

# Themelios: a model-checked reimplementation of Kubernetes

Andrew Jeffery, Richard Mortier

18th April 2024 @UKSys

# Themelios: Breaking it down

A model-checked reimplementation of Kubernetes



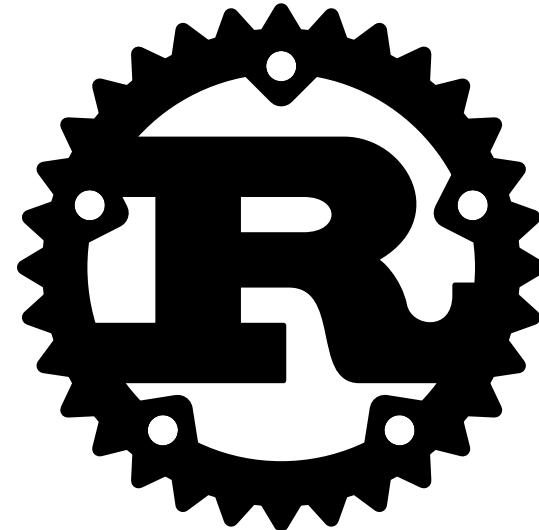
UNIVERSITY OF  
CAMBRIDGE

[andrew.jeffery@cst.cam.ac.uk](mailto:andrew.jeffery@cst.cam.ac.uk)

2

# Themelios: Breaking it down

A model-checked reimplementation of Kubernetes



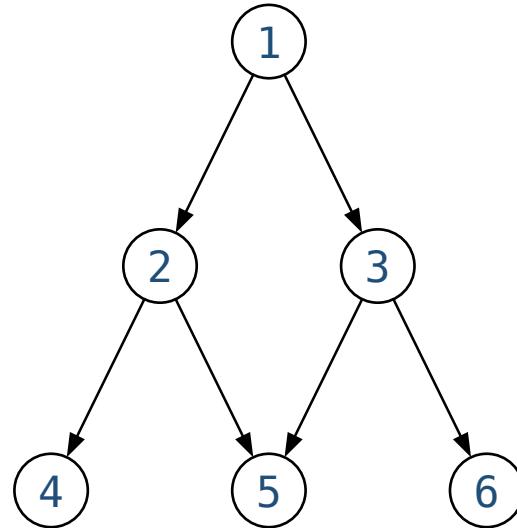
UNIVERSITY OF  
CAMBRIDGE

[andrew.jeffery@cst.cam.ac.uk](mailto:andrew.jeffery@cst.cam.ac.uk)

