



VANDERBILT
UNIVERSITY

STRATUM: A Serverless Framework for the Lifecycle Management of Machine Learning-based Data Analytics Tasks

- 2019 USENIX Conference on Operational Machine Learning
- MAY 20, 2019

❖ *Anirban Bhattacharjee*

❖ Yogesh Barve

❖ Shweta Khare

❖ Shunxing Bao

❖ Aniruddha Gokhale

❖ Thomas Damiano

Acknowledgments

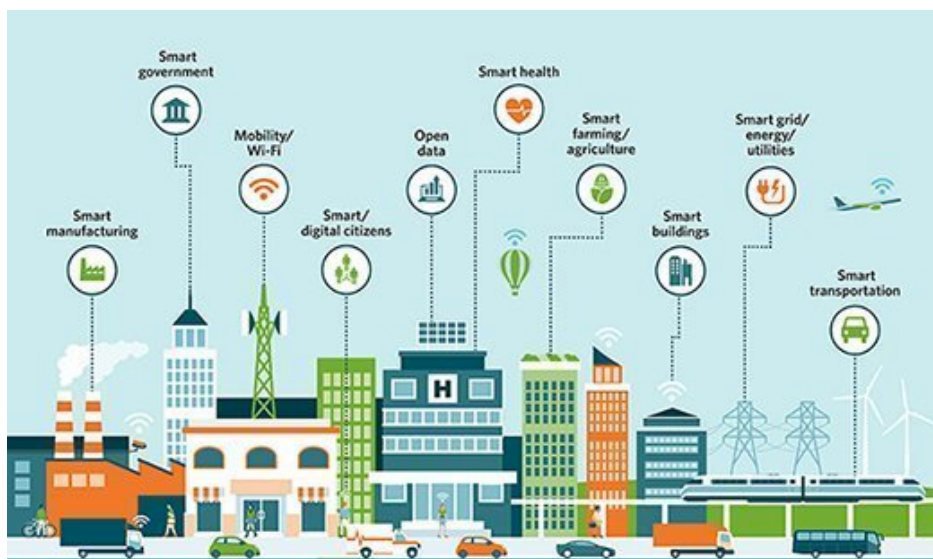
This work was supported in part by NSF US Ignite CNS 1531079, AFOSR DDDAS FA9550-18-1-0126 and AFRL/Lockheed Martin StreamlinedML program. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of NSF, AFOSR or AFRL.



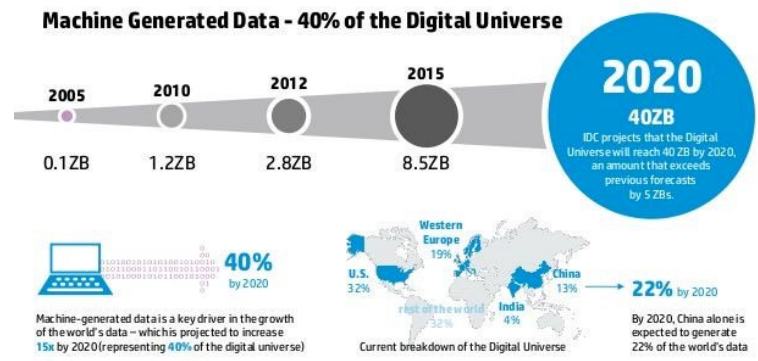
Data Analytics Trends

The world is changing and accelerating

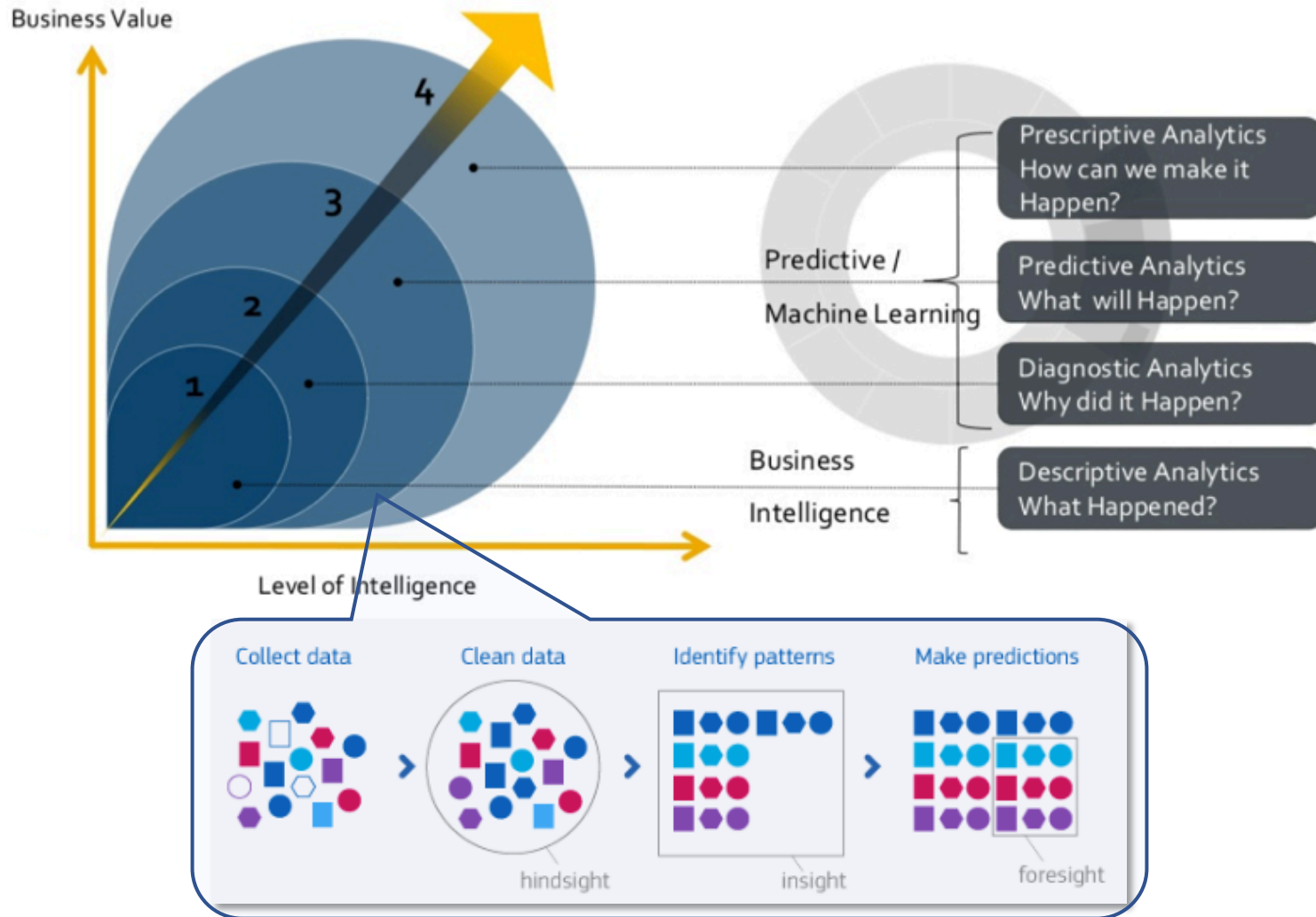
Internet of Things (IoT) applications, such as cognitive assistance, voice assistance, patient health monitoring, and connected vehicles are increasingly using the cloud and edge analytics.



