



# HOW TO CLUSTER & FAILOVER (ALMOST) ANYTHING

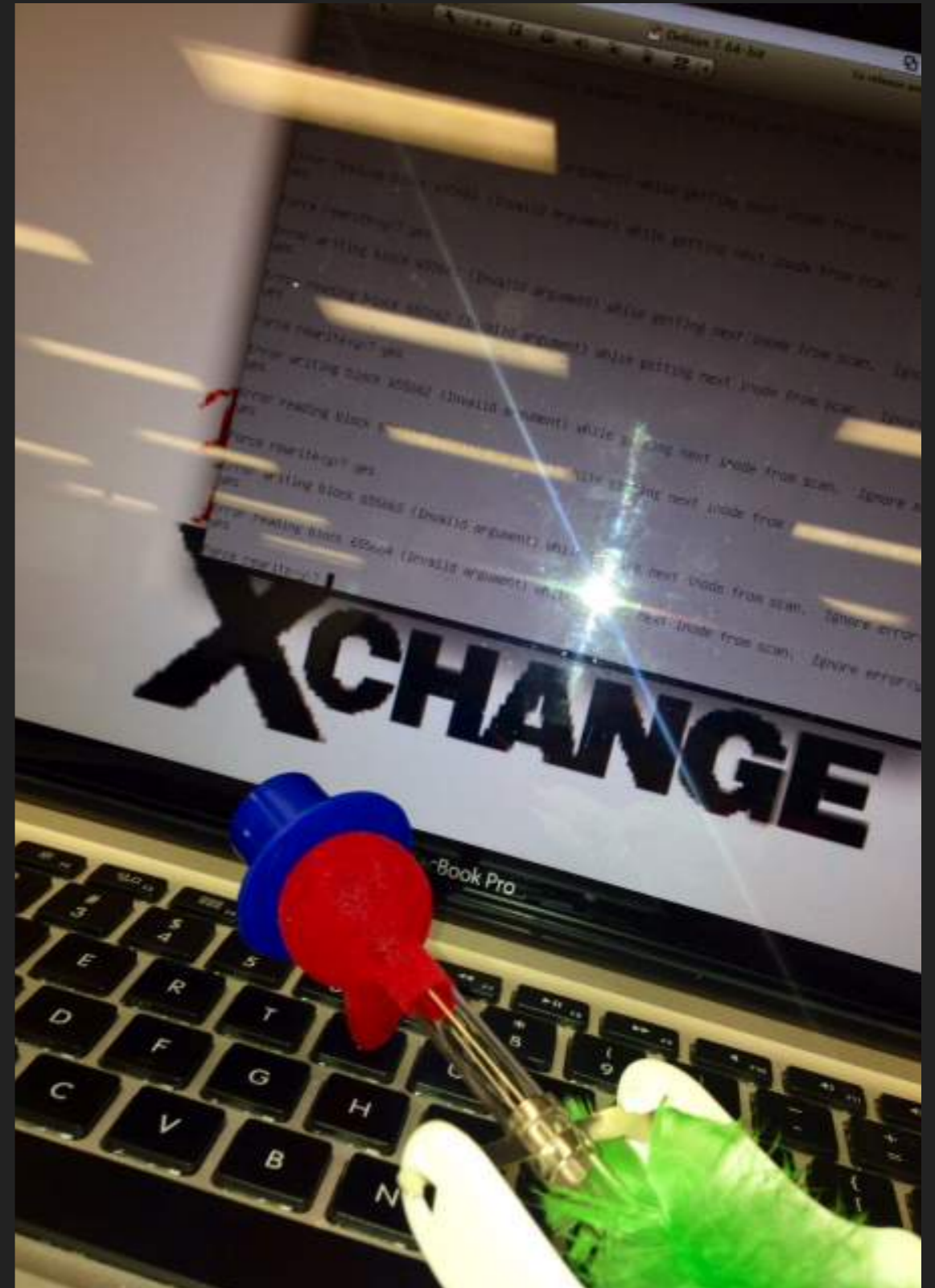
---

## AN INTRO TO PACEMAKER & COROSYNC

Sam McLeod | <https://smcleod.net> | [@s\\_mcleod](#)

## WHO AM I?

- ▶ Ops Lead @Infoxchange
- ▶ 2014 - Researched state of 'enterprise' storage
- ▶ 2014 - 'Enterprise' storage sucks
- ▶ 2014/2015 - Designed & implemented high performance SSD SANs using standard OSS technologies





