tags, to help generalizing to **unseen tags** in **unseen domains**.

Flight Service 1

Adding Intent Description

```
{
  "name": "SearchFlight",
  "description": "Find a flight itinerary
between cities for a given date",
  "required_slots": .....
},
```

Adding Slot Description, Type, Possible values

```
{
  "name": "is_direct",
  "description": "Whether the flight
directly arrive without any stop",
  "is_categorical": True,
  "possible_values": ["true", "false"]
},
```

Flight Service 2

```
{
  "name": "FindFlight",
  "description": "Search for flights to a
destinaion",
  "required_slots": .....
},
```

```
{
  "name": "num_stops",
  "description": "Number of layovers
in the flight",
  "is_categorical": true,
  "possible_values": ["0", "1", "2"]
},
```

Schema-guided Dialogue State Tracking

Given **service schema** with description and dialogue history, predict **dialog states** after each user turn.

Schema

Services

```
"service_name": "Restaurants_1"
"description": "A leading provider for restaurant search and reservations"
```

Intents

```
"name": "FindRestaurants"
"description": "Find a restaurant of a
particular cuisine in a city"
"required_slots": ["cuisine", "city"],
"optional_slots": {
  "has_live_music": "dontcare",
  "serves_alcohol": "dontcare"
},
```

Slots

```
"name": "restaurant_name"
"description": "Name of the restaurant"
```

...

```
"name": "has_live_music"
"description": "Boolean flag indicating if the
restaurant has live music"
"is_categorical": true
"possible_values": ["True", "False"]
```

Schema-guided Dialogue State Tracking

Dialogue

User: I am looking for **SFO** food in **Asian**

System: I found 10 restaurants in **San Francisco**.

User: Any popular one? By the way, I also want to **buy a drink**.

System: There is a nice restaurant called Butterfly Restaurant.

User: Do they have **live music**? **Where are they located**?

System: They do not have live music. They are at 33 The Embarcadero.

Dialogue State Tracking

User Turn 1		User Turn 2		User Turn 3	
intent	FindRestaurants	intent	FindRestaurants	intent	FindRestaurants
req_slots	N/A	req_slots	N/A	req_slots	has_live_music, street_address
slot_values		slot_values		slot_values	
city	SFO	city	SFO, San Francisco	city	SFO, San Francisco
cuisine	Asian	cuisine	Asian	cuisine	Asian
		cuisine	alcohol	cuisine	alcohol
			TRUE		TRUE

Datasets

Datasets	Splits	Dialog	Domains	Services	Zero-shot Domains	Zero-shot Services	Function Overlap	Collecting Method
SG-DST	Train	16142	16	26	-	-	Across-domain Within-domain	M2M
	Dev	2482	16	17	1	8		
	Test	4201	18	21	3	11		
MULTWOZ 2.2	Train	9617	3	3	-	-	Across-domain	H2H
	Dev	2455	5	5	2	2		
	Test	2969	8	8	5	5		

SG-DST has more overlapping functionalities than MultiWOZ 2.2

Challenges: Three Comparative Studies

Q1: How to encode the dialog and schema?

- For each turn, matching the same dialog history with all schema descriptions **multiple times**.
- **Sentence-pair(SNL)** and **Token-level classification(QA)**

