

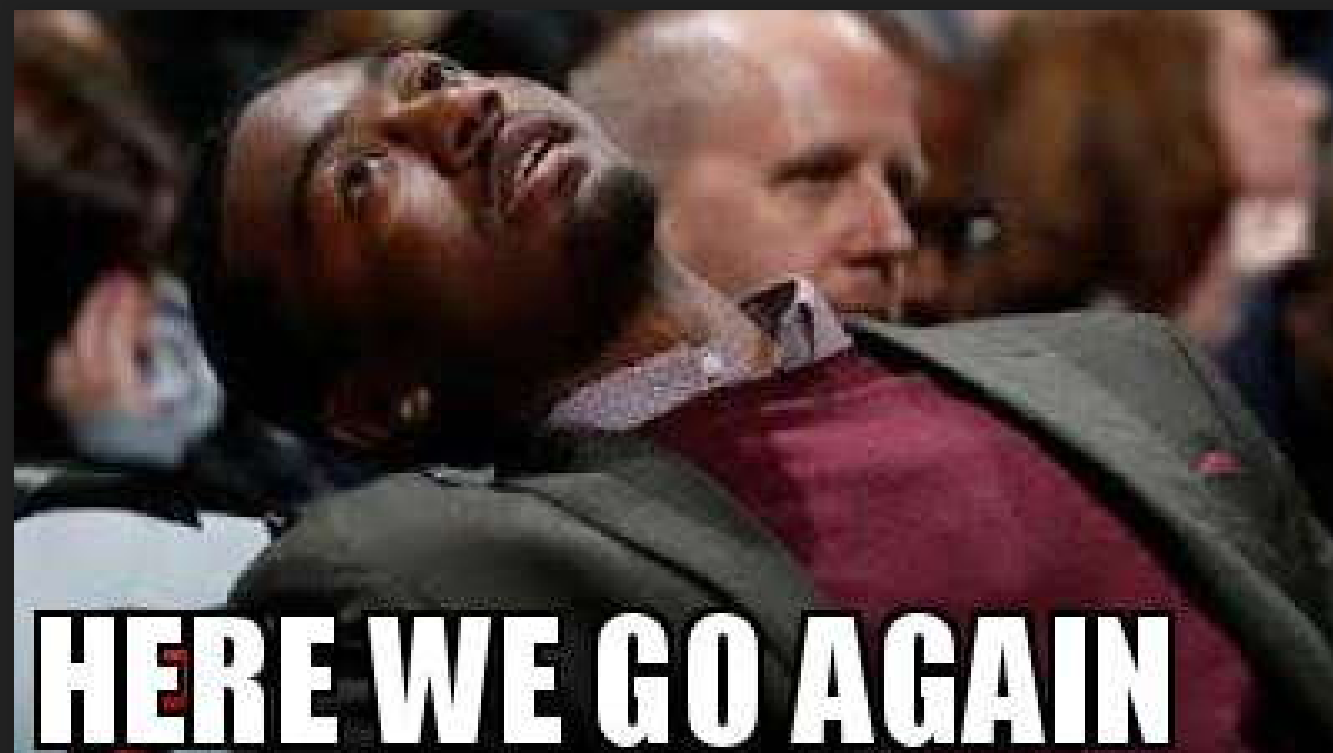
# PAF

...Another brick in the HA wall

# WHO AM I ?

- Jehan-Guillaume de Rorthais
- aka. ioguix
- using PostgreSQL since 2007
- involved in PostgreSQL community since 2008
- @dalibo since 2009

# High Availability



# MENU

- Quick intro about HA
- Quick intro about Pacemaker
- Why PAF ?
- PAF abilities

# HA IN SHORT

- Part of the Business Continuity Planning
- In short: double everything
- should it include automatic failover ?

# AUTO FAILOVER: TECH

Hard to achieve:

- how to detect a real failure?
- why the master doesn't answer?
- is it under high load? switched off?
- is it a network issue? hick up?
- how to avoid split brain?

# BUILDING AUTO FAILOVER

- many issues to understand
- solutions: quorum, fencing, watchdog, ...
- complex setup
- complex maintenances
- document, document, document
- test, test, test

If you don't have time, don't do auto failover (almost).

