

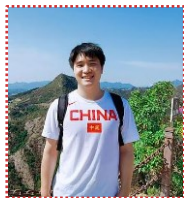
Table Pre-training via Learning a Neural SQL Executor



Qian Liu¹



Bei Chen²



Jiaqi Guo³



Morteza Ziyadi²



Zeqi Lin²



Weizhu Chen²



Jian-Guang Lou²



北京航空航天大学
BEIHANG UNIVERSITY



Microsoft



西安交通大学
XI'AN JIAOTONG UNIVERSITY



10th International Conference on
Learning Representations (ICLR 2022)

Background: Table-based Question Answering

Year	City	Country	Nations
1896	Athens	Greece	14
1900	Paris	France	24
1904	St. Louis	USA	12
...
2004	Athens	Greece	201
2008	Beijing	China	204



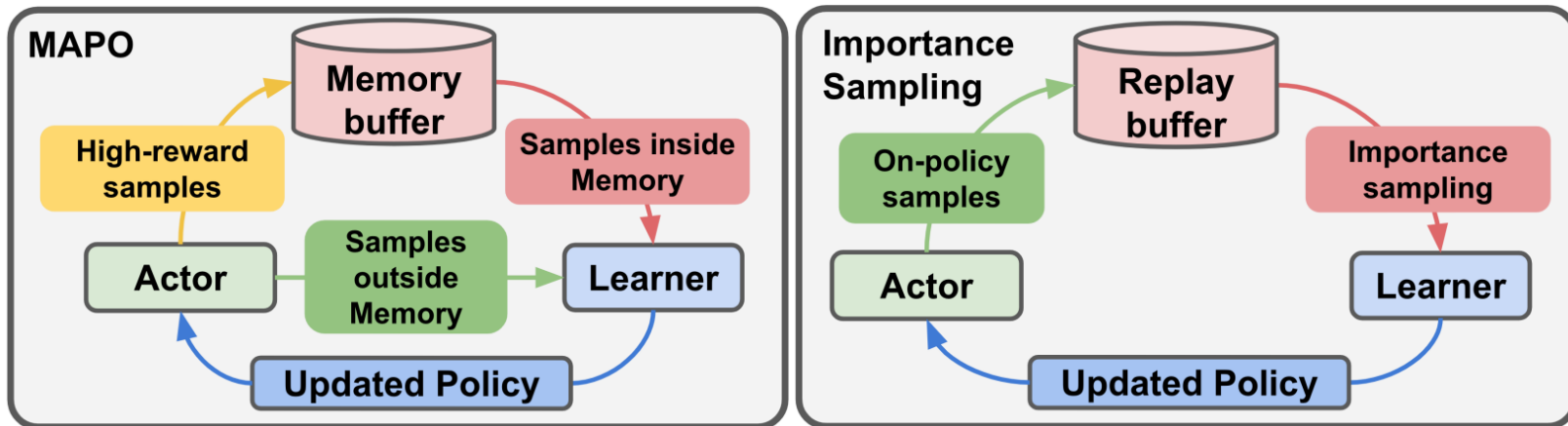
Greece held its last Summer Olympics in which year?

2004



Previous Work: Reinforcement Learning

Obtain rewards by comparing execution results of sampled SQL queries with golden answers to train a text-to-SQL semantic parser. *Hard to scale* to complex scenarios.



[Chen et al. 2018]

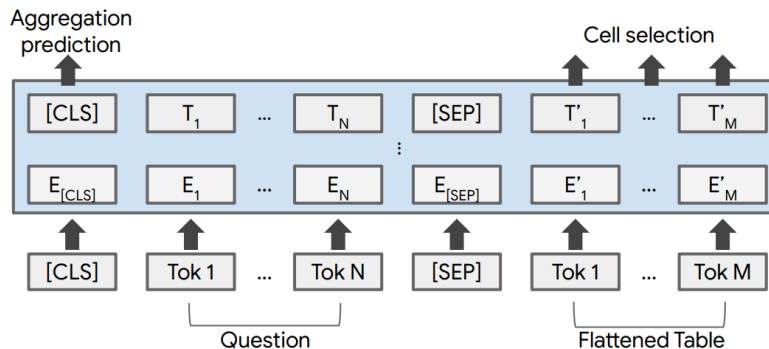
Previous Work: Table Parsing

Predict answer by selecting table cell values and optionally applying an aggregation operator to the selected region. **Flexibility is limited.**

op	$P_s(op)$	compute(op, P_s, T)
NONE	0	-
COUNT	0.1	.9 + .9 + .2 = 2
SUM	0.8	.9×37 + .9×31 + .2×15 = 64.2
AVG	0.1	64.2 ÷ 2 = 32.1

