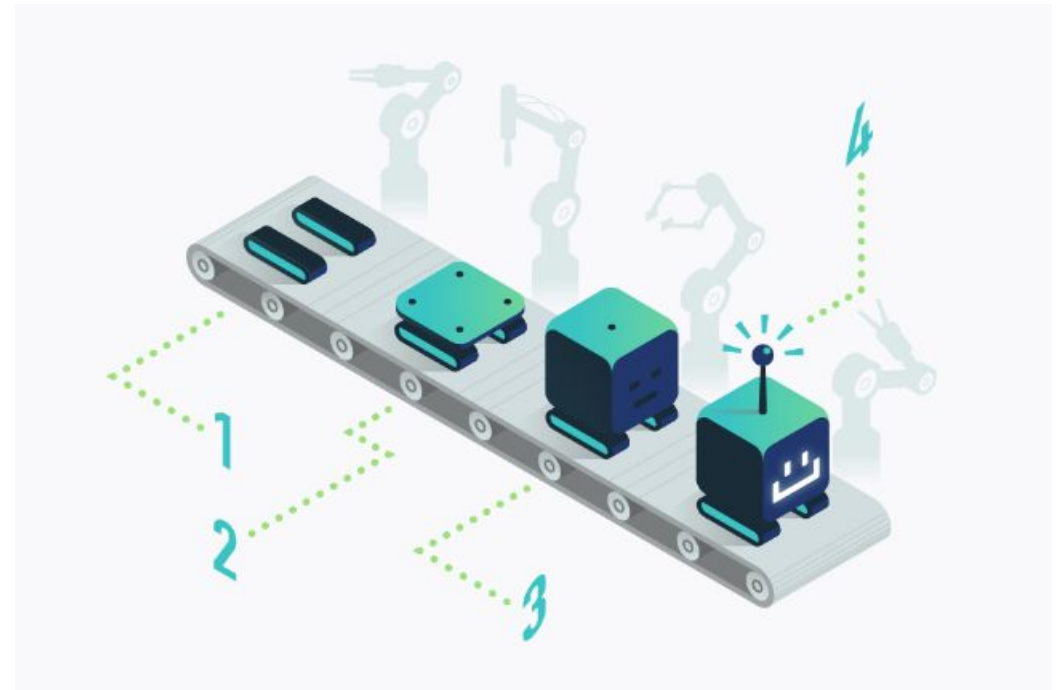# SystemML:
# Declarative Machine Learning on Spark