Smart IoT devices generates data in volume and velocity, which needs to be analyzed to get valuable insight.



Big Data Value Model

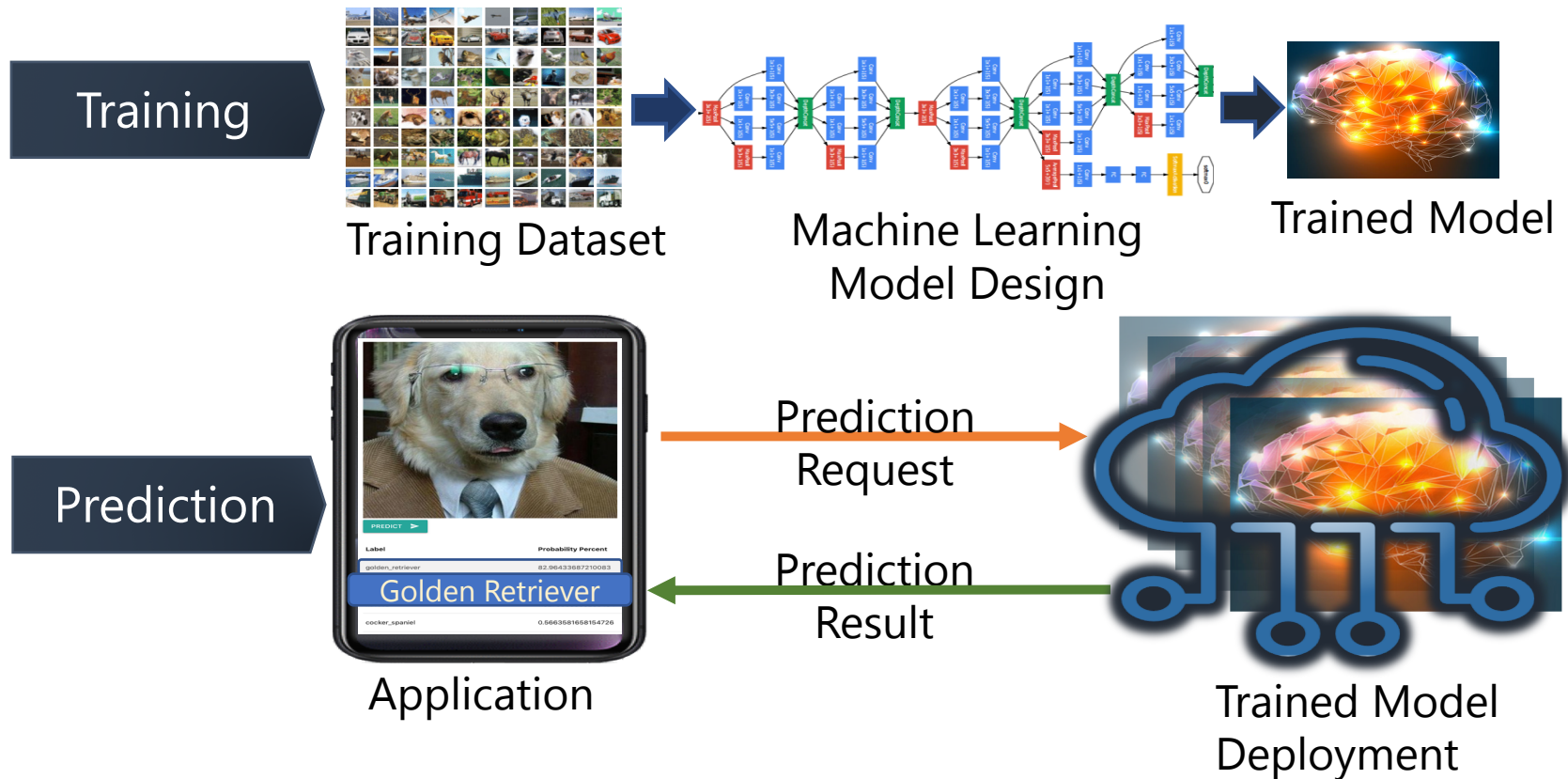




ML-based Predictive Analytics

- **Predictive analytics**

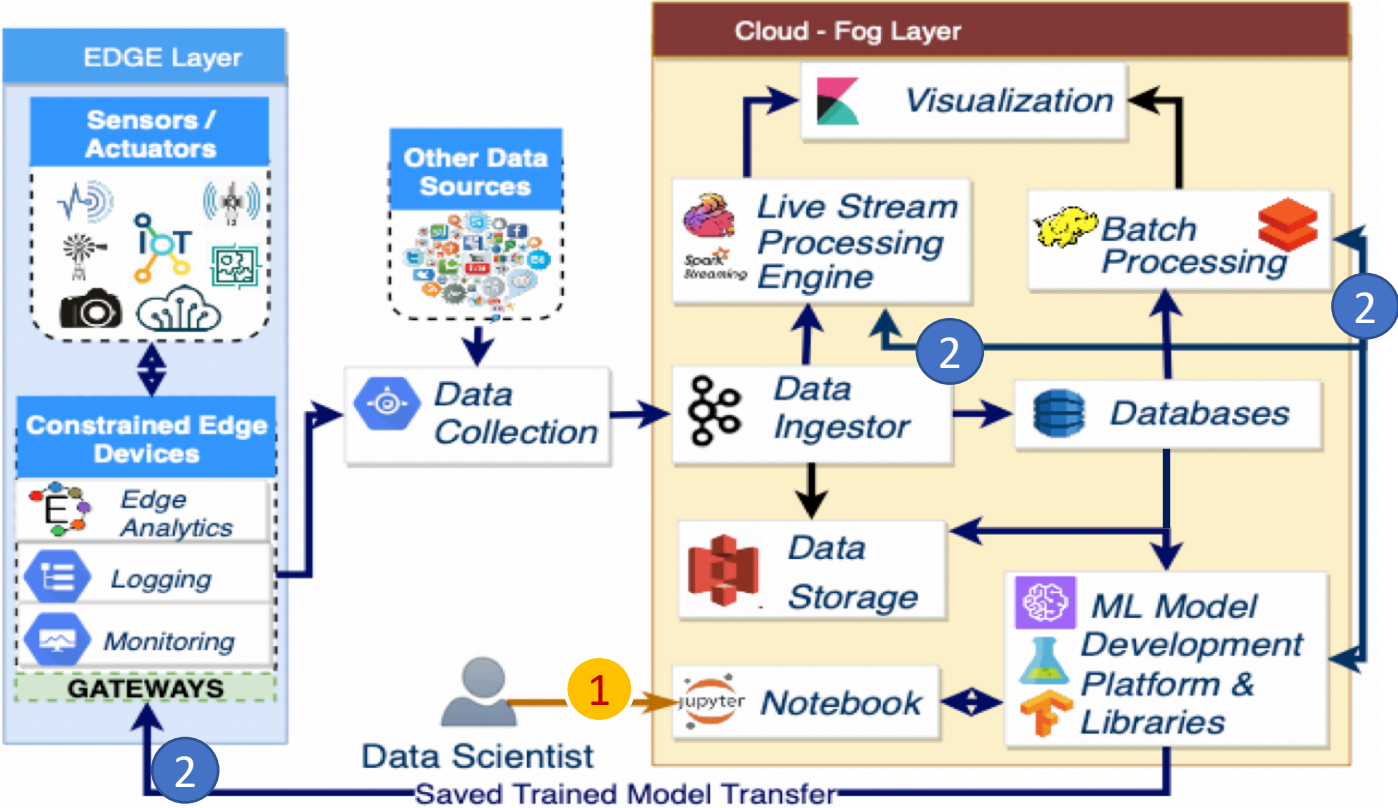
- It uses various statistical modelling and machine learning techniques to analyze past data and predict the future outcomes.



