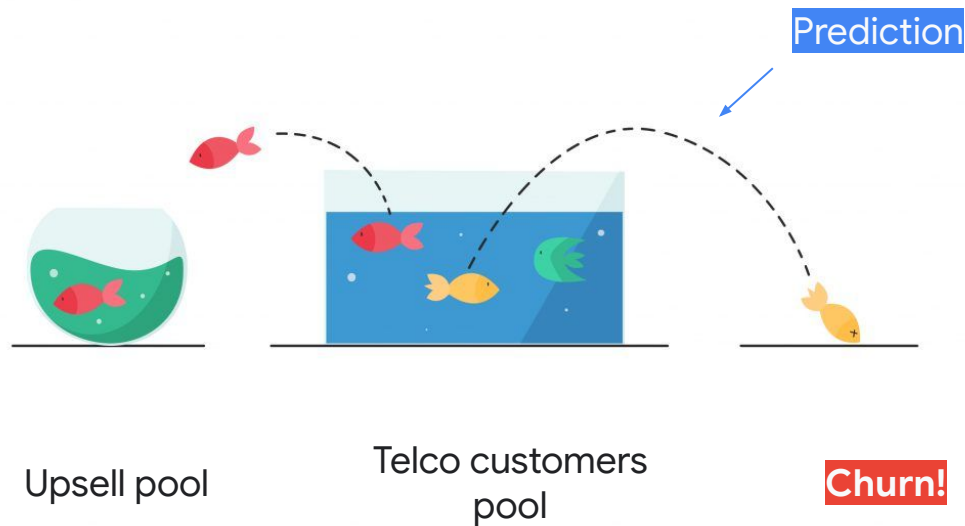


The **job** to be done: A **churn prediction** model



The **job** to be done: Build a data product

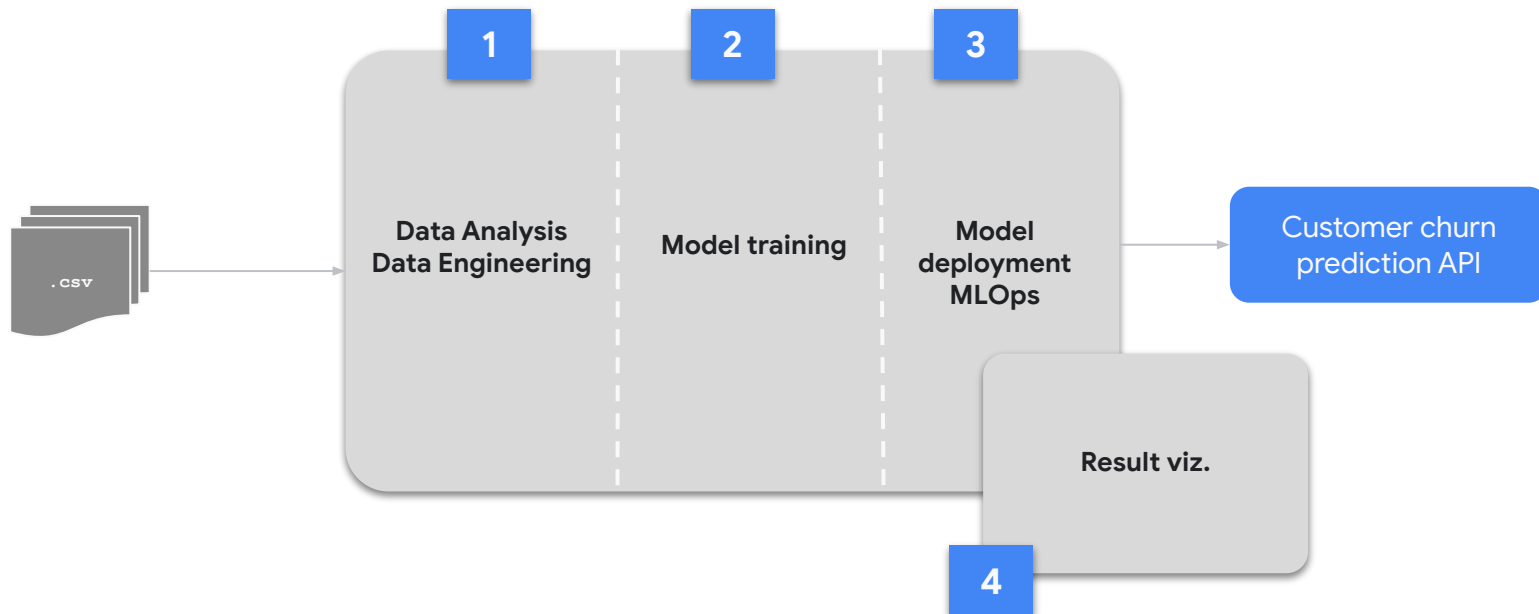
High value telco
customer list variables
and **churn tag (Y/N)**



Data system

Customer churn
prediction API

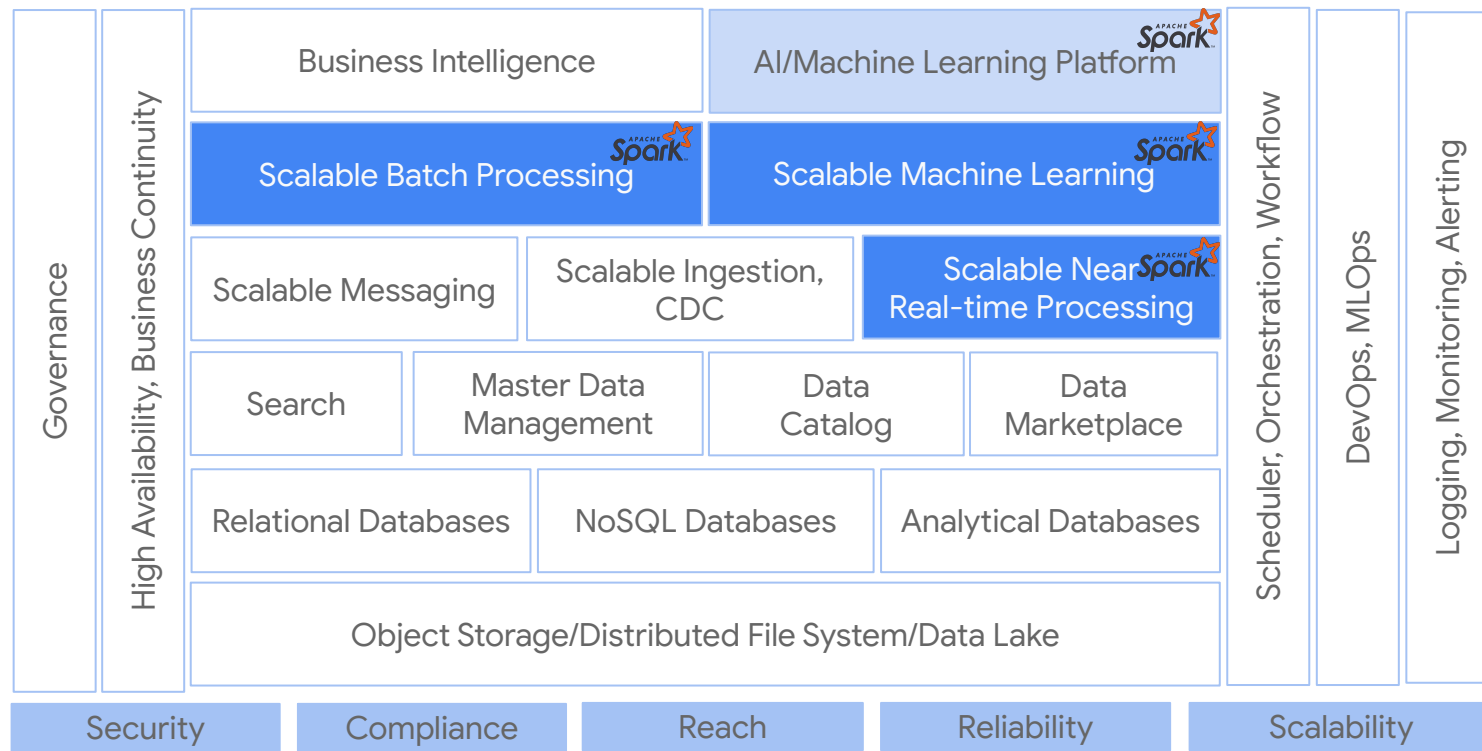
The **job** to be done: Product dev stages



The toolkit - SPARK



The toolkit



The personas



Data Analyst
Business Analyst



Data Engineer
ML Engineer



Data Scientist

The trigger - You've got email

[URGENT] Can you help me with this request?

