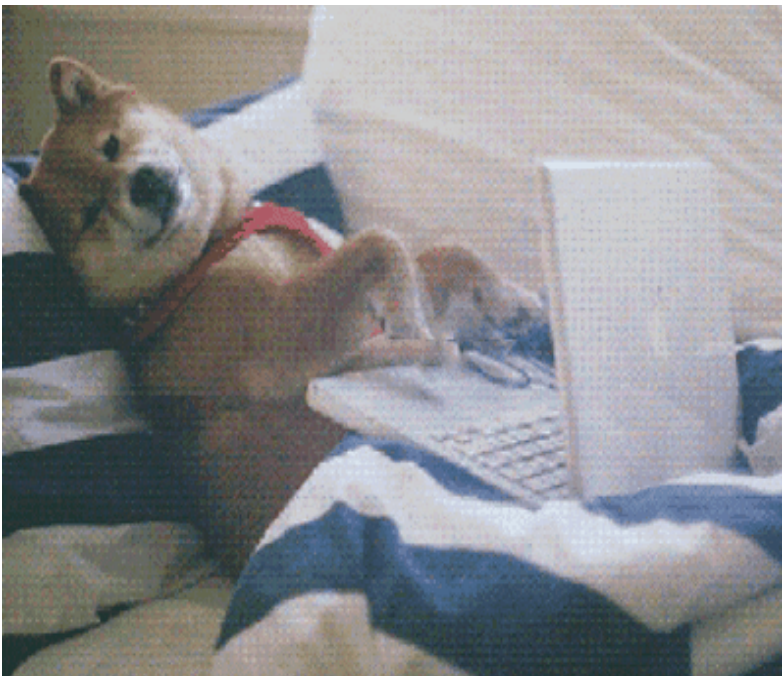


ASKING QUESTIONS TO YOUR DATABASE WITH LLMs

Nahuel Defossé  nahuel.defosse@ibm.com
IBM Research Africa

ABOUT MYSELF



- 🐍 Pythonista with 18 years of experience
 - Co-organized SciPy
 - PyCon Argentina and
- 🚜 Worked as CTO in H
- 🧪 Software Engineer a
- 🛰️ Worked in Foundat
- Geospatial application
- ... Currently working o
- core features to differ

TOMORROW








Learn about [Multimodal Geospatial Foundation Models](#) by Wanjiru, Beldine and Reggie

FOLLOW ALONG (OR AT )



INTRO

In this talk we're gonna show how to use F






-  Connect to a database and execute t
-  Convert natural language questions
-  Create a workflow
-  More advanced nodes
-  Lessons learned


TODO: Put some drawing & rephrase

CONNECT TO A DATABASE AND THE QUERIES


PUBLIC DATASETS USED IN TEXT TO SQL

These datasets define:

-  Natural language questions
-  Expected SQL
-  Database schema & content
-  Evaluation metrics
-  Leaderboard

-  Spide
-  Spi
-  BIRD
-  Arche

BIRD



BIRD-SQL

A Big Bench for Large-Scale Database Grounded Text-to-SQLs

About BIRD

Page Views 49409

BIRD (**B**ig **B**ench for **L**a**R**ge-scale **D**atabase **G**rounded **T**ext-to-**S**QL **E**valuation) represents a pioneering, cross-domain dataset that examines the impact of extensive database contents on text-to-SQL parsing. BIRD contains over **12,751** unique question-SQL pairs, **95** big databases with a total size of **33.4 GB**. It also covers more than **37** professional domains, such as blockchain, hockey, healthcare and education, etc.

Paper

Code

Train Set

🔥 Dev Set

Leaderboard - Execution Accuracy (EX)

	Model	Code	Score
	Human Performance		
	Data Engineers + DB Students		
🏆 1	ExSL + granite-20b-code		20.5
May 14, 2024	IBM Research AI		
🥈 2	Byte-SQL		3.0
May 31, 2024	ByteDance ByteBrain		
🥉 3	CHESS		1.0
May 21, 2024	Stanford		
	[Taleei et al.'24]		
4	MCS-SQL + GPT-4		0.0
Jan 14, 2024	Dunamu		

IBM Research in the lederboard 2024-06-02

BIRD MINI-DEV

- It consist of 500 queries classified as **simple** and **challenging** 

```
1 from datasets import load_dataset, DownloadConfig
2
3 # Load the dataset
4 dataset = load_dataset(
5     "birdsql/bird_mini_dev",
6     download_config=DownloadConfig(disable_tqdm=True)
7 )
8
9 # Contents
10 print("Database types: ", *dataset.keys())
11 sqlite_df = (dataset["mini_dev_sqlite"]
12     .to_pandas()
13     .pipe(lambda df: df.drop(columns=['question_id'])
14     display(sqlite_df.head(5)))
```