# QUORUM

- resources run in the cluster partition hosting the greater number of nodes
- useful on network split
- ...or when you require a minimal number of node alive
- based on vote



# FENCING

- ability to poweroff or reboot any node of your cluster
- the definitive solution to know the real state of an unresponsive node
- hardware fencing (smart PDU, UPS, IPMI)
- IO fencing (SAN, network)
- virtual fencing (libvirt, xen, vbox, ...)
- software: do not rely on it (eg. ssh)
- meatware

Really, do it. Do not think you are safe without it.

# WATCHDOG

- feed your local dog or it will kill your node
- either hardware or software (cf. softdog)
- self-fencing (suicide) on purpose
- auto-self-fencing when node is unresponsive

# Pacemaker

Will assimilate your resource...

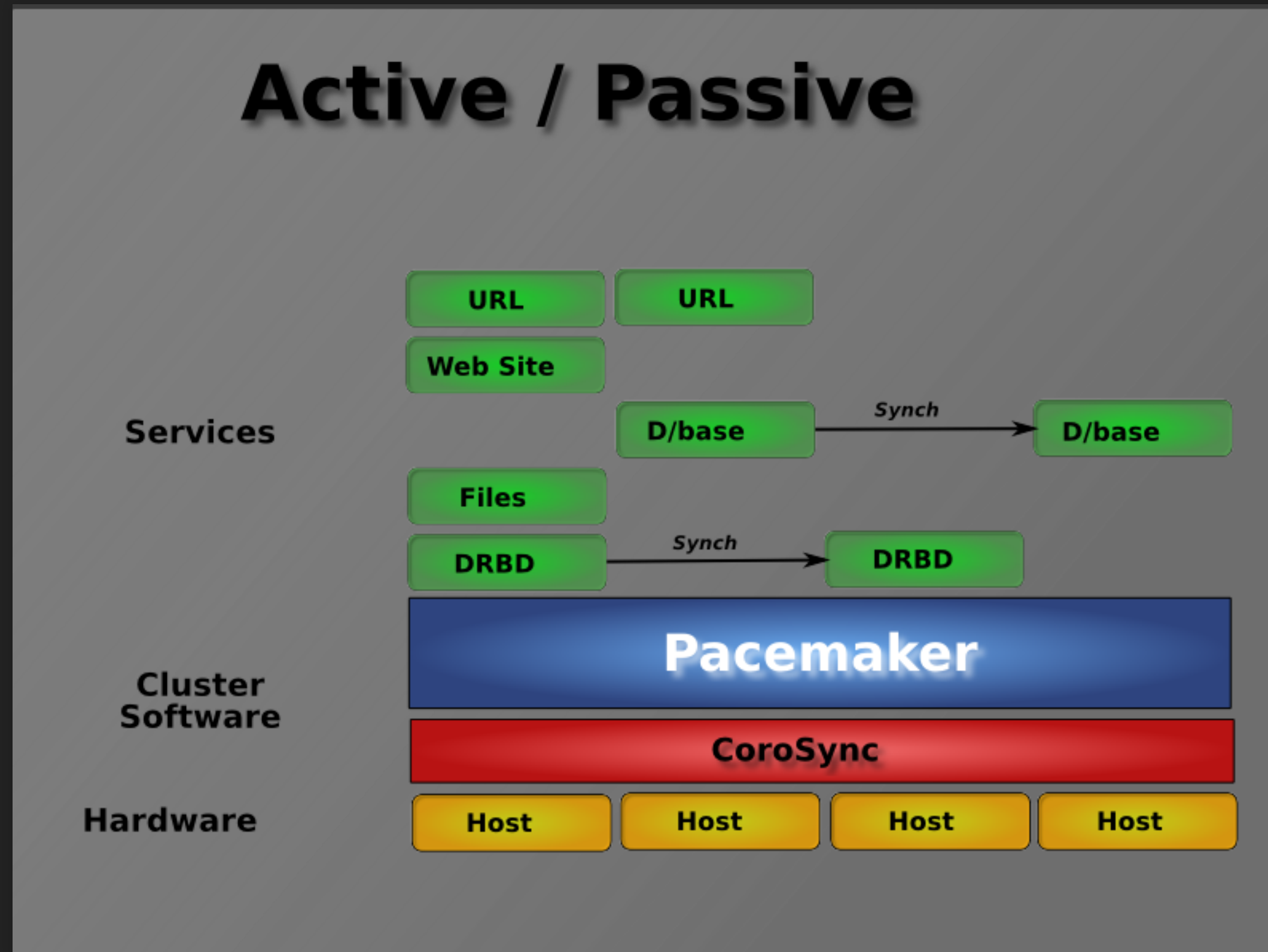
Resistance is futile.

