

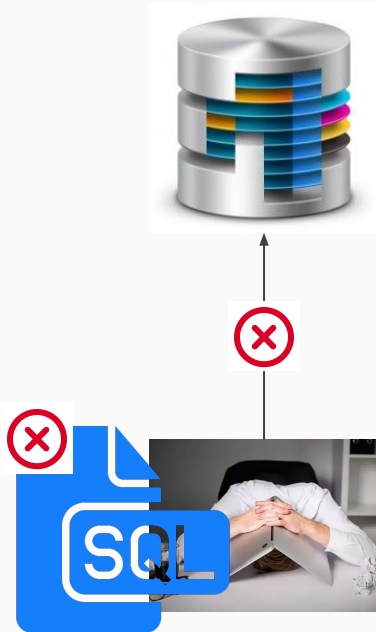
# Text to SQL

Automated ever-improving ad-hoc SQL query generation

**Atakan Okan**

A decorative light blue triangle is located in the bottom right corner of the slide, pointing towards the top right.

# Business Problem



# Business Problem



# Defining the Problem

JUST ADDRESSING:

SELECT X  
FROM Y  
WHERE Z > W

**SQLNet: GENERATING STRUCTURED QUERIES FROM  
NATURAL LANGUAGE WITHOUT REINFORCEMENT  
LEARNING**

Xiaojun Xu\*  
Shanghai Jiao Ton

**Robust Text-to-SQL Generation with Execution-Guided Decoding**

Chenglong Wang,<sup>1\*</sup> Kedar Tatwawadi,<sup>2\*</sup> Marc Brockschmidt,<sup>3</sup> Po-Sen Huang,<sup>3†</sup>

Yi Mao,<sup>3</sup> Oleksandr Polozov,<sup>3</sup> Rishabh Singh<sup>4†</sup>

<sup>1</sup>University of Washington <sup>2</sup>Stanford University <sup>3</sup>Microsoft Research <sup>4</sup>Google Brain

{clwang@cs.washington.edu, kedar@stanford.edu, {mabrocks, pshuang, maoyi, polozov}@microsoft.com, rising@google.com}

**DialSQL: Dialogue Based Structured Query Generation**

Izzeddin Gonen and Sarah Young and Xu Song and Yifan Yan

Department of INCSQL: TRAINING INCREMENTAL TEXT-TO-SQL  
{izzeddin, sarah}@cs.stanford.edu PARSERS WITH NON-DETERMINISTIC ORACLES

Tianze Shi \*  
Stanford University

Kedar Tatwawadi \*  
Stanford University

**Natural Language to Structured Query Generation via Meta-Learning**

Po-Sen Huang\*, Chenglong Wang†, Rishabh Singh\*, Wen-tau Yih†, Xiaodong He\*◊

\*Microsoft Research †University of Washington

**TypeSQL: Knowledge-based Type-Aware Neural Text-to-SQL Generation**

Tao Yu

Yale University

tao.yu@yale.edu

Zifan Li

Yale University

zifan.li@yale.edu

Zilin Zhang

Yale University

zilin.zhang@yale.edu

Achieving 90% accuracy in WikiSQL

Wonseok Hwang

wonseok.hwang@navercorp.com

Jinyeong Yim

jinyeong.yim@navercorp.com

Seunghyun Park

seung.park@navercorp.com

Minjoon Seo

minjoon.seo@navercorp.com

Clova AI Research, NAVER Corp., Seongnam, Korea

January 9, 2019\*

ragomir Radev

Yale University

r.radev@yale.edu

# Natural Language to Structured Language



**STRUCTURED  
QUERY  
LANGUAGE**

Our SQL generation process should use its structural rules!

# SyntaxSQLNet: How it works

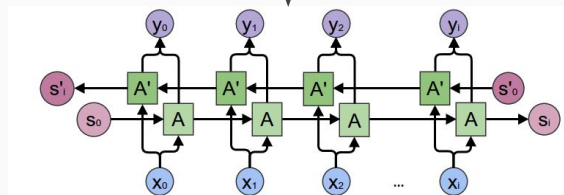
**QUESTION**

What are the maximum and minimum budget of the departments?

# SyntaxSQLNet: How it works

## QUESTION

What are the maximum and minimum budget of the departments?



## ENCODING

Bidirectional LSTM with Attention

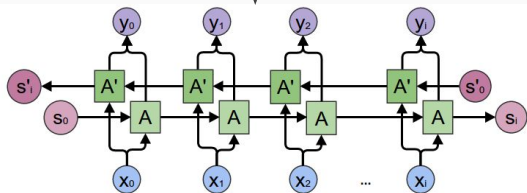
QUESTION EMBEDDING (CONTEXT VECTOR)

# SyntaxSQLNet: How it works

RECURSIVE

## QUESTION

What are the maximum and minimum budget of the departments?



