

SystemDS: Introduction and Overview

Rohit Lokwani

Sravan Hande

Tharun Kumar Reddy Karasani



Agenda

- 1 Introduction to SystemDS
- 2 Architecture Discussion
- 3 Installation and Integration
- 4 Performance Benchmarking
- 5 System Demo



SystemDS

SystemDS is an open source ML system for the end-to-end data science lifecycle from data integration, cleaning, and feature engineering, over efficient, local and distributed ML model training, to deployment and serving

- 1
- 2
- 3

Customizability

Algorithm customizability via R-like and Python-like languages

Scalability

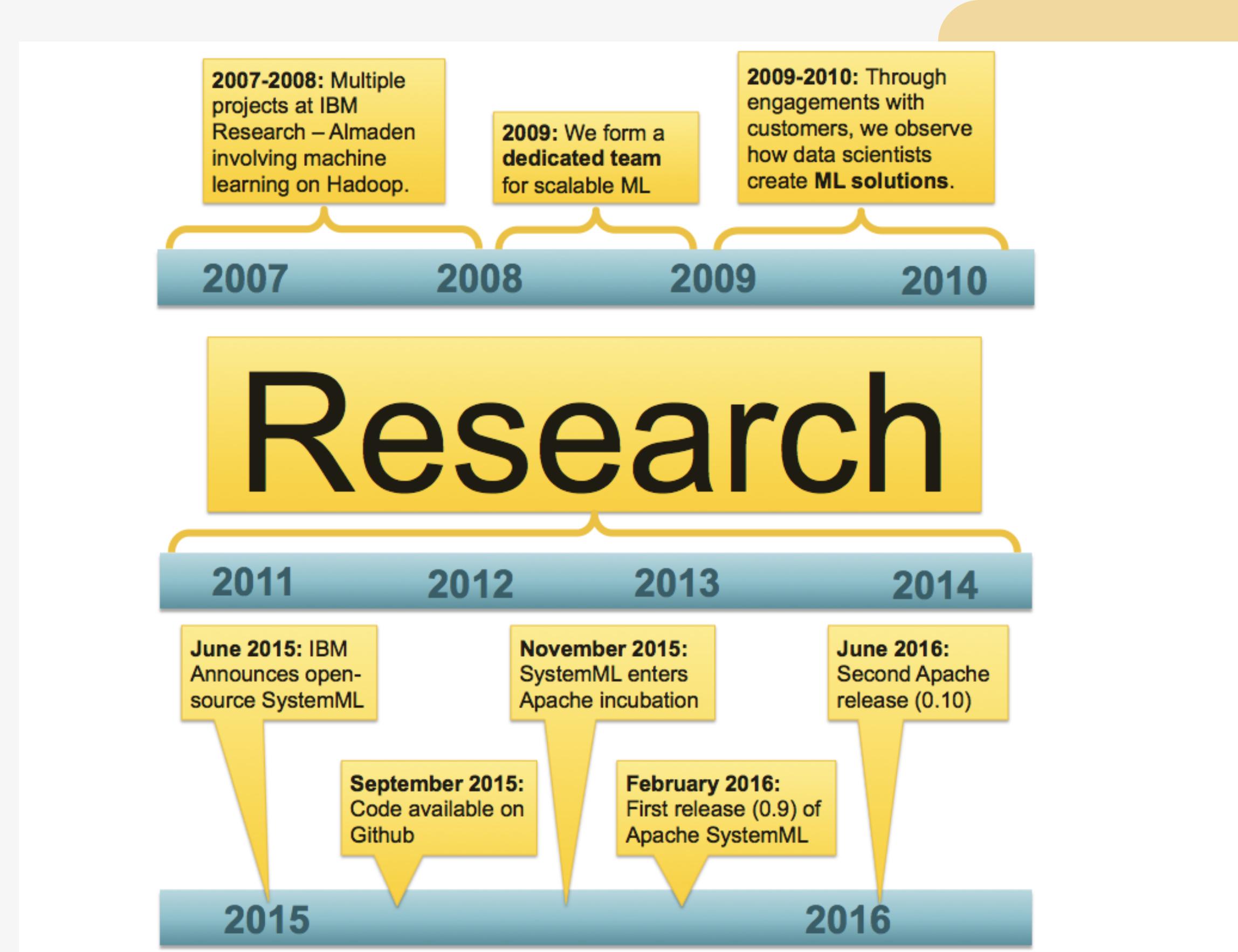
Multiple execution nodes including Spark MLContext, Batch, Standalone

