**Name:** TAPAS: Weakly Supervised Table Parsing via Pre-training

**Paper link:** https://arxiv.org/abs/2004.02349

A. **Introduction**
   ● Authors introduce an approach to answer NL questions from tables without generating logical forms (as is done in **semantic parsing models**)
   ● **Weakly supervised (learning with incomplete, inexact and inaccurate supervision signal) pre training** on text segments and tables crawled from Wikipedia
   ● Use **BERT architecture** to encode data in the table. Uses **transfer learning**
   ● Selects relevant cells + applies relevant **AGGREGATION** operations (SUM, MAX, MIN etc)

B. **Description**:
   a. Issues with Semantic Parsing models: They rely on **generated logical form** of the input user query. This logical form comes with issues like **maintaining a logical formalism** and **label bias problem.**
   b. **Contribution**:
      i. **ARCHITECTURE** : BERT architecture + table specific additional embeddings + 2 classification layers (one for cell selection and one for applying aggregation)
      ii. **PRETRAINING** : BERT's MLM (Masked Language Modelling) on millions of Wikipedia text and tables. Textual and tabular information is masked and the model is made to predict it based on context.
      iii. **FINETUNING**: End-to-end fully differentiable training that allows TaPaS to be trained under **weak supervision** which allows:
         1. **Only cell selection**
         2. **Cell selection + Aggregation**
   c. **Details:**
      i. **ARCHITECTURE**
         1. **EMBEDDINGS**

