

OPEN QUESTION ANSWERING OVER TABLES AND TEXT

Anonymous authors

Paper under double-blind review

ABSTRACT

In open question answering (QA), the answer to a question is produced by retrieving and then analyzing documents that might contain answers to the question. Most open QA systems have considered only retrieving information from unstructured text. Here we consider for the first time open QA over *both* tabular and textual data and present a new large-scale dataset *Open Table-and-Text Question Answering* (OTT-QA) to evaluate performance on this task. Most questions in OTT-QA require multi-hop inference across tabular data and unstructured text, and the evidence required to answer a question can be distributed in different ways over these two types of input, making evidence retrieval challenging—our baseline model using an iterative retriever and BERT-based reader achieves an exact match score less than 10%. We then propose two novel techniques to address the challenge of retrieving and aggregating evidence for OTT-QA. The first technique is to use “early fusion” to group multiple highly relevant tabular and textual units into a fused block, which provides more context for the retriever to search for. The second technique is to use a cross-block reader to model the cross-dependency between multiple retrieved evidence with global-local sparse attention. Combining these two techniques improves the score significantly, to above 27%.

1 INTRODUCTION

Open question answering considers the problem of retrieving documents from a fixed corpus with a *retriever*, and then analyzes retrieved evidence to provide answers to a given question with a *reader*. Prior open question answering systems focused only on retrieving and reading free-form passages or documents. However, a significant amount of real-world information is stored in other forms, such as semi-structured web tables due to its compact representation to aggregate related information. For example, tables are often used to hold large quantities of related facts, especially numeric facts, such as ‘Career Statistics for Lebron James’. This type of detailed information is found much less frequently in unstructured text. Tables are also commonly used for collections of homogeneous entities or recurring events, like ‘List of Periodic Comets’ or ‘List of Champions League Winners since 1966’. Hence tabular information serves as an excellent complement to textual data, especially in the open setting. Despite these advantages, no previous studies have exploited the millions of web tables to augment their open QA system.

In this paper, we describe the first study to jointly exploit tables and text for open-domain question answering. For this purpose, we construct a new dataset, Open Table-and-Text Question Answering (OTT-QA). OTT-QA is built on the HybridQA dataset (Chen et al., 2020), and like HybridQA, OTT-QA questions are multi-hop questions which require aggregating information from both tables and text to answer. However, unlike HybridQA, OTT-QA requires the system to *retrieve* relevant tables and text — in contrast, in HybridQA, the ground truth tables and textual passages required for each question are given. To produce OTT-QA’s questions, we begin by re-annotating the questions from HybridQA to ‘decontextualize’ them—i.e., we make questions suitable for the open-domain setting so that unique answers can be determined from the question alone, without needing context from the provided text and tables. We then add new questions to remove potential biases. After these steps, OTT-QA contains 45K human-annotated questions that require retrieving and aggregating information over tables and text from the whole Wikipedia. Examples from OTT-QA are depicted in Figure 1. Note the table and passages contain non-overlapping information, and both of them must be understood to answer the question. For example, the question has a low lexical overlap with

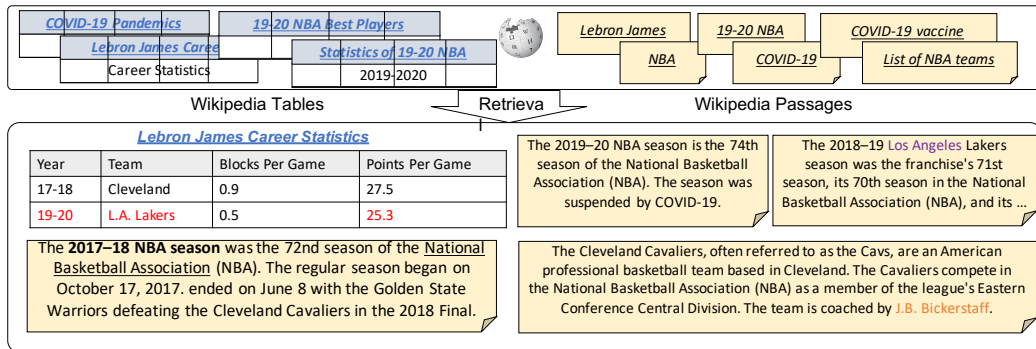


Figure 1: The problem setting: A OTT-QA model needs to retrieve from two candidate pools and then perform multi-hop reasoning to find answers.

the passage about the 'Lakers', and it needs the table as the bridge to retrieve this passage. Such cross-modality multi-hop retrieval features OTT-QA. More examples are displayed in Appendix.

OTT-QA is distinguished from the existing QA datasets in two aspects. Existing table-based QA datasets (Pasupat & Liang, 2015; Yu et al., 2018; Chen et al., 2020) operates in the closed setting without requiring any retrieval, whereas most existing open QA datasets (Joshi et al., 2017; Yang et al., 2018) require only text retrieval, not table retrieval. One dataset, Natural Questions (Kwiatkowski et al., 2019) includes some tabular information in its corpus, but the tables are nearly always of a restricted type (infobox tables with only a single row). In contrast, OTT-QA models require retrieving *both* tabular data and text, and unlike the Natural Questions dataset, requires information fusion from text and tables in non-trivial ways.

