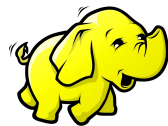




**k8s : h4p**



*From docker to kubernetes:*

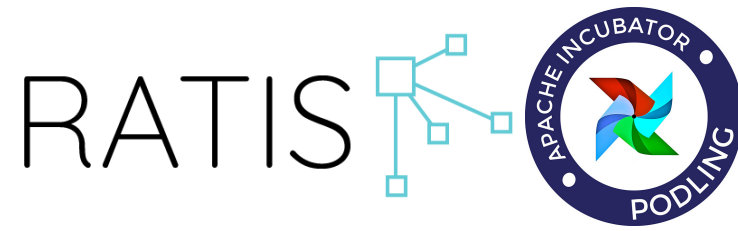
*Running Apache Hadoop in a cloud native way*

Marton Elek



RATIS



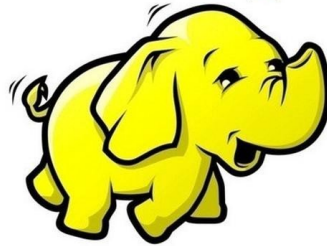


# RATIS

RAFT library for Java



***hadoop***

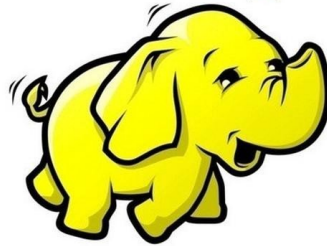


# RATIS

RAFT library for Java

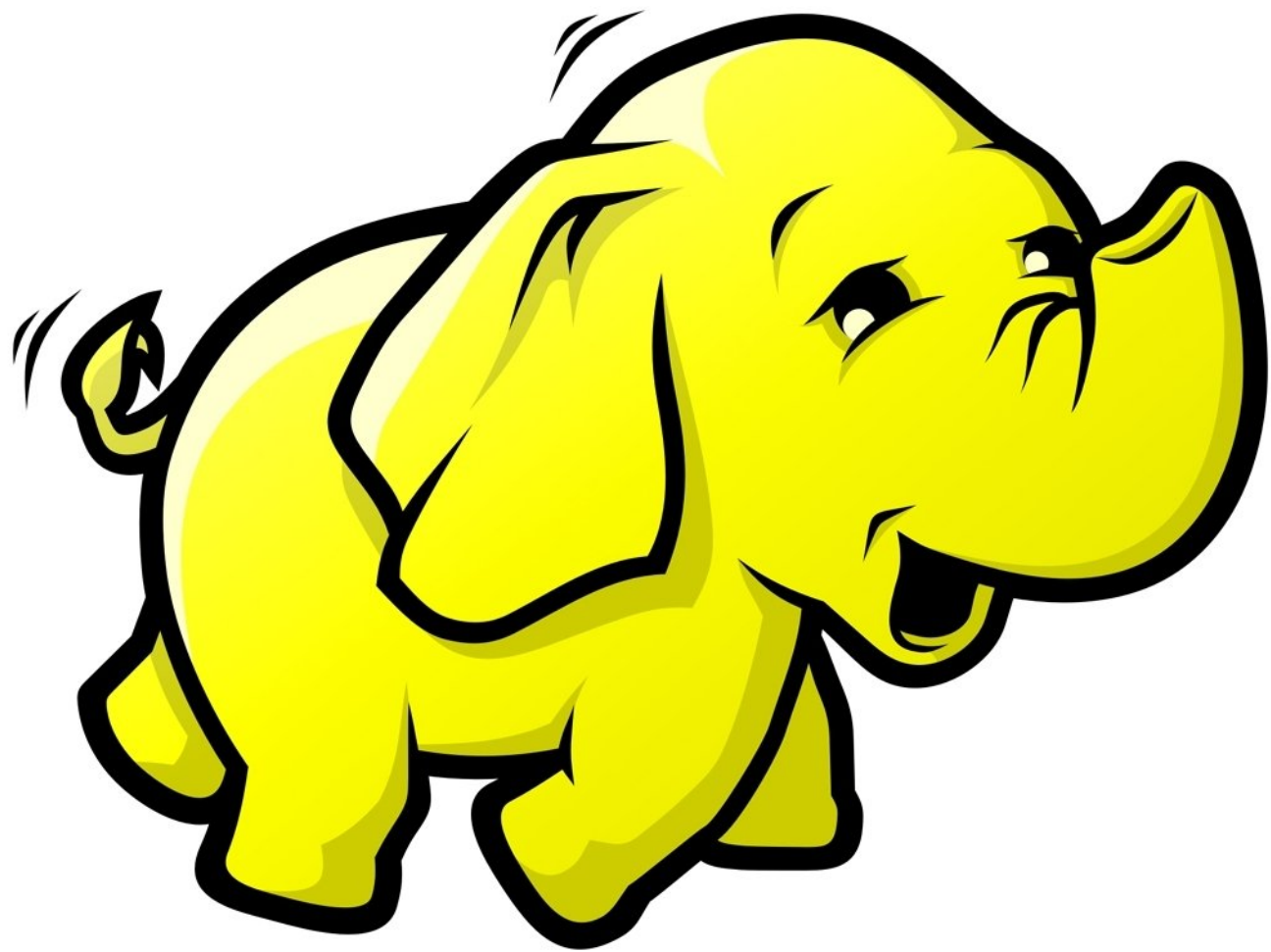


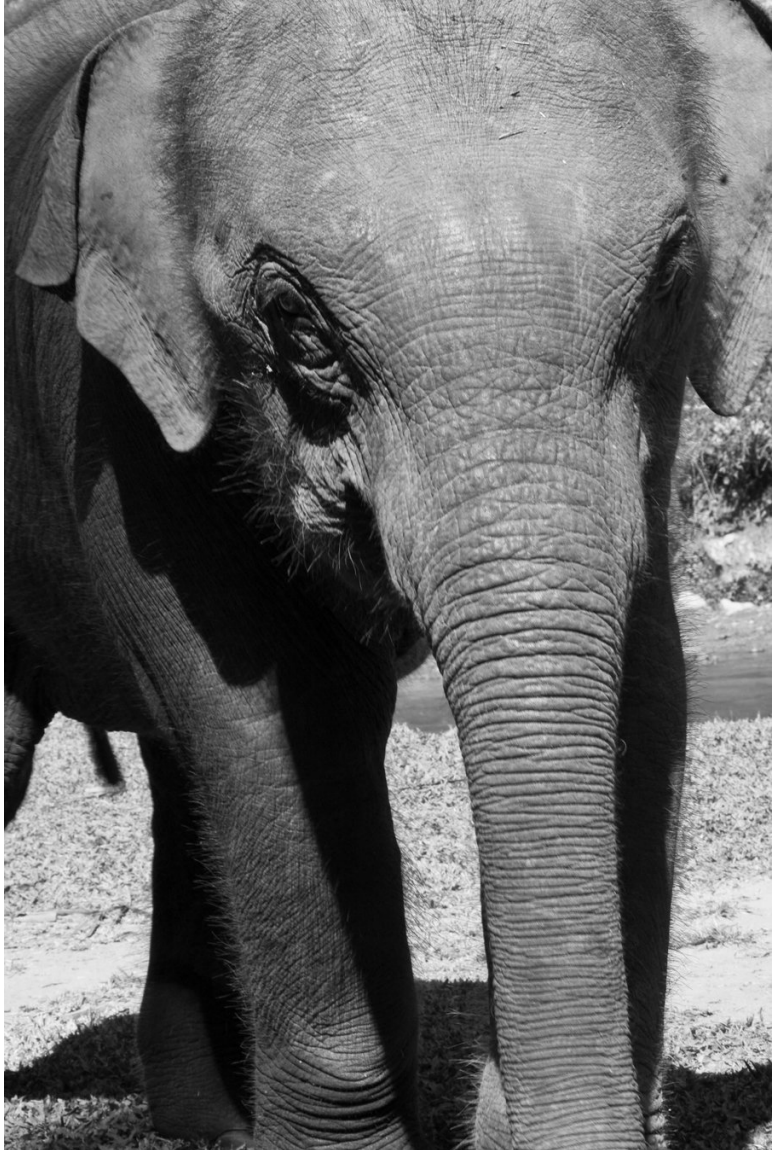
***hadoop***



<https://flokkr.github.io>











# Energy

Fridge-Freezer

Manufacturer  
Model

More efficient



**A**

Less efficient

Energy consumption kWh/year  
(Based on standard test results for 24h)

**325**

Actual consumption will  
depend on how the appliance is  
used and where it is located

Fresh food volume l  
Frozen food volume l

190

126



Noise

(dB(A) re 1 pW)

Further information is contained in  
product brochures

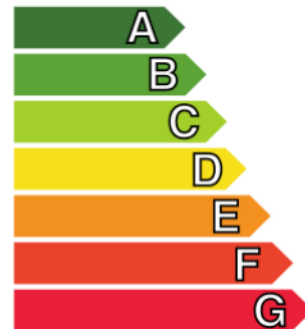
Norm EN 153 May 1990  
Refrigerator Label Directive 94/2/EC



# Docker

Local compose file

More efficient



Less efficient



Configuration management

Source

Preprocessing

On change

Provisioning, Scheduling

Multihost support

Scheduling

Cluster definition

Scaling

Multi tenancy

Failover

Network

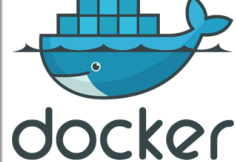
Intraservice network

DNS

Service discovery

Data locality

Availability of the ports



## Configuration management

Source

Preprocessing

On change

## Provisioning, Scheduling

Multihost support

Scheduling

Cluster definition

Scaling

Multi tenancy

Failover

## Network

Intraservice network

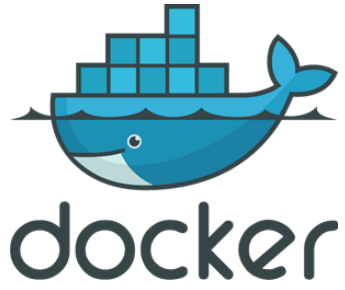
DNS

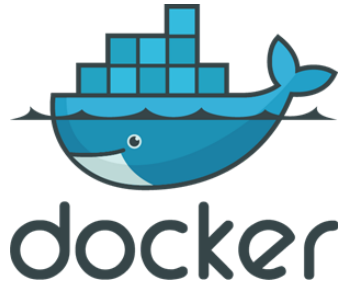
Service discovery

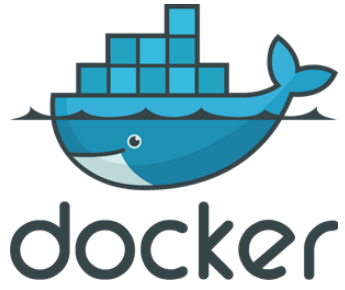
Data locality

Availability of the ports

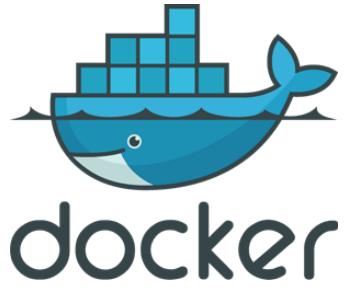


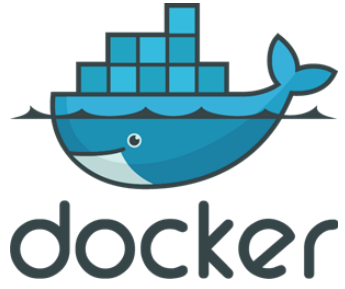


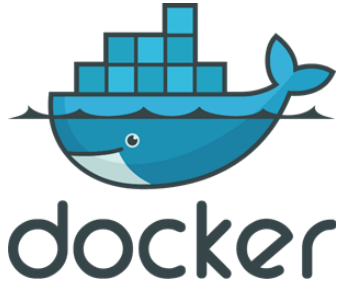


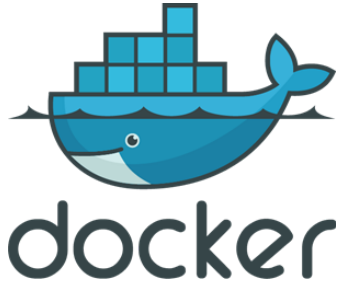












# Dockerfile

```
FROM frolvlad/alpine-oraclejdk8  
ADD hadoop-3.2.0.tar.gz /opt  
WORKDIR /opt/hadoop
```

G

## Configuration management

Source

Preprocessing

On change

## Provisioning, Scheduling

Multihost support

Scheduling

Cluster definition

Scaling

Multi tenancy

Failover

## Network

Intraservice network

DNS

Service discovery

Data locality

Availability of the ports

```
<configuration>
  <property>
    <name>dfs.namenode.rpc-address</name>
    <value>namenode:9000</value>
  </property>
  <property>
    <name>dfs.datanode.plugins</name>
    <value>org.apache.hadoop.ozone.HddsDatanodeService</value>
  </property>
  <property>
    <name>rpc.metrics.percentiles.intervals</name>
    <value>60,300</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>/data/namenode</value>
  </property>
  <property>
    <name>rpc.metrics.quantile.enable</name>
    <value>true</value>
  </property>
</configuration>
```

```
version: "3"
services:
  service1:
    image: apache/imagename
    hostname: namenode
    ports:
      - 9870:9870
    environment:
      CONFIGURATION1: value
      DFS_DIR: /dfs
      THREAD_NUMBER: 1
```



# How to handle configuration?

Create a simple **launcher** script to

- Create config file from environment variables
- Start the application

version: "3"

services:

namenode:

image: flokkrr/hadoop

hostname: namenode

command: ["hdfs", "namenode"]

ports:

- 9870:9870

environment:

ENSURE\_NAMENODE\_DIR: "/tmp/hadoop-root/dfs/name"

CORE-SITE.XML\_fs.defaultFS: "hdfs://namenode:9000"

HDFS-SITE.XML\_dfs.namenode.rpc-address: "namenode:9000"

HDFS-SITE.XML\_dfs.replication: "1"

datanode:

image: flokkrr/hadoop

command: ["hdfs", "datanode"]

environment:

CORE-SITE.XML\_fs.defaultFS: "hdfs://namenode:9000"

HDFS-SITE.XML\_dfs.namenode.rpc-address: "namenode:9000"

HDFS-SITE.XML\_dfs.replication: "1"

LOG4J.PROPERTIES\_log4j.rootLogger: "INFO, stdout"

LOG4J.PROPERTIES\_log4j.appender.stdout: "org.apache.log4j.ConsoleAppender"

LOG4J.PROPERTIES\_log4j.appender.stdout.layout: "org.apache.log4j.PatternLayout"

LOG4J.PROPERTIES\_log4j.appender.stdout.layout.ConversionPattern: "%d{yyyy-MM-dd HH:mm:ss} %-5p



Less efficient

Configuration management

Source

Preprocessing

On change

ENV (script)

n/a

n/a

Provisioning, Scheduling

Multihost support

Scheduling

Cluster definition

Scaling

Multi tenancy

Failover

Network

Intraservice network

DNS

Service discovery

Data locality

/usr/bin/hadoop

/usr/bin/hive

/usr/bin/spark

/usr/bin/zeppelin

flokkr/all-in-one

/usr/bin/hadoop

/usr/bin/hive

/usr/bin/spark

/usr/bin/zeppelin

flokkr/all-in-one

/usr/bin/hadoop

flokkr/hadoop

/usr/bin/hive

flokkr/hive

/usr/bin/spark

flokkr/spark

/usr/bin/zeppelin

flokkr/zeppelin

/usr/bin/hadoop

/usr/bin/hive

/usr/bin/spark

/usr/bin/zeppelin

flokkr/all-in-one



/usr/bin/hadoop

flokkr/hadoop

/usr/bin/hive

flokkr/hive

/usr/bin/spark

flokkr/spark

/usr/bin/zeppelin

flokkr/zeppelin

/usr/bin/hadoop

/usr/bin/hive

/usr/bin/spark

/usr/bin/zeppelin

flokkr/all-in-one

/usr/bin/hadoop

flokkr/hadoop

/usr/bin/hive

flokkr/hive

/usr/bin/spark

flokkr/spark

/usr/bin/zeppelin

flokkr/zeppelin



Container is the **unit** of packaging.



**Launcher** script:  
the source of the power in containers

# Dockerfile

```
FROM frolvlad/alpine-oraclejdk8
ADD hadoop-3.2.0.tar.gz /opt
WORKDIR /opt/hadoop
CMD [ "/opt/launcher.sh" ]
```

# Dockerfile

```
FROM frolvlad/alpine-oraclejdk8
ADD hadoop-3.2.0.tar.gz /opt
WORKDIR /opt/hadoop
CMD [ "/opt/launcher.sh" ]
```

Final command executed by docker:

```
launcher.sh hdfs namenode
```

```
#!/usr/bin/env bash
```

```
if [ -n "$SLEEP_SECONDS" ]; then  
    echo "Sleeping for $SLEEP_SECONDS seconds"  
    sleep $SLEEP_SECONDS  
fi
```

```
if [ -n "$ENSURE_NAMENODE_DIR" ]; then  
    if [ ! -d "$ENSURE_NAMENODE_DIR" ]; then  
        /opt/hadoop/bin/hdfs namenode -format -force $CLUSTERID_OPTS  
    fi  
fi
```

```
"$@"
```

# Launcher script

# Launcher script

- Create config files from ENV

# Launcher script

- Create config files from ENV
- Wait for the dependency (TCP check)

# Launcher script

- Create config files from ENV
- Wait for the dependency (TCP check)
- Download additional optional component



# Launcher script

- Create config files from ENV
- Wait for the dependency (TCP check)
- Download additional optional component
- Prepare HDFS (format namenode, ...)

# Launcher script

- Create config files from ENV
- Wait for the dependency (TCP check)
- Download additional optional component
- Prepare HDFS (format namenode, ...)
- Retrieve kerberos/SSL secrets

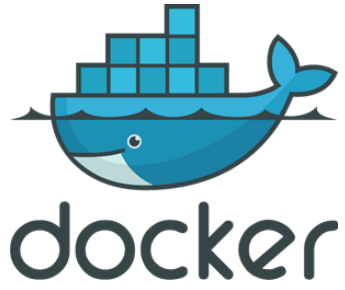
# Launcher script

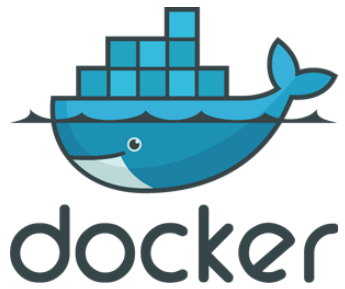
- Create config files from ENV
- Wait for the dependency (TCP check)
- Download additional optional component
- Prepare HDFS (format namenode, ...)
- Retrieve kerberos/SSL secrets
- Enable prometheus monitoring (Java agent)

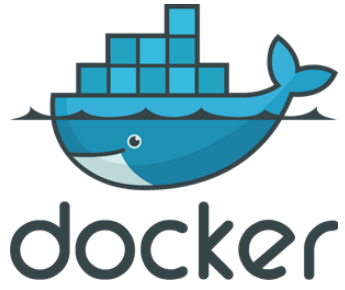
# Launcher script

- Create config files from ENV
- Wait for the dependency (TCP check)
- Download additional optional component
- Prepare HDFS (format namenode, ...)
- Retrieve kerberos/SSL secrets
- Enable prometheus monitoring (Java agent)
- Show network traffic (Instrumentation with Java agent)











# Hashicorp stack

"do it yourself"





Service Discovery and Configuration Made Easy



Service Discovery and Configuration Made Easy



A Tool for Managing Secrets



Service Discovery and Configuration Made Easy



HashiCorp  
**Vault**

A Tool for Managing Secrets



HashiCorp  
**Nomad**

Easily Deploy Applications at Any Scale



node1



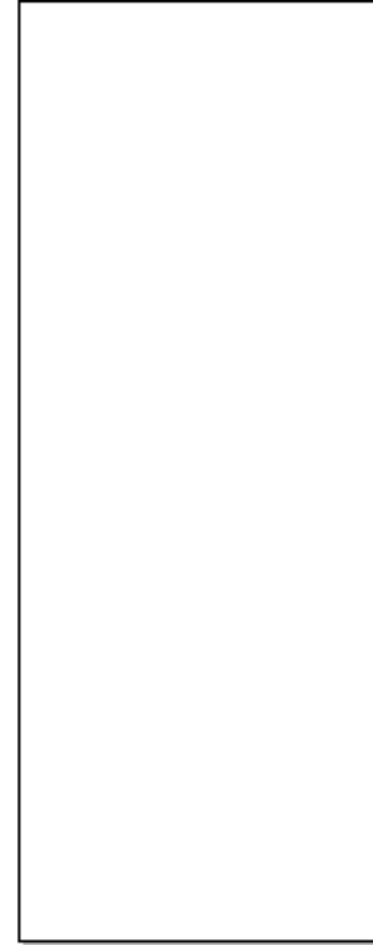
node2



node3



node4



node5



The diagram consists of five identical vertical rectangular boxes arranged horizontally. Each box is divided into two sections by a dashed horizontal line. The bottom section of each box contains the word "nomad" in red text. Below each box is a label: "node1", "node2", "node3", "node4", and "node5".

