

# Scaling Machine Learning with Apache Spark

Holly Smith  
Niall Turbitt

# About

Holly Smith

Senior Consultant at Databricks

- Professional Services and Training
- Experience
  - Credit Risk & Decisioning
  - Mobile Banking
  - Forecasting & Optimisation
- BSc Mathematics University of Greenwich

**DATA+AI SUMMIT EUROPE**



**#DataTeams #DataAISummit**



# About

Niall Turbitt

Senior Data Scientist at Databricks

- Professional Services and Training
- Experience
  - e-Commerce
  - Supply Chain and Logistics
  - Recommender Systems & Personalisation
- MS Statistics University College Dublin
- BA Mathematics & Economics Trinity College Dublin

**DATA+AI SUMMIT EUROPE**



**#DataTeams #DataAISummit**



# Outline

- Motivation
- Spark Architecture Recap
- Paradigms of ML on Spark:
  - Training & Tuning
    - Spark MLlib
    - Pandas Function APIs
    - Hyperopt
  - Inference
    - Pandas UDFs

# Motivation

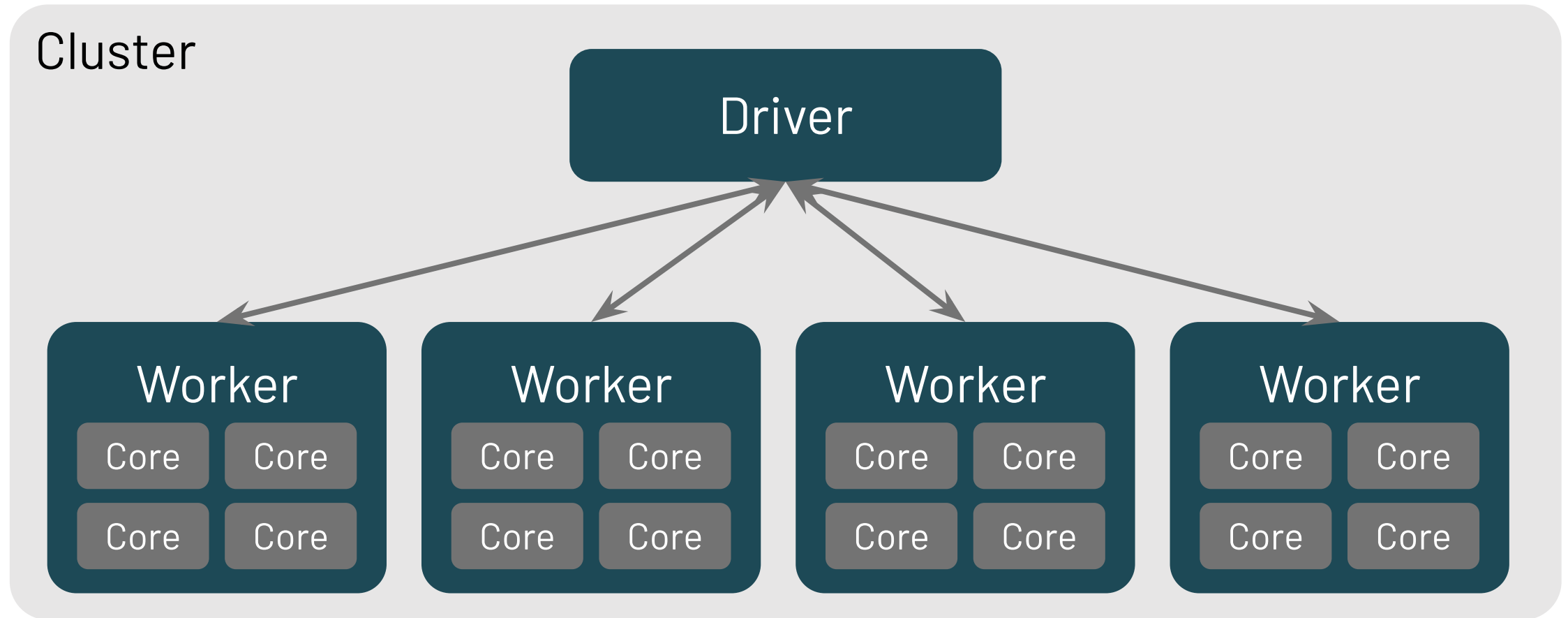
- Confusion around where and when to use Spark for Machine Learning
- Spark is powerful: harness the full potential of Spark for Machine Learning
- Fast moving environment
  - Use new APIs and techniques before they become mainstream

# Motivation

- Confusion around where and when to use Spark for Machine Learning
- Spark is powerful: harness the full potential of Spark for Machine Learning
- Fast moving environment
  - Use new APIs and techniques before they become mainstream



# Refresher: Spark Architecture



= a single node



# ML on Spark

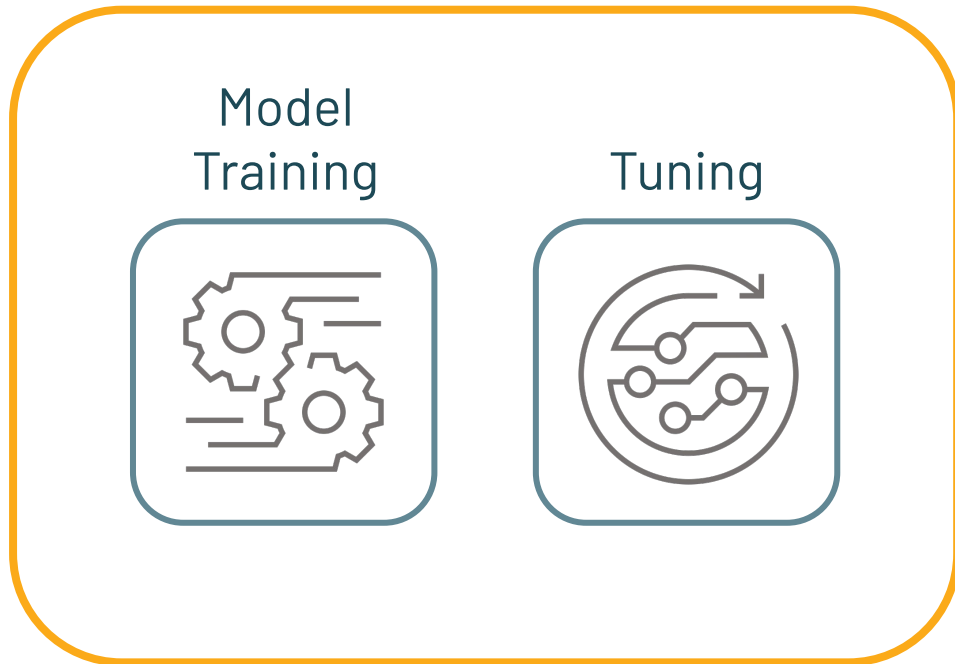
**DATA+AI SUMMIT EUROPE**

**#DataTeams #DataAISummit**



# The ML Lifecycle

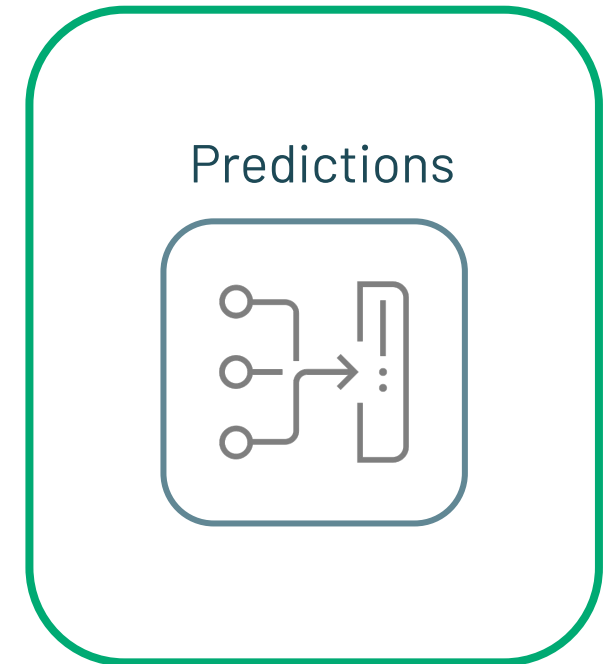
## Training Phase



## Model Artefact

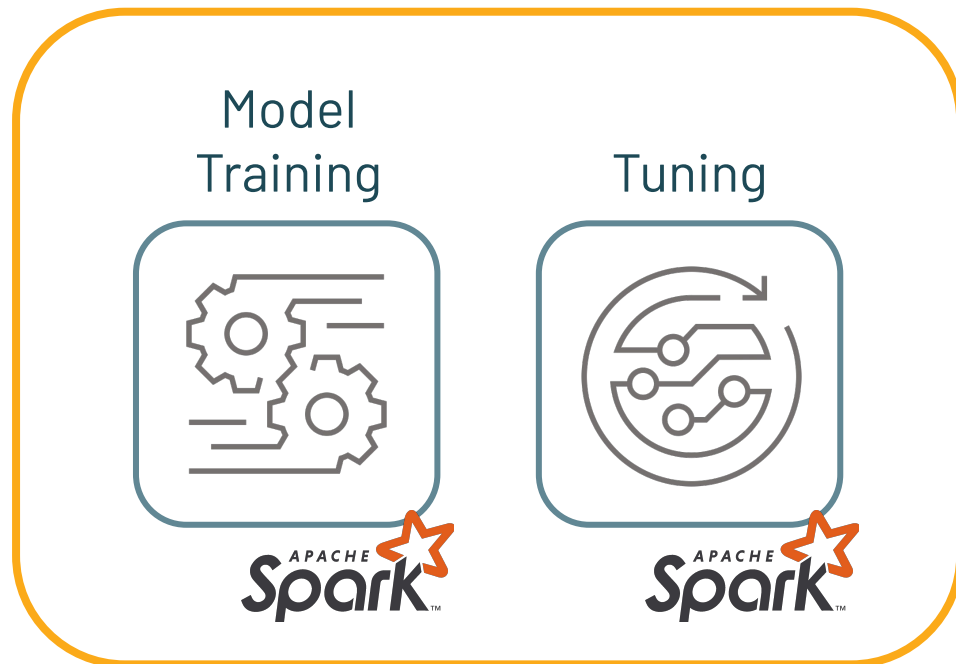


## Inference



# The ML Lifecycle

## Training Phase



## Model Artefact



## Inference





# ML Training & Tuning on Spark

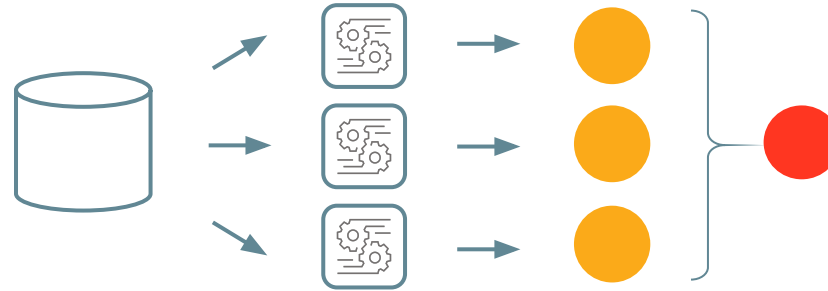
Distributed  
ML Library

Data Parallel  
Training

Training Data

Nodes

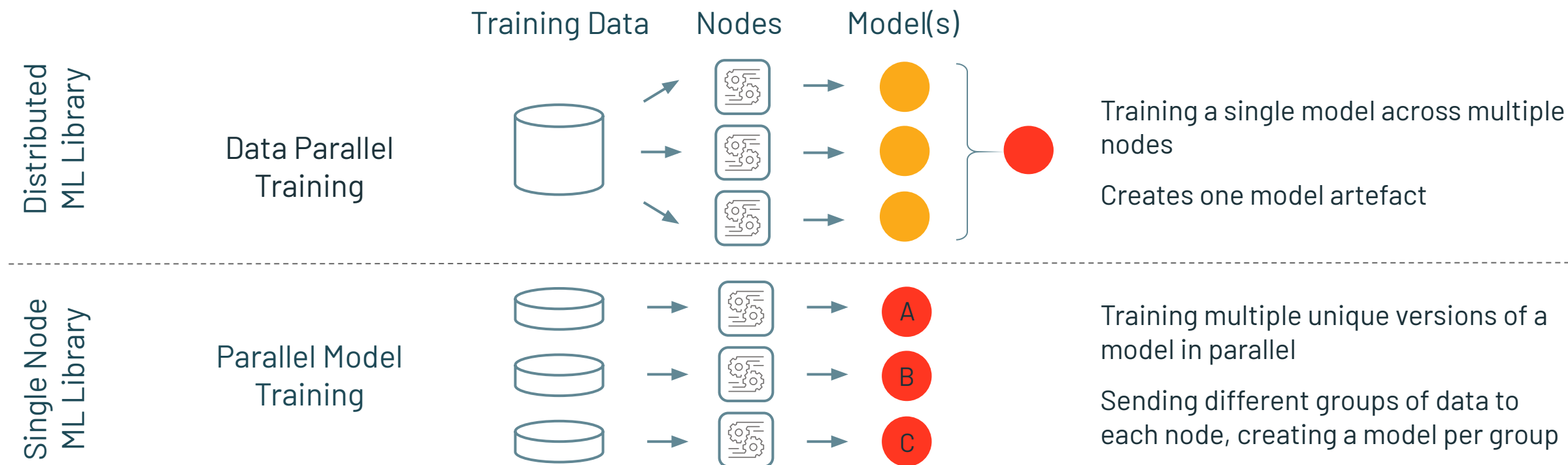
Model(s)



Training a single model across multiple  
nodes

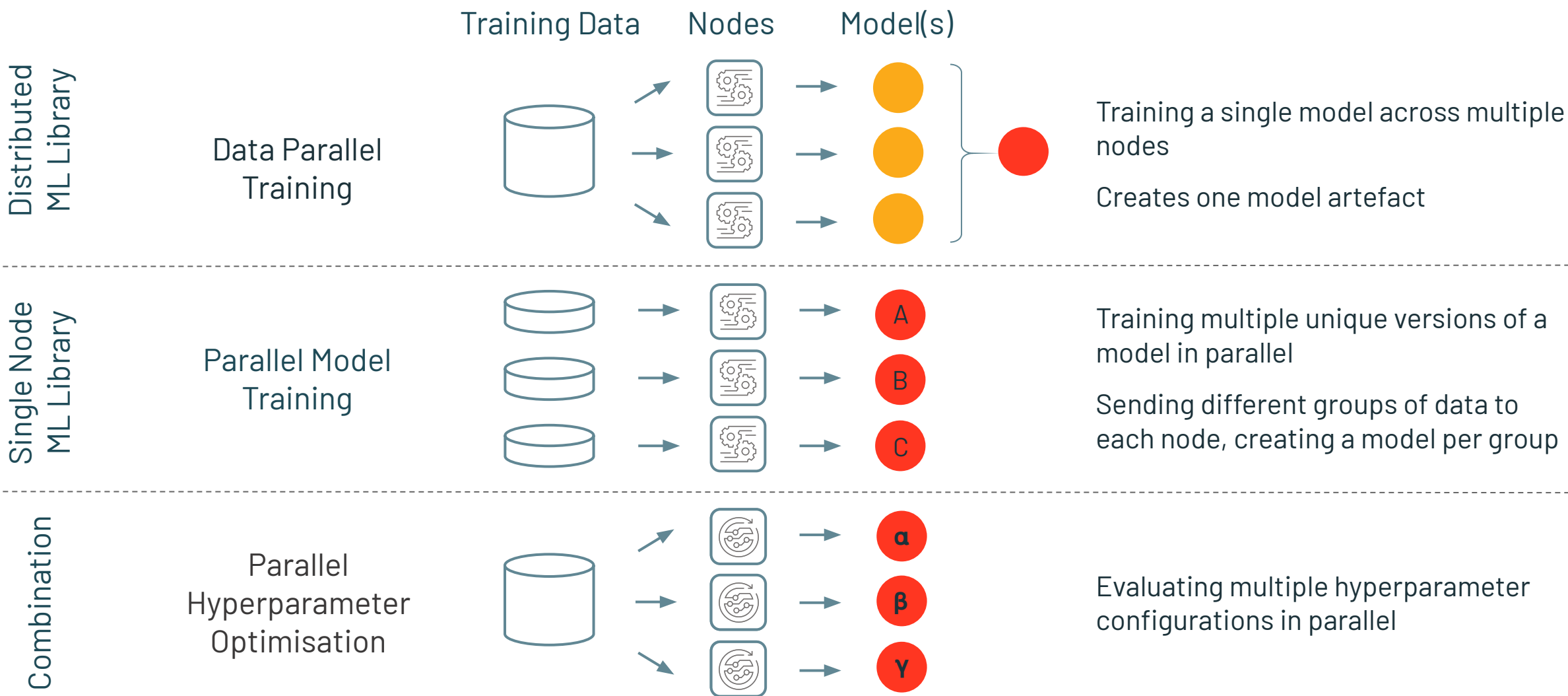
Creates one model artefact

# ML Training & Tuning on Spark





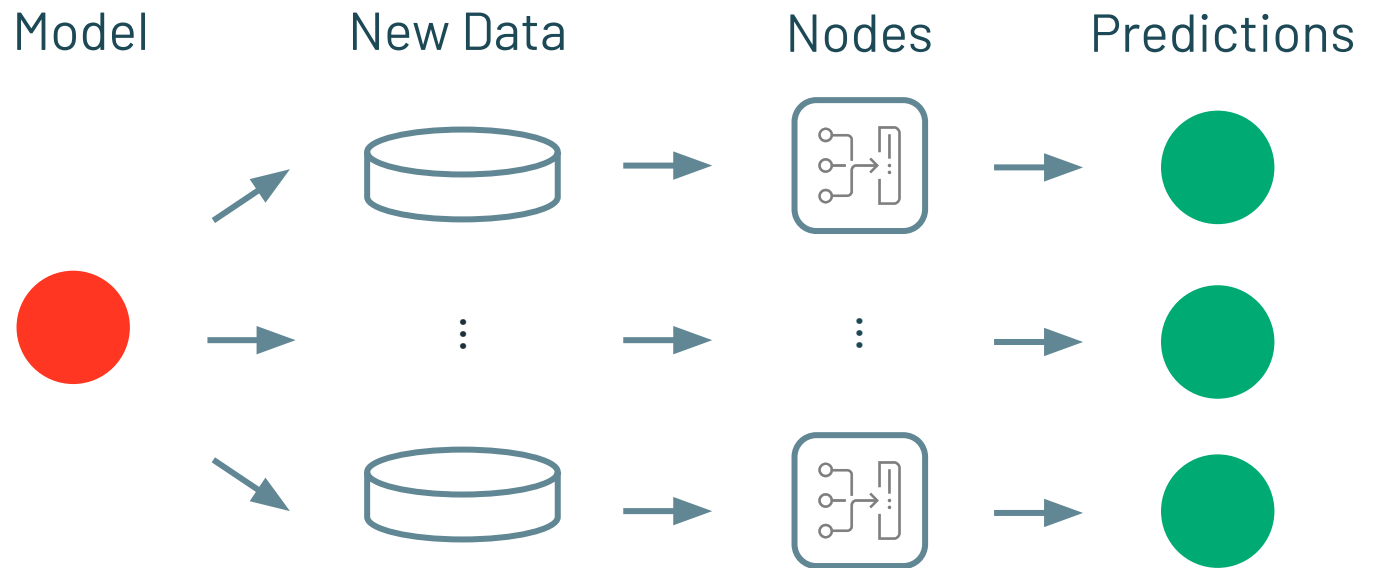
# ML Training & Tuning on Spark



# ML Inference on Spark

For both distributed and single node ML libraries:

1. Take trained model
2. Distribute new instances
3. Apply model in parallel





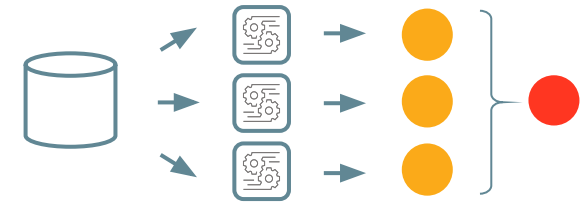
# ML Project Considerations

- Data Dependent
- Compute Resources Available
- Single machine vs distributed computing
- Inference: Deployment Requirements

	Throughput	Latency	Example
<b>Batch</b>	High	Hours to days	Customer churn prediction
<b>Streaming</b>	Medium	Seconds to minutes	Predictive maintenance
<b>Real-time</b>	Low	Milliseconds	Fraud detection

# Spark MLlib

Parallelising Single-Model Training



Distributed ML Library

- Spark's Machine Learning Library
  - ML algorithms
  - Featurization
  - Pipelines
- MLlib vs sklearn
- A note on terminology:

What is meant by "MLlib"

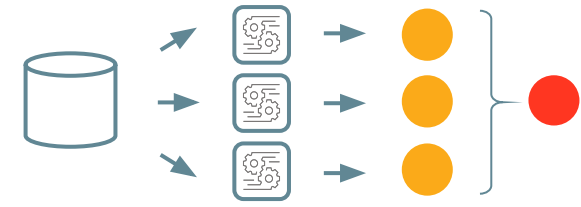
Spark.mllib RDD based API  
**Maintenance Mode**

Spark.ml Dataframe based  
API  
**Recommended**



# Spark MLlib

Parallelising Single-Model Training



Distributed ML Library

- Spark's Machine Learning Library
  - ML algorithms
  - Featurization
  - Pipelines
- MLlib vs sklearn
- A note on terminology:

```
from pyspark.ml.regression import LinearRegression

train_df = spark.read...
test_df = spark.read...

lr = LinearRegression().fit(train_df)
predictions = lr.transform(test_df)
```

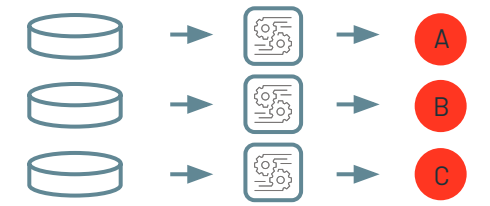
What is meant by "MLlib"

Spark.mllib RDD based API  
**Maintenance Mode**

Spark.ml Dataframe based API  
**Recommended**

# Pandas Function API – Grouped Map

Parallelising training of independent models

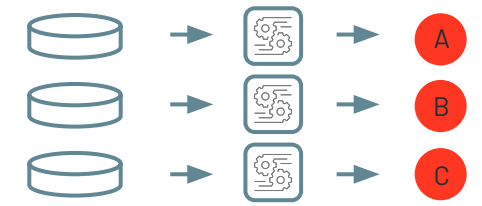


Single Node ML Library

- `DataFrame.groupby().applyInPandas()` **NEW**
- Directly apply a Python native function against a Spark DataFrame as if each group is a Pandas DataFrame
- “split-apply-combine” pattern:
  - Split data into groups
  - Apply function on each group
  - Combine results into new Spark DataFrame

# Pandas Function API – Grouped Map

Parallelising training of independent models



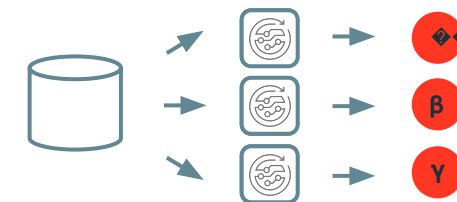
Single Node ML Library

- `DataFrame.groupby().applyInPandas()` **NEW**
- Directly apply a Python native function against a Spark DataFrame as if each group is a Pandas DataFrame
- “split-apply-combine” pattern:
  - Split data into groups
  - Apply function on each group
  - Combine results into new Spark DataFrame

```
def train_model(pd.DataFrame) -> pd.DataFrame:  
  
    # fit single-node model  
    sklearn_model.fit()  
    ...  
  
    return pandas_df  
  
spark_df.groupby("device_id")  
        .applyInPandas(train_model,  
                        schema=return_schema)
```

# Hyperopt

## Parallelising Hyperparameter Optimisation



Combination

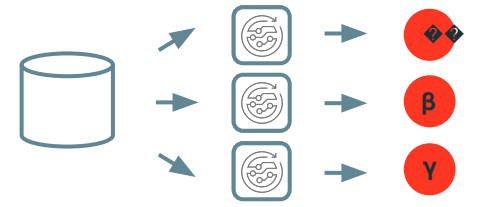
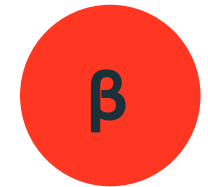
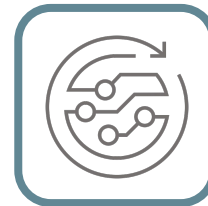
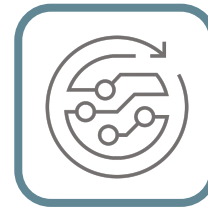
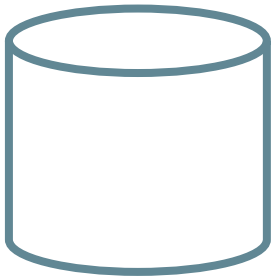
- Open source hyperparameter optimisation package
  - Enables either serial or parallel optimisation over provided search spaces
  - Can tune both distributed and single node libraries
  - HOWEVER: distributed training and distributed tuning don't mix
- Bayesian based approach
  - Adaptively selects new hyperparameter settings to explore based on prior results
  - Enables exploration of the hyperparameter space in an intelligent way
  - Allows a wider search space with more hyperparameters



# Hyperopt

Parallelising Hyperparameter Optimisation

- Distributed Hyperopt with a single node library
  - SparkTrials

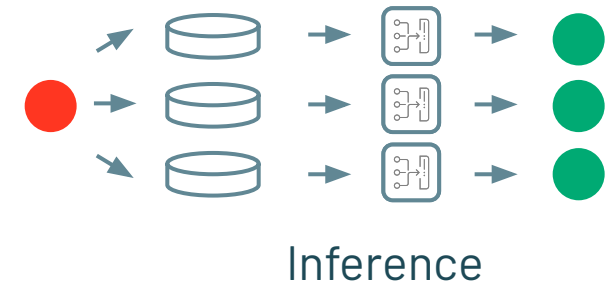


Combination

# Pandas Scalar Iterator UDF

## Distributing inference

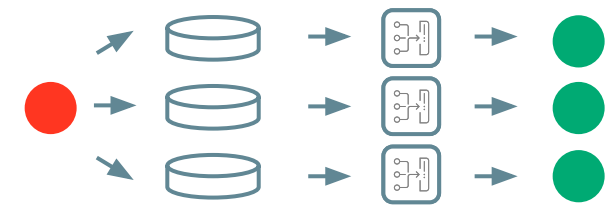
- Pandas UDFs can accept an iterator of `pandas.Series` or `pandas.DataFrame`
- Spark DataFrame is split into batches and the function called for each batch
- Iterator negates the need to repeatedly load the same model for every batch in the same Python worker process



# Pandas Scalar Iterator UDF

## Distributing inference

- Pandas UDFs can accept an iterator of `pandas.Series` or `pandas.DataFrame`
- Spark DataFrame is split into batches and the function called for each batch
- Iterator negates the need to repeatedly load the same model for every batch in the same Python worker process



Inference

```
@pandas_udf
def predict_udf(iterator):

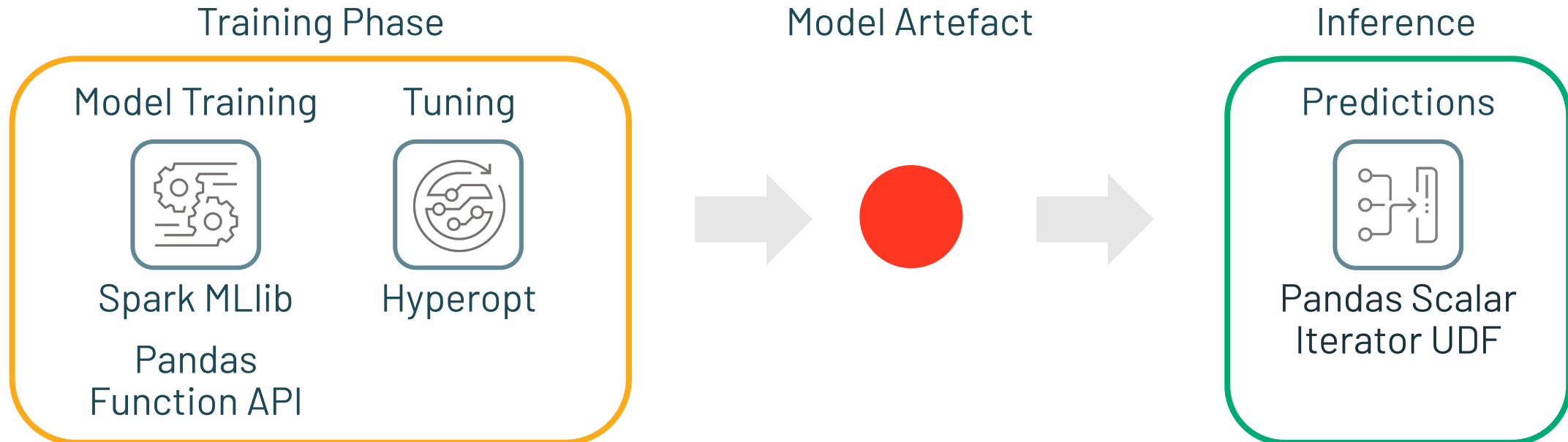
    # load model
    model = ...

    for features in iterator:
        yield pd.Series(model.predict(features))

spark_df.withColumn("prediction",
                    predict_udf(*input_cols))
```

# Conclusion

Distributing workloads allows you to scale, either by using libraries that are multi or single node to suit your project.





# Resources

Notebook: [bit.ly/scaling\\_ml\\_spark\\_2020](https://bit.ly/scaling_ml_spark_2020)

Pandas UDF Blog post – [bit.ly/Pandas\\_UDF](https://bit.ly/Pandas_UDF)

Docs:

- MLlib – [bit.ly/ML\\_lib](https://bit.ly/ML_lib)
- Hyperopt – [bit.ly/hyperopt\\_spark](https://bit.ly/hyperopt_spark)
- Pandas Grouped Map – [bit.ly/grouped\\_map](https://bit.ly/grouped_map)

# Feedback

Your feedback is important to us.  
Don't forget to rate  
and review the sessions.