ML-based Predictive Analytics

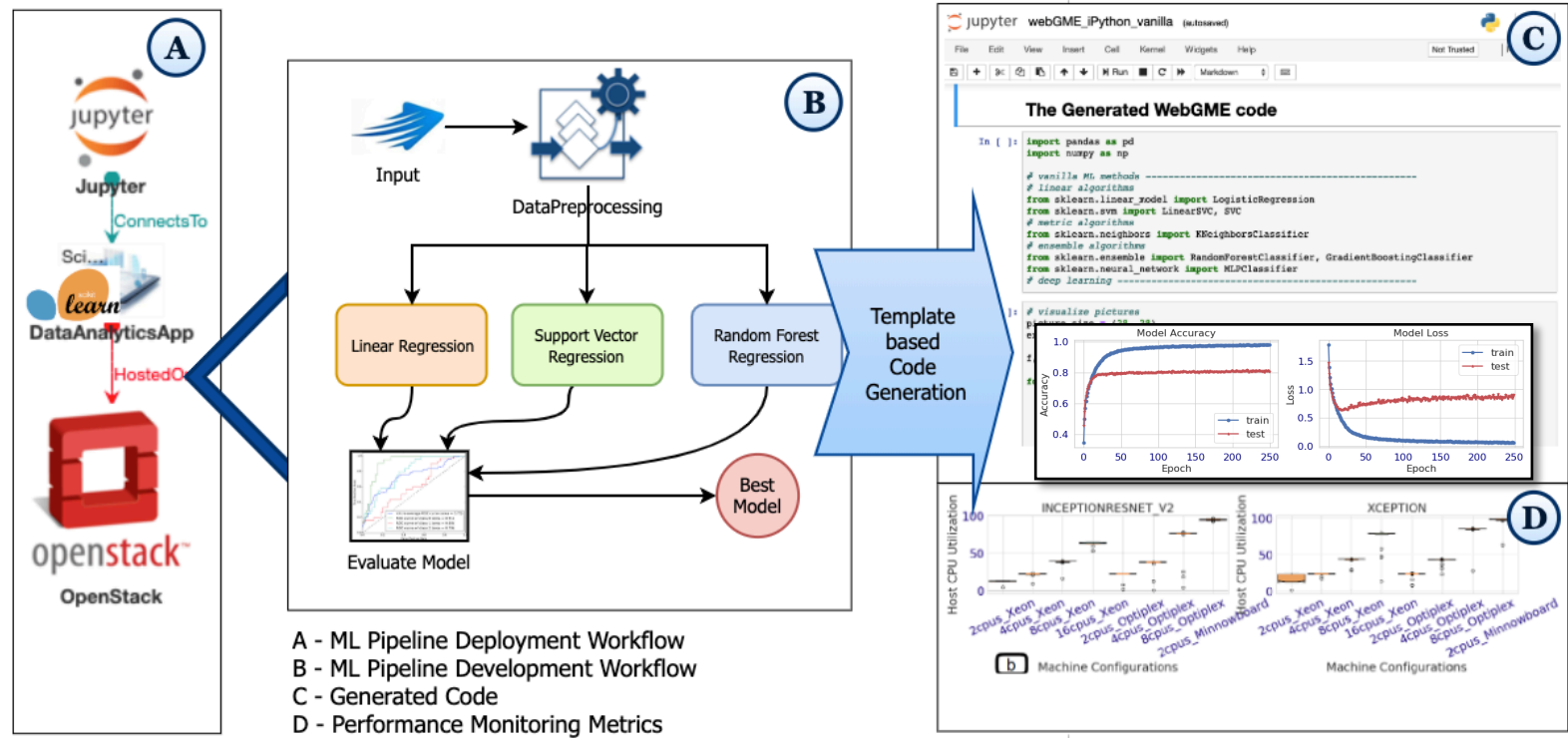
- **Predictive analytics**

- It uses various statistical modelling and machine learning techniques to analyze past data and predict the future outcomes.



Requirement

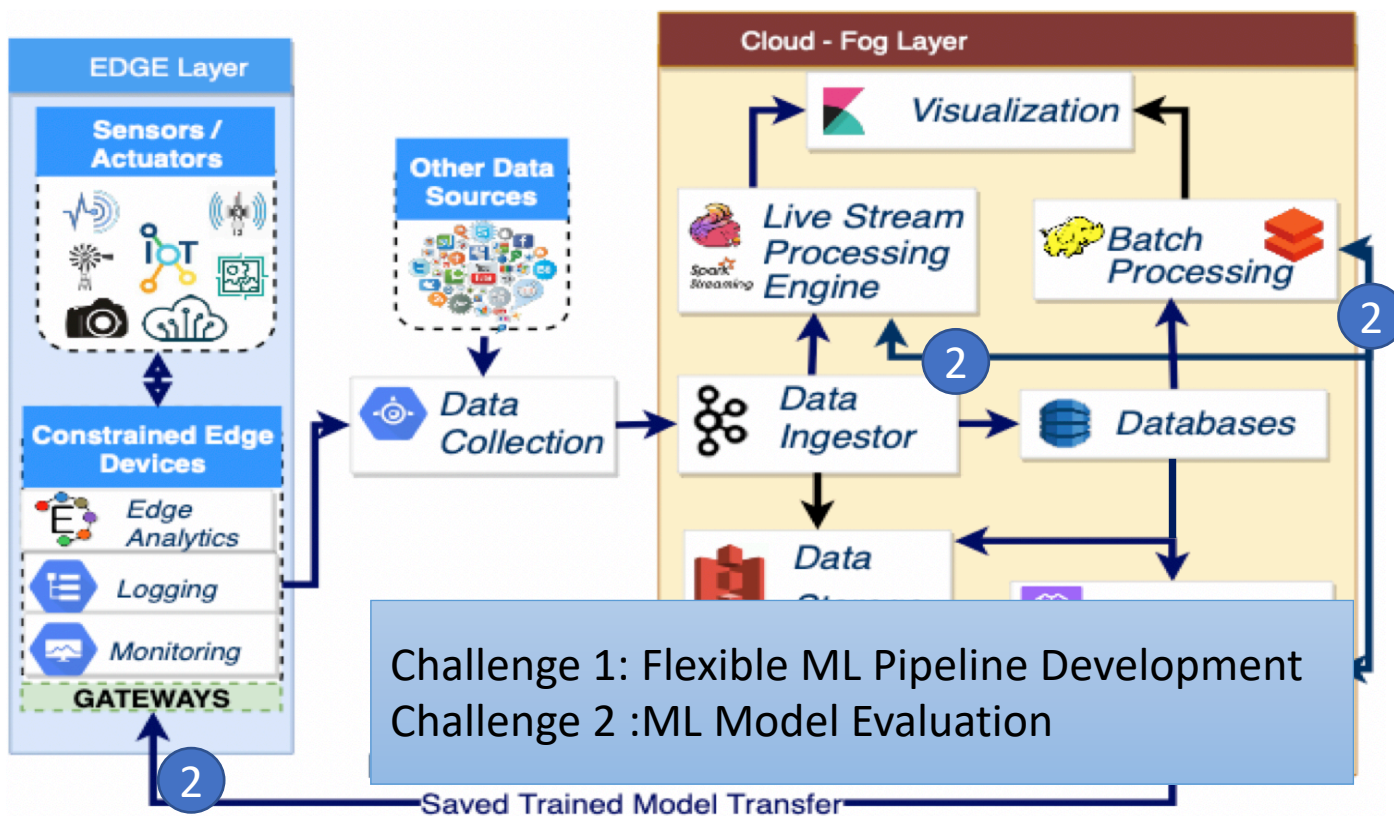
- Automation of Machine Learning (ML) Model Development and Deployment
 - Alleviate ML developers from writing the code from scratch
 - ML Library and Framework Agnostic



ML-based Predictive Analytics

- **Predictive analytics**

- It uses various statistical modelling and machine learning techniques to analyze past data and predict the future outcomes.



