# Productionizing H2O Models with Apache Spark

Jakub Háva, Michal Malohlava; H2O.ai

**#ML4SAIS**

# Who are we?

- **Michal**
  - Chief Architect of Platforms at H2O.ai
  - Creator of Sparkling Water
  - Ph.D at Charles University (CZ), PostDoc at Purdue Uni (US)
- **Kuba**
  - Senior Software engineer at H2O.ai - Core Sparkling Water
  - Master's at Charles University (CZ)
  - Implemented high-performance cluster monitoring tool for JVM based languages (JNI, JVMTI, instrumentation)
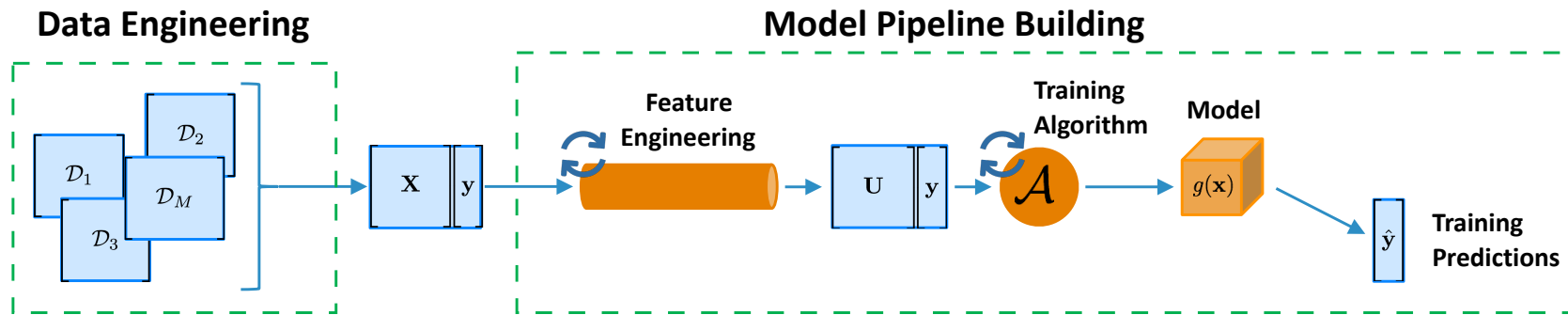
# Machine Learning (ML) Lifecycle
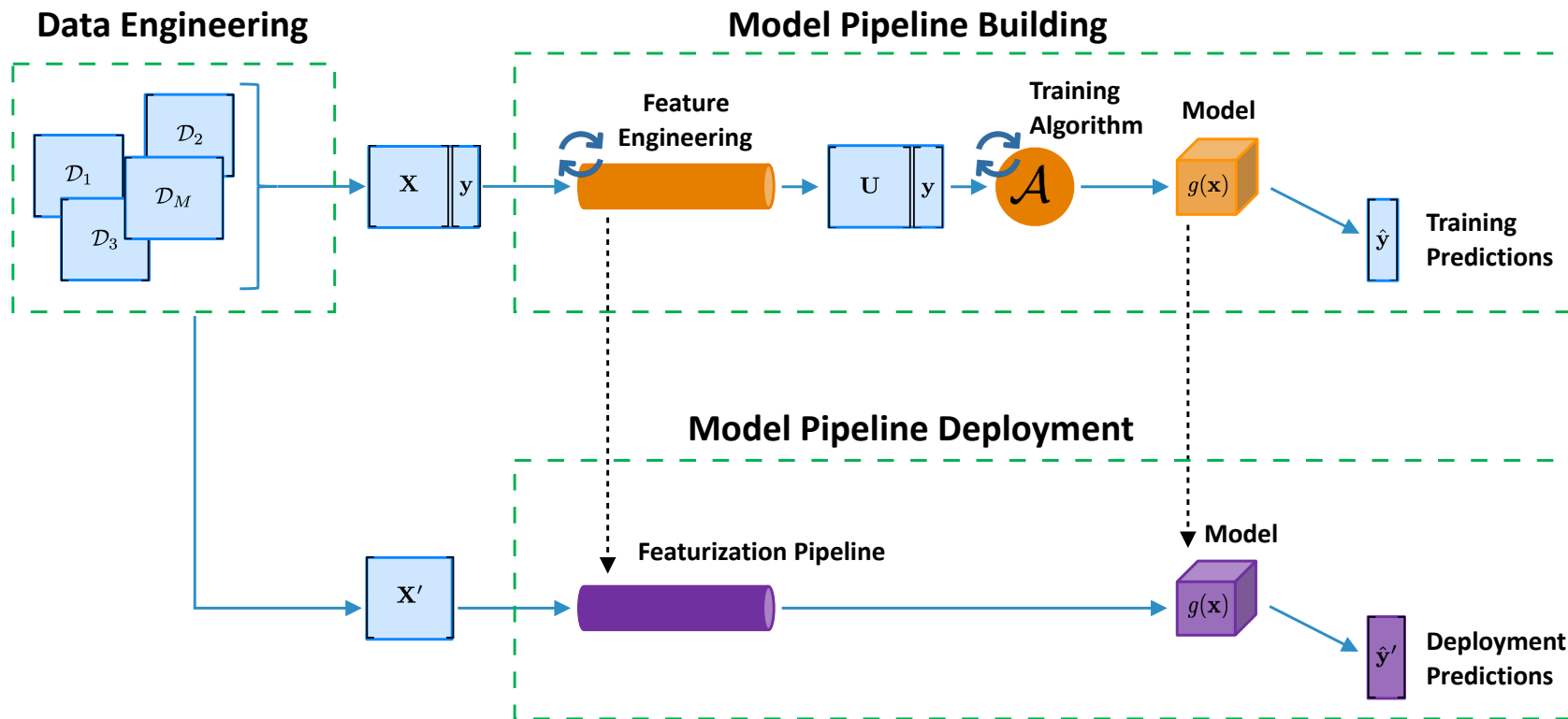
#ML4SAIS