Database types: mini_dev_mysql mini_dev_pg
mini_dev_sqlite

	db_id	question	evidence	
0	debit_card_specializing	What is the ratio of customers who pay in EUR ...	ratio of customers who pay in EUR against cust...	SELE CAST 'EUR
1	debit_card_specializing	In 2012, who had the least consumption in LAM?	Year 2012 can be presented as Between 201201 A...	SELE FROM INNE
2	debit_card_specializing	What was the average monthly consumption of cu...	Average Monthly consumption = AVG(Consumption)...	SELE AVG(FROM
3	debit_card_specializing	What was the difference in gas consumption bet...	Year 2012 can be presented as Between 201201 A...	SELE SUM 'CZK
4	debit_card_specializing	Which year recorded the most consumption of ga...	The first 4 strings of the Date values in the ...	SELE 4) FR

DOWNLOADING BIRD DATABASES

```
1 uvx gdown https://drive.google.com/file/d/13VLWIwp
```



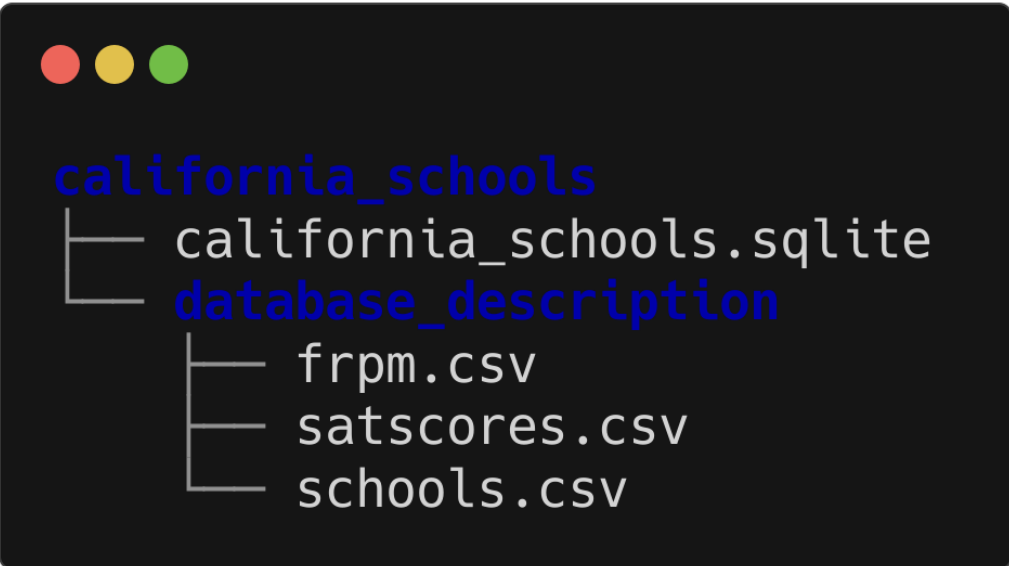
```
→ uvx gdown https://drive.google.com/file/d/13VLWIwpw5E3d5DUkMvzw7hvHE67a4
Downloading...
From (original): https://drive.google.com/uc?id=13VLWIwpw5E3d5DUkMvzw7hvHE
From (redirected): https://drive.google.com/uc?id=13VLWIwpw5E3d5DUkMvzw7hv
id=29863c54-8181-4f9e-b437-e75026526da9
To: /Users/nahueldefosse/workspace/slides/pycon-2025-insights-from-sql-llm
100%|████████████████████████████████████████████████████████████████████████████████| 800M/800M [
```

Extracting the archive (3.3GiB)

```
1 unzip minidev_703.zip
```

PICKING THE EXAMPLE DATABASE `californi`

In `minidev/MINIDEV/dev_databases/california_schools/` we find the  



```
california_schools
├── california_schools.sqlite
└── database_description
    ├── frpm.csv
    ├── satscores.csv
    └── schools.csv
```

```
1 from pathlib import Path
2 base = Path('./minidev/MINIDEV/dev_databases/calif
```