Challenges[1/3]

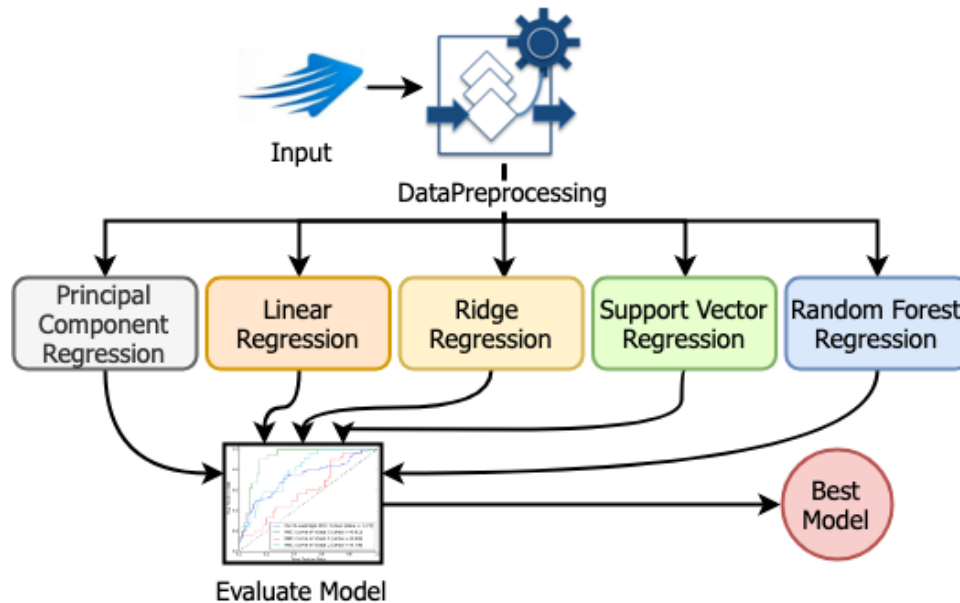
Flexible ML Pipeline Development



- A diverse set of ML algorithms -
 - classification (logistic regression, naive Bayes), regression, decision trees, random forests, and gradient-boosted trees, recommendation (ALS), clustering (K-means, GMMs), and many others.
- A diverse set of different ML libraries and frameworks –
 - Scikit-learn, Spark MLlib, TensorFlow etc.
- ML pipeline capabilities needs to be captured, and abstracted in the metamodel.
 - Attributes of ML algorithms, data preprocessing strategies, evaluation methods etc.

Challenges[2/3]

ML Model Evaluation



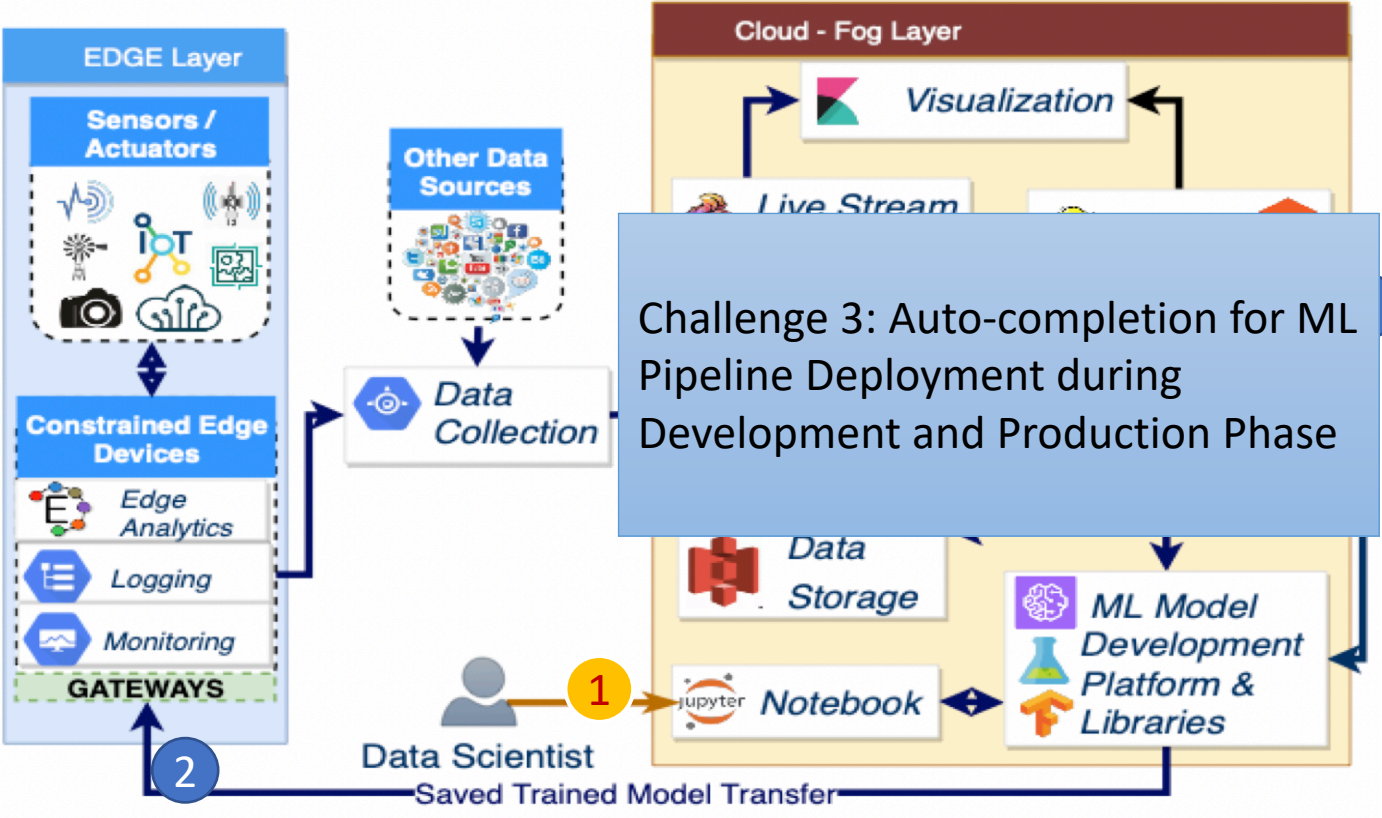
- After training the ML models with the diverse set of ML algorithms, the **best model** for the dataset needs to be selected.
 - save it for prediction jobs.
- The model can be evaluated based on **different scoring methods** such as accuracy, f1 score, precision, r2 score, mean square error which is captured in a metamodel.
- To speed up the training process, the ML models with different algorithms can be distributed.

ML Trained Model with all the software dependencies needs to be encapsulated in a container on the specific hardware.

