```
1 from google.colab import drive
2 drive.mount('/content/drive')
3 !cp -r /content/drive/MyDrive/labs/project4/* .
4 !pip install -r requirements.txt
5 # restart the runtime
6 import os
7 os._exit(00)
```

```
      Downloading websocket_client-1.3.3-py3-none-any.whl (54 kB)
         |████████████████████████████████| 54 kB 2.6 MB/s
      Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.7/dist-packages (from r
      Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.7/dist-packages (
      Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in /usr/local/lib/pyth
      Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.7/dist-packages (f
      Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-packages (from impo
      Requirement already satisfied: decorator in /usr/local/lib/python3.7/dist-packages (from ipyt
      Requirement already satisfied: simplegeneric>0.8 in /usr/local/lib/python3.7/dist-packages (f
      Requirement already satisfied: traitlets>=4.2 in /usr/local/lib/python3.7/dist-packages (from
      Requirement already satisfied: prompt-toolkit<2.0.0,>=1.0.4 in /usr/local/lib/python3.7/dist-
      Requirement already satisfied: pygments in /usr/local/lib/python3.7/dist-packages (from ipyth
      Requirement already satisfied: pickleshare in /usr/local/lib/python3.7/dist-packages (from ip
      Requirement already satisfied: pexpect in /usr/local/lib/python3.7/dist-packages (from ipytho
      Requirement already satisfied: six>=1.9.0 in /usr/local/lib/python3.7/dist-packages (from pro
      Requirement already satisfied: wcwidth in /usr/local/lib/python3.7/dist-packages (from prompt
      Requirement already satisfied: MarkupSafe>=0.23 in /usr/local/lib/python3.7/dist-packages (fr
      Requirement already satisfied: jupyter-core in /usr/local/lib/python3.7/dist-packages (from n
      Requirement already satisfied: bleach in /usr/local/lib/python3.7/dist-packages (from nbconve
      Requirement already satisfied: pandocfilters>=1.4.1 in /usr/local/lib/python3.7/dist-packages
      Requirement already satisfied: testpath in /usr/local/lib/python3.7/dist-packages (from nbcor
      Requirement already satisfied: entrypoints>=0.2.2 in /usr/local/lib/python3.7/dist-packages (
      Requirement already satisfied: mistune<2,>=0.8.1 in /usr/local/lib/python3.7/dist-packages (f
      Requirement already satisfied: defusedxml in /usr/local/lib/python3.7/dist-packages (from nbc
      Requirement already satisfied: fastjsonschema in /usr/local/lib/python3.7/dist-packages (from
      Requirement already satisfied: jsonschema>=2.6 in /usr/local/lib/python3.7/dist-packages (fro
      Requirement already satisfied: attrs>=17.4.0 in /usr/local/lib/python3.7/dist-packages (from
      Requirement already satisfied: importlib-resources>=1.4.0 in /usr/local/lib/python3.7/dist-pa
      Requirement already satisfied: pyrsistent!=0.17.0,!=0.17.1,!=0.17.2,>=0.14.0 in /usr/local/li
      Requirement already satisfied: webencodings in /usr/local/lib/python3.7/dist-packages (from b
      Requirement already satisfied: python-dateutil>=2.7.3 in /usr/local/lib/python3.7/dist-packag
      Requirement already satisfied: pytz>=2017.3 in /usr/local/lib/python3.7/dist-packages (from p
      Requirement already satisfied: ptyprocess>=0.5 in /usr/local/lib/python3.7/dist-packages (fro
      Building wheels for collected packages: func-timeout, wget
        Building wheel for func-timeout (setup.py) ... done
        Created wheel for func-timeout: filename=func_timeout-4.3.5-py3-none-any.whl size=15098 sha
        Stored in directory: /root/.cache/pip/wheels/68/b5/a5/67c4364c354e141f5a1bd3ec568126f77877a
        Building wheel for wget (setup.py) ... done
        Created wheel for wget: filename=wget-3.2-py3-none-any.whl size=9675 sha256=d4b52b1482461c4
        Stored in directory: /root/.cache/pip/wheels/a1/b6/7c/0e63e34eb06634181c63adacca38b79ff8f35
      Successfully built func-timeout wget
      Installing collected packages: websocket-client, pyyaml, torch, tokenizers, PyPDF2, pdfkit, h
        Attempting uninstall: pyyaml
          Found existing installation: PyYAML 3.13
          Uninstalling PyYAML-3.13:
            Successfully uninstalled PyYAML-3.13
        Attempting uninstall: torch
          Found existing installation: torch 1.12.0+cu113
          Uninstalling torch-1.12.0+cu113:
            Successfully uninstalled torch-1.12.0+cu113
        Attempting uninstall: torchtext
          Found existing installation: torchtext 0.13.0
          Uninstalling torchtext-0.13.0:
            Successfully uninstalled torchtext-0.13.0
      ERROR: pip's dependency resolver does not currently take into account all the packages that a
      torchvision 0.13.0+cu113 requires torch==1.12.0, but you have torch 1.10.2 which is incompati
```

```
1 # Please do not change this cell because some hidden tests might depend on it.
2 import os
3
4 # Otter grader does not handle ! commands well, so we define and use our
5 # own function to execute shell commands.
6 def shell(commands, warn=True):
7     """Executes the string `commands` as a sequence of shell commands.
8
9         Prints the result to stdout and returns the exit status.
10        Provides a printed warning on non-zero exit status unless `warn`
11        flag is unset.
12     """
13     file = os.popen(commands)
14     print (file.read().rstrip('\n'))
15     exit_status = file.close()
16     if warn and exit_status != None:
17         print(f"Completed with errors. Exit status: {exit_status}\n")
18     return exit_status
19
20 shell("""
21 ls requirements.txt >/dev/null 2>&1
22 if [ ! $? = 0 ]; then
23  rm -rf .tmp
24  git clone https://github.com/cs236299-2022-spring/project4.git .tmp
25  mv .tmp/requirements.txt ./
26  rm -rf .tmp
27 fi
28 pip install -q -r requirements.txt
29 """)
```

```
1 # Initialize Otter
2 import otter
3 grader = otter.Notebook()
```

Unsupported Cell Type. Double-Click to inspect/edit the content.

# ▾ 236299 - Introduction to Natural Language Processing

## Project 4: Semantic Interpretation – Question Answering

The goal of semantic parsing is to convert natural language utterances to a meaning representation such as a *logical form* expression or a *SQL query*. In the previous project segment, you built a parsing system to reconstruct parse trees from the natural-language queries in the ATIS dataset. However, that only solves an intermediary task, not the end-user task of obtaining answers to the queries.

In this final project segment, you will go further, building a semantic parsing system to convert English queries to SQL queries, so that by consulting a database you will be able to answer those questions. You

will implement both a rule-based approach and an end-to-end sequence-to-sequence (seq2seq) approach. Both algorithms come with their pros and cons, and by the end of this segment you should have a basic understanding of the characteristics of the two approaches.

## Goals

1. Build a semantic parsing algorithm to convert text to SQL queries based on the syntactic parse trees from the last project.
2. Build an attention-based end-to-end seq2seq system to convert text to SQL.
3. Improve the attention-based end-to-end seq2seq system with self-attention to convert text to SQL.
4. Discuss the pros and cons of the rule-based system and the end-to-end system.
5. (Optional) Use the state-of-the-art pretrained transformers for text-to-SQL conversion.

This will be an extremely challenging project, so we recommend that you start early.

## ▾ Setup

```
 1 import copy
 2 import datetime
 3 import math
 4 import re
 5 import sys
 6 import warnings
 7
 8 import wget
 9 import nltk
10 import sqlite3
11 import torch
12 import torch.nn as nn
13 import torchtext.legacy as tt
14
15 from cryptography.fernet import Fernet
16 from func_timeout import func_set_timeout
17 from torch.nn.utils.rnn import pack_padded_sequence as pack
18 from torch.nn.utils.rnn import pad_packed_sequence as unpack
19 from tqdm import tqdm
20 from transformers import BartTokenizer, BartForConditionalGeneration
```

```
 1 # Set random seeds
 2 seed = 1234
 3 torch.manual_seed(seed)
 4 # Set timeout for executing SQL
 5 TIMEOUT = 3 # seconds
 6
 7 # GPU check: Set runtime type to use GPU where available
 8 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
 9 print (device)
```

    cuda

```
 1 ## Download needed scripts and data
 2 os.makedirs('data', exist_ok=True)
```

```
3 os.makedirs('scripts', exist_ok=True)
4 source_url = "https://raw.githubusercontent.com/nlp-236299/data/master"
5
6 # Grammar to augment for this segment
7 if not os.path.isfile('data/grammar'):
8   wget.download(f"{source_url}/ATIS/grammar_distrib4.crypt", out="data/")
9
10  # Decrypt the grammar file
11  key = b'bfksTY2BJ5VKKK9xZb1PDDLaGkdu7KCDFYfVePSEfGY='
12  fernet = Fernet(key)
13  with open('./data/grammar_distrib4.crypt', 'rb') as f:
14    restored = Fernet(key).decrypt(f.read())
15  with open('./data/grammar', 'wb') as f:
16    f.write(restored)
17
18 # Download scripts and ATIS database
19 wget.download(f"{source_url}/scripts/trees/transform.py", out="scripts/")
20 wget.download(f"{source_url}/ATIS/atis_sqlite.db", out="data/")
```

```
'data//atis_sqlite.db'
```

```
1 # Import downloaded scripts for parsing augmented grammars
2 sys.path.insert(1, './scripts')
3 import transform as xform
```

## ▾ Semantically augmented grammars

In the first part of this project segment, you'll be implementing a rule-based system for semantic interpretation of sentences. Before jumping into using such a system on the ATIS dataset – we'll get to that soon enough – let's first work with some trivial examples to get things going.

The fundamental idea of rule-based semantic interpretation is the rule of compositionality, that *the meaning of a constituent is a function of the meanings of its immediate subconstituents and the syntactic rule that combined them*. This leads to an infrastructure for specifying semantic interpretation in which each syntactic rule in a grammar (in our case, a context-free grammar) is associated with a semantic rule that applies to the meanings associated with the elements on the right-hand side of the rule.

### Example: arithmetic expressions

As a first example, let's consider an augmented grammar for arithmetic expressions, familiar from lab 3-1. We again use the function `xform.parse_augmented_grammar` to parse the augmented grammar. You can read more about it in the file `scripts/transform.py`.

```
1 arithmetic_grammar, arithmetic_augmentations = xform.parse_augmented_grammar(
2     """
3     ## Sample grammar for arithmetic expressions
4
5     S -> NUM                        : lambda Num: Num
6       | S OP S                      : lambda S1, Op, S2: Op(S1, S2)
7
8     OP -> ADD                       : lambda Op: Op
9        | SUB
10       | MULT
11       | DIV
```

```
12
13     NUM -> 'zero'                           : lambda: 0
14          | 'one'                            : lambda: 1
15          | 'two'                            : lambda: 2
16          | 'three'                          : lambda: 3
17          | 'four'                           : lambda: 4
18          | 'five'                           : lambda: 5
19          | 'six'                            : lambda: 6
20          | 'seven'                          : lambda: 7
21          | 'eight'                          : lambda: 8
22          | 'nine'                           : lambda: 9
23          | 'ten'                            : lambda: 10
24
25     ADD -> 'plus' | 'added' 'to'            : lambda: lambda x, y: x + y
26     SUB -> 'minus'                          : lambda: lambda x, y: x - y
27     MULT -> 'times' | 'multiplied' 'by'     : lambda: lambda x, y: x * y
28     DIV -> 'divided' 'by'                   : lambda: lambda x, y: x / y
29     """
30 )
```

Recall that in this grammar specification format, rules that are not explicitly provided with an augmentation (like all the `OP` rules after the first `OP -> ADD`) are associated with the textually most recent one (`lambda Op: Op`).

The `parse_augmented_grammar` function returns both an NLTK grammar and a dictionary that maps from productions in the grammar to their associated augmentations. Let's examine the returned grammar.

```
1 for production in arithmetic_grammar.productions():
2   print(f"{repr(production):25}    {arithmetic_augmentations[production]}")
```

```
S -> NUM                   <function <lambda> at 0x7f18eef5e050>
S -> S OP S                <function <lambda> at 0x7f18eef5e0e0>
OP -> ADD                  <function <lambda> at 0x7f18eef5e200>
OP -> SUB                  <function <lambda> at 0x7f18eef5e320>
OP -> MULT                 <function <lambda> at 0x7f18eef5e440>
OP -> DIV                  <function <lambda> at 0x7f18eef5e560>
NUM -> 'zero'              <function <lambda> at 0x7f18eef5e680>
NUM -> 'one'               <function <lambda> at 0x7f18eef5e7a0>
NUM -> 'two'               <function <lambda> at 0x7f18eef5e8c0>
NUM -> 'three'             <function <lambda> at 0x7f18eef5e9e0>
NUM -> 'four'              <function <lambda> at 0x7f18eef5eb00>
NUM -> 'five'              <function <lambda> at 0x7f18eef5ec20>
NUM -> 'six'               <function <lambda> at 0x7f18eef5ed40>
NUM -> 'seven'             <function <lambda> at 0x7f18eef5ee60>
NUM -> 'eight'             <function <lambda> at 0x7f18eef5ef80>
NUM -> 'nine'              <function <lambda> at 0x7f18eef5f0e0>
NUM -> 'ten'               <function <lambda> at 0x7f18eef5f200>
ADD -> 'plus'              <function <lambda> at 0x7f18eef5f3b0>
ADD -> 'added' 'to'        <function <lambda> at 0x7f18eef5f560>
SUB -> 'minus'             <function <lambda> at 0x7f18eef5f710>
MULT -> 'times'            <function <lambda> at 0x7f18eef5f8c0>
MULT -> 'multiplied' 'by'  <function <lambda> at 0x7f18eef5fa70>
DIV -> 'divided' 'by'      <function <lambda> at 0x7f18eef5fc20>
```
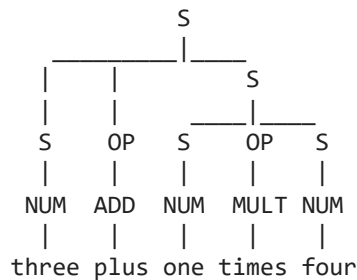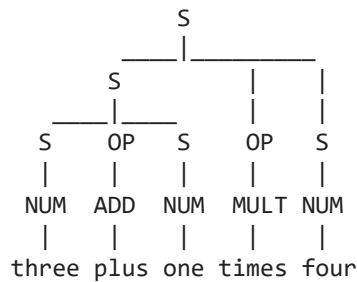
We can parse with the grammar using one of the built-in NLTK parsers.

```
1 arithmetic_parser = nltk.parse.BottomUpChartParser(arithmetic_grammar)
2 parses = [p for p in arithmetic_parser.parse('three plus one times four'.split())]
```

```
3 for parse in parses:
4   parse.pretty_print()
```

```
                  S
           _____|_____
          S                |      |
       ___|___             |      |
      S      OP     S      OP     S
      |      |      |      |      |
    NUM    ADD    NUM   MULT   NUM
      |      |      |      |      |
   three  plus   one   times  four
```

```
                  S
        _____|____
       |     |         S
       |     |      ___|____
       S    OP     S     OP     S
       |    |      |     |      |
     NUM   ADD   NUM   MULT   NUM
       |    |      |     |      |
    three plus   one  times  four
```

Now let's turn to the augmentations. They can be arbitrary Python functions applied to the semantic representations associated with the right-hand-side nonterminals, returning the semantic representation of the left-hand side. To interpret the semantic representation of the entire sentence (at the root of the parse tree), we can use the following pseudo-code:

```
to interpret a tree:
  interpret each of the nonterminal-rooted subtrees
  find the augmentation associated with the root production of the tree
    (it should be a function of as many arguments as there are nonterminals on the right-hand side)
  return the result of applying the augmentation to the subtree values
```

(The base case of this recursion occurs when the number of nonterminal-rooted subtrees is zero, that is, a rule all of whose right-hand side elements are terminals.)

Suppose we had such a function, call it `interpret`. How would it operate on, for instance, the tree `(S (S (NUM three)) (OP (ADD plus)) (S (NUM one)))`?

```
interpret (S (S (NUM three)) (OP (ADD plus)) (S (NUM one)))
    |->interpret (S (NUM three))
    |        |->interpret (NUM three)
    |        |        |->(no subconstituents to evaluate)
    |        |        |->apply the augmentation for the rule NUM -> three to the empty set of values
    |        |        |        (lambda: 3) () ==> 3
    |        |        \==> 3
    |        |->apply the augmentation for the rule S -> NUM to the value 3
    |        |        (lambda NUM: NUM)(3) ==> 3
    |        \==> 3
    |->interpret (OP (ADD plus))
    |        |...
    |        \==> lambda x, y: x + y
    |->interpret (S (NUM one))
    |        |...
```

```
|        \==> 1
|->apply the augmentation for the rule S -> S OP S to the values 3, (lambda x, y: x + y), and 1
|        (lambda S1, Op, S2: Op(S1, S2))(3, (lambda x, y: x + y), 1) ==> 4
\==> 4
```

Thus, the string "three plus one" is semantically interpreted as the value 4.

We provide the `interpret` function to carry out this recursive process, copied over from lab 4-2:

```
1 def interpret(tree, augmentations):
2   syntactic_rule = tree.productions()[0]
3   semantic_rule = augmentations[syntactic_rule]
4   child_meanings = [interpret(child, augmentations)
5                     for child in tree
6                     if isinstance(child, nltk.Tree)]
7   return semantic_rule(*child_meanings)
```

Now we should be able to **evaluate** the arithmetic example from above.

```
1 interpret(parses[0], arithmetic_augmentations)
```

```
16
```

And we can even write a function that parses and interprets a string. We'll have it **evaluate** each of the possible parses and print the results.

```
1 def parse_and_interpret(string, grammar, augmentations):
2   parser = nltk.parse.BottomUpChartParser(grammar)
3   parses = parser.parse(string.split())
4   for parse in parses:
5     parse.pretty_print()
6     print(parse, "==>", interpret(parse, augmentations))
```

```
1 parse_and_interpret("three plus one times four", arithmetic_grammar, arithmetic_augmentations)
```

```
                 S
          ____|_____
         S            |     |
     ___|___          |     |
    S   OP   S        OP    S
    |   |    |        |     |
   NUM  ADD NUM     MULT  NUM
    |   |    |        |     |
  three plus one   times four

(S
  (S (S (NUM three)) (OP (ADD plus)) (S (NUM one)))
  (OP (MULT times))
  (S (NUM four))) ==> 16
                 S
          _____|____
        |    |         S
        |    |      ___|___
        S   OP    S   OP    S
        |    |    |    |     |
       NUM  ADD  NUM MULT  NUM
```

```
       |    |    |    |    |
     three plus one times four

    (S
      (S (NUM three))
      (OP (ADD plus))
      (S (S (NUM one)) (OP (MULT times)) (S (NUM four)))) ==> 7
```

Since the string is syntactically ambiguous according to the grammar, it is semantically ambiguous as well.

