# MLbase: A Distributed Machine-learning System

Tim Kraska[1],
Rean Griffith[2],
Ameet Talwalkar[3],
John Duchi[3],
Michael J. Franklin[3],
Michael Jordan[3]

[1]Brown University [2]VMware [3]AMPLab, UC Berkeley

Introduction
Architecture
Optimization
Conclusion

MLbase
Motivation

# Outline

Introduction
Architecture
Optimization
Conclusion

MLbase
Motivation

# MLbase

- A distributed machine learning system.
- Focus on user experience.
  - End-user
  - ML researchers

Introduction
Architecture
Optimization
Conclusion

MLbase
**Motivation**

## Motivation

- Explosion of abundant data.
- Data in no longer confined to academic researchers.
- Extracting value from such big data is a growing concern.
- The complexity of existing algorithms is overwhelming.
- Layman users may not understand the trade-offs between different learning techniques and parameterization.

Introduction
**Architecture**
Optimization
Conclusion

Workflow
Architecture

# Outline

1. **Introduction**

2. **Architecture**

3. **Optimization**

4. **Conclusion**

Introduction
Architecture
Optimization
Conclusion

Workflow
Architecture

## Workflow

- User specify simple machine learning task.
- Parse into a logical learning plan (LLP).
    - LLP is the most general workflow to perform the ML task.
- An optimizer translate LLP into physical learning plan (PLP).
    - PLP specifies exactly the parameters and the data sets to be used.
- Distribute PLP onto the worker nodes.
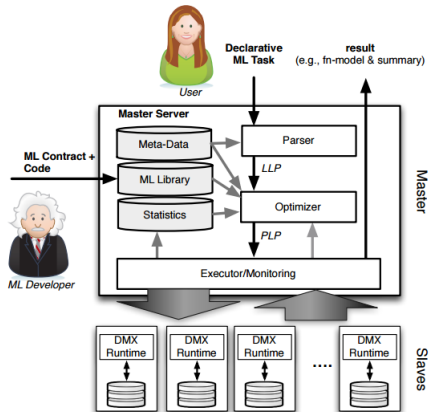- Return a learned model that can be immediately used.

Introduction
Architecture
Optimization
Conclusion

Workflow
Architecture

# Architecture



**Figure 1: MLbase Architecture**
MLbase

Introduction
Architecture
**Optimization**
Conclusion

LLP
PLP
Return result
Good extensibility

# Outline

Introduction
Architecture
**Optimization**
Conclusion

**LLP**
PLP
Return result
Good extensibility

## LLP

- Many operations are mapped 1-to-1 to PLP which can translate earlier.
- Insert down-sample stage.

Introduction
Architecture
**Optimization**
Conclusion

LLP
PLP
Return result
Good extensibility

# LLP



(fn-model, summary)

Introduction
Architecture
**Optimization**
Conclusion
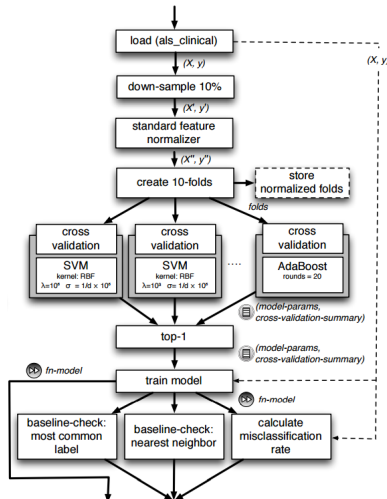
LLP
**PLP**
Return result
Good extensibility

# PLP

- Transform LLP to concrete parameters and learning algorithm.
- Estimate execution time based on statistical models.
  - Normalized data sets for SVM and AdaBoost.
- Use 10-folds cross validation.

Introduction
Architecture
**Optimization**
Conclusion

LLP
PLP
Return result
Good extensibility

# PLP



(3) Optimized Plan

Introduction
Architecture
**Optimization**
Conclusion

LLP
PLP
**Return result**
Good extensibility

# Return first result, find best in the background

- Become more interactive.
- Reduce the risk of stopping too early.

Introduction
Architecture
**Optimization**
Conclusion

LLP
PLP
Return result
**Good extensibility**

# Good extensibility

- A set of high-level operators to enable ML researchers to implement a wide range ML algorithms without deep systems knowledge.

$$\textbf{while } \text{Not}(condition\ for\ completeness) \textbf{ do}$$
$$\theta = U\left(\theta, \frac{1}{|X|} \sum_{x \in X} G(x, \theta)\right)$$
$$\textbf{end while}$$

## Outline

1. **Introduction**

2. **Architecture**

3. **Optimization**

4. **Conclusion**

## Conclusion

- This system aiming to make ML more accessible to non-experts.
- The core of MLbase is its optimizer,which transforms the hight-level ML tasks into the executable ML codes.
- MLbase quickly returns a first quality answer to the user.
- MLbase can improve the result in the background.
- ML developers can constantly add new ML techniques to the system.