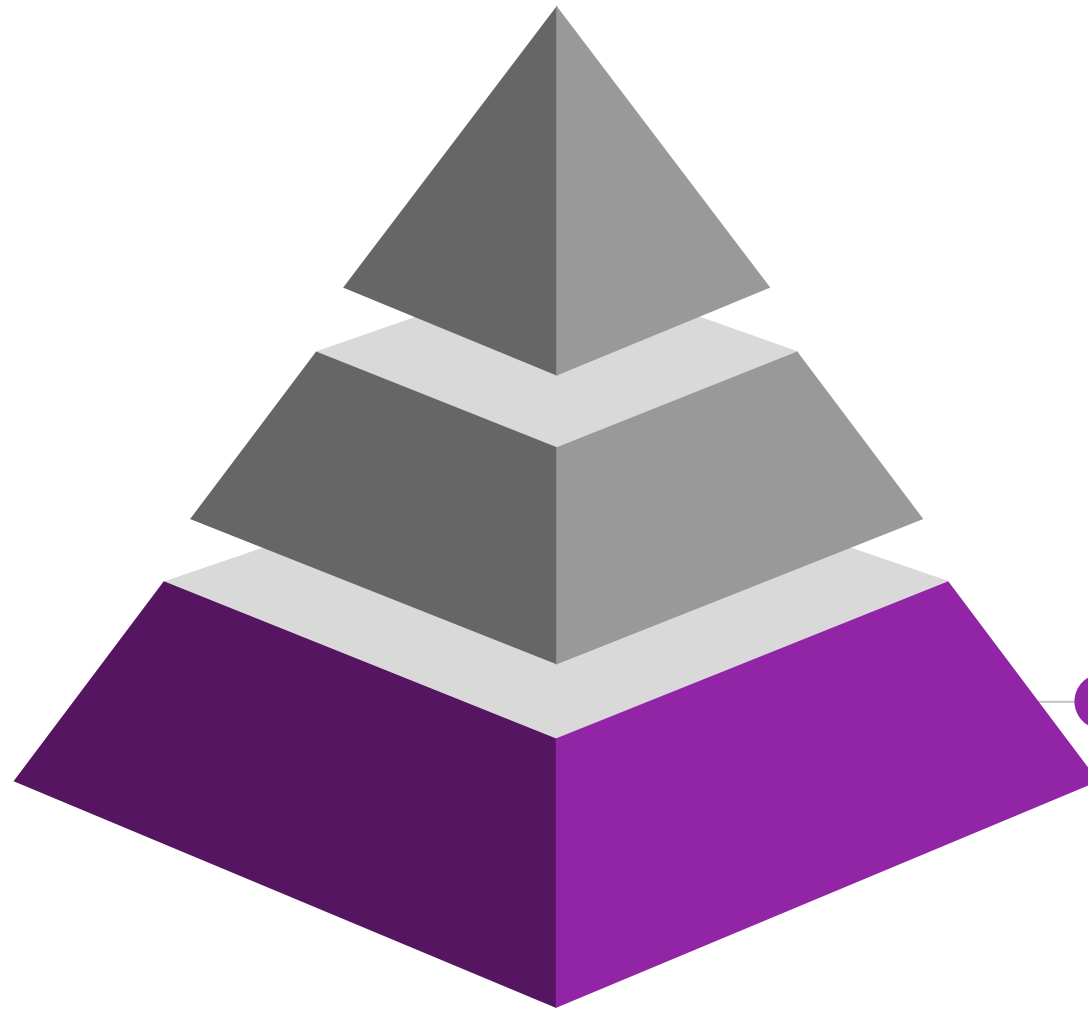05/03/19

Presented by: Juan Carrillo
Candidate for MASc. in Computer Software
Department of Electrical & Computer Engineering
University of Waterloo

UNIVERSITY OF
WATERLOO

# Agenda

1. Introduction

2. SystemML core features

3. Experiments

4. Conclusions

5. Discussion

UNIVERSITY OF
**WATERLOO**

**1** **Introduction**

UNIVERSITY OF
**WATERLOO**

1.  Introduction

# Machine Learning for Big Data Analytics

UNIVERSITY OF
**WATERLOO**

# 1. Introduction

# The problem, and the SystemML approach



**Usual workflow**

Data Scientist → R or Python → Systems Programmer → Scala → Spark → Results

💥 Time consuming
💥 Error prone

**SystemML approach**

Data Scientist → R or Python → DML SystemML → Spark → Results

⚡ Accelerates model development
⚡ Simplifies deployment

Source: Spark Summit. Inside Apache SystemML

SystemML: Declarative Machine Learning on Spark

UNIVERSITY OF WATERLOO

# 1. Introduction

# SystemML background

| **2010** | **2015** | 2017 | 2018 |
|---|---|---|---|
| **Creation** | **Open-source** | **Top Level Project** | **Current release 1.2** |
| **By researchers at the IBM Almaden Research Center** | **Spark Summit in San Francisco** | Apache Software Foundation Board | Deep learning functions Ultra-sparse data |

