# apache Ignite

In-Memory Performance
Durability of Disk

# Scalable Machine and Deep Learning with Apache Ignite

Denis Magda
Apache Ignite PMC Chair
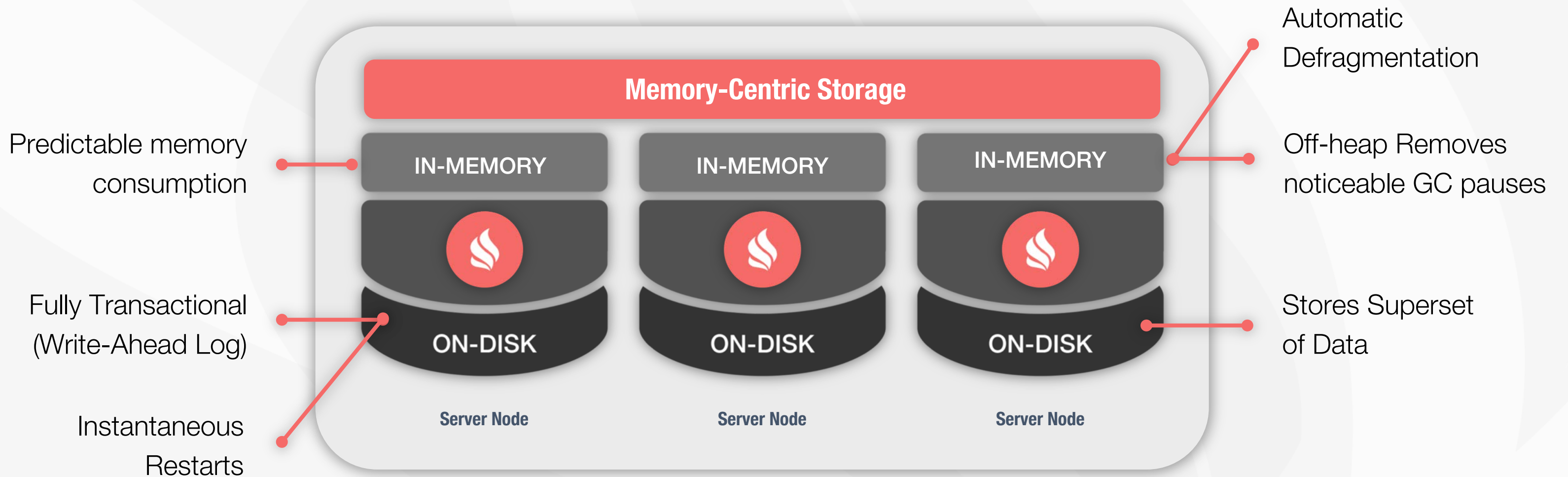GridGain Director of Product Management

# Agenda

- Why Machine Learning at Scale?
- Ignite Machine Learning
- Genetic Algorithms
- TensorFlow Integration
- Demo
- Q&A

# Why Machine Learning at Scale?

1. Models trained and deployed in different systems
   - Move data out for training
   - Wait for training to complete
   - Redeploy models in production

