# Quilt

**Ethan J. Jackson**, Aurojit Panda, Kevin Lin,

Johann Schleier-Smith, Nicholas Sun, Luise Valentin,
Yuen Mei Wan, Scott Shenker

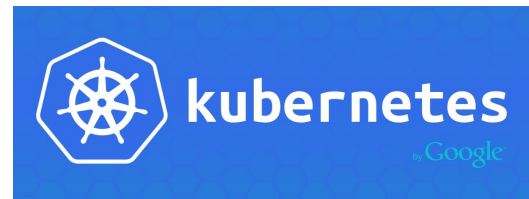quilt.io

# **Everything has an API**

# Compute

# Network

# DevOps

1. Choose a Compute API

2. Choose a Network API

3. Write a Deployment Script

# Deployment Script

Simple right?

# spark-ec2.py

- Official Spark Script

- 1528 Lines of Code

- Incomprehensible

# Network Security

- Status Quo
  - Secure the Perimeter

- A Better Way
  - Tight East-West Firewall
  - Increased script complexity

# Portability

# Quilt

Automated Deployment

# Quilt DSL: Stitch

- Declarative Application Specification
- Lisp Dialect
- Declaration Includes:
  - Application Network and Compute
  - Infrastructure

# Example: Wordpress

# WordPress

```
1  (import "haproxy")
2  (import "memcached")
3  (import "mysql")
4  (import "spark")
5  (import "wordpress")
6  (import "zookeeper")
7
8  (let ((db (mysql.New "db" 2))
9        (memcd (memcached.New "memcd" 3))
10       (wp (wordpress.New "wp" 8 db memcd))
11       (hap (haproxy.New "hap" 2 wp))
12       (zk (zookeeper.New "zk" 3))
13       (spark (spark.New "spark" 2 4 zk)))
14    (connect 7077 (hmapValues spark)
15                  (hmapValues db))
16    (connect 80 "public" hap))
```

# wordpress.New

```
1 (define (New name n db memcd)
2   (let ((dk (makeList n (docker image)))
3         (labelNames (strings.Range name n))
4         (wp (map label labelNames dk)))
5     (configure wp db memcd)
6     (connect 3306 wp (hmapGet db "master"))
7     (connect 3306 wp (hmapGet db "slave"))
8     (connect 11211 wp memcd)
9     wp))
```

(wordpress.New "wp" 8 db memcd)

# WordPress

```
1  (import "haproxy")
2  (import "memcached")
3  (import "mysql")
4  (import "spark")
5  (import "wordpress")
6  (import "zookeeper")
7
8  (let ((db (mysql.New "db" 2))
9        (memcd (memcached.New "memcd" 3))
10       (wp (wordpress.New "wp" 8 db memcd))
11       (hap (haproxy.New "hap" 2 wp))
12       (zk (zookeeper.New "zk" 3))
13       (spark (spark.New "spark" 2 4 zk)))
14    (connect 7077 (hmapValues spark)
15                  (hmapValues db))
16    (connect 80 "public" hap))
```

# WordPress



```
1 (import "haproxy")
2 (import "memcached")
3 (import "mysql")
4 (import "spark")
5 (import "wordpress")
6 (import "zookeeper")
7
8 (let ((db (mysql.New "db" 2))
9       (memcd (memcached.New "memcd" 3))
10      (wp (wordpress.New "wp" 8 db memcd))
11      (hap (haproxy.New "hap" 2 wp))
12      (zk (zookeeper.New "zk" 3))
13      (spark (spark.New "spark" 2 4 zk)))
14   (connect 7077 (hmapValues spark)
15                 (hmapValues db))
16   (connect 80 "public" hap))
```

# Infrastructure

```
1 (define cfg
2    (list (provider "Amazon") (region "us-west-1")
3          (ram 32 64) (cpu 4 8) (sshkey "elided")))
4
5 (makeList 3 (machine (role "Master") cfg))
6 (makeList 32 (machine (role "Worker") cfg))
```

# Infrastructure

```
1 (define cfg
2   (list (provider "Amazon") (region "us-west-1")
3           (ram 32 64) (cpu 4 8) (sshkey "elided")))
4
5 (makeList 3 (machine (role "Master") cfg))
6 (makeList 32 (machine (role "Worker") cfg))
```

# Infrastructure

```
1 (define cfg
2    (list (provider "Amazon") (region "us-west-1")
3          (ram 32 64) (cpu 4 8) (sshkey "elided")))
4
5 (makeList 3 (machine (role "Master") cfg))
6 (makeList 32 (machine (role "Worker") cfg))
```