SystemML: Declarative Machine Learning on Spark

UNIVERSITY OF
**WATERLOO**

**SystemML core features** ②

SystemML: Declarative Machine Learning on Spark

UNIVERSITY OF WATERLOO

2. SystemML core features

# Optimizer integration



Source: Spark Summit. Inside Apache SystemML

UNIVERSITY OF
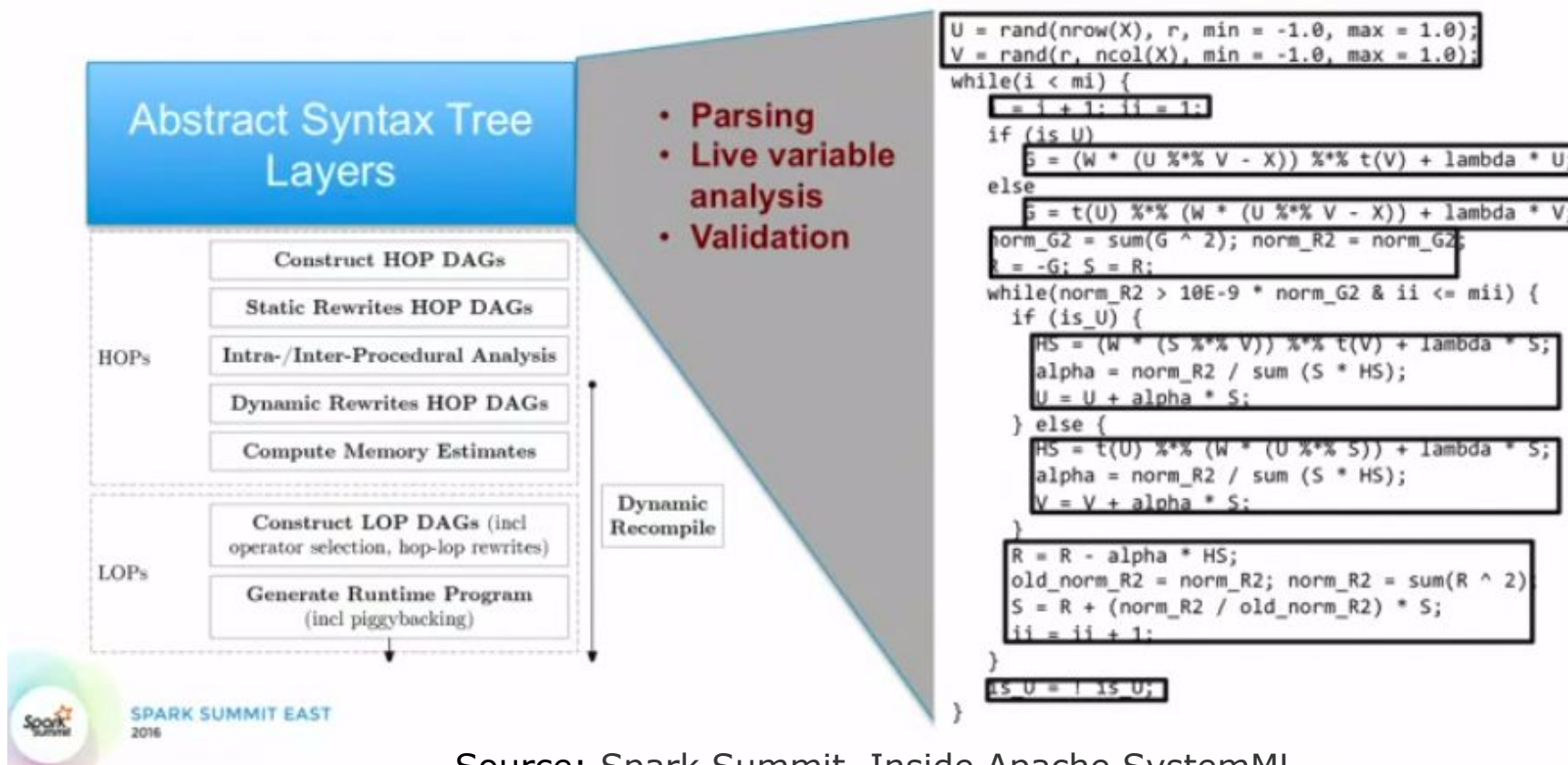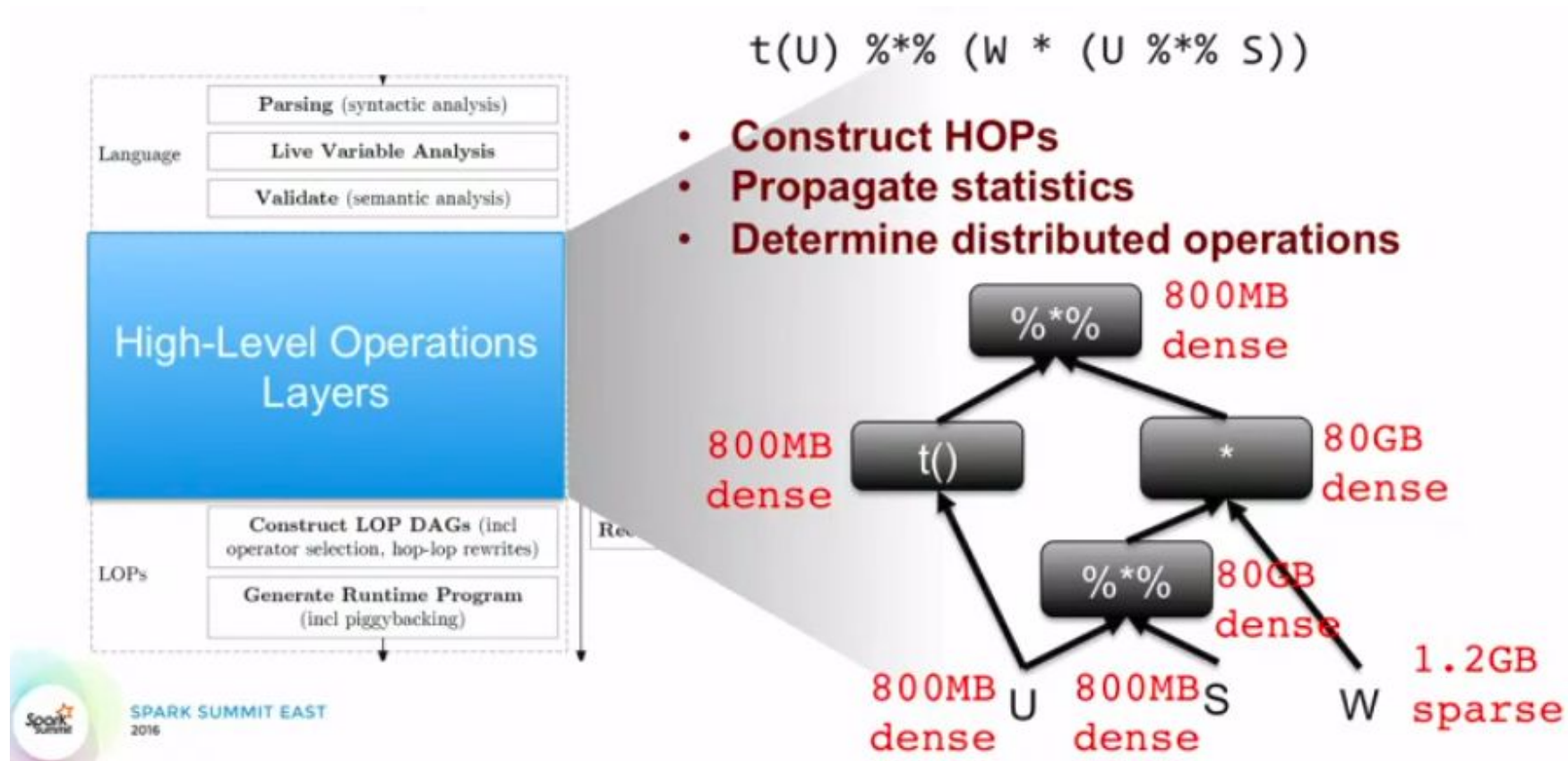WATERLOO