CREATING THE Engine

Creating an Engine instance connected to database.

```
1 from sqlalchemy import create_engine
2 db_path = base / 'california_schools.sqlite'
3 db_url = f"sqlite:/// {db_path}"
4 engine = create_engine(db_url)
5 engine
```

```
Engine(sqlite:///minidev/MINIDEV/dev_databases/california_schools.sqlite)
```


ER

frpm
A-Z CDSCode
A-Z Academic Year
A-Z County Code
123 District Code
A-Z School Code
A-Z County Name
A-Z District Name
A-Z School Name
A-Z District Type
A-Z School Type
A-Z Educational Option Type
A-Z NSLP Provision Status
123 Charter School (Y/N)
A-Z Charter School Number
A-Z Charter Funding Type
123 IRC
A-Z Low Grade
A-Z High Grade
123 Enrollment (K-12)
123 Free Meal Count (K-12)
123 Percent (%) Eligible Free (K-12)
123 FRPM Count (K-12)
123 Percent (%) Eligible FRPM (K-12)
123 Enrollment (Ages 5-17)
123 Free Meal Count (Ages 5-17)
123 Percent (%) Eligible Free (Ages 5-17)
123 FRPM Count (Ages 5-17)
123 Percent (%) Eligible FRPM (Ages 5-17)
123 2013-14 CALPADS Fall 1 Certification Status

schools
A-Z CDSCode
A-Z NCESDist
A-Z NCESSchool
A-Z StatusType
A-Z County
A-Z District
A-Z School
A-Z Street
A-Z StreetAbr
A-Z City
A-Z Zip
A-Z State
A-Z MailStreet
A-Z MailStrAbr
A-Z MailCity
A-Z MailZip
A-Z MailState
A-Z Phone
A-Z Ext
A-Z Website
A-Z OpenDate
A-Z ClosedDate
123 Charter
A-Z CharterNum
A-Z FundingType
A-Z DOC
A-Z DOCType
A-Z SOC
A-Z SOCType
A-Z EdOpsCode
A-Z EdOpsName
A-Z EILCode
A-Z EILName
A-Z GSoffered
A-Z GSServed

satscores
A-Z cds
A-Z rtype
A-Z sname
A-Z dname
A-Z cname
123 enroll12
123 NumTstTakr
123 AvgScrRead
123 AvgScrMath
123 AvgScrWrite
123 NumGE1500

GETTING ONE SIMPLE QUESTION

```
1 simple_queries_df = (sqlite_df
2     .pipe(lambda df: df[df.db_id == 'california_scho
3     .pipe(lambda df: df[df.difficulty == 'simple']))
4
5 simple_queries_df.head(2).set_index("db_id")
```

		question	evidence
db_id			
california_schools	How many schools with an average score in Math...	Exclusively virtual refers to Virtual = 'F'	SELECT T2.Scho
california_schools	Please list the codes of the schools with a to...	Total enrollment can be represented by `Enroll...	SELECT schools

GETTING ONE SIMPLE QUESTION

```
1 question_sql_df = (simple_queries_df
2     .pipe(lambda df: df[["question", "SQL"]])
3     .reset_index())
```

Let's take a look at the `question` and `SQL`

```
1 question, query = question_sql_df.loc[
2     0,
3     ["question", "SQL"]]
4 display(question, query)
```

'How many schools with an average score in Math greater than 500 on the state test are exclusively virtual?'

"SELECT COUNT(DISTINCT T2.School) FROM satscores AS T1
T2 ON T1.cds = T2.CDSCode WHERE T2.Virtual = 'F' AND T1.AvgScore > 500"

EXECUTE THE QUERIES

Now we run the SQL column captured in the variable `query` through S as a DataFrame.






```
1 from sqlalchemy import text
2 with engine.connect() as conn:
3     result = conn.execute(text(query))
4     res_df = pd.DataFrame(result.fetchall()) # 🐼
5     display(res_df)
```

	COUNT(DISTINCT T2.School)
0	4

Question: How many schools with an average score in Math greater than 400 in the SAT test are exclusively virtual?