## Some grammar specification conveniences

Before going on, it will be useful to have a few more conveniences in writing augmentations for rules. First, since the augmentations are arbitrary Python expressions, they can be built from and make use of other functions. For instance, you'll notice that many of the augmentations at the leaves of the tree took no arguments and returned a constant. We can define a function `constant` that returns a function that ignores its arguments and returns a particular value.

```
1 def constant(value):
2   """Return `value`, ignoring any arguments"""
3   return lambda *args: value
```

Similarly, several of the augmentations are functions that just return their first argument. Again, we can define a generic form `first` of such a function:

```
1 def first(*args):
2   """Return the value of the first (and perhaps only) subconstituent,
3     ignoring any others"""
4   return args[0]
```

We can now rewrite the grammar above to take advantage of these shortcuts.

> In the call to `parse_augmented_grammar` below, we pass in the global environment, extracted via a `globals()` function call, via the named argument `globals`. This allows the `parse_augmented_grammar` function to make use of the global bindings for `constant`, `first`, and the like when evaluating the augmentation expressions to their values. You can check out the code in `transform.py` to see how the passed in `globals` bindings are used. To help understand what's going on, see what happens if you don't include the `globals=globals()`.

```
1 arithmetic_grammar_2, arithmetic_augmentations_2 = xform.parse_augmented_grammar(
2     """
3     ## Sample grammar for arithmetic expressions
4
5     S -> NUM                              : first
6        | S OP S                           : lambda S1, Op, S2: Op(S1, S2)
7
8     OP -> ADD                             : first
9        | SUB
10       | MULT
11       | DIV
12
```

```
13    NUM -> 'zero'                            : constant(0)
14        | 'one'                              : constant(1)
15        | 'two'                              : constant(2)
16        | 'three'                            : constant(3)
17        | 'four'                             : constant(4)
18        | 'five'                             : constant(5)
19        | 'six'                              : constant(6)
20        | 'seven'                            : constant(7)
21        | 'eight'                            : constant(8)
22        | 'nine'                             : constant(9)
23        | 'ten'                              : constant(10)
24
25    ADD -> 'plus' | 'added' 'to'        : constant(lambda x, y: x + y)
26    SUB -> 'minus'                       : constant(lambda x, y: x - y)
27    MULT -> 'times' | 'multiplied' 'by'  : constant(lambda x, y: x * y)
28    DIV -> 'divided' 'by'                : constant(lambda x, y: x / y)
29    """,
30    globals=globals())
```

Finally, it might make our lives easier to write a template of augmentations whose instantiation depends on the right-hand side of the rule.

We use a reserved keyword `_RHS` to denote the right-hand side of the syntactic rule, which will be replaced by a **list** of the right-hand-side strings. For example, an augmentation `numeric_template(_RHS)` would be as if written as `numeric_template(['zero'])` when the rule is `NUM -> 'zero'`, and `numeric_template(['one'])` when the rule is `NUM -> 'one'`. The details of how this works can be found at `scripts/transform.py`.

This would allow us to use a single template function, for example,

```
1 def numeric_template(rhs):
2   """Ignore the subphrase meanings and lookup the first right-hand-side symbol
3     as a number"""
4   return constant({'zero':0, 'one':1, 'two':2, 'three':3, 'four':4, 'five':5,
5           'six':6, 'seven':7, 'eight':8, 'nine':9, 'ten':10}[rhs[0]])
```

and then further simplify the grammar specification:

```
1 arithmetic_grammar_3, arithmetic_augmentations_3 = xform.parse_augmented_grammar(
2     """
3     ## Sample grammar for arithmetic expressions
4
5     S -> NUM                             : first
6        | S OP S                          : lambda S1, Op, S2: Op(S1, S2)
7
8     OP -> ADD                            : first
9        | SUB
10       | MULT
11       | DIV
12
13    NUM -> 'zero'  | 'one'   | 'two'     : numeric_template(_RHS)
14        | 'three' | 'four'  | 'five'
15        | 'six'   | 'seven' | 'eight'
16        | 'nine'  | 'ten'
17
18    ADD -> 'plus' | 'added' 'to'         : constant(lambda x, y: x + y)
```
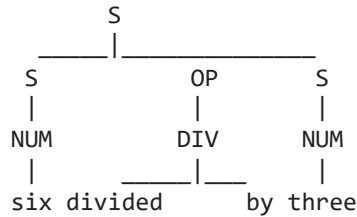
```
19    SUB -> 'minus'                        : constant(lambda x, y: x - y)
20    MULT -> 'times' | 'multiplied' 'by'   : constant(lambda x, y: x * y)
21    DIV -> 'divided' 'by'                 : constant(lambda x, y: x / y)
22    """,
23    globals=globals())
```

```
1 parse_and_interpret("six divided by three", arithmetic_grammar_3, arithmetic_augmentations_3)
```

```
             S
      _____|_____
      S             OP        S
      |             |         |
     NUM           DIV       NUM
      |         ____|___      |
    six divided      by three

    (S (S (NUM six)) (OP (DIV divided by)) (S (NUM three))) ==> 2.0
```

## ▾ Example: *Green Eggs and Ham* revisited

This stuff is tricky, so it's useful to see more examples before jumping in the deep end. In this simple GEaH fragment grammar, we use a larger set of auxiliary functions to build the augmentations.

```
1 def forward(F, A):
2   """Forward application: Return the application of the first
3       argument to the second"""
4   return F(A)
5
6 def backward(A, F):
7   """Backward application: Return the application of the second
8       argument to the first"""
9   return F(A)
10
11 def second(*args):
12   """Return the value of the second subconstituent, ignoring any others"""
13   return args[1]
14
15 def ignore(*args):
16   """Return `None`, ignoring everything about the constituent. (Good as a
17       placeholder until a better augmentation can be devised.)"""
18   return None
```

Using these, we can build and test the grammar.
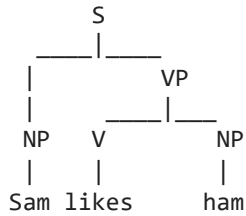
```
1 geah_grammar_spec = """
2   ## Productions
3   S -> NP VP            : backward
4   VP -> V NP            : forward
5
6   ## Lexicon
7   V -> 'likes'         : constant(lambda Object: lambda Subject: f"like({Subject}, {Object})")
8   NP -> 'Sam' | 'sam'  : constant(_RHS[0])
9   NP -> 'ham'
10  NP -> 'eggs'
11 """
```

```
1 geah_grammar, geah_augmentations = xform.parse_augmented_grammar(geah_grammar_spec,
2                                                        globals=globals())
```

```
1 parse_and_interpret("Sam likes ham", geah_grammar, geah_augmentations)
```

```
            S
        ____|____
       |         VP
       |      ___|___
      NP     V       NP
      |      |       |
     Sam   likes    ham

   (S (NP Sam) (VP (V likes) (NP ham))) ==> like(Sam, ham)
```

## Semantics of ATIS queries

Now you're in a good position to understand and add augmentations to a more comprehensive grammar, say, one that parses ATIS queries and generates SQL queries.

In preparation for that, we need to load the ATIS data, both NL and SQL queries.

## Loading and preprocessing the corpus

To simplify things a bit, we'll only consider ATIS queries whose question type (remember that from project segment 1?) is `flight_id`. We download training, development, and test splits for this subset of the ATIS corpus, including corresponding SQL queries.

```
1 # Acquire the datasets - training, development, and test splits of the
2 # ATIS queries and corresponding SQL queries
3 wget.download(f"{source_url}/ATIS/test_flightid.nl", out="data/")
4 wget.download(f"{source_url}/ATIS/test_flightid.sql", out="data/")
5 wget.download(f"{source_url}/ATIS/dev_flightid.nl", out="data/")
6 wget.download(f"{source_url}/ATIS/dev_flightid.sql", out="data/")
7 wget.download(f"{source_url}/ATIS/train_flightid.nl", out="data/")
8 wget.download(f"{source_url}/ATIS/train_flightid.sql", out="data/")
```

```
   'data//train_flightid.sql'
```

Let's take a look at the data: the NL queries are in `.nl` files, and the SQL queries are in `.sql` files.

```
1 shell("head -1 data/dev_flightid.nl")
2 shell("head -1 data/dev_flightid.sql")
```

```
   what flights are available tomorrow from denver to philadelphia
   SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 , c
```

## Corpus preprocessing

We'll use `torchtext` to process the data. We use two `Field`s: `SRC` for the questions, and `TGT` for the SQL

```
1 ## Tokenizer
2 tokenizer = nltk.tokenize.RegexpTokenizer('\d+|st\.|[\w-]+|\$[\d\.]+|\S+')
3 def tokenize(string):
4   return tokenizer.tokenize(string.lower())
5
6 ## Demonstrating the tokenizer
7 ## Note especially the handling of `"11pm"` and hyphenated words.
8 print(tokenize("Are there any first-class flights from St. Louis at 11pm for less than $3.50?"))
```

```
['are', 'there', 'any', 'first-class', 'flights', 'from', 'st.', 'louis', 'at', '11', 'pm', 'fo
```

```
 1 SRC = tt.data.Field(include_lengths=True,          # include lengths
 2                     batch_first=False,             # batches will be max_len x batch_size
 3                     tokenize=tokenize,             # use our tokenizer
 4                     )
 5 TGT = tt.data.Field(include_lengths=False,
 6                     batch_first=False,             # batches will be max_len x batch_size
 7                     tokenize=lambda x: x.split(),  # use split to tokenize
 8                     init_token="<bos>",            # prepend <bos>
 9                     eos_token="<eos>")             # append <eos>
10 fields = [('src', SRC), ('tgt', TGT)]
```

> Note that we specified `batch_first=False` (as in lab 4-4), so that the returned batched
> tensors would be of size `max_length x batch_size`, which facilitates seq2seq
> implementation.

Now, we load the data using `torchtext`. We use the `TranslationDataset` class here because our task is essentially a translation task: "translating" questions into the corresponding SQL queries. Therefore, we also refer to the questions as the *source* side (`SRC`) and the SQL queries as the *target* side (`TGT`).

```
 1 # Make splits for data
 2 train_data, val_data, test_data = tt.datasets.TranslationDataset.splits(
 3     ('_flightid.nl', '_flightid.sql'), fields, path='./data/',
 4     train='train', validation='dev', test='test')
 5
 6 MIN_FREQ = 3
 7 SRC.build_vocab(train_data.src, min_freq=MIN_FREQ)
 8 TGT.build_vocab(train_data.tgt, min_freq=MIN_FREQ)
 9
10 print (f"Size of English vocab: {len(SRC.vocab)}")
11 print (f"Most common English words: {SRC.vocab.freqs.most_common(10)}\n")
12
13 print (f"Size of SQL vocab: {len(TGT.vocab)}")
14 print (f"Most common SQL words: {TGT.vocab.freqs.most_common(10)}\n")
15
16 print (f"Index for start of sequence token: {TGT.vocab.stoi[TGT.init_token]}")
17 print (f"Index for end of sequence token: {TGT.vocab.stoi[TGT.eos_token]}")
```

```
    Size of English vocab: 421
    Most common English words: [('to', 3478), ('from', 3019), ('flights', 2094), ('the', 1550), ('o

    Size of SQL vocab: 392
    Most common SQL words: [('=', 38876), ('AND', 36564), (',', 22772), ('airport_service', 8314),
```

```
Index for start of sequence token: 2
Index for end of sequence token: 3
```

Next, we batch our data to facilitate processing on a GPU. Batching is a bit tricky because the source and target will typically be of different lengths. Fortunately, `torchtext` allows us to pass in a `sort_key` function. By sorting on length, we can minimize the amount of padding on the source side, but since there is still some padding, we need to handle them with pack and unpack later on in the seq2seq part (as in lab 4-5).

```
 1 BATCH_SIZE = 16 # batch size for training/validation
 2 TEST_BATCH_SIZE = 1 # batch size for test, we use 1 to make beam search implementation easier
 3
 4 train_iter, val_iter = tt.data.BucketIterator.splits((train_data, val_data),
 5                                                       batch_size=BATCH_SIZE,
 6                                                       device=device,
 7                                                       repeat=False,
 8                                                       sort_key=lambda x: len(x.src),
 9                                                       sort_within_batch=True)
10 test_iter = tt.data.BucketIterator(test_data,
11                                    batch_size=TEST_BATCH_SIZE,
12                                    device=device,
13                                    repeat=False,
14                                    sort=False,
15                                    train=False)
```

Let's look at a single batch from one of these iterators.

```
 1 batch = next(iter(train_iter))
 2 train_batch_text, train_batch_text_lengths = batch.src
 3 print (f"Size of text batch: {train_batch_text.shape}")
 4 print (f"Third sentence in batch: {train_batch_text[:, 2]}")
 5 print (f"Length of the third sentence in batch: {train_batch_text_lengths[2]}")
 6 print (f"Converted back to string: {' '.join([SRC.vocab.itos[i] for i in train_batch_text[:, 2]])}")
 7
 8 train_batch_sql = batch.tgt
 9 print (f"Size of sql batch: {train_batch_sql.shape}")
10 print (f"Third SQL in batch: {train_batch_sql[:, 2]}")
11 print (f"Converted back to string: {' '.join([TGT.vocab.itos[i] for i in train_batch_sql[:, 2]])}"
```

```
    Size of text batch: torch.Size([12, 16])
    Third sentence in batch: tensor([  9,   7,   5,   4,   3,  11,   2,  20,  33, 267,  18, 415],
            device='cuda:0')
    Length of the third sentence in batch: 12
    Converted back to string: show me the flights from boston to pittsburgh leaving wednesdays and
    Size of sql batch: torch.Size([163, 16])
    Third SQL in batch: tensor([  2,  14,  31,  11,  13,  12,  16,   6,   7,  22,   6,   8,  23,
              7,  29,   6,   8,  30,   6,  33,  40,   6,  33, 101,  15,  21,   4,
             18,   5,  19,   4,  17,   5,  20,   4,  52,   5,   9,  24,   4,  25,
              5,  26,   4,  27,   5,  28,   4,  59,   5,   9,  34,   4,  36,   5,
             37,   4, 248,   5,  34,   4,  99,   5,  98,   4, 221,  10,  10,   3,
              1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,
              1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,
              1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,
              1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,
              1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,
              1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,   1,
              1,   1,   1,   1,   1,   1,   1,   1,   1], device='cuda:0')
    Converted back to string: <bos> SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airpo
```

Alternatively, we can directly iterate over the raw examples:

```
1 for example in train_iter.dataset[:1]:
2   train_text_1 = ' '.join(example.src) # detokenized question
3   train_sql_1 = ' '.join(example.tgt)  # detokenized sql
4   print (f"Question: {train_text_1}\n")
5   print (f"SQL: {train_sql_1}")
```

Question: list all the flights that arrive at general mitchell international from various citie

SQL: SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport airport_1 , airport_serv

## ▾ Establishing a SQL database for evaluating ATIS queries

The output of our systems will be SQL queries. How should we determine if the generated queries are correct? We can't merely compare against the gold SQL queries, since there are many ways to implement a SQL query that answers any given NL query.

Instead, we will execute the queries – both the predicted SQL query and the gold SQL query – on an actual database, and verify that the returned responses are the same. For that purpose, we need a SQL database server to use. We'll set one up here, using the Python `sqlite3` module.

```
1 @func_set_timeout(TIMEOUT)
2 def execute_sql(sql):
3   conn = sqlite3.connect('data/atis_sqlite.db')  # establish the DB based on the downloaded data
4   c = conn.cursor()                              # build a "cursor"
5   c.execute(sql)
6   results = list(c.fetchall())
7   c.close()
8   conn.close()
9   return results
```
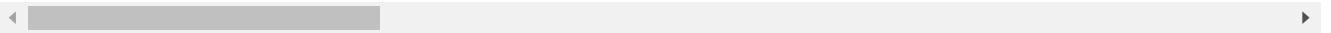
To run a query, we use the cursor's `execute` function, and retrieve the results with `fetchall`. Let's get all the flights that arrive at General Mitchell International – the query `train_sql_1` above. There's a lot, so we'll just print out the first few.

```
1 predicted_ret = execute_sql(train_sql_1)
2
3 print(f"""
4 Executing: {train_sql_1}
5
6 Result: {len(predicted_ret)} entries starting with
7
8 {predicted_ret[:10]}
9 """)
```

Executing: SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport airport_1 , airpor

Result: 534 entries starting with

```
[(107929,), (107930,), (107931,), (107932,), (107933,), (107934,), (107935,), (107936,), (10793
```

For your reference, the SQL database we are using has a database schema described at
https://github.com/jkkummerfeld/text2sql-data/blob/master/data/atis-schema.csv, and is consistent with
the SQL queries provided in the various `.sql` files loaded above.

