

ModelOps

**A programming model for reusable,
platform-independent, and composable AI workflows**

Waldemar Hummer, Vinod Muthusamy

IBM Research AI

IBM Programming Languages Day

T.J. Watson Research Center

December 10, 2018

<https://ibm.biz/plday2018>

Contact: whummer@ibm.com

Team:

Yorktown: Waldemar Hummer, Anupama Murthi,
Kaoutar El Maghraoui, Benjamin Herta, Darrell Reimer,
Punleuk Oum, Gaodan Fang

Austin: Vinod Muthusamy

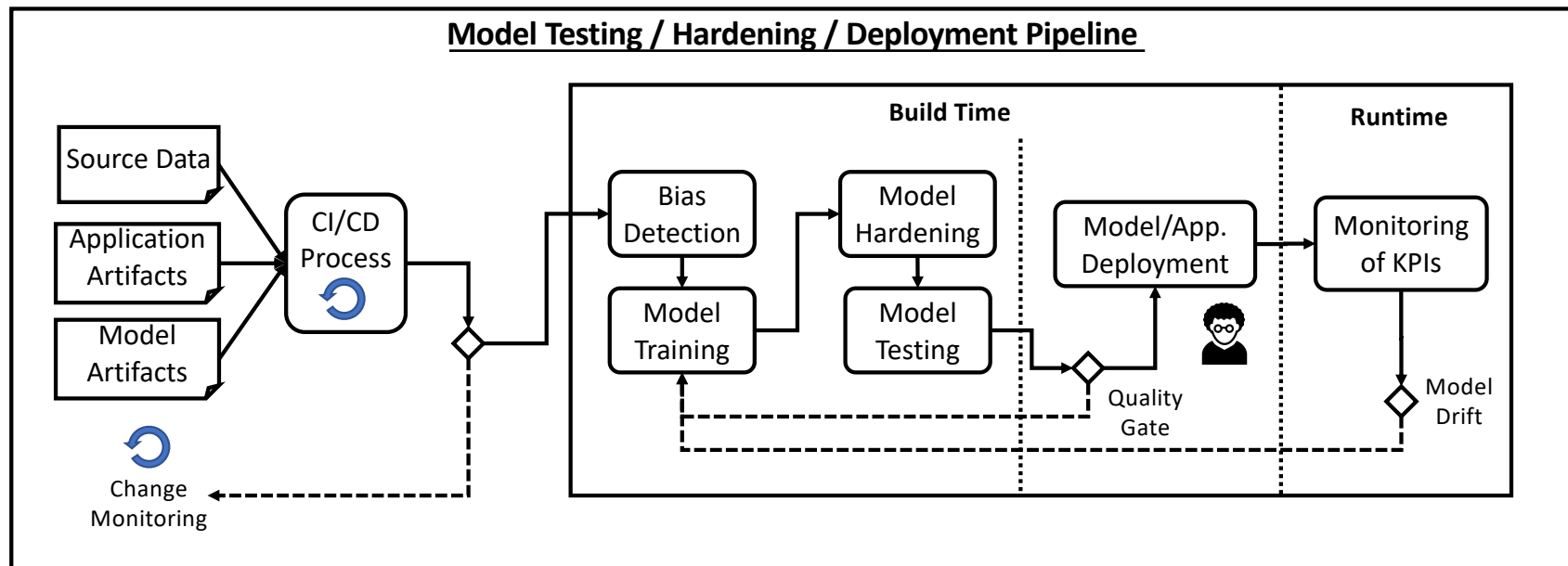
Cambridge: Scott Boag

AI Models become engrained in the Software Lifecycle

- What's the difference to traditional DevOps?

Classical Application Lifecycle	AI Application Lifecycle
Requires dev / ops skills	Involves more diverse skill sets
Relatively short running	Long-running, resource intensive
Human speed (low change frequency)	Continuous (re-)training
Few versions of software artifacts	Huge number of models
Linear evolution of artifacts	Specialized models coexist
Configurations applied at runtime	Parameters tuned at training time
Codebase changes trigger new builds	Data/code changes trigger model re-training
Deterministic testing	Statistical/probabilistic testing
Monitoring of application performance & KPIs	Monitoring of model accuracy, drift, and KPIs

AI Operations – Pipeline Scenario



- **AI models** introduce **risk** to business applications
- Training and deployment pipelines can become quite **complex**
- Pipelines often **hand-crafted**
 - Hard to maintain
 - Hard to optimize or reason over

Core Problems Addressed

- **Reusability**

- Pipelines are often hand-crafted (shell scripts, Makefiles, Python scripts, ...)
- Yet, there are common patterns: artifacts could be shared across users and teams

- **Composability**

- Despite the common patterns: hardly a one-size-fits-all solution
- Ability to compose pipelines and other features from building blocks is critical

- **Platform Independence**

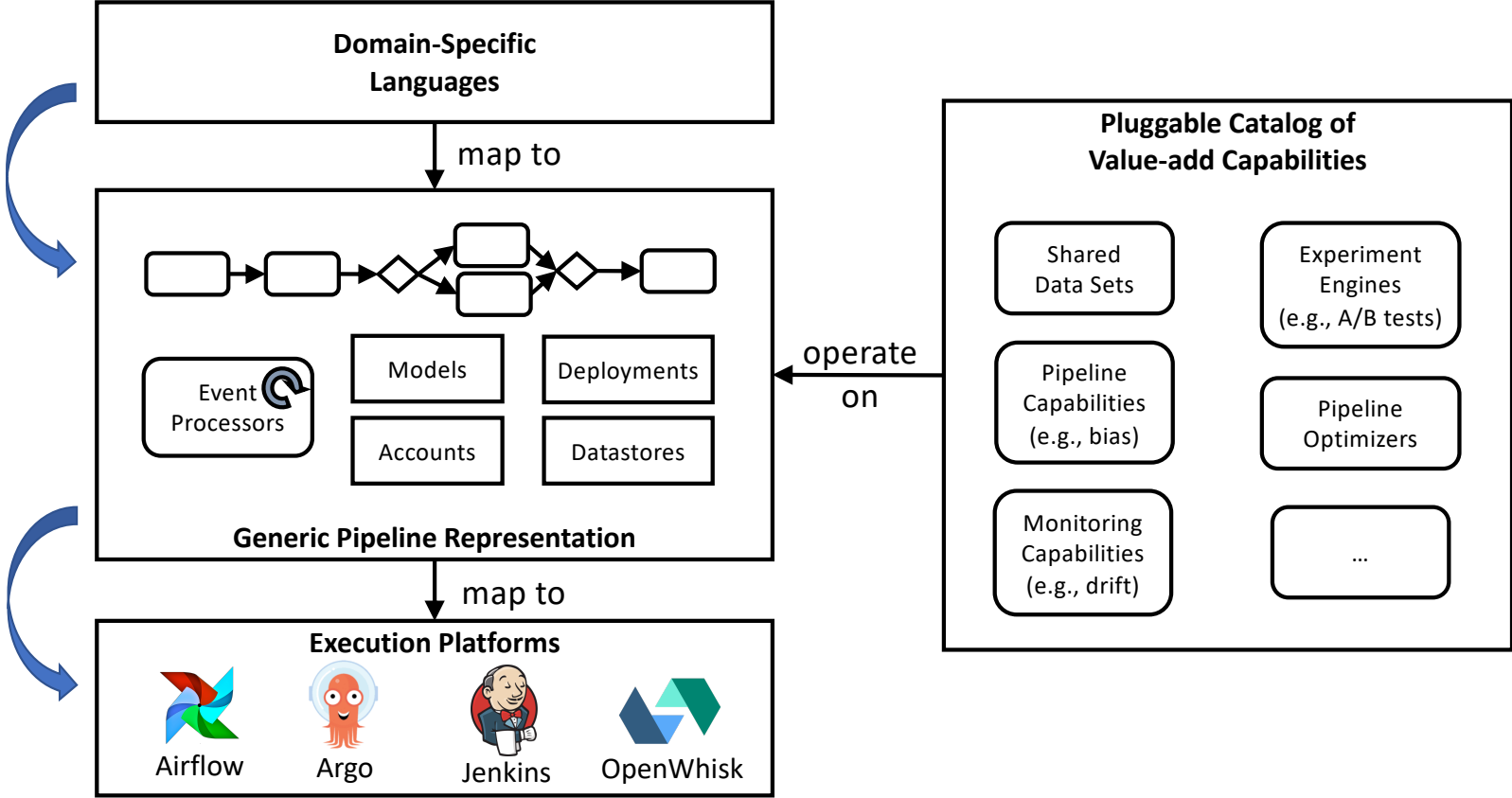
- Different Cloud vendors use different APIs – yet very similar concepts
 - E.g., “train model”, “deploy model” → map to different API calls in WML / AWS / GCP
- Parts of pipelines can run on various different environments
 - On local machines, on smart phones, in the cloud, on edge devices, HPC clusters, ...