**Joel Shea** @realist

2★

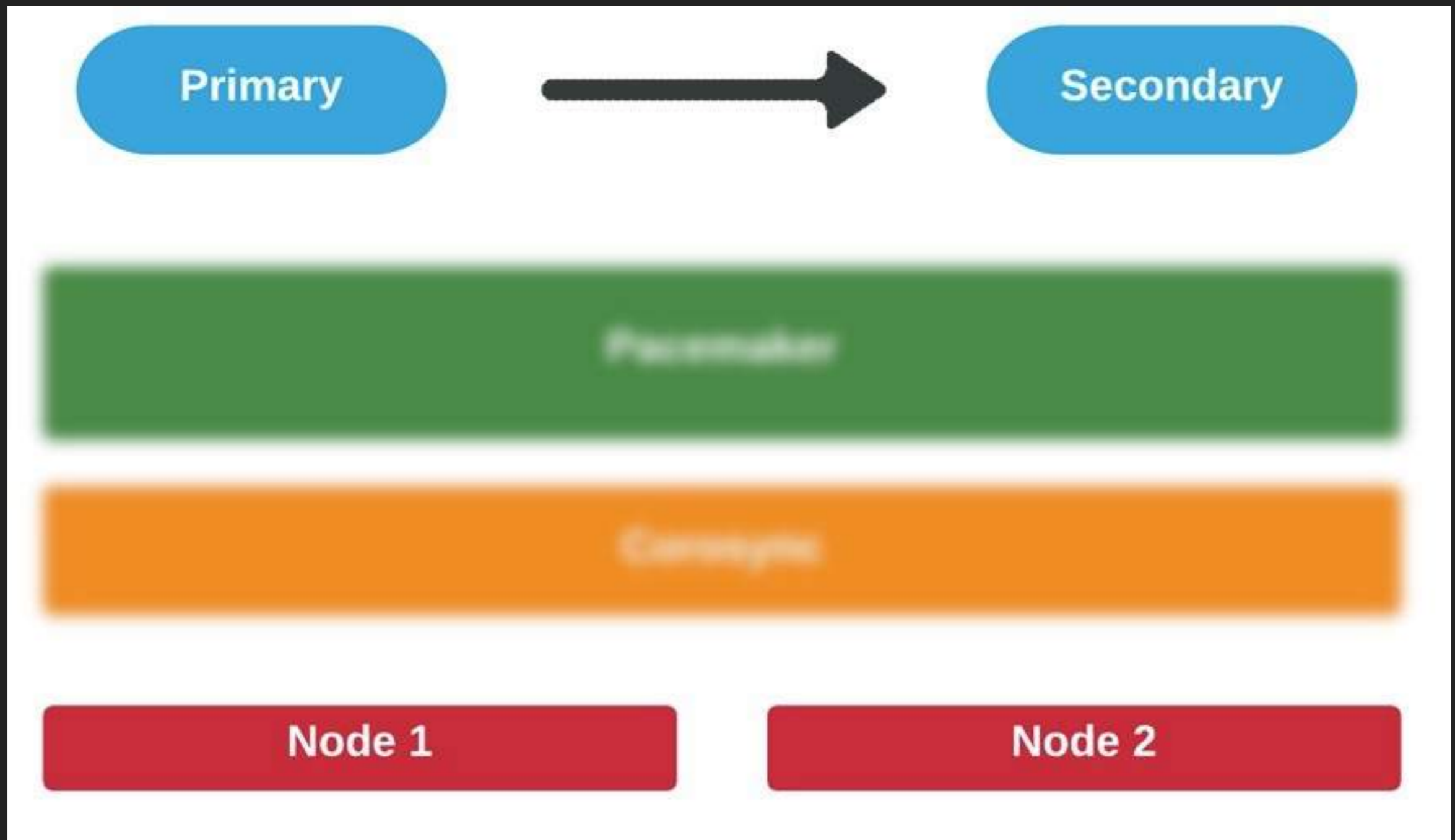
. @s\_mcleod isn't full of shit when talking redundancy, he even has TP failover at home #devopsdays [pic.twitter.com/nA7FRihmeO](https://pic.twitter.com/nA7FRihmeO)



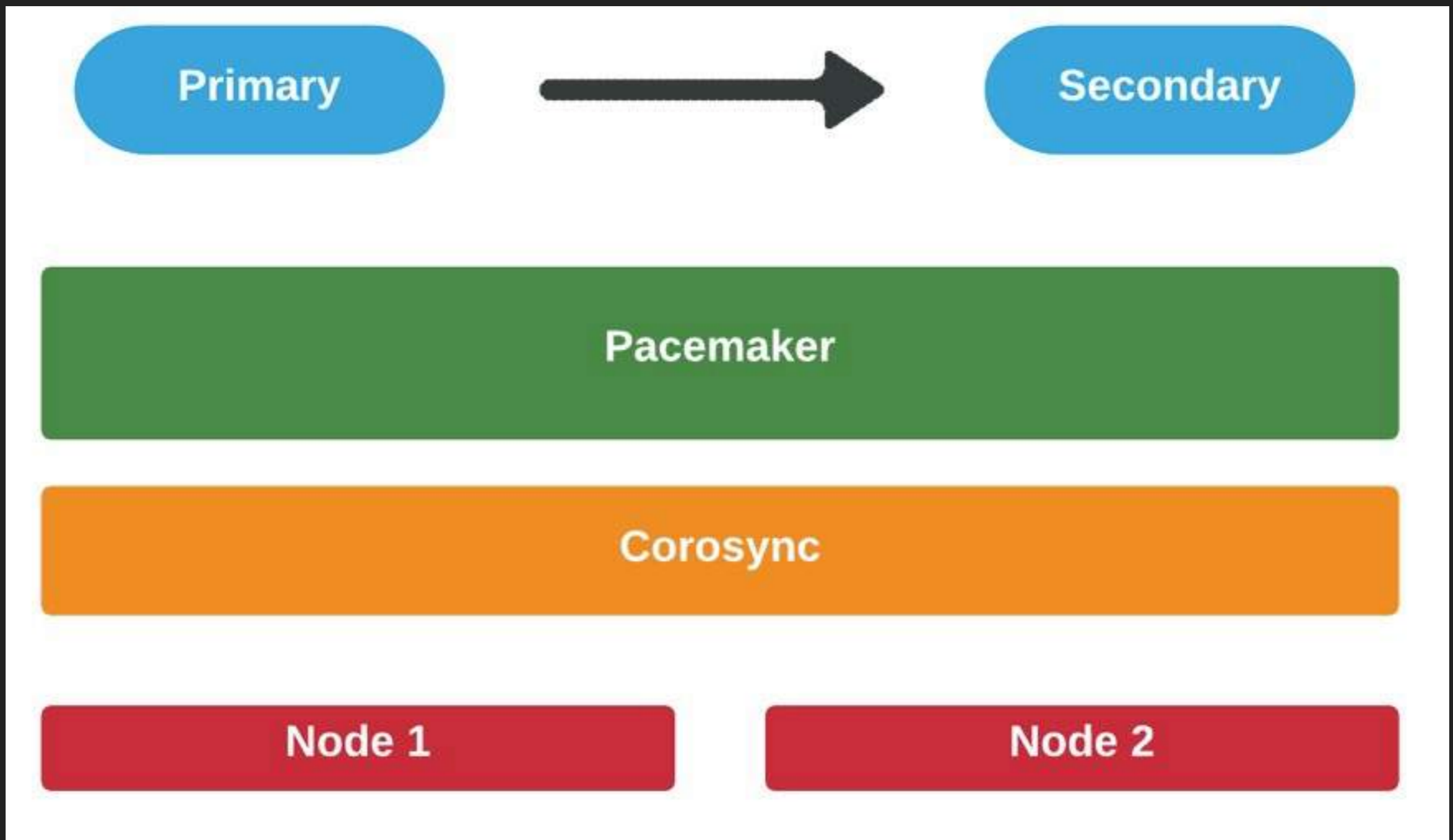
## WHAT CAN YOU CLUSTER WITH COROSYNC+PACEMAKER?

