

# HAUTE DISPONIBILITÉ AVEC POSTGRESQL ET PATRONI

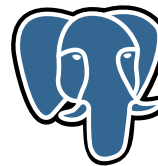
Jean-Christophe Arnu

Capitôle Du Libre 2024

# ENCHANTÉ, JE ME PRÉSENTE !

## JEAN-CHRISTOPHE ARNU

- Création d'une association autour des logiciels libres à Toulouse 1996
- PostgreSQL depuis 1998
- Un des membres fondateurs de PostgreSQLFr
- Organisateur du premier PGDay.fr à Toulouse en 2008
- Consultant PostgreSQL chez LOXODATA



Tux: par Larry Ewing, Simon Budig, Anja Gerwinski, ©  
GNU: par Aurelio A. Heckert, Free Art Licence  
PostgreSQL: par Daniel Lundin, PostgreSQL License  
Toulouse: par "The blazon project" wikipedia CC-BY-SA

# QUOI ?

- ~~Mais au fait c'est quoi la haute dispo ?~~
- Quels problèmes la haute dispo cherche à résoudre ?  
**Spoiler-alert : c'est barbant, c'est de la théorie**
  - Ok PostgreSQL? Il le gère comment ?
  - **PATRONI** et la haute disponibilité.

# LA HAUTE DISPONIBILITÉ ? KEZAKO ?

Fonctionne tout le temps, de manière optimale:

- On redonde (Alim, réseaux, serveurs ...)
- On distribue (plusieurs zones (AZ), plusieurs DC)

Les indicateurs:

- On ne veut aucun « downtime » : RTO=0
- On ne veut perdre aucune donnée : RPO=0
- ... ?

Que veut-on au juste ?

## LES MYTHES

*« Je veux que ça fonctionne de manière performante, sans interruption de service, sans perte de donnée, quoi qu'il arrive »*

*« Je veux que ça fonctionne 99.999% du temps »\**

Nous sommes dans une conférence technique, nous allons voir pourquoi ce n'est pas possible.

# VOYONS VOIR (SUR UN AN)...

99% :3 JOURS 14H 24 MINUTES

```
select justify_interval('1 year'::interval*0.01)
justify_interval
-----
3 days 14:24:00
```

99.9% : 8H 38 MINUTES

```
select justify_interval('1 year'::interval*0.001)
justify_interval
-----
08:38:24
```

99.99%: ~ 52 MINUTES

```
select justify_interval('1 year'::interval*0.0001)
justify_interval
-----
00:51:50.4
```

99.999% : 5 MINUTES

```
select justify_interval('1 year'::interval*0.00001)
justify_interval
-----
00:05:11.04
```

## PLEIN DE 9 ?

**\* MAIS NE FAIT-ON JAMAIS DE MAINTENANCE DE L'APPLI,  
DES BASES (UPGRADES), L'OS, ... ?**

**POURQUOI NE PAS APPLIQUER LES MÊMES RÈGLES  
À LA BASE ET AUX APPLIS ?**

## LE MYTHE DU RPO=0

*« Je ne veux perdre aucune donnée »*



## LE MYTHE DU RPO=0

D'accord, mais en cas de désastre, on perd :

- Le cache système (fsync)
- Le cache disque (batteries ?)
- L'hyperviseur ou le CP
- Le cache réseau (pour le stockage réseau)

## LE MYTHE DU RPO=0

**QUE PERDONS-NOUS RÉELLEMENT ?**

**EN BASE DE DONNÉES : DES TRANSACTIONS**