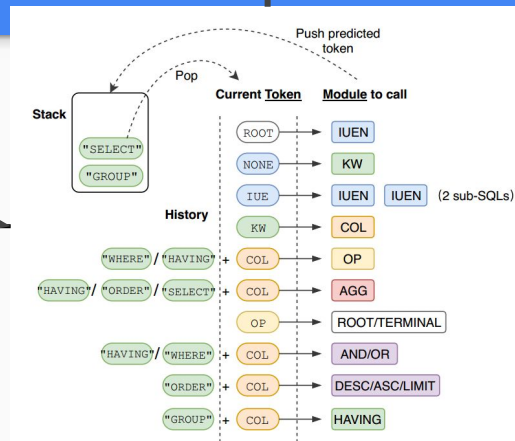
## ENCODING

Bidirectional LSTM with Attention

QUESTION EMBEDDING (CONTEXT VECTOR)

## STACK

ROOT



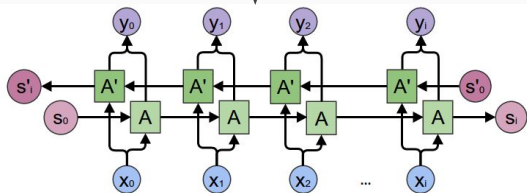


# SyntaxSQLNet: How it works

RECURSIVE

## QUESTION

What are the maximum and minimum budget of the departments?



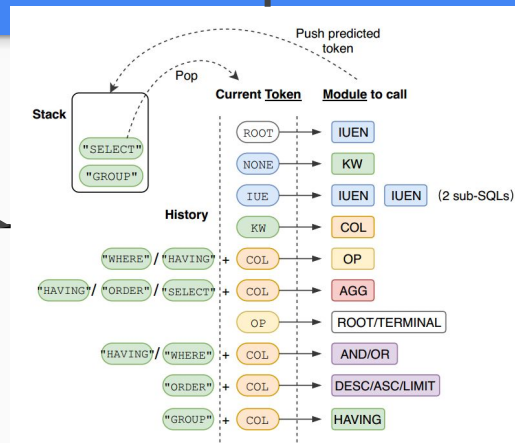
## ENCODING

Bidirectional LSTM with Attention

QUESTION EMBEDDING (CONTEXT VECTOR)

## STACK

ROOT



## GENERATED SET

{'sql': {'select': [('department', 'budget in billions', 5), 'min', ('department', 'budget in billions', 5), 'max']}}

## PARSED QUERY

select min(Budget\_in\_Billions),max(Budget\_in\_Billions) from department

# Challenges

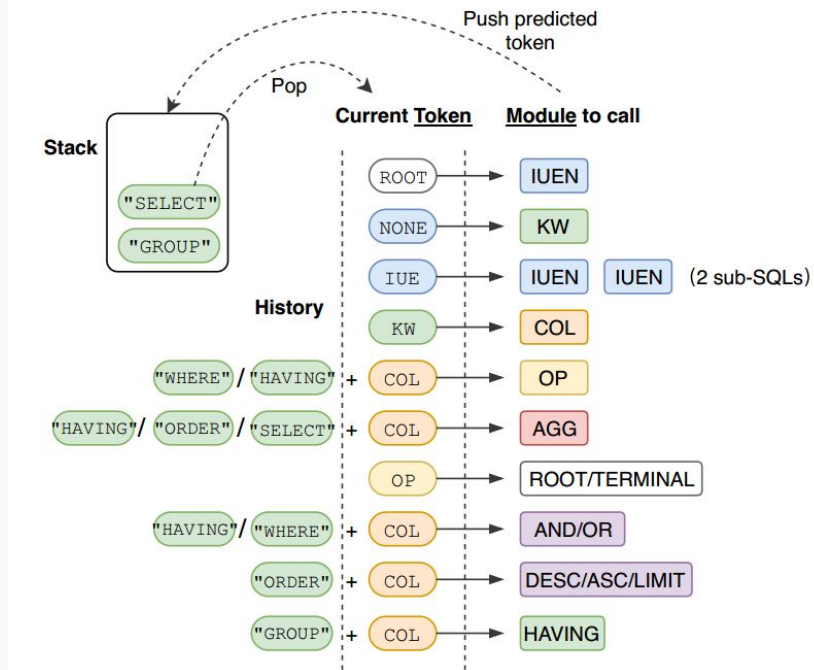
USER SUPPLIED QUESTION	CORRECT QUERIES	SYNTAXSQLNET
What are the maximum and minimum budget of the departments?	<b>select min(Budget_in_Billions), max(Budget_in_Billions) from department</b>	<b>select min(Budget_in_Billions), max(Budget_in_Billions) from department</b>
What is the department ID of the state department?	<b>select Department_ID from department where Name = 'State'</b>	<b>select Department_ID from department where Name = 'terminal'</b>
What are the names of the departments that were founded after 1800?	<b>select Name from department where Creation &gt; 1800</b>	<b>select T1.Name from department as T1 join management as T2 on T1.Department_ID = T2.department_ID join head as T3 on T2.head_ID = T3.head_ID where T3.age &gt; 'terminal'</b>

