

## CHAPTER 6

# Data to Text

This Chapter covers data-to-text generation, which is conceptually the converse of the text-to-data material covered in Chapter 4.

## 6.1 INTRODUCTION

The term **Natural Language Generation** (NLG) is heavily used in the NLP literature but it means different things in different contexts. It can mean predicting (generating) the next word or sentence in a sequence using a Hidden Markov Model (HMM), a Recurrent Neural Network (RNN), or a pretrained language model such as CTRL [Keskar et al., 2019], GPT [Radford et al., 2018], and GPT-3 [Brown et al., 2020b].

Natural language generation can be done with or without a given context. When the output depends on the context, sometimes the term “controllable generation” is used. The context can be in the form of a user query, a prompt, a user profile, or a document, among other things.

However, in a more traditional sense, the term NLG describes the process of converting a *semantic representation* (which could be a tree or a database record, among other things) to natural, running text. In this chapter we will talk about both of these types of generations, referred to as unconditional and conditional, with a focus on conditional data-to-text generation, where the input is a structured representation (e.g., a table) and the output is plain text.

### 6.1.1 TRADITIONAL GENERATION

$$\mathcal{F}_4 \rightarrow \begin{bmatrix} cat: & s \\ prot: & \boxed{1} \quad \left[ n: "john" \right] \\ verb: & \left[ v: "like" \right] \\ goal: & \boxed{1} \end{bmatrix}$$

Figure 6.1: Sample FUF input for the sentence “John likes himself.” (Image from ?).

### 6.1.2 DATA-TO-TEXT GENERATION

FUF [Elhadad, 1988] is one of the early table-to-text systems, developed in late 1980s, that used functional unification grammars to describe a set of key-value features. The features could be nested (e.g., the value of a key can be a feature structure). Also two keys (e.g., *prot* and *goal* in Figure 6.1) can have the same values. Another early example of data-to-text generation appears in Konstas and Lapata [2012]. The input comes from the ATIS airline reservation dataset [Dahl et al., 1994]. It is represented in a structured data form, with nested attributes, e.g., the day of the departure is the 9th of August (Figure 6.2). One more example from the same time period is Konstas and Lapata [2013]. They use WeatherGov and WinHelp (Figure 6.3), two standard datasets, and generate multisentential text that describes the database records.

	Flight	Day Number	Month	Condition	Search
Database:	<b>from</b> <b>to</b> denver boston	<b>number</b> <b>dep/ar</b> 9 departure	<b>month</b> <b>dep/ar</b> august departure	<b>arg1</b> <b>arg2</b> <b>type</b> arrival_time 1600 <	<b>type</b> <b>what</b> query flight
$\lambda$ -expression:	$\lambda x. \text{flight}(x) \wedge \text{from}(x, \text{denver}) \wedge \text{to}(x, \text{boston}) \wedge \text{day\_number}(x, 9) \wedge \text{month}(x, \text{august}) \wedge \text{less\_than}(\text{arrival\_time}(x), 1600)$				
Text:	Give me the flights leaving Denver August ninth coming back to Boston before 4pm.				

Figure 6.2: Structured data as input, along with a lambda expression representing its meaning, and the matching output text. (Image from Konstas and Lapata [2012]).

Database Records	Database Records
temp(time:6-21, min:9, mean:15, max:21) wind-spd(time:6-21, min:15, mean:20, max:30) sky-cover(time:6-9, percent:25-50) sky-cover(time:9-12, percent:50-75) wind-dir(time:6-21, mode:SSE) gust(time:6-21, min:20, mean:30, max:40)	desktop(cmd:lclick, name:start, type:button) start(cmd:lclick, name:settings, type:button) start-target(cmd:lclick, name:control panel, type:button) win-target(cmd:dblclick, name:users and passwords, type:item) contMenu(cmd:lclick, name:advanced, type:tab) action-contMenu(cmd:lclick, name:advanced, type:button)
Output Text	Output Text
Cloudy, with a high around 20. South southeast wind between 15 and 30 mph. Gusts as high as 40 mph.	Click start, point to settings, and then click control panel. Double-click users and passwords. On the advanced tab, click advanced.

Figure 6.3: WeatherGov and WinHelp examples. (Image from Konstas and Lapata [2013]).

Other papers use as input attribute-value pairs extracted from Wikipedia biographical infoboxes, e.g., Figure 6.4. In this scenario, the generated output has to describe the entity in question using grammatical sentences, covering the most important facts in the input.

The Neural Wikipedian paper [Vougiouklis et al., 2017] and its follow up paper [Vougiouklis et al., 2020] introduce another interesting dataset that aligns sets of RDF triples to verbal summaries. An example data point from this dataset is shown in Figure 6.5. The text summary is a sentence that succinctly covers multiple facts presented in the input triples.

(a)	<table border="1"> <tr> <td><b>Born</b></td><td>Robert Joseph Flaherty February 16, 1884 Iron Mountain, Michigan, U.S.</td></tr> <tr> <td><b>Died</b></td><td>July 23, 1951 (aged 67) Dummerston, Vermont, U.S.</td></tr> <tr> <td><b>Cause of death</b></td><td>Cerebral thrombosis</td></tr> <tr> <td><b>Occupation</b></td><td>Filmmaker</td></tr> <tr> <td><b>Spouse(s)</b></td><td>Frances Johnson Hubbard</td></tr> </table>	<b>Born</b>	Robert Joseph Flaherty February 16, 1884 Iron Mountain, Michigan, U.S.	<b>Died</b>	July 23, 1951 (aged 67) Dummerston, Vermont, U.S.	<b>Cause of death</b>	Cerebral thrombosis	<b>Occupation</b>	Filmmaker	<b>Spouse(s)</b>	Frances Johnson Hubbard	(b) Robert Joseph Flaherty, (February 16, 1884 – July 23, 1951) was an American film-maker who directed and produced the first commercially successful feature-length documentary film, <i>Nanook of the North</i> (1922). The film made his reputation and nothing in his later life fully equalled its success, although he continued the development of this new genre of narrative documentary, e.g., with <i>Moana</i> (1926), set in the South Seas, and <i>Man of Aran</i> (1934), filmed in Ireland's Aran Islands. He is considered the “father” of both the documentary and the ethnographic film. Flaherty was married to writer Frances H. Flaherty from 1914 until his death in 1951. Frances worked on several of her husband's films, and received an Academy Award nomination for Best Original Story for <i>Louisiana Story</i> (1948).
<b>Born</b>	Robert Joseph Flaherty February 16, 1884 Iron Mountain, Michigan, U.S.											
<b>Died</b>	July 23, 1951 (aged 67) Dummerston, Vermont, U.S.											
<b>Cause of death</b>	Cerebral thrombosis											
<b>Occupation</b>	Filmmaker											
<b>Spouse(s)</b>	Frances Johnson Hubbard											
(c)	(c) Robert Joseph Flaherty, (February 16, 1884 – July 23, 1951) was an American film-maker. Flaherty was married to Frances H. Flaherty until his death in 1951.											

Figure 6.4: Sample Wikipedia biographical box, along with the matching text and verbalization (Image from Perez-Beltrachini and Lapata [2018]).

Triples	Atlas_Shrugged <b>literaryGenre</b> Science_fiction Atlas_Shrugged <b>country</b> United_States John_Galt <b>series</b> Atlas_Shrugged Atlas_Shrugged <b>publicationYear</b> "1957" Atlas_Shrugged <b>author</b> Ayn_Rand
Text Summary	Atlas Shrugged is a science fiction novel by Ayn Rand.

Figure 6.5: Text summary aligned to a set of five RDF triples. Each of the triples includes a predicate (verb) in bold face and two arguments connected by that predicate. (Example from Vougiouklis et al. [2020]).

### 6.1.3 ABSTRACT MEANING REPRESENTATION (AMR) FOR TEXT GENERATION

One representation suitable for meaning to text generation is AMR (Abstract Meaning Representation), previously discussed in the text-to-data chapter.

In natural language understanding, there is not much leeway in representing the content of the input sentence using a meaning representation. In contrast, in natural language generation, the expected output sentence for a given input is underspecified. For example, the AMR representation in Figure 6.6 indicates that a “break” event has taken or is taking place, and that its main arguments are “dog” and “window”. This input can be translated to natural language sentences in multiple ways, e.g., “the dog broke the window”, or “the dog will break the window”, or “the window was broken by the dog”, or even “the dog broke it.”

```
(b / break-01
:ARG0 (d / dog)
:ARG1 (w / window))
```

**Figure 6.6:** An example of meaning input underspecification. The input covers only some of the information needed to generate a full sentence. The rest of the features required to generate such a sentence have to come from other sources, e.g., discourse constraints, user preferences, etc.

In other words, natural language generation depends significantly on a series of choices made during the different stages mentioned earlier, e.g., content selection, planning, lexicalization, and surface realization. Such choices may be made in an arbitrary way, or they may depend on additional constraints. Such constraints can be based on knowledge about the time of the events (to determine the grammatical tense of the sentences) or the discourse structure (e.g., whether the dog or the window is the theme of the passage, or whether to replace one or more of the entities that have already been mentioned with a pronoun). A good survey that covers such topics is [Gatt and Krahmer \[2017\]](#).

Using AMR as a representation for NLG has been very popular over the recent years. The SemEval 2017 Task 9<sup>1</sup> was on generating text from Abstract Meaning Representations. It included two tasks, one of which is on AMR to English generation using (domain-independent) news and discussion text.

### 6.1.4 NEURAL GENERATION

Many neural methods for natural language generation are based on an encoder-decoder model. Encoding was traditionally done using a recurrent neural network (RNN) that converts the input text into a formal meaning representation, which is then decoded.

<sup>1</sup><https://alt.qcri.org/semeval2017/task9/>

```
s / say-01
    :ARG0 (s2 / service
        :mod (e / emergency)
        :location (c / city :wiki "London"
            :name (n / name :op1 "London")))
    :ARG1 (s3 / send-01
        :ARG1 (p / person :quant 11)
        :ARG2 (h / hospital)
        :mod (a / altogether)
        :purpose (t / treat-03
            :ARG1 p
            :ARG2 (w / wound-01
                :ARG1 p
                :mod (m / minor)))))
```

Figure 6.7: AMR input representation. This structure can be generated as “The London emergency services said that altogether 11 people had been sent to hospital for treatment due to minor wounds.”

This is essentially the same machinery that we described in Chapter 4, with the output being plain text this time.

At each time step  $t$ , the model computes a vector of scores for each item in the output vocabulary and converts them to probabilities using the softmax function. During decoding, the algorithm decides what token to choose as the next item in the output. A simple way to perform decoding is to use *teacher forcing* (or maximum likelihood training) to produce, at each step, the most likely token, given the tokens generated so far. Beam decoding is a greedy algorithm that considers at each time step the  $k$  most likely candidates. Argmax decoding is an extreme case of a greedy algorithm with a beam size of  $k = 1$ .

Greedy methods tend to become repetitive and are usually avoided in favor of sampling methods such as top- $k$  sampling or nucleus sampling. The amount of randomness can be regulated by adjusting the temperature parameter of the softmax. A survey of decoding methods in NLG is [Zarrieß et al. \[2021\]](#).

## 6.2 DOMAIN SPECIFIC TABLE-TO-TEXT

Research contributions in table-to-text generation have been split into two thrusts - generic domain and domain-specific. A number of tasks have been created over the years in domain-specific table-to-text generation. We will now describe some of them.

### 6.2.1 SRST

The First Multilingual Surface Realisation Shared Task (SR'18) [Mille et al. \[2018a,b\]](#) was presented at an ACL workshop, as a follow up to a Surface Realisation Shared Task in 2011 [Belz et al. \[2011\]](#). It includes a shallow track (in ten languages) with universal dependency structures from CONLL where the lemmatized words appear in a shuffled order, as well as a deep track (in three languages) where all function words, as well as most of the morphological information, are additionally removed. Figure 6.8 shows a sample input representation. A shallow input representation (from Track 1 of SR'18) is shown in Figure 6.9 and the corresponding deep input representation (from Track 2 of SR'18) is shown in Figure 6.10. The Third Multilingual Surface Realisation Shared Task (SR'20) [\[Mille et al., 2020\]](#) was introduced in 2020. In this version of the task, each of the two tracks (shallow and deep) had two variants, one in which only the data given as input for the track could be used for training and a second one that allowed the use of any source of data for training.

### 6.2.2 E2E

The **e2e** dataset<sup>2</sup>, described in [Dušek et al. \[2018, 2020\]](#), [Novikova et al. \[2017\]](#), is limited to the restaurant domain and is intended to exhibit high syntactic complexity and different discourse phenomena.

<sup>2</sup><http://www.macs.hw.ac.uk/InteractionLab/E2E/>

**134 6. DATA TO TEXT**

1	The	the	DET	DT	Definite=Def PronType=Art	2	det
2	third	third	ADJ	JJ	Degree=Pos NumType=Ord	5	nsubj-pass
3	was	be	AUX	VBD	Mood=Ind Number=Sing Person=3 Tense=Past VerbForm=Fin	5	aux
4	being	be	AUX	VBG	VerbForm=Ger	5	aux-pass
5	run	run	VERB	VBN	Tense=Past VerbForm=Part Voice=Pass	0	root
6	by	by	ADP	IN	-	8	case
7	the	the	DET	DT	Definite=Def PronType=Art	8	det
8	head	head	NOUN	NN	Number=Sing	5	obl
9	of	of	ADP	IN	-	12	case
10	an	a	DET	DT	Definite=Ind PronType=Art	12	det
11	investment	investment	NOUN	NN	Number=Sing	12	compound
12	firm	firm	NOUN	NN	Number=Sing	8	nmod
13	.	.	PUNCT	.	-	5	punct

**Figure 6.8:** A universal dependency structure for English from SRST'18. (Image from [Mille et al. \[2018b\]](#)).

1	the	-	DET	DT	Definite=Def PronType=Art	2	det
2	third	-	ADJ	JJ	Degree=Pos NumType=Ord	3	nsubj-pass
3	run	-	VERB	VBN	Tense=Past VerbForm=Part Voice=Pass	0	root
4	be	-	AUX	VBD	Mood=Ind Number=Sing Person=3 Tense=Past VerbForm=Fin	3	aux
5	be	-	AUX	VBG	VerbForm=Ger	3	aux-pass
6	head	-	NOUN	NN	Number=Sing	3	obl
7	.	-	PUNCT	.	-	3	punct
8	by	-	ADP	IN	-	6	case
9	the	-	DET	DT	Definite=Def PronType=Art	6	det
10	firm	-	NOUN	NN	Number=Sing	6	nmod
11	an	-	DET	DT	Definite=Ind PronType=Art	10	det
12	investment	-	NOUN	NN	Number=Sing	10	compound
13	of	-	ADP	IN	-	10	case

**Figure 6.9:** Shallow input version of the previous example (Track 1). (Image from [Mille et al. \[2018b\]](#)).

1	third	-	ADJ	-	Degree=Pos	2	A2
2	run	-	VERB	-	Tense=Past Aspect=Progr	0	ROOT
3	head	-	NOUN	-	Number=Sing Definiteness=Def	2	A1
4	firm	-	NOUN	-	Number=Sing Definiteness=Indef	3	A2
5	investment	-	NOUN	-	Number=Sing	4	AM

**Figure 6.10:** Deep input derived from the previous example (Track 2). (Image from [Mille et al. \[2018b\]](#)).

Figure 6.11 shows a sample data point, including a meaning representation for the restaurant domain and a single corresponding reference sentence. The domain ontology is shown in Figure 6.12. An example with multiple reference sentences is included in 6.13. The e2e dataset includes 6,039 meaning representations and 51,426 reference sentences.

