

HiTab: A Hierarchical Table Dataset for Question Answering and Natural Language Generation

Anonymous ACL submission

Abstract

Tables are often created with hierarchies, but existing works on table reasoning mainly focus on flat tables and neglect hierarchical tables. Hierarchical tables challenge table reasoning by complex hierarchical indexing, as well as implicit relationships of calculation and semantics. We present a new dataset, HiTab, to study question answering (QA) and natural language generation (NLG) over hierarchical tables. HiTab is a cross-domain dataset constructed from a wealth of statistical reports and Wikipedia pages, and has unique characteristics: (1) nearly all tables are hierarchical, and (2) questions are not proposed by annotators from scratch, but are revised from real and meaningful sentences authored by analysts. (3) to reveal complex numerical reasoning in analyses, we provide fine-grained annotations of quantity and entity alignment. Experiment results show that HiTab presents a strong challenge for existing baselines and a valuable benchmark for future research. Targeting hierarchical structure, we devise an effective hierarchy-aware logical form for symbolic reasoning over tables. Furthermore, we leverage entity and quantity alignment to explore partially supervised training in QA and conditional generation in NLG, and largely reduce spurious predictions in QA and meaningless descriptions in NLG.

1 Introduction

In recent years, there are a flurry of works on reasoning over semi-structured tables, e.g., answering questions over tables (Yu *et al.*, 2018; Pasupat and Liang, 2015) and generating fluent and faithful text from tables (Lebret *et al.*, 2016; Parikh *et al.*, 2020). But they mainly focus on simple flat tables and neglect complex tables, e.g., hierarchical tables. A table is regarded as hierarchical if its header exhibits a multi-level structure (Lim and Ng, 1999;

¹<https://www.nsf.gov/statistics/2019/nsf19319/>

	A	B	C	D	E	F	G
1	TABLE 3. Primary source and mechanism of support for full-time master's and doctoral students in science and engineering: 2017						
2		All full-time graduate students		Master's		Doctoral	
3	Source and mechanism	Total	Percent	All	Percent	All	Percent
4	All full-time	433,916	100.0	209,221	100.0	224,695	100.0
5	Self-support	161,641	37.3	139,373	66.6	22,268	9.9
6	All sources of support	272,275	62.7	69,848	33.4	202,427	90.1
7	Federal	65,999	15.2	10,736	5.1	55,263	24.6
8	Department of Agricu	2,361	0.5	938	0.4	1,423	0.6
9	Department of Defens	8,089	1.9	2,568	1.2	5,521	2.5
16	Other	9,098	2.1	3,462	1.7	5,636	2.5
17	Institutional	182,135	42.0	52,319	25.0	129,816	57.8
18	Other U.S. source	19,432	4.5	5,136	2.5	14,296	6.4
19	Foreign	4,709	1.1	1,657	0.8	3,052	1.4
20	All mechanisms of support	272,275	62.7	69,848	33.4	202,427	90.1
21	Fellowships	39,368	9.1	5,687	2.7	33,681	15.0
22	Traineeships	10,945	2.5	1,497	0.7	9,448	4.2
23	Research assistantships	103,586	23.9	19,702	9.4	83,884	37.3
24	Teaching assistantships	84,499	19.5	22,171	10.6	62,328	27.7
25	Other mechanisms	33,877	7.8	20,791	9.9	13,086	5.8

- Teaching assistantships were most commonly reported as the primary mechanism of support for master's students (11%).

Figure 1: A hierarchical table and accompanied descriptions in a National Science Foundation report.¹

Chen and Cafarella, 2014; Wang *et al.*, 2020). Hierarchical tables are widely used, especially in data products, statistical reports, and research papers in government, finance, and science-related domains.

Hierarchical tables challenge QA and NLG due to: (1) **Hierarchical indexing.** Hierarchical headers, such as D2:G3 and A4:A25 in Figure 1, are informative and intuitive for readers, but make cell selection much more compositional than flat tables, requiring multi-level and bi-dimensional indexing. For example, to select the cell E5 (“66.6”), one needs to specify two top header cells, “Master’s” and “Percent”, and two left header cells, “All full-time” and “Self-support”. (2) **Implicit calculation relationships among quantities.** In hierarchical tables, it is common to insert aggregated rows and columns without explicit indications, e.g., total (columns B,D,F and rows 4,6,7,20) and proportion (columns C,E,G, which challenge precise numerical inference. (3) **Implicit semantic relationships among entities.** There are various cross-row, cross-column, and cross-level entity relationships, but lack explicit indications, e.g., “source” and “mecha-

nism” in A2 describe A6:A19 and A20:A25 respectively, and D2 (“Master’s”) and F2 (“Doctoral”) can be jointly described by a virtual entity, “Degree”. How to identify semantic relationships and link entities correctly is also a challenge.

In this paper, we aim to build a dataset for hierarchical table QA and NLG. But without sufficient data analysts, it’s hard to ensure questions and descriptions are meaningful and diverse (Gururangan *et al.*, 2018; Poliak *et al.*, 2018). Fortunately, large amounts of statistical reports are public from a variety of organizations (StatCan; NSF; Census; CDC; BLS; IMF), containing rich hierarchical tables and textual descriptions. Take Statistics Canada (StatCan) for example, it consists of 6,039 reports in 27 domains authored by over 1,000 professions. Importantly, since both tables and sentences are authored by domain experts, sentences are natural and reflective of real understandings of tables.

To this end, we propose a new dataset, HiTab, for QA and NLG on hierarchical tables. (1) All sentence descriptions of hierarchical tables are carefully extracted and revised by human annotators. (2) It shows that annotations of fine-grained and lexical-level entity linking significantly help table QA (Lei *et al.*, 2020; Shi *et al.*, 2020), motivating us to align entities in text with table cells. In addition to entity, we believe aligning quantities (Ibrahim *et al.*, 2019), especially composite quantities (computed by multiple cells), is also important for table reasoning, so we annotate underlying numerical relationships between quantities in text and table cells, as Table 1 shows. (3) Since real sentences in statistical reports are natural, diverse, and reflective of real understandings of tables, we devise a process to construct QA pairs based on existing sentence descriptions instead of asking annotators to propose questions from scratch.

HiTab presents a strong challenge to state-of-the-art baselines. For the QA task, MAPO (Liang *et al.*, 2018) only achieves 29.2% accuracy due to the ineffectiveness of the logical form customized for flat tables. To leverage the hierarchy for table reasoning, we devise a hierarchy-aware logical form for table QA, which shows high effectiveness. We propose partially supervised training given annotations of linked mentions and formulas, which helps models to largely reduce spurious predictions and achieve 45.1% accuracy. For the NLG task, models also have difficulties in understanding deep hierarchies and generate complex analytical texts.

We explore controllable generation (Parikh *et al.*, 2020), showing that conditioning on both aligned cells and calculation types helps models to generate meaningful texts.

2 Dataset Construction and Analysis

We design an annotation process with six steps. To well-handle the annotation complexity, we recruit 18 students or graduates (13 females and 5 males) in computer science, finance, and English majors from top universities, and provide them with comprehensive online training, documents, and QAs. Labeling totally spends 2,400 working hours, and ethical considerations can be found in Section 8.

2.1 Hierarchical Table Collection

We select two representative organizations, Statistics Canada (StatCan) and National Science Foundation (NSF), that are rich of statistical reports. Different from (Census; CDC; BLS; IMF) that only provide PDF reports where table hierarchies are hard to extract precisely (Schreiber *et al.*, 2017), StaCan and NSF also provide HTML reports, in which cell information such as text and formats can be extracted in precise using HTML tags.

First, we crawl English HTML statistical reports published in recent five years from StatCan (1,083 reports in 27 well-categorized domains) and NSF (208 reports from 11 organizations in science foundation domain). We merge StatCan and NSF and get a total of 28 domains. In addition, ToTTo contains a small proportion (5.03%) of hierarchical tables, so we include them to cover more domains from Wikipedia. To keep the balance between statistical reports and Wikipedia pages, we only randomly include 40% (1,851) of tables in ToTTo. Next, we transform HTML tables to spreadsheet tables using a preprocessing script. Since spreadsheet formula is easy to write, execute, and check, the spreadsheet is naturally a great annotation tool to align quantities and answer questions. To enable correct formula execution, we normalize quantities in data cells by excluding surrounding superscripts, internal commas, etc. Super small or large tables are filtered out (Appendix A.1 gives more details).