- **→ Ability to Optimize and Reason over AI Pipelines**

- E.g., Task pruning, task co-location, task scheduling

Approach

ModelOps: Towards a programming model for Operationalizing AI



Generic AI Pipeline Representation

- Metamodel of approx. 25 entities
 - Pipelines (DAG of tasks with dependencies)
 - Models (AI Classifier)
 - Event Triggers / Subscribers
 - Datastores
 - Accounts
 - Environments
 - Code Plugins
 - ...
- Entities implemented as Python classes
 - Attributes defined in JSON Schema
 - Stored in markup files (e.g., YAML, JSON)
- Language features
 - Plugin Imports
 - Variable placeholders (lazy evaluation)
 - Control flow (pipelines DAG)

```
1  accounts:
2    - name: account_wml
3      type: ibm/wml
4      password: '{{credentials.wml.password}}'
5      username: '{{credentials.wml.username}}'
6      instance_id: '{{credentials.wml.instance_id}}'
7  datastores:
8    - name: training_data
9      type: cos
10     bucket: trainint-data
11     account: s3_account
12  models:
13    - name: model1
14      type: wml
15      command: 'python train.py'
16      framework_name: tensorflow
17      source: training_data
18      target: model_results
19  subscribers:
20    - _type: plugins.templates.DriftDetector
21      name: drift_detector_model1
22  pipelines:
23    - name: train-deploy-model
24      tasks:
25        - _type: modelops.tasks.TrainModel
26          name: TrainModel
27          model: model1
28        - _type: modelops.tasks.StoreModel
29          name: StoreModel
30          model: model1
31        - _type: modelops.tasks.DeployModel
32          name: DeployModel
33          model: model1
```


Extensibility via Plugins

- Extensions as Plugins
 - Pipeline tasks
 - E.g. “Train Model”
 - Subscribers
 - E.g., “Detect Drift”
 - Transformers
 - (See following slides)
- Plugins can be referenced in the config
 - Automatically loaded into Python path at runtime

```
class PipelineTask(EmbeddedEntity):
    """ Represents a task in a pipeline. """

    def execute(self, **params):
        """ Main method that starts the execution of this task. Implemented by subclasses. """
        raise NotImplementedError()

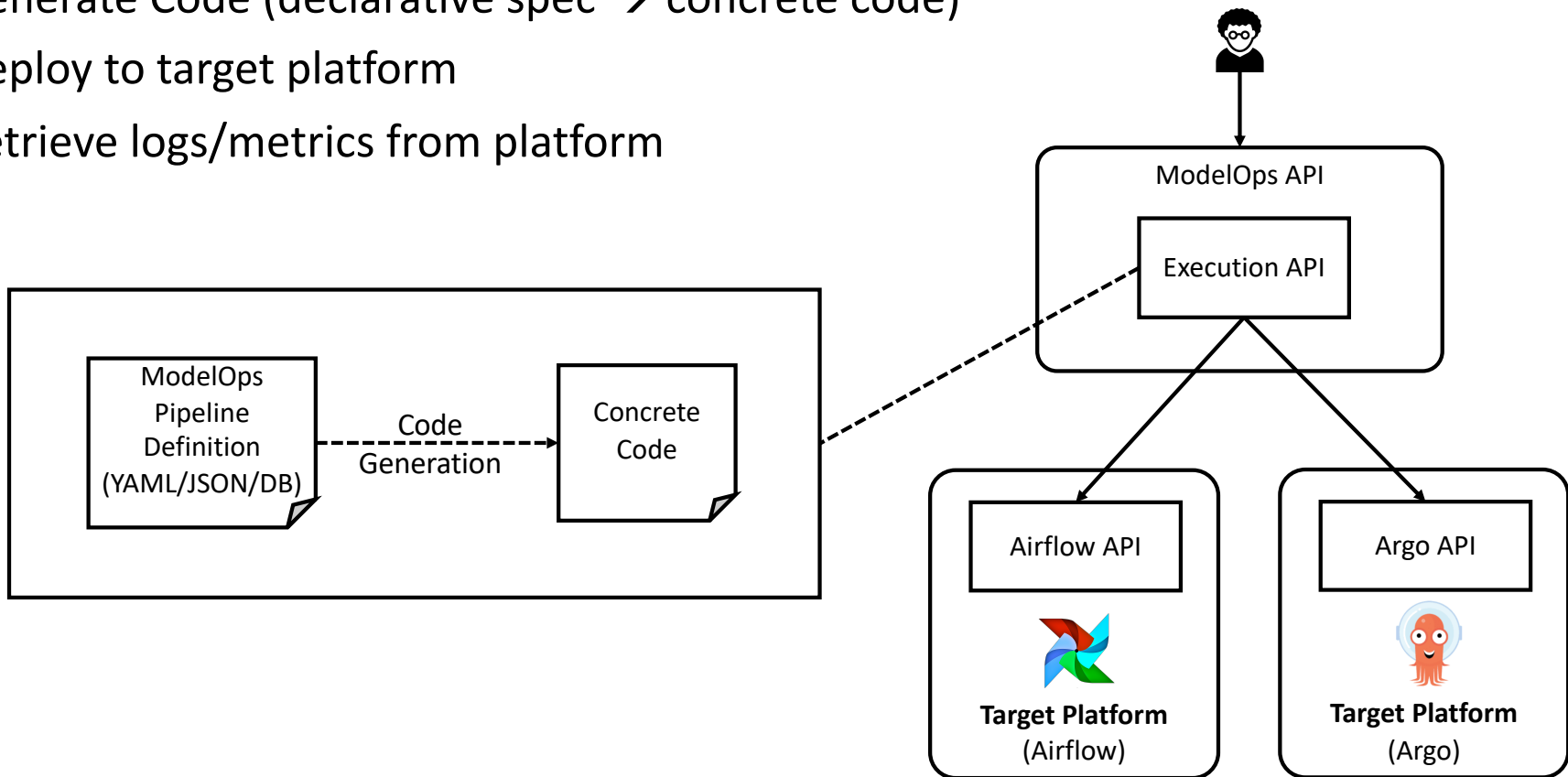
    def finished(self, execute_result, **params):
        """ Return the final result of this task, or False if this task has not finished yet.
        This method receives as parameter the return value of the "execute" method.
        By default, we assume that "execute" runs synchronously, and hence this method
        returns a truthy value by default. Subclasses need to overwrite this method with
        specific logic that polls for task completion. """
        return execute_result or True

    def poll_interval(self, **params):
        """ Interval in seconds to apply when polling for results. If this value is 0
        (default), then this is a synchronous task that is finished after "execute(..)"
        returns (i.e., no polling required). If this value is a negative number,
        then let the task executor pick a suitable polling interval. """
        return -1

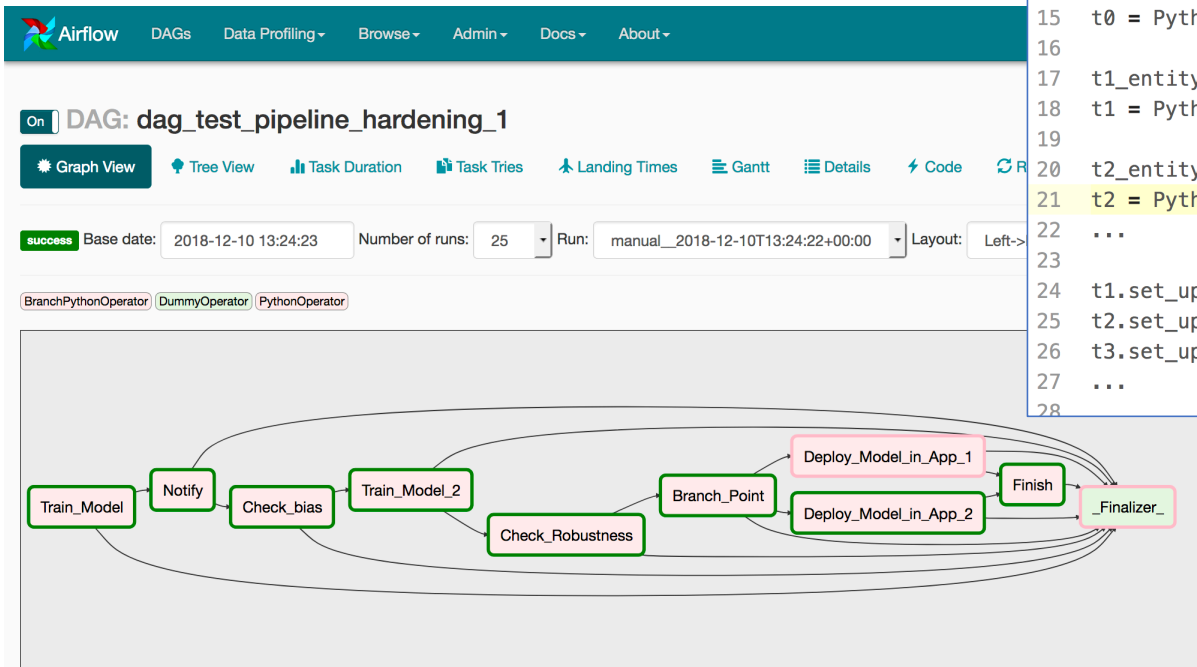
    def get_timeout(self, **params):
        """ Return the timeout interval (in seconds) after which this task execution should
        be considered failed. Returns the value of the "timeout" property of this task
        instance (if present), or 30 minutes by default. Subclasses can overwrite this method.
        If a negative or non-integer value is provided, let the task executor choose a timeout. """
        return self.timeout if is_positive_number(self.timeout) else self.DEFAULT_TIMEOUT
```