Wrong column selection

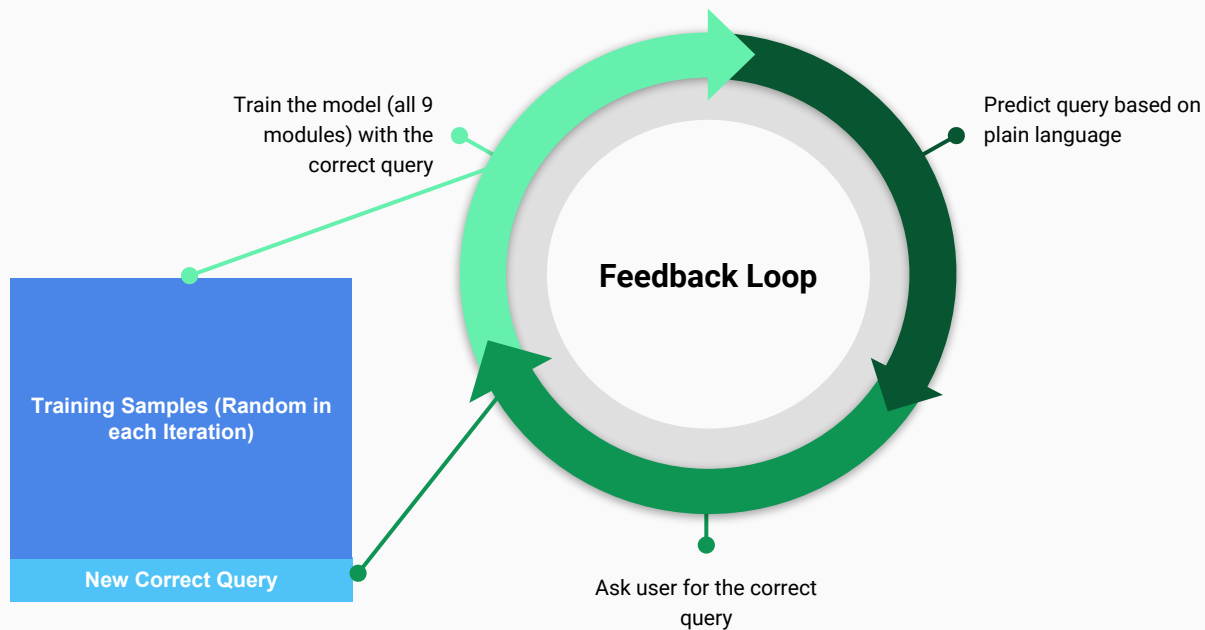
Question: Simple query  
Model prediction: complex query  
Correct answer: simple query

# Other Challenges

- Decoder: 9 Modules -> Trained
  - LONG TRAINING TIME
- Each module has different datasets created by different parsers



# Human-in-the-loop Training / Active Learning

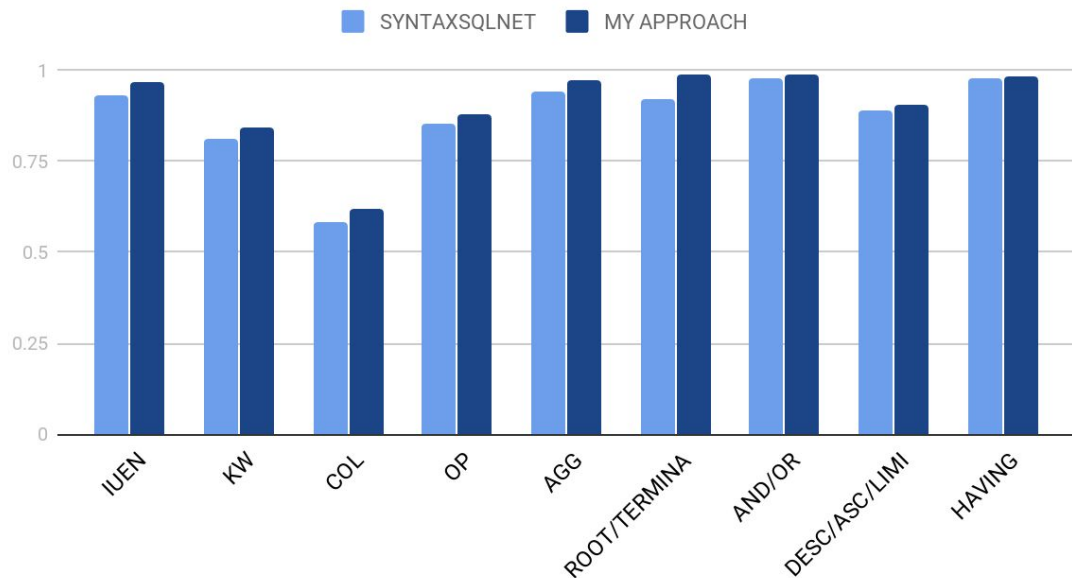


# Improvement

USER SUPPLIED QUESTION	CORRECT QUERIES	SYNTAXSQLNET	MY APPROACH
What are the maximum and minimum budget of the departments?	<b>select min(Budget_in_Billions), max(Budget_in_Billions) from department</b>	<b>select min(Budget_in_Billions), max(Budget_in_Billions) from department</b>	<b>select min(Budget_in_Billions), max(Budget_in_Billions) from department</b>
What is the department ID of the state department?	<b>select Department_ID from department where Name = 'State'</b>	<b>select Department_ID from department where Name = 'terminal'</b>	<b>select Department_ID from department where Name = 'terminal'</b>
What are the names of the departments that were founded after 1800?	<b>select Name from department where Creation &gt; 1800</b>	<b>select T1.Name from department as T1 join management as T2 on T1.Department_ID = T2.department_ID join head as T3 on T2.head_ID = T3.head_ID where T3.age &gt; 'terminal'</b>	<b>select Name from department where Creation &gt; 1800</b>

# Side by Side Comparison

Binary Accuracy For Each Module



## Text to SQL

Uses SyntaxSQLNet Implementation

**Please enter your question:**

**Please enter the database to be queried:**

Generate SQL Query!

