# Microsoft

# Magpie

## Python at Speed and Scale using Cloud Backends

Alekh Jindal, K. Venkatesh Emani, Maureen Daum, Olga Poppe, Brandon Haynes, Anna Pavlenko, Ayushi Gupta, Karthik Ramachandra, Carlo Curino, Andreas Mueller, Wentao Wu, Hiren Patel

Gray Systems Lab, Microsoft
Azure Data, Microsoft
Microsoft Research
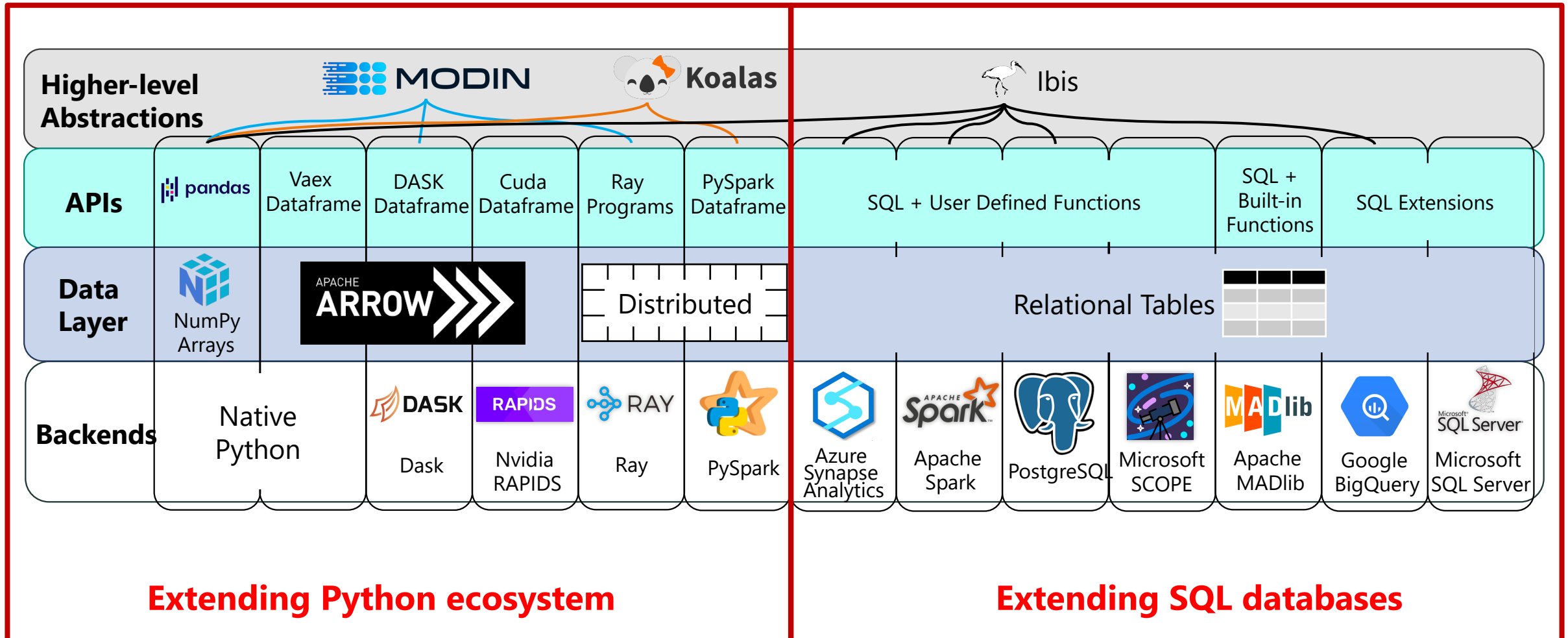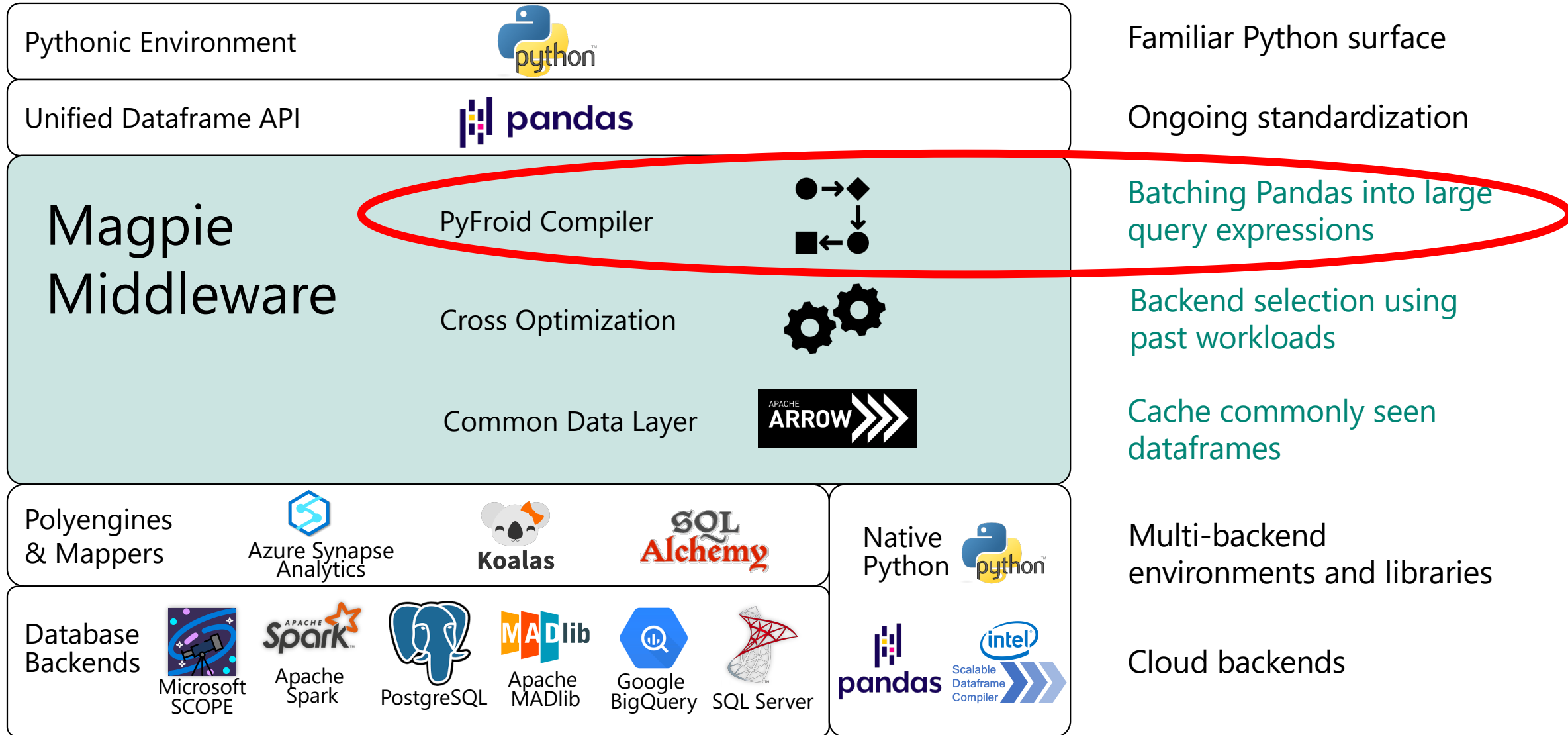University of Washington