Pipeline Executor

1. Generate Code (declarative spec → concrete code)
2. Deploy to target platform
3. Retrieve logs/metrics from platform



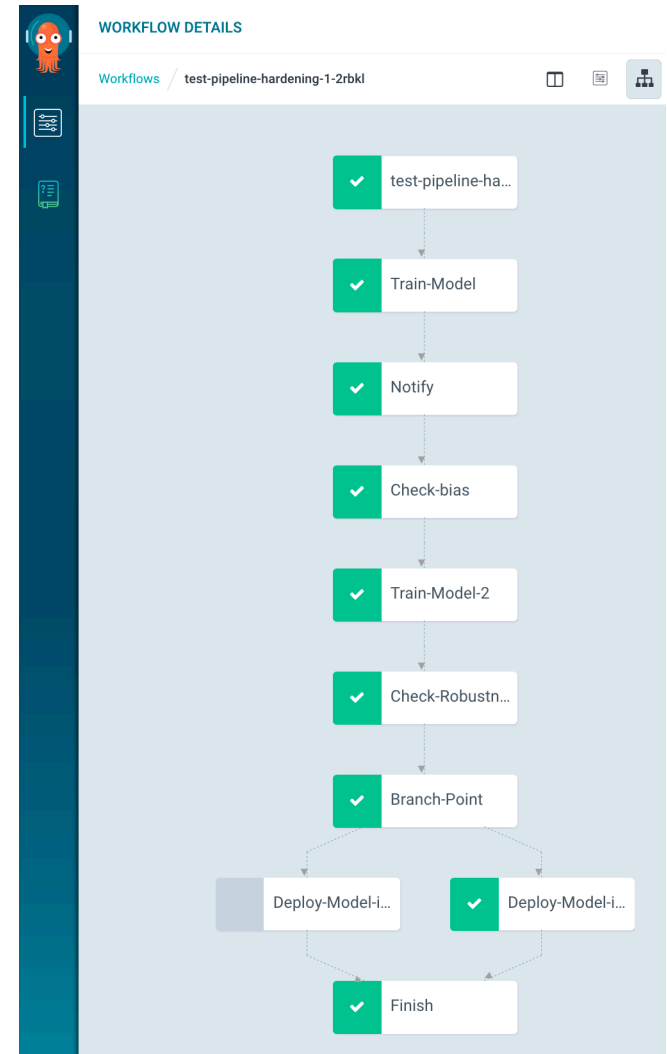
Generated Pipeline Code – Airflow



```
1 import json
2 from airflow import DAG
3 ...
4
5 def run_task(task_entity, **params):
6     from modelops import main
7     ...
8
9 args = {
10     'start_date': datetime.datetime.fromtimestamp(1544448249.35)
11 }
12 dag = DAG(dag_id='dag_test_pipeline_hardening_1', default_args=args)
13
14 t0_entity = json.loads('{"_type": "modelops.tasks.TrainModel", "e
15 t0 = PythonOperator(task_id='Train_Model', python_callable=run_task)
16
17 t1_entity = json.loads('{"_type": "modelops.tasks.PostMessage", "
18 t1 = PythonOperator(task_id='Notify', python_callable=run_task, op
19
20 t2_entity = json.loads('{"_type": "modelops.tasks.DataBiasCheck",
21 t2 = PythonOperator(task_id='Check_bias', python_callable=run_task,
22 ...
23
24 t1.set_upstream(t0)
25 t2.set_upstream(t1)
26 t3.set_upstream(t2)
27 ...
28
```