Rank	...	Days	P_s
1	...	37	0.9
2	...	31	0.9
3	...	17	0
4	...	15	0.2
...	0

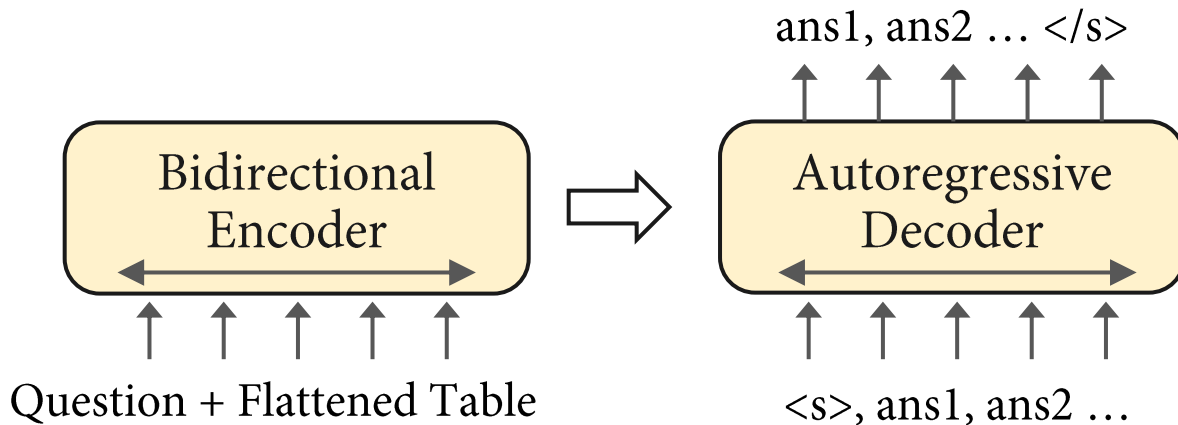
$$s_{\text{pred}} = .1 \times 2 + .8 \times 64.2 + .1 \times 32.1 = 54.8$$



[Herzig et al. 2020]

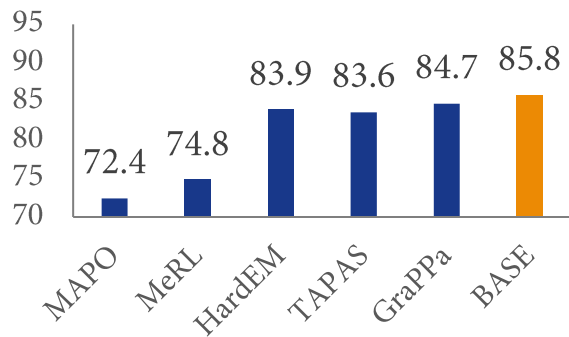
Our Proposal: Generative Language Model

We formulate the task of table-based question answering as answer generation, and leverages generative language models (e.g., BART) to output autoregressively.

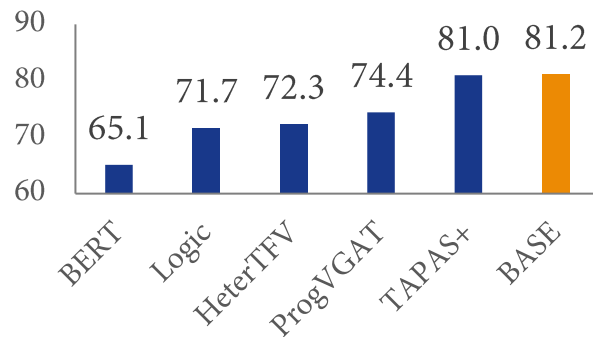


Preliminary Result: Models Are Data-Hungry

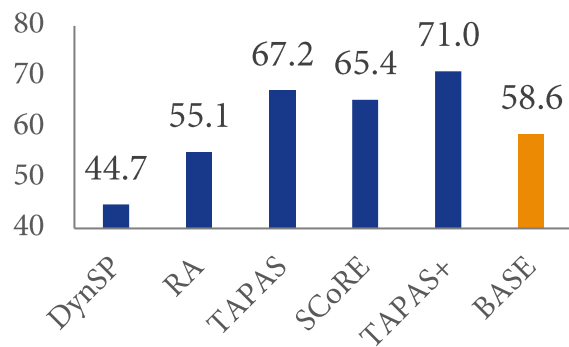
WikiSQL (Weak)



