# Managing Containers with Helix
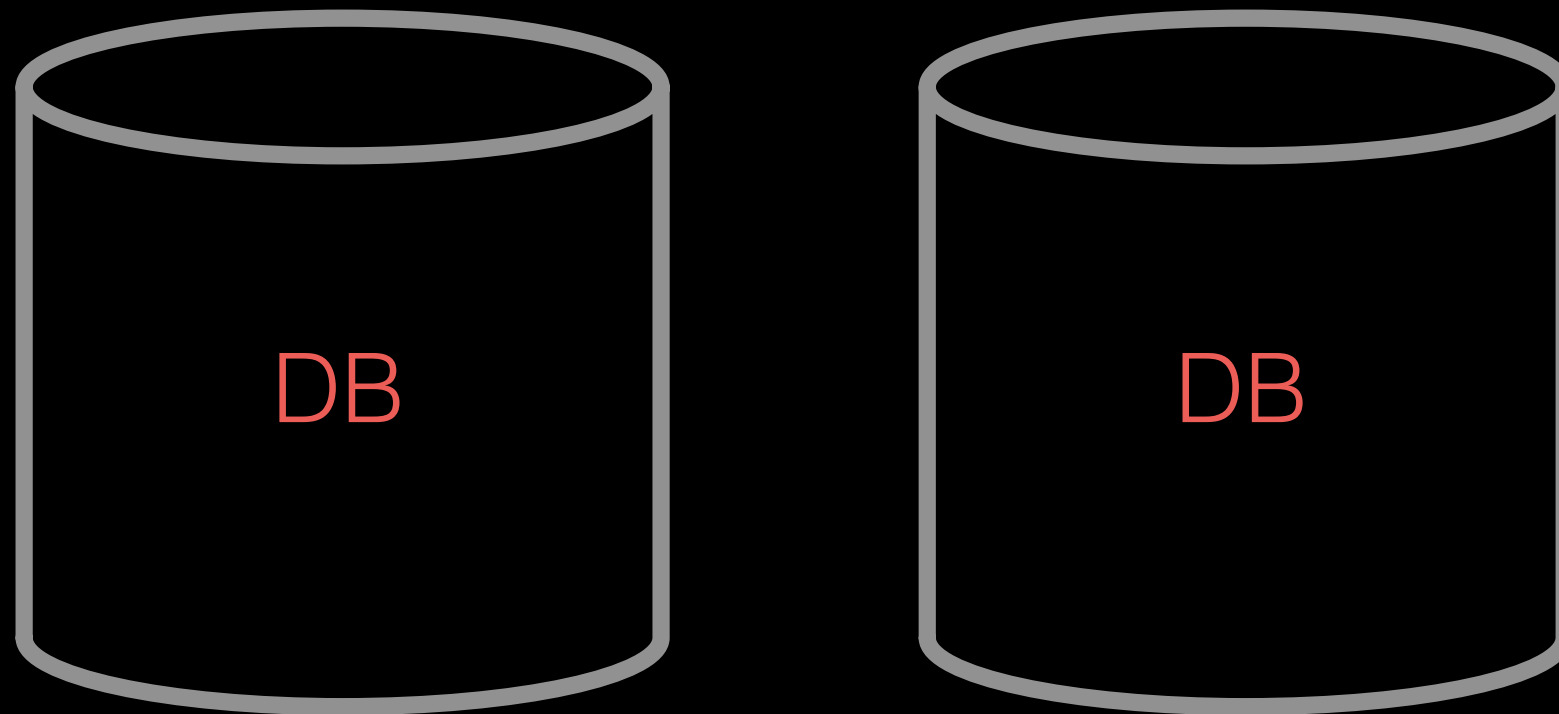
**Kanak Biscuitwala**

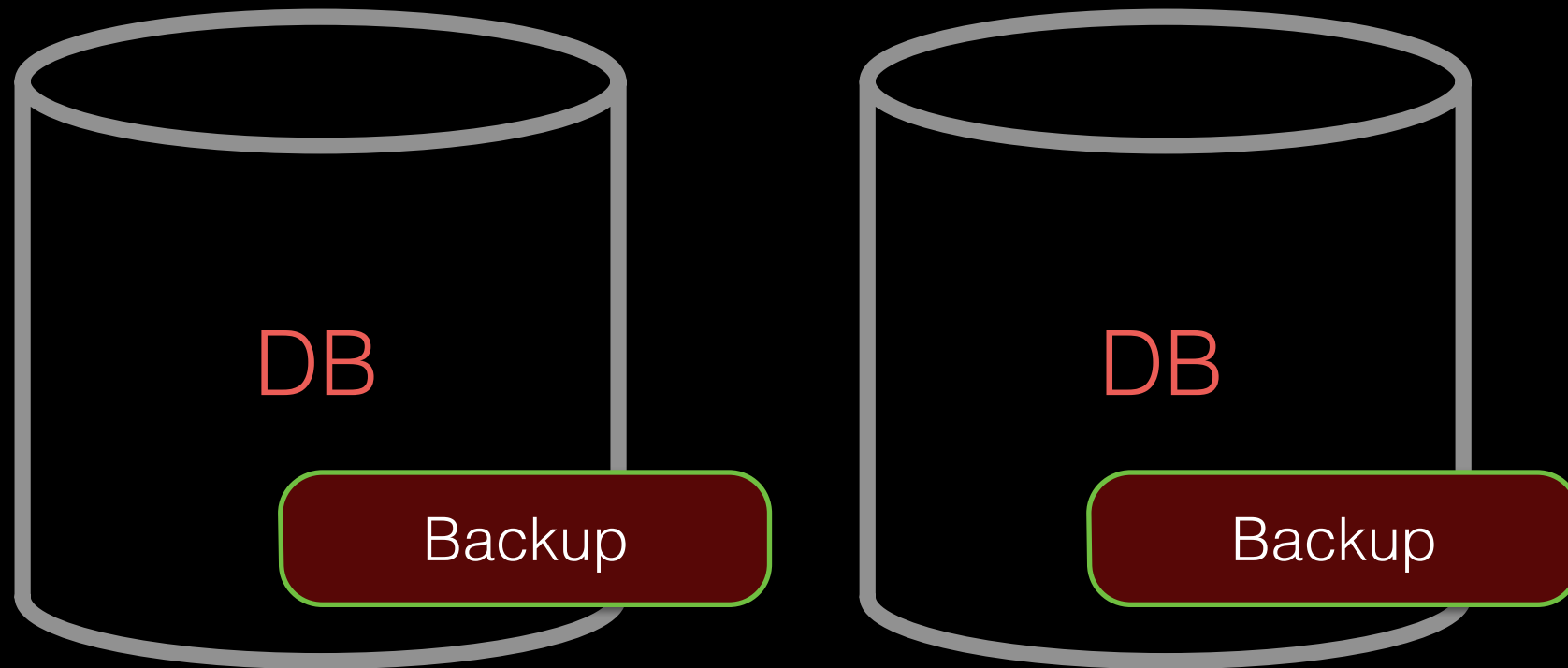**Jason Zhang**

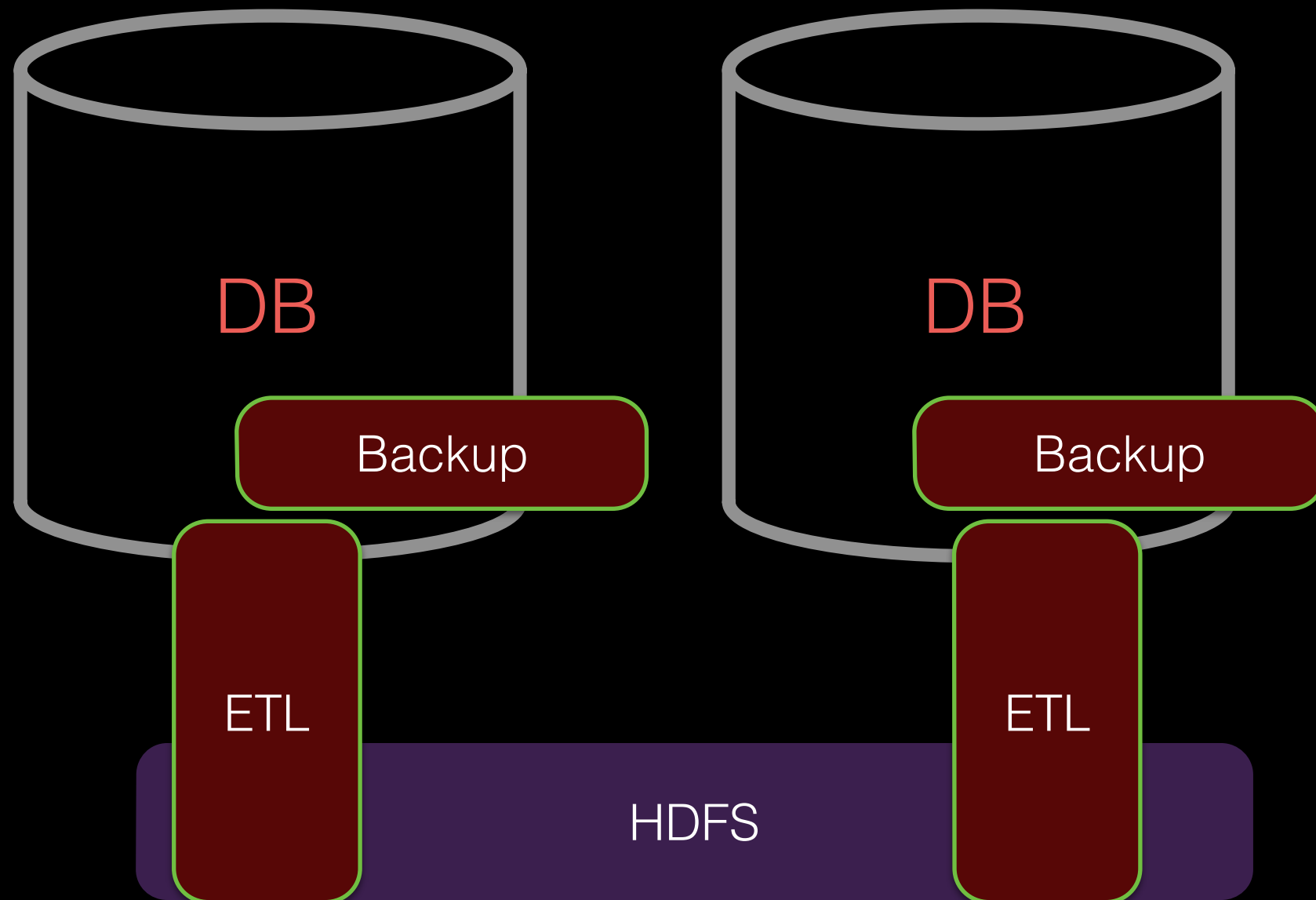Apache Helix Committers @ LinkedIn

helix.apache.org
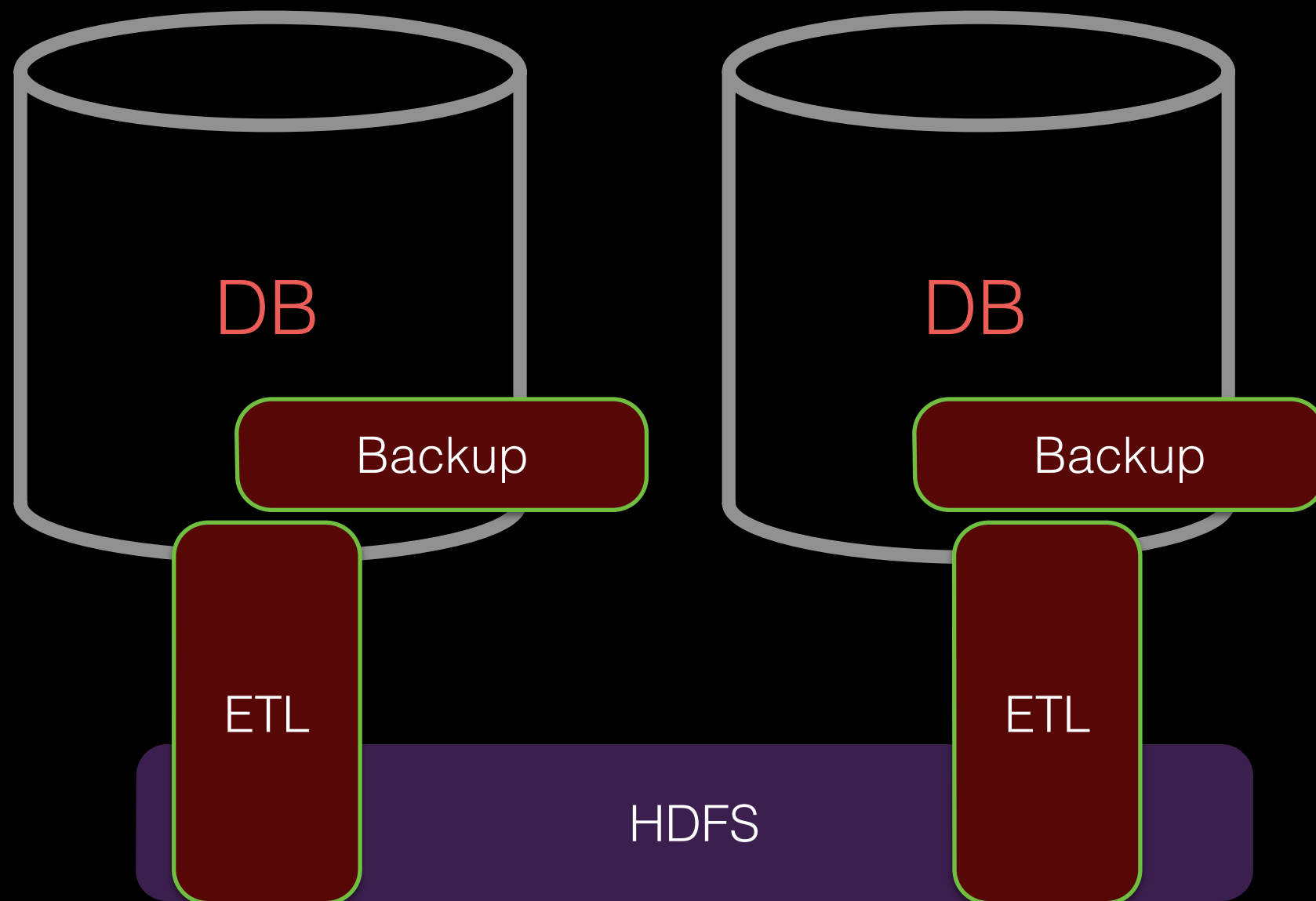
@apachehelix

# Intersection of Job Types

# Intersection of Job Types

# Intersection of Job Types

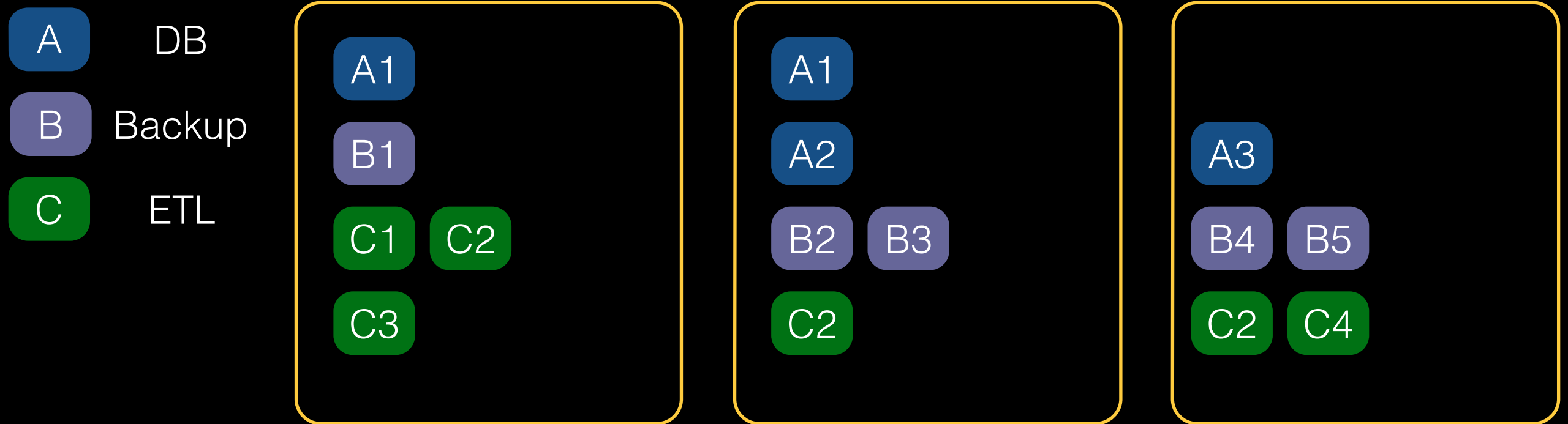# Intersection of Job Types



Long-running and batch jobs running together!

# Cloud Deployment

| A | online |
| B | nearline |
| C | batch |

