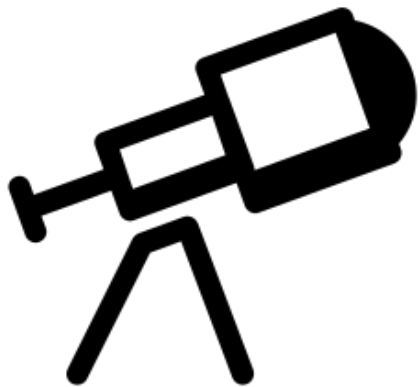

TypeSQL: From Natural Text to SQL Queries

— *DataBase Systems 2020* —
Apostolos Papatheodorou

Slides Structure

- *TypeSQL introduction*
- *Implementation*
- *Evaluation & Plan*



Presentation: Instead of Introduction 1/2



Question for a relation



Corresponding
Sql query



Type Sql

- a) **Slot filling approach**
- b) **Knowledge based & Type aware**

i.e. Assign each query word a type(Column, Number, KB entity)

- c) **Search DB rows for better accuracy (+9% improvement)**

Presentation: Basic idea (2/2)

What are Slots ???

```
SELECT $AGGR $SEL_COL  
WHERE $COND_COL $OP $COND_VAL  
(AND $COND_COL $OP $COND_VAL)*
```

Query:

Was J. Biden v. president of US in 2019 ?

(none) (person) (column) (country) (year)

Three steps Formula

- 1) **Preprocessing: type recognition**
- 2) **Bi-directional LSTM for**
 - a) **Type-word encoding**
 - b) **Column's name encoding**
- 3) **Slot Prediction (more LSTMs)**

Implementation: Encoding (1/3)

Two Bi-LSTMs for

- a) Quest: word-type \rightarrow H-QT
- b) Column schema \rightarrow H-COL

Output

Slot Predictions

```
SELECT $AGGR $SEL_COL  
WHERE $COND_COL $OP $COND_VAL  
AND $COND_COL $OP $COND_VAL)*
```

Implementation of

Three LSTMs
for prediction

\$AGG

\$OP, \$COND_VAL,

\$SEL_COL, \$COND_COL, \$COND#

(5 LSMTs all together)

$$\alpha_{QT/COL} = \text{softmax}(\mathbf{H}_{COL} \mathbf{W}_{ct} \mathbf{H}_{QT}^T)$$
$$\mathbf{H}_{QT/COL} = \alpha_{QT/COL} \mathbf{H}_{QT}$$

Implementation: Preprocessing (2/3)

Query:

Was J. Biden v. president of US in 2019 ?

(none) (person) (column) (country) (year)

- 1) Tokenization k-grams, $k=[2-6]$ for type assignment.

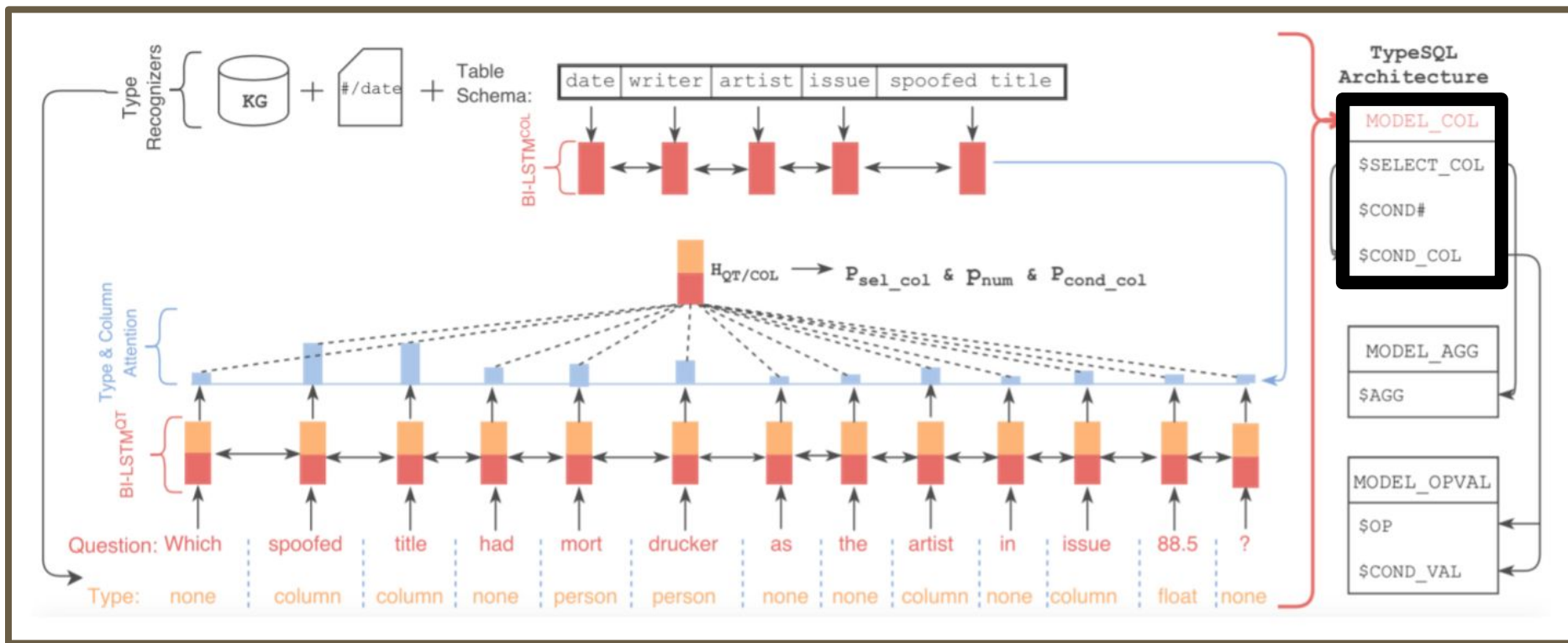
MyQuest: Random choice of k value?