OTT-QA poses great challenges to both the retriever and reader in open QA. Retrievers for OTT-QA need to consider two information formats, making the search space larger. Even worse, as questions in OTT-QA often require multi-hop inference, one round of retrieval is often not enough. Readers for OTT-QA also need to aggregate a significant amount of knowledge-intensive information, compared to other reader models: a single table in OTT-QA has an average length of over 300 words. Moreover, readers are often expected to process multiple retrieved units due to the uncertainty in retrieval, which makes it difficult to design strong reader models (Devlin et al., 2019; Liu et al., 2019) with a limited length budget of 512 subword tokens.

The baseline system that we propose to address these challenges uses an iterative retriever (Das et al., 2018; Sun et al., 2019; Qi et al., 2019; Min et al., 2019; Ding et al., 2019; Asai et al., 2019) and a BERT reader (Devlin et al., 2019). The iterative retriever explores multiple evidence documents iteratively, interacting with the candidate pool to gradually reformulate the query. Beam search is used to find multiple subsets of documents that may contain all the required evidence, and each subset is then fed to the BERT reader to predict the answer span. The highest-scored prediction is chosen as the answer. The iterative retriever needs to re-encode the query with a big transformer and re-search over millions of candidates, such a procedure (especially dense) can be computationally expensive. Furthermore, the BERT reader fails to have a global overview of the retrieved document set, which might lead to bad local prediction.

We propose a more sophisticated system that addresses these challenges with two novel strategies: namely *fusion* retrieval and *cross-block* reading. The **fusion retriever** first pre-aligns the table segments to their highly related passages, using entity linking. Then, the aligned table segments and passages are grouped as a *fused block*, which contains aggregated information from two modalities; hence, compared to the previous documents, it contains richer context to benefit the following retrieval. We view the fused block as the basic unit to be retrieved, and instead of performing multiple runs of retrieval iteratively, the fusion retriever is used once to retrieve the top K fused blocks; however, due to errors in fusion and retrieval, the retrieved top-1 fused block might not contain the necessary information. We thus also propose a **cross-block reader** based on a sparse-attention based transformer architecture (Ainslie et al., 2020; Zaheer et al., 2020), which can process extremely long sequences efficiently. We use the cross-block reader to read all the top- K retrieved fused blocks jointly. Both strategies have proven effective compared to the baseline system: the best model combining the two strategies improves the accuracy of the baseline system by a huge margin.

2 BACKGROUND

The aim of an open QA system is to extract an answer to a question q from a given large corpus. Most open QA models are retriever-reader models, which extract answers in two steps: retrieval and reading. In the *retrieval* step, a retrieval model f is used to retrieve a set of passages from the text corpus. In the *reading* step, the reader is then used to extract the answer from them.

Retrieval Function There are two commonly-used types of retrieval function f : sparse retrievers and dense retrievers. Our sparse retriever uses a unigram-based BM-25 score to retrieve an evidence unit b from the candidate pool \mathbb{B} . Our dense retrieval function is a dual-encoder model (Bromley et al., 1994), and we follow (Lee et al., 2019; Guu et al., 2020) for the dual encoder design. The query and the passage are encoded with separate Transformers. As in (Devlin et al., 2019), the vector corresponding to the first token, $[\text{CLS}]$, is used as a “pooled” representation of the sequence. The dense retrieval function is the dot product between $h_q = \text{BERT}_q(q)[\text{CLS}]$ and $h_b = \text{BERT}_B(b)[\text{CLS}]$ for each evidence block b in the candidate corpus—i.e., the scoring function is $f(q, b) = h_q^T h_b$, which can be viewed as finding the nearest neighbor in vector space. In the multi-hop open QA setting (Yang et al., 2018), an iterative retrieval function (Sun et al., 2019; Min et al., 2019; Ding et al., 2019) is proposed, which defines the retrieval process as an auto-regressive formula. Our iterative retriever function is denoted as $f([q, b_1, \dots, b_{j-1}], b_j)$, which appends the previous $j - 1$ rounds of retrieval to the original q in the j -th round of retrieval. Beam search is used in test time.

Single-Block Reader Due to the uncertainty in retrieval, the top-1 document might not always contain the answer. Existing models normally retrieve the top- k documents and feed them to the reader for span selection. The standard reader (Chen et al., 2017; Joshi et al., 2017) aims to extract a span from each of the retrieved blocks b_i and assign a confidence $f(q, b_i) f_{read}(a|q, b_i)$ to it, with $f(q, b_i)$ indicating the retrieval probability and $f_{read}(a|q, b_i)$ denoting the span selection probability by reader. Multiple answers $\{a_1, \dots, a_k\}$ are ranked with this confidence score and the highest scored answer span \hat{a} is the final answer. Note that the reader needs to run k times, once for each of the top- k retrievals. We refer to this model as the *single-block reader* and use it as our baseline.

HybridQA HybridQA (Chen et al., 2020), a closed-domain QA dataset, is the most related to ours. During the annotation of HybridQA, a table T and its relevant passages $\{P_1, \dots, P_N\}$ (surrounding text and hyperlinked passage) are given to a crowd worker to write questions which necessarily require both the passage and table information to answer. The original dataset contains 72K multi-hop questions paired with 13K tables with their paired passages. During training/testing time, the ground-truth tables and passages are given to a model, HYBRIDER, to find the final answer. HYBRIDER also serves as an important baseline in our paper.