SQL: SELECT COUNT(DISTINCT T2.School) FROM satscores AS T1 INNER JOIN schools AS T2 ON T1.cds = T2.CDSCode WHERE T2.Virtual = 'F' AND T1.AvgScrMath > 400

DATABASE SCHEMA WITH

LangChain (  ) community   provides a way to provide a database schema that can retrieve some schema information

```
1 !uv pip install langchain-community
```

Audited **1 package** in 5ms

```
1 from langchain_community.utilities import SQLDatabase
2 db = SQLDatabase(engine=engine)
3
4 display(db.get_usable_table_names())
```

```
['frpm', 'satscores', 'schools']
```

As we can see, the table names may not be understandable 🤔

```
1 print(db.get_table_info())
```

```
CREATE TABLE frpm (  
    "CDSCode" TEXT NOT NULL,  
    "Academic Year" TEXT,  
    "County Code" TEXT,  
    "District Code" INTEGER,  
    "School Code" TEXT,  
    "County Name" TEXT,  
    "District Name" TEXT,  
    "School Name" TEXT,  
    "District Type" TEXT,  
    "School Type" TEXT,  
    "Educational Option Type" TEXT,  
    "NSLP Provision Status" TEXT,  
    "Charter School (Y/N)" INTEGER,  
    "Charter School Number" TEXT,  
    "Charter Funding Type" TEXT,  
    "IRC" INTEGER,  
    "Low Grade" TEXT,  
    "High Grade" TEXT,  
    "Enrollment (K-12)" REAL,  
    "Free Meal Count (K-12)" REAL,  
    "Percent (%) Eligible Free (K-12)"  
REAL,  
    "FRPM Count (K-12)" REAL,  
    "Percent (%) Eligible FRPM (K-12)"  
REAL,
```

```

    "Enrollment (Ages 5-17)" REAL,
    "Free Meal Count (Ages 5-17)" REAL,
    "Percent (%) Eligible Free (Ages
5-17)" REAL,
    "FRPM Count (Ages 5-17)" REAL,
    "Percent (%) Eligible FRPM (Ages
5-17)" REAL,
    "2013-14 CALPADS Fall 1 Certification
Status" INTEGER,
    PRIMARY KEY ("CDSCode"),
    FOREIGN KEY("CDSCode") REFERENCES
schools ("CDSCode")
)
```

/*

3 rows from frpm table:

CDSCode	Academic Year	County Code
District Code	School Code	County Name
District Name	School Name	District Type
School Type	Educational Option	Type NSLP
Provision Status	Charter School	(Y/N)
Charter School Number	Charter Funding	Type
IRC	Low Grade	High Grade
Enrollment (K-12)	Free Meal Count	
(K-12)	Percent (%) Eligible Free	(K-12)
FRPM Count (K-12)	Percent (%) Eligible	
FRPM (K-12)	Enrollment (Ages 5-17)	
Free Meal Count (Ages 5-17)	Percent (%)	
Eligible Free (Ages 5-17)	FRPM Count (Ages 5-17)	

5-17) Percent (%) Eligible FRPM (Ages 5-17)
2013-14 CALPADS Fall 1 Certification Status
01100170109835 2014-2015 01 10017
0109835 Alameda Alameda County Office of
Education FAME Public Charter County
Office of Education (COE) K-12 Schools
(Public) Traditional None 1
0728 Directly funded 1 K 12
1087.0 565.0 0.519779208831647 715.0
0.657773689052438 1070.0 553.0
0.516822429906542 702.0
0.65607476635514 1
01100170112607 2014-2015 01 10017
0112607 Alameda Alameda County Office of
Education Envision Academy for Arts &
Technology County Office of Education (COE)
High Schools (Public) Traditional None
1 0811 Directly funded 1 9
12 395.0 186.0 0.470886075949367
186.0 0.470886075949367 376.0 182.0
0.484042553191489 182.0
0.484042553191489 1
01100170118489 2014-2015 01 10017
0118489 Alameda Alameda County Office of
Education Aspire California College
Preparatory Academy County Office of
Education (COE) High Schools (Public)
Traditional None 1 1049
Directly funded 1 9

```
134.0    0.549180327868853    175.0
0.717213114754098    230.0    128.0
0.556521739130435    168.0
0.730434782608696    1
*/
```

```
CREATE TABLE satscores (
    cds TEXT NOT NULL,
    rtype TEXT NOT NULL,
    sname TEXT,
    dname TEXT,
    cname TEXT,
    enroll12 INTEGER NOT NULL,
    "NumTstTskr" INTEGER NOT NULL,
    "AvgScrRead" INTEGER,
    "AvgScrMath" INTEGER,
    "AvgScrWrite" INTEGER,
    "NumGE1500" INTEGER,
    PRIMARY KEY (cds),
    FOREIGN KEY(cds) REFERENCES schools
("CDSCode")
)
```

```
/*
```




3 rows from satscores table:

cds	rtype	sname	dname	cname
enroll12		NumTstTskr		AvgScrRead
AvgScrMath		AvgScrWrite		NumGE1500