# Basic ML Lifecycle

**Data Engineering**

**Model Pipeline Building**

# Basic ML Lifecycle

# Example Implementations

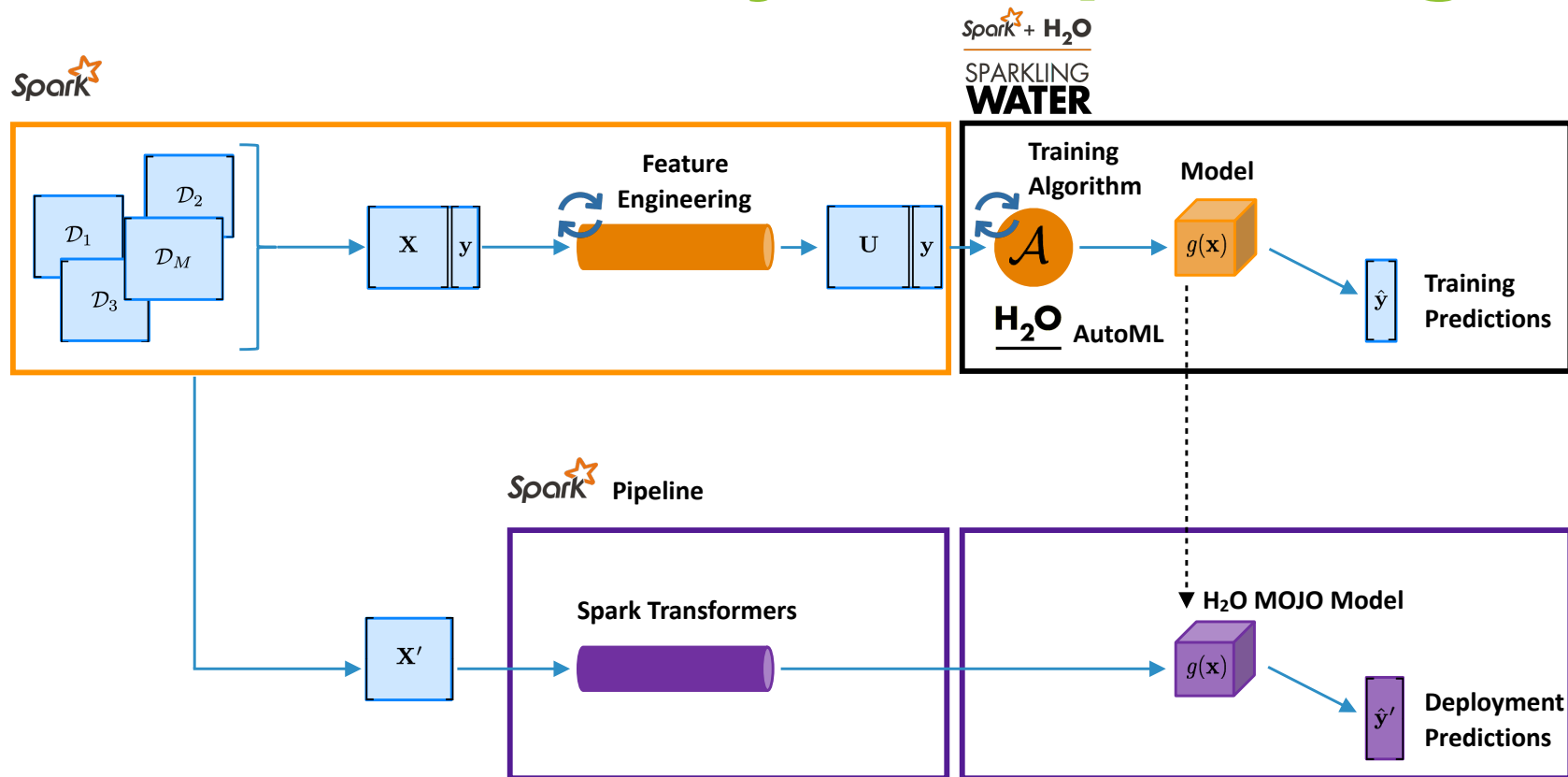| | Model Building | | | Model Deployment | |
|---|---|---|---|---|---|
| **Data Engineering** | **Feature Engineering** | **Training Algorithm** | **Deployment Pipeline** | **Model** |
| Spark | | H2O | Spark | H2O MOJO |
| Spark | H2O Driverless AI | | Spark | H2O Driverless AI MOJO |