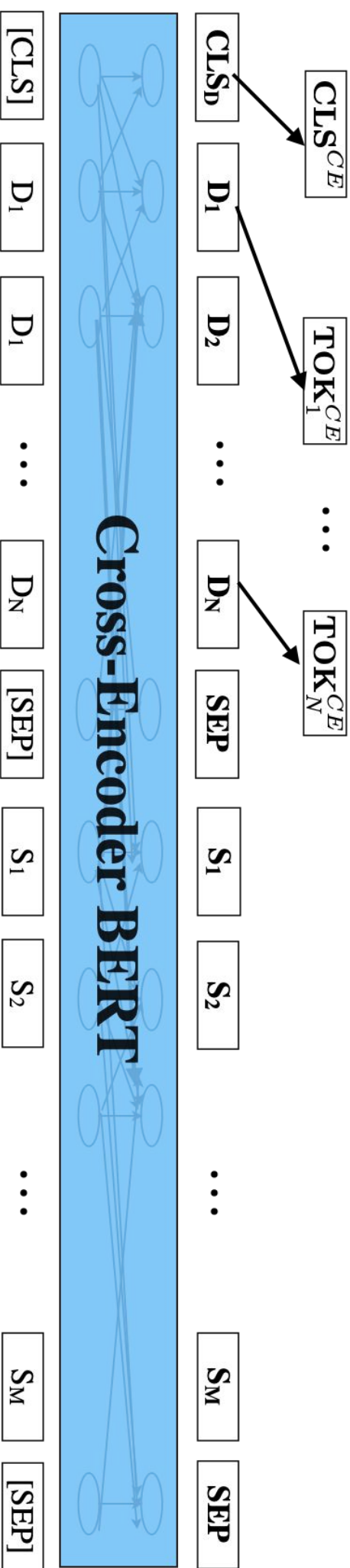
Q2: How do different supplementary trainings help?

- **Zero-shot learning** for unseen services

Q3: How the model performs on various description styles?

- Unseen service may have heterogeneous styles.

Q1: How to encode the dialog and schema? (Cross-Encoder)



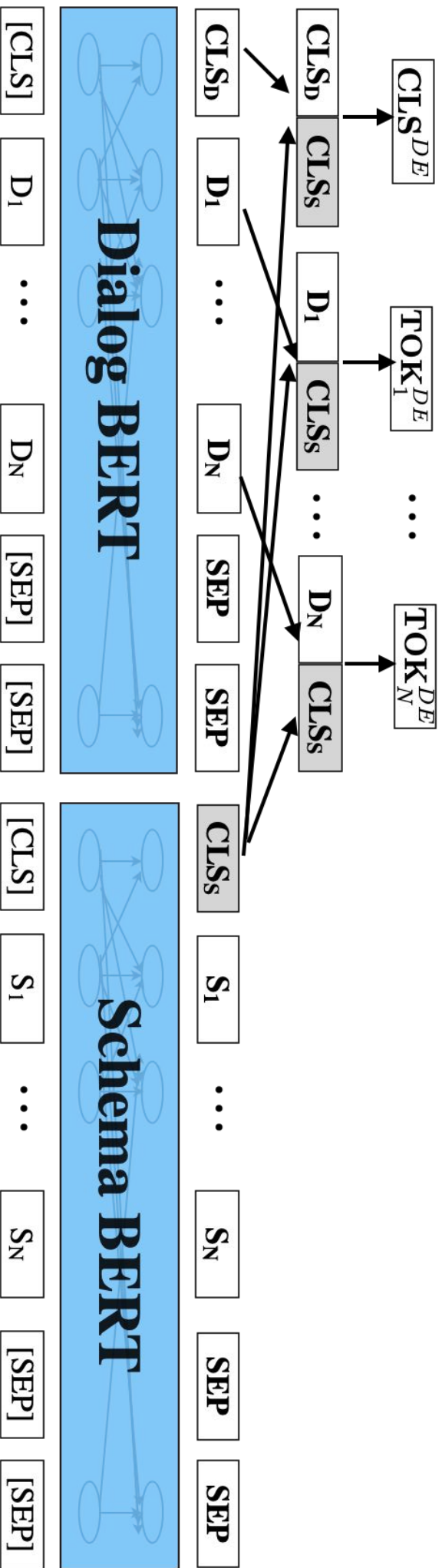
Pros:

Accurate, each representation is contextualized via full-attention

Cons: a lot of recomputing, slow

- a. Dialog encoded **multiple times** within the same turn
- b. Schema encoded **multiple times** across different turns.