Based on k-grams Search

- a) ***Columns from schema***
- b) ***Entities in Freebase KB***
- c) ***numbers, dates, years***

D) IF RELATION AVAILABLE USE IT
(+9% improvement)

Paper's paradigm



Implementation: Final Models(3/3)

MODEL_COL: (3 Slots)

\$SELECT_COL

\$COND#

\$COND_Col



MODEL_COL- \$SELECT_COL

$$s = V^{sel} \tanh(W_c^{sel} H_{COL}^T + W_{qt}^{sel} H_{QT/COL}^T)$$

$$P_{sel.col} = \text{softmax}(s)$$



MODEL_COL- \$SELECT_COND#

$$c = V^{col} \tanh(W_c^{col} H_{COL}^T + W_{qt}^{col} H_{QT/COL}^T + W_{qt}^{scol} H_{QT/SCOL}^T)$$

$$P_{cond.col} = \text{softmax}(c)$$

Evaluation & Plan

	Dev			Test		
	Acc _{agg}	Acc _{sel}	Acc _{where}	Acc _{agg}	Acc _{sel}	Acc _{where}
Seq2SQL (Zhong et al., 2017)	90.0%	89.6%	62.1%	90.1%	88.9%	60.2%
SQLNet (Xu et al., 2017)	90.1%	91.5%	74.1%	90.3%	90.9%	71.9%
TypeSQL (ours)	90.3%	93.1%	78.5%	90.5%	92.2%	77.8%
TypeSQL+TC (ours)	90.3%	93.5%	92.8%	90.5%	92.1%	87.9%

Recap & Action Plan

For accomplishing the ultimate goal: **NLP text** to **SQL query**

**Three LSTMs
for prediction**

Evaluation & Plan

	Dev			Test		
	Acc _{agg}	Acc _{sel}	Acc _{where}	Acc _{agg}	Acc _{sel}	Acc _{where}
Seq2SQL (Zhong et al., 2017)	90.0%	89.6%	62.1%	90.1%	88.9%	60.2%
SQLNet (Xu et al., 2017)	90.1%	91.5%	74.1%	90.3%	90.9%	71.9%
TypeSQL (ours)	90.3%	93.1%	78.5%	90.5%	92.2%	77.8%
TypeSQL+TC (ours)	90.3%	93.5%	92.8%	90.5%	92.1%	87.9%

Recap & Action Plan

For accomplishing the ultimate goal: **NLP text** to **SQL query**

Two Basic Competitors:

Seq2SQL & **SQLNet**

*Three LSTMs
for prediction*

\$AGG

\$OP, \$COND_VAL,

\$SEL_COL, \$COND_COL, \$COND#

+2 LSMTs During preprocessing

Evaluation & Plan

	Dev			Test		
	Acc _{agg}	Acc _{sel}	Acc _{where}	Acc _{agg}	Acc _{sel}	Acc _{where}
Seq2SQL (Zhong et al., 2017)	90.0%	89.6%	62.1%	90.1%	88.9%	60.2%
SQLNet (Xu et al., 2017)	90.1%	91.5%	74.1%	90.3%	90.9%	71.9%
TypeSQL (ours)	90.3%	93.1%	78.5%	90.5%	92.2%	77.8%
TypeSQL+TC (ours)	90.3%	93.5%	92.8%	90.5%	92.1%	87.9%

Recap & Action Plan

For accomplishing the ultimate goal: **NLP text** to **SQL query**

Two Basic Competitors:

Seq2SQL & **SQLNet**

Implementation details

- ❖ WikiSQL Dataset
- ❖ Pennington et al & *Wieting Glmbel*
- ❖ Tuning: **dimens**, **dropout rate**, etc.
- ❖ **Adam** optimizer
- ❖ *My experiments on optimizers & tuning*

Thanks for watching