**Box 1:** A1, B1, C1, C2, C3

**Box 2:** A1, A2, B2, B3, C2

**Box 3:** A3, B4, B5, C2, C4

Applications with diverse requirements running
together in a datacenter

HELIX

# Cloud Deployment

| A | DB |
|---|---|
| B | Backup |
| C | ETL |

**Box 1:**
A1
B1
C1  C2
C3

**Box 2:**
A1
A2
B2  B3
C2

**Box 3:**
A3
B4  B5
C2  C4

Applications with diverse requirements running together in a datacenter

HELIX

# Processes on Machines

Machine    Process    VM    Container

# Processes on Machines

Process

No Isolation

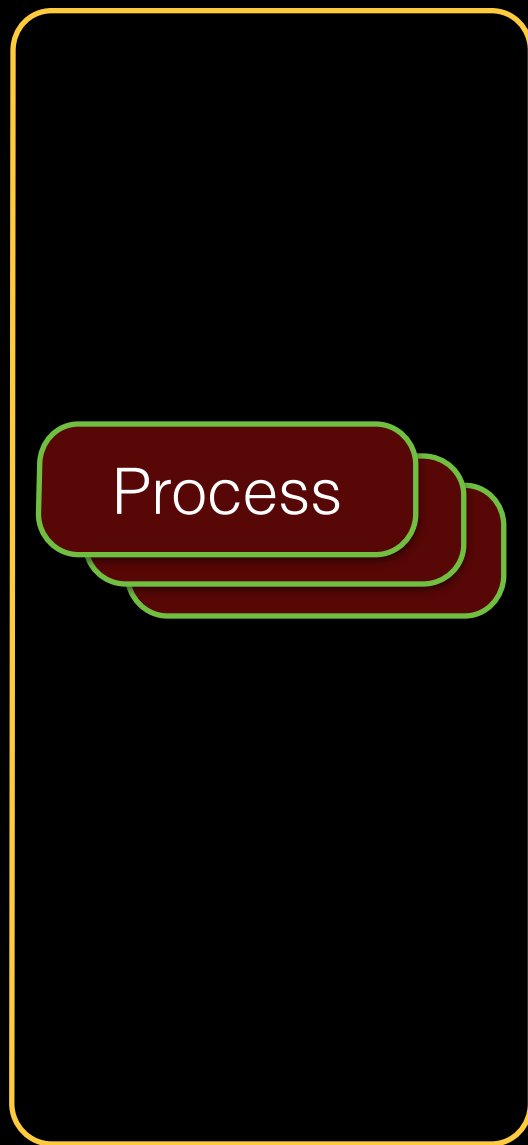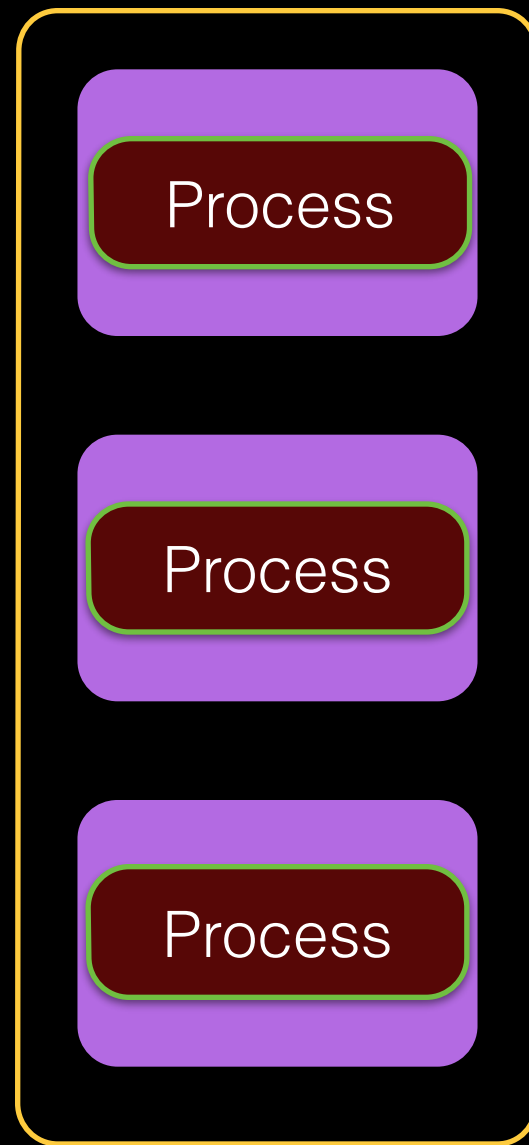Machine   Process   VM   Container

