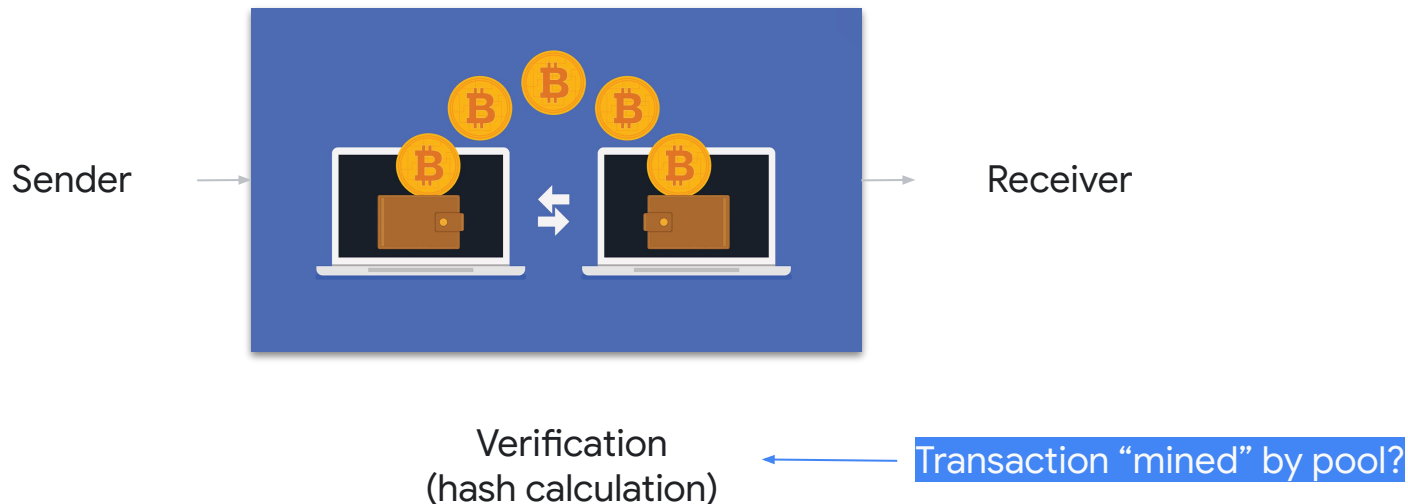


## The **job** to be done: A **mining pool classifier** model



# The **job** to be done: A **mining pool classifier** model

## Address classification

Blockchain transaction history can be aggregated by address and used to analyze user behavior. To motivate further exploration, we present a simple classifier that can detect Bitcoin [mining pools](#). As a brief historical note, mining pools were created when the difficulty of mining Bitcoin reached such a level that rewards could be expected only once every few years. Miners began to pool their resources to earn a smaller share of rewards more consistently and in proportion to their contribution to the pool in which they were mining.

First, we constructed 26 feature vectors to characterize incoming and outgoing transaction flows to each address. Next, we trained the model using labels derived from transaction signatures. Many large mining pools identify themselves in the signature of blocks' [Coinbase](#) transactions. Parsing these signatures, we labelled 10,000 addresses as belonging to known mining pools. One million other addresses were included in the dataset as “non-miners.” The query used to generate our features and labels can be seen [here](#), and the source code for this analysis can be found in a Kaggle notebook [here](#).