3 TASK AND DATASET

In OTT-QA, the retrieval corpus consists of a set of table candidates \mathbb{B}_T and a set of passage candidates \mathbb{B}_P . The task is to answer question q by extracting answer strings from blocks $b \in \mathbb{B}_T \cup \mathbb{B}_P$, where b can be either textual and tabular data. We adopt the standard exact match (EM) and F1 scores (Yang et al., 2018) for evaluation. Different from HybridQA, OTT-QA’s table candidates are web tables *without* hyperlinks provided. This decision was made to make the problem setting more general, as otherwise systems that solve OTT-QA could only be applied to high-quality data in Wikipedia. However, in OTT-QA, we provide hyperlinks in the training subset, but not dev/test set. Removing hyperlinks in tables makes the overall task much more challenging, but makes the final systems applicable to more general domains. Thus, an OTT-QA model needs to jointly retrieve both tables and text, without abusing gold hyperlinks, and then aggregate them to find the answer.

Candidate Pool For our table collection \mathbb{B}_T , we extracted all Wikipedia regular tables with their metadata including page title, page section title, and section text. The metadata, denoted T_M , is essential for de-contextualization. We obtain a table corpus containing over 400k high-quality tables with an average length of 320 words including metadata. For the text passage collection \mathbb{B}_P , we crawl English Wikipedia dump pages and filter out noisy pages. We follow HybridQA (Chen et al., 2020) and only keep a maximum of 12 sentences in the introduction section as the passage. We obtain a corpus containing over 5 million passages, with an average of 94 words.

Notation We define each table as a matrix T , which consists of cells $T_{i,j}$ with i specifying the row, and j specifying the column. Each cell $T_{i,j}$ could be a number, date, phrase or even sentence due to

its semi-structured nature. However, a single complete table with structured representation (Herzig et al., 2020) can easily exceed the 512-token limit, which poses great challenges to the downstream reader to process top- K retrieval. Hence we propose to decompose each table T into multiple rows R_i , which are combined with the headers, metadata, and global max/min information from the original table as a **table segment**. The table segment is used as the basic retrieval block in our paper. This decomposition procedure increases candidate \mathbb{B}_T from 400k to 5 million, making the retrieval problems even more fine-grained and more challenging. Our table segment representation is described in Appendix subsection B.1. In summary, we build a candidate pool of 5 million table segments \mathbb{B}_T and a pool of 5 million passages \mathbb{B}_P . We denote as \mathbb{B} as our full candidate pool, which our model needs to find the block b (a table segment or a passage) containing the answer span.

3.1 QUESTION AND ANSWER ANNOTATIONS

Our question and answer pairs are built upon the existing HybridQA (Chen et al., 2020) dataset, with several significant changes. First, crowd workers ‘decontextualize’ the questions so that they are not under-specified or context-dependent, and thus suitable for the open setting. Second, we add more questions to the development/test set to remove possible annotation bias. During annotation, we adopt strict quality control¹ and more details are described in Appendix section A.1.

Decontextualization Most questions in HybridQA are contextualized with a given table and several passages, with corresponding questions written by crowd workers. Often, the crowd-sourced questions assume the context. For example, the questions might contain the words "the players" because the given table is about "Netherlands players". We thus needed ‘de-contextualize’ (Parikh et al., 2020) the original context-dependent questions, so they could serve as standalone questions, specific enough to imply a unique answer relative to the corpus. To discourage excessive unwanted modification, we enforce a two-step annotation procedure, as depicted in Figure 2. In the first phase, the worker is only allowed to insert **minimum** words or phrases (or replace pronouns) into the questions based on the information presented by Wikipedia Title, Section Title, and Section Text to make the question have a unique answer. After this step, we often potentially obtain overly-complicated questions that are artificial and unnatural. Therefore, we manually selected the worst 25% questions and sent them back to make them more concise and natural.

Title: Netherlands at the European Track Championships					
Section Title: European Track Championships (elite) 2010-current – Medalists					
schema	Medal	Championship	Name	Event	Ranking
content	Silver	2010 Pruszków	Tim Veldt	Men's omnium	2nd
	Bronze	2011 Apeldoorn	Kirsten Wild	Women's omnium	3rd

Original: Which city does the player winning the silver medal in Men's Omnium come from?	↓	Context Insertion
Step1: Which city does the Netherlands player winning the Men's Omnium silver medal in ETC after 2010 come from?	↓	Naturalization
Final: Which city does the Netherlands Men's Omnium silver medalist after 2010 in ETC come from?		

Figure 2: The ‘de-contextualization’ annotation phase of OTT-QA. In the first step, the annotator is restricted to add phrases from the context. In the second step, the annotator is specifically requested to make the sentence more concise and natural.

Additional Evaluation Examples As all the questions from HybridQA are based on the 13k tables from the HybridQA set, no questions are asked about the newly crawled 400k tables. This potentially generates unwanted statistical biases or artifacts for the model to exploit, and potentially biases the final evaluation results. Therefore, we randomly sampled another 1100 tables from the newly crawled tables, and follow the original annotation process used by HybridQA to re-collect 2200 new questions. These new questions were mainly used in the dev/test set. Below we refer to the subset of tables used by original HybridQA as the *in-domain* tables.