- ▶ Floating IPs
- ▶ Filesystem / Block Storage
- ▶ SystemD / LSB Services
- ▶ Database Servers
- ▶ VMs (Failover and Load Balancing)
- ▶ Network Services

# ACTIVE/STANDBY CLUSTERING



# ACTIVE/STANDBY CLUSTERING





RESOURCE MANAGER

---

**PACEMAKER**

**Pacemaker is the thing that starts and stops services** (like your database or mail server) and contains logic for ensuring both that they're running, and that they're only running in one location (to avoid data corruption).

**Andrew Beekhof**

<http://blog.clusterlabs.org/blog/2010/pacemaker-heartbeat-corosync-wtf/>

Sam McLeod | <https://smcleod.net> | [@s\\_mcleod](https://twitter.com/s_mcleod)



MESSAGING LAYER

---

# COROSYNC

**Think of Corosync as dbus but between nodes.**

Somewhere that any node can throw messages on and know that they'll be received by all its peers.

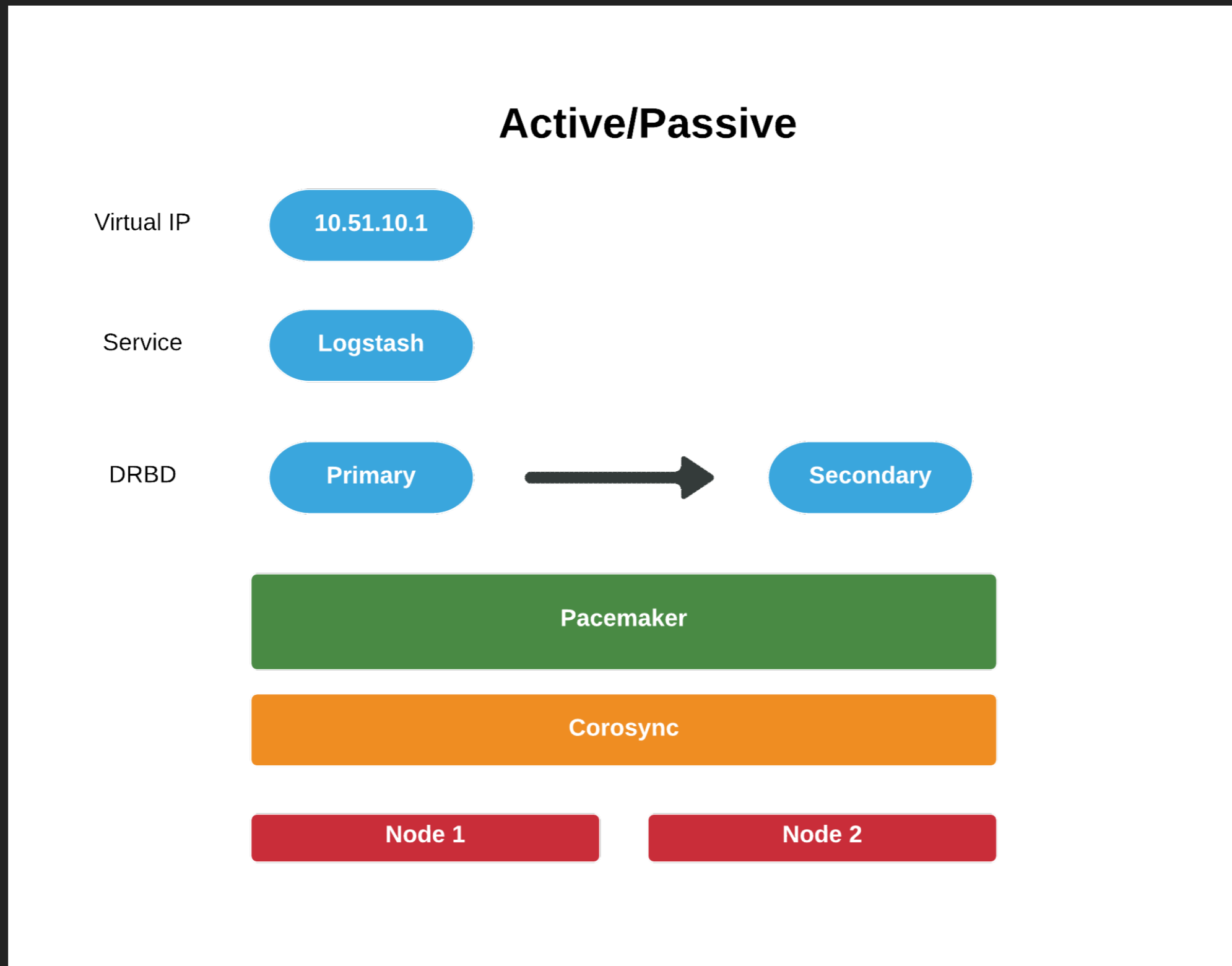
This bus also ensures that everyone agrees who is (and is not) connected to the bus and tells Pacemaker when that list changes.

**Andrew Beekhof**

<http://blog.clusterlabs.org/blog/2010/pacemaker-heartbeat-corosync-wtf/>

Sam McLeod | <https://smcleod.net> | [@s\\_mcleod](https://twitter.com/s_mcleod)

# FLOATING IP, APPLICATION SERVICE & FILESYSTEM CLUSTER





OPEN CLUSTER FRAMEWORK (OCF)

---

# RESOURCES

To avoid teaching Pacemaker about every possible service that people might want to make highly available which we call Resource Agents. **Any series of command-line actions can be easily turned into a resource agent by adding them to an existing template.**

**Andrew Beekhof**

<http://blog.clusterlabs.org/blog/2010/pacemaker-heartbeat-corosync-wtf/>

Sam McLeod | <https://smcleod.net> | [@s\\_mcleod](https://twitter.com/s_mcleod)

## [HTTPS://GITHUB.COM/CLUSTERLABS/RESOURCE-AGENTS](https://github.com/clusterlabs/resource-agents)

### # pcs resource standards

ocf  
lsb  
service  
systemd  
stonith

### # pcs resource agents

CTDB  
Delay  
Dummy  
Filesystem  
IPAddr2  
LVM  
MailTo

### # pcs resource describe IPAddr2

ocf:heartbeat:IPAddr2 – Manages virtual IPv4 and IPv6 addresses  
(Linux specific version)