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

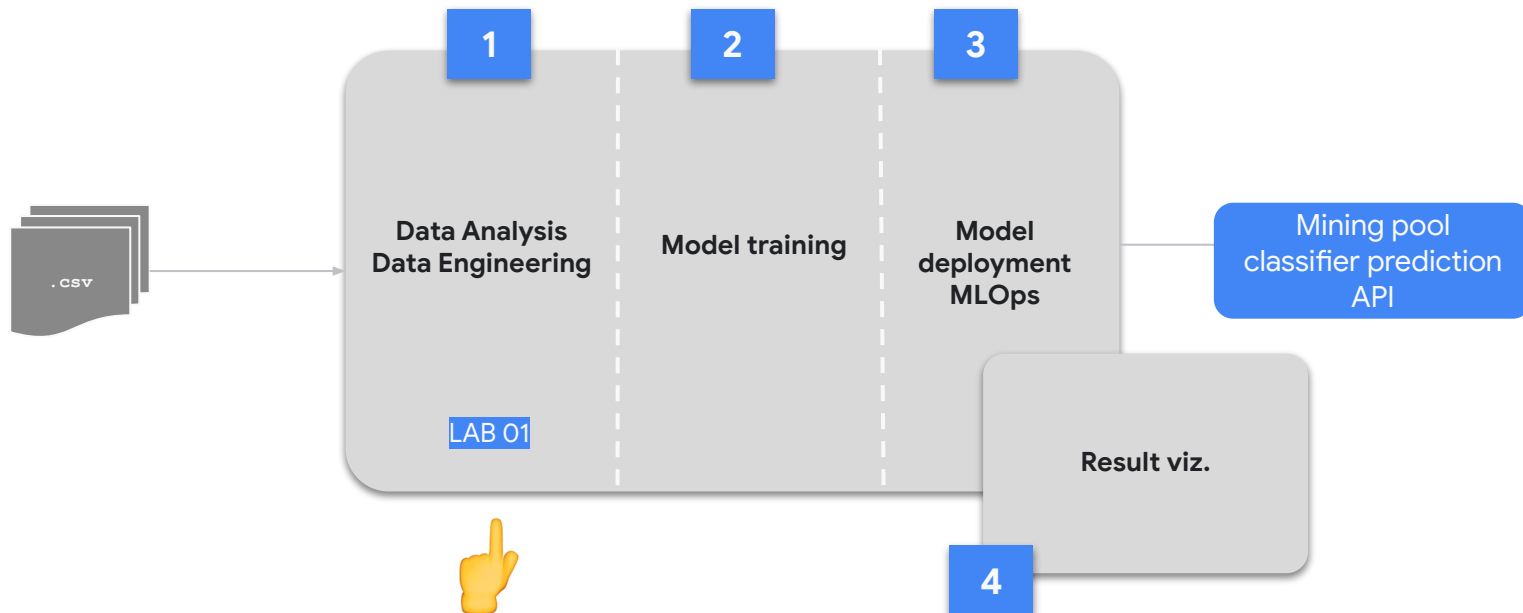
Blockchain transaction



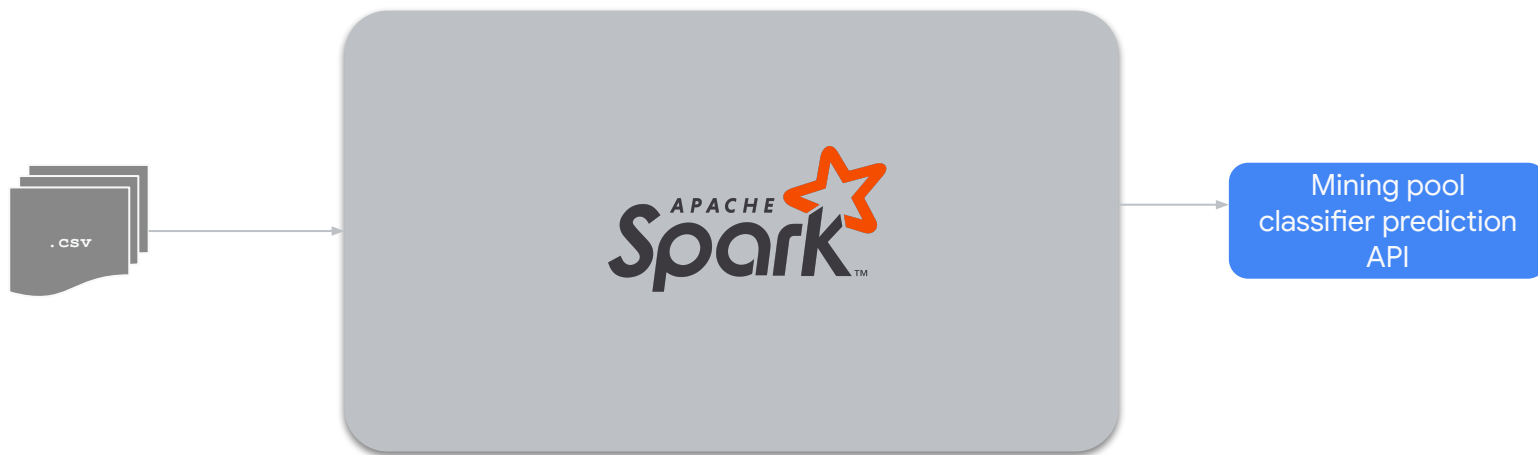
**Data system**

Mining pool  
classifier prediction  
API

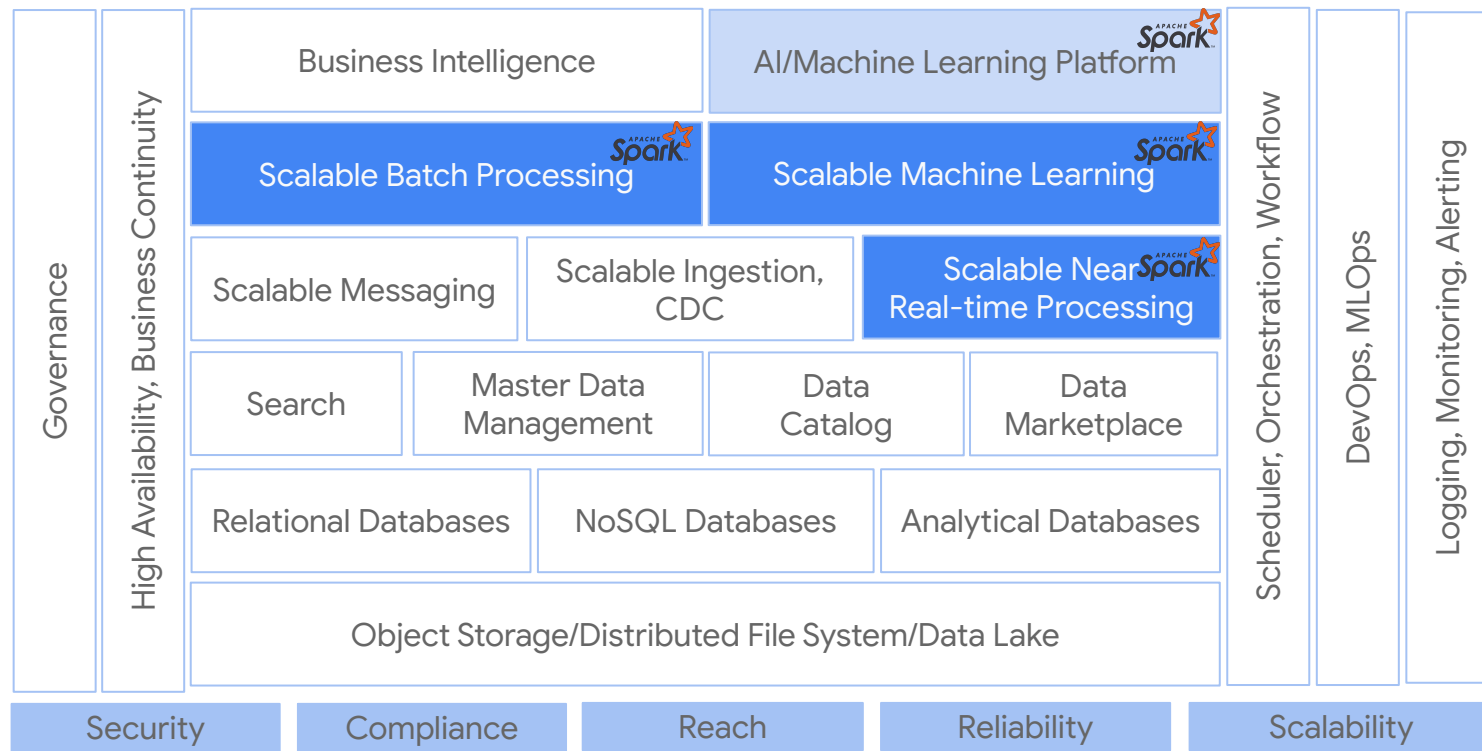
# The **job** to be done: Product dev stages



# 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

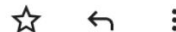
Inbox x



**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 mining pool classifier. After discussing with the ML team, we have identified a valuable public dataset with bitcoin transaction information, however data is not ready for ML training.