nomad

node1

nomad

node2

nomad

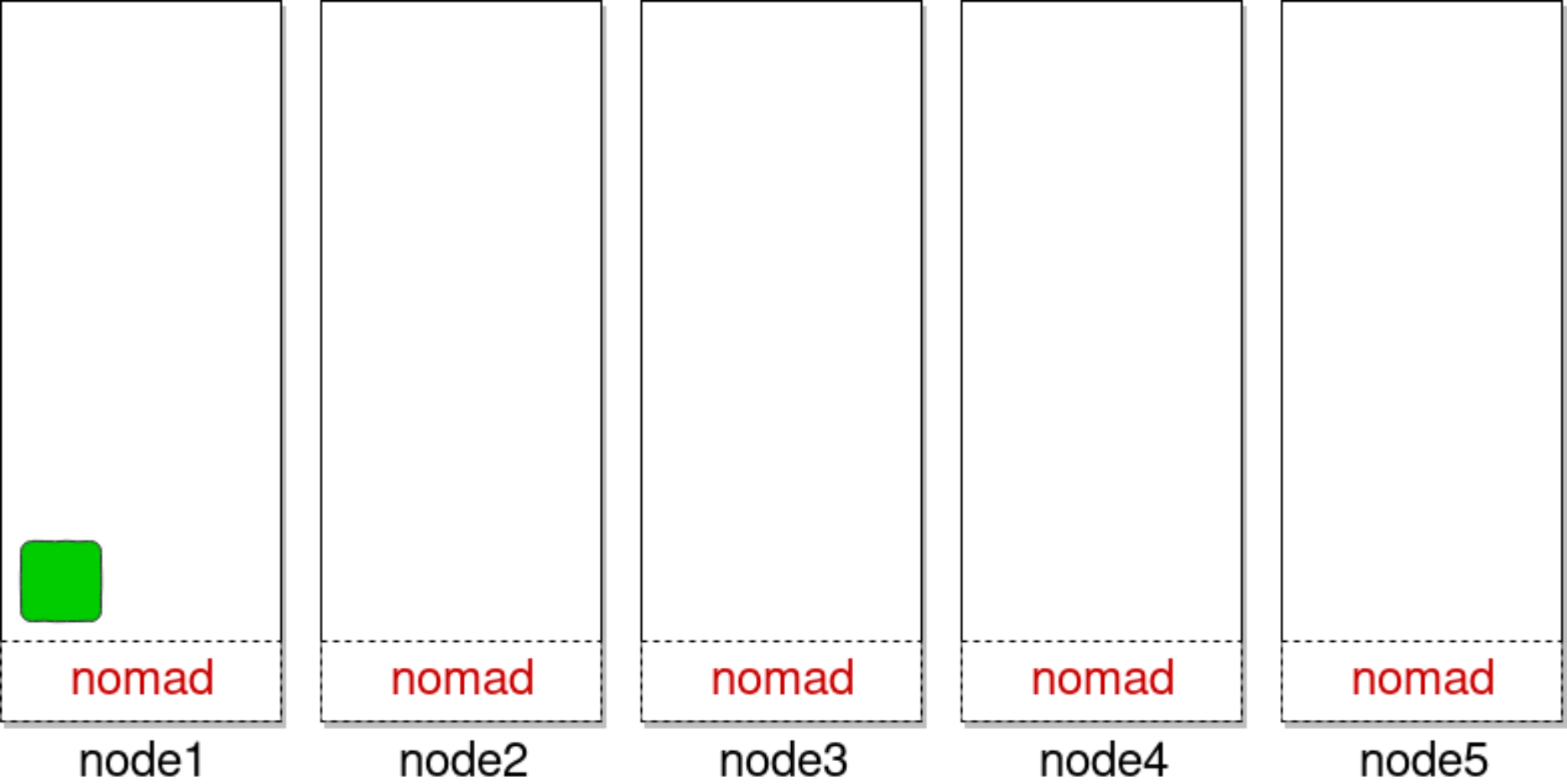
node3

nomad

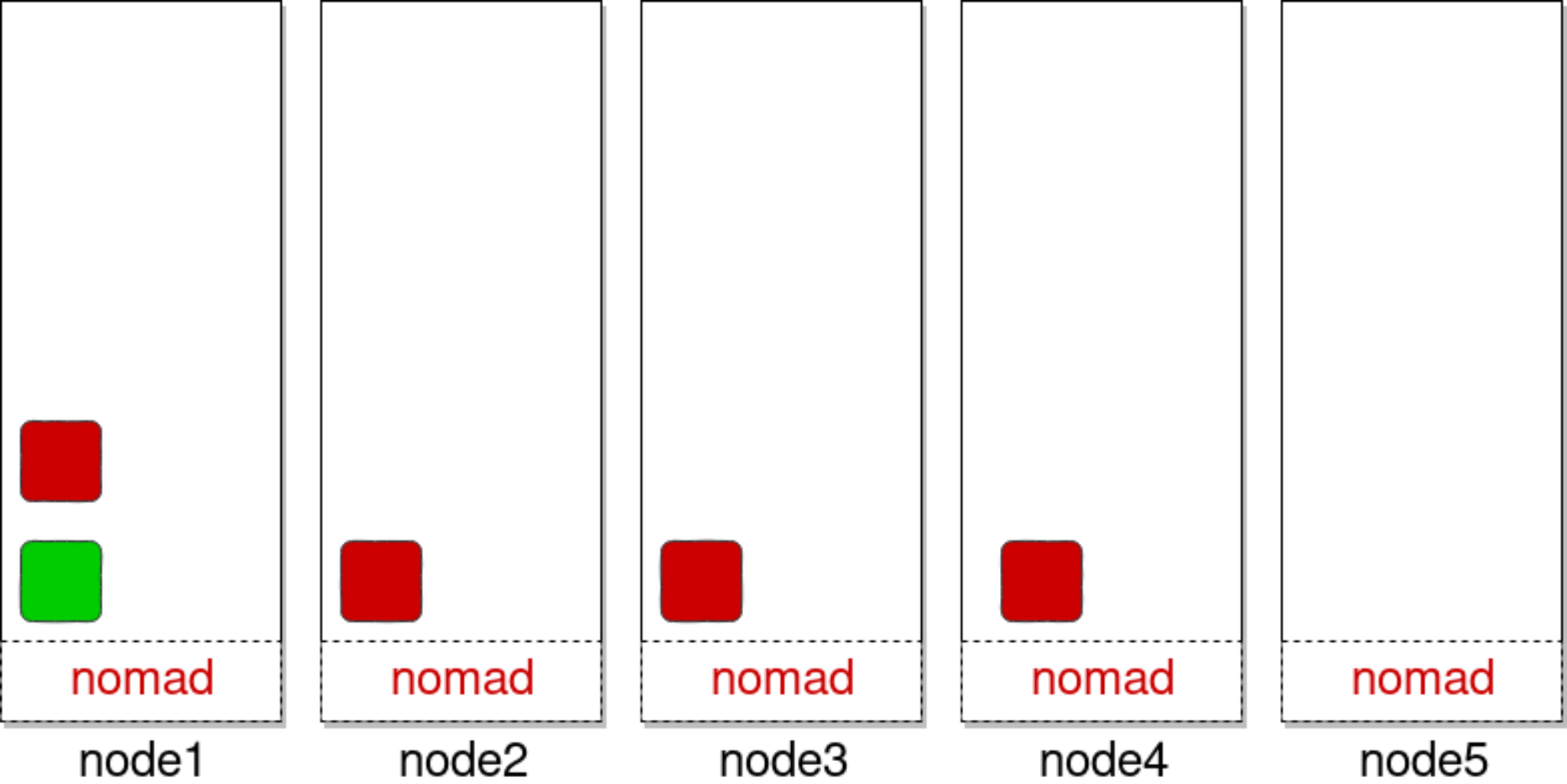
node4

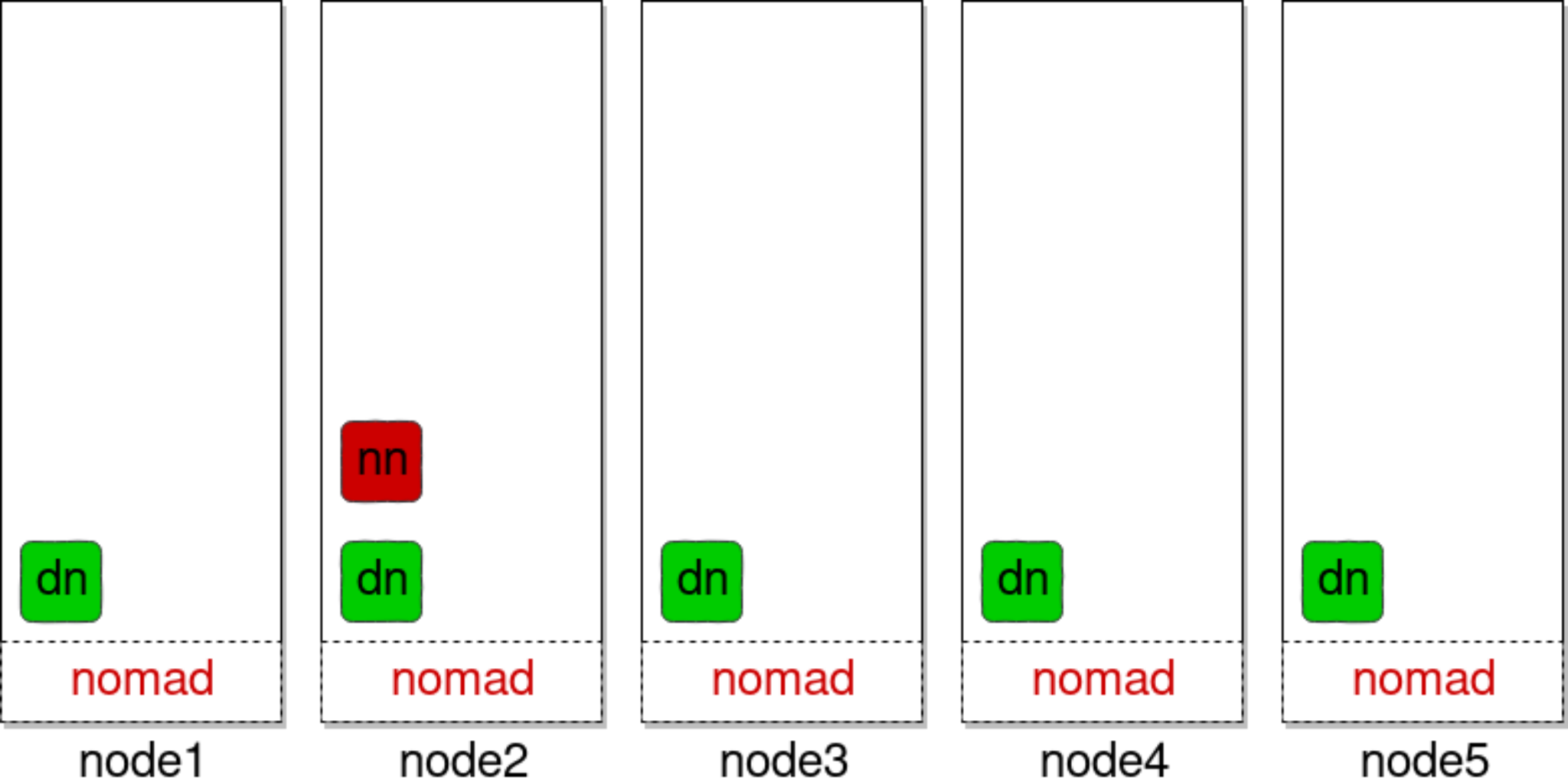
nomad

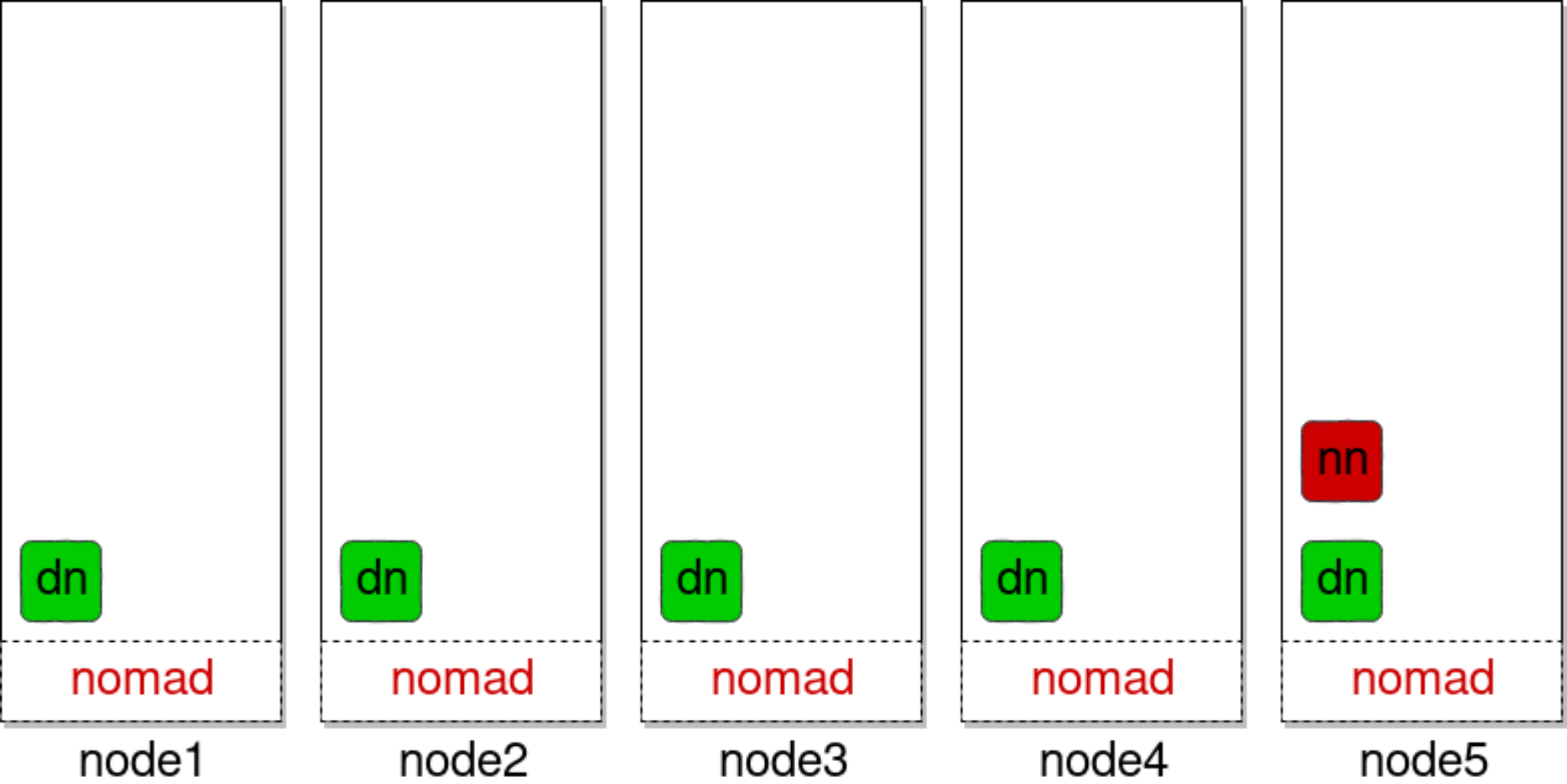
node5

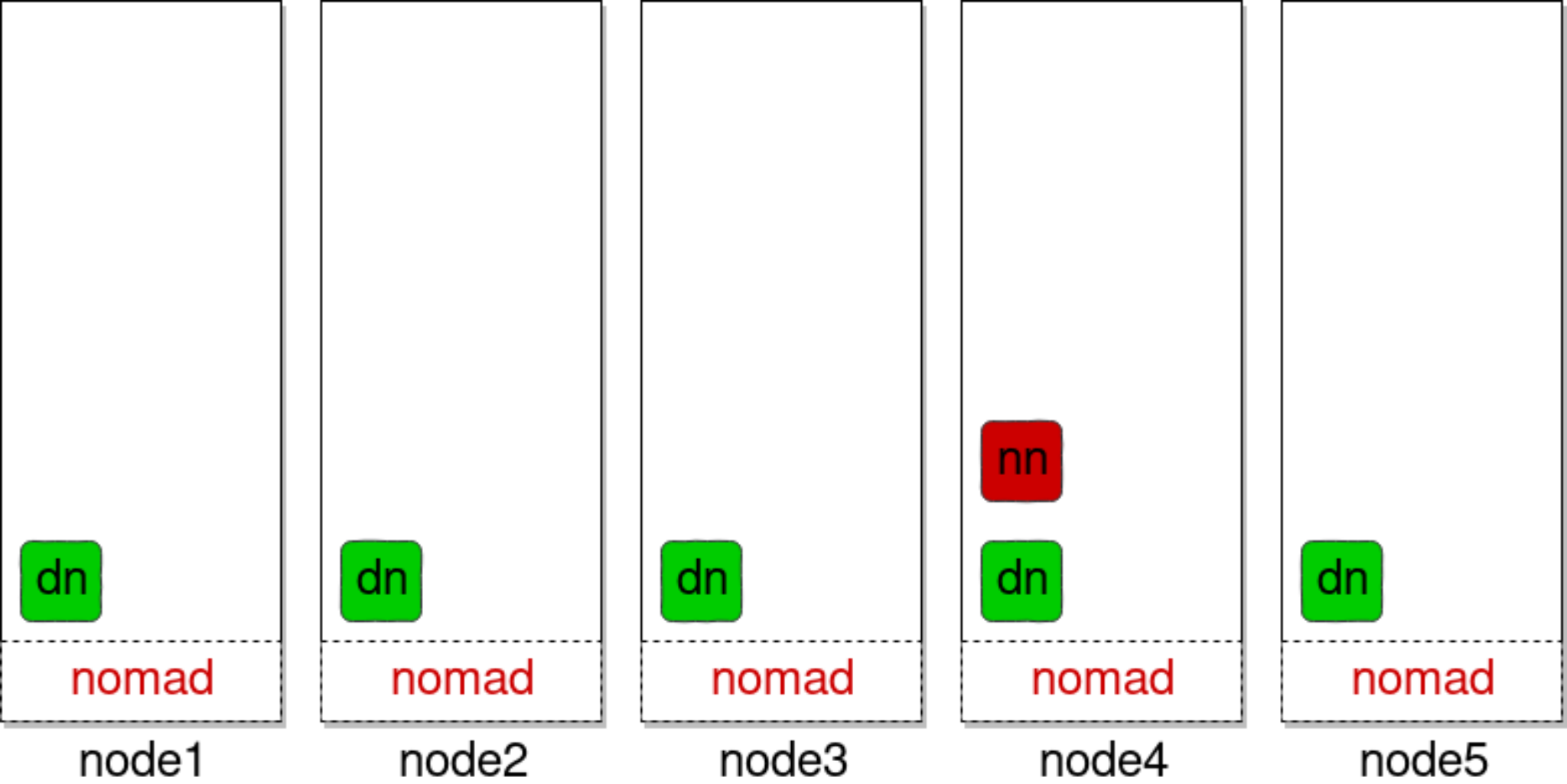




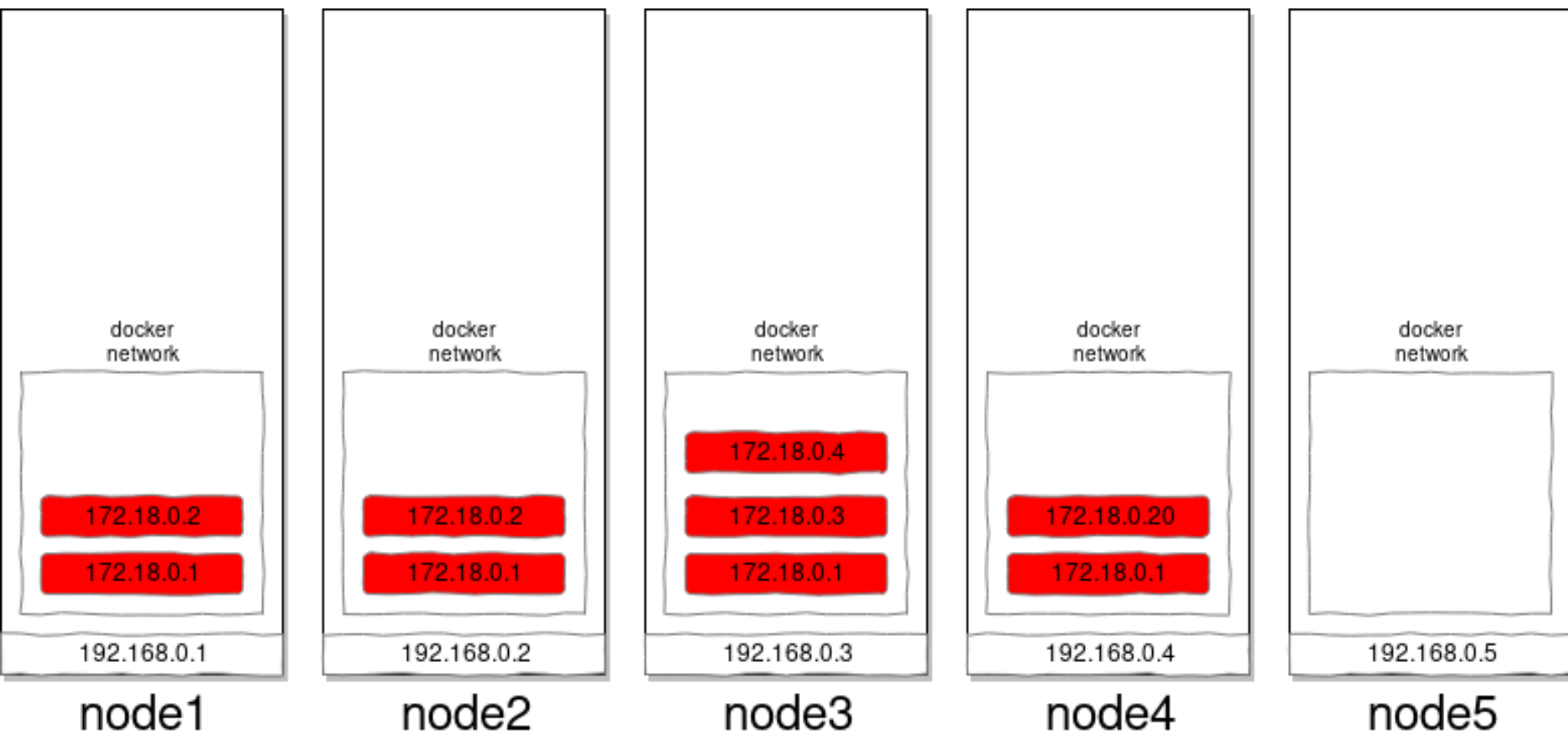




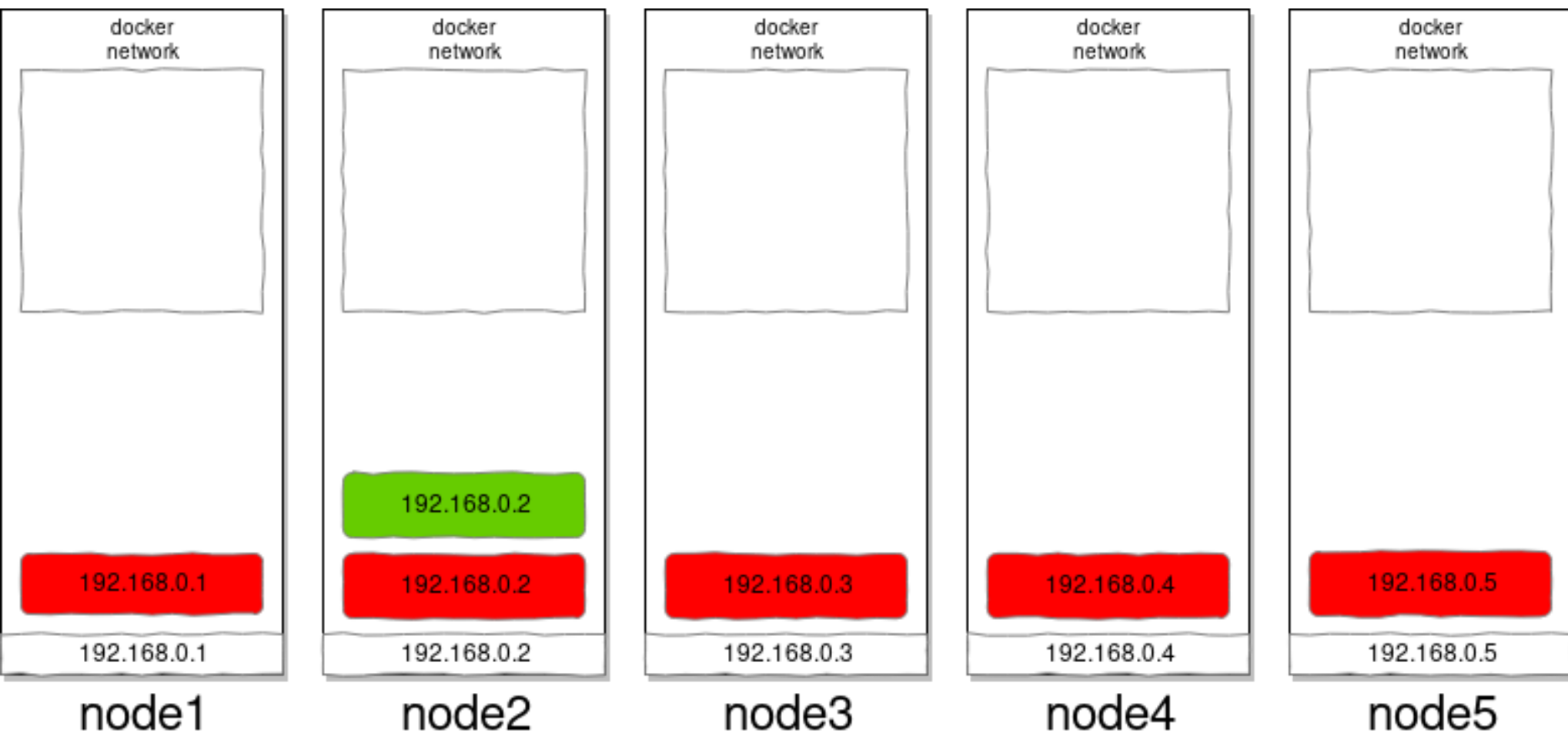




# Docker network



# Host network



## Configuration management

Source

Preprocessing

On change

## Provisioning, Scheduling

Multihost support

Scheduling

Cluster definition

Scaling

Multi tenancy

Failover

## Network

Intraservice network

DNS

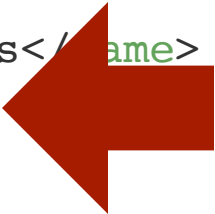
Service discovery

Data locality

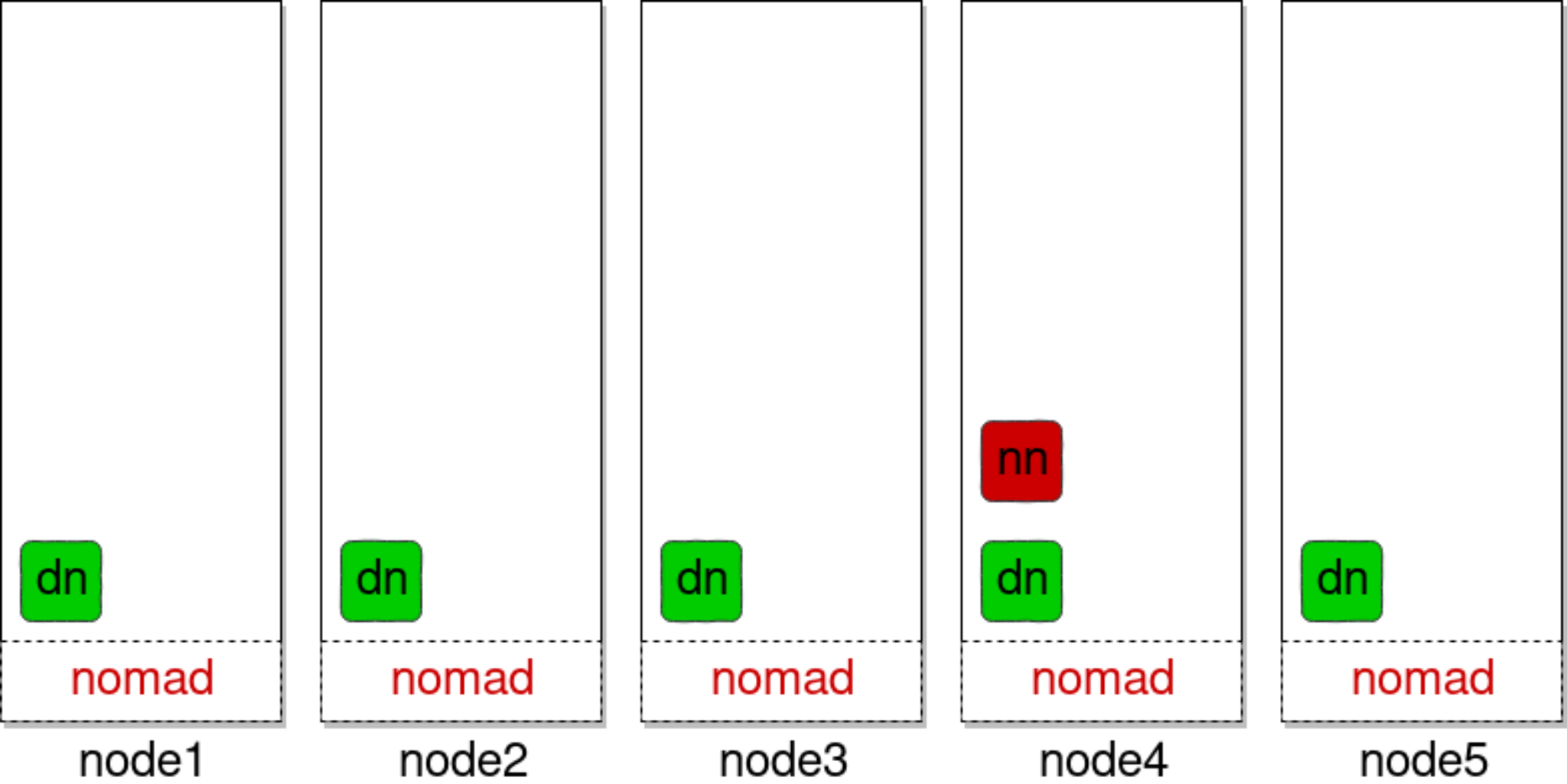