(T.N.: 'service' eq 'resource')

# PACEMAKER IN SHORT

- is a "Cluster Resource Manager"
- support fencing, quorum and watchdog
- multi-resource, dependencies, resources order, constraints, rules, etc
- **Resource Agents** are the glue between the CRM and the services
- RA can be stateless or multi-state
- RA API: script OCF, upstart, systemd, LSB

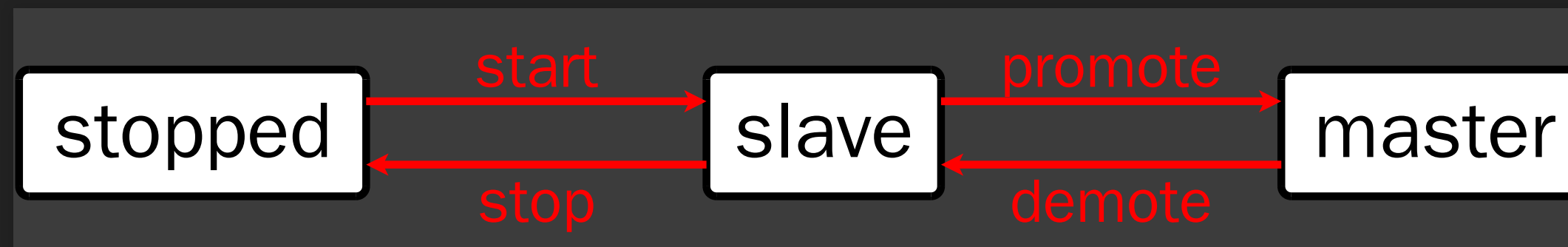
# PACEMAKER ARCHITECTURE





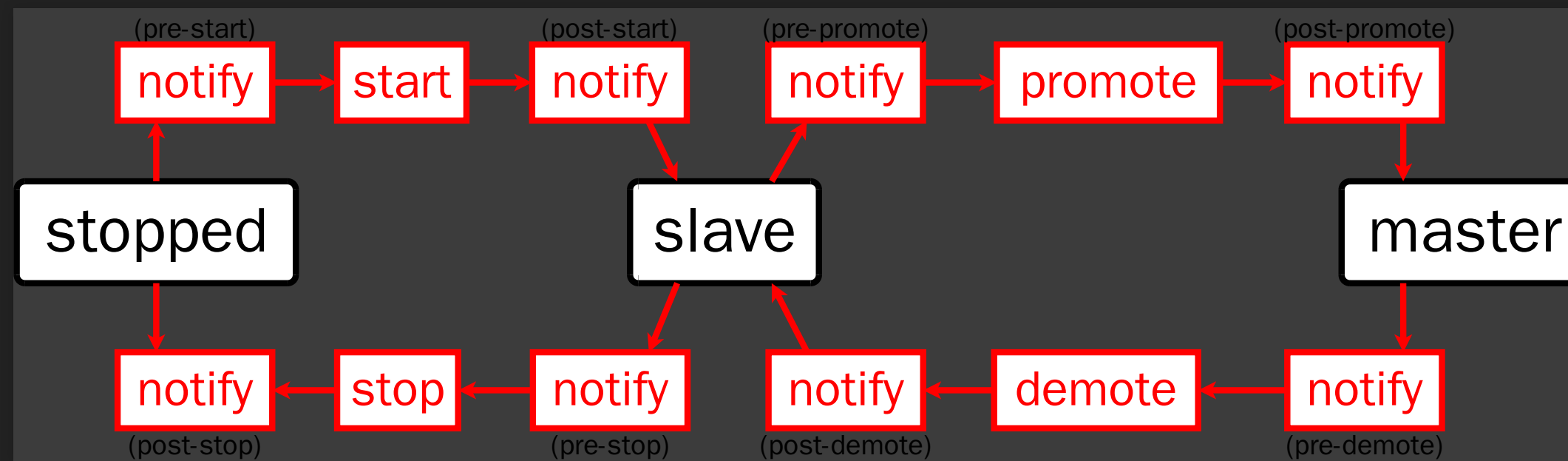
# CRM MECHANISM

- kind of automate
- 4 states: stopped, started, slave or master
- the CRM compute transitions between two states
- only **ONE** CRMd is handling the whole cluster: the DC
- minimal actions API (eg. systemd): start, stop, monitor(status)
- extended actions API (OCF): start, stop, promote, demote, monitor, **notify**
- for a multi-state resource:



# NOTIFY ACTION

- only available with OCF resource agents
- triggered before and after actions
- triggered on **ALL** node
- action wait for all pre-notify feedback to run
- next actions wait for all post-notify feedback to run
- allows the resource agent to run specific service code



# NOTIFY DATAS

Datas available to the RA during the notify actions:

```
active    => [ ],
inactive  => [
    { rsc => 'pgsql:2', uname => 'srv1' },
    { rsc => 'pgsql:0', uname => 'srv2' },
    { rsc => 'pgsql:1', uname => 'srv3' }
],
master    => [ ],
slave     => [ ],
promote   => [ { rsc => 'pgsql:0', uname => 'srv1' } ],
demote    => [ ],
start     => [
    { rsc => 'pgsql:0', uname => 'srv1' },
    { rsc => 'pgsql:1', uname => 'srv3' },
    { rsc => 'pgsql:2', uname => 'srv2' }
],
stop      => [ ],
type      => 'pre'
```



# MASTER SCORE

- set preference on slave to promote
- highest score is promoted to master
- a slave must have a positive score to be promoted
- no promotion if no master score anywhere