We need to build a data pipeline to transform 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!

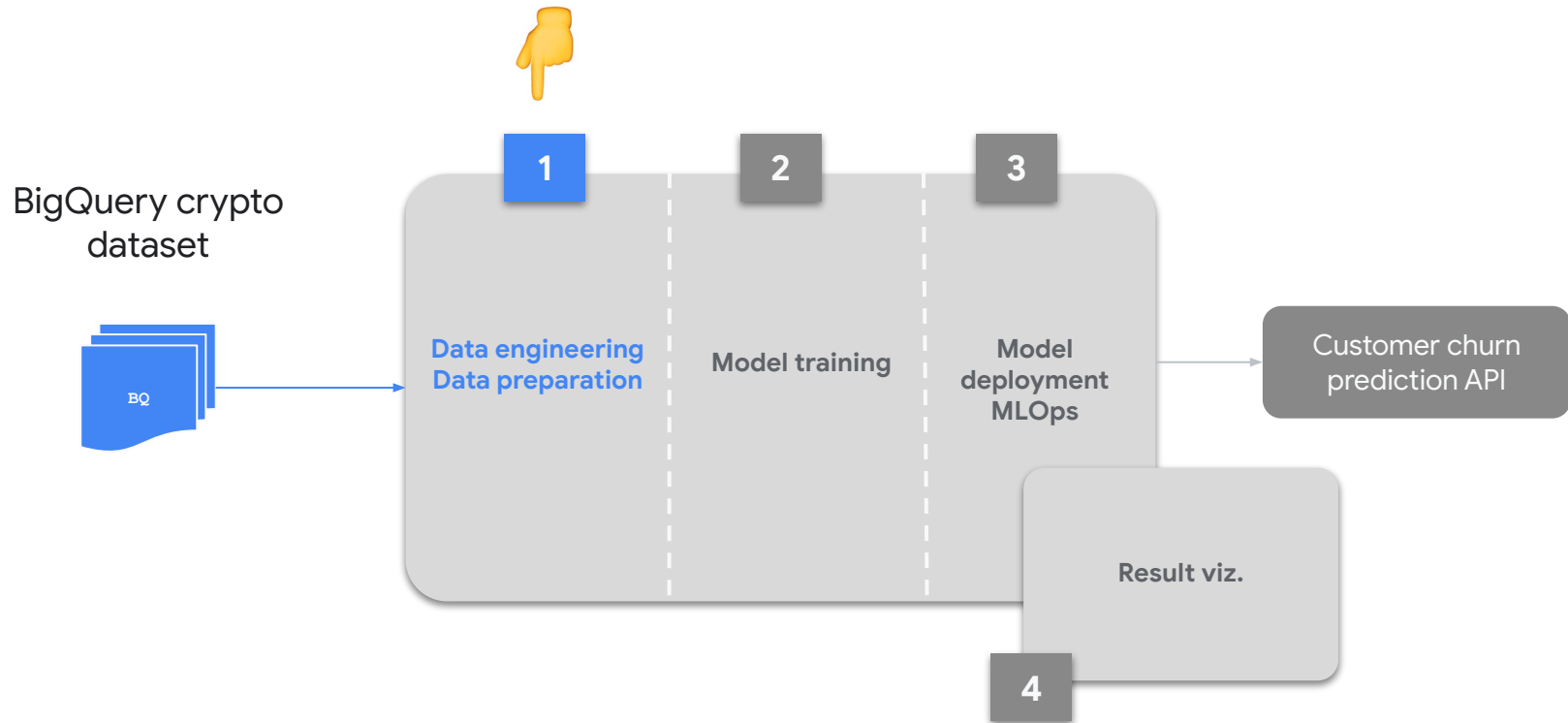
[Message clipped] [View entire message](#)





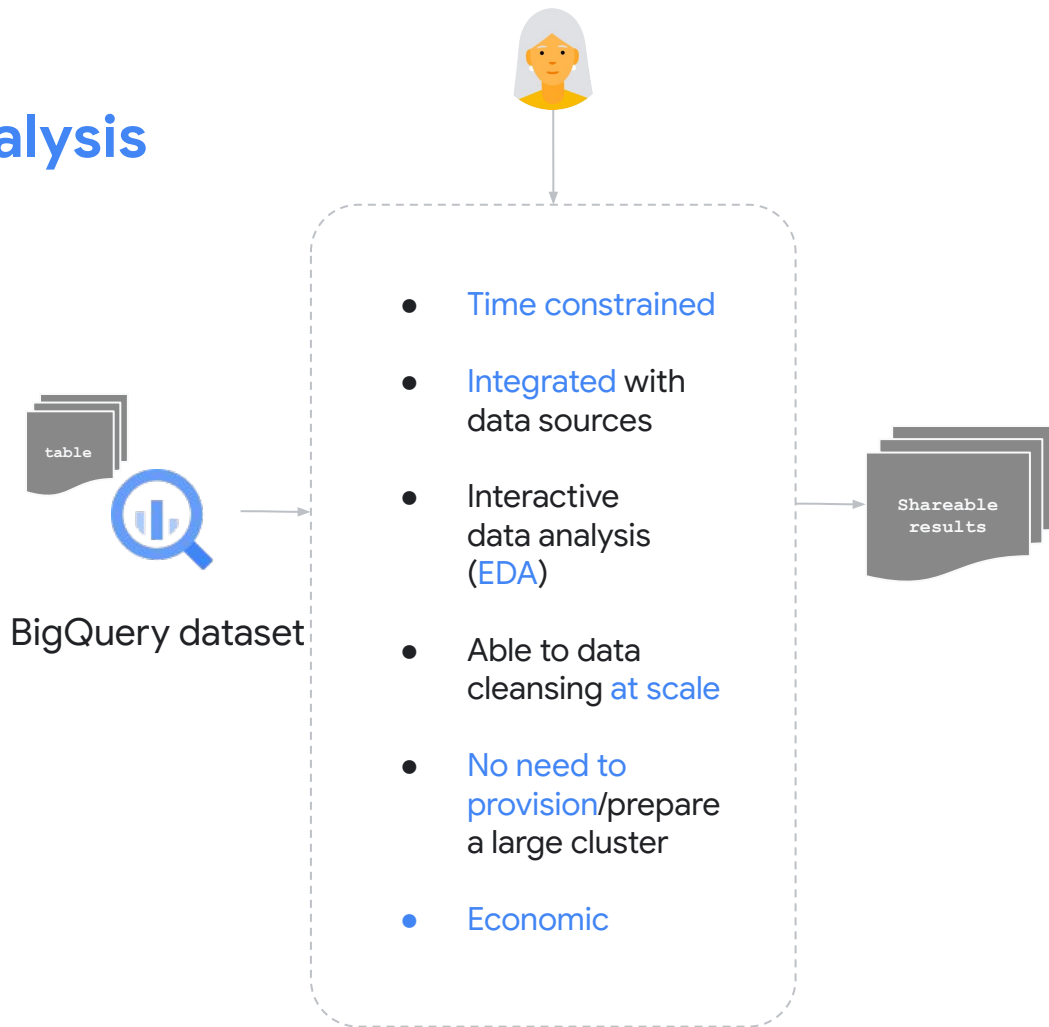
Where's the data?