| Table | |
|---|---|
| col1 | col2 |
| 0 | 1 |
| 2 | 3 |

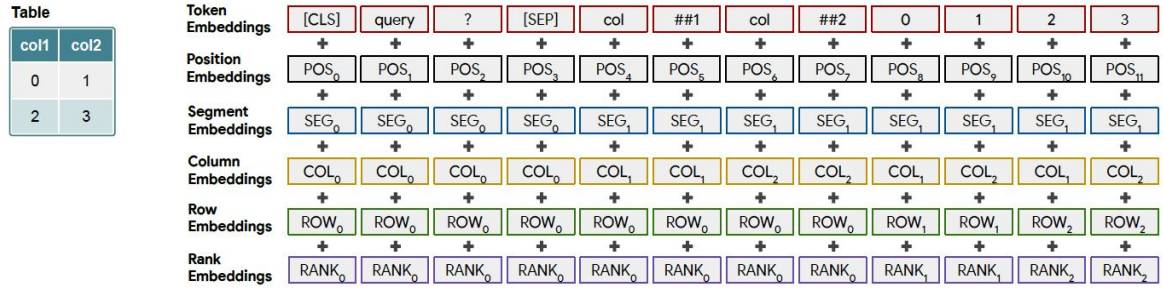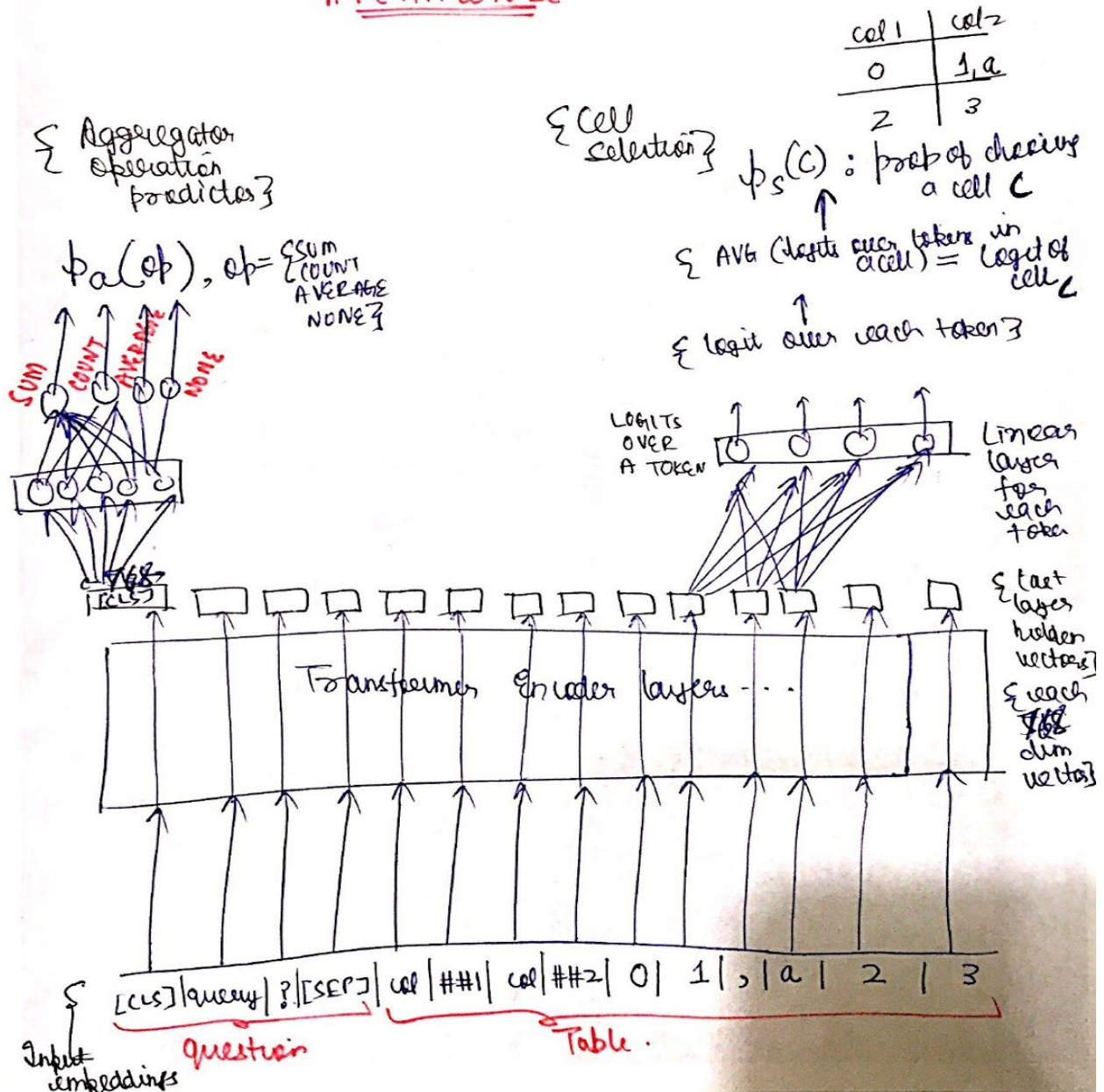| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Token Embeddings** | [CLS] | query | ? | [SEP] | col | ##1 | col | ##2 | 0 | 1 | 2 | 3 |
| | + | + | + | + | + | + | + | + | + | + | + | + |
| **Position Embeddings** | $POS_0$ | $POS_1$ | $POS_2$ | $POS_3$ | $POS_4$ | $POS_5$ | $POS_6$ | $POS_7$ | $POS_8$ | $POS_9$ | $POS_{10}$ | $POS_{11}$ |
| | + | + | + | + | + | + | + | + | + | + | + | + |
| **Segment Embeddings** | $SEG_0$ | $SEG_0$ | $SEG_0$ | $SEG_0$ | $SEG_1$ | $SEG_1$ | $SEG_1$ | $SEG_1$ | $SEG_1$ | $SEG_1$ | $SEG_1$ | $SEG_1$ |
| | + | + | + | + | + | + | + | + | + | + | + | + |
| **Column Embeddings** | $COL_0$ | $COL_0$ | $COL_0$ | $COL_0$ | $COL_1$ | $COL_1$ | $COL_2$ | $COL_2$ | $COL_1$ | $COL_2$ | $COL_1$ | $COL_2$ |
| | + | + | + | + | + | + | + | + | + | + | + | + |
| **Row Embeddings** | $ROW_0$ | $ROW_0$ | $ROW_0$ | $ROW_0$ | $ROW_0$ | $ROW_0$ | $ROW_0$ | $ROW_0$ | $ROW_1$ | $ROW_1$ | $ROW_2$ | $ROW_2$ |
| | + | + | + | + | + | + | + | + | + | + | + | + |
| **Rank Embeddings** | $RANK_0$ | $RANK_0$ | $RANK_0$ | $RANK_0$ | $RANK_0$ | $RANK_0$ | $RANK_0$ | $RANK_0$ | $RANK_1$ | $RANK_1$ | $RANK_2$ | $RANK_2$ |