# The Python and The Cloud



- De-facto for ad-hoc analysis
- Pandas dataframes highly popular
- Performance is a challenge

- Hyper-scale performance
- Several SQL processing backends
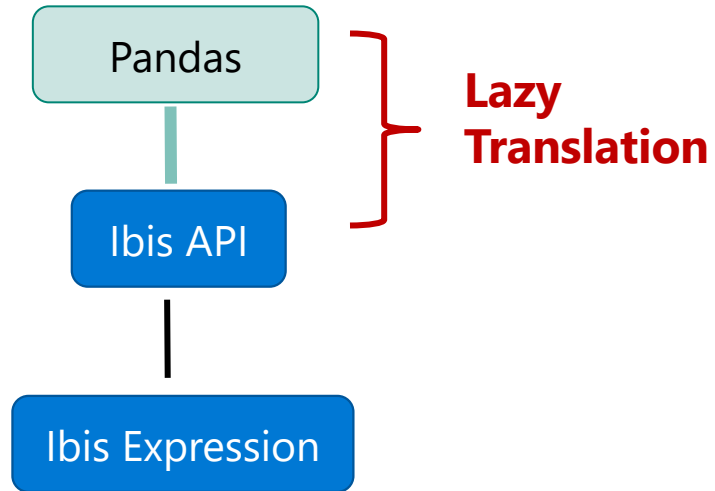- Enterprise data already on cloud

# The current landscape ... is a fragmented jungle!



**Extending Python ecosystem**

**Extending SQL databases**

# Magpie

| | | |
|---|---|---|
| Pythonic Environment | | Familiar Python surface |
| Unified Dataframe API | **pandas** | Ongoing standardization |

**Magpie Middleware**

| | | |
|---|---|---|
| PyFroid Compiler | | Batching Pandas into large query expressions |
| Cross Optimization | | Backend selection using past workloads |
| Common Data Layer | APACHE ARROW | Cache commonly seen dataframes |

| | | | | |
|---|---|---|---|---|
| Polyengines & Mappers | Azure Synapse Analytics | **Koalas** | SQL Alchemy | Native Python |

Multi-backend environments and libraries

| | | | | | |
|---|---|---|---|---|---|
| Database Backends | Microsoft SCOPE | Apache Spark | PostgreSQL | Apache MADlib | Google BigQuery | SQL Server |

pandas · intel Scalable Dataframe Compiler
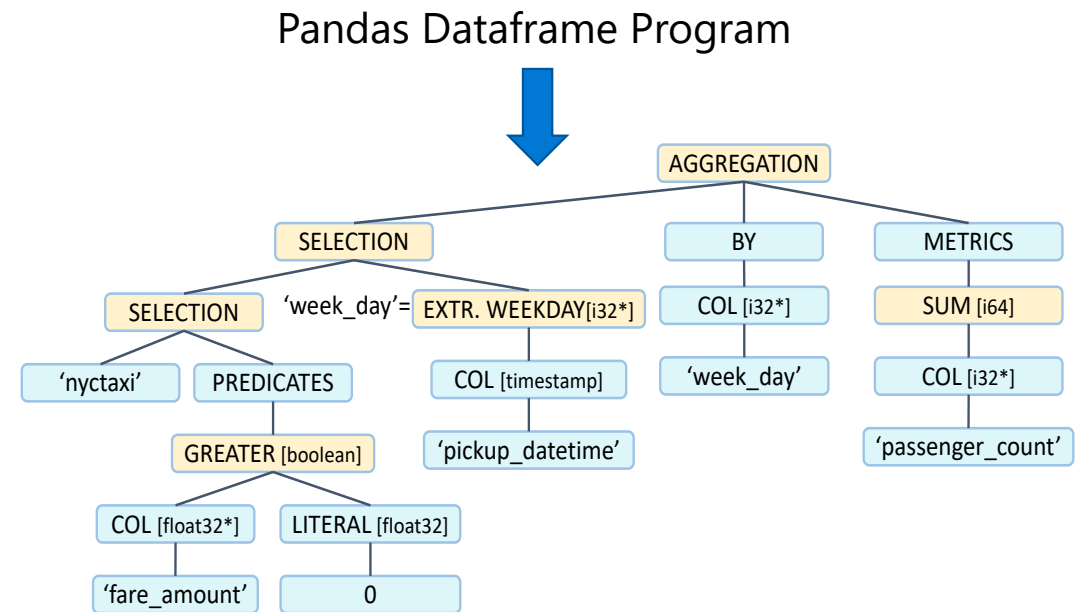
Cloud backends

