# EasyML: Ease the Process of Machine Learning with Data Flow

Jun Xu

Institute of Computing Technology, Chinese Academy of Sciences
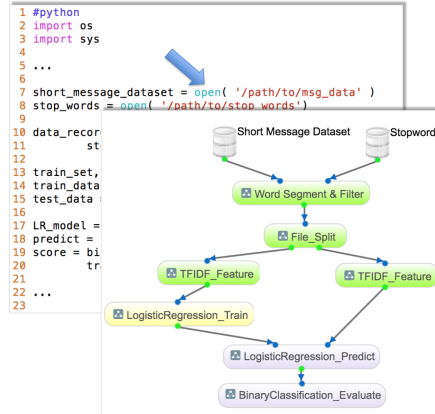Oct. 28, 2017

# Applying Machine Learning is not Easy



Collaborating and sharing

Share the data, algorithms, and experience



UI

Simple UI is helpful



Mobility

Can access the service everywhere

The barrier comes not only from the advanced algorithms, but also from the complex process of using the algorithms!

# Large Scale Machine Learning System@ICT

**Easy ML: interactive graphical UI**

Designer: ML task creation, editing, submitting and management

Monitor: task monitoring, result visualization, and task reusing

Focus of EasyML

**LIB: scalable machine learning algorithms**

Conventional ML algorithms

DL/RL algorithms for ranking & matching

Data pre-processing, ETL, model evaluation …
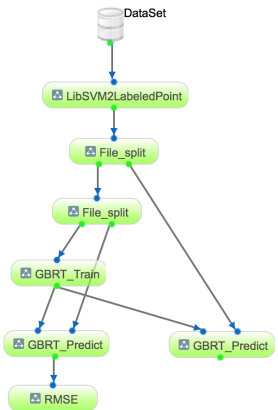
**Distributed Computing**