# Processes on Machines

Process

Process

Process

Process

No Isolation

VM-based Isolation

Machine | Process | VM | Container

# Processes on Machines

**No Isolation**

Process

**VM-based Isolation**

Process

Process

Process

**Container-based Isolation**

Process

Process

Machine   Process   VM   Container

HELIX

# Processes on Machines

- Run as individual processes
  - Poor isolation or poor utilization
- Virtual machines
  - Better isolation
  - Xen, Hyper-V, ESX, KVM
- Containers
  - cgroup
  - YARN, Mesos
  - Super lightweight, dynamic based on application requirements

# Processes on Machines

Virtualization and containerization significantly improve process isolation and open up possibilities for efficient utilization of physical resources
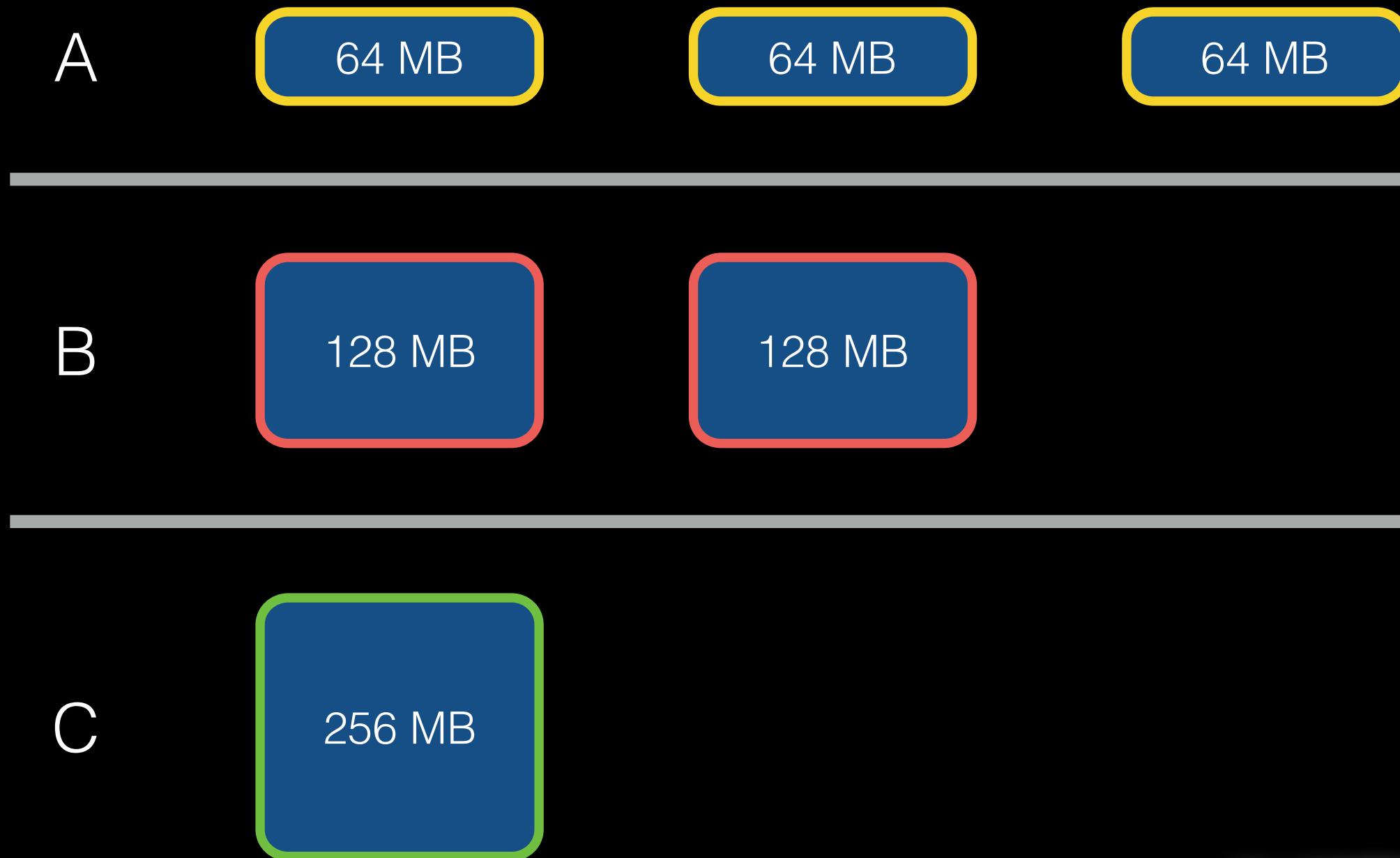
# Container-Based Solution
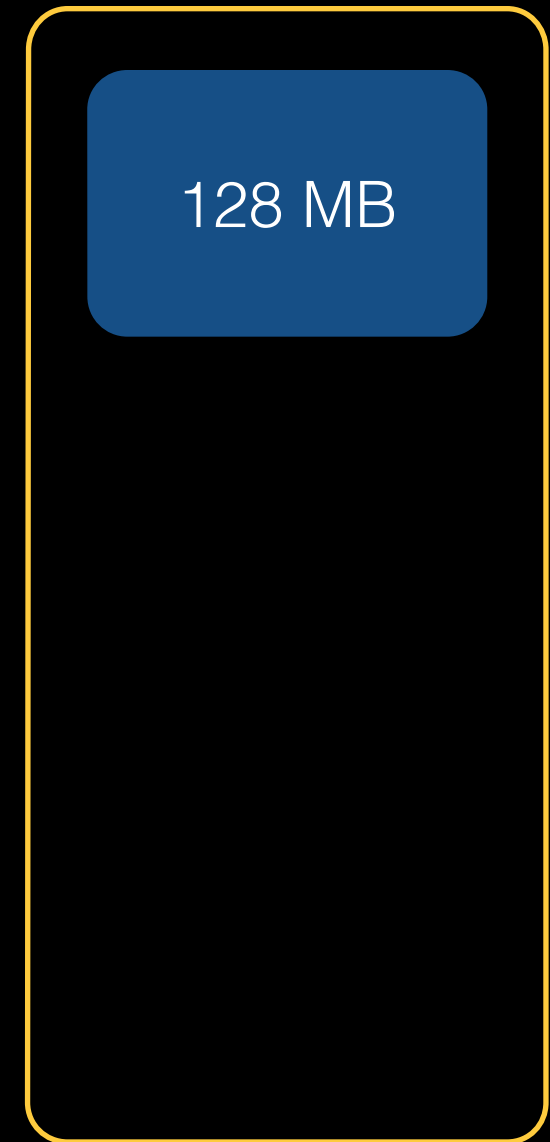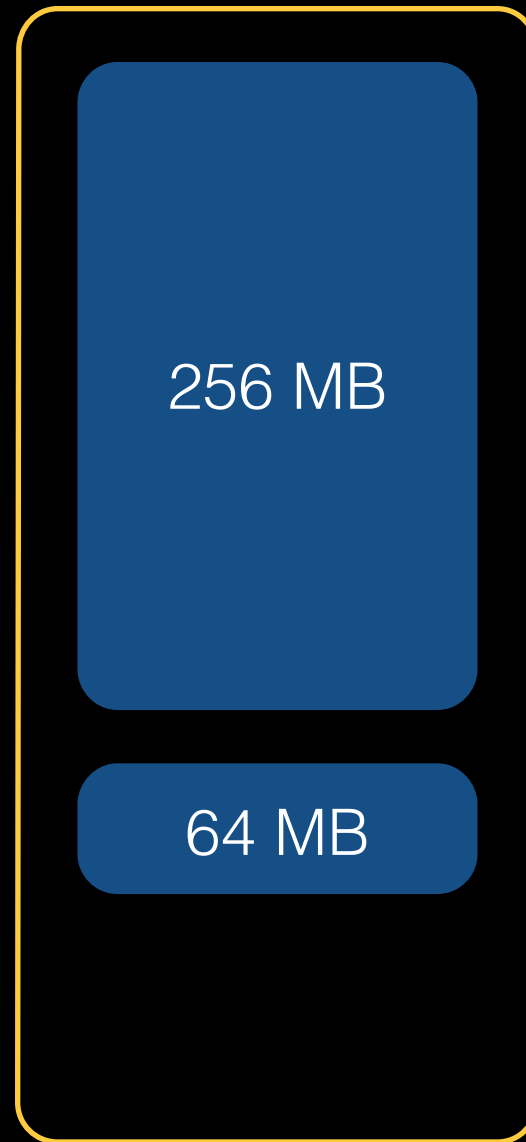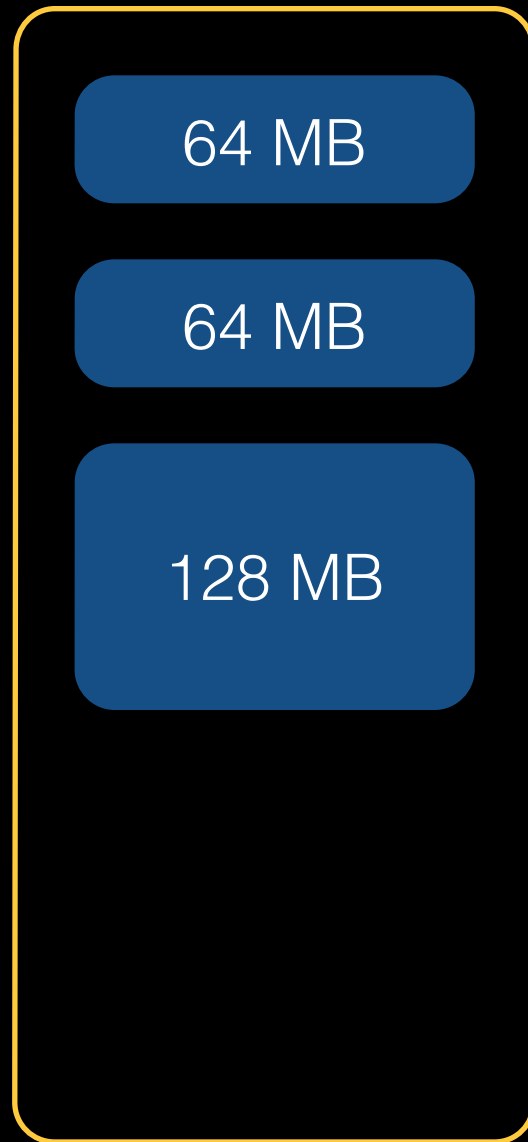
# Container-Based Solution