# Batching Pandas

The number of taxi trips per weekday over the NYC Taxi dataset
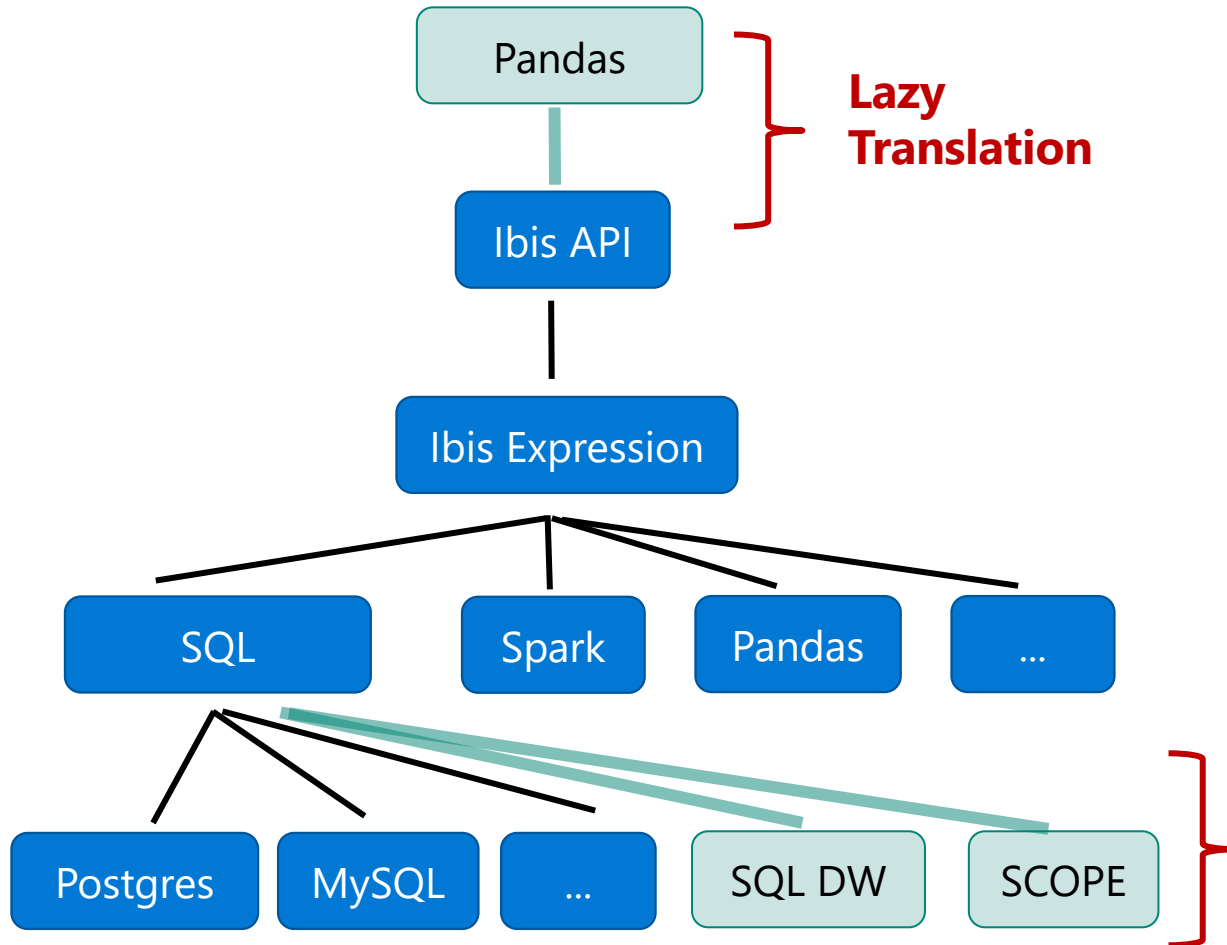
```
1  import pyfroid.pandas as pd                          # vs import pandas as pd
2  df = pd.read_sql('nyctaxi', con)                     # fetch data
3  df = df[df.fare_amount > 0]                          # filter bad rows
4  df['day'] = df.pickup_datetime.dt.dayofweek          # add features
5  df = df.groupby(['day'])['passenger_count'].sum()    # aggregation
6  print(df)                                            # use dataframe
```
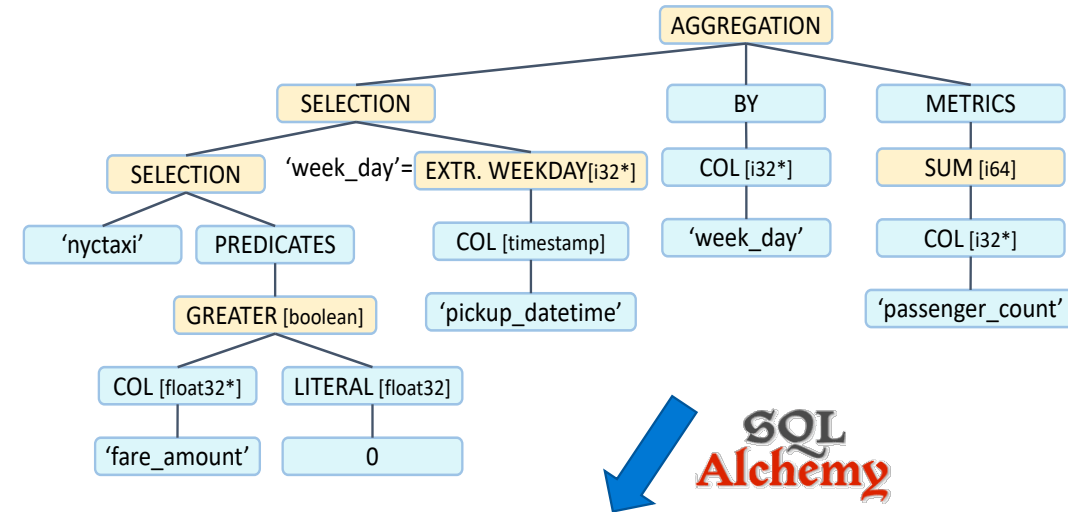
Pandas Dataframe Program

Intermediate Representation

Blue parts: already in IBIS, Green parts: our contributions

# Pushing Data Science down



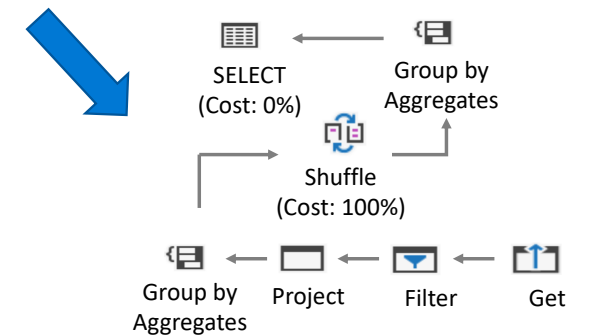**Lazy Translation**

Pandas → Ibis API

Ibis Expression → SQL, Spark, Pandas, ...