Check on BigQuery  
Good luck!



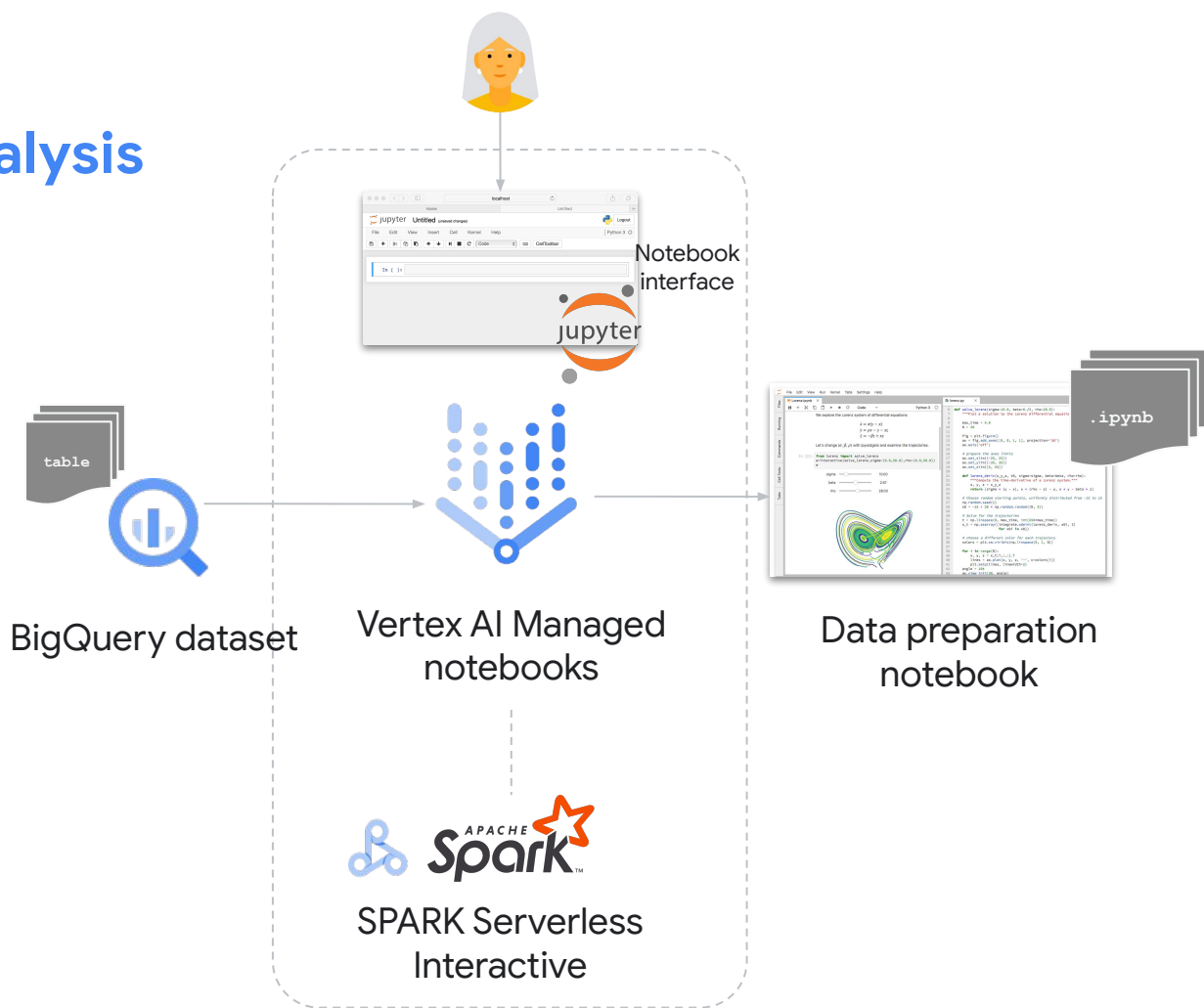
1

# 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

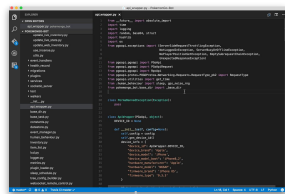


Prepared data

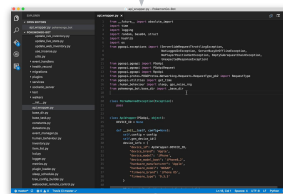
1

# Data Engineering

IDE interface



IDE interface



Production  
Code  
refactor

Pipeline  
code

```
def parse_arguments():  
    """Read arguments from a command line."""  
    parser = argparse.ArgumentParser(description='Arguments for period use --command')  
    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(args.v), stream=sys.stdout)  
    return args  
  
def main():  
    pass  
  
if __name__ == '__main__':  
    args = parse_arguments()  
    main()
```

data\_preparation.py



Data preparation  
notebook

Testing  
Error handling  
Argument Parsing  
Performance  
Security  
Libraries  
...