ML-based Predictive Analytics

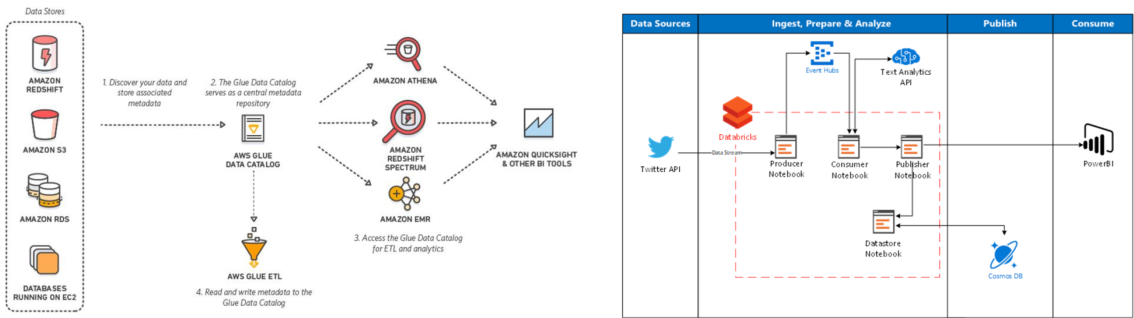
- **Predictive analytics**

- It uses various statistical modelling and machine learning techniques to analyze past data and predict the future outcomes.



Challenges[3/3]

Auto-completion for ML Pipeline Deployment



```

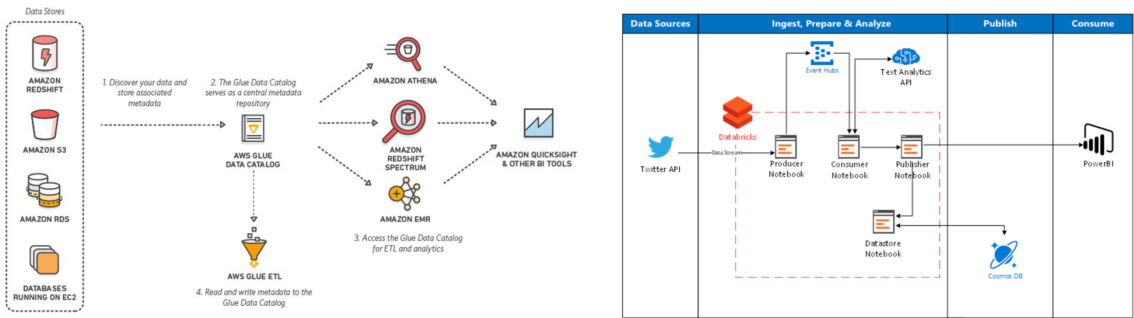
- name: Install PHP-FPM
  apt:
    name: "{{ item }}"
    state: latest
    update_cache: yes
    with_items: "{{ ubuntu_php_pkgs }}"

- name: Copy the templates to their respective destination
  template:
    src: "{{ item.src }}"
    dest: "{{ item.dest }}"
    owner: root
    group: root
  with_items:
    - { src: 'www.conf.j2', dest: '/etc/php/7.0/fpm/pool.d/www.conf' }
    - { src: 'php.ini.j2', dest: '/etc/php/7.0/fpm/php.ini' }
  notify:
    - restart apache2
  
```



Challenges[3/3]

Auto-completion for ML Pipeline Deployment



```

- name: Install PHP-FPM
  apt:
    name: "{{ item }}"
    state: latest
    update_cache: yes
    with_items: "{{ ubuntu_php_pkgs }}"
- name: Copy the templates to their respective destination
  templates
  src:
  dest:
  owner:
  group:
  with_items:
  - {}
  notify:
  - {}

```

Service Providers Concerns

How to deploy and maintain the application components with ease to increase productivity and usability while reducing the time-to-market?



