



Kubernetes & Google Container Engine Overview

DevOps Meetup Singapore

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

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

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

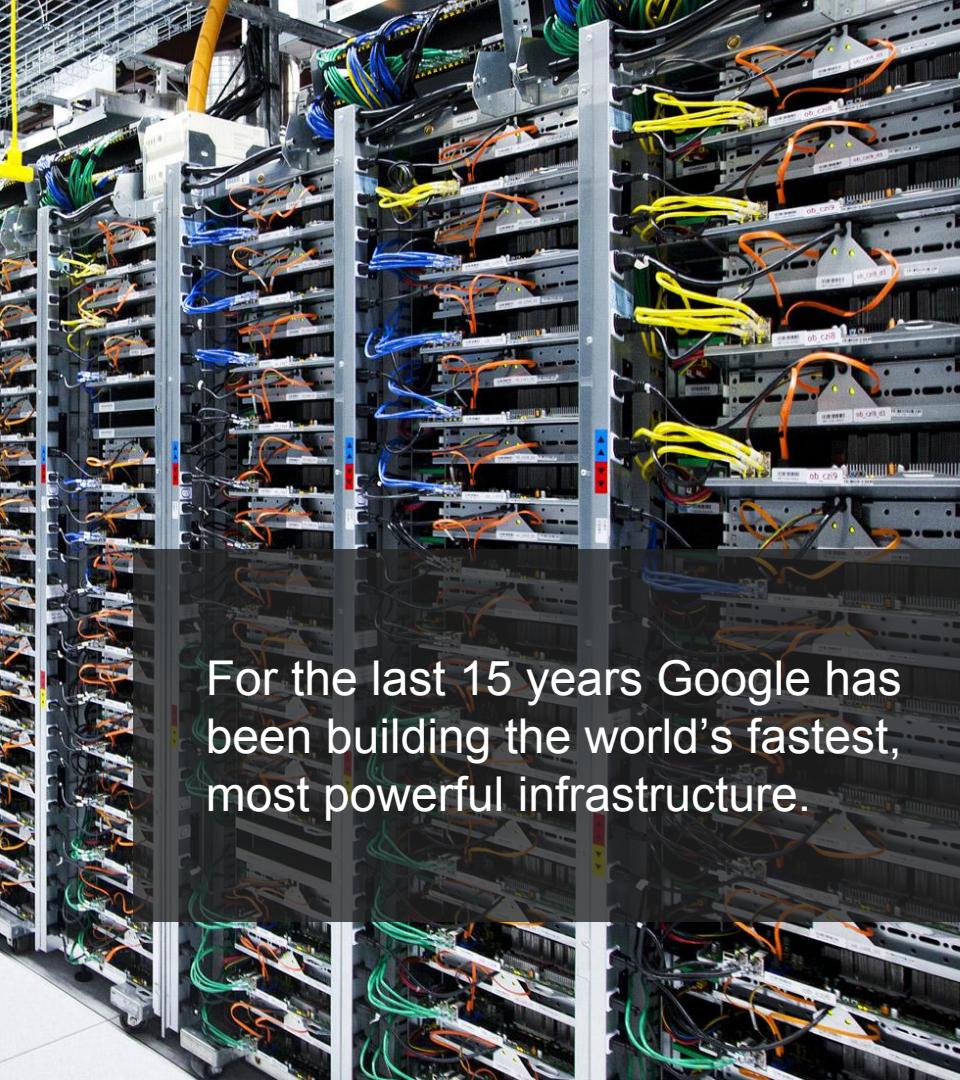
351

352



Ian Lewis
Developer Advocate - Google Cloud Platform
Tokyo, Japan

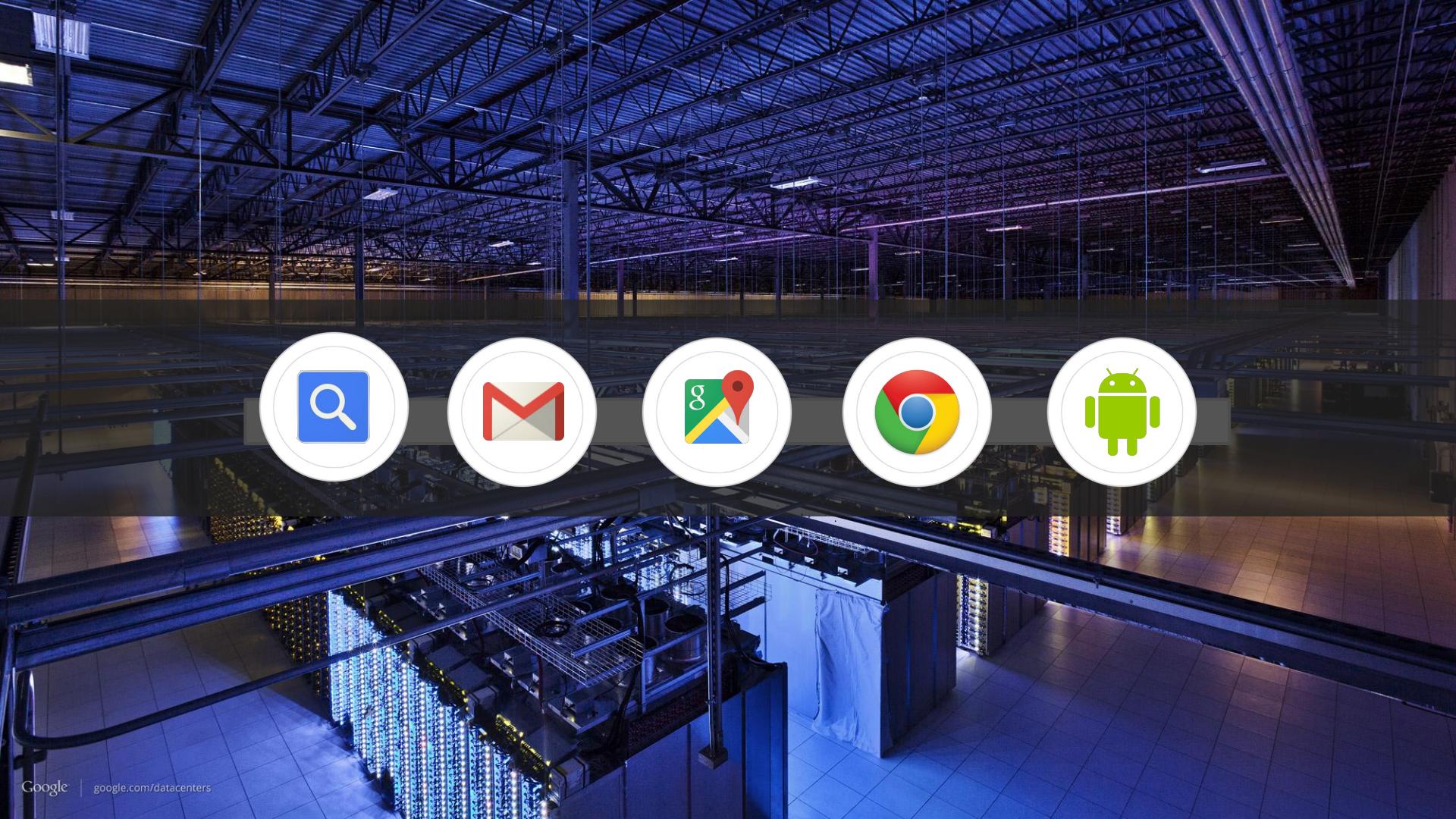
@IanMLewis



For the last 15 years Google has
been building the world's fastest,
most powerful infrastructure.

A satellite image of the Earth at night, showing city lights as white and yellow dots against the dark oceans. The image is overlaid with a network of glowing blue lines, representing Google's world-spanning backbone network. The network density is higher in North America and Europe, with a prominent line running from the US through the Atlantic Ocean to Europe.

Google's World Spanning Backbone Network





Google-Grade Networking

33 Countries
70 Edge Locations

The most of any Cloud Provider



Google Cloud Platform



Monitoring

Big
Data



Storage



Compute



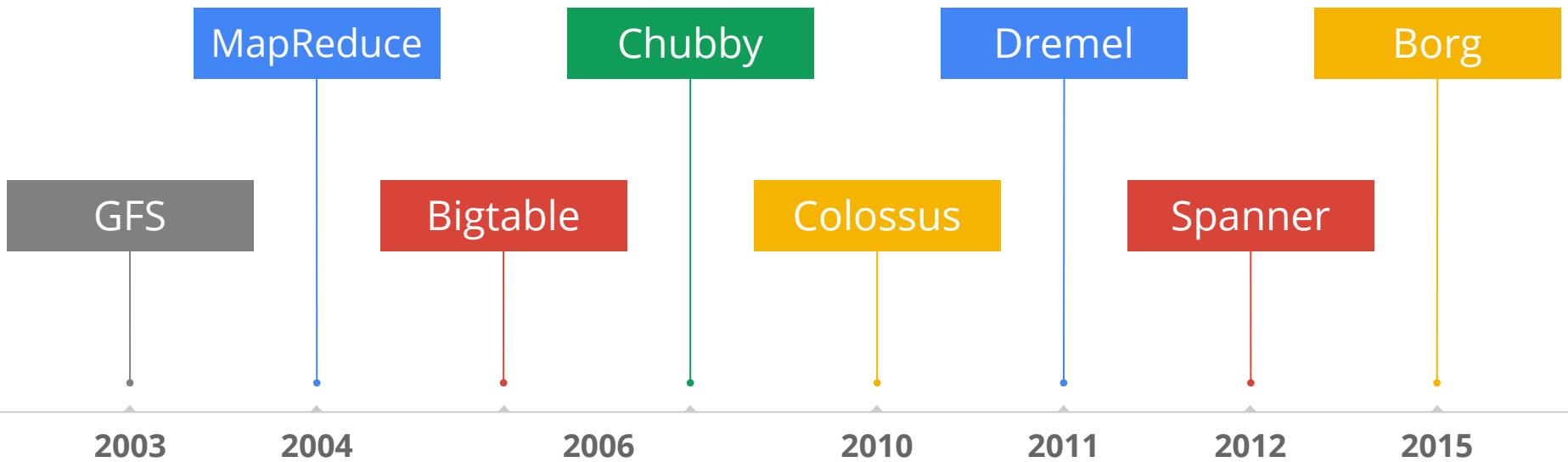
Network



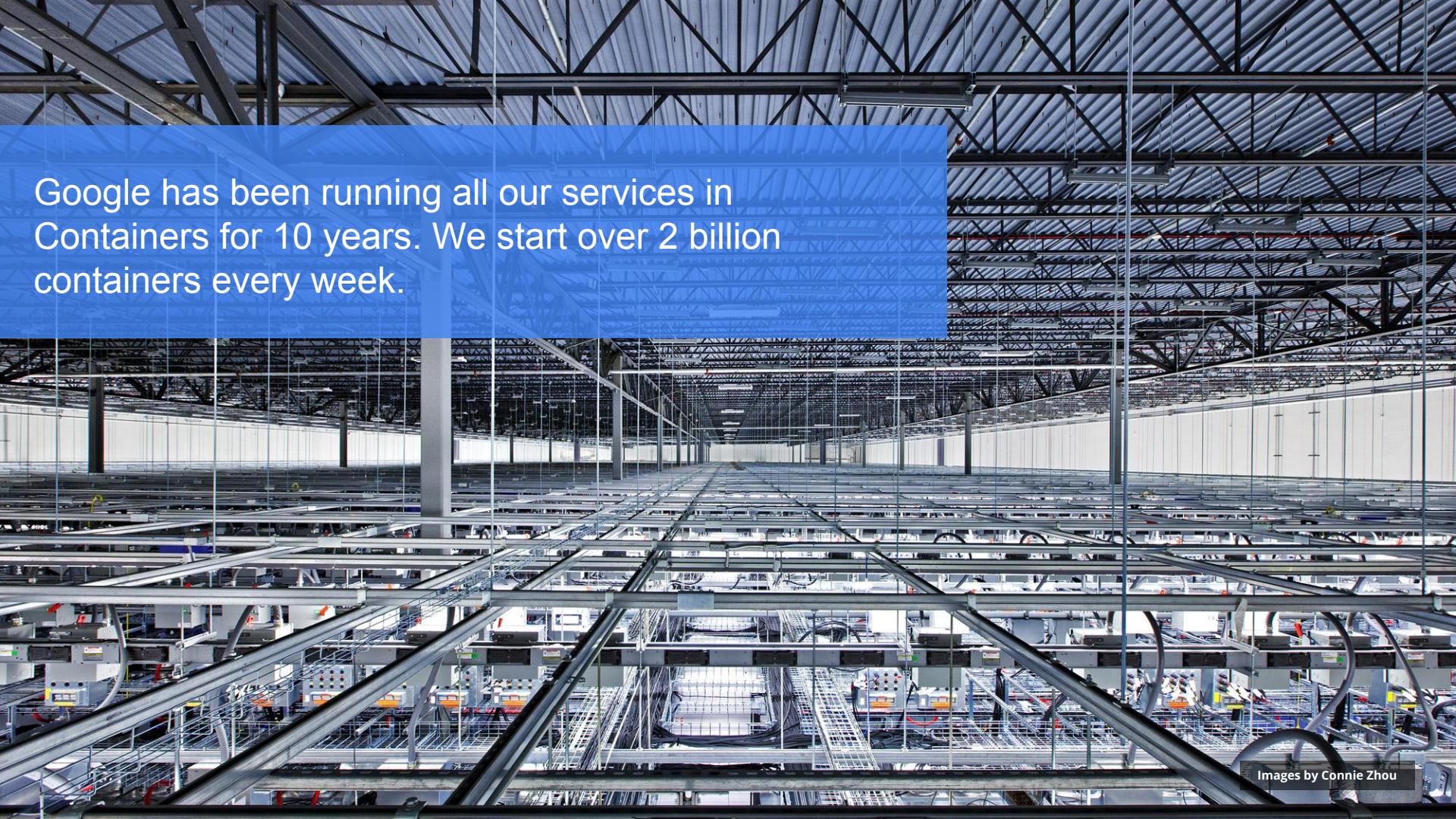
Development



Mobile



Google has been running all our services in Containers for 10 years. We start over 2 billion containers every week.



Large-scale cluster management at Google with Borg

Abhishek Verma[†] Luis Pedrosa[‡] Madhukar Korupolu

David Oppenheimer Eric Tune John Wilkes

Google Inc.

<http://research.google.com/pubs/pub43438.html>

Abstract

Google's Borg system is a cluster manager that runs hundreds of thousands of jobs, from many thousands of different applications, across a number of clusters each with up to tens of thousands of machines.

It achieves high utilization by combining admission control, efficient task-packing, over-commitment, and machine sharing with process-level performance isolation. It supports high-availability applications with runtime features that minimize fault-recovery time, and scheduling policies that reduce the probability of correlated failures. Borg simplifies life for its users by offering a declarative job specification language, name service integration, real-time job monitoring, and tools to analyze and simulate system behavior.

We present a summary of the Borg system architecture

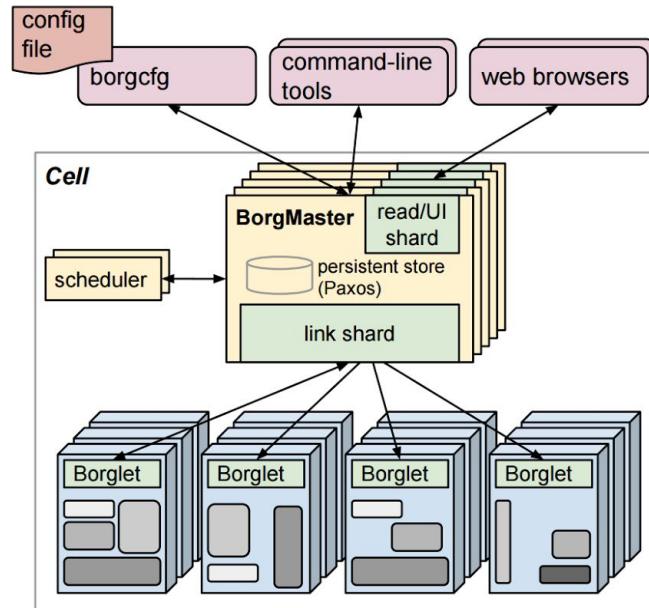




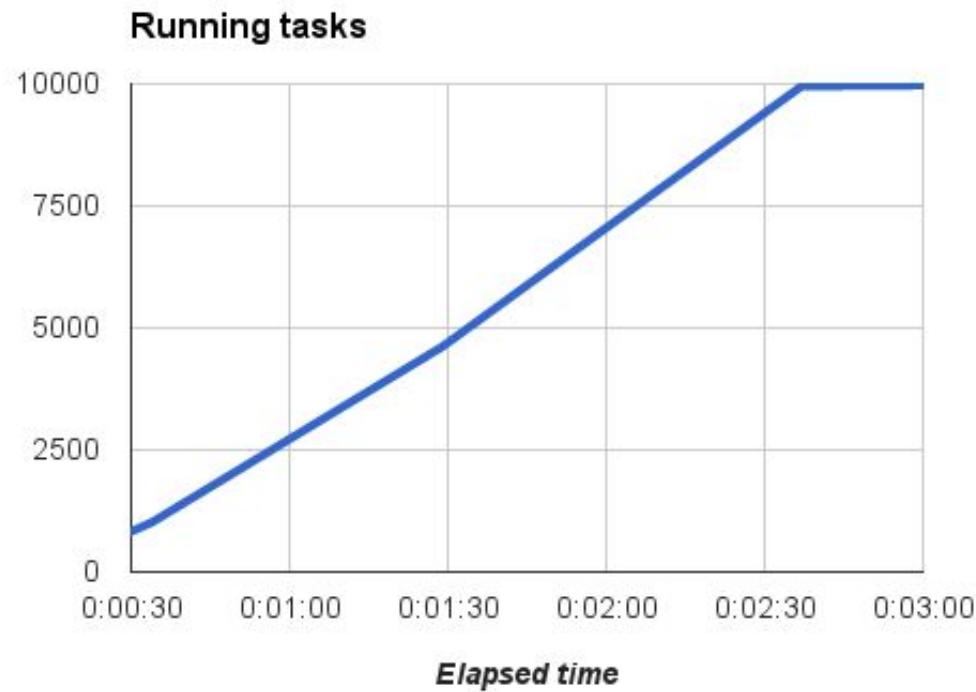
Image by Connie Zhou

Developer View

```
job hello_world = {
    runtime = { cell = 'ic' }                      // Cell (cluster) to run in
    binary = '.../hello_world_webserver'           // Program to run
    args = { port = '%port%' }                     // Command line parameters
    requirements = {                               // Resource requirements
        ram = 100M
        disk = 100M
        cpu = 0.1
    }
    replicas = 10000     // Number of tasks
}
```



Developer View



Developer View

What just
happened?

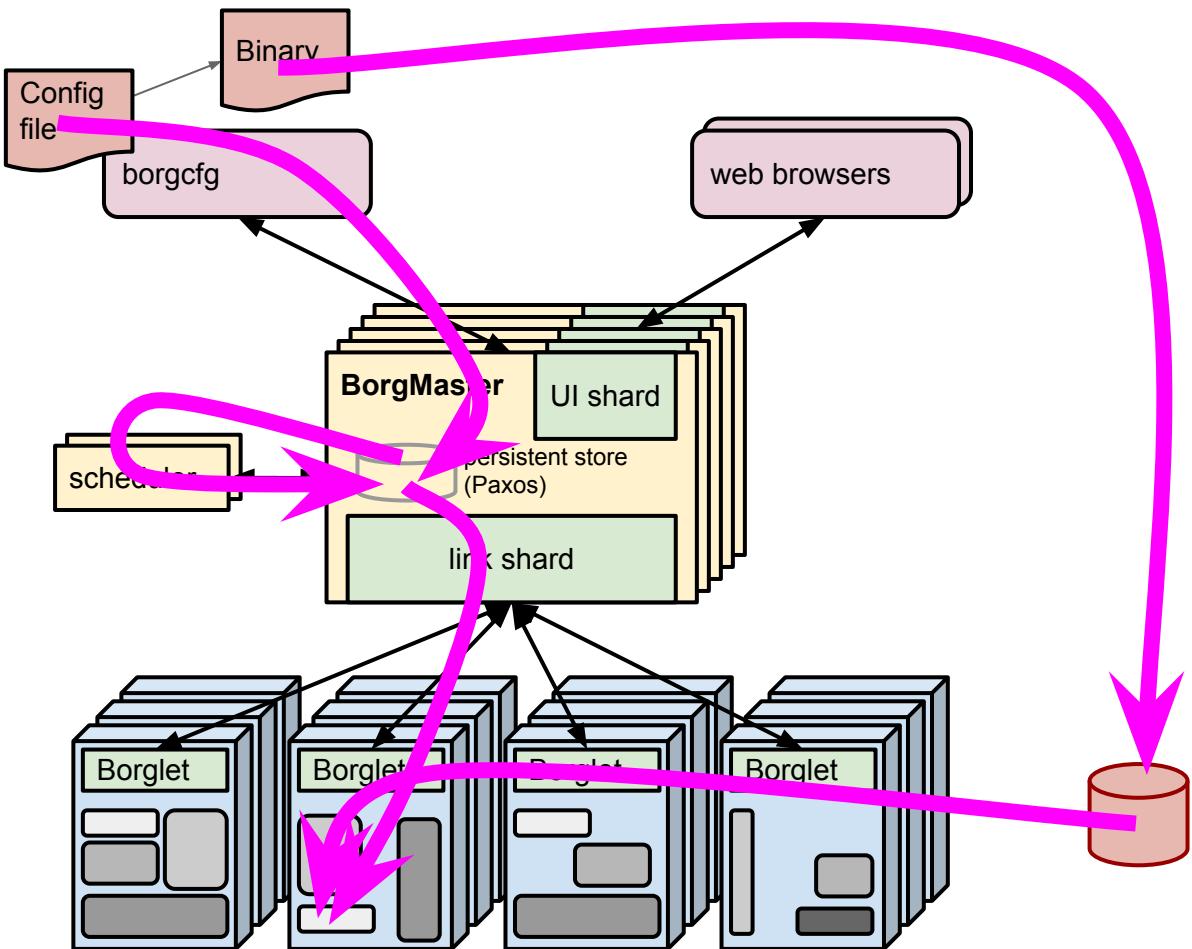




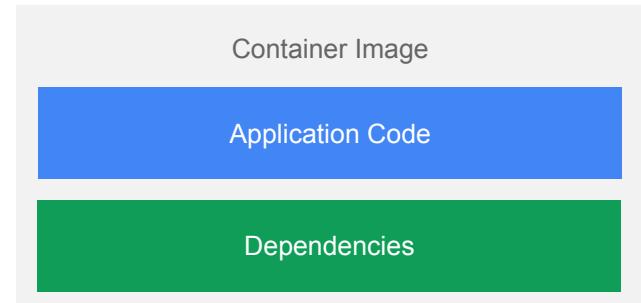
Image by Connie Zhou

What are Containers?