Distant Supervision Signals For the in-domain tables ($\approx 8k$), the cell-wise hyperlinks are provided in OTT-QA as a potential signal for supervision. We use $H_{i,j} = \{b_1, b_2, \dots \in \mathbb{B}_P\}$ to denote the hyperlinks in cell $T_{i,j}$. Since in HybridQA the oracle fine-grained answer span is not explicitly

¹The collection is conducted on an established crowdsourcing platform with annotators from countries with English as the native language. The annotators were required to meet the requirements: 1) a native speaker in an English-speaking country 2) having an approval rate of over 95% and 3) having at least 500 approved jobs.

annotated, we approximate this by traversing the table and hyperlinked passages to find all exact matches. This process contains some noise—a manual study reveals that it roughly contains 15% error. We use this ‘weakly-supervised’ fine-grained information to train our models. We denote the ‘approximate’ block of the answer span for answer a as b_a , and use it to train our model.

3.2 DATASET STATISTICS

After annotation, we sampled roughly 2K questions from the in-domain HybridQA dataset, and then mix them with the newly collected out-domain questions to construct our dev and test sets. Finally, we have 41,469 questions in the training set, 2,214 questions in the dev set, and 2,158 questions in the test set. We conduct more in-detailed analysis over the reasoning types and show them in the Appendix A.3, a remarkable difference from original HybridQA is that a proportion of questions actually have multiple plausible inference chains in the open-domain setting.

4 MODEL

Our model for OTT-QA is a retriever-reader model with new designs for both retriever and reader. As discussed briefly above, we propose to use a fusion retriever instead of using a standard iterative retrieve, and we also propose to use cross-block readers to replace a standard single-block reader.

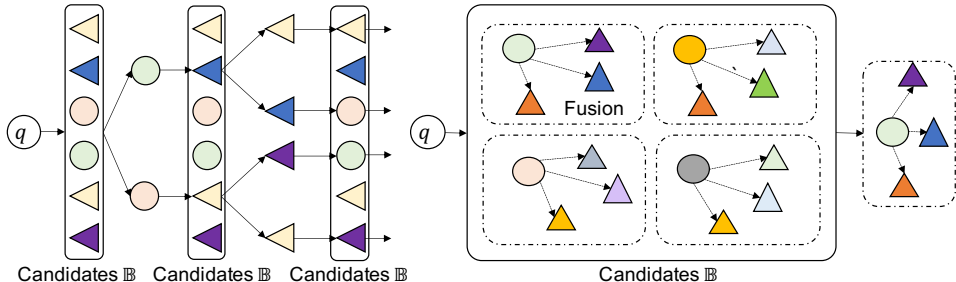


Figure 3: **Left:** Iterative 3-step retrieval (baseline). **Right:** Fusion 1-step retrieval. Circles denote table segments and triangles denotes passages.

4.1 FUSION RETRIEVER

Iterative retrieval (Figure 3, Left) has the following issues. First, iterative retrieval training often requires having supervision signals for every retrieval step to reach good performance, which is not available in OTT-QA. The iterative retrieval also suffers from the problem of error propagation, as early mistakes can propagate to later retrieval stages. Finally, the computation cost for applying a dual-encoder for iterative retrieval is very high, as for every stage, the query embedding has to be re-encoded to include the entire retrieval history.

We propose an alternative strategy to replace multi-step retrieval, namely fusion retrieval (Figure 3, Right). In the fusion retriever, we first use an ‘early fusion’ strategy to group relevant heterogeneous data before retrieval. The fusion procedure groups several highly-relevant blocks from different modalities as a self-contained group (fused block), which provides more clues for the retriever to utilize. Early fusion is very important for retrieving table segments, which often have incomplete context by themselves. The early fusion process aims to fuse a table segment and relevant passages into a group. Here we propose to fuse entities mentioned in a table segment to the appropriate passages for those entities; this is similar to document expansion based on a traditional entity linking step. The problem is challenging due to the mismatch between the lexical forms from the table (which for brevity are often abbreviated) and the relevant passage titles. For example, a cell in the table of "NCAA Division I Men’s Football Tournament" contains the term "Penn State". Directly matching "Penn State" against the passage corpus will lead to "Penn State University" rather than the ground-truth hyperlinked entity, named "Penn State Nittany Lions football". Therefore, we propose an additional augmentation step, which takes in a table segment block b_T and generates a sequence of augmented queries q_1, q_2, \dots, q_n token by token to make the queries more similar to the passage title. The augmented queries are then used to search

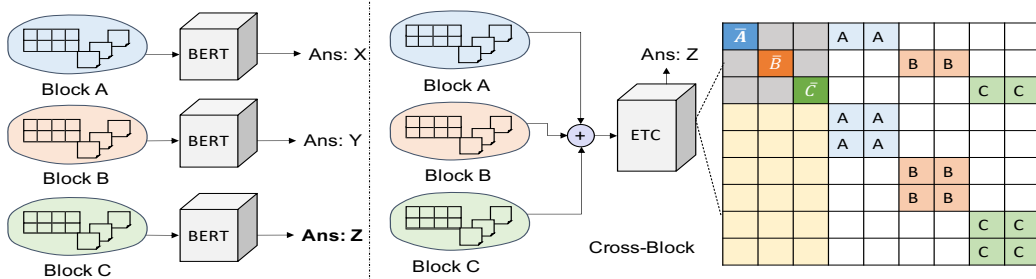


Figure 4: **Left:** Single-block reader (baseline). **Right:** Cross-block reader, with \bar{A} denoting the global state corresponding to local block A.