Map-Reduce

Spark

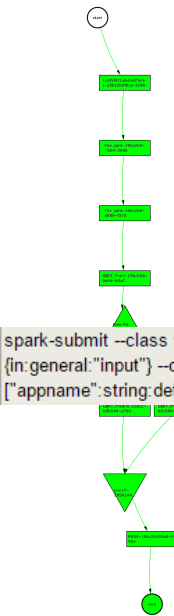TensorFlow

**Data Storage and Management**

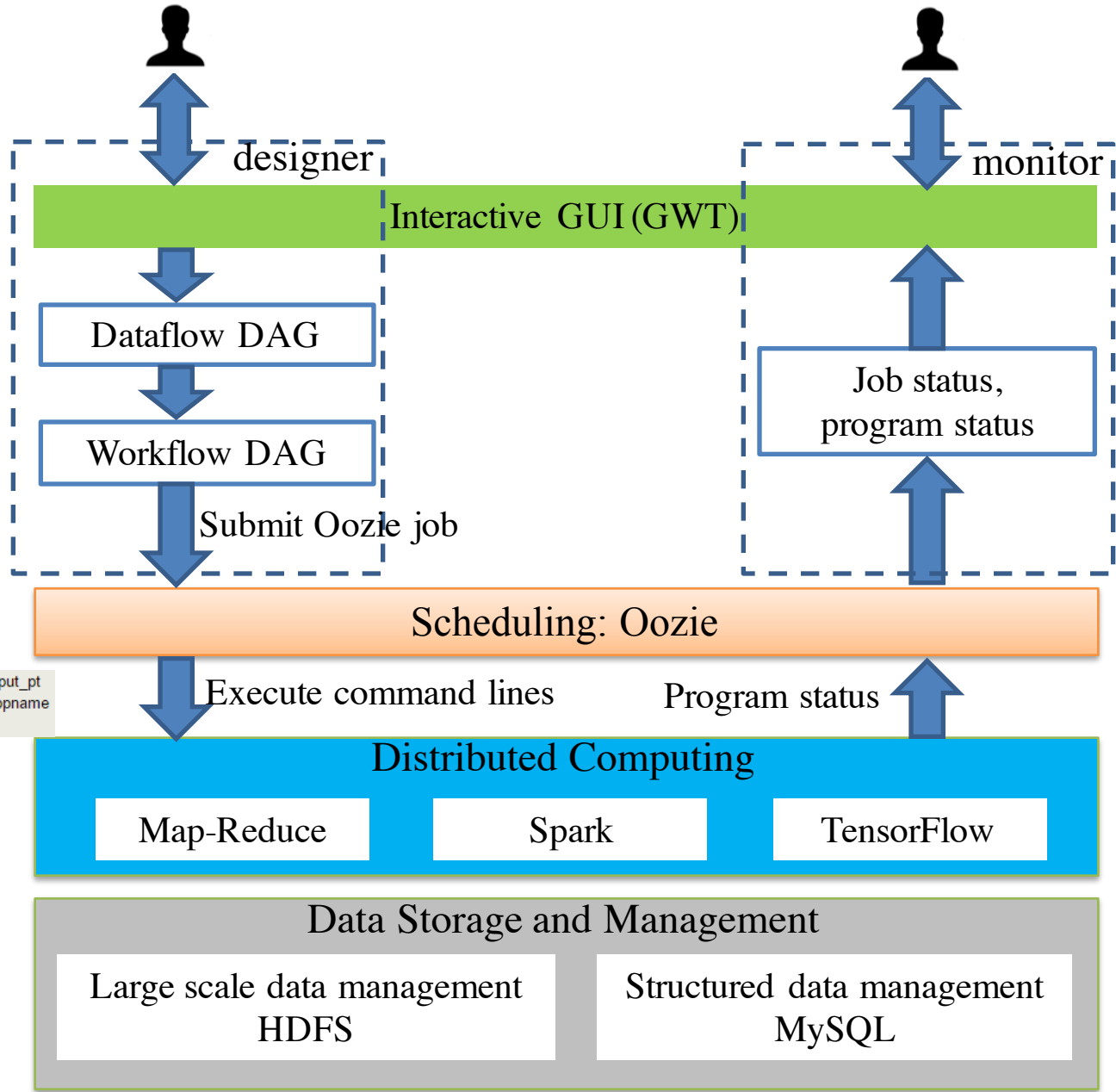Large scale data management HDFS

Structured data management MySQL

# Design of Easy Machine Learning



Node: program / data
Edge: dataflow

```
spark-submit --class word.WordCount wordcount.jar --input_pt
{in:general:"input"} --output_pt {out:general:"output"} --appname
["appname":string:default,"wordcount"]
```

Node: program / start /
end / fork / join
Edge: dependency

designer

monitor

## Interactive GUI (GWT)

Dataflow DAG

Workflow DAG

Submit Oozle job

Job status,
program status

## Scheduling: Oozie

Execute command lines

Program status

## Distributed Computing

| Map-Reduce | Spark | TensorFlow |

## Data Storage and Management

Large scale data management
HDFS

Structured data management
MySQL

# Key Features — Resource Management



Managing programs, data, and tasks

Uploading new algorithms

- Managing the algorithms, data, and tasks
- Uploading algorithms and data

# Key Features — Task Design



- Creating the task DAG (usually by cloning and editing an existing task) with drag-and-drop manner

- Setting the parameters for each node

- Submitting for execution

# Key Features — Task Monitoring



Task/node status monitoring

Data/results visualization

- Monitoring status of tasks and nodes
- Checking / downloading the outputs
- Visualizing the data / models

# Key Features — Task Reusing



- Editing (appending nodes, deleting nodes, and changing parameters) and re-submiting

# Deploy as Web Service
## http://159.226.40.104:18080/dev



Brower          EasyML service          Hadoop/Spark/TensorFlow  cluster

- Advantages
  - **Sharing**: share data/programs/tasks among users
  - **Collaborating**: working together for one task
  - **Mobility**: accessing with web browsers anywhere
  - **Open**:  ETL for data import/export; can run third-party algorithms

# Source Shared at Github

## https://github.com/ICT-BDA/EasyML

- Top 1 Java project at Github trending for one week

- 1400 + stars and ~300 forks

- CIKM 2016 best demo candidate
  [Guo et al., CIKM '16]