SQL → Postgres, MySQL, ..., SQL DW, SCOPE

**Cloud backends**

**Blue** parts: already in IBIS, **Green** parts: our contributions

AGGREGATION
- SELECTION
  - SELECTION
    - 'nyctaxi'
    - PREDICATES
      - GREATER [boolean]
        - COL [float32*]
          - 'fare_amount'
        - LITERAL [float32]
          - 0
  - 'week_day'= EXTR. WEEKDAY[i32*]
    - COL [timestamp]
      - 'pickup_datetime'
- BY
  - COL [i32*]
    - 'week_day'
- METRICS
  - SUM [i64]
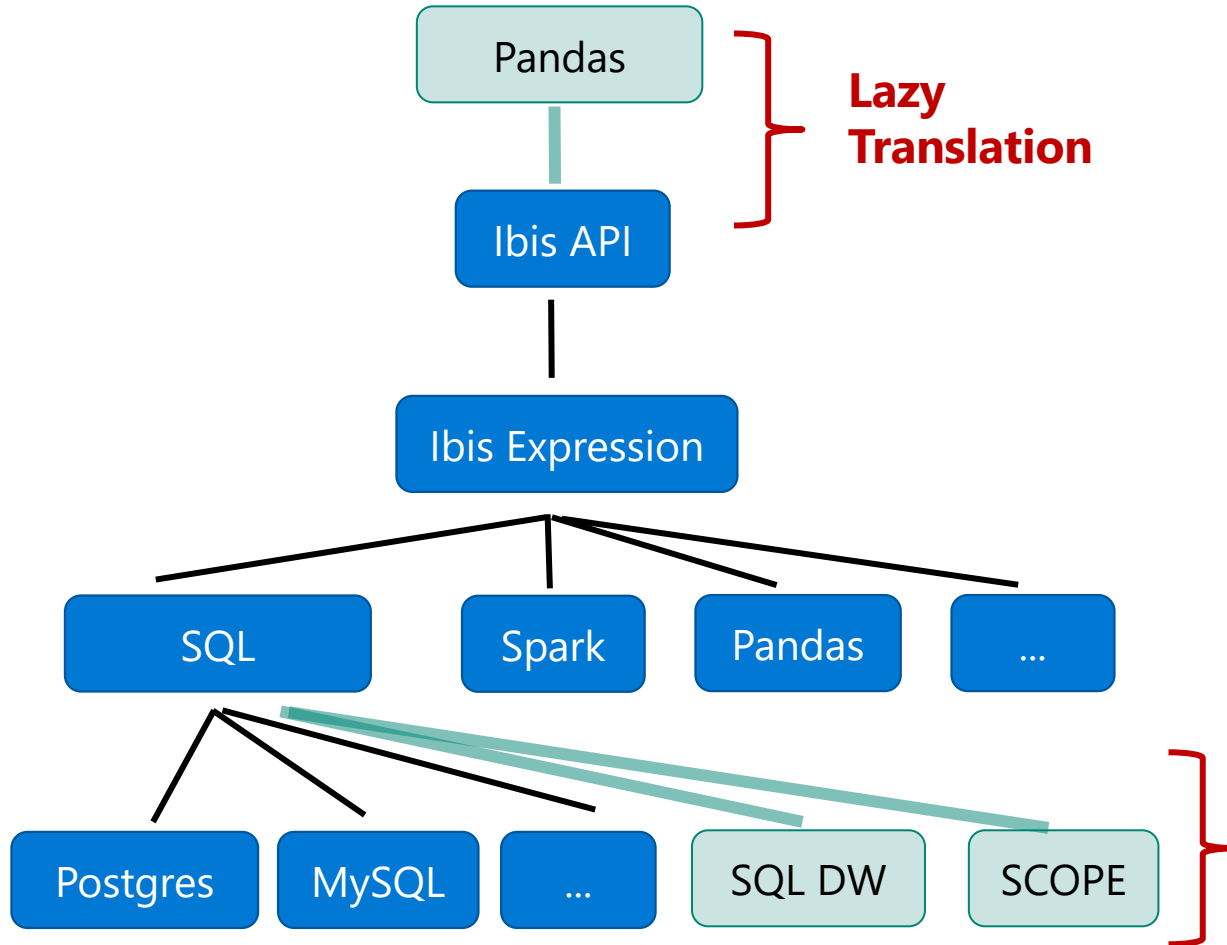    - COL [i32*]
      - 'passenger_count'

SQL Alchemy

SELECT DATEPART(WEEKDAY, pickup_datetime) AS day,
       SUM(passenger_count)
FROM nyctaxi WHERE fare_amount > 0
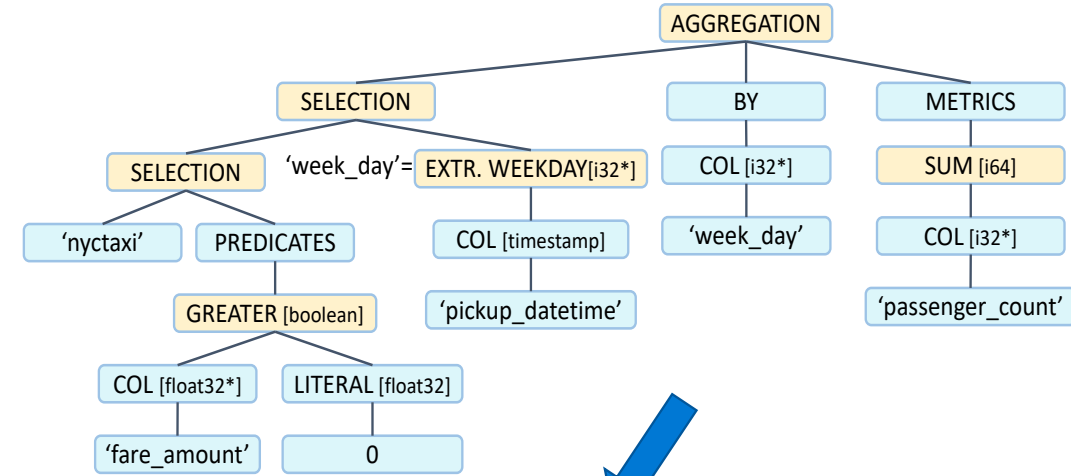GROUP BY DATEPART(WEEKDAY, pickup_datetime)

T-SQL Statement

SELECT (Cost: 0%) ← Group by Aggregates

Shuffle (Cost: 100%)