# RESOURCE COLOCATION

You can specify that two (or more) resources must be started at the same location (node), in doing so you also define their ordering.

```
# pcs constraint colocation add virtualIP with NFS INFINITY
```

```
# pcs constraint colocation add DRBD with RAID INFINITY
```

# RESOURCE ORDERING

Explicit ordering can be set on resources regardless of location

```
# pcs constraint order set virtualIP NFS
```



PACEMAKER + OCF

---

# RESOURCE MONITORING

## RESOURCE MONITORING

The frequency that resources are monitored for health as well as how long to wait for them to start / stop and what to do if they fail to perform their action can be added as operational configuration.

```
# pcs resource create myvip IPAddr2 ip=10.51.10.150 cidr_netmask=23 \  
op monitor interval=30s timeout=30s \  
op start interval=0 timeout=15s on-fail=restart \  
op stop interval=0 timeout=15s on-fail=block
```

AVOIDING SPLIT BRAIN

---

# FENCING & STONITH

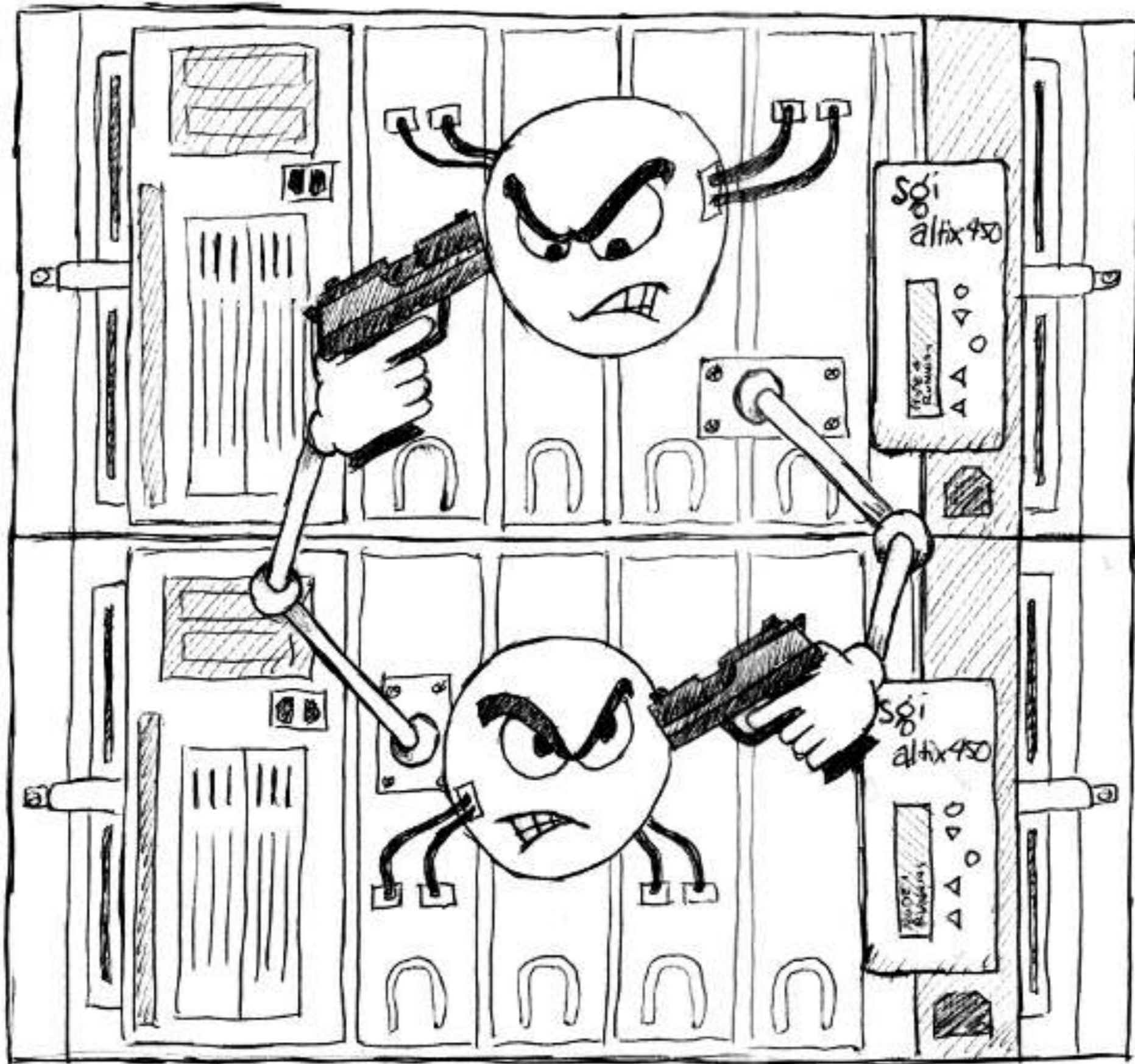
**STONITH = SHOOT THE OTHER NODE IN THE HEAD**

What if a node can't ensure the state of its peer?

The cluster is at risk of split brain.

The peers state is ensured, usually by remote power off or restart.