Trending in open source
See what the GitHub community is most excited about this

| Repositories | Developers | | Trending: this week |

**ICT-BDA / EasyML**
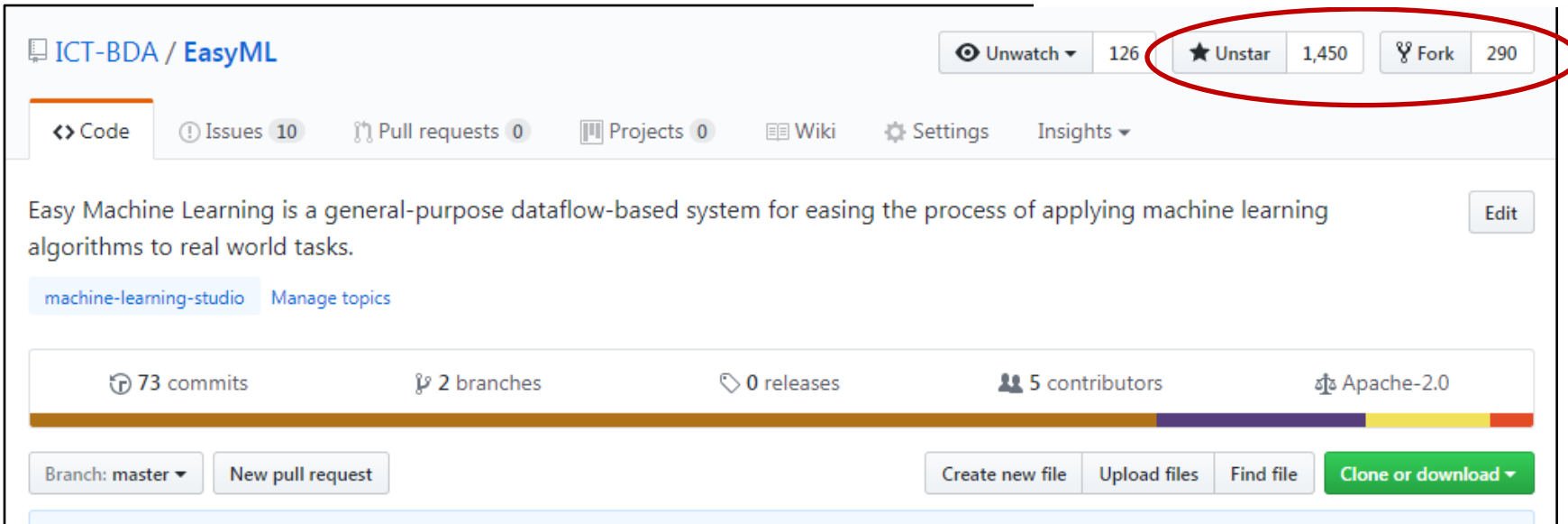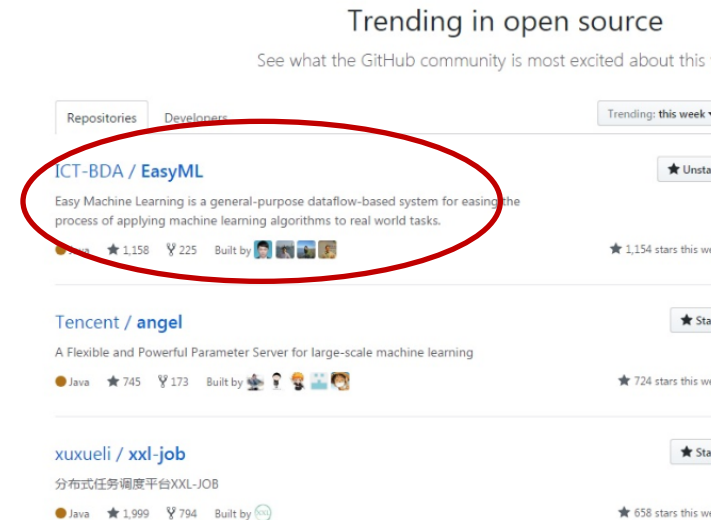Easy Machine Learning is a general-purpose dataflow-based system for easing the process of applying machine learning algorithms to real world tasks.
Java ★ 1,158 ⑂ 225 Built by
★ 1,154 stars this week

**Tencent / angel**
A Flexible and Powerful Parameter Server for large-scale machine learning
Java ★ 745 ⑂ 173 Built by
★ 724 stars this we

**xuxueli / xxl-job**
分布式任务调度平台XXL-JOB
Java ★ 1,999 ⑂ 794 Built by
★ 658 stars this we

---

📕 ICT-BDA / **EasyML**

⊙ Unwatch ▾ 126 | ★ Unstar 1,450 | ⑂ Fork 290

<> Code | ⊙ Issues 10 | ⋔ Pull requests 0 | ▥ Projects 0 | ▦ Wiki | ⚙ Settings | Insights ▾

Easy Machine Learning is a general-purpose dataflow-based system for easing the process of applying machine learning algorithms to real world tasks.

Edit

machine-learning-studio    Manage topics

| ⊙ 73 commits | ⑂ 2 branches | ◌ 0 releases | ⚊ 5 contributors | ⚖ Apache-2.0 |

Branch: master ▾ | New pull request | | Create new file | Upload files | Find file | Clone or download ▾

# Related Systems

- Stanford CodaLab
  - A collaborative platform for reproducible research
- Rapid Miner Studio
  - Interactive GUI and integrated machine learning algorithms
- Microsoft Azure machine learning
  - GUI-based IDE for constructing and operationalizing machine learning workflow on Azure
- Alibaba 御膳房 / DT PAI
- The Fourth Paradigm Prophet

# Summary

- Ease the process of using machine learning
  - Machine learning process as dataflow DAG
  - Interactive GUI for designing and managing ML tasks
  - Deployed as a web service
    - Resource management
    - Task design
    - Task monitoring
    - Result reusing
  - Source code @Github
    https://github.com/ICT-BDA/EasyML

# EasyML Team

- Xinjie Chen, ICT CAS
- Zhaohui Li, ICT CAS
- Tianyou Guo, Sogou Inc.
- Jianpeng Hou, Google China
- Ping Li, Tencent Wechat
- Jiashuo Cao, CUIT
- Dong Huang, UCAS
- Xiaohui Yan
- Xueqi Cheng, ICT CAS

# Thanks!

junxu@ict.ac.cn

http://www.bigdatalab.ac.cn/~junxu