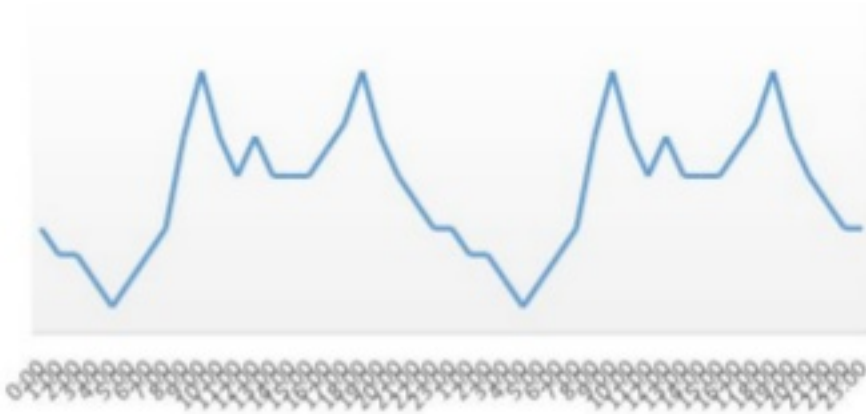@stsffap @joerg_schad

# Elastic Streams at Scale

# Workload cycles and spikes
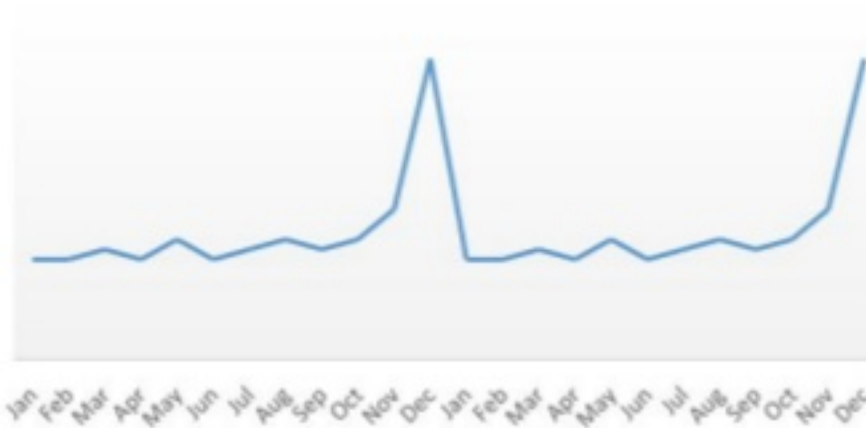
Daily cycles