Availability of the ports

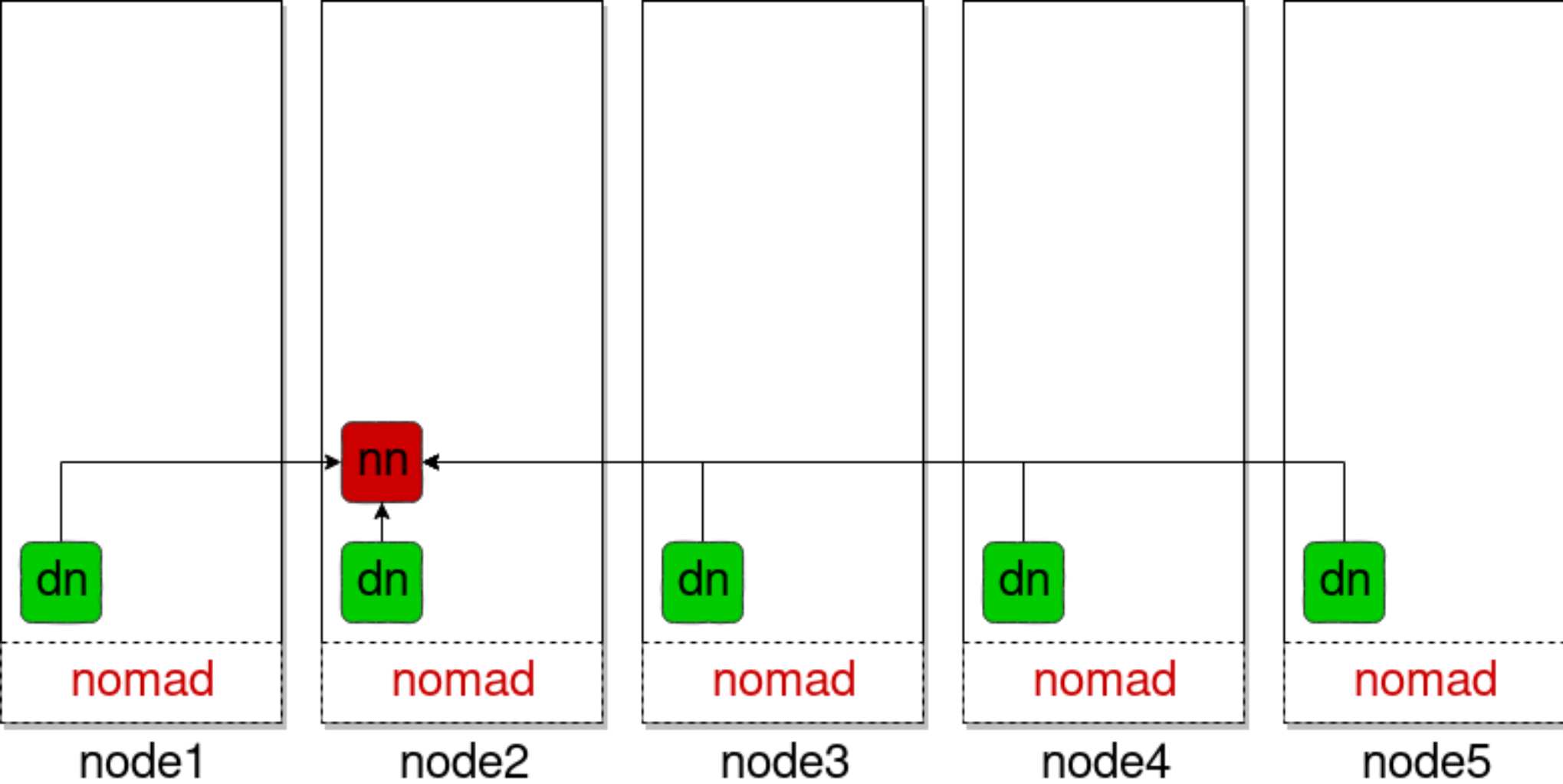
docker host  
host dns  
dns  
yes  
host

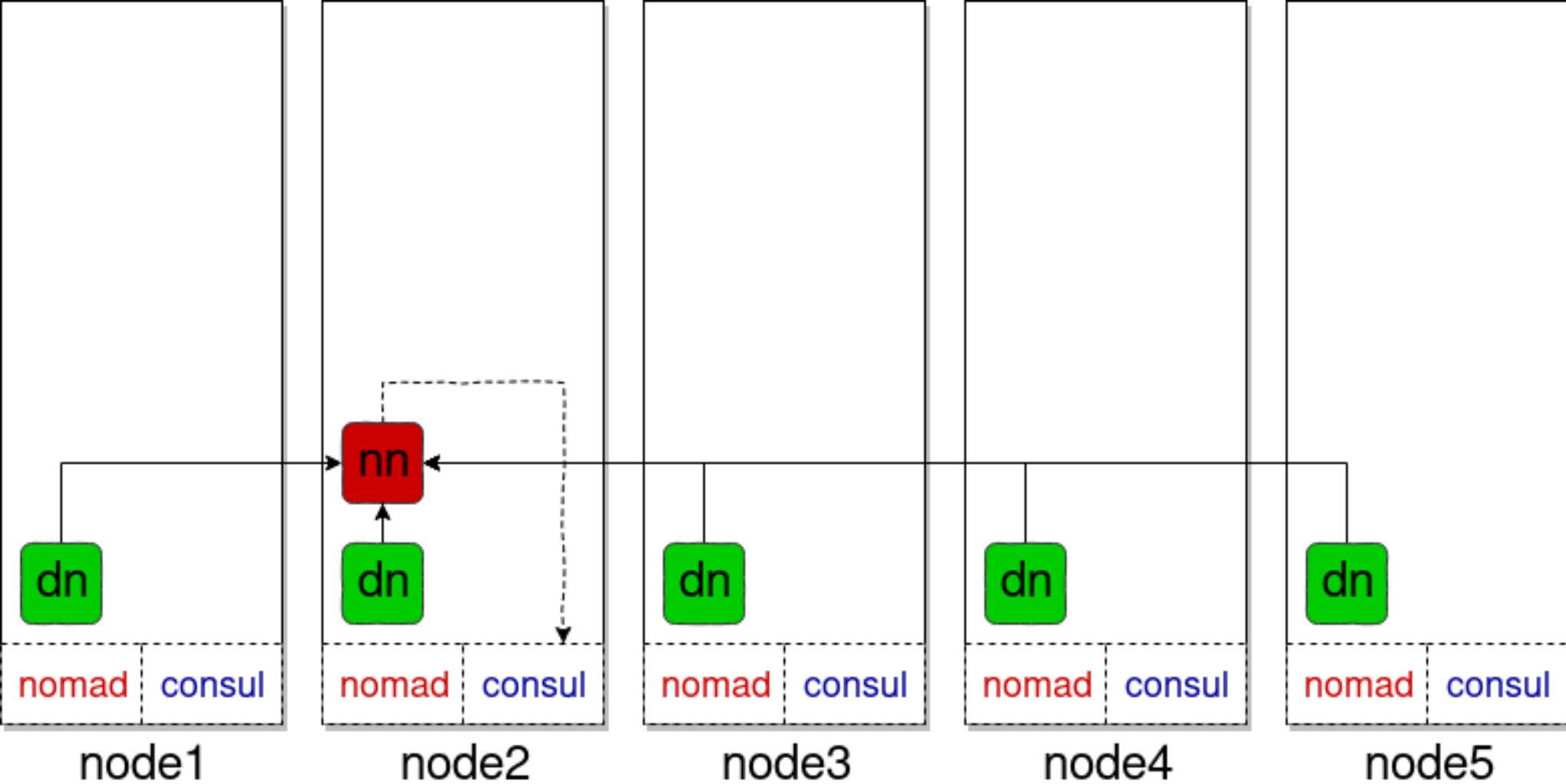
```
<configuration>
  <property>
    <name>dfs.namenode.rpc-address</name>
    <value>namenode:9000</value>
  </property>
  <property>
    <name>dfs.datanode.plugins</name>
    <value>org.apache.hadoop.ozone.HddsDatanodeService</value>
  </property>
  <property>
    <name>rpc.metrics.percentiles.intervals</name>
    <value>60,300</value>
  </property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>/data/namenode</value>
  </property>
  <property>
    <name>rpc.metrics.quantile.enable</name>
    <value>true</value>
  </property>
</configuration>
```