for nearest neighbors in the passage corpus \mathbb{B}_P using BM25 as the final entity linking step, which is depicted in Appendix. The query augmentation is implemented with a GPT-2 model (Radford et al., 2019), fine-tuned on the supervised pairs of (table segment, hyperlink) from the in-domain tables. Each b_T is fed to find its companions b_P^1, \dots, b_P^n , they are collectively called b_F .

We follow the standard dual-encoder setting (section 2) and the only difference is that we replace the input of the block encoder with $h_b = \text{BERT}_B([b_T, b_P^1, \dots, b_P^n])$, which captures the cross-attention between the table and the text within a block. The fused embedding contains richer context from both modalities to complement each other. The retriever only needs to retrieve once from the candidate pool, which dramatically decreases the complexity compared to the existing iterative retrievers.

To enhance the neural retrieval system to retrieve fused blocks, we apply the Inverse Cloze Task (ICT) (Lee et al., 2019) pretraining task on the corpus of fused blocks. ICT is a way to generate pseudo-training data for dense retrieval. Unlike standard document-wise ICT, our fused block contains both table segments and multiple passages. Given a fused block b_F , we generate the pseudo-query in the following way: 1) we first corrupt the table segment by randomly dropping half of the words from the table metadata and cells to obtain a partial table segment \hat{b}_T . 2) We then randomly sample a sentence \hat{b}_P from the fused passage. We combine \hat{b}_T and \hat{b}_P as a pseudo query \hat{q} and pair it with the original fused block b_F as pre-training data. The pre-training data is applied to enhance the dual encoder’s ability to select lexically matched documents. After pre-training, the retriever is fine-tuned on OTT-QA. Finally, at inference time, the retriever is used to retrieve the top K fused blocks for a question, which are then passed to the reader for answer prediction.

4.2 CROSS-BLOCK READER

The reader typically needs to process the top- k retrieved blocks returned by the retriever to extract the best answer, as the top-1 block might not contain enough evidence to answer the question. As demonstrated in Figure 4, the cross-block reader aims to address this issue by using cross attention between different blocks to model their dependencies. To obtain the cross-block reader, we take the pre-trained long-range sparse attention transformer (ETC) (Ainslie et al., 2020), which can accept up to 4096 tokens as input, and then fine-tune the model on the distant supervision data. During training, the ground truth (fused) blocks are mixed with hard negative blocks from the retriever. We take the top- k retrieval results to fill the 4096 token space (roughly 15 fused blocks).

Cross-attention between blocks allows a much more powerful way to aggregate information across the k retrieved blocks compared to the single-block reader, especially when the blocks are fused. This is feasible because of the design of the sparse attention structure in ETC, which can constrain the attention of each token to its neighboring tokens within a local radius in its local block. Such sparse attention can decrease the attention computation complexity from quadratic $\mathcal{O}(N^2)$ to linear $\mathcal{O}(N|R|)$, where $|R|$ is the local radius (where $N = 4096$ and $|R| = 84$ in our experiments). To allow cross-block interaction, ETC assigns a global state for each local block in the long sequence, and blocks can attention to each other through multiple layers of such global-local structures.

Retriever	Dev-Sparse		Dev-Dense		Test-Best	
	EM	F1	EM	F1	EM	F1
Model						
HYBRIDER (Top-1) (Chen et al., 2020)	8.7	10.9	8.9	11.3	8.4	10.6
HYBRIDER (best Top-K) (Chen et al., 2020)	9.9	12.2	10.3	13.0	9.7	12.8
Iterative-Retrieval + Single-Block Reader	9.8	13.3	7.9	11.1	9.6	13.1
Fusion-Retrieval + Single-Block Reader	14.3	17.8	13.8	17.2	13.4	16.9
Iterative-Retrieval + Cross-Block Reader	17.1	20.7	14.4	18.5	16.9	20.9
Fusion-Retrieval + Cross-Block Reader	27.7	31.8	28.1	32.5	27.2	31.5
† Table-only Retrieval + Cross-Block Reader	4.6	6.9	4.9	7.2	4.4	7.0
† Text-only Retrieval + Cross-Block Reader	8.2	12.4	8.9	12.8	8.8	12.1
† Oracle Link + Fusion-Retrieval + Cross-Block Reader	35.8	40.1	35.2	39.9	35.0	39.5
† Oracle Table + Link (w/o Retrieval) + HYBRIDER	44.1	50.8	44.1	50.8	43.0	49.8

Table 1: **Main Results.** We conduct experiments with both sparse and dense retrievers using the dev set, and then select the best setting to report the test set results (as indicated by the word "Best"). Fusion-Retriever and Cross-Block Reader are combined to obtain the highest score, the scores are median value over 3 runs with different random seeds. † refers to ablation study results.