<b>MR</b>	name[The Wrestlers], priceRange[cheap], customerRating[low]
<b>Reference</b>	The wrestlers offers competitive prices, but isn't highly rated by customers.

Figure 6.11: A meaning representation and the corresponding reference sentence.  
(Image from Dušek et al. [2018]).

Attribute	Data Type	Example value
name	verbatim string	<i>The Eagle, ...</i>
eatType	dictionary	<i>restaurant, pub, ...</i>
familyFriendly	boolean	<i>Yes / No</i>
priceRange	dictionary	<i>cheap, expensive, ...</i>
food	dictionary	<i>French, Italian, ...</i>
near	verbatim string	<i>Zizzi, Cafe Adriatic, ...</i>
area	dictionary	<i>riverside, city center, ...</i>
customerRating	dictionary	<i>1 of 5 (low), 4 of 5 (high), ...</i>

Figure 6.12: Domain ontology for the e2e task. (Image from Dušek et al. [2018]).

```
name[The Eagle],  
eatType[coffee shop],  
food[French],  
priceRange[moderate],  
customerRating[3/5],  
area[riverside],  
kidsFriendly[yes],  
near[Burger King]
```

Figure 6.13: Example from the e2e dataset: “The three star coffee shop, The Eagle, gives families a mid-priced dining experience featuring a variety of wines and cheeses. Find The Eagle near Burger King.”

One of the early approaches to the e2e task is the neural template induction approach of Wiseman et al. [2018]. An example is shown in Figure 6.14. The templates

136 6. DATA TO TEXT

are learned automatically and include alternative phrasings, e.g., “providing/serving/of-  
ering”. Once generated, these templates are then filled in with the missing words for  
the specific example, e.g., as in Figure 6.15.

## **Source Entity:** Cotto

type[coffee shop], rating[3 out of 5],  
food[English], area[city centre],  
price[moderate], near[The Portland Arms]

## System Generation:

Cotto is a coffee shop serving English food in the moderate price range. It is located near The Portland Arms. Its customer rating is 3 out of 5.

(a) An example from the e2e dataset.

## Neural Template:

The \_\_\_\_\_ is a \_\_\_\_\_ providing \_\_\_\_\_  
... is an expensive \_\_\_\_\_ serving \_\_\_\_\_  
...  
food in the price range ...  
cuisine with a price bracket ...  
foods and has a ... pricing It's  
... The place is  
...  
located in the Its customer rating is  
located near ... Their customer rating is  
near ... Customers have rated it  
...  
...

(b) An induced neural template.

Figure 6.14: E2e examples. (Images from Wiseman et al. [2018]).

- |    |                                 |   |                                    |  |   |                                    |                                    |                           |  |
|----|---------------------------------|---|------------------------------------|--|---|------------------------------------|------------------------------------|---------------------------|--|
| 1. | The Eagle<br>Zizzi              | provides<br>serves                            | Indian<br>Chinese<br>English       | food<br>cuisine<br>Food                | in the<br>with a<br>and has a                         | high<br>moderate<br>average        | price range<br>customer rating .   | It is<br>They are<br>It's | near<br>located in the<br>located near |
|    | ...<br>riverside<br>city centre | ...<br>Its customer rating is                 | ...<br>1 out of 5                  |  |   | ...                                | ...                                | ...                       | ...                                    |
|    | Cafe Sicilia                    | ...<br>The price range is                     | ...<br>average<br>high             |  |   |                                    |                                    |                           |  |
| 2. | Located near<br>Near            | The Portland Arms<br>riverside<br>city centre | is a family friendly               | Italian<br>fast food<br>French         | restaurant called<br>place called<br>restaurant named | The Waterman<br>Cocum<br>Loch Fyne |                                    |                           |  |
|    | ...<br>...                      | ...<br>there is a                             | ...<br>...                         | ...<br>...                             | ...<br>...  | ...                                |                                    |                           |  |
| 3. | A<br>An<br>A family friendly    | Italian<br>fast food<br>French                | restaurant<br>pub<br>coffee shop   | is called<br>named                     | The Waterman<br>Cocum<br>Loch Fyne                    |                                    |                                    |                           |  |
|    | ...<br>...                      | ...<br>...                                    | ...<br>...                         | ...<br>...                             | ...   |                                    |                                    |                           |  |
| 4. | Located near<br>Near            | The Portland Arms<br>riverside<br>city centre | The Eagle<br>Golden Curry<br>Zizzi | is a family friendly<br>is an          | cheap<br>family-friendly<br>family friendly           | Italian<br>fast food<br>French     | restaurant<br>pub<br>coffee shop . |                           |  |
|    | ...<br>...                      | ...<br>...                                    | ...<br>...                         | ...<br>...                             | ...<br>...  |                                    |                                    |                           |  |
| 5. | A<br>An<br>A family friendly    | Italian<br>fast food<br>French                | restaurant<br>pub<br>coffee shop   | near<br>located in the<br>located near | riverside<br>city centre<br>Cafe Sicilia              | is called<br>named                 | The Waterman<br>Cocum<br>Loch Fyne |                           |  |
|    | ...<br>...                      | ...<br>...                                    | ...<br>...                         | ...<br>...                             | ...   |                                    |                                    |                           |  |

Figure 6.15: Filled out templates. (Image from Wiseman et al. [2018]).

### 6.2.3 WEBNLG

The **WebNLG** challenge [Gardent et al. \[2017\]](#), introduced in 2017, is based on mapping RDF triples to plain text. The WebNLG dataset includes 25,298 (data, text) pairs and 9,674 sets of up to seven RDF triples extracted from DBpedia. The texts verbalize these sets of RDF triples. The examples are all domain-specific in 15 domains such as University, Building, Airport, City, Athlete, Politician, Astronaut, Artist, etc. In the

example shown in Figure 6.16, three RDF triples about the same person are combined into a single sentence.

```
(JOHN E BLAHA BIRTHDATE 1942 08 26)
(JOHN E BLAHA BIRTHPLACE SAN ANTONIO)
(JOHN E BLAHA OCCUPATION FIGHTER PILOT)
```

**Figure 6.16:** WebNLG example from Gardent et al. [2017] corresponding to the sentence “John E Blaha, born in San Antonio on 1942-08-26, worked as a fighter pilot.”.

A follow up paper [Castro Ferreira et al., 2018] expands<sup>3</sup> the corpus to German and also introduces some new tasks such as Discourse Ordering, Lexicalization (word choice), and Referring Expression (e.g., pronoun) Generation. An example of an input set with five RDF triples is shown in Figure 6.17. Figure 6.18 shows the mapping between the tags and entities for the template “AGENT-1 is located in BRIDGE-1 , PATIENT-1 and serves the city of PATIENT-2 . BRIDGE-1 is part of PATIENT-3 and PATIENT-4.”[Castro Ferreira et al., 2018].

Subject	Predicate	Object
Appleton International Airport	location	Greenville, Wisconsin
Greenville, Wisconsin	isPartOf	Ellington, Wisconsin
Greenville, Wisconsin	isPartOf	Menasha (town), Wisconsin
Greenville, Wisconsin	country	United States
Appleton International Airport	cityServed	Appleton, Wisconsin

**Figure 6.17:** Input set of five RDF triples that can be verbalized as “The Appleton International Airport is located in Greenville, Wisconsin, United States and serves the city of Appleton, Wisconsin. Greenville is part of the town of Menasha and Ellington, Wisconsin.”. (Example from Castro Ferreira et al. [2018]).

Tag	Entity
AGENT-1	Appleton International Airport
BRIDGE-1	Greenville, Wisconsin
PATIENT-1	United States
PATIENT-2	Appleton, Wisconsin
PATIENT-3	Menasha (town), Wisconsin
PATIENT-4	Ellington, Wisconsin

**Figure 6.18:** Tags and entities for the delexicalized/wikified templates. (Example from Castro Ferreira et al. [2018]).

A more recent<sup>4</sup> version of WebNLG expands the generation task to other languages such as Russian. An example with six triples is shown in Figure 6.19. This

<sup>3</sup><https://github.com/ThiagoCF05/webnlg>

<sup>4</sup>[https://webnlg-challenge.loria.fr/challenge\\_2020](https://webnlg-challenge.loria.fr/challenge_2020)

## 138 6. DATA TO TEXT

challenge includes three tasks: RDF-to-Text (Generation), Text-to-RDF, and Automatic Evaluation for WebNLG.

```
<entry category="Company" eid="Id3" shape="(X (X) (X) (X) (X) (X))" shape_type="sibling" size="6">
    <modifiedtripleset>
        <mtriple>Chinabank | foundingDate | 1920-08-16</mtriple>
        <mtriple>Chinabank | numberOfLocations | 295</mtriple>
        <mtriple>Chinabank | foundationPlace | Manila</mtriple>
        <mtriple>Chinabank | type | Public_company</mtriple>
        <mtriple>Chinabank | foundationPlace | Insular_Government_of_the_Philippine_Islands</mtriple>
        <mtriple>Chinabank | location | Philippines</mtriple>
    </modifiedtripleset>
</entry>
```

**Figure 6.19:** WebNLG example with six input triples. The output sentences are (1) “Chinabank was founded on August 16, 1920 in the Insular Government of the Philippine Islands In Manila. Located in the Philippines, it is a publicly traded company with 295 branches.”, (2) “Chinabank was founded on August 16, 1920 in Manila, in the Insular Government of the Philippine Islands. It is a publicly traded company, located in the Philippines with 295 branches.”, and (3) “Publicly traded Chinabank, founded in Manila, Philippines, at the time of the Insular Government of the Philippine Islands, on August 16, 1920, operates 295 banking centers.”.

### 6.2.4 WIKIBIO

**WikiBio** was introduced by [Lebret et al. \[2016\]](#). It includes more than 700K sentences from Wikipedia biographies. Figure 6.20a shows an example from the WikiBio dataset. Figure 6.20b shows the attention scores from slot to word. Finally, Figure 6.21 displays three examples of the outputs of their table-conditioned neural language model (TableNLM).

Some of the early work on WikiBio<sup>5</sup> is [Liu et al. \[2017\]](#). In this paper, the authors propose a structure-aware seq2seq architecture that uses a field-gating encoder which associates an LSTM unit with the corresponding field value. During decoding, dual attention to both the words and fields connects the generated description and the table. Figure 6.22 shows a sample Wikipedia infobox and the corresponding field representation. Figure 6.23 shows the structure-aware seq2seq architecture. Another example from [Liu et al. \[2017\]](#) is shown in Figure 6.24.

The follow up paper by [Sha et al. \[2018\]](#) introduces an order-planning generation model for Wikipedia data (Figure 6.25). They use a self-adaptive gate Figure 6.26 that takes into account both content-based and link-based attention (Figure 6.28). A sample output is shown in Figure 6.27.

<sup>5</sup><https://github.com/tyliupku/wiki2bio>

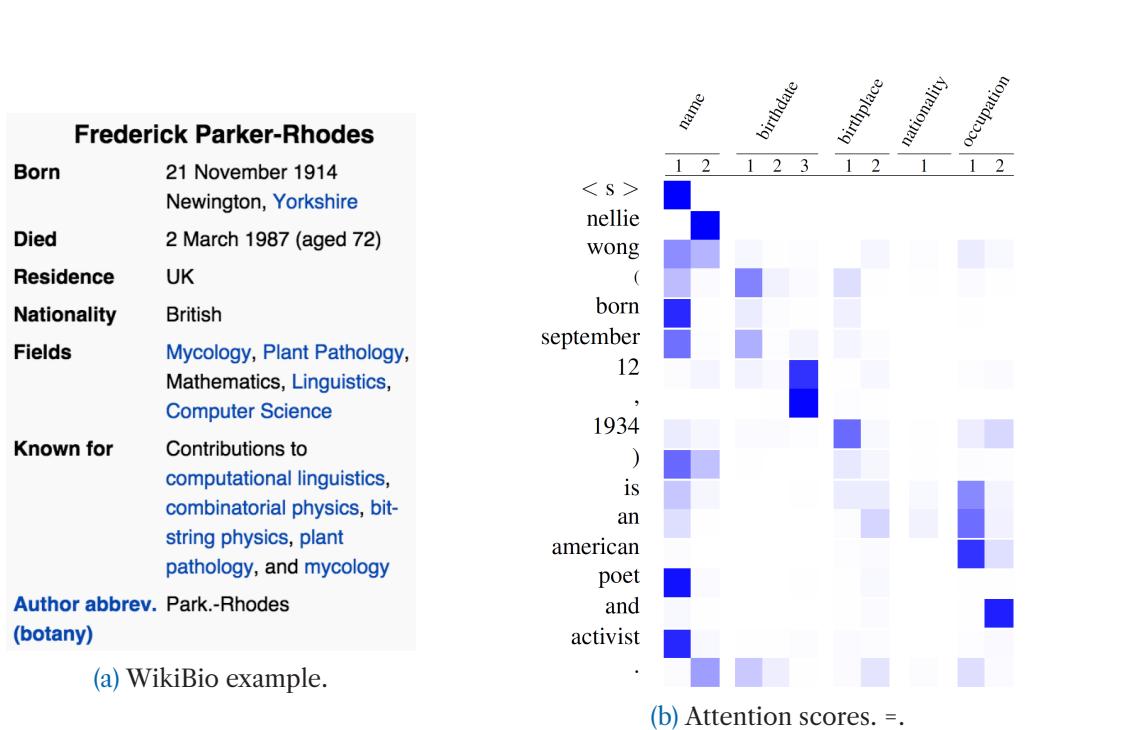


Figure 6.20: Wikibio examples. (Images from Lebret et al. [2016]).

Model	Generated Sentence
Reference	frederick parker-rhodes (21 march 1914 – 21 november 1987) was an english linguist, plant pathologist, computer scientist, mathematician, mystic, and mycologist.
Baseline (Template KN)	frederick parker-rhodes ( born november 21 , 1914 – march 2 , 1987 ) was an english cricketer .
Table NLM +Local (field, start)	frederick parker-rhodes ( 21 november 1914 – 2 march 1987 ) was an australian rules footballer who played with carlton in the victorian football league ( vfl ) during the XXXXs and XXXXs .
+ Global (field)	frederick parker-rhodes ( 21 november 1914 – 2 march 1987 ) was an english mycology and plant pathology , mathematics at the university of uk .
+ Global (field, word)	frederick parker-rhodes ( 21 november 1914 – 2 march 1987 ) was a british computer scientist , best known for his contributions to computational linguistics .

Figure 6.21: Reference and sample system outputs for the WikiBio input shown in Figure 6.20a. (Image from Lebret et al. [2016]).

word	Field embedding
name	George Mikell
birthdate	Jurgis Mikelaitis
birthplace	4 April 1929 (age 88)
nationality	Bildeniai, Lithuania
occupation	Lithuanian, Australian
years active	Actor, writer
known for	1957–present
The	The Guns of Navarone
Great	The Great Escape
Escape	

Figure 6.22: Wikipedia infobox and its field representation for George Mikell. (Image from Liu et al. [2017]).