2. SystemML core features

# Optimizer integration



Source: Spark Summit. Inside Apache SystemML

## 2. SystemML core features

# Optimizer integration
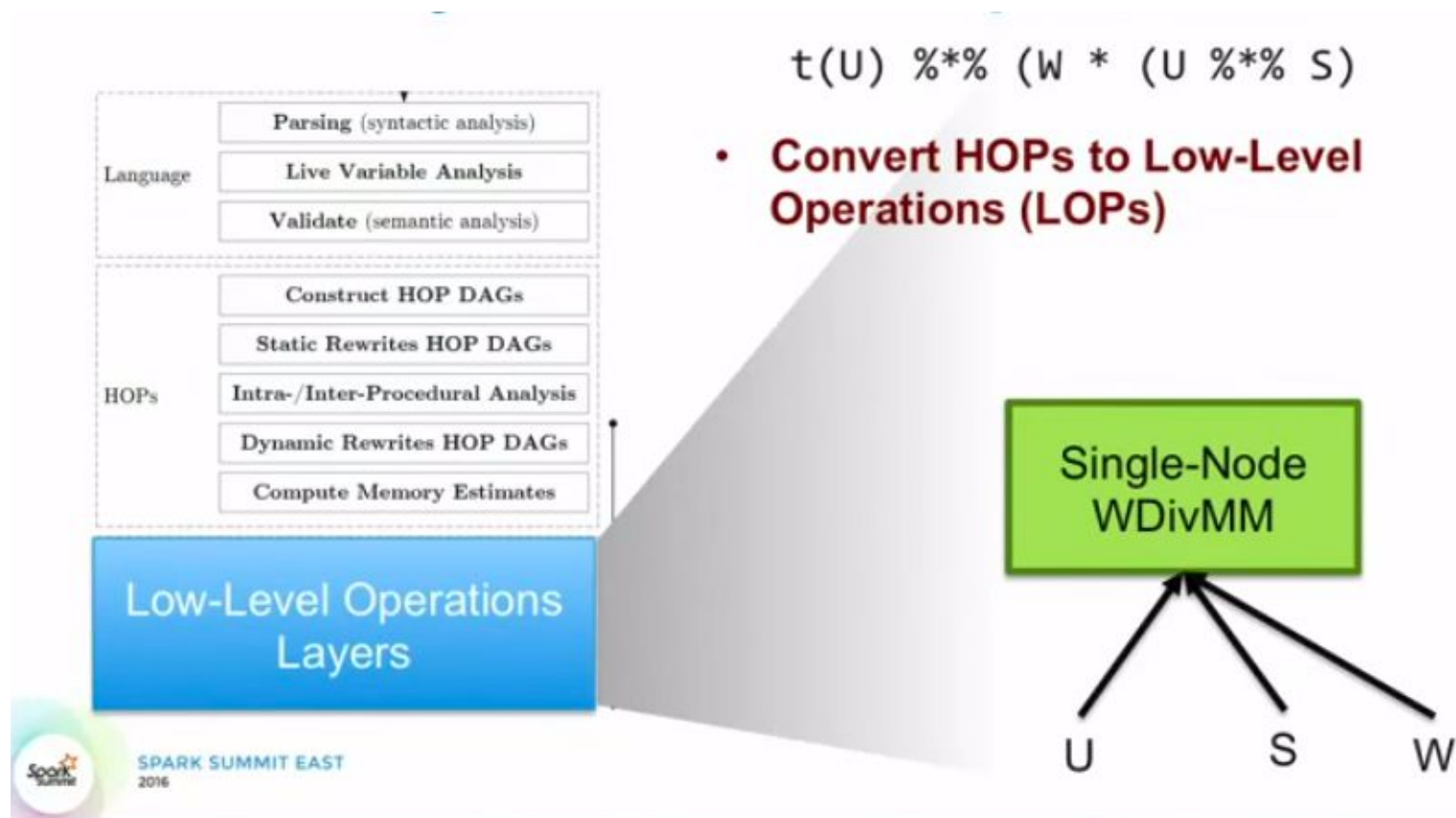


$$t(U) \;\%*\%\; (W * (U \;\%*\%\; S))$$

- **Convert HOPs to Low-Level Operations (LOPs)**

Language
- Parsing (syntactic analysis)
- Live Variable Analysis
- Validate (semantic analysis)

HOPs
- Construct HOP DAGs
- Static Rewrites HOP DAGs
- Intra-/Inter-Procedural Analysis
- Dynamic Rewrites HOP DAGs
- Compute Memory Estimates