2.2 Sentence Extraction and Revision

In this step, annotators manually go through statistical reports and extract sentence descriptions for each table. Sentences consisting of multiple semantic-independent sub-sentences will be care-

Original	After revision	Entity & quantity alignment	Question-answering conversion
Two-thirds (67%) of master's students and only one-tenth (10%) of doctoral students were self-supported (table 3).	Two-thirds (67%) of master's students and only one-tenth (10%) of doctoral students were self-supported.	two-thirds (67%) → =E5% master's → =D2 one-tenth (10%) → =G5% self-supported → =A5	What are the percentages of master's students and doctoral students who are self-supported? =E5, =G5
Teaching assistantships were most commonly reported as the primary mechanism of support for master's students (11%).	Teaching assistantships were most commonly reported as the primary mechanism of support for master's students (11%).	teaching assistantships → =A24 mechanism of support → =A20 master's → =D2 11% → =E24%	Which is the primary mechanism of support for master's students? =XLOOKUP(MAX(E21:E24), E21:E24, A21:A24)
For doctoral students, the proportion of support from research assistantships is 10 points higher than that from teaching assistantships.	For doctoral students, the proportion of support from research assistantships is 10 points higher than that from teaching assistantships.	doctoral → =F2 proportion → =E3 research assistantships → =A23 10 points → =G23-G24 teaching assistantships → =A24	For doctoral students, what is the difference between the proportions of research assistantships and teaching assistantships? =G23-G24

Table 1: Examples of the annotation process. All sentences describe the table in Figure 1.

fully split into multiple ones. Annotators are instructed to eliminate redundancy and ambiguity in sentences through revisions including decontextualization and phrase deletion like (Parikh *et al.*, 2020). Fortunately, most sentences in statistical reports are clean and fully supported by table data, so few revisions are needed to get high-quality text.

2.3 Entity and Quantity Alignment

In this phase, annotators are instructed to align mentions in text with corresponding cells in tables. It has two parts, entity alignment and quantity alignment, as shown in Table 1. For entity alignment, we record the mappings from entity mentions in text to corresponding cells. Single-cell quantity mentions can be linked similar with entity mentions, but composite quantity mentions are calculated from two or more cells through operators like *max/sum/div/diff* (Table 2). The spreadsheet formula is powerful and easy-to-use for tabular data calculation, so we use the formula to record the calculations process of composite quantities in text, e.g., ‘10 points higher’ (=G23-G24). Although quantities are often rounded in descriptions, we neglect rounding and refer to precise quantities in table cells.

2.4 Converting Sentences to QA Pairs

Existing QA datasets instruct annotators to propose questions from scratch, but it’s hard to guarantee the meaningfulness and diversity of proposed questions. In HiTab, we simply revise declarative sentences to QA pairs. For each sentence, annotators need to identify a target key part to question about (according to the underlying logic), then convert

Operators	Formula template (ranges are placeholders)
opposite, percent	=-A5, =B2%
kth-argmax/argmin	=XLOOKUP(SMALL(D1:D3, k), D1:D3, A1:A3)
pair-argmax/argmin	=IF(B1>B2, A1, A2)
sum, average	=SUM(D2:D4), =AVERAGE(D2:D4)
max, count	=MAX(D2:D4), =COUNT(D2:D4)
product, diff, div	=D3*D4, =D3-D4, =D3/D4

Table 2: Example operators and formula templates.

it to the QA form. All questions are answered by formulas that reflect the numerical inference process. For example, the ‘XLOOKUP’ operator is frequently used to retrieve the header cells of superlatives, as shown in Table 1. To keep sentences as natural as they are, we do not encourage unnecessary sentence modification during the conversion. If an annotator finds multiple ways to question regarding a sentence, she only needs to choose one way that best reflects the overall meaning.

2.5 Regular Inspections and the Final Review

We ask two most experienced annotators to perform regular inspections and the final review. (1) In the labeling process, they regularly sample annotations (about 10%) from all annotators to give timely feedback on labeling issues. (2) Finally, they review all annotations and fix labeling errors. Also, to assist the final review, we write a script to automatically identify spelling issues and formula issues. To double check the labeling quality before the final review, we study the agreement of annotators by collecting and comparing annotations on a randomly sampled 50 tables from two annotators. It shows 0.89 and 0.82 for quantity and entity alignment in Fleiss Kappa respectively, which are regarded as “almost perfect agreement” (Landis and Koch, 1977), and 64.5 in BLEU-4 after sentence revision, which also indicates high agreement.

2.6 Hierarchy Extraction

We follow existing work (Lim and Ng, 1999; Chen and Cafarella, 2014; Wang *et al.*, 2020) and use the tree structure to model hierarchical headers. Since cell formats such as merging, indentation, and font bold, are commonly used to present hierarchies, we adapt heuristics in (Wang *et al.*, 2020) to extract top and left hierarchical trees, which has high accuracy. We go through 100 randomly sampled tables in HiTab, 94% of them are precisely extracted. Figure 7 in Appendix shows an illustration.

Dataset	Tables	Data source			Fine-grained alignment		QA and NLG tasks				
		Table	Question or sentence	Real sentences revised per table	Entity	Quantity	QA	NLG	Questions	Words per question	Sentences
WTQ (Pasupat and Liang, 2015)	2,108	Wikipedia	Post-created	-	-	-	-	22,033	10.0	-	
WikiSQL (Zhong et al., 2017)	26,521	Wikipedia	Post-created	-	-	-	-	80,654	11.7	-	
Spider (Yu et al., 2018)	1,020	College data, WikiSQL	Post-created	-	-	-	-	10,181	13.2	-	
HybridQA (Chen et al., 2020b)	13,000	Wikipedia	Post-created	-	-	-	-	69,611	18.9	-	
TAT-QA (Zhu et al., 2021)	2,757	Financial reports (PDF)	Post-created	-	-	-	-	16,552	12.5	-	
FinQA (Chen et al., 2021)	2,776	Financial reports (PDF)	Post-created	-	-	-	-	8,281	16.6	-	
DART (Nan et al., 2020)	5,623	WTQ, WikiSQL, ...	Post-created	-	-	-	Yes	-	-	82,191	
LogicNLG (Chen et al., 2020a)	7,392	Wikipedia	Post-created	-	-	-	Yes	-	-	37,015	
ToTTo (Parikh et al., 2020)	83,141	Wikipedia	Pre-existing	1.4	-	-	Yes	-	-	120,000	
NumericNLG (Suadaa et al., 2021)	1,300	Scientific papers (ACL)	Pre-existing	3.8	-	-	Yes	-	-	4,756	
HiTab	3,597	Stat. reports, Wiki.	Pre-existing	5.0 (reports)	Yes	Yes	Yes	10,686	16.5	10,686	

Table 3: Dataset statistics and comparison.

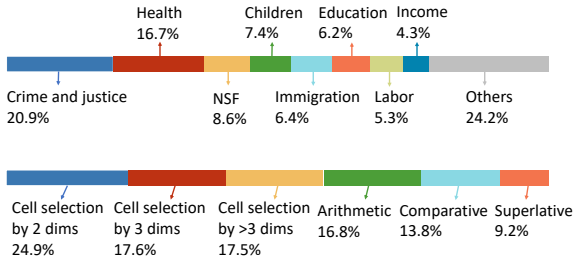


Figure 2: Distribution of domains and operations in StatCan and NSF. *Cell selection by k dims* means that header cells in k levels are used in cell selection.

	A	B	C	D	E
1	Table 4: Quantity and contribution of selected beverages to nutrient intake by year				
3		Total beverages		Skim, 1% or 2% milk	
4		2004	2015	2004	2015
5	19 to 50 years, male				
6	Quantity (grams)	2,458.0	2,279.0	164.0	97
7	Proportion of energy (%)	18.1	15.6	2.9	2.0
8	Proportion of vitamin C (%)	46.6	41.5	1.2	0.2
15	19 to 50 years, female				
16	Quantity (grams)	2,169.0	1,813.0	152.0	87.0
17	Proportion of energy (%)	16.2	12.9	3.6	2.4
18	Proportion of vitamin C (%)	41.9	33.4	1.2	0.1
25	51 to 70 years				
35	71 years or older				
51	What is the percentage points change in daily energy intake from total beverage among adults aged 19 to 50 between 2004 and 2015 ?				
53	$=(B6*B7\%+B16*B17\%)/(B6+B16)-(C6*C7\%+C16*C17\%)/(C6+C16)$				

Figure 3: A meaningful but challenging case in HiTab.