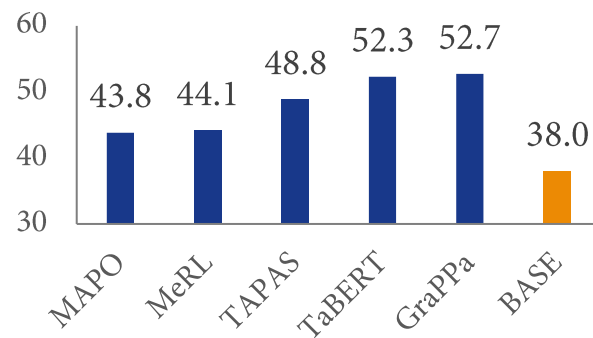
TabFact



SQA



WikiTableQuestions



Motivation: Program as Proxy

Programming Context : Database

Year	City	Country	Nations
1896	Athens	Greece	14
1900	Paris	France	24
1904	St. Louis	USA	12
...
2004	Athens	Greece	201
2008	Beijing	China	204

Program:

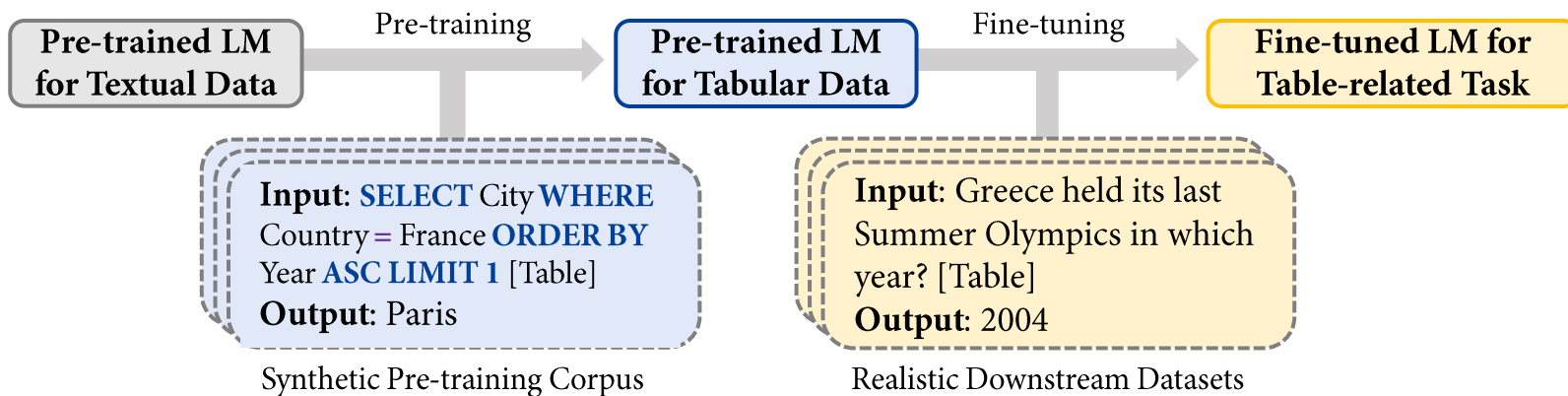
```
SELECT City WHERE  
Country = France ORDER BY  
Year ASC LIMIT 1
```

Executor

Paris

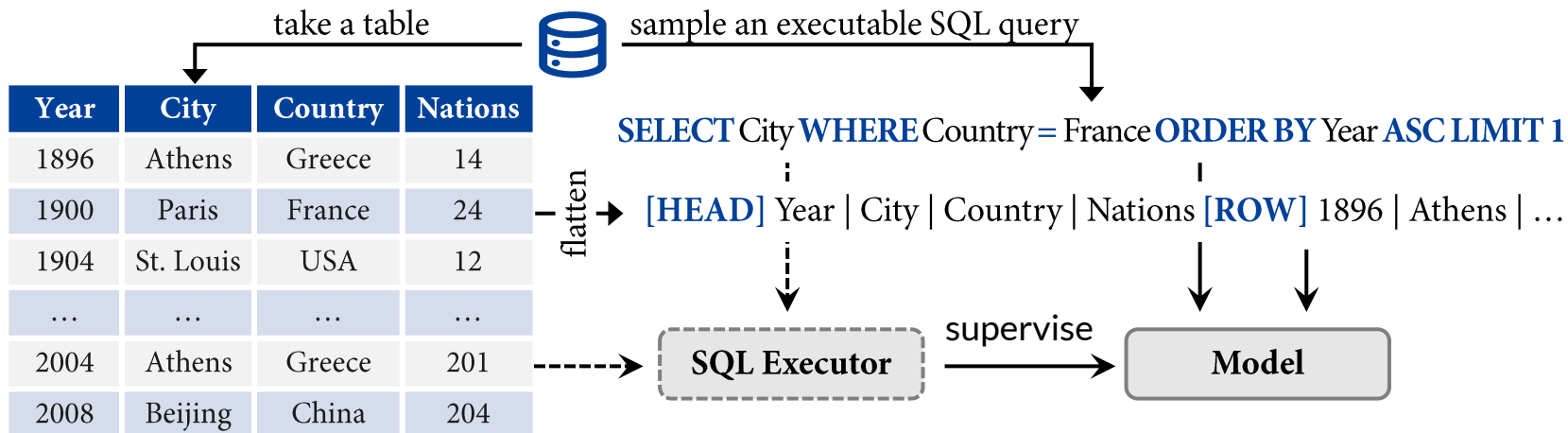
Our Proposal: SQL Execution Pre-training