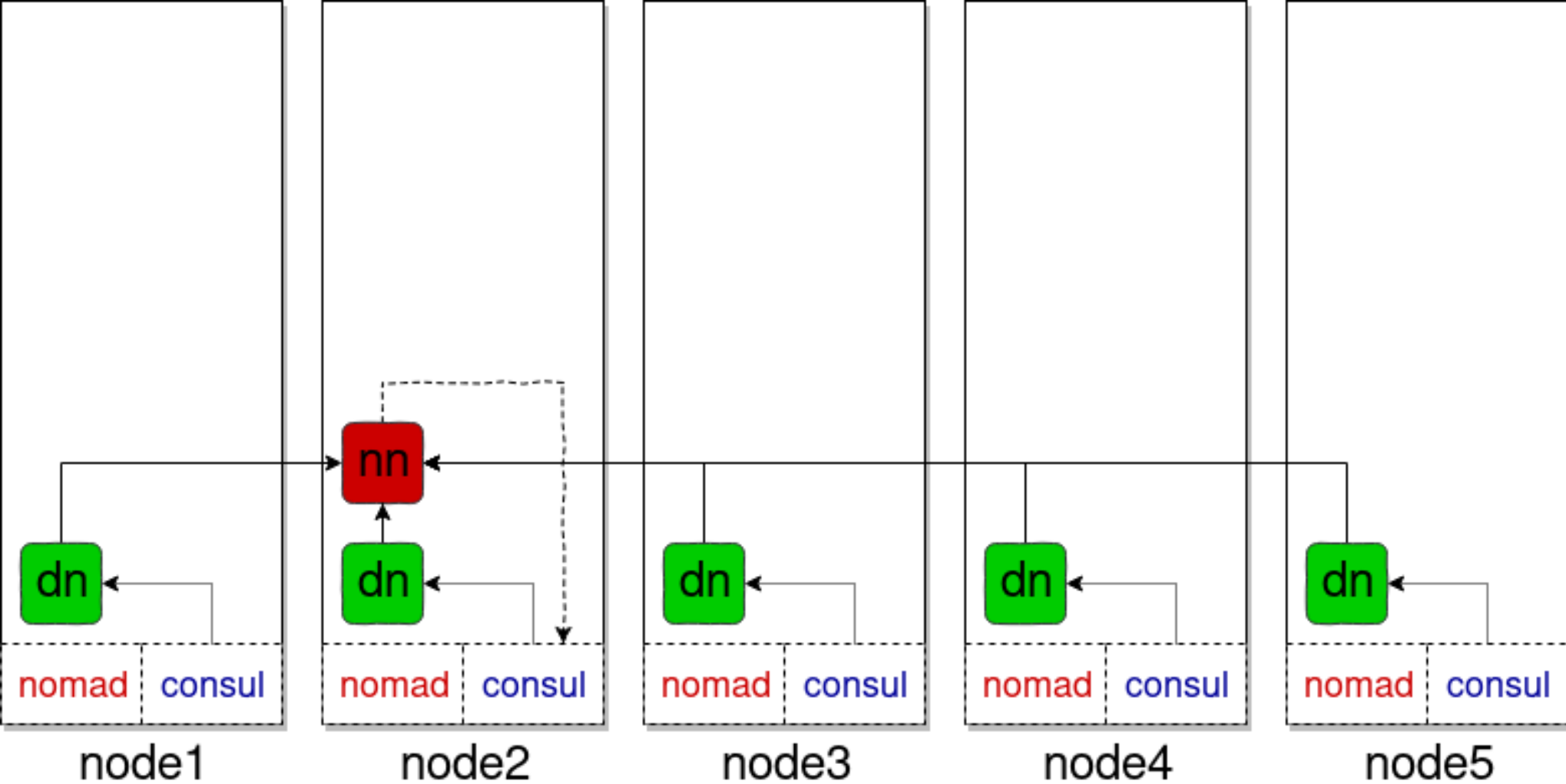


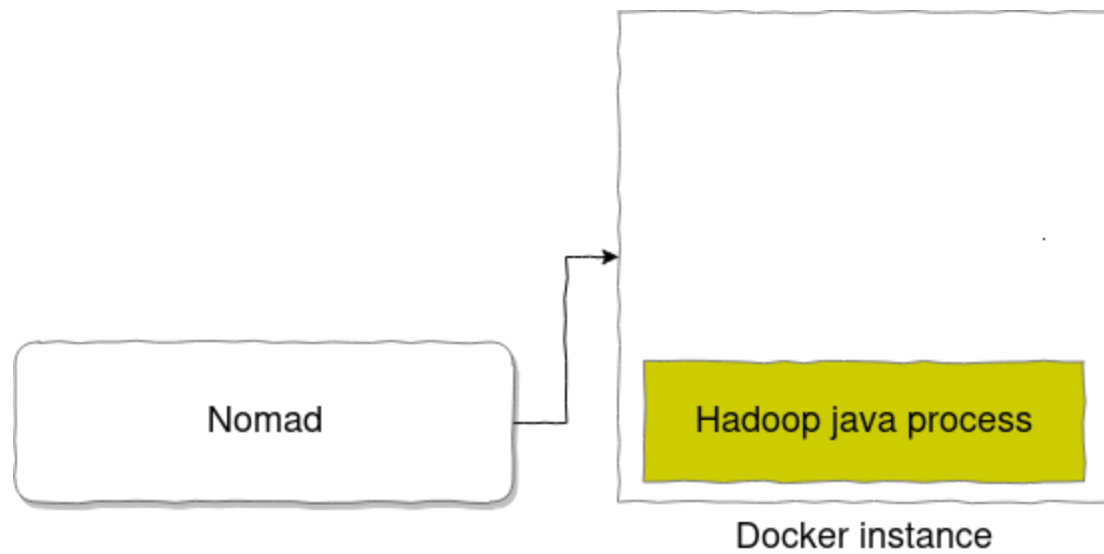


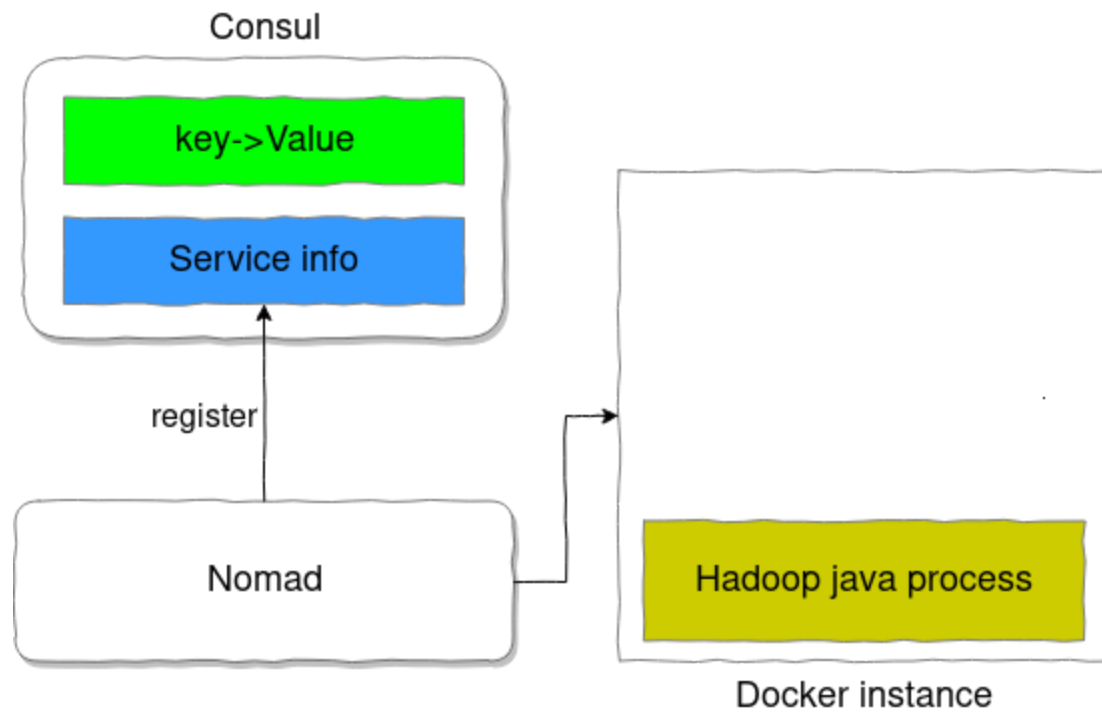


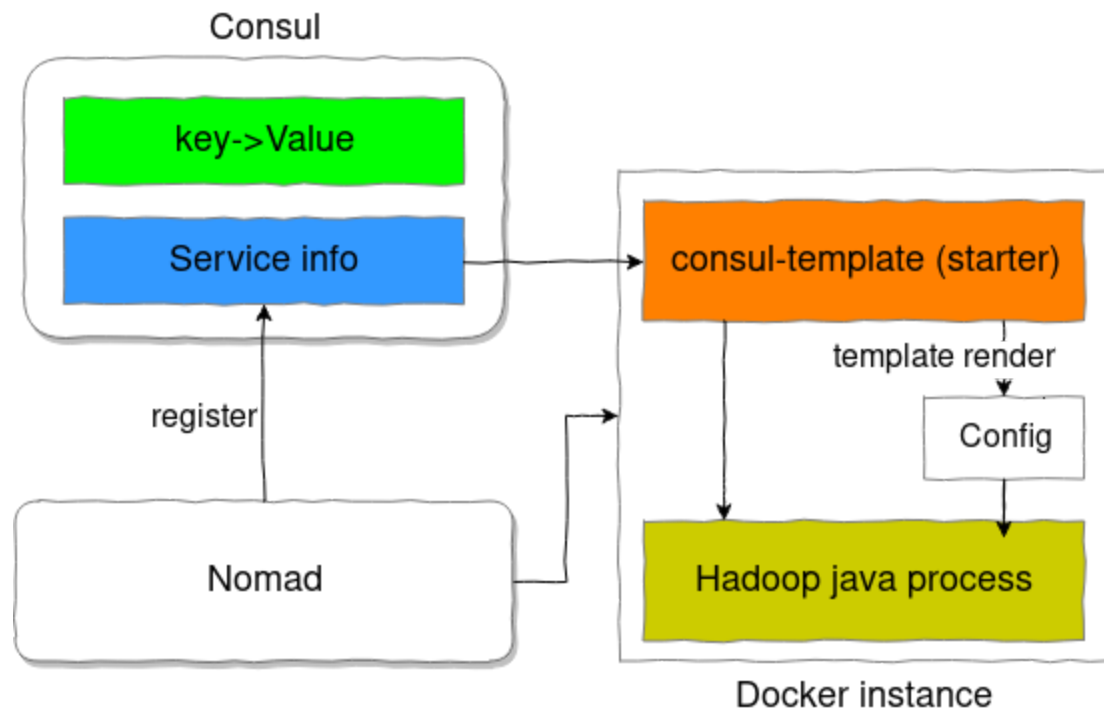


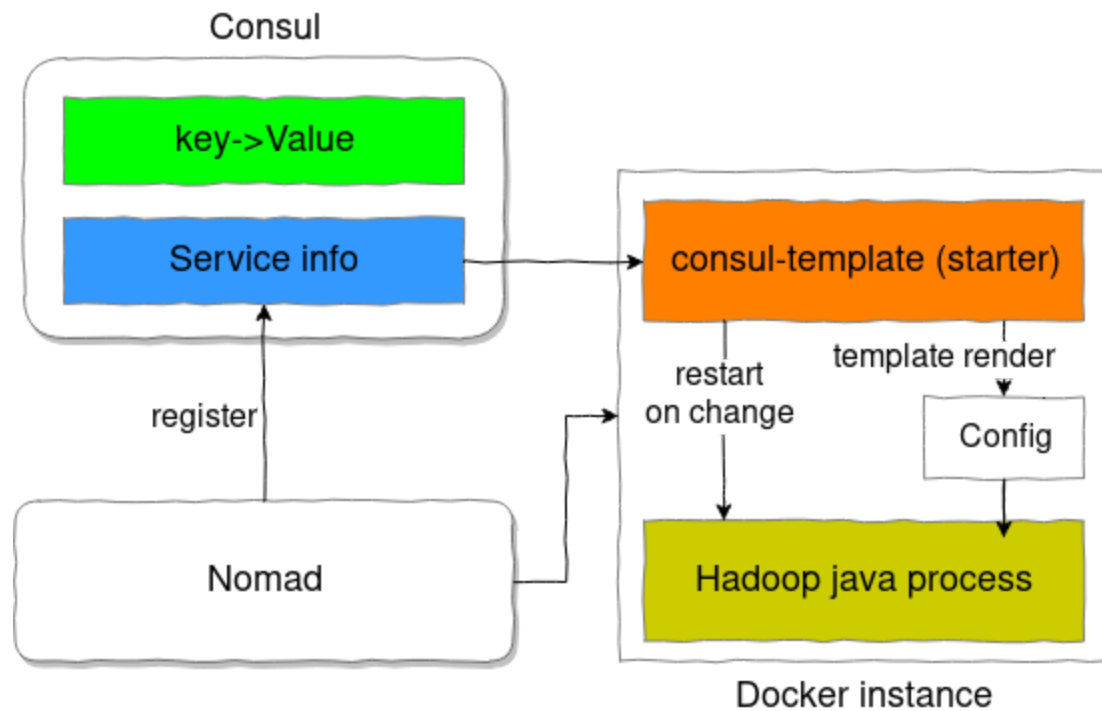




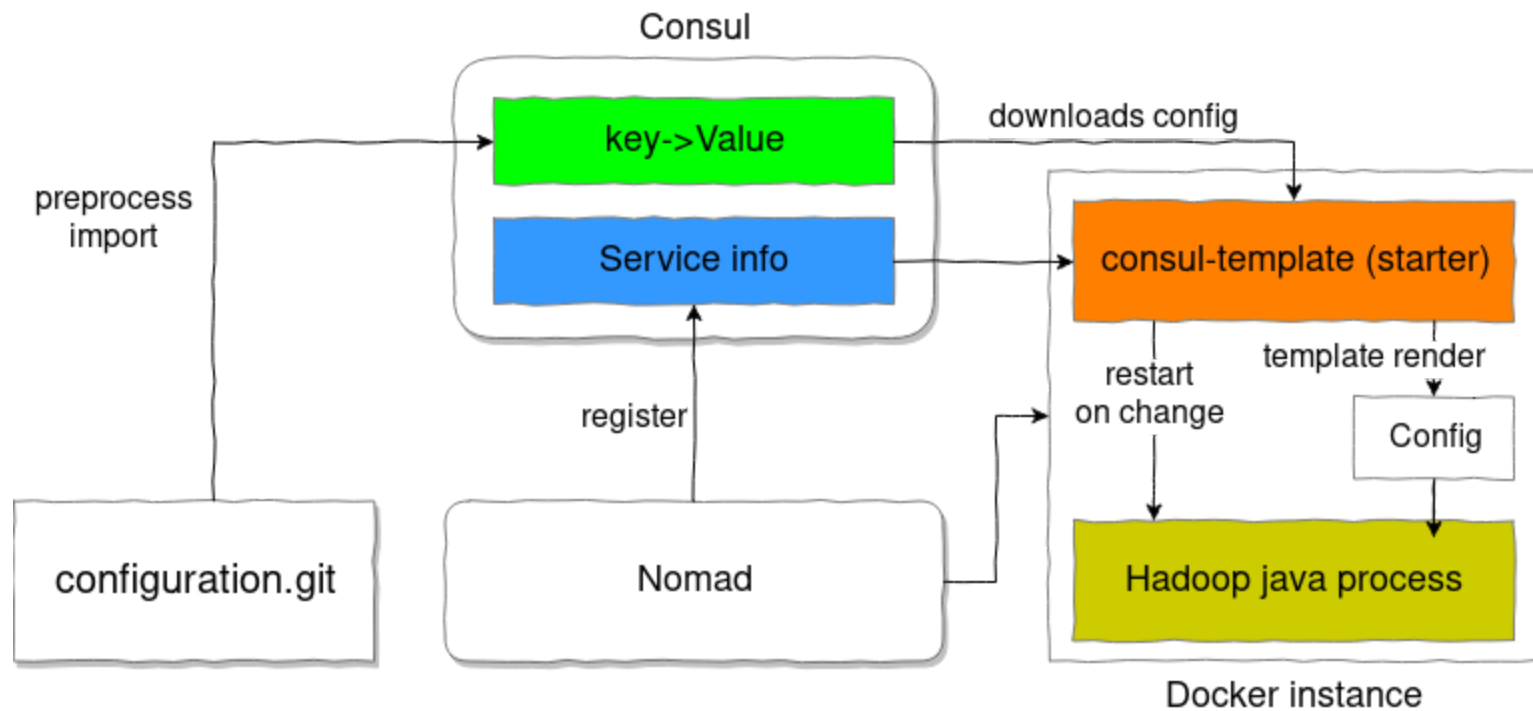














# DIY



## Nodes



## Healthy Nodes



```
2018-06-04 21:24:04 INFO DataNode:422 - Successfully sent block report 0x763c45c65afe8d18, containing 1 storage report(s), of which we sent 1. The reports had 0 total blocks and used 1 RPC(s). This took 12 msec to generate and 75 msec for RPC and MN processing. Got back one command: FinalizeCommand/5.
```

```
2018-06-04 21:24:04 INFO DataNode:759 - Got finalize command for block pool BP-643505683-127.0.0.1-1528147013324
```

```
2018-06-04 21:24:01 INFO DataNode:422 - Successfully sent block report 0x900e6d4b9e81dba1, containing 1 storage report(s), of which we sent 1. The reports had 0 total blocks and used 1 RPC(s). This took 7 msec to generate and 49 msec for RPC and MN processing. Got back one command: FinalizeCommand/5.
```

```
2018-06-04 21:24:01 INFO DataNode:759 - Got finalize command for block pool BP-643505683-127.0.0.1-1528147013324
```



### Configuration management

Source

Preprocessing

On change

Consul  
Yes (script)  
Restart

### Provisioning, Scheduling

Multihost support

Scheduling

Cluster definition

Scaling

Multi tenancy

Failover

host netw  
Nomad  
.nomad  
redploy  
no  
yes

### Network

Intraservice network

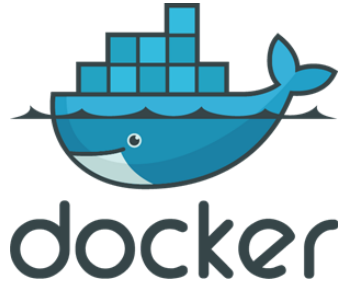
DNS

Service discovery

Data locality

Availability of the ports

host netw  
yes  
consul  
yes  
host



# Kubernetes

"out of the box"