2.7 Dataset Statistics and Comparison

Table 3 shows a comprehensive comparison of related datasets. HiTab is not among the largest ones, but (1) it is the first dataset to study table reasoning over hierarchical tables (accounting for 98.1% tables in HiTab); (2) it is annotated with fine-grained entity and quantity alignment; (3) compared with TAT-QA, FinQA, and NumericNLG that are single-domain, HiTab is cross-domain; (4) the number of real descriptions per table (5.0) in statistical reports (HiTab) is much richer than 1.4 in Wikipedia (ToTTo) and 3.8 in scientific papers, contributing more analytical aspects per table.

Figure 2 analyzes this dataset by domains and operations: domains are diverse, covering 28 domains from statistical reports (fully listed in Appendix A.2) and other open domains from Wikipedia; a large proportion of questions involves complex cell selection and numerical operations.

3 Hierarchical Table QA

Table QA is essential for table understanding, document retrieval, ad-hoc search, *etc.* Hierarchical tables are quite common in these scenarios like in webpages and reports, while current Table QA tasks and methods focus on simple flat tables.

Problem Statement Hierarchical Table QA is defined as follows: given a hierarchical table t and a question x in natural language, output answer y . The question-answer pair should be fully supported by the table. Our dataset $D = \{(x_i, t_i, y_i)\}$, $i \in [1, N]$ is a set of N question-table-answer triples.

Table QA is usually formulated as a semantic parsing problem (Pasupat and Liang, 2015; Liang et al., 2017), where a parser converts questions into logical forms, and an executor executes it to produce the answer. However, existing logical forms for Table QA (Pasupat and Liang, 2015; Liang et al., 2017; Yin et al., 2020) are customized for flat or database tables. The three challenges mentioned in Section 1 make QA more difficult on hierarchical tables, *i.e.*, hierarchical indexing, implicit calculation and semantic relationships.

3.1 Hierarchy-aware Logical Forms

To this end, we propose a hierarchy-aware logical form that exploits table hierarchies to mitigate these challenges. Specifically, we define *region* as the operating object, and propose two functions for hierarchical region selection.

Definitions Given tree hierarchies of tables extracted in Section 2.6, we define *header* as a header cell (e.g., A7(“Federal”) in Figure 1), and *level* as a level in the left/top tree (e.g., A5, A6, A20 are on the same level). Existing logical forms on tables treat rows as operating objects and columns as attributes, and thus can not perform arithmetic operations on cells in the same row. However, a row in hierarchical tables is not necessarily a subject or record, thus operations can be applied on cells in the same row. Motivated by this, we define *region* as our operating object, which is a data region in table

indexed by both left and top headers (e.g., B6:C19 is a rectangular region indexed by A6:B2). The logical form execution process is divided into two phases: region selection and region operation.

Region Selection We design two functions (*filter_tree* h) and (*filter_level* l) to do region selection, where h is a header, l is a level. Functions can be stringed up: the subsequent function applies on the return region of the previous function. (*filter_tree* h) selects a sub-tree region according to a header cell h : if h is a leaf header (e.g., A8), the selected region should be the row/column indexed by h (row 8); if h is a non-leaf header (e.g., A7), the selected region should be the rows/columns indexed by both h and its children headers (row 7-16). (*filter_level* l) selects a sub-tree from the input tree according to a level l and return the sub-region indexed by headers on level l . These two functions mitigate aforementioned three challenges: (1) hierarchical indexing is achieved by applying these two functions sequentially; (2) with *filter_level*, data with different calculation types (e.g., rows 4-5) will not be co-selected, thus not incorrectly operated together; (3) level-wise semantics can be captured by aggregating header cell semantics (e.g., embeddings) on this level. Some logical form execution examples are shown in Appendix B.2.

Region Operation Operators are applied on the selected region to produce the answer. We define 19 operators, mostly following MAPO (Liang et al., 2018), and further include some operators (e.g., *difference rate*) for hierarchical tables. Complete logical form functions are shown in Appendix B.1.

3.2 Experimental Setup

3.2.1 Baselines

We present baselines in two branches. One is logical form-based semantic parsing, and the other is end-to-end table parsing without logical forms.

Neural Symbolic Machine (Liang et al., 2017) is a powerful semantic parsing framework consisting of a programmer to generate programs from NL and save intermediate results, and a computer to execute programs. We replace the LSTM encoder with BERT (Devlin et al., 2018), and implement a lisp interpreter for our logical forms as executor. Table is linearized by placing headers in level order, which is shown in detail in Figure 7.

TaPas (Herzig et al., 2020) is a state-of-the-art end-to-end table parsing model without generating logical forms. Its power to select cells and reason over

tables is gained from its pretraining on millions of tables. To fit TaPas input, we convert hierarchical tables into flat ones following WTQ (Pasupat and Liang, 2015). Specifically, we unmerge the cells spanning many rows/columns on left/top headers and duplicate the contents into unmerged cells. The first top header row is specified as column names.

3.2.2 Weak Supervision

In weak supervision, the model is trained with QA pairs, without golden logical forms. For NSM, we compare three widely-studied learning paradigms.

MML (Dempster et al., 1977) maximizes marginal likelihood of observed programs. **REINFORCE** (Williams, 1992) maximizes the reward of on-policy samples. **MAPO** (Liang et al., 2018) learns from programs both inside and outside buffer and samples efficiently by systematic exploration.

All methods require consistent programs for learning or warm start. We randomly search 15000 programs per sample before training. The pruning rules are shown in Appendix B.5. Finally, 6.12 consistent programs are found for each sample.

For TaPas, we use the pre-trained version and follow its weak supervised training process on WTQ.

3.2.3 Partial Supervision

Given labeled entity links, quantity links, and calculations (from the formula), we further explore to guide training in a *partially supervised* way. These three annotations indicate selected headers, region, and operators in QA. For NSM, we exploit them to prune spurious programs, *i.e.*, incorrect programs that accidentally produce correct answers, in two ways. (1) When searching consistent programs, besides producing correct answers, programs are required to satisfy at least two constraints. In this way, the average consistent programs reduces from 6.12 to 2.13 per sample. (2) When training, satisfying each condition will add 0.2 to the original binary 0/1 reward. Sampled programs with reward $r \geq 1.4$ are added to the program buffer.

For TaPas, we additionally provide answer coordinates and calculation types in training following its WikiSQL setting.

3.2.4 Evaluation Metrics

We use *Execution Accuracy* (EA) as our metric following (Pasupat and Liang, 2015), measuring the percentage of samples with correct answers. We also report *Spurious Program Rate* to study the percentage that incorrect logical forms produce cor-

Weak Supervision			
Method	Dev	Test	%Spurious
MAPO <i>w.</i> original logical form	31.9	29.2	-
TaPas <i>w/o.</i> logical form	39.7	38.9	-
MML <i>w.</i> h.a. logical form	38.9	36.7	22.7
REINFORCE <i>w.</i> h.a. logical form	42.7	38.4	39.3
MAPO <i>w.</i> h.a. logical form	43.5	40.7	19.0
Partial Supervision			
TaPas <i>w/o.</i> logical form	41.2	40.1	-
MML <i>w.</i> h.a. logical form	45.4	45.1	10.3
REINFORCE <i>w.</i> h.a. logical form	44.0	39.7	23.9
MAPO <i>w.</i> h.a. logical form	44.8	44.3	10.7

Table 4: QA execution accuracy (*EA*) on dev/test and spurious program rate of 150 samples on dev. *h.a.* stands for *hierarchy-aware*.

rect answer. Since we do not have golden logical forms, we manually annotate logical forms for 150 random samples in dev set for evaluation.

3.2.5 Implementations

We split 3,597 tables into train (70%), dev (15%) and test (15%) with no overlap. We download pre-trained models from huggingface². For NSM, we utilize ‘bert-base-uncased’, and fine-tune 20K steps on HiTab. Beam size is 5 for both training and inference. To test MAPO original logical form, we convert flatten tables as we do for TaPas. For TaPas, we adopt the PyTorch (Paszke et al., 2019) version in huggingface. We utilize ‘tapas-base’, and fine-tune 40 epochs on HiTab. All experiments are conducted on a server with four V100 GPUs.

3.3 Results

Table 4 summarizes our evaluation results.