5 EXPERIMENTS

All of our code is based on Tensorflow (Abadi et al., 2016). For the retriever part, the sparse retriever is built on top of DrQA (Chen et al., 2017) with unigram features, and the dense retriever is built with BERT. The single-block retriever is based on BERT-uncased, and the cross-block reader is based on ETC (Ainslie et al., 2020). Both of them consist of 12 layers with a hidden size of 768, the minor differences in the relative positional embedding used in ETC. All the models are trained with a learning rate of 1e-5 optimized by AdamW (Loshchilov & Hutter, 2019). We use in-batch negatives (Lee et al., 2019) to train our dense retrievers. A more detailed implementation of the baseline iterative retriever is described in Appendix subsection B.2. In fusion retriever, we use the ‘fused’ block containing the ‘approximate’ answer block b_a as the positive instance. In iterative retriever, since the auto-regressive model $f(b_j|q, b_1, \dots, b_{j-1})$ requires fine-grained inference chain for step-wise supervision, which is not given in OTT-QA. We apply lexical match based heuristics to synthesize inference chains as weakly supervised training data (described in Appendix). For both all the dense retrievers, we pre-train with 10K steps using the generated pseudo query and then fine-tune them another 10K step using a batch size of 2048. For the cross-block reader, we fine-tune with a batch size of 64. Both are using 16 cloud TPUs.

Main Results In our experiments, we experiment with different types of retriever and reader models under both sparse and dense setting, the details are described as follows:

- **HYBRIDER:** this model, designed for closed domain HybridQA questions, is one baseline. Since this model requires a ground truth table with its hyperlinks to do modularized reasoning, we use BM25 to retrieve the most relevant table and passages to reconstruct an ‘approximated’ input for this model. We experiment with top-1,2,3,4 cases where we use the answer with the highest confidence as the final result. We also directly feed the ground-truth table and hyperlinks to HYBRIDER, which roughly estimates an upper limit of this task.
- **Iterative-Retriever (Sparse):** We use a 2-step iterative retriever: in the first step, we apply the question to retrieve the top-10 table segments and top-10 passages. In the second step, we use each retrieved table segment to retrieve its related top-5 passages and concatenate each retrieved passage title with the original question to retrieve the top-5 table segments. We merge and calculate the retrieval score of each unique block and rank them by their score. For the single-block reader, we split the retrieved blocks into 512-token chunks and feed them to the BERT reader. For the cross-block reader, we truncate the top 4096 subword tokens and only feed these tokens to reader.
- **Iterative-Retriever (Dense):** We use a 3-step iterative retriever. In the first step, we encode the question and retrieve the top-8 blocks (either table segment or passage); in the second step, we concatenate the previous retrieved block and the question to re-encode the query vector to further retrieve top-4 blocks; similarly, the last step retrieves top-2 blocks.
- **Fusion Retriever (Iterative):** We use a sparse retriever to directly retrieve the top-15 fused blocks

based on bag-of-words BM25 score, and then split it into individual table segments and passage blocks. Since passage could be associated with multiple fused blocks, we merge duplicate blocks and use their summed score. Finally, we rank each block based on its merged retrieval score and truncate the first 4096 subword tokens for the next step.

- Fusion Retriever (Iterative): We use a dual-encoder dense retriever to directly retrieve the top-15 fused blocks, and then follow the same procedure as above. Without specifying the dense retriever uses ICT for pre-training by default.
- Fusion Retriever w/o ICT and w/o GPT-2: these two ablation studies are aimed to show the effectiveness of our proposed ICT pre-training and query augmentation.

The main results are presented in Table 1. First, we can observe that best HYBRIDER top-2 can only achieve a comprised exact match of 9.9% while the oracle HYBRIDER can obtain a score of 44%, which reflects the difficulties of the hybrid retrieval in our dataset. We restrain the retriever to only retrieve table and text to answer the questions and report their results in Table 1, even with the strong cross-block reader, the model only obtains 10% EM. These experiments demonstrate the necessity to integrate information from both forms in OTT-QA.

By combining the standard iterative retriever and single-block reader, we can slightly improve the score can to roughly 10%. By replacing the iterative retriever with the proposed sparse fusion retriever, the EM score can reach 14%, a 4.5% absolute improvement. By replacing the single-block with the proposed cross-block reader, the EM score can reach 17% , a 7% absolute improvement. However, by combining the two strategies, the final EM score can reach 28%, with an 18% absolute improvement, which is greater than the sum of individual improvements. The observation suggests the two components can affect each other in a positive way. We conjecture that the fusion retriever is more likely to retrieve mutually-supportive blocks in a group, which makes the multi-hop reasoning across different blocks easier for the following cross-block reader. In comparison, the iterative retriever retrieves isolated table segments and passages separately, which can easily miss out on the bridging evidence for building the complete reasoning chain. Thus, the cross-block reader cannot maximize its advantage in reasoning across blocks.

By removing the ICT pre-training and query augmentation, we observe that the Dev-EM score drops to 24.6%. By removing the GPT-2 query augmentation, the Dev-EM performance drops to 22.1%. These two results indicate the effectiveness of the proposed two strategies. By replacing the predicted hyperlinks with the oracle links, the fusion model performance can increase by 7% EM. This indicates that there is still plenty of room to improve for the table-passage fusion model.