System Requirements

A    64 MB     64 MB     64 MB

B    128 MB     128 MB

C    256 MB

HELIX

# Container-Based Solution

Allocation

**Machine**
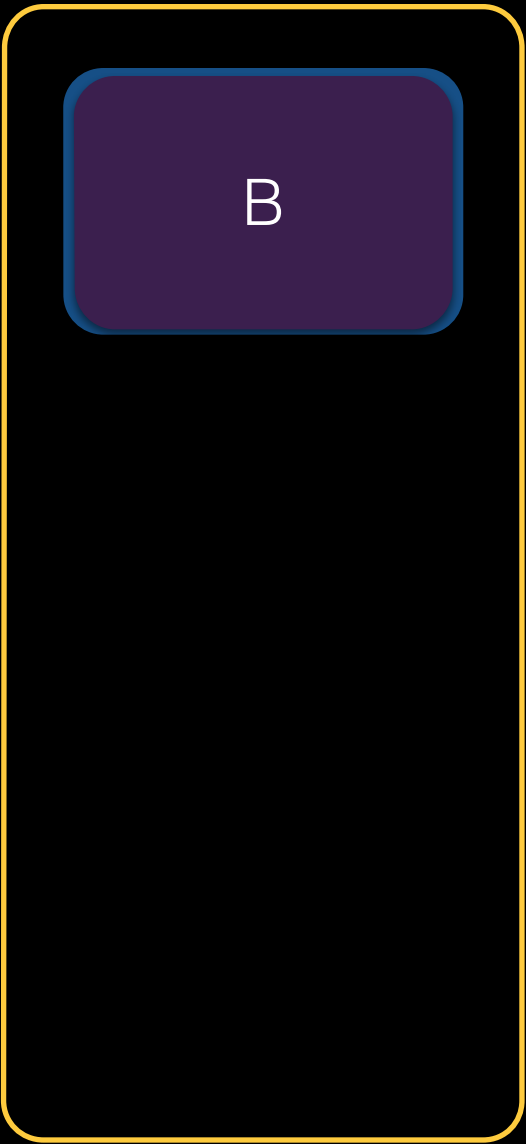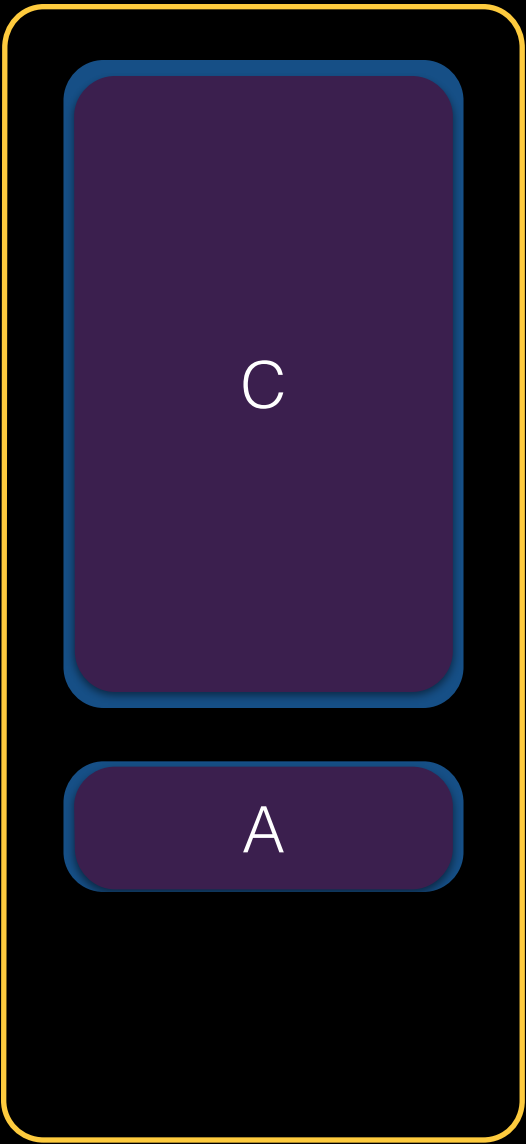
**Container**

64 MB

64 MB

128 MB

256 MB

64 MB

128 MB

HELIX

# Container-Based Solution

Allocation

Machine

Process    Container

A

A

B

C

A

B

# Container-Based Solution

Allocation

Machine

Process    Container

| | | |
|---|---|---|
| A | C | B |
| A | | |
| B | A | |

Containerization is powerful!

HELIX

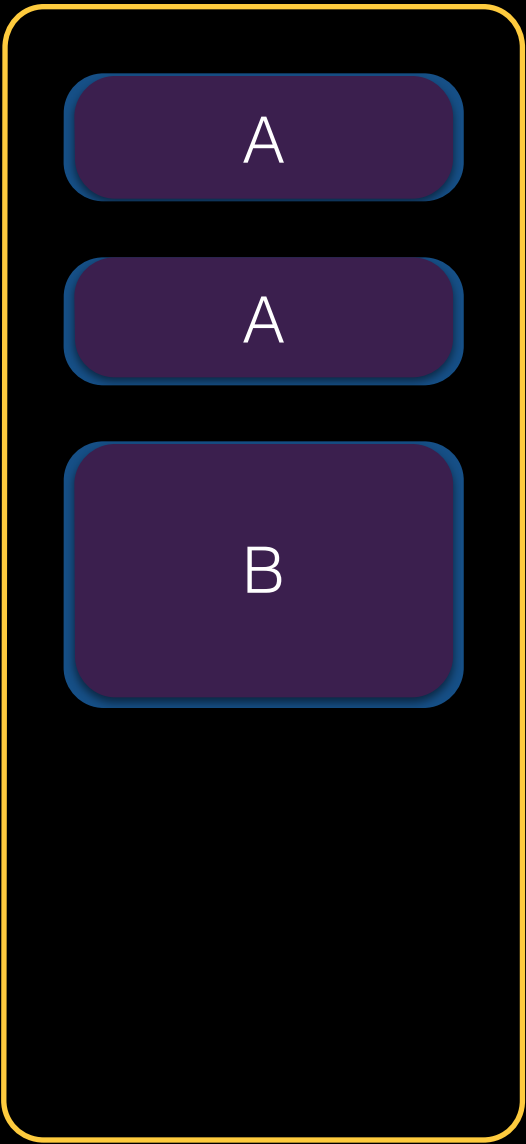# Container-Based Solution

Allocation

Machine

Process

Container

A

A

B
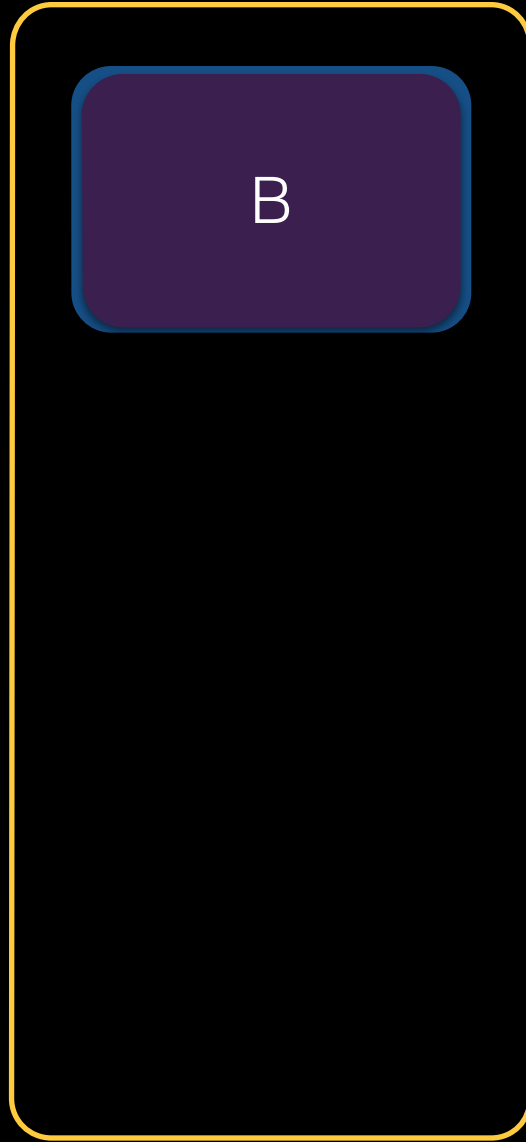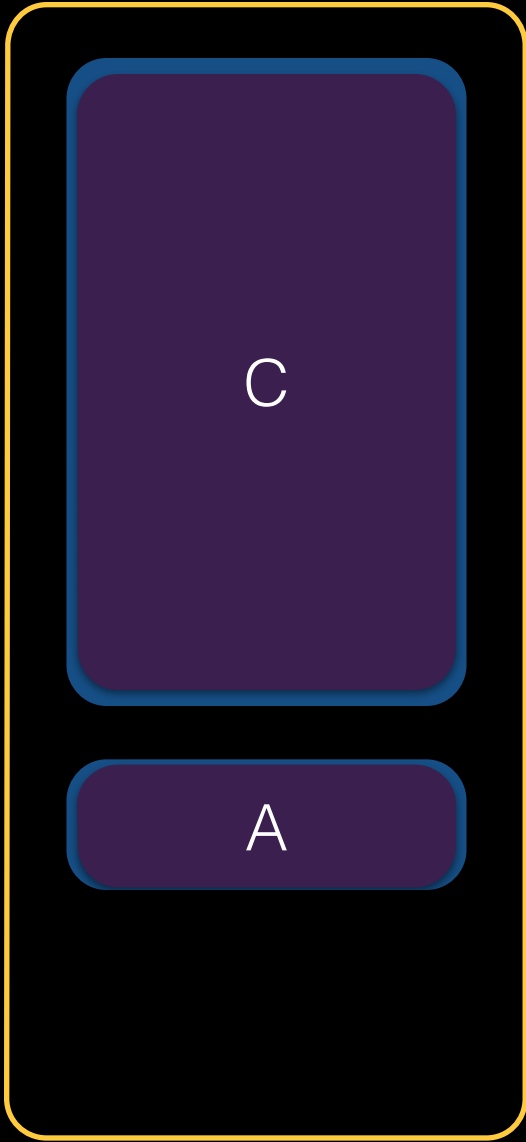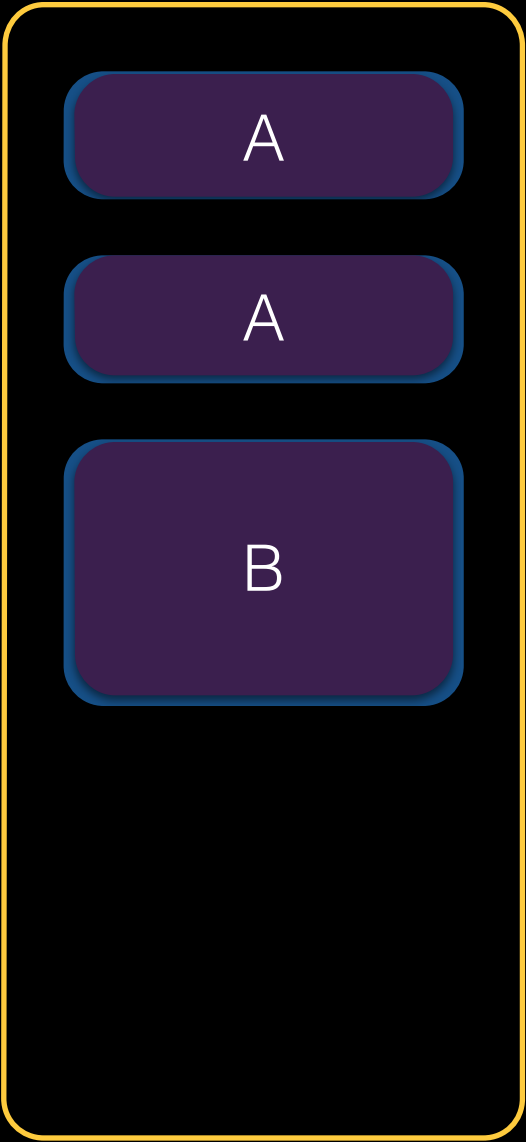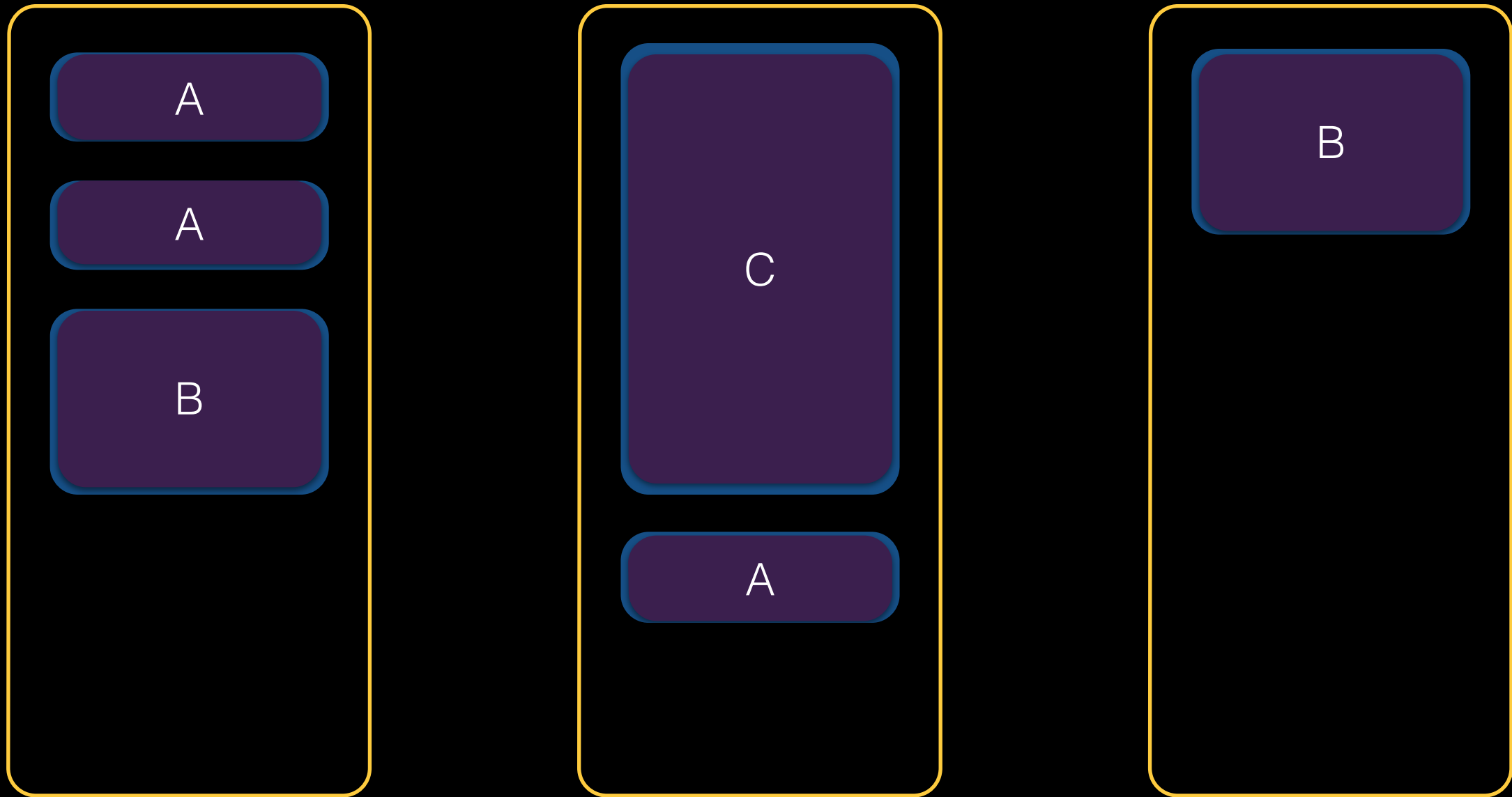
C