Weak Supervision First, MAPO with our hierarchy-aware logical form outperforms that using its original logical form by a large margin 11.5%, indicating the necessity of designing a logical form leveraging hierarchies. Second, MAPO achieves the best *EA* (40.7%) with the lowest spurious rate (19%). But >50% questions are answered incorrectly, proving QA on HiTab is challenging. Third, though TaPas benefits from pretraining on tables, it performs worse than the best logical form-based method without table pretraining. Detailed level-wise results are shown in Appendix B.4.

Partial Supervision From Table 4, we can conclude the effectiveness of partial supervision in two aspects. First, it improves *EA*. The model learns how to deal with more cases given high-quality programs. Second, it largely lowers *%Spurious*. The model learns to generate correct programs instead of some tricks. MML, whose performance highly

depends on the quality of searched programs, benefits the most (36.7% to 45.1%), indicating partial supervision improves the quality of consistent programs by pruning spurious ones. However, TaPas does not gain much improvements from partial supervision, which we will discuss in error analysis.

Error Analysis For TaPas, 98.7% of success cases are cell selections, which means TaPas benefits little from partial supervision. This may be caused by: (1) TaPas does not support some common operators on hierarchical table like *difference*; (2) the coarse-to-fine cell selection strategy first selects columns then cells, but cells in different columns may also aggregate in hierarchical tables.

For MAPO under partial supervision, we analyze 100 error cases. Error cases fall into four categories: (1) entity missing (23%): the header to *filter* is not mentioned in question, where a common case is omitted *Total*; model failure, including (2) failing to select correct regions (38%) and (3) failing to generate correct operations (20%); (4) out of coverage (19%): question types unsolvable with the logical form, which is explained in Appendix B.1.

Spurious programs occur mostly in two patterns. In cell selection, there may exist multiple data cells with correct answers (e.g., G9,G16 in Figure 1), while only one is golden. In superlatives, the model can produce the target answer by operating on different regions (e.g., in both region B21:B25 and B23:B25, B23 is the largest).

4 Hierarchical Table to Text

4.1 Problem Statement

Some works formulate table-to-text as a summarization problem (Lebret et al., 2016; Wiseman et al., 2017). However, since a full table often contains quite rich information, there lack explicit signals on what to generate and renders the task unconstrained and the evaluation difficult. On the other hand, some recent works propose *control-table* generation to enable more specific and logical generation: (1) LogicNLG generates a sentence conditioned on a logical form guiding symbolic operations over given cells, but writing correct logical forms as conditions is challenging for common users who are more experienced to write natural language directly, thus restricting the application to real scenario; (2) ToTTo generates a sentence given a table as well as a set of highlighted cells. In ToTTo’s formulation, the condition of cell selection is much easier to specify than the logical

²<https://huggingface.co/transformers/>

form, but it neglects symbolic operations which are critical for generating some analytical sentences concerning numerical reasoning in HiTab.

We place HiTab as a middle-ground of ToTTo and LogicNLG to make the task more controllable than ToTTo and closer to real application than LogicNLG. In our setting, given a table, the model generates a sentence conditioned on a group of selected cells (similar to ToTTo) and operators (much easier to be specified than logical forms). Although we use two strong conditions to guide symbolic operations over cells, there still leaves a considerable amount of content planning to be done by the model, such as retrieving contextual cells in a hierarchical table given selected cells, identifying how operators are applied on given cells, and composing sentences in a faithful and logical manner.

We now define our task as: given a hierarchical table T , highlighted cells C , and specified operators O , generating a faithful description S . The dataset $H = (T_i, S_i), i \in [1, N]$ is a set of N table-description instances. Description S_i is a sentence about a table T_i and involves a series of operations $O_i = [O_{i1}, O_{i2}, \dots, O_{in}]$ on certain table cells $C_i = [c_{i1}, c_{i2}, \dots, c_{im}]$.

4.2 Controlled Generation

4.2.1 With Highlighted Cells

An entity or quantity in text can be supported by table cells if it is directly stated in cell contents, or can be logically inferred by them. Different from only taking data cells as highlighted cells (Parikh *et al.*, 2020), we also take header cells as highlighted cells, and it is usually the case for superlative ARG-type operations on a specific header level in hierarchical tables, e.g., “Teaching assistantships” is retrieved by ARGMAX in Figure 1. In our dataset, highlighted cells are extracted from annotations of the entity and quantity alignment.

4.2.2 With Operators

Highlighted cells can tell the target for text generation, but is not sufficient, especially for analytical descriptions involving cell operations in HiTab. So we introduce to use operators as extra control. It contributes to text clarity and meaningfulness in two ways. 1) It clarifies the numerical reasoning intent on cells. For example, given the same set of data cells, applying SUM, AVERAGE, or COUNT conveys different meanings thus should yield different texts. 2) Operation results on highlighted

cells can be used as additional input sources. Existing seq2seq models are not powerful enough to do arithmetic operations (Thawani *et al.*, 2021), e.g., adding up a group of numbers, and it greatly limits their ability to generate correct numbers in sentences. Explicitly pre-computing calculation results is a promising alternative way to mitigate this gap in seq2seq models.

4.2.3 Sub Table Selection and Serialization

Sub Table Selection Under controls of selected cells and operators, we devise a heuristic to retrieve all contextual cells as a sub table. (1) we start with highlighted cells extracted from our entity and quantity alignment, then use the extracted table hierarchy to group the selected cells into the top header, the left header, and the data region. (2) based on the extracted table hierarchy, we use the source set of top and left header cells to include their indexed data cells, and we also use the source set of data cells to include corresponding header cells. (3) we leverage the table hierarchy to include their parent header cells to construct a full set of headers. In the end, we take the union of of them as the result of sub table selection.

Serialization On each sub table, we do a row-turn traversal on linked cells and concatenate their cell strings using [SEP] tokens. Operator tokens and calculation results are also concatenated with the input sequence. We also experimented with other serialization methods, such as header-data pairing or template-based method, yet none reported superiority over the simple concatenation. Appendix C.1 gives an illustration.

4.3 Experiments

We conduct experiments by fine-tuning four state-of-the-art text generation methods on HiTab.

Pointer Generator (See *et al.*, 2017) A LSTM-based seq2seq model with copy mechanism. While originally designed for text summarization, it is also used in data-to-text (Gehrmann *et al.*, 2018).

BERT-to-BERT (Rothe *et al.*, 2020) A transformer encoder-decoder model (Vaswani *et al.*, 2017) initialized with BERT (Devlin *et al.*, 2018).

BART (Lewis *et al.*, 2019) A pre-trained denoising autoencoder with standard Transformer-based architecture and shows effectiveness in NLG.

T5 (Raffel *et al.*, 2019) A transformer-based pre-trained model. It converts all textual language problems into text-to-text and proves to be effective.

Model	Cell Highlight		Cell & Calculation	
	BLEU-4	PARENT	BLEU-4	PARENT
Pointer-Generator	5.8	8.8	9.0	10.8
BERT-to-BERT	11.4	16.7	11.7	15.4
BART	17.9	28.0	23.8	31.4
T5	19.5	35.7	26.6	36.9

Table 5: Results of hierarchical-table-to-text.

4.3.1 Evaluation Metrics

We use two automatic metrics, BLEU and PARENT. BLEU (Papineni *et al.*, 2002) is broadly used to evaluate text generation. PARENT (Dhingra *et al.*, 2019) is proposed specifically for data-to-text evaluation that additionally aligns n-grams from the reference and generated texts to the source table.

4.3.2 Experiment Setup

Samples are split into train (70%), dev (15%), and test (15%) sets just the same as the QA task. The maximum length of input/output sequence is set to 512/64. Implementation details of all baselines are given in Appendix C.2.

4.3.3 Experiment Result and Analysis

As shown in Table 5, **first**, from an overall point of view, both metrics are not scored high. This well proves the difficulty of HiTab. It could be caused by the hierarchical structure, as well as statements with logical and numerical complexity. **Second**, by comparing two controlled scenarios (cell highlights & both cell highlights and operators), we see that add operators to conditions greatly help models to generate descriptions with higher scores, showing the effectiveness of our augmented conditional generation setting. **Third**, results on two controlled scenarios across baselines are quite consistent. Replacing the traditional LSTM with transformers shows large increasing. Leveraging seq2seq-like pretraining yields a rise of +6.5 BLEU and +11.3 PARENT. Lastly, between pretrained transformers, T5 reports higher scores over BART, probably for T5 is more extensively tuned during pre-training.