Linker/Retriever Results To understand the results more, we evaluate the standalone table-passage entity linking accuracy and retriever recall.

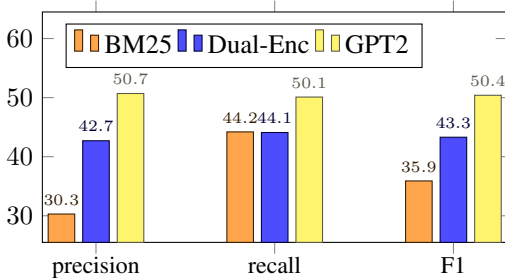


Figure 5: Entity linker performance (F1).

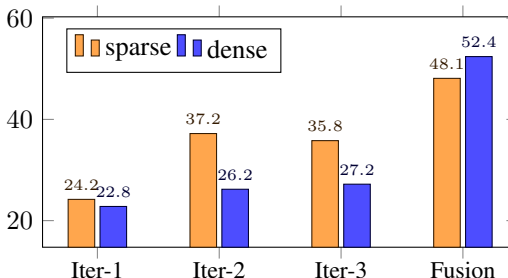


Figure 6: Retriever performance (HITS@4K).

We consider the following linking models: a) BM25 model, which directly uses the cell value to retrieve passages based on their titles without query augmentation, b) a Dual-Encoder model, which encodes the cell value and meta information into a query vector to compute dot-product over all the passage candidate to retrieve, c) a GPT-2 model, which first augments the cell value by the context and then uses BM25. We demonstrate our findings in Figure 5, and evaluate with table-segment-wise F1 score. We observe that directly using BM25 leads to compromised precision of 30.3%, which is mainly due to the lack of context information. By using a dual-encoder retriever, the precision can be improved to 42%. However, many table segments have either zero or multiple linked passages and can be better modeled by an auto-regressive retrieval process.

We use HITS@4K is used to measure the retriever performance, which indicates the chance of ground truth block existing in the retrieved 4096 subword tokens. The results are reported in Figure 6. We vary the steps of iterative retrievers to show the necessity of multi-hop retrieval in OTT-QA. We observe that the 1-step retrieval has the lowest recall because the answer block in OTT-QA normally has a lower lexical overlap with the query. Adding the second retrieval step can greatly improve the recall, but adding the third retrieval hop has very little impact. In contrast to the iterative retriever, the fusion retriever can consistently improve the performance over the iterative setting for both sparse and dense setting. The sparse setting can rise from 35.8% to 48.1% indicating the advantage of ‘early’ fusion. The dense retriever’s improvement is more dramatic (from 27.2% to 52.4%). We believe this is because the iterative retriever heavily relies on noisy synthetic inference chain data, while the fusion retriever does not require such a fine-grained supervision signal, thus less prone to noise. To better understand the retriever, we conduct detailed error analysis in Appendix C.

6 RELATED WORK

Table Retrieval: Tables are pervasive on the Web, there have been some studies on mining web tables to answer open-domain questions (Sun et al., 2016; Chakrabarti et al., 2020). In Sun et al. (2016), the authors have proposed a pipeline framework to first detect the topic entity and then generate a candidate chain, finally ranking chains to predict the answer cell. In Chakrabarti et al. (2020), the authors investigate different similarity matching features to retrieve tables from the web. Our paper is significantly different from these two studies in two aspects: 1) the previous papers use private small-scale datasets while we collect a large-scale dataset and release it for public use, 2) the previous studies are restricted to only using tables as evidence, while our paper considers a more realistic and challenging setting with both table and text corpus. Tables have been a ubiquitous information representation form to express semi-structured information. There has been a long-standing effort to utilize tables in natural language processing applications (Pasupat & Liang, 2015; Zhong et al., 2017; Yu et al., 2018; Parikh et al., 2020; Chen et al., 2019). However, these existing tasks are restricted to in-domain cases without requiring any retrieval, and our paper is the first to investigate retrieving web tables for downstream tasks. Another pair of related works are TAPAS (Herzig et al., 2020) and TABERT (Yin et al., 2020), which investigate joint pre-training over textual and tabular data. Our method draws inspiration from these models, and also uses special tokens and embeddings to encode spatial and logical operations inside tables.

Long Range Transformer: Recently, many transformer variants to resolve the $\mathcal{O}(n^2)$ attention cost have been proposed including Sparse Attention (Child et al., 2019), Reformer (Kitaev et al., 2020), Routing Transformer (Roy et al., 2020), Longformer (Beltagy et al., 2020) and ETC (Ainslie et al., 2020). These different transformer models apply hierarchical architecture, local-sensitive hashing, global-local state to decrease the attention complexity to nearly linear. Our cross-block reader is based on ETC (Ainslie et al., 2020), but unlike prior works that process one long document for QA, our task requires reading multiple blocks containing both structured and unstructured data. To handle the long sequence of retrieved documents in open-domain question answering, Fusion-in-Decoder (Izacard & Grave, 2020) has been proposed to replace the extractive model with an encoder-decoder generative model. The long sequence of passages are split and encoded independently to decrease the computation complexity, but the decoder still uses full attention over the tens of thousands of encoded vectors to generate the answer token by token. Such full-attention can decrease the decoding speed by an order of magnitude, while our sparse-attention-based cross-block reader can still maintain the same speed as the standard BERT model.