A

B

Containerization is powerful!

But do processes always fit so nicely?

HELIX

# Container-Based Solution

Over-Utilization

Machine

Process    Container

256 MB

# Container-Based Solution

Over-Utilization

Machine

Process

Container

Process 1
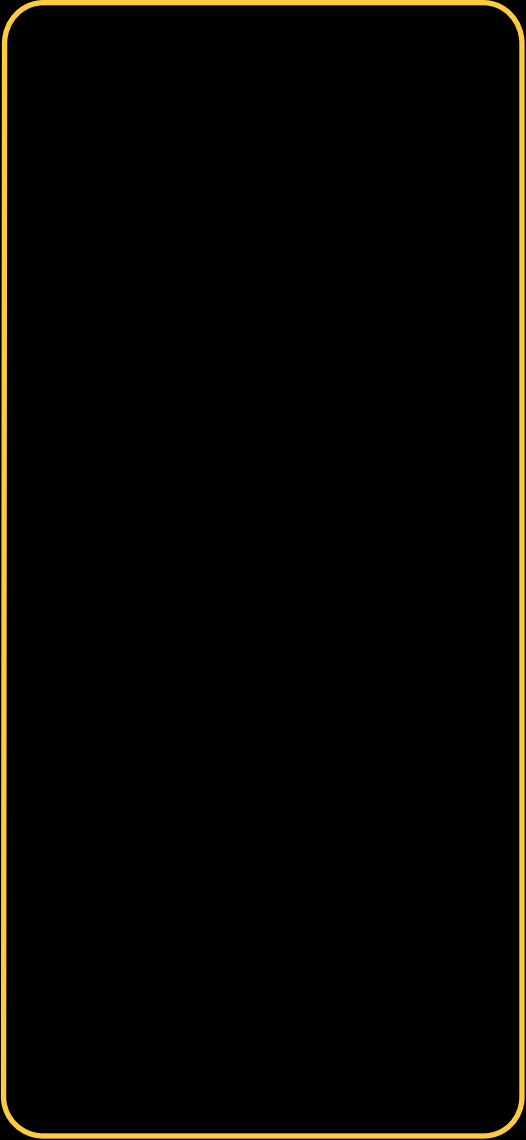
HELIX

# Container-Based Solution

Over-Utilization

Machine

Process
Container

Process 1

Outcome: Preemption and relaunch

HELIX

# Container-Based Solution

Over-Utilization

Machine

Process    Container

384 MB
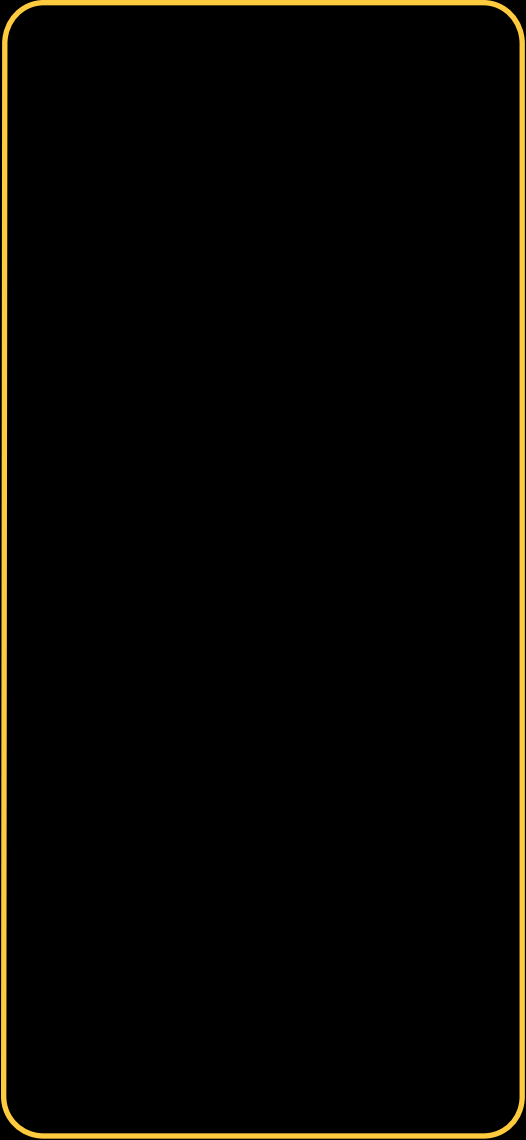
Outcome: Preemption and relaunch

HELIX

# Container-Based Solution

Over-Utilization

Machine

Process    Container

Process 1

Outcome: Preemption and relaunch

HELIX