Group by Aggregates ← Project ← Filter ← Get

SQL DW Execution Plan

# Pushing Data Science down



**Blue** parts: already in IBIS, **Green** parts: our contributions

SCOPE Execution Plan
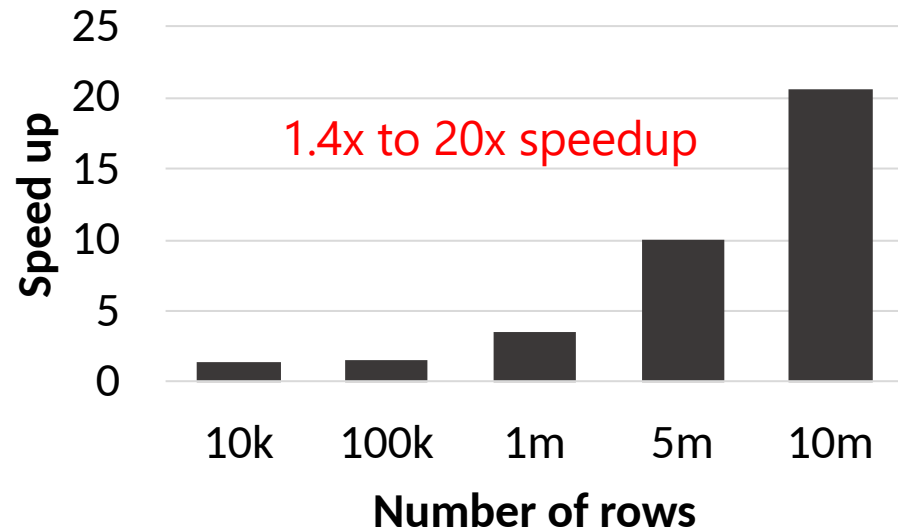
# Impact: speed-up using SQL DW

Growing input size

Growing query complexity

1.4x to 20x speedup

7x to 380x speedup
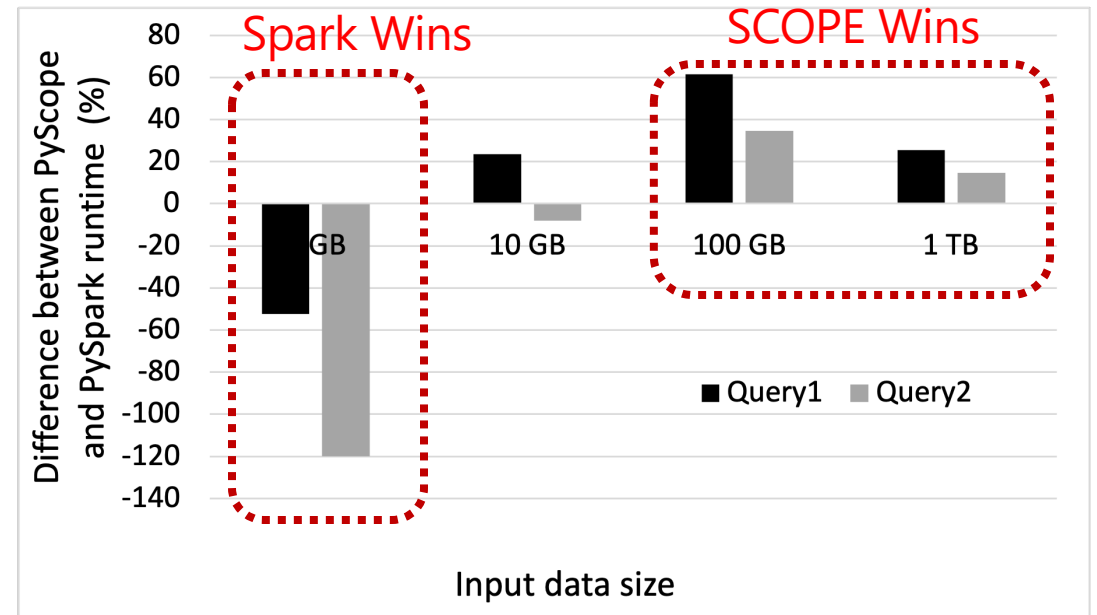
# Impact: scale-out using SCOPE



Scale data science to big data!

SCOPE vs Spark

# Magpie

| | | | |
|---|---|---|---|
| Pythonic Environment |  python | | Familiar Python surface |
| Unified Dataframe API |  pandas | | Ongoing standardization |

**Magpie Middleware**

| | | |
|---|---|---|
| PyFroid Compiler | ●→◆ ■←● | Batching Pandas into large query expressions |
| Cross Optimization | ⚙⚙ | Backend selection using past workloads |
| Common Data Layer | APACHE ARROW ⟩⟩⟩ | Cache commonly seen dataframes |

| | | | | | |
|---|---|---|---|---|---|
| Polyengines & Mappers | Azure Synapse Analytics | **Koalas** | SQL Alchemy | Native Python python | Multi-backend environments and libraries |
| Database Backends | Microsoft SCOPE  Apache Spark  PostgreSQL  Apache MADlib  Google BigQuery  SQL Server | | | pandas  intel Scalable Dataframe Compiler | Cloud backends |

# Backend Selection



**Lazy Translation**

**Cost-based optimization**

**Cloud backends**

Pandas → Ibis API → Ibis Expression → Backend Selection → SQL, Spark, Pandas, ...