2. Scalability
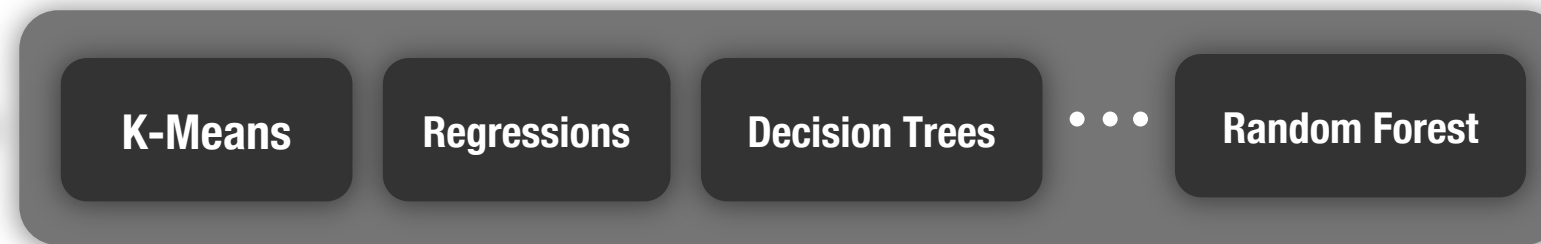   - Data exceed capacity of single server
   - Burden for developers
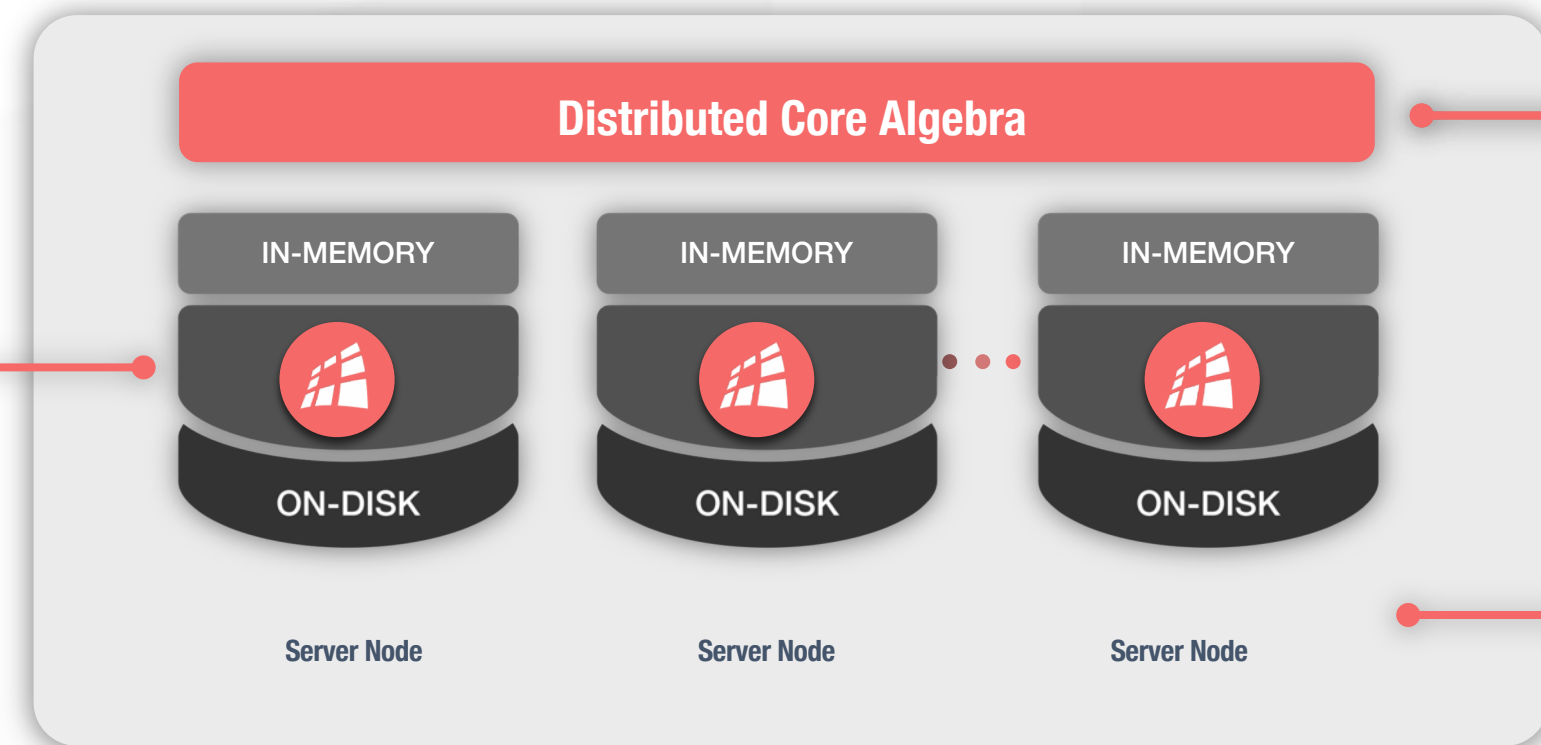
# Memory-Centric Storage

# Machine Learning

**R**  **C++**  **Python**  **Java**  **Scala**  **REST** — Multi-Language Support

Distributed Algorithms — **K-Means**  **Regressions**  **Decision Trees**  • • •  **Random Forest**