# Container-Based Solution

Under-Utilization

Machine

Process  Container

384 MB

128 MB

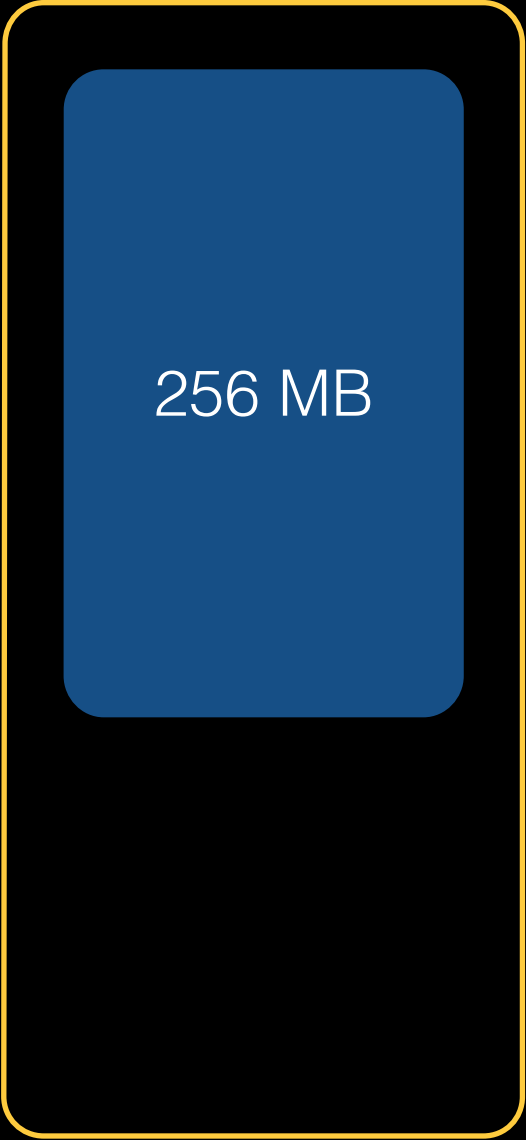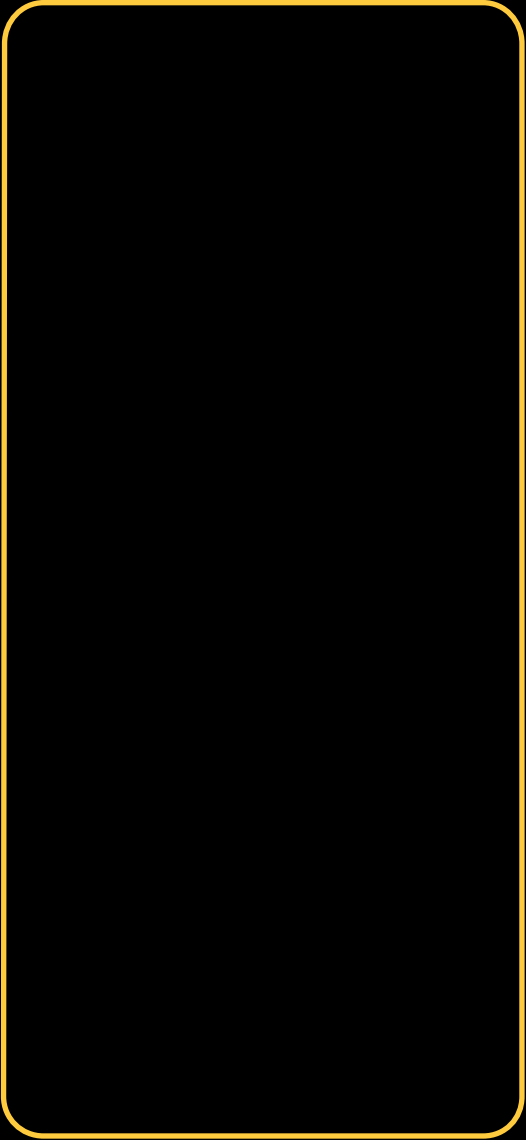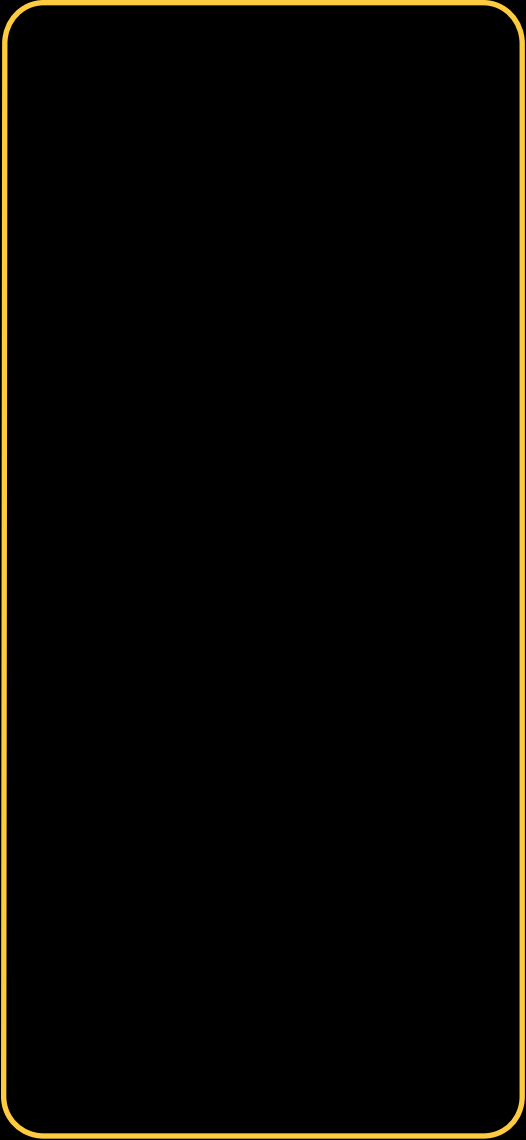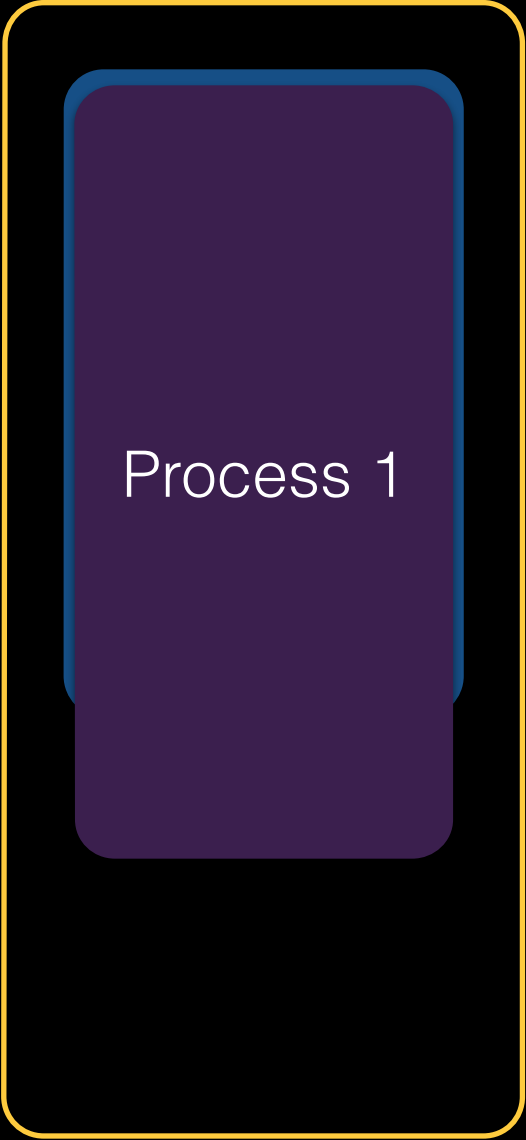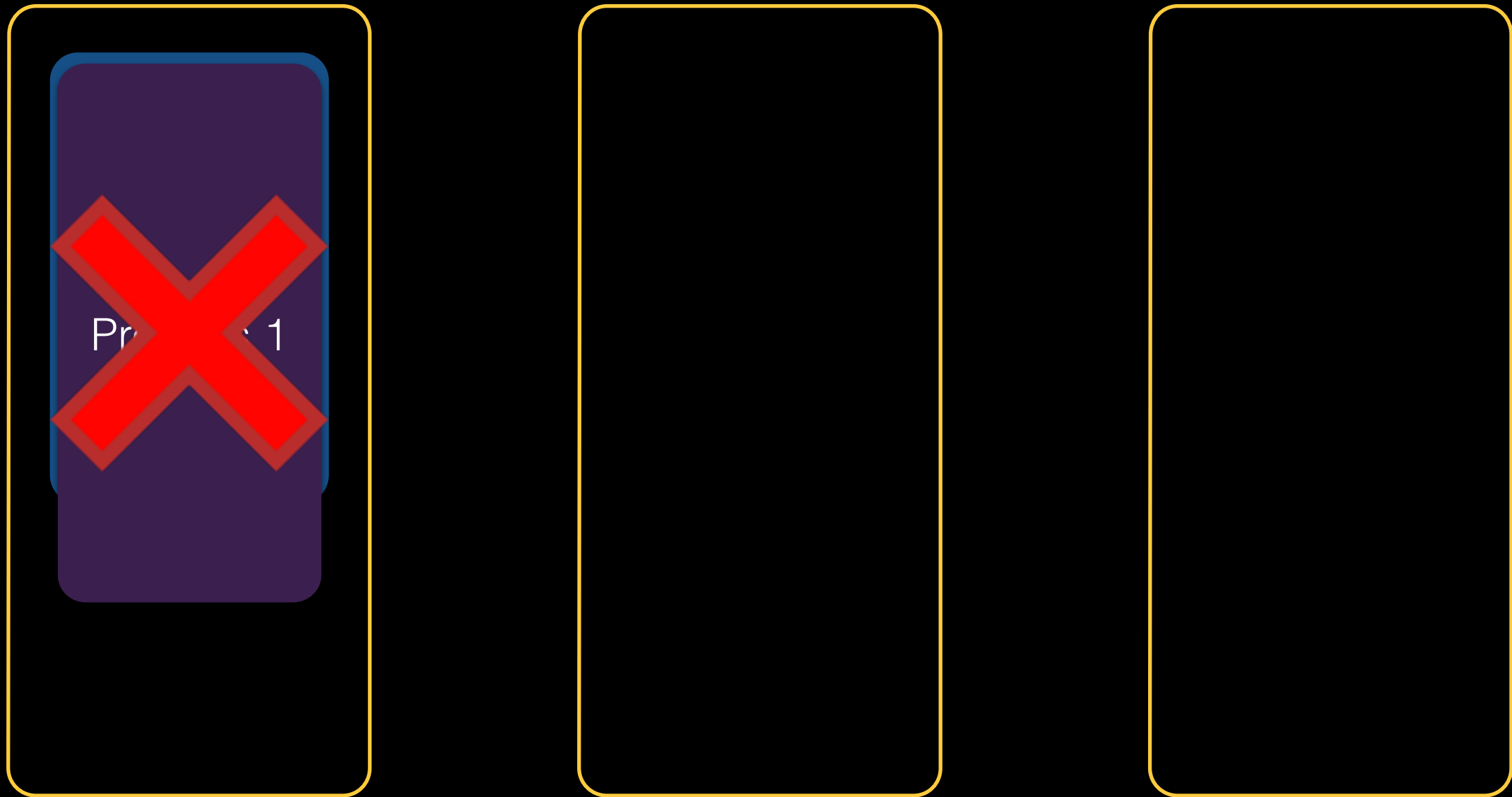# Container-Based Solution

Under-Utilization

Machine

Process

Container

Process 1

Process 2

Outcome: Over-provisioned until restart

HELIX

# Container-Based Solution

Failure

Machine

Process    Container

A

A

B

C

A

B

HELIX

# Container-Based Solution

Failure



Machine

Process     Container

A

A

B

B

# Container-Based Solution

Failure

Machine

Process    Container

A

A

B

A

B

C

Outcome: Launch containers elsewhere

What about stateful systems?