RAW data



APACHE  
**Spark**

Serverless  
Batch

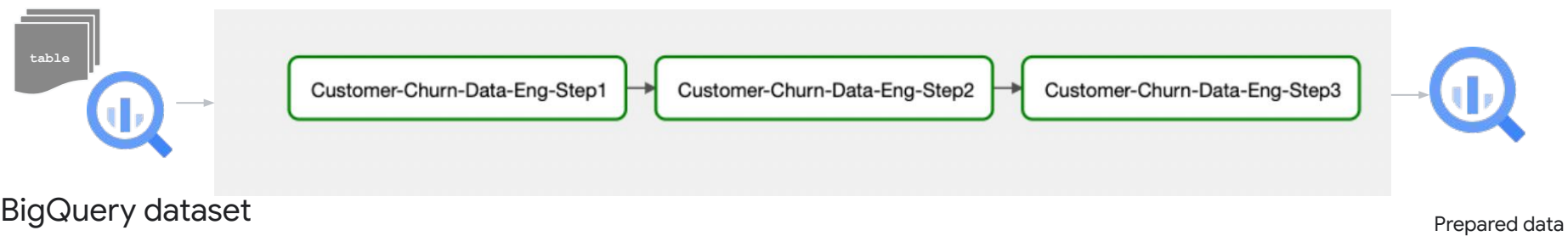


Composer  
dataproc  
serverless  
operator

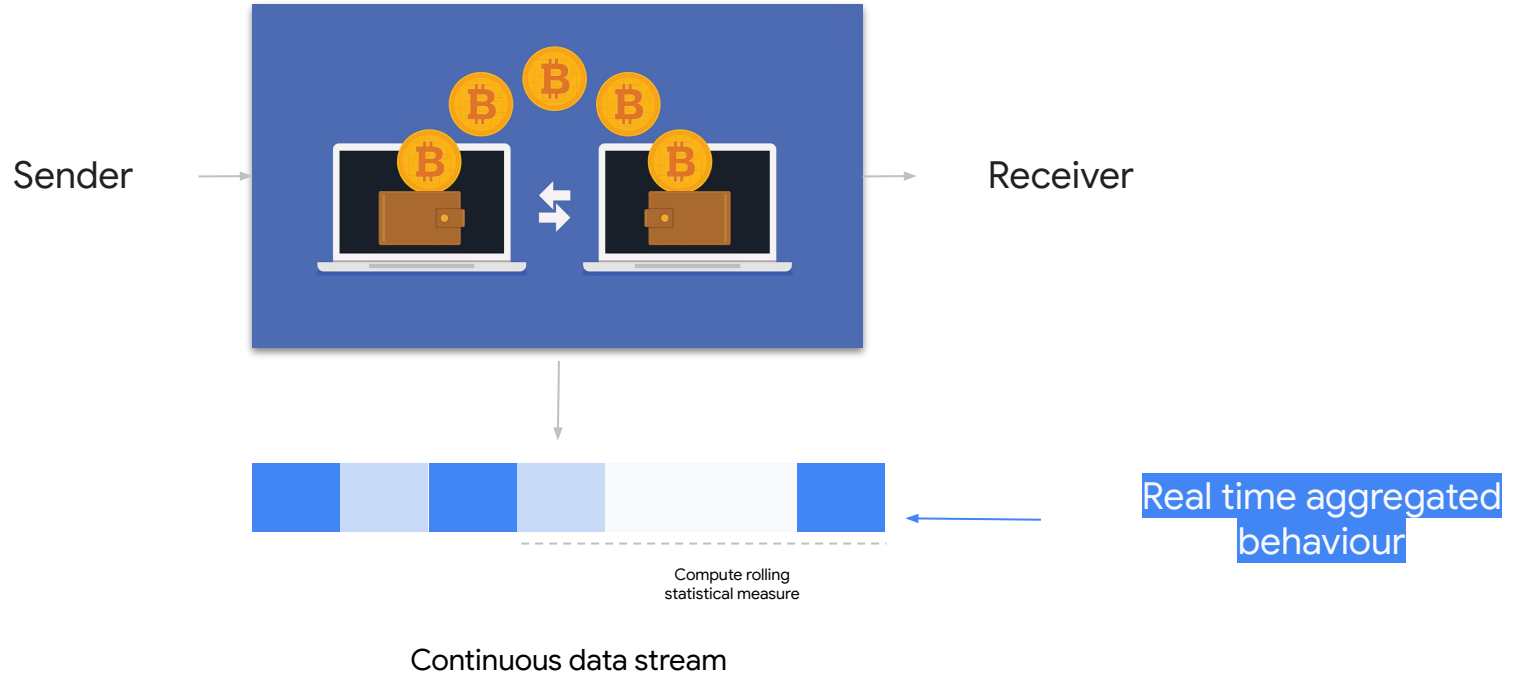


Prepared data

# Whats next?



# The **job** to be done: A **real time analytics dashboard**





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

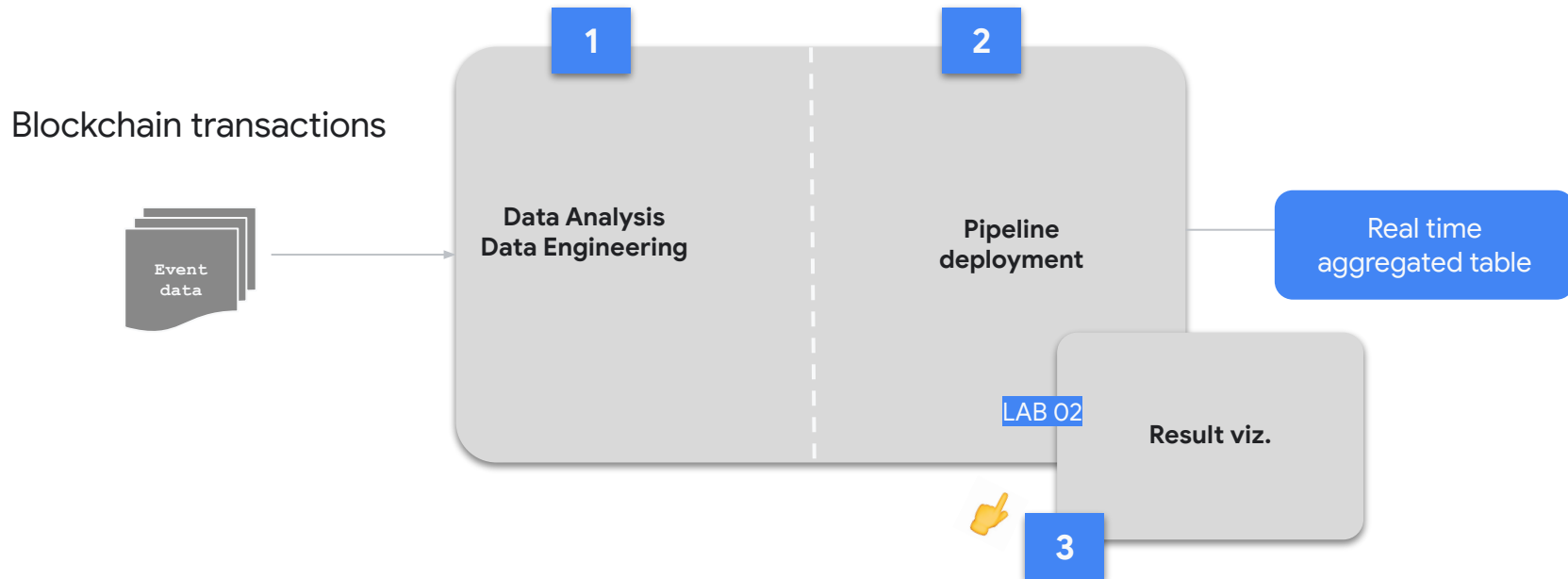
Blockchain transactions



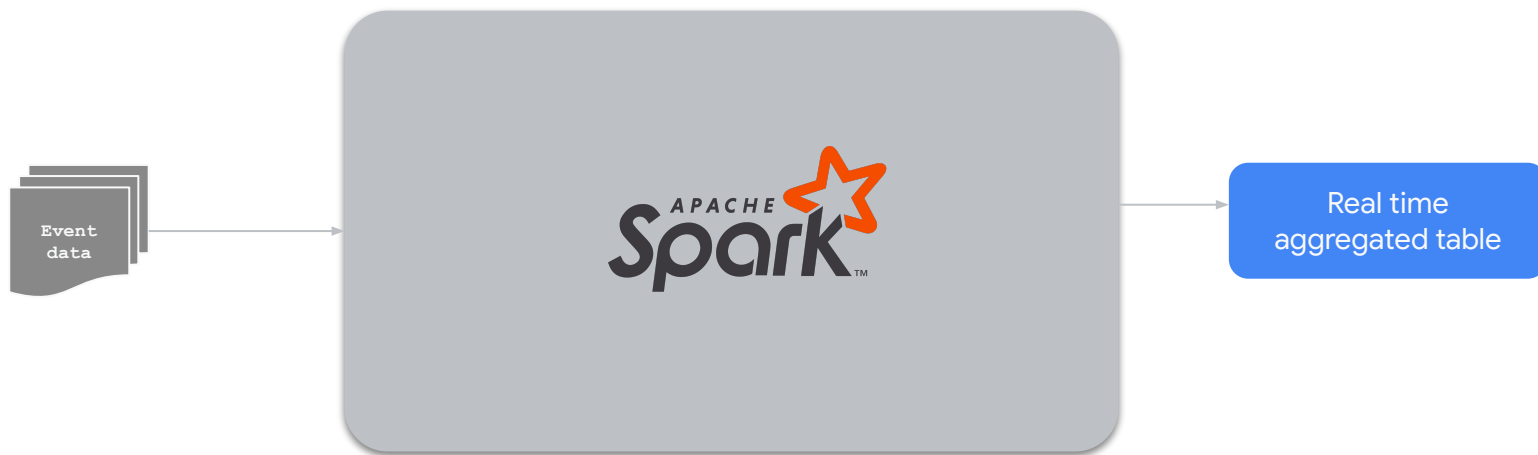
**Data system**

Real time  
aggregated table

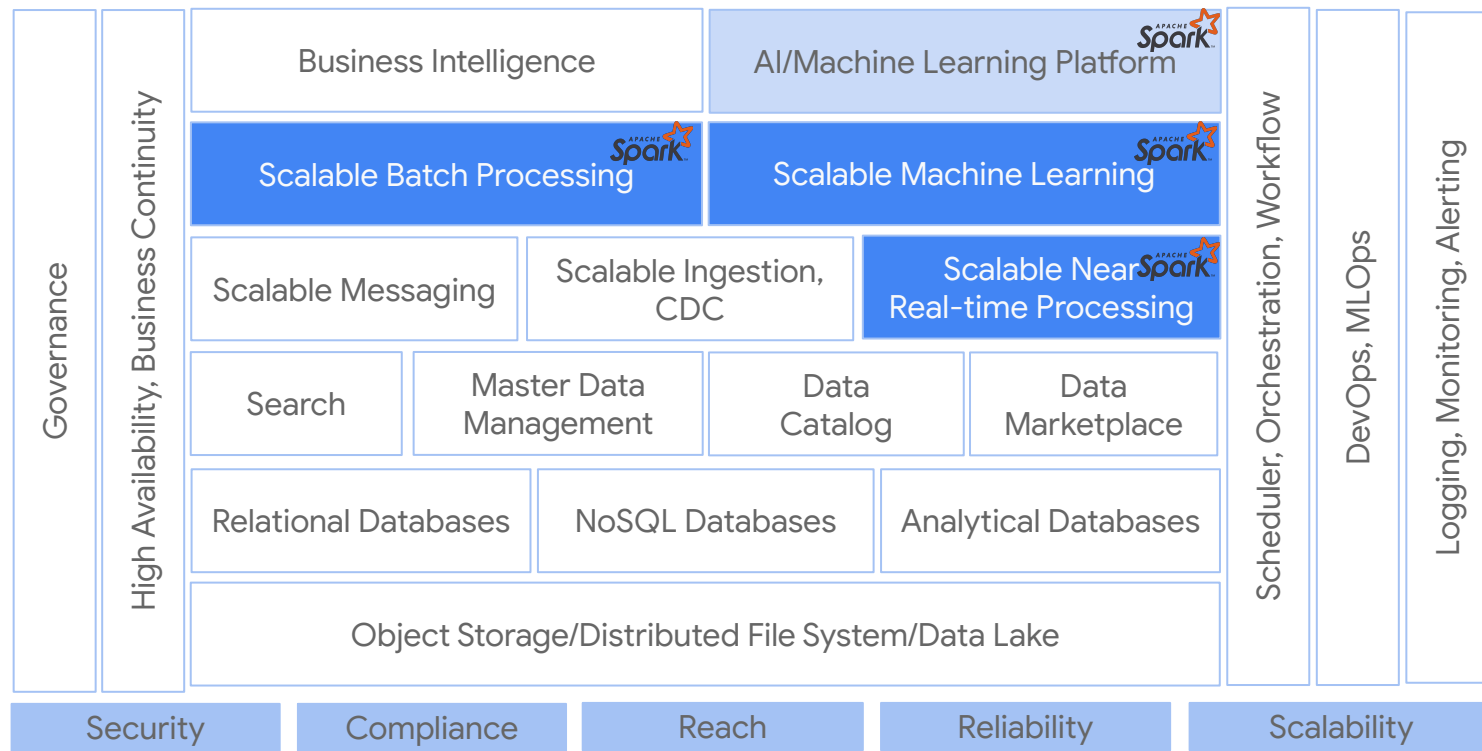
# The **job** to be done: Product dev stages



# The toolkit - SPARK



# The toolkit



# The personas



**Data Analyst**  
**Business Analyst**

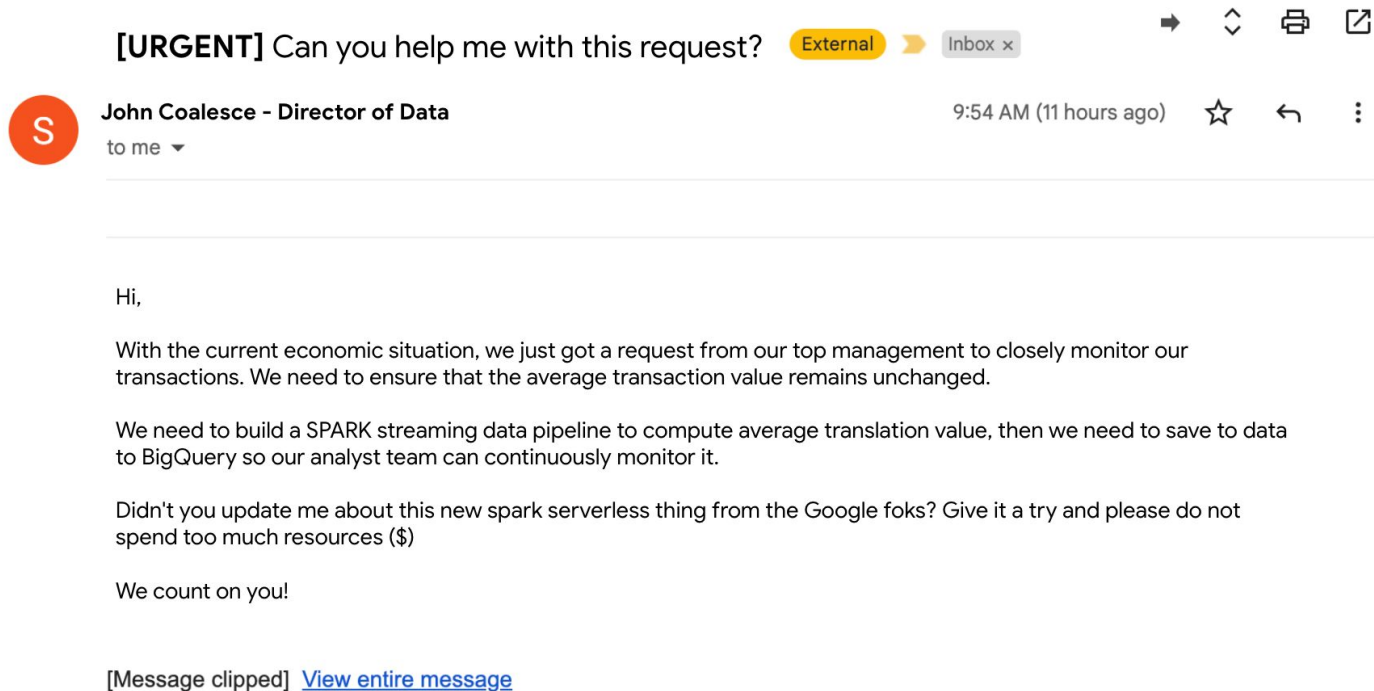


**Data Engineer**  
**ML Engineer**

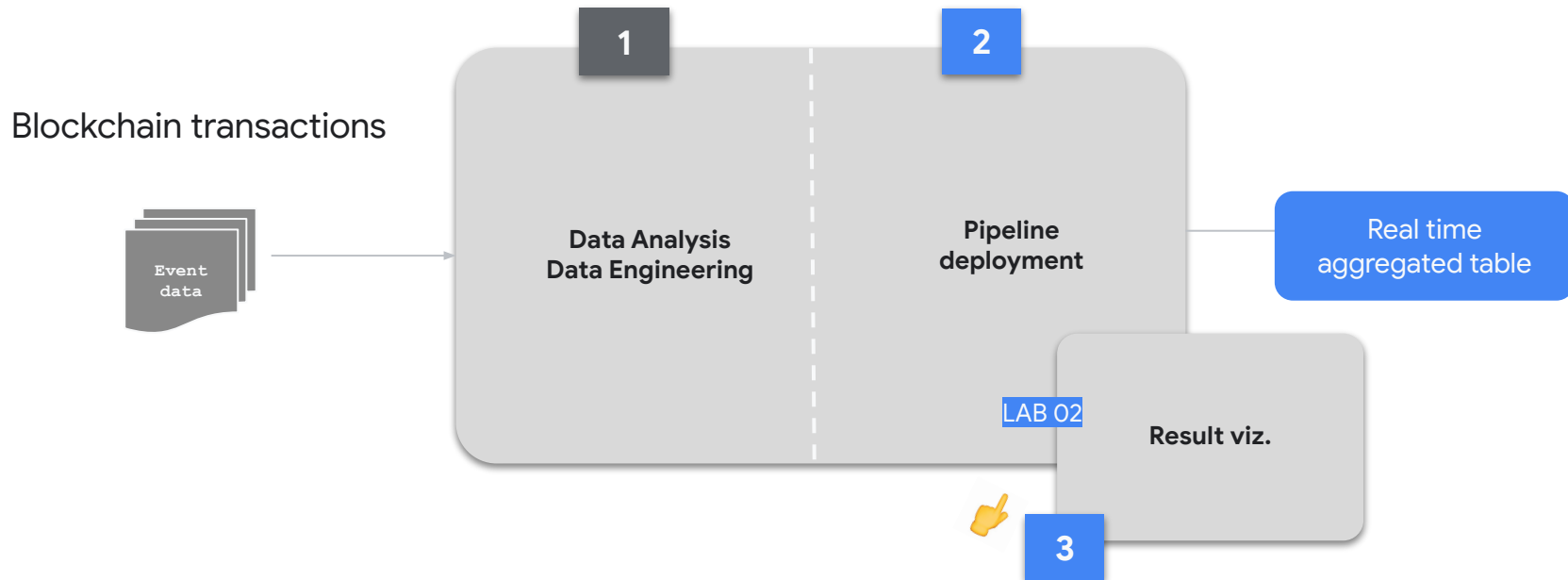


**Data Scientist**

# The trigger - You've got email

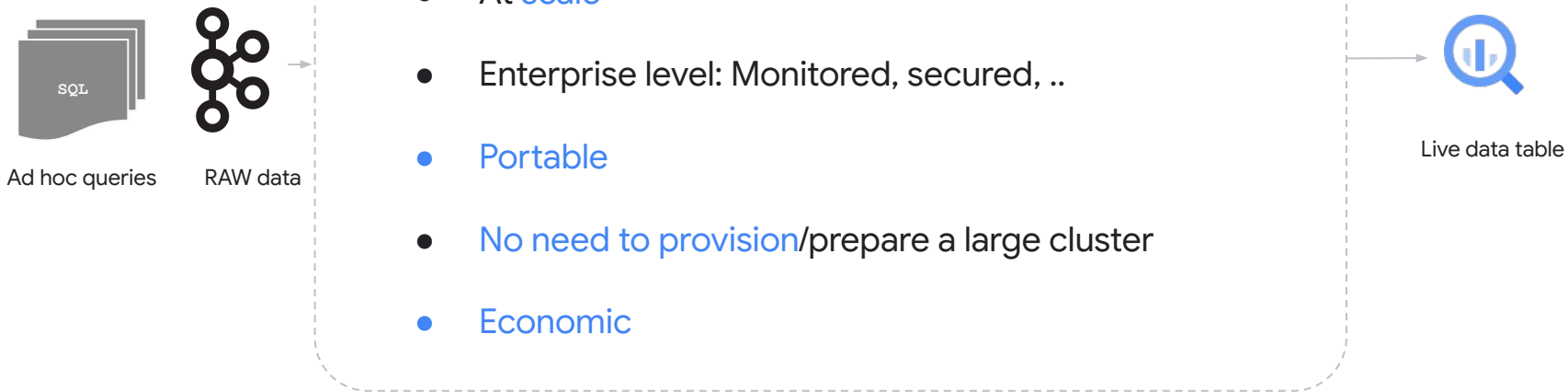


# The **job** to be done: Product dev stages



## 1

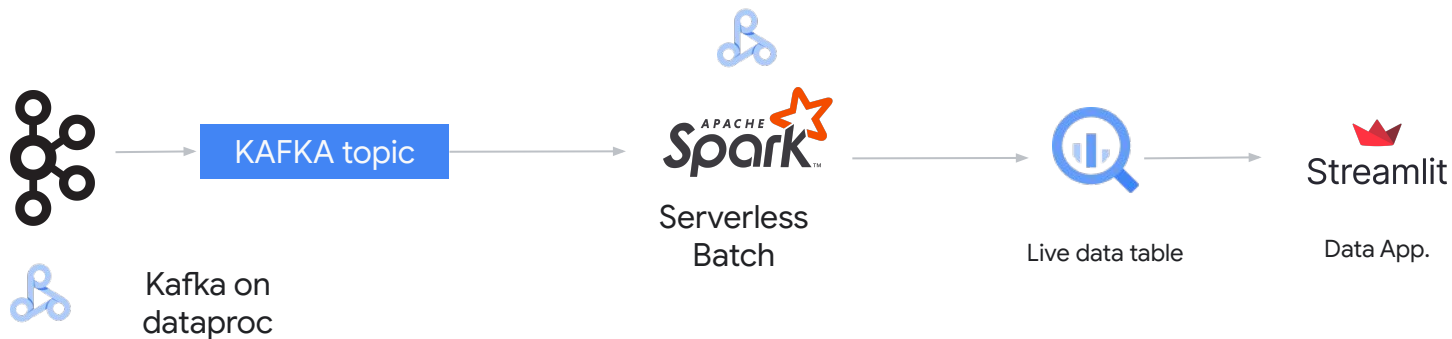
# Pipeline deployment





1

# Pipeline deployment



# SPARK - A first class citizen in Google Cloud