STRATUM: ML Pipeline Automation



FRAMEWORK DESIGN



EVALUATION RESULTS

STRATUM: ML Pipeline Automation

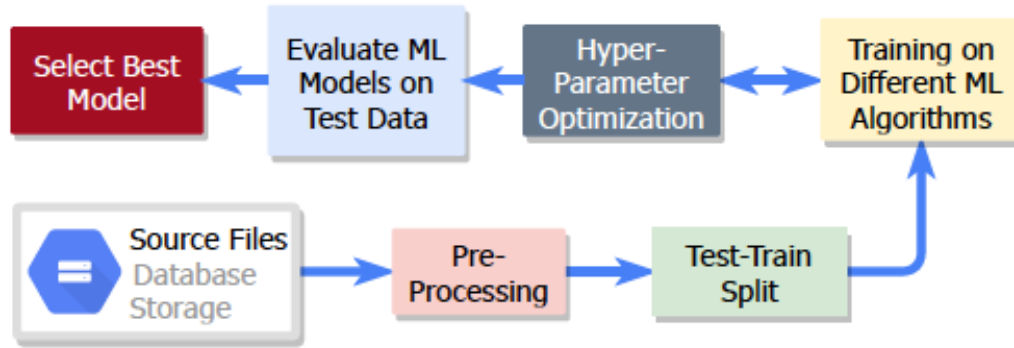


FRAMEWORK DESIGN



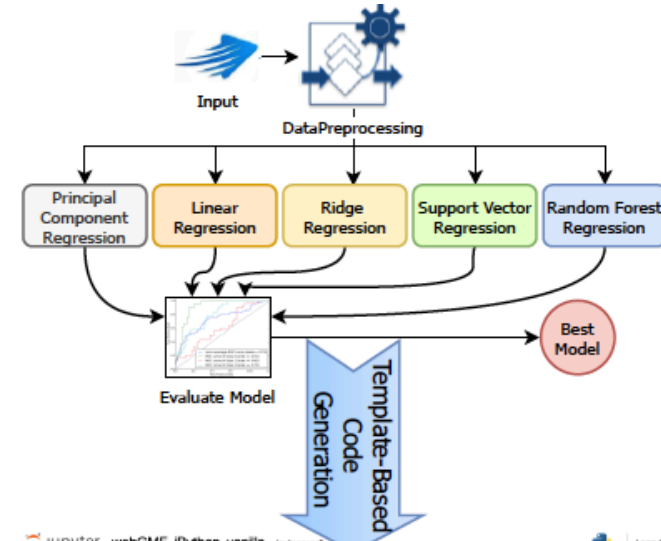
EVALUATION RESULTS

ML Model Evaluation



Sample Machine Learning Pipeline

- The framework provide the GUI for ML pipeline construction, model evaluation, and hyper-parameter tuning capabilities, which forms the basis for continuous evaluation.
- We also integrated Jupyter Notebook (notebook-based environment) to provide data-scientists the ability to train their models interactively.
- The ML execution pipeline needs to be bind with a specific library or framework such as Scikit-learn or TensorFlow.



```

jupyter webGME_iPython_vanilla auto saved
File Edit View Insert Cell Kernel Widgets Help Not Trained Python 3.0
The Generated WebGME code
In [ ]: import pandas as pd
import numpy as np

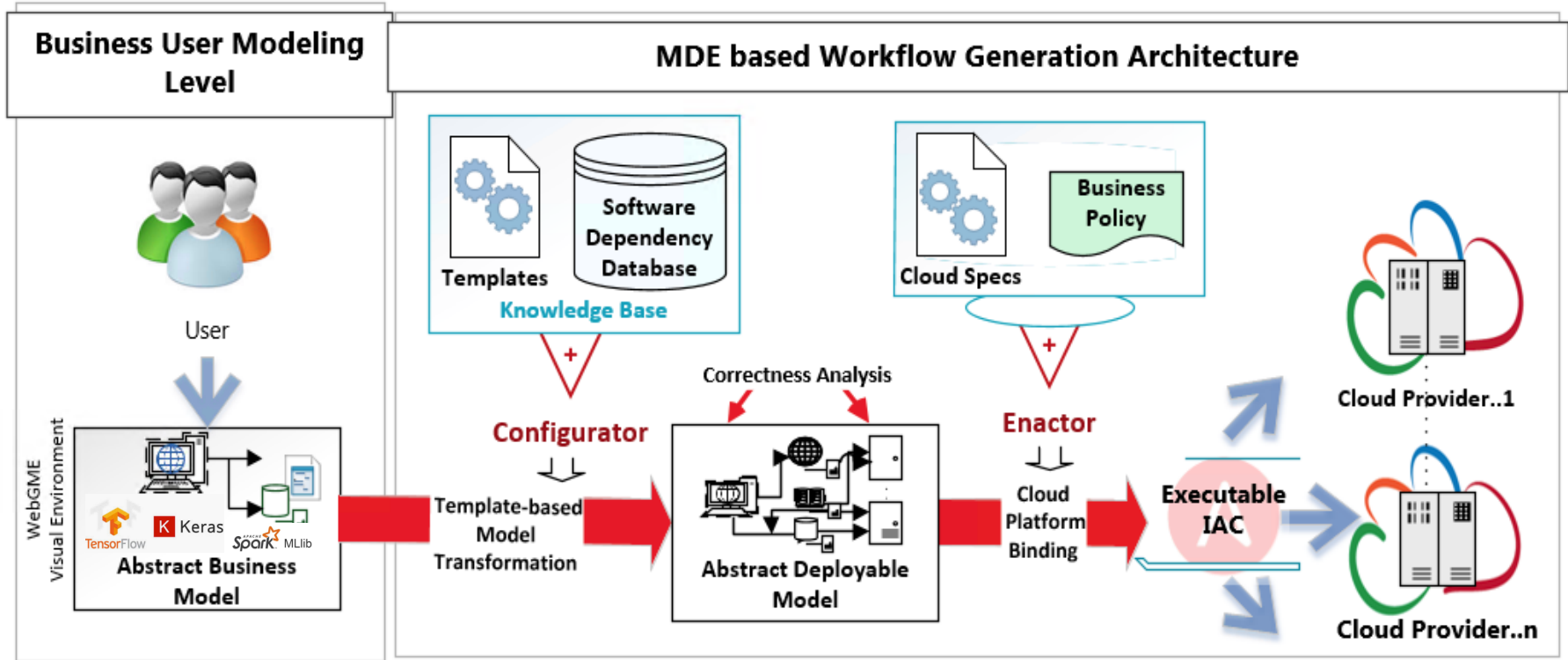
# vanilla ML methods -----
# linear algorithms
from sklearn.linear_model import LogisticRegression
from sklearn.svm import LinearSVC, SVC
# metric algorithms
from sklearn.neighbors import KNeighborsClassifier
# ensemble algorithms
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
from sklearn.neural_network import MLPClassifier
# deep learning -----

In [ ]: # visualize pictures
picture_size = (28, 28)
examples = [train[train.label == k].sample(1, random_state=42).values for k in range(10)]
f, axarr = plt.subplots(1, 10, squeeze=False, figsize=(12, 1.2))

for i, e in enumerate(examples):
    # reshape 1D to 2D array - a picture
    img = e[:, 1:].reshape(picture_size).astype(float)
    # then draw it
    axarr[0, i].set_title(i)
    axarr[0, i].imshow(img, cmap='gray', interpolation='bicubic')
  
```

Sample generative capabilities of Stratum

DSML for STRATUM Deployment Framework



- **Abstraction of Model:** Deployers provision by selecting only business-relevant components.
- **Configurator:** Transforms abstract service components to Ansible-specific automation tasks using DSML.
- **Enactor:** Generates IAC by integrating automation code, cloud-specs & inter-component connection types.
- **Knowledge Base:** Software dependencies for service component types are stored in RDBMS table.