Low-Level Operations Layers

Single-Node WDivMM

U     S     W

SPARK SUMMIT EAST 2016

Source: Spark Summit. Inside Apache SystemML

UNIVERSITY OF WATERLOO

## 2. SystemML core features

# Runtime integration
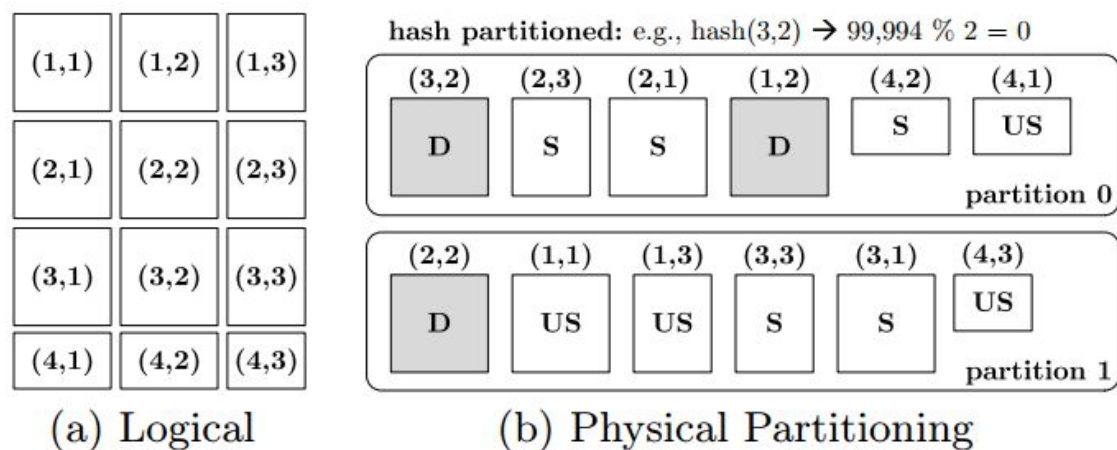