SQL → Postgres, MySQL, ..., SQL DW, SCOPE

Blue parts: already in IBIS, Green parts: our contributions
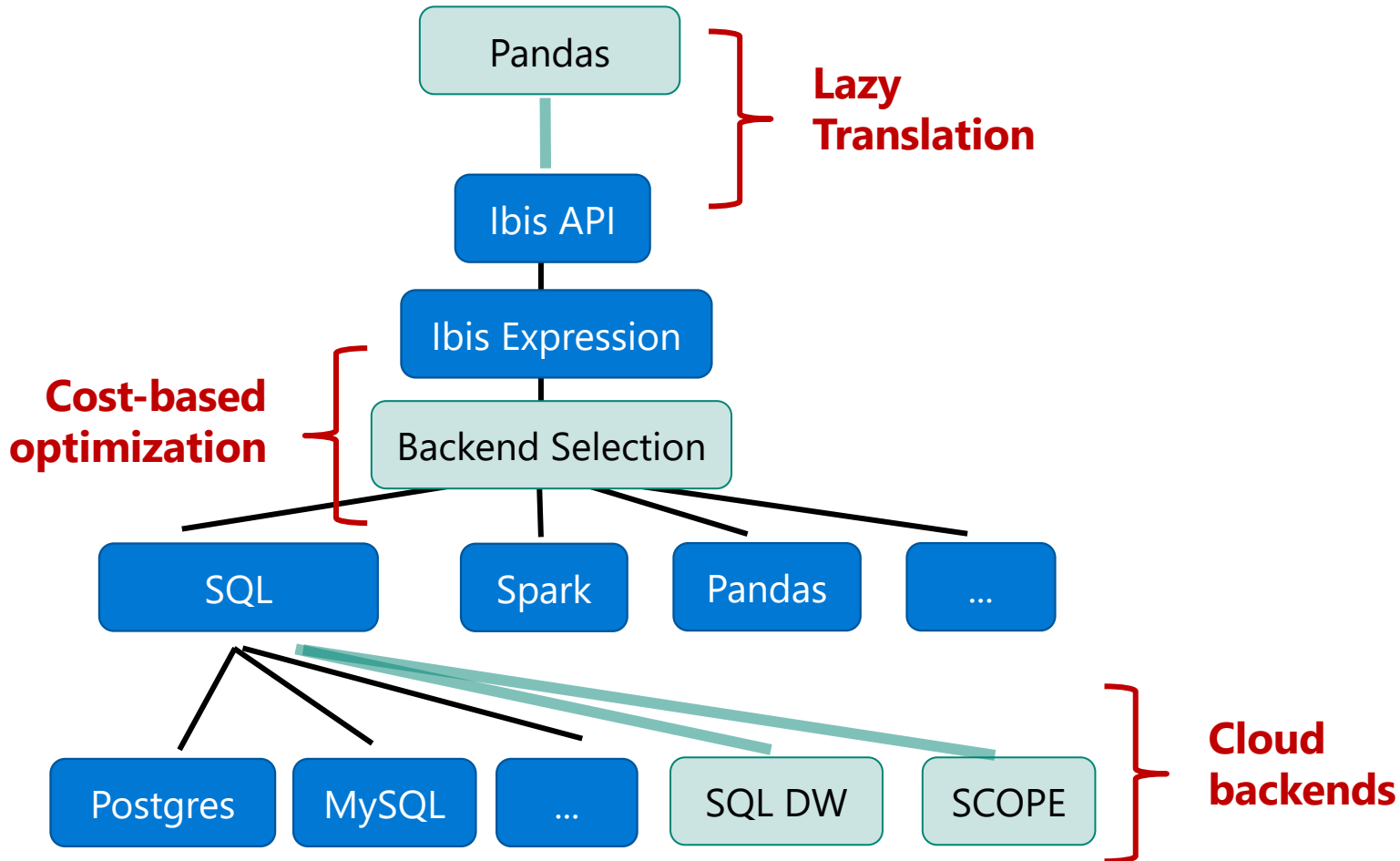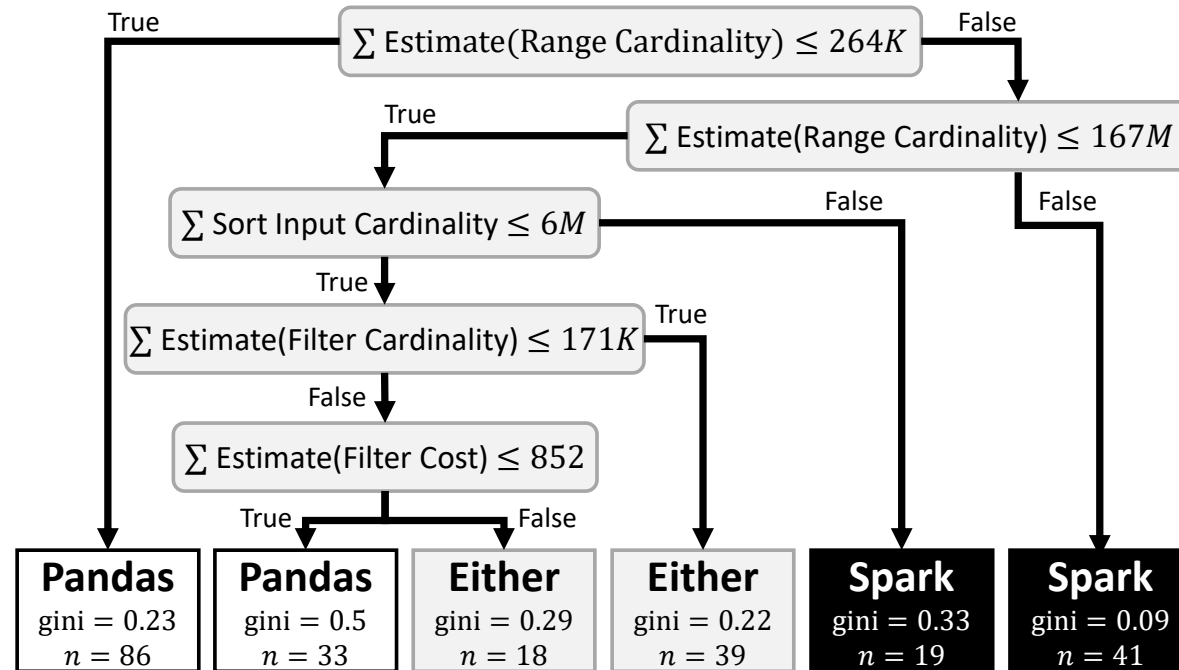
- Leverage past workloads from cloud backends to learn a decision tree
- At compile time:
  - User provides the list of available backends
  - Compile the plan into a common representation
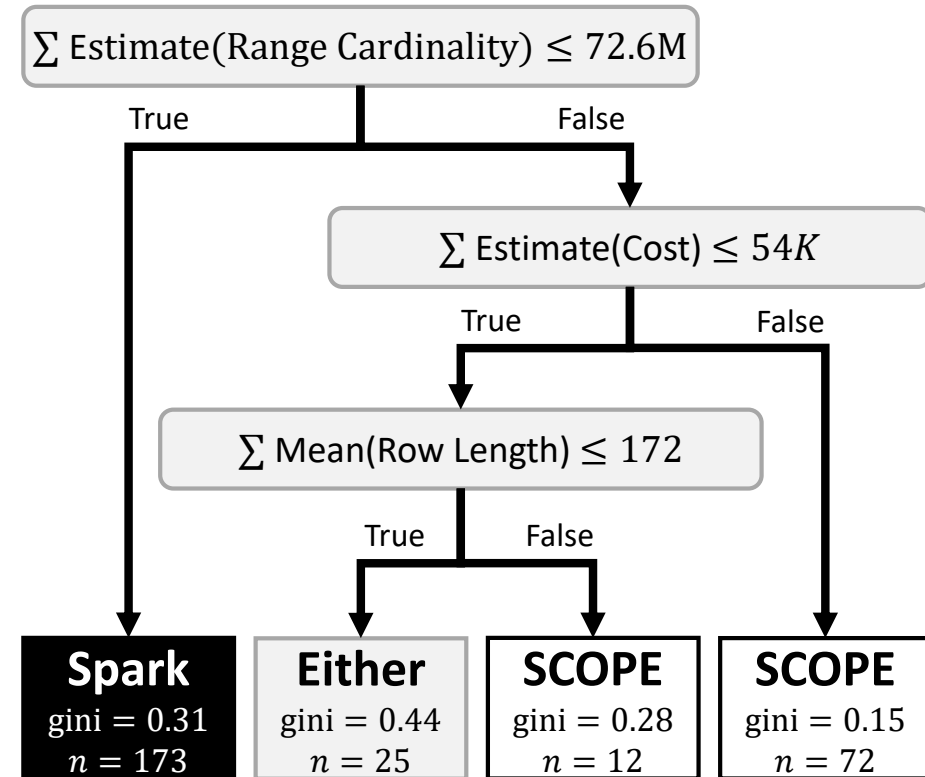  - Infer best backend using the decision tree