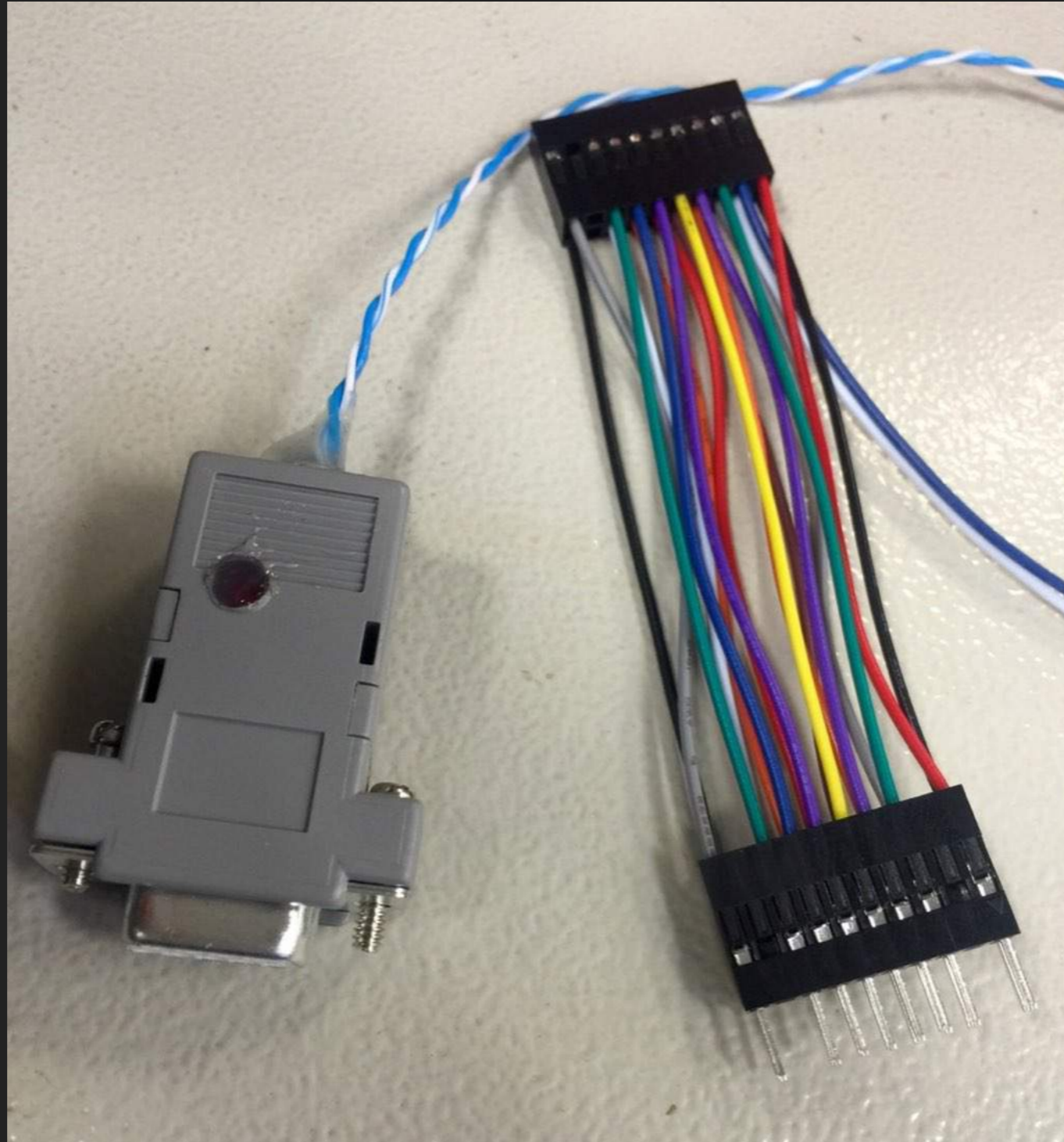
DON'T ANYBODY MOVE ...

# [HTTPS://GITHUB.COM/CLUSTERLABS/RESOURCE-AGENTS](https://github.com/clusterlabs/resource-agents)

### # pcs stonith list

fence\_apc - Fence agent for APC over telnet/ssh  
fence\_bladecenter - Fence agent for IBM BladeCenter  
fence\_brocade - Fence agent for HP Brocade over telnet/ssh  
fence\_cisco\_mds - Fence agent for Cisco MDS  
fence\_cisco\_ucs - Fence agent for Cisco UCS  
fence\_compute - Fence agent for nova compute nodes  
fence\_drac5 - Fence agent for Dell DRAC CMC/5  
fence\_hpblade - Fence agent for HP BladeSystem  
fence\_idrac - Fence agent for IPMI  
fence\_ifmib - Fence agent for IF MIB  
fence\_ilo - Fence agent for HP iLO  
fence\_ipmilan - Fence agent for IPMI  
fence\_kdump - Fence agent for use with kdump  
fence\_rhevm - Fence agent for RHEV-M REST API  
fence\_scsi - Fence agent for SCSI persistent reservation  
fence\_virt - Fence agent for virtual machines  
fence\_vmware\_soap - Fence agent for VMWare over SOAP API  
fence\_xvm - Fence agent for virtual machines







CLUSTER DEMO

---

# PRODUCTION FAILOVER

rcd\_serial (stonith\_pm-san6)  
running on: pm-san5

rcd\_serial (stonith\_pm-san5)  
running on: pm-san6

online  
pm-san5

pacemaker:ping (ping\_gateway)  
Clone Set (ping\_gateway-clone)  
running on:  
pm-san5 / ping: 1  
pm-san6 / ping: 1

online  
pm-san6

linbit:drbd (r0)  
Master/Slave Set (drbd\_r0)  
primary on:  
pm-san5  
secondary on:  
pm-san6

linbit:drbd (r1)  
Master/Slave Set (drbd\_r1)  
primary on:  
pm-san6  
secondary on:  
pm-san5

linbit:drbd (r2)  
Master/Slave Set (drbd\_r2)  
primary on:  
pm-san6  
secondary on:  
pm-san5

IPAddr2 (ip\_r0 / 10.51.40.75)  
running on: pm-san5

IPAddr2 (ip\_r1 / 10.51.40.85)  
running on: pm-san6

IPAddr2 (ip\_r2 / 10.51.40.95)  
running on: pm-san6

iSCSITarget (iscsi\_target\_r0)  
running on: pm-san5

iSCSITarget (iscsi\_target\_r1)  
running on: pm-san6

iSCSITarget (iscsi\_target\_r2)  
running on: pm-san6

iSCSILogicalUnit (iscsi\_lun\_r0)  
running on: pm-san5

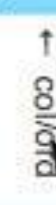
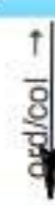
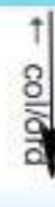
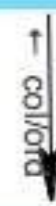
iSCSILogicalUnit (iscsi\_lun\_r1)  
running on: pm-san6

iSCSILogicalUnit (iscsi\_lun\_r2)  
running on: pm-san6

anything (iscsi\_conf\_r0)  
running on: pm-san5

anything (iscsi\_conf\_r1)  
running on: pm-san6

anything (iscsi\_conf\_r2)  
running on: pm-san6





CLUSTER DEMO

---

# TWO NODE STORAGE CLUSTER



# Two Node Cluster

Design  
Performance  
Failover Demo

Pause video as required



CLUSTER CREATION

---

**COROSYNC + PACEMAKER**

Virtual IP

10.51.10.150

Service

Logstash

Pacemaker

Corosync

dev-samm-01

dev-samm-02

Virtual IP

10.51.10.150

Service

Logstash

Pacemaker

Corosync

dev-sa

Bang!