Containers encapsulate application code and all dependencies.

Applications can be deployed less on the infrastructure where it runs.

- In traditional IT environments, applications needed specific infrastructure. Dependencies needed to be installed beforehand.
- Containers incorporate applications and their dependencies so deployment to development, test, and production can be made easier.
- Don't need to be dependent on on-premise, private or public cloud environments.



Why Containers?



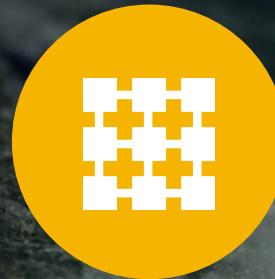
Fast

Simple and Fast compared to VMs. Can be started in just a few milliseconds.



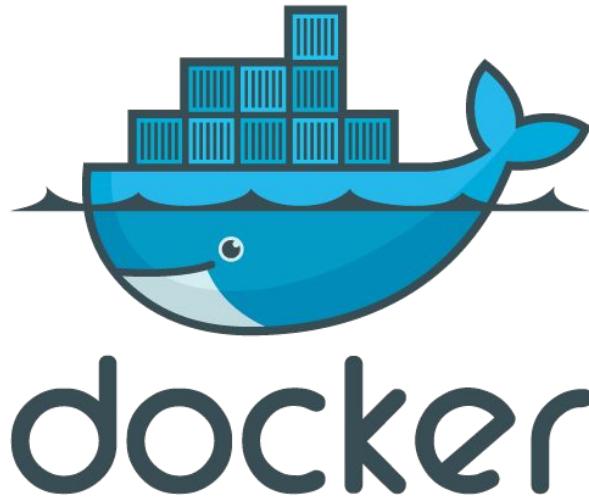
Portable

Can be run in a many environments.

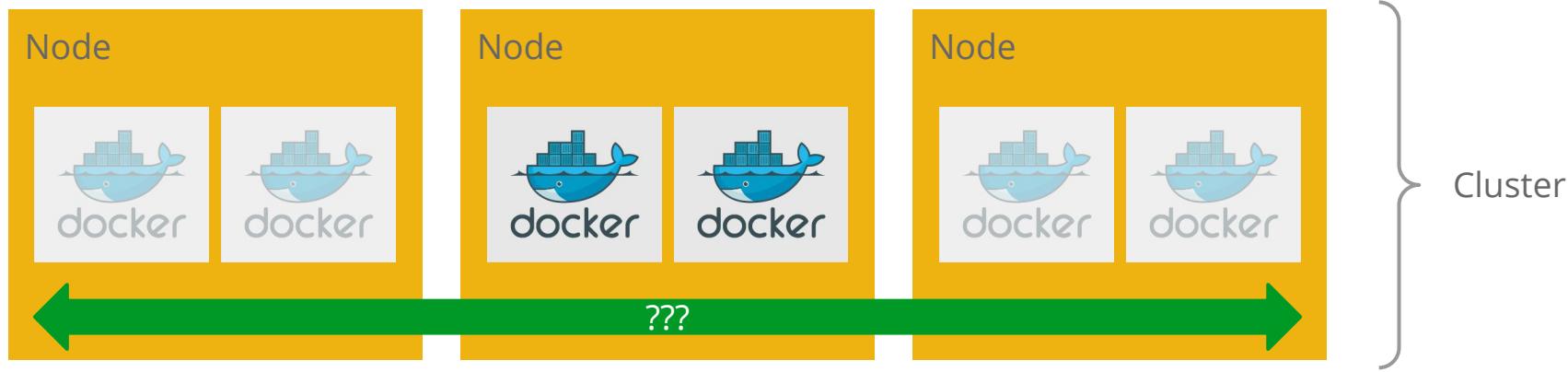


Efficiency

Low overhead. Resources used by containers can be limited.



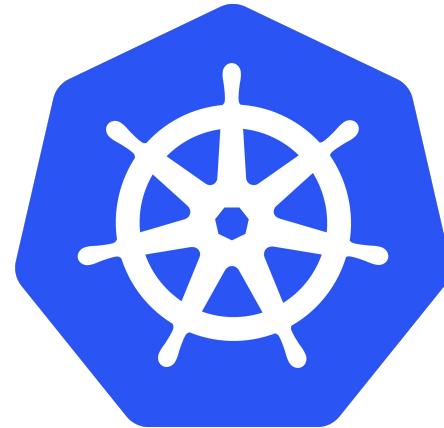
Container Management

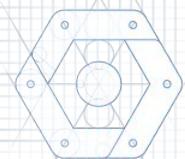


- How to deploy to multiple nodes?
- How to deal with node failures?
- How to deal with container failures?
- How do you update your applications?

Kubernetes

κυβερνήτης: *Greek for “pilot” or “helmsman of a ship”*
the open source cluster manager from Google





CNCF(Cloud Native Computing Foundation)



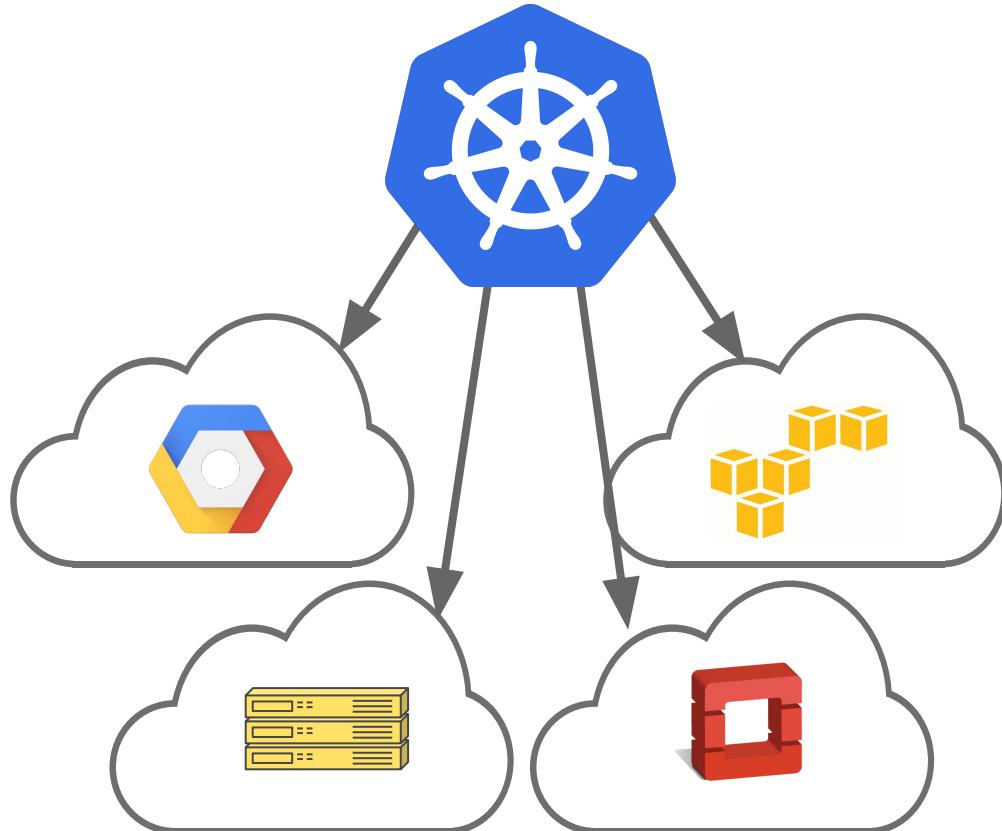
Workload portability

Goal: Avoid vendor lock-in

Runs in many environments, including
“bare metal” and “your laptop”

The API and the implementation are
100% open

The whole system is modular and
replaceable



Workload portability

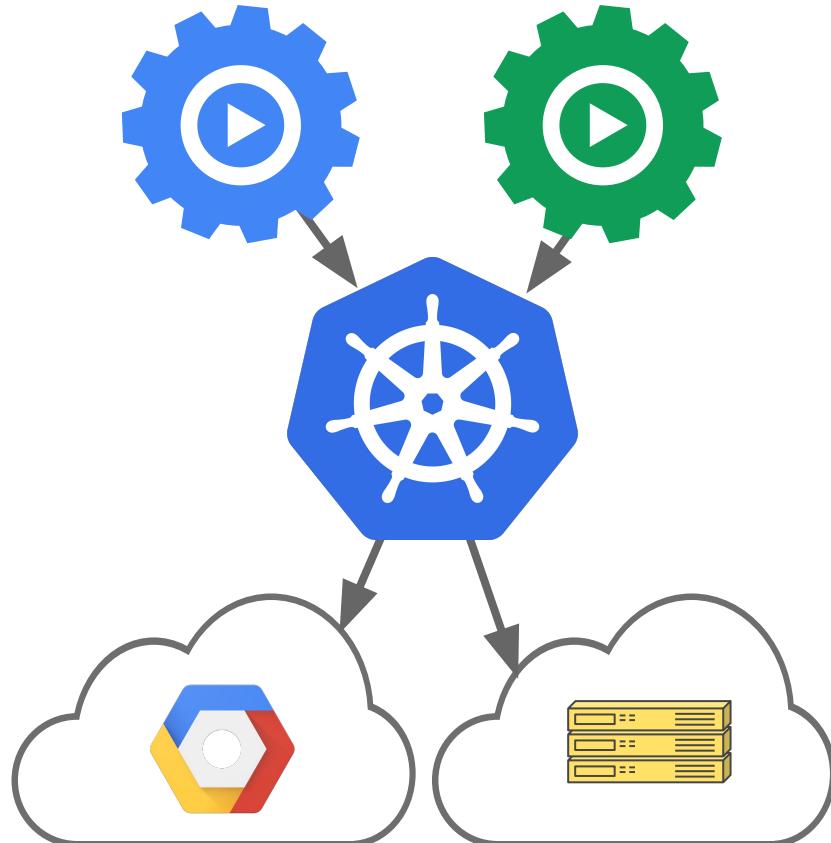
Goal: Write once, run anywhere^{*}

Don't force apps to know about
concepts that are
cloud-provider-specific

Examples of this:

- Network model
- Ingress
- Service load-balancers
- PersistentVolumes

** approximately*



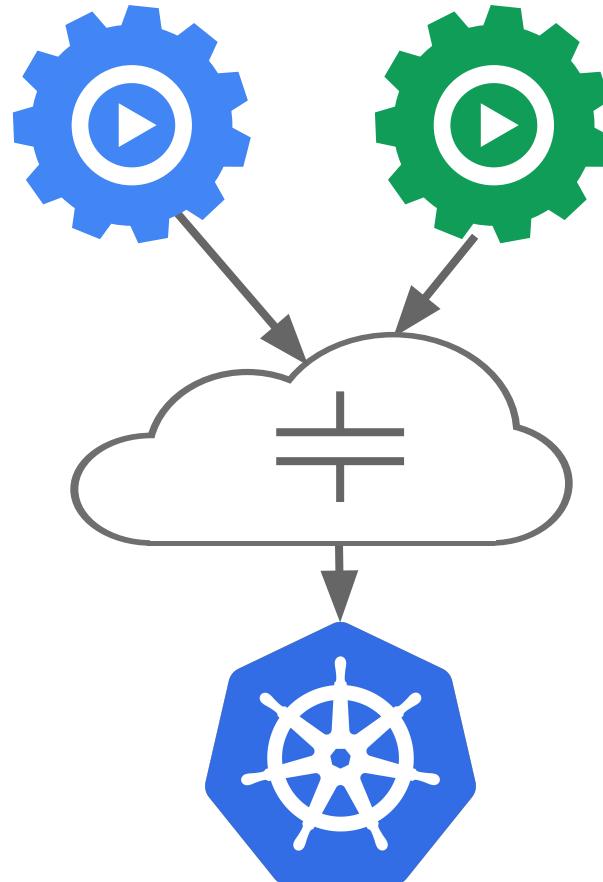
Workload portability

Goal: Avoid coupling

Don't force apps to know about concepts that are Kubernetes-specific

Examples of this:

- Namespaces
- Services / DNS
- Downward API
- Secrets / ConfigMaps



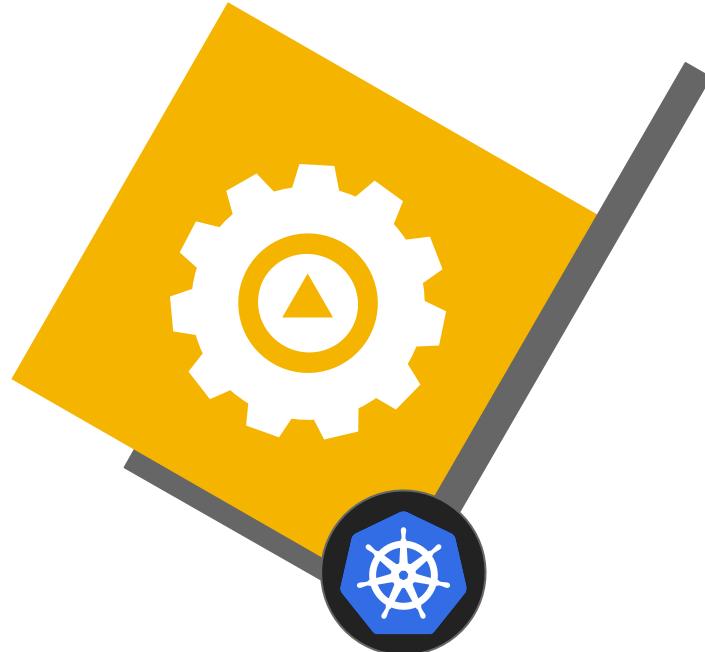
Workload portability

Result: Portability

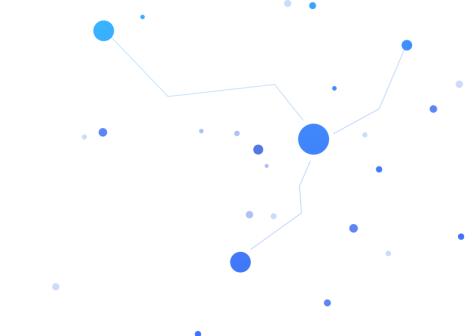
Build your apps on-prem, lift-and-shift into cloud when you are ready

Don't get stuck with a platform that doesn't work for you

Put your app on wheels and move it whenever and wherever you need



Container Engine



Pods

Small group of containers & volumes

Tightly coupled

The atom of scheduling & placement

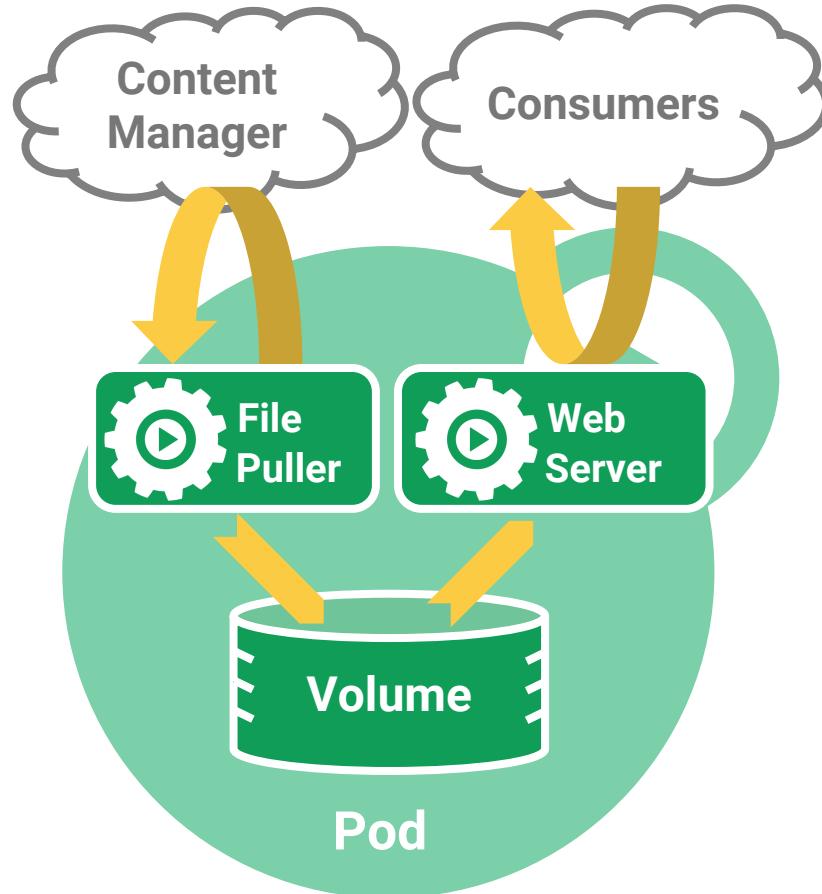
Shared namespace

- share IP address & localhost
- share IPC, etc.