External

Inbox x



John Coalesce - Director of Data

to me ▾

9:54 AM (11 hours ago)



Hi,

As you know, one of our key OKRs this Q is to build a churn prediction model. After discussing with the ML team, we have been able to build a data extraction pipeline with key customer information, however data seems to be very messy and doesn't meet the quality requirements for ML training.

We need to build a data pipeline to sanitize the records, please, work on this as P0!
Remember that our SPARK production cluster has some heavy workloads these days during quarter closure.

Didn't you update me about this new spark serverless thing from the Google folks? Give it a try and please do not spend too much resources (\$)

We count on you!

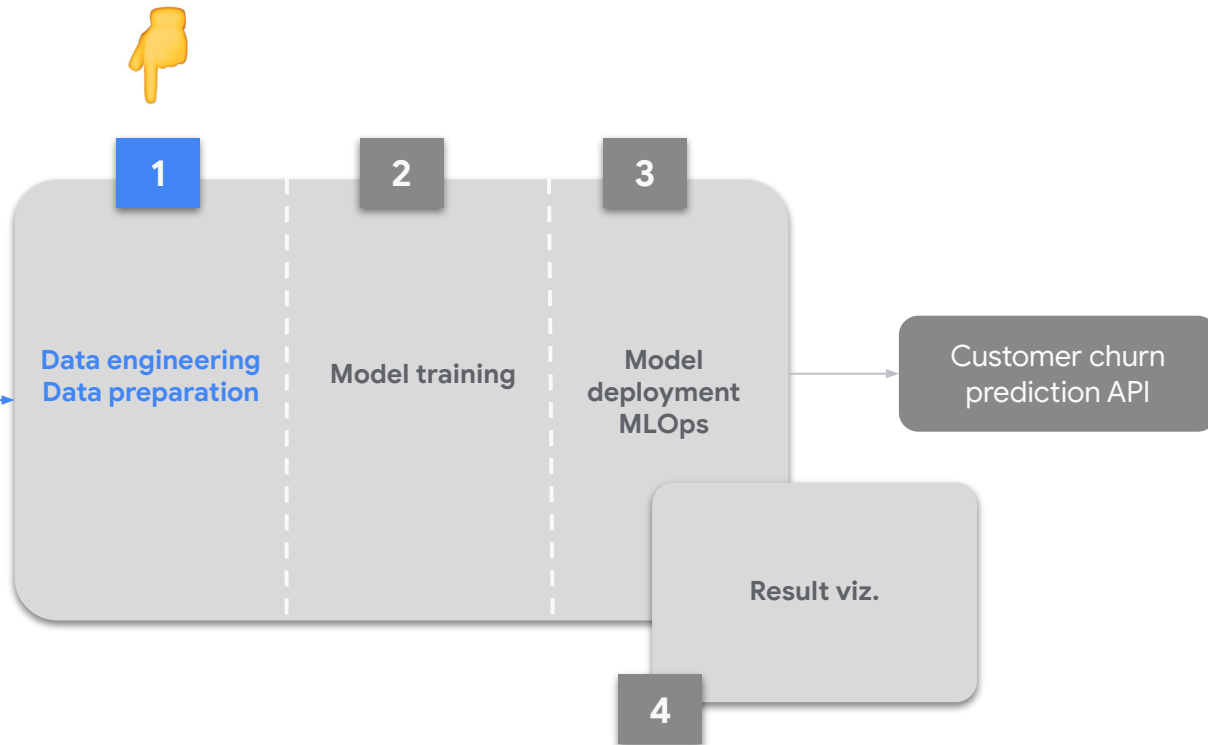
[Message clipped] [View entire message](#)



Where's the data?

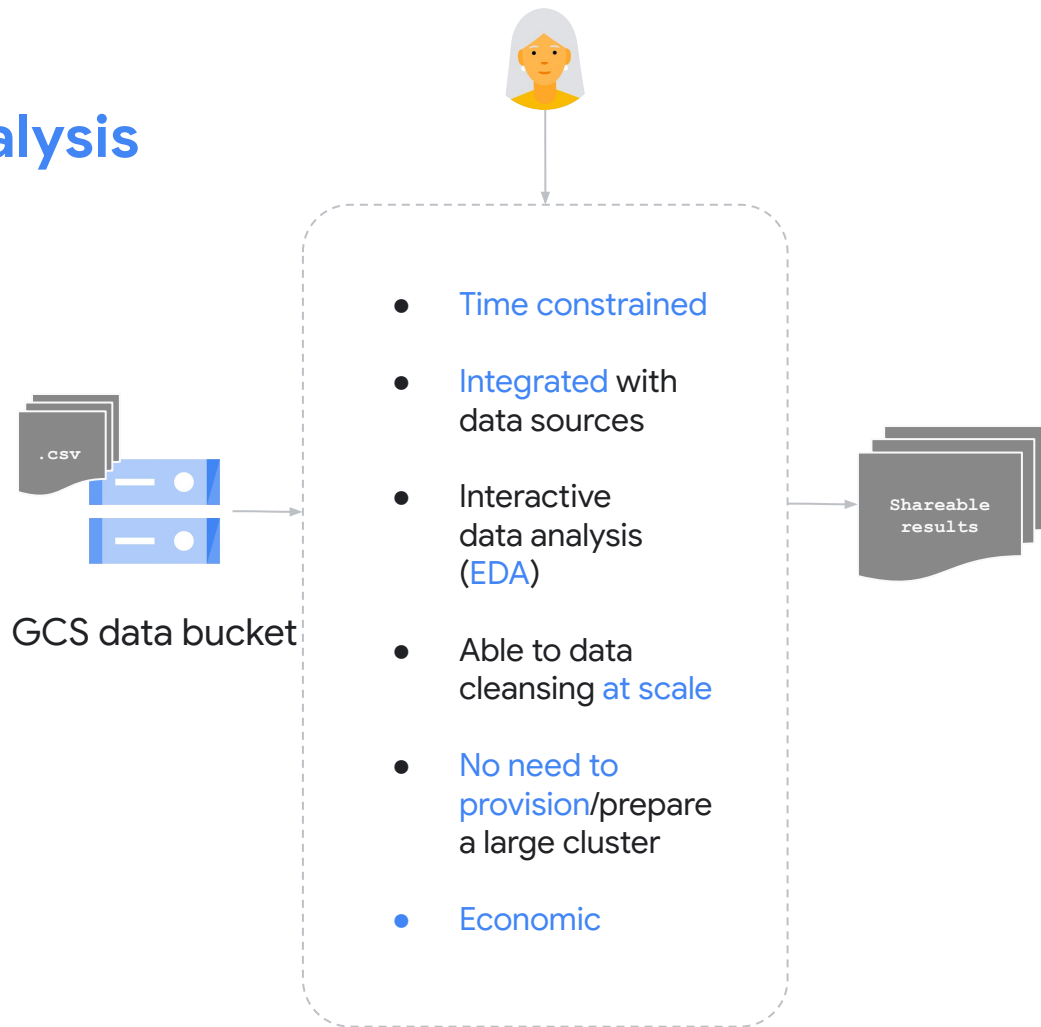
Left it on a bucket
Good luck!

