



ArgoCon

CI/CD for Data Pipelines with Argo Workflows

J.P. Zivalich, Pipekit

#argocon2022

@pipekitio



AGENDA

- **Intro & goals for the talk**
- **Background: What is CI/CD for data pipelines?**
- **Problems CI/CD for data pipelines solves**
- **Solutions using Argo Events and Argo Workflows**
- **Demo & How-to**
- **Next steps**
- **Q&A**

Intro & goals for the talk

GOALS FOR THE TALK

- Show a development setup of CI/CD for Data Pipelines using Argo Workflows
- Understand how to test WorkflowTemplates
- Learn some strategies for versioning WorkflowTemplates

ABOUT ME



J.P. Zivalich
CTO & Founder @ [Pipekit](#)

Contributer to the [Argo Workflows](#) project,
ex-Docker Hub maintainer

Focus: Data engineering architecture,
product development



#argocon2022

Background: What is CI/CD for data pipelines?

BACKGROUND: CI/CD FOR DATA PIPELINES

Data engineering commonly pushes to prod without testing

- People push changes to their data transforms and hope for the best
- Ideally should validate these transforms when making changes
- Rollbacks are currently difficult - no concept of versioning for pipeline components

Forward-thinkers are applying CI/CD concepts to data pipelines

- Factoring out transforms and other critical pieces into components
- Running tests on these components during pull requests and other events
- Versioning these components using semantic versioning

PROBLEMS IT SOLVES

◆ Wasted cloud spend

- Data pipelines can run anywhere from several minutes up to over a month
- Hardware intensive compute instances can be very expensive
- A failure occurring late in the data pipeline can result in lots of wasted resources

◆ Data scientist time

- Learn earlier in the development cycle if bugs are being introduced

BACKGROUND: WORKFLOW TEMPLATES

➡ **WorkflowTemplates and ClusterWorkflowTemplates are Argo Workflows's native reuse component**

Why factor items into WorkflowTemplates?

- Factoring to WorkflowTemplates you get...
 - DRY-code
 - Easier troubleshooting
- The following are good candidates to make components
 - Data transforms
 - Set-up/tear-down of K8s resources like Dask/Spark deployments
 - Utilities like cloning git repositories, etc.

TESTING WORKFLOW TEMPLATES

→ **To test, use the inputs/outputs of WorkflowTemplates**

WorkflowTemplates are essentially functions

- Each can take inputs and generate outputs
- For easy testing/components, ensure tested WorkflowTemplates are pure functions
 - The function return values are the same for identical arguments
 - The function has no side effects
- This allows us to compare inputs and actual outputs vs expected outputs