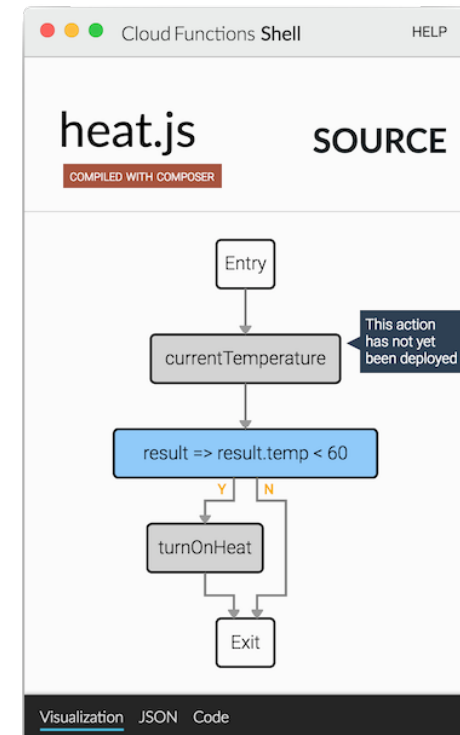
Generated Pipeline Code – Argo

```
1 name: test-pipeline-hardening-1
2 dag:
3   tasks:
4     - arguments: {parameters: [{name: task, value: eyJfdHwcy5wbHV
5       name: Train-Model
6       template: run-task
7     - arguments: {parameters: [{name: task, value: eyJfdHyNrLLBv
8       dependencies: [Train-Model]
9       name: Notify
10      template: run-task
11     - arguments: {parameters: [{name: task, value: eyJfcy5wJpYXM
12       dependencies: [Notify]
13       name: Check-bias
14       template: run-task
15     - arguments: {parameters: [{name: task, value: eyJfcy5wbLmFw
16       dependencies: [Check-bias]
17       name: Train-Model-2
18       template: run-task
19     - arguments: {parameters: [{name: task, value: eyJfdHlwZSI6I
20       dependencies: [Train-Model-2]
21       name: Check-Robustness
22       template: run-task
23     - arguments: {parameters: [{name: task, value: eyJfluZV9icmF
24       dependencies: [Check-Robustness]
25       name: Branch-Point
26       template: run-task
```



Generated Pipeline Code – OpenWhisk Composer

```
1 const composer = require('@ibm-functions/composer/composer')
2 const conductor = require('@ibm-functions/composer/conductor')
3 const pipeline_name = 'test_pipeline_hardening_1'
4 const wsk = conductor({ignore_certs: true})
5
6 const composition = composer.sequence(
7   composer.sequence("test_pipeline_hardening_1_Train_Model",
8     composer.dowhile_nosave((tmp) => (tmp),
9       composer.sequence("test_pipeline_hardening_1_Train_Model_poll", (r)
10     ),
11   composer.sequence("test_pipeline_hardening_1_Notify",
12     composer.dowhile_nosave((tmp) => (tmp),
13       composer.sequence("test_pipeline_hardening_1_Notify_poll", (r) => (
14     ),
15   composer.sequence("test_pipeline_hardening_1_Check_bias",
16     composer.dowhile_nosave((tmp) => (tmp),
17       composer.sequence("test_pipeline_hardening_1_Check_bias_poll", (r)
18     ),
19   composer.sequence("test_pipeline_hardening_1_Train_Model_2",
20     composer.dowhile_nosave((tmp) => (tmp),
21       composer.sequence("test_pipeline_hardening_1_Train_Model_2_poll", (
22     ),
23   composer.sequence("test_pipeline_hardening_1_Check_Robustness",
24     composer.dowhile_nosave((tmp) => (tmp),
25       composer.sequence("test_pipeline_hardening_1_Check_Robustness_poll"
26     ),
27   composer.if("test_pipeline_hardening_1_Branch_Point",
```



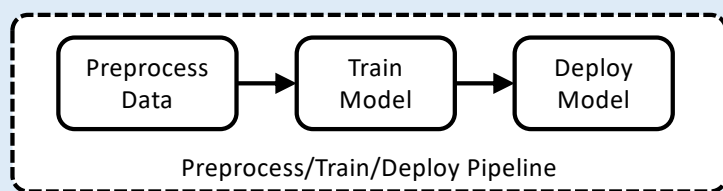
Example Image, taken from:
<https://www.ibm.com/blogs/bluemix/2017/10/serverless-composition-ibm-cloud-functions>

Composable Pipeline Templates

Pipeline Templates

→ Parameterizable templates of common pipeline patterns

- **Example:** Preprocess/Train/Deploy Pipeline
 - Easily bootstrap configurations with default values
 - Users can customize and fine-tune the configuration



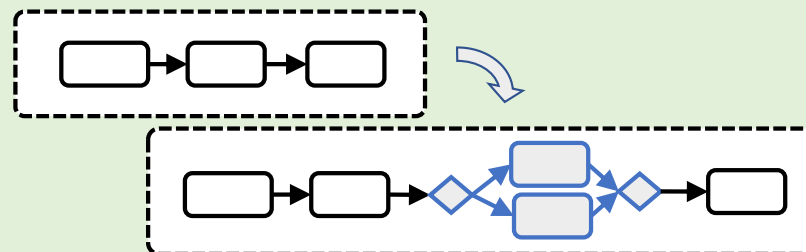
Extensible Catalog of Common Patterns

- Composability becomes critical
- Leverage techniques from BPM, service composition, and configuration management

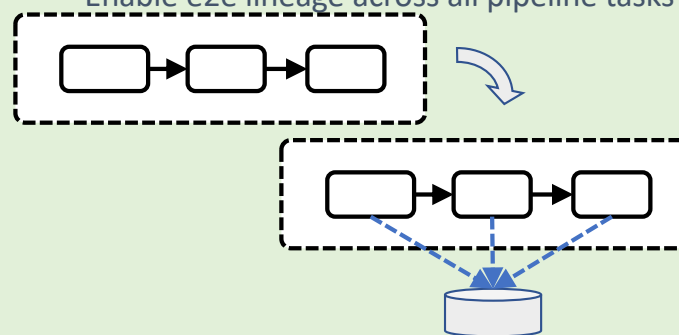
Pipeline Transformers

→ Enrich existing configurations with additional features

- **Example:** Add safeguards for deployment
 - A/B testing with feedback and human in the loop



- **Example:** Plug in cross-cutting features for trusted AI
 - Enable e2e lineage across all pipeline tasks



Simple Demo – Pipeline Templates

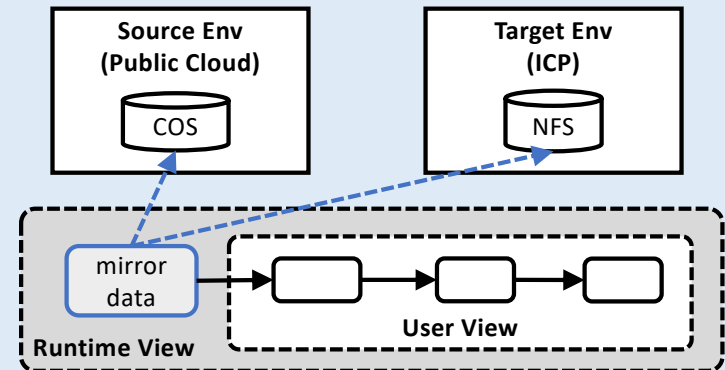
```
bash-3.2$  
bash-3.2$ █
```

Reusable domain abstractions

Use domain knowledge to build smarts into the pipeline

Example:

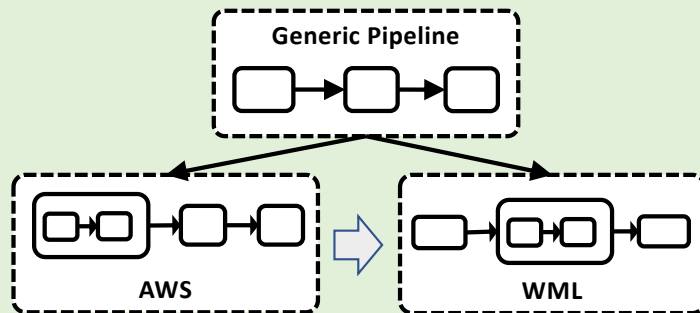
- Seamlessly move between Public Cloud and ICP
- Automatically enrich and adjust the pipeline at runtime
- Managing all artifacts required to run the pipeline



Platform-independent pipeline tasks

Example:

- Generic pipelines allow easy on-ramp into our platform



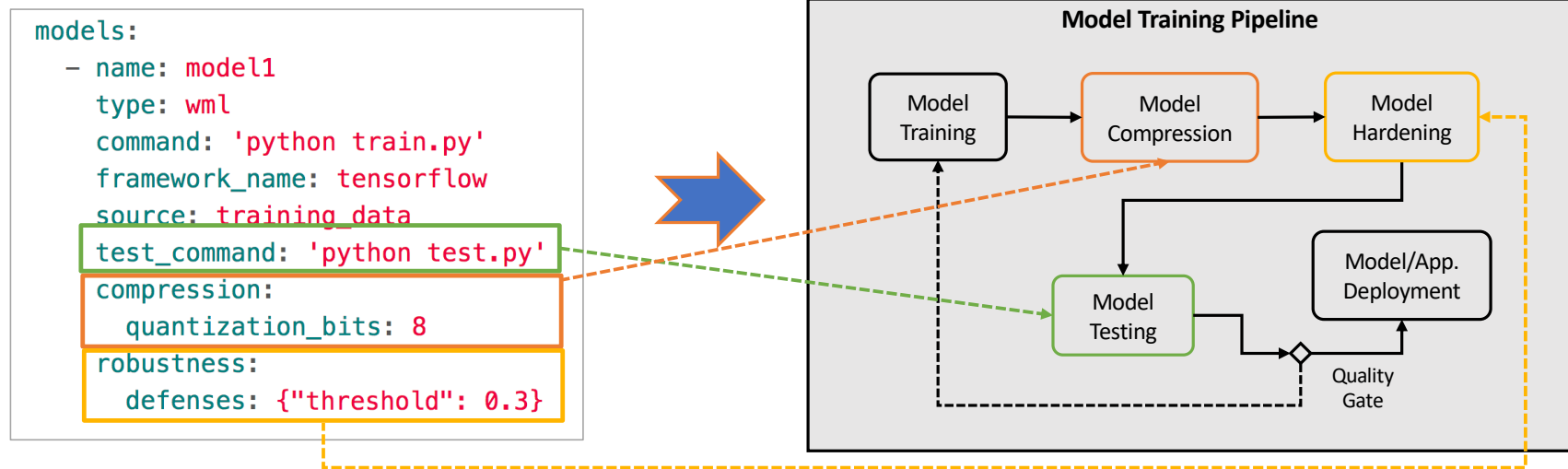
Rollout to mobile and distributed deployments

Use cases:

- VR models on mobile devices with automated retraining loop
- Patient-specific models for predicting hypoglycemia based on real-time data

Testing and Fine-Tuning Classifiers

- **Context:** Model training entails multiple **specialized stages** (e.g., data bias checks[1], model robustness checks [1], model compression)
- **Problem:** Each stage requires custom configuration, input-output mapping, etc, which can be tedious
- **Solution:** Annotate the model entity with desired features – pipeline tasks will get inserted automatically



[1] <https://github.com/IBM/AIF360>

[2] <https://github.com/IBM/adversarial-robustness-toolbox>

Staged Deployments

- **Context:** Models are deployed into business-critical applications
- **Problem:** Need a controlled way to deploy, compare, and update model versions
- **Solution:** Hide new model deployments behind feature flags, gradually roll out the change

```
experiments:  
  - name: a_b_test  
    settings:  
      _type: modelops.routing.ABTest  
      model: model1  
      routing:  
        type: fixed  
        new_model: 0.2  
        old_model: 0.8  
    deployment:  
      _type: modelops.routing.DeployToFlagr  
      endpoint: {{flagr_url}}
```

A/B Testing experiment configuration



Segment ID: 2

Description: Staged Model Deployment Feature Flag

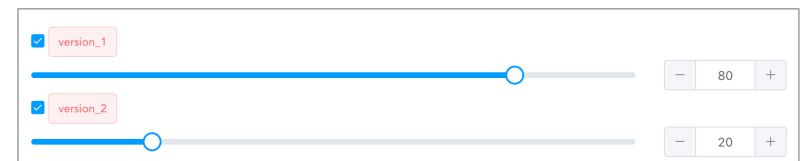
Constraints (match ALL of them): No constraints

Property == Value, e.g. "CA", ["CA", "NY"]

Distribution edit

version_1: 80%

version_2: 20%



Deployed feature flag for different model versions

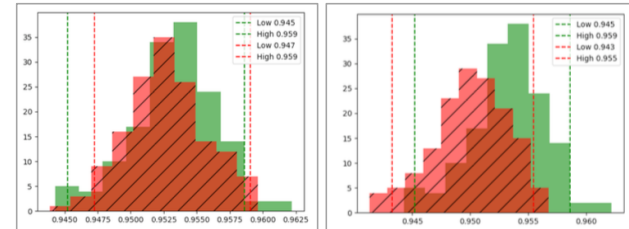
Model Drift Detection

– **Context:** Models operate in **dynamic environments**

→ Data and conditions may change over time

– **Problem:** Models may result in incorrect or **inaccurate predictions** (denoted model or data “drift”)

– **Solution:** Monitor the runtime traffic and raise an alert if drift is detected



```
subscribers:
```

- name: drift-detector
- _type: modelops.drift.DriftDetector
- model: model1
- window_size: 500
- p_value_threshold: 0.05
- actions:
 - notify_slack
 - update_ensemble_classifier