STRATUM: ML Pipeline Automation



FRAMEWORK DESIGN



EVALUATION RESULTS

Usability of STRATUM Framework

The screenshot displays the STRATUM framework interface. At the top, the breadcrumb path is 'GME > erudite > master > TestMLModel'. The interface is divided into several sections:

- Visualizer Selector:** A sidebar on the left with a 'Composition' tab. It contains a 'Meta' section with options like 'Set membership', 'Crosscut', and 'Graph view'. Below this are icons for 'DataPreprocessing', 'Documentation', 'EvaluateModel', 'Input', 'MLAlgorithms', and 'Visualize'.
- Workflow Diagram:** A central diagram showing an 'Input' node leading to a 'DataPreprocessing' node. From 'DataPreprocessing', arrows point to five model nodes: 'Principal Component Regression', 'Linear Regression', 'Ridge Regression', 'Support Vector Regression', and 'Random Forest Regression'. These nodes all feed into an 'Evaluate Model' node, which then points to a 'Best Model' node.
- Property Editor:** A panel on the right titled 'PROPERTY EDITOR' for a 'RandomForestRegressor'. It shows various attributes and their values:

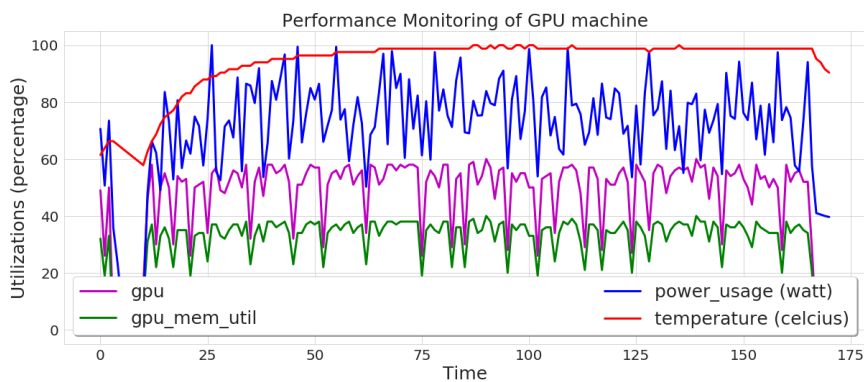
Attribute	Value
GUID	5dcccfb8-eaa4-fc91-3514-2ae
ID	/x
Meta type	RandomForestRegressor
Platform	scikitLearn
bootstrap	TRUE
criterion	mse
max_depth	0
max_features	auto
min_samples_split	2
n_estimators	10
n_jobs	1
name	RandomForestRegressor
- ROC Curve:** A plot at the bottom showing the True Positive Rate (Y-axis) versus the False Positive Rate (X-axis). The plot is titled 'ROC curve' and contains five curves with their respective Area Under the Curve (AUC) values:

Model	AUC
RT + LR	0.672
RF	0.946
RF + LR	0.960
GBT	0.969
GBT + LR	0.971

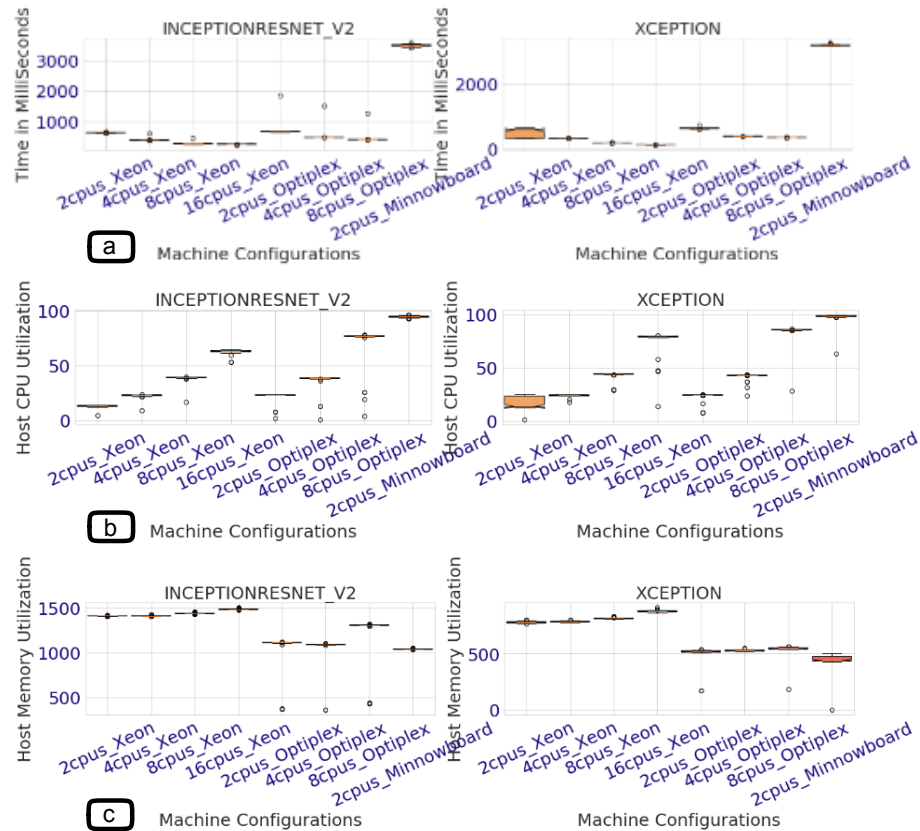
A yellow banner at the bottom of the screenshot contains the text: **Flexible ML Pipeline Development and Model Evaluation**.

- The ML model can be developed by dragging and dropping the build blocks (from box 1→2).
- All the attributes of the selected ML algorithms such as max_depth, criteria need to be specified by the user (box 3).
- The Erudite model transformer can distribute different jobs with different ML techniques over a cluster of connected machines.
- It aids the developer to select the best model or ensemble of models based on the user's choice of evaluation methods.

Performance Monitoring



Sample Deep Learning Training on GPU Machine



Performance Monitoring of the prediction services (a)The execution latency of InceptionResnetV2 and Xception model on different ML containers with variable configurations, (b) Host CPU utilization of the ML containers (c) Host Memory utilization of ML containers (in MB).

Summary



- We presented a model-driven engineering and generative programming approach for automated development of ML pipeline.
- We integrated a monitoring framework to analyze the performance of ML pipeline during training and prediction phase.
- We proposed a ML pipeline deployment methodology across cloud-fog-edge spectrum.



<https://doc-vu.github.io/Stratum/>