**Distributed Core Algebra** — Dense and Sparse Algebra

IN-MEMORY  IN-MEMORY  IN-MEMORY

Large Scale Parallelization

ON-DISK  ON-DISK  • • •  ON-DISK — No ETL

Server Node  Server Node  Server Node

GridGain

# Record to Node Mapping

**Key** → **Partition** → 

Server Node

ON-DISK

# Caches and Partitions

**Cache**

| Partition 1 | Partition 2 |
|---|---|
| K1, V1 | K5, V5 |
| K3, V3 | K7,V7 |
| K2, V2 | K6, V6 |
| K4, V4 | K8, V8 |
| | K9, V9 |

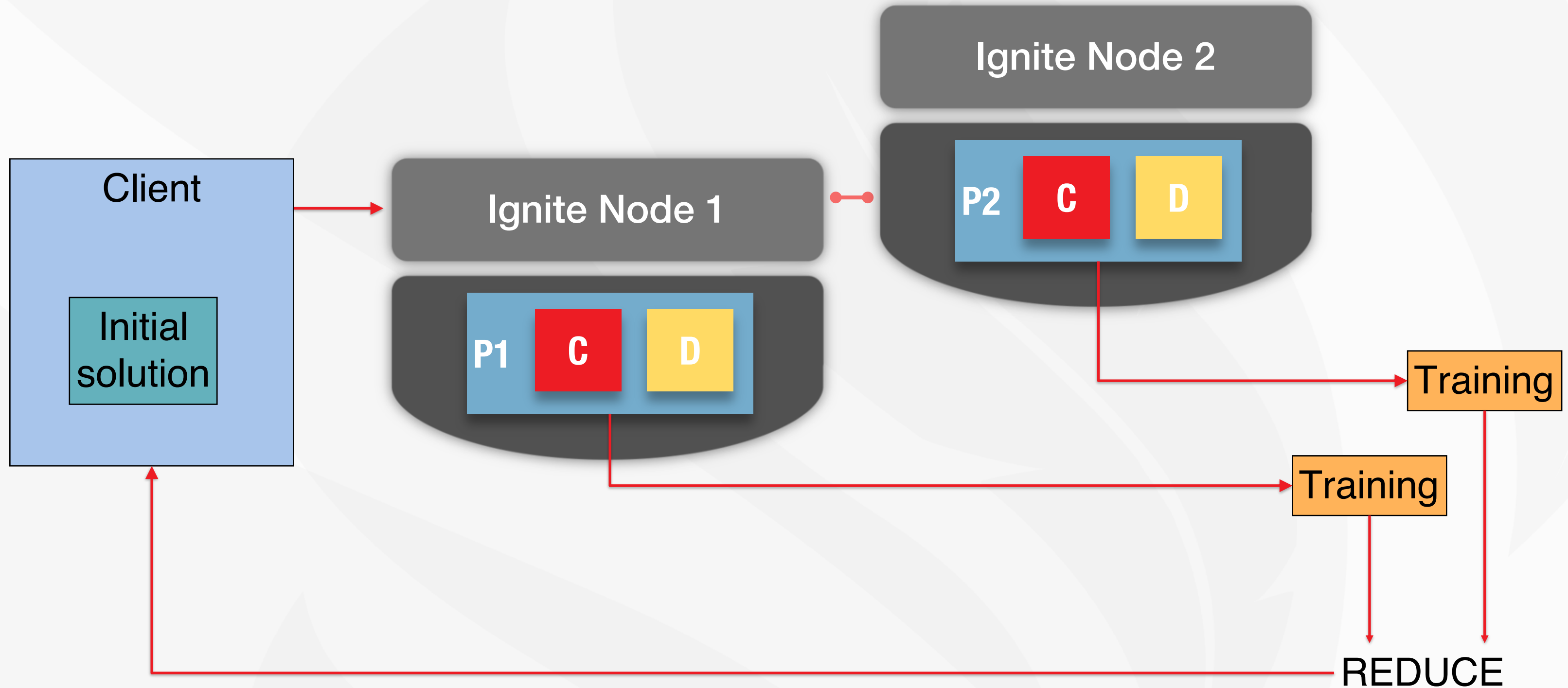**Partition 1**

**Partition 2**
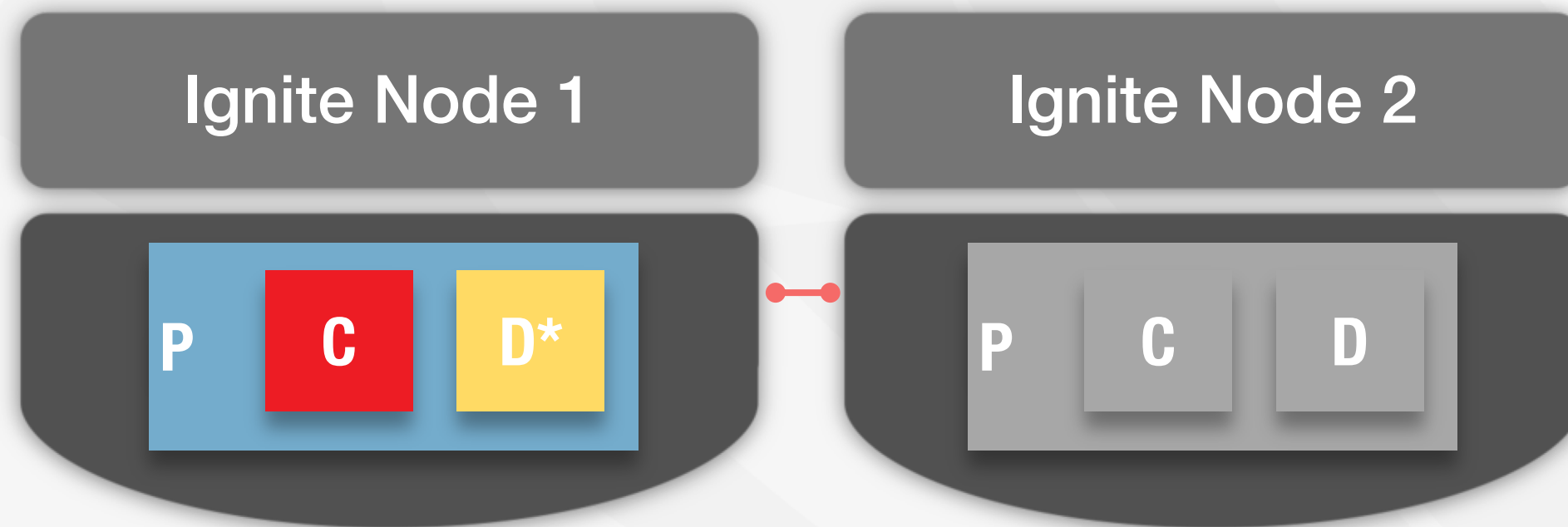
Partition-Based Dataset

# Training Failover



P = Partition
C = Partition Context
D = Partition Data
D* = Local ETL

# Continuous Learning

**Data** → **Training**

**Solution**

Iteration

# Algorithms and Applicability

|  | Classification | Regression |
|---|---|---|
| Description | Identify to which category a new observation belongs, on the basis of a training set of data | Modeling the relationship between a scalar dependent variable y and one or more explanatory variables x |
| Applicability | spam detection, image recognition, credit scoring, disease identification | drug response, stock prices, supermarket revenue |
| Algorithms | nearest neighbor, decision tree classification, neural network | linear regression, decision tree regression, nearest neighbor, neural network |

# Algorithms and Applicability

| | Clustering | Preprocessing |
|---|---|---|
| Description | Grouping a set of objects in such a way that objects in the same group are more similar to each other than to those in other groups | Feature extraction and normalization |
| Applicability | customer segmentation, grouping experiment outcomes, grouping shopping items | transform input data, such as text, for use with machine learning algorithms |
| Algorithms | k-means | Normalization preprocessor |