## Text to SQL

Uses SyntaxSQLNet Implementation

Please enter your question:

Please enter the database to be queried:

Generate SQL Query!

## Generated Query:

Question:

What are the maximum and minimum budget of the departments?

DB Name:

department\_management

Copy paste this query into your SQL query UI:

**select min(Budget\_in\_Billions),max(Budget\_in\_Billions) from department**



## Text to SQL

Uses SyntaxSQLNet Implementation

Please enter your question:

-

Please enter the database to be queried:

-

Generate SQL Query!

## Generated Query:

Question:

What are the maximum and minimum budget of the departments?

DB Name:

department\_management

Copy paste this query into your SQL query UI:

**`select min(Budget_in_Billions),max(Budget_in_Billions) from department`**

## Query Feedback:

Please enter your question:

-

Please enter the database to be queried:

-

Please enter the correct query:

Enter Correct Query Here

Train the model!

## Database Tables:

department

Department_ID	Name	Creation	Ranking	Budget_in_Billions	Num_Employees
01	State	1789	1	9.96	30266.0
12	Treasury	1789	2	11.10	115897.0
23	Defense	1947	3	439.30	3000000.0

head

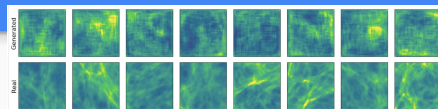
head_ID	name	born_state	age
01	Tiger Woods	Alabama	67.0
12	Sergio García	California	68.0
23	K. J. Choi	Alabama	69.0

management

department_ID	head_ID	temporary_acting
02	5	Yes
115	4	Yes
22	6	Yes

# Atakan Okan

- Replicating numerical simulations using MMD-GANs at Flatiron Institute CCA
- NYU MS in Data Science
  - Data Science
  - Deep Learning
    - Computer Vision
    - Natural Language Processing
    - Generative Adversarial Networks
- Junior Data Scientist at Blisce Venture Capital
  - Data Infrastructure/Databases
  - Business Intelligence
  - Natural Language Processing



Q&A

# APPENDIX

# Mathematical Notations and Explanations

$$\mathbf{H}_{1/2} = \mathbf{softmax}(\mathbf{H}_1 \mathbf{W} \mathbf{H}_2^\top) \mathbf{H}_1.$$

$$\mathcal{P}(\mathbf{U}) = \mathbf{softmax}(\mathbf{V} \mathbf{tanh}(\mathbf{U})).$$

H\_Q = LSTM Hidden State on question embedding

H\_HS = LSTM Hidden State on path history

H\_COL = LSTM Hidden State on column embeddings

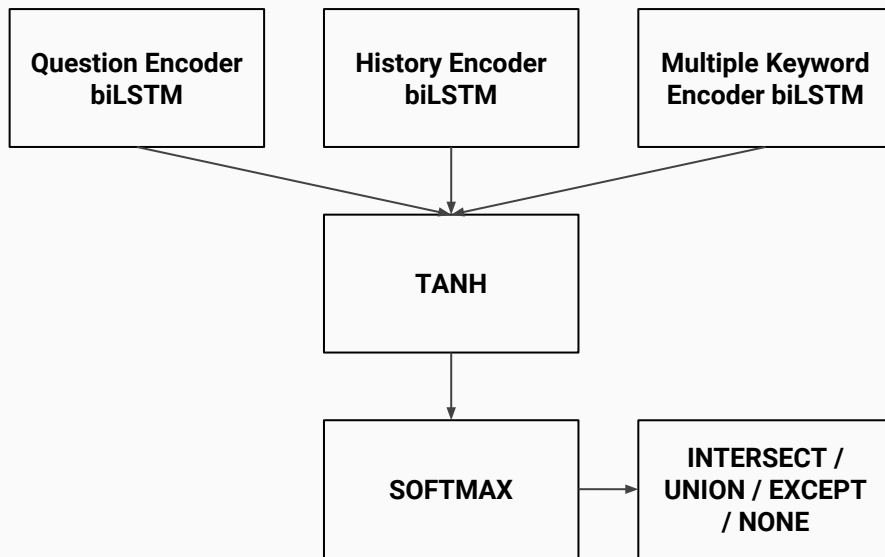
H\_MKW = LSTM hidden state on multiple keywords