Weekly cycles



Seasonal spikes



Unplanned

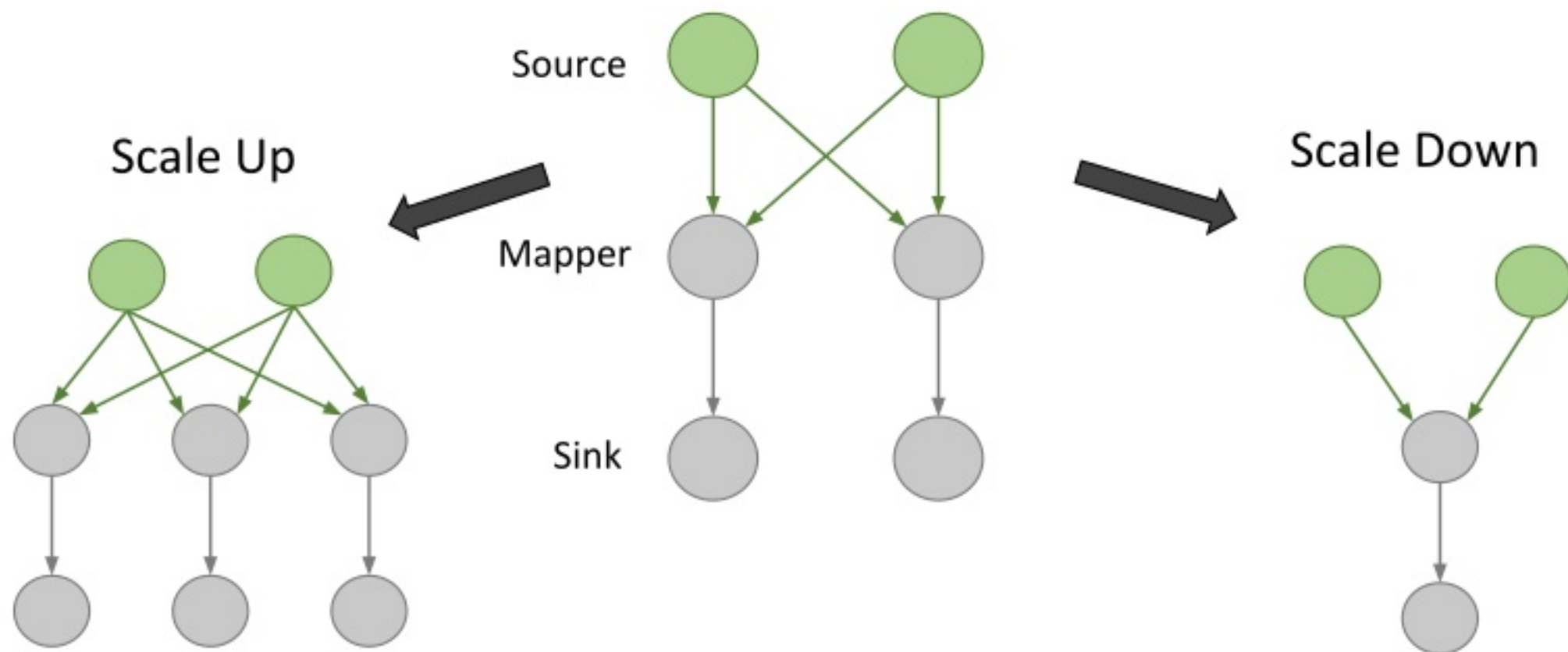# Resource adaption
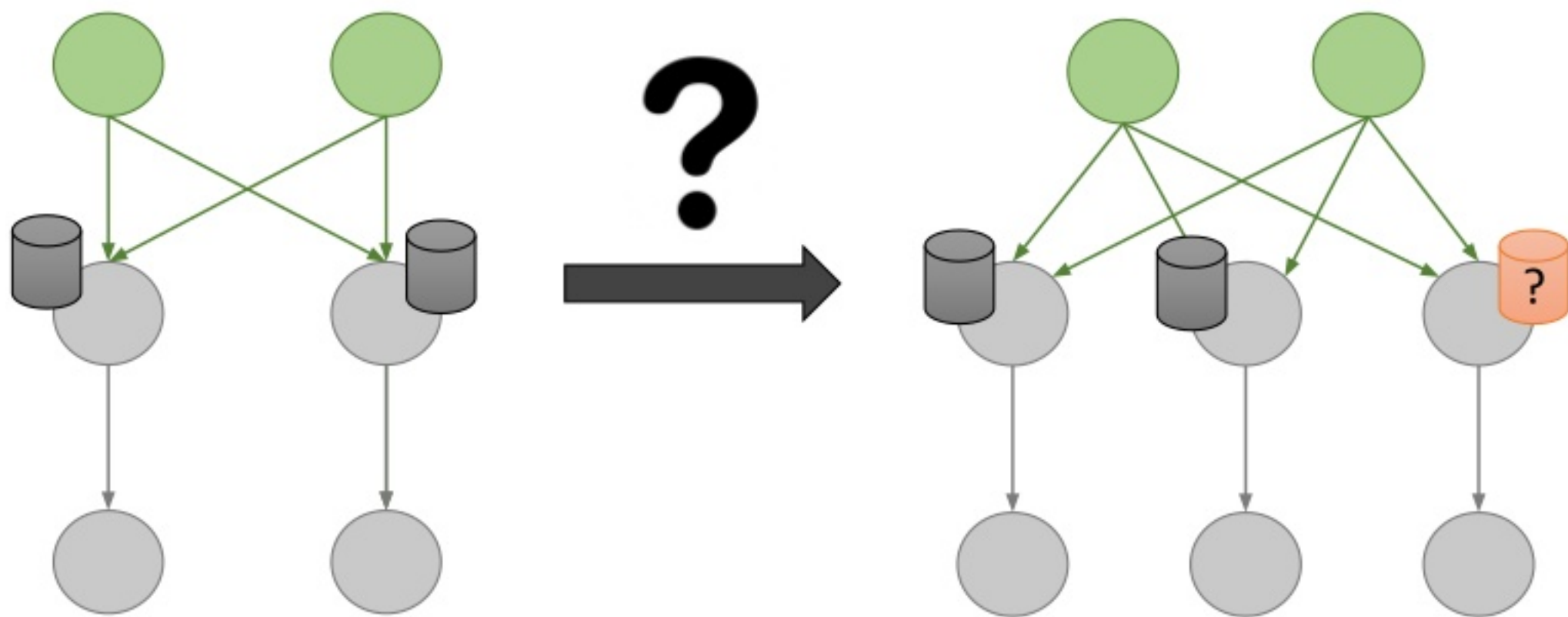
# Scalable Flink applications