Figure 2: Encoding of the question "query?" and a simple table using the special embeddings of TAPAS. The previous answer embeddings are omitted for brevity.

2. **CELL SELECTION + AGGREGATOR SELECTOR:** We choose the relevant cells + appropriate aggregator using the architecture shown below to get the inference/ output

CELL SELECTION and AGGREGATION SELECTOR
↳ each cell modelled as Bernoulli R.V. $p_s(c) > 0.5$ is selected

**ARCHITECTURE**

| col1 | col2 |
|------|------|
| 0 | 1,a |
| 2 | 3 |

{ Aggregator operation predictor }

$p_a(op)$, $op = \{$ SUM, COUNT, AVERAGE, NONE $\}$

SUM  COUNT  AVERAGE  NONE

768 [CLS]

{ Cell selection }  $p_s(c)$ : prob of choosing a cell $c$

{ AVG (logits over tokens in a cell) = logit of cell $c$ }

{ logit over each token }

LOGITS OVER A TOKEN

Linear layer for each token

{ last layer hidden vectors }
{ each 768 dim vector }

Transformer Encoder layers ...

{ [CLS] |query| ?|[SEP]| col |##1| col |##2| 0| 1|,|a| 2 | 3 |

question          Table.

Input embeddings

ii.  **PRE TRAINING:**
1.  **EXTRACTED DATA:** 6.2M tables (max 500 cells) and 21.3M relevant text snippets (table caption, article title, article description,

segment title, text of segment) from Wiki to form extracted text-table pairs.

2. **CONVERT DATA TO PRETRAINING EXAMPLES**: **len(tokenized text + table cells) <= 128** Randomly choose a **word piece snippet of length 8 to 16** from the associated text. To fit the table, we start by only adding the **first word of each column name and cell**. We then keep adding words turn-wise until we reach the word piece budget. **For every table we generate 10 different snippets in this way**

3. **OBJECTIVE:** MLM pre training objective with whole word + whole cell masking

4. Allows the model to learn correlations between text and table , and between cells of column and the header.

5. Used as initialization for the table parsing task.

    iii.    **FINETUNING**:



FINETUNING    Training set of N examples: $\{(x_i, T_i, y_i)\}_{i=1}^{i=N}$ where,

$x_i$: utterance , $T_i$: table , $y_i$: set of denotations.

Aim is to find a program/model $z$ that when $x_i$ is performed against $T_i$ via $z$, $y_i$ is returned.

PREPROCESSING: Translate $y$ to tuple $(C, s)$ where $C$: cell coordinates
$s$: scalar populated only when $y$ is SCALAR.

Cell selection: $s$ is not populated, C is populated.

Scalar answer: $s$ is populated , $C$ is empty.

e.g. $x_i$: Who is the highest ranking wrestler , $y_i = \{RIC flair\}$.
Preprocessing converts $y_i \rightarrow (\downarrow, \quad)$ empty.
coordinates of cell

$x_i$: What is average time for top 2 wrestlers , $y_i = 3426$
Preprocessing converts $y_i \rightarrow (\quad , 3426)$
empty

Task-1 Cell selection: $\{C \neq \phi, u = \phi\}$

★ Train model to choose column 'col' that has highest # of cells in $C$.