# Linear Regression

- Ordinary Least Squares
- Linear Regression Trainer
  - **QR Decomposition**
  - **Gradient Descent**

```
// y = bx + a
LinearRegressionModel model = trainer.train(trainSet);
double prediction = model.predict(sampleObject);

// Prepare trainSet
...

// QR Decomposition
LinearRegressionQRTrainer trainer = new LinearRegressionQRTrainer();
LinearRegressionModel mdl = trainer.train(trainSet);

// Gradient Descent
LinearRegressionSGDTrainer trainer = new LinearRegressionSGDTrainer(
    1000, 1e-6);
LinearRegressionModel mdl = trainer.train(trainSet);
```

# Decision Trees

- Data stored by features
- Related data on same node
- Features
  - Continuous
  - Categorical

```java
// Train the model
DecisionTreeModel mdl = trainer.train(
    new BiIndexedCacheColumnDecisionTreeTrainerInput(
        cache, new HashMap<>(), ptsCnt, featCnt));

// Estimate the model on the test set
IgniteTriFunction<Model<Vector, Double>,
    Stream<IgniteBiTuple<Vector, Double>>,
    Function<Double, Double>,
    Double> mse = Estimators.errorsPercentage();

Double accuracy = mse.apply(mdl, testMnistStream.map(
    v -> new IgniteBiTuple<>(v.viewPart(0, featCnt), v.getX(featCnt))),
    Function.identity());

System.out.println(">>> Errs percentage: " + accuracy);
```
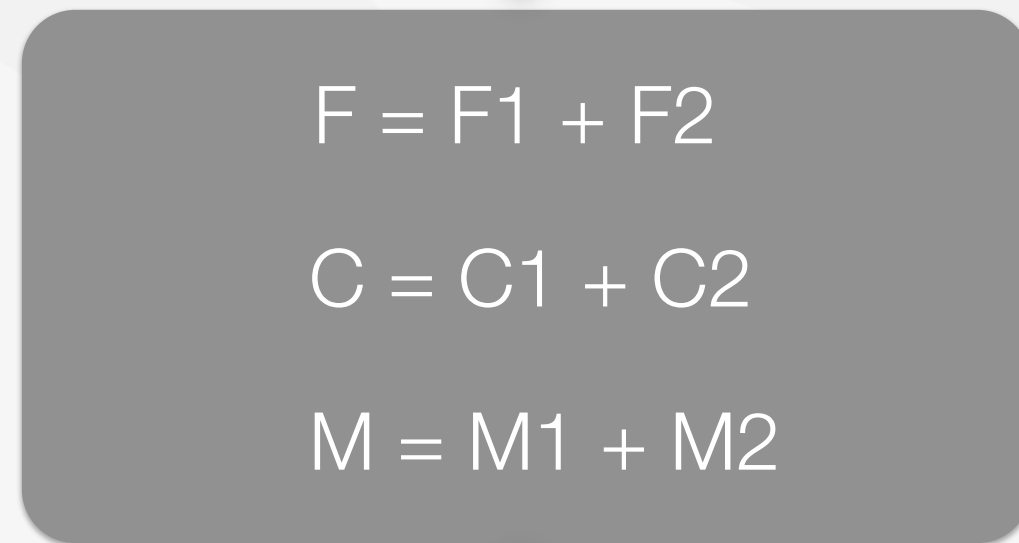
# Demo: Fraud Detection

# Genetic Algorithms

Biological Evolution
Simulation

F = F1 + F2

C = C1 + C2

M = M1 + M2

F = Fitness Calculation
C = Crossover
M = Mutation

Collocated Computation

F1, C1, M1

F2, C2, M2

**Chromosome and Genes Cluster**

IN-MEMORY

ON-DISK

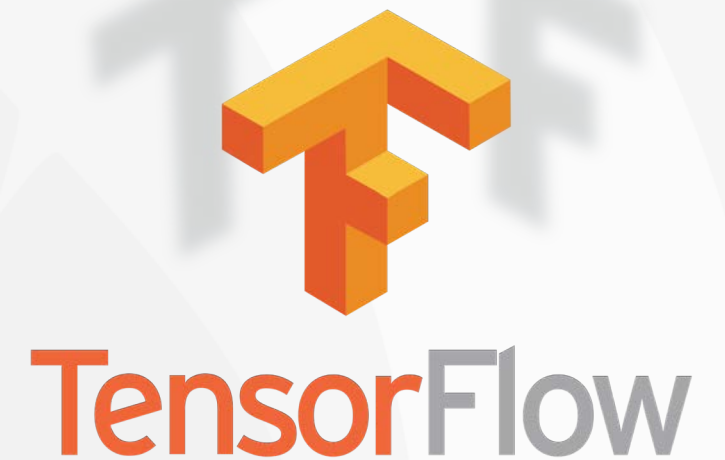IN-MEMORY

ON-DISK

**Ignite Cluster**

# TensorFlow Integration: Benefits

- Ignite as distributed data source
  - Perfect fit for distributed TF training

- Less ETL
  - TF nodes deployed together with Ignite nodes
  - In-machine data movement only