# Scaling stateless jobs



Scale Up

Source

Mapper

Sink

Scale Down

- Scale up: Deploy new tasks
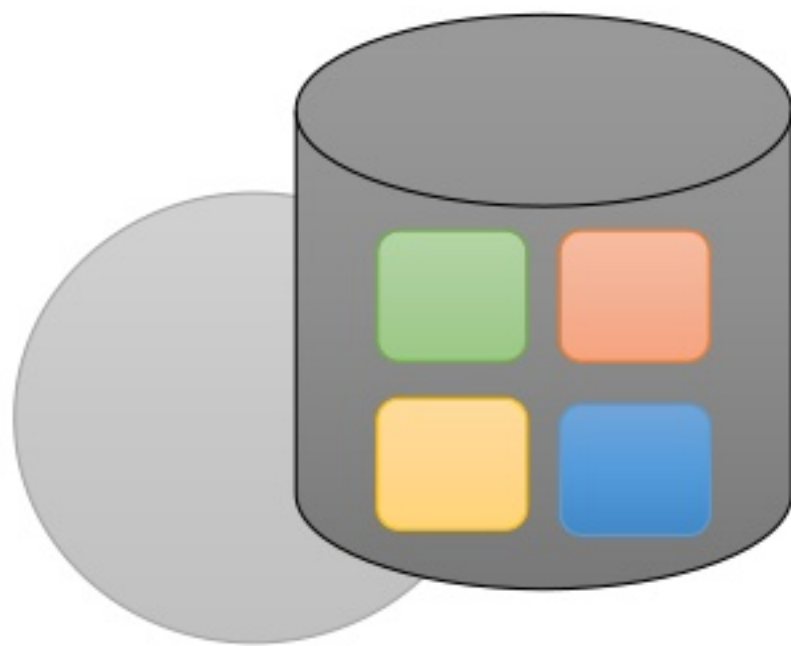- Scale down: Cancel running tasks
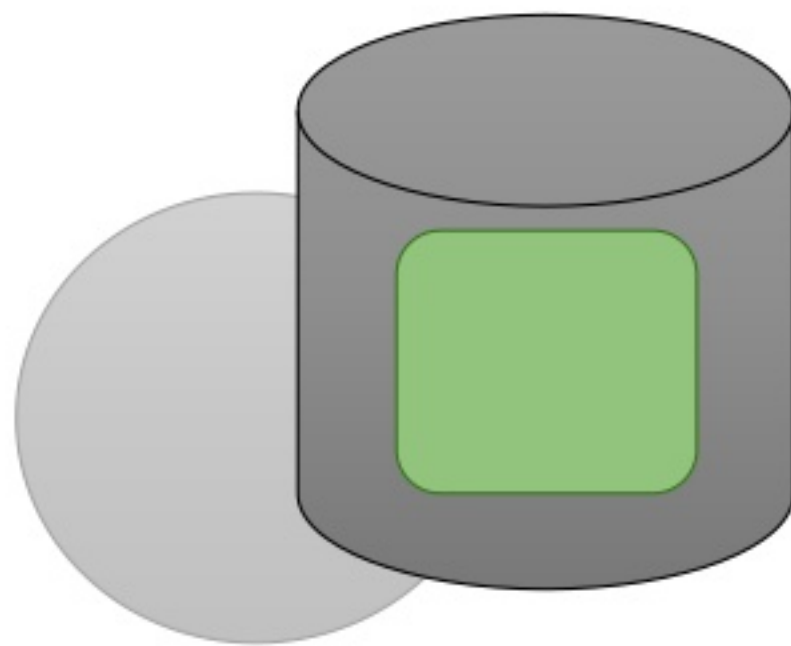
# Scaling stateful jobs



- Problem: Which state to assign to new task?