2

# Themelios: Breaking it down

A model-checked reimplementation of Kubernetes



# Kubernetes isn't the *only* one



Kubernetes



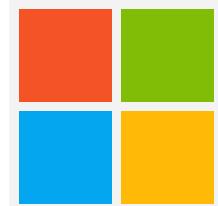
Mesos



HashiCorp  
**Nomad**



Twine



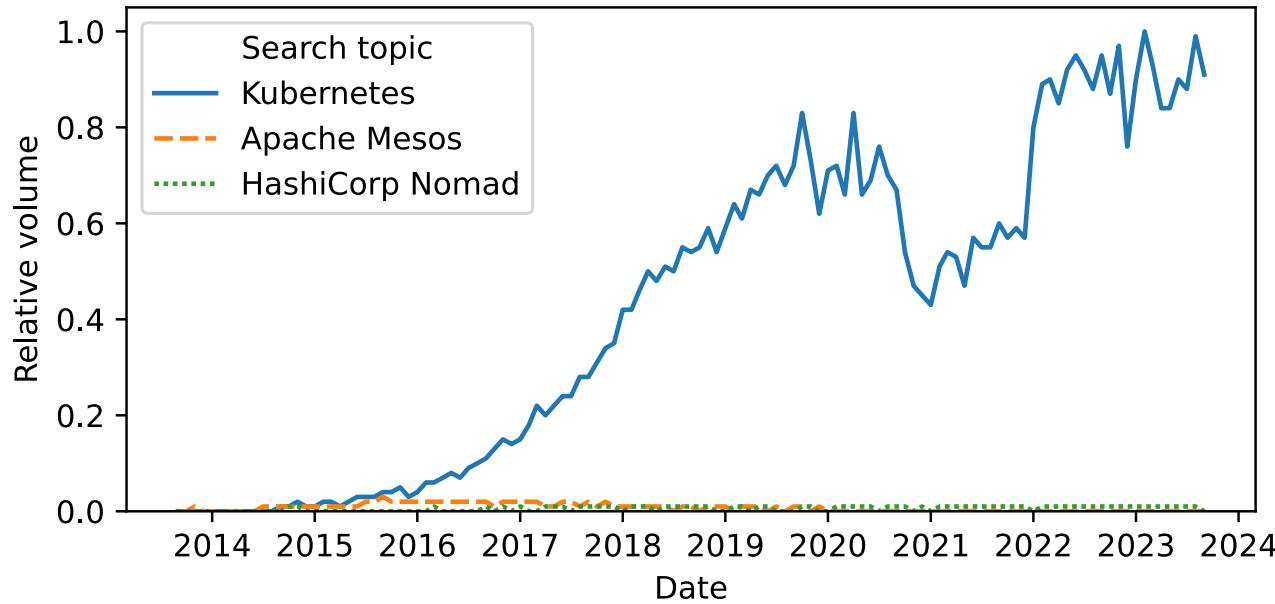
Autopilot



Borg



# But Kubernetes is *the* one



# But what is going on?

Kubernetes is vulnerable to stale reads, violating critical pod safety guarantees #59848

 Open

**smarterclayton** opened this issue on Feb 14, 2018 · 89 comments

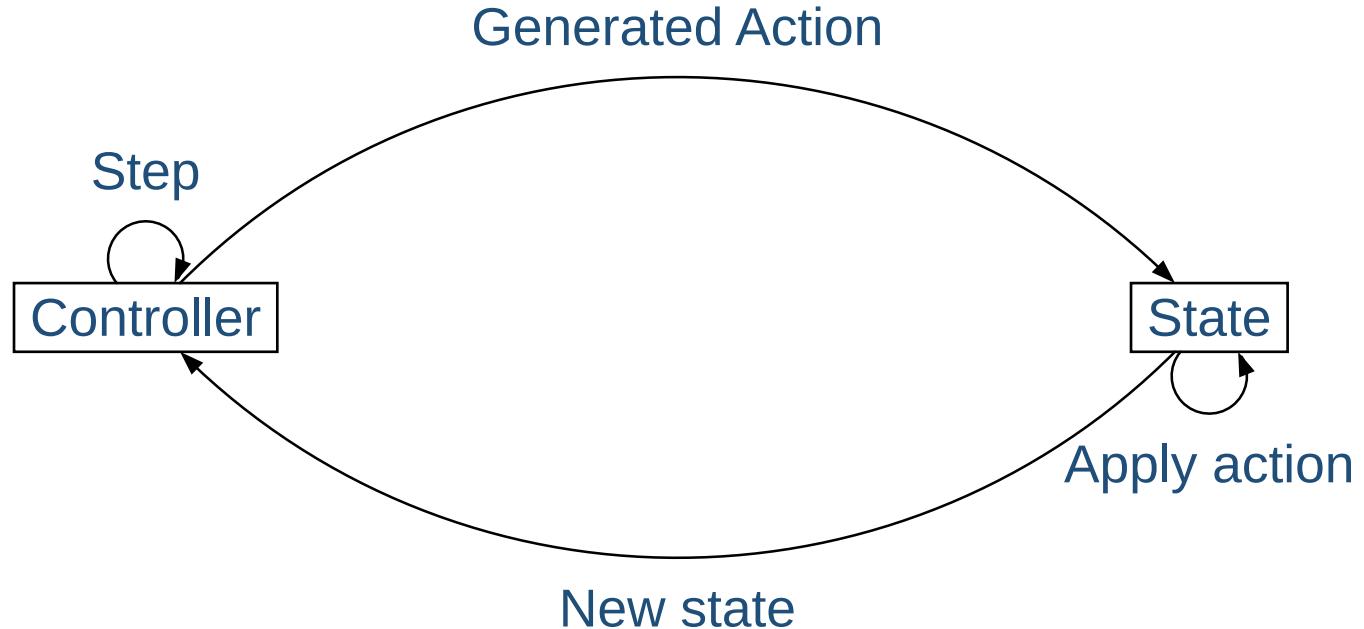


UNIVERSITY OF  
CAMBRIDGE

andrew.jeffery@cst.cam.ac.uk

5

# The model



# The model — controllers

```
trait Controller {  
    fn step(&mut self, // local state  
           gs: &GlobalState) // etcd state  
           -> Option<Action>; // what to do  
}  
  
// typically managed by an API server  
impl GlobalState {  
    fn apply(&mut self, action: Action) {  
        // ...  
    }  
}
```