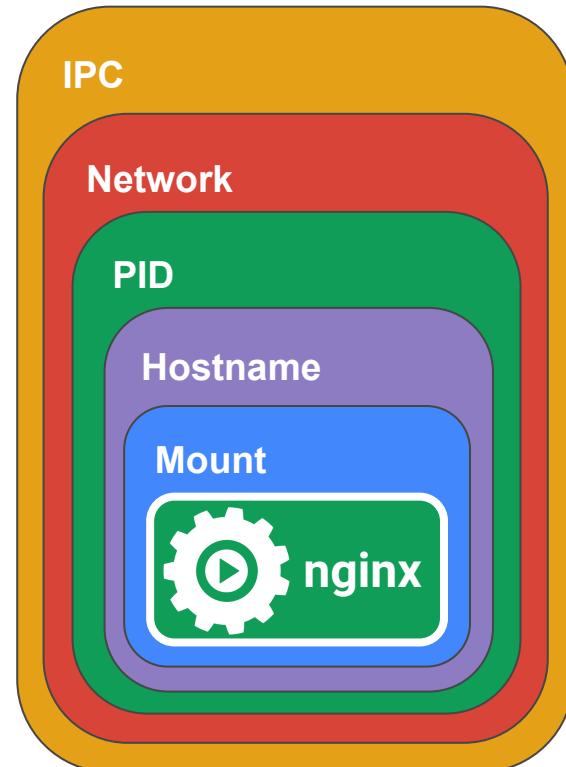
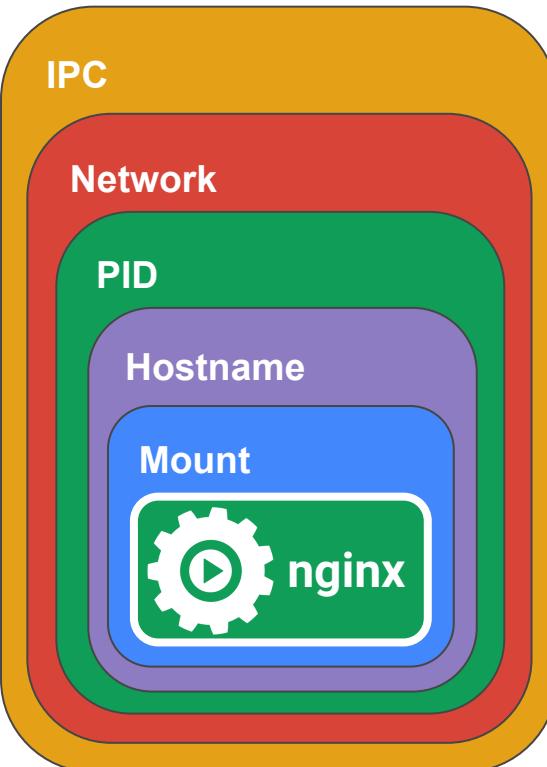
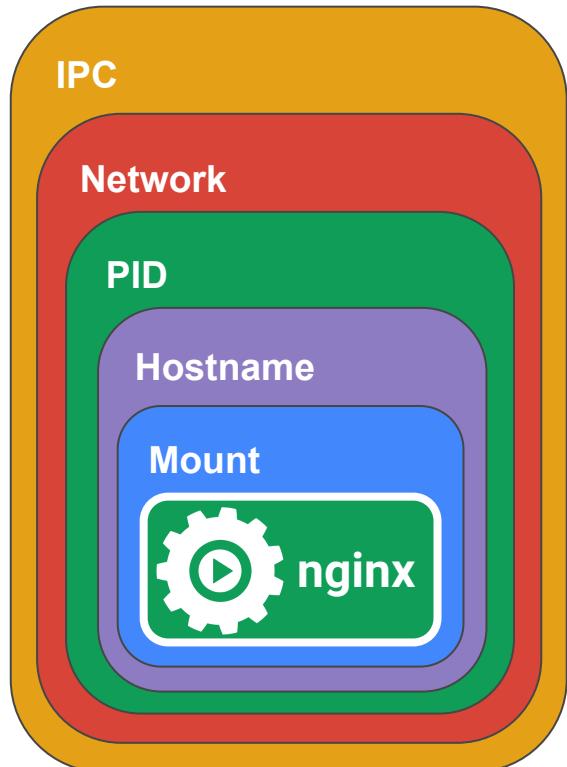
Managed lifecycle