7 CONCLUSION

We focus on the problem of performing open question answering over tables and text in this paper. One interesting question we would like to ask in the future is: can we extend open question answering system to more modalities? Some questions can be better answered by images and other resources, but the task can be drastically more challenging by including more modalities, as we have learned from this paper. Finally, we believe the techniques we proposed might be useful for other open-QA setting, especially the comparisons between iterative retriever and fusion retriever.

REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. In *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, pp. 265–283, 2016.
- Joshua Ainslie, Santiago Ontanon, Chris Alberti, Philip Pham, Anirudh Ravula, and Sumit Sanghai. Etc: Encoding long and structured data in transformers. *Proceedings of EMNLP 2020*, 2020.
- Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi, Richard Socher, and Caiming Xiong. Learning to retrieve reasoning paths over wikipedia graph for question answering. In *International Conference on Learning Representations*, 2019.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. Signature verification using a” siamese” time delay neural network. In *Advances in neural information processing systems*, pp. 737–744, 1994.
- Kaushik Chakrabarti, Zhimin Chen, Siamak Shakeri, and Guihong Cao. Open domain question answering using web tables. *arXiv preprint arXiv:2001.03272*, 2020.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1870–1879, 2017.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyu Zhou, and William Yang Wang. Tabfact: A large-scale dataset for table-based fact verification. In *International Conference on Learning Representations*, 2019.
- Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan Xiong, Hong Wang, and William Wang. Hybridqa: A dataset of multi-hop question answering over tabular and textual data. *Proceedings of Findings of EMNLP 2020*, 2020.
- Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. *arXiv preprint arXiv:1904.10509*, 2019.
- Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, and Andrew McCallum. Multi-step retriever-reader interaction for scalable open-domain question answering. In *International Conference on Learning Representations*, 2018.
- Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. Autoregressive entity retrieval. *arXiv preprint arXiv:2010.00904*, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, 2019.
- Bhuvan Dhingra, Manzil Zaheer, Vidhisha Balachandran, Graham Neubig, Ruslan Salakhutdinov, and William W Cohen. Differentiable reasoning over a virtual knowledge base. In *International Conference on Learning Representations*, 2019.
- Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. Cognitive graph for multi-hop reading comprehension at scale. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 2694–2703, 2019.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm: Retrieval-augmented language model pre-training. *Proceedings of ICML 2020*, 2020.
- Jonathan Herzig, Paweł Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Martin Eisenschlos. Tapas: Weakly supervised table parsing via pre-training. *ACL 2020*, 2020.

- Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2020.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1601–1611, 2017.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Ledell Wu, Sergey Edunov, Danqi Chen, and Wentaou Yih. Dense passage retrieval for open-domain question answering. *EMNLP 2020*, 2020.
- Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. *ICLR*, 2020.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. Latent retrieval for weakly supervised open domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6086–6096, 2019.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *ICLR 2019*, 2019.
- Sewon Min, Danqi Chen, Luke Zettlemoyer, and Hannaneh Hajishirzi. Knowledge guided text retrieval and reading for open domain question answering. *arXiv preprint arXiv:1911.03868*, 2019.
- Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. Totto: A controlled table-to-text generation dataset. *arXiv preprint arXiv:2004.14373*, 2020.
- Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1470–1480, 2015.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vassilis Plachouras, Tim Rocktäschel, et al. Kilt: a benchmark for knowledge intensive language tasks. *arXiv preprint arXiv:2009.02252*, 2020.
- Peng Qi, Xiaowen Lin, Leo Mehr, Zijian Wang, and Christopher D Manning. Answering complex open-domain questions through iterative query generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2590–2602, 2019.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI Blog*, 1(8):9, 2019.
- Stephen Robertson and Hugo Zaragoza. *The probabilistic relevance framework: BM25 and beyond*. Now Publishers Inc, 2009.
- Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *arXiv preprint arXiv:2003.05997*, 2020.

- Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William Cohen. Open domain question answering using early fusion of knowledge bases and text. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 4231–4242, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1455. URL <https://www.aclweb.org/anthology/D18-1455>.
- Haitian Sun, Tania Bedrax-Weiss, and William Cohen. PullNet: Open domain question answering with iterative retrieval on knowledge bases and text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2380–2390, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1242. URL <https://www.aclweb.org/anthology/D19-1242>.
- Huan Sun, Hao Ma, Xiaodong He, Wen-tau Yih, Yu Su, and Xifeng Yan. Table cell search for question answering. In *Proceedings of the 25th International Conference on World Wide Web*, pp. 771–782, 2016.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380, 2018.
- Wen-tau Yih, Kristina Toutanova, John C Platt, and Christopher Meek. Learning discriminative projections for text similarity measures. In *Proceedings of the fifteenth conference on computational natural language learning*, pp. 247–256, 2011.
- Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Sebastian Riedel. Tabert: Pretraining for joint understanding of textual and tabular data. *ACL 2020*, 2020.
- Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, et al. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 3911–3921, 2018.
- Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. *arXiv preprint arXiv:2007.14062*, 2020.
- Victor Zhong, Caiming Xiong, and Richard Socher. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv preprint arXiv:1709.00103*, 2017.