## Distributed Matrix Representation



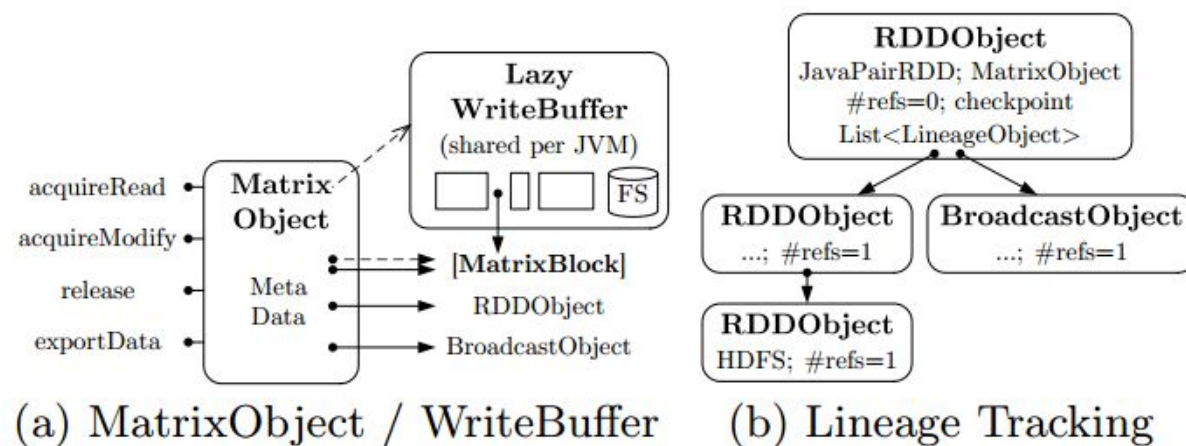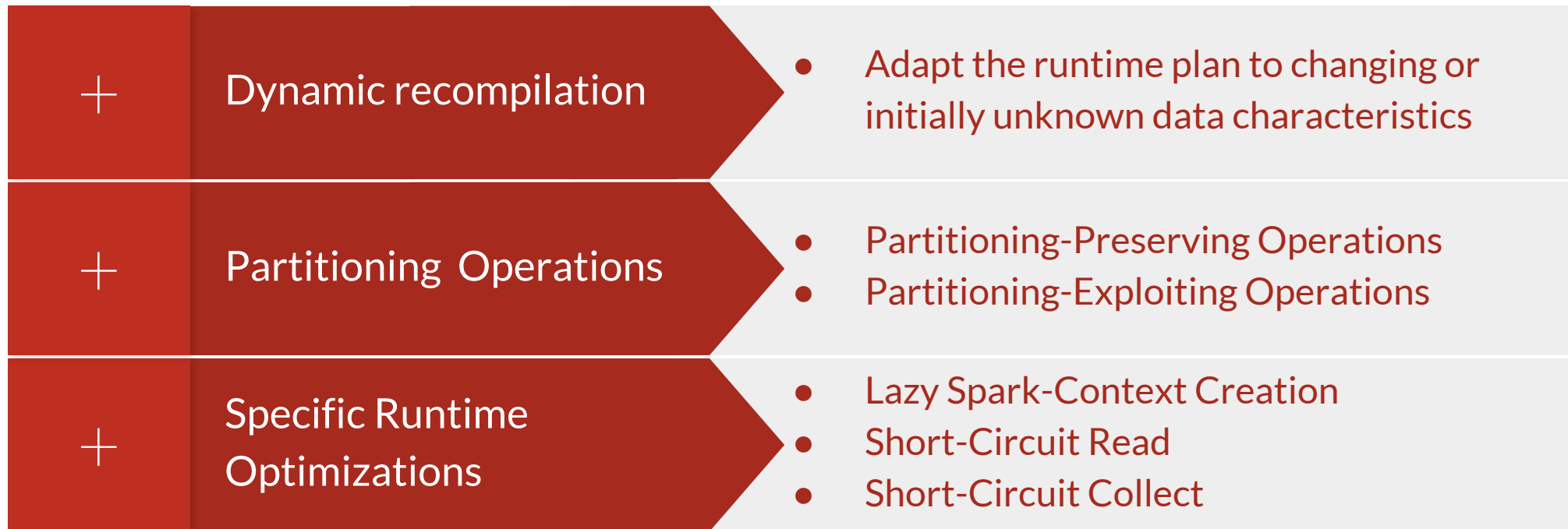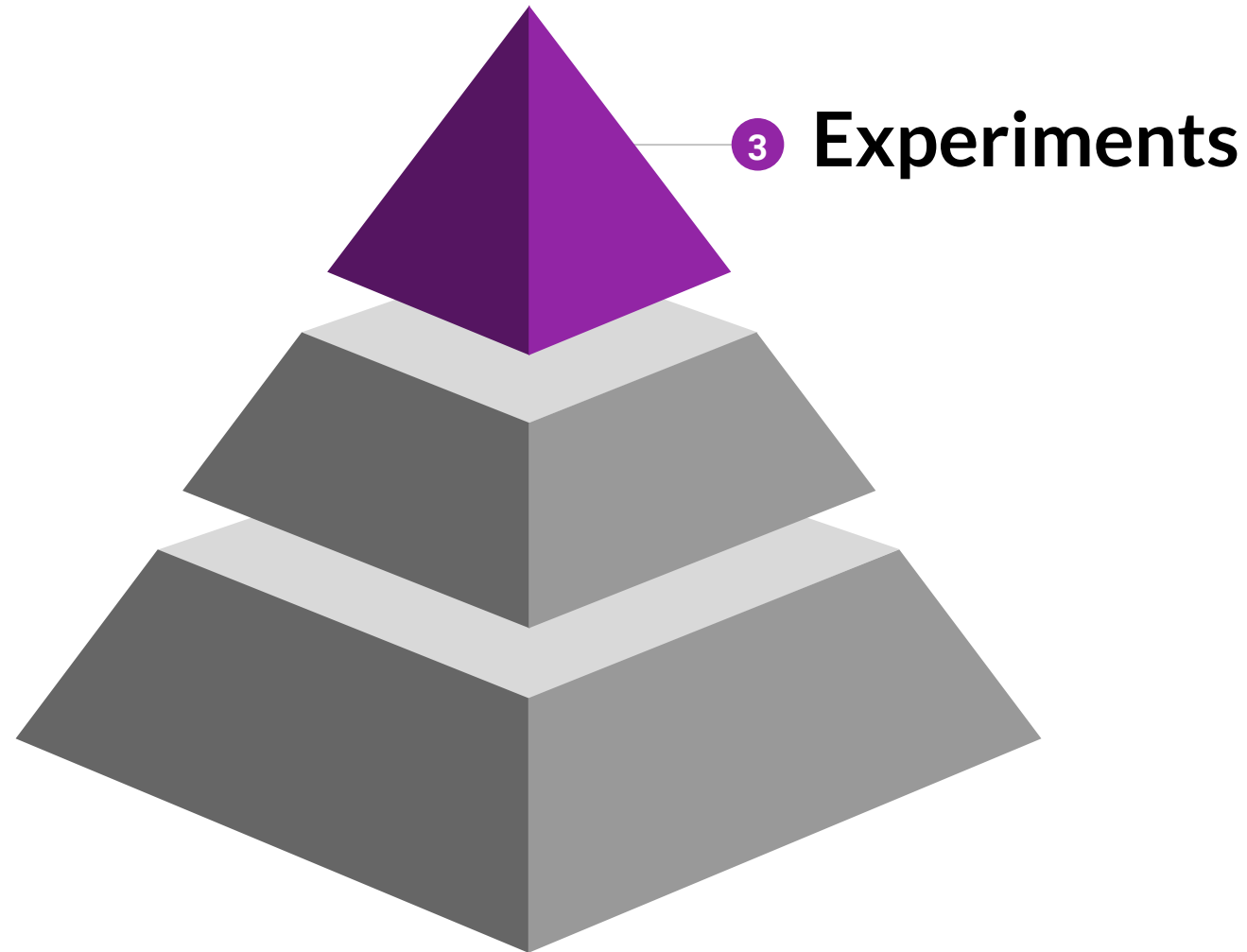Figure 2: Distributed Matrix Representation.