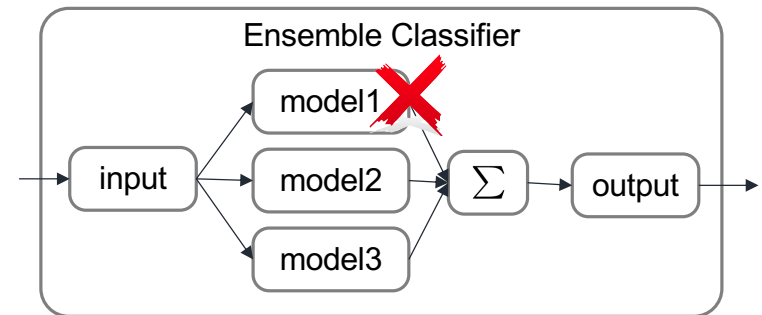
Drift detector configuration



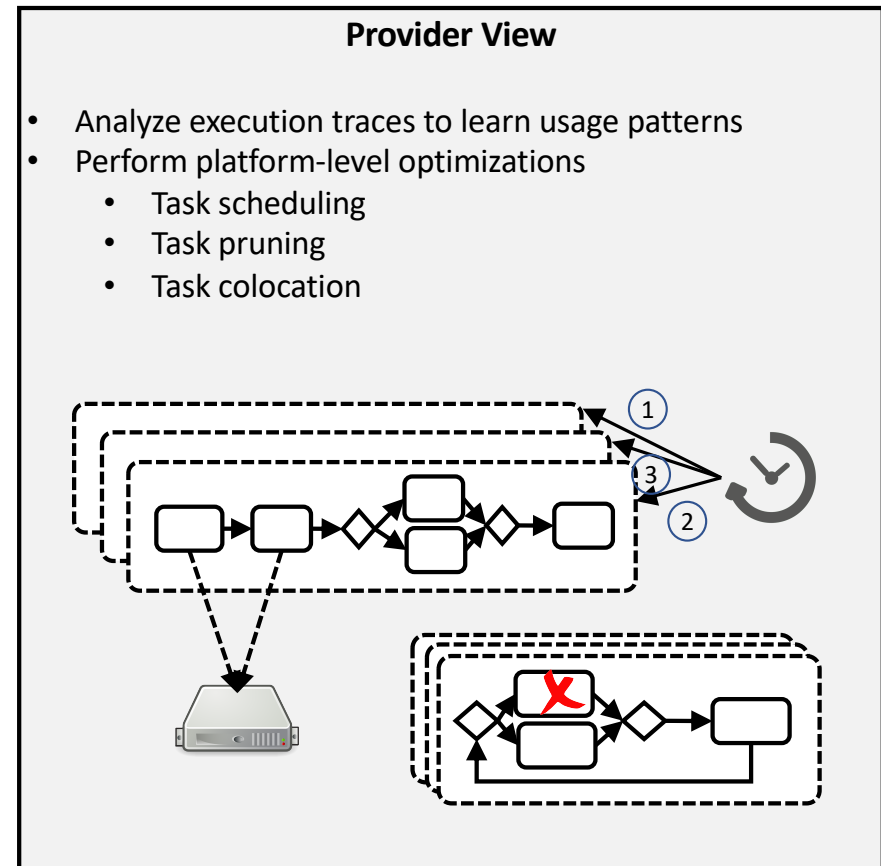
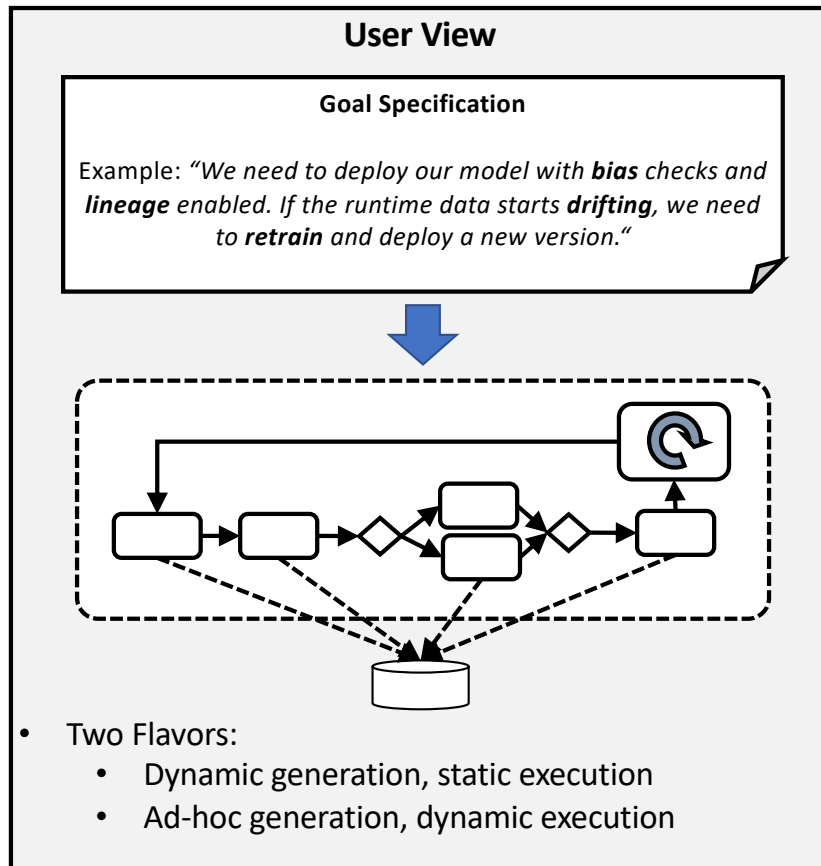
Today

incoming-webhook APP 1:24 PM
Drift detected for model "model1" with p-value 0.034125128

+ Message #modelops-alerts @ 😊



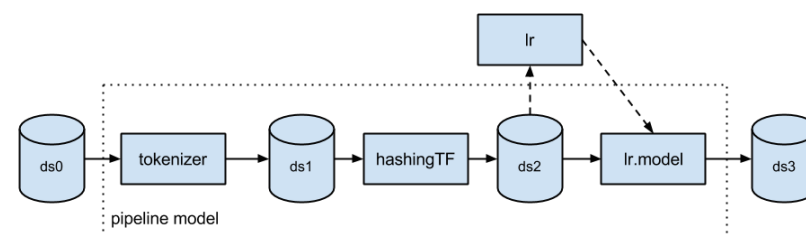
Future Work: Dynamic pipeline generation and large-scale optimization



Related Work

Pipelines in different Machine Learning environments

- ML Pipelines in Spark using MLlib [Meng'16]
- ML Pipelines in scikit-learn [Pedregosa'11]
 - Transformers and estimators



Automated Machine Learning

- Efficient and robust automated machine learning [Feurer'15]
 - AutoML as a Combined Algorithm Selection and Hyperparameter optimization (CASH) problem

Production-grade machine learning platforms

- TFX: A TensorFlow-Based Production-Scale Machine Learning Platform [Baylor'17]
- Michelangelo: Uber's Machine Learning Platform [Hermann'17]

ModelOps – Future Directions

- **Domain Abstractions**
 - Cross-platform model train&deploy (hybrid/multi-cloud)
 - Model versions, staged deployments, A/B testing (e.g., large-scale transfer learning pipelines)
 - Multi-level pipelines in edge scenarios
- **Usability / UX**
 - User study on UX and configuration formats
 - Configuration CLI – (“modelops init”)
- **Large-Scale Pipeline Scheduling & Optimization**
 - Looking at ModelOps pipelines from **provider’s point of view**
 - Ad-hoc pipeline creation and scheduling under resource constraints
 - Simulation, experimentation, and analytics environment to evaluate different strategies
- **Extended Pipelines for Online and Reinforcement Learning**
 - Explore the computational model for stateful checkers, bandits, RL agents

ModelOps

**A programming model for reusable,
platform-independent, and composable AI workflows**

Waldemar Hummer, Vinod Muthusamy

IBM Research AI

IBM Programming Languages Day

T.J. Watson Research Center

December 10, 2018

<https://ibm.biz/plday2018>

Contact: whummer@ibm.com

Team:

Yorktown: Waldemar Hummer, Anupama Murthi,
Kaoutar El Maghraoui, Benjamin Herta, Darrell Reimer,
Punleuk Oum, Gaodan Fang