Pre-training a model to mimic the behavior of a symbolic execution engine.



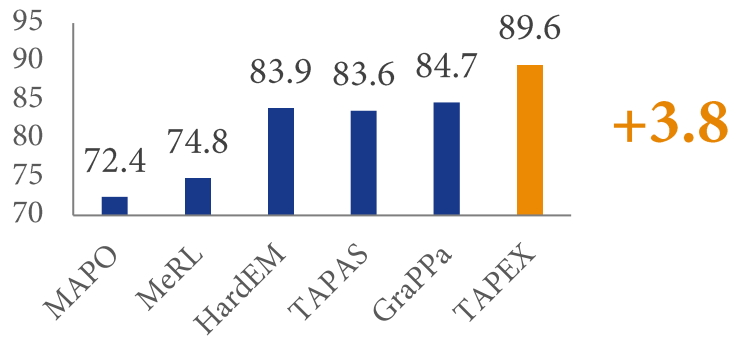
Our Proposal: SQL Execution Pre-training

If we train a model to mimic the SQL query execution procedure over databases, we believe it **learns latent programmatic operations** from the execution engine.

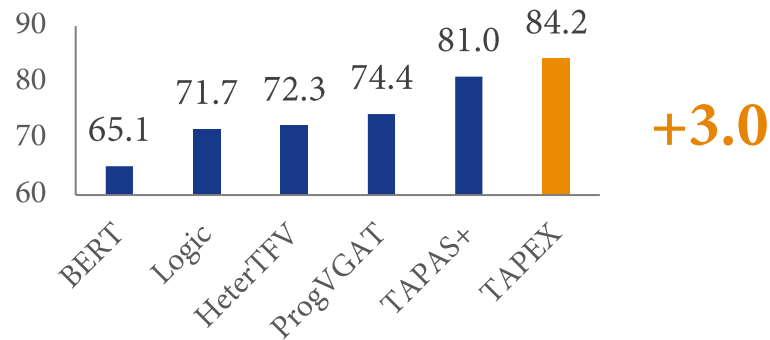


Experimental Result: SOTA Across Benchmarks

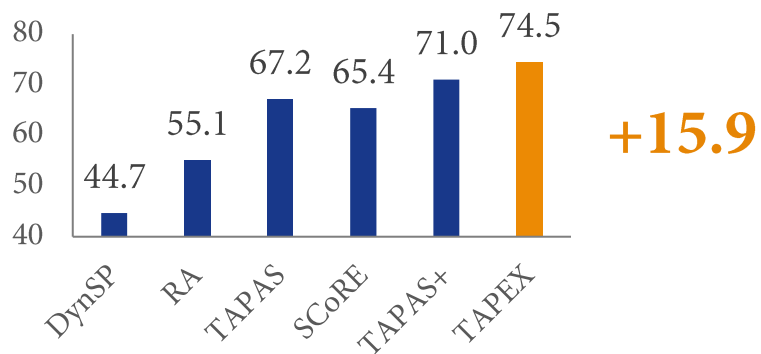
WikiSQL (Weak)