dev-samm-02

```
yum install  
pcs corosync pacemaker fence-agents-all
```

## /ETC/COROSYNC/COROSYNC.CONF

```
totem {
  version: 2
  secauth: off
  cluster_name: dev-samm

  interface {
    ringnumber: 0
    bindnetaddr: 10.51.11.231
    broadcast: yes
  }
}

quorum {
  provider: corosync_votequorum
  two_node: 1
}
```

```
nodelist {
  node {
    ring0_addr: 10.51.11.231
    nodeid: 1
  }
  node {
    ring0_addr: 10.51.11.232
    nodeid: 2
  }
}
```

---

## /ETC/COROSYNC/SERVICE.D/PCMK

```
service {
  name: pacemaker
  ver: 1
}
```

# START SERVICES

```
systemctl start corosync
```

```
systemctl start pacemaker
```

- ▶ Starting Corosync does not start Pacemaker, it just sets up the quorum and messaging interfaces needed by the rest of the stack.
- ▶ The Pacemaker daemon must be started after Corosync

## CREATE A RESOURCE FOR A FLOATING IP

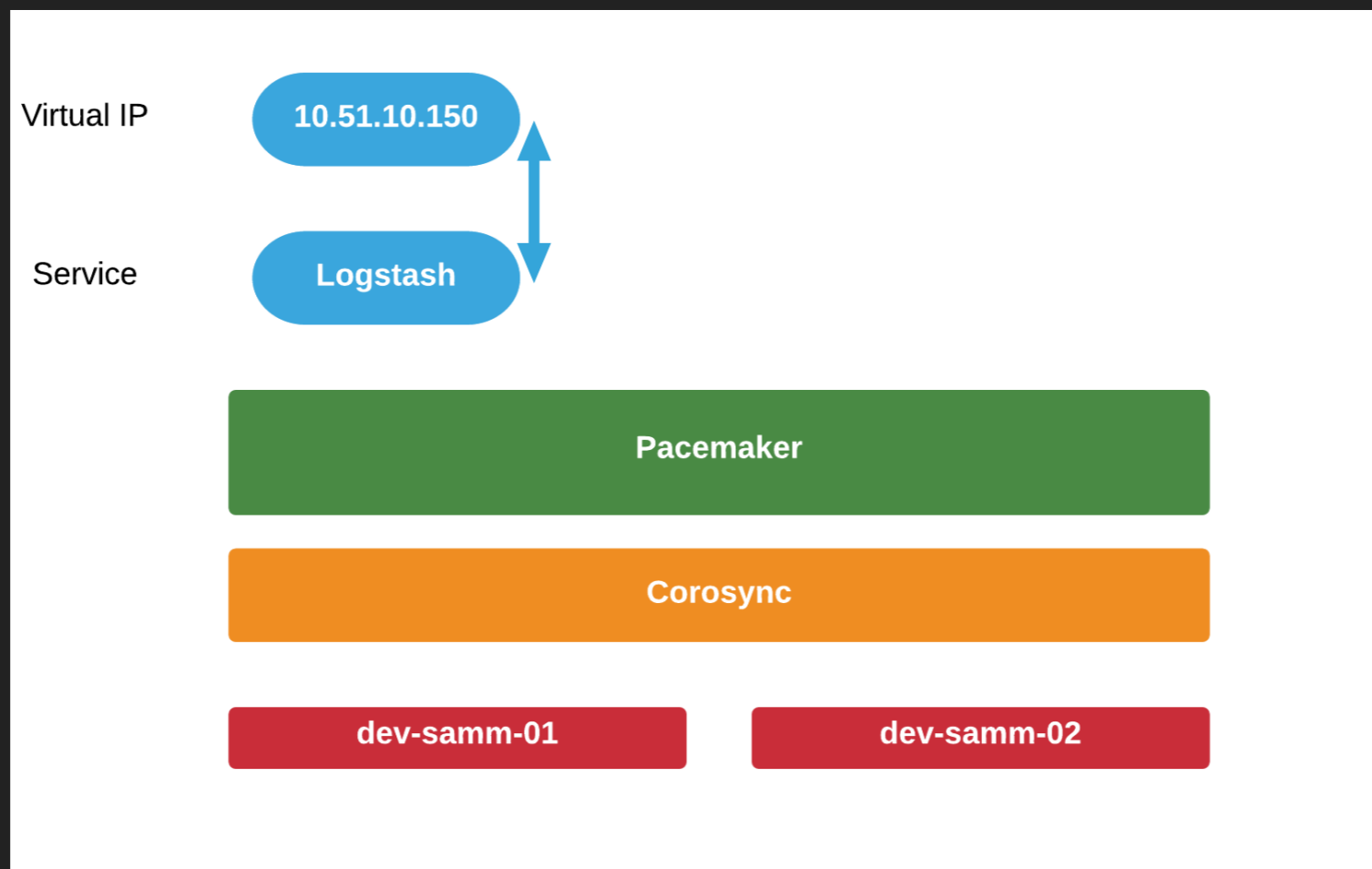
```
# pcs resource create myvip IPAddr2 ip=10.51.10.150  
cidr_netmask=23
```

## CREATE A RESOURCE FOR A SYSTEMD SERVICE

```
# pcs resource create logstash systemd:logstash
```

# CREATE A CLUSTER CONSTRAINT THAT COLOCATES THE SERVICE AND VIRTUAL IP

```
# pcs constraint colocation add myvip with logstash INFINITY
```