Telco customer list with
churn information



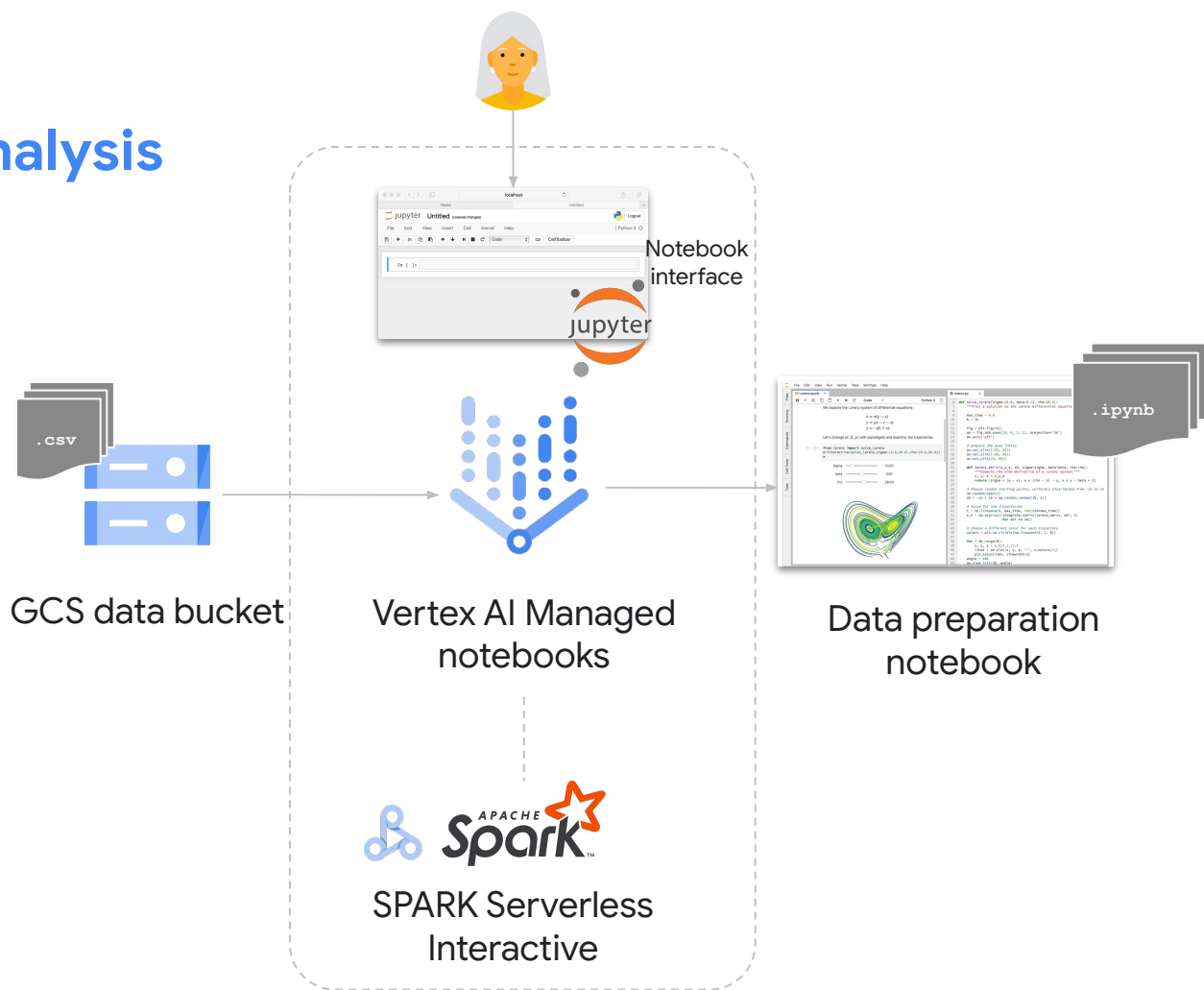
1

Data Analysis



1

Data Analysis



1

Data Engineering



Data preparation notebook



RAW data

- Repeatable process
- At scale
- Enterprise level: Monitored, secured, ..
- Portable
- No need to provision/prepare a large cluster
- Economic

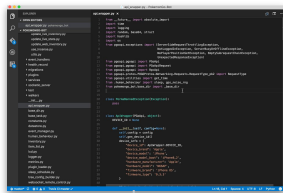


Prepared data

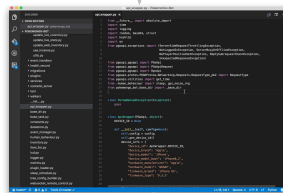
1

Data Engineering

IDE interface



IDE interface



Production
Code
refactor

Pipeline
code

```
def parse_arguments():
    """Read arguments from a command line."""
    parser = argparse.ArgumentParser(
        formatter_class=argparse.ArgumentDefaultsHelpFormatter,
        description="Script to parse arguments and return them as a dictionary."
    )
    parser.add_argument(
        "-v", "--verbose",
        type=int,
        help="Verbosity of logging: 0 - critical, 1 - warning, 2 - info, 3 - debug"
    )
    return parser.parse_args()

def main():
    args = parse_arguments()
    logging.basicConfig(
        level=logging.DEBUG if args.verbose > 2 else logging.INFO,
        format='%(asctime)s - %(name)s - %(levelname)s: %(message)s'
    )
    logging.info("Starting the script")
    # Your code here
    return 0

if __name__ == '__main__':
    main()
```

data_preparation.py

Data preparation
notebook

Testing
Error handling
Argument Parsing
Performance
Security
Libraries
...