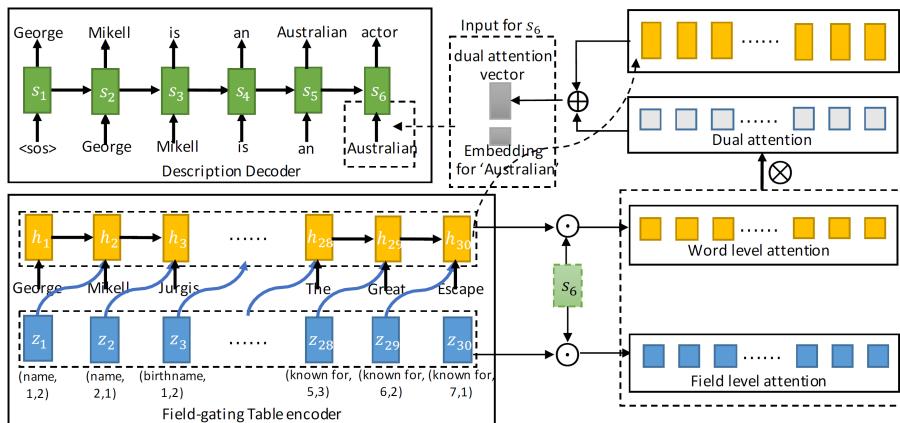


Figure 6.23: Structure-aware seq2seq architecture for generating the description for George Mikell in Figure 6.22 (Image from Liu et al. [2017]).

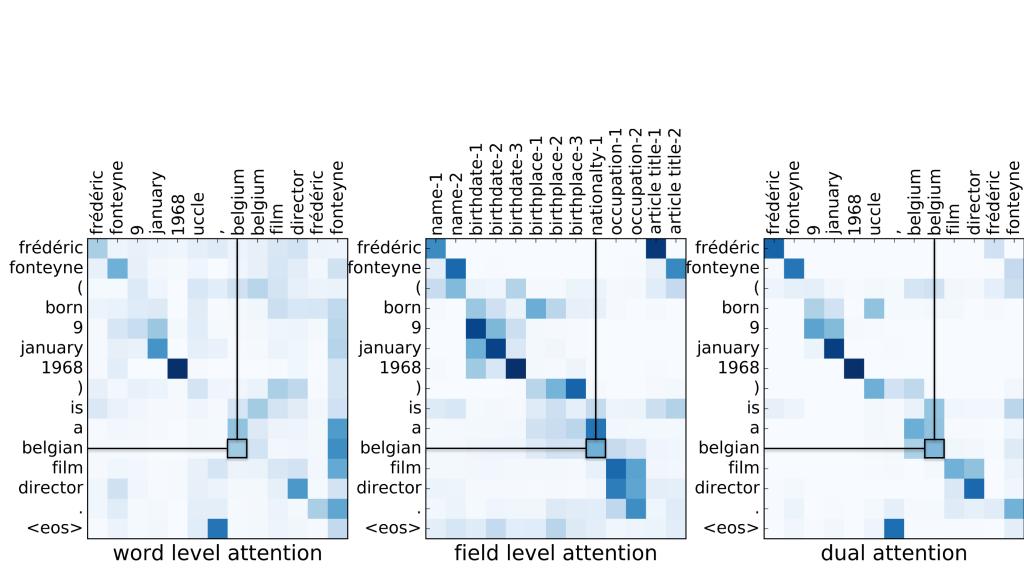


Figure 6.24: Dual attention. (Image from Liu et al. [2017]).

**Table:**

ID	Field	Content
1	Name	<i>Arthur Ignatius Conan Doyle</i>
2	Born	<i>22 May 1859 Edinburgh, Scotland</i>
3	Died	<i>7 July 1930 (aged 71) Crowborough, England</i>
4	Occupation	<i>Author, writer, physician</i>
5	Nationality	<i>British</i>
6	Alma mater	<i>University of Edinburgh Medical School</i>
7	Genre	<i>Detective fiction fantasy</i>
8	Notable work	<i>Stories of Sherlock Holmes</i>

**Text:** Sir Arthur Ignatius Conan Doyle (22 May 1859 – 7 July 1930) was a British writer best known for his detective fiction featuring the character Sherlock Holmes.

Figure 6.25: A matching pair of a Wikipedia infobox and its description. (Image from Sha et al. [2018]).

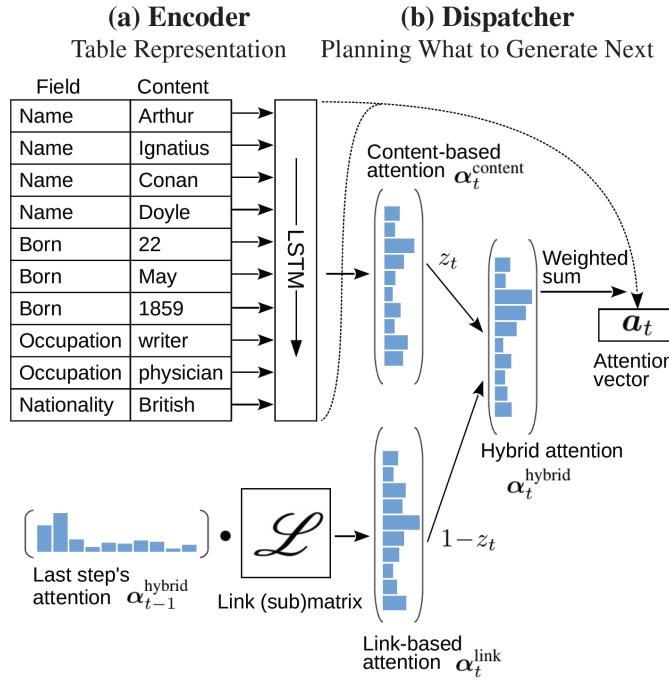
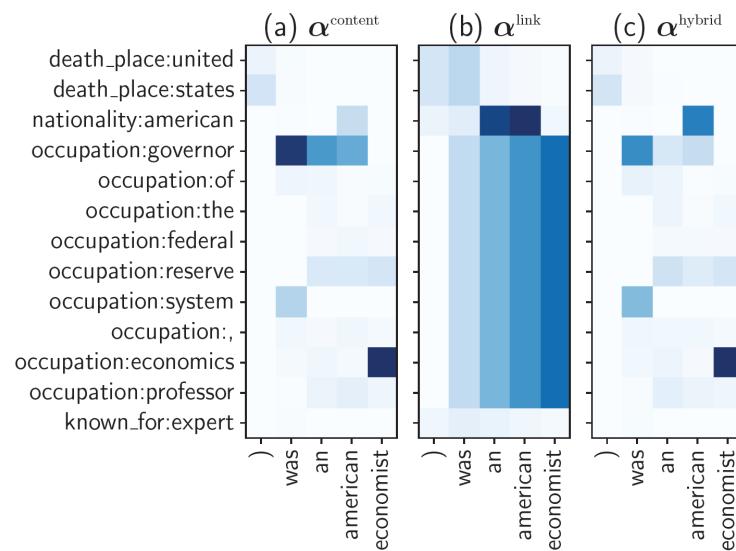


Figure 6.26: The encoder and the dispatcher of the system. (Image from Sha et al. [2018]).

<table border="1"> <tr><td><b>Name</b></td><td>Emmett John Rice</td></tr> <tr><td><b>Birth date</b></td><td>December 21, 1919</td></tr> <tr><td><b>Birth place</b></td><td>Florence, South Carolina, United States</td></tr> <tr><td><b>Death date</b></td><td>March 10, 2011 (aged 91)</td></tr> <tr><td><b>Death place</b></td><td>Camas, Washington, United States</td></tr> <tr><td><b>Nationality</b></td><td>American</td></tr> <tr><td><b>Occupation</b></td><td>Governor of the Federal Reserve System, Economics Professor</td></tr> <tr><td><b>Known for</b></td><td>Expert in the Monetary System of Developing Countries, Father to Susan E. Rice</td></tr> </table>	<b>Name</b>	Emmett John Rice	<b>Birth date</b>	December 21, 1919	<b>Birth place</b>	Florence, South Carolina, United States	<b>Death date</b>	March 10, 2011 (aged 91)	<b>Death place</b>	Camas, Washington, United States	<b>Nationality</b>	American	<b>Occupation</b>	Governor of the Federal Reserve System, Economics Professor	<b>Known for</b>	Expert in the Monetary System of Developing Countries, Father to Susan E. Rice	<p><b>Reference</b></p> <p>emmett john rice ( december 21 , 1919 – march 10 , 2011 ) was a former governor of the federal reserve system , a Cornell university economics professor , expert in the monetary systems of developing countries and the father of the current national security advisor to president barack obama , susan e . rice .</p> <hr/> <p><b>Content-based attention</b></p> <p>emmett john rice ( december 21 , 1919 – march 10 , 2011 ) was an economist , author , public official and the former american governor of the federal reserve system , the first african american UNK .</p> <hr/> <p><b>Hybrid attention</b></p> <p>emmett john rice ( december 21 , 1919 – march 10 , 2011 ) was an american economist , author , public official and the former governor of the federal reserve system , expert in the monetary systems of developing countries .</p>
<b>Name</b>	Emmett John Rice																
<b>Birth date</b>	December 21, 1919																
<b>Birth place</b>	Florence, South Carolina, United States																
<b>Death date</b>	March 10, 2011 (aged 91)																
<b>Death place</b>	Camas, Washington, United States																
<b>Nationality</b>	American																
<b>Occupation</b>	Governor of the Federal Reserve System, Economics Professor																
<b>Known for</b>	Expert in the Monetary System of Developing Countries, Father to Susan E. Rice																

Figure 6.27: Sample outputs for the given infobox. (Image from Sha et al. [2018]).



**Figure 6.28:** Attention map between words and fields. (a) Content-based attention. (b) Link-based attention. (c) Hybrid attention (Image from Sha et al. [2018]).

## 6.2.5 ROTOWIRE

The **Rotowire** dataset<sup>6</sup> was introduced in [Wiseman et al. \[2017\]](#). It includes summaries of basketball games aligned to the matching box-and-line scores. A sample data point is shown in Figure 6.29. An output of a system<sup>7</sup> for input records describing a basketball game is shown in Figure 6.30.

TEAM	WIN	LOSS	PTS	FG_PCT	RB	AS ...
Heat	11	12	103	49	47	27
Hawks	7	15	95	43	33	20
PLAYER	AS	RB	PT	FG	FGA	CITY ...
Tyler Johnson	5	2	27	8	16	Miami
Dwight Howard	4	17	23	9	11	Atlanta
Paul Millsap	2	9	21	8	12	Atlanta
Goran Dragic	4	2	21	8	17	Miami
Wayne Ellington	2	3	19	7	15	Miami
Dennis Schroder	7	4	17	8	15	Atlanta
Rodney McGruder	5	5	11	3	8	Miami
Thabo Sefolosha	5	5	10	5	11	Atlanta
Kyle Korver	5	3	9	3	9	Atlanta
...						

(a) Rotowire Table.

The Atlanta Hawks defeated the Miami Heat , 103 - 95 , at Philips Arena on Wednesday . Atlanta was in desperate need of a win and they were able to take care of a shorthanded Miami team here . Defense was key for the Hawks , as they held the Heat to 42 percent shooting and forced them to commit 16 turnovers . Atlanta also dominated in the paint , winning the rebounding battle , 47 - 34 , and outscoring them in the paint 58 - 26. The Hawks shot 49 percent from the field and assisted on 27 of their 43 made baskets . This was a near wire - to - wire win for the Hawks , as Miami held just one lead in the first five minutes . Miami ( 7 - 15 ) are as beat - up as anyone right now and it 's taking a toll on the heavily used starters . Hassan Whiteside really struggled in this game , as he amassed eight points , 12 rebounds and one blocks on 4 - of - 12 shooting ...

(b) The corresponding summary.

Figure 6.29: RotoWire examples. (Images from [Wiseman et al. \[2017\]](#)).

A related paper is [Puduppully et al. \[2018\]](#). Figure 6.31 shows the generation model, and the next two figures, Figure 6.32 and Figure 6.33, show a sample template and a system output. This work was followed by [Puduppully et al. \[2019\]](#) and [Puduppully and Lapata \[2021\]](#).

## 6.2.6 WIKITABLET

**WikiTablet**<sup>8</sup> was introduced by [Chen et al. \[2021b\]](#). This dataset is based on Wikipedia sections and the corresponding tables and matching metadata. An example is shown in Figure 6.34.

<sup>6</sup><https://github.com/harvardnlp/boxscore-data>

<sup>7</sup><https://github.com/harvardnlp/data2text>

<sup>8</sup><https://github.com/mingdachen/WikiTableT>

The Utah Jazz ( 38 - 26 ) defeated the Houston Rockets ( 38 - 26 ) 117 - 91 on Wednesday at Energy Solutions Arena in Salt Lake City . The Jazz got out to a quick start in this one , out - scoring the Rockets 31 - 15 in the first quarter alone . Along with the quick start , the Rockets were the superior shooters in this game , going 54 percent from the field and 43 percent from the three - point line , while the Jazz went 38 percent from the floor and a meager 19 percent from deep . The Rockets were able to out - rebound the Rockets 49 - 49 , giving them just enough of an advantage to secure the victory in front of their home crowd . The Jazz were led by the duo of Derrick Favors and James Harden . Favors went 2 - for - 6 from the field and 0 - for - 1 from the three - point line to score a game - high of 15 points , while also adding four rebounds and four assists ....

Figure 6.30: RotoWire output. (Image from Wiseman et al. [2017]).

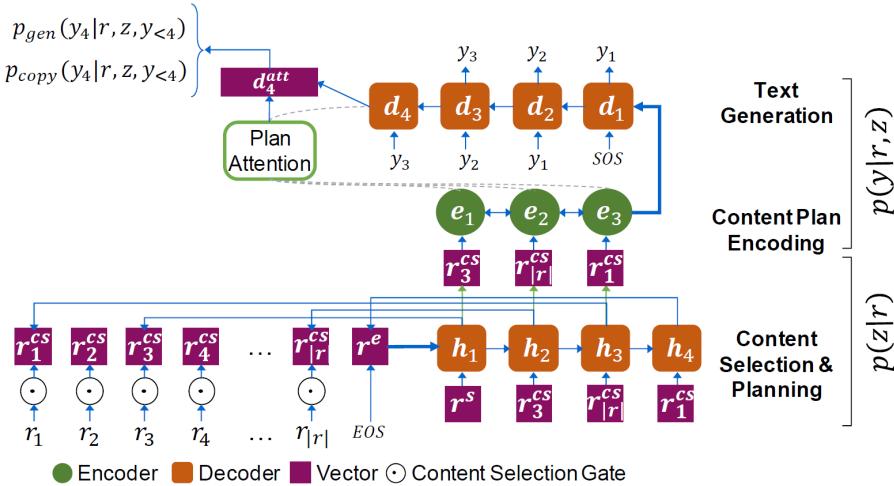


Figure 6.31: Generation model. (Image from Puduppully et al. [2018]).

The Washington Wizards (8-13) defeated the Denver Nuggets (8-15) 92-85. Bradley Beal scored 26 points (8-15 FG, 4-7 3PT, 6-6 FT) to go with 3 rebounds. Nikola Jokic scored 17 points (6-10 FG, 0-0 3PT, 5-7 FT) to go with 11 rebounds. Markieff Morris scored 15 points (5-12 FG, 0-0 3PT, 5-5 FT) to go with 3 rebounds. John Wall scored 15 points (5-14 FG, 0-4 3PT, 5-6 FT) to go with 7 rebounds. Danilo Gallinari scored 14 points (3-11 FG, 1-8 3PT, 7-9 FT) to go with 4 rebounds. Jusuf Nurkic scored 13 points (6-6 FG, 0-0 3PT, 1-2 FT) to go with 7 rebounds. The Washington Wizards' next game will be at home against the Dallas Mavericks, while the Denver Nuggets will travel to play the Bulls.

Figure 6.32: Sample template. (Image from Puduppully et al. [2018]).