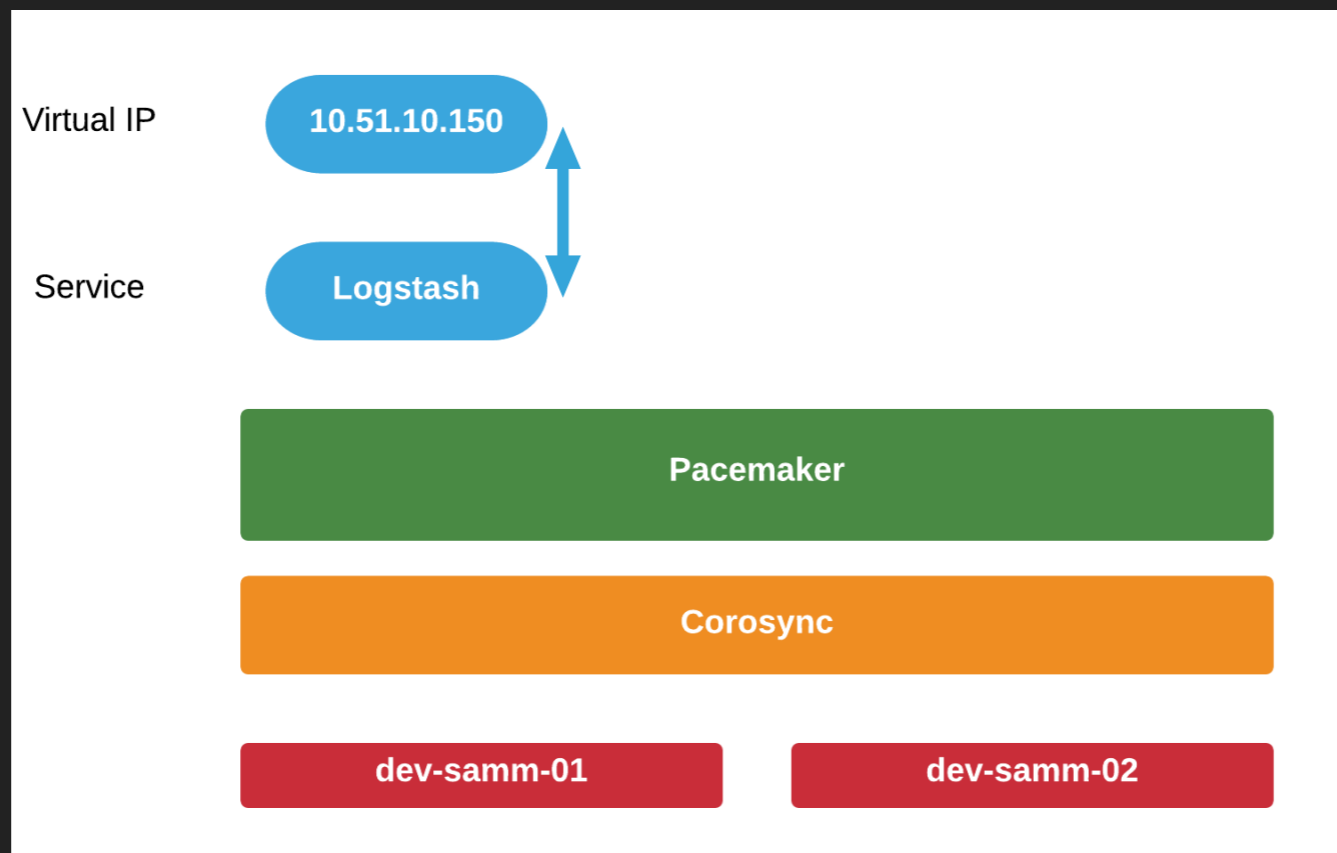
# VIEW THE CLUSTER STATUS

```
# pcs status
```

```
Online: [ dev-samm-01, dev-samm-02 ]
```

```
myvip (ocf::heartbeat:IPaddr2): Started dev-samm-01
```

```
logstash (systemd:logstash): Started dev-samm-01
```



Reminder

**YOU MUST USE FENCING  
IN PRODUCTION!**

# TEST BASIC FAILOVER

```
# pcs cluster standby dev-samm-01
```

```
# pcs status
```

```
Node dev-samm-01: standby
```

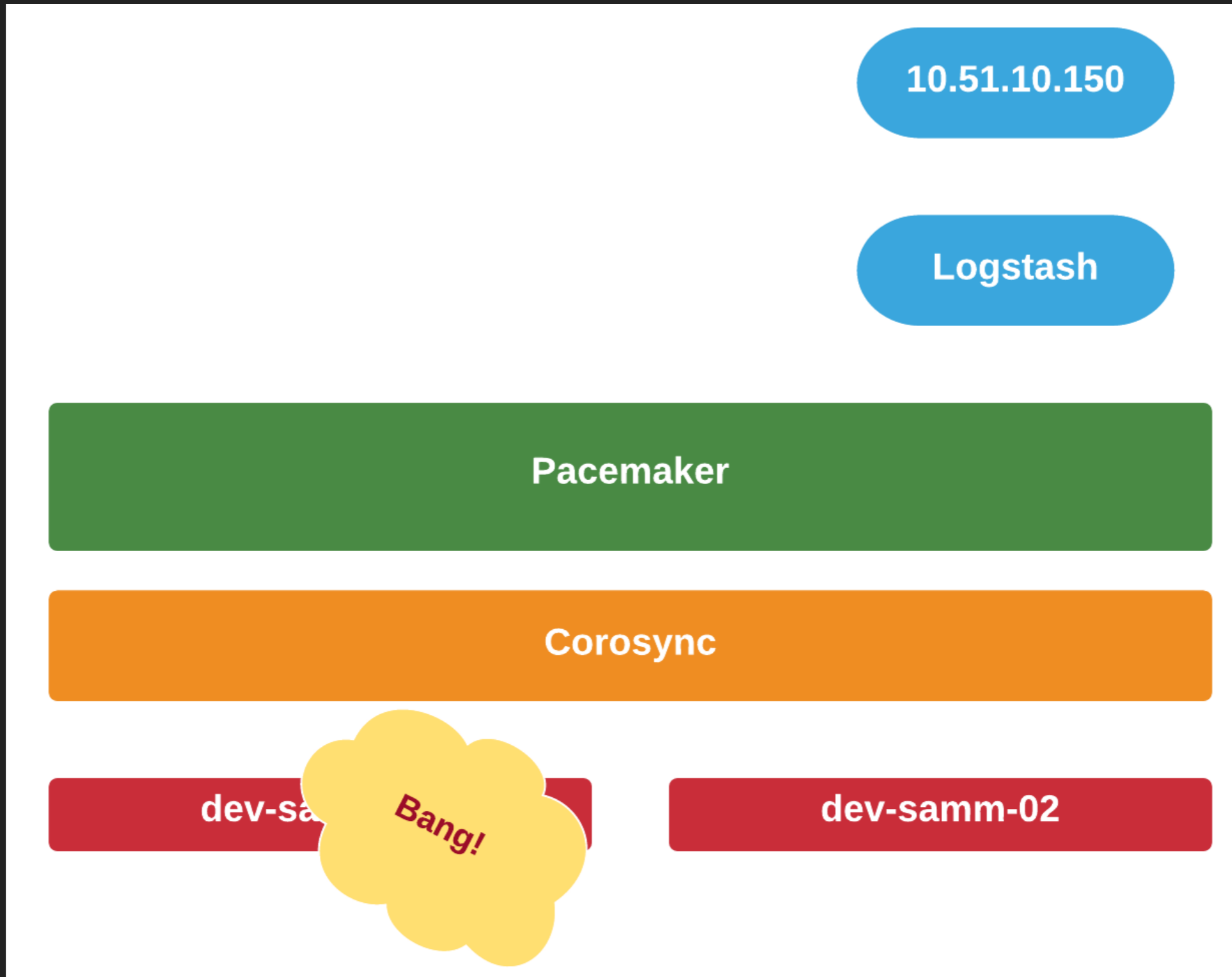
```
Online: [ dev-samm-02 ]
```

```
Full list of resources:
```

```
myvip (ocf::heartbeat:IPaddr2): Started dev-samm-02
```