- set by the resource agent and/or the admin

Started

Stopped



# PostgreSQL Automatic Failover

# HISTORY

- pgconf.eu 2012 talk on Pacemaker/pgsql
- had a hard time to build a PoC and document
- discussion with Magnus about demote
- (other small projects around this before)
- PAF started in 2015
- lots of questions to Pacemaker's devs
- authors: Maël Rimbault, Me
- some contributors and feedbacks (Thanks!)

# WHY ?

The existing RA:

- achieve multiple architectures (stateless and multistate)
- implementation details to understand (lock file)
- only failover (no role swapping or recovery)
- hard and heavy to manage (start/stop order, etc)
- hard to setup
- fake Pacemaker state because of demote, mess in the code
- old code...

# GOALS

- keep Pacemaker: it does most of the job for us
- focus on our expertise: PostgreSQL
- stick to the OCF API and Pacemaker behavior, embrace them
- keep a **SIMPLE** RA setup
- support **ONLY** multi-state
- support **ONLY** Streaming Replication
- **REQUIRE** Streaming Replication and Hot Standby
- ease of administration
- keep the code clean and documented
- support PostgreSQL 9.3 and after

# VERSIONS

Two versions to catch them (almost) ALL!

- **1.x:** up to EL6 and Debian 7
- ...or Pacemaker 1.12/corosync 1.x
- **2.x:** from EL7 and Debian 8
- ... or Pacemaker 1.13/Corosync 2.x

# GUTS

- written in perl
- demote = stop + start (= slave)
- slave election during failover
- detect various kind of transitions thanks to notify (recover and move)

# PAF CONFIGURATION

- system\_user
- bindir
- datadir (oops, 1.1 only)
- pgdata
- pgghost
- pgport
- recovery\_template
- start\_opts



# OLD CONFIGURATION

Compare with historical `pgsql` RA:

- pgctl
- start\_opt
- ctl\_opt
- psql
- pgdata
- pgdba
- pgghost
- pgport
- pglibs
- monitor\_user

# OLD CONFIGURATION (2)

Encore?