01100170000000	D	None	Alameda			
County Office of Education			Alameda	398		
88	418	418	417	14		
01100170109835	S		FAME Public Charter			
Alameda County Office of Education						
Alameda	62	17	503	546	505	
9						
01100170112607	S		Envision Academy for			
Arts & Technology			Alameda County Office of			
Education		Alameda	75	71	397	
387	395	5				

*/

```
CREATE TABLE schools (  
  "CDSCode" TEXT NOT NULL,  
  "NCESDist" TEXT,  
  "NCESSchool" TEXT,  
  "StatusType" TEXT NOT NULL,  
  "County" TEXT NOT NULL,  
  "District" TEXT NOT NULL,  
  "School" TEXT,  
  "Street" TEXT,  
  "StreetAbr" TEXT,  
  "City" TEXT,  
  "Zip" TEXT,  
  "State" TEXT,  
  "MailStreet" TEXT,  
  "MailStrAbr" TEXT,
```

Asking question to  with LLMs -  PyCon  2025

```
"MailCity" TEXT,  
"MailZip" TEXT,  
"MailState" TEXT,  
"Phone" TEXT,  
"Ext" TEXT,  
"Website" TEXT,  
"OpenDate" DATE,  
"ClosedDate" DATE,  
"Charter" INTEGER,  
"CharterNum" TEXT,  
"FundingType" TEXT,  
"DOC" TEXT NOT NULL,  
"DOCType" TEXT NOT NULL,  
"SOC" TEXT,  
"SOCType" TEXT,  
"EdOpsCode" TEXT,  
"EdOpsName" TEXT,  
"EILCode" TEXT,  
"EILName" TEXT,  
"GSoffered" TEXT,  
"GSserved" TEXT,  
"Virtual" TEXT,  
"Magnet" INTEGER,  
"Latitude" REAL,  
"Longitude" REAL,  
"AdmFName1" TEXT,  
"AdmLName1" TEXT,  
"AdmEmail1" TEXT,  
"AdmFName2" TEXT,
```

/*

CDSCode	NCESDist	NCESSchool		
StatusType	County	District		
School	Street	StreetAbr	City	Zip
State	MailStreet	MailStrAbr		
MailCity	MailZip	MailState	Phone	
Ext	Website	OpenDate	ClosedDate	
Charter	CharterNum	FundingType	DOC	
DOCType	SOC	SOCType	EdOpsCode	
EdOpsName	EILCode	EILName	GSoffered	
GSserved	Virtual	Magnet	Latitude	
Longitude	AdmFName1	AdmLName1		
AdmEmail1	AdmFName2	AdmLName2		
AdmEmail2	AdmFName3	AdmLName3		
AdmEmail3	LastUpdate			
01100170000000	0691051	None	Active	
Alameda	Alameda	County Office of Education		
None	313	West Winton Avenue	313	West
Winton Ave.	Hayward	94544-1136	Asking question to with CAMs - 2	

313 West Winton Avenue 313 West Winton Ave.
Hayward 94544-1136 CA (510)
887-0152 None www.acoe.org None
None None None None 00
County Office of Education (COE) None
None None None None None None
None None None 37.658212
-122.09713 L Karen Monroe
lkmonroe@acoe.org None None None
None None None 2015-06-23
01100170109835 0691051 10546 Closed
Alameda Alameda County Office of Education
FAME Public Charter 39899 Balentine
Drive, Suite 335 39899 Balentine Dr.,
Ste. 335 Newark 94560-5359 CA
39899 Balentine Drive, Suite 335 39899
Balentine Dr., Ste. 335 Newark 94560-5359
CA None None None 2005-08-29
2015-07-31 1 0728 Directly
funded 00 County Office of Education
(COE) 65 K-12 Schools (Public)
TRAD Traditional ELEMHIGH
Elementary-High Combination K-12 K-12
P 0 37.521436 -121.99391
None None None None None None
None None None 2015-09-01
01100170112607 0691051 10947 Active
Alameda Alameda County Office of Education
Envision Academy for Arts & Technology with LMs 15152 PyCon 2025

Webster Street 1515 Webster St.
Oakland 94612-3355 CA 1515 Webster
Street 1515 Webster St. Oakland
94612 CA (510) 596-8901 None
www.envisionacademy.org/ 2006-08-28
None 1 0811 Directly funded 00
County Office of Education (COE) 66
High Schools (Public) TRAD Traditional
HS High School 9-12 9-12 N
0 37.80452 -122.26815 Laura
Robell laura@envisionacademy.org None
None None None None None
2015-06-18
*/

CONVERT NATURAL LANGUAGE QUESTIONS INTO SQL

LLMs are quite capable of writing functional SQL queries, from the [code](#) or [coder](#) ones, to specific ones for SQL generation.