## 146 6. DATA TO TEXT

The Washington Wizards defeated the visiting Denver Nuggets 92-85 at Verizon Center on Monday. The Wizards (8-13) came into this game winners of five of their last eight games, but the Wizards (8-15) jumped out to a 10-point lead at the end of the first quarter. Bradley Beal led the way for the Wizards with a game-high 26 points on 8-of-15 shooting from the field. John Wall shot 5-of-14 from the field on his way to 15 points, to go along with seven rebounds, five assists, three steals and one block. Jusuf Nurkic chipped in 13 points, seven rebounds and one assist. Jameer Nelson filled out the stat sheet with 10 points, eight assists, four rebounds and two steals in 36 minutes. As a team, it was a forgettable shooting night for the Nuggets, as the team shot just 46 percent from the field. Next up, the Nuggets play the second game of a back-to-back when they host the Denver Nuggets on Wednesday, while the Wizards host the Portland Trail Blazers on Friday.

Figure 6.33: Output. (Image from Puduppully et al. [2018]).

Journey to the Center of the Earth (also called Jules Verne's Journey to the Center of the Earth) is a 1959 American science fiction adventure film in color by De Luxe, distributed by 20th Century Fox. The film, produced by Charles Brackett and directed by Henry Levin, stars James Mason, Pat Boone, and Arlene Dahl. Bernard Herrmann wrote the film score, and the film's storyline was adapted by Charles Brackett from the 1864 novel of the same name by Jules Verne.

Section Data		Article Data	
Attribute	Value	Attribute	Value
musical composition	20th Century Fox	instance of	film
PERSON	Jules Verne	director	Henry Levin
dependence syndrome	alcoholic	composer	Bernard Herrmann
film genre	adventure film	released	1959, 12, 16
business	Deluxe Entertainment Services Group, Inc.	genre	science fiction film
based on		genre	fantasy film
A Journey to the Center of the Earth		starring	James Mason, Pat Boone, Arlene Dahl
Title Data		Document title	Journey to the Center of the Earth (1959 film)
Section title		Section title	Introduction

Figure 6.34: An example from WikiTableT. (Image from Chen et al. [2021b]).

## 6.3 DOMAIN INDEPENDENT TABLE-TO-TEXT

In recent years, the focus in data-to-text research has shifted significantly toward domain-independent dataset. In this section, we will discuss the recent work in this area.

### 6.3.1 TOTTO

**ToTTo**<sup>9</sup> is an open-domain table-to-text dataset with 120,000 training examples for controlled generation from Wikipedia tables. The annotators were instructed to generate a one-sentence description of a region of a given table. Figure 6.35 illustrates the annotation process which consists of (1) showing an original Wikipedia text and table and asking an annotator to highlight the cells mentioned in text, (2) deleting phrases from original text not supported by the highlighted cells, (3) replacing pronouns with named entities from the table, and (4) polishing the produced text by another annotator. Another example from the dataset is shown in Figure 6.36, and some statistics about the dataset, compared to the data-to-text datasets that were available at the time, is given in Figure 6.37.

**Table Title:** Gabriele Becker  
**Section Title:** International Competitions  
**Table Description:** None

Year	Competition	Venue	Position	Event	Notes
<b>Representing Germany</b>					
1992	World Junior Championships	Seoul, South Korea	10th (semis)	100 m	11.83
1993	European Junior Championships	San Sebastián, Spain	7th	100 m	11.74
			3rd	4x100 m relay	44.60
1994	World Junior Championships	Lisbon, Portugal	12th (semis)	100 m	11.66 (wind: +1.3 m/s)
			2nd	4x100 m relay	44.78
1995	World Championships	Gothenburg, Sweden	7th (q-finals)	100 m	11.54
			3rd	4x100 m relay	43.01

**Original Text:** After winning the German under-23 100 m title, she was selected to run at the 1995 World Championships in Athletics both individually and in the relay.

**Text after Deletion:** she at the 1995 World Championships in both individually and in the relay.

**Text After Decontextualization:** Gabriele Becker competed at the 1995 World Championships in both individually and in the relay.

**Final Text:** Gabriele Becker competed at the 1995 World Championships both individually and in the relay.

**Figure 6.35:** A datapoint from ToTTo showing the annotation process with the text summarizing the highlighted cells. (Image from Parikh et al. [2020]).

### 6.3.2 DART

**DART**<sup>10</sup> is very similar to ToTTo. It is also based on selected areas from a table and the matching human-written verbalizations of these areas. DART was introduced in Nan et al. [2021]. Some statistics about the dataset in comparison with other related datasets are shown in Figure 6.38. Data is collected from both human annotations and automatic

<sup>9</sup><https://github.com/google-research-datasets/totto>

<sup>10</sup><https://github.com/Yale-LILY/dart>

## 148 6. DATA TO TEXT

**Table Title:** Montpellier  
**Section Title:** Climate  
**Table Description:** None

Month	Climate data for Montpellier (1981–2010 averages)												
	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Year
<b>Record high °C (°F)</b>	21.2 (70.2)	22.5 (72.5)	27.4 (81.3)	30.4 (86.7)	35.1 (95.2)	37.2 (99.0)	<b>37.5 (99.5)</b>	36.8 (98.2)	36.3 (97.3)	31.8 (89.2)	27.1 (80.8)	22.0 (71.6)	37.5 (99.5)
<b>Average high °C (°F)</b>	11.6 (52.9)	12.8 (55.0)	15.9 (60.6)	18.2 (64.8)	22.0 (71.6)	26.4 (79.5)	29.3 (84.7)	28.9 (84.0)	25.0 (77.0)	20.5 (68.9)	15.3 (59.5)	12.2 (54.0)	19.9 (67.8)
<b>Daily mean °C (°F)</b>	7.2 (45.0)	8.1 (46.6)	10.9 (51.6)	13.5 (56.3)	17.3 (63.1)	21.2 (70.2)	24.1 (75.4)	23.7 (74.7)	20.0 (74.7)	16.2 (68.0)	11.1 (61.2)	8.0 (46.4)	15.1 (59.2)
<b>Average low °C (°F)</b>	2.8 (37.0)	3.3 (37.9)	5.9 (42.6)	8.7 (47.7)	12.5 (54.5)	16.0 (60.8)	18.9 (66.0)	18.5 (65.3)	15.0 (59.0)	11.9 (53.4)	6.8 (44.2)	3.7 (38.7)	10.4 (50.7)
<b>Record low °C (°F)</b>	-15 (5)	<b>-17.8 (0.0)</b>	-9.6 (14.7)	-1.7 (28.9)	0.6 (33.1)	5.4 (41.7)	8.4 (47.1)	8.2 (46.8)	3.8 (38.8)	-0.7 (30.7)	-5 (23) (9.7)	-12.4 (0.0)	-17.8 (0.0)
<b>Average precipitation mm (inches)</b>	55.6 (2.19)	51.8 (2.04)	34.3 (1.35)	55.5 (2.19)	42.7 (1.68)	27.8 (1.09)	16.4 (0.65)	34.4 (1.35)	80.3 (3.16)	96.8 (3.81)	66.8 (2.63)	66.7 (2.63)	629.1 (24.77)
<b>Average precipitation days</b>	5.5	4.4	4.7	5.7	4.9	3.6	2.4	3.6	4.6	6.8	6.1	5.6	57.8
<b>Average snowy days</b>	0.6	0.7	0.3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.1	0.7	2.4
<b>Average relative humidity (%)</b>	75	73	68	68	70	66	63	66	72	77	75	76	70.8
<b>Mean monthly sunshine hours</b>	142.9	168.1	220.9	227.0	263.9	312.4	339.7	298.0	241.5	168.6	148.8	136.5	2,668.2
Source #1: Météo France													
Source #2: Infoclimat.fr (humidity and snowy days, 1961–1990)													

**Target sentence:** Extreme temperatures of Montpellier have ranged from -17.8 °C recorded in February and up to 37.5 °C (99.5 °F) in July.

**Figure 6.36:** A datapoint from ToTTo, with “interesting reference language”. (Image from Parikh et al. [2020]).

Dataset	Train Size	Domain	Target Quality	Target Source	Content Selection
Wikibio (Lebret et al., 2016)	583K	Biographies	Noisy	Wikipedia	Not specified
Rotowire (Wiseman et al., 2017)	4.9K	Basketball	Noisy	Rotowire	Not specified
WebNLG (Gardent et al., 2017b)	25.3K	15 DBpedia categories	Clean	Annotator Generated	Fully specified
E2E (Novikova et al., 2017)	50.6K	Restaurants	Clean	Annotator Generated	Partially specified
LogicNLG (Chen et al., 2020)	28.5K	Wikipedia (open-domain)	Clean	Annotator Generated	Columns via entity linking
<b>TOTTo</b>	<b>120K</b>	<b>Wikipedia (open-domain)</b>	<b>Clean</b>	<b>Wikipedia (Annotator Revised)</b>	<b>Annotator highlighted</b>

**Figure 6.37:** Comparing data-to-text datasets. (Image from Parikh et al. [2020]).

annotations. For the former, as shown in Figure 6.39, the annotation is done by (1) collecting parent-child relations between columns (top panel), (2) constructing an ontology from the annotations and selecting the nodes that correspond to the highlighted cells (middle panel), and (3) extracting triples. For the latter, WikiSQL questions are mapped to declarative sentences and the cells are highlighted based on the provided answers and/or SQL queries. This is turning text-to-data annotations into a data-to-text dataset. An example from DART that uses the table title is shown in Figure 6.40, and one more example is shown in Figure 6.41.

	Input Unit	Examples	Vocab Size	Words per SR	Sents per SR	Tables
WikiTableText	Row	13,318	—	13.9	1.0	4,962
LogicNLG	Table	37,015	122K	13.8	1.0	7,392
ToTTo	Highlighted Cells	136,161	136K	17.4	1.0	83,141
DART	Triple Set	82,191	33.2K	21.6	1.5	5,623

Figure 6.38: Data-to-text dataset statistics. (Image from Nan et al. [2021]).

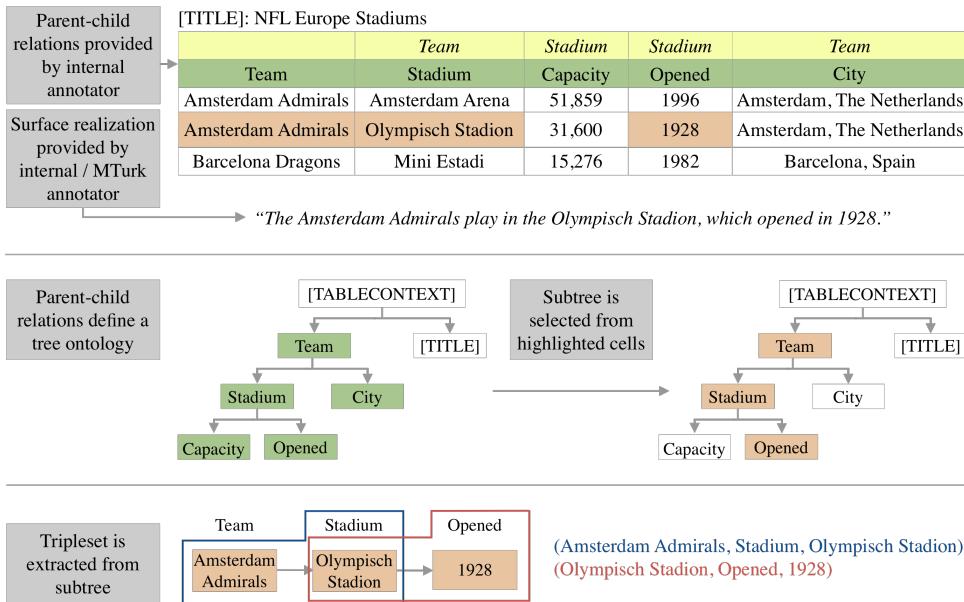
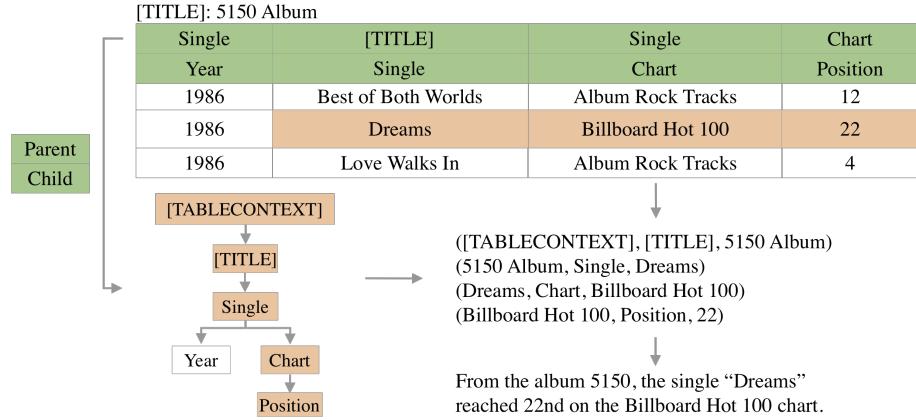


Figure 6.39: Human annotation interface for DART. (Image from Nan et al. [2021]).

### 6.3.3 FETAQA

FetaQA was introduced by Nan et al. [2022b]. As the name suggests, this is a table question-answering dataset with a focus on complex reasoning than simple schema comprehension. The dataset includes 10K Wikipedia-based tuples that include a table,

## 150 6. DATA TO TEXT



**Figure 6.40:** An example from DART that uses the table title. (Image from Nan et al. [2021]).

### **Input triples:**

<H> Andrew Rayel <R> associated Band/associated Musical Artist <T> Christian Burns  
 <H> Andrew Rayel <R> associated Band/associated Musical Artist <T> Jonathan Mendelsohn

### **Reference:**

andrew rayel , is associated with musical artist jonathan mendelsohn and christian burns .

### **train on WebNLG - BART-base output:**

christian mendelsohn and andrew rayel are both associated with the same band , christian burns .

### **train on DART - BART-base output:**

andrew rayel is associated with christian burns and jonathan mendelsohn .

**Figure 6.41:** A sample DART data point. (Image from Nan et al. [2021]).

### 6.3. DOMAIN INDEPENDENT TABLE-TO-TEXT 151

a question, a free-form answer, and the supporting table cells. The data is sampled from ToTTo based on some constraints on table size and the number of highlighted rows. Answers were selected from table-grounded sentences in the dataset and the questions were collected from human annotators. The annotators were asked to write questions that were answered by the ToTTo sentence. Examples from FeTaQA are shown in Figure 6.42. The FeTaQA annotation interface is shown in Figure 6.43.

(a) Page Title: German submarine U-60 (1939)				
Date	Ship	Nationality	Tonnage (GRT)	Fate
19 December 1939	City of Kobe	United Kingdom	4,373	Sunk (Mine)
13 August 1940	Nils Gorthon	Sweden	1,787	Sunk
31 August 1940	Volendam	Netherlands	15,434	Damaged
3 September 1940	Ulva	United Kingdom	1,401	Sunk

Q: How destructive was the U-60?  
A: U-60 sank three ships for a total of 7,561 GRT and damaged another one of 15,434 GRT.

(b) Page Title: High-deductible health plan				
Year	Minimum deductible (single)	Minimum deductible (family)	Maximum out-of-pocket (single)	Maximum out-of-pocket (family)
2016	\$1,300	\$2,600	\$6,550	\$13,100
2017	\$1,300	\$2,600	\$6,550	\$13,100
2018	\$1,350	\$2,700	\$6,650	\$13,300

Q: What is the high-deductible health plan's latest maximum yearly out-of-pocket expenses?  
A: In 2018, a high-deductible health plan's yearly out-of-pocket expenses can't be more than \$6,650 for an individual or \$13,300 for a family.