Further, to study the generation difficulty concerning **table hierarchy**, we respectively evaluate samples at different hierarchical depths, *i.e.*, table’s maximum depths in top and left header trees. In groups of 2, 3, 4+ depth, BLEU scores 31.7, 26.5, and 21.3; PARENT scores 40.9, 36.5, and 31.6. The reason could be that, as table headers grow deeper, data indexing is more compositional, so it’s harder for baselines to identify entity relationships and compose logical sentences.

Method	Test Acc.	
	BLEU	PARENT
MAPO <i>w.</i> partial supervision	32.6	
T5 <i>w.</i> cell & calculation	16.9	28.8

Table 6: Results of cross-domain evaluation.

5 Related Work

Table-to-Text Existing datasets are restricted in flat tables or specific subjects (Liang *et al.*, 2009; Chen and Mooney, 2008; Wiseman *et al.*, 2017; Novikova *et al.*, 2016; Banik *et al.*, 2013; Lebrete *et al.*, 2016; Moosavi *et al.*, 2021). The most related table-to-text dataset to HiTab is ToTTo (Parikh *et al.*, 2020), in which complex tables are also included. There are two main differences between HiTab and ToTTo: (1) in ToTTo, hierarchical tables only account for a small proportion (5%), and there are no indication and usage of table hierarchies. (2) in addition to cell highlights, Hitab conditions on operators that reflect symbolic operations on cells.

Table QA mainly focuses on DB tables (Wang *et al.*, 2015; Yu *et al.*, 2018; Zhong *et al.*, 2017) and semi-structured flat tables (Pasupat and Liang, 2015; Sun *et al.*, 2016). Recently, there are some datasets on domain-specific table QA (Chen *et al.*, 2021; Zhu *et al.*, 2021) and jointly QA over tables and texts (Chen *et al.*, 2020b; Zhu *et al.*, 2021), but hierarchical tables still have not been studied in depth. HiTab explores QA on hierarchical tables.

6 Discussion

HiTab also presents cross-domain and complicated-calculation challenges. (1) To explore cross-domain generalizability, we randomly split train/dev/test by domains for three times and present the average results of our best methods in Table 6. We found decreases in all metrics in QA and NLG. (2) Figure 3 shows a case that challenges existing methods: performing complicated calculations needs to jointly consider quantity relationships, header semantics, and hierarchies.

7 Conclusion

We present a new dataset, HiTab, that simultaneously supports QA and NLG on hierarchical tables. Importantly, we provide fine-grained annotations on entity and quantity alignment. We introduce baselines and conduct comprehensive experiments. Results suggest that HiTab can serve as a challenging and valuable benchmark for future research.

8 Ethical Considerations

This work presents HiTab, a free and open English dataset for the research community to study table question-answering and table-to-text over hierarchical tables. Our dataset contains well-processed tables, annotations (QA pairs, target text, and bidirectionally mappings between entities and quantities in text and the corresponding cells in table), recognized table hierarchies, and source code. Data in HiTab are collected from two public organizations, StatCan and NSF. Both of them allow sharing and redistribution of their public reports, so there is no privacy issue. We collect tables and accompanied descriptive sentences from StatCan and NSF. We also include hierarchical tables in Wikipedia. We recruit 18 students or graduates in computer science, finance, and English majors from top universities(13 females and 5 males). Each student is paid \$7.8 per hour (above the average local payment of similar jobs), totally spending 2,400 hours. We finally get 3,597 tables and 10,686 well-annotated sentences. The details for our data collection and characteristics are introduced in Section 2.

References

Eva Banik, Claire Gardent, and Eric Kow. The kbgen challenge. In *the 14th European Workshop on Natural Language Generation (ENLG)*, pages 94–97, 2013.

BLS. U.s. bureau of labor statistics. <https://www.bls.gov> Accessed July 4, 2021.

CDC. Centers for disease control and prevention. <https://www.cdc.gov> Accessed July 4, 2021.

Census. Census bureau. <https://www.census.gov>. Accessed July 4, 2021.

Zhe Chen and Michael Cafarella. Integrating spreadsheet data via accurate and low-effort extraction. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1126–1135, 2014.

David L Chen and Raymond J Mooney. Learning to sportscast: a test of grounded language acquisition. In *Proceedings of the 25th international conference on Machine learning*, pages 128–135, 2008.

Wenhu Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and William Yang Wang. Logical natural language generation from open-domain tables. *arXiv preprint arXiv:2004.10404*, 2020.

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. *arXiv preprint arXiv:2004.07347*, 2020.

Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, et al. Finqa: A dataset of numerical reasoning over financial data. *arXiv preprint arXiv:2109.00122*, 2021.

Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1):1–22, 1977.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William W Cohen. Handling divergent reference texts when evaluating table-to-text generation. *arXiv preprint arXiv:1906.01081*, 2019.

Sebastian Gehrmann, Falcon Z Dai, Henry Elder, and Alexander M Rush. End-to-end content and plan selection for data-to-text generation. *arXiv preprint arXiv:1810.04700*, 2018.

Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R Bowman, and Noah A Smith. Annotation artifacts in natural language inference data. *arXiv preprint arXiv:1803.02324*, 2018.

Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Mueller, Francesco Piccinno, and Julian Eisenschlos. Tapas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, 2020.

Yusra Ibrahim, Mirek Riedewald, Gerhard Weikum, and Demetrios Zeinalipour-Yazti. Bridging quantities in tables and text. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pages 1010–1021. IEEE, 2019.

IMF. International monetary fund. <https://www.imf.org>. Accessed July 4, 2021.

J Richard Landis and Gary G Koch. The measurement of observer agreement for categorical data. *biometrics*, pages 159–174, 1977.

Rémi Lebret, David Grangier, and Michael Auli. Neural text generation from structured data with application to the biography domain. *arXiv:1603.07771*, 2016.