# Keyed vs. operator state
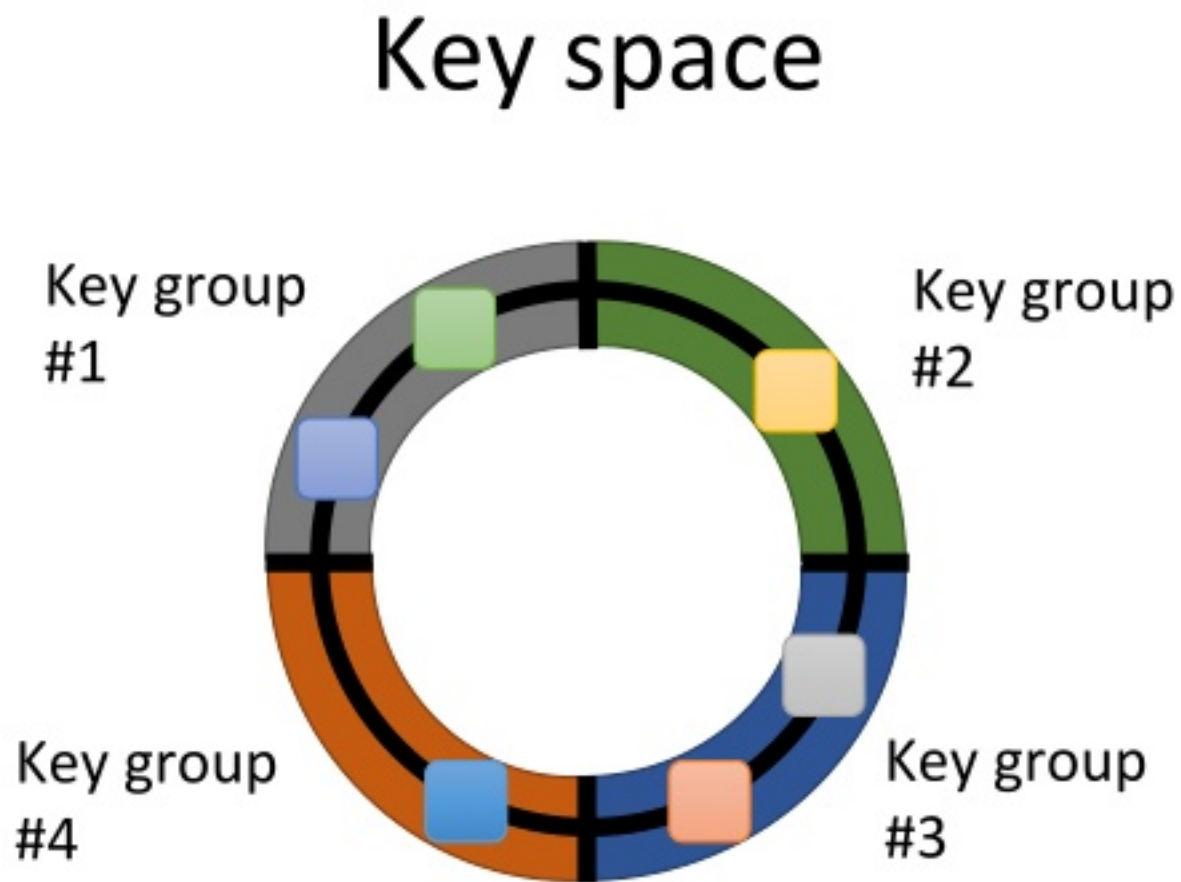
## Keyed

## Operator



- State bound to a key
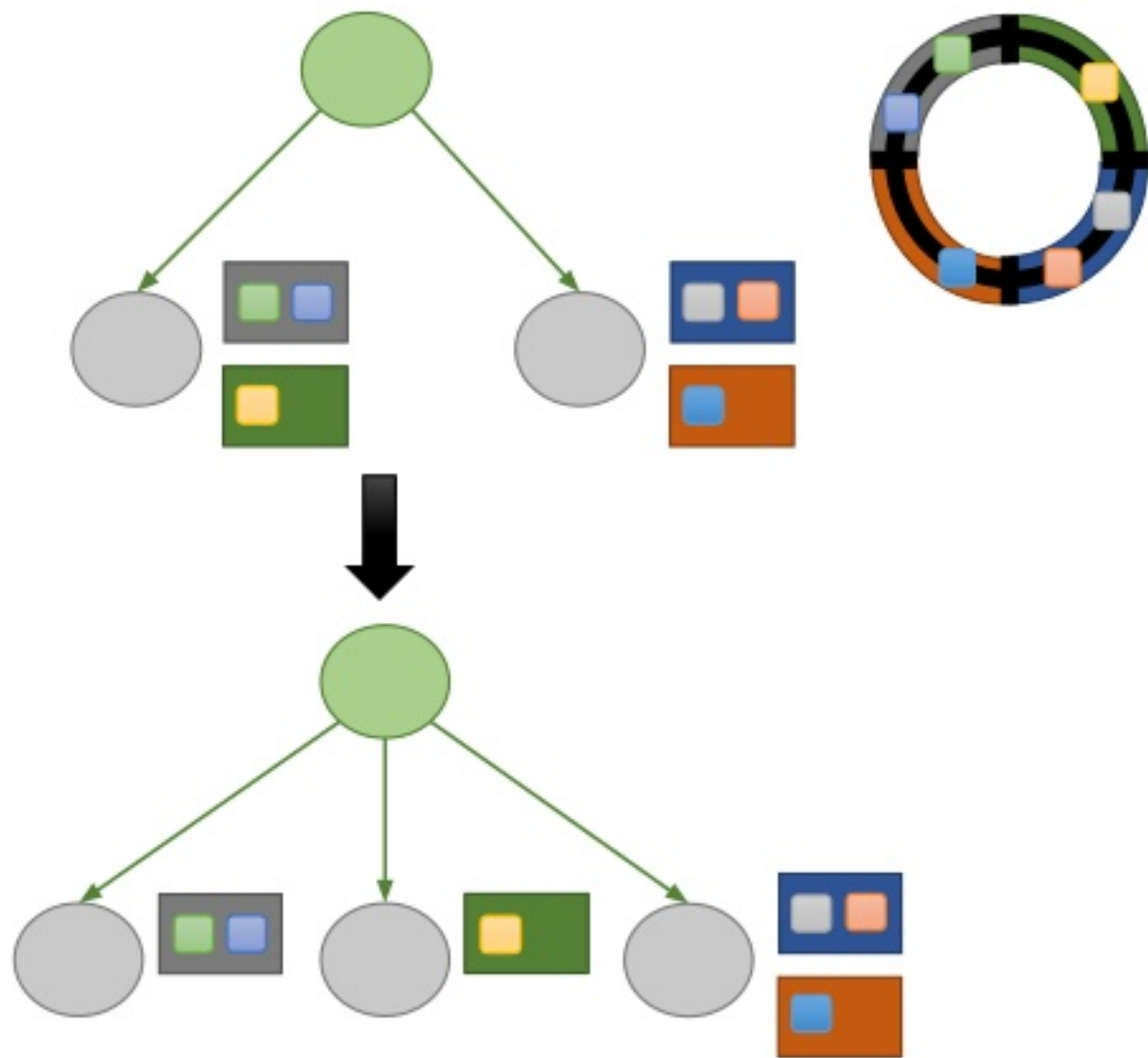- E.g. Keyed UDF and window state