## Buffer Pool Integration



Figure 3: Buffer Pool Integration.

UNIVERSITY OF
WATERLOO

2. SystemML core features

# Runtime integration

| | Dynamic recompilation | • Adapt the runtime plan to changing or initially unknown data characteristics |
|---|---|---|
| + | | |
| + | Partitioning Operations | • Partitioning-Preserving Operations<br>• Partitioning-Exploiting Operations |
| + | Specific Runtime Optimizations | • Lazy Spark-Context Creation<br>• Short-Circuit Read<br>• Short-Circuit Collect |

SystemML: Declarative Machine Learning on Spark

UNIVERSITY OF
WATERLOO

**3 Experiments**

UNIVERSITY OF
**WATERLOO**
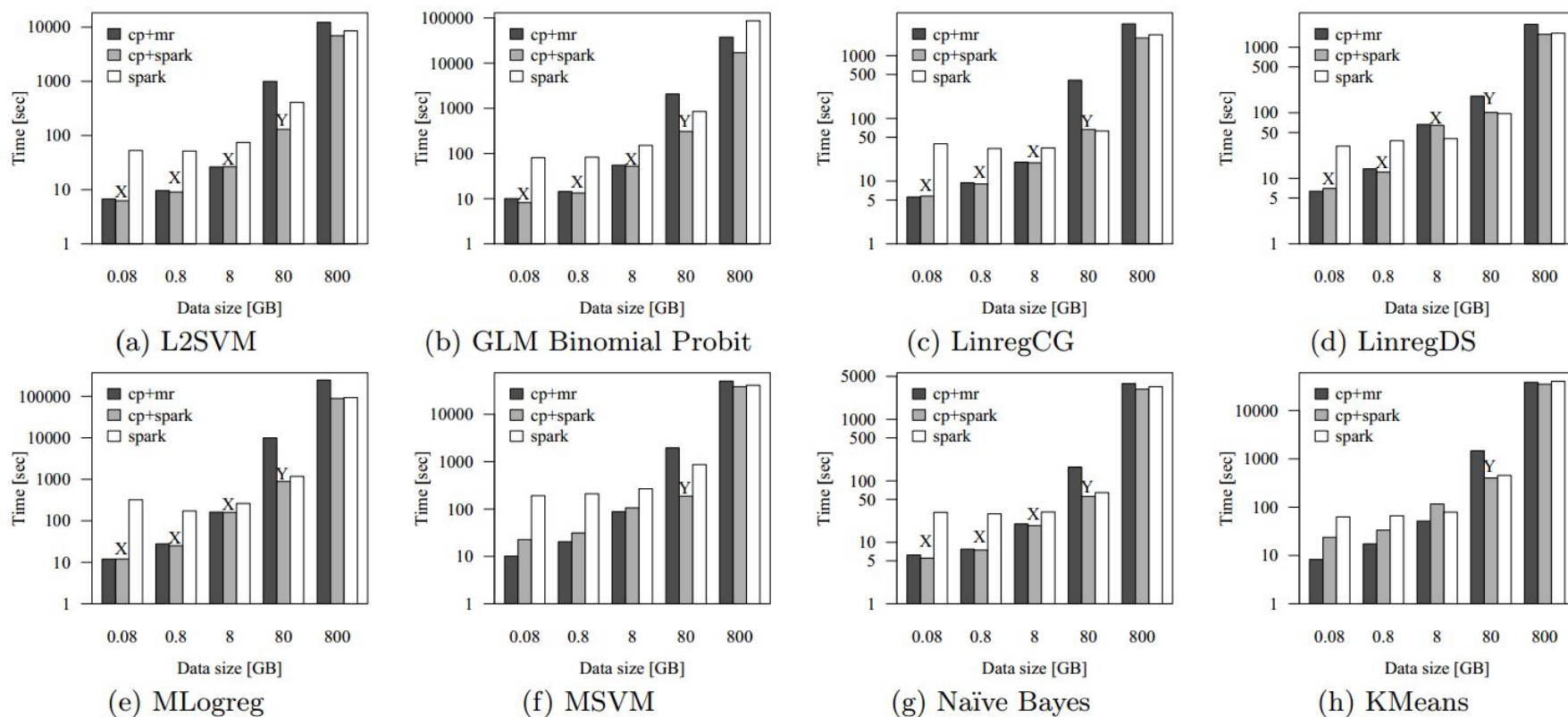
## 3. Experiments