H\_KW = LSTM hidden state on keyword embeddings

# IUEN

INTERSECT, UNION, EXCEPT, NONE

Only one will be selected via softmax

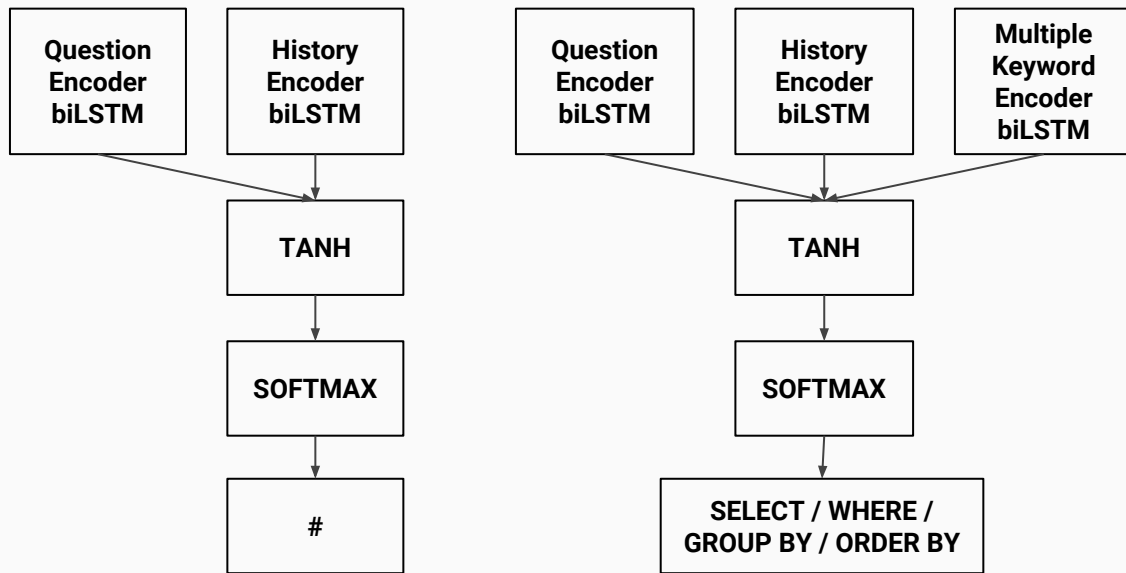


$$P_{\text{IUEN}} = \mathcal{P} \left( \mathbf{w}_1 \mathbf{H}_{\text{Q/MKW}}^{\top} + \mathbf{w}_2 \mathbf{H}_{\text{HS/MKW}}^{\top} + \mathbf{w}_3 \mathbf{H}_{\text{MKW}}^{\top} \right)$$

# KW

KEYWORD: SELECT, WHERE, GROUP  
BY & ORDER BY

1. Predict the number of keywords in SQL
2. Then predict keywords from the list



$$P_{KW}^{\text{num}} = \mathcal{P} \left( \mathbf{W}_1^{\text{num}} \mathbf{H}_{Q/KW}^{\text{num} \top} + \mathbf{W}_2^{\text{num}} \mathbf{H}_{HS/KW}^{\text{num} \top} \right)$$

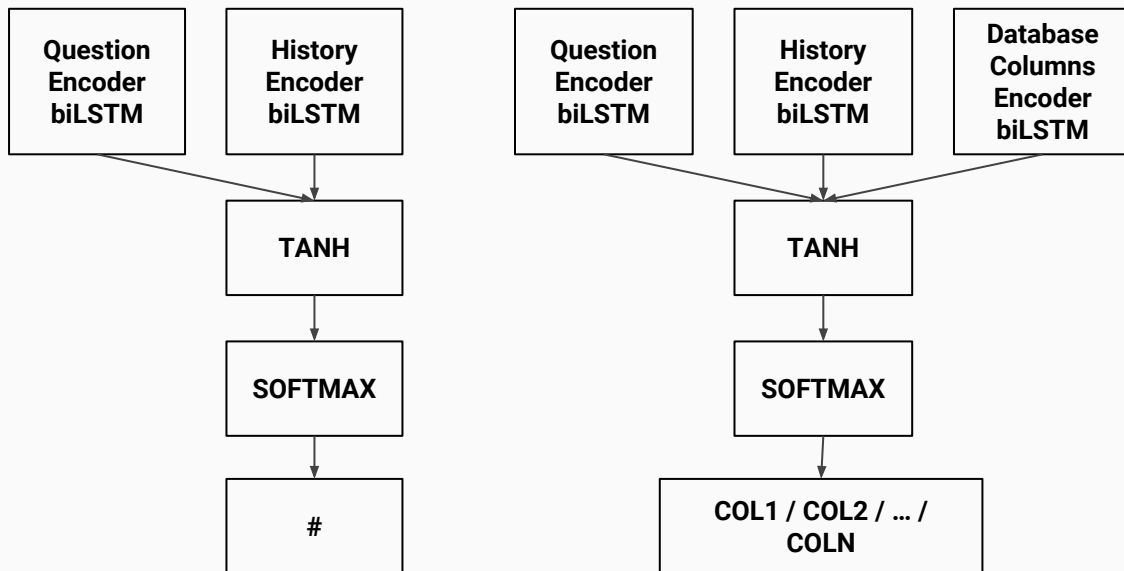
$$P_{KW}^{\text{val}} = \mathcal{P} \left( \mathbf{W}_1^{\text{val}} \mathbf{H}_{Q/KW}^{\text{val} \top} + \mathbf{W}_2^{\text{val}} \mathbf{H}_{HS/KW}^{\text{val} \top} + \mathbf{W}_3^{\text{val}} \mathbf{H}_{KW}^{\text{val} \top} \right)$$