- bound to a node, restart in place
- can die, cannot be reborn with same ID

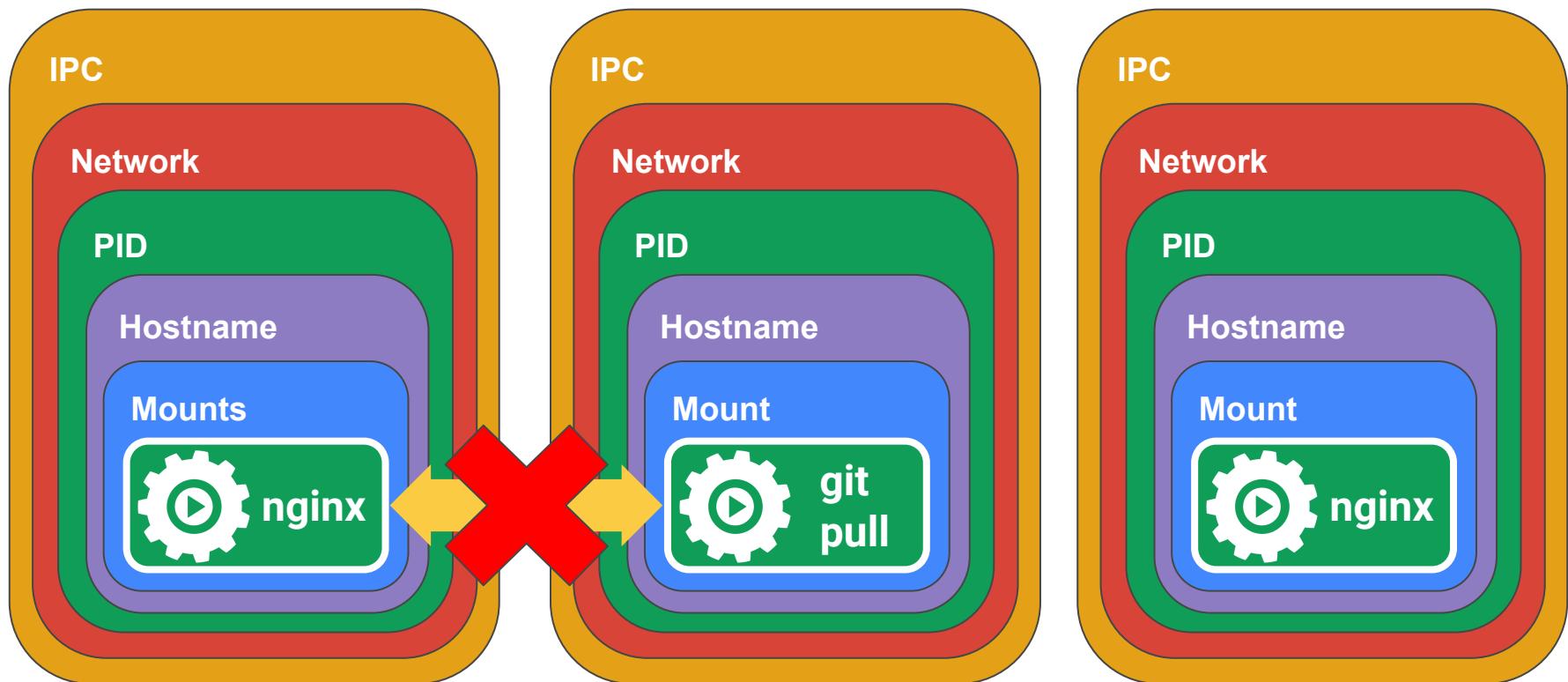
Example: data puller & web server



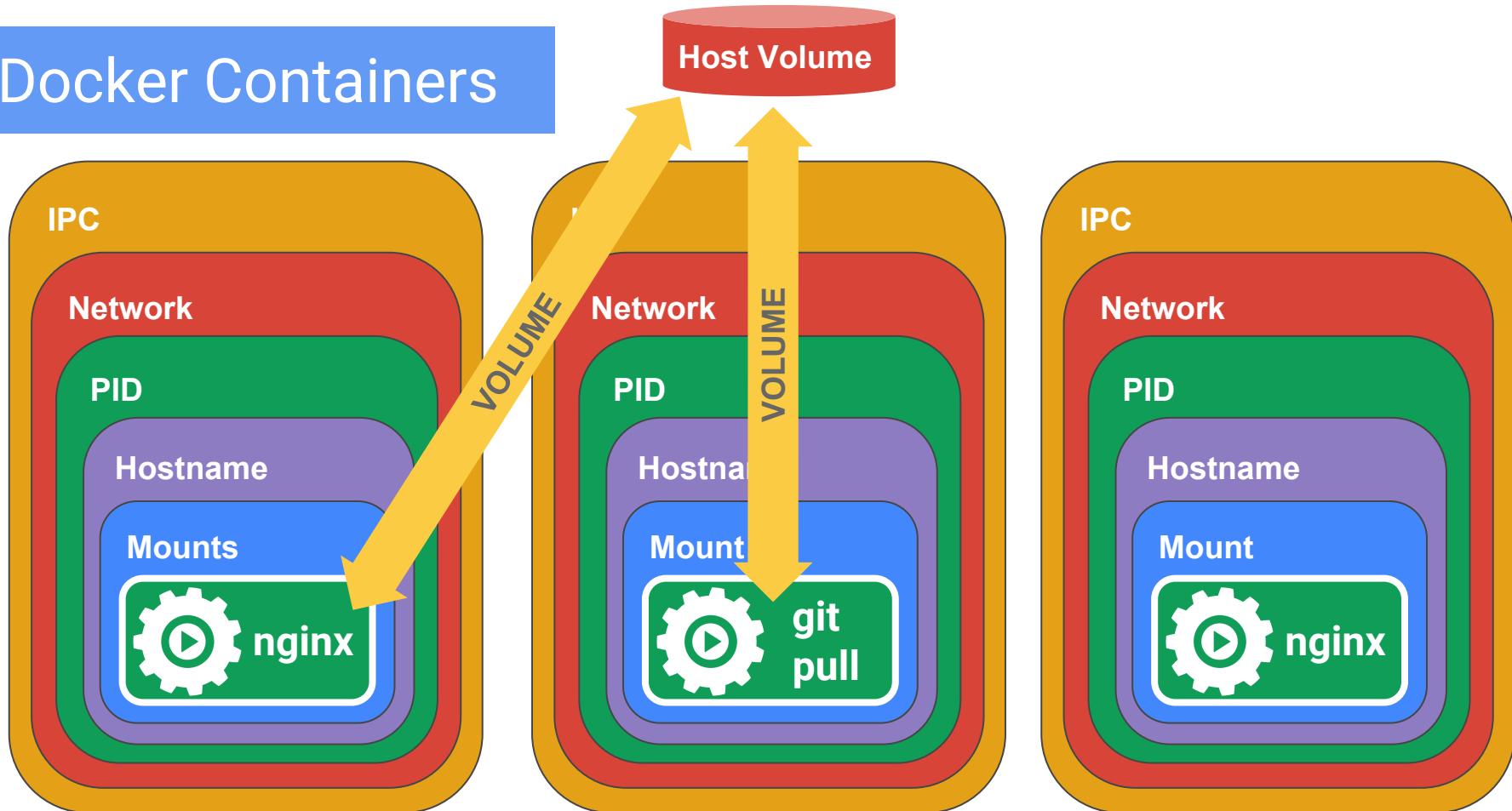
Docker Containers



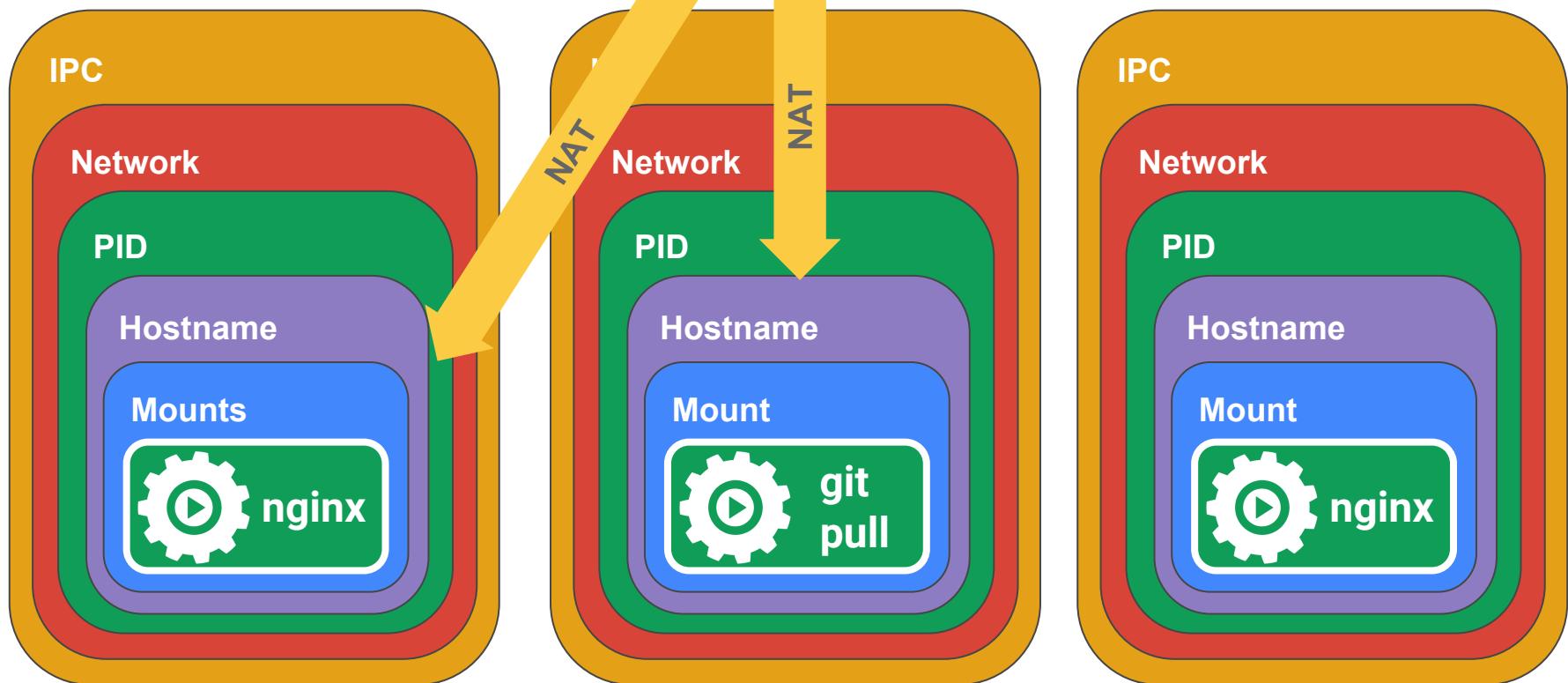
Docker Containers



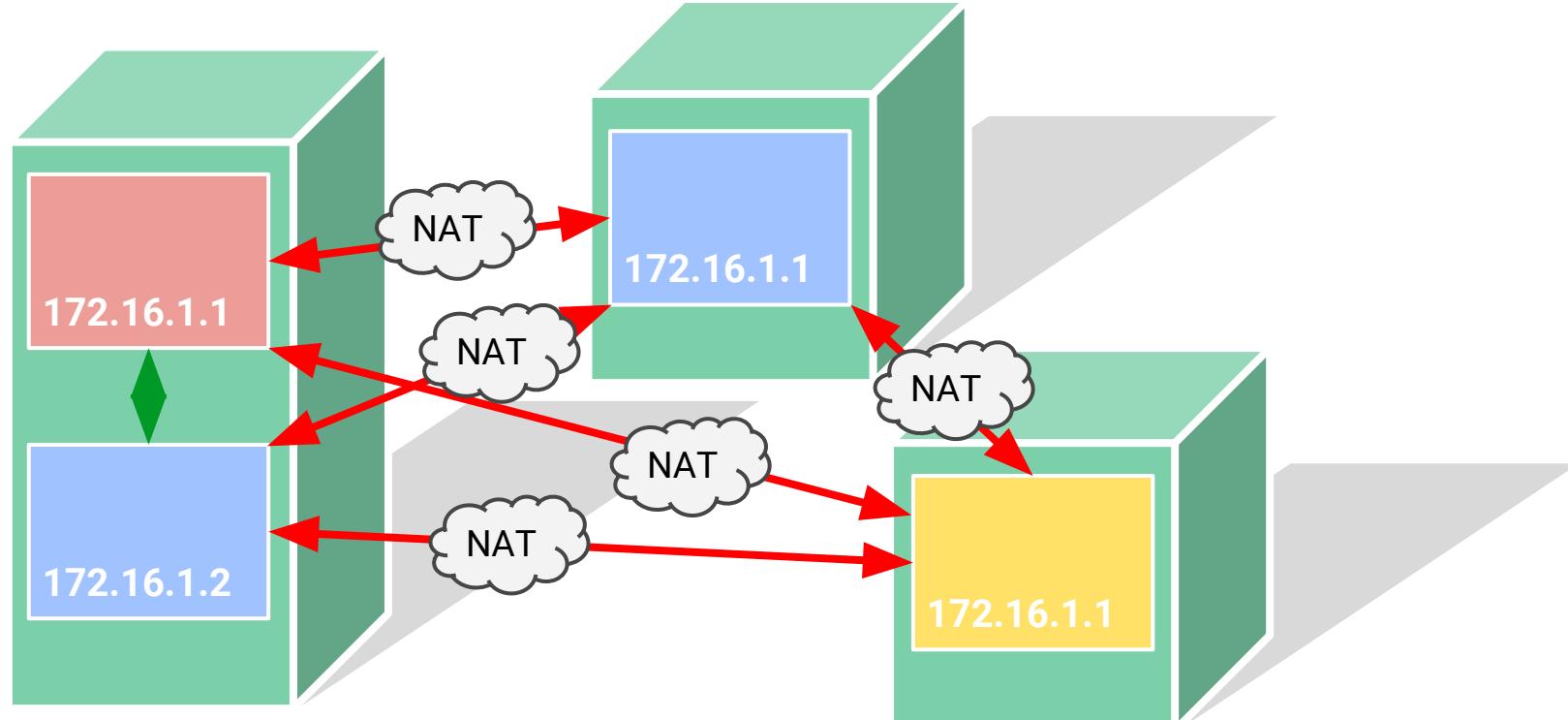
Docker Containers



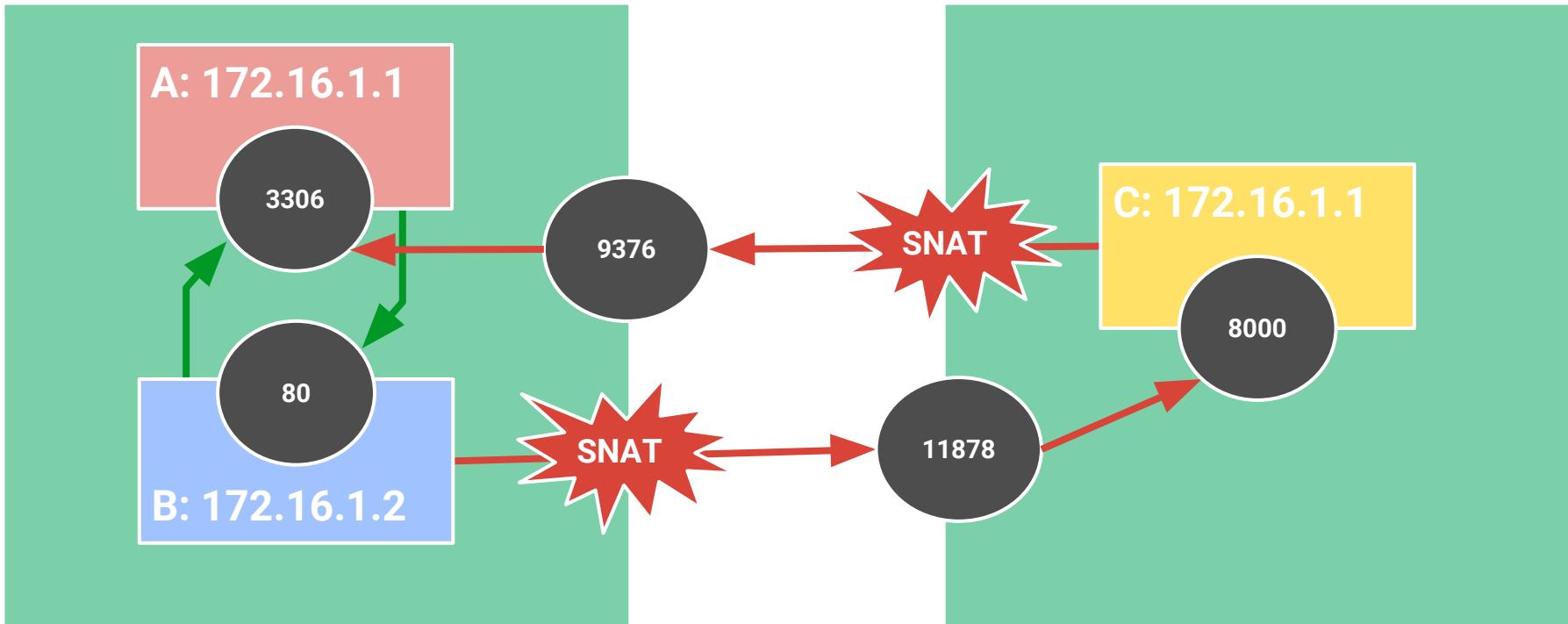
Docker Containers



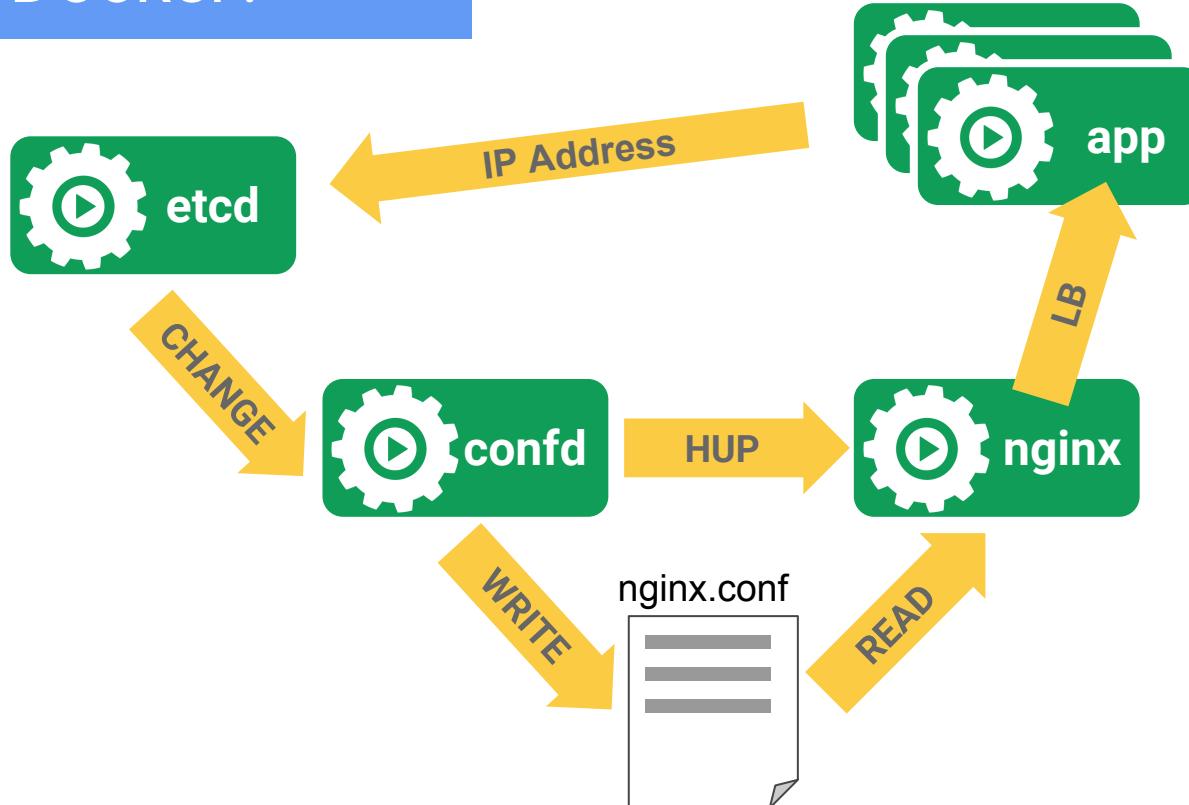
Docker networking



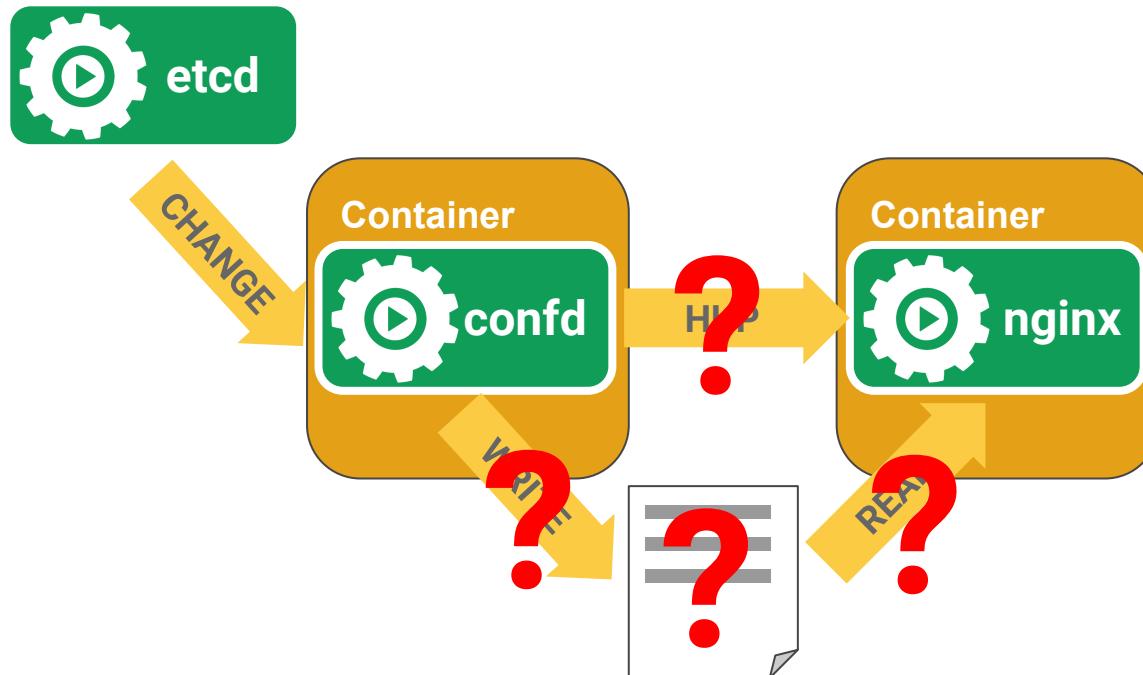
Port mapping



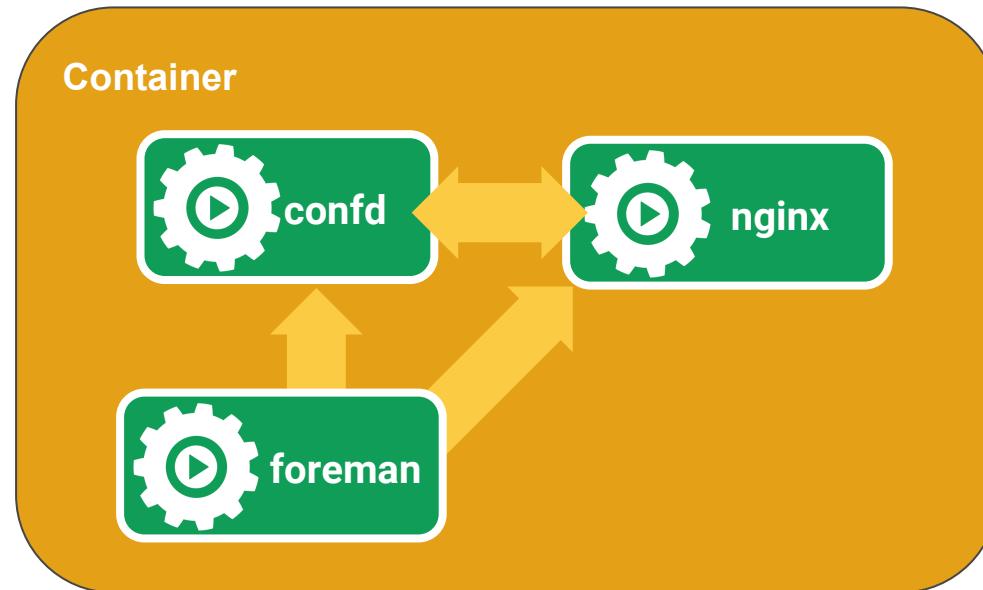
Pods & Docker?



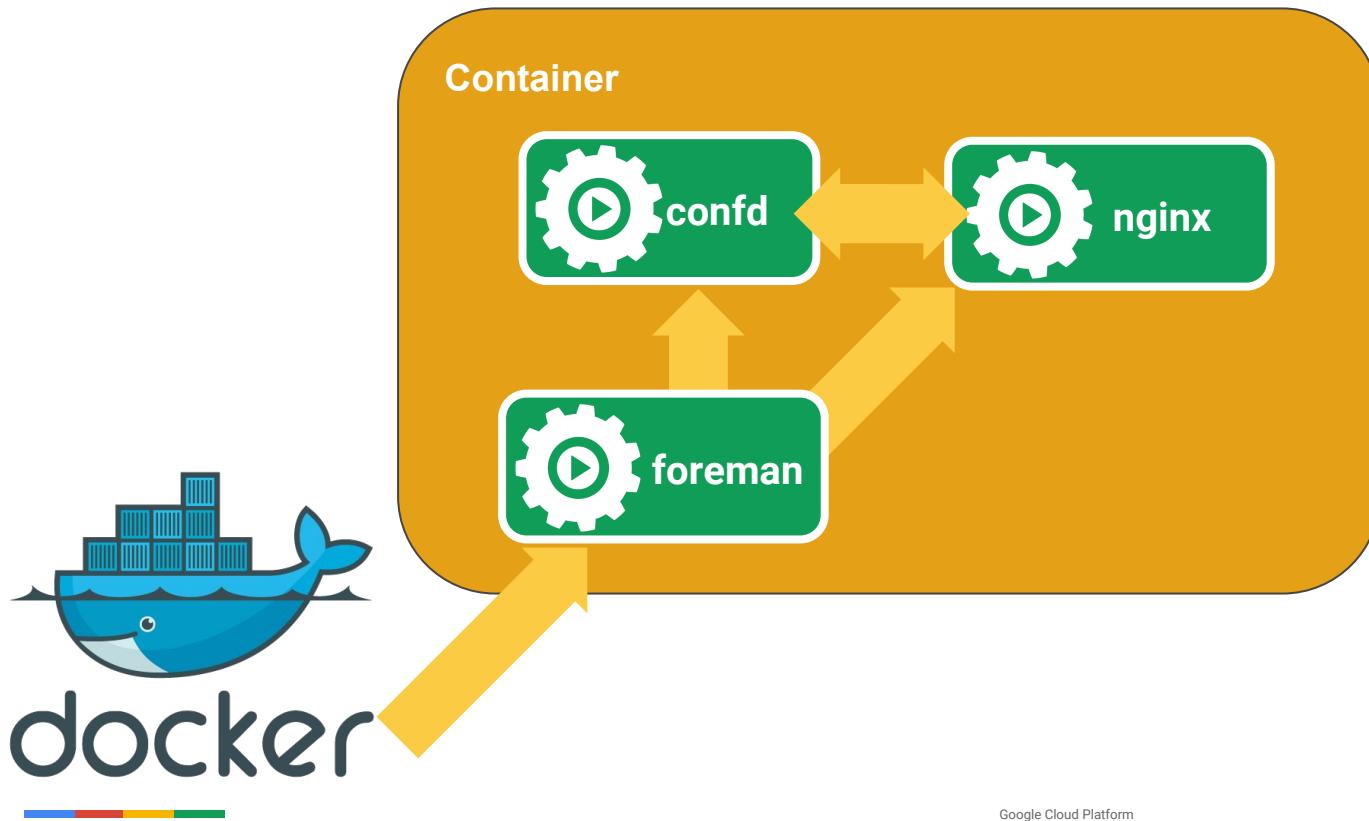
Pods & Docker?



Pods & Docker?



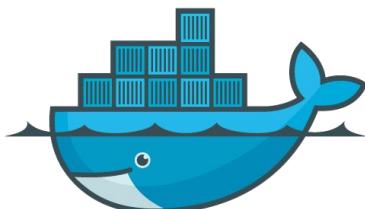
Pods & Docker?



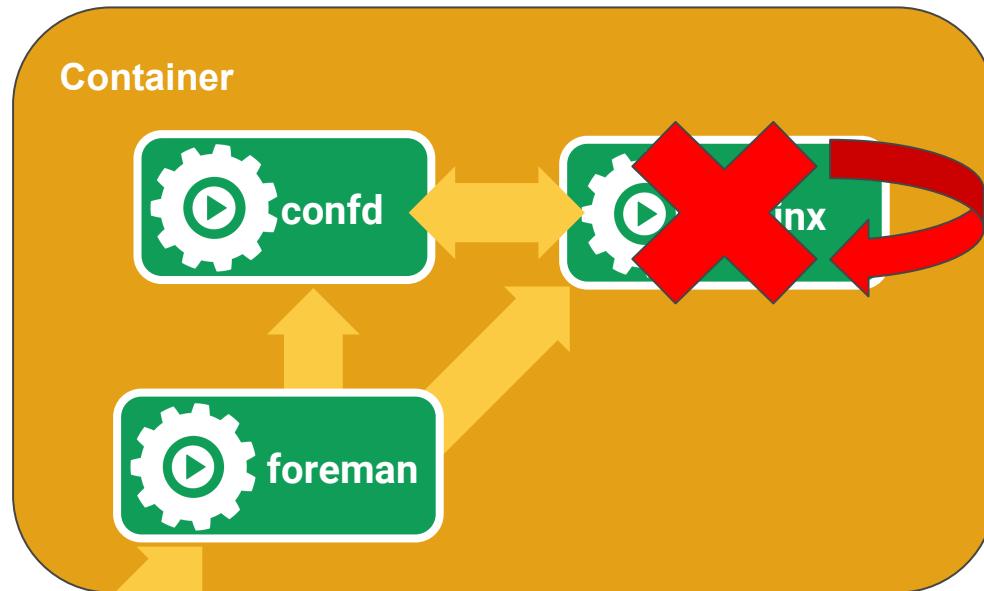
Pods & Docker?

Crash-Restart
Loop

Everything's
A-OK!!



docker

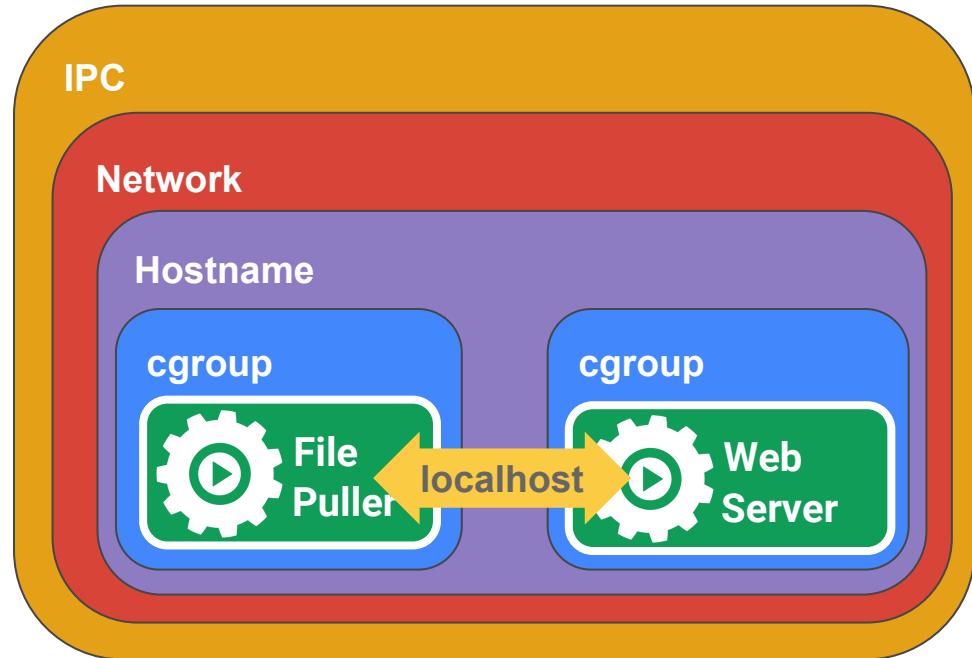


Pods

docker ...

--net=container:id

--ipc=container:id

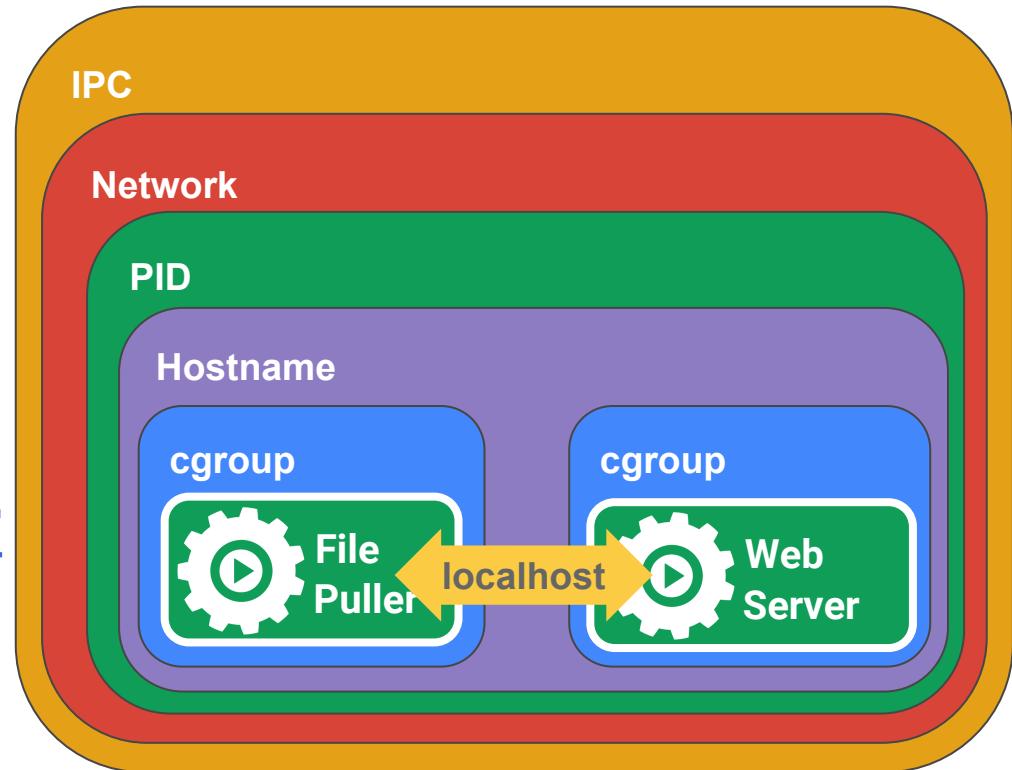


Pods (TODO)

docker ...

--net=container:id
--ipc=container:id
--pid=container:id

<https://github.com/docker/docker/issues/10163>



Kubernetes networking

IPs are **cluster-scoped**

- vs docker default private IP

Pods can reach each other directly

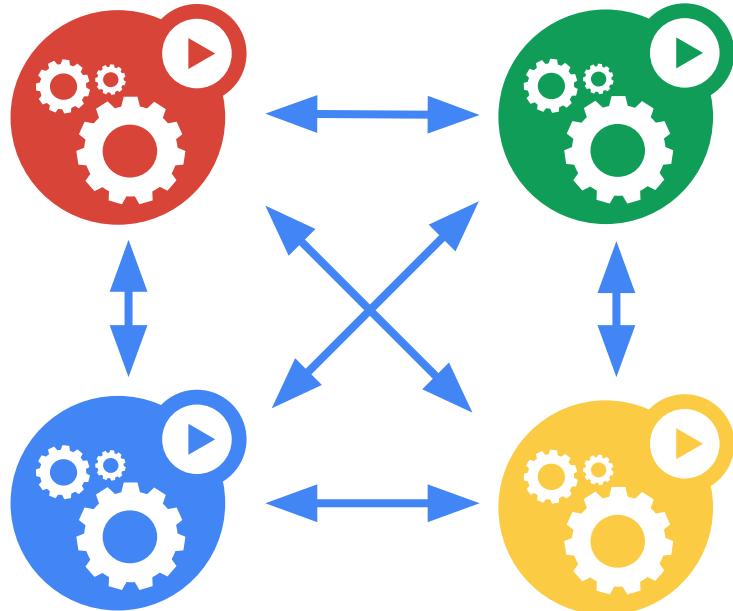
- even across nodes

No brokering of port numbers

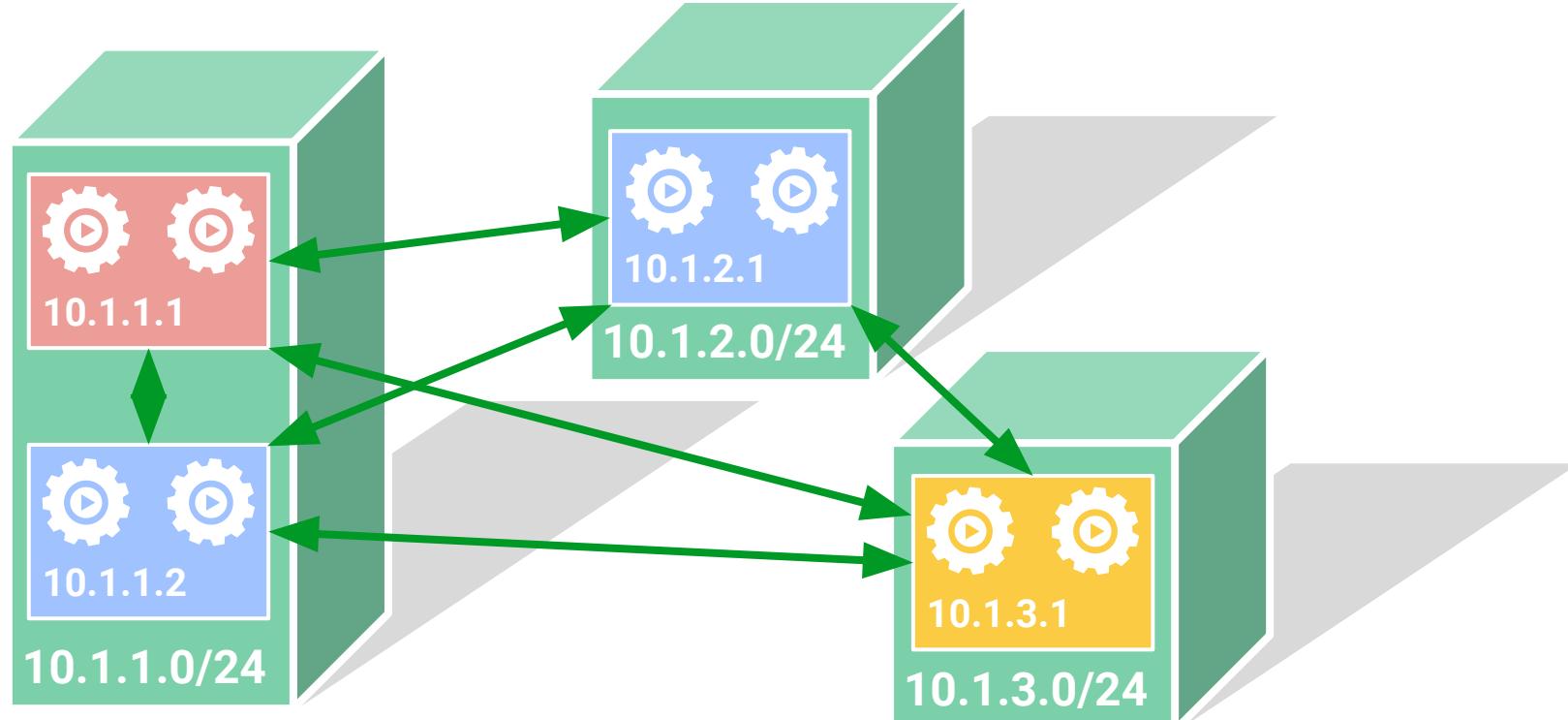
- too complex, why bother?