- monitor\_password
- monitor\_sql
- config
- pgdb
- logfile
- socketdir
- **stop\_escalate**
- rep\_mode
- node\_list
- restore\_command

# OLD CONFIGURATION (3)

Not done yet...

- archive\_cleanup\_command
- recovery\_end\_command
- master\_ip
- repuser
- primary\_conninfo\_opt
- restart\_on\_promote
- replication\_slot\_name
- tmpdir
- xlog\_check\_count
- crm\_attr\_timeout

# OLD CONFIGURATION (4)

Promise, the last ones:

- stop\_escalate\_in\_slave
- check\_wal\_receiver



# Features

The following demos considers:

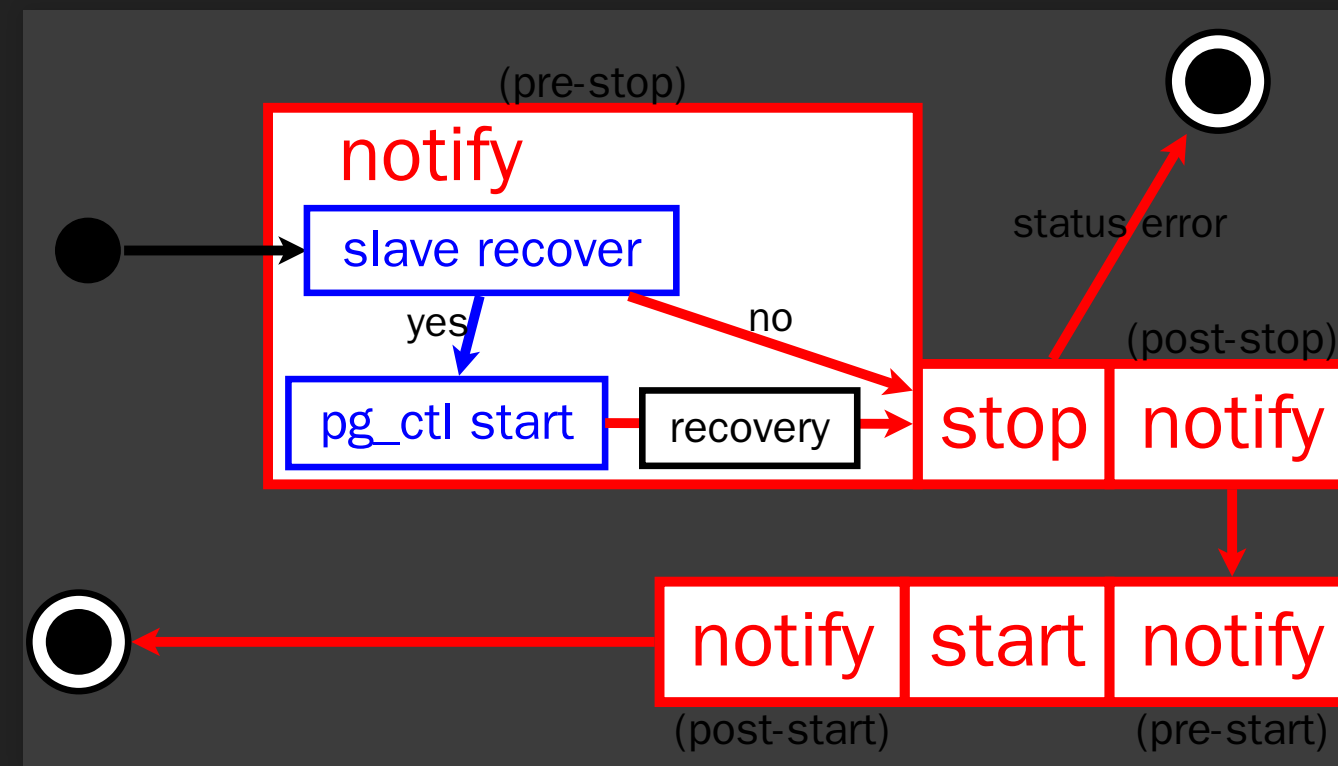
- one master & two slaves
- a secondary IP address following the master role:  
192.168.122.50
- a really simple recovery.conf template file:

```
standby_mode = on  
primary_conninfo = 'host=192.168.122.50 application_name=$(hostname -s)'  
recovery_target_timeline = 'latest'
```