# Scenario 1: Pandas vs PySpark

- Question:
  - When to switch to a cluster?
  - Or to local execution?
- Decision tree:
  - 84% accuracy on test set
  - On Pandas:
    - 84% median improvement
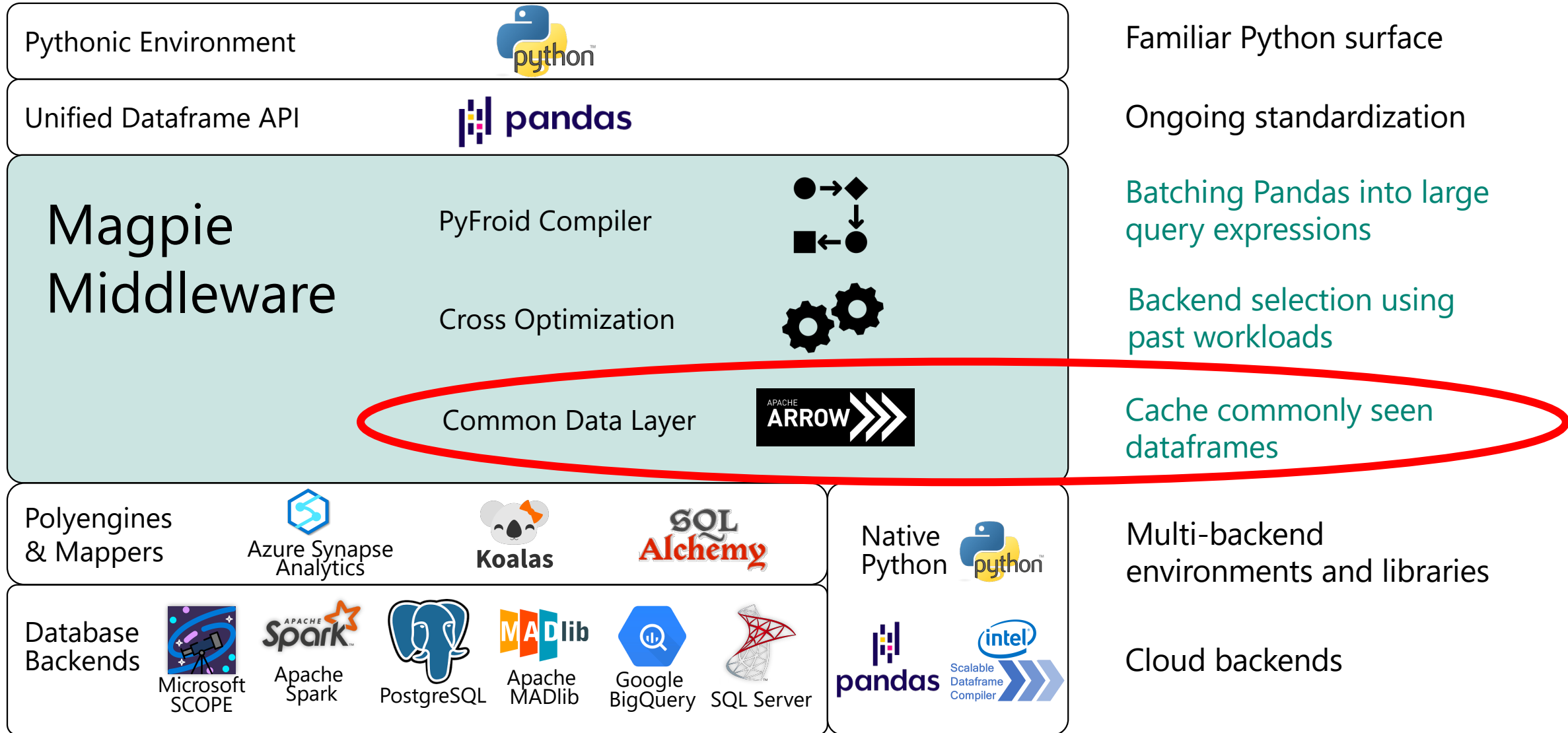    - Up to 99% improvement
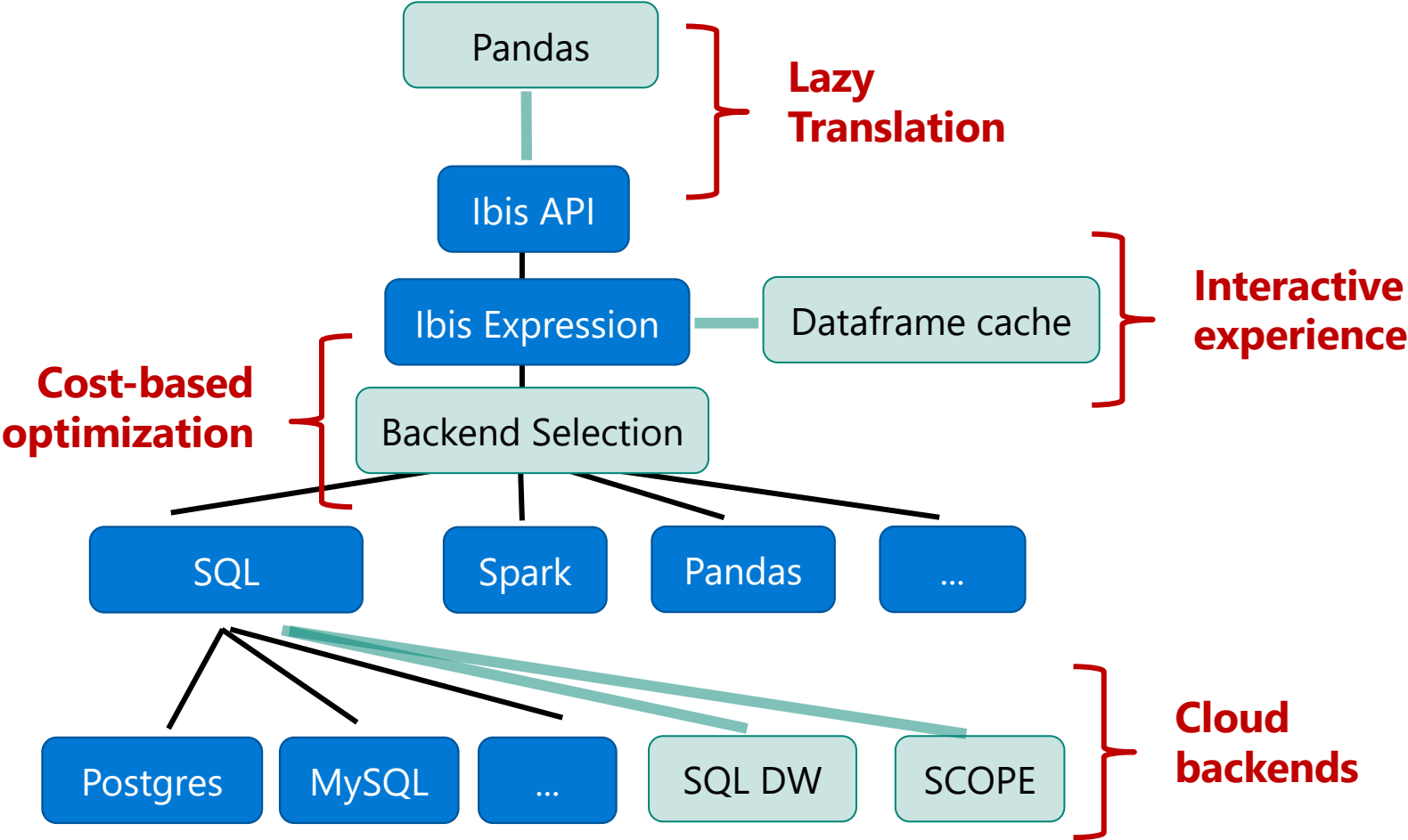
# Scenario 2: PyScope vs PySpark