## ▾ Rule-based parsing and interpretation of ATIS queries

First, you will implement a rule-based semantic parser using a grammar like the one you completed in the
third project segment. We've placed an initial grammar in the file `data/grammar`. In addition to the helper
functions defined above (`constant`, `first`, etc.), it makes use of some other simple functions. We've
included those below, but you can (and almost certainly should) augment this set with others that you
define as you build out the full set of augmentations.

```
1 def upper(term):
2   return '"' + term.upper() + '"'
3
4 def weekday(day):
5   return f"flight.flight_days IN (SELECT days.days_code FROM days WHERE days.day_name = '{day.uppe
6
7 def month_name(month):
8   return {'JANUARY' : 1,
9           'FEBRUARY' : 2,
10          'MARCH' : 3,
11          'APRIL' : 4,
12          'MAY' : 5,
13          'JUNE' : 6,
14          'JULY' : 7,
15          'AUGUST' : 8,
16          'SEPTEMBER' : 9,
17          'OCTOBER' : 10,
18          'NOVEMBER' : 11,
19          'DECEMBER' : 12}[month.upper()]
20
21 def airports_from_airport_name(airport_name):
22   return f"(SELECT airport.airport_code FROM airport WHERE airport.airport_name = {upper(airport_n
23
24 def airports_from_city(city):
25   return f"""
26     (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code IN
27       (SELECT city.city_code FROM city WHERE city.city_name = {upper(city)}))
28   """
29
30 def null_condition(*args, **kwargs):
31   return 1
32
33 def depart_around(time):
34   return f"""
35     flight.departure_time >= {add_delta(miltime(time), -15).strftime('%H%M')}
36     AND flight.departure_time <= {add_delta(miltime(time), 15).strftime('%H%M')}
37   """.strip()
38
39 def add_delta(tme, delta):
```

```
40      # transform to a full datetime first
41      return (datetime.datetime.combine(datetime.date.today(), tme) +
42              datetime.timedelta(minutes=delta)).time()
43
44 def miltime(minutes):
45   return datetime.time(hour=int(minutes/100), minute=(minutes % 100))
46
```

```
1 ## Added augmentation ##
2 def flight_dest(dest):
3     return f"""flight.to_airport IN {dest}"""
4
5 def flight_src(src):
6     return f"""flight.from_airport IN {src}"""
7
8 def arrive_before(time):
9     return f"""flight.arrival_time < {time}"""
10
11 def depart_before(time):
12     return f"""flight.departure_time < {time}"""
13
14 def arrive_after(time):
15     return f"""flight.arrival_time > {time}"""
16
17 def depart_after(time):
18     return f"""flight.departure_time > {time}"""
19
20 def arrive_around(time):
21   return f"""
22     flight.arrival_time >= {add_delta(miltime(time), -15).strftime('%H%M')}
23     AND flight.arrival_time <= {add_delta(miltime(time), 15).strftime('%H%M')}
24     """.strip()
25
26 def date_from_DMY(day_num = None, month = None, year = None):
27   constraints = []
28
29   if day_num is not None:
30     constraints.append(f"""DAY_NUMBER = {day_num}""")
31   if month is not None:
32     constraints.append(f"""MONTH_NUMBER = {month}""")
33   if year is not None:
34     constraints.append(f"""YEAR = {year}""")
35   constraints_s = " AND ".join(constraints)
36   return f"""
37   flight.flight_days IN
38     (SELECT DAYS_CODE FROM days WHERE DAY_NAME IN
39       (SELECT DAY_NAME FROM date_day WHERE {constraints_s}))
40   """
41
42 def union_dates(day_num1 = None, month1 = None, year1 = None, day_num2 = None, month2 = None, year
43
44     return f"""{date_from_DMY(day_num1, month1, year1)} OR  {date_from_DMY(day_num2, month2, year2
45
46 def and_(a, b):
47   return f""" {a} AND {b} """
48
49 def airline_code(code):
50   return f"flight.airline_code = '{code}'"
51
52 def S_NP(NP):
```

```
53    return f"""SELECT DISTINCT flight.flight_id FROM flight WHERE {NP}"""
54
```

We can build a parser with the augmented grammar:

```
1 atis_grammar, atis_augmentations = xform.read_augmented_grammar('data/grammar', globals=globals())
2 atis_parser = nltk.parse.BottomUpChartParser(atis_grammar)
```

We'll define a function to return a parse tree for a string according to the ATIS grammar (if available).

```
 1 def parse_tree(sentence):
 2   """Parse a sentence and return the parse tree, or None if failure."""
 3   try:
 4     parses = list(atis_parser.parse(tokenize(sentence)))
 5     if len(parses) == 0:
 6       return None
 7     else:
 8       return parses[0]
 9   except:
10     return None
```

We can check the overall coverage of this grammar on the training set by using the `parse_tree` function to determine if a parse is available. The grammar that we provide should get about a 44% coverage of the training set.

```
 1 # Check coverage on training set
 2 parsed = 0
 3 with open("data/train_flightid.nl") as train:
 4   examples = train.readlines()[:]
 5 for sentence in tqdm(examples):
 6   if parse_tree(sentence):
 7     parsed += 1
 8   else:
 9     next
10
11 print(f"\nParsed {parsed} of {len(examples)} ({parsed*100/(len(examples)):.2f}%)")
```

```
    100%|██████████| 3651/3651 [00:19<00:00, 185.86it/s]
    Parsed 1609 of 3651 (44.07%)
```

## ▾ Goal 1: Construct SQL queries from a parse tree and **evaluate** the results

It's time to turn to the first major part of this project segment, implementing a rule-based semantic parsing system to answer flight-ID-type ATIS queries.

Recall that in rule-based semantic parsing, each syntactic rule is associated with a semantic composition rule. The grammar we've provided has semantic augmentations for some of the low-level phrases – cities, airports, times, airlines – but not the higher level syntactic types. You'll be adding those.

In the ATIS grammar that we provide, as with the earlier toy grammars, the augmentation for a rule with $n$ nonterminals and $m$ terminals on the right-hand side is assumed to be called with $n$ positional arguments (the values for the corresponding children). The `interpret` function you've already defined should therefore work well with this grammar.

Let's run through one way that a semantic derivation might proceed, for the sample query "flights to boston":

```
1 sample_query = "flights to boston"
2 print(tokenize(sample_query))
3 sample_tree = parse_tree(sample_query)
4 sample_tree.pretty_print()
```

```
['flights', 'to', 'boston']
                  S
                  |
              NP_FLIGHT
                  |
              NOM_FLIGHT
                  |
               N_FLIGHT
         _____|_____
        |                   PP
        |                   |
        |                PP_PLACE
        |             _____|_____
    N_FLIGHT         |             N_PLACE
        |            |               |
   TERM_FLIGHT     P_PLACE       TERM_PLACE
        |            |               |
     flights         to            boston
```

Given a sentence, we first construct its parse tree using the syntactic rules, then compose the corresponding semantic rules bottom-up, until eventually we arrive at the root node with a finished SQL statement. For this query, we will go through what the possible meaning representations for the subconstituents of "flights to boston" might be. But this is just one way of doing things; other ways are possible, and you should feel free to experiment.

Working from bottom up:

1. The `TERM_PLACE` phrase "boston" uses the composition function template
   `constant(airports_from_city(' '.join(_RHS)))`, which will be instantiated as
   `constant(airports_from_city(' '.join(['boston'])))` (recall that `_RHS` is replaced by the right-hand side of the rule). The meaning of `TERM_PLACE` will be the SQL snippet

   ```
   SELECT airport_service.airport_code
   FROM airport_service
   WHERE airport_service.city_code IN
    (SELECT city.city_code
      FROM city
      WHERE city.city_name = "BOSTON")
   ```

   (This query generates a list of all of the airports in Boston.)

2. The `N_PLACE` phrase "boston" can have the same meaning as the `TERM_PLACE`.

3. The `P_PLACE` phrase "to" might be associated with a function that maps a SQL query for a list of airports to a SQL condition that holds of flights that go to one of those airports, i.e., `flight.to_airport IN (...)`.

4. The `PP_PLACE` phrase "to boston" might apply the `P_PLACE` meaning to the `TERM_PLACE` meaning, thus generating a SQL condition that holds of flights that go to one of the Boston airports:

```
flight.to_airport IN
 (SELECT airport_service.airport_code
  FROM airport_service
  WHERE airport_service.city_code IN
      (SELECT city.city_code
       FROM city
       WHERE city.city_name = "BOSTON"))
```

5. The `PP` phrase "to Boston" can again get its meaning from the `PP_PLACE`.

6. The `TERM_FLIGHT` phrase "flights" might also return a condition on flights, this time the "null condition", represented by the SQL truth value `1`. Ditto for the `N_FLIGHT` phrase "flights".

7. The `N_FLIGHT` phrase "flights to boston" can conjoin the two conditions, yielding the SQL condition

```
flight.to_airport IN
 (SELECT airport_service.airport_code
  FROM airport_service
  WHERE airport_service.city_code IN
      (SELECT city.city_code
       FROM city
       WHERE city.city_name = "BOSTON"))
AND 1
```

which can be inherited by the `NOM_FLIGHT` and `NP_FLIGHT` phrases.

8. The `S` phrase "flights to boston" can use the condition provided by the `NP_FLIGHT` phrase to select all flights satisfying the condition with a SQL query like

```
SELECT DISTINCT flight.flight_id
FROM flight
WHERE flight.to_airport IN
      (SELECT airport_service.airport_code
       FROM airport_service
       WHERE airport_service.city_code IN
           (SELECT city.city_code
            FROM city
            WHERE city.city_name = "BOSTON"))
    AND 1
```

This SQL query is then taken to be a representation of the meaning for the NL query "flights to boston", and can be executed against the ATIS database to retrieve the requested flights.

Now, it's your turn to add augmentations to `data/grammar` to make this example work. The augmentations that we have provided for the grammar make use of a set of auxiliary functions that we defined above. You should feel free to add your own auxiliary functions that you make use of in the grammar.

```
1 #TODO: add augmentations to `data/grammar` to make this example work
2 atis_grammar, atis_augmentations = xform.read_augmented_grammar('data/grammar', globals=globals())
3 atis_parser = nltk.parse.BottomUpChartParser(atis_grammar)
4 predicted_sql = interpret(sample_tree, atis_augmentations)
5 print("Predicted SQL:\n\n", predicted_sql, "\n")
```

```
    Predicted SQL:

     SELECT DISTINCT flight.flight_id FROM flight WHERE  1 AND flight.to_airport IN
         (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
           (SELECT city.city_code FROM city WHERE city.city_name = "BOSTON"))
```

## ▾ Verification on some examples

With a rule-based semantic parsing system, we can generate SQL queries given questions, and then execute those queries on a SQL database to answer the given questions. To **evaluate** the performance of the system, we compare the returned results against the results of executing the ground truth queries.

We provide a function `verify` to compare the results from our generated SQL to the ground truth SQL. It should be useful for testing individual queries.

```
 1 def verify(predicted_sql, gold_sql, silent=True):
 2     """
 3     Compare the correctness of the generated SQL by executing on the
 4     ATIS database and comparing the returned results.
 5     Arguments:
 6         predicted_sql: the predicted SQL query
 7         gold_sql: the reference SQL query to compare against
 8         silent: print outputs or not
 9     Returns: True if the returned results are the same, otherwise False
10     """
11     # Execute predicted SQL
12     try:
13       predicted_result = execute_sql(predicted_sql)
14     except BaseException as e:
15       if not silent:
16         print(f"predicted sql exec failed: {e}")
17       return False
18     if not silent:
19       print("Predicted DB result:\n\n", predicted_result[:10], "\n")
20
21     # Execute gold SQL
22     try:
23       gold_result = execute_sql(gold_sql)
24     except BaseException as e:
```

```
25      if not silent:
26        print(f"gold sql exec failed: {e}")
27      return False
28    if not silent:
29      print("Gold DB result:\n\n", gold_result[:10], "\n")
30
31    # Verify correctness
32    if gold_result == predicted_result:
33      return True
```

Let's try this methodology on a simple example: "flights from phoenix to milwaukee". we provide it along with the gold SQL query.

```
1 def rule_based_trial(sentence, gold_sql):
2   print("Sentence: ", sentence, "\n")
3   tree = parse_tree(sentence)
4   print("Parse:\n\n")
5   tree.pretty_print()
6
7   predicted_sql = interpret(tree, atis_augmentations)
8   print("Predicted SQL:\n\n", predicted_sql, "\n")
9
10  if verify(predicted_sql, gold_sql, silent=False):
11    print ('Correct!')
12  else:
13    print ('Incorrect!')
```

```
1 # Run this cell to reload augmentations after you make changes to `data/grammar`
2 atis_grammar, atis_augmentations = xform.read_augmented_grammar('data/grammar', globals=globals())
3 atis_parser = nltk.parse.BottomUpChartParser(atis_grammar)
```

```
1 #TODO: add augmentations to `data/grammar` to make this example work
2 # Example 1
3 example_1 = 'flights from phoenix to milwaukee'
4 gold_sql_1 = """
5   SELECT DISTINCT flight_1.flight_id
6   FROM flight flight_1 ,
7        airport_service airport_service_1 ,
8        city city_1 ,
9        airport_service airport_service_2 ,
10       city city_2
11  WHERE flight_1.from_airport = airport_service_1.airport_code
12        AND airport_service_1.city_code = city_1.city_code
13        AND city_1.city_name = 'PHOENIX'
14        AND flight_1.to_airport = airport_service_2.airport_code
15        AND airport_service_2.city_code = city_2.city_code
16        AND city_2.city_name = 'MILWAUKEE'
17  """
18
19 rule_based_trial(example_1, gold_sql_1)
```
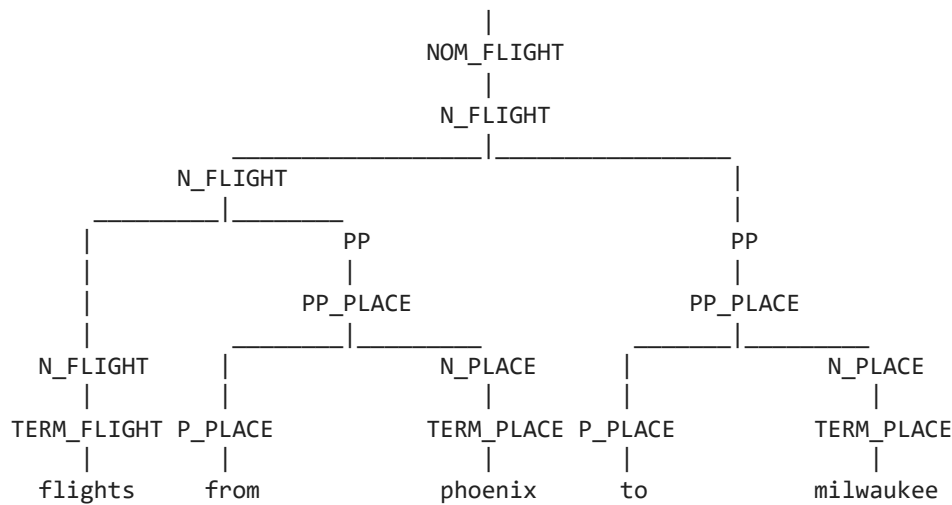
```
    Sentence:  flights from phoenix to milwaukee

    Parse:


                            S
                            |
                         NP_FLIGHT
```

```
                                    |
                                NOM_FLIGHT
                                    |
                                 N_FLIGHT
                     _____|_____
                N_FLIGHT                              |
             _____|_____                          |
            |              PP                          PP
            |               |                          |
            |            PP_PLACE                    PP_PLACE
            |          _____|_____              _____|_____
        N_FLIGHT     |           N_PLACE         |           N_PLACE
            |        |              |            |              |
       TERM_FLIGHT P_PLACE    TERM_PLACE P_PLACE          TERM_PLACE
            |        |              |            |              |
         flights    from        phoenix         to         milwaukee
```

Predicted SQL:

```
 SELECT DISTINCT flight.flight_id FROM flight WHERE   1 AND flight.from_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "PHOENIX"))
    AND flight.to_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "MILWAUKEE"))
```

Predicted DB result:

```
 [(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,), (304881,), (3106
```

Gold DB result:

```
 [(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,), (304881,), (3106
```

Correct!

To make development faster, we recommend starting with a few examples before running the full evaluation script. We've taken some examples from the ATIS dataset including the gold SQL queries that they provided. Of course, yours (and those of the project segment solution set) may differ.

```
 1 #TODO: add augmentations to `data/grammar` to make this example work
 2 # Example 2
 3 example_2 = 'i would like a united flight'
 4 gold_sql_2 = """
 5   SELECT DISTINCT flight_1.flight_id
 6   FROM flight flight_1
 7   WHERE flight_1.airline_code = 'UA'
 8   """
 9
10 rule_based_trial(example_2, gold_sql_2)
```

Sentence:  i would like a united flight

Parse:

```
                                    S
              _____|_____
                                    |
                                    |
```

```
                                    PREIGNORE
          _____|_____                              ____
         |                     PREIGNORE         |                           ADJ
         |            _____|_____  |                            |
         |           |                 PREIGNORE  |                       ADJ_AIRLIN
         |           |            _____|_____                        |
         |           |           |                PREIGNORE    TERM_AIRLI
         |           |           |                   |             |
  PREIGNORESYMBOL PREIGNORESYMBOL    PREIGNORESYMBOL  PREIGNORESYMBOL TERM_AIRBRA
         |           |                   |                 |             |
         i         would                like               a          united
```

Predicted SQL:

 SELECT DISTINCT flight.flight_id FROM flight WHERE  flight.airline_code = 'UA' AND 1

Predicted DB result:

 [(100094,), (100099,), (100145,), (100158,), (100164,), (100167,), (100169,), (100203,), (1002

Gold DB result:

 [(100094,), (100099,), (100145,), (100158,), (100164,), (100167,), (100169,), (100203,), (1002

Correct!
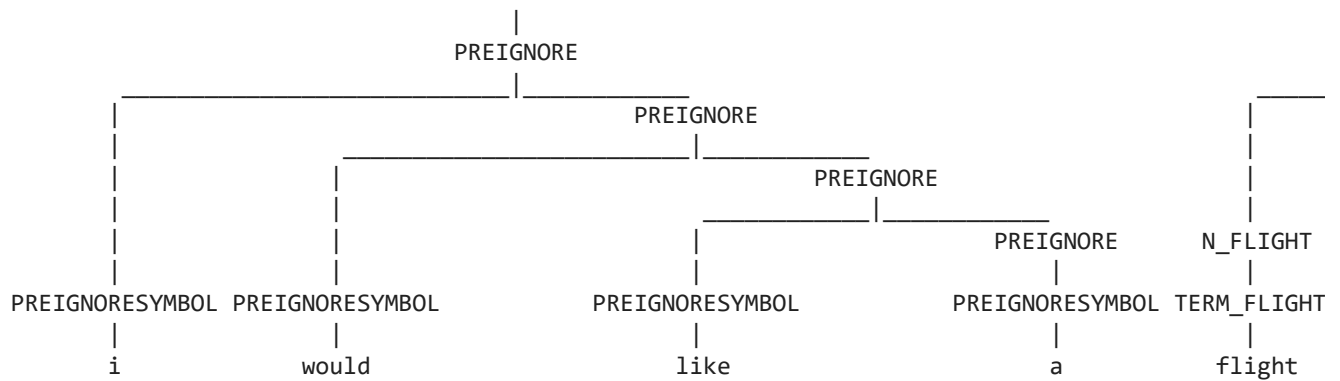
```
 1 #TODO: add augmentations to `data/grammar` to make this example work
 2 # Example 3
 3 example_3 = 'i would like a flight between boston and dallas'
 4 gold_sql_3 = """
 5   SELECT DISTINCT flight_1.flight_id
 6   FROM flight flight_1 ,
 7        airport_service airport_service_1 ,
 8        city city_1 ,
 9        airport_service airport_service_2 ,
10        city city_2
11   WHERE flight_1.from_airport = airport_service_1.airport_code
12        AND airport_service_1.city_code = city_1.city_code
13        AND city_1.city_name = 'BOSTON'
14        AND flight_1.to_airport = airport_service_2.airport_code
15        AND airport_service_2.city_code = city_2.city_code
16        AND city_2.city_name = 'DALLAS'
17   """
18
19 # Note that the parse tree might appear wrong: instead of
20 # `PP_PLACE -> 'between' N_PLACE 'and' N_PLACE`, the tree appears to be
21 # `PP_PLACE -> 'between' 'and' N_PLACE N_PLACE`. But it's only a visualization
22 # error of tree.pretty_print() and you should assume that the production is
23 # `PP_PLACE -> 'between' N_PLACE 'and' N_PLACE` (you can verify by printing out
24 # all productions).
25 rule_based_trial(example_3, gold_sql_3)
```