- a monitor action every 15s

# STANDBY RECOVER

Transition: **stop** -> **start**



# Slave recover demo:

```
root@srv1:~#
root@srv1:~#

root@srv2:~#
Online: [ srv1 srv2 srv3 ]

Full list of resources:

fence_vm_srv1 (stonith:fence_virsh): Started srv2
fence_vm_srv2 (stonith:fence_virsh): Started srv1
fence_vm_srv3 (stonith:fence_virsh): Started srv1
Master/Slave Set: pgsql-ha [pgsql]
Masters: [ srv1 ]
Slaves: [ srv2 srv3 ]
pgsql-master-ip (ocf::heartbeat:IPaddr2): Started srv1

Node Attributes:
* Node srv1:
+ master-pgsql : 1001
* Node srv2:
+ master-pgsql : 1000
* Node srv3:
+ master-pgsql : 990

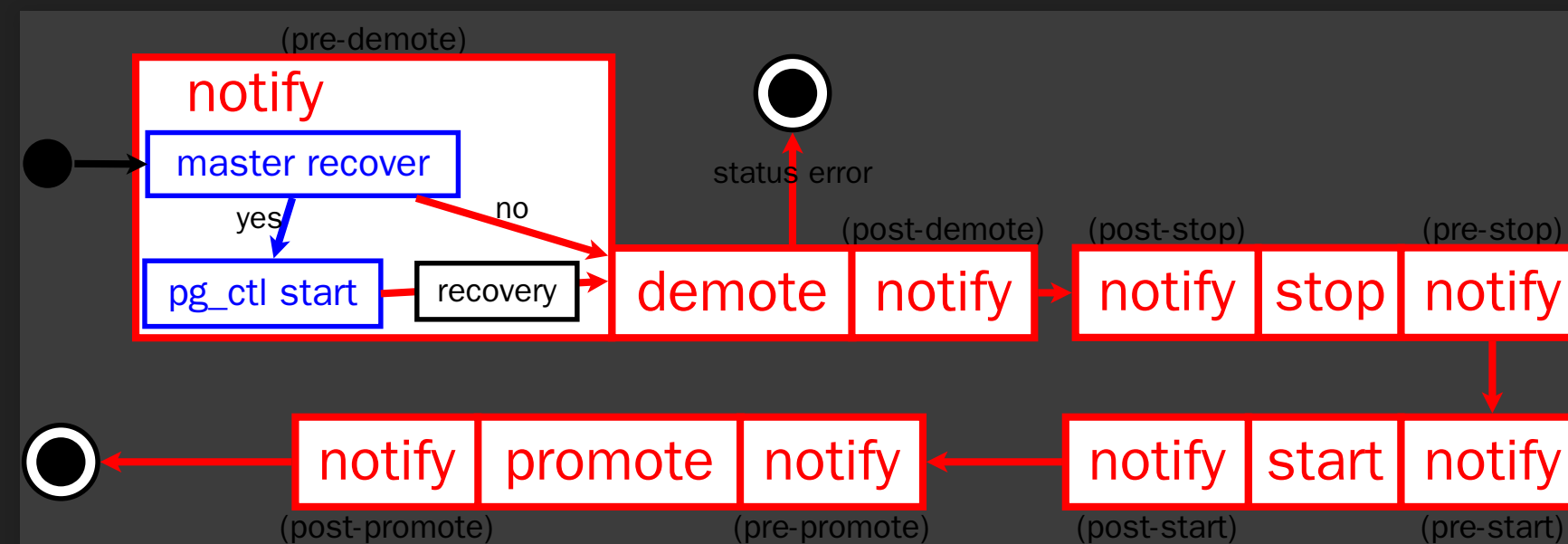
Migration Summary:
* Node srv2:
* Node srv3:
* Node srv1:
[ ]

root@srv3:~#
Every 1,0s: ps f -o cmd= -u postgres      Fri Sep  9 16:27:00 2016

/usr/pgsql-9.4/bin/postgres -D /var/lib/pgsql/9.4/data
\_ postgres: logger process
\_ postgres: startup process  recovering 00000018000000000000000019
\_ postgres: checkpointer process
\_ postgres: writer process
\_ postgres: stats collector process
\_ postgres: wal receiver process  streaming 0/19FB63F0
```