· Cosmos has both SCOPE and Spark engines now

· Question: which one to use for data science?

· Decision tree
  · 87% accuracy on test set
  · On Spark:
    · Median improvement 85%
    · Up to 98% improvement

# Magpie

| | | |
|---|---|---|
| Pythonic Environment |  | Familiar Python surface |
| Unified Dataframe API | pandas | Ongoing standardization |

| Magpie Middleware | PyFroid Compiler | | Batching Pandas into large query expressions |
|---|---|---|---|
| | Cross Optimization | | Backend selection using past workloads |
| | Common Data Layer | APACHE ARROW | Cache commonly seen dataframes |

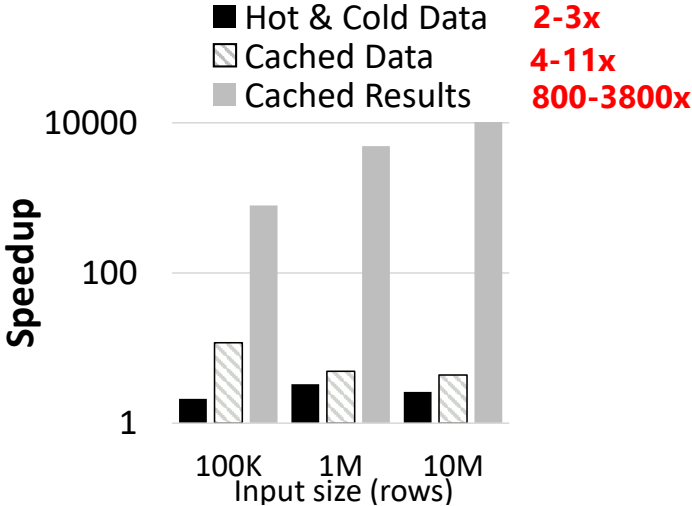| | | | |
|---|---|---|---|
| Polyengines & Mappers | Azure Synapse Analytics    Koalas    SQL Alchemy | Native Python | Multi-backend environments and libraries |
| Database Backends | Microsoft SCOPE    Apache Spark    PostgreSQL    Apache MADlib    Google BigQuery    SQL Server | pandas    intel Scalable Dataframe Compiler | Cloud backends |

# Common Data Layer



Blue parts: already in IBIS, Green parts: our contributions

- Dataframe cache
  - Generate unique signatures
  - Store repeated dataframes in ArrowFlight server
  - Skip accessing the backend in case of cache hit

# Summary

| Pythonic Environment | Lingua franca for many analyses |
|---|---|

| Unified Dataframe API | Increasingly getting standardized |
|---|---|

**Magpie Middleware**

**Pandas Without Regret!**
**Write once, execute anywhere**
**Abstracting Data Processing Complexity**

| Polyengines & Mappers | From polystores to polyengines |
|---|---|

| Database Backends | Hyperscale performance<br>Data already in the cloud |
|---|---|

Native Python

Gray Systems Lab

https://azuredata.microsoft.com/

**Hiring Summer Interns!**