For example, some include:

- granite-20b-coder
- granite-34b-coder
- granite-20b-coder
- granite-20b-coder

[More info on the](#)



PROMPTS FOR SQL GENERATION

LLMs don't know the  structure of our database, so they can hallucinate about it, or create some flat out wrong queries.

We have to provide **extra** information about the database schema in the instructions.

For this we will use a **prompt** string with placeholders { }.

Some research papers from our team from last year:



[Weakly Supervised Detection of Hallucinations in LLMs](#)



[Localizing Persona Representations in LLMs](#)

PROMPTS

```
1 system_message = """
2 Given an input question, create a syntactically co
3 run to help find the answer. Unless the user speci
4 specific number of examples they wish to obtain, a
5 at most {top_k} results. You can order the results
6 return the most interesting examples in the databa
7
8 Never query for all the columns from a specific ta
9 few relevant columns given the question.
10
11 Pay attention to use only the column names that yo
12 description. Be careful to not query for columns t
13 pay attention to which column is in which table.
14
15 Only use the following tables:
16 {table_info}
17 """
```

source 

CREATING A PROMPT

Now we construct a list of messages. These are `dicts` which `system`, and a `content`.

```
1 def generate_messages(question, dialect="SQL", top_k=5, table_info=""):
2     # Create a ChatPromptTemplate
3     messages = [
4         {
5             "role": "system",
6             "content": system_message.format(
7                 dialect=dialect,
8                 top_k=top_k,
9                 table_info=table_info
10            )},
11         {
12             "role": "user",
13             "content": question
14         }
15     ]
16
17     return messages
18
19 messages = generate_messages(
20     question=question,
21     dialect=db.dialect, top_k=10,
22     table_info=db.get_table_info())
23 messages
```

CREATING A PROMPT

```
[{'role': 'system',
  'content': '\nGiven an input question, create a syntactically correct sqlit
answer. Unless the user specifies in his question a\nspecific number of examp
limit your query to\nat most 10 results. You can order the results by a relev
interesting examples in the database.\n\nNever query for all the columns from
the\nfew relevant columns given the question.\n\nPay attention to use only th
the schema\ndescription. Be careful to not query for columns that do not exist
column is in which table.\n\nOnly use the following tables:\n\nCREATE TABLE f
\n\t"Academic Year" TEXT, \n\t"County Code" TEXT, \n\t"District Code" INTEGER
\n\t"County Name" TEXT, \n\t"District Name" TEXT, \n\t"School Name" TEXT, \n\t
Type" TEXT, \n\t"Educational Option Type" TEXT, \n\t"NSLP Provision Status" T
INTEGER, \n\t"Charter School Number" TEXT, \n\t"Charter Funding Type" TEXT, \n
TEXT, \n\t"High Grade" TEXT, \n\t"Enrollment (K-12)" REAL, \n\t"Free Meal Cou
Eligible Free (K-12)" REAL, \n\t"FRPM Count (K-12)" REAL, \n\t"Percent (%) EL
\n\t"Enrollment (Ages 5-17)" REAL, \n\t"Free Meal Count (Ages 5-17)" REAL, \n
5-17)" REAL, \n\t"FRPM Count (Ages 5-17)" REAL, \n\t"Percent (%) Eligible FRP
CALPADS Fall 1 Certification Status" INTEGER, \n\tPRIMARY KEY ("CDSCode"), \n
schools ("CDSCode")\n\n\n/*\n3 rows from frpm table:\nCDSCode\tAcademic Year
Code\tSchool Code\tCounty Name\tDistrict Name\tSchool Name\tDistrict Type\tSc
Type\tNSLP Provision Status\tCharter School (Y/N)\tCharter School Number\tCha
```

THE system MESSAGE

Given an input question, create a syntactically correct sqlite query to run to help find the answer. Unless the user specifies in his question a specific number of examples they wish to obtain, always limit your query to at most **10** results. You can order the results by a relevant column to return the most interesting examples in the database.

Never query for all the columns from a specific table, only ask for a the few relevant columns given the question.

Pay attention to use only the column names that you can see in the schema description. Be careful to not query for columns that do not exist. Also, pay attention to which column is in which table.