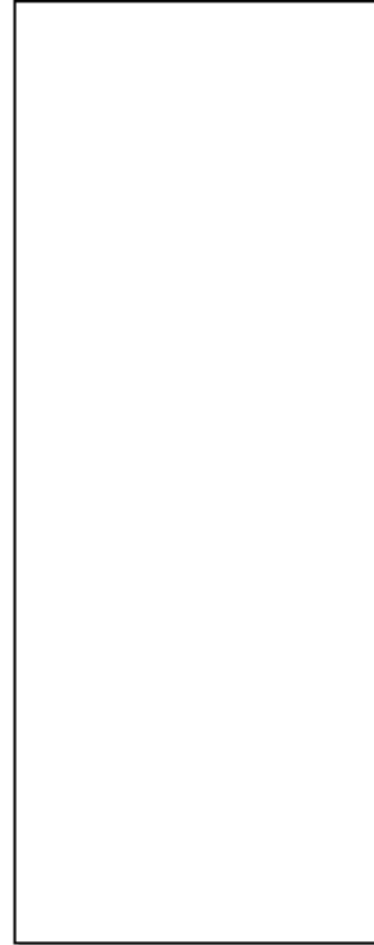
node1



node2



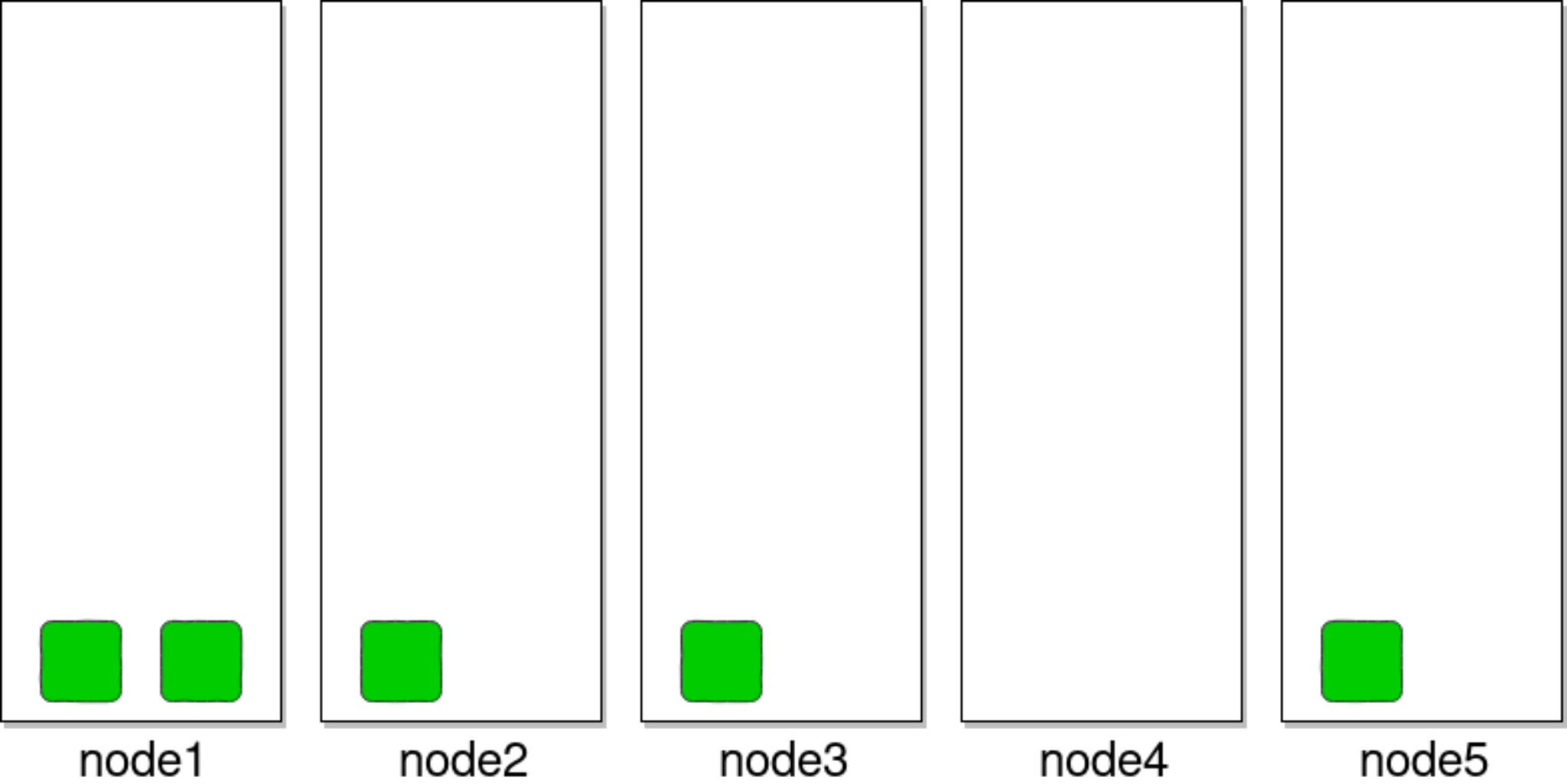
node3



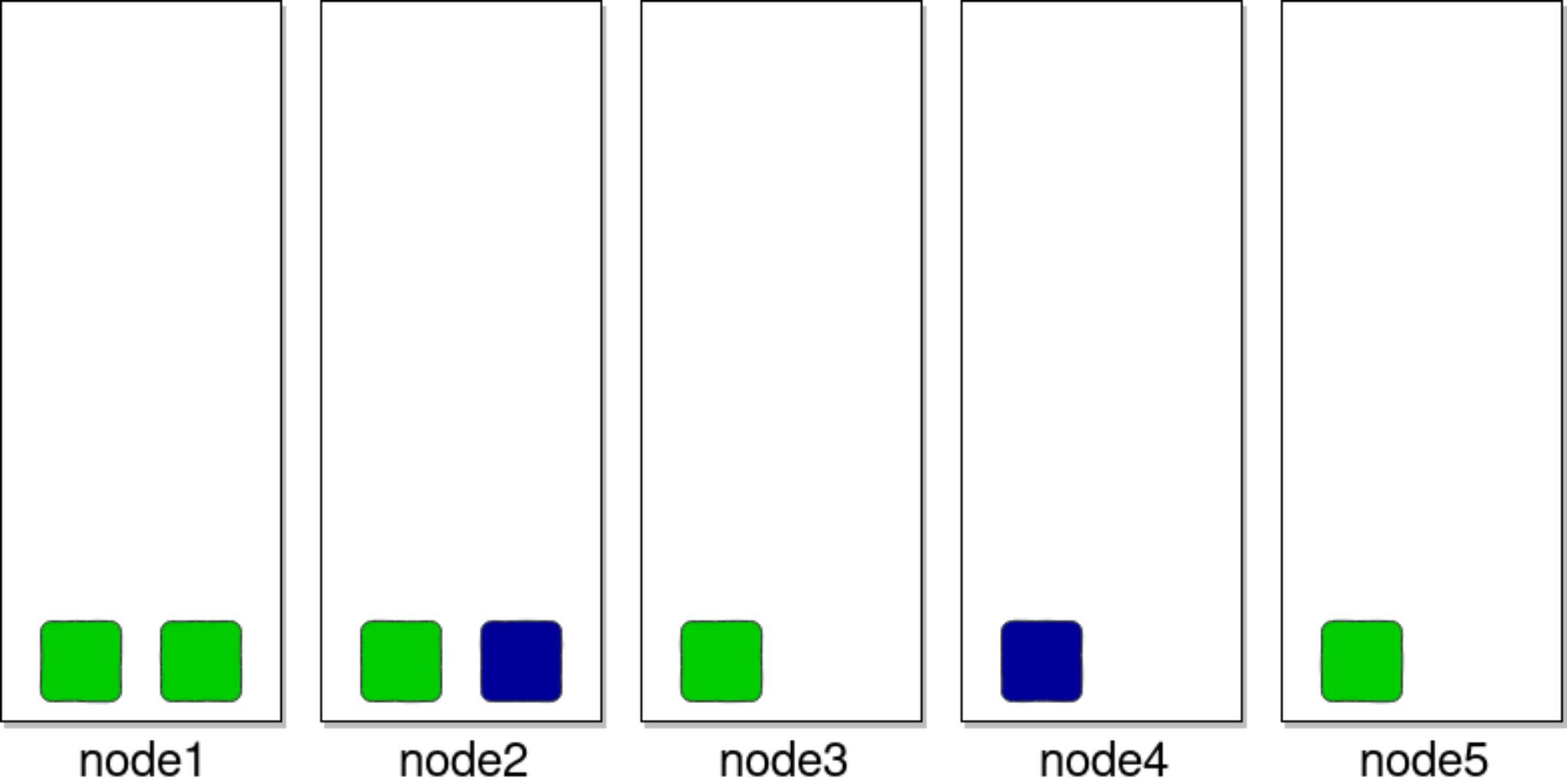
node4



node5

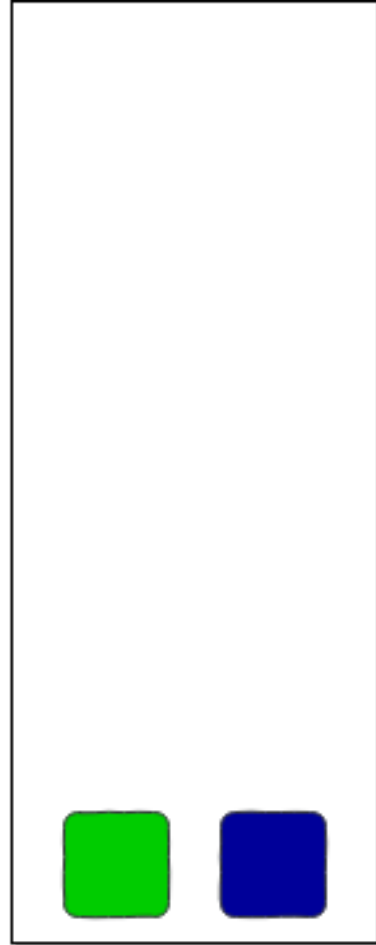








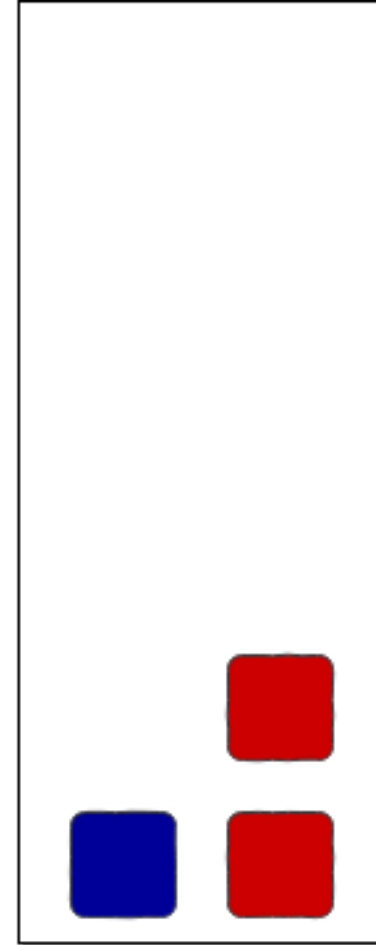
node1



node2



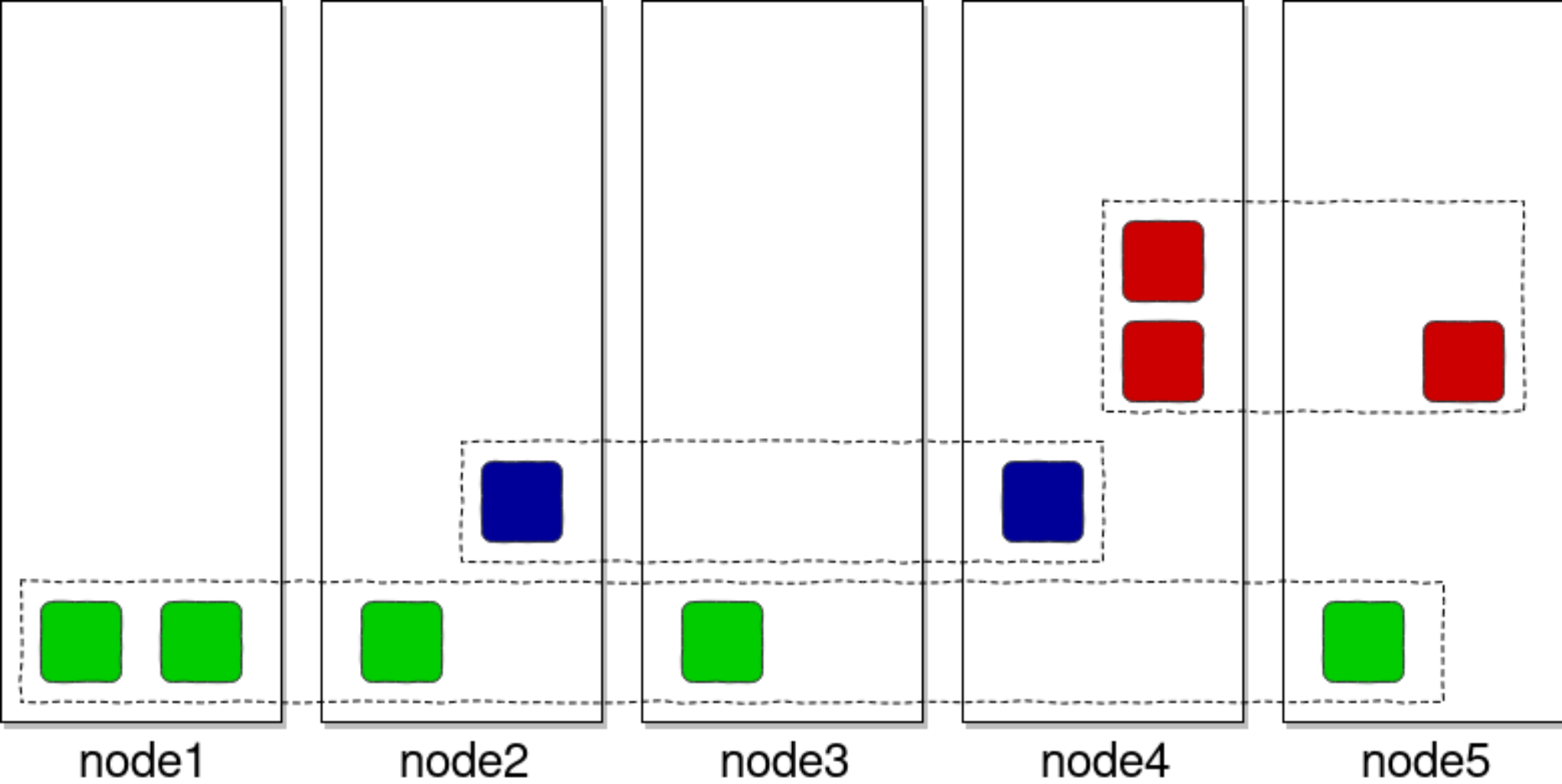
node3



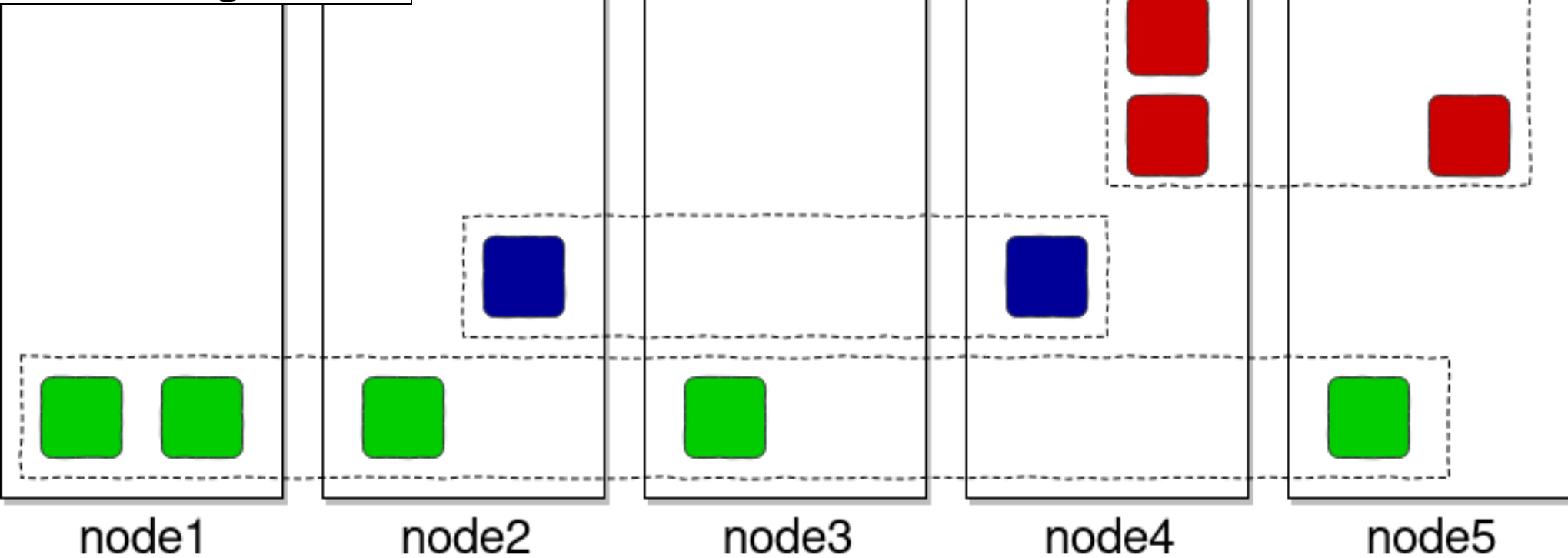
node4

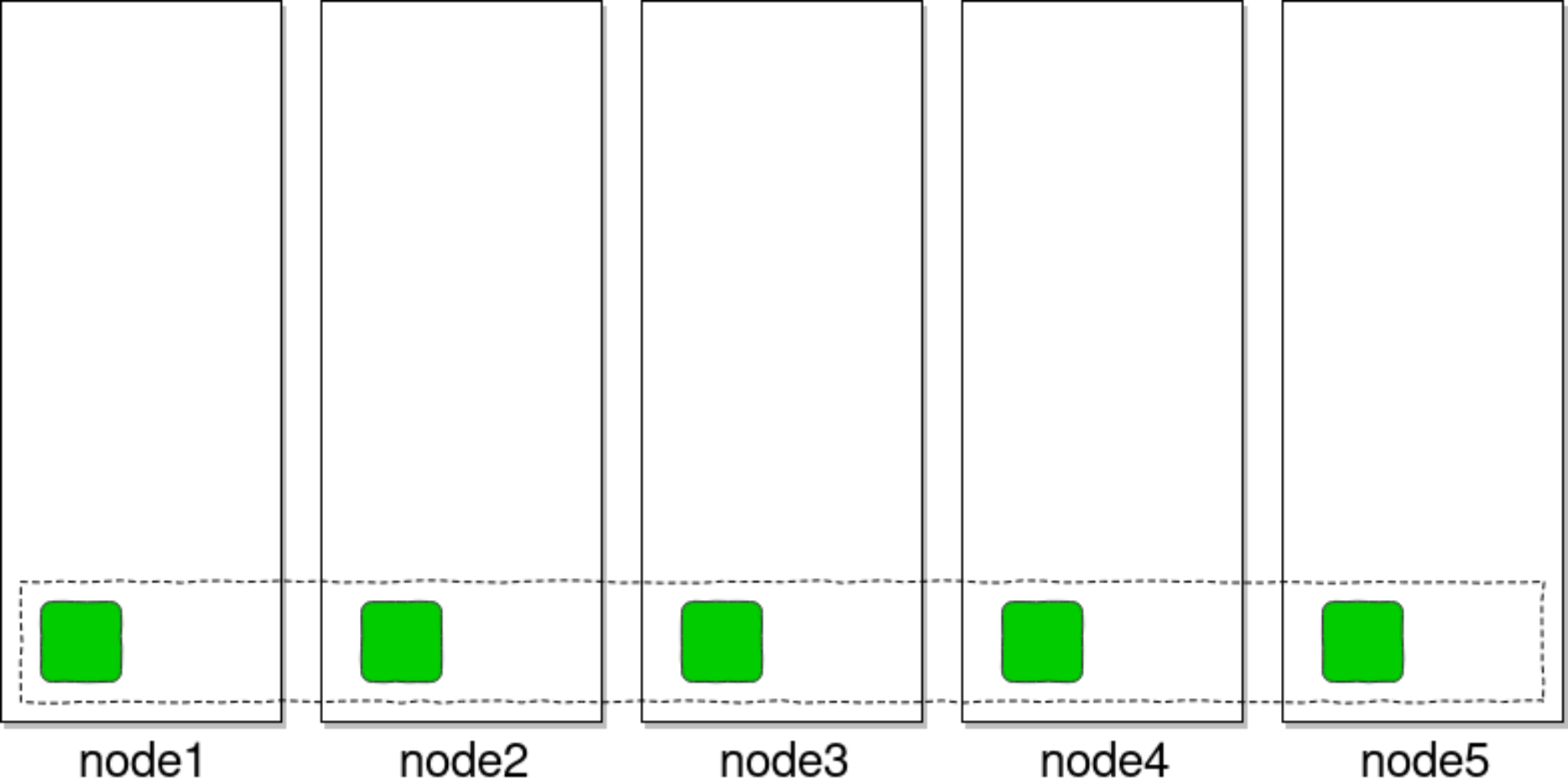


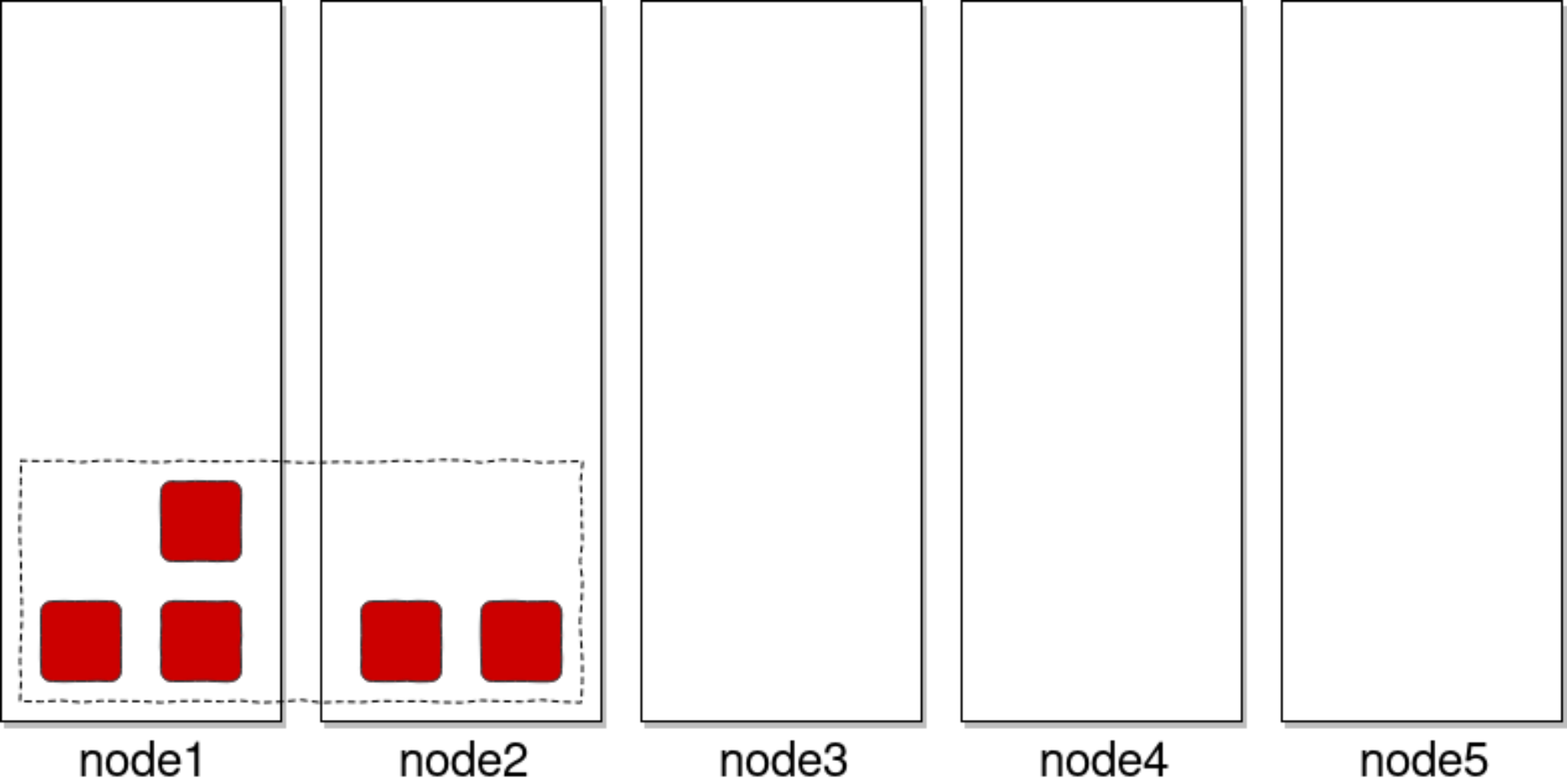
node5



+Network!  
+Storage!!  
(volume, secrets,  
configs)