Austin: Vinod Muthusamy

Cambridge: Scott Boag

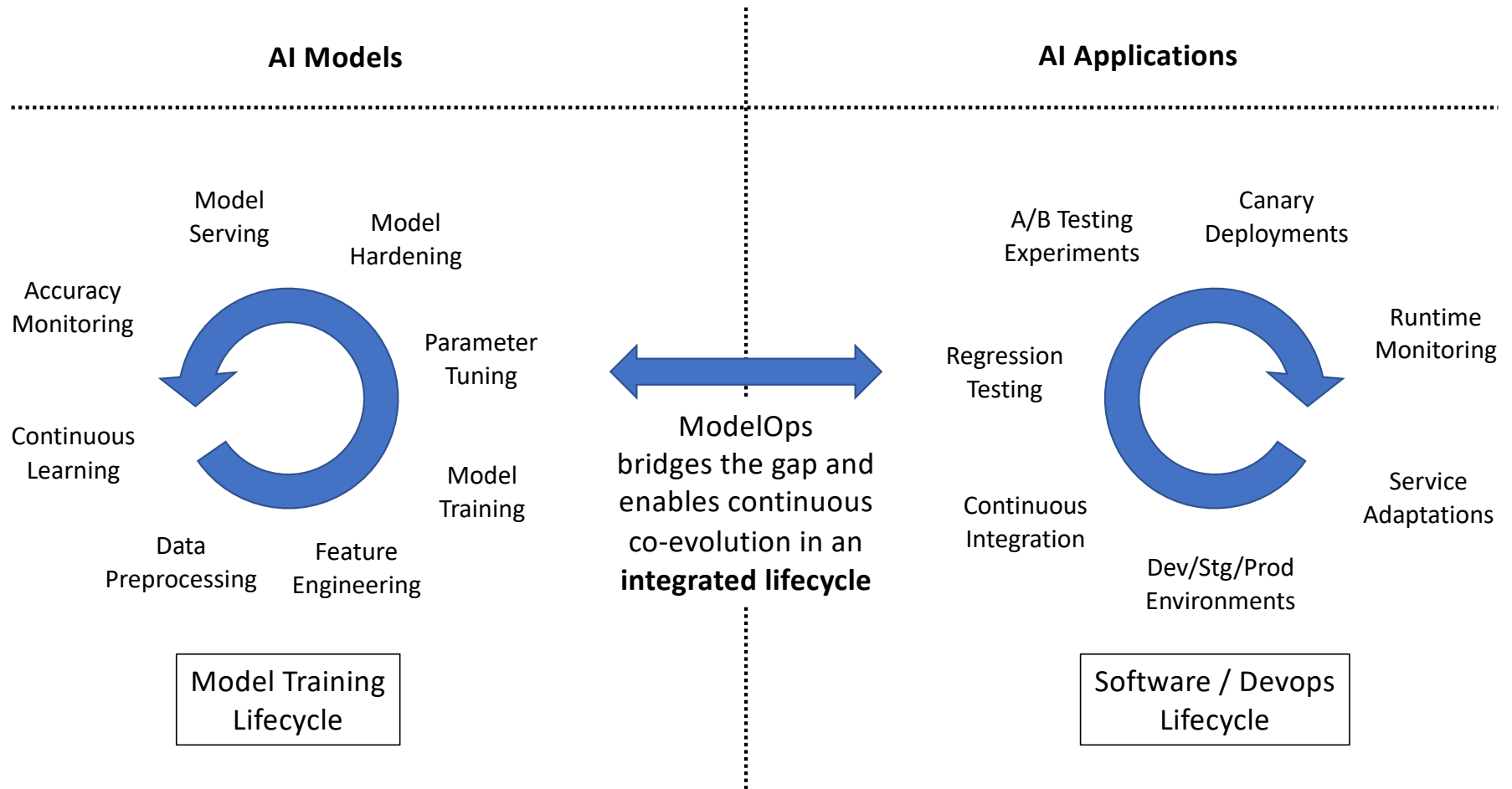
Discussion

References

- Meng, X., Bradley, J., Yavuz, B., Sparks, E., Venkataraman, S., Liu, D., ... & Xin, D. (2016). Mllib: Machine learning in apache spark. *The Journal of Machine Learning Research*, 17(1), 1235-1241.
- Feurer, M., Klein, A., Eggenberger, K., Springenberg, J., Blum, M., & Hutter, F. (2015). Efficient and robust automated machine learning. In *Advances in Neural Information Processing Systems* (pp. 2962-2970).
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Vanderplas, J. (2011). Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct), 2825-2830
- Baylor, D., Breck, E., Cheng, H. T., Fiedel, N., Foo, C. Y., Haque, Z., ... & Koo, C. Y. (2017, August). Tfx: A tensorflow-based production-scale machine learning platform. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1387-1395). ACM.
- Hermann, J., & Del Balso, M. (2017). Meet Michelangelo: Uber's machine learning platform. URL <https://eng.uber.com/michelangelo>