# COL

## COLUMN

PREDICTS THE TABLE COLUMNS

1. Predict the number of columns
2. Predict which ones



$$P_{\text{COL}}^{\text{num}} = \mathcal{P} \left( \mathbf{W}_1^{\text{num}} \mathbf{H}_{\text{Q/COL}}^{\text{num} \top} + \mathbf{W}_2^{\text{num}} \mathbf{H}_{\text{HS/COL}}^{\text{num} \top} \right)$$
$$P_{\text{COL}}^{\text{val}} = \mathcal{P} \left( \mathbf{W}_1^{\text{val}} \mathbf{H}_{\text{Q/COL}}^{\text{val} \top} + \mathbf{W}_2^{\text{val}} \mathbf{H}_{\text{HS/COL}}^{\text{val} \top} + \mathbf{W}_3^{\text{val}} \mathbf{H}_{\text{COL}}^{\text{val} \top} \right)$$



# OP

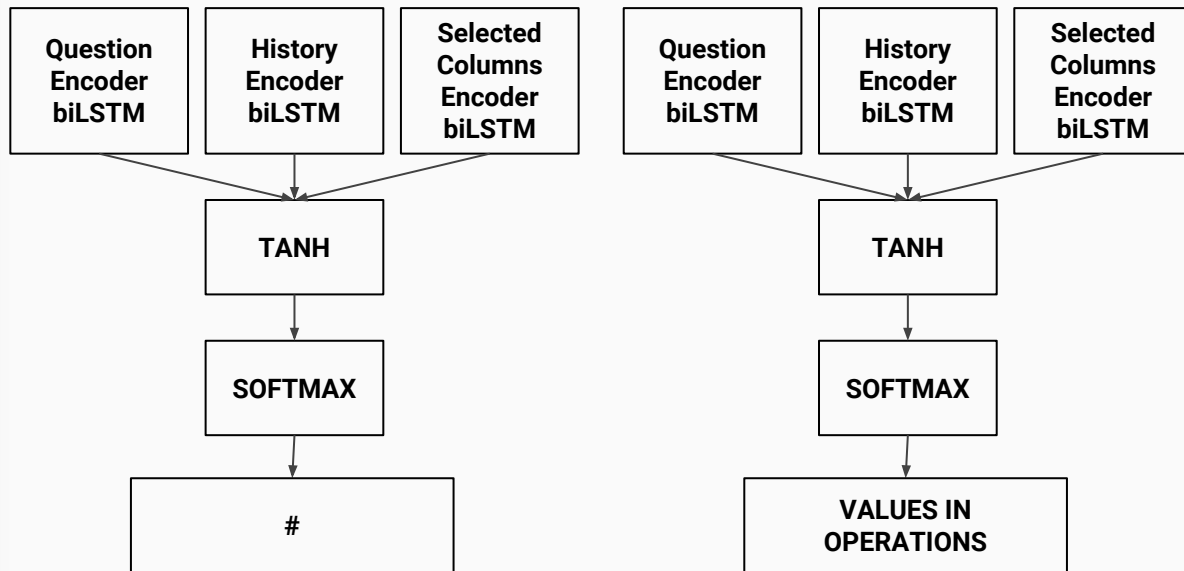
## OPERATIONS

= / > / < / >= / <= / != / LIKE / NOT IN /  
IN / BETWEEN

For each predicted column that is in the WHERE clause:

1. Predict the number of operators on each column
2. Predict which operation(s) from the list

H\_CS = Embedding of one of the predicted columns from COL module



$$P_{OP}^{num} = \mathcal{P} \left( \mathbf{W}_1^{num} \mathbf{H}_{Q/CS}^{num \top} + \mathbf{W}_2^{num} \mathbf{H}_{HS/CS}^{num \top} + \mathbf{W}_3^{num} \mathbf{H}_{CS}^{num \top} \right)$$
$$P_{OP}^{val} = \mathcal{P} \left( \mathbf{W}_1^{val} \mathbf{H}_{Q/CS}^{val \top} + \mathbf{W}_2^{val} \mathbf{H}_{HS/CS}^{val \top} + \mathbf{W}_3^{val} \mathbf{H}_{CS}^{val \top} \right)$$