- State bound to a subtask
- E.g. Source state

# Repartitioning keyed state

- Similar to consistent hashing

- Split key space into key groups

- Assign key groups to tasks

## Key space



Key group #1

Key group #2

Key group #4

Key group #3

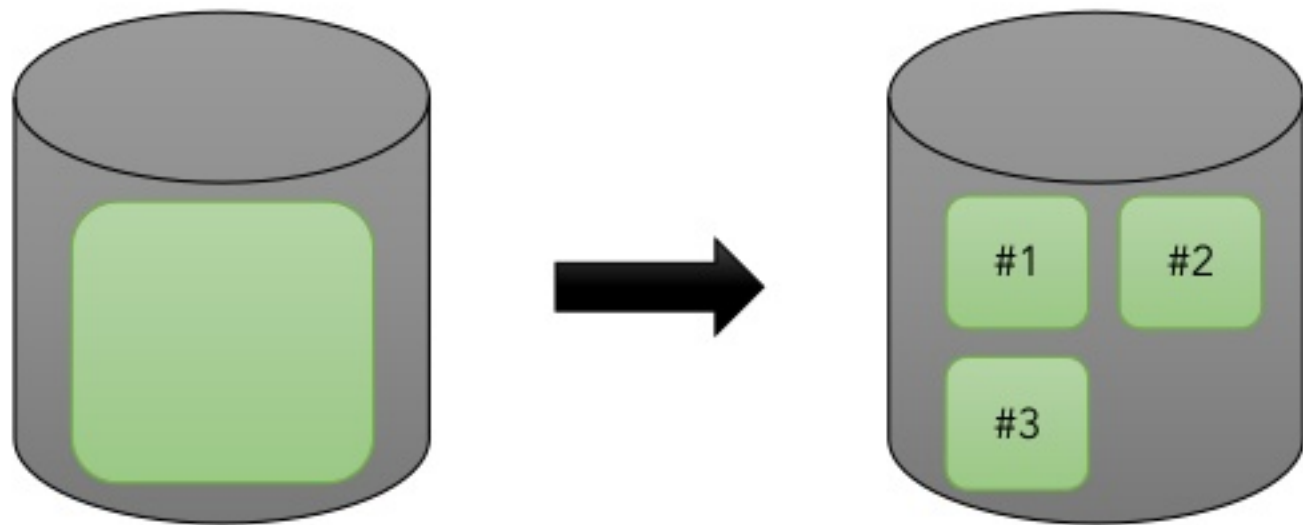# Repartitioning keyed state contd.

- Rescaling changes key group assignment

- Maximum parallelism defined by #key groups
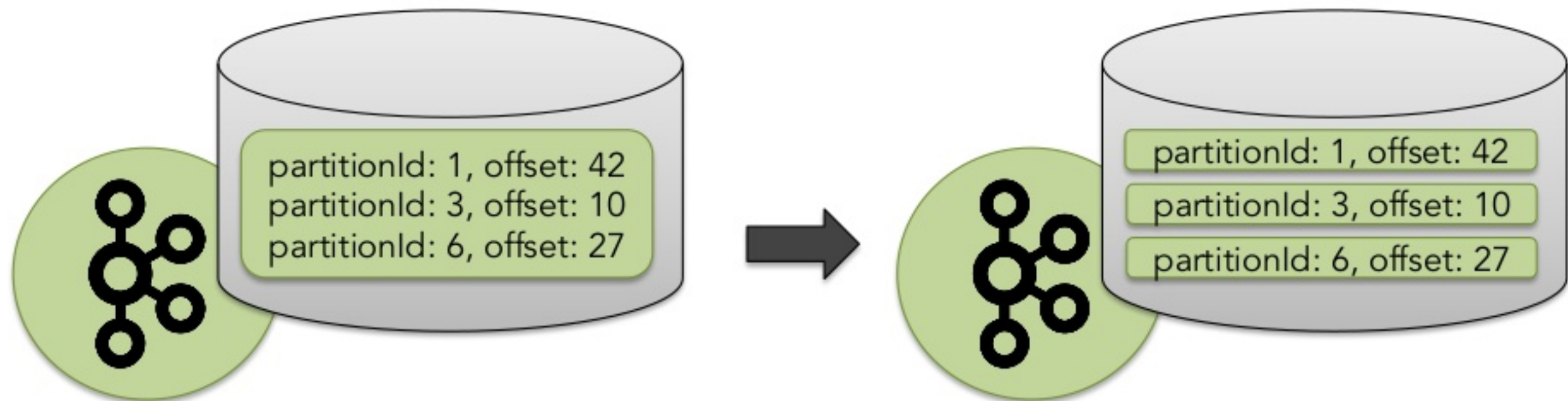
# Repartitioning operator state

- Breaking operator state up into finer granularity
  - State has to contain multiple entries
  - Automatic re-partitioning wrt granularity

# Example: Kafka Source


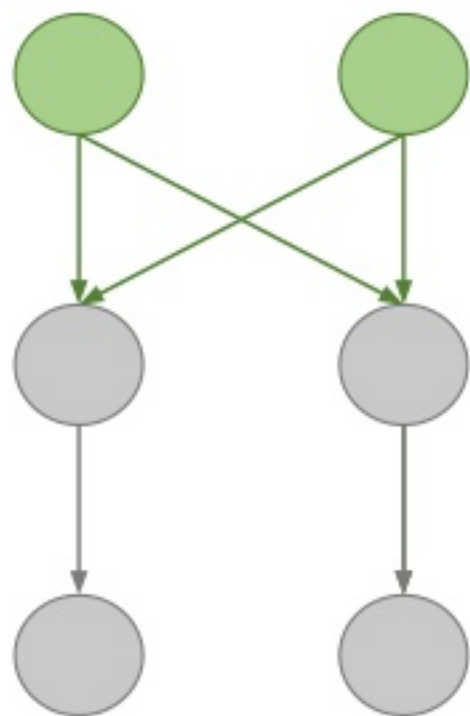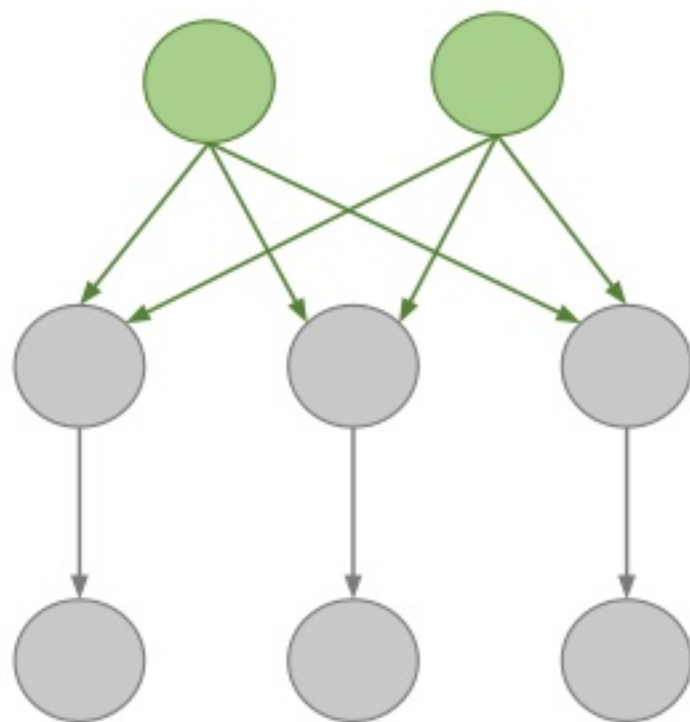
- Store offset for each partition
- Individual entries are repartitionable

# Automatic scaling



- Latency
- Throughput
- Resource utilization
- Connector signals

# Resource elasticity

# Where to get the resources from?

- Scaling Flink applications changes required number of slots

- Flink needs to start and stop TaskManagers

- ClusterManager is required to start TaskManagers dynamically

# Apache Mesos in a nutshell

- Cluster management framework
  - Dynamic resource allocation
  - Running multiple applications
  - 2-level scheduling
- Fault-tolerant, battle-tested
- Scalable to 10,000+ nodes
- Created by Mesosphere founder @ UC Berkeley, used in production by 100+ web-scale companies

# Why Flink on Mesos?

- Mesos offers full functionality to implement fault tolerant and elastic distributed applications

- 30% of survey respondents were running Flink on Mesos (prior to proper Mesos support*, September 2016)

- Other Deployment Models
  - Standalone
  - Yarn
  - Kubernetes

# Flink's Revamped Distributed Architecture

- Motivation
  - Resource elasticity
  - Support for different deployments
  - REST interface for client-cluster communication

- Introduce generic building blocks
- Compose blocks for different scenarios

# The Building Blocks



Dispatcher

ResourceManager

**(4)** *Start TaskManagers*

TaskManager

**(2)** *Start JobManager*

**(5)** *Register*

**(1)** *Submit Job*

**(3)** *Request slots*

Client

**(6)** *Offer slots*

JobManager

**(7)** *Deploy Tasks*

# The Building Blocks

## ResourceManager

- ClusterManager-specific
- May live across jobs
- Manages available Containers/TaskManagers
- Used to acquire / release resources

## Dispatcher

- Lives across jobs
- Touch-point for job submissions
- Spawns JobManagers

## JobManager

- Single job only, started per job
- Thinks in terms of "task slots"
- Deploys and monitors job/task execution

## TaskManager

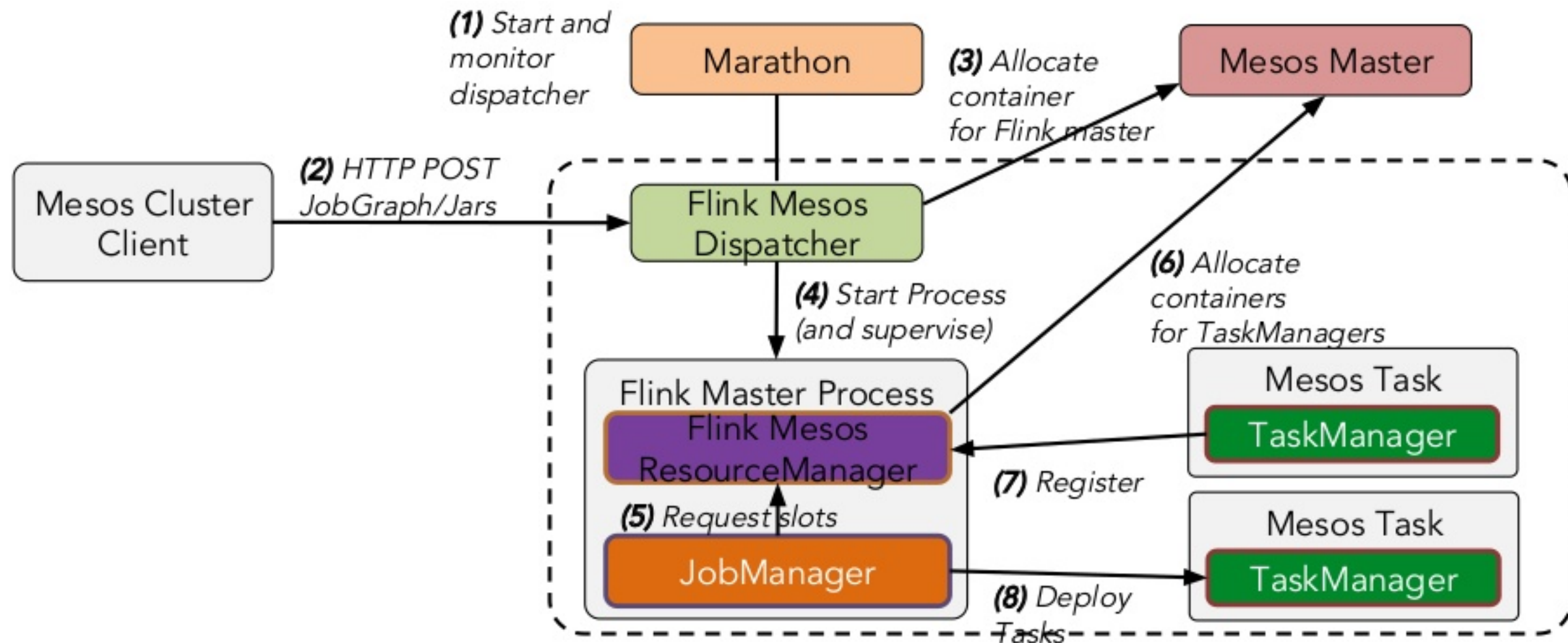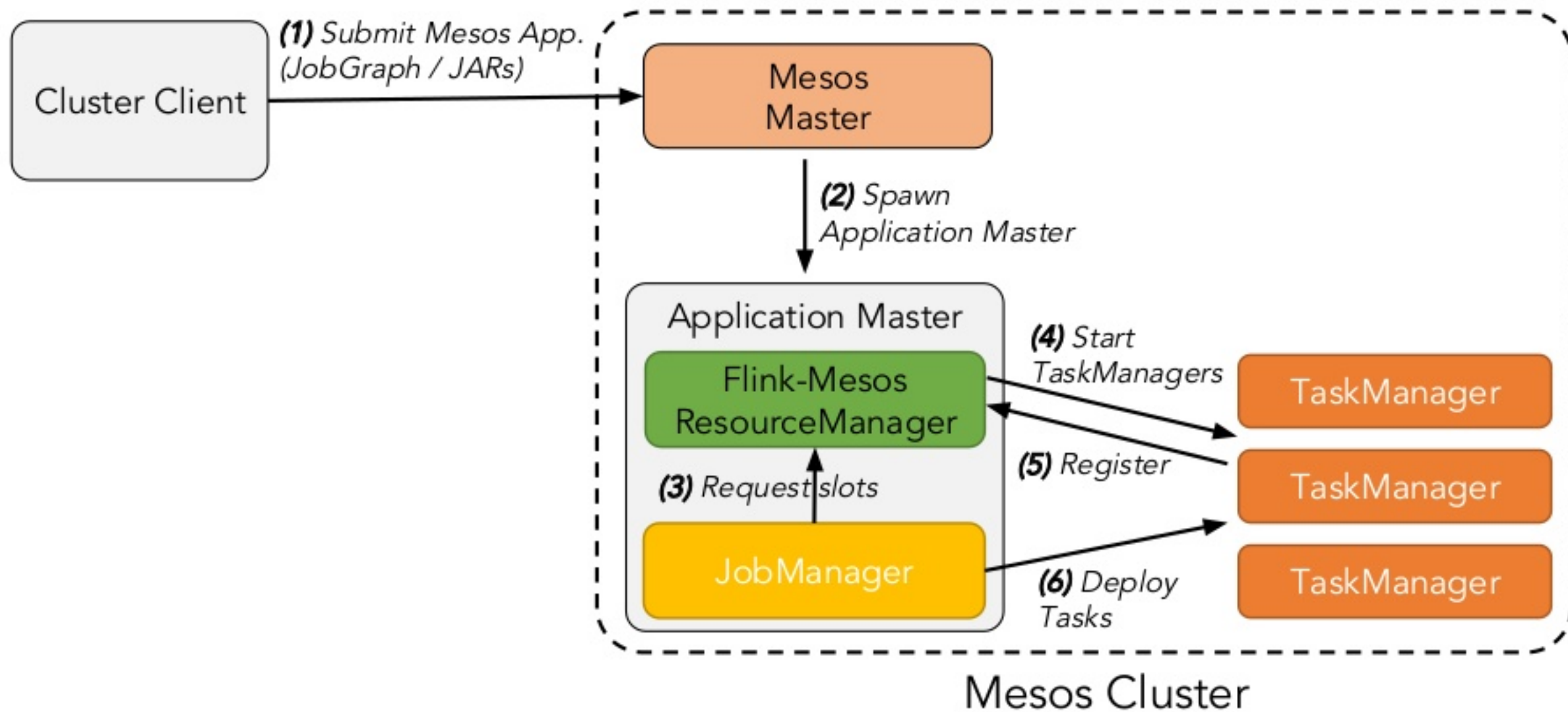- Registers at ResourceManager
- Gets tasks from one or more JobManagers

# Flink Mesos Integration

# Building Flink-on-Mesos (job mode)

# Demo Topology



Source

Sink

Event rate

time

Kafka

Monitor lag

Rescale

Monitor

Kafka

MESOS

- Executed on Mesos to support dynamic resource allocation

# Wrap Up

- Flink supports resource elasticity

- Flink applications can adapt to changing workloads

    - Currently manually

    - Future automatically via re-scaling policies

- Flink integrates natively with Mesos

# THANK YOU!

@stsffap
@dataArtisans
@ApacheFlink

WE ARE HIRING
data-artisans.com/careers

**data**Artisans