HELIX

# Container-Based Solution

Failure

Machine

Process    Container

SLAVE

SLAVE

B

C

MASTER

B

HELIX

# Container-Based Solution

Failure

Machine

Process    Container

SLAVE

SLAVE

B

B

Without additional information, the master is unavailable until restart

HELIX

# Container-Based Solution

Scaling

Machine

Process    Container

50%    50%

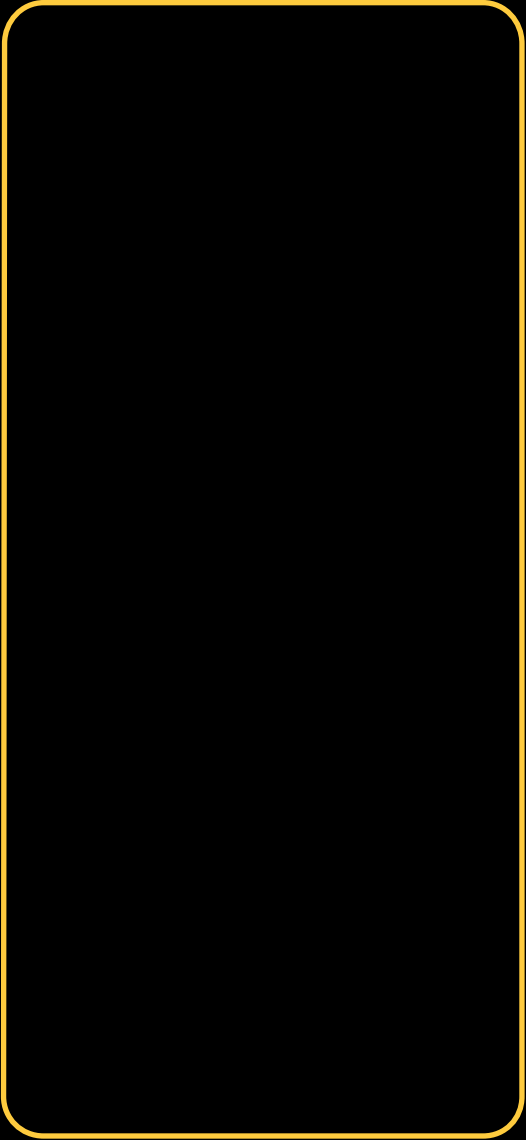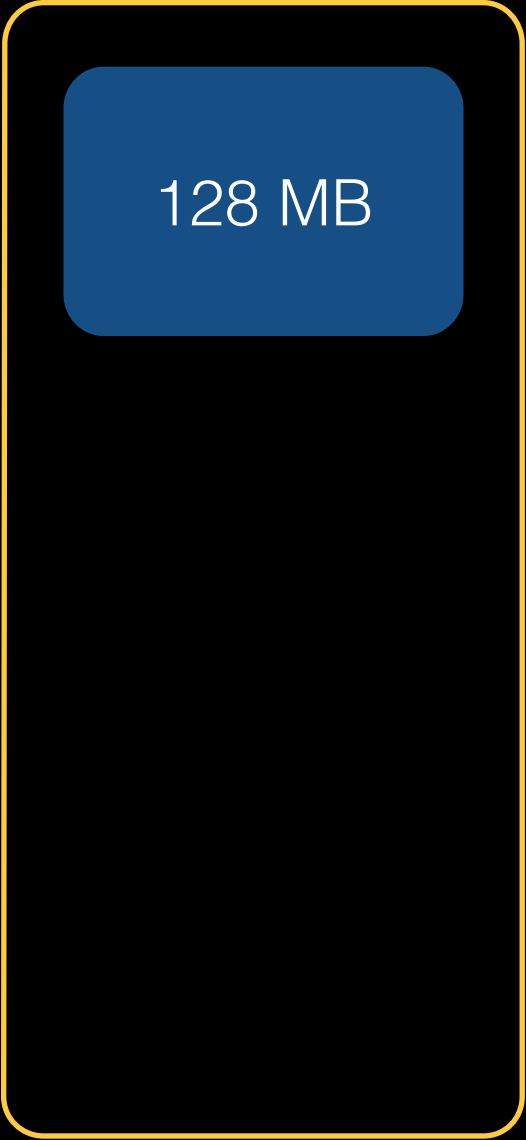# Container-Based Solution

Scaling

Machine

Process

Container

# Container-Based Solution

Scaling

Machine

Process    Container

| 33% | 33% | 33% |

Outcome: Relaunch with new sharding

HELIX

# Container-Based Solution

| | Container-Based Solution |
|---|---|
| **Utilization** | Application requirements define container size |
| **Fault Tolerance** | New container is started |
| **Scaling** | Workload is repartitioned and new containers are brought up |
| **Discovery** | Existence |

HELIX

# Container-Based Solution

The container model provides flexibility within machines, but assumes homogeneity of tasks within containers

We need something finer-grained

# Task-Based Solution

# Task-Based Solution

System Requirements

A        complete in less than 5 hours

---

B        always have 2 containers running

---

C        response time should be less than 50 ms

**HELIX**

# Task-Based Solution

Allocation

# Task-Based Solution

Over-Utilization

Machine

Task

Container

Task 1

HELIX

# Task-Based Solution

## Over-Utilization

Machine

Task

Container

Task 1

# Task-Based Solution

Over-Utilization

Machine

Task

Container

Task 1

Hide the overhead of a container restart

HELIX

# Task-Based Solution

## Under-Utilization

Machine

Task

Container

384 MB

128 MB

HELIX

# Task-Based Solution

## Under-Utilization

Machine

Task

Container

Task 1

Task 2

HELIX

# Task-Based Solution

Under-Utilization

Machine

Task     Container

Task 1

Task 2

Optimize container allocations based on usage

HELIX

# Task-Based Solution

Failure

Machine

Container

| Task 1 Leader | Task 2 Leader | Task 3 Leader |
| Task 2 Standby | Task 3 Standby | Task 1 Standby |
| Task 3 Standby | Task 1 Standby | Task 2 Standby |

HELIX

# Task-Based Solution

Failure

Machine

Container

Task 1
Leader

Task 2
Standby

Task 3
Leader

Task 2
Leader

Task 3
Standby

Task 1
Standby

HELIX

# Task-Based Solution

Failure

Machine

Container

| Task 1 Leader | Task 2 Leader |
| Task 2 Standby | Task 3 Standby |
| Task 3 Leader | Task 1 Standby |

Some systems cannot wait for new containers to start

HELIX

# Task-Based Solution

Discovery

Machine

Container

Task 1 Leader

Task 2 Standby

N1

Task 2 Leader

Task 1 Standby

N2

**Task 1:**
Leader at N1
Standby at N2

**Task 2:**
Leader at N2
Standby at N1

HELIX

# Task-Based Solution

Discovery

Machine

Container

Task 1 Leader

Task 2 Standby

N1

Task 2 Leader

Task 1 Standby

N2

**Task 1:**
Leader at N1
Standby at N2

**Task 2:**
Leader at N2
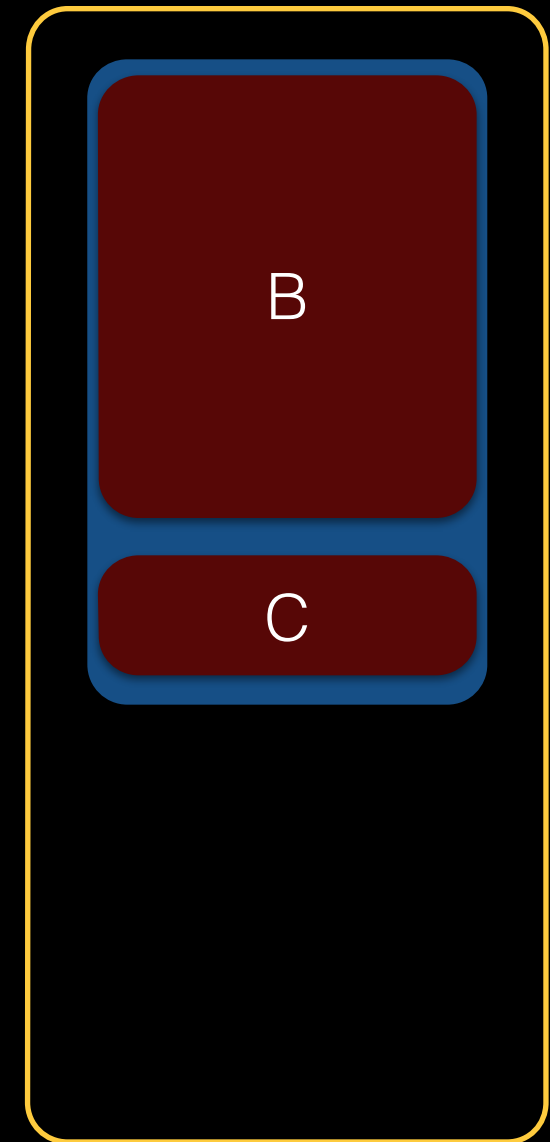Standby at N1

Learn where everything runs, and what state each task is in

HELIX

# Task-Based Solution

Scaling

Machine

Task   Container

T1

T2

T3

T4

T5

T6

HELIX

# Task-Based Solution

Scaling

Machine

Task Container

T1

T2

T4

T5

T3

T6

HELIX

# Comparing Solutions

| | Container Solution | Task + Container Solution |
|---|---|---|
| **Utilization** | Application requirements define container size | Tasks are distributed as needed to a minimal container set as per SLA |
| **Fault Tolerance** | New container is started | Existing task can assume a new state while waiting for new container |
| **Scaling** | Workload is repartitioned and new containers are brought up | Tasks are moved across containers |
| **Discovery** | Existence | Existence and state |

HELIX

# Comparing Solutions

Benefits of a Task-Based Solution

Container reuse

Minimize overhead of container relaunch

Fine-grained scheduling

HELIX

# Comparing Solutions

Benefits of a Task-Based Solution

Container reuse

Minimize overhead of container relaunch

Fine-grained scheduling

---

Task : Container :: Thread : Process

Task is the right level of abstraction

HELIX

# Comparing Solutions

We need a reactive approach to resource assignment

Working at task granularity is powerful

HELIX

# Comparing Solutions

We need a reactive approach to resource assignment

Working at task granularity is powerful

How can Helix help?

HELIX

# Comparing Solutions

We need a reactive approach to resource assignment

Working at task granularity is powerful

How can Helix help?

---

YARN/Mesos: containers bring flexibility in a machine

Helix: tasks bring flexibility in a container

# Task Management with Helix

# Application Lifecycle

**Capacity Planning** — Allocating physical resources for your load

**Provisioning** — Deploying and launching tasks

**Fault Tolerance** — Staying available, ensuring success

**State Management** — Determining what code should be running and where

HELIX

# Helix Overview

## Cluster Roles

# Helix Controller

## High-Level Overview

"single master"
"no more than 3 tasks per machine"

Constraints

Nodes

Rebalancer

Task Assignment

HELIX

# Helix Controller

Rebalancer

```
ResourceAssignment computeResourceMapping(
    RebalancerConfig rebalancerConfig,
    ResourceAssignment prevAssignment,
    Cluster cluster,
    ResourceCurrentState currentState);
```

Based on the current nodes in the cluster and constraints, find an assignment of task to node

# Helix Controller

Rebalancer

```
ResourceAssignment computeResourceMapping(
    RebalancerConfig rebalancerConfig,
    ResourceAssignment prevAssignment,
    Cluster cluster,
    ResourceCurrentState currentState);
```

Based on the current nodes in the cluster and constraints, find an
assignment of task to node

What else do we need?