Only use the following tables:

CREATE TABLE frpm (Asking question to  with LLMs -  PyCon  2025

"GPGCode" TEXT NOT NULL

CALLING THE LLM WITH THE PROMPT

`litellm` is a client for multiple LLM providers

```
1 !uv add litellm --quiet
```

To run inference, we just call the `completion` function:

```
1 import litellm
2
3 # model = "ollama_chat/granite-code:20b"
4 model = "watsonx/ibm/granite-3-2-8b-instruct"
5
6 response = litellm.completion(
7     model=model,
8     messages=messages,
9 )
```

```
1 print(response)
```

```
ModelResponse(  
    id='chatcpl-f79a3df0-2506-4e7e-aa2b-b65d  
21f5d680---5976a0e7cf57f032504389f213b02431--  
-e4f953c6-f9cd-4161-bccf-1f1fd819ea7b',  
    created=1758486227,  
    model='watsonx/ibm/granite-3-2-8b-instruc  
t',  
    object='chat.completion',  
    system_fingerprint=None,  
    choices=[  
        Choices(  
            finish_reason='stop',  
            index=0,  
            message=Message(  
                content='SELECT COUNT(*) FROM  
schools\nINNER JOIN frpm ON schools.CDSCode =  
frpm.CDSCode\nINNER JOIN satscores ON  
schools.CDSCode = satscores.cds\nWHERE  
satscores."AvgScrMath" > 400\nAND  
schools.Virtual = \'Y\'\nAND schools.Magnet =  
0\nAND frpm.CharterSchool(Y/N) = 0\nLIMIT  
10;',  
                role='assistant',  
                tool_calls=None,  
                function_call=None,
```


STRUCTURED OUTPUT

Now that we get the SQL, we're going to ask the LLM to return using a Pydantic model.

```
1 from pydantic import BaseModel, Field
2 class SQLOutput(BaseModel):
3     sql: str = Field(description="The SQL query")
4     explanation: str = Field(
5         description="The reasoning for the query construction")
6
7 # Optional
8 # litellm.enable_json_schema_validation = True
9
10 response = litellm.completion(
11     model=model,
12     messages=messages,
13     response_format=SQLOutput,
14 )
15 output = SQLOutput.model_validate_json(
16     response.choices[0].message.content)
17 output
```

STRUCTURED OUTPUT

```
SQLOutput(sql="SELECT COUNT(*) FROM schools INNER JOIN satscores ON schools.CDSCode = satscores.CDSCode WHERE schools.Virtual = 'Y' AND satscores.AvgScrMath > 400;", explanation="The query counts the number of schools from the 'schools' table that have a 'Virtual' column value of 'Y' and whose average score in Math (AvgScrMath) in the 'satscores' table is greater than 400. It uses an INNER JOIN to combine data from both tables based on the common 'CDSCode' column.")
```

Now we can use the LLM like we were getting

```
1 output.sql
```

```
"SELECT COUNT(*) FROM schools INNER JOIN satscores ON  
satscores.cds WHERE schools.Virtual = 'Y' AND satscore
```

```
1 output.explanation
```

```
"The query selects the count of virtual schools from t  
have a 'Virtual' column value of 'Y' and an 'AvgScrMat  
Math) in the 'satscores' table greater than 400. It us  
data from both tables based on the common 'CDSCode' co
```

Result.fetchall() takes **1** positional argument
but **2** were given

CREATE A WORKFLOW

CHAINING GENERATION AND

🐦 When our code starts to become large and hard to handle 🧘, **LangGraph** is a great tool to provide a reusable organization. We can add it with `langgraph`.

We will build a small pipeline where to run and execution.

Each node in the pipeline will be a function that takes a state object and returns another state object.

STATE FOR SQL WORKFLOW

LangGraph uses state that is propagated through the graph.
This state can be defined with a TypedDict or Pydantic BaseModel.

```
1 from typing import Optional, Dict, List
2 from dataclasses import dataclass, field
3 from sqlalchemy import Engine
4
5 @dataclass
6 class State:
7     question: str = ""
8     engine: Optional[Engine] = None
9     database: Optional[SQLDatabase] = None
10    messages: List[Dict[str, str]] = field(default_factory=list)
11    sql: Optional[str] = None
12    results: Optional[str] = None
```

CREATING NODES

The initial node will receive the user input state.

```
1 def init(user_input: Dict[str, str]) -> State:
2     engine = create_engine(user_input["database"])
3     database = SQLiteDatabase(engine=engine)
4     return State(
5         question=user_input["question"],
6         database=database,
7     )
```