This is a fundamental requirement

- can be L3 routed
- can be underlaid (cloud)
- can be overlayed (SDN)



Kubernetes networking



ConfigMaps

Goal: manage app configuration

- ...without making overly-brITTLE container images

[12-factor](#) says config comes from the environment

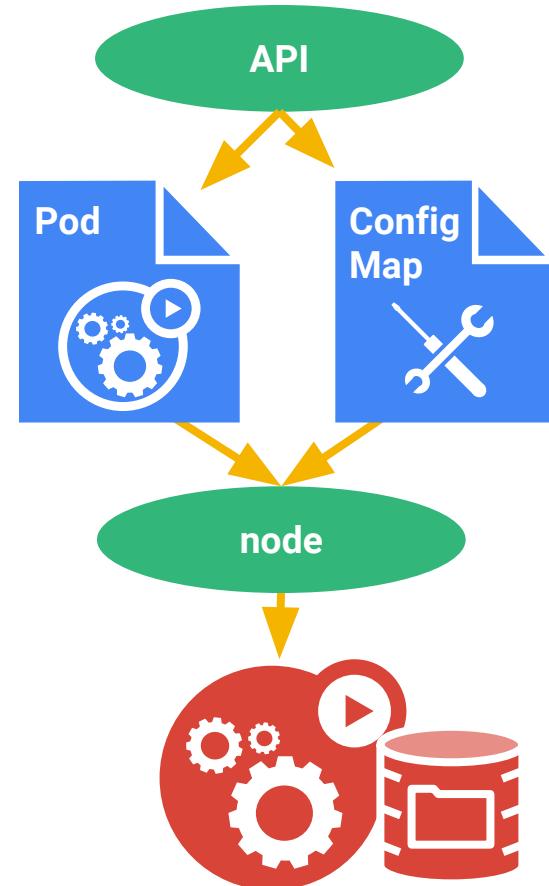
- Kubernetes is the environment

Manage config via the Kubernetes API

Inject config as a virtual volume into your Pods

- late-binding, live-updated (atomic)
- also available as env vars

Status: GA in Kubernetes v1.2



Secrets

Goal: grant a pod access to a secured *something*

- don't put secrets in the container image!

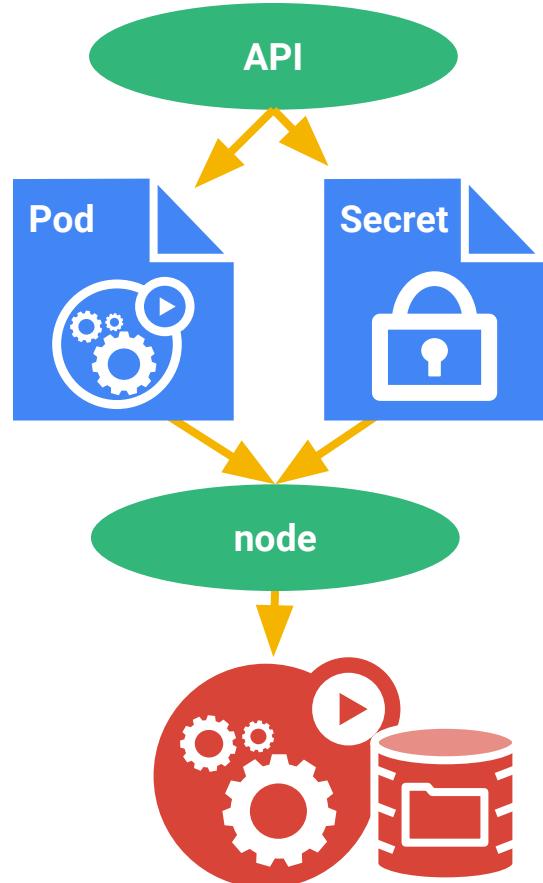
[12-factor](#) says config comes from the environment