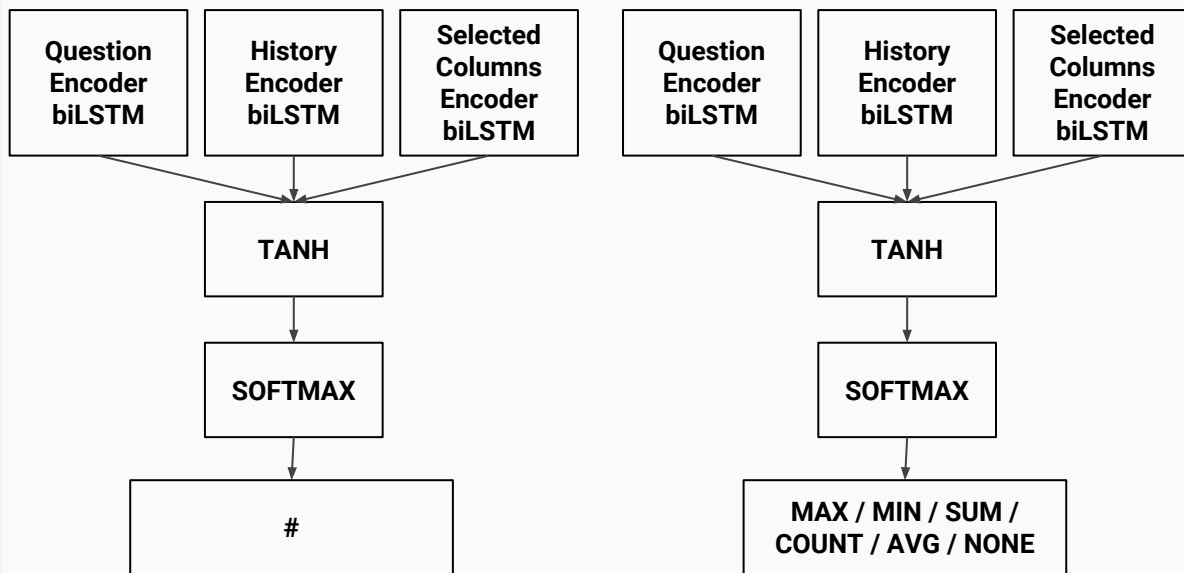
# AGG

## AGGREGATOR

MAX / MIN / SUM / COUNT / AVG /  
NONE

For each predicted column in COL  
module:

1. Predict the number of  
aggregators on that column
2. Predict which aggregators



$$P_{AGG}^{num} = \mathcal{P} \left( \mathbf{W}_1^{num} \mathbf{H}_{Q/CS}^{num \top} + \mathbf{W}_2^{num} \mathbf{H}_{HS/CS}^{num \top} + \mathbf{W}_3^{num} \mathbf{H}_{CS}^{num \top} \right)$$
$$P_{AGG}^{val} = \mathcal{P} \left( \mathbf{W}_1^{val} \mathbf{H}_{Q/CS}^{val \top} + \mathbf{W}_2^{val} \mathbf{H}_{HS/CS}^{val \top} + \mathbf{W}_3^{val} \mathbf{H}_{CS}^{val \top} \right)$$

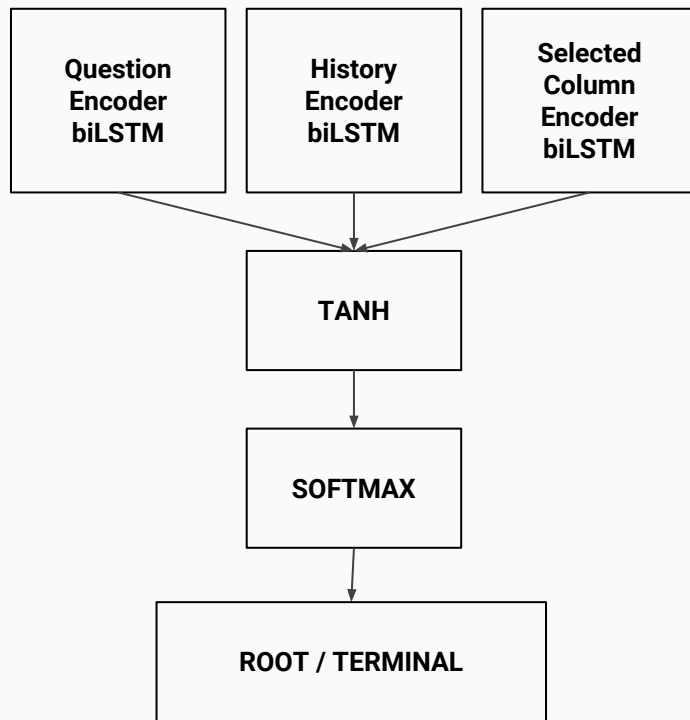
# ROOT/TERMINAL

ROOT OF THE SUBQUERY //  
TERMINAL VALUE

ENABLES NESTED QUERIES

For each of the predicted column from the COL module that is in the WHERE clause:

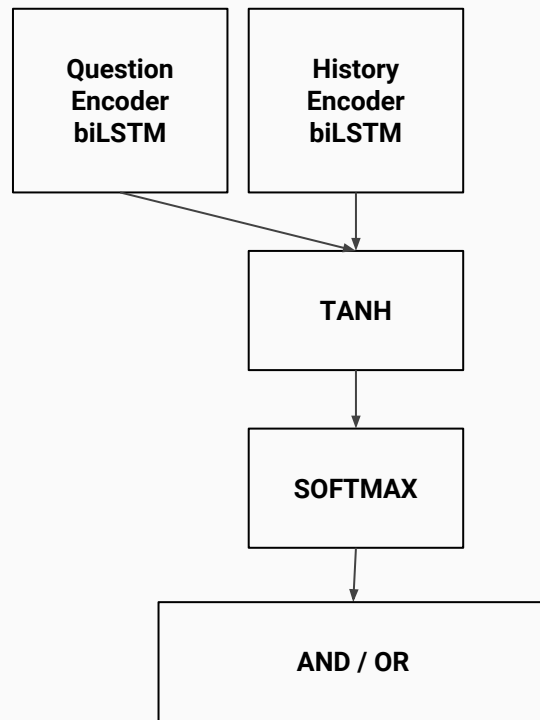
1. Call OP module
2. Predict whether the next decoding step is a “ROOT” node or a value terminal node



$$P_{RT} = \mathcal{P} \left( \mathbf{W}_1 \mathbf{H}_{Q/CS}^\top + \mathbf{W}_2 \mathbf{H}_{HS/CS}^\top + \mathbf{W}_3 \mathbf{H}_{CS}^\top \right)$$