762	Wenqiang Lei, Weixin Wang, Zhixin Ma, Tian Gan,	Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann,	816
763	Wei Lu, Min-Yen Kan, and Tat-Seng Chua. Re-	Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and	817
764	examining the role of schema linking in text-to-sql.	Dipanjan Das. Totto: A controlled table-to-text gen-	818
765	In <i>Proceedings of the 2020 Conference on Empirical</i>	eration dataset. <i>arXiv preprint arXiv:2004.14373</i> ,	819
766	<i>Methods in Natural Language Processing (EMNLP)</i> ,	2020.	820
767	pages 6943–6954, 2020.		
768	Mike Lewis, Yinhan Liu, Naman Goyal, Mar-	Panupong Pasupat and Percy Liang. Compositional	821
769	jan Ghazvininejad, Abdelrahman Mohamed, Omer	semantic parsing on semi-structured tables. <i>arXiv</i>	822
770	Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart:	<i>preprint arXiv:1508.00305</i> , 2015.	823
771	Denosing sequence-to-sequence pre-training for		
772	natural language generation, translation, and com-	Adam Paszke, Sam Gross, Francisco Massa, Adam	824
773	prehension. <i>arXiv preprint arXiv:1910.13461</i> ,	Lerer, James Bradbury, Gregory Chanan, Trevor	825
774	2019.	Killeen, Zeming Lin, Natalia Gimelshein, Luca	826
775	Percy Liang, Michael I Jordan, and Dan Klein. Learn-	Antiga, et al. Pytorch: An imperative style, high-	827
776	ing semantic correspondences with less supervision.	performance deep learning library. <i>Advances in neu-</i>	828
777	In <i>Proceedings of the Joint Conference of the 47th</i>	<i>ral information processing systems</i> , 32:8026–8037,	829
778	<i>Annual Meeting of the ACL and the 4th International</i>	2019.	830
779	<i>Conference on Natural Language Processing</i>		
780	<i>of the AFNLP</i> , pages 91–99, 2009.	Adam Poliak, Jason Naradowsky, Aparajita Haldar,	831
781	Chen Liang, Jonathan Berant, Quoc Le, Kenneth D For-	Rachel Rudinger, and Benjamin Van Durme. Hy-	832
782	bus, and Ni Lao. Neural symbolic machines: Learn-	pothesis only baselines in natural language infer-	833
783	ing semantic parsers on freebase with weak supervi-	ence. <i>arXiv preprint arXiv:1805.01042</i> , 2018.	834
784	sion. In <i>Proceedings of the 55th Annual Meeting of</i>		
785	<i>the Association for Computational Linguistics (Vol-</i>	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	835
786	<i>ume 1: Long Papers)</i> , volume 1, pages 23–33, 2017.	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	836
787	Chen Liang, Mohammad Norouzi, Jonathan Berant,	Wei Li, and Peter J Liu. Exploring the limits of trans-	837
788	Quoc V Le, and Ni Lao. Memory augmented pol-	fer learning with a unified text-to-text transformer.	838
789	icy optimization for program synthesis and semantic	<i>arXiv:1910.10683</i> , 2019.	839
790	parsing. In <i>Advances in Neural Information Process-</i>		
791	<i>ing Systems</i> . Curran Associates, Inc., 2018.	Sascha Rothe, Shashi Narayan, and Aliaksei Severyn.	840
792	Seung-Jin Lim and Yiu-Kai Ng. An automated ap-	Leveraging pre-trained checkpoints for sequence	841
793	proach for retrieving hierarchical data from html ta-	generation tasks. <i>Transactions of the Association for</i>	842
794	bles. In <i>Proceedings of the eighth international con-</i>	<i>Computational Linguistics</i> , 8:264–280, 2020.	843
795	<i>ference on Information and knowledge management</i> ,		
796	pages 466–474, 1999.	Sebastian Schreiber, Stefan Agne, Ivo Wolf, Andreas	844
797	Nafise Sadat Moosavi, Andreas Rücklé, Dan Roth, and	Dengel, and Sheraz Ahmed. Deepdesrt: Deep learn-	845
798	Iryna Gurevych. Learning to reason for text genera-	ing for detection and structure recognition of tables	846
799	tion from scientific tables. <i>arXiv:2104.08296</i> , 2021.	in document images. In <i>2017 14th IAPR interna-</i>	847
800	Linyong Nan, Dragomir Radev, Rui Zhang, Amrit	<i>tional conference on document analysis and recog-</i>	848
801	Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xian-	<i>nition (ICDAR)</i> , volume 1, pages 1162–1167. IEEE,	849
802	gru Tang, Aadit Vyas, Neha Verma, Pranav Krishna,	2017.	850
803	et al. Dart: Open-domain structured data record to	Abigail See, Peter J Liu, and Christopher D	851
804	text generation. <i>arXiv preprint arXiv:2007.02871</i> ,	Manning. Get to the point: Summarization	852
805	2020.	with pointer-generator networks. <i>arXiv preprint</i>	853
806	Jekaterina Novikova, Oliver Lemon, and Verena Rieser.	<i>arXiv:1704.04368</i> , 2017.	854
807	Crowd-sourcing nlg data: Pictures elicit better data.	Tianze Shi, Chen Zhao, Jordan Boyd-Graber, Hal	855
808	<i>arXiv preprint arXiv:1608.00339</i> , 2016.	Daumé III, and Lillian Lee. On the potential of	856
809	NSF. National science foundation. https://www.nsf.gov .	lexico-logical alignments for semantic parsing to sql	857
810	Accessed July 4, 2021.	queries. <i>arXiv:2010.11246</i> , 2020.	858
811	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	StatCan. Statistics canada. https://www150.statcan.gc.ca .	859
812	Jing Zhu. Bleu: a method for automatic evalua-	Accessed July 4, 2021.	860
813	tion of machine translation. In <i>Proceedings of the</i>	Lya Hulliyayatus Suadaa, Hidetaka Kamigaito, Kotaro	861
814	<i>40th annual meeting of the Association for Compu-</i>	Funakoshi, Manabu Okumura, and Hiroya Taka-	862
815	<i>tational Linguistics</i> , pages 311–318, 2002.	mura. Towards table-to-text generation with numer-	863
		ical reasoning. In <i>Proceedings of the 59th Annual</i>	864
		<i>Meeting of the Association for Computational Lin-</i>	865
		<i>guistics and the 11th International Joint Conference</i>	866
		<i>on Natural Language Processing (Volume 1: Long</i>	867
		<i>Papers)</i> , pages 1451–1465, 2021.	868

869 Huan Sun, Hao Ma, Xiaodong He, Wen-tau Yih, Yu Su,
870 and Xifeng Yan. Table cell search for question an-
871 swering. In *Proceedings of the 25th International*
872 *Conference on World Wide Web*, pages 771–782,
873 2016.

874 Avijit Thawani, Jay Pujara, Pedro A Szekely, and Filip
875 Ilievski. Representing numbers in nlp: a survey and
876 a vision. *arXiv preprint arXiv:2103.13136*, 2021.

877 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
878 Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz
879 Kaiser, and Illia Polosukhin. Attention is all you
880 need. *arXiv preprint arXiv:1706.03762*, 2017.

881 Yushi Wang, Jonathan Berant, and Percy Liang. Build-
882 ing a semantic parser overnight. In *Proceedings*
883 *of the 53rd Annual Meeting of the Association for*
884 *Computational Linguistics and the 7th International*
885 *Joint Conference on Natural Language Processing*
886 *(Volume 1: Long Papers)*, pages 1332–1342, 2015.

887 Zhiruo Wang, Haoyu Dong, Ran Jia, Jia Li, Zhiyi
888 Fu, Shi Han, and Dongmei Zhang. Structure-aware
889 pre-training for table understanding with tree-based
890 transformers. *arXiv:2010.12537*, 2020.

891 Ronald J Williams. Simple statistical gradient-
892 following algorithms for connectionist reinforce-
893 ment learning. *Machine learning*, 8(3-4):229–256,
894 1992.

895 Sam Wiseman, Stuart M Shieber, and Alexander M
896 Rush. Challenges in data-to-document generation.
897 *arXiv preprint arXiv:1707.08052*, 2017.

898 Pengcheng Yin, Graham Neubig, Wen-tau Yih, and Se-
899 bastian Riedel. Tabert: Pretraining for joint under-
900 standing of textual and tabular data. *arXiv preprint*
901 *arXiv:2005.08314*, 2020.

902 Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga,
903 Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingn-
904 ing Yao, Shanelle Roman, et al. Spider: A large-
905 scale human-labeled dataset for complex and cross-
906 domain semantic parsing and text-to-sql task. *arXiv*
907 *preprint arXiv:1809.08887*, 2018.

908 Yuchen Zhang, Panupong Pasupat, and Percy Liang.
909 Macro grammars and holistic triggering for efficient
910 semantic parsing. *arXiv preprint arXiv:1707.07806*,
911 2017.

912 Victor Zhong, Caiming Xiong, and Richard Socher.
913 Seq2sql: Generating structured queries from
914 natural language using reinforcement learning.
915 *arXiv:1709.00103*, 2017.

916 Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao
917 Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and
918 Tat-Seng Chua. Tat-qa: A question answering
919 benchmark on a hybrid of tabular and textual content
920 in finance. *arXiv preprint arXiv:2105.07624*, 2021.

A More Details on Dataset

A.1 Dataset Preprocessing

We filter tables using these constraints: (1) number of rows and columns are more than 2 and less than 64; (2) cell strings have no more than one non-ASCII character and 20 tokens; (3) hierarchies are successfully parsed via the method in 2.6. (4) hierarchies have no more than four levels. Finally, 85% tables meet all constraints.

A.2 Domain Distribution

The full 29 domains of sample distribution in HiTab are shown in Figure 4.

A.3 Annotation Interface

The annotation interface looks like Figure 8. Since spreadsheet formula is easy to write, execute, and check, the spreadsheet is naturally a great annotation tool. Annotators can use the Excel formula conveniently for cell linking and calculation in entity alignment and answering questions.

B Hierarchical Table QA

B.1 Logical Form Function List

We list our logical form functions in Table 9.

Union selection is required for comparative and arithmetic operations. It is achieved by allowing variable number of headers in *filter_tree*, where “variable” is one or two in practice.

In our implementation, a function by default takes the selected region of last function as input region to prune search space. We use grammars to filter left headers before top headers, and a (*filter_level*) is necessary after filtering one direction of tree even when only the leaf level is available. And we deactivate order relation functions (e.g., *eq* function) and the order argument *k* in *argmax/argmin* because there are few questions in these types and activating them will largely increase number of spurious programs when searching.

The logical form coverage after deactivation is 78.3% in 300 iterations of random exploration. Some typical question types that can not be covered are: (1) scale conversion, e.g., 0.984 to 98.4%, (2) operating data indexed by different levels of headers, e.g., proportion of total, (3) complex composite operations, e.g., Figure 3.