Sentence:  i would like a flight between boston and dallas

Parse:

```
                                                                       S
                          _____|_____
                         |
                         |
                         |
```

```
                                    |
                                 PREIGNORE
        _____|_____
        |                             PREIGNORE                                    _____
        |              _____|_____                         |
        |              |                      PREIGNORE                              |
        |              |              _____|_____                        |
        |              |              |                     PREIGNORE       N_FLIGHT
        |              |              |              _____|_____        |
        |              |              |              |              PREIGNORE     |
        |              |              |              |                  |         |
 PREIGNORESYMBOL PREIGNORESYMBOL PREIGNORESYMBOL PREIGNORESYMBOL TERM_FLIGHT
        |              |              |              |                  |         |
        i            would          like            a              flight
```

Predicted SQL:

```
 SELECT DISTINCT flight.flight_id FROM flight WHERE  1 AND  flight.to_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "DALLAS"))
   AND flight.from_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "BOSTON"))
```

Predicted DB result:

`[(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,), (103178,), (1031`

Gold DB result:

`[(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,), (103178,), (1031`
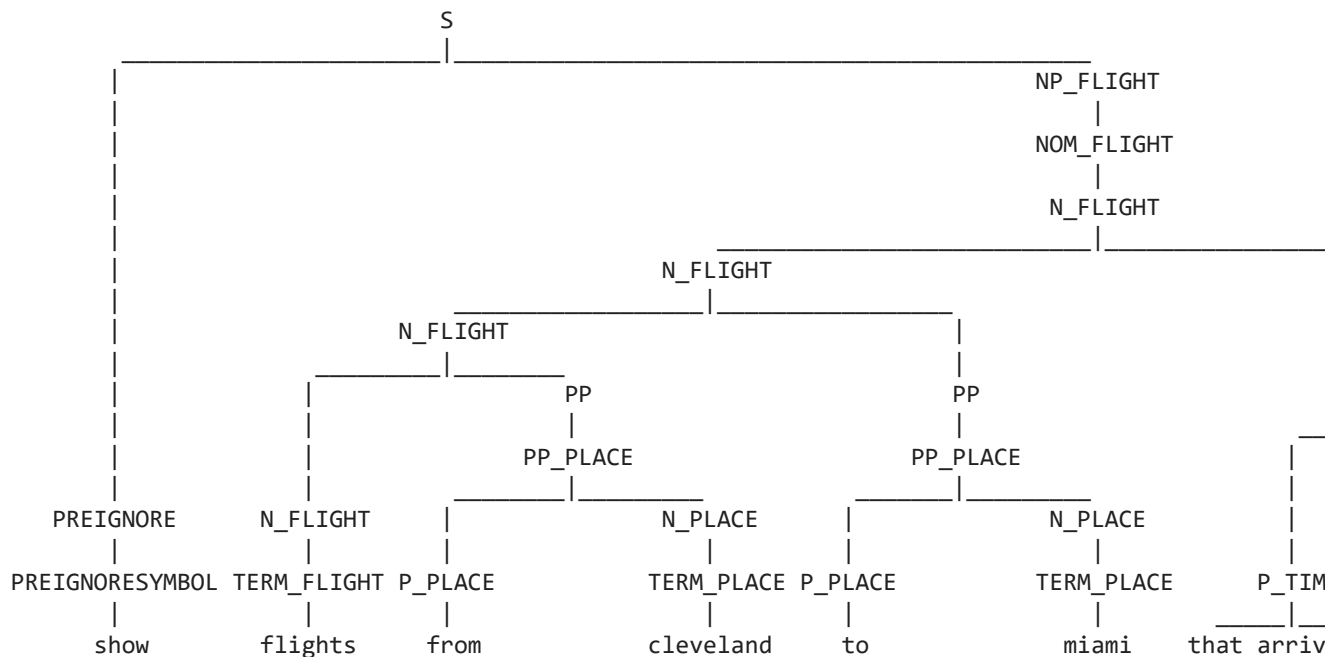
Correct!

```
1 #TODO: add augmentations to `data/grammar` to make this example work
2 # Example 4
3 example_4 = 'show me the united flights from denver to baltimore'
4 gold_sql_4 = """
5   SELECT DISTINCT flight_1.flight_id
6   FROM flight flight_1 ,
7        airport_service airport_service_1 ,
8        city city_1 ,
9        airport_service airport_service_2 ,
10       city city_2
11  WHERE flight_1.airline_code = 'UA'
12       AND ( flight_1.from_airport = airport_service_1.airport_code
13            AND airport_service_1.city_code = city_1.city_code
14            AND city_1.city_name = 'DENVER'
15            AND flight_1.to_airport = airport_service_2.airport_code
16            AND airport_service_2.city_code = city_2.city_code
17            AND city_2.city_name = 'BALTIMORE' )
18
19  """
20
21 rule_based_trial(example_4, gold_sql_4)
```

Sentence:  show me the united flights from denver to baltimore

Parse:

```
                                                        S
                    _____|_____
                    |                                          NP_FLIGHT
```

```
                                                          |
                                                      NOM_FLIGHT
                                           _____|_____
                                          |                             |
                                          |                             |
                                          |                             |
                                          |                             |
                                          |                           _____
                                          |                          |
                                          |                        N_FLIGHT
                                          |               _____|_____
                   PREIGNORE             ADJ             |
          _____|_____     |             |
         |                   |          ADJ_AIRLINE      |
         |                 PREIGNORE       |             |
         |          _____|_____  |             |     _____
         |         |               PREIGNORE  TERM_AIRLINE N_FLIGHT |
         |         |                   |       |        |        |
    PREIGNORESYMBOL PREIGNORESYMBOL   PREIGNORESYMBOL TERM_AIRBRAND TERM_FLIGHT P_PLACE
         |              |                   |       |        |        |
       show            me                  the    united  flights   from
```

Predicted SQL:

```
  SELECT DISTINCT flight.flight_id FROM flight WHERE  flight.airline_code = 'UA' AND   1 AND fli
     (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
       (SELECT city.city_code FROM city WHERE city.city_name = "DENVER"))
     AND flight.to_airport IN
     (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
       (SELECT city.city_code FROM city WHERE city.city_name = "BALTIMORE"))
```

Predicted DB result:

```
  [(101231,), (101233,), (305983,)]
```

Gold DB result:

```
  [(101231,), (101233,), (305983,)]
```

Correct!

```python
 1  #TODO: add augmentations to `data/grammar` to make this example work
 2  # Example 5
 3  example_5 = 'show flights from cleveland to miami that arrive before 4pm'
 4  gold_sql_5 = """
 5    SELECT DISTINCT flight_1.flight_id
 6    FROM flight flight_1 ,
 7         airport_service airport_service_1 ,
 8         city city_1 ,
 9         airport_service airport_service_2 ,
10         city city_2
11    WHERE flight_1.from_airport = airport_service_1.airport_code
12         AND airport_service_1.city_code = city_1.city_code
13         AND city_1.city_name = 'CLEVELAND'
14         AND ( flight_1.to_airport = airport_service_2.airport_code
15              AND airport_service_2.city_code = city_2.city_code
16              AND city_2.city_name = 'MIAMI'
17              AND flight_1.arrival_time < 1600 )
18    """
19
20  rule_based_trial(example_5, gold_sql_5)
```

Sentence:  show flights from cleveland to miami that arrive before 4pm

Parse:

```
                                                       S
         _____|_____
        |                                                                   NP_FLIGHT
        |                                                                       |
        |                                                                   NOM_FLIGHT
        |                                                                       |
        |                                                                    N_FLIGHT
        |                                                          _____|_____
        |                                                     N_FLIGHT
        |                                          _____|_____
        |                        N_FLIGHT         |                             |
        |                    _____|_____       |                             |
        |                   |              PP      |                            PP
        |                   |              |        |                            |
        |                   |           PP_PLACE    |                        PP_PLACE              __
        |                   |         _____|_____    |                      _____|_____            |
     PREIGNORE          N_FLIGHT    |         N_PLACE  |                   |        N_PLACE         |
        |                   |        |          |      |                   |          |            |
  PREIGNORESYMBOL      TERM_FLIGHT P_PLACE   TERM_PLACE P_PLACE         TERM_PLACE            P_TIM
        |                   |        |          |      |                   |          |        ___|__
       show             flights    from     cleveland  to               miami          that arriv
```

Predicted SQL:

```
 SELECT DISTINCT flight.flight_id FROM flight WHERE     1 AND flight.from_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "CLEVELAND"))
    AND flight.to_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "MIAMI"))
    AND flight.arrival_time < 1600
```

Predicted DB result:

```
 [(107698,), (301117,)]
```
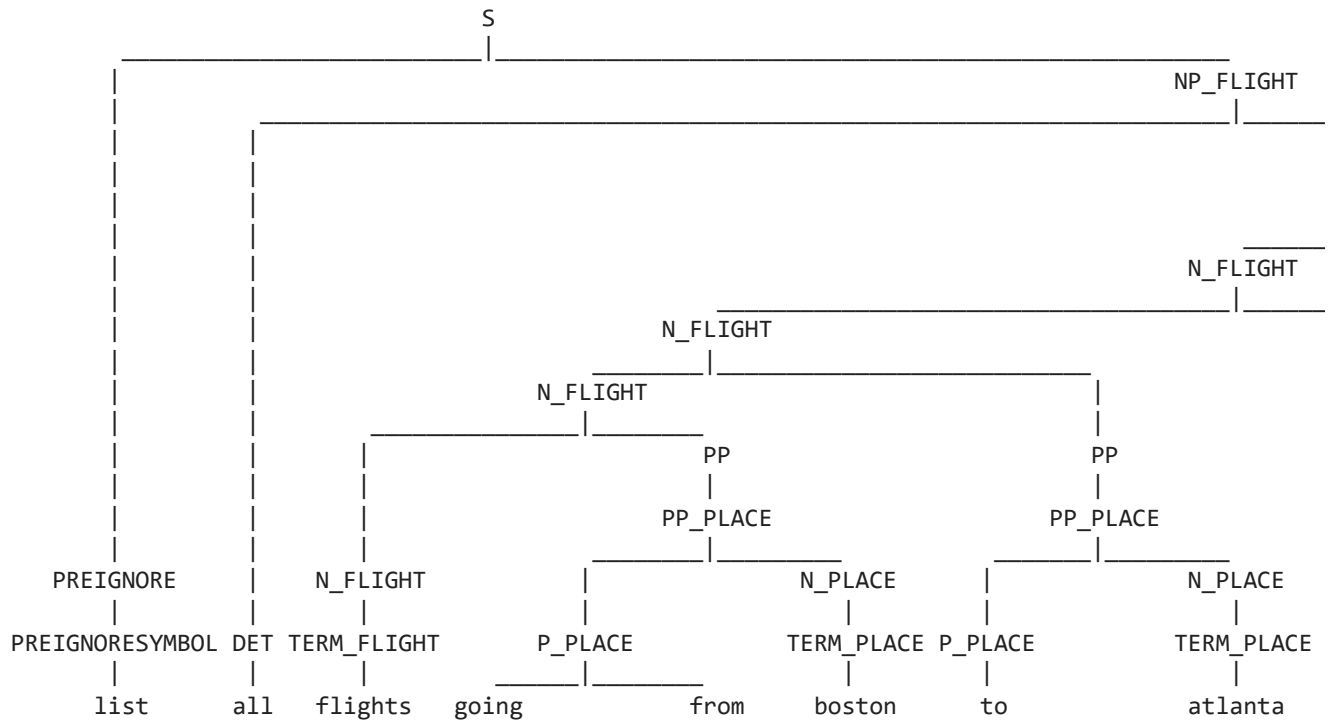
Gold DB result:

```
 [(107698,), (301117,)]
```

Correct!

```
 1 #TODO: add augmentations to `data/grammar` to make this example work
 2 # Example 6
 3 example_6 = 'okay how about a flight on sunday from tampa to charlotte'
 4 gold_sql_6 = """
 5   SELECT DISTINCT flight_1.flight_id
 6   FROM flight flight_1 ,
 7        airport_service airport_service_1 ,
 8        city city_1 ,
 9        airport_service airport_service_2 ,
10        city city_2 ,
11        days days_1 ,
12        date_day date_day_1
13   WHERE flight_1.from_airport = airport_service_1.airport_code
14       AND airport_service_1.city_code = city_1.city_code
15       AND city_1.city_name = 'TAMPA'
16       AND ( flight_1.to_airport = airport_service_2.airport_code
17           AND airport_service_2.city_code = city_2.city_code
18           AND city_2.city_name = 'CHARLOTTE'
19           AND flight_1.flight_days = days_1.days_code
```

```
20                  AND days_1.day_name = date_day_1.day_name
21                  AND date_day_1.year = 1991
22                  AND date_day_1.month_number = 8
23                  AND date_day_1.day_number = 27 )
24   """
25
26 # You might notice that the gold answer above used the exact date, which is
27 # not easily implementable. A more implementable way (generated by the project
28 # segment 4 solution code) is:
29 gold_sql_6b = """
30   SELECT DISTINCT flight.flight_id
31   FROM flight
32   WHERE ((((1
33           AND flight.flight_days IN (SELECT days.days_code
34                                      FROM days
35                                      WHERE days.day_name = 'SUNDAY')
36           )
37          AND flight.from_airport IN (SELECT airport_service.airport_code
38                                      FROM airport_service
39                                      WHERE airport_service.city_code IN (SELECT city.city_code
40                                                                          FROM city
41                                                                          WHERE city.city_name =
42          AND flight.to_airport IN (SELECT airport_service.airport_code
43                                    FROM airport_service
44                                    WHERE airport_service.city_code IN (SELECT city.city_code
45                                                                        FROM city
46                                                                        WHERE city.city_name = "CH
47   """
48
49 rule_based_trial(example_6, gold_sql_6b)
```

Sentence:  okay how about a flight on sunday from tampa to charlotte

Parse:

```
                                  _____
                                 |
                                 |
                                 |
                                 |
                                 |
                                 |
                                 |
                                 |
                             PREIGNORE
               _____|_____
              |                              PREIGNORE                           ____
              |                _____|_____                |
              |               |                             PREIGNORE             |
              |               |                    _____|_____        |
              |               |                   |                  PREIGNORE   N_FLIGHT
              |               |                   |                      |          |
      PREIGNORESYMBOL PREIGNORESYMBOL      PREIGNORESYMBOL      PREIGNORESYMBOL TERM_FLIGHT
              |               |                   |                      |          |
            okay            how                about                     a        flight
```

Predicted SQL:

```
SELECT DISTINCT flight.flight_id FROM flight WHERE    1 AND flight.flight_days IN (SELECT days
   (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
     (SELECT city.city_code FROM city WHERE city.city_name = "TAMPA"))
     AND flight.to_airport IN
```

```
        (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
          (SELECT city.city_code FROM city WHERE city.city_name = "CHARLOTTE"))


    Predicted DB result:

     [(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]

    Gold DB result:

     [(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]

    Correct!
```

```python
1  #TODO: add augmentations to `data/grammar` to make this example work
2  # Example 7
3  example_7 = 'list all flights going from boston to atlanta that leaves before 7 am on thursday'
4  gold_sql_7 = """
5    SELECT DISTINCT flight_1.flight_id
6    FROM flight flight_1 ,
7        airport_service airport_service_1 ,
8        city city_1 ,
9        airport_service airport_service_2 ,
10       city city_2 ,
11       days days_1 ,
12       date_day date_day_1
13   WHERE flight_1.from_airport = airport_service_1.airport_code
14        AND airport_service_1.city_code = city_1.city_code
15        AND city_1.city_name = 'BOSTON'
16        AND ( flight_1.to_airport = airport_service_2.airport_code
17              AND airport_service_2.city_code = city_2.city_code
18              AND city_2.city_name = 'ATLANTA'
19              AND ( flight_1.flight_days = days_1.days_code
20                    AND days_1.day_name = date_day_1.day_name
21                    AND date_day_1.year = 1991
22                    AND date_day_1.month_number = 5
23                    AND date_day_1.day_number = 24
24                    AND flight_1.departure_time < 700 ) )
25   """
26
27 # Again, the gold answer above used the exact date, as opposed to the
28 # following approach:
29 gold_sql_7b = """
30   SELECT DISTINCT flight.flight_id
31   FROM flight
32   WHERE ((1
33           AND ((((1
34                 AND flight.from_airport IN (SELECT airport_service.airport_code
35                                             FROM airport_service
36                                             WHERE airport_service.city_code IN (SELECT city.city
37                                                                                 FROM city
38                                                                                 WHERE city.city_
39               AND flight.to_airport IN (SELECT airport_service.airport_code
40                                         FROM airport_service
41                                         WHERE airport_service.city_code IN (SELECT city.city_cc
42                                                                             FROM city
43                                                                             WHERE city.city_nan
44             AND flight.departure_time <= 0700)
45           AND flight.flight_days IN (SELECT days.days_code
46                                      FROM days
```

```
47                                                     WHERE days.day_name = 'THURSDAY'))))
48  """
49
50 rule_based_trial(example_7, gold_sql_7b)
```