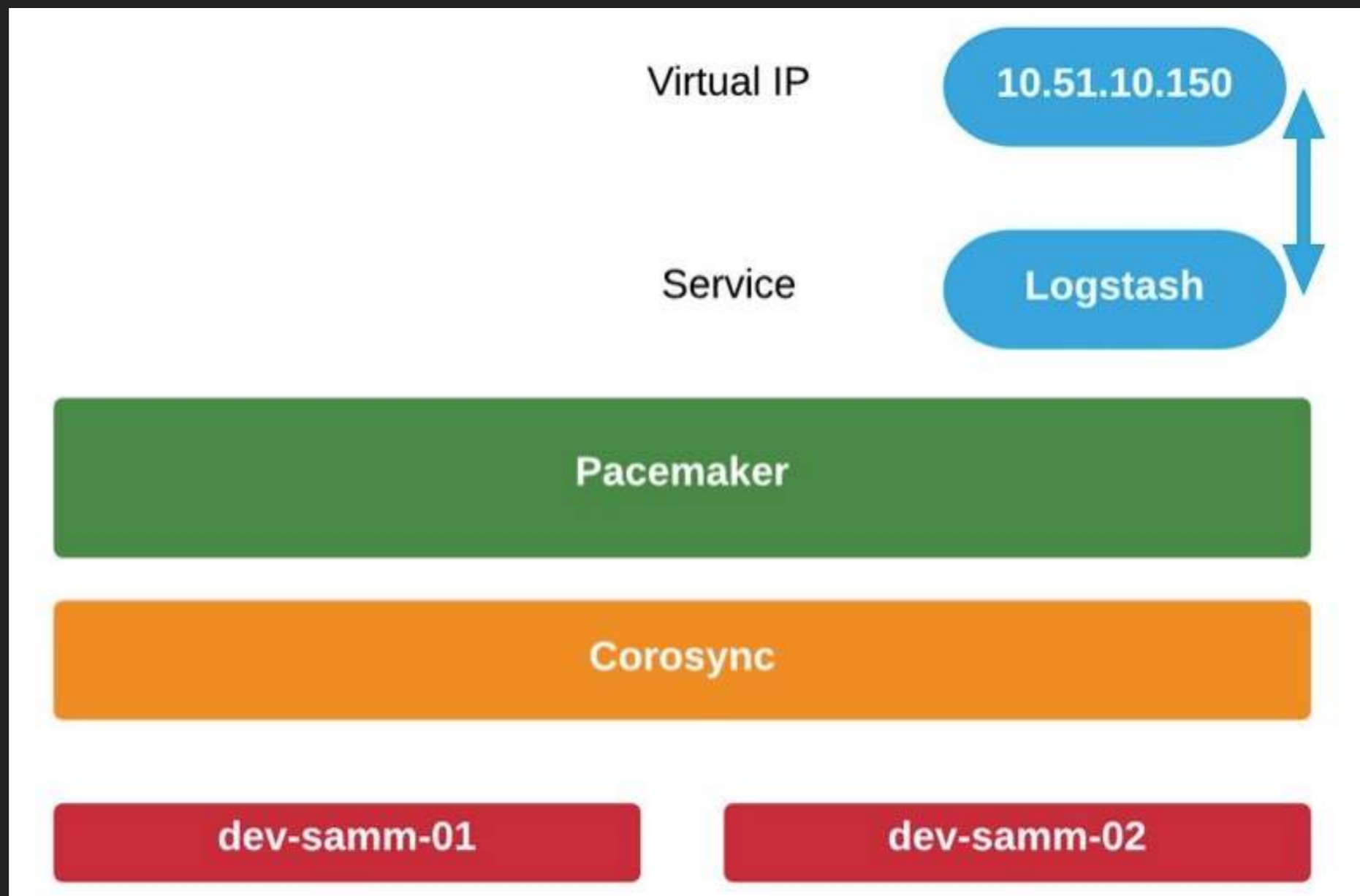
```
logstash (systemd:logstash): Started dev-samm-02
```

# RESOURCE CREATION - RESOURCE FAILOVER



# TEST BASIC FAILOVER

```
# pcs cluster online dev-samm-01
```



# WHERE DO YOU START?

CLUSTERS FROM SCRATCH - PACEMAKER 1.1 FOR COROSYNC 2.X AND PCS

---

**CLUSTERLABS.COM/DOC**

THANK YOU

---

## USEFUL LINKS

- ▶ <http://clusterlabs.com/doc>
- ▶ [https://alteeve.ca/w/High-Availability\\_Clustering\\_in\\_the\\_Open\\_Source\\_Ecosystem](https://alteeve.ca/w/High-Availability_Clustering_in_the_Open_Source_Ecosystem)
- ▶ <https://github.com/ClusterLabs/pacemaker/blob/master/doc/pcs-crmsh-quick-ref.md>
- ▶ <https://github.com/ClusterLabs/resource-agents>
- ▶ #clusterlabs on FreeNode



Link to these slides

Sam McLeod | <https://smcleod.net> | [@s\\_mcleod](#)