- Kubernetes is the environment

Manage secrets via the Kubernetes API

Inject secrets as virtual volumes into your Pods

- late-binding, tmpfs - never touches disk
- also available as env vars



PersistentVolumes

A higher-level storage abstraction

- insulation from any one cloud environment

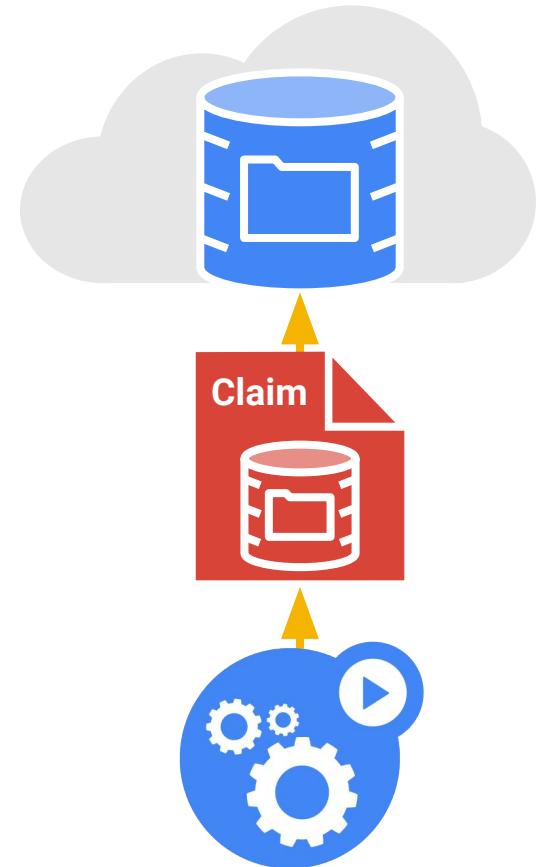
Admin provisions them, users claim them

- **NEW: auto-provisioning (alpha in v1.2)**

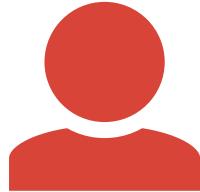
Independent lifetime from consumers

- lives until user is done with it
- can be handed-off between pods

Dynamically “scheduled” and managed, like nodes and pods



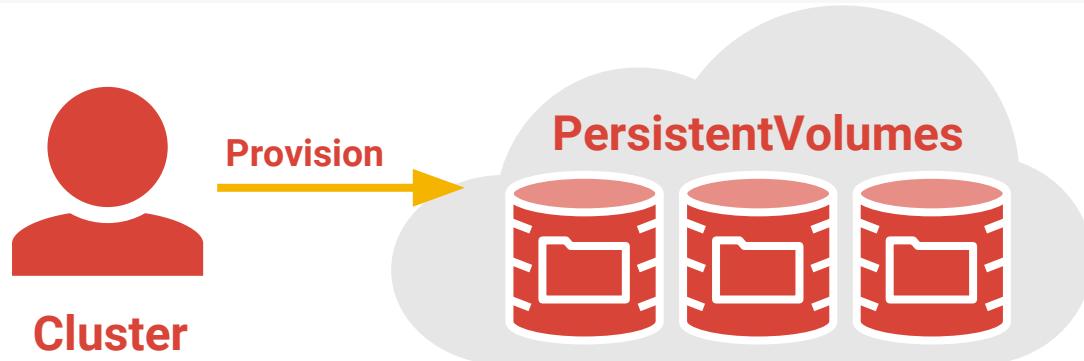
PersistentVolumes



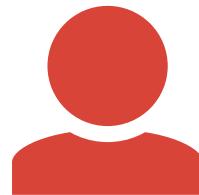
**Cluster
Admin**



PersistentVolumes



PersistentVolumes

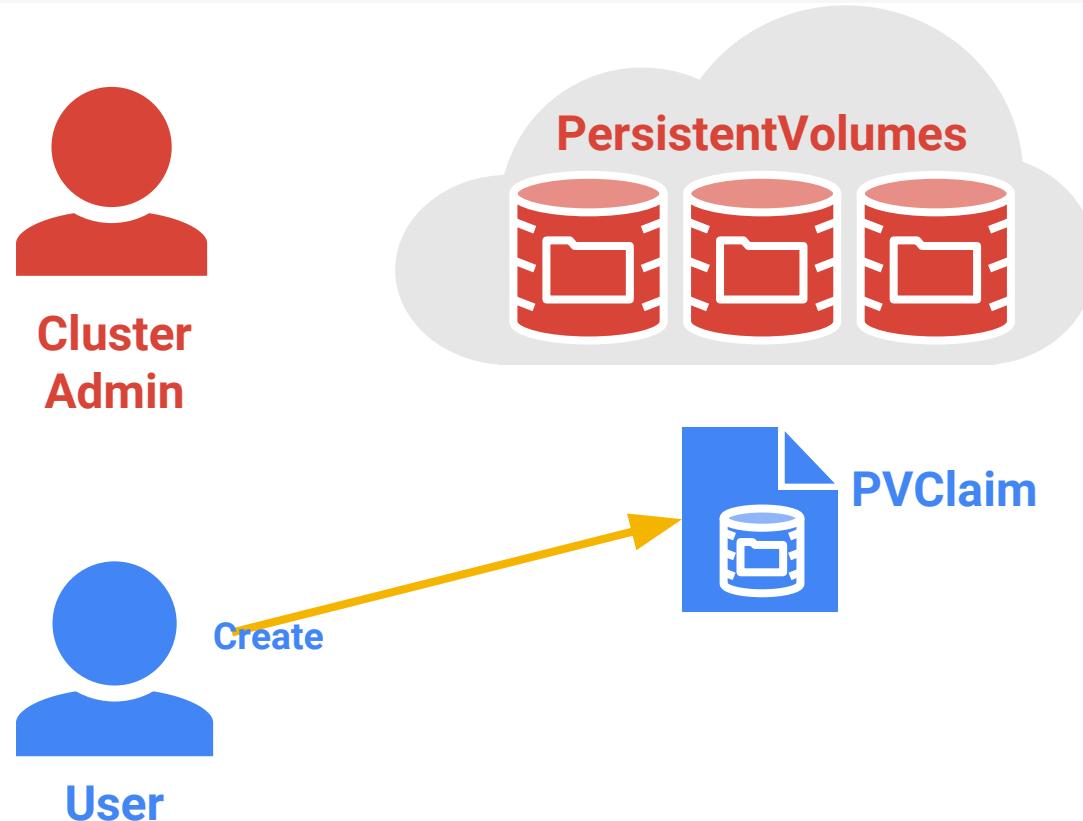


**Cluster
Admin**

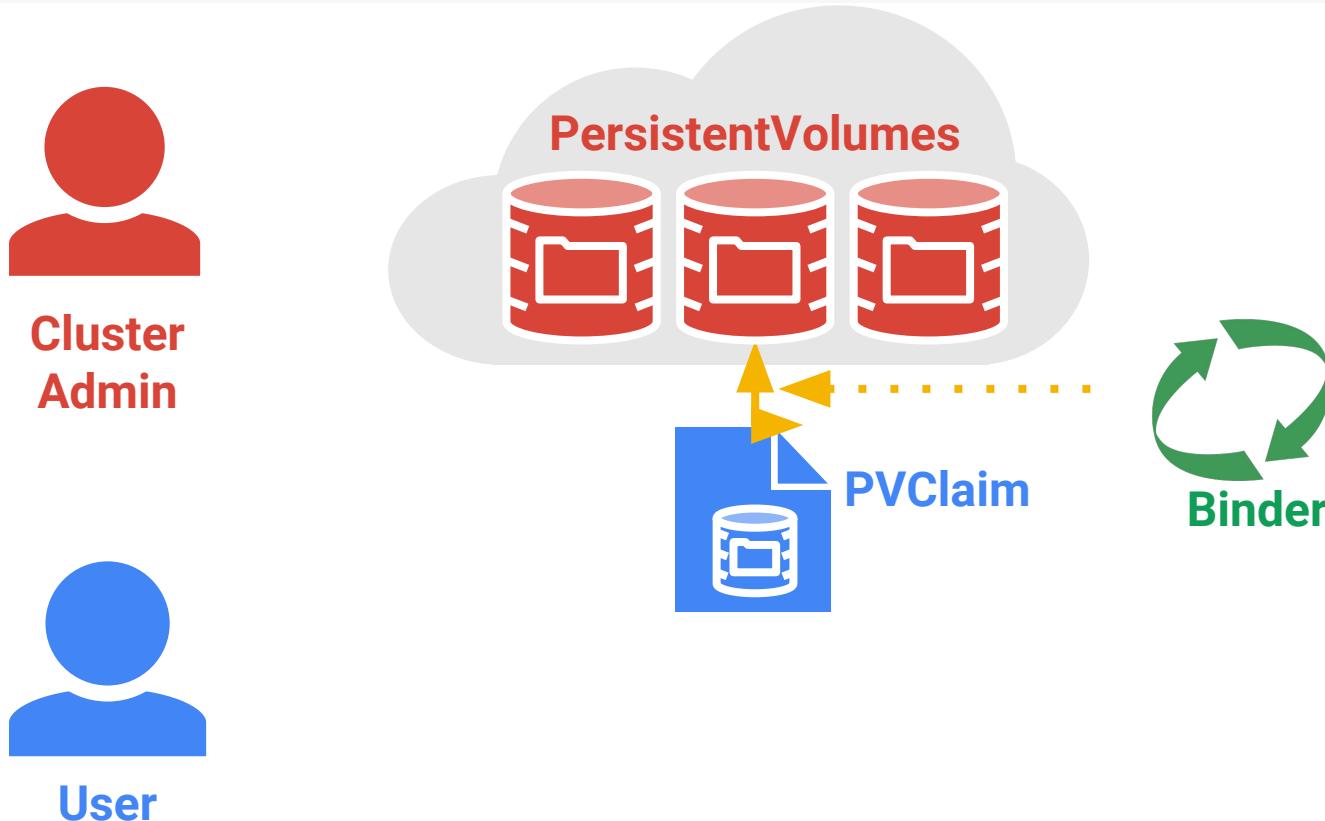


User

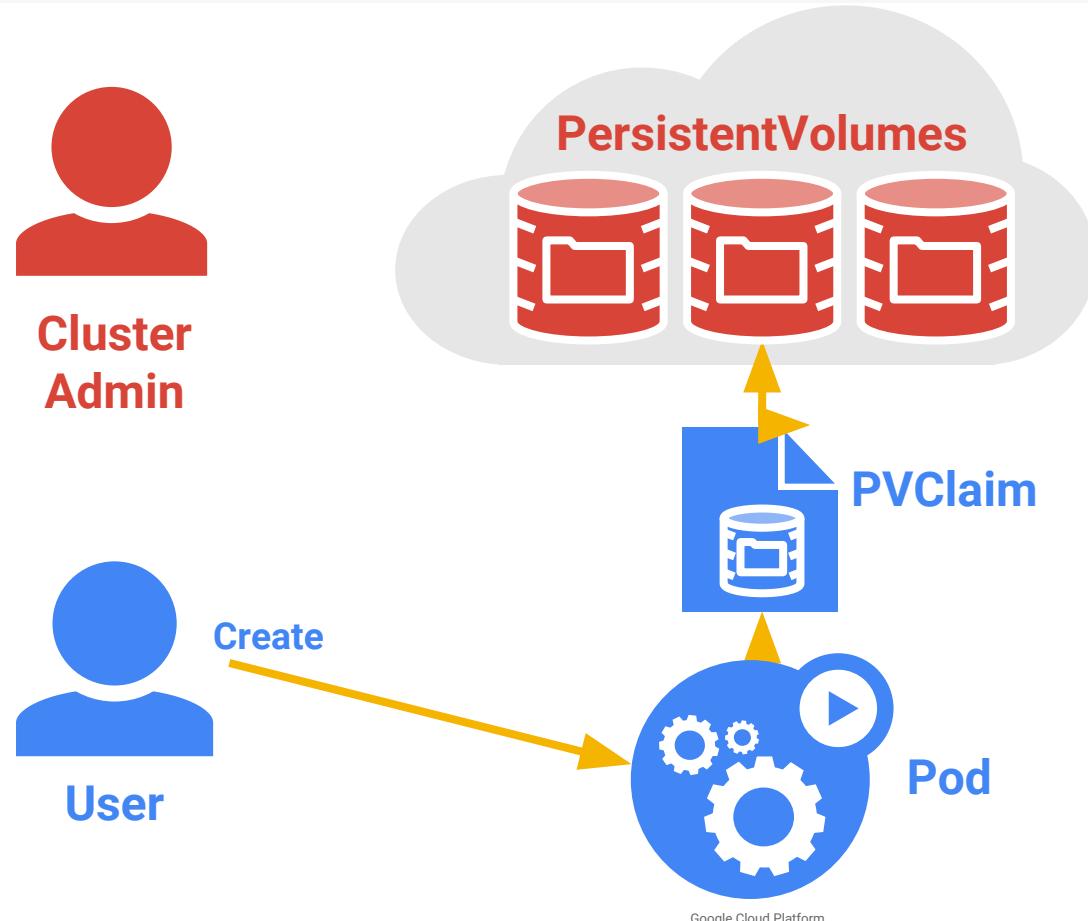
PersistentVolumes



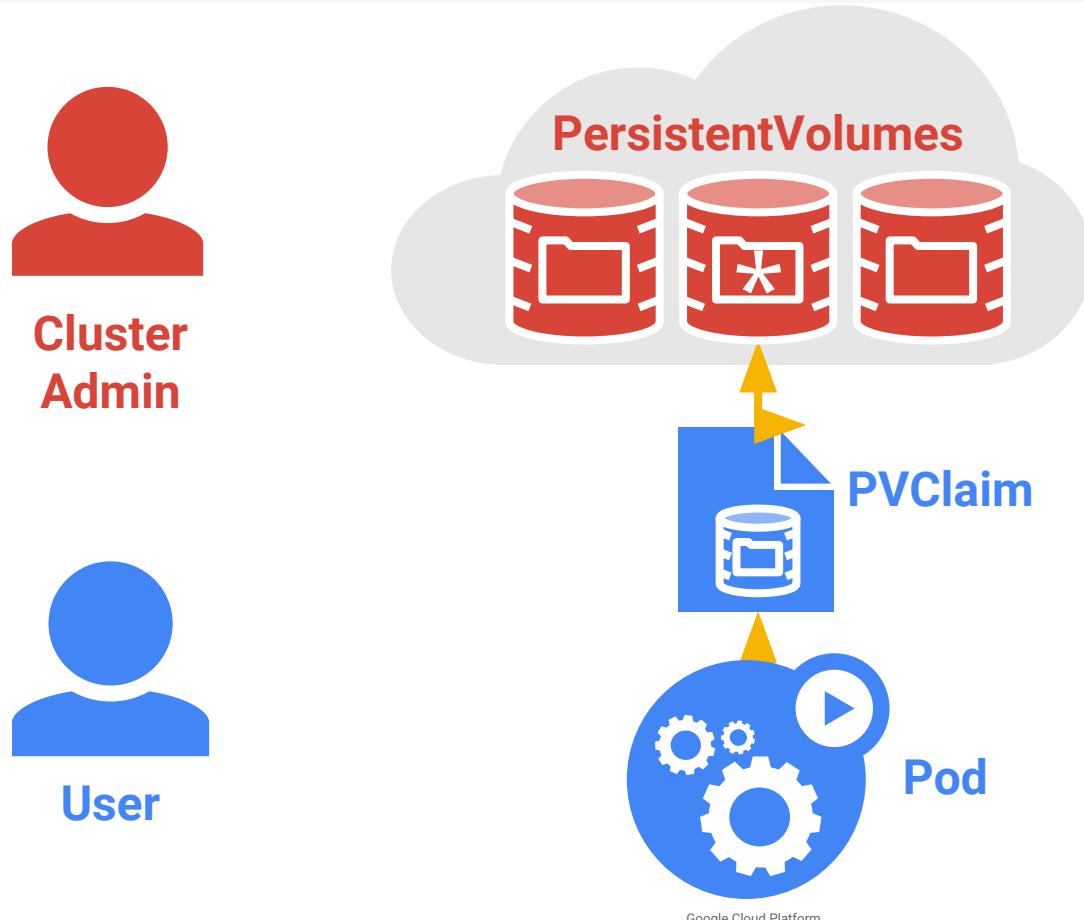
PersistentVolumes



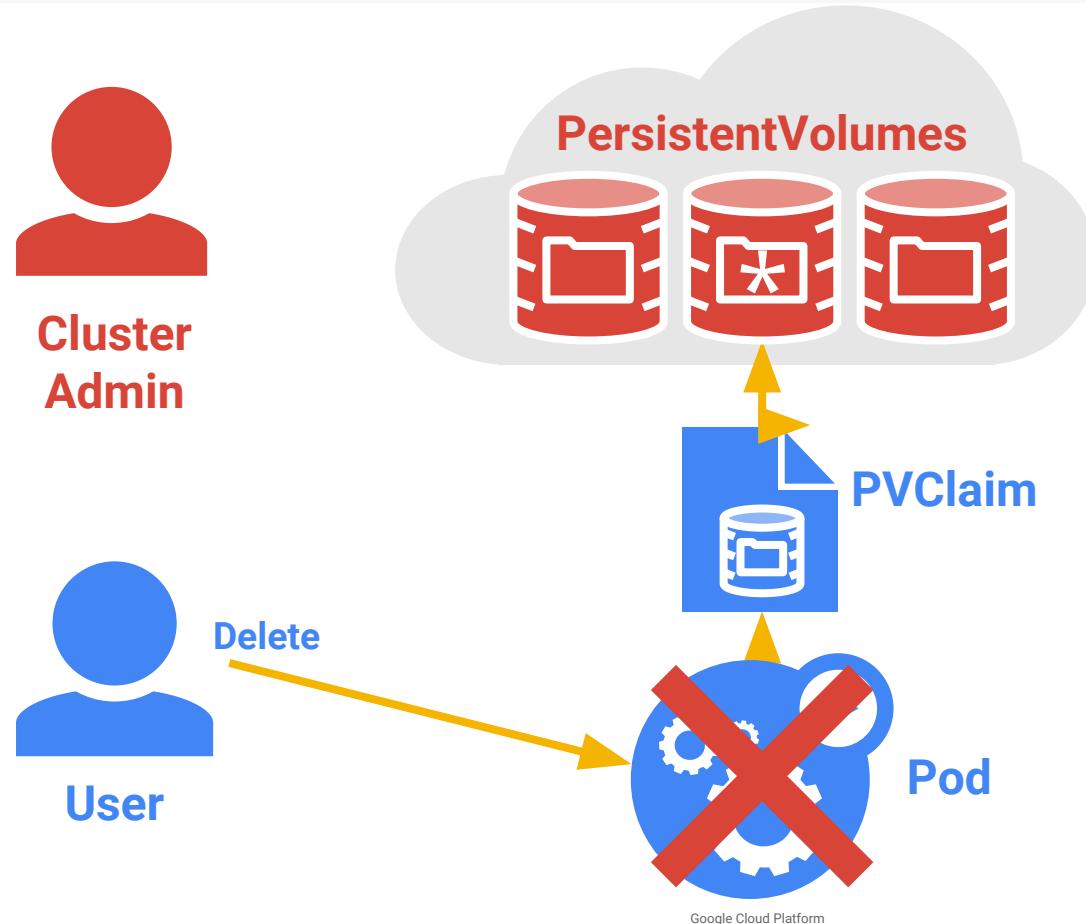
PersistentVolumes



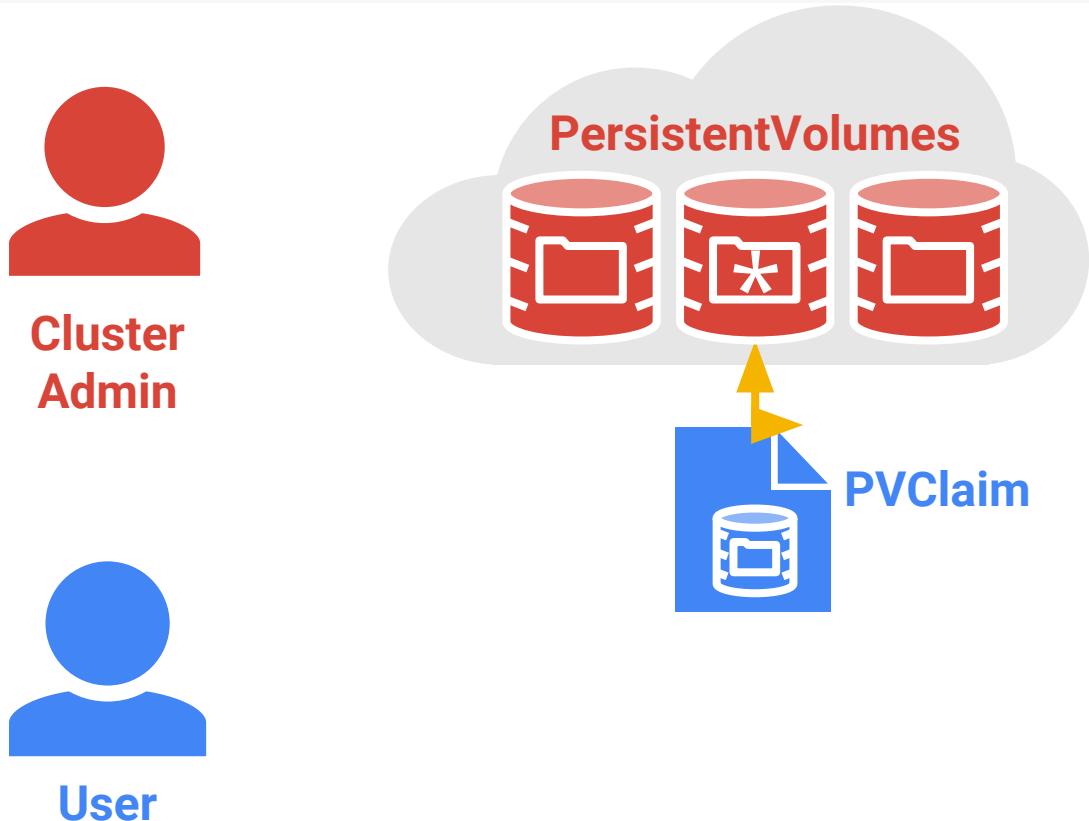
PersistentVolumes



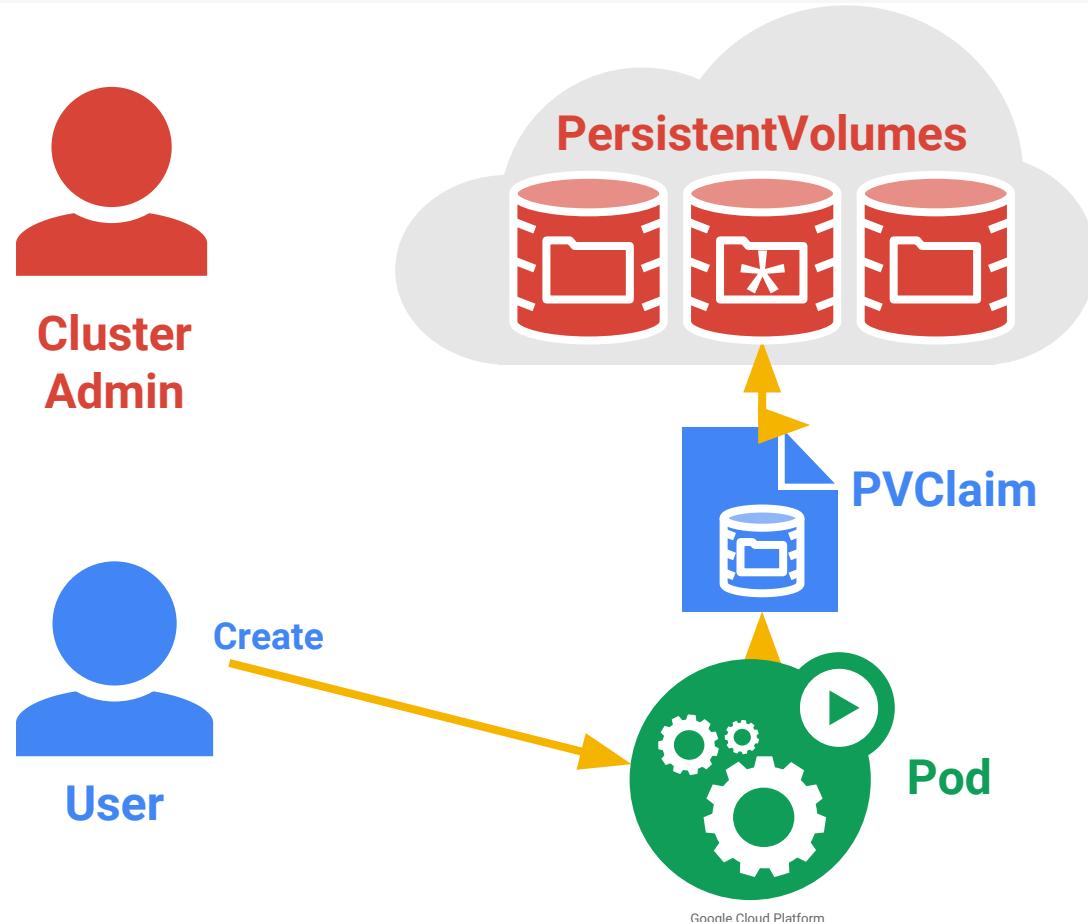
PersistentVolumes



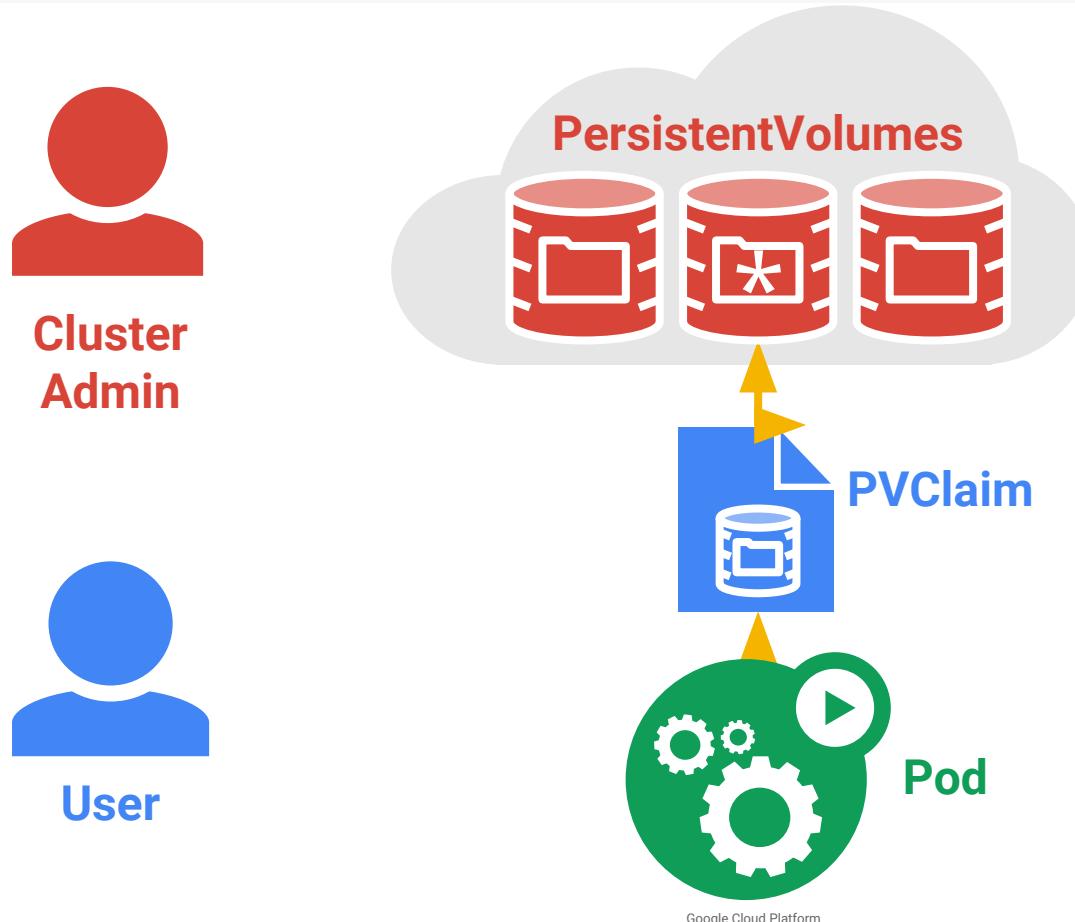
PersistentVolumes



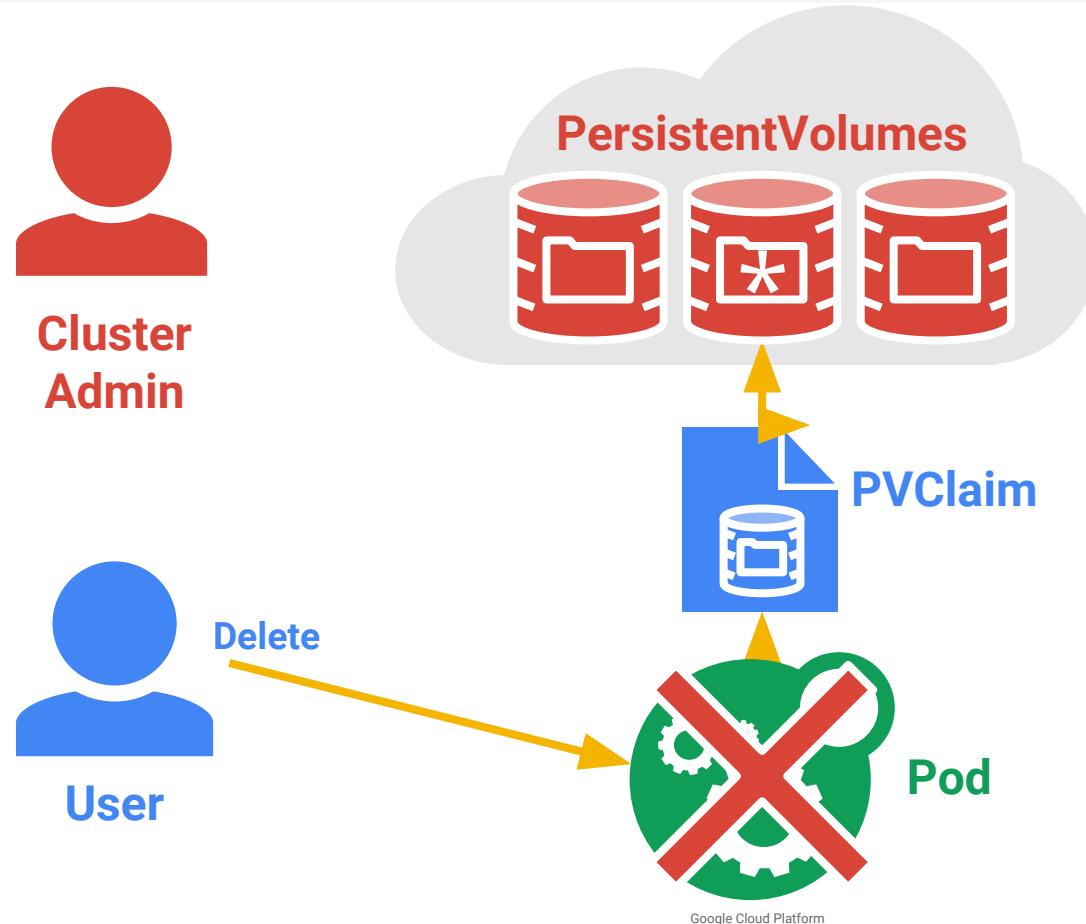
PersistentVolumes



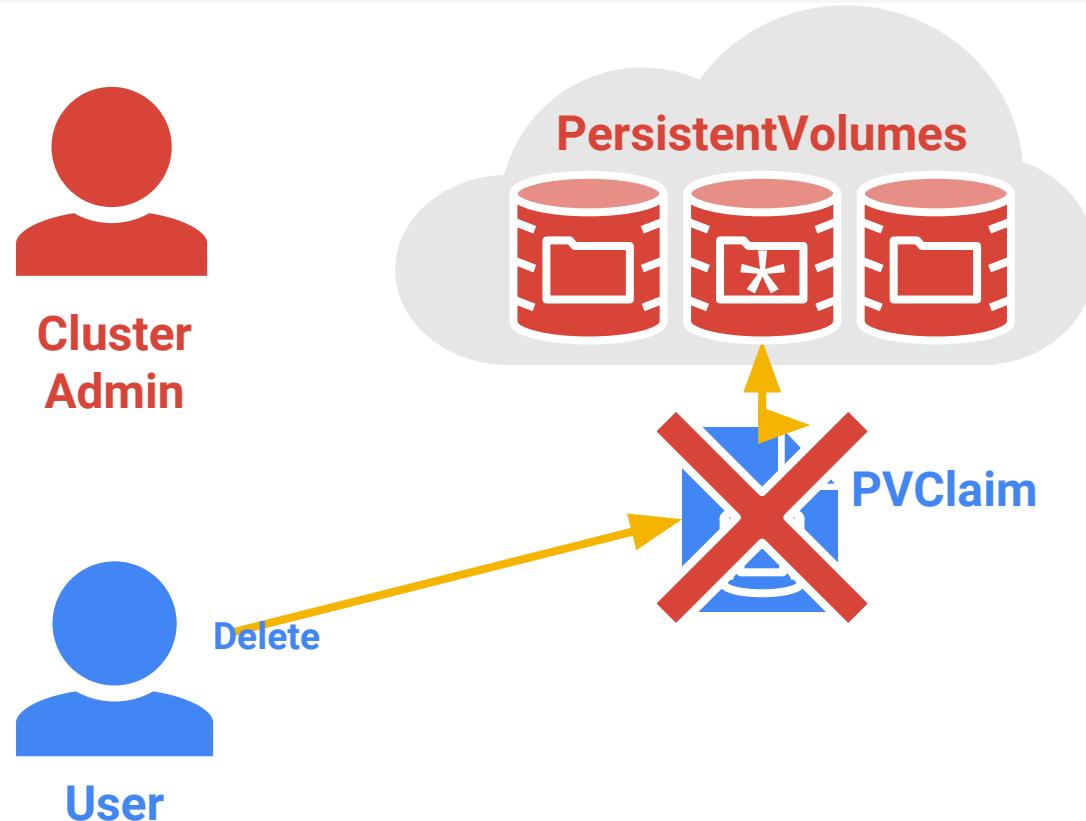
PersistentVolumes



PersistentVolumes



PersistentVolumes

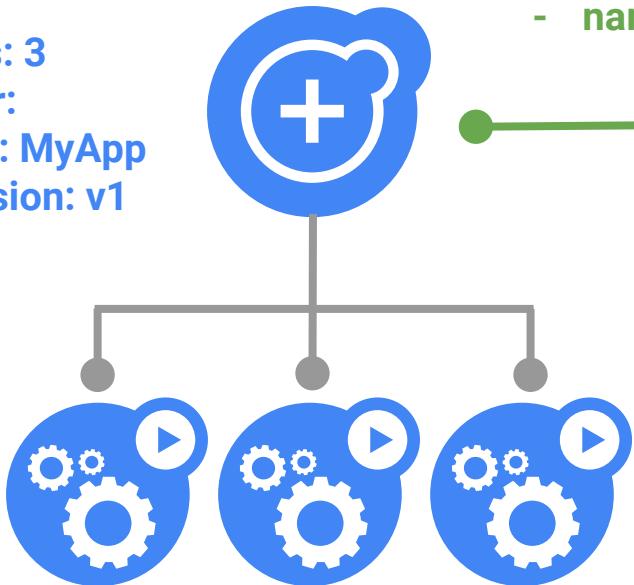


PersistentVolumes

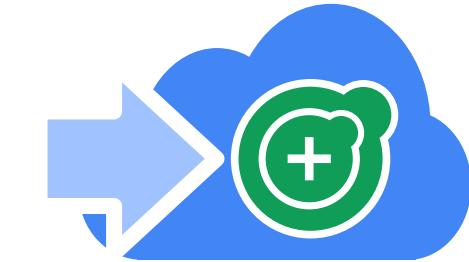


Deployments

ReplicaSet
- replicas: 3
- selector:
 - app: MyApp
 - version: v1



Deployment
- name: MyApp



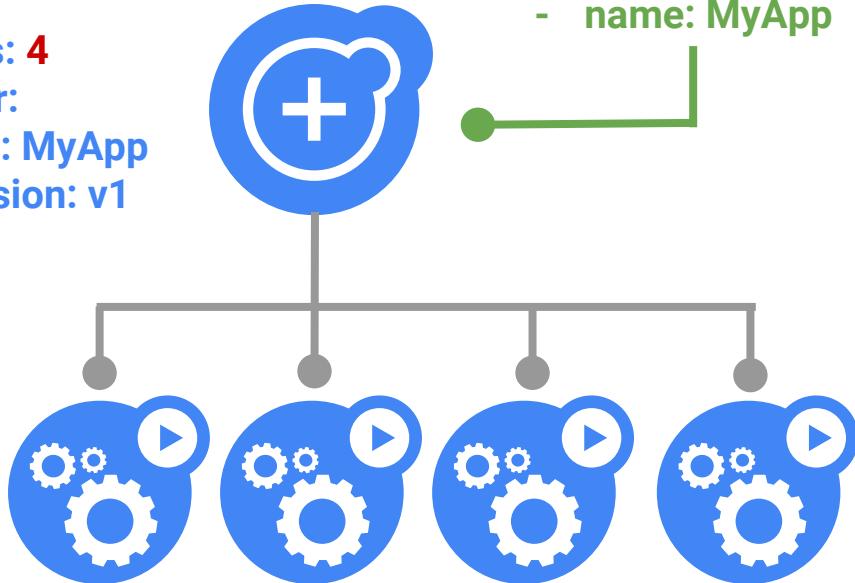
kubectl create ...



Deployments

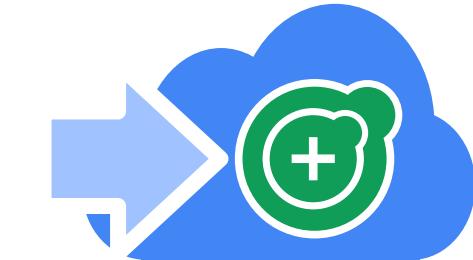
ReplicaSet

- replicas: 4
- selector:
 - app: MyApp
 - version: v1



Deployment

- name: MyApp



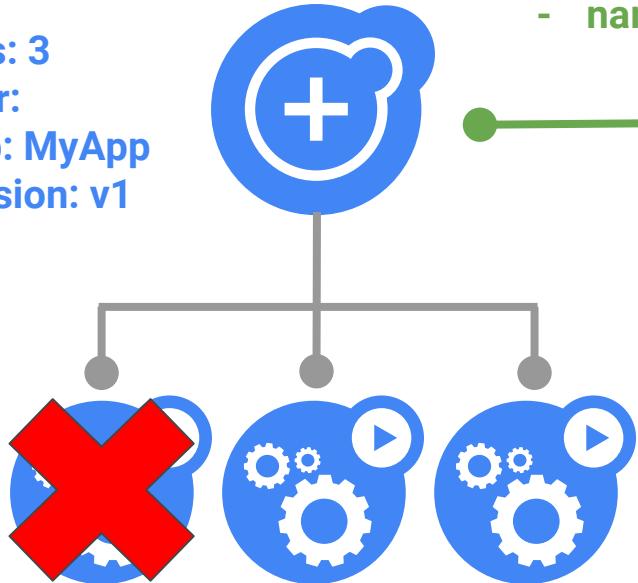
kubectl create ...