Composer
dataproc
serverless
operator



dbt python
models



BigQuery Spark
stored
procedure

RAW data

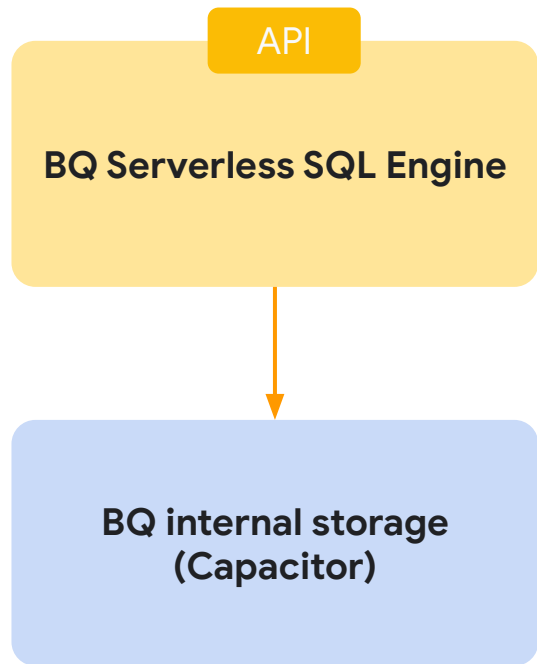


Serverless
Batch

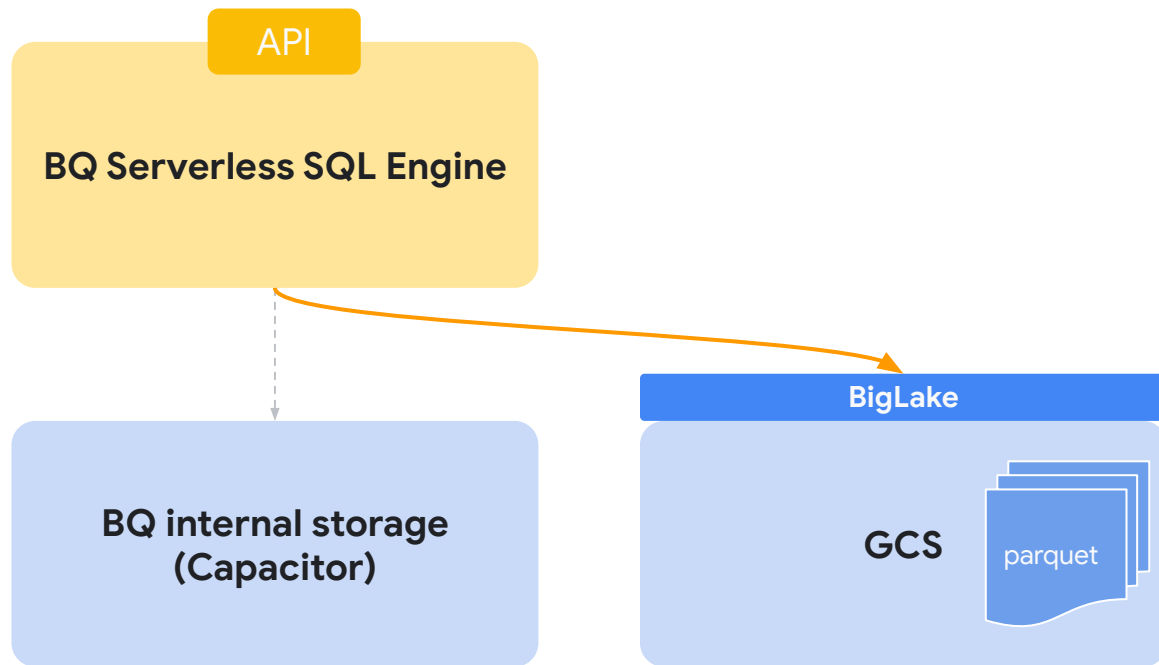


Prepared data

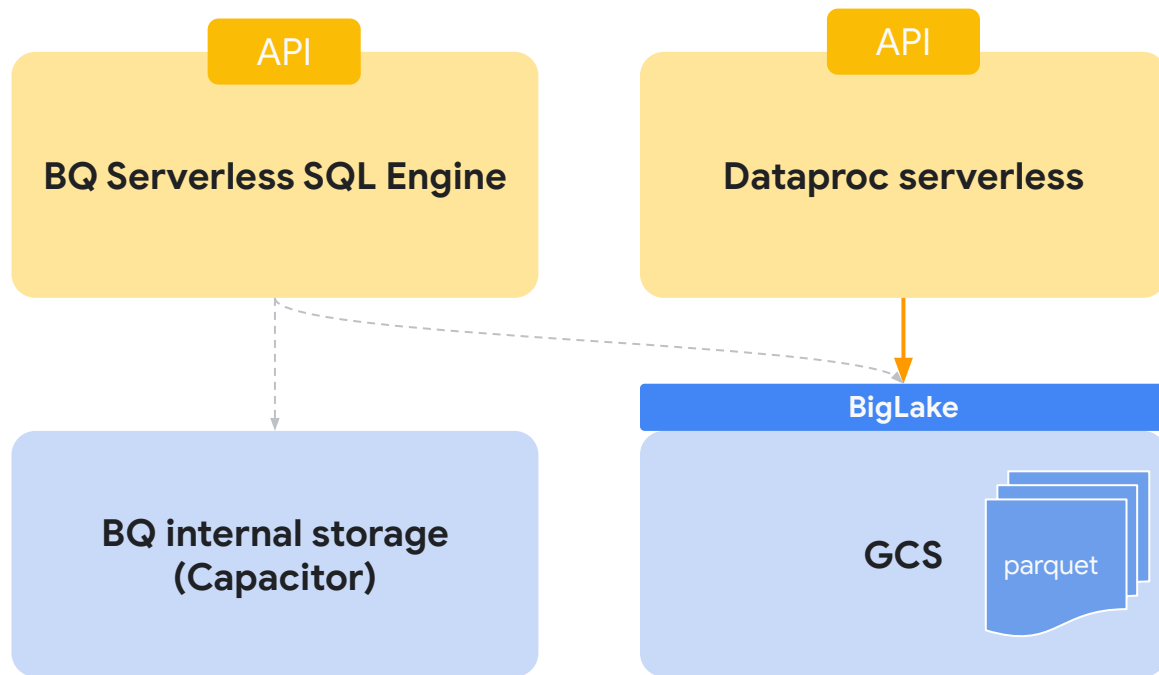
{Interlude} Interoperability



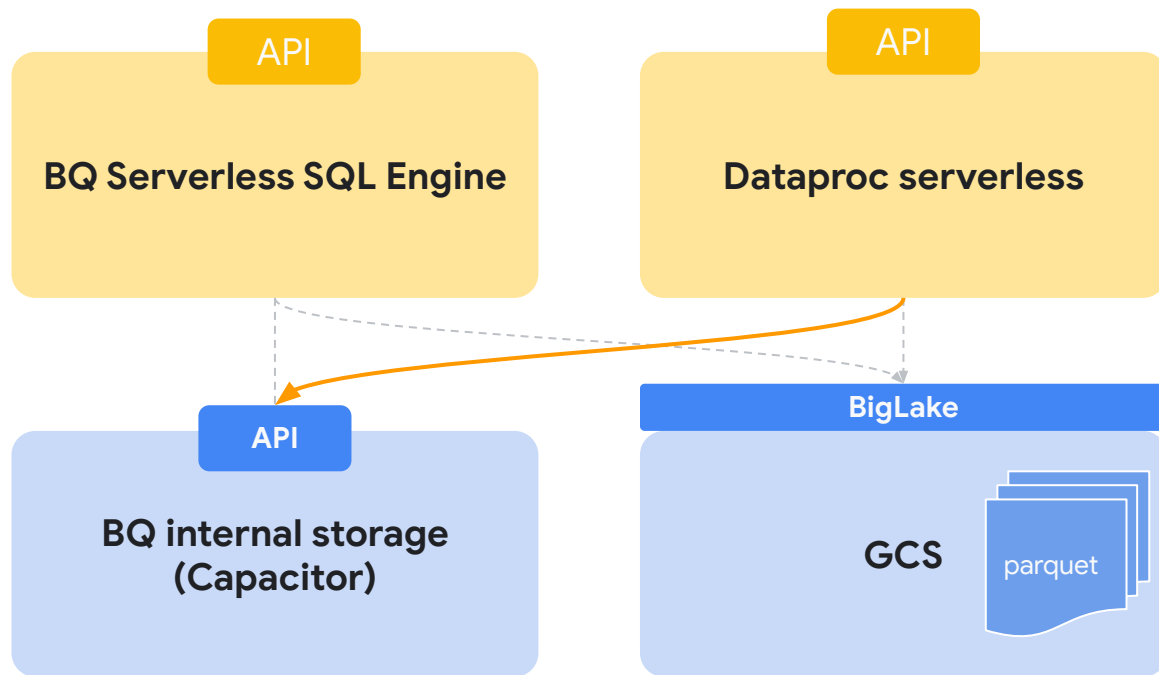
{Interlude} Interoperability



{Interlude} Interoperability

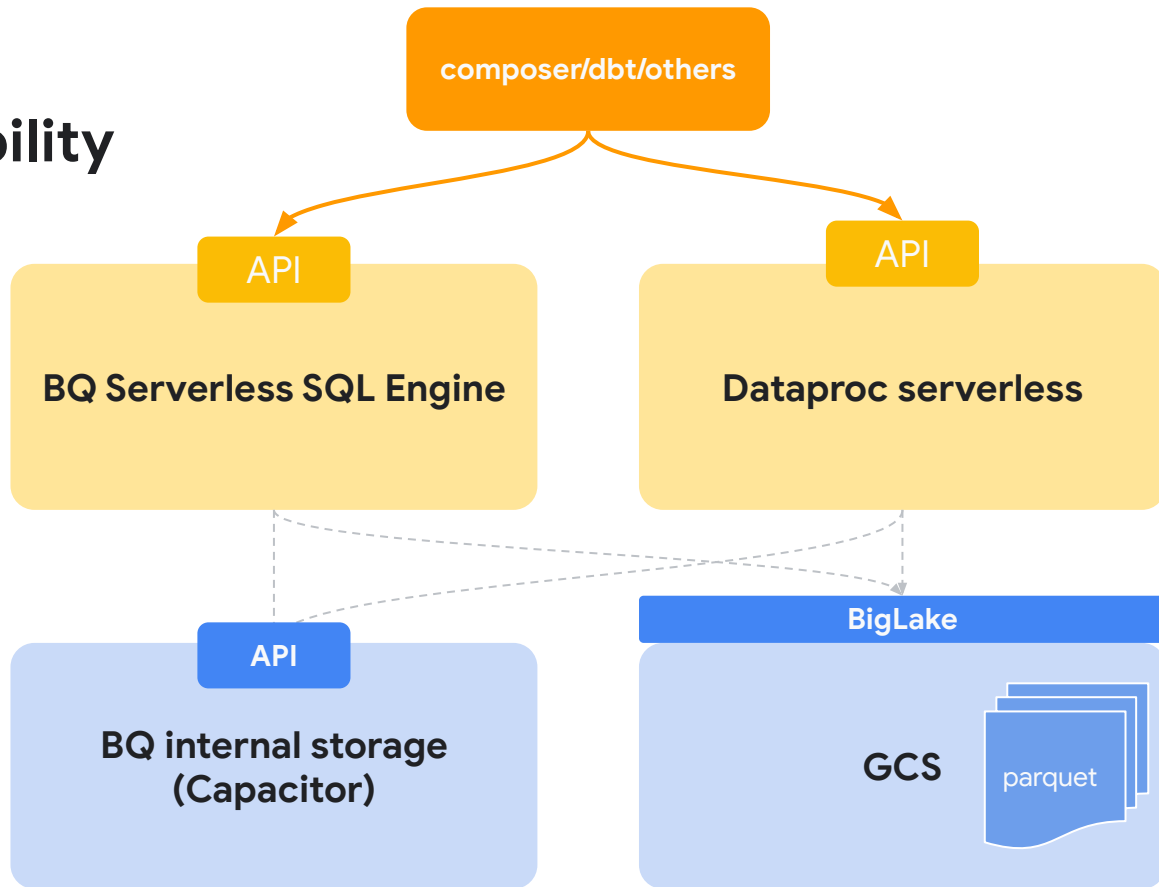


{Interlude} Interoperability

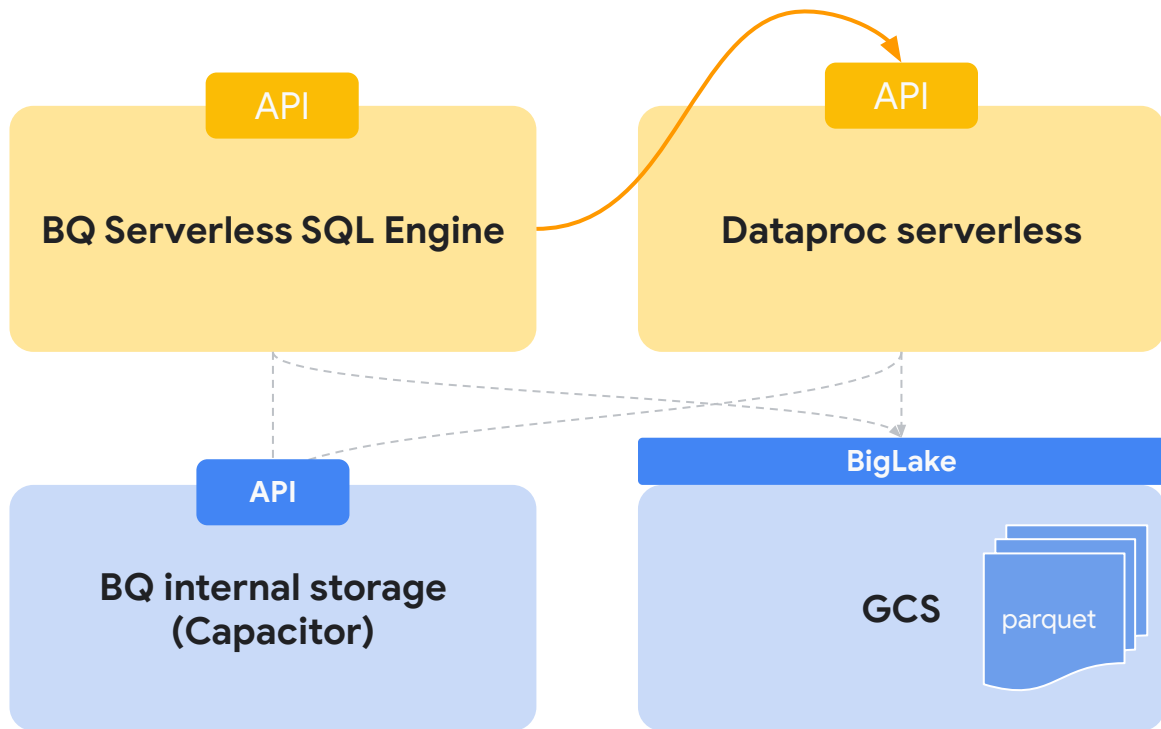


{Interlude}

Interoperability



{Interlude} Interoperability



The trigger

[URGENT] Model churn dataset is ready to go!

External

Inbox x



John Overfitted - Director of Data Science