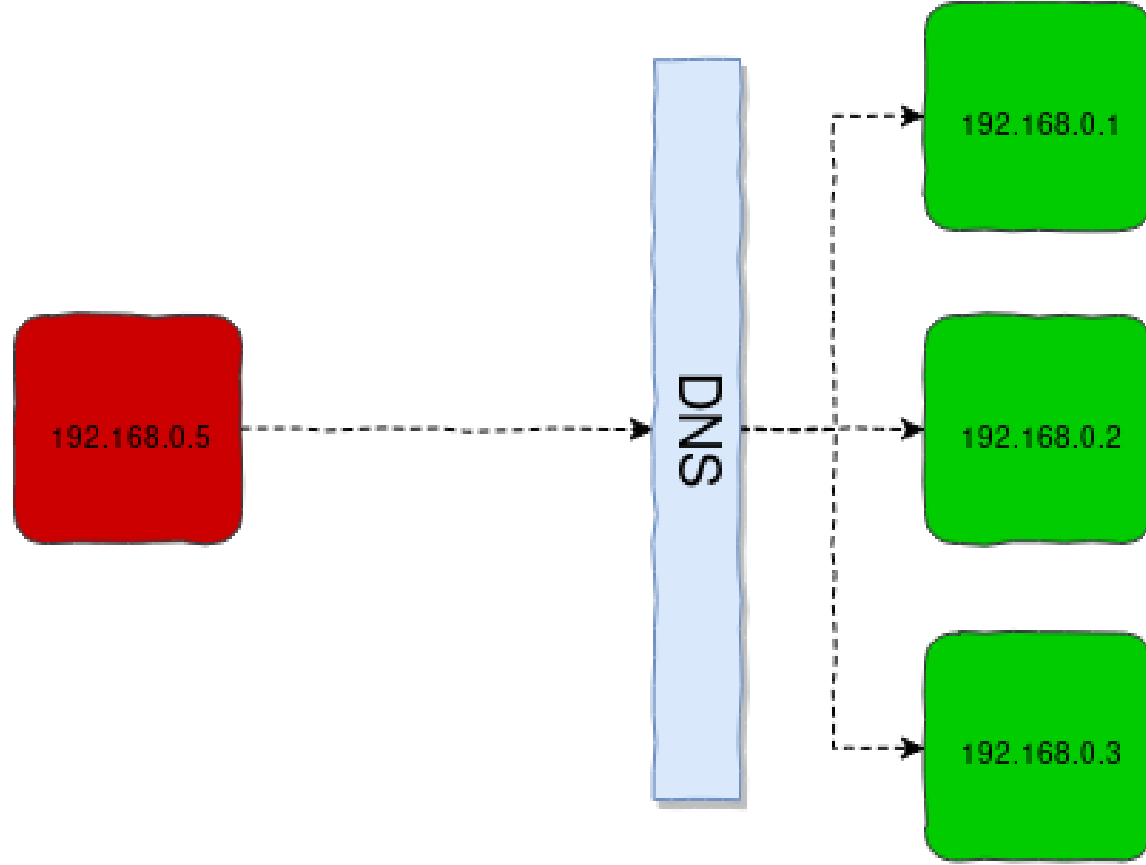
192.168.0.1

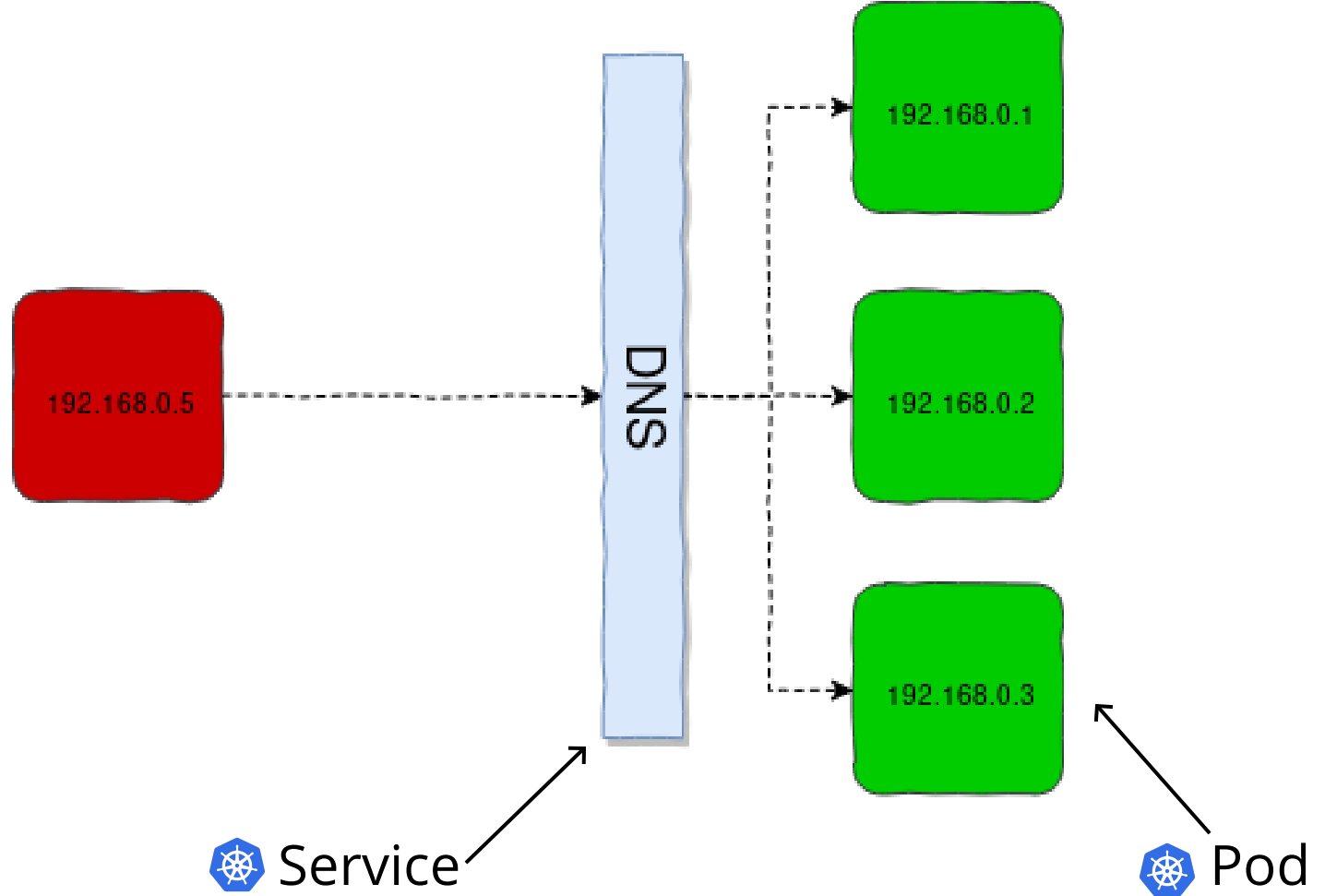
192.168.0.2

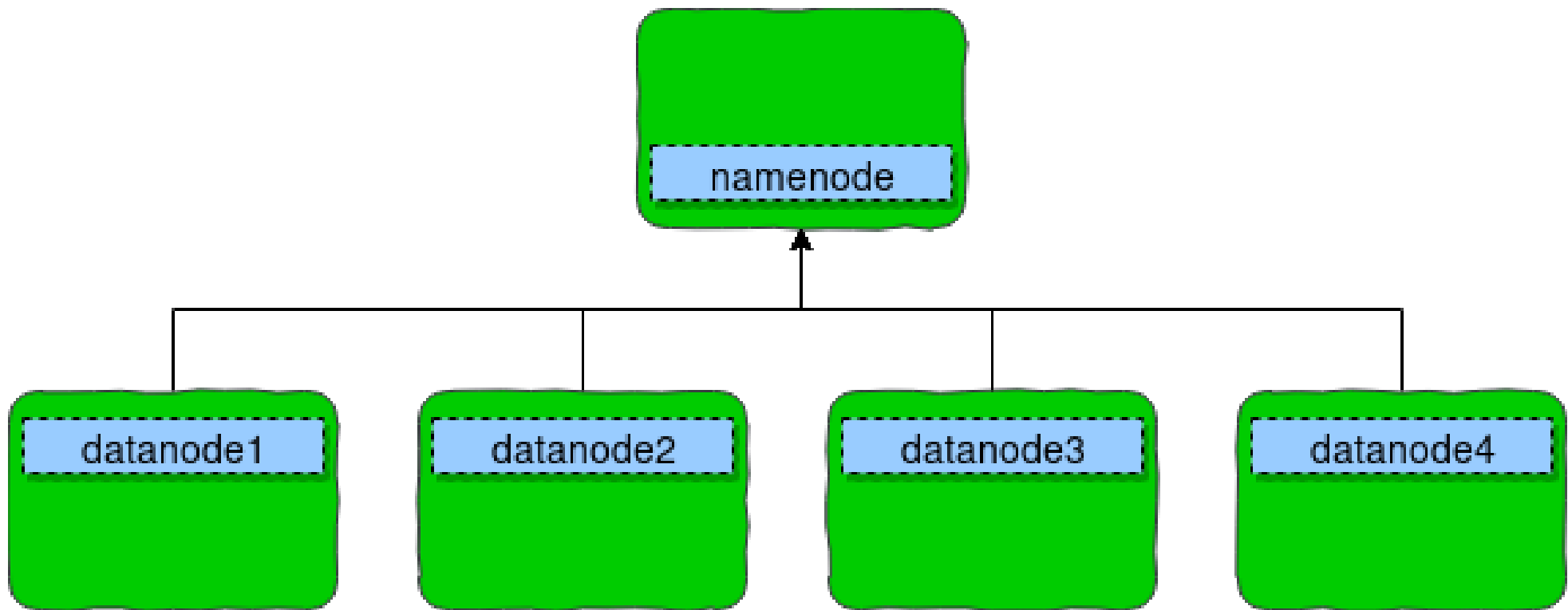
192.168.0.3

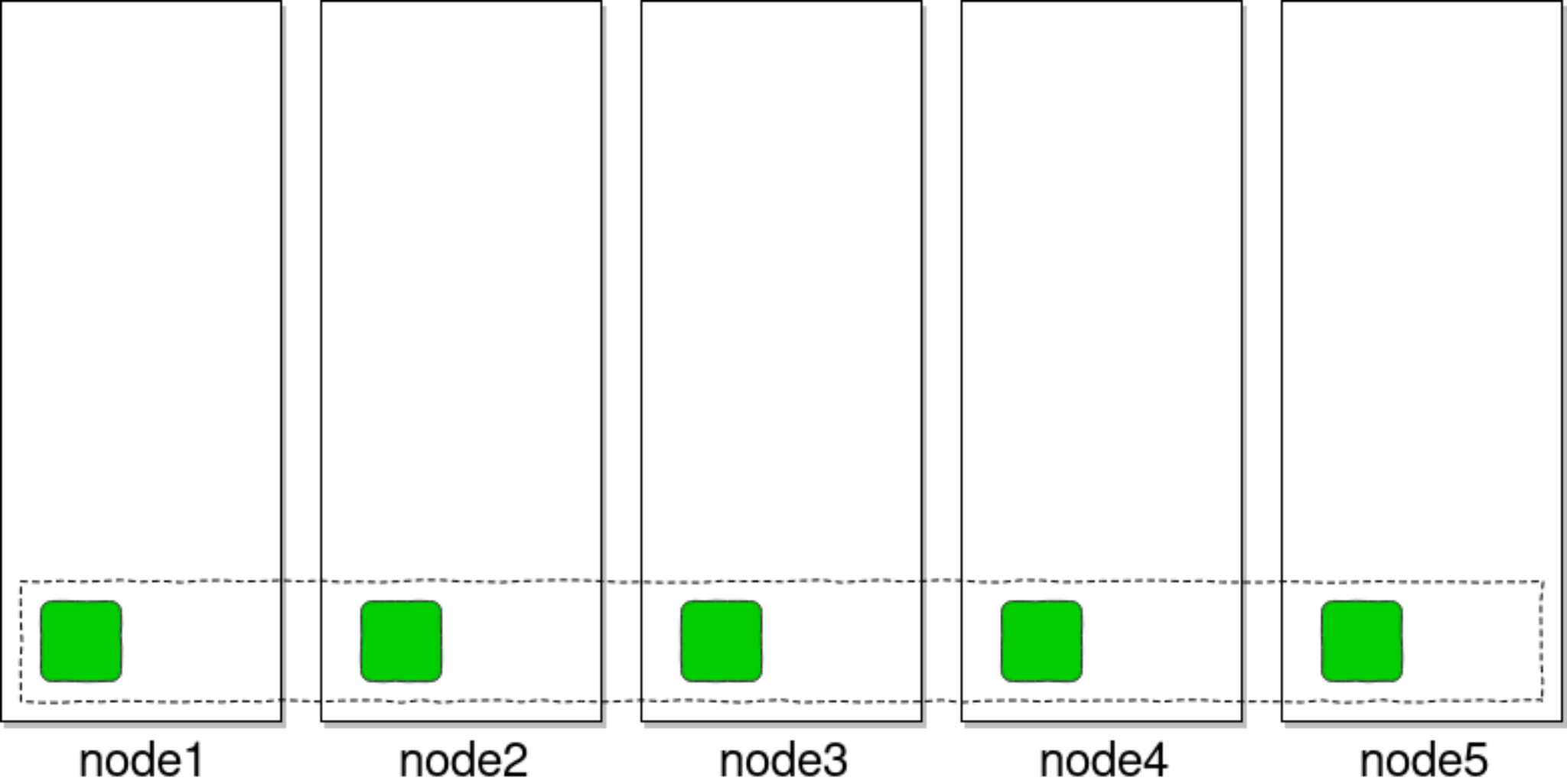


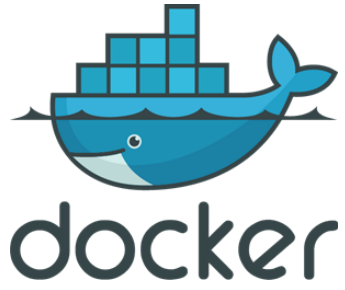


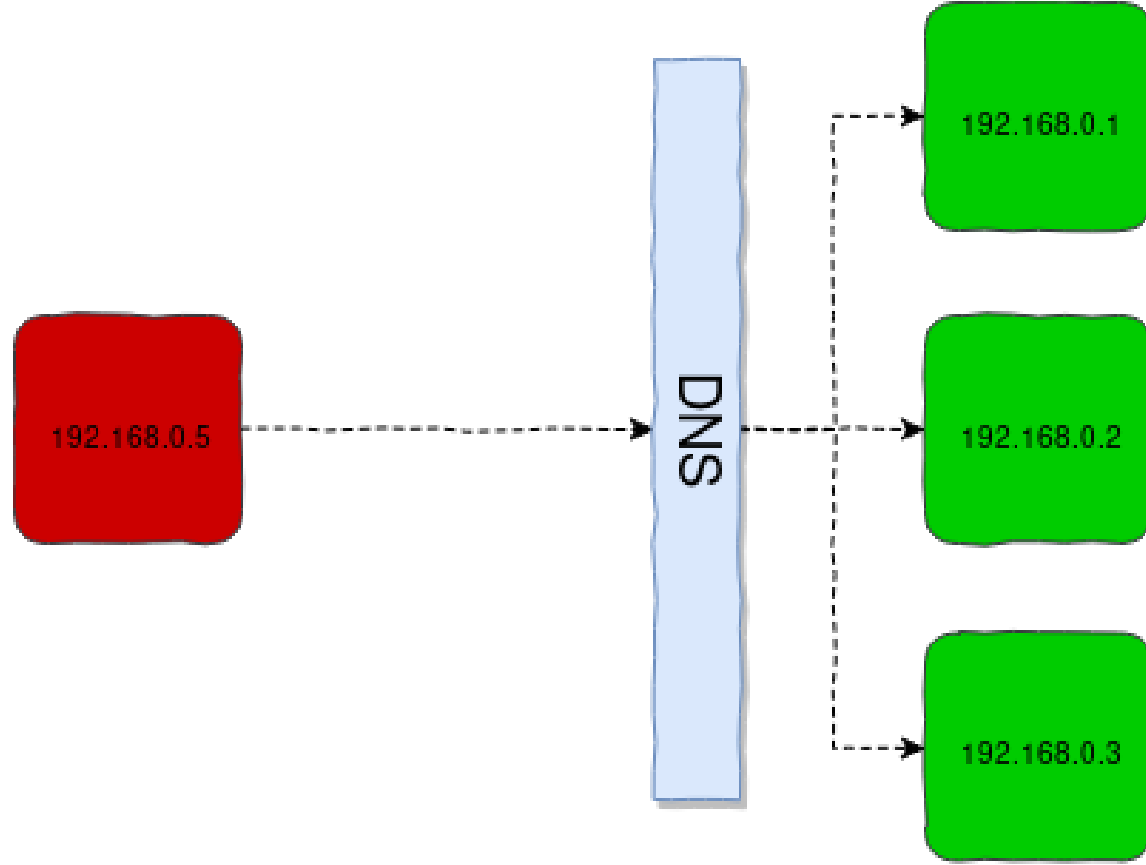


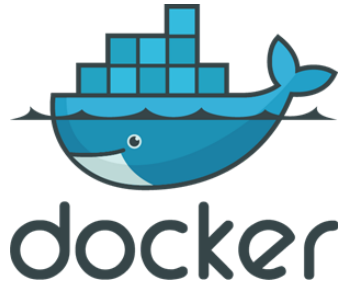












# Benefits of Hadoop + k8s?





Benefits of  
Hadoop + k8s?

Ecosystem  
Flexibility



# Example:

## Monitor Hadoop with Prometheus

[Enable query history](#)

rate(Hadoop\_KeySpaceManager\_NumKeyCommits[10m])

Load time: 187ms  
Resolution: 7s  
Total time series:

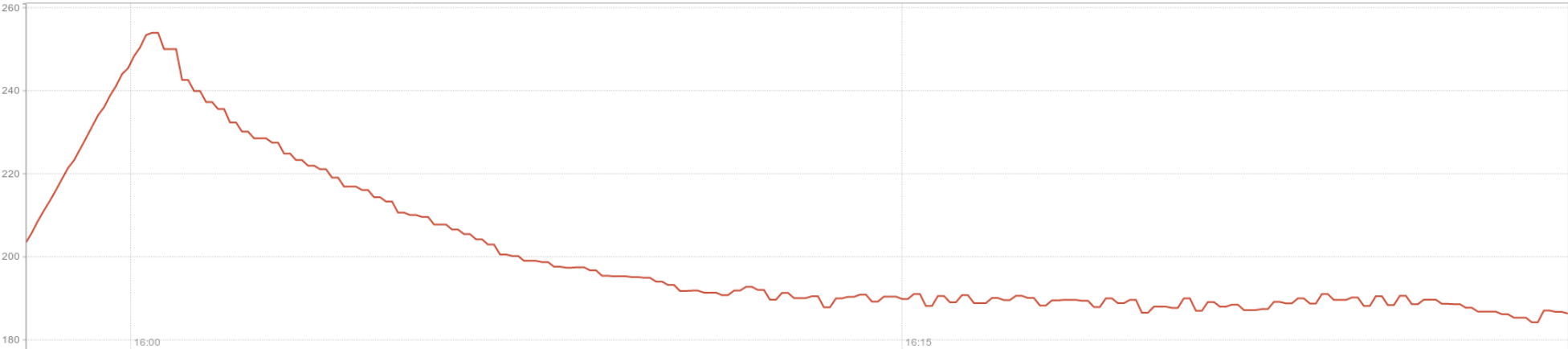
Execute

java\_lang\_OperatingSyste

Graph

Console

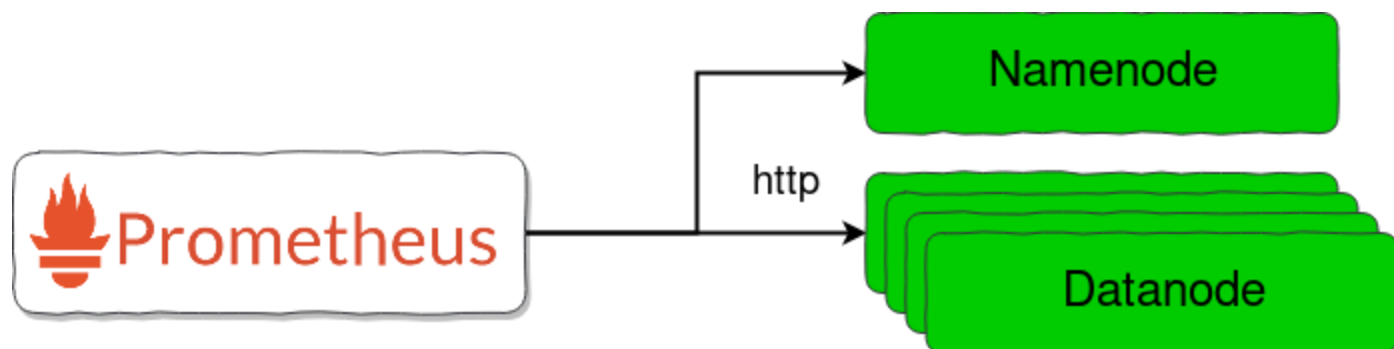
- 30m + ◀ Until ▶ Res. (s) stacked

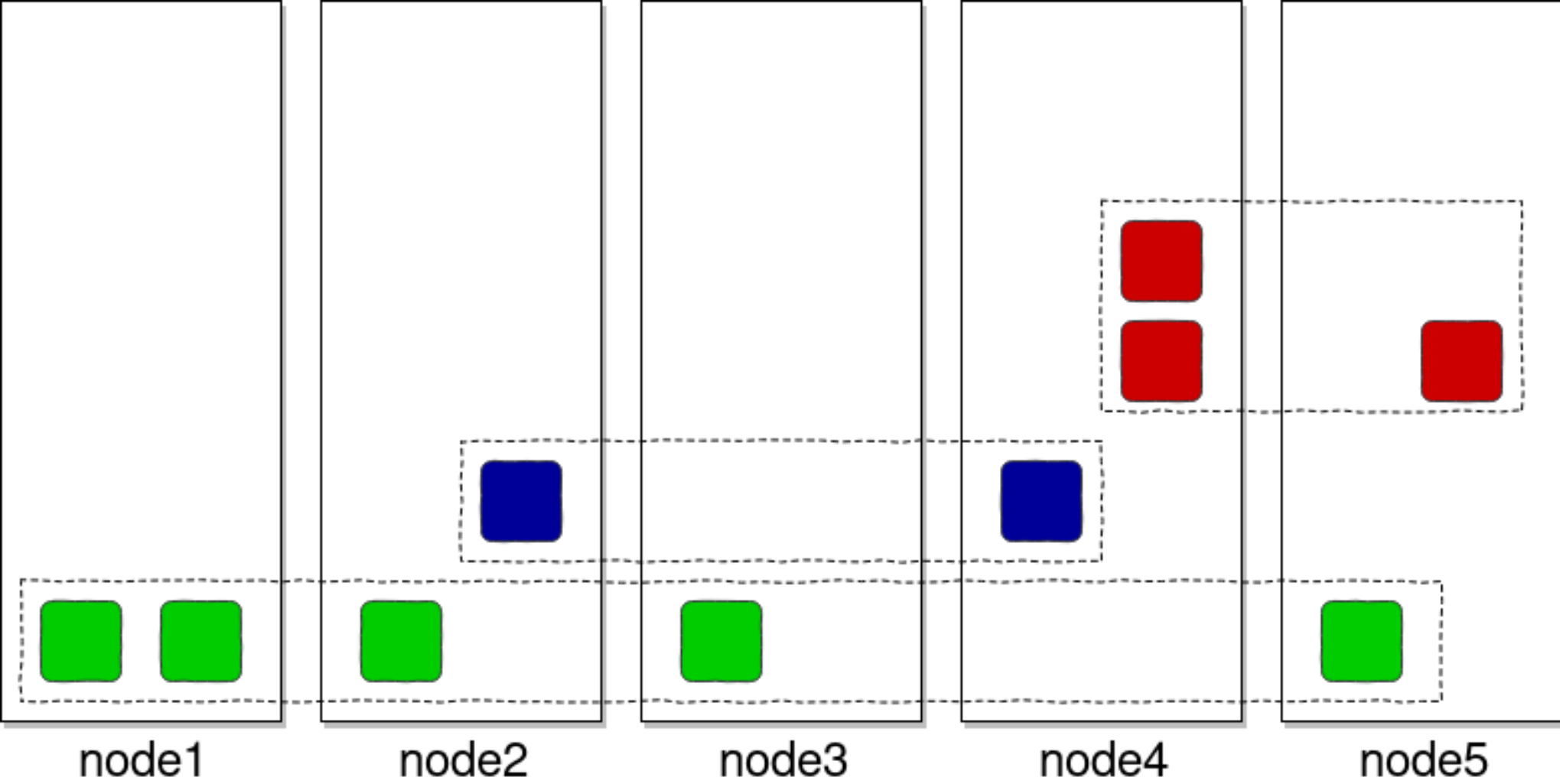


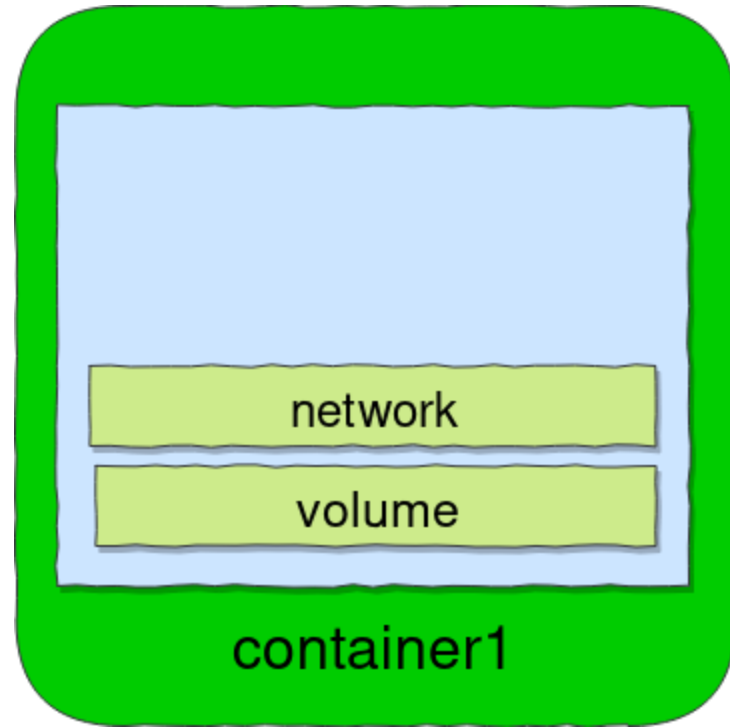
◻ (instance="c3e21ade589b:38271" job="jmxexporter", name="KSMMetrics")

[Remove Graph](#)

Add Graph

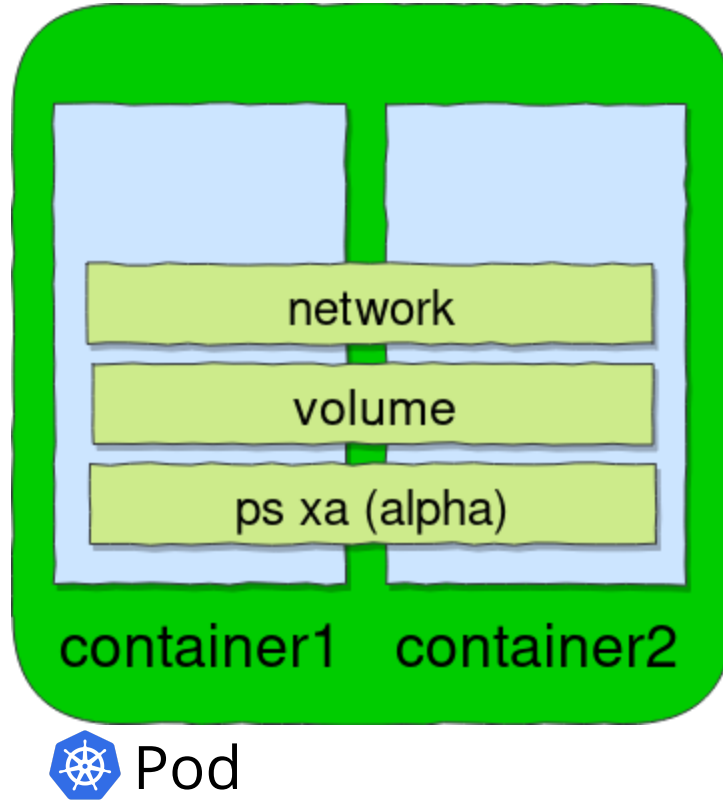




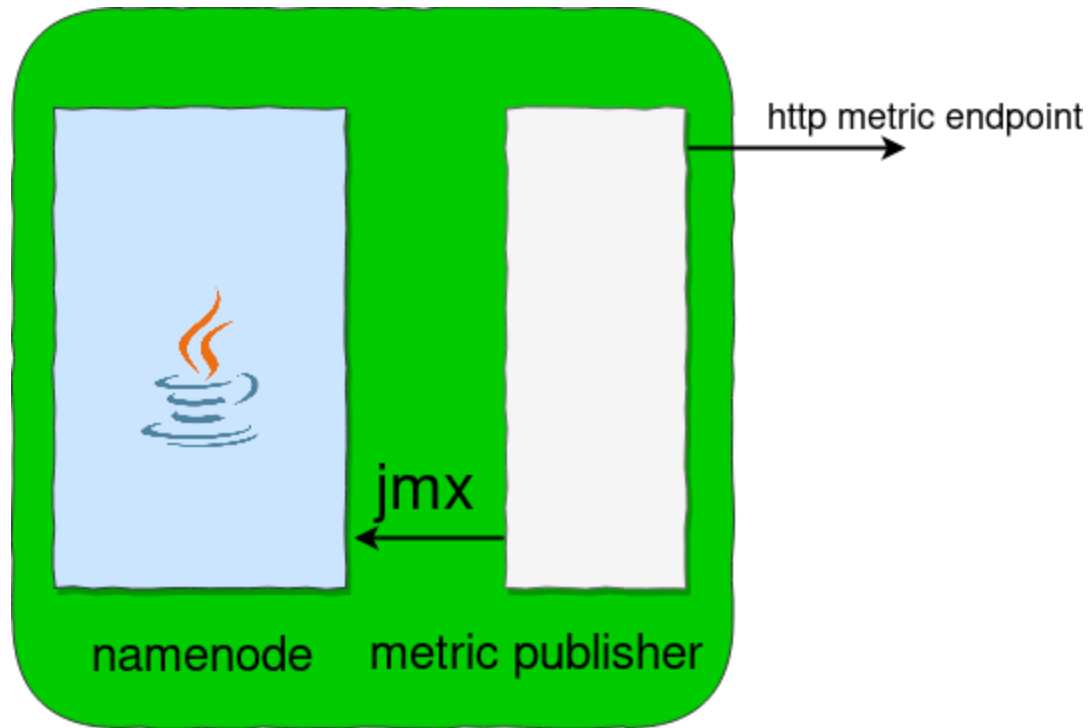


 Pod

# Sidecar pattern



# Sidecar for monitoring over jmx





```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  name: ozone-hdfs-namenode
spec:
  serviceName: ozone2-hdfs-namenode
  replicas: 1
  template:
    metadata:
      labels:
        app: ozone
    spec:
      containers:
        - name: hdfs-namenode
          image: flokkir/ozone:2.1.0
          args: ["hdfs", "namenode"]
```

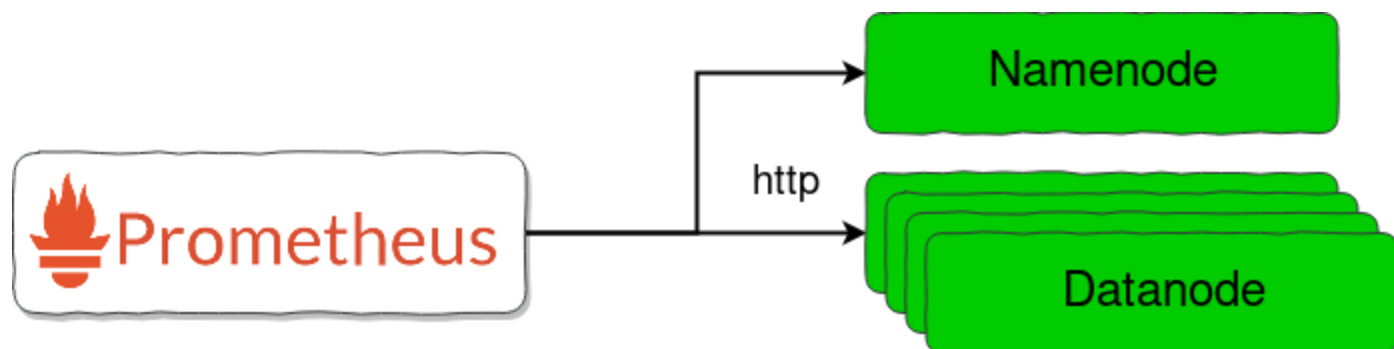
```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  name: ozone-hdfs-namenode
spec:
  serviceName: ozone2-hdfs-namenode
  replicas: 1
  template:
    metadata:
      labels:
        app: ozone
    spec:
      containers:
        - name: hdfs-namenode
          image: flokkkr/ozone:2.1.0
          args: ["hdfs", "namenode"]
```

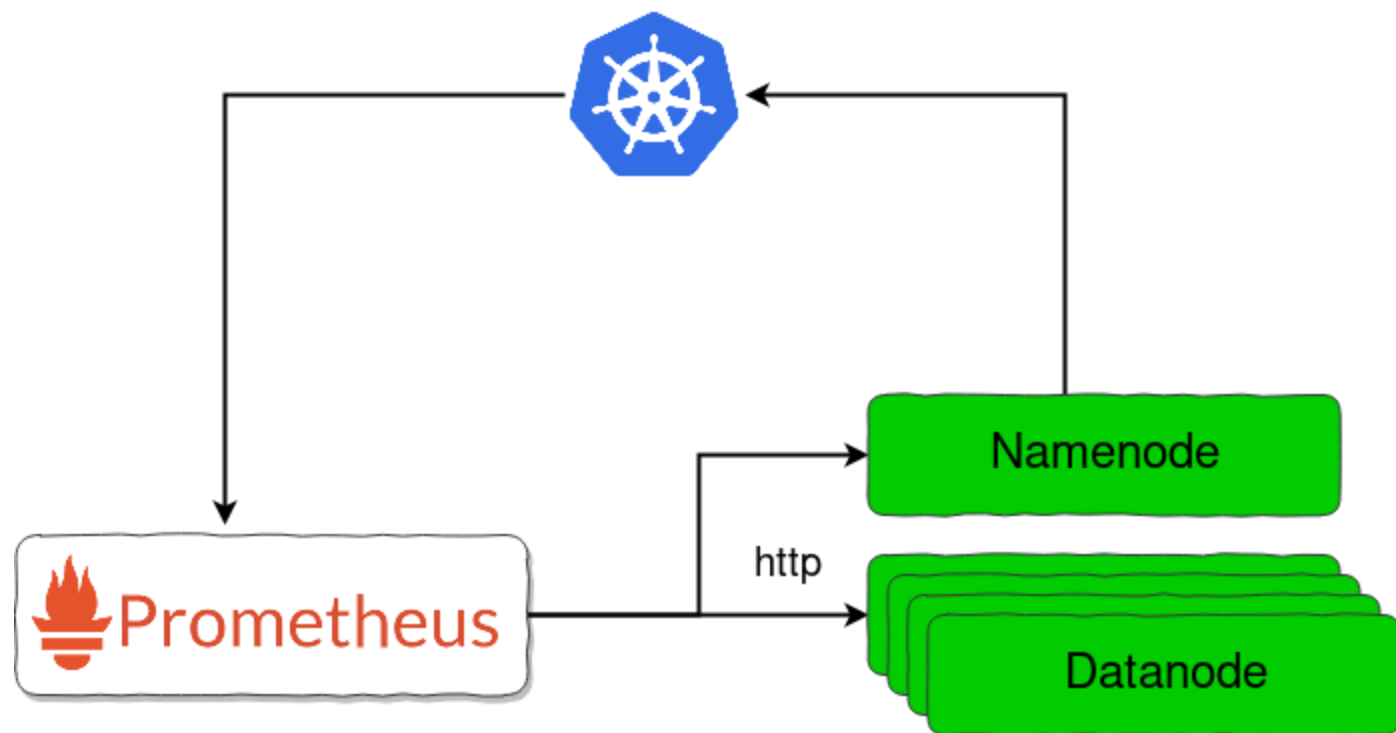
# Flexibility

```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  name: ozone-hdfs-namenode
spec:
  serviceName: ozone2-hdfs-namenode
  replicas: 1
  template:
    metadata:
      labels:
        app: ozone
    spec:
      shareProcessNamespace: true
      containers:
        - name: hdfs-namenode
          image: flokkkr/ozone:2.1.0
          args: ["hdfs", "namenode"]
        - name: jmxpromo
          image: flokkkr/jmxpromo-sidecar
```

```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  name: ozone-hdfs-namenode
spec:
  serviceName: ozone2-hdfs-namenode
  replicas: 1
  template:
    metadata:
      labels:
        app: ozone
    spec:
      shareProcessNamespace: true
      containers:
        - name: hdfs-namenode
          image: flokkkr/ozone:2.1.0
          args: ["hdfs", "namenode"]
        - name: jmxpromo
          image: flokkkr/jmxpromo-sidecar
```