Sentence:  list all flights going from boston to atlanta that leaves before 7 am on thursday

Parse:

```
                                                      S
        _____|_____
       |                                                                                NP_FLIGHT
       |                                  _____|_____
       |                                 |                                                       |
       |                                 |                                                        _____
       |                                 |                                                      N_FLIGHT
       |                                 |                                          _____|_____
       |                                 |                                N_FLIGHT  |
       |                                 |                        _____|_____ |
       |                                 |           N_FLIGHT     |                                 |   |
       |                                 |      _____|_____ |                               PP   |
       |                                 |     |               PP  |                               |    |
       |                                 |     |               |   PP_PLACE                     PP_PLACE
       |                                 |     |            PP_PLACE  _____|_____            _____|_____
       |            PREIGNORE            |  N_FLIGHT        |       N_PLACE     |              N_PLACE
       |               |                 |     |           |          |        |                 |
   PREIGNORESYMBOL  DET TERM_FLIGHT   P_PLACE       TERM_PLACE  P_PLACE   TERM_PLACE
       |               |    |              _____|_____       |        |              |
      list            all flights        going            from    boston       to        atlanta
```

Predicted SQL:

```
 SELECT DISTINCT flight.flight_id FROM flight WHERE      1 AND flight.from_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "BOSTON"))
    AND flight.to_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "ATLANTA"))
    AND flight.departure_time < 700  AND flight.flight_days IN (SELECT days.days_code FROM days
```

Predicted DB result:

  [(100014,)]

Gold DB result:

  [(100014,)]

Correct!

```
1 #TODO: add augmentations to `data/grammar` to make this example work
2 # Example 8
3 example_8 = 'list the flights from dallas to san francisco on american airlines'
4 gold_sql_8 = """
5   SELECT DISTINCT flight_1.flight_id
6   FROM flight flight_1 ,
7        airport_service airport_service_1 ,
8        city city_1 ,
9        airport_service airport_service_2 ,
```

```
10        city city_2
11  WHERE flight_1.airline_code = 'AA'
12        AND ( flight_1.from_airport = airport_service_1.airport_code
13              AND airport_service_1.city_code = city_1.city_code
14              AND city_1.city_name = 'DALLAS'
15              AND flight_1.to_airport = airport_service_2.airport_code
16              AND airport_service_2.city_code = city_2.city_code
17              AND city_2.city_name = 'SAN FRANCISCO' )
18  """
19
20  rule_based_trial(example_8, gold_sql_8)
```

Sentence:  list the flights from dallas to san francisco on american airlines

Parse:

```
                                                                          S
                                                      _____|_____
                           |                                                     NP_F
                           |
                           |                                                     NOM_
                           |
                           |                                                     N_F
                           |                                             _____
                           |                                            N_FLIGHT
                           |                               _____|_____
                           |                              N_FLIGHT
                           |                        _____|_____
                           |                       |                 PP
                           |                       |                 |
                  PREIGNORE                        |              PP_PLACE                    PP_
          _____|_____                      |           _____|_____           _____
         |                   |       PREIGNORE  N_FLIGHT       |              N_PLACE       |
         |                   |          |          |          |           |               |
  PREIGNORESYMBOL     PREIGNORESYMBOL TERM_FLIGHT P_PLACE             TERM_PLACE P_PLACE
         |                   |          |          |                       |          |
         |                   |          |          |                       |          |
        list                the      flights      from                   dallas      to        s
```

Predicted SQL:

```
 SELECT DISTINCT flight.flight_id FROM flight WHERE     1 AND flight.from_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "DALLAS"))
    AND flight.to_airport IN
    (SELECT airport_service.airport_code FROM airport_service WHERE airport_service.city_code I
      (SELECT city.city_code FROM city WHERE city.city_name = "SAN FRANCISCO"))
    AND flight.airline_code = 'AA'
```

Predicted DB result:

 [(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,), (111091,), (1110

Gold DB result:

 [(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,), (111091,), (1110

Correct!

## ▾ Systematic evaluation on a test set

We can perform a more systematic evaluation by checking the accuracy of the queries on an entire test set for which we have gold queries. The `evaluate` function below does just this, calculating precision, recall, and F1 metrics for the test set. It takes as argument a "predictor" function, which maps token sequences to predicted SQL queries. We've provided a predictor function for the rule-based model in the next cell (and a predictor for the seq2seq system below when we get to that system).

The rule-based system does not generate predictions for all queries; many queries won't parse. The precision and recall metrics take this into account in measuring the efficacy of the method. The recall metric captures what proportion of *all of the test examples* for which the system generates a correct query. The precision metric captures what proportion of *all of the test examples for which a prediction is generated* for which the system generates a correct query. (Recall that F1 is just the geometric mean of precision and recall.)

Once you've made some progress on adding augmentations to the grammar, you can **evaluate** your progress by seeing if the precision and recall have improved. For reference, the solution code achieves precision of about 66% and recall of about 28% for an F1 of 39%.

```python
1 def evaluate(predictor, dataset, num_examples=0, silent=True):
2   """Evaluate accuracy of `predictor` by executing predictions on a
3   SQL database and comparing returned results against those of gold queries.
4
5   Arguments:
6     predictor:    a function that maps a token sequence (provided by torchtext)
7                   to a predicted SQL query string
8     dataset:      the dataset of token sequences and gold SQL queries
9     num_examples: number of examples from `dataset` to use; all of
10                   them if 0
11    silent: if set to False, will print out logs
12  Returns: precision, recall, and F1 score
13  """
14  # Prepare to count results
15  if num_examples <= 0:
16    num_examples = len(dataset)
17  example_count = 0
18  predicted_count = 0
19  correct = 0
20  incorrect = 0
21
22  # Process the examples from the dataset
23  for example in tqdm(dataset[:num_examples]):
24    example_count += 1
25    # obtain query SQL
26    predicted_sql = predictor(example.src)
27    if predicted_sql == None:
28      continue
29    predicted_count += 1
30    # obtain gold SQL
31    gold_sql = ' '.join(example.tgt)
32
33    # check that they're compatible
34    if verify(predicted_sql, gold_sql,silent):
35      correct += 1
36    else:
37      incorrect += 1
38
39  # Compute and return precision, recall, F1
40  precision = correct / predicted_count if predicted_count > 0 else 0
```

```
41  recall = correct / example_count
42  f1 = (2 * precision * recall) / (precision + recall) if precision + recall > 0 else 0
43  return precision, recall, f1
```

```
 1 def rule_based_predictor(tokens):
 2   query = ' '.join(tokens)    # detokenized query
 3   tree = parse_tree(query)
 4   if tree is None:
 5     return None
 6   try:
 7     predicted_sql = interpret(tree, atis_augmentations)
 8   except Exception as err:
 9     return None
10   return predicted_sql
```

```
1 precision, recall, f1 = evaluate(rule_based_predictor, test_iter.dataset, num_examples=0)
2 print(f"precision: {precision:3.2f}")
3 print(f"recall:    {recall:3.2f}")
4 print(f"F1:        {f1:3.2f}")
```

```
    100%|██████████| 332/332 [00:07<00:00, 46.28it/s]precision: 0.66
    recall:    0.28
    F1:        0.39
```

# End-to-End Seq2Seq Model

In this part, you will implement a seq2seq model **with attention mechanism** to directly learn the translation from NL query to SQL. You might find labs 4-4 and 4-5 particularly helpful, as the primary difference here is that we are using a different dataset.

**Note:** We recommend using GPUs to train the model in this part (one way to get GPUs is to use Google Colab and clicking Menu -> Runtime -> Change runtime type -> GPU), as we need to use a very large model to solve the task well. For development we recommend starting with a smaller model and training for only 1 epoch.

# Goal 2: Implement a seq2seq model (with attention)

In lab 4-5, you implemented a neural encoder-decoder model with attention. That model was used to convert English number phrases to numbers, but one of the biggest advantages of neural models is that we can easily apply them to different tasks (such as machine translation and document summarization) by using different training datasets.

$$P(y_1|y_0) \qquad P(y_2|y_{<2}) \quad P(y_3 = \langle eos \rangle | y_{<3})$$



Implement the class `AttnEncoderDecoder` to convert natural language queries into SQL statements. You may find that you can reuse most of the code you wrote for lab 4-5. A reasonable way to proceed is to implement the following methods:

- **Model**

    1. `__init__` : an initializer where you create network modules.

    2. `forward` : given source word ids of size `(max_src_len, batch_size)`, source lengths of size `(batch_size)` and decoder input target word ids `(max_tgt_len, batch_size)`, returns logits `(max_tgt_len, batch_size, V_tgt)`. For better modularity you might want to implement it by implementing two functions `forward_encoder` and `forward_decoder`.

- **Optimization**

    3. `train_all` : compute loss on training data, compute gradients, and update model parameters to minimize the loss.

    4. `evaluate_ppl` : **evaluate** the current model's perplexity on a given dataset iterator, we use the perplexity value on the validation set to select the best model.

- **Decoding**

    5. `predict` : Generates the target sequence given a list of source tokens using beam search decoding. Note that here you can assume the batch size to be 1 for simplicity.

```
1 ## Attention utility function from lab4-5
2
3 def attention(batched_Q, batched_K, batched_V, mask=None):
4     """
5     Performs the attention operation and returns the attention matrix
6     `batched_A` and the context matrix `batched_C` using queries
7     `batched_Q`, keys `batched_K`, and values `batched_V`.
8
9     Arguments:
10        batched_Q: (q_len, bsz, D)
11        batched_K: (k_len, bsz, D)
12        batched_V: (k_len, bsz, D)
13        mask: (bsz, q_len, k_len). An optional boolean mask *disallowing*
14              attentions where the mask value is *`False`*.
15    Returns:
16        batched_A: the normalized attention scores (bsz, q_len, k_ken)
17        batched_C: a tensor of size (q_len, bsz, D).
18    """
19    # Check sizes
20    D = batched_Q.size(-1)
```

```python
21   bsz = batched_Q.size(1)
22   q_len = batched_Q.size(0)
23   k_len = batched_K.size(0)
24   assert batched_K.size(-1) == D and batched_V.size(-1) == D
25   assert batched_K.size(1) == bsz and batched_V.size(1) == bsz
26   assert batched_V.size(0) == k_len
27
28   if mask is not None:
29     assert mask.size() == torch.Size([bsz, q_len, k_len])
30
31   batched_Q = batched_Q.transpose(0,1)
32   batched_K = batched_K.transpose(0,1)
33   batched_K = batched_K.transpose(1,2)
34   A = torch.bmm(batched_Q,batched_K)
35   if mask is not None:
36     A = A.masked_fill(~mask, float('-inf'))
37   batched_A = torch.softmax(A, dim=-1)
38   batched_V = batched_V.transpose(0,1)
39   batched_C = torch.bmm(batched_A,batched_V)
40   batched_C = batched_C.transpose(0,1)
41   # Verify that things sum up to one properly.
42   assert torch.all(torch.isclose(batched_A.sum(-1),
43                                  torch.ones(bsz, q_len).to(device)))
44   return batched_A, batched_C
```

```python
1  class AttnEncoderDecoder(nn.Module):
2   def __init__(self, src_field, tgt_field, hidden_size=64, layers=3):
3     """
4     Initializer. Creates network modules and loss function.
5     Arguments:
6         src_field: src field
7         tgt_field: tgt field
8         hidden_size: hidden layer size of both encoder and decoder
9         layers: number of layers of both encoder and decoder
10    """
11    super().__init__()
12    self.src_field = src_field
13    self.tgt_field = tgt_field
14
15    # Keep the vocabulary sizes available
16    self.V_src = len(src_field.vocab.itos)
17    self.V_tgt = len(tgt_field.vocab.itos)
18
19    # Get special word ids
20    self.padding_id_src = src_field.vocab.stoi[src_field.pad_token]
21    self.padding_id_tgt = tgt_field.vocab.stoi[tgt_field.pad_token]
22    self.bos_id = tgt_field.vocab.stoi[tgt_field.init_token]
23    self.eos_id = tgt_field.vocab.stoi[tgt_field.eos_token]
24
25    # Keep hyper-parameters available
26    self.embedding_size = hidden_size
27    self.hidden_size = hidden_size
28    self.layers = layers
29
30    # Create essential modules
31    self.word_embeddings_src = nn.Embedding(self.V_src, self.embedding_size)
32    self.word_embeddings_tgt = nn.Embedding(self.V_tgt, self.embedding_size)
33
34    # RNN cells
35    self.encoder_rnn = nn.LSTM(
```

```python
36       input_size   = self.embedding_size,
37       hidden_size  = hidden_size // 2, # to match decoder hidden size
38       num_layers   = layers,
39       bidirectional = True              # bidirectional encoder
40     )
41     self.decoder_rnn = nn.LSTM(
42       input_size   = self.embedding_size,
43       hidden_size  = hidden_size,
44       num_layers   = layers,
45       bidirectional = False             # unidirectional decoder
46     )
47
48     # Final projection layer
49     self.hidden2output = nn.Linear(2*hidden_size, self.V_tgt) # project the concatenation to logit
50
51     # Create loss function
52     self.loss_function = nn.CrossEntropyLoss(reduction='sum',
53                                              ignore_index=self.padding_id_tgt)
54
55   def forward_encoder(self, src, src_lengths):
56     """
57     Encodes source words `src`.
58     Arguments:
59         src: src batch of size (max_src_len, bsz)
60         src_lengths: src lengths of size (bsz)
61     Returns:
62         memory_bank: a tensor of size (src_len, bsz, hidden_size)
63         (final_state, context): `final_state` is a tuple (h, c) where h/c is of size
64                                 (layers, bsz, hidden_size), and `context` is `None`.
65     """
66     emb_src = self.word_embeddings_src(src)
67     src_lengths = src_lengths.tolist()
68     packed_src = pack(emb_src, src_lengths)
69     packed_output_rnn, (h, c) = self.encoder_rnn(packed_src)
70     swap_h = h.transpose(0, 1)
71     swap_c = c.transpose(0, 1)
72     join_h = swap_h.reshape(-1, int(swap_h.shape[1]/2), swap_h.shape[2]*2)
73     join_c = swap_c.reshape(-1, int(swap_c.shape[1]/2), swap_c.shape[2]*2)
74     h = join_h.transpose(0,1)
75     c = join_c.transpose(0,1)
76     h = h.contiguous()
77     c = c.contiguous()
78     memory_bank,_ = unpack(packed_output_rnn)
79     final_state= (h,c)
80     context = None
81     return memory_bank, (final_state, context)
82
83   def forward_decoder(self, encoder_final_state, tgt_in, memory_bank, src_mask):
84     """
85     Decodes based on encoder final state, memory bank, src_mask, and ground truth
86     target words.
87     Arguments:
88         encoder_final_state: (final_state, None) where final_state is the encoder
89                              final state used to initialize decoder. None is the
90                              initial context (there's no previous context at the
91                              first step).
92         tgt_in: a tensor of size (tgt_len, bsz)
93         memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder outputs
94                      at every position
95         src_mask: a tensor of size (src_len, bsz): a boolean tensor, `False` where
96                   src is padding (we disallow decoder to attend to those places).
```

```python
97      Returns:
98          Logits of size (tgt_len, bsz, V_tgt) (before the softmax operation)
99      """
100     max_tgt_length = tgt_in.size(0)
101
102     # Initialize decoder state, note that it's a tuple (state, context) here
103     decoder_states = encoder_final_state
104
105     all_logits = []
106     for i in range(max_tgt_length):
107       logits, decoder_states, attn = \
108         self.forward_decoder_incrementally(decoder_states,
109                                            tgt_in[i],
110                                            memory_bank,
111                                            src_mask,
112                                            normalize=False)
113       all_logits.append(logits)          # list of bsz, vocab_tgt
114     all_logits = torch.stack(all_logits, 0) # tgt_len, bsz, vocab_tgt
115     return all_logits
116
117   def forward(self, src, src_lengths, tgt_in):
118     """
119     Performs forward computation, returns logits.
120     Arguments:
121         src: src batch of size (max_src_len, bsz)
122         src_lengths: src lengths of size (bsz)
123         tgt_in:  a tensor of size (tgt_len, bsz)
124     """
125     src_mask = src.ne(self.padding_id_src) # max_src_len, bsz
126     # Forward encoder
127     memory_bank, encoder_final_state = self.forward_encoder(src, src_lengths)  # return memory_bar
128     # Forward decoder
129     logits = self.forward_decoder(encoder_final_state, tgt_in, memory_bank, src_mask)
130     return logits
131
132   def forward_decoder_incrementally(self, prev_decoder_states, tgt_in_onestep,
133                                     memory_bank, src_mask,
134                                     normalize=True):
135     """
136     Forward the decoder for a single step with token `tgt_in_onestep`.
137     This function will be used both in `forward_decoder` and in beam search.
138     Note that bsz can be greater than 1.
139     Arguments:
140         prev_decoder_states: a tuple (prev_decoder_state, prev_context). `prev_context`
141                              is `None` for the first step
142         tgt_in_onestep: a tensor of size (bsz), tokens at one step
143         memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder outputs
144                      at every position
145         src_mask: a tensor of size (src_len, bsz): a boolean tensor, `False` where
146                   src is padding (we disallow decoder to attend to those places).
147         normalize: use log_softmax to normalize or not. Beam search needs to normalize,
148                    while `forward_decoder` does not
149     Returns:
150         logits: log probabilities for `tgt_in_token` of size (bsz, V_tgt)
151         decoder_states: (`decoder_state`, `context`) which will be used for the
152                         next incremental update
153         attn: normalized attention scores at this step (bsz, src_len)
154     """
155     prev_decoder_state, prev_context = prev_decoder_states
156     tgt_embeddings = self.word_embeddings_tgt(tgt_in_onestep[None,...])
157     # Forward decoder RNN
```

```python
158      input_decoder = tgt_embeddings + (prev_context if prev_context is not None else 0)
159
160      decoder_outs, decoder_state = self.decoder_rnn(input_decoder, prev_decoder_state)
161
162      src_mask = src_mask.transpose(0,1)
163      src_mask = torch.unsqueeze(src_mask,1)
164      attn, attn_context = attention(decoder_outs, memory_bank, memory_bank, src_mask)
165      concated = torch.cat((decoder_outs, attn_context),dim=2)
166      logits = self.hidden2output(concated)
167      decoder_states = (decoder_state, attn_context)
168      if normalize:
169        logits = torch.log_softmax(logits, dim=-1)
170      return logits, decoder_states, attn
171
172    def evaluate_ppl(self, iterator):
173      """Returns the model's perplexity on a given dataset `iterator`."""
174      # Switch to eval mode
175      self.eval()
176      total_loss = 0
177      total_words = 0
178      for batch in iterator:
179        # Input and target
180        src, src_lengths = batch.src
181        tgt = batch.tgt # max_length_sql, bsz
182        tgt_in = tgt[:-1] # remove <eos> for decode input (y_0=<bos>, y_1, y_2)
183        tgt_out = tgt[1:] # remove <bos> as target        (y_1, y_2, y_3=<eos>)
184        # Forward to get logits
185        logits = self.forward(src, src_lengths, tgt_in)
186        # Compute cross entropy loss
187        loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.view(-1))
188        total_loss += loss.item()
189        total_words += tgt_out.ne(self.padding_id_tgt).float().sum().item()
190      return math.exp(total_loss/total_words)
191
192    def train_all(self, train_iter, val_iter, epochs=10, learning_rate=0.001):
193      """Train the model."""
194      # Switch the module to training mode
195      self.train()
196      # Use Adam to optimize the parameters
197      optim = torch.optim.Adam(self.parameters(), lr=learning_rate)
198      best_validation_ppl = float('inf')
199      best_model = None
200      # Run the optimization for multiple epochs
201      for epoch in range(epochs):
202        total_words = 0
203        total_loss = 0.0
204        for batch in tqdm(train_iter):
205          # Zero the parameter gradients
206          self.zero_grad()
207          # Input and target
208          src, src_lengths = batch.src # text: max_src_length, bsz
209          tgt = batch.tgt # max_tgt_length, bsz
210          tgt_in = tgt[:-1] # Remove <eos> for decode input (y_0=<bos>, y_1, y_2)
211          tgt_out = tgt[1:] # Remove <bos> as target        (y_1, y_2, y_3=<eos>)
212          bsz = tgt.size(1)
213          # Run forward pass and compute loss along the way.
214          logits = self.forward(src, src_lengths, tgt_in)
215          loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.view(-1))
216          # Training stats
217          num_tgt_words = tgt_out.ne(self.padding_id_tgt).float().sum().item()
218          total_words += num_tgt_words
```

```
219             total_loss += loss.item()
220             # Perform backpropagation
221             loss.div(bsz).backward()
222             optim.step()
223
224         # Evaluate and track improvements on the validation dataset
225         validation_ppl = self.evaluate_ppl(val_iter)
226         self.train()
227         if validation_ppl < best_validation_ppl:
228             best_validation_ppl = validation_ppl
229             self.best_model = copy.deepcopy(self.state_dict())
230         epoch_loss = total_loss / total_words
231         print (f'Epoch: {epoch} Training Perplexity: {math.exp(epoch_loss):.4f} '
232                 f'Validation Perplexity: {validation_ppl:.4f}')
233
234     def predict(self, tokens, max_T, K=1):
235         beam_searcher = BeamSearcher(self)
236         ## Adjust tokens to fit for BeamSearcher
237         tokens = [self.src_field.vocab.stoi[i] for i in tokens]
238         tokens = torch.IntTensor(tokens)
239         tokens = torch.unsqueeze(tokens, 1).to(device)
240         ## Adjust src_length to fit pack() later
241         src_lengths = torch.IntTensor([len(tokens)]).to(device)
242         src = tokens
243         prediction, _ = beam_searcher.beam_search(src, src_lengths, K, max_T=max_T)
244         # Convert to string
245         prediction = ' '.join([TGT.vocab.itos[token] for token in prediction])
246         prediction = prediction.lstrip('<bos>').rstrip('<eos>').strip()
247         return prediction
```

```
1  ## Beam search utility
2
3  class Beam():
4    """
5    Helper class for storing a hypothesis, its score and its decoder hidden state.
6    """
7    def __init__(self, decoder_state, tokens, score):
8      self.decoder_state = decoder_state
9      self.tokens = tokens
10     self.score = score
11
12 class BeamSearcher():
13   """
14   Main class for beam search.
15   """
16   def __init__(self, model):
17     self.model = model
18     self.bos_id = model.bos_id
19     self.eos_id = model.eos_id
20     self.padding_id_src = model.padding_id_src
21     self.V = model.V_tgt
22
23
24   def beam_search(self, src, src_lengths, K, max_T):
25     """
26     Performs beam search decoding.
27     Arguments:
28         src: src batch of size (max_src_len, 1)
29         src_lengths: src lengths of size (1)
30         K: beam size
```

```python
31              max_T: max possible target length considered
32          Returns:
33              a list of token ids and a list of attentions
34          """
35          finished = []
36          all_attns = []
37          # Initialize the beam
38          self.model.eval()
39          memory_bank, encoder_final_state = self.model.forward_encoder(src, src_lengths)
40          init_beam =  Beam(encoder_final_state,[torch.LongTensor(1).fill_(self.bos_id).to(device)], sco
41          beams = [init_beam]
42
43          with torch.no_grad():
44            for t in range(max_T): # main body of search over time steps
45
46              # Expand each beam by all possible tokens y_{t+1}
47              all_total_scores = []
48              for beam in beams:
49                y_1_to_t, score, decoder_state = beam.tokens, beam.score, beam.decoder_state
50                y_t = y_1_to_t[-1]
51                src_mask = src.ne(self.padding_id_src)
52                y_t_tensor = torch.ones(1, dtype=torch.long, device=device) * y_t
53                logits, decoder_state, attn = self.model.forward_decoder_incrementally(decoder_state,
54                                                    y_t_tensor, memory_bank, src_mask, normalize=True)
55
56                if attn is not None:
57                  attn = attn.reshape(1, -1)
58                total_scores = score + logits
59                # ours ^^^
60                all_total_scores.append(total_scores)
61                all_attns.append(attn) # keep attentions for visualization
62                beam.decoder_state = decoder_state # update decoder state in the beam
63              all_total_scores = torch.stack(all_total_scores) # (K, V) when t>0, (1, V) when t=0
64
65              # Find K best next beams
66              all_scores_flattened = all_total_scores.view(-1) # K*V when t>0, 1*V when t=0
67              topk_scores, topk_ids = all_scores_flattened.topk(K, 0)
68              beam_ids = topk_ids.div(self.V, rounding_mode='floor')
69              next_tokens = topk_ids - beam_ids * self.V
70              new_beams = []
71              for k in range(K):
72                beam_id = beam_ids[k]        # which beam it comes from
73                y_t_plus_1 = next_tokens[k] # which y_{t+1}
74                score = topk_scores[k]
75                beam = beams[beam_id]
76                decoder_state = beam.decoder_state
77                y_1_to_t = beam.tokens
78                new_beam = Beam(decoder_state, y_1_to_t + [y_t_plus_1], score)  # ours
79                new_beams.append(new_beam)
80              beams = new_beams
81
82              # Set aside completed beams
83              new_beams = []
84              for beam in beams:
85                if beam.tokens[-1] == self.eos_id:
86                  finished.append(beam)
87                else:
88                  new_beams.append(beam)
89              beams = new_beams
90
91              # Break the loop if everything is completed
```

```
92          if len(beams) == 0:
93              break
94
95      # Return the best hypothesis
96      if len(finished) > 0:
97        finished = sorted(finished, key=lambda beam: -beam.score)
98        return finished[0].tokens, all_attns
99      else: # when nothing is finished, return an unfinished hypothesis
100       return beams[0].tokens, all_attns
```