9:54 AM (11 hours ago)



to me ▾

Hi,

Surprisingly enough, the data engineering team just update us: the churn prediction model dataset seems to be ready to go. It's up to our team now to create a proper classification model, but remember, we are working against the clock here!

The data engineering team seems to have use some new magic from Google Cloud that really accelerated their time-to-market, spark serverless I believe its called.

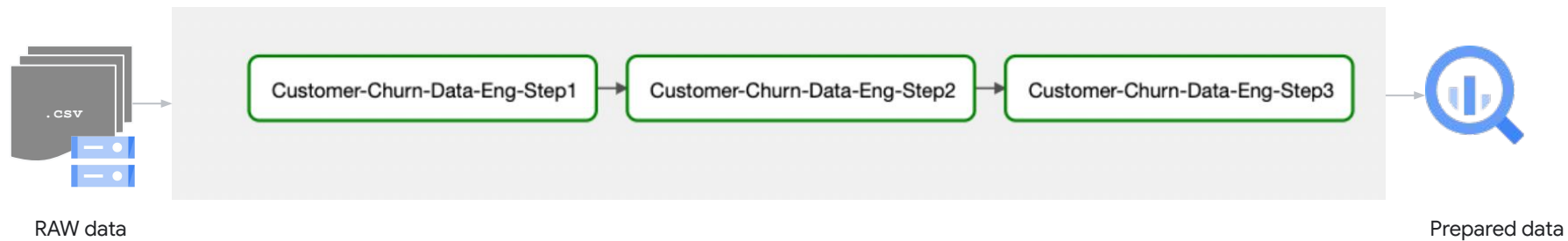
Would that fit our use case as well?

We count on you!

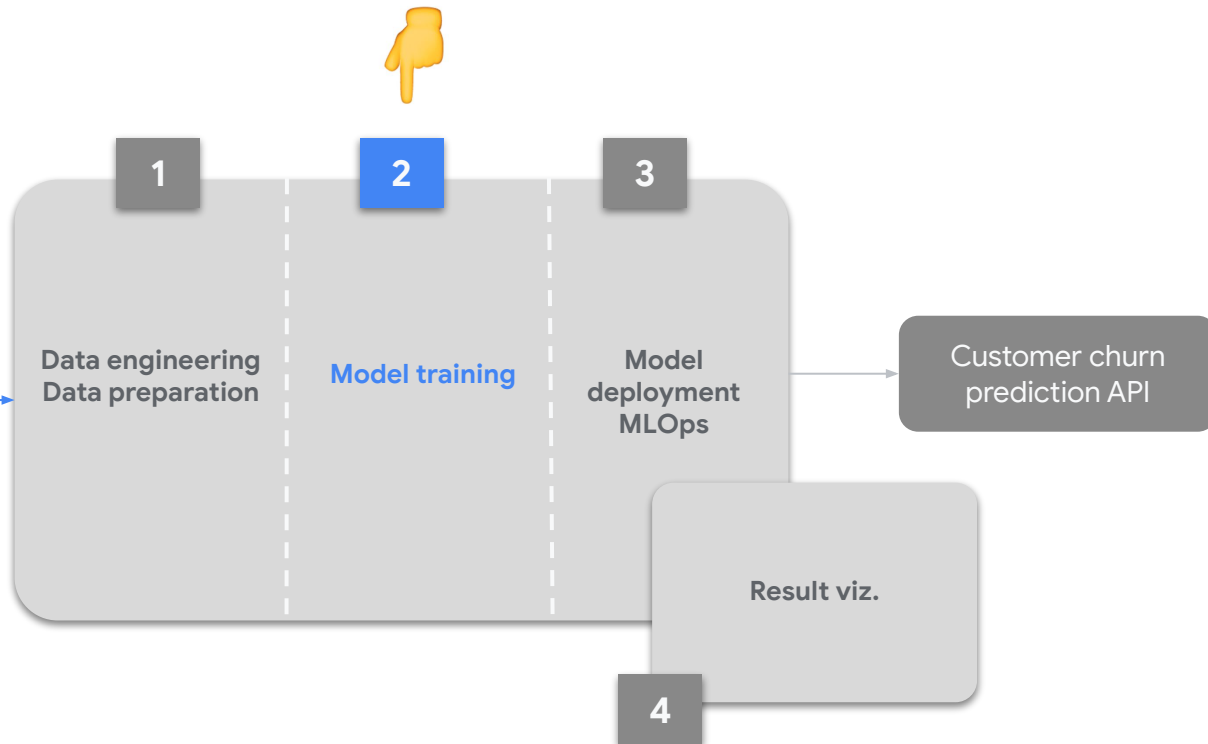
P.S: Have a look at the consumption, we are above the monthly budget already

[Message clipped] [View entire message](#)

Whats next?