(c) Page Title: 1964 United States presidential election in Illinois			
Party	Candidate	Votes	%
Democratic	Lyndon B. Johnson (inc.)	2,796,833	59.47%
Republican	Barry Goldwater	1,905,946	40.53%
Write-in		62	0.00%
Total votes		4,702,841	100.00%

Q: How did Lyndon B. Johnson fare against his opponent in the Illinois presidential election?  
A: Lyndon B. Johnson won Illinois with 59.47% of the vote, against Barry Goldwater, with 40.53% of the vote.

(d) Page Title: Joshua Jackson			
Year	Title	Role	Notes
1998–2003	Dawson's Creek	Pacey Witter	124 episodes
2000	The Simpsons	Jesse Grass	Voice; Episode: "Lisa the Tree Hugger"
2001	Cubix	Brian	Voice

Q: Did Joshua Jackson ever star in The Simpsons?  
A: In 2000, Joshua Jackson starred in The Simpsons, voicing the character of Jesse Grass in the episode "Lisa the Tree Hugger".

Figure 6.42: FetaQA examples. (Image from Nan et al. [2022b]).

Page Title: German submarine U-60 (1939)  
Section Title: Summary of raiding History  
Table Section Text: None  
Src url: [http://en.wikipedia.org/wiki/German\\_submarine\\_U-60\\_\(1939\)](http://en.wikipedia.org/wiki/German_submarine_U-60_(1939))

Edit Cells Disable Coloring Edit Sentences Save Changes

Date	Ship	Nationality	Tonnage (GRT)	Fate
19 December 1939	City of Kobe	United Kingdom	4,373	Sunk (Mine)
13 August 1940	Nils Gorthon	Sweden	1,787	Sunk
31 August 1940	Volendam	Netherlands	15,434	Damaged
3 September 1940	Ulva	United Kingdom	1,401	Sunk

Return Previous Page Next Page

Sentence(s):

1. U-60 sank three ships for a total of 7,561 GRT and damaged another one of 15,434 GRT.

Figure 6.43: FetaQA Annotation interface. (Image from Nan et al. [2022b]).

Two benchmark models are suggested for the task, as shown in Figure 6.44: (1) a pipeline model that first parses a question and the table that has the answer into some denotations before generating an answer using a data-to-text model; and (2) an encoder-decoder model based on large pretrained models.

152 6. DATA TO TEXT

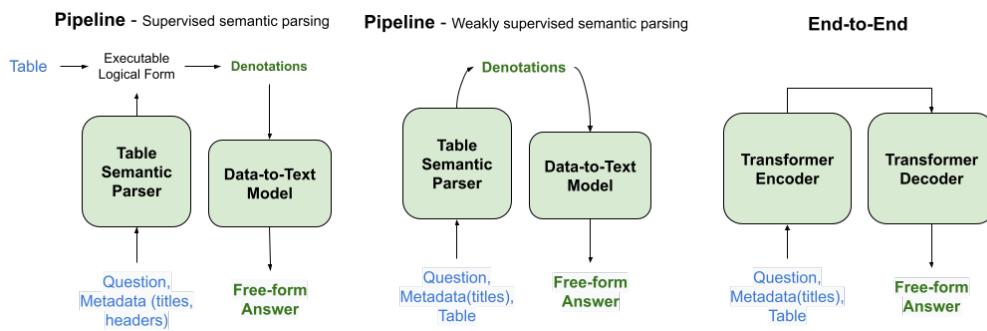


Figure 6.44: FetaQA model diagrams. (Image from Nan et al. [2022b]).

### 6.3.4 TABFACT

**TabFact**<sup>11</sup> [Chen et al., 2019] introduces fact verification for tables. The dataset includes 16,000 Wikipedia tables and 118,000 manually annotated statements in natural language, labeled as either entailed or refuted. The authors describe two models, Table-BERT and Latent Program Algorithm (LPA). The former is used to linearize and encode the tables and statements into continuous vectors. Then LPA parses the statements and executes the programs against the tables. An example is shown in Figure 6.45.

United States House of Representatives Elections, 1972				
District	Incumbent	Party	Result	Candidates
California 3	John E. Moss	democratic	re-elected	John E. Moss (d) 69.9% John Rakus (r) 30.1%
California 5	Phillip Burton	democratic	re-elected	Phillip Burton (d) 81.8% Edlo E. Powell (r) 18.2%
California 8	George Paul Miller	democratic	lost renomination democratic hold	Pete Stark (d) 52.9% Lew M. Warden , Jr. (r) 47.1%
California 14	Jerome R. Waldie	republican	re-elected	Jerome R. Waldie (d) 77.6% Floyd E. Sims (r) 22.4%
California 15	John J. Mcfall	republican	re-elected	John J. Mcfall (d) unopposed

Entailed Statement

- 1. John E. Moss and Phillip Burton are both re-elected in the house of representative election.
- 2. John J. Mcfall is unopposed during the re-election.
- 3. There are three different incumbents from democratic.

Refuted Statement

- 1. John E. Moss and George Paul Miller are both re-elected in the house of representative election.
- 2. John J. Mcfall failed to be re-elected though being unopposed.
- 3. There are five candidates in total, two of them are democrats and three of them are republicans.

Figure 6.45: TabFact example. (Image from Chen et al. [2019]).

### 6.3.5 LOGICNLG

**LogicNLG** was introduced in [Chen et al., 2020b]. This innovative project focuses on generating statements that can be logically entailed from the data in open-domain semi-structured tables. A sample data point is shown in Figure 6.46. A comparison with other similar datasets is shown in Figure 6.47. Sample outputs can be seen in Figure 6.48. Finally, Figure 6.49 illustrates the evaluation of semantic parsing in the paper.

### 6.3.6 GEM

The **GEM** benchmark [Gehrmann et al., 2021] is a more recent task that includes many datasets, such as the aforementioned WebNLG, DART, and ToTTo, as well as the summarization tasks MLSum [Scialom et al., 2020], XSum [Narayan et al., 2018], and WikiLingua [Ladhak et al., 2020]. Figure 6.50 lists these datasets. A multilingual version, GEM v.2 was, recently released [Gehrmann et al., 2022] and it covers 40 datasets in 51 languages.

## 6.4 PRETRAINING FOR TABLES

When dealing with hybrid data (e.g., tables and text), it is important to learn joint representations of the different modalities (see Figure 6.51 for the structure of a

<sup>11</sup><https://github.com/wenhuchen/Table-Fact-Checking>

## 154 6. DATA TO TEXT

Medal Table from Tournament				
Nation	Gold Medal	Silver Medal	Bronze Medal	Sports
Canada	3	1	2	Ice Hockey
Mexico	2	3	1	Baseball
Colombia	1	3	0	Roller Skating

Surface-level Generation				
<b>Sentence:</b> Canada has got 3 gold medals in the tournament.				
<b>Sentence:</b> Mexico got 3 silver medals and 1 bronze medal.				
Logical Natural Language Generation				
<b>Sentence:</b> Canada obtained 1 more gold medal than Mexico.				
<b>Sentence:</b> Canada obtained the most gold medals in the game.				

Figure 6.46: Sample from LogicNLG. (Image from Chen et al. [2020b]).

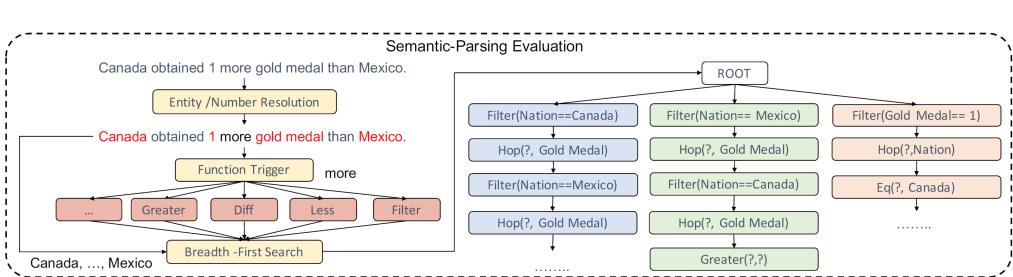
	Vocab	Examples	Vocab/Sent	Tables	Domain	Source	Inference	Schema
WEATHERGOV	394	22.1K	0.01	22.1K	Weather	Crawled	No	Known
WikiBIO	400K	728K	0.54	728K	Biography	Crawled	No	Limited
ROTOWIRE	11.3K	4.9K	0.72	4.9K	NBA	Annotated	Few	Known
LOGICNLG	122K	37.0K	<b>3.31</b>	7.3K	<b>Open</b>	Annotated	<b>Rich</b>	<b>Unlimited</b>

Figure 6.47: Statistics about LogicNLG and related datasets. (Image from Chen et al. [2020b]).

player	country	year (s) won	total	to par
larry nelson	united states	1981 , 1987	152	+ 8
jack nicklaus	united states	1963 , 1971 , 1973 1975 , 1980	152	+ 8
lee Trevino	united states	1974 , 1984	152	+ 8
hubert green	united states	1985	153	+ 9
lanny wadkins	united states	1977	155	+ 11
dave stockton	united states	1970 , 1976	157	+ 13

larry nelson , jack nicklaus , and lee Trevino all shot 8 strokes over par  
 larry nelson , lee Trevino , and dave stockton each won two pga championships in the 1970s - 1980s  
 jack nicklaus had more pga championship wins than larry nelson and lee Trevino combined  
 dave stockton shot five strokes worse than larry nelson , jack nicklaus , and lee Trevino  
 three golfers shot worse than 8 strokes over par

Figure 6.48: Statements generated as part of LogicNLG. (Image from Chen et al. [2020b]).



**Figure 6.49:** Semantic Parsing Evaluation for LogicNLG. (Image from [Chen et al. \[2020b\]](#)).

Dataset	Communicative Goal	Language(s)	Size	Input Type
CommonGEN (Lin et al., 2020)	Produce a likely sentence which mentions all of the source concepts.	en	67k	Concept Set
Czech Restaurant (Dušek and Jurčíček, 2019)	Produce a text expressing the given intent and covering the specified attributes.	cs	5k	Meaning Representation
DART (Radev et al., 2020)	Describe cells in a table, covering all information provided in triples.	en	82k	Triple Set
E2E clean (Novikova et al., 2017) (Dušek et al., 2019)	Describe a restaurant, given all and only the attributes specified on the input.	en	42k	Meaning Representation
MLSum (Scialom et al., 2020)	Summarize relevant points within a news article	*de/es	*520k	Articles
Schema-Guided Dialog (Rastogi et al., 2020)	Provide the surface realization for a virtual assistant	en	*165k	Dialog Act
ToTTo (Parikh et al., 2020)	Produce an English sentence that describes the highlighted cells in the context of the given table.	en	136k	Highlighted Table
XSum (Narayan et al., 2018)	Highlight relevant points in a news article	en	*25k	Articles
WebNLG (Gardent et al., 2017)	Produce a text that verbalises the input triples in a grammatical and natural way.	en/ru	50k	RDF triple
WikiAuto + Turk/ASSET (Jiang et al., 2020) (Alva-Manchego et al., 2020)	Communicate the same information as the source sentence using simpler words and grammar.	en	594k	Sentence
WikiLingua (Ladhak et al., 2020)	Produce high quality summaries of an instructional article.	*en/es/ru/tr/vi	*175k	Article

**Figure 6.50:** Datasets included in GEM v.1. (Image from [Gehrmann et al. \[2021\]](#)).

## 156 6. DATA TO TEXT

Wikipedia table). Many recent language models for tables have focused on just that. These transformer-based tabular models (TaLMs), usually built on top of pre-trained models such as BERT and T5, not only learn a joint representation of content and structure but also inherit some of the semantic and text understanding features of the underlying pre-trained models.

Many of TaLMs have adopted an encoder architecture to learn a contextual representation of tables. These include TaPas [Herzig et al., 2020], TaBERT [Yin et al., 2020], TURL Deng et al. [2020], GraPPa [Yu et al., 2020], TABBIE Iida et al. [2021], TUTA Wang et al. [2021], TableFormer Yang et al. [2022] and others. Some TaLMs have used an encode-decoder architecture to better support text generation tasks such as table-to-text. These include KGPT Chen et al. [2020c], RPT Tang et al. [2021], TaPEX Liu et al. [2021], UnifiedSKG Xie et al. [2022], STTP Xing and Wan [2021], etc. TableGPT<sup>12</sup> adopts a decoder architecture and is fine-tuned on GPT-2.

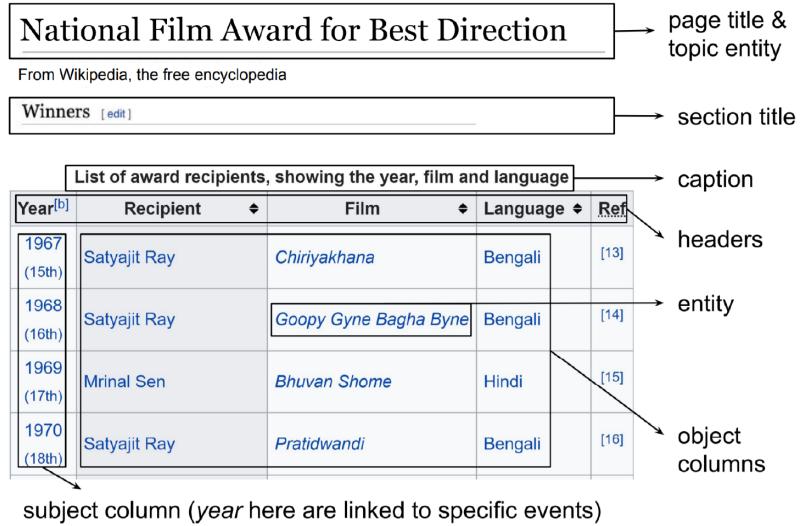


Figure 6.51: Wikipedia Table Structure. (Image from Deng et al. [2020]).

We will look at models for pre-training from tables and text such as TaPas [Herzig et al., 2020], TaBERT [Yin et al., 2020], GraPPa [Yu et al., 2020], and others. See Dong et al. [2022] for a partial list of such model architectures.

### 6.4.1 TURL

TURL<sup>12</sup> [Deng et al., 2020] introduces a pre-training and finetuning approach to relational tables. The pre-training stage is used to learn deep contextualized representations of the tables. Different table components are separately encoded and are fused together. A structure-aware transformer encoder is used to model the table rows and

<sup>12</sup><https://github.com/sunlab-osu/TURL>

columns. A Masked Entity Recovery objective is used to pre-train the system, allowing each element to attend to other table elements on the same row or column. The system is evaluated on six table understanding tasks, including cell filling, row population, entity linking, column type annotation and relation extraction. The TURL architecture is shown in Figure 6.52.

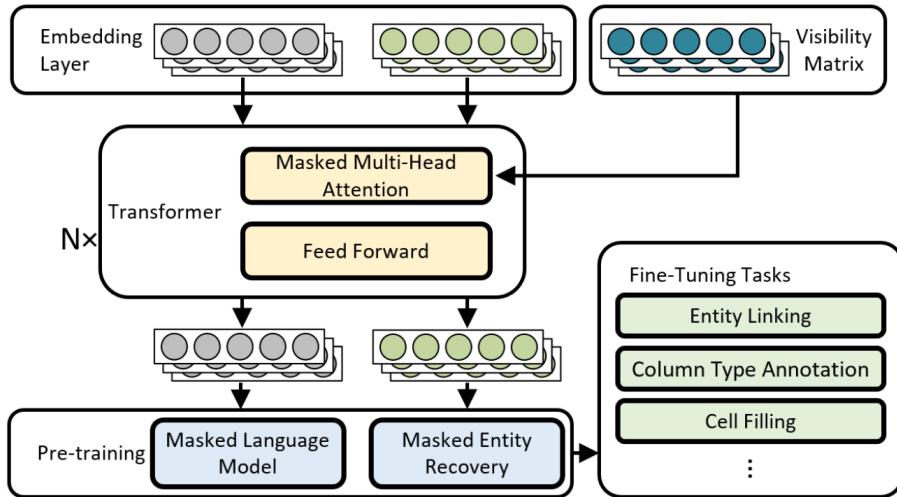


Figure 6.52: Architecture of TURL. (Image from Deng et al. [2020]).