Q1: How should the dialog and schema be encoded? (Dual-Encoder)



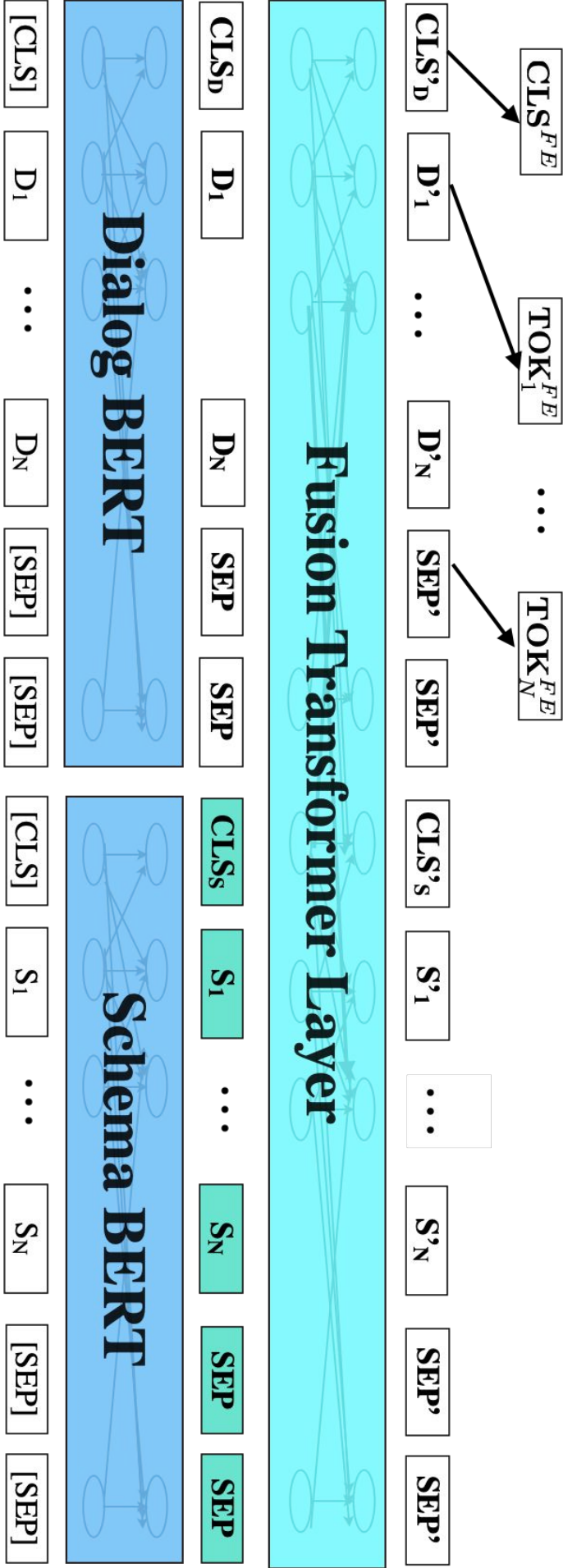
Pros:

Encoding dialog history and schema independently, can be precomputed once and cached. Fast inference.

Cons:

- Local self-attention
- inaccurate

Q1: How to encode the dialog and schema?



Pros: moderate inference
Still independent precomputing.
but a thin full-attention fusion layer for better performance