# MASTER RECOVER

Transition: demote -> stop -> start -> promote

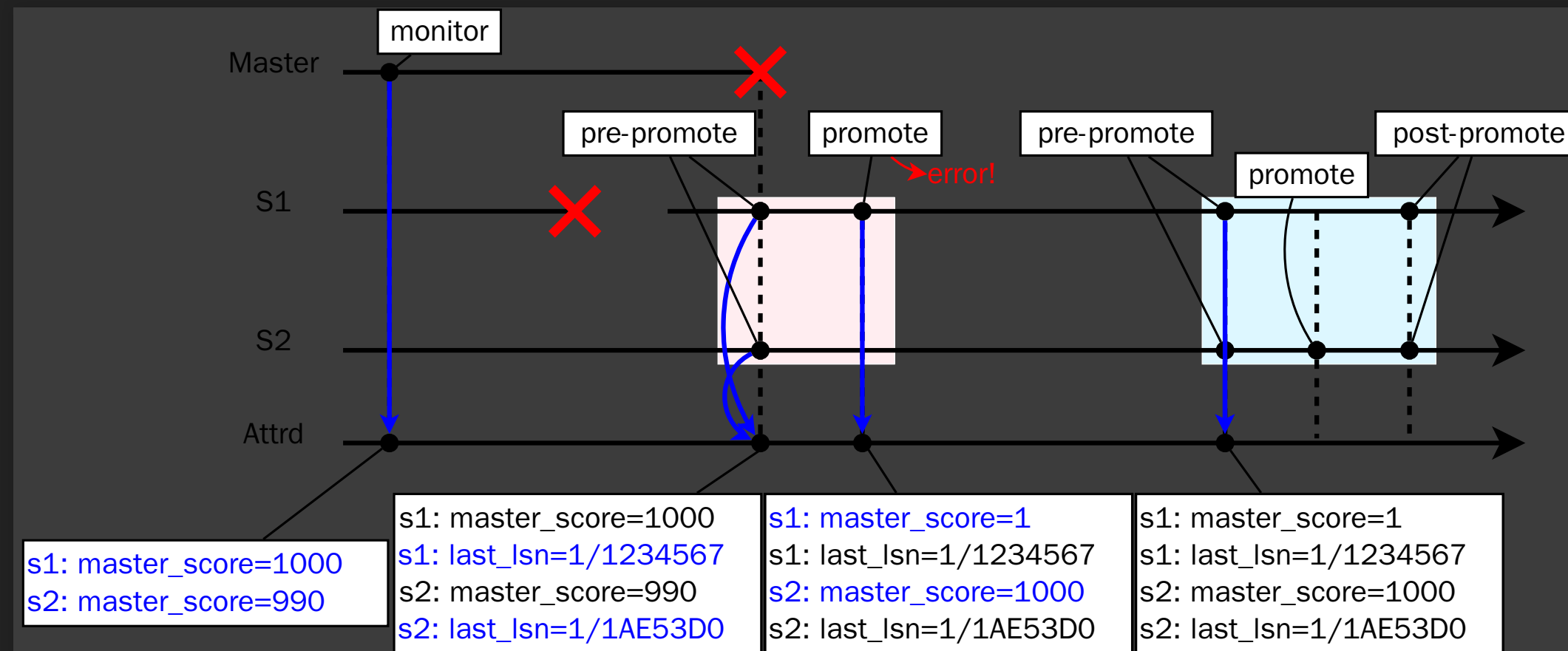




# Master recover demo:

```
root@srv1:~#  
  
root@srv2:~#  
Online: [ srv1 srv2 srv3 ]  
Full list of resources:  
fence_vm_srv1 (stonith:fence_virsh): Started srv2  
fence_vm_srv2 (stonith:fence_virsh): Started srv1  
fence_vm_srv3 (stonith:fence_virsh): Started srv1  
Master/Slave Set: pgsql-ha [pgsql]  
Masters: [ srv1 ]  
Slaves: [ srv2 srv3 ]  
pgsql-master-ip (ocf::heartbeat:IPaddr2): Started srv1  
  
Node Attributes:  
* Node srv1:  
+ master-pgsql : 1001  
* Node srv2:  
+ master-pgsql : 1000  
* Node srv3:  
+ master-pgsql : 990  
  
Migration Summary:  
* Node srv2:  
* Node srv3:  
* Node srv1:  
[  
  
root@srv3:~#  
Every 1,0s: ps f -o cmd= -u postgres      Fri Sep  9 17:46:41 2016  
  
/usr/pgsql-9.4/bin/postgres -D /var/lib/pgsql/9.4/data  
\_ postgres: logger process  
\_ postgres: startup process  recovering 0000001A0000000000000001D  
\_ postgres: checkpoint process  
\_ postgres: writer process  
\_ postgres: stats collector process  
\_ postgres: wal receiver process streaming 0/1D1FAEB8
```

# FAILOVER & ELECTION



# Failover demo:

```
root@srv1:~# echo c > /proc/sysr
root@srv2:~# tail -f /var/log/cluster/corosync.log|grep -o "DEBUG: pgsql_promote.*"

Online: [ srv1 srv2 srv3 ]

Full list of resources:

fence_vm_srv1 (stonith:fence_virsh): Started srv2
fence_vm_srv2 (stonith:fence_virsh): Started srv3
fence_vm_srv3 (stonith:fence_virsh): Started srv2
Master/Slave Set: pgsql-ha [pgsql]
Masters: [ srv1 ]
Slaves: [ srv2 srv3 ]
pgsql-master-ip (ocf::heartbeat:IPaddr2): Started srv1

Node Attributes:
* Node srv1:
+ master-pgsql : 1001