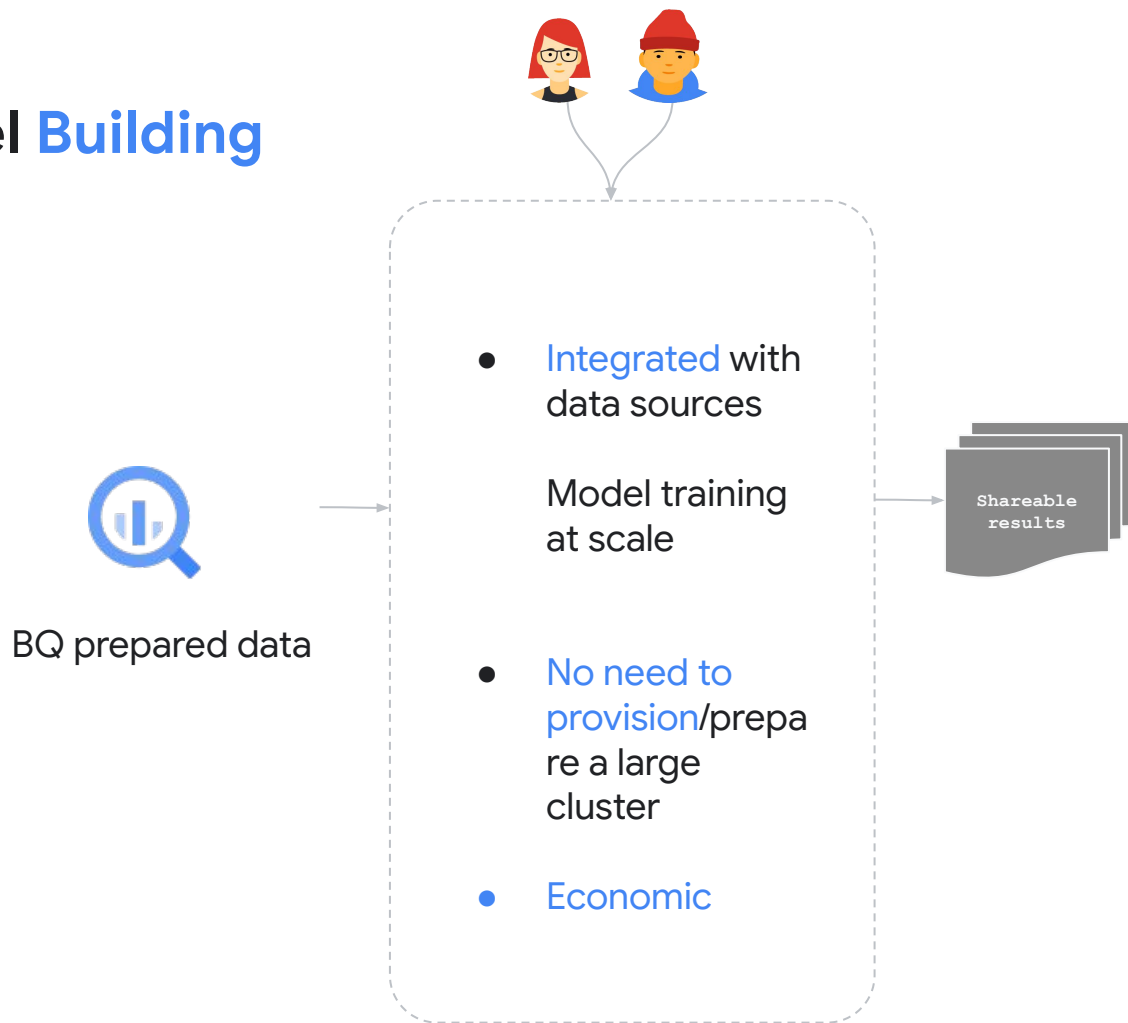


Telco customer list with
churn information



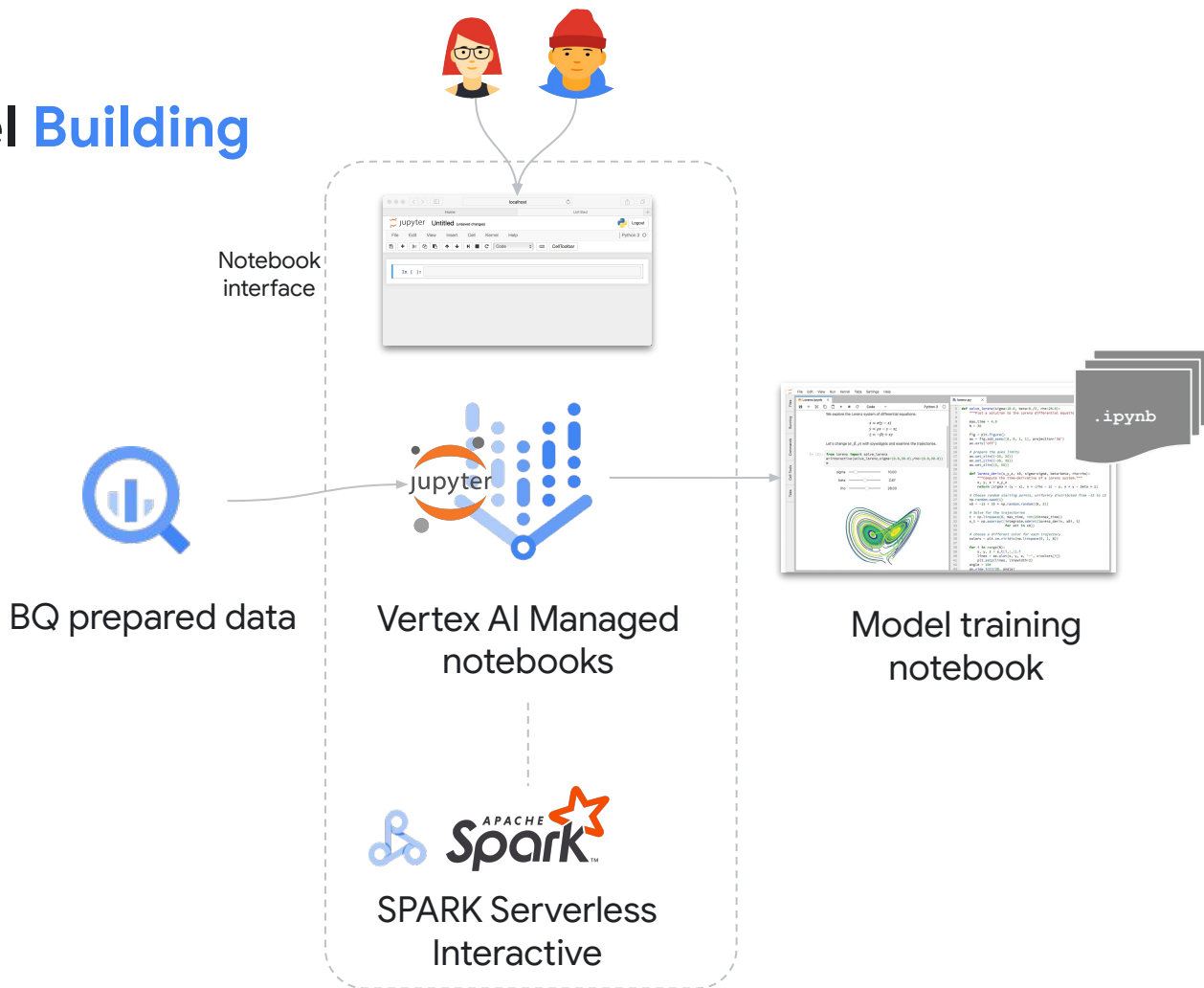
2

Model Building



2

Model Building



The trigger

[URGENT] Model churn to prod!

External



Inbox x



Maria DriftSkew - Director of ML Engineering

to me ▾

9:54 AM (11 hours ago)



Hi,

Have you heard the news? The DS team has just released their churn prediction model! No delays this time!