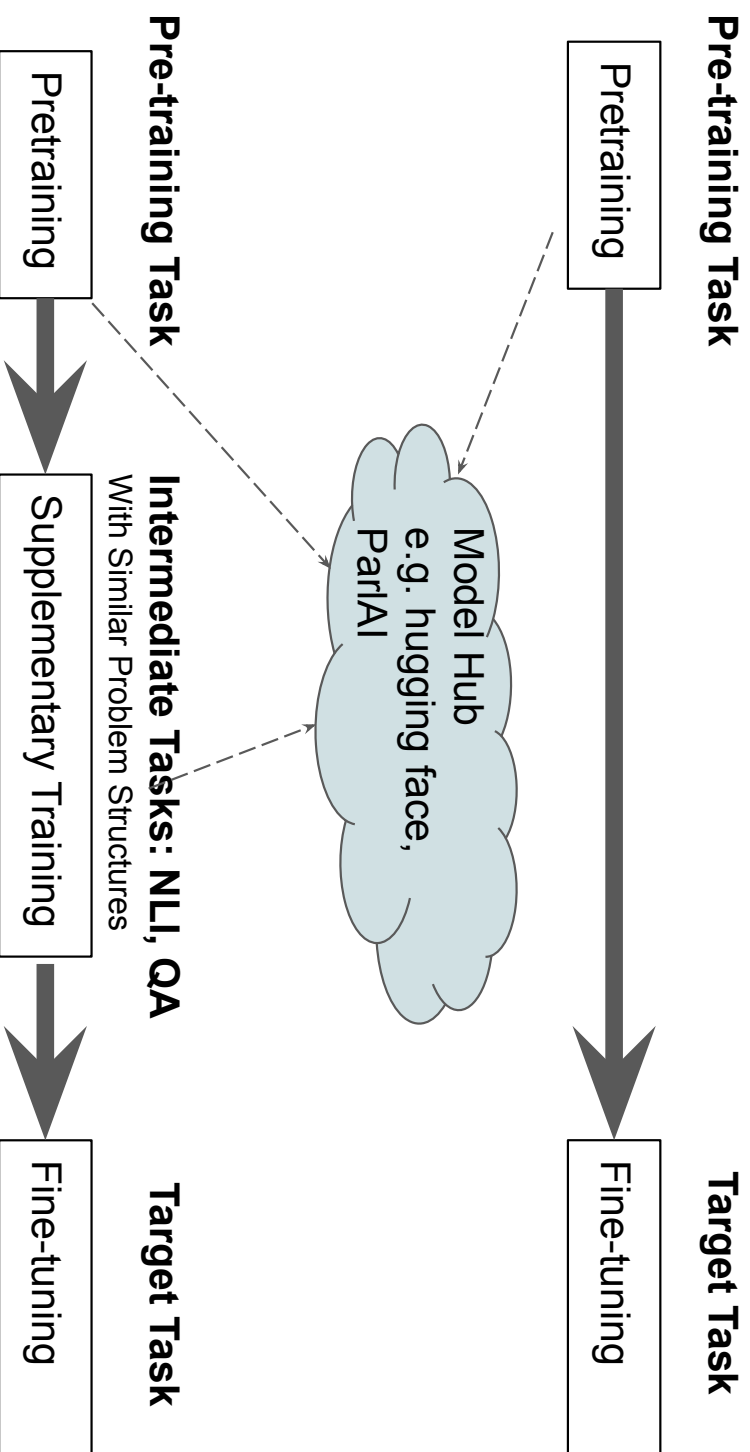
Cons:
Moderate accuracy

Q1: How to encode the dialog and schema(for 4 subtasks) ?

	SG-DST					MULTIWOZ 2.2				
Method/Task	Acc	F1	Joint Acc			Joint Acc				
	Intent	Req	Cat	NonCat	All	Cat	NonCat	All		
Seen Services										
Dual-Encoder	94.51	99.62	87.92	47.77	43.20	79.20	79.34	65.64		
Fusion-Encoder	94.90	99.69	88.94	48.78	58.52	81.37	80.58	67.43		
Cross-Encoder	95.55	99.59	93.68	91.85	87.58	85.99	81.02	71.93		
Unseen Services										
Dual-Encoder	89.73	95.20	42.44	31.62	19.51	56.92	50.82	31.83		
Fusion-Encoder	90.47	95.95	48.79	35.91	22.85	57.01	52.23	33.64		
Cross-Encoder	93.84	98.26	71.55	74.13	54.54	59.85	59.62	38.46		

By caching the token embedding instead of the single CLS embedding, a simple partial-attention ***Fusion-Encoder*** can achieve much better performance than ***Dual-Encoder***, while still infers two times faster than ***Cross-Encoder***

Q2: How do different supplementary trainings help?



Q2: How do different supplementary trainings help?

SG-DST							
	Intent		Req		Cat		NonCat
	seen	unseen	seen	unseen	seen	unseen	seen unseen
Δ_{SNLI}	+0.02	+0.68	+0.38	-0.38	-2.87	-1.23	-0.1 -6.25
Δ_{SQuAD}	-0.17	-1.32	-0.01	-0.33	-3.02	-5.17	-1.79 +3.25

- **SNLI only** helps for Intent (emphasizing the whole sentence entailment), although Req and Cat are also sentence-pair classification tasks.
- **SQuAD** consistently helps for non-categorical slot identification tasks, due to span-based retrieving
- Supplementary training helps more on unseen services

Q3: How the model performs on various description styles?