Backup

ModelOps – Mission Statement

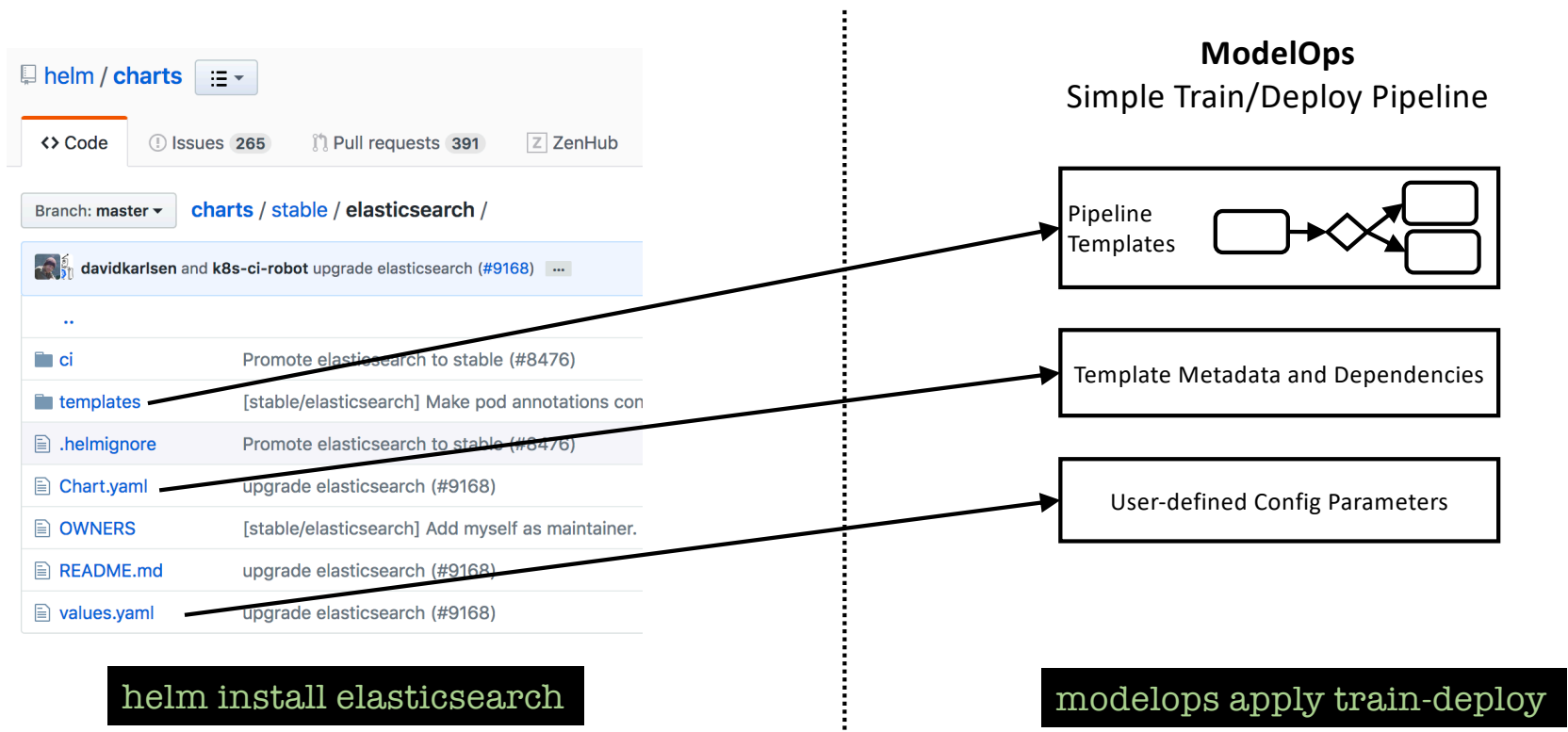


Generic AI Pipeline Representation - Python

```
1  from modelops.plugins.slack import PostMessage as SlackNotification
2
3  # add single datastore to list
4  datastores + DataStore('account_wml', type='ibm/wml', ...)
5
6  # define pipelines
7  pipeline = Pipeline('dsl-demo-flow')
8
9  # define tasks
10 hardening_task = Harden('harden_model', ...)
11 notification_task = SlackNotification('notify_training', ...)
12 deploy_task = Deploy(...)
13 train_task = Train(...)
14
15 # add task
16 pipeline + train_task
17
18 # define a pipeline branch/fork
19 pipeline | (hardening_task, notification_task)
20
21 # define a pipeline join
22 pipeline & ('harden_model', 'notify_training', deploy_task)
23
24 ...
```

Pipeline plugins from the community

- Lifecycle capabilities can be **crowdsourced**, based on the pluggable framework
- Analogy to Helm Charts

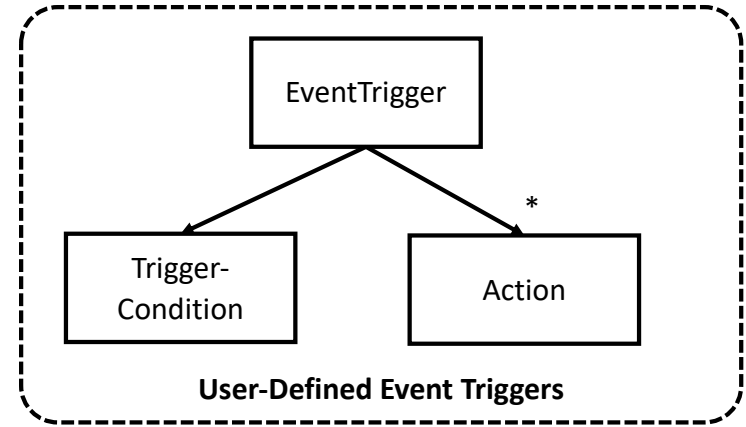
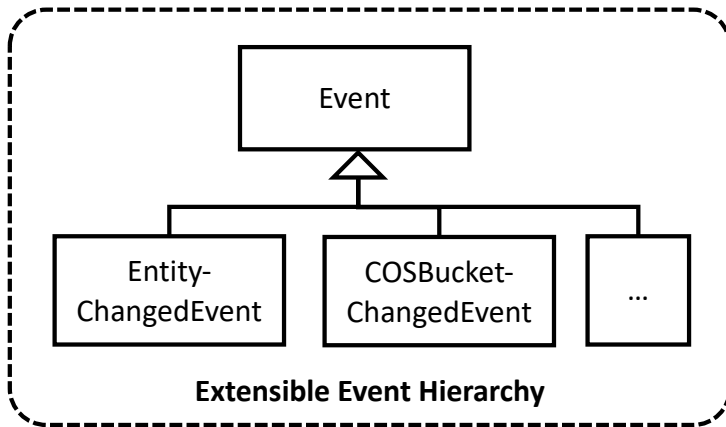


Environment Switches

```
name: test_manifest
environments:
  - name: dev
    imports:
      - type: account
        name: wml_account_dev
        localname: wml_account
  - name: stg
    imports:
      - type: account
        name: wml_account_stg
        localname: wml_account
accounts:
  - name: wml_account_dev
    username: user_dev
  - name: wml_account_stg
    username: user_stg
```

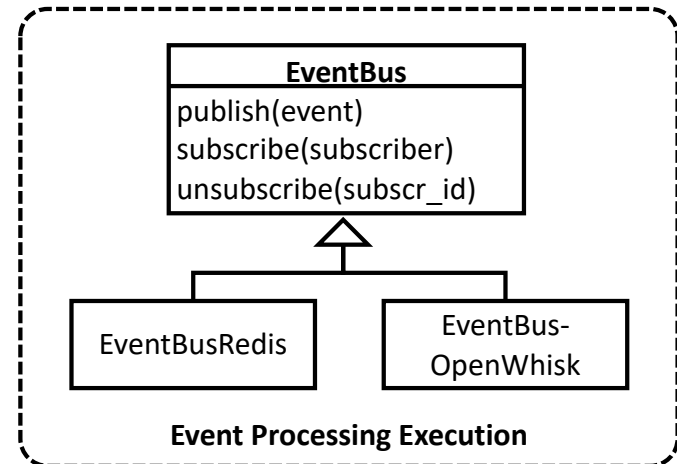
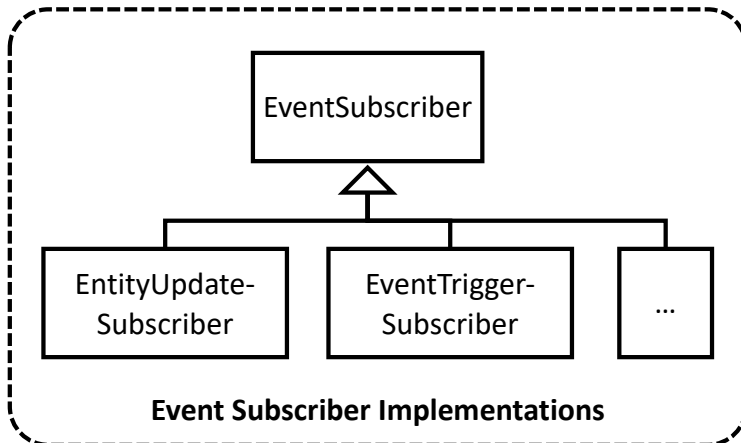
- Import entities into an environment namespace
 - Create a local name alias
 - wml_account_dev → wml_account
 - wml_account_stg → wml_account
 - Local name “wml_account” can then be used in other parts of the configuration
- Activate different environments:
 - **ENV=dev** modelops run
 - **ENV=stg** modelops run

ModelOps – Lifecycle Event Processing



User Facing

Platform Internal



ModelOps – Event Trigger Example

```
triggers:  
  - name: wml-e2e-demo-trigger  
    condition:  
      _type: modelops.eventing.triggers.BucketChangedCondition  
      datastore: wml-e2e-demo-store  
      file_pattern: train.csv  
    actions:  
      - type: execute-pipeline  
        pipeline: wml-e2e-demo-train
```

<https://github.ibm.com/ModelOps/modelops-demo-titanic/blob/master/modelops.yml>