VERSIONING WORKFLOW TEMPLATES

◆ Semantic versioning: Great to have, not in vanilla Argo ◆

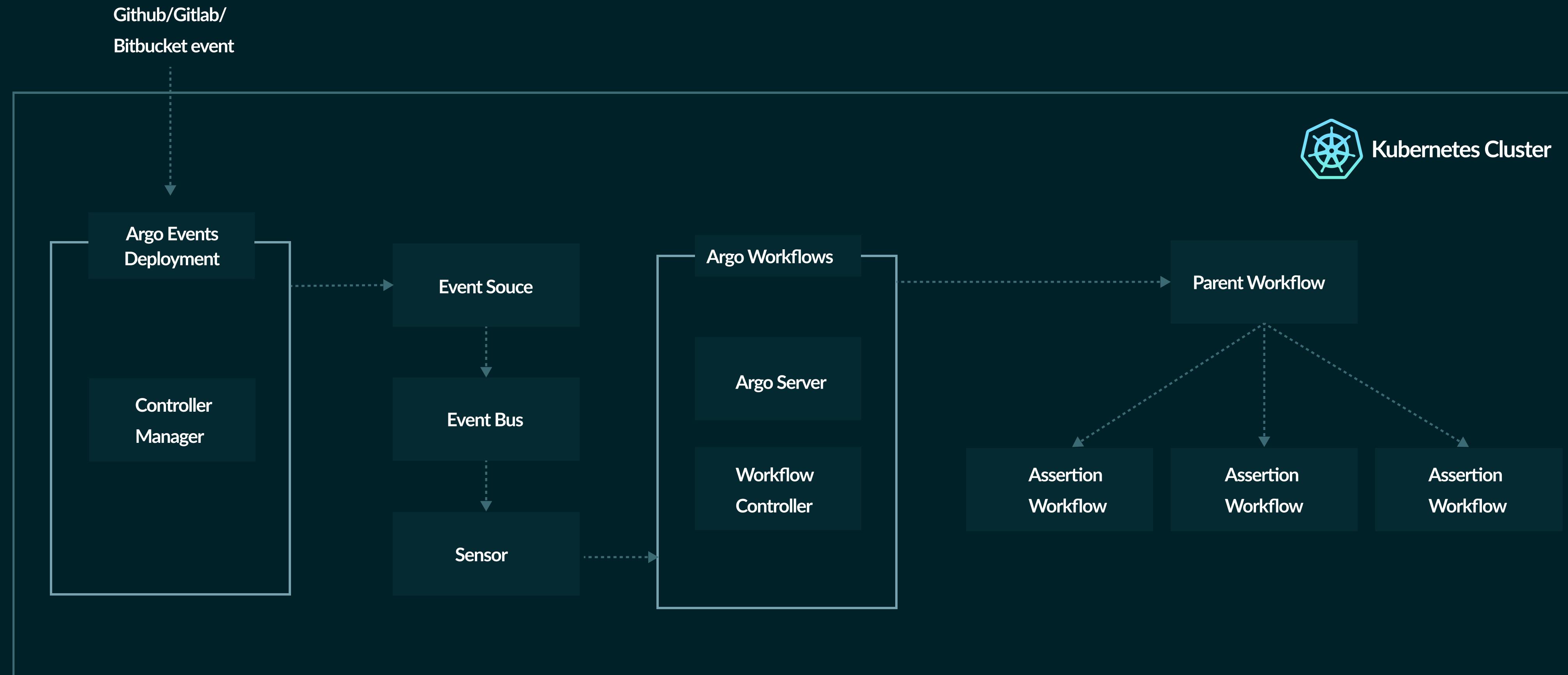
- Semantic versioning allows for structured promotion of components and easy rollbacks
- We've seen two ways of implementing it for WorkflowTemplates in Argo Workflows:
 1. Appending the version to the name with dashes, ie template-12-3-9
 2. Adding a label or annotation denoting the version

Solution



#argocon2022

CI/CD FOR DATA ARCHITECTURE WITH ARGO EVENTS AND WORKFLOWS



#argocon2022

ARGO EVENTS

→ One-time configuration of event source and event bus.
Update the sensor as needed with test cases.

- **Event Source:**
 - Configure your version control provider's webhook
- **Sensor:**
 - Map the version control provider's webhook to an Argo Workflow
 - Process the webhook and extract information such as the PR title, sha, etc.

ARGO WORKFLOWS

- Apply the WorkflowTemplate(s) from a given PR, write assertions, and validate them using a Workflow of Workflows pattern
 - Use a Workflow to call WorkflowTemplates that clone the repository and apply the WorkflowTemplates to be tested
 - To run newly created/updated WorkflowTemplates, use a Workflow of Workflows pattern
 - WorkflowTemplates are incorporated into a given Workflow at compile time
 - Specify the test cases using **withItems** and loop over the given Workflow of Workflows

Demo & How-to



#argocon2022

DEMO & HOW-TO

→ In this demo, we will...

- Create an Argo Event source and sensor that reads Github PRs and runs Argo Workflows
- Test changes made to a WorkflowTemplate in the PR

→ Prerequisites:

- K8s cluster running locally (K3d, Docker Desktop, Minikube, etc.)
- Github repository
- Argo Workflows v3.3 installed

DEMO & HOW-TO

Sample WorkflowTemplate:

```
1 1---  
2 apiVersion: argoproj.io/v1alpha1  
3 kind: WorkflowTemplate  
4 metadata:  
5   name: doubler  
6 spec:  
7   entrypoint: double  
8   templates:  
9     - name: double  
10    inputs:  
11      parameters:  
12        - name: input  
13    script:  
14      command:  
15        - python  
16      image: "python:alpine3.6"  
17      # a simple transform that doubles an input  
18      source: "print(int(\"{{inputs.parameters.input}}") * 2)"  
19  -
```