* Node srv2:
+ master-pgsql : 1000
* Node srv3:
+ master-pgsql : 990

Migration Summary:
* Node srv2:
* Node srv3:
* Node srv1:
```

# YOU THINK IT'S OVER ?



# CONTROLLED SWITCHOVER

- only with 2.0
- the designated standby checks itself
- it cancel the promotion if the previous master will not be able to catch up with it.

# Controlled switchover demo:

```
root@srv1:~#  
  
root@srv2:~#  
Online: [ srv1 srv2 srv3 ]  
Full list of resources:  
fence_vm_srv1 (stonith:fence_virsh): Started srv2  
fence_vm_srv2 (stonith:fence_virsh): Started srv1  
fence_vm_srv3 (stonith:fence_virsh): Started srv1  
Master/Slave Set: pgsql-ha [pgsql]  
Masters: [ srv1 ]  
Slaves: [ srv2 srv3 ]  
pgsql-master-ip (ocf::heartbeat:IPaddr2): Started srv1  
  
Node Attributes:  
* Node srv1:  
+ master-pgsql : 1001  
* Node srv2:  
+ master-pgsql : 1000  
* Node srv3:  
+ master-pgsql : 990  
  
Migration Summary:  
* Node srv2:  
* Node srv3:  
* Node srv1:  
[ ]  
  
root@srv3:~#  
Every 1,0s: ps f -o cmd= -u postgres Fri Sep 9 18:28:45 2016  
/usr/pgsql-9.4/bin/postgres -D /var/lib/pgsql/9.4/data  
\_ postgres: logger process  
\_ postgres: startup process recovering 0000002C0000000000000029  
\_ postgres: checkpoint process  
\_ postgres: writer process  
\_ postgres: stats collector process  
\_ postgres: wal receiver process streaming 0/29088FB0
```

# Whishlist

- recovery.conf as GUC
- live demote
- pgbench handling of errors

# Where?

- site officiel: <http://dalibo.github.io/PAF/>
- code: <https://github.com/dalibo/PAF>
- packages: <https://github.com/dalibo/PAF/releases>
- support: <https://github.com/dalibo/PAF/issues>
- mailing list: pgsql-general



**PLZ, ASK ME QUESTIONS**