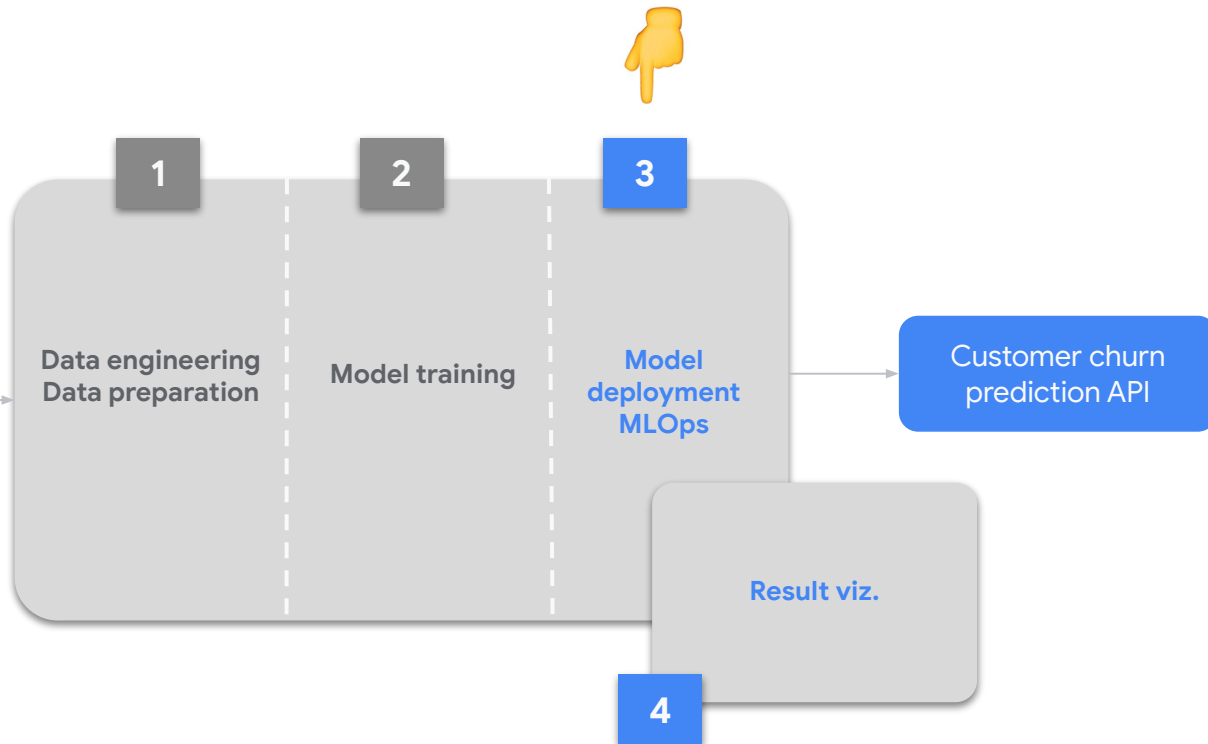
So, it's our turn now, we need to build the a ML pipeline to ensure the model is up to date!

Honestly, never deployed a SPARK MLlib model before, good luck!

P.S: Feel free to use another pipeline tool other than Composer

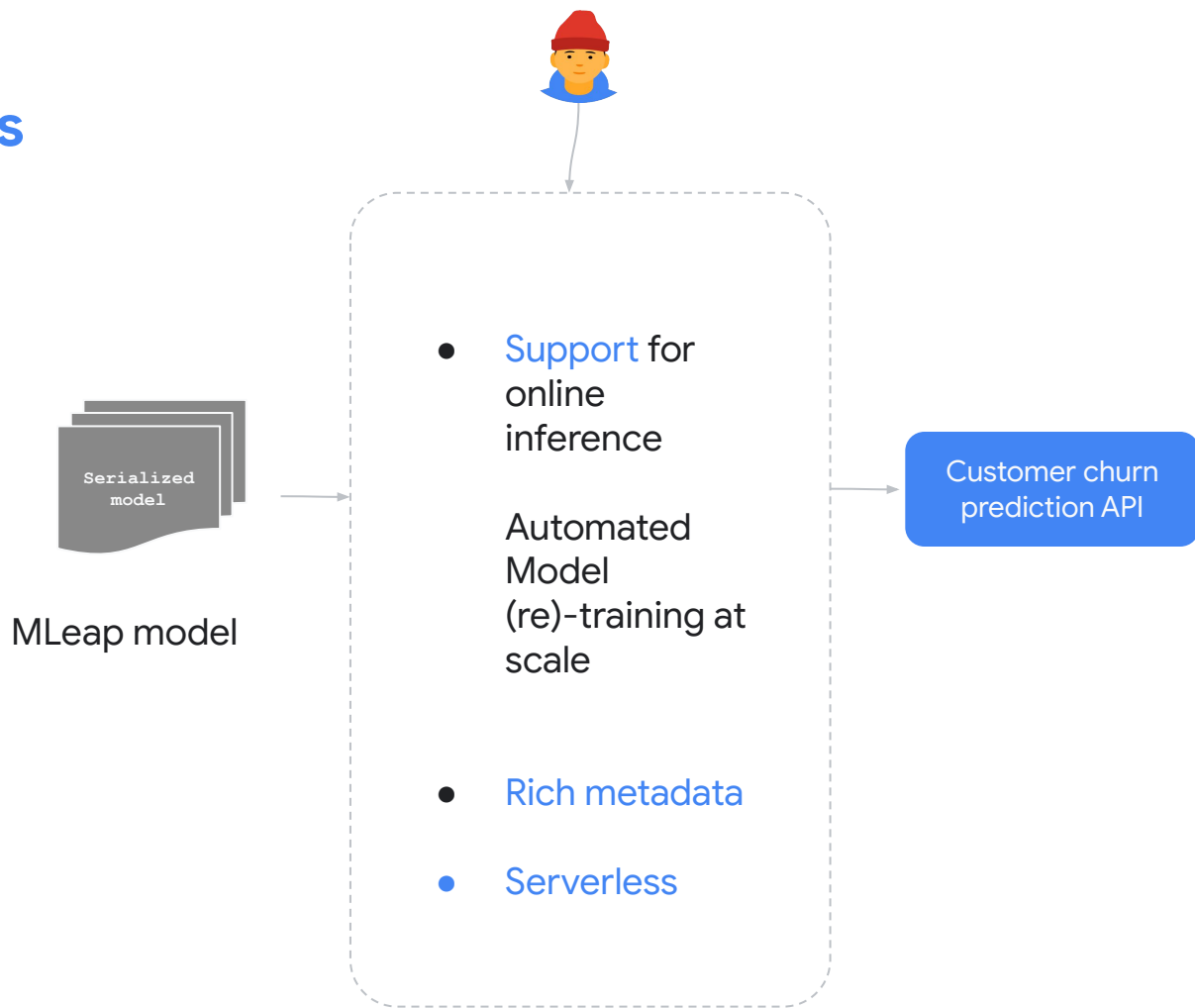
[Message clipped] [View entire message](#)

Telco customer list with
churn information



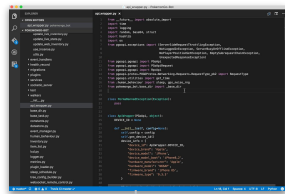
3

ML Ops

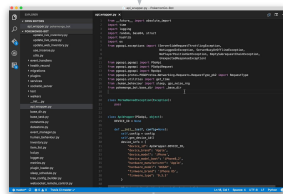


3 ML Ops

IDE interface



IDE interface



Production
Code
refactor

Testing
Model serialization
Model deployment
Error handling
Argument Parsing
Performance
Security
Libraries
...