Automatic Optimization

Based on data and cluster characteristics to ensure efficiency and scalability

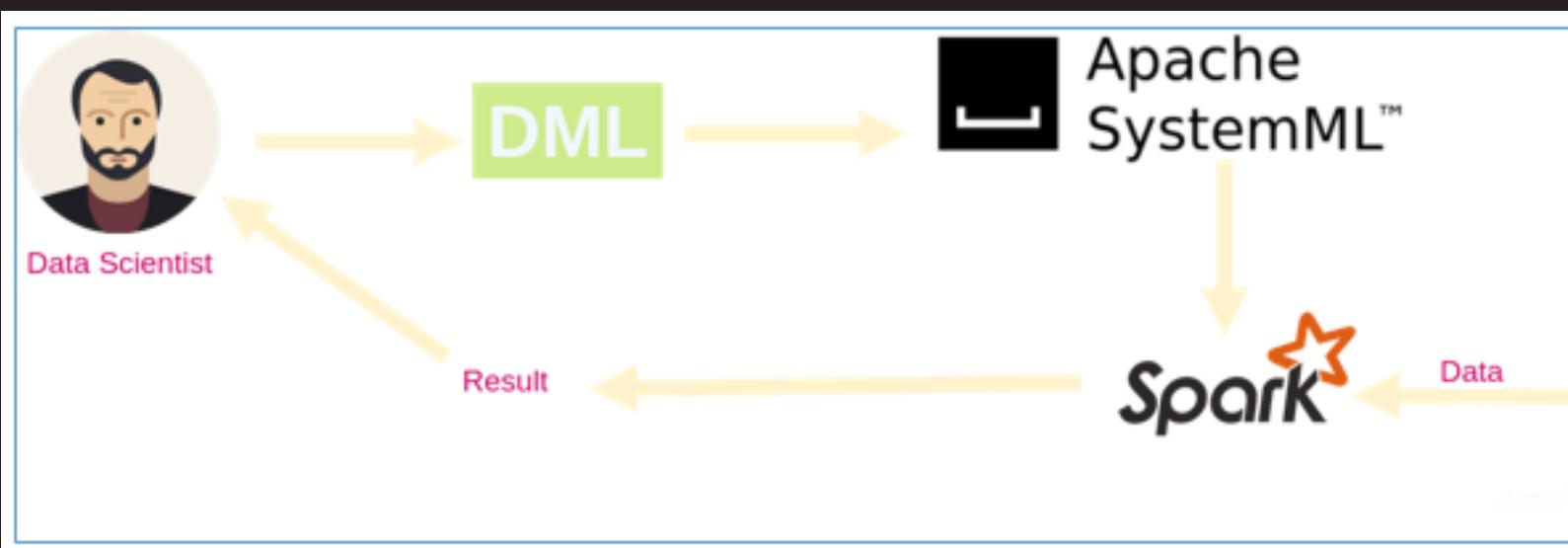
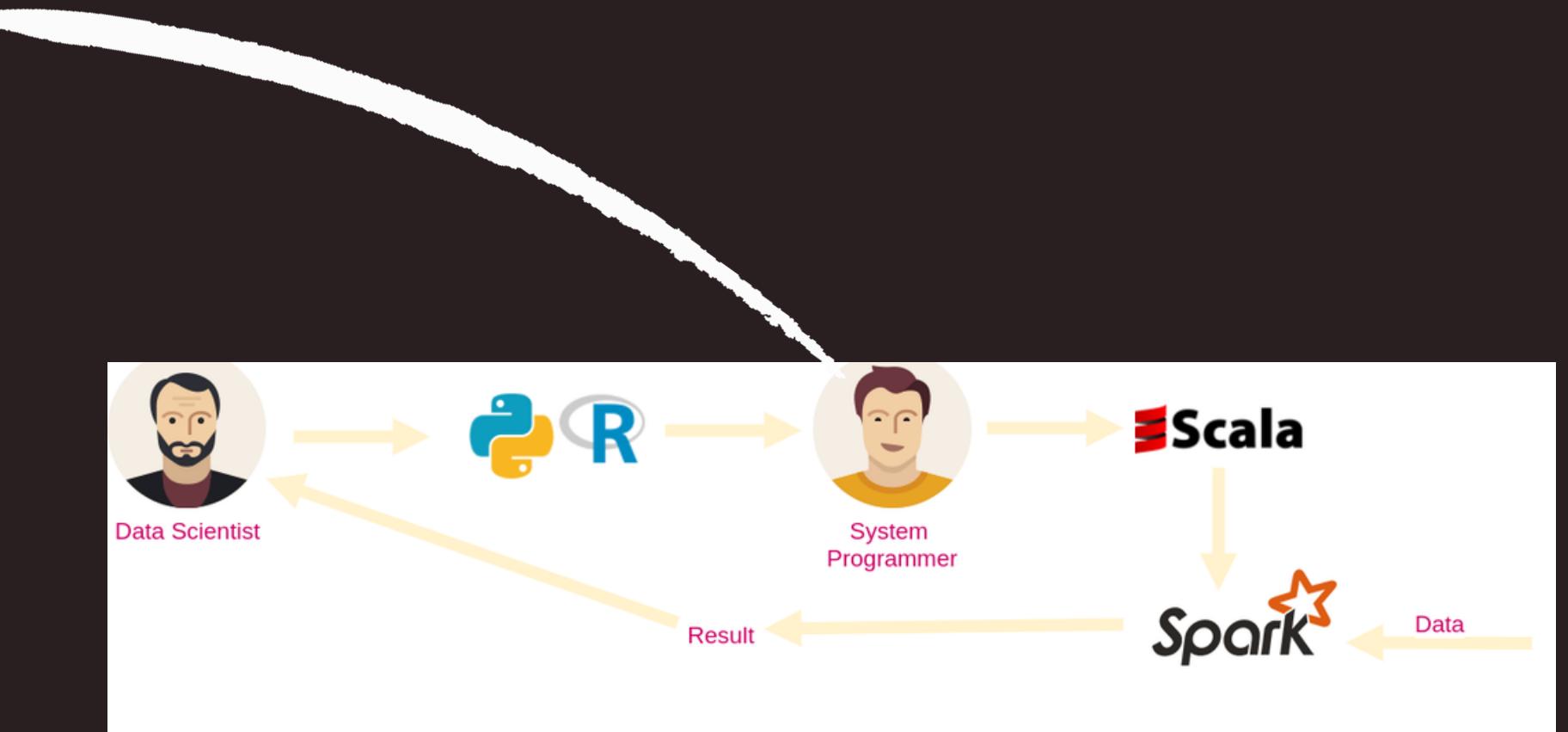
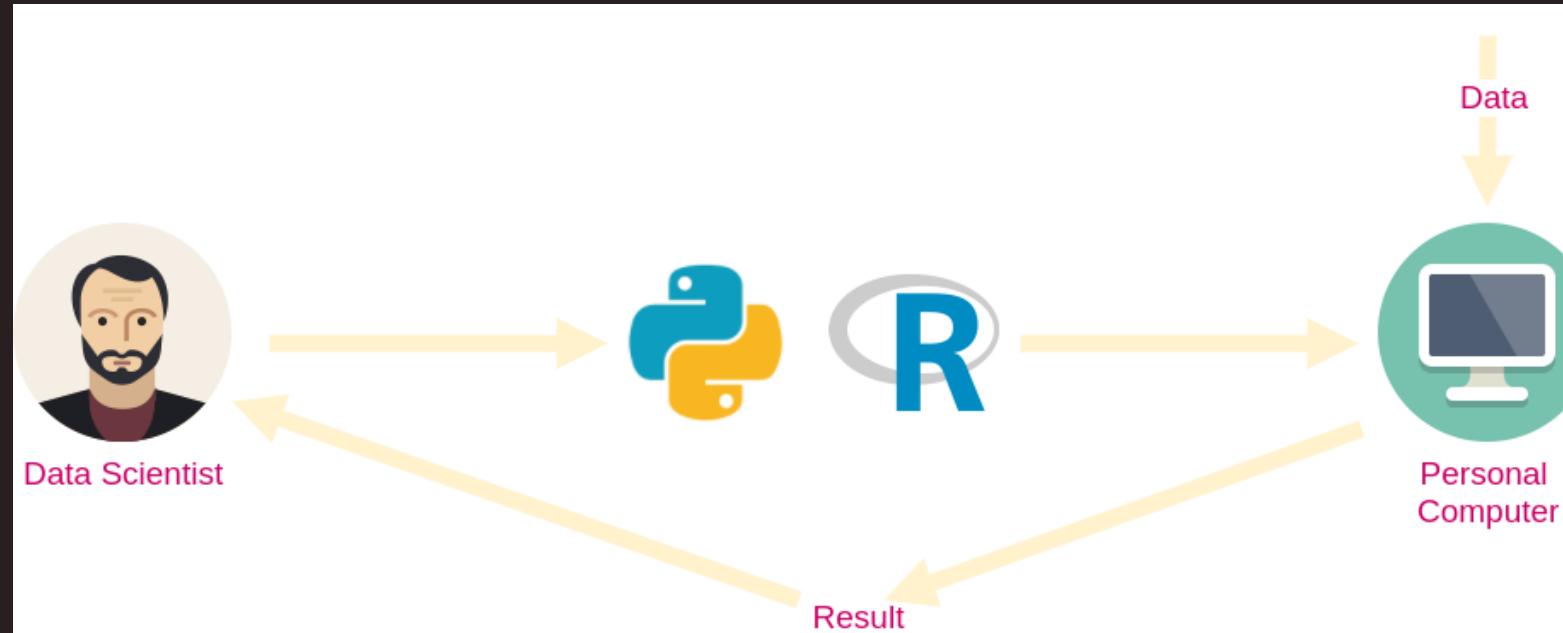