We provide the recommended hyperparameters for the final model in the script below, but you are free to tune the hyperparameters or change any part of the provided code.

> For quick debugging, we recommend starting with smaller models (by using a very small `hidden_size`), and only a single epoch. If the model runs smoothly, then you can train the full model on GPUs.

```
1 ## Already trained a good model, so we'll just load it.
2 ## Meaning, no training anymore
3
4 model = AttnEncoderDecoder(SRC, TGT,
5   hidden_size    = 1024, ##1024
6   layers         = 1,
7 ).to(device)
8 model.load_state_dict(torch.load("./model1_params",map_location=torch.device('cpu')))
9
10 ## The code below was used to train a model and later store its' params.
11 ## The code above is used to load the params for the model.
12
13 # EPOCHS = 20 # epochs; we recommend starting with a smaller number like 1, will be 50
14 # LEARNING_RATE = 1e-4 # learning rate
15
16 # # Instantiate and train classifier
17 # model = AttnEncoderDecoder(SRC, TGT,
18 #   hidden_size    = 1024, ##1024
19 #   layers         = 1,
20 # ).to(device)
21
22 # model.train_all(train_iter, val_iter, epochs=EPOCHS, learning_rate=LEARNING_RATE)
23 # model.load_state_dict(model.best_model)
24
25 # Evaluate model performance, the expected value should be < 1.2
26 print (f'Validation perplexity: {model.evaluate_ppl(val_iter):.3f}')
```

```
    Validation perplexity: 1.094
```

With a trained model, we can convert questions to SQL statements. We recommend making sure that the model can generate at least reasonable results on the examples from before, before evaluating on the full test set.

```
1 def seq2seq_trial(sentence, gold_sql):
2   print("Sentence: ", sentence, "\n")
3   tokens = tokenize(sentence)
4
5   predicted_sql = model.predict(tokens, K=1, max_T=400)
6   print("Predicted SQL:\n\n", predicted_sql, "\n")
```

```
 7
 8   if verify(predicted_sql, gold_sql, silent=False):
 9     print ('Correct!')
10   else:
11     print ('Incorrect!')
```

```
1 seq2seq_trial(example_1, gold_sql_1)
```

Sentence:  flights from phoenix to milwaukee

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

Predicted DB result:

 [(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,), (304881,), (3106

Gold DB result:

 [(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,), (304881,), (3106

Correct!

```
1 seq2seq_trial(example_2, gold_sql_2)
```

Sentence:  i would like a united flight

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

predicted sql exec failed: no such column: airport_service_2.airport_code
Incorrect!

```
1 seq2seq_trial(example_3, gold_sql_3)
```

Sentence:  i would like a flight between boston and dallas

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

Predicted DB result:

 [(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,), (103178,), (1031

Gold DB result:

 [(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,), (103178,), (1031

Correct!

```
1 seq2seq_trial(example_4, gold_sql_4)
```

Sentence:  show me the united flights from denver to baltimore

Predicted SQL:

```
     SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

     Predicted DB result:

      [(101231,), (101233,), (305983,)]

     Gold DB result:

      [(101231,), (101233,), (305983,)]

     Correct!
```

```
 1 seq2seq_trial(example_5, gold_sql_5)

     Sentence:  show flights from cleveland to miami that arrive before 4pm

     Predicted SQL:

      SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

     Predicted DB result:

      [(107698,), (301117,)]

     Gold DB result:

      [(107698,), (301117,)]

     Correct!
```

```
 1 seq2seq_trial(example_6, gold_sql_6b)

     Sentence:  okay how about a flight on sunday from tampa to charlotte

     Predicted SQL:

      SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

     Predicted DB result:

      [(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]

     Gold DB result:

      [(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]

     Correct!
```

```
 1 seq2seq_trial(example_7, gold_sql_7b)

     Sentence:  list all flights going from boston to atlanta that leaves before 7 am on thursday

     Predicted SQL:

      SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

     Predicted DB result:

      [(100014,)]

     Gold DB result:
```

```
    [(100014,)]
```

Correct!

```
1 seq2seq_trial(example_8, gold_sql_8)
```

Sentence:  list the flights from dallas to san francisco on american airlines

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1 , airport_service airport_service_1 ,

Predicted DB result:

 [(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,), (111091,), (1110

Gold DB result:

 [(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,), (111091,), (1110

Correct!

## ▾ Evaluation

Now we are ready to run the full evaluation. A proper implementation should reach more than 35% precision/recall/F1.

```
1 def seq2seq_predictor(tokens):
2   prediction = model.predict(tokens, K=1, max_T=400)
3   return prediction
```

```
1 precision, recall, f1 = evaluate(seq2seq_predictor, test_iter.dataset, num_examples=0)
2 print(f"precision: {precision:3.2f}")
3 print(f"recall:    {recall:3.2f}")
4 print(f"F1:        {f1:3.2f}")
```

```
100%|███████████| 332/332 [01:50<00:00,  3.01it/s]precision: 0.39
recall:    0.39
F1:        0.39
```

## ▾ Goal 3: Implement a seq2seq model (with cross attention and self attention)

In the previous section, you have implemented a seq2seq model with attention. The attention mechanism used in that section is usually referred to as "cross-attention", as at each decoding step, the decoder attends to encoder outputs, enabling a dynamic view on the encoder side as decoding proceeds.

Similarly, we can have a dynamic view on the decoder side as well as decoding proceeds, i.e., the decoder attends to decoder outputs at previous steps. This is called "self attention", and has been found very useful in modern neural architectures such as transformers.

Augment the seq2seq model you implemented before with a decoder self-attention mechanism as class `AttnEncoderDecoder2`. A model diagram can be found below:



At each decoding step, the decoder LSTM first produces an output state $o_t$, then it attends to all previous output states $o_1, \ldots, o_{t-1}$ (decoder self-attention). You need to special case the first decoding step to not perform self-attention, as there are no previous decoder states. The attention result is added to $o_t$ itself and the sum is used as $q_t$ to attend to the encoder side (encoder-decoder cross-attention). The rest of the model is the same as encoder-decoder with attention.

```python
1  ## Beam search utility for New model.
2
3  class Beam2(Beam):
4      """
5      Helper class for storing a hypothesis, its score and its decoder hidden state.
6      """
7      def __init__(self, decoder_state, tokens, score):
8          self.decoder_state = decoder_state
9          self.tokens = tokens
10         self.score = score
11         self.prev_outs = None
12
13  class BeamSearcher2(BeamSearcher):
14      """
15      Main class for beam search.
16      """
17      def beam_search(self, src, src_lengths, K, max_T):
18          """
19          Performs beam search decoding.
20          Arguments:
21              src: src batch of size (max_src_len, 1)
22              src_lengths: src lengths of size (1)
23              K: beam size
24              max_T: max possible target length considered
25          Returns:
26              a list of token ids and a list of attentions
27          """
28          finished = []
29          all_attns = []
30          # Initialize the beam
31          self.model.eval()
```

```
32      memory_bank, encoder_final_state = self.model.forward_encoder(src, src_lengths)
33      init_beam =  Beam2(encoder_final_state,[torch.LongTensor(1).fill_(self.bos_id).to(device)], sc
34      beams = [init_beam]
35      with torch.no_grad():
36        for t in range(max_T): # main body of search over time steps
37          # Expand each beam by all possible tokens y_{t+1}
38          all_total_scores = []
39          for beam in beams:
40            y_1_to_t, score, decoder_state = beam.tokens, beam.score, beam.decoder_state
41            y_t = y_1_to_t[-1]
42            src_mask = src.ne(self.padding_id_src)
43
44            y_t_tensor = torch.ones(1, dtype=torch.long, device=device) * y_t
45            logits, decoder_state, attn, new_out = self.model.forward_decoder_incrementally(decoder_
46                                                    y_t_tensor, memory_bank, src_mask, beam.prev_outs, r
47            beam.prev_outs = torch.cat((beam.prev_outs, new_out), dim=0) if beam.prev_outs is not No
48            if attn is not None:
49              attn = attn.reshape(1, -1)
50            total_scores = score + logits
51            all_total_scores.append(total_scores)
52            all_attns.append(attn) # keep attentions for visualization
53            beam.decoder_state = decoder_state # update decoder state in the beam
54          all_total_scores = torch.stack(all_total_scores) # (K, V) when t>0, (1, V) when t=0
55
56          # Find K best next beams
57          all_scores_flattened = all_total_scores.view(-1) # K*V when t>0, 1*V when t=0
58          topk_scores, topk_ids = all_scores_flattened.topk(K, 0)
59          beam_ids = topk_ids.div(self.V, rounding_mode='floor')
60          next_tokens = topk_ids - beam_ids * self.V
61          new_beams = []
62          for k in range(K):
63            beam_id = beam_ids[k]        # which beam it comes from
64            y_t_plus_1 = next_tokens[k] # which y_{t+1}
65            score = topk_scores[k]
66            beam = beams[beam_id]
67            decoder_state = beam.decoder_state
68            y_1_to_t = beam.tokens
69            new_beam = Beam2(decoder_state, y_1_to_t + [y_t_plus_1], score)
70            new_beams.append(new_beam)
71          beams = new_beams
72
73          # Set aside completed beams
74          new_beams = []
75          for beam in beams:
76            if beam.tokens[-1] == self.eos_id:
77              finished.append(beam)
78            else:
79              new_beams.append(beam)
80          beams = new_beams
81
82          # Break the loop if everything is completed
83          if len(beams) == 0:
84              break
85
86    # Return the best hypothesis
87    if len(finished) > 0:
88      finished = sorted(finished, key=lambda beam: -beam.score)
89      return finished[0].tokens, all_attns
90    else: # when nothing is finished, return an unfinished hypothesis
91      return beams[0].tokens, all_attns
```

```python
1  ## Our latest implementation
2
3  class AttnEncoderDecoder2(nn.Module):
4    def __init__(self, src_field, tgt_field, hidden_size=64, layers=3):
5      """
6      Initializer. Creates network modules and loss function.
7      Arguments:
8          src_field: src field
9          tgt_field: tgt field
10         hidden_size: hidden layer size of both encoder and decoder
11         layers: number of layers of both encoder and decoder
12      """
13     super().__init__()
14     self.src_field = src_field
15     self.tgt_field = tgt_field
16
17     # Keep the vocabulary sizes available
18     self.V_src = len(src_field.vocab.itos)
19     self.V_tgt = len(tgt_field.vocab.itos)
20
21     # Get special word ids
22     self.padding_id_src = src_field.vocab.stoi[src_field.pad_token]
23     self.padding_id_tgt = tgt_field.vocab.stoi[tgt_field.pad_token]
24     self.bos_id = tgt_field.vocab.stoi[tgt_field.init_token]
25     self.eos_id = tgt_field.vocab.stoi[tgt_field.eos_token]
26
27     # Keep hyper-parameters available
28     self.embedding_size = hidden_size
29     self.hidden_size = hidden_size
30     self.layers = layers
31
32     # Create essential modules
33     self.word_embeddings_src = nn.Embedding(self.V_src, self.embedding_size)
34     self.word_embeddings_tgt = nn.Embedding(self.V_tgt, self.embedding_size)
35
36     # RNN cells
37     self.encoder_rnn = nn.LSTM(
38       input_size    = self.embedding_size,
39       hidden_size   = hidden_size // 2, # to match decoder hidden size
40       num_layers    = layers,
41       bidirectional = True              # bidirectional encoder
42     )
43     self.decoder_rnn = nn.LSTM(
44       input_size    = self.embedding_size,
45       hidden_size   = hidden_size,
46       num_layers    = layers,
47       bidirectional = False             # unidirectional decoder
48     )
49
50     # Final projection layer
51     self.hidden2output = nn.Linear(2*hidden_size, self.V_tgt) # project the concatenation to logit
52
53     # Create loss function
54     self.loss_function = nn.CrossEntropyLoss(reduction='sum',
55                                              ignore_index=self.padding_id_tgt)
56
57    def forward(self, src, src_lengths, tgt_in):
58        """
59        Performs forward computation, returns logits.
60        Arguments:
61            src: src batch of size (max_src_len, bsz)
```

```
62            src_lengths: src lengths of size (bsz)
63            tgt_in:  a tensor of size (tgt_len, bsz)
64         """
65         src_mask = src.ne(self.padding_id_src) # max_src_len, bsz
66         # Forward encoder
67         memory_bank, encoder_final_state = self.forward_encoder(src, src_lengths)
68         # Forward decoder
69         logits = self.forward_decoder(encoder_final_state, tgt_in, memory_bank, src_mask)
70         return logits
71
72     def forward_encoder(self, src, src_lengths):
73        """
74        Encodes source words `src`.
75        Arguments:
76            src: src batch of size (max_src_len, bsz)
77            src_lengths: src lengths of size (bsz)
78        Returns:
79            memory_bank: a tensor of size (src_len, bsz, hidden_size)
80            (final_state, context): `final_state` is a tuple (h, c) where h/c is of size
81                                    (layers, bsz, hidden_size), and `context` is `None`.
82        """
83        emb_src = self.word_embeddings_src(src)
84        src_lengths = src_lengths.tolist()
85        packed_src = pack(emb_src, src_lengths)
86        packed_output_rnn, (h, c) = self.encoder_rnn(packed_src)
87        swap_h = h.transpose(0, 1)
88        swap_c = c.transpose(0, 1)
89        join_h = swap_h.reshape(-1, int(swap_h.shape[1]/2), swap_h.shape[2]*2)
90        join_c = swap_c.reshape(-1, int(swap_c.shape[1]/2), swap_c.shape[2]*2)
91        h = join_h.transpose(0,1)
92        c = join_c.transpose(0,1)
93        h = h.contiguous()
94        c = c.contiguous()
95        memory_bank,_ = unpack(packed_output_rnn)
96        final_state= (h,c)
97        context = None
98        return memory_bank, (final_state, context)
99
100    def forward_decoder(self, encoder_final_state, tgt_in, memory_bank, src_mask):
101       """
102       Decodes based on encoder final state, memory bank, src_mask, and ground truth
103       target words.
104       Arguments:
105           encoder_final_state: (final_state, None) where final_state is the encoder
106                                final state used to initialize decoder. None is the
107                                initial context (there's no previous context at the
108                                first step).
109           tgt_in: a tensor of size (tgt_len, bsz)
110           memory_bank: a tensor of size (src_len, bsz, hidden_size), encoder outputs
111                        at every position
112           src_mask: a tensor of size (src_len, bsz): a boolean tensor, `False` where
113                     src is padding (we disallow decoder to attend to those places).
114       Returns:
115           Logits of size (tgt_len, bsz, V_tgt) (before the softmax operation)
116       """
117       max_tgt_length = tgt_in.size(0)
118
119       # Initialize decoder state
120       decoder_states = encoder_final_state
121
122       all_logits = []
```

```python
123      prev_outs = None
124      for i in range(max_tgt_length):
125        self_attn_mask = None
126        logits, decoder_states, attn, new_out = \
127        self.forward_decoder_incrementally(decoder_states,
128                                           tgt_in[i],
129                                           memory_bank,
130                                           src_mask,
131                                           prev_outs,
132                                           self_attn_mask,
133                                           normalize=False)
134      prev_outs = torch.cat((prev_outs, new_out), dim=0) if prev_outs is not None else new_out
135      all_logits.append(logits)              # list of bsz, vocab_tgt
136
137    all_logits = torch.stack(all_logits, 0) # tgt_len, bsz, vocab_tgt
138    return all_logits
139
140  def forward(self, src, src_lengths, tgt_in):
141    """
142    Performs forward computation, returns logits.
143    Arguments:
144        src: src batch of size (max_src_len, bsz)
145        src_lengths: src lengths of size (bsz)
146        tgt_in:  a tensor of size (tgt_len, bsz)
147    """
148    src_mask = src.ne(self.padding_id_src) # max_src_len, bsz
149    # Forward encoder
150    memory_bank, encoder_final_state = self.forward_encoder(src, src_lengths)  # return memory_ban
151    # Forward decoder
152    logits = self.forward_decoder(encoder_final_state, tgt_in, memory_bank, src_mask)
153    return logits
154
155  def forward_decoder_incrementally(self, prev_decoder_states, tgt_in_onestep, memory_bank,
156                                    src_mask, prev_outs ,self_mask=None, normalize=True):
157
158    prev_decoder_state, prev_context = prev_decoder_states
159    tgt_embeddings = self.word_embeddings_tgt(tgt_in_onestep[None,...])
160    # Forward decoder RNN
161    input_decoder = (tgt_embeddings + prev_context) if prev_context is not None else tgt_embedding
162
163    decoder_outs, decoder_state = self.decoder_rnn(input_decoder, prev_decoder_state)
164    out = decoder_outs
165    if prev_outs is not None:
166      self_attn, self_context = attention(out, prev_outs, prev_outs, self_mask)
167      out = out + self_context
168    src_mask = src_mask.transpose(0,1)
169    src_mask = torch.unsqueeze(src_mask,1)
170    attn, attn_context = attention(out, memory_bank, memory_bank, src_mask)
171
172    concated = torch.cat((decoder_outs, attn_context),dim=2)
173
174    logits = self.hidden2output(concated)
175
176    decoder_states = (decoder_state, attn_context)
177    if normalize:
178      logits = torch.log_softmax(logits, dim=-1)
179    return logits, decoder_states, attn, decoder_outs
180
181  def evaluate_ppl(self, iterator):
182    """Returns the model's perplexity on a given dataset `iterator`."""
183    # Switch to eval mode
```

```python
184        self.eval()
185        total_loss = 0
186        total_words = 0
187        for batch in iterator:
188          # Input and target
189          src, src_lengths = batch.src
190          tgt = batch.tgt # max_length_sql, bsz
191          tgt_in = tgt[:-1] # remove <eos> for decode input (y_0=<bos>, y_1, y_2)
192          tgt_out = tgt[1:] # remove <bos> as target        (y_1, y_2, y_3=<eos>)
193          # Forward to get logits
194          logits = self.forward(src, src_lengths, tgt_in)
195          # Compute cross entropy loss
196          loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.view(-1))
197          total_loss += loss.item()
198          total_words += tgt_out.ne(self.padding_id_tgt).float().sum().item()
199        return math.exp(total_loss/total_words)
200
201
202      def train_all(self, train_iter, val_iter, epochs=10, learning_rate=0.001):
203        """Train the model."""
204        # Switch the module to training mode
205        self.train()
206        # Use Adam to optimize the parameters
207        optim = torch.optim.Adam(self.parameters(), lr=learning_rate)
208        best_validation_ppl = float('inf')
209        best_model = None
210        # Run the optimization for multiple epochs
211        for epoch in range(epochs):
212          total_words = 0
213          total_loss = 0.0
214          for batch in tqdm(train_iter):
215            # Zero the parameter gradients
216            self.zero_grad()
217            # Input and target
218            src, src_lengths = batch.src # text: max_src_length, bsz
219            tgt = batch.tgt # max_tgt_length, bsz
220            tgt_in = tgt[:-1] # Remove <eos> for decode input (y_0=<bos>, y_1, y_2)
221            tgt_out = tgt[1:] # Remove <bos> as target        (y_1, y_2, y_3=<eos>)
222            bsz = tgt.size(1)
223            # Run forward pass and compute loss along the way.
224            logits = self.forward(src, src_lengths, tgt_in)
225            loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.view(-1))
226            # Training stats
227            num_tgt_words = tgt_out.ne(self.padding_id_tgt).float().sum().item()
228            total_words += num_tgt_words
229            total_loss += loss.item()
230            # Perform backpropagation
231            loss.div(bsz).backward()
232            optim.step()
233
234          # Evaluate and track improvements on the validation dataset
235          validation_ppl = self.evaluate_ppl(val_iter)
236          # Switch the module to back to training mode since evaluate() changed it to eval mode
237          self.train()
238          if validation_ppl < best_validation_ppl:
239            best_validation_ppl = validation_ppl
240            self.best_model = copy.deepcopy(self.state_dict())
241          epoch_loss = total_loss / total_words
242          print (f'Epoch: {epoch} Training Perplexity: {math.exp(epoch_loss):.4f} '
243                 f'Validation Perplexity: {validation_ppl:.4f}')
244
```

```
245
246
247   def predict(self, tokens, max_T, K=1):
248     beam_searcher = BeamSearcher2(self)
249     ## Adjust tokens to fit for the beam_searcher function
250     tokens = [self.src_field.vocab.stoi[i] for i in tokens]
251     tokens = torch.IntTensor(tokens)
252     tokens = torch.unsqueeze(tokens, 1).to(device)
253     ## Adjust src_length to fit for pack() later
254     src_lengths = torch.IntTensor([len(tokens)]).to(device)
255     src = tokens
256     prediction, _ = beam_searcher.beam_search(src, src_lengths, K, max_T=max_T)
257     # Convert to string
258     prediction = ' '.join([TGT.vocab.itos[token] for token in prediction])
259     prediction = prediction.lstrip('<bos>').rstrip('<eos>').strip()
260     return prediction
```

```
1 ## Load the already trained model, instead of training a model again.
2
3
4 model2 = AttnEncoderDecoder2(SRC, TGT,
5   hidden_size   = 1024, ##1024
6   layers        = 1,
7 ).to(device)
8 model2.load_state_dict(torch.load("./best_self_attn_model_params",map_location=torch.device('cpu')
9
10 ## The code below was used to train a model and later store its' params.
11 ## The code above is used to load the params for the model.
12
13
14 # EPOCHS = 20 # epochs, we recommend starting with a smaller number like 1
15 # LEARNING_RATE = 1e-4 # learning rate
16
17 # # Instantiate and train classifier
18 # model2 = AttnEncoderDecoder2(SRC, TGT,
19 #   hidden_size   = 1024,
20 #   layers        = 1,
21 # ).to(device)
22
23 # model2.train_all(train_iter, val_iter, epochs=EPOCHS, learning_rate=LEARNING_RATE)
24 # model2.load_state_dict(model2.best_model)
25
26 # Evaluate model performance, the expected value should be < 1.2
27 print (f'Validation perplexity: {model2.evaluate_ppl(val_iter):.3f}')
```

```
Validation perplexity: 1.102
```

## Evaluation

Now we are ready to run the full evaluation. A proper implementation should reach more than 35% precision/recall/F1.

```
1 def seq2seq_predictor2(tokens):
2   prediction = model2.predict(tokens, K=1, max_T=400)
3   return prediction
```

```
1 precision, recall, f1 = evaluate(seq2seq_predictor2, test_iter.dataset, num_examples=0)
```

```
2 print(f"precision: {precision:3.2f}")
3 print(f"recall:    {recall:3.2f}")
4 print(f"F1:        {f1:3.2f}")
```

```
    100%|██████████| 332/332 [01:56<00:00,  2.84it/s]precision: 0.38
    recall:     0.38
    F1:         0.38
```

## ▾ Goal 4: Use state-of-the-art pretrained transformers

The most recent breakthrough in natural-language processing stems from the use of pretrained transformer models. For example, you might have heard of pretrained transformers such as GPT-3 and BERT. (BERT is already used in Google search.) These models are usually trained on vast amounts of text data using variants of language modeling objectives, and researchers have found that finetuning them on downstream tasks usually results in better performance as compared to training a model from scratch.

In the previous part, you implemented an LSTM-based sequence-to-sequence approach. To "upgrade" the model to be a state-of-the-art pretrained transformer only requires minor modifications.

The pretrained model that we will use is BART, which uses a bidirectional transformer encoder and a unidirectional transformer decoder, as illustrated in the below diagram (image courtesy https://arxiv.org/pdf/1910.13461):



We can see that this model is strikingly similar to the LSTM-based encoder-decoder model we've been using. The only difference is that they use transformers instead of LSTMs. Therefore, we only need to change the modeling parts of the code, as we will see later.

First, we download and load the pretrained BART model from the transformers package by Huggingface. Note that we also need to use the "tokenizer" of BART, which is actually a combination of a tokenizer and a mapping from strings to word ids.

```
1 pretrained_bart = BartForConditionalGeneration.from_pretrained('facebook/bart-base')
2 bart_tokenizer = BartTokenizer.from_pretrained('facebook/bart-base')
```

Below we demonstrate how to use BART's tokenizer to convert a sentence to a list of word ids, and vice versa.

```
1 # BART uses a predefined "tokenizer", which directly maps a sentence
2 # to a list of ids
3 def bart_tokenize(string):
4   return bart_tokenizer(string)['input_ids'][:1024] # BART model can process at most 1024 tokens
5
6 def bart_detokenize(token_ids):
7     return bart_tokenizer.decode(token_ids, skip_special_tokens=True)
8
9 ## Demonstrating the tokenizer
10 question = 'Are there any first-class flights from St. Louis at 11pm for less than $3.50?'
11
12 tokenized_question = bart_tokenize(question)
13 print('tokenized:', tokenized_question)
14
15 detokenized_question = bart_detokenize(tokenized_question)
16 print('detokenized:', detokenized_question)
```

```
tokenized: [0, 13755, 89, 143, 78, 12, 4684, 4871, 31, 312, 4, 3217, 23, 365, 1685, 13, 540, 87
detokenized: Are there any first-class flights from St. Louis at 11pm for less than $3.50?
```

We need to reprocess the data using our new tokenizer. Note that here we set `batch_first` to `True`, since that's the expected input shape of the transformers package.

```
1 SRC_BART = tt.data.Field(include_lengths=True,    # include lengths
2                          batch_first=True,        # batches will be batch_size x max_len
3                          tokenize=bart_tokenize,  # use bart tokenizer
4                          use_vocab=False,         # bart tokenizer already converts to int ids
5                          pad_token=bart_tokenizer.pad_token_id
6                          )
7 TGT_BART = tt.data.Field(include_lengths=False,
8                          batch_first=True,        # batches will be batch_size x max_len
9                          tokenize=bart_tokenize,  # use bart tokenizer
10                         use_vocab=False,         # bart tokenizer already converts to int ids
11                         pad_token=bart_tokenizer.pad_token_id
12                         )
13 fields_bart = [('src', SRC_BART), ('tgt', TGT_BART)]
14
15 # Make splits for data
16 train_data_bart, val_data_bart, test_data_bart = tt.datasets.TranslationDataset.splits(
17     ('_flightid.nl', '_flightid.sql'), fields_bart, path='./data/',
18     train='train', validation='dev', test='test')
19
20 BATCH_SIZE = 1 # batch size for training/validation
21 TEST_BATCH_SIZE = 1 # batch size for test, we use 1 to make beam search implementation easier
22
23 train_iter_bart, val_iter_bart = tt.data.BucketIterator.splits((train_data_bart, val_data_bart),
24                                                 batch_size=BATCH_SIZE,
25                                                 device=device,
26                                                 repeat=False,
```

```
27                                                     sort_key=lambda x: len(x.src),
28                                                     sort_within_batch=True)
29 test_iter_bart = tt.data.BucketIterator(test_data_bart,
30                                     batch_size=1,
31                                     device=device,
32                                     repeat=False,
33                                     sort=False,
34                                     train=False)
```

Token indices sequence length is longer than the specified maximum sequence length for this mod

Let's take a look at the batch. Note that the shape of the batch is `batch_size x max_len`, instead of `max_len x batch_size` as in the previous part.

```
 1 batch = next(iter(train_iter_bart))
 2 train_batch_text, train_batch_text_lengths = batch.src
 3 print (f"Size of text batch: {train_batch_text.shape}")
 4 print (f"First sentence in batch: {train_batch_text[0]}")
 5 print (f"Length of the third sentence in batch: {train_batch_text_lengths[0]}")
 6 print (f"Converted back to string: {bart_detokenize(train_batch_text[0])}")
 7
 8 train_batch_sql = batch.tgt
 9 print (f"Size of sql batch: {train_batch_sql.shape}")
10 print (f"First sql in batch: {train_batch_sql[0]}")
11 print (f"Converted back to string: {bart_detokenize(train_batch_sql[0])}")
```

```
Size of text batch: torch.Size([1, 16])
First sentence in batch: tensor([    0, 12196,    16,     5,    78,  2524,    71,   316,  5996,
        14784,  1054,    13,  3069,  2802,     2], device='cuda:0')
Length of the third sentence in batch: 16
Converted back to string: what is the first flight after 12 noon from washington for denver
Size of sql batch: torch.Size([1, 338])
First sql in batch: tensor([    0, 49179,   211, 11595,  2444,  7164,  2524,  1215,   134,
        15801,  1215,   808, 11974,  2524,  2524,  1215,   134,  2156,  3062,
         1215, 11131,  3062,  1215, 11131,  1215,   134,  2156,   343,   343,
         1215,   134,  2156,  3062,  1215, 11131,  3062,  1215, 11131,  1215,
          176,  2156,   343,   343,  1215,   176, 29919,  2524,  1215,   134,
            4, 17272,  2013,  2407,  1215,   958,  5457,    36, 44664, 18335,
           36,  2524,  1215,   134,     4, 17272,  2013,  2407,  1215,   958,
         4839, 11974,  2524,  2524,  1215,   134,  2156,  3062,  1215, 11131,
         3062,  1215, 11131,  1215,   134,  2156,   343,   343,  1215,   134,
         2156,  3062,  1215, 11131,  3062,  1215, 11131,  1215,   176,  2156,
          343,   343,  1215,   176, 29919,  2524,  1215,   134,     4,  7761,
         1215,  2456,  3427,  5457,  3062,  1215, 11131,  1215,   134,     4,
         2456,  3427,  1215, 20414,  4248,  3062,  1215, 11131,  1215,   134,
            4, 14853,  1215, 20414,  5457,   343,  1215,   134,     4, 14853,
         1215, 20414,  4248,   343,  1215,   134,     4, 14853,  1215, 13650,
         5457,   128,  5762,   108,  4248,    36,  2524,  1215,   134,     4,
          560,  1215,  2456,  3427,  5457,  3062,  1215, 11131,  1215,   176,
            4,  2456,  3427,  1215, 20414,  4248,  3062,  1215, 11131,  1215,
          176,     4, 14853,  1215, 20414,  5457,   343,  1215,   176,     4,
        14853,  1215, 20414,  4248,   343,  1215,   176,     4, 14853,  1215,
        13650,  5457,   128, 28082,  9847,   108,  4248,  2524,  1215,   134,
            4, 17272,  2013,  2407,  1215,   958,  8061, 23777,  4839,  4839,
         4248,    36,  2524,  1215,   134,     4,  7761,  1215,  2456,  3427,
         5457,  3062,  1215, 11131,  1215,   134,     4,  2456,  3427,  1215,
        20414,  4248,  3062,  1215, 11131,  1215,   134,     4, 14853,  1215,
        20414,  5457,   343,  1215,   134,     4, 14853,  1215, 20414,  4248,
          343,  1215,   134,     4, 14853,  1215, 13650,  5457,   128,  5762,
          108,  4248,    36,  2524,  1215,   134,     4,   560,  1215,  2456,
         3427,  5457,  3062,  1215, 11131,  1215,   176,     4,  2456,  3427,
```

```
            1215, 20414,  4248,  3062,  1215, 11131,  1215,   176,     4, 14853,
            1215, 20414,  5457,   343,  1215,   176,     4, 14853,  1215, 20414,
            4248,   343,  1215,   176,     4, 14853,  1215, 13650,  5457,   128,
           28082,  9847,   108,  4248,  2524,  1215,   134,     4, 17272,  2013,
            2407,  1215,   958,  8061, 23777,  4839,  4839,     2],
         device='cuda:0')
   Converted back to string: SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_serv
```

Now we are ready to implement the BART-based approach for the text-to-SQL conversion problem. In the below `BART` class, we have provided the constructer `__init__`, the `forward` function, and the `predict` function. Your job is to implement the main optimization `train_all`, and `evaluate_ppl` for evaluating validation perplexity for model selection.