ML Pipeline
code

```
def parse_arguments():  
    """Read arguments from a command line"""  
    parser = argparse.ArgumentParser(description='ML Pipeline')  
    parser.add_argument('-v', '--verbose', type=int, default=0,  
                        help='Verbosity of logging: 0 - critical, 1 - warning, 2 - info, 3 - debug')  
  
    args = parser.parse_args()  
    verbose = 0 if logging.CRITICAL, 1: logging.WARNING, 2: logging.INFO, 3: logging.DEBUG  
    logging.basicConfig(format='%(asctime)s: %(message)s', level=verbose, datefmt='%Y-%m-%d %H:%M:%S')  
  
    return args  
  
def main():  
    pass  
  
if __name__ == '__main__':  
    args = parse_arguments()  
    main()
```

model_trainer.py



Vertex AI
Pipelines

Prepared data

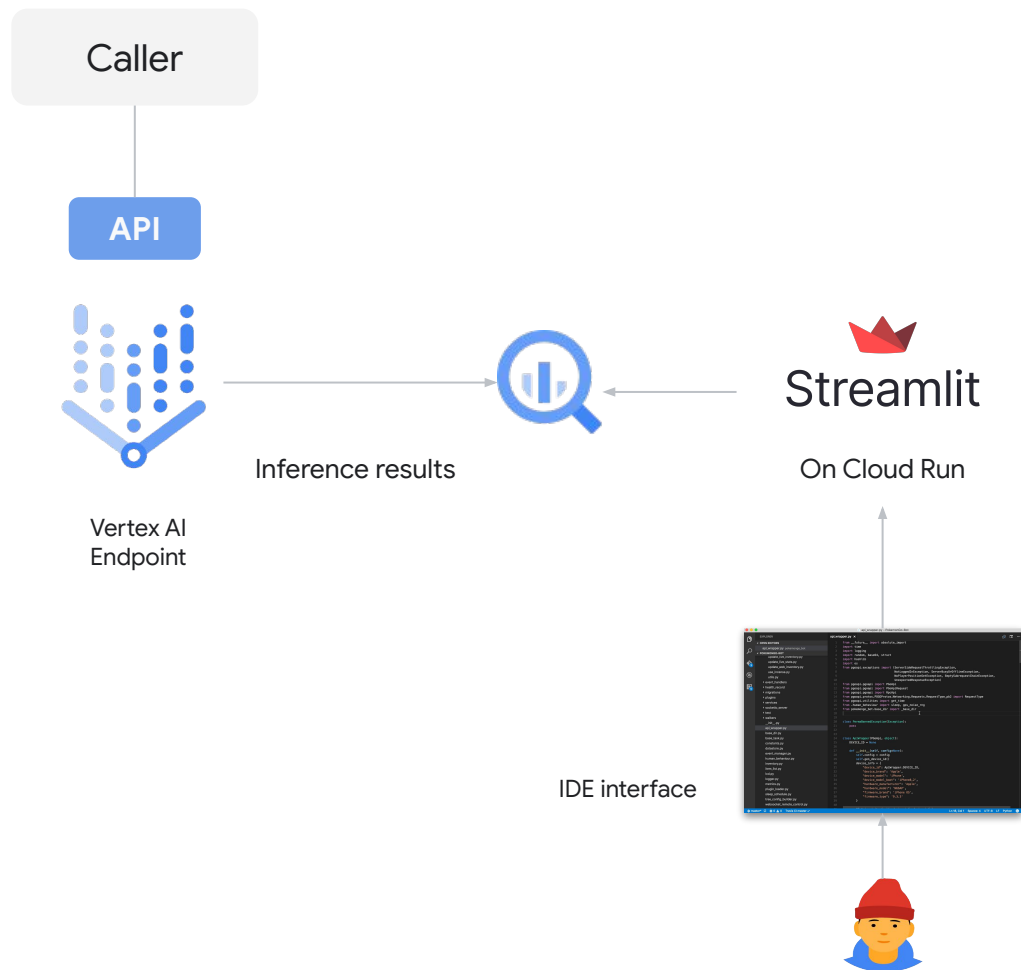


Serverless
Batch

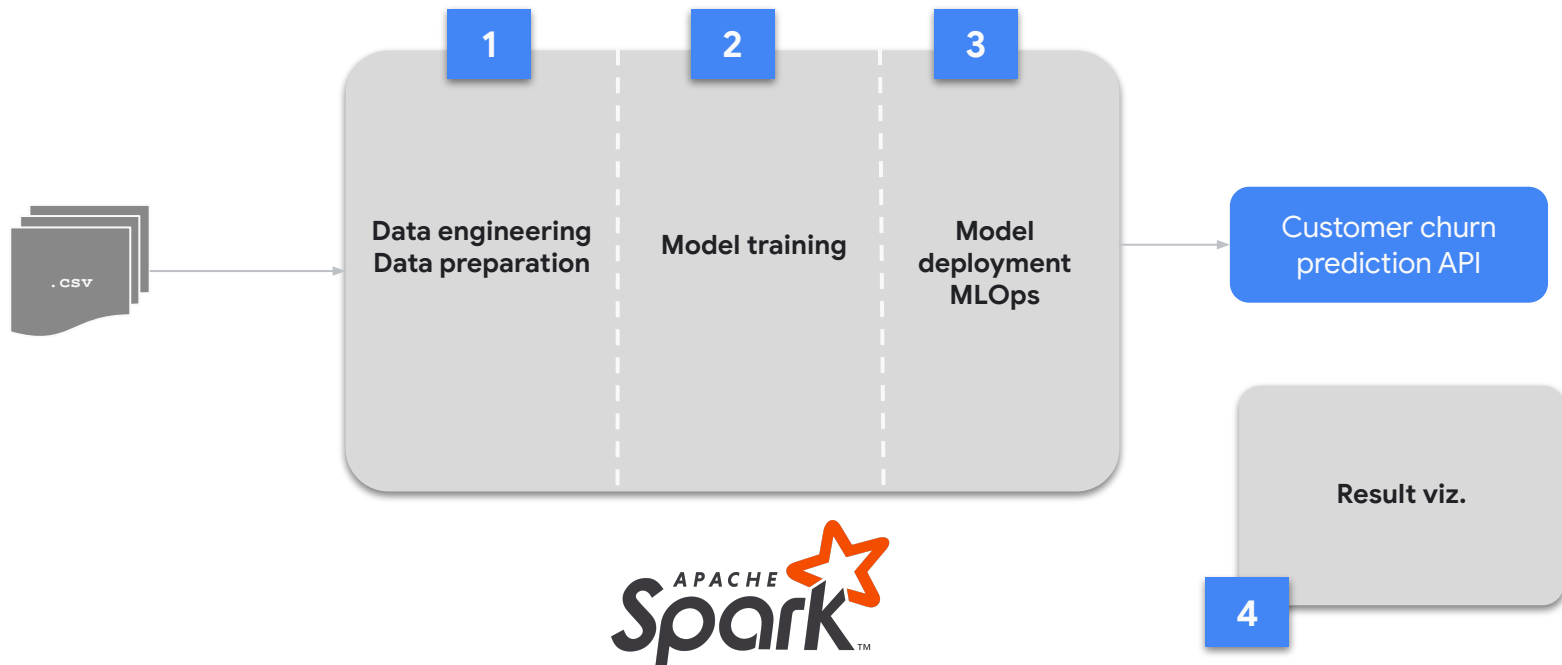
Model
endpoint
REST

4

Data Viz

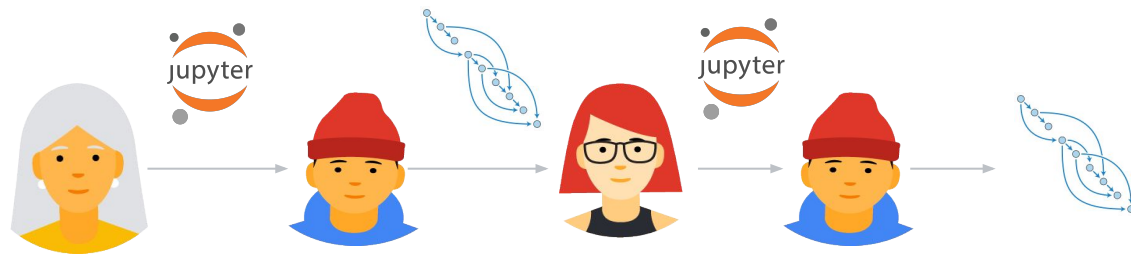


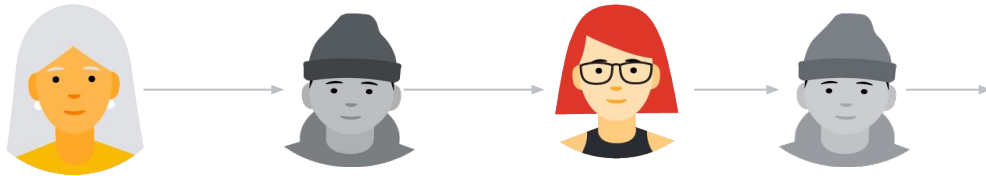
All in all



SPARK - A first class citizen in Google Cloud









Using decades
of software
engineering
best practices



Scheduling
Jupyter
notebooks
in production