#### 6.4.2 TUTA

Wang et al. [2021] describe a unified pre-training architecture for tables. In order to understand a table, the paper introduces a tree-based structure used to describe the spatial and hierarchical information in tables, followed by tree-based attention and position embeddings. TUTA<sup>13</sup> is evaluated on five datasets for cell type classification and table type classification. An example from TUTA is shown in Figure 6.53. The TUTA architecture is shown in Figure 6.54.

#### 6.4.3 TAPAS

TAPAS [Herzig et al., 2020] is based on the idea of adding table layout information to the pretrained model. It modifies the BERT architecture so that it can be pretrained on tables and the matching text segments. It also uses positional embeddings to encode the table structure. It also uses weak supervision in order to use the tables without the need of intermediate logical forms. The denotations are predicted by a selection of table cells. An example of the TaPaS model is shown in Figure 6.55. Figure 6.56 shows an internal representation, while Figure 6.57 shows a sample table from TaPaS and its corresponding questions.

<sup>13</sup>[https://github.com/microsoft/TUTA\\_table\\_understanding/](https://github.com/microsoft/TUTA_table_understanding/)

## 158 6. DATA TO TEXT

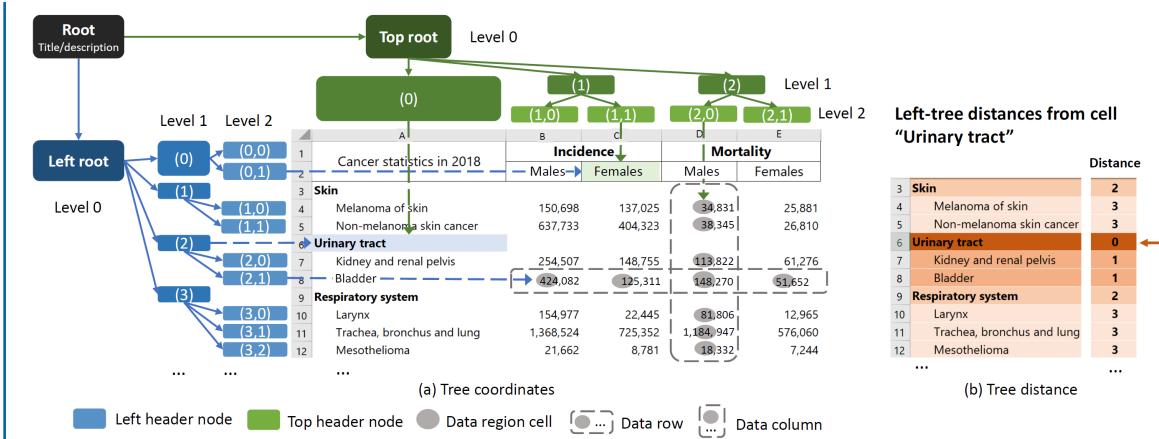


Figure 6.53: Tree coordinates and tree distance for tables, as used in TUTA. (Image from Wang et al. [2021]).

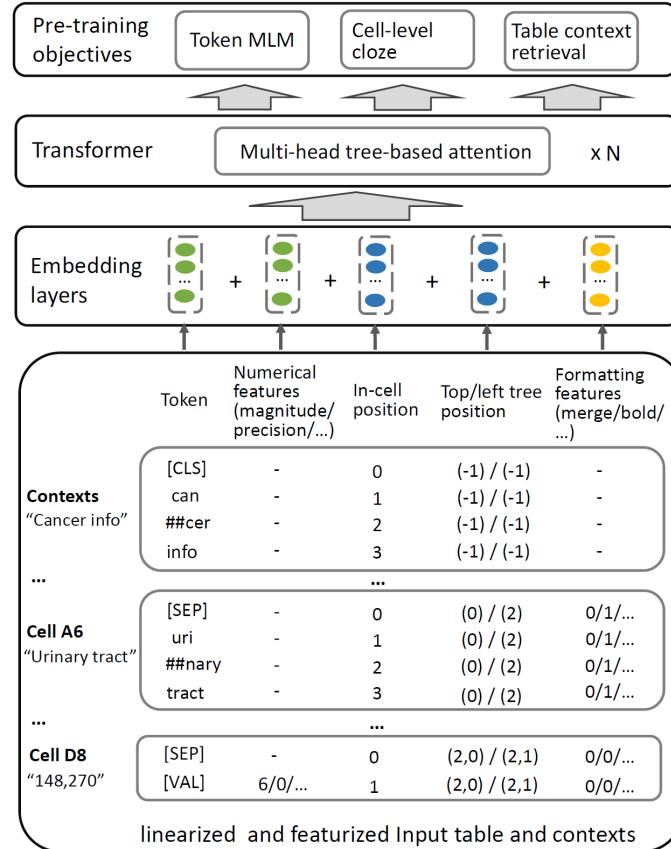


Figure 6.54: The architecture of TUTA. (Image from Wang et al. [2021]).

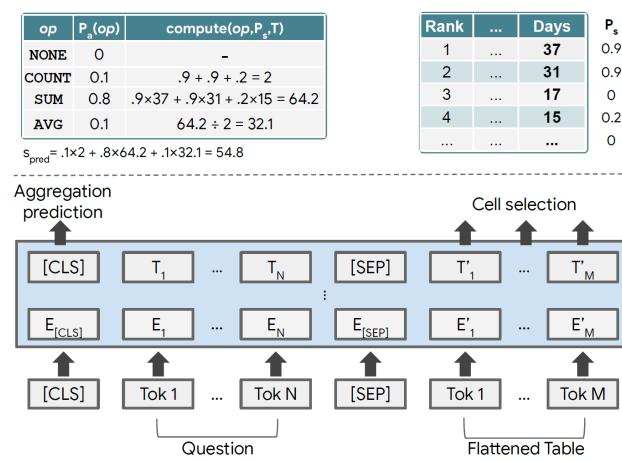


Figure 6.55: TaPaS model. (Image from Herzog et al. [2020]).

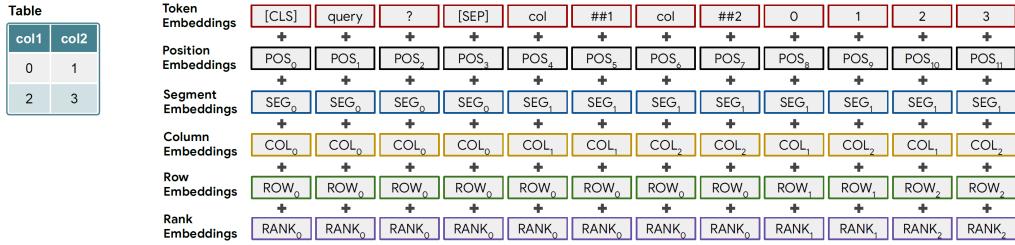


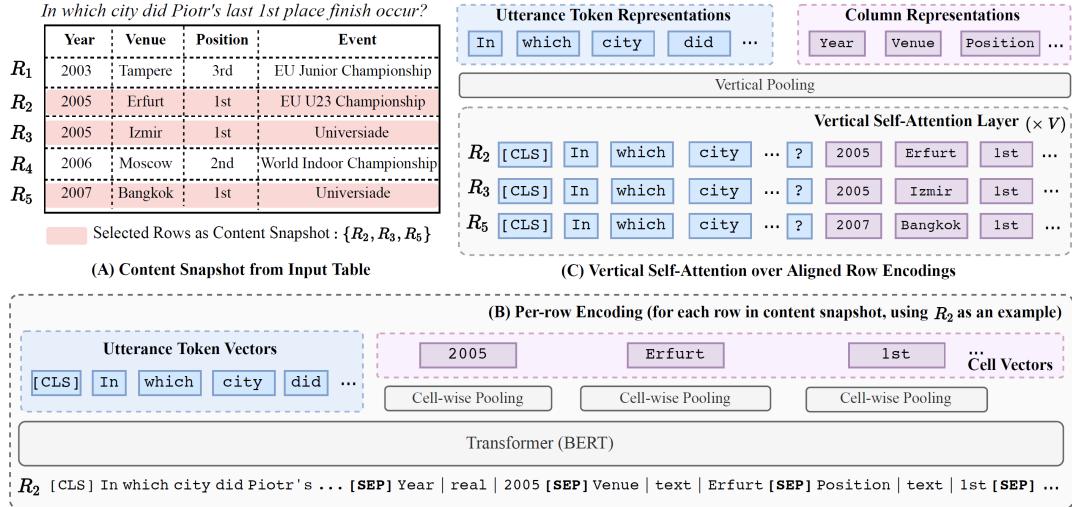
Figure 6.56: Encoding of the question “query?” in TaPaS. (Image from Herzog et al. [2020]).

Table				Example questions											
Rank	Name	No. of reigns	Combined days	#	Question			Answer			Example Type				
1	Lou Thesz	3	3,749	1	Which wrestler had the most number of reigns?	Ric Flair								Cell selection	
2	Ric Flair	8	3,103	2	Average time as champion for top 2 wrestlers?	AVG(3749,3103)=3426								Scalar answer	
3	Harley Race	7	1,799	3	How many world champions are there with only one reign?	COUNT(Dory Funk Jr., Gene Kiniski)=2								Ambiguous answer	
4	Dory Funk Jr.	1	1,563	4	What is the number of reigns for Harley Race?	7								Ambiguous answer	
5	Dan Severn	2	1,559		Which of the following wrestlers were ranked in the bottom 3?	{Dory Funk Jr., Dan Severn, Gene Kiniski}								Cell selection	
6	Gene Kiniski	1	1,131		Out of these, who had more than one reign?	Dan Severn								Cell selection	

Figure 6.57: A table from TaPaS and some of the corresponding questions. (Image from Herzog et al. [2020]).

#### 6.4.4 TABERT

The semantic parser used by **TaBert** [Yin et al., 2020] was built on top of Bert and then trained on 26 million linearized tables and their contexts in English. It achieved state of the art on the WikiTableQuestions dataset and also did very well on the Spider dataset. Figure 6.58 shows how TaBert represents utterances and tables.



**Figure 6.58:** TaBert representations of utterances and tables from WikiTableQuestions. (Image from Yin et al. [2020]).

#### 6.4.5 GRAPPA

**Grappa** Yu et al. [2020] uses a synchronous context-free grammar to generate synthetic question-SQL pairs. It can learn compositional inductive bias through techniques commonly used for semantic parsing, namely masked language modeling (MLM) on joint table+language datasets. Combined with semantic parsers, Grappa achieves SOTA on four tasks: full supervised, weakly supervised x WikisQL, WikiTableQuestions. Figure 6.59 shows the Grappa pre-training approach. The next Figure 6.60 shows some of the non-terminals and production rules using the SCFG for generating synthetic training data for Grappa.

Grappa is different from related work such as TaBERT and TaPas. These two earlier systems are trained on millions of web tables aligned with noisy nearby content. In contrast, Grappa is trained on synthetic data plus a smaller but cleaner set of existing text+table datasets. The four table-based semantic parsing and question answering datasets used in Grappa are shown in Figure 6.61. The top two are fully supervised, while the other two are weakly-supervised. Figure 6.62 shows examples of the inputs and annotations for the four semantic parsing tasks.

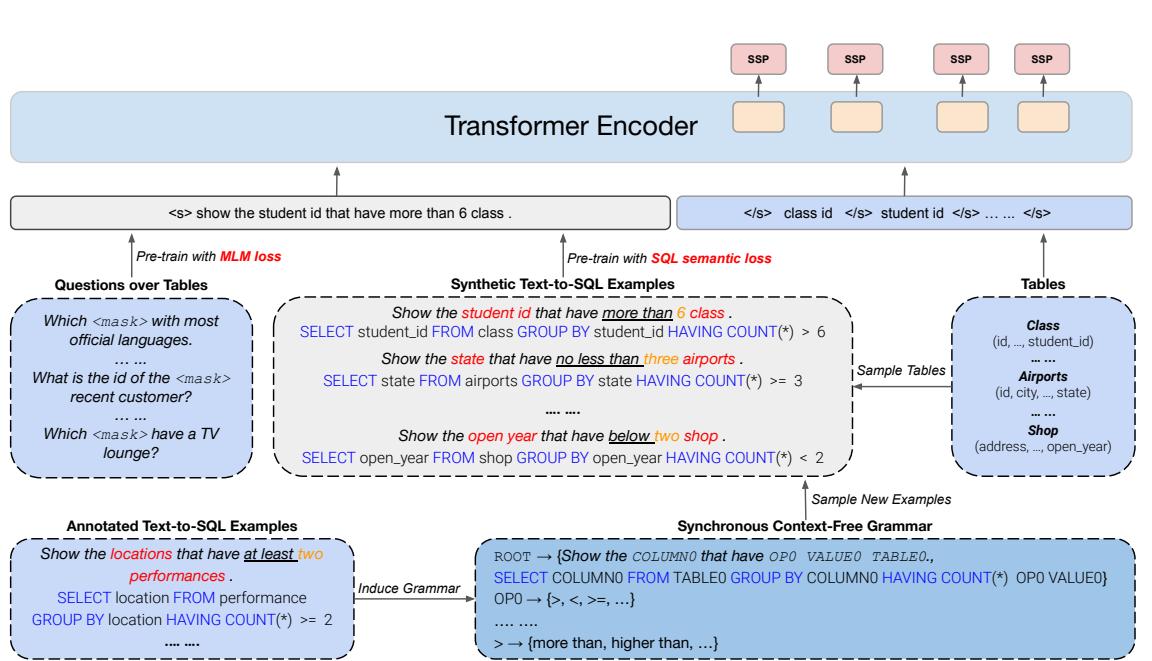


Figure 6.59: Grappa pre-training approach. (Image from Yu et al. [2020]).

Non-terminals	Production rules
TABLE → $t_i$	1. ROOT → ‘‘For each COLUMN0 , return how many times TABLE0 with COLUMN1 OP0 VALUE0 ?’’, SELECT COLUMN0 , COUNT ( * ) WHERE COLUMN1 OP0 VALUE0 GROUP BY COLUMN0 ’’
COLUMN → $c_i$	2. ROOT → ‘‘What are the COLUMN0 and COLUMN1 of the TABLE0 whose COLUMN2 is OP0 AGG0 COLUMN2 ?’’, SELECT COLUMN0 , COLUMN1 WHERE COLUMN2 OP0 ( SELECT AGG0 ( COLUMN2 ) ) ’’
VALUE → $v_i$	
AGG → ⟨ MAX, MIN, COUNT, AVG, SUM ⟩	
OP → ⟨ =, ≤, ≠, …, LIKE, BETWEEN ⟩	
SC → ⟨ ASC, DESC ⟩	
MAX → ‘‘maximum’’, ‘‘the largest’’…	
≤ → ‘‘no more than’’, ‘‘no above’’…	
…	

Figure 6.60: Examples of non-terminals and production rules in the SCFG for Grappa. (Image from Yu et al. [2020]).

Task & Dataset	# Examples	Resource	Annotation	Cross-domain
SPIDER Yu et al. (2018b)	10,181	database	SQL	✓
Fully-sup. WIKISQL Zhong et al. (2017)	80,654	single table	SQL	✓
WIKITABLEQUESTIONS Pasupat & Liang (2015)	2,2033	single table	answer	✓
Weakly-sup. WIKISQL Zhong et al. (2017)	80,654	single table	answer	✓

Figure 6.61: The four table-based semantic parsing and question answering datasets used in Grappa. (Image from Yu et al. [2020]).