# AND/OR

For each condition column predicted from the COL module with number bigger than 1



$$P_{AO} = \mathcal{P} \left( \mathbf{W}_1 \mathbf{H}_Q^\top + \mathbf{W}_2 \mathbf{H}_{HS}^\top \right)$$

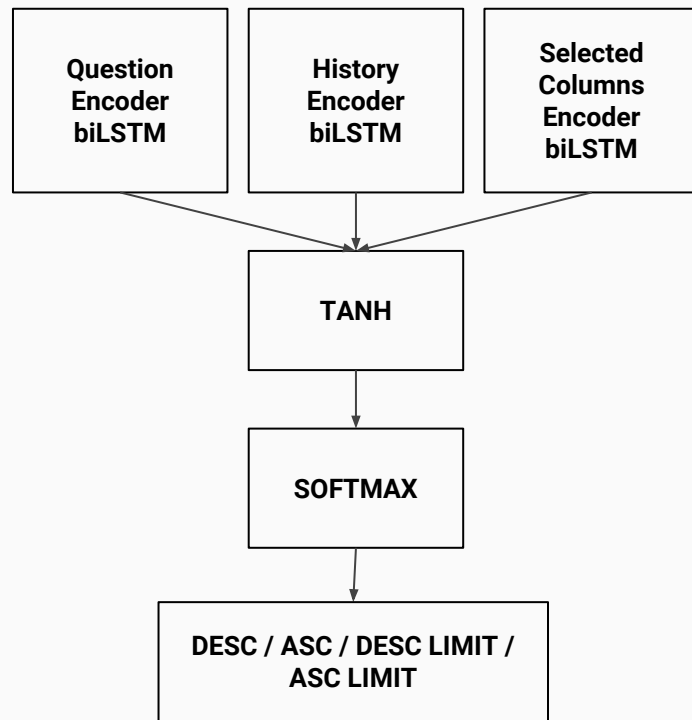
# DESC/ASC/LIMIT

KEYWORDS ASSOCIATED WITH  
'ORDER BY'

CALLED ONLY IF ORDER BY IS  
PREDICTED BEFORE

DESC/ASC/DESC LIMIT/ASC LIMIT

For each predicted column from the  
COL module that is in the ORDER BY  
clause

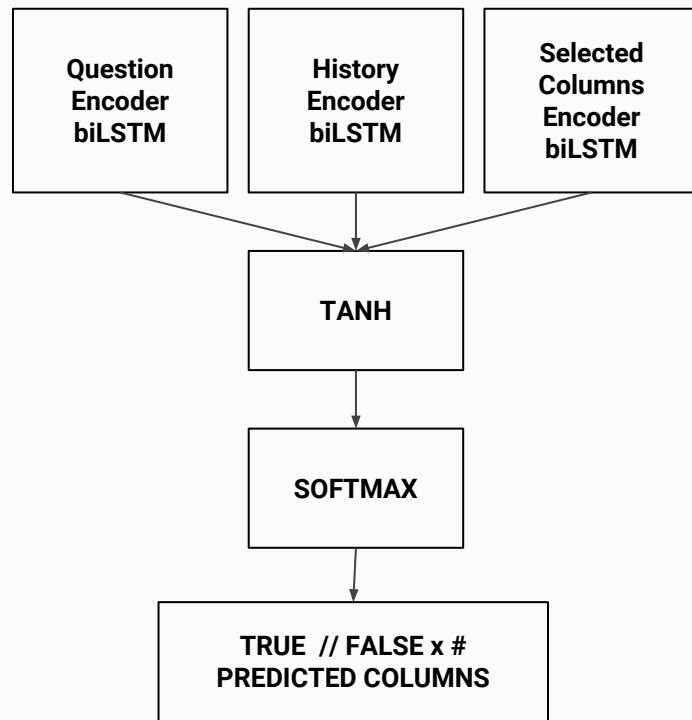


$$P_{\text{DAL}} = \mathcal{P} \left( \mathbf{W}_1 \mathbf{H}_{\text{Q/CS}}^{\top} + \mathbf{W}_2 \mathbf{H}_{\text{HS/CS}}^{\top} + \mathbf{W}_3 \mathbf{H}_{\text{CS}}^{\top} \right)$$

# HAVING

PREDICTS WHETHER 'ORDER BY'  
CONTAINS 'HAVING' ELEMENT

For each predicted column from the  
COL module that is in the GROUP BY  
clause, predict whether it is in the  
HAVING clause



$$P_{\text{HAVING}} = \mathcal{P} \left( \mathbf{W}_1 \mathbf{H}_{Q/CS}^{\top} + \mathbf{W}_2 \mathbf{H}_{HS/CS}^{\top} + \mathbf{W}_3 \mathbf{H}_{CS}^{\top} \right)$$