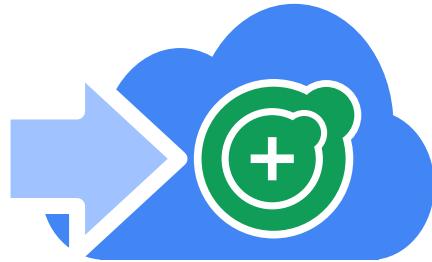
Deployments

ReplicaSet

- replicas: 3
- selector:
 - app: MyApp
 - version: v1



Deployment
- name: MyApp



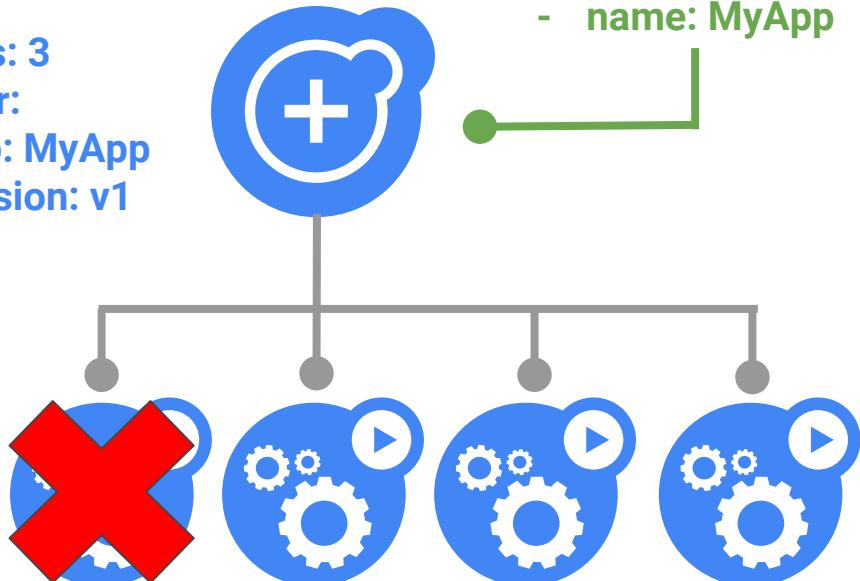
kubectl create ...



Deployments

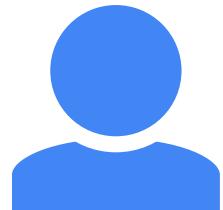
ReplicaSet

- replicas: 3
- selector:
 - app: MyApp
 - version: v1



Deployment
- name: MyApp

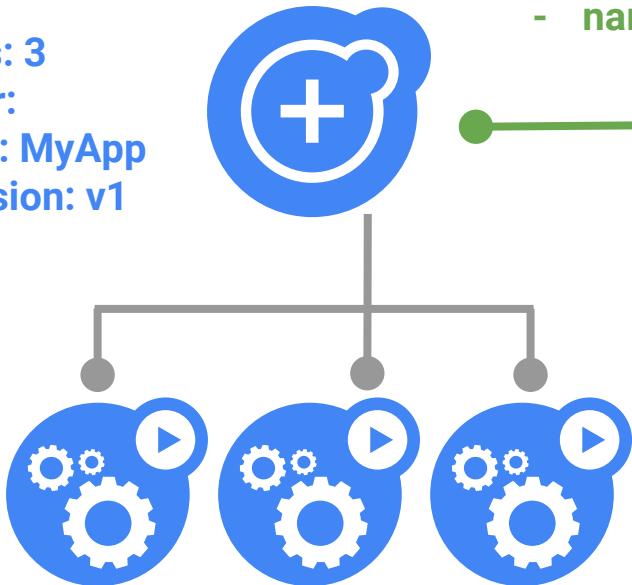
kubectl create ...



Rolling Updates

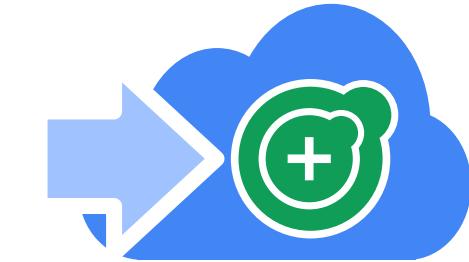
ReplicaSet

- replicas: 3
- selector:
 - app: MyApp
 - version: v1

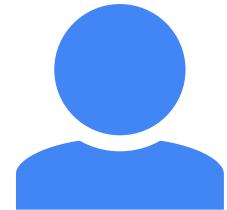


Deployment

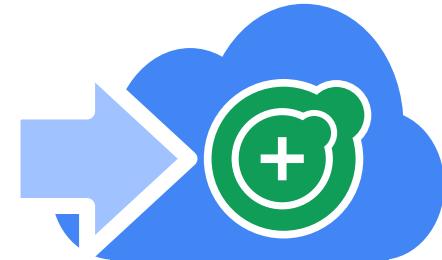
- name: MyApp



kubectl apply ...

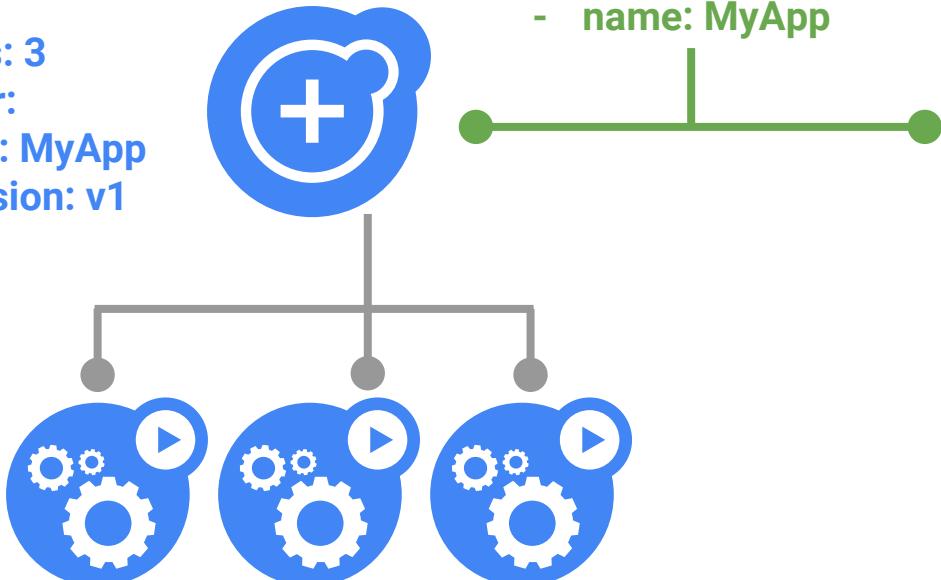


Rolling Updates



Deployment
- name: MyApp

ReplicaSet
- replicas: 3
- selector:
 - app: MyApp
 - version: v1



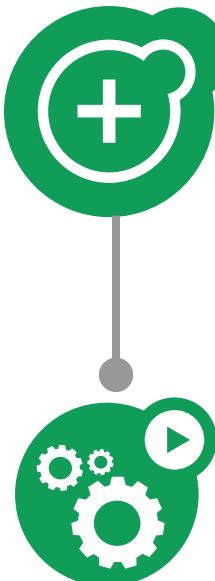
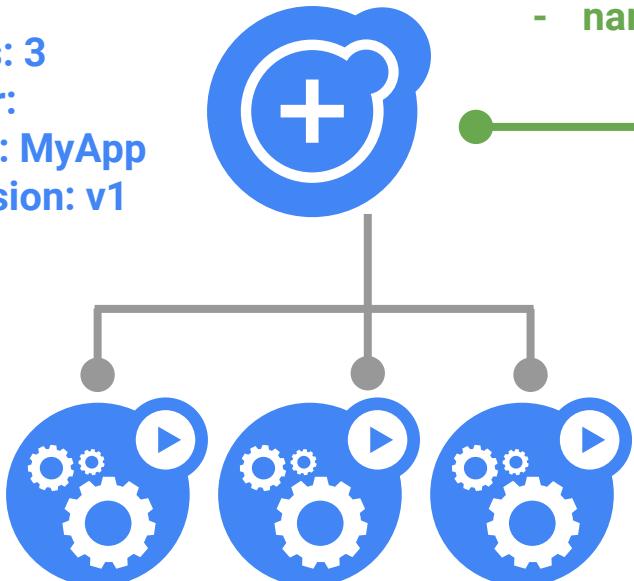
ReplicaSet
- replicas: 0
- selector:
 - app: MyApp
 - version: v2

Rolling Updates



Deployment
- name: MyApp

ReplicaSet
- replicas: 3
- selector:
 - app: MyApp
 - version: v1



ReplicaSet
- replicas: 1
- selector:
 - app: MyApp
 - version: v2

Rolling Updates

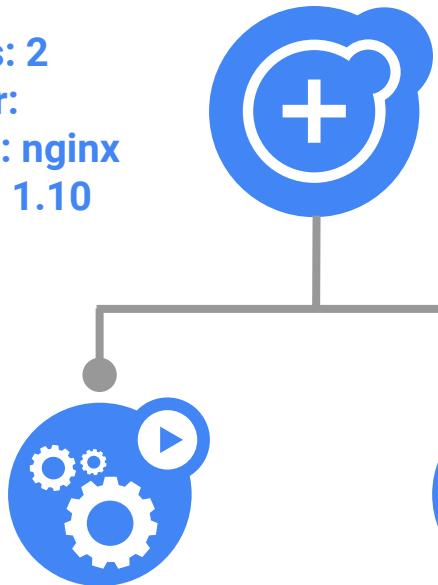


Deployment

- ## - app: nginx

ReplicaSet

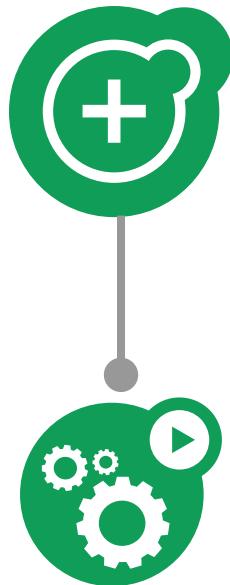
- replicas: 2
 - selector:
 - app: nginx
 - ver: 1.10



www.nationalgeographic.com

ReplicaSet

- replicas: 1
 - selector:
 - app: nginx
 - ver: 1.11

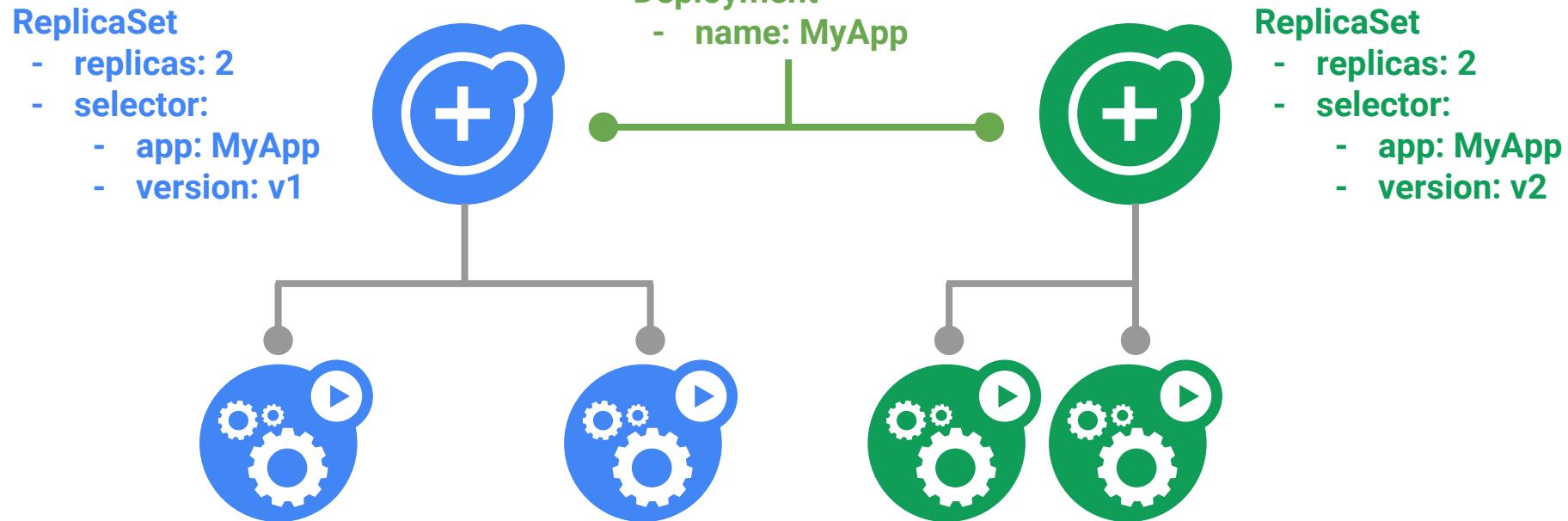


Rolling Updates



Deployment

- name: MyApp



Rolling Updates

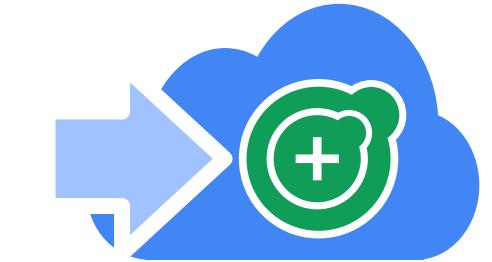
ReplicaSet

- replicas: 1
- selector:
 - app: MyApp
 - version: v1



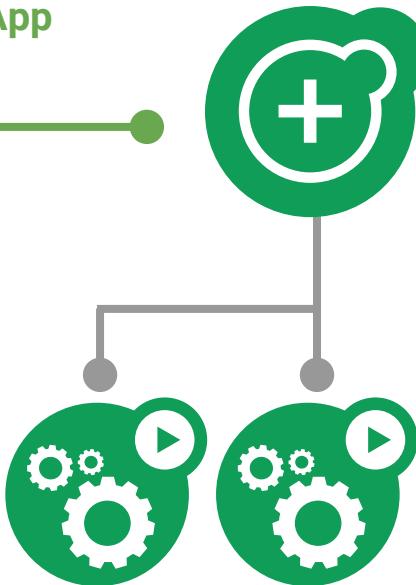
Deployment

- name: MyApp



ReplicaSet

- replicas: 2
- selector:
 - app: MyApp
 - version: v2



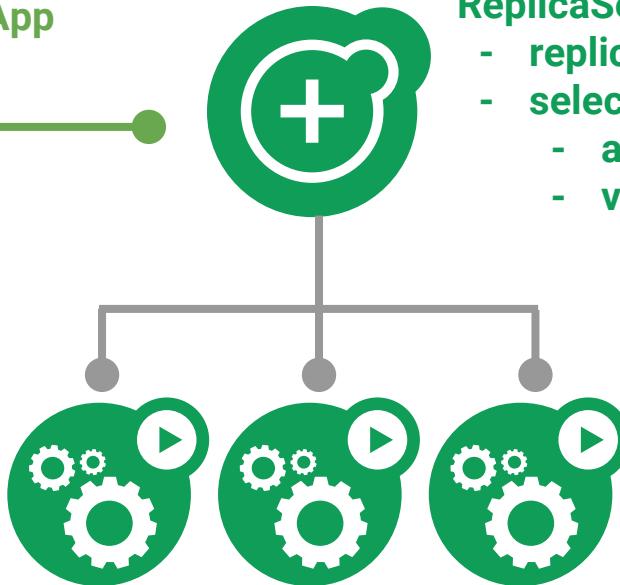
Rolling Updates

ReplicaSet
- replicas: 1
- selector:
 - app: MyApp
 - version: v1



Deployment
- name: MyApp

ReplicaSet
- replicas: 3
- selector:
 - app: MyApp
 - version: v2



Rolling Updates

ReplicaSet

- replicas: 0
- selector:
 - app: MyApp
 - version: v1



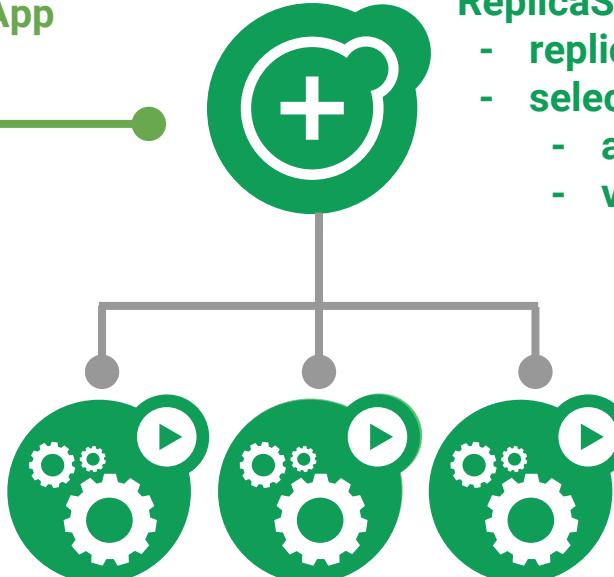
Deployment

- name: MyApp

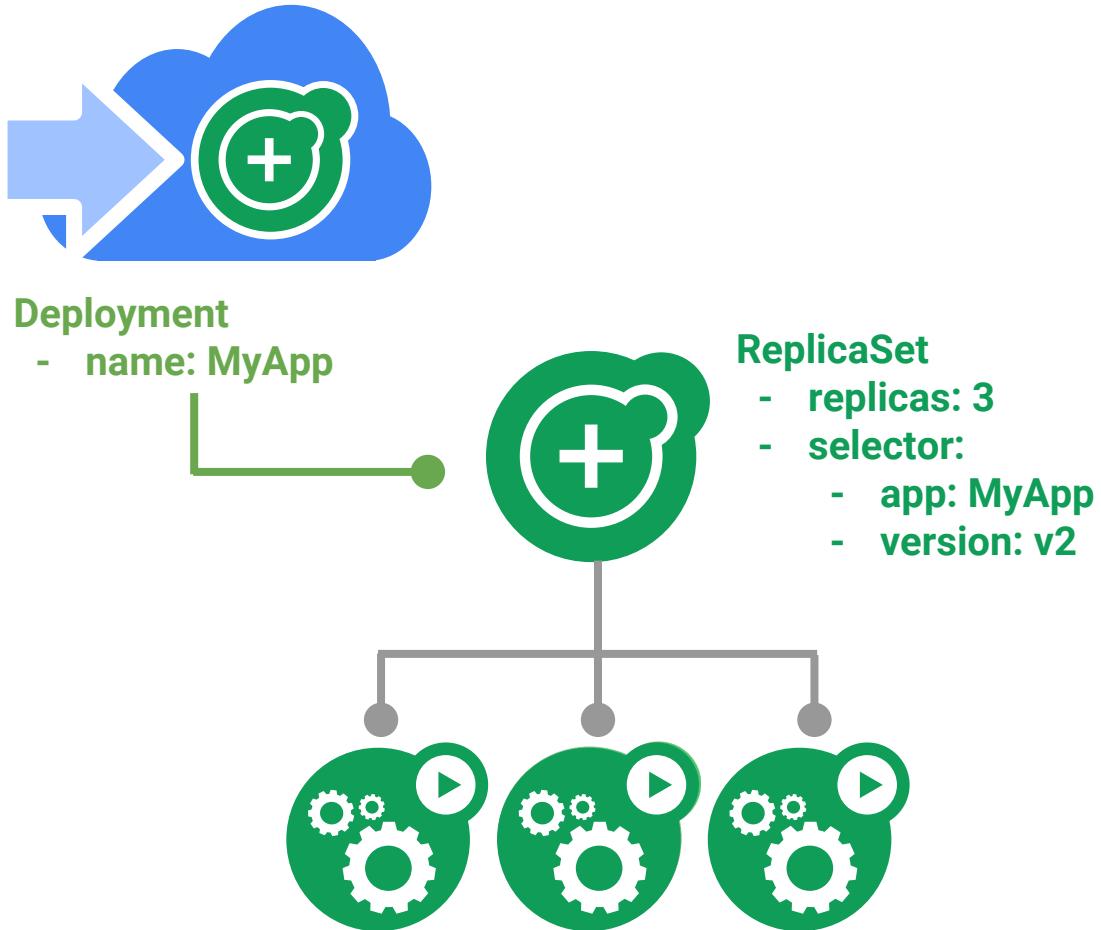


ReplicaSet

- replicas: 3
- selector:
 - app: MyApp
 - version: v2



Rolling Updates



Services

A group of pods that **work together**

- grouped by a selector

Defines access policy

- “load balanced” or “headless”