Background:

1. To compatible with previous tag-based DST system, many previous papers show **simply adding question format to those tags** may help.
 - a. Is name-based description enough?
 - b. Does question format helps?
2. **Unseen services may use different description style**
 - a. **heterogeneous evaluation?**

style	Intent Description		Slot Description	
<i>Identifier</i>	intent_1		slot_4	
<i>NameOnly</i>	CheckBalance		account_type	
<i>Q-Name</i>	Is the user intending to CheckBalance?		What is the value of acctount_type ?	
<i>Orig</i>	Check the amount of money in a user's bank account		The account type of the user	
<i>Q-Orig</i>	Does the user want to check the amount of money in the bank account ?		What is the account type of the user ?	
<i>Orig-Para</i>	Check the balance of the user's bank account		Type of the user account	

Q3: How the model performs on various description styles? (homogeneous)

Style\Task	SG-DST				MULTIWOZ 2.2	
	Intent	Req	Cat	NonCat	Cat	NonCat
<i>Identifier</i>	61.16	91.48	62.47	30.19	34.25	52.28
<i>NameOnly</i>	94.24	98.84	74.01	75.63	53.72	56.18
<i>Q-Name</i>	93.31	98.86	74.36	74.86	54.19	56.17
<i>Orig</i>	93.01	98.55	74.51	75.76	52.19	57.20
<i>Q-Orig</i>	93.42	98.51	76.64	76.60	53.61	57.80

Is named-based description enough?

- Most name are meaningful, and perform **not bad**, especially on Intent/Req subtasks
- Rich description outperforms the name-based on **NonCat**, but inconsistent on other tasks.

Q3: How the model performs on various description styles? (homogeneous)

Style/Dataset	SG-DST		MULTIWOZ 2.2	
	seen	unseen	seen	unseen
<i>Orig</i>	-1.79	+3.25	-2.21	+4.27
<i>Q-Orig</i>	-2.01	+8.84	-1.28	+3.06
<i>NameOnly</i>	-1.49	-0.11	+0.58	+1.77
<i>Q-Name</i>	-2.98	+1.04	-0.32	+1.25

Is question-format helpful?

- It generally helps on Cat/NonCat
 - Adding it to rich description will benefit more from SQuAD2
- supplementary training on unseen. However, not on MultiWOZ.

Q3: How the model performs on various description styles? (Heterogeneous)

Style\Task	SG-DST							
	Intent(Acc)		Req(F1)		Cat(Joint Acc)		NonCat(Joint Acc)	
	mean	Δ	mean	Δ	mean	Δ	mean	Δ
<i>NameOnly</i>	82.47	-11.47	96.92	-1.64	61.37	-5.54	56.53	-14.68
<i>Q-Name</i>	93.27	+0.58	97.88	-0.76	68.55	+2.63	62.92	-6.30
<i>Orig</i>	79.47	-12.70	97.42	-0.74	68.58	-0.3	66.72	-3.11
<i>Q-Orig</i>	84.57	-8.24	96.70	-1.45	68.40	-2.89	56.17	-15.00
<i>NameOnly</i>	para	Δ	para	Δ	para	Δ	para	Δ
	92.22	-1.74	97.69	-0.87	67.39	-0.7	67.17	-4.04
<i>Orig</i>	91.54	-0.63	98.42	+0.26	71.74	+2.86	67.68	-2.16

What if unseen service in different description styles?

- For unseen styles, all tasks suffer from inconsistencies, though to varying degrees
- For paraphrased styles, richer description are relatively more robust than named-based descriptions.

Takeaways

1. Cross-Encoder > Fusion-Encoder > Dual-Encoder in accuracy, while opposite on inference speed.
2. To support low-resource unseen services, we quantified the gain via supplementary training on different subtasks.
3. Simple named-based description are actually meaningful, and they perform not bad, but not as robust as rich description in most cases.
4. All subtasks suffers from inconsistencies when using heterogeneous description on unseen services, which requires future work on more robust cross-style schema-guided dialog modeling.

Q&A?

Thanks