★ Loss = Loss in       + Loss in        + Loss applying
           Choosing column   choosing        NONE operator
                             Cells within column
           $L_{cols}$    +    $L_{cells}$    +    $L_{agg}$

$$= \frac{1}{|column|} \sum_{co \in column\ list} CE\left(\phi_{col}^{(co)}, 1_{co=col}\right) + \frac{1}{|cells(col)|} \sum_{\substack{c \in cells \\ (col)}} CE\left(\phi_s^{(c)}, 1_{c \in C}\right)$$

$$- 1 \log p_a \binom{op_o =}{NONE}$$

---

Task-2 Scalar answer ($u$)  $\{C = \phi, u \neq \phi\}$

★ $y$ is a scalar which does not appear in table. Hence it is assumed to be an o/p of aggregation.

★ OP: COUNT, AVERAGE, SUM.

$\longrightarrow$ calculated using SOFT IMPLEM.

★ $s_{pred} = \sum_{i=1}^{} \hat{p}_a(op_i) \times compute(op_i, p_s, T)$

$\dfrac{p_a(op_i)}{\sum_{i=1} p_a(op_i)}$  (except $op = NONE$)

★ Calculate HUBER loss B/W $s$ and $s_{pred}$. ($L_{scalar}$)
   Also calculate $L_{agg} = -\log\left(\sum_{i=1}^{} p_a(op_i)\right)$  (except $op_o = NONE$)

$$Loss = L_{scalar} + \beta L_{agg}$$

---

| $op$ | $compute(op, p_s, T)$ |
|---|---|
| COUNT | $\sum_{c \in T} p_s^{(c)}$ |
| SUM | $\sum_{c \in T} p_s^{(c)} \cdot T[c]$ |
| AVERAGE | $\dfrac{compute(\text{SUM}, p_s, T)}{compute(\text{COUNT}, p_s, T)}$ |

Task-3 Ambiguous answer $\{C \neq \phi, s \neq \phi\}$

 ↳ Value $\iota$ is available in table.

 ↳ Let model decide whether it wants cell selection based on threshold $S$.

  OR

  Scalar answer

$0 < S < 1$

$$p_a(\phi_0) = p_a(NONE) \geq S \quad, \quad CELL\ SELECTION$$
$$ELSE \quad , \quad SCALAR\ ANSWER$$

iv. **EXPERIMENTATION**:
  1. Experiment on semantic parsing datasets : **WIKITQ, SQA, WIKISQL**
  2. Metric: Denotational accuracy and report median for 5 independent runs.
  3. From ablation studies, it is concluded that **model pretraining** and **column/row input embeddings** are very important

v. **LIMITATION:**
  1. Fails to capture very large tables which cannot fit in memory
  2. Fails to capture questions which have multiple aggregations.
  3. Fails to handle cases when the denotation is text that does not appear in the table

## C. My Assessment/ Analysis
  a. Authors present a new idea of directly extracting answers from tables without intermediate logical form.
  b. Uses Transfer Learning (BERT-large architecture), pre-training on WIKI and fine tuning on specific dataset for TABULAR DATA
  c. Combination of 2 tasks: Cell selection and aggregation.
  d. Authors very cleverly calculate SOFT DIFFERENTIABLE ESTIMATION for each operator instead of pinpointing one single operator. I.e. s_pred is calculated using all the operators.
  e. Can be fine tuned /trained via **weak supervision** where supervision signal is either **Cell selection** OR **Scalar answer.** i.e. **y => (C,s)** via code during preprocessing. The training for 2 tasks as mentioned in b. is done using **(C,s) tuple**

## D. LIMITATIONS/ CONFUSIONS for me

a. Very simplistic assumption of considering cells as independent Bernoulli Random Variables.
b. Adding INDUCTIVE BIAS of choosing a column and then cells within it might not be applicable for all scenarios.
c. Does not seem to work well on complex queries across multiple columns and aggregations.
d. Aggregations are limited to AVERAGE, SUM, MAXIMUM, NONE. Not sure how to extend this.