Azure

Central US

# Geographical Distribution

# Geographical Distribution



Australia - Amazon

| wp | hap | memcd | → | MySQL | ← | spark | zk |

Iowa - Microsoft

| wp | hap | memcd |

Belgium - Google

| memcd | hap | wp |

```
1  (define cfg (list (ram 32 64) (cpu 4 8)
2                     (sshkey "<elided>")))
3
4  (define db (mysql.New "db" 2))
5  (define zk (zookeeper.New "zk" 3))
6  (define spark (spark.New "spark" 2 4 zk))
7  (connect 7077 (hmapValues spark) (hmapValues db))
8
9  (define (makeLoc prvd rgn)
10   (list (provider prvd) (region rgn)))
11
12 (define (makePod name)
13   (let ((memcd (memcached.New (+ name "-mem") 1))
14         (wp (wordpress.New (+ name "-wp")
15                            2 db memcd))
16         (hap (haproxy.New (+ name "-hap") 1 wp)))
17     (connect 80 "public" hap)
18     (list memcd wp hap)))
19
20 (define (deploy pod loc)
21   (makeList 16 (machine (role "Worker") cfg loc))
22   (place (machineRule "on" loc) pod))
23
24 (deploy (makePod "gce")
25         (makeLoc "Google" "europe-west1-b"))
26
27 (deploy (makePod "azure")
28         (makeLoc "Azure" "Central US"))
29
30 (let ((loc (makeLoc "Amazon" "ap-southeast-2"))
31       (nodes (append (makePod "aws") zk
32                      (hmapValues db)
33                      (hmapValues spark))))
34   (machine (role "Master") cfg loc)
35   (deploy nodes loc))
```

# Stitch

# Stitch

- Lisp (Scheme)
  - Variables
  - Arithmetic
  - Functions
  - Modules
- Domain Specific Primitives

# Stitch — Primitives

- Application Primitives
  - "docker", "label", "connect", "place", "setEnv"
- Infrastructure Primitives
  - "machine"
  - "role", "provider", "region", "ram", "cpu", "size"

# Stitch — Primitives

```
(label "spark-master" (docker "quilt/spark" "start-master.sh"))

(label "spark-worker"
       (makeList 10 (docker "quilt/spark" "start-worker.sh"
                            "spark://spark-master.di:7077")))

// Spark workers listen on random ports. Must open up everything.
(connect (list 1000 65535)
         (list "spark-master" "spark-worker")
         "spark-worker")

(connect 7077 "spark-worker" "spark-master")
```

# Stitch — Primitives



```
(label "spark-master" (docker "quilt/spark" "start-master.sh"))

(label "spark-worker"
        (makeList 10 (docker "quilt/spark" "start-worker.sh"
                                "spark://spark-master.di:7077")))

// Spark workers listen on random ports. Must open up everything.
(connect (list 1000 65535)
         (list "spark-master" "spark-worker")
         "spark-worker")

(connect 7077 "spark-worker" "spark-master")
```

# Quilt Architecture

# Goals

- Simple
- Robust
- Portable

# Quilt Architecture

# Infrastructure Controller

- Import Infrastructure Spec
- Update Cluster
- Cloud Provider Plugins
  - Amazon EC2
  - Google Compute Engine
  - Microsoft Azure

# Cloud Provider

- Boot, Stop, List
- Network Reachability
- **Application Agnostic**

```go
type Provider interface {
    Connect(namespace string) error

    List() ([]Machine, error)

    Boot([]Machine) error

    Stop([]Machine) error

    SetACLs(acls []string) error

    Disconnect()

    ChooseSize(ram dsl.Range, cpu dsl.Range,
            maxPrice float64) string
}
```

# Quilt Cluster

- Virtual Machines Running …
- Application Containers
- Open Virtual Network
  - SDN Overlay
- **Infrastructure Agnostic**

Cluster

| Worker | Worker | Worker | Worker |
|--------|--------|--------|--------|

| Master | Master | Master |
|--------|--------|--------|

# Unsolved Problems

- Application Configuration
- Container Security
- State
- External Services

# Related Work

# Related Work

- Container Orchestrators
  - Kubernetes, Docker Swarm, Mesos, Nomad
  - No explicit application specification
  - No tight network firewall
- Quilt is a policy layer above these systems

# Related Work

- Docker Compose / Kubernetes Helm
  - Declare Groups of Containers to Boot
- Static Data Serialization Format
  - Poor modularity
- Missing network graph

# Future Work

# Stitch: New Domains

- Security policy
  - Key Management
  - User Management
- Data
- Application Configuration

# Stitch Analysis

- Verification
  - Stitch specifies app *entirely*
  - Simpler to verify than deployed systems
- Reachability
- Availability

# Summary

- Portable Application Deployment
- Strict Network Security
- Modular, Shareable, Reusable Specifications
- In Future — Formal Analysis

# Thank you

quilt.io

ejj@eecs.berkeley.edu