> Hint: you can use almost the same `train_all` and `evaluate_ppl` function you implemented before, but here a major difference is that due to setting `batch_first=True`, the batched source/target tensors are of size `batch_size x max_len`, as opposed to `max_len x batch_size` in the LSTM-based approach, and you need to make changes in `train_all` and `evaluate_ppl` accordingly.

```python
1 ## Our latest implementation
2
3 #TODO - finish implementing the `BART` class.
4 class BART(nn.Module):
5   def __init__(self, tokenizer, pretrained_bart):
6     """
7     Initializer. Creates network modules and loss function.
8     Arguments:
9         tokenizer: BART tokenizer
10        pretrained_bart: pretrained BART
11    """
12    super(BART, self).__init__()
13
14    self.V_tgt = len(tokenizer)
15
16    # Get special word ids
17    self.padding_id_tgt = tokenizer.pad_token_id
18    self.bos_id = tokenizer.bos_token_id
19
20    # Create essential modules
21    self.bart = pretrained_bart
22
23    # Create loss function
24    self.loss_function = nn.CrossEntropyLoss(reduction="sum",
25                                              ignore_index=self.padding_id_tgt)
26
27  def forward(self, src, src_lengths, tgt_in):
28    """
29    Performs forward computation, returns logits.
30    Arguments:
31        src: src batch of size (batch_size, max_src_len)
32        src_lengths: src lengths of size (batch_size)
33        tgt_in:  a tensor of size (tgt_len, bsz)
34    """
35    # BART assumes inputs to be batch-first
36    # This single function is forwarding both encoder and decoder (w/ cross attn),
37    # using `input_ids` as encoder inputs, and `decoder_input_ids`
```

```
38        # as decoder inputs.
39        tgt_in = torch.unsqueeze(tgt_in, 0)
40        logits = self.bart(input_ids=src,
41                           decoder_input_ids=tgt_in,
42                            use_cache=False
43                           ).logits
44        return logits
45
46     def evaluate_ppl(self, iterator):
47        """Returns the model's perplexity on a given dataset `iterator`."""
48        # Switch to eval mode
49        self.eval()
50        total_loss = 0
51        total_words = 0
52        for batch in iterator:
53          # Input and target
54          src, src_lengths = batch.src  # bsz,max_len_src
55          tgt = batch.tgt #  bsz,max_length_sql
56          tgt_in_without_bos = tgt[0][:-1] # Remove <eos> for decode input (y_0=<bos>, y_1, y_2)
57          tgt_in = torch.cat((torch.LongTensor([self.bos_id]).to(device), tgt_in_without_bos), dim=-1)
58          tgt_out = tgt
59          # Forward to get logits
60          logits = self.forward(src, src_lengths, tgt_in)
61          # Compute cross entropy loss
62          loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.view(-1))
63          total_loss += loss.item()
64          total_words += tgt_out.ne(self.padding_id_tgt).float().sum().item()
65        return math.exp(total_loss/total_words)
66
67     def train_all(self, train_iter, val_iter, epochs=10, learning_rate=0.001):
68        """Train the model."""
69        # Switch the module to training mode
70        self.train()
71        # Use Adam to optimize the parameters
72        optim = torch.optim.Adam(self.parameters(), lr=learning_rate)
73        best_validation_ppl = float('inf')
74        best_model = None
75        # Run the optimization for multiple epochs
76        for epoch in range(epochs):
77          total_words = 0
78          total_loss = 0.0
79          for batch in tqdm(train_iter):
80            # Zero the parameter gradients
81            self.zero_grad()
82            # Input and target
83            src, src_lengths = batch.src # text: bsz, max_src_length
84            tgt = batch.tgt # bsz, max_tgt_length
85
86            ## Current best solution, achieved 48% within 2 epochs and PP of 1.02
87            ## Remove eos and insert bos in tgt_in
88            ## tgt_out remain the same as tgt
89            tgt_in_without_eos = tgt[0][:-1] # Remove <eos> for decode input (y_0=<bos>, y_1, y_2)
90            tgt_in = torch.cat((torch.LongTensor([self.bos_id]).to(device), tgt_in_without_eos), dim=-
91            tgt_out = tgt
92            bsz = tgt.size(1)
93            # Run forward pass and compute loss along the way.
94            logits = self.forward(src, src_lengths, tgt_in)
95            loss = self.loss_function(logits.view(-1, self.V_tgt), tgt_out.view(-1))
96            # Training stats
97            num_tgt_words = tgt_out.ne(self.padding_id_tgt).float().sum().item()
98            total_words += num_tgt_words
```

```
99          total_loss += loss.item()
100         # Perform backpropagation
101         loss.div(bsz).backward()
102         optim.step()
103
104      # Evaluate and track improvements on the validation dataset
105      validation_ppl = self.evaluate_ppl(val_iter)
106      self.train()
107      if validation_ppl < best_validation_ppl:
108        best_validation_ppl = validation_ppl
109        self.best_model = copy.deepcopy(self.state_dict())
110      epoch_loss = total_loss / total_words
111      print (f'Epoch: {epoch} Training Perplexity: {math.exp(epoch_loss):.4f} '
112             f'Validation Perplexity: {validation_ppl:.4f}')
113
114
115  def predict(self, tokens, K=1, max_T=400):
116    """
117    Generates the target sequence given the source sequence using beam search decoding.
118    Note that for simplicity, we only use batch size 1.
119    Arguments:
120        tokens: a list of strings, the source sentence.
121        max_T: at most proceed this many steps of decoding
122    Returns:
123        a string of the generated target sentence.
124    """
125    string = ' '.join(tokens) # first convert to a string
126    # Tokenize and map to a list of word ids
127    inputs = torch.LongTensor(bart_tokenize(string)).to(device).view(1, -1)
128    # The `transformers` package provides built-in beam search support
129    prediction = self.bart.generate(inputs,
130                                     num_beams=K,
131                                     max_length=max_T,
132                                     early_stopping=True,
133                                     no_repeat_ngram_size=0,
134                                     decoder_start_token_id=0,
135                                     use_cache=True)[0]
136    return bart_detokenize(prediction)
```

The code below will kick off training, and **evaluate** the validation perplexity. You should expect to see a value very close to 1.

```
1 ## Load the best model params we trained already, to save computation time
2
3 bart_model = BART(bart_tokenizer,
4                   pretrained_bart
5 ).to(device)
6 bart_model.load_state_dict(torch.load("./bart_best_model_params",map_location=torch.device('cpu'))
7
```

```
<All keys matched successfully>
```

```
1 ## Trained a model already, so no training this time.
2
3 EPOCHS = 2 # epochs, we recommend starting with a smaller number like 1
4 LEARNING_RATE = 1e-5 # learning rate
5
6 # Instantiate and train classifier
```

```
 7 bart_model = BART(bart_tokenizer,
 8                    pretrained_bart
 9 ).to(device)
10
11 bart_model.train_all(train_iter_bart, val_iter_bart, epochs=EPOCHS, learning_rate=LEARNING_RATE)
12 bart_model.load_state_dict(bart_model.best_model)
13
14 # Evaluate model performance, the expected value should be < 1.2
15 print (f'Validation perplexity: {bart_model.evaluate_ppl(val_iter_bart):.3f}')
```

```
    100%|██████████| 3651/3651 [07:04<00:00,  8.61it/s]
    Epoch: 0 Training Perplexity: 1.3632 Validation Perplexity: 1.0746
    100%|██████████| 3651/3651 [07:05<00:00,  8.59it/s]
    Epoch: 1 Training Perplexity: 1.0846 Validation Perplexity: 1.0454
    Validation perplexity: 1.045
```

As before, make sure that your model is making reasonable predictions on a few examples before evaluating on the entire test set.

```
 1 def bart_trial(sentence, gold_sql):
 2   print("Sentence: ", sentence, "\n")
 3   tokens = tokenize(sentence)
 4
 5   predicted_sql = bart_model.predict(tokens, K=1, max_T=300)
 6   print("Predicted SQL:\n\n", predicted_sql, "\n")
 7
 8   if verify(predicted_sql, gold_sql, silent=False):
 9     print ('Correct!')
10   else:
11     print ('Incorrect!')
```

```
 1 bart_trial(example_1, gold_sql_1)
```

```
    Sentence:  flights from phoenix to milwaukee

    Predicted SQL:

     SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci

    Predicted DB result:

     [(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,), (304881,), (3106

    Gold DB result:

     [(108086,), (108087,), (301763,), (301764,), (301765,), (301766,), (302323,), (304881,), (3106

    Correct!
```

```
 1 bart_trial(example_2, gold_sql_2)
```

```
    Sentence:  i would like a united flight

    Predicted SQL:

     SELECT DISTINCT flight_1.flight_id FROM flight flight_1 WHERE flight_1.airline_code = 'UA' AND

    Predicted DB result:
```
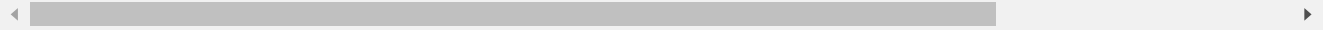
```
 [(104617,), (111035,), (111122,), (111123,)]
```

Gold DB result:

```
 [(100094,), (100099,), (100145,), (100158,), (100164,), (100167,), (100169,), (100203,), (1002
```

Incorrect!

---

`1 bart_trial(example_3, gold_sql_3)`

Sentence:  i would like a flight between boston and dallas

Predicted SQL:

```
 SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci
```
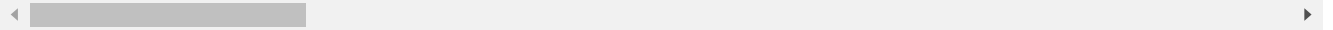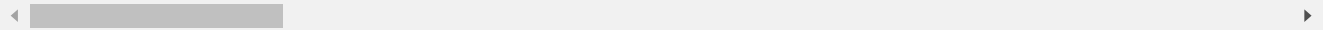
Predicted DB result:

```
 [(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,), (103178,), (1031
```

Gold DB result:

```
 [(103171,), (103172,), (103173,), (103174,), (103175,), (103176,), (103177,), (103178,), (1031
```

Correct!

---

`1 bart_trial(example_4, gold_sql_4)`

Sentence:  show me the united flights from denver to baltimore

Predicted SQL:

```
 SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci
```

Predicted DB result:

```
 [(101231,), (101233,), (305983,)]
```

Gold DB result:

```
 [(101231,), (101233,), (305983,)]
```

Correct!

---

`1 bart_trial(example_5, gold_sql_5)`

Sentence:  show flights from cleveland to miami that arrive before 4pm

Predicted SQL:

```
 SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci
```

Predicted DB result:

```
 [(107698,), (301117,)]
```

Gold DB result:

```
 [(107698,), (301117,)]
```

Correct!

```
1  bart_trial(example_6, gold_sql_6b)
```

Sentence:   okay how about a flight on sunday from tampa to charlotte

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci

Predicted DB result:

 [(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]

Gold DB result:

 [(101860,), (101861,), (101862,), (101863,), (101864,), (101865,), (305231,)]

Correct!

```
1  bart_trial(example_7, gold_sql_7b)
```

Sentence:   list all flights going from boston to atlanta that leaves before 7 am on thursday

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci

Predicted DB result:

 [(100014,), (100015,), (100016,), (100017,), (100018,), (100019,), (304692,), (307330,), (1000

Gold DB result:

 [(100014,)]

Incorrect!

```
1  bart_trial(example_8, gold_sql_8)
```

Sentence:   list the flights from dallas to san francisco on american airlines

Predicted SQL:

 SELECT DISTINCT flight_1.flight_id FROM flight flight_1, airport_service airport_service_1, ci

Predicted DB result:

 [(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,), (111091,), (1110

Gold DB result:

 [(108452,), (108454,), (108456,), (111083,), (111085,), (111086,), (111090,), (111091,), (1110

Correct!

## ▾ Evaluation

The code below will **evaluate** on the entire test set. You should expect to see precision/recall/F1 greater than 40%.

```
1 def seq2seq_predictor_bart(tokens):
2   prediction = bart_model.predict(tokens, K=4, max_T=400)
3   return prediction
```

```
1 precision, recall, f1 = evaluate(seq2seq_predictor_bart, test_iter.dataset, num_examples=0)
2 print(f"precision: {precision:3.2f}")
3 print(f"recall:    {recall:3.2f}")
4 print(f"F1:        {f1:3.2f}")
```

```
100%|██████████| 332/332 [13:50<00:00,  2.50s/it]precision: 0.43
recall:    0.43
F1:        0.43
```

# ▾ Discussion

## ▾ Goal 5: Compare the pros and cons of rule-based and neural approaches.

Compare the pros and cons of the rule-based approach and the neural approaches with relevant examples from your experiments above. Concerning the accuracy, which approach would you choose to be used in a product? Explain.

Our observations regarding pros and cons are as follows:

1. The precision of the rule based approach is higher compared to the neural approach. On the other hand, the recall of the neural approach is higher. In general, all the models achieve the same $F1$ score.

2. The rule based approach doesn't have any training phase, thus not require any training. However, We do need to struct the derivation rules by hand. This stage requires prior knowledge and understanding of the language and the domain of the problem, which is not required for the neural approach.

3. The neural approach require a lot of tagged data in order to achieve good results, while the rule based doesn't need any.

4. Inference time: We note the rule-based approach inference is much faster compared to the neural approach (about x10). We think it's because the rule based require on average $O(log_2(N))$ while using CNF grammar, while the neural approach is implemeneted here using LSTM, which is iterative (Thus at least $O(N)$) and require expensive matrices multiplication.

5. The rule-based approach is tailored for our domain (English to SQL flight queries), while the neural approach will probably work on other domains (after retraining the model).

6. As we saw with the BART model, neural approach enable us to use knowledge transfer, meaning usiong a pre-trained model and fine-tune it to our domain.

In the light of all the above we would choose a neural model for our product, since its developing and deployment will be much easier and cheaper, and we could use it for multiple tasks. A rule-based approach

could be appropriate for Real-Time applications or embedded systems where large models cannot be stored, inference has to be performed quickly, or clouds cannot be accessed remotely

## ▾ Debrief

**Question:** We're interested in any thoughts you have about this project segment so that we can improve it for later years, and to inform later segments for this year. Please list any issues that arose or comments you have to improve the project segment. Useful things to comment on might include the following:

- Was the project segment clear or unclear? Which portions?
- Were the readings appropriate background for the project segment?
- Are there additions or changes you think would make the project segment better?

but you should comment on whatever aspects you found especially positive or negative.

*This project was way too long and extensive. Should have been shorter, easier and released sooner, so we could finish it up before exam period, instead of working on it while prepeare for exams.*

## Instructions for submission of the project segment

This project segment should be submitted to Gradescope at https://rebrand.ly/project4-submit-code and https://rebrand.ly/project4-submit-pdf, which will be made available some time before the due date.

Project segment notebooks are manually graded, not autograded using otter as labs are. (Otter is used within project segment notebooks to synchronize distribution and solution code however.) **We will not run your notebook before grading it.** Instead, we ask that you submit the already freshly run notebook. The best method is to "restart kernel and run all cells", allowing time for all cells to be run to completion. You should submit your code to Gradescope at the code submission assignment at https://rebrand.ly/project4-submit-code. Make sure that you are also submitting your `data/grammar` file as part of your solution code as well.

We also request that you **submit a PDF of the freshly run notebook**. The simplest method is to use "Export notebook to PDF", which will render the notebook to PDF via LaTeX. If that doesn't work, the method that seems to be most reliable is to export the notebook as HTML (if you are using Jupyter Notebook, you can do so using `File -> Print Preview`), open the HTML in a browser, and print it to a file. Then make sure to add the file to your git commit. Please name the file the same name as this notebook, but with a `.pdf` extension. (Conveniently, the methods just described will use that name by default.) You can then perform a git commit and push and submit the commit to Gradescope at https://rebrand.ly/project4-submit-pdf.

## End of project segment 4

✓ 13m 50s    completed at 6:22 PM    ● ✕