SystemDS Timeline





Efficiency and Scalability

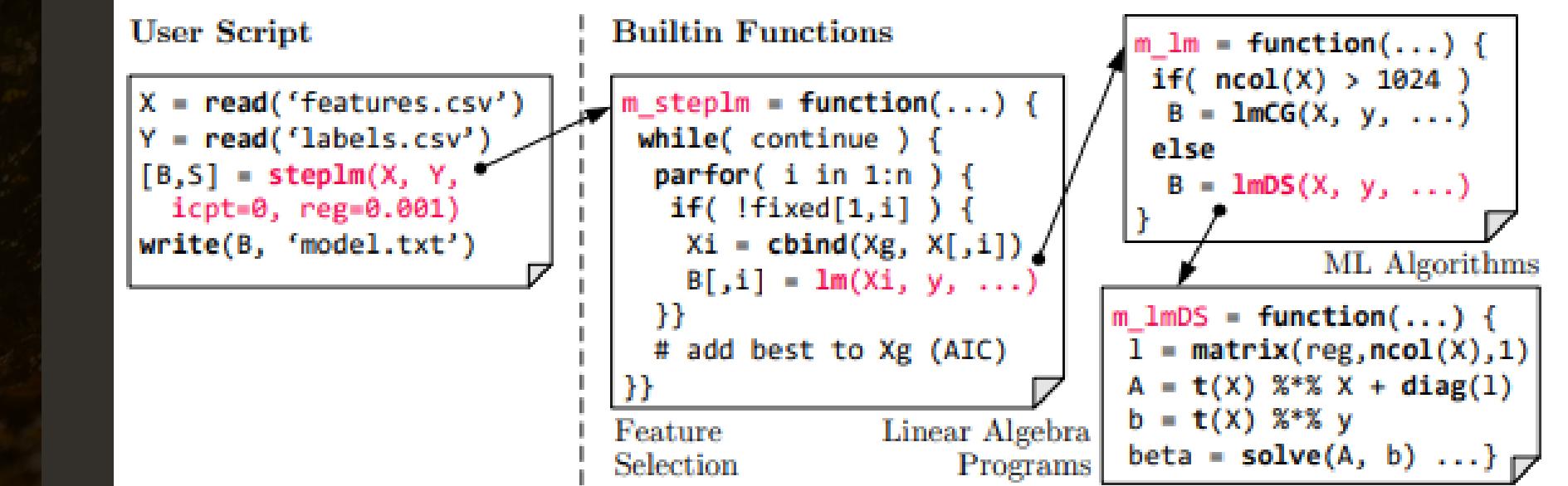
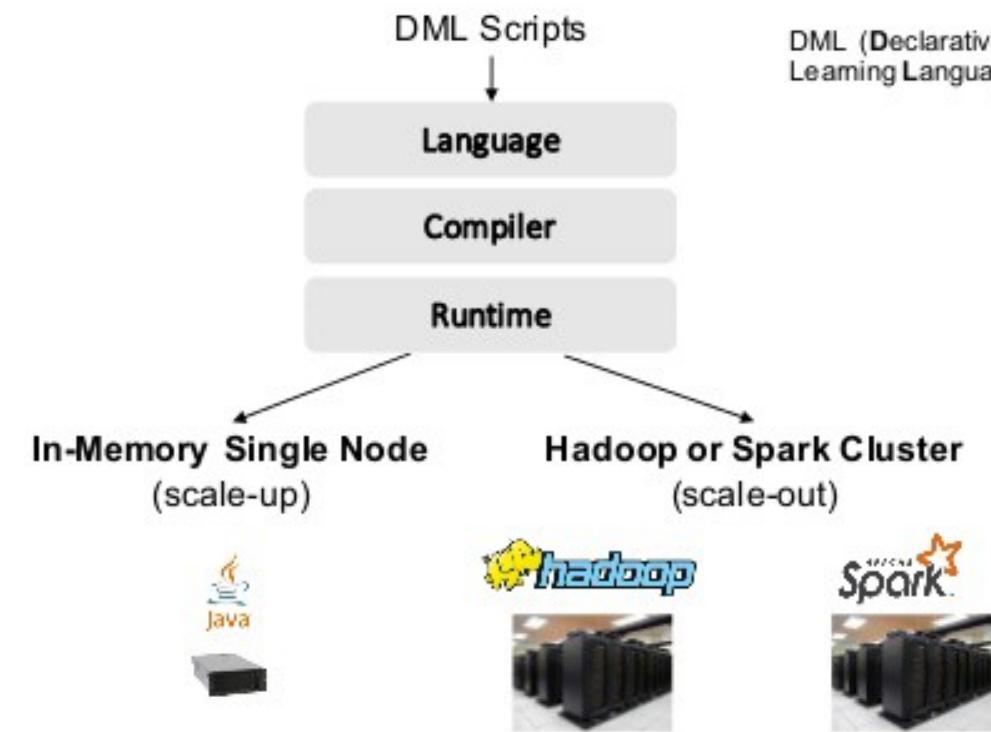