Question	Logical Forms
Cell Selection	(filter_tree 2012)
Q: What is the GDP of China in 2012?	(filter_tree china) (filter_level LEFT_2) (filter_tree gdp) (filter_level TOP_1)
Superlative	(filter_tree 2012)
Q: Which country has the highest GDP in 2012?	(filter_level LEFT_2) (filter_tree gdp) (filter_level TOP_1) (argmax 1)
Arithmetic	(filter_tree 2013)
Q: How much more is U.S. GDP higher than China in 2013?	(filter_tree u.s. china) (filter_level LEFT_2) (filter_tree gdp) (filter_level TOP_1) (difference)

Table 7: Examples of our logical form. The table to be questioned is in Fig. 7. *LEFT_1* is a symbol for the first level on the left.

B.2 Examples of Logical Form Execution

Take the table in Figure 7 as input table, we demonstrate three types of questions with complete logical forms in Table 7.

B.3 Table Linearization

We linearize the question and table according to Figure 7.

The input is concatenation of question and table. Table is linearized by putting headers in level order. Each level is led by a [*LEVEL*] token to gather current level embedding. The first [*LEVEL*] token stands for level zero of left. Each header is linearized as *name* | *type*. *name* is the tokenized header string. *type* is the entity type parsed by Stanford CoreNLP, which includes “string”, “number”, “datetime” in our case. Headers with the same *name* will gather token embeddings by mean pooling.

B.4 More Experiment Results

In Figure 5, we present level-wise accuracy of HiTab QA with MAPO and our hierarchy-aware logical form. The *Level* in table means sum of left header levels and top header levels. The QA accuracy degrades when table level increases when table structure becomes more complex, except for tables level = 2, i.e., tables with no hierarchies. The reason level = 2 performs comparatively worse is that only 1.9% tables with hierarchies are seen in HiTab, and thus number of training samples for level = 2 is relatively small.

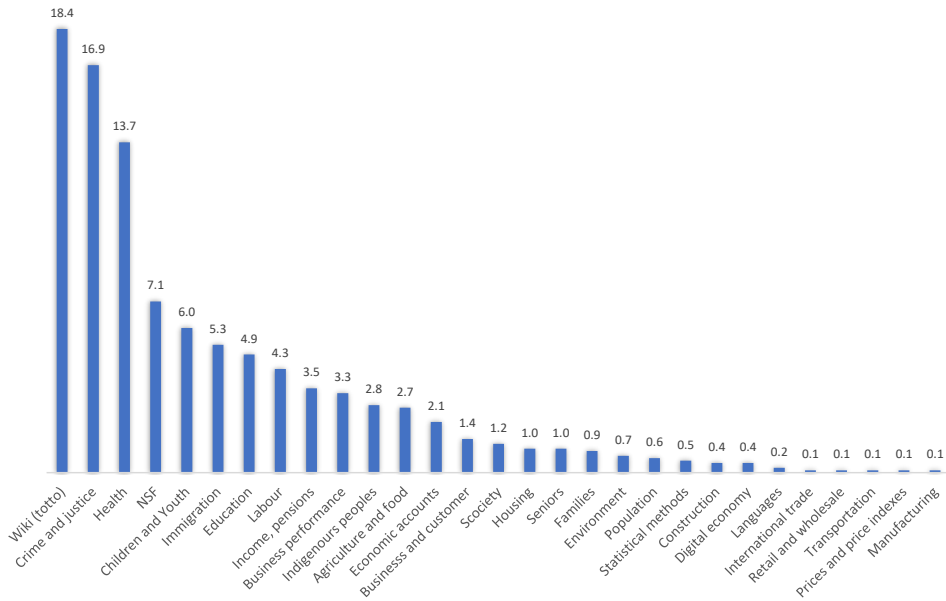


Figure 4: Proportion of samples in different 29 domains.

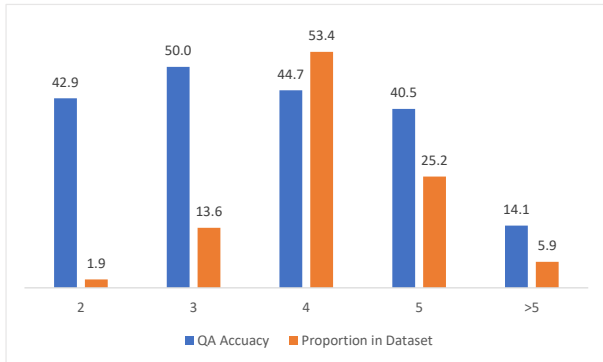


Figure 5: Level-wise QA accuracy and proportion of samples with MAPO and hierarchy-aware logical form.

Function	Trigger Words
argmax	JJR, JJS, RBR, RBS, top,
argmin	first, bottom, and last.
max	JJS, RBS
min	
average	average, mean
sum	all, combine, total, sum
count	how, many, total, number
difference	difference, more, than,
difference_rate	change, compare, JJR
difference_rate_rev	RBR.
proportion	times, percent,
proportion_rev	percentage, fraction

Table 8: Trigger Words for Functions

B.5 Pruning Rules in Searching

We use trigger words and POS tags for some functions in random exploration, which is inspired by (Zhang *et al.*, 2017; Liang *et al.*, 2018). Functions are allowed to be selected only when triggers appear in the question. Triggers are listed in Table 8.

C Hierarchical Table to Text

C.1 Illustration on controllable generation in hierarchical table to text.

Please find the illustration shown in Figure 6.

C.2 Baseline Implementation Details

We perform optimized tuning for baselines using the following settings.

Pointer Generator (See *et al.*, 2017) A LSTM-based seq2seq model with copy mechanism. The model uses two-layer bi-directional LSTMs for the encoder with 300-dim word embeddings and 300 hidden units. We perform fine-tuning using batch size 2, learning rate 0.05, and beam size 5.

BERT-to-BERT (Rothe *et al.*, 2020) A transformer encoder-decoder model (Vaswani *et al.*, 2017) where the encoder and decoder are both

	A	B	C	D	E	F	G
1	TABLE 3. Primary source and mechanism of support for full-time master's and doctoral students in science and engineering: 2017						
2		All full-time graduate students		Master's		Doctoral	
3	Source and mechanism	Total	Percent	All	Percent	All	Percent
4	All full-time	433,916	100.0	209,221	100.0	224,695	100.0
5	Self-support	161,641	37.3	139,373	66.6	22,268	9.9
6	All sources of support	272,275	62.7	69,848	33.4	202,427	90.1
7	Federal	65,999	15.2	10,736	5.1	55,263	24.6
8	Department of Agricu	2,361	0.5	938	0.4	1,423	0.6
9	Department of Defens	8,089	1.9	2,568	1.2	5,521	2.5
16	Other	9,098	2.1	3,462	1.7	5,636	2.5
17	Institutional	182,135	42.0	52,319	25.0	129,816	57.8
18	Other U.S. source	19,432	4.5	5,136	2.5	14,296	6.4
19	Foreign	4,709	1.1	1,657	0.8	3,052	1.4
20	All mechanisms of support	272,275	62.7	69,848	33.4	202,427	90.1
21	Fellowships	39,368	9.1	5,687	2.7	33,681	15.0
22	Traineeships	10,945	2.5	1,497	0.7	9,448	4.2
23	Research assistantships	103,586	23.9	19,702	9.4	83,884	37.3
24	Teaching assistantships	84,499	19.5	22,171	10.6	62,328	27.7
25	Other mechanisms	33,877	7.8	20,791	9.9	13,086	5.8

Target text:

For doctoral students, the proportion of support from research assistantships is 10 points higher than that from teaching assistantships.

Highlighted cells:

From entity alignment: Doctoral, percent, research assistantships, teaching assistantships. From quantity alignment: 37.3, 27.7

Operators:

DIFF

Input sequence after sub table selection and serialization:

[SEP] source and mechanism [SEP] doctoral [SEP] percent [SEP] all mechanisms of support [SEP] research assistantships [SEP] 37.3 [SEP] teaching assistantships [SEP] 27.7 [SEP] DIFF [SEP] 9.6

Figure 6: An illustration on controllable generation.

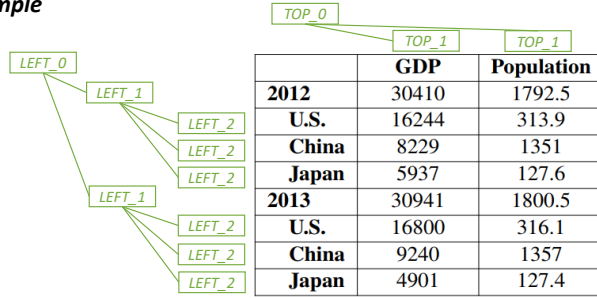
1017 initialized with BERT (Devlin *et al.*, 2018) by
 1018 loading the checkpoint named ‘bert-base-uncased’
 1019 provided by the huggingface/transformers repos-
 1020 itory. We perform fine-tuning using batch-size 2
 1021 and learning rate $3e^{-5}$.