# H2O + Spark = Sparkling Water

#ML4SAIS

# H2O + Spark

- H2O
  - Machine Learning Library
  - Distributed Algorithms
  - For ML experts

- Sparkling Water
  - Integrates H2O & Spark Ecosystems
  - Transparent for Spark users
  - Based on Spark pipelines & H2O

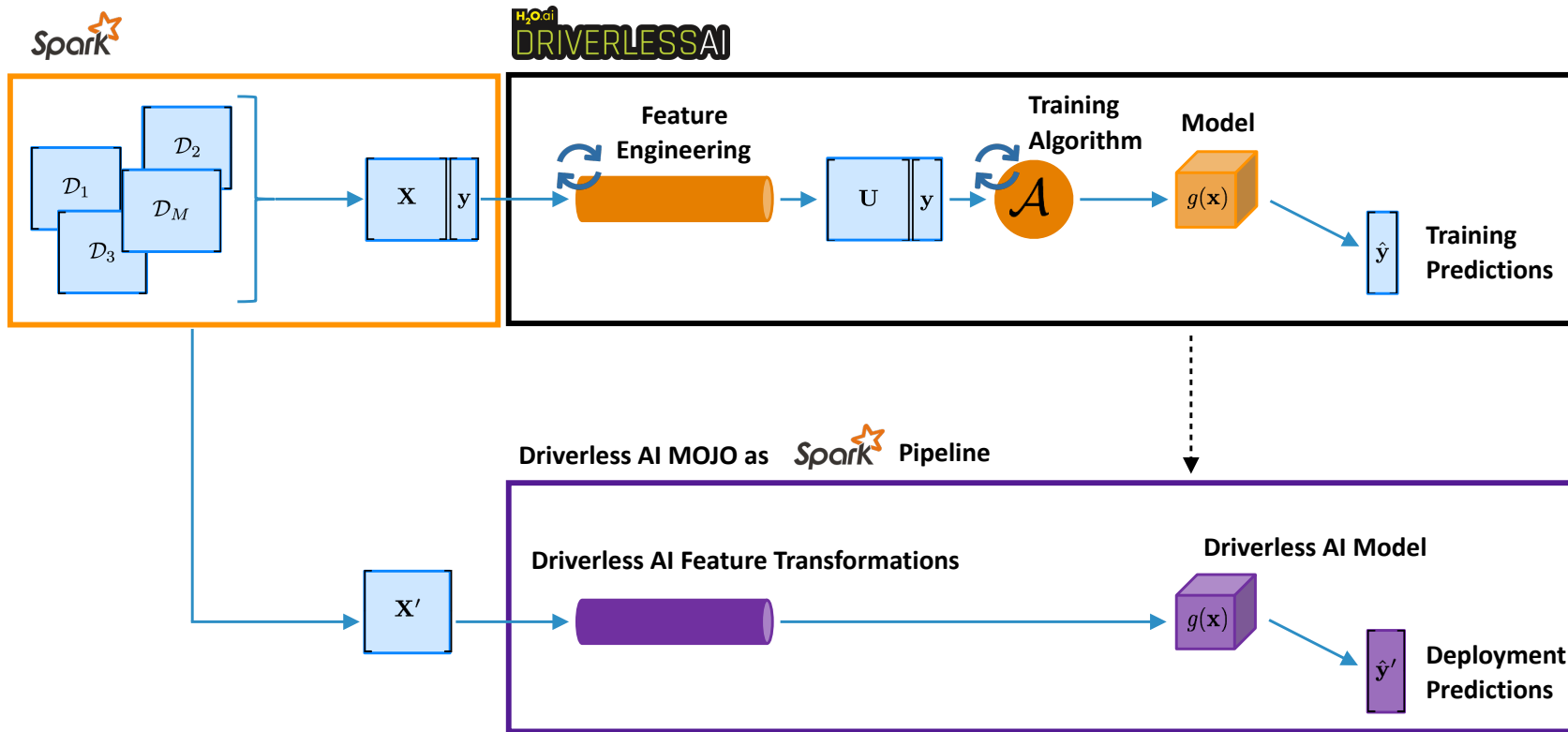# Basic ML Lifecycle: Sparkling Water

# Demo:
# Spark Pipeline

# H2O Driverless AI

# H2O Driverless AI

- What if I'm not expert ?
  - H2O Driverless AI
- H2O Driverless AI
  - No expert knowledge required
  - Automatic **Feature Engineering** & ML

# Basic ML Lifecycle: Driverless AI

# Demo: Driverless AI as Spark Pipeline

# H2O.ai    Experiment

## TRAINING DATA

DATASET
train.csv

| ROWS | COLUMNS | DROPPED COLS | VALIDATION DATASET | TEST DATASET |
|------|---------|--------------|--------------------|--------------|
| 24K | 25 | -- | -- | -- |

TARGET COLUMN
default payment next

FOLD COLUMN
--

WEIGHT COLUMN
--

TIME COLUMN
[OFF]

| TYPE | COUNT | UNIQUE | TARGET FREQ |
|------|-------|--------|-------------|
| int | 23999 | 2 | 18630 |

---

### What do these settings mean?

ACCURACY ■▬▬▬▬▬
- Training data size: **4,000 rows, 25 cols** (sampled)
- Feature evolution: **XGBoost, 1/3 validation split, 2 reps**