# End-to-End Performance



Figure 4: End-to-End Performance of Different Algorithms with Different Execution Modes.
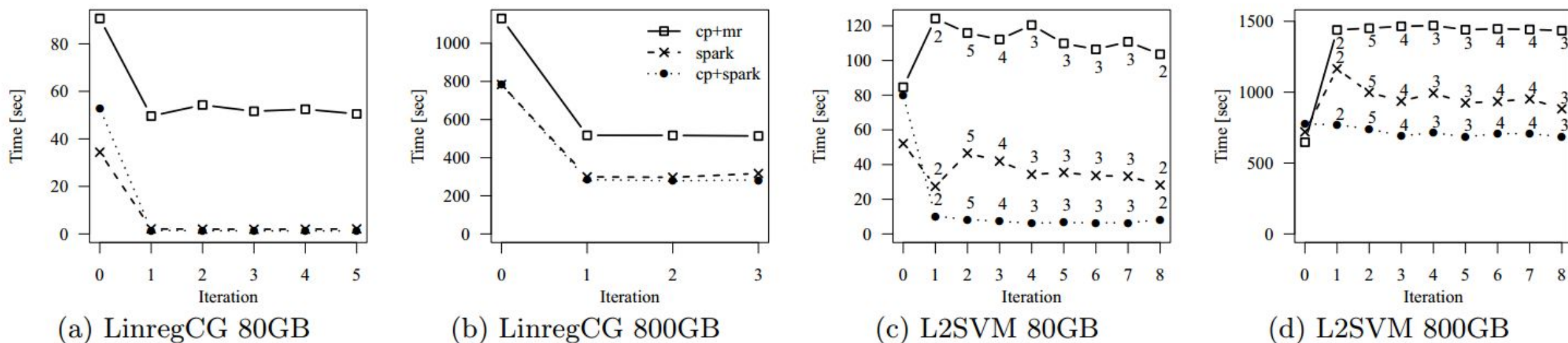
3. Experiments

# Runtime per Iteration



Figure 5: Runtime per Iteration of LinregCG and L2SVM with Different Execution Modes.

# Conclusions 4

UNIVERSITY OF
WATERLOO

4. Conclusions

## Takeaways and paper contributions

✓ Importance of DML as a high-level language to improve interoperability and scalability of Machine Learning models on Spark

✓ Multiple layers of abstraction and optimizations make SystemML a powerful tool for accelerating the development of Machine Learning models over Big Data

✓ Experimental evaluation on multiple ML models and datasets

SystemML: Declarative Machine Learning on Spark

UNIVERSITY OF
WATERLOO

# Thanks for your attention

UNIVERSITY OF
**WATERLOO**

# Discussion ⑤

UNIVERSITY OF
**WATERLOO**

5. Discussion

# Research

1. Optimizer. How to optimize ML models over data streams?
2. Runtime. In dynamic recompilation, what could be unknown data characteristics?
3. Experiments. How SystemML might perform for the KNN algorithm?

# Industry

5. Current capabilities compared to other tools such as Numpy, Scikit Learn, or TensorFlow?
6. Adoption in the current ML and Big Data user base?
7. SystemML in Cloud computing infrastructure. Beyond IBM?

SystemML: Declarative Machine Learning on Spark

UNIVERSITY OF
**WATERLOO**