1022 **BART** (Lewis *et al.*, 2019) BART is a pre-
 1023 trained denoising autoencoder for seq2seq lan-
 1024 guage modeling. It uses standard Transformer-
 1025 based architecture and shows effectiveness in NLG.
 1026 We align model configuration with the BASE ver-
 1027 sion of BART, and use the model ‘facebook/bart-
 1028 base’ in huggingface/transformers. During fine-
 1029 tuning, we use a batch size of 8 and a learning rate
 1030 of $2e^{-4}$.

1031 **T5** (Raffel *et al.*, 2019) T5 is also a transformer-
 1032 based pre-training LM. It trains extensively on text-
 1033 to-text tasks and scores high on generation tasks.
 1034 We use the pre-trained model ‘t5-base’ in hugging-
 1035 face/transformers. For fine-tuning, we set batch
 1036 size to 8 and learning rate to $2e^{-4}$.

1037 We use a beam size of 5 to search decoded out-
 1038 puts (sequence lengths range from 8 to 60 tokens)

Example



Q: What is the GDP of China in 2012?
A: 8229

Model

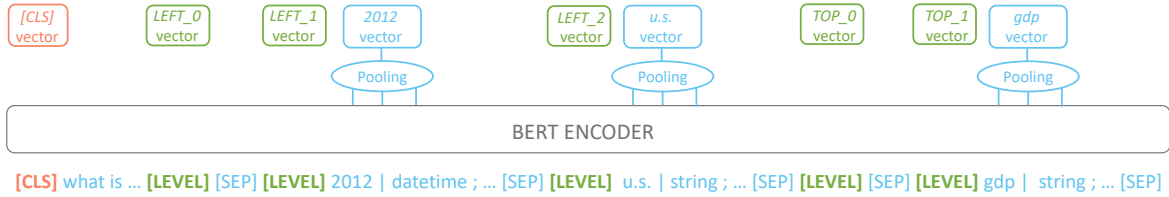


Figure 7: An QA example table with hierarchy and its linearized input to the encoder. Each level in the hierarchical header starts with a *LEVEL* token to learn a level representation. *LEFT_k* means the *k*th level in the left tree. Each header cell has a unique header cell representation.

Function	Arguments	Returns	Description
(filter_tree h)	h : a header	a region	Select a region indexed by sub-tree of the given header in the given region.
(filter_level l)	l : a level	a region	Select a region indexed by headers on the given level in the given region.
(argmax k) (argmin k)	k : a number	a list of headers	Find the header(s) with k-th largest/smallest value in the region. [Input region should have one row or one column of data]
(max l) (min l) (sum l) (average l)	l : a level	a region	Maximum/minimum/sum/average of the given region, grouping by headers of the given level, <i>i.e.</i> , data values aggregate according to their header strings on the given level.
(count l)	l : a level	a number	Count the number of headers on the given level of given region.
(difference) (proportion) (proportion_rev) (difference_rate) (difference_rate_rev)		a number	Absolute difference, proportion and difference rate of given two elements <i>a</i> and <i>b</i> in region. <i>rev</i> means changing order of operands. <i>e.g.</i> , <i>proportion</i> applies <i>b/a</i> and <i>proportion_rev</i> applies <i>a/b</i> . [Input region should have two data elements]
(greater_than n) (greater_eq_than n) (less_than n) (less_eq_than n) (eq n) (not_eq n)	n : a number	a list of headers	Find the header(s) with data value(s) that have certain order relation with the given number. [Input region should have one row or one column of data]
(opposite)		a number	Take opposite value of data in a given region. [Input region should have one data element]

Table 9: Function list of hierarchy-aware logical form

	A	B	C	D	E	F	G	H
1	Table 3: Sex and marital status by FOLS of workers in the agricultural sector aged 15 years and over, three agricultural regions of New Brunswick, 2011							
5	percent							
6	Sex							
7	Female	35.3	28	41.8	30.6	35.9	26.6	
8	Male	64.7	72	58.2	69.4	64.1	73.4	
9	Marital Status							
10	Single	18.7	24.8	26.9	26.1	25.6	32.8	
11	Married	51.3	53.9	47.8	56.7	54.4	57.8	
12	Common-Law	17.1	10.9	22.4	6.4	11.8	7.8	
13	Separated, divorced, or widowed	12.8	10.6	0	10.8	7.7	0	
14								
15								
16								
17	table descriptive sentence id:	178						
18	table descriptive sentence:	Regardless of the region or language, male workers outnumbered female workers in the agricultural sector in 2011.						
19								
20	sub-sentence (complete & fix gra	Regardless of the region or language, male workers outnumbered female workers in the agricultural sector in 2011.						
21	sub-sentence after deletion & decontextualization:	Regardless of the region or language, male workers outnumbered female workers in the agricultural sector in 2011.						
22	key part to be questioned:	male workers						
23	schema linking phrases:	region	language	female workers	agricultural sector	in 2011		
24	schema linking positions:	=B3	French-language worker	Female	Table 3: Sex and marital	Table 3: Sex and marital status by FOLS of workers in the a		
25	question rewrite:	Which group of people has more workers in the agricultural sector in 2011, regardless of the region or language? Male or female?						
26	answer (formula):	Male						
27	aggregation type:	pair-argmax						
28								
29	table descriptive sentence id:	179						
30	table descriptive sentence:	Compared with English-language workers, there were fewer men among French-language agricultural workers.						
31								
32	sub-sentence (complete & fix gra	Compared with English-language workers, there were fewer men among French-language agricultural workers.						
33	sub-sentence after deletion & decontextualization:	Compared with English-language workers, there were fewer men among French-language agricultural workers.						
34	key part to be questioned:	French-language agricultural workers						
35	schema linking phrases:	English-language worke men						
36	schema linking positions:	English-language worke Male						
37	question rewrite:	Which sector has fewer male agricultural workers? English-language workers or French-language workers?						
38	answer (formula):	French-language workers						
39	aggregation type:	pair-argmax						
40								
41	table descriptive sentence id:	180						
42	table descriptive sentence:	In 2011, the majority of New Brunswick's agricultural workers, both English-language and French-language workers, were married.						
43								
44	sub-sentence (complete & fix gra	In 2011, the majority of New Brunswick's agricultural workers, both English-language and French-language workers, were married.						
45	sub-sentence after deletion & decontextualization:	In 2011, the majority of New Brunswick's agricultural workers, both English-language and French-language workers, were married.						
46	key part to be questioned:	married						
47	schema linking phrases:	in 2011	New Brunswick	agricultural workers	English-language	French-language workers		
48	schema linking positions:	Table 3: Sex and marita	Table 3: Sex and marital	Table 3: Sex and marita	English-language worke	French-language workers		
49	question rewrite:	What is the marital status for the majority of New Brunswick's agricultural workers, both English-language and French-language workers in 2011?						
50	answer (formula):	Married						
51	aggregation type:	argmax						
52								
53	table descriptive sentence id:	181						
54	table descriptive sentence:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
55								
56	sub-sentence (complete & fix grammar):	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
57	sub-sentence after deletion & decontextualization:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
58	key part to be questioned:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
59	schema linking phrases:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
60	schema linking positions:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
61	question rewrite:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
62	answer (formula):	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
63	aggregation type:	As a general rule, French-language workers were less likely to be married or single than their English-language colleagues.						
64								
65	table descriptive sentence id:	182						
66	table descriptive sentence:	In 2011, French-language workers were more likely to be in a common-law relationship than their English-language colleagues.						
67								
68	sub-sentence (complete & fix gra	In 2011, French-language workers were more likely to be in a common-law relationship than their English-language colleagues.						
69	sub-sentence after deletion & decontextualization:	In 2011, French-language workers were more likely to be in a common-law relationship than their English-language colleagues.						
70	key part to be questioned:	French-language workers						
71	schema linking phrases:	in 2011	English-language collea	common-law relationship				
72	schema linking positions:	Table 3: Sex and marita	English-language worke	Common-Law				
73	question rewrite:	In 2011, which sector of workers were more likely to be in a common-law relationship? French-language workers or English-language workers?						
74	answer (formula):	French-language workers						
75	aggregation type:	pair-argmax						

Figure 8: Annotation interface in Excel.