# The model — state

State

```
{  
  "pods": {},  
  
  "nodes": {  
    "node-0": []  
  },  
}
```

Actions

# The model — state

State

```
{  
  "pods": {  
    "llama2": {  
      "image": "hf/llama:2"  
      "node": null,  
    }  
  },  
  "nodes": {  
    "node-0": []  
  },  
}
```

Actions

- CreatePod("llama2")



UNIVERSITY OF  
CAMBRIDGE

# The model — state

State

```
{  
  "pods": {  
    "llama2": {  
      "image": "hf/llama:2"  
      "node": "node-0",  
    }  
  },  
  "nodes": {  
    "node-0": []  
  },  
}
```

Actions

- CreatePod("llama2")
- SchedulePod("llama2", "node-0")



UNIVERSITY OF  
CAMBRIDGE

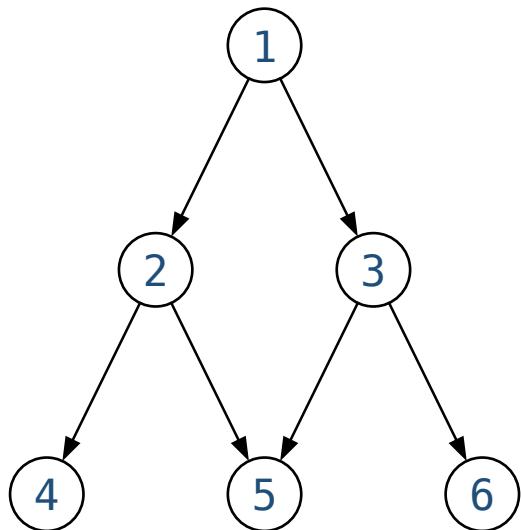
# The model — state

| State  | Actions  |
|--|--|
| <pre>{<br/>  "pods": {<br/>    "llama2": {<br/>      "image": "hf/llama:2"<br/>      "node": "node-0",<br/>    }<br/>  },<br/>  "nodes": {<br/>    "node-0": ["llama2"]<br/>  },<br/>}</pre> | <ul style="list-style-type: none"><li>• CreatePod("llama2")</li><li>• SchedulePod("llama2", "node-0")</li><li>• RunPod("node-0", "llama2")</li></ul> |

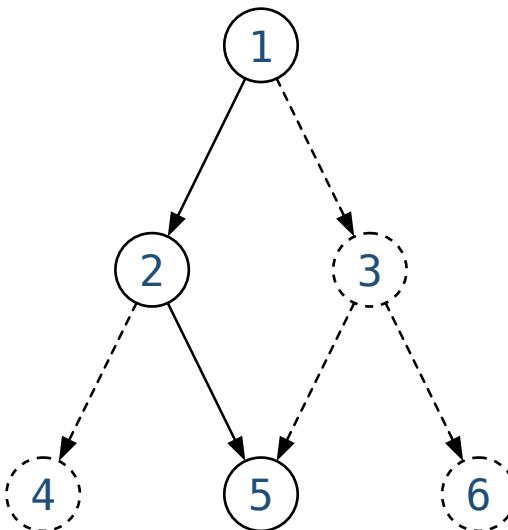


# What can we do with it? Check it!

Exhaustive



Simulation



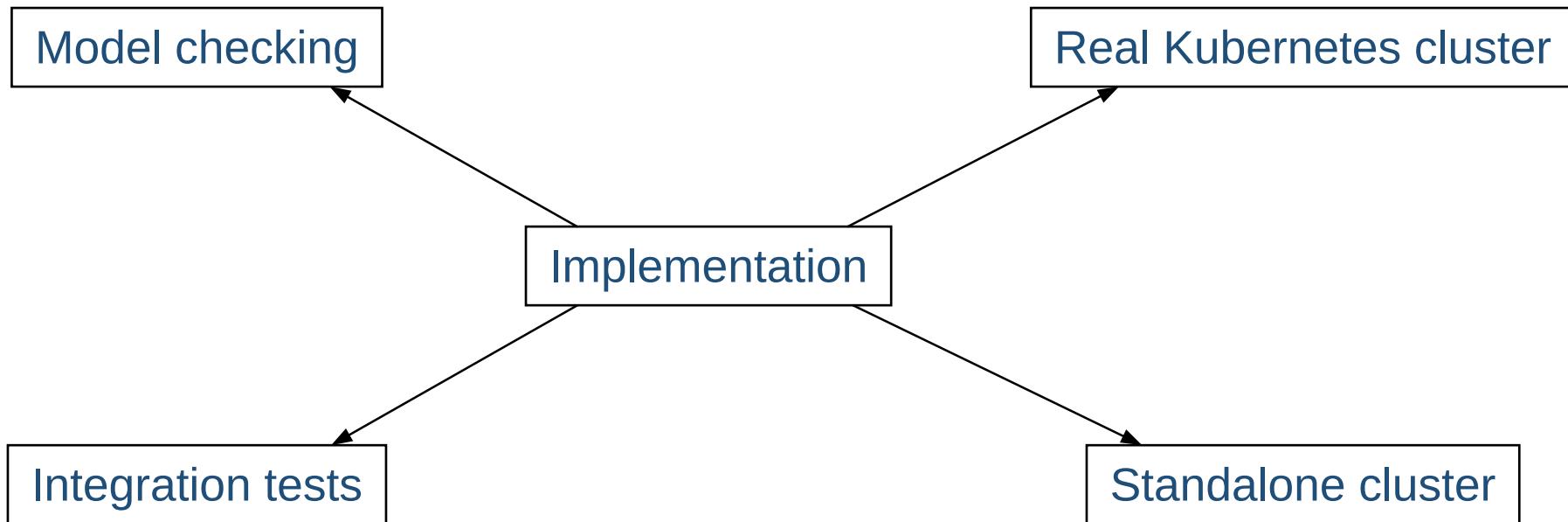
For the stale reads property:

Pods running on Nodes **always** have unique names across the cluster



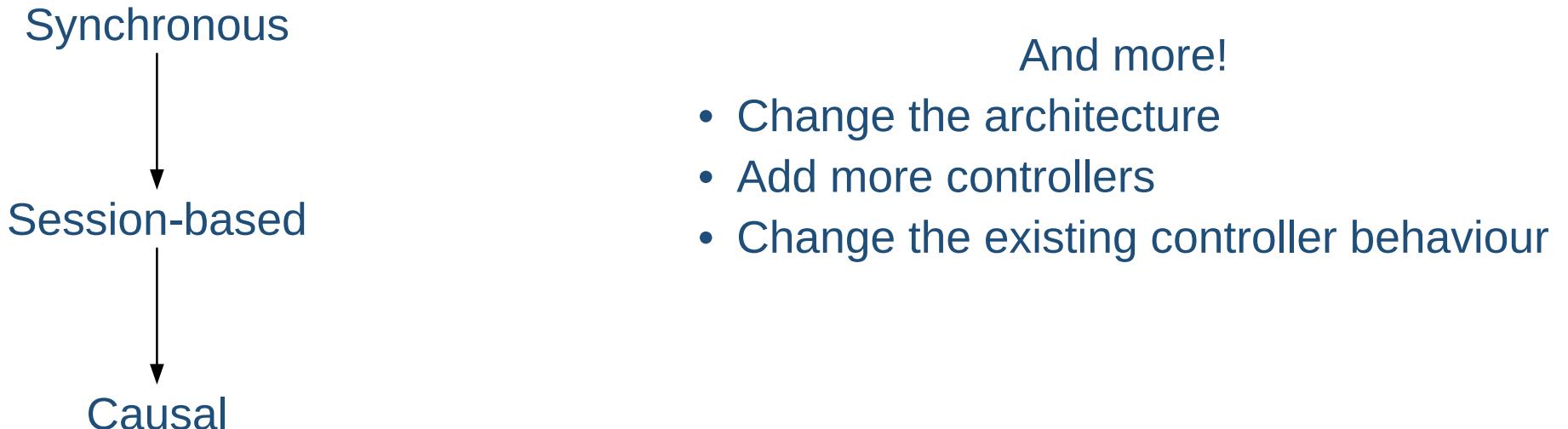
# What can we do with it? Run it!

Single, checked implementation is used everywhere



# What can we do with it? Adapt it!

Weaken the consistency model  
(swap out etcd)



# And how fast does it check?

| Consistency        | Total states explored | $\mu\text{s per state}$ |
|--------------------|-----------------------|-------------------------|
| Synchronous        | 24,758,530            | 2.464                   |
| Monotonic Session  | 23,819,028            | 2.561                   |
| Resettable Session | 23,461,999            | 2.600                   |
| Causal             | 10,100,172            | 6.040                   |

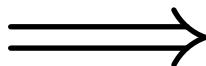
Best results of checking runs of 60s, max depth of 100



# Let's start building strong foundations



Kubernetes



Themelios



UNIVERSITY OF  
CAMBRIDGE

[andrew.jeffery@cst.cam.ac.uk](mailto:andrew.jeffery@cst.cam.ac.uk)

13