Platform Independence

Ease of programming

High-Level SystemML Architecture





SystemDS + Spark Overview

01 Consistency

Immutable RDDs, use lineage graphs to recover

02 Scalability

Multiple execution nodes

03 Replication

Replicate RDDs with no or some transformations

04 Partitioning

Logical partitioning in RDDs

02

03

04

05

06

05 Ease of programming

Language flexibility, minimum LOC

06 Documentation

Lack of detailed documentation

SystemDS on Spark

Distributed processing

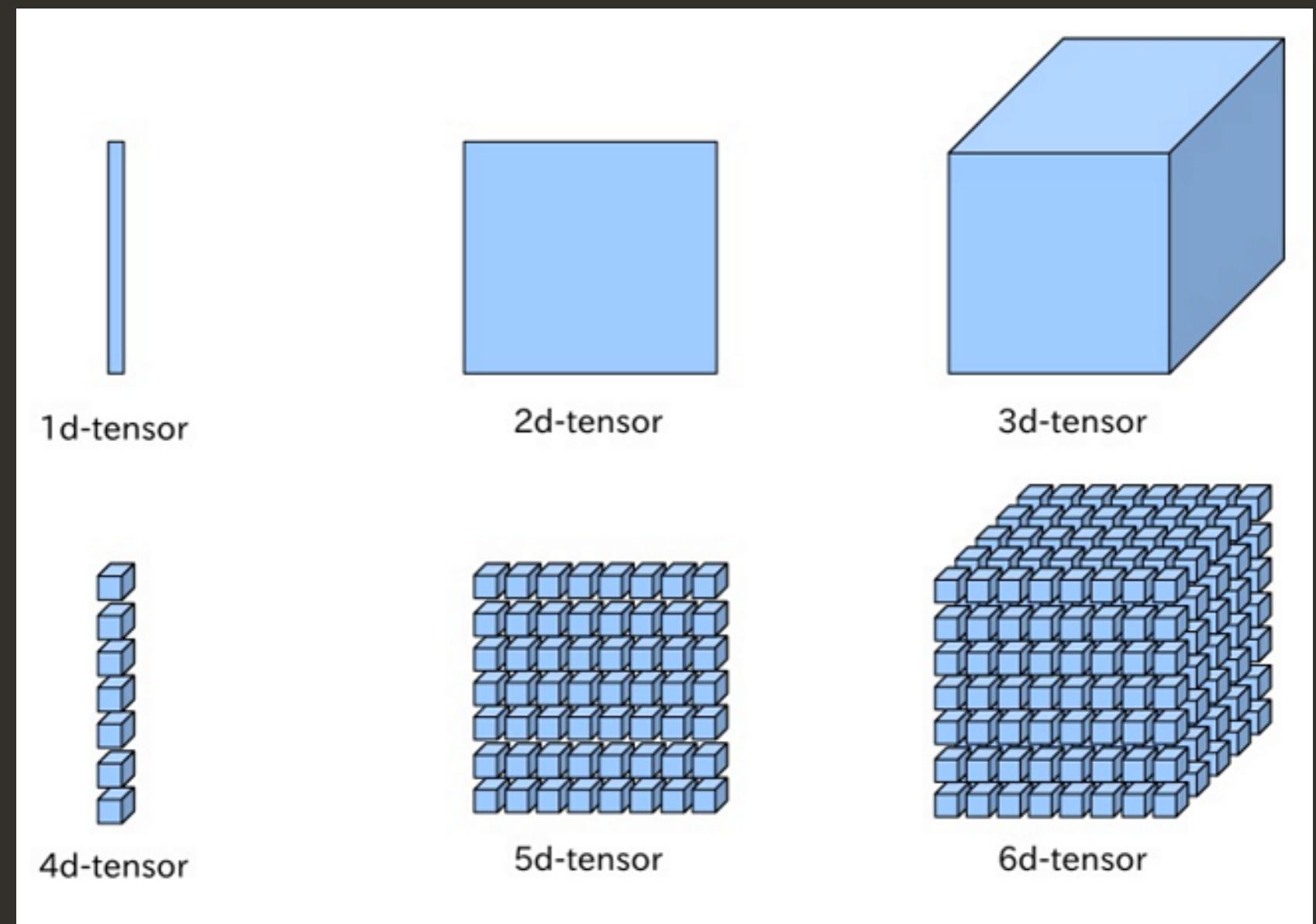
Utilizes spark RDD

In-memory Processing

Write your topic or idea

Data Model

- Basic Tensor Blocks (CPU)
- Data Tensor Blocks (GPU)
- Distributed Tensor Blocks (Spark)
- Federated Tensors



User Interface



CLI and DML

Briefly elaborate on what you want to discuss.



Java API

Briefly elaborate on what you want to discuss.



Python API

Briefly elaborate on what you want to discuss.

Installation

- Dependancies (Java8)
- Python & Jupyter
- SystemDS pip package
- Required GPU libraries



https://systemds.apache.org/docs/2.2.1/api/python/getting_started/install.html

Dataset - Airbnb

- 4 Cities
- Listings, Availability, Reviews
and Neighbourhood data
- 350MB
- Data Analysis and Regression

Backend Support

Local

GPU

Distributed

Discuss about support for SPARK and other distributed frameworks

Federated

Discuss about federated environment

Benchmark Comparison

Baseline Comparison

TensorFlow2 Comparison

System Demo

Demo of SystemDS in distributed and federated environments

Conclusion / Summary

Add a little bit of body text

References

- <https://systemds.apache.org/>
- <https://github.com/apache/systemds>
- https://en.wikipedia.org/wiki/Apache_SystemDS
- <https://apache.github.io/systemds/>
- <https://apache.github.io/systemds/site/algorithms-reference>

Questions?

Thank You!