Assertion WorkflowTemplate:

```
12   - name: run-doubler-and-assert  
13     steps:  
14       # invoke the doubler template with a given input  
15       - - name: run-doubler  
16         templateRef:  
17           name: doubler  
18           template: double  
19           arguments:  
20             parameters:  
21               - name: input  
22                 value: "{{workflow.parameters.input}}"  
23       - - name: assert  
24         template: assert-equal  
25         arguments:  
26           parameters:  
27               - name: input  
28                 value: "{{workflow.parameters.input}}"  
29               - name: actual  
30                 value: "{{steps.run-doubler.outputs.result}}"  
31               - name: expected  
32                 value: "{{workflow.parameters.expected}}"  
33   - name: assert-equal  
34     inputs:  
35       parameters:  
36         - name: input  
37         - name: actual  
38         - name: expected  
39     script:  
40       image: python:alpine3.6  
41       command: [python]  
42       # ensure the actual value matched the expected value  
43       source: |  
44         assert {{inputs.parameters.actual}} == {{inputs.parameters.expected}},  
45         f"for {{inputs.parameters.input}}, expected {{inputs.parameters.expected}}, got {{inputs.parameters.actual}}"  
46  -
```

DEMO & HOW-TO

Test cases:

```
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
    - name: assert-doubler-
      steps:-
        - - name: checkout-pr-
          templateRef:-
            name: git-checkout-pr-
            template: git-checkout-pr-
            arguments:-
              parameters:-
                - name: app_repo-
                  value: talk-demos-
                - name: git_branch-
                  value: "{{workflow.parameters.short_sha}}"-_
        - - name: get-template-
          template: get-template-
        - - name: apply-template-
          template: apply-template-
          arguments:-
            parameters:-
              - name: workflowtemplate-
                value: "{{steps.get-template.outputs.result}}"-_
        - - name: run-doubler-workflow-
          template: run-doubler-workflow-
          arguments:-
            parameters:-
              - name: input-
                value: "{{item.input}}"-_
              - name: expected-
                value: "{{item.expected}}"-_
            withItems: # the test cases to be ran-
              - { input: '2', expected: '4' }-
              - { input: '4', expected: '8' }-
              - { input: '0', expected: '0' }-
```

Workflow of Workflows:

```
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
```

- **name:** run-doubler-workflow-
- inputs:**-
 - **parameters:**-
 - **name:** input-
 - **name:** expected-
- resource:**-
 - **action:** create-
 - **successCondition:** status.phase == Succeeded-
 - **failureCondition:** status.phase in (Failed, Error)-
- manifest:** |-
 - **apiVersion:** argoproj.io/v1alpha1-
 - **kind:** Workflow-
 - **metadata:**-
 - **generateName:** assert-doubler-
 - **spec:**-
 - **arguments:**-
 - **parameters:**-
 - **name:** input-
 - **value:** "{{inputs.parameters.input}}"-
 - **name:** expected-
 - **value:** "{{inputs.parameters.expected}}"-
 - **# invoke the assertion template with given arguments-**
 - **workflowTemplateRef:**-
 - **name:** assert-doubler-



ArgoCon

#argocon2022

DEMO & HOW-TO

The output of the parent Workflow and children in the UI:

The screenshot shows the Argo Workflows UI interface. On the left is a sidebar with icons for Home, Workflows, Pipelines, Triggers, Metrics, CronWorkflows, and Help. The main area is titled "Workflows / default". It features a search bar and a "SUBMIT NEW WORKFLOW" button. A sidebar on the left lists filters for "NAMESPACE" (set to "default"), "LABELS", "WORKFLOW TEMPLATE", and "CRON WORKFLOW". The main table displays four workflow entries:

	NAME	NAMESPACE	STARTED	FINISHED	DURATION	PROGRESS	MESSAGE	DETAILS
<input type="checkbox"/>	assert-doubler-4tpkd	default	56m ago	56m ago	41s	2/2	-	SHOW ▾
<input type="checkbox"/>	assert-doubler-hp7k7	default	56m ago	56m ago	41s	2/2	-	SHOW ▾
<input type="checkbox"/>	assert-doubler-wgnd5	default	56m ago	56m ago	31s	2/2	-	SHOW ▾
<input type="checkbox"/>	github-doubler-test-4-54ca501	default	57m ago	55m ago	1m	6/6	-	SHOW ▾

At the bottom are navigation buttons for "FIRST PAGE" and "NEXT PAGE >". A dropdown menu shows "500 results per page".

Next steps



#argocon2022

NEXT STEPS

➡ What does this unlock for us?

- Audit trail of workflows run with different component versions
- Assertions that correspond to cases found in production runs
- Clear guidelines for promoting WorkflowTemplates to higher environments

Resources & References



#argocon2022

RESOURCES & REFERENCES

- ➡ Visit our repo: <https://github.com/pipekit/talk-demos/tree/main/argocon-demos/2022-ci-cd-for-data>
- ➡ Visit our blog: <https://pipekit.io/blog/>
- ➡ Follow [@pipekitio](#) on Twitter



ArgoCon