- Final pipeline: **XGBoost, 4-fold CV**

TIME ■▬▬▬▬▬▬▬
- Feature evolution: **8 individuals**, up to **500 iterations**
- Early stopping: After **50** iterations of no improvement

INTERPRETABILITY ■■▬▬▬
- Feature pre-pruning strategy: None
- Monotonicity constraints: disabled
- Feature engineering search space (where applicable):
['Clustering', 'Date', 'FrequencyEncoding', 'Identity',
'Interactions', 'TargetEncoding', 'Text', 'TruncatedSVD',
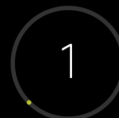'WeightOfEvidence']

XGBoost models to train:
- Feature evolution: **4024**
- Final pipeline: **1**

Estimated max. total memory usage:
- Feature engineering: **8.0MB**
- GPU XGBoost: **1.2GB**

Estimated runtime: **20 minutes**

## EXPERIMENT SETTINGS  HELP

| ( 1 ) | ( 10 ) | ( 3 ) |
|-------|--------|-------|
| ACCURACY | TIME | INTERPRETABILITY |

| CLASSIFICATION | REPRODUCIBLE | ENABLE GPUS |

SCORER

| GINI |
| MCC |
| F05 |
| F1 |
| F2 |
| ACCURACY |
| LOGLOSS |
| AUC |
| AUCPR |

LAUNCH EXPERIMENT

# Driverless AI Pipeline

# Governed ML Lifecycle

#ML4SAIS

# Governed ML Lifecycle

# Materials



https://bit.ly/2sxowxD

# Thank you!

Sparkling Water enables deployment of H2O ML models with Spark Pipelines