Task	Question	Table/Database	Annotation
SPIDER	Find the first and last names of the students who are living in the dorms that have a TV Lounge as an amenity.	database with 5 tables e.g. student, dorm_amenity, ...	<pre>SELECT T1.FNAME, T1.LNAME FROM STUDENT AS T1 JOIN LIVES_IN AS T2 ON T1.STUDID=T2.STUD WHERE T2.DORMID IN ( SELECT T3.DORMID FROM HAS_AMENITY AS T3 JOIN DORM_AMENITY AS T4 ON T3.AMENID=T4.AMENID WHERE T4.AMENITY.NAME='TV LOUNGE')</pre>
Fully-sup. WIKISQL	How many CFL teams are from York College?	a table with 5 columns e.g. player, position, ...	<pre>SELECT COUNT CFL TEAM FROM CFLDRAFT WHERE COLLEGE='YORK'</pre>
WIKITABLEQUESTIONS	In what city did Piotr's last 1st place finish occur?	a table with 6 columns e.g. year, event, ...	"Bangkok, Thailand"
Weakly-sup. WIKISQL	How many CFL teams are from York College?	a table with 5 columns e.g. player, position, ...	2

**Figure 6.62:** Examples from the four semantic parsing tasks. (Image from Yu et al. [2020]).

### 6.4.6 TABBIE

Iida et al. [2021] introduced **Tabbie**. Its goal is different from previous work such as TaPas, TaBERT, and Grappa. It is intended to predict corrupted cell values and therefore is not trained on text. Tabbie’s corrupted cell prediction objective is based on ELECTRA. It uses two transformers, one for rows and another one for columns, as can be seen in Figure 6.63. Tabbie is evaluated on three tasks that measure the level of semantic understanding of tables (in the absence of associated text), namely column population, row population, and column type prediction.

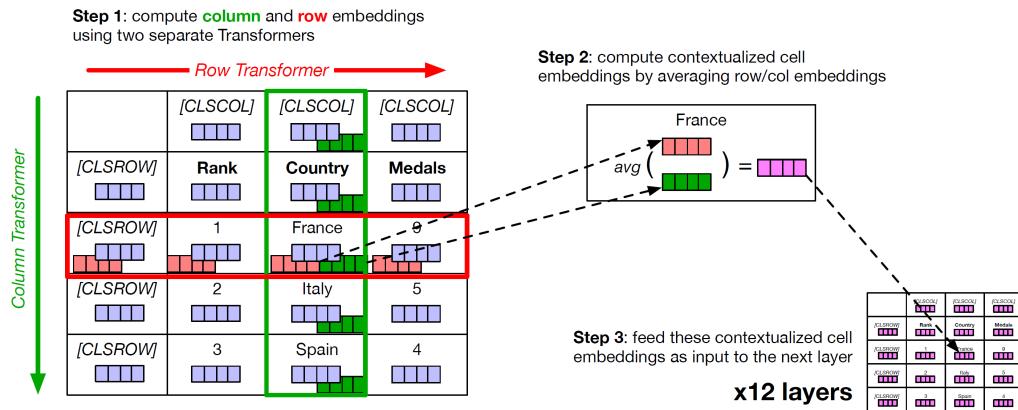


Figure 6.63: Row and column transformers in Tabbie. (Image from Iida et al. [2021]).

### 6.4.7 OTHER PAPERS

**TableFormer** was introduced by Yang et al. [2022]. It represents a structure-aware encoding architecture that transforms table structure information into attention. This structure is invariant to both row and column order. It performed well on the SQA, WikiTableQuestions, and TabFact datasets, and did especially well in settings in which the rows or columns are shuffled. The architecture of TableFormer is shown in Figure 6.65.

Another paper by Herzig et al. [2021], a follow up to TaPaS, is about open domain question answering over tables via dense retrieval. They use a retriever trained to deal with tabular data. Due to the lack of datasets in this domain, they create a new task based on a subset of the Natural Questions [Kwiatkowski et al., 2019] dataset.

**TaPEX**<sup>14</sup> [Liu et al., 2021] trains a neural SQL program executor using a synthetic corpus that includes SQL queries and their execution outputs. The method achieves strong results on WikiSQL, WikiTableQuestions, SQA, and TabFact. An example of TaPEX is shown in Figure 6.67.

Another paper examines pretraining and finetuning strategies for data-to-text tasks is Kale and Rastogi [2020]. They show that text-to-text pretraining as in the

<sup>14</sup><https://github.com/microsoft/Table-Pretraining>

Nation	Gold	Silver	Bronze
Great Britain	2	1	2
<b>Spain</b>	1	2	0
Ukraine	0	2	0

**Question:** Which nation received 2 silver medals?

**Gold Answer:** Spain, Ukraine

**TAPAS:** Spain

**TableFormer:** Spain, Ukraine

**TableFormer w/o a proposed structural bias:**

Spain

Figure 6.64: An example that shows how TableFormer does better than TAPAS. (Image from Yang et al. [2022]).

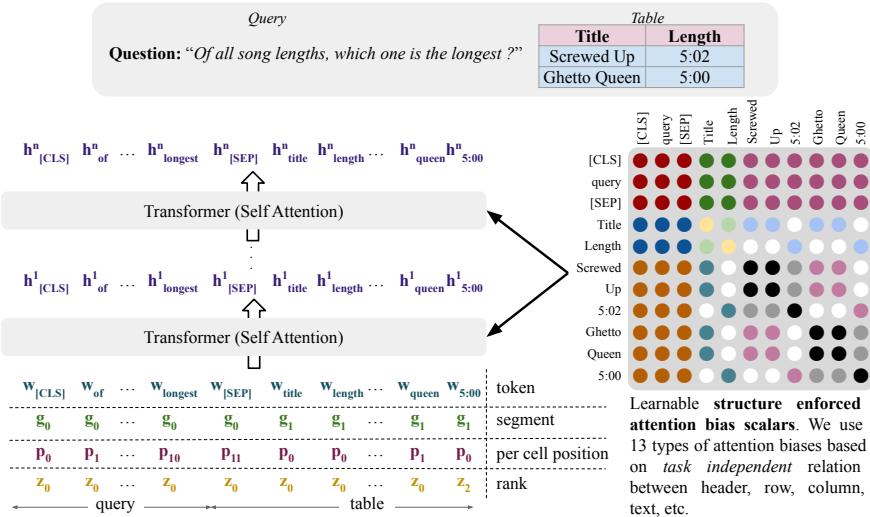


Figure 6.65: TableFormer architecture. (Image from Yang et al. [2022]).

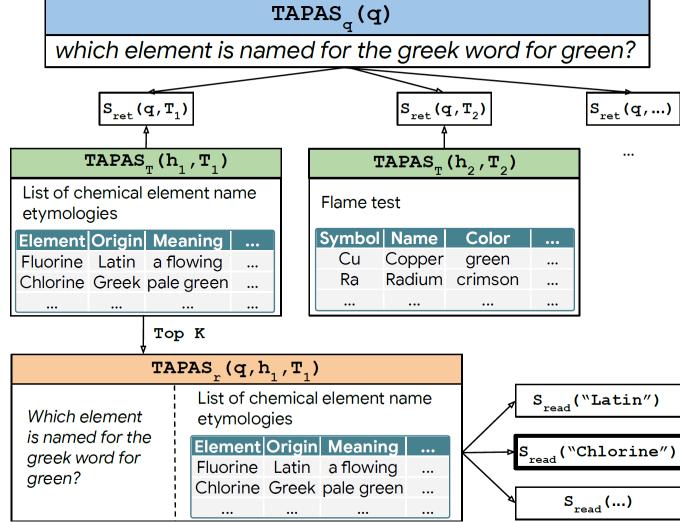


Figure 6.66: A dense table retriever. (Image from Herzig et al. [2021]).

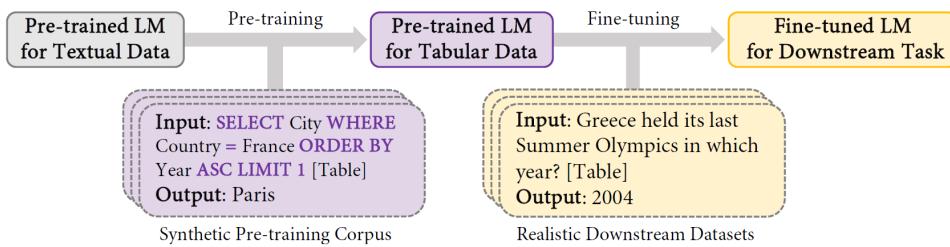


Figure 6.67: TaPEx Example. (Image from Liu et al. [2021]).

## 166 6. DATA TO TEXT

T5 model Raffel et al. [2020] can outperform end-to-end neural architectures specifically designed for data-to-text generation. They obtain strong results on out-of-domain datasets. The dataset sizes are shown in Figure 6.68 and examples of each are shown in Figure 6.69 (WebNLG), and Figure 6.70 (ToTTo).

Dataset	Train	Dev	Test
WebNLG	18.1K	2.2k	4.9k
ToTTo	120K	7.7k	7.7k
Multiwoz	56.8K	7.3k	7.3k

Figure 6.68: Dataset sizes. (Image from Kale and Rastogi [2020]).

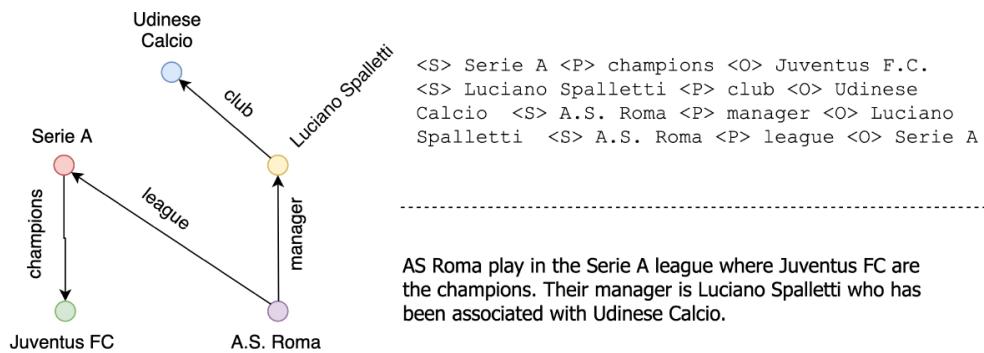


Figure 6.69: WebNLG example. (Image from Kale and Rastogi [2020]).

Table Title: Cristhian Stuani				
Section Title: International goals				
No.	Date	Venue	Opponent	Result
2	13 November 2013	Amman International Stadium, Amman, Jordan	Jordan	5-0

```
<page_title> Cristhian Stuani </page_title>
<section_title> International goals </section_title>
<table> <cell> 2. <col_header> No. </col_header> </cell>
<cell> 13 November 2013 <col_header> Date </col_header>
</cell> <cell> Amman International Stadium, Amman, Jordan <col_header> Venue </col_header> </cell>
<cell> Jordan <col_header> Opponent </col_header> </cell>
<cell> 5-0 <col_header> Result </col_header> </cell>
</table>
```

On 13 November 2013 Cristhian Stuani netted the second in a 5–0 win in Jordan.

Figure 6.70: ToTTo example. (Image from Kale and Rastogi [2020]).

A recent paper<sup>15</sup> [Wang et al., 2022b] makes the claim that models don't really have to be designed with tables in mind. They look at the task of table retrieval and demonstrate that a simple, text-based model can achieve results comparable to these of table-specific models. They compare a generic dense passage retriever, which when

<sup>15</sup><https://github.com/zorazrw/nqt-retrieval>

fine-tuned on linearized tables, performs favorably compared to a dense table retriever on the NQtable dataset which is a table-focused subset of the Natural Questions dataset.

## 6.5 OTHER RECENT WORK

**HybridQA**<sup>16</sup> [Chen et al., 2020d] is a large-scale QA dataset that involves reasoning on heterogeneous information, both Wikipedia tables and free-form text linked to the entities in the tables. Figure 6.71 and Figure 6.72 show examples of HybridQA datapoints.

The 2016 Summer Olympics officially known as the Games of the XXXI Olympiad (Portuguese : Jogos da XXXI Olimpíada) and commonly known as Rio 2016 , was an international multi-sport event .....

Name	Year	Season	Flag bearer
XXXI	<a href="#">2016</a>	Summer	<a href="#">Yan Naing Soe</a>
XXX	<a href="#">2012</a>	Summer	<a href="#">Zaw Win Thet</a>
XXIX	<a href="#">2008</a>	Summer	<a href="#">Phone Myint Tayzar</a>
XXVIII	<a href="#">2004</a>	Summer	<a href="#">Hla Win U</a>
XXVII	<a href="#">2000</a>	Summer	<a href="#">Maung Maung Nge</a>
XX	<a href="#">1972</a>	Summer	<a href="#">Win Maung</a>

Yan Naing Soe ( born [31 January 1979](#) ) is a Burmese judoka . He competed at the 2016 Summer Olympics in the men 's 100 kg event , ..... He was the flag bearer for Myanmar at the Parade of Nations .

Zaw Win Thet ( born [1 March 1991](#) in Kyonpyaw , Pathein District , Ayeyarwady Division , Myanmar ) is a Burmese runner who .....

Myint Tayzar Phone ( Burmese : မြင်တော်ဖျား ) born [July 2 , 1978](#) ) is a sprint canoer from Myanmar who competed in the late 2000s .

.....

Win Maung ( born [12 May 1949](#) ) is a Burmese footballer . He competed in the men 's tournament at the 1972 Summer Olympics ...

Q: In which year did the judoka bearer participate in the Olympic opening ceremony? A: 2016

Q: Which event does the does the XXXI Olympic flag bearer participate in? A: men's 100 kg event

Q: Where does the Burmese judoka participate in the Olympic opening ceremony as a flag bearer? A: Rio

Q: For the Olympic event happening after 2014, what session does the Flag bearer participate? A: Parade of Nations

Q: For the XXXI and XXX Olympic event, which has an older flag bearer? A: XXXI

Q: When does the oldest flag Burmese bearer participate in the Olympic ceremony? A: 1972

Figure 6.71: Annotated question answering pairs in HybridQA. (Image from Chen et al. [2020d]).

Another recent dataset is **OTT-QA**<sup>17</sup> [Chen et al., 2020a]. The acronym stands for Open Table-and-Text Question Answering. Answering such questions requires multi-hop inferences across tables and text. The methods used in the paper include “arly fusion” that combines related table and text unites into a single block, and “cross-block reader” which applies global-local sparse attention. The best results are obtained when the two methods used in combination. An example is included in Figure 6.73. Some of the matching questions are shown in Figure 6.74.

**Turning Tables** [Yoran et al., 2022] is a recent paper that adds reasoning to semi-structured tables. They describe a method for the automatic generation of question-paragraph pairs. These pairs require multiple reasoning skills such as fact composition and number comparison. Figure 6.75 and Figure 6.76 include examples.

**HiTab**<sup>18</sup> [Cheng et al., 2022] moves forward from QA and NLG for flat tables to hierarchical tables (e.g., Figure 6.77). Figure 6.78 shows a HiTab example.