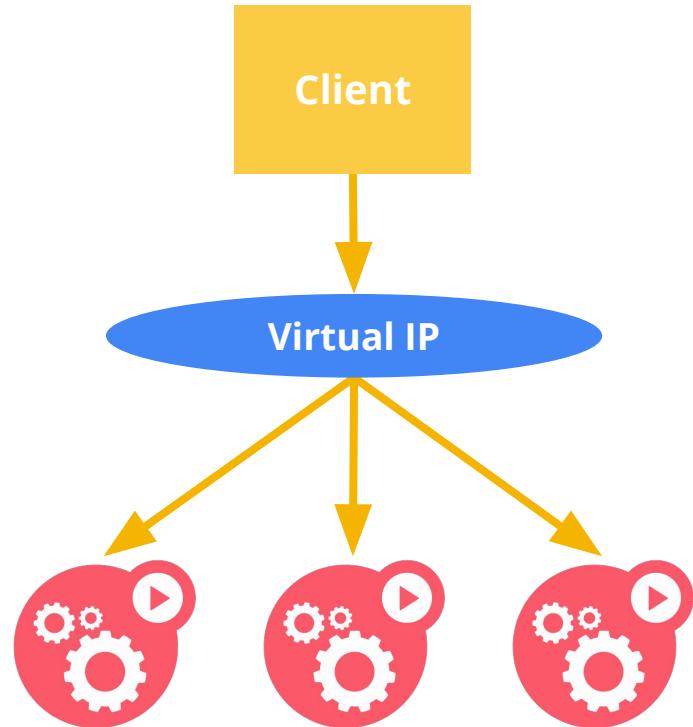
Gets a stable **virtual IP** and port

- sometimes called the service *portal*
- also a DNS name

VIP is managed by *kube-proxy*

- watches all services
- updates iptables when backends change

Hides complexity - ideal for non-native apps



Labels

Arbitrary metadata

Attached to any API object

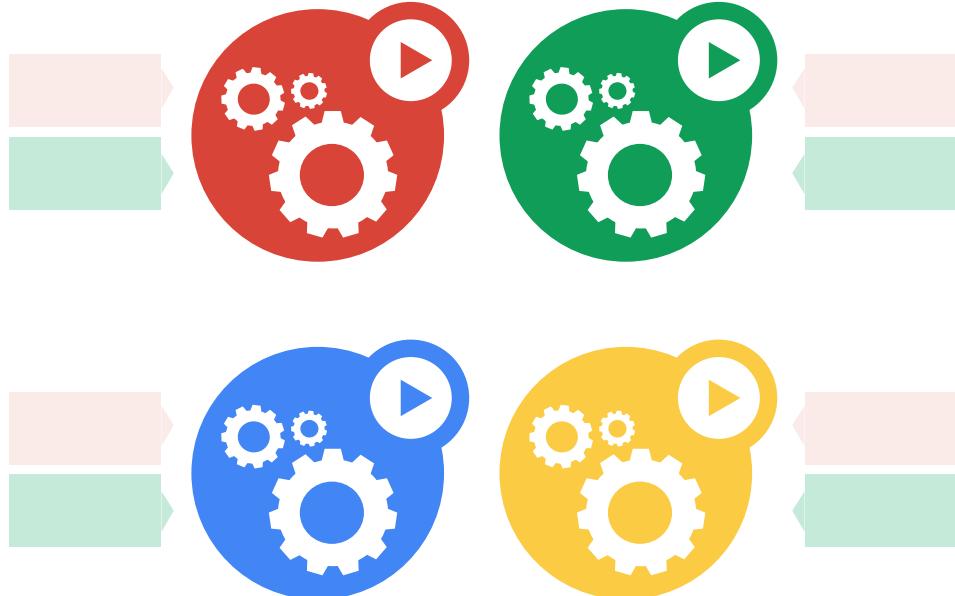
Generally represent **identity**

Queryable by **selectors**

- think SQL 'select ... where ...'

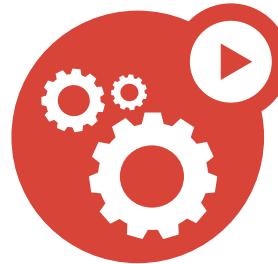
The **only** grouping mechanism

- pods under a ReplicationController
- pods in a Service
- capabilities of a node (constraints)



Selectors

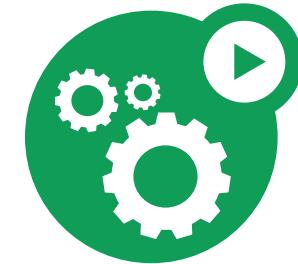
App: MyApp
Phase: prod
Role: FE



App: MyApp
Phase: test
Role: FE



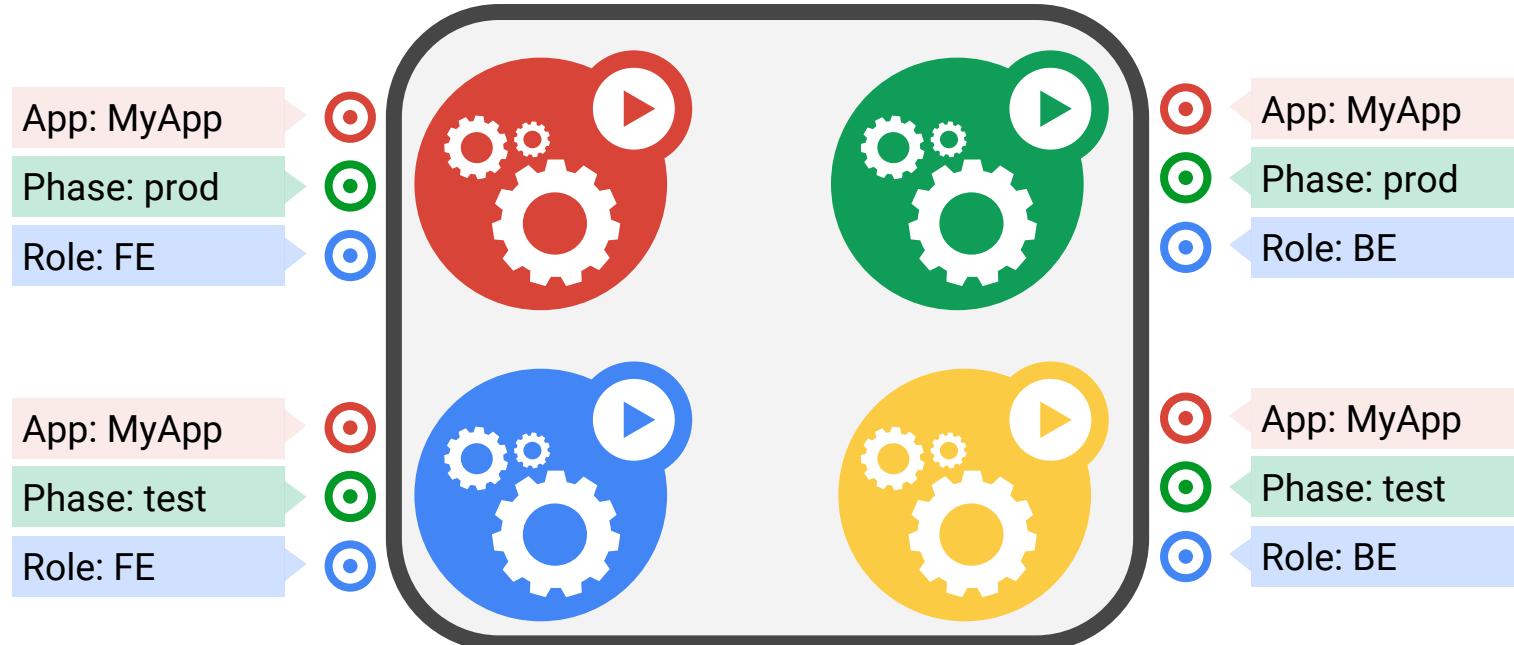
App: MyApp
Phase: prod
Role: BE



App: MyApp
Phase: test
Role: BE

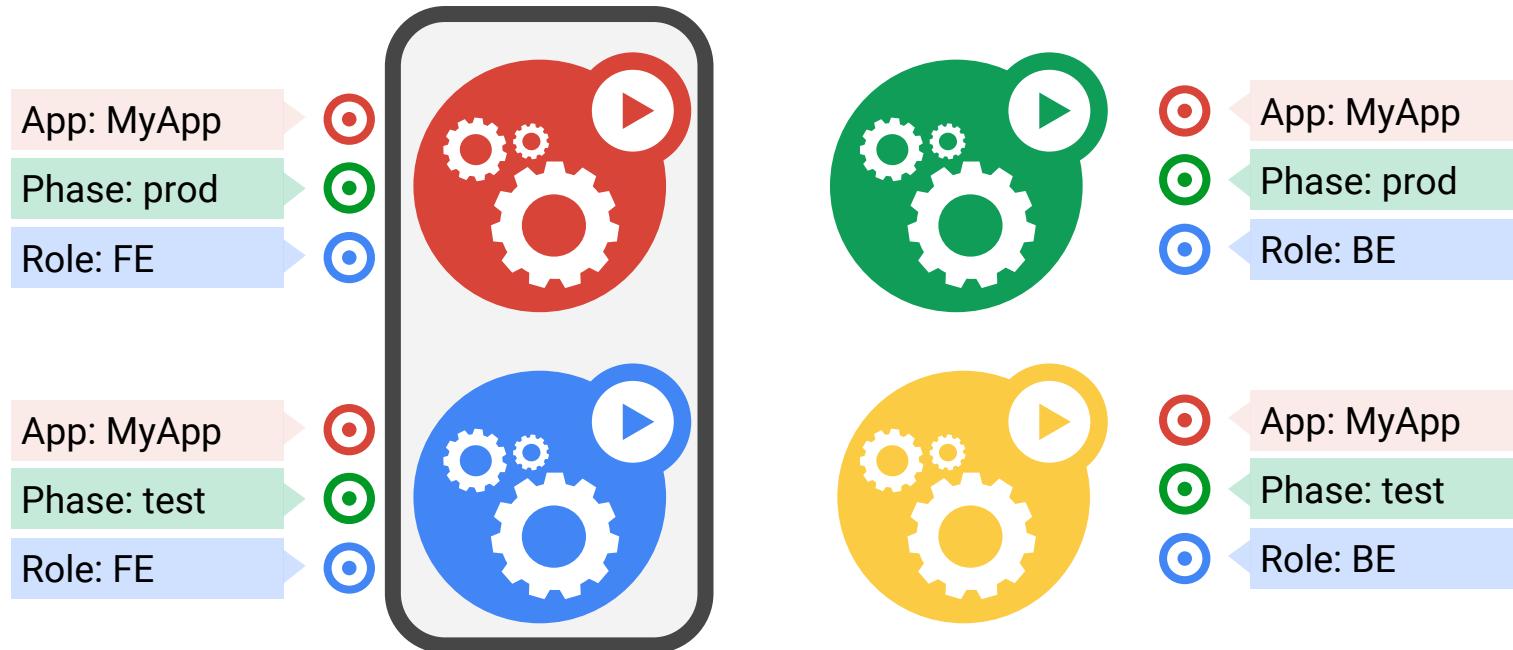


Selectors



App = MyApp

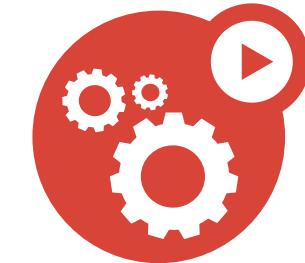
Selectors



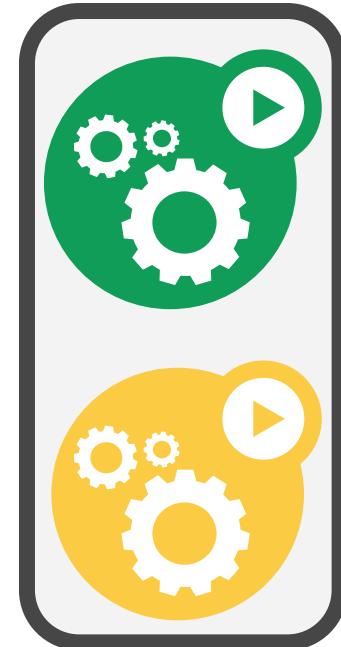
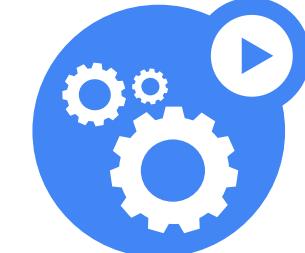
App = MyApp, Role = FE

Selectors

App: MyApp
Phase: prod
Role: FE



App: MyApp
Phase: test
Role: FE

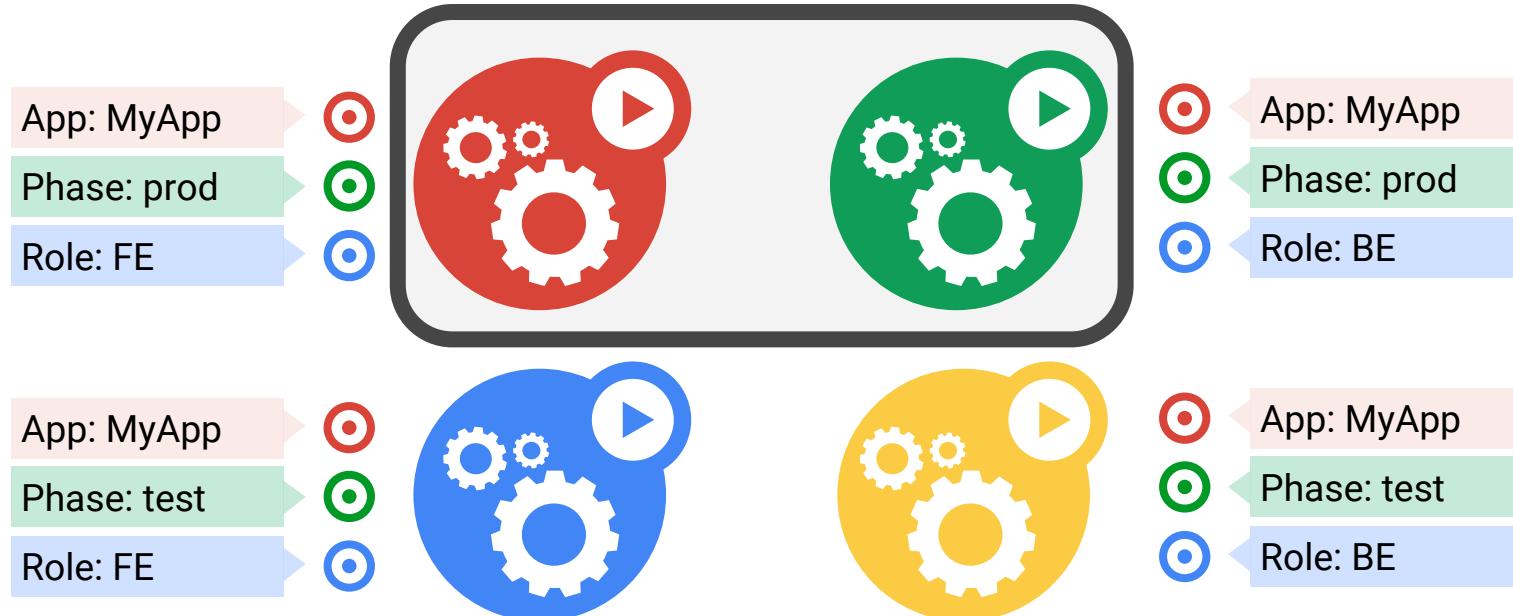


App: MyApp
Phase: prod
Role: BE

App: MyApp
Phase: test
Role: BE

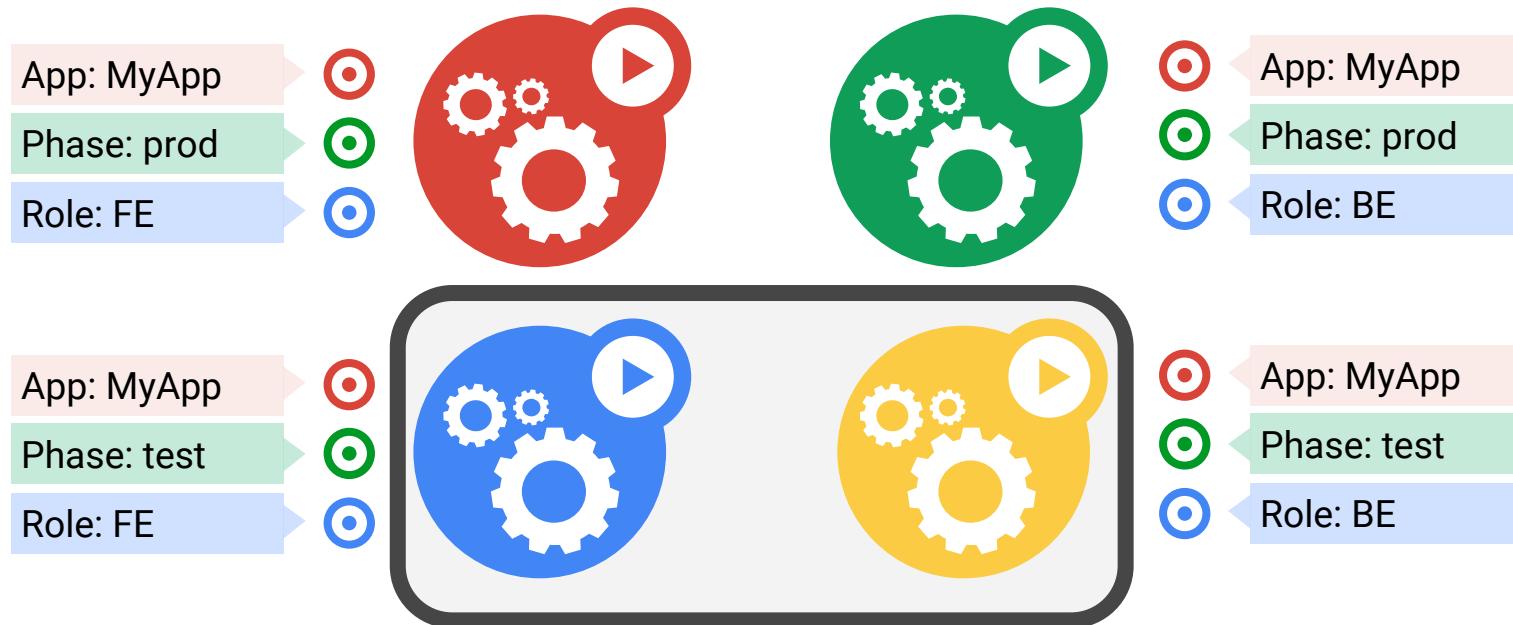
App = MyApp, Role = BE

Selectors



App = MyApp, Phase = prod

Selectors



App = MyApp, Phase = test

Jobs

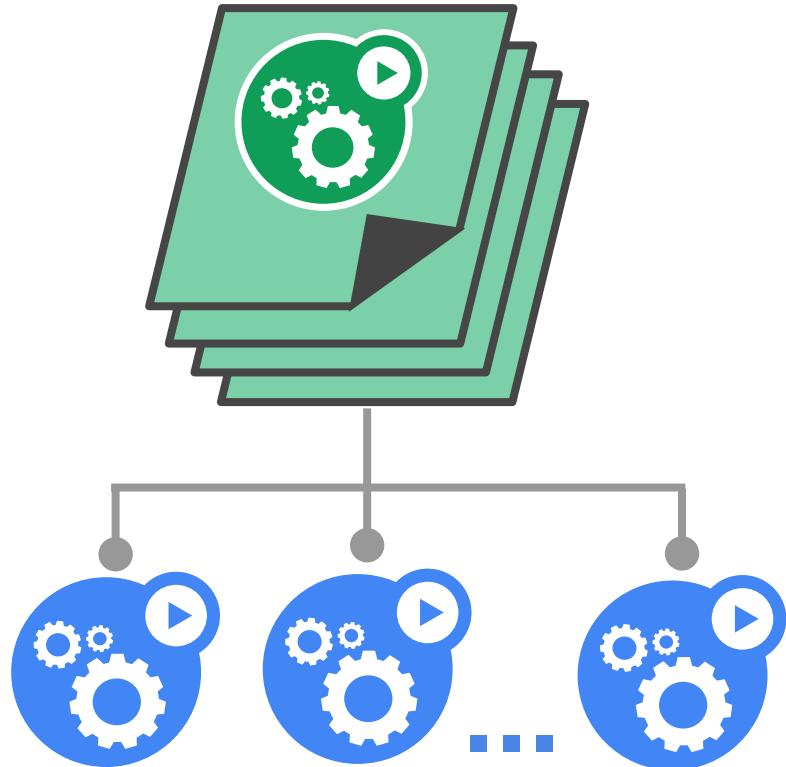
Run-to-completion, as opposed to run-forever

- Express parallelism vs. required completions
- Workflow: restart on failure
- Build/test: don't restart on failure

Aggregates success/failure counts

Built for batch and big-data work

Status: GA in Kubernetes v1.2



Namespaces

Problem: I have too much stuff!

- name collisions in the API
- poor isolation between users
- don't want to expose things like Secrets

Solution: Slice up the cluster

- create new Namespaces as needed
 - per-user, per-app, per-department, etc.
- part of the API - NOT private machines
- most API objects are namespaced
 - part of the REST URL path
- Namespaces are just another API object
- One-step cleanup - delete the Namespace
- Obvious hook for policy enforcement (e.g. quota)





+



Join **Kubernetes** on Slack.

935 users online now of **8062** registered.

GET MY INVITE

or [sign in](#).

slack.kubernetes.io

Thank You