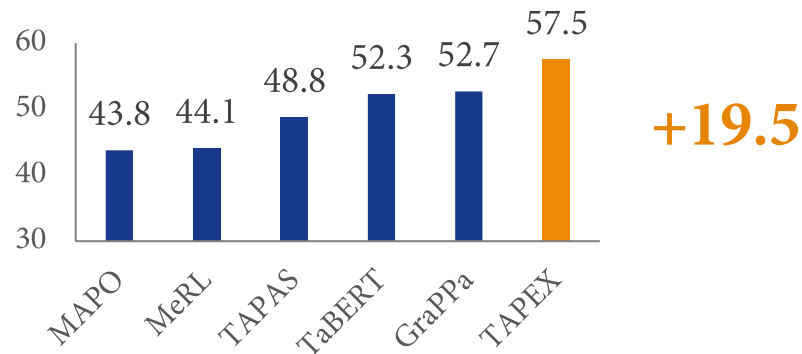
TabFact



SQA

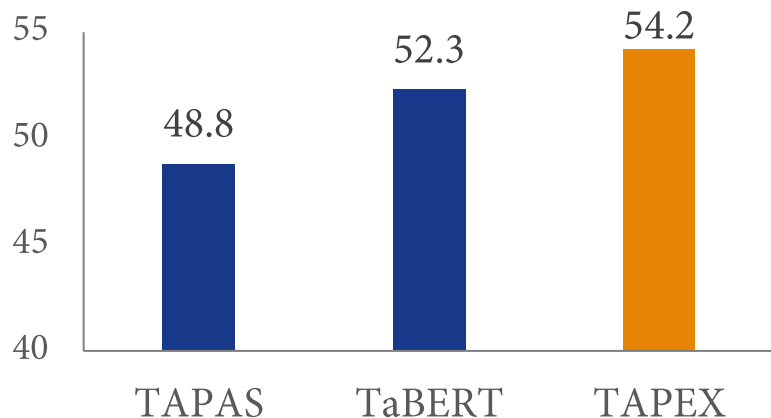


WikiTableQuestions

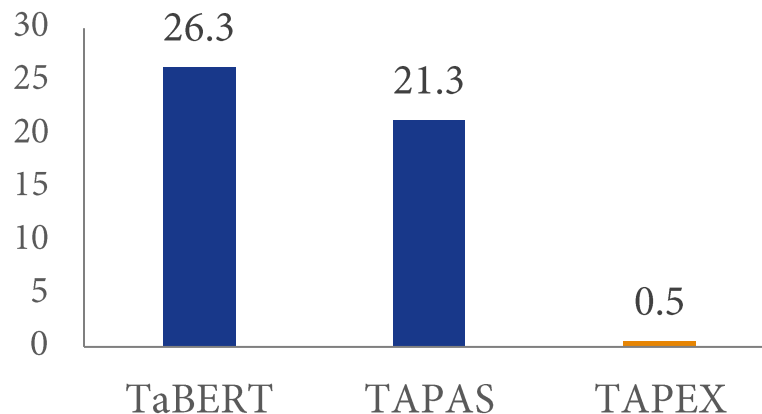


Experimental Result: Cost-Effective Pre-training

Fine-tuning Performance



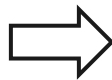
Pre-training Corpus (Million)



Compared with TaBERT, **2%** of corpus yields **2%** improvement!

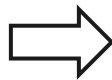
Take Away: Pre-training without Real Data

When performing continual pre-training, instead of mining a large noisy web corpus, we can also try to synthesize an accurate and small corpus.



Take Away: Pre-training without Language Modeling

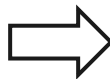
When performing continual pre-training, instead of performing the general-purpose language modeling, we can also try to simulate the specialized skill.



Take Away: Pre-training without Natural Language

When performing continual pre-training for natural language tasks, instead of natural language, we can also try to leverage programs.

We introduce a new language representation model called **BERT**, which stands for **Bidirectional Encoder Representations from Transformers**. Unlike recent language representation models (Peters et al., 2018a; Radford et al., 2018), BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers. As a result, the pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks, such as question answering and language inference, without substantial task-specific architecture modifications.



```
30 def start_link(name, opts \\ [] do
31   GenServer.start_link(__MODULE__, {name}, opts)
32 end
33
34 def init({name}) do
35   require Logger
36   Logger.log :debug, "Started channel #{name}!"
37   :pg2.join(:channels, self)
38   :ets.insert(:channels, {name, self})
39   users = :ets.new(:users, [:set, :protected])
40   {:ok, {name, users, []}}
41 end
42
43 def handle_call({:send, message}, _from, {name, _users, _buffer} = state) do
44   Kaguya.Util.sendPM(name, message)
45   {:reply, :ok, state}
46 end
47
48 def handle_call({:rename_user, {old_nick, new_nick}}, _from, {_name, users, _buffer} = state) do
49   case :ets.lookup(users, old_nick) do
50     [{^old_nick, user}] ->
51       new_user = %{user | nick: new_nick}
52       :ets.delete(users, old_nick)
53       :ets.insert(users, {new_nick, new_user})
54     [] -> :ok
55   end
56   {:reply, :ok, state}
57 end
```