# INTROS

pipekit

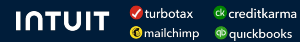INTUIT · turbotax · creditkarma · mailchimp · quickbooks

## Alan Clucas

Senior Software Engineer @ Pipekit.io

Argo Workflows maintainer

## Julie Vogelman

Staff Software Engineer @ Intuit

Argo Workflows maintainer

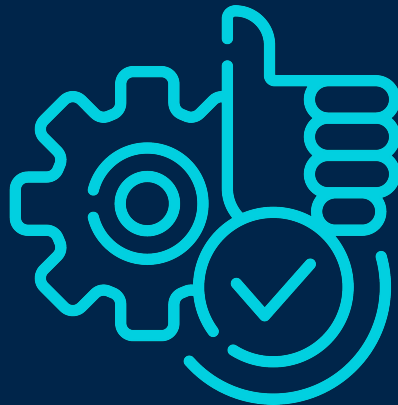ArgoCon

# GOALS

**Saving time**
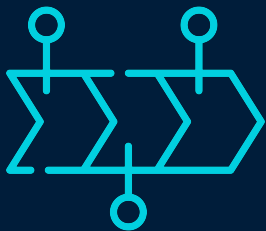
**Saving Cost**

**Don't repeat work.**

**If we know the answer already, use that...**

# WHAT THINGS

**Individual steps/
tasks between
workflows**

**Templates**

**(DAG, steps)**

# Terminology

**Inputs and outputs**

**Parameters and artifacts**

**Memoization and Work Avoidance**

# Introducing memoization

```yaml
apiVersion: argoproj.io/v1alpha1
kind: Workflow
metadata:
    generateName: memoized-workflow-
spec:
    entrypoint: whalesay
    templates:
      - name: whalesay
        memoize:
            key: "{{inputs.parameters.message}}"
            maxAge: "10s"
            cache:
                configMap:
                    name: whalesay-cache
```

ArgoCon

# What happens?

- Argo looks up keys in a cache

- If they match...

  - Remembers the outputs and provides them instead of running the step

| DURATION | | 17s |
|---|---|---|
| PROGRESS | | 1/1 |
| MEMOIZATION | KEY | hi-there-world |
| | CACHE NAME | whalesay-cache |
| | HIT? | NO |

| DURATION | | 0s |
|---|---|---|
| PROGRESS | | 1/1 |
| MEMOIZATION | KEY | hi-there-world |
| | CACHE NAME | whalesay-cache |
| | HIT? | YES |

ArgoCon

# This also works for output artifacts

# You can memoize whole templates

Workflow skipping an entire DAG template:

# ConfigMap

Created for you by the controller

Usually limited to 1MiB

Contains 'human readable' JSON

ArgoCon

# What if my step outputs to somewhere other than an artifact or parameter?

dna-sequencer-45czr

workflow.parameters.person="Sally"

Step A:
Compute DNA

Step B

Sally's DNA

DB

dna-sequencer-18ab4

workflow.parameters.person="Sally"

Step A:
???

Step B

# What if my step outputs to somewhere other than an artifact or parameter?

# What if my step outputs to somewhere other than an artifact or parameter?

Memoization works for this in latest

Alternative: use a technique called "Work Avoidance" instead

# Avoiding work sounds pretty good 😅
## ...but what is it?

Pod logic checks to see if the data exists before recomputing

Or can use a separate marker file to indicate data written (in PVC, DB, etc) and check for that

dna-sequencer-45czr   workflow.parameters.person="Sally"   Step A: Compute DNA   Step B

Sally's DNA

DB

Is Sally's DNA already there?

dna-sequencer-18ab4   workflow.parameters.person="Sally"   Step A: nothing for me to do!   Step B

# Any other time I can't use memoization?

If your data's too big for a ConfigMap (> 1 MiB)

Use "Work Avoidance" instead

# Caching

There are only two hard things in Computer Science: cache invalidation and naming things.

-- Phil Karlton

ArgoCon

# Guaranteed safe

*Pure* steps

- Inputs - all as key
    - Outputs are derived only from inputs
- Has no external interactions
    - Nothing is read from the outside world
    - Changes nothing about the outside world, no side effects

ArgoCon

# Otherwise – choose keys wisely



**CAUTION**

THIS MACHINE
HAS NO BRAIN
USE YOUR OWN

Will skip step when key(s) match

And time since storing < maxAge

ArgoCon

# The End

Come have a chat with Alan about memoization this week.

@Alan Clucas on CNCF Slack

@Joibel on GitHub

Further reading:

https://argoproj.github.io/argo-workflows/memoization/

https://argoproj.github.io/argo-workflows/work-avoidance/

https://s.pipekit.io/chat-argo-wf

ArgoCon