GENERATING THE PROMPT

```
1 def make_prompt(state: State) -> State:
2     messages = [
3         {
4             "role": "system",
5             "content": system_message.format(
6                 dialect=state.database.dialect,
7                 top_k=top_k,
8                 table_info=state.database.get_table_info
9             )},
10        {
11            "role": "user",
12            "content": state.question,
13        }
14    ]
15    return state
```

CALLING THE LLM

The initial node will receive the user input state.

```
1 def generate(state: State) -> State:
2
3     response = litellm.completion(
4         model=model,
5         messages=messages,
6         response_format=SQLOutput,
7     )
8     output = SQLOutput.model_validate_json(response.json())
9     print(output)
10    engine = create_engine(user_input["database"])
11    database = SQLiteDatabase(engine=engine)
12    return State(
13        question=user_input["question"],
14        database=database,
15    )
```

EXECUTING THE SQL

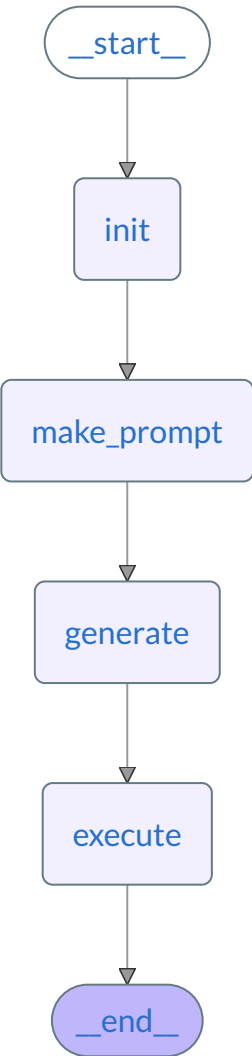
The initial node will receive the user input state.

```
1 def execute(state: State) -> State:
2     with state.engine.connect() as con:
3         state.result = conn.execute(text(state.sql))
4
5     return State(
6         question=user_input["question"],
7         database=database,
8     )
```

GRAPH CREATION

```
1 from langgraph.graph import StateGraph, START, END
2
3 graph_builder = StateGraph(State)
4 graph_builder.add_node("init", init)
5 graph_builder.add_node("make_prompt", make_prompt)
6 graph_builder.add_node("generate", generate)
7 graph_builder.add_node("execute", execute)
8 graph_builder.add_edge(START, "init")
9 graph_builder.add_edge("init", "make_prompt")
10 graph_builder.add_edge("make_prompt", "generate")
11 graph_builder.add_edge("generate", "execute")
12 graph = graph_builder.compile()
```

VISUALIZING THE GRAPH



LESSONS LEARNED

MOVE FROM FUNCTION TO CLASSES

- The workflow state can be a `dict` or `TypedDict`, but `BaseModel` is more convenient.

```
1 class State(TypedDict):
2     connection: ...
3     messages: ...
4     schema: ...
```

- A simple function that works with the state.

```
1 def create_sql(state: State) -> State:
2     # Do something
3     return State(...)
```

- This can also be used to pass multiple states

```
1 class CreateSQL:
2     def __init__(self, state: State):
3         self.state = state
4
5     def __call__(self, state: State) -> State:
6         self.state = state
7         # do something
8         return self.state
```

SQL VALIDATION

- We can pass the SQL to SQLAlchemy engine and check if the syntax is correct
- SQLGlot `parse_one()` provides as an AST
 - Finding dangerous operations (DML)
 - For any updates to the query, working better than string substitution.

ENHANCING THE CONTEXT

- Tables with a high number of columns can make the context larger
 - RAG at the column level is a common way to improve efficiency and accuracy
- Retrieving the values from the DB can help to produce better queries
- Providing some examples in the context, such as **question** and **SQL** pairs also improves results

DYNAMIC CONTEXT

- Other types of examples of well understood concepts that require some mapping.
 - For example, the term Q1 can mean, and the context should be supplied dynamically.
- Function nodes, when we want to modify the behavior of another approach is to create
 - WASM?

EXTENSIBILITY

- Offload node responsibilities into external services (e.g. DB execution)
- Pipelines defined as YAML.
- Each node is registered in a collection of nodes. Each node in the pipeline accepts the `__init__` method.
- Use Python's built in `entry-point` system. This is a good approach to separate pro-

THANK YOU