<sup>16</sup><https://github.com/wenhuchen/HybridQA>

<sup>17</sup><https://github.com/wenhuchen/OTT-QA>

<sup>18</sup><https://github.com/microsoft/HiTab>

## 168 6. DATA TO TEXT

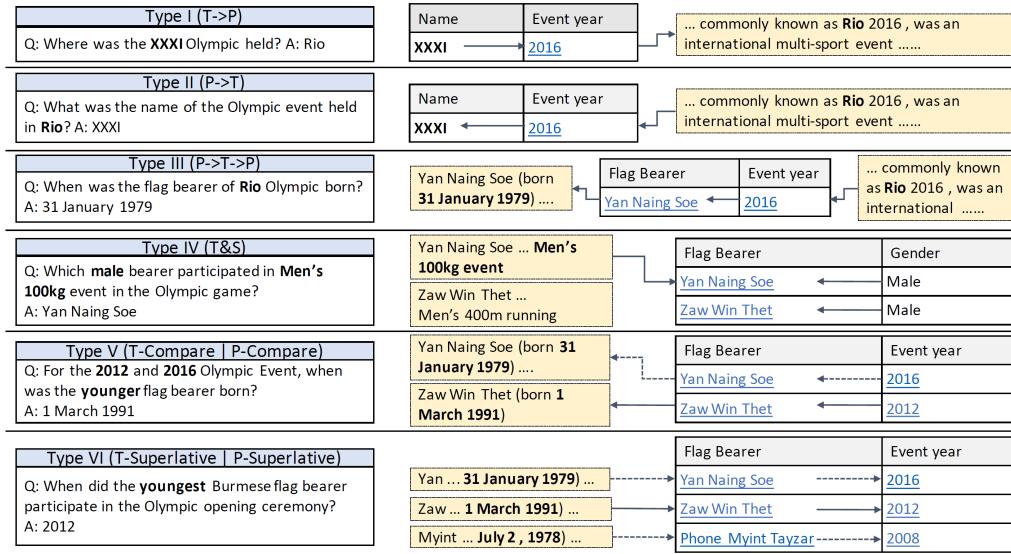


Figure 6.72: Multihop questions in HybridQA. (Image from Chen et al. [2020d]).

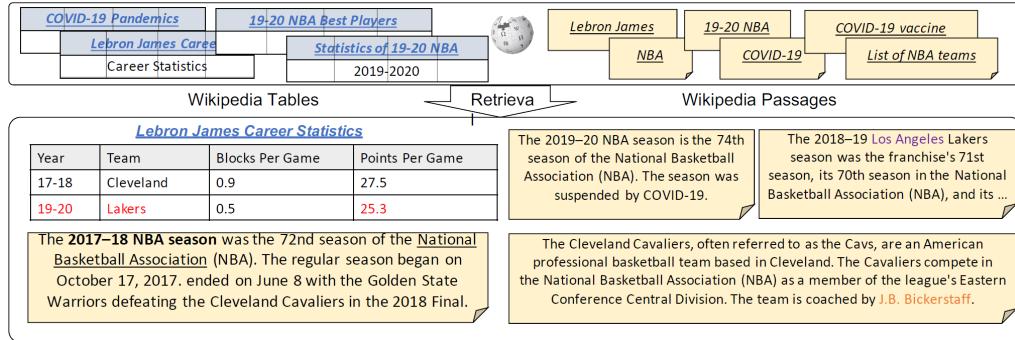


Figure 6.73: Multi-hop reasoning over two candidate pools. (Image from Chen et al. [2020a]).

Q1: How many points per game did Lebron James get in the COVID-19 NBA Season?

A1: COVID-19 -> 19-20 Season -> 25.3

Q2: Who is the coach of the team that Lebron James played in to achieve his highest score in his career?

A2: 27.5 -> Cleveland -> J. B. Bickerstaff

Q3: What suspends the NBA season during which Lebron James has an average points per game of 25.3?

A3: 25.3 -> 19-20 Season -> COVID-19

Q4: For the season suspended by COVID-19 and the season defeated by Warriors in Final, which season has Lebron obtained more points?

A4: 25.3 < 27.5 -> 17-18

Figure 6.74: More sample questions. (Image from Chen et al. [2020a]).

Round	Date	Opponent	Venue	Result	Attendance
R3	31 October 1990	Portsmouth	H	0–0	16,699
R3R	6 November 1990	Portsmouth	A	3–2	16,085
R4	28 November 1990	Oxford United	A	2–1	9,789
QF	16 January 1991	Tottenham Hotspur	H	0–0	34,178
QFR	23 January 1991	Tottenham Hotspur	A	3–0	33,861
SF 1st Leg	24 February 1991	Sheffield Wednesday	H	0–2	34,074
SF 2nd Leg	27 February 1991	Sheffield Wednesday	A	1–3	34,669



**Composition:** q: What was the Result when the Round was R4? c: The Date when the Round was R4 was 28 November 1990. The Result when the Date was 28... a: 2-1

**Comparison:** q: Which Round had a higher Attendance: QF or QFR? c: The Attendance when the Round was QF was 34,178. The Attendance when the Round was QFR... a: QF

**Date Difference:** q: The Opponent was Portsmouth how much time before the Opponent was Sheffield Wednesday? c: The Date when the Opponent... a: 3 months and 18 days

**Figure 6.75:** An example of a table and automatically generated question-context answer triples. (Image from Yoran et al. [2022]).

EG	Template	Question
2/3-hop Composition	What was the col:1(s) when the col:2 was val:2 in table-title of page-title?	"What was the Play(s) when the Author was William Shakespeare in Notable works of Lidia Zamkow?"
Conjunction	What was the col:1 when the col:2 was val:2 and the col:3 was val:3 in table-title of page-title?	"What was the Common name when the Family was Picidae and the Distribution was Okinawa in List of species of List of endemic birds of Japan?"
Quantifiers Only	Is val:1 the only col:1 that has col:2 val:2 in table-title of page-title?	"Is Jean Philippe the only Artist that has Language French in Results of Eurovision Song Contest 1959?"
Quantifiers Every/Most	In table-title of page-title, does [OPERATOR] col:1 have col:2 val:2?	"In Coal of List of Mines in South Africa, does every Mine have Owner Exxaro?"
Num. Comparison	In table-title of page-title, which col:1 had [OPERATOR] col:2: val:1 or val:1?	"In Administration of Mueang Nonthaburi District, which Name had a higher population: Suan Yai or Bang Khen?"
Temp. Comparison	In table-title of page-title, what happened [OPERATOR]: the col:1 was val:1 or the col:2 was val:2?	"In Awards and nominations of Alexandre Pires, what happened earlier: the Category was Pop New Artist or the Category was Album of the Year?"
Num. Yes/No Comparison	In table-title of page-title did val:1 have [OPERATOR] col:2 than val:1?	"In Top employers of Chula Vista, California, did Walmart have more Employees than Target?"
Temp. Yes/No Comparison	The col:1 was val:1 [OPERATOR] the col:2 was val:2 in table-title of page-title?	"The Referee was Salim Oussassi more recently than when the Referee was Rachid Medjiba in 1980 to 1999 of Algerian Cup Final referees?"
Temp./Num. Superlatives	In table-title of page-title, which col:1 has the [OPERATOR] col:2?	"In List of graphic novels of Minx (comics), which Title has the earliest Release date?"
Arithmetic Superlatives	In table-title of page-title, what was the [OPERATOR] col:1 when the col:2 was val:2?	"In By rocket of 1961 in spaceflight, what was the highest Successes when the Remarks was Maiden flight?"
Counting	How many col:1 have col:2 val:2 in table-title of page-title?	"How many Elections have Candidate John Kufuor in Presidential elections of New Patriotic Party?"
Arithmetic Addition	In table-title of page-title, what was the total number of col:1 when the col:2 was val:2?	"In Assists table of 2010-11 La Liga, what was the total number of Assists when the Club was Villarreal?"
Date Difference	In table-title of page-title, how much time had passed between when the col:1 was val:1 and when the col:2 was val:2?	"In Notable events   Concerts of Candlestick Park, how much time had passed between when the Artist was Paul McCartney and when the Artist was The Beatles?"

**Figure 6.76:** Question templates. (Image from Yoran et al. [2022]).

170 6. DATA TO TEXT

	A	B	C	D	E	F	G
Source and mechanism	All full-time graduate students		Master's		Doctoral		
	Total	Percent	All	Percent	All	Percent	
All full-time	433,916	100.0	209,221	100.0	224,695	100.0	
Self-support	161,641	37.3	139,373	66.6	22,268	9.9	
All sources of support	272,275	62.7	69,848	33.4	202,427	90.1	
Federal	65,999	15.2	10,736	5.1	55,263	24.6	
Department of Agricu	2,361	0.5	938	0.4	1,423	0.6	
Department of Defens	8,089	1.9	2,568	1.2	5,521	2.5	
Other	9,098	2.1	3,462	1.7	5,636	2.5	
Institutional	182,135	42.0	52,319	25.0	129,816	57.8	
Other U.S. source	19,432	4.5	5,136	2.5	14,296	6.4	
Foreign	4,709	1.1	1,657	0.8	3,052	1.4	
All mechanisms of support	272,275	62.7	69,848	33.4	202,427	90.1	
Fellowships	39,368	9.1	5,687	2.7	33,681	15.0	
Traineeships	10,945	2.5	1,497	0.7	9,448	4.2	
Research assistantships	103,586	23.9	19,702	9.4	83,884	37.3	
Teaching assistantships	84,499	19.5	22,171	10.6	62,328	27.7	
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 6.77: Sample hierarchical table from HiTab. (Image from [Cheng et al. \[2022\]](#)).

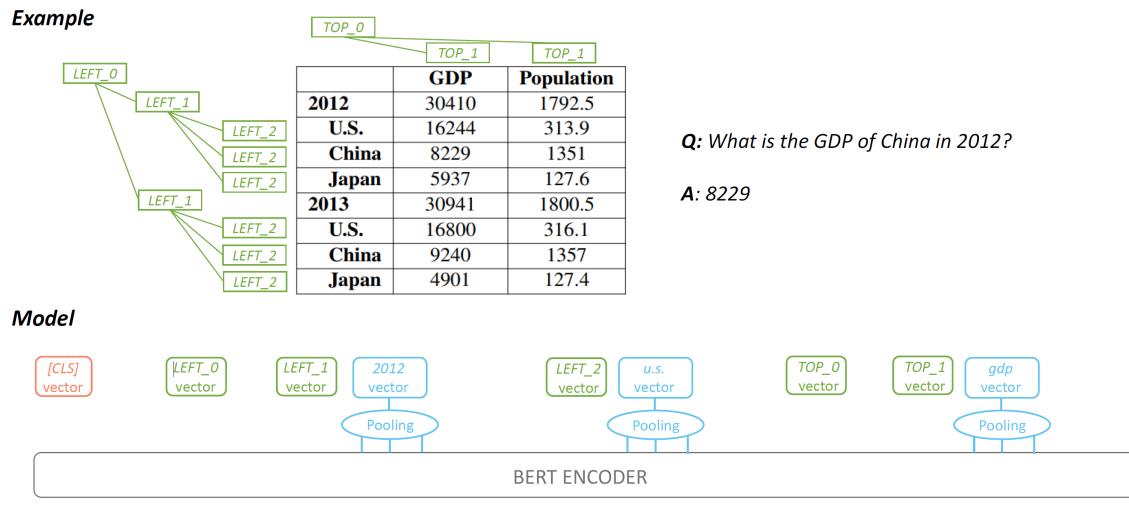


Figure 6.78: HiTab example. (Image from Cheng et al. [2022]).

**R2D2** Nan et al. [2022a] is another recent paper. It focuses on unfaithful data-to-text generation by training a system as both generator and faithfulness discriminator. The resulting R2D2 system achieves state-of-the-art results on FeTaQA, LogicNLG, and ToTTo.

Figure 6.79: Example from R2D2. (Image from Nan et al. [2022a]).

**Logic2Text**<sup>19</sup> Chen et al. [2020e] addresses the problem of generating descriptions of tables that require inference across records. They formulate such a task as generation from logical forms. The dataset includes 10,753 descriptions, paired with the corresponding logical forms. Figure 6.80 provides an illustration.

## 6.6 ENCODING DATABASE INFORMATION IN PLAIN TEXT.

The last work that we want to mention is **NeuralDB**<sup>20</sup> Thorne et al. [2021a,b]. The idea here is to encode all the data into plain text instead of a relational schema. This is a novel approach to answering database queries over plain text representations of tabular data. Their approach scales to databases containing thousands of facts. An example is shown in Figure 6.81.

## 6.7 SUMMARY

In this chapter we looked at generation of natural language text from tabular data. This is still a very active research area and new relevant papers get published on a weekly basis.

For more information, we refer the reader to five recent surveys on text generation Gatt and Krahmer [2017], Narayan and Gardent [2020] Garbacea and Mei [2020], Celikyilmaz et al. [2020], Jin et al. [2020a], Iqbal and Qureshi [2022], Dong et al. [2021], Yu et al. [2022], and Erdem et al. [2022]. A survey of table pretraining is Dong et al. [2022]. Another survey is Borisov et al. [2021].

<sup>19</sup><https://github.com/czyssrs/Logic2Text>

<sup>20</sup><https://github.com/facebookresearch/NeuralDB>

table caption: opec				
country	region	joined opec	population (july 2012)	area (km square)
algeria	africa	1969	37367226	2381740
angola	africa	2007	18056072	1246700
iraq	middle east	1960	31129225	437072
libya	africa	1962	5613380	1759540
nigeria	africa	1971	170123740	923768
...	...	...	...	...

**Surface-level NLG**

**Description:** angola, from the region africa, joined opec in 2007, with an population of 18056072 in 2012.

**Description:** algeria, from the region africa, joined opec in 1969, with an population of 37367226 in 2012.

**Logical-level NLG with logical forms ( our dataset )**

**logical form:** eq { count { filter\_eq { all\_rows ; region ; africa } } ; 4 } = True

```

graph TD
    eq((eq)) --> count((count))
    eq --> four((4))
    count --> filterEq((filter_eq))
    filterEq --> allRows((all_rows))
    filterEq --> region((region))
    filterEq --> africa((africa))
  
```

**Description:** In 2012 in opec, there were 4 member countries from africa.

**logical form:** and { eq { hop { argmax { all\_rows ; joined opec } ; region } ; africa } ; eq { hop { argmax { all\_rows ; joined opec } ; country } ; angola } } = True

```

graph TD
    and((and)) --> eq1((eq))
    and --> eq2((eq))
    eq1 --> africa((africa))
    eq1 --> hop1((hop))
    hop1 --> region1((region))
    hop1 --> argmax1((argmax))
    argmax1 --> allRows1((all_rows))
    argmax1 --> joinedOpec1((joined opec))
    eq2 --> hop2((hop))
    hop2 --> argmax2((argmax))
    argmax2 --> allRows2((all_rows))
    argmax2 --> joinedOpec2((joined opec))
    hop2 --> country((country))
    country --> angola((angola))
  
```

**Description:** In 2012 in opec, angola, from africa, was the latest country to join.

**Figure 6.80:** A sample pair of a table and a description. The bottom part of the figure shows outputs that are based on the logical form. (Image from [Chen et al. \[2020e\]](#)).

**Facts:** (4 of 50 shown)

Nicholas lives in Washington D.C. with Sheryl.  
Sheryl is Nicholas's spouse.  
Teuvo was born in 1912 in Ruskala.  
In 1978, Sheryl's mother gave birth to her in Huntsville.

**Queries:**

Does Nicholas's spouse live in Washington D.C.?  
(Boolean Join) —> TRUE

Who is Sheryl's husband?  
(Lookup) —> Nicholas

Who is the oldest person in the database?  
(Max) —> Teuvo

Who is Sheryl's mother?  
(Lookup) —> NULL

Figure 6.81: A database encoded in plain text by NeuralDB. (Image from Thorne et al. [2021a]).

Some links to additional information can be found in Figure 6.82.

<a href="http://1">http://1</a>
<a href="http://2">http://2</a>

**Figure 6.82:** Additional links for Data to Text.

Acknowledgements ?