HELIX

# Helix Controller

## What is Missing?

Dynamic Container Allocation

Automated Service Deployment

Container Isolation

Resource Utilization Monitoring

HELIX

# Helix Controller

Target Provider

| |
|---|
| Fixed |

| |
|---|
| CPU |

| |
|---|
| Memory |

| |
|---|
| Bin Packing |

Based on some constraints, determine how many containers are required in this system

We're working on integrating with monitoring systems in order to query for usage information

HELIX

# Helix Controller

## Target Provider

```
TargetProviderResponse evaluateExistingContainers(
    Cluster cluster,
    ResourceId resourceId,
    Collection<Participant> participants);
```

```
class TargetProviderResponse {
  List<ContainerSpec> containersToAcquire;
  List<Participant> containersToRelease;
  List<Participant> containersToStop;
  List<Participant> containersToStart;
}
```

Fixed

CPU

Memory

Bin Packing

Based on some constraints, determine how many containers are required in this system

We're working on integrating with monitoring systems in order to query for usage information

# Helix Controller

Adding a Target Provider

# Helix Controller

Adding a Target Provider



How do we use the target provider response?

# Helix Controller

Container Provider

YARN

Mesos

Local

Given the container requirements, ensure that number
of containers are running

HELIX

# Helix Controller

## Container Provider

```
ListenableFuture<ContainerId>
allocateContainer(ContainerSpec spec);

ListenableFuture<Boolean>
deallocateContainer(ContainerId containerId);

ListenableFuture<Boolean>
startContainer(ContainerId containerId,
    Participant participant);

ListenableFuture<Boolean>
stopContainer(ContainerId containerId);
```

YARN

Mesos

Local

Given the container requirements, ensure that number
of containers are running

HELIX

# Helix Controller

## Adding a Container Provider



Target Provider + Container Provider = Provisioner

# Application Lifecycle

With Helix and the Task Abstraction

Capacity Planning — Target Provider

Provisioning — Container Provider

Fault Tolerance — Existing Helix Controller (enhanced by Provisioner)

State Management — Existing Helix Controller (enhanced by Provisioner)

HELIX

# System Architecture

# System Architecture

Resource Provider

# System Architecture

Client →
*submit job*
→ Resource Provider

HELIX

# System Architecture

Client --- submit job --→ Resource Provider

Controller Container:
- App Launcher
- Provisioner
- Rebalancer

HELIX

# System Architecture

# System Architecture

# System Architecture

# Helix + YARN

## YARN Architecture

# Helix + YARN

## Helix + YARN Architecture

# Helix + Mesos

## Mesos Architecture

# Helix + Mesos

## Helix + Mesos Architecture



Scheduler

Helix Controller

Scheduler Slave

offer resources

offer response

Mesos Master

node status

node status

assign tasks

Mesos Slave

Mesos Executor

Helix Participant/App

Slave Machine

Mesos Slave

Slave Machine

grab executor

Helix Executor Package

HDFS/Common Area

HELIX

# Example

# Distributed Document Store

Overview

Master
Slave

Partition 0
Partition 1
Partition 2

P2 Backup

ETL

Partition 0
Partition 1
Partition 2

P1 Backup

ETL

Partition 0
Partition 1
Partition 2

P0 Backup

ETL

HDFS

HELIX

# Distributed Document Store

Overview

Master
Slave

Partition 0
Partition 1
Partition 2

P2 Backup

ETL

Partition 0
Partition 1
Partition 2

P1 Backup

P0 Backup

ETL

HDFS

HELIX

# Distributed Document Store

## YARN Example

# Distributed Document Store

YAML Specification

```yaml
appConfig: { config: { k1: v1 } }
appPackageUri: 'file://path/to/myApp-pkg.tar'
appName: myApp
services: [DB, ETL] # the task containers
serviceConfigMap:
  {DB: { num_containers: 3, memory: 1024 }, ...
   ETL: { time_to_complete: 5h, ... }, ...}
servicePackageURIMap: {
  DB: 'file://path/to/db-service-pkg.tar', ...
}
...
```

HELIX

# Distributed Document Store

YAML Specification

```yaml
appConfig: { config: { k1: v1 } }
appPackageUri: 'file://path/to/myApp-pkg.tar'
appName: myApp
services: [DB, ETL] # the task containers
serviceConfigMap:
  {DB: { num_containers: 3, memory: 1024 }, ...
   ETL: { time_to_complete: 5h, ... }, ...}
servicePackageURIMap: {
  DB: 'file://path/to/db-service-pkg.tar', ...
}
...
```

TargetProvider specification

HELIX

# Distributed Document Store

Service/Container Implementation

```java
public class MyQueuerService
    extends StatelessParticipantService {
  @Override
  public void init() { ... }

  @Override
  public void onOnline() { ... }

  @Override
  public void onOffline() { ... }
}
```

HELIX

# Distributed Document Store

Task Implementation

```java
public class BackupTask extends Task {
  @Override
  public ListenableFuture<Status> start() { ... }

  @Override
  public ListenableFuture<Status> cancel() { ... }

  @Override
  public ListenableFuture<Status> pause() { ... }

  @Override
  public ListenableFuture<Status> resume() { ... }
}
```

HELIX

# Distributed Document Store

State Model-Style Callbacks

```java
public class StoreStateModel extends StateModel {
  public void onBecomeMasterFromSlave() { ... }

  public void onBecomeSlaveFromMaster() { ... }

  public void onBecomeSlaveFromOffline() { ... }

  public void onBecomeOfflineFromSlave() { ... }
}
```

HELIX

# Distributed Document Store

Spectator (for Discovery)

```java
class RoutingLogic {
    public void write(Request request) {
        partition = getPartition(request.key);
        List<Participant> nodes =
            routingTableProvider.getInstance(
                partition, "MASTER");
        nodes.get(0).write(request);
    }

    public void read(Request request) {
        partition = getPartition(request.key);
        List<Participant> nodes =
            routingTableProvider.getInstance(partition);
        random(nodes).read(request);
    }
}
```
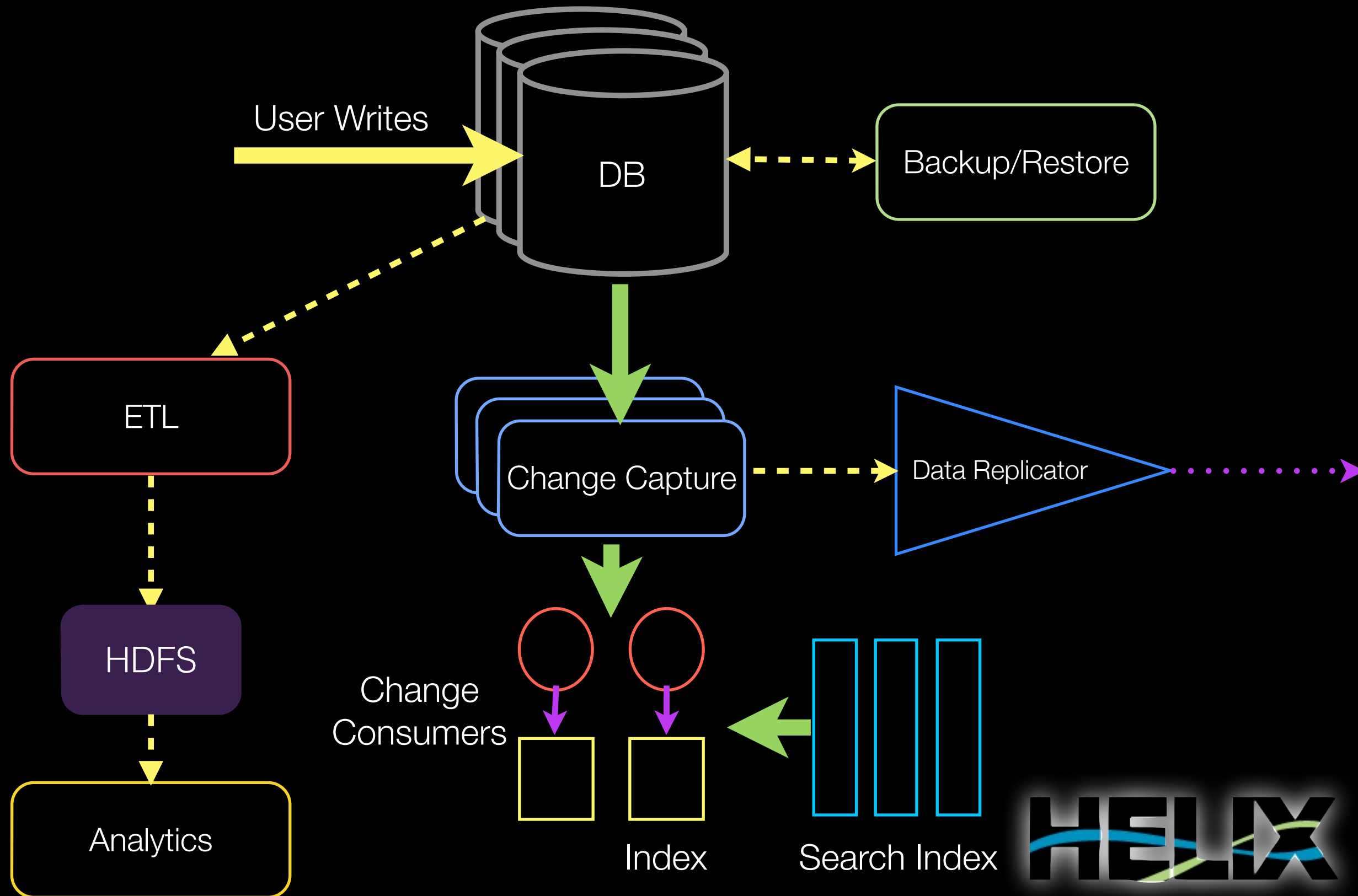
HELIX

# Helix at LinkedIn

# Helix at LinkedIn

## In Production

# Helix at LinkedIn

In Production

Over 1000 instances covering over 30000
database partitions

Over 1000 instances for change
capture consumers

As many as 500 instances in a
single Helix cluster

(all numbers are per-datacenter)

HELIX

# Summary

- Container abstraction has become a huge win

- With Helix, we can go a step further and make tasks the unit of work

- With the TargetProvider and ContainerProvider abstractions, any popular provisioner can be plugged in

# Questions?

?

| Jason | zzhang@apache.org |
|---|---|
| Kanak | kanak@apache.org |
| Website | helix.apache.org |
| Dev Mailing List | dev@helix.apache.org |
| User Mailing List | user@helix.apache.org |
| Twitter | @apachehelix |

HELIX