# Flexibility

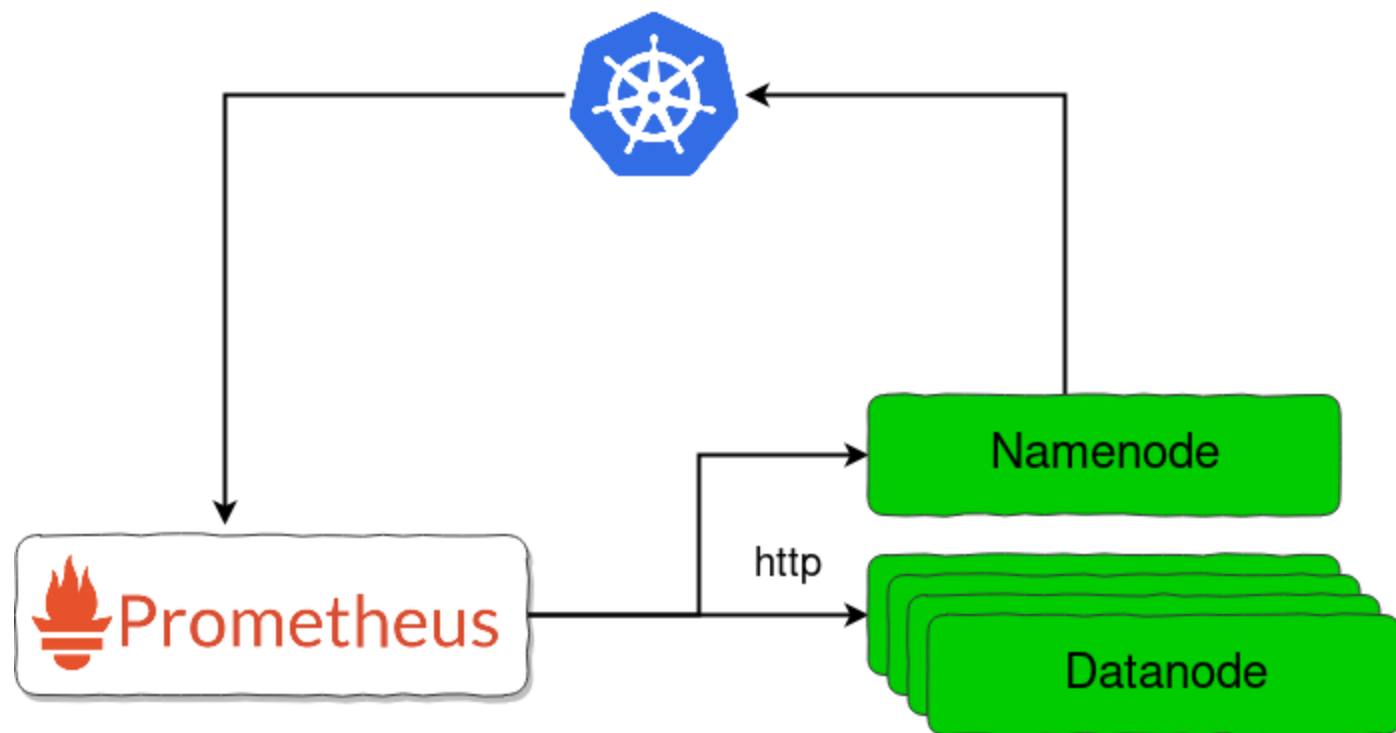




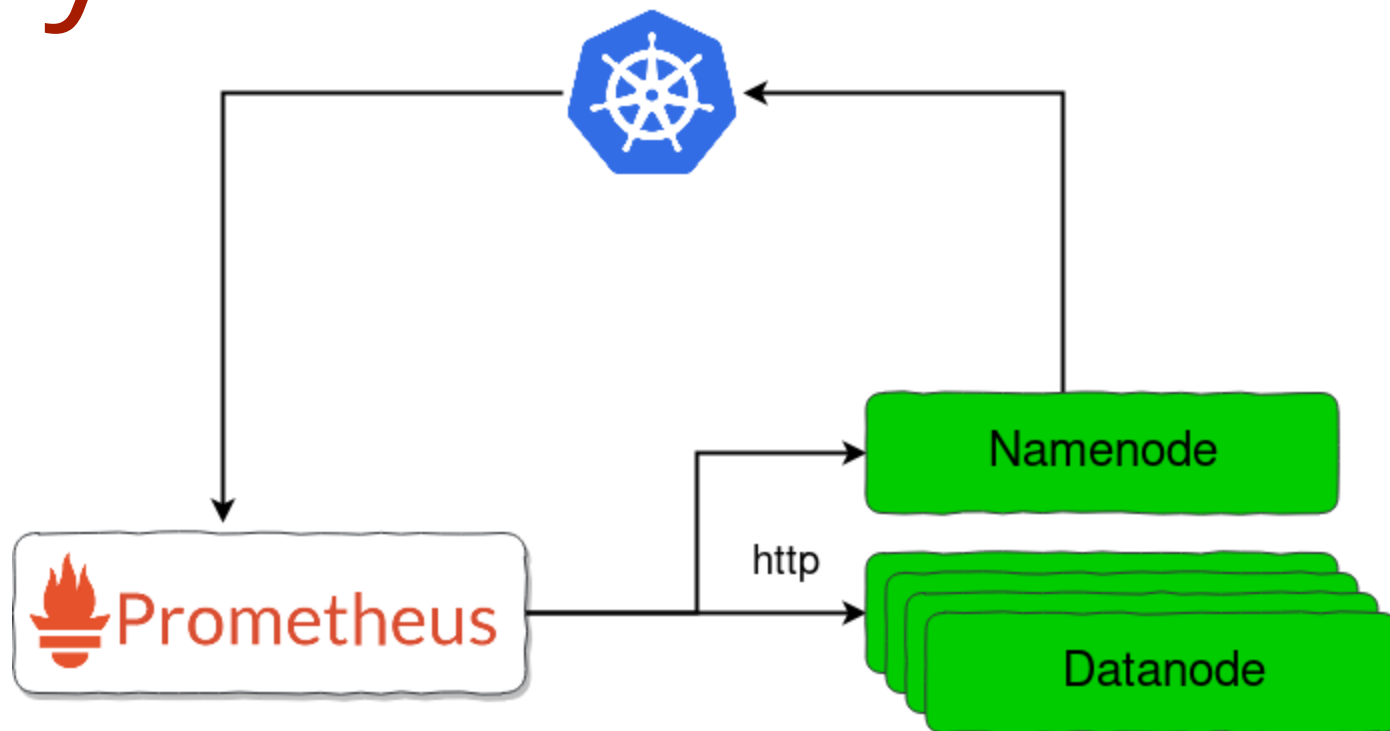
```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  name: ozone-hdfs-namenode
spec:
  serviceName: ozone2-hdfs-namenode
  replicas: 1
  template:
    metadata:
      labels:
        app: ozone
    spec:
      containers:
        - name: hdfs-namenode
          image: flokkkr/ozone:2.1.0
          args: [ "hdfs", "namenode" ]
```


```
apiVersion: apps/v1beta1
kind: StatefulSet
metadata:
  name: ozone-hdfs-namenode
spec:
  serviceName: ozone2-hdfs-namenode
  replicas: 1
  template:
    metadata:
      labels:
        app: ozone
    annotations:
      prometheus.io/scrape: "true"
      prometheus.io/port: "28942"
  spec:
    containers:
      - name: hdfs-namenode
        image: flokkkr/ozone:2.1.0
        args: [ "hdfs", "namenode" ]
```



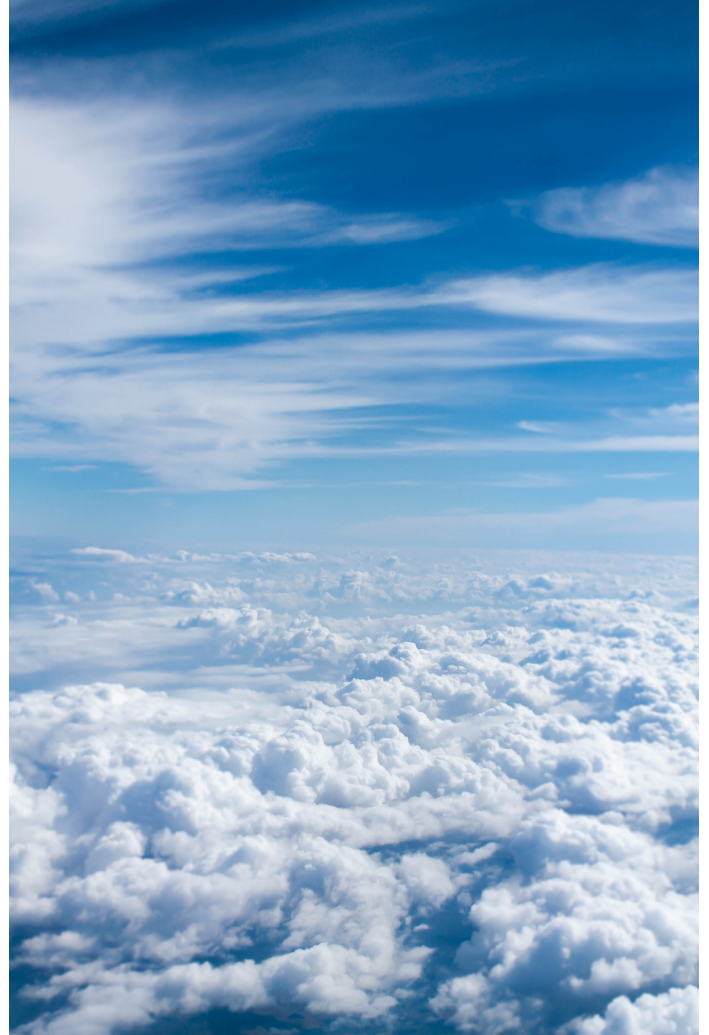


# Ecosystem



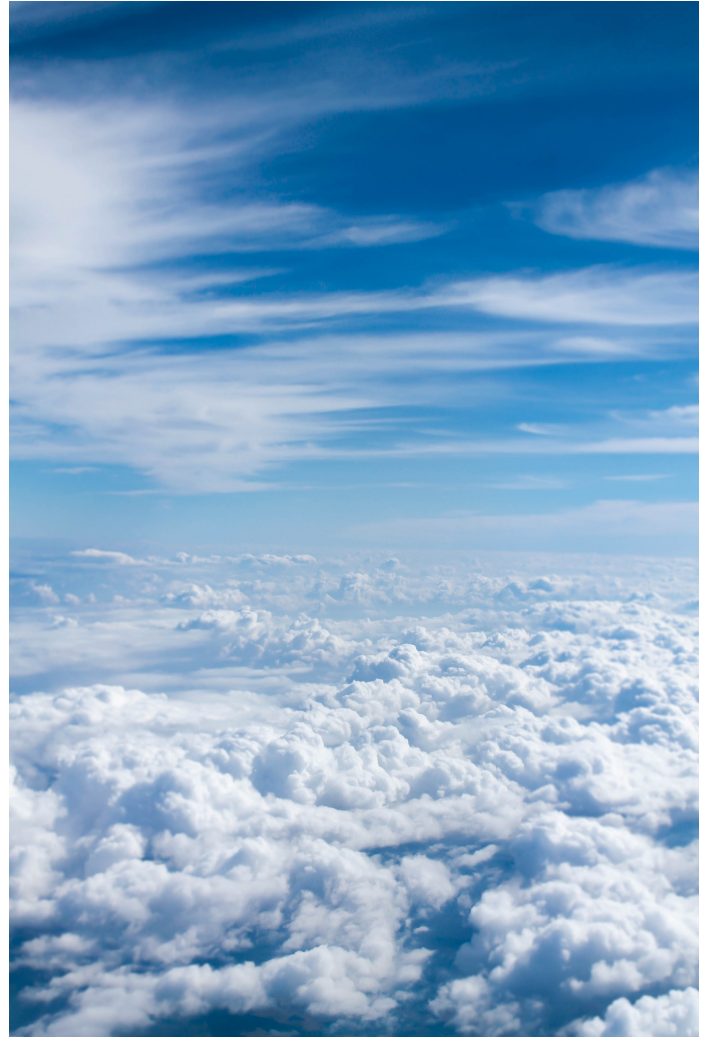
<div> <div>F</div> <div>G</div> </div>	
<b>Configuration management</b> Source Preprocessing On change	configmap helm n/a
<b>Provisioning, Scheduling</b> Multihost support Scheduling Cluster definition Scaling Multi tenancy Failover	CNI k8s helm, yaml yes namespaces yes
<b>Network</b> Intrасervice network DNS Service discovery Data locality Availability of the ports	CNI statefulset DNS no service/ingress

**Is Hadoop  
cloud native?**



# Is Hadoop cloud native?

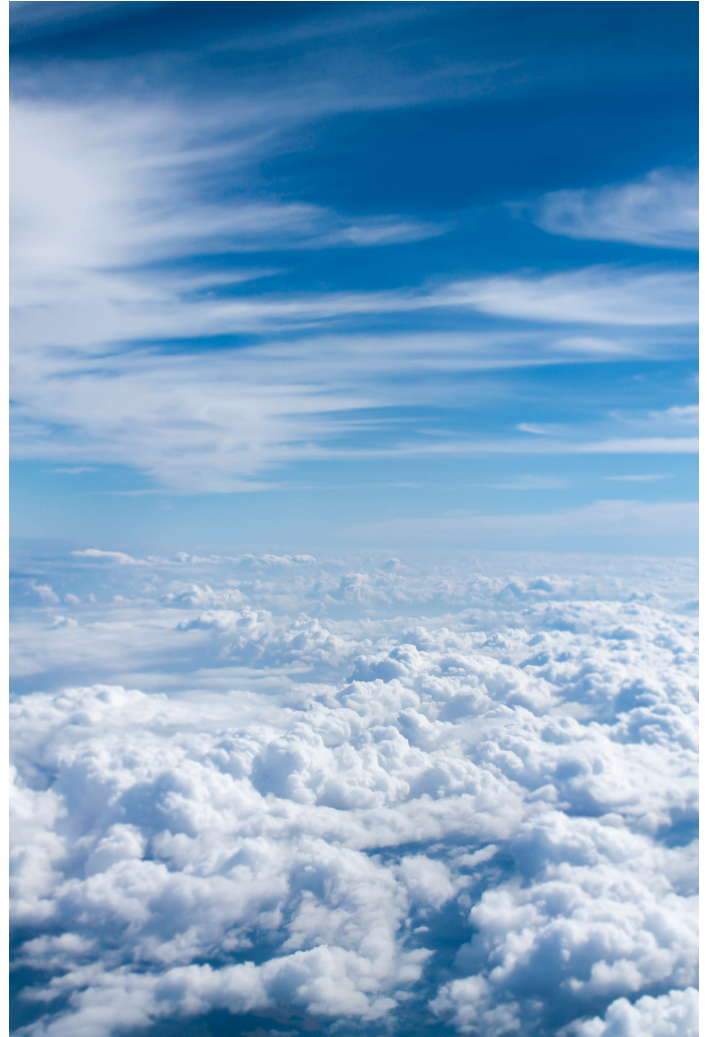
Yes, with small modifications



# Is Hadoop cloud native?

Yes, with small modifications

- DNS/IP handling. Should work
  - without DNS
  - with changing DNS



# Is Hadoop cloud native?

Yes, with small modifications

- DNS/IP handling. Should work
  - without DNS
  - with changing DNS
- More flexible configuration loading

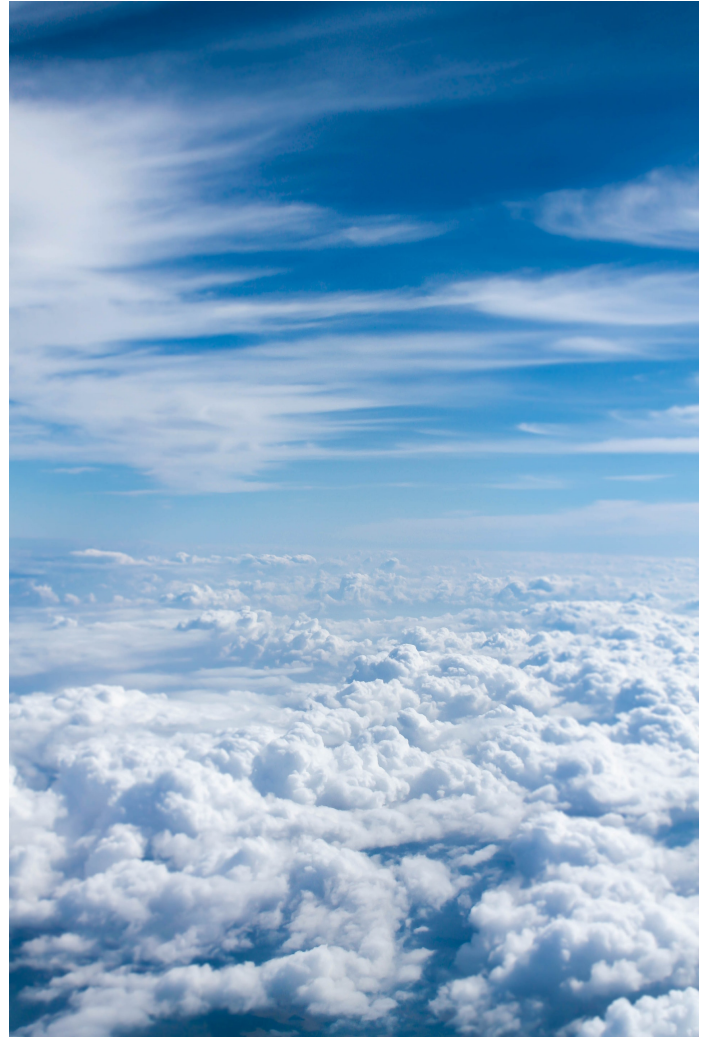




# Is Hadoop cloud native?

Yes, with small modifications

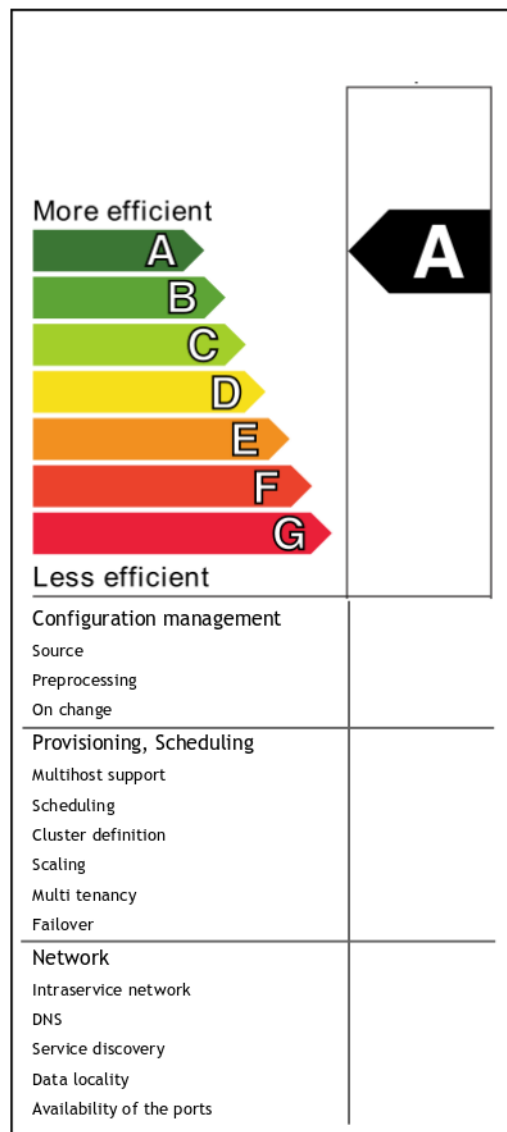
- DNS/IP handling. Should work
  - without DNS
  - with changing DNS
- More flexible configuration loading
- Reverse proxy friendly UI





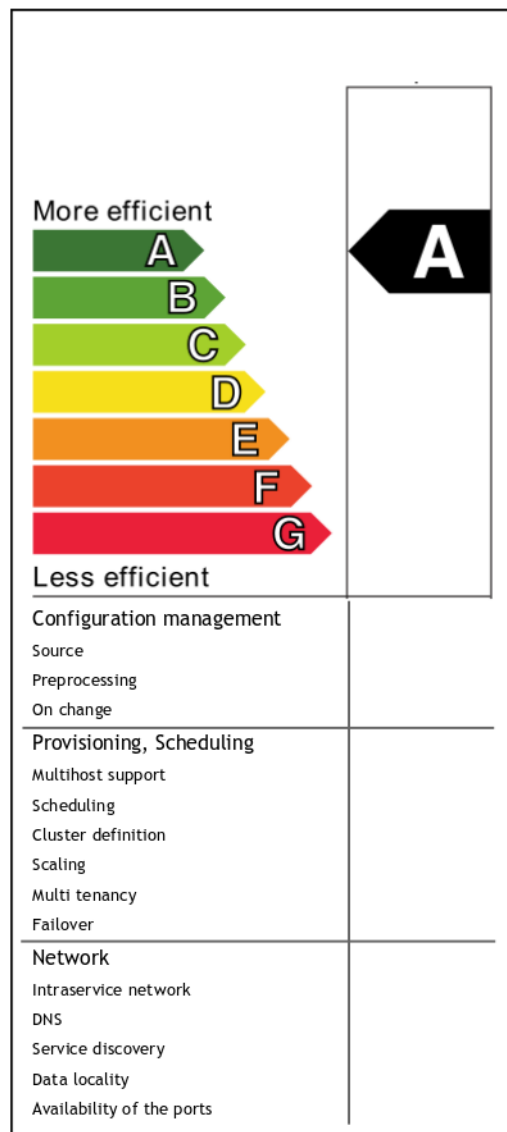
<p>More efficient</p> <p>A</p> <p>B</p> <p>C</p> <p>D</p> <p>E</p> <p>F</p> <p>G</p> <p>Less efficient</p>	<p>A</p>
<p>Configuration management</p> <p>Source</p> <p>Preprocessing</p> <p>On change</p>	
<p>Provisioning, Scheduling</p> <p>Multihost support</p> <p>Scheduling</p> <p>Cluster definition</p> <p>Scaling</p> <p>Multi tenancy</p> <p>Failover</p>	
<p>Network</p> <p>Intraservice network</p> <p>DNS</p> <p>Service discovery</p> <p>Data locality</p> <p>Availability of the ports</p>	

# Summary



# Summary

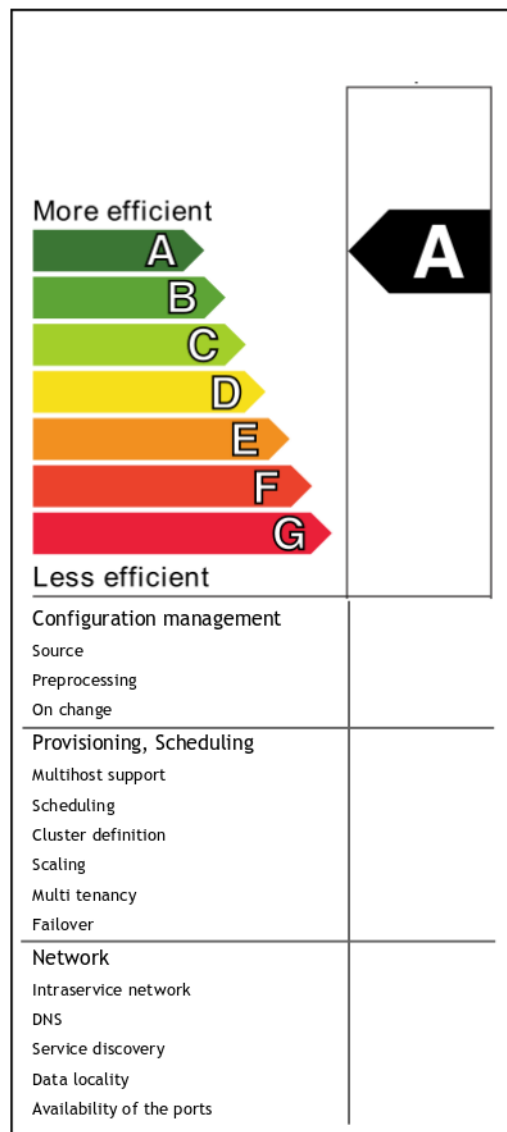
Don't buy without checking the label



# Summary

Don't buy without checking the label

Containerization can help a lot to manage our *Bigdata* clusters



# Summary

Don't buy without checking the label

Containerization can help a lot to manage our *Bigdata* clusters

Hadoop is first class citizen of cloud-native/containerized environments\*

# Q&A

**Márton Elek** @anzix

<https://flokkr.github.io> (bigdata + containers project)

<https://github.com/flokkr> (source)

[elek@apache.org](mailto:elek@apache.org)

# Image credits



Yan Pritzker (CC)



Carrie Cizauskas (CC)