- TF tasks execution in-place in Ignite
  - Roadmap

# TensorFlow Integration: Main Features

- Distribution of user tasks written in Python

- Automatic creation and maintenance of TF cluster

- Minimization of ETL costs

- Fault tolerance for both Ignite and TF instances

# Demo: TensorFlow and Ignite

# Summary: Apache Ignite Benefits

- Massive scalability
  - Horizontal + Vertical
  - RAM + Disk

- Zero-ETL
  - Train models and run algorithms in place

- Fault tolerance and continuous learning
  - Partition-based dataset

# Resources

- Apache Ignite ML Documentation:
  - https://apacheignite.readme.io/docs

- ML Blogging Series:
  - [Genetic Algorithms with Apache Ignite](#)
  - [Introduction to Machine Learning with Apache Ignite](#)
  - [Using Linear Regression with Apache Ignite](#)
  - [Using k-NN Classification with Apache Ignite](#)
  - [Using K-Means Clustering with Apache Ignite](#)
  - [Using Apache Ignite's Machine Learning for Fraud Detection at Scale](#)

# Among Top 5 Apache Projects

Over 1M downloads per year

## Top 5 Developer Mailing Lists

1. **Ignite**
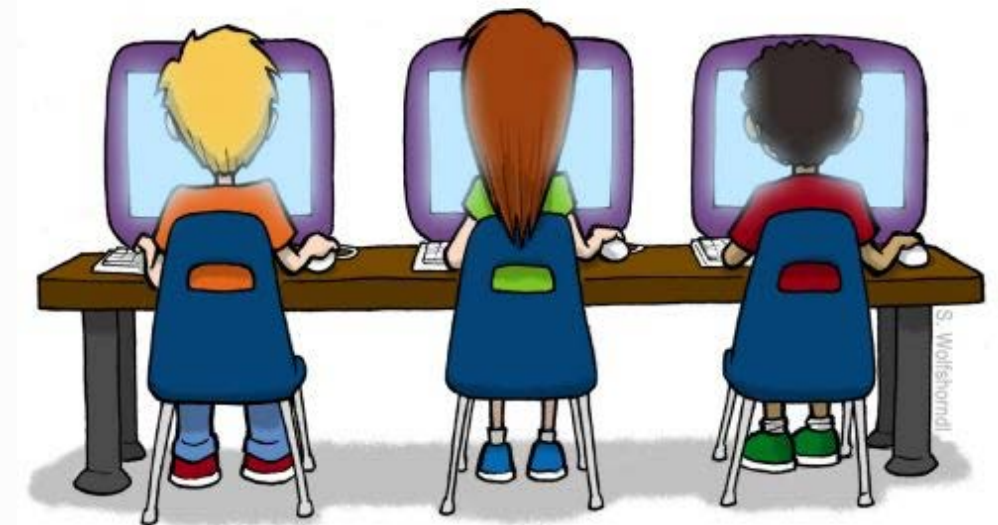2. Kafka
3. Tomcat
4. Beam
5. James

## Top 5 User Mailing Lists

1. Lucene/Solr
2. **Ignite**
3. Flink
4. Kafka
5. Cassandra

## Top 5 by Commits

1. Hadoop
2. Ambari
3. Camel
4. **Ignite**
5. Beam

GridGain

# Apache Ignite – We're Hiring ;)

- Very Active Community

- Great Way to Learn Distributed Computing

- How To Contribute:
  - https://ignite.apache.org/

# Any Questions?

Thank you for joining us. Follow the conversation.
http://ignite.apache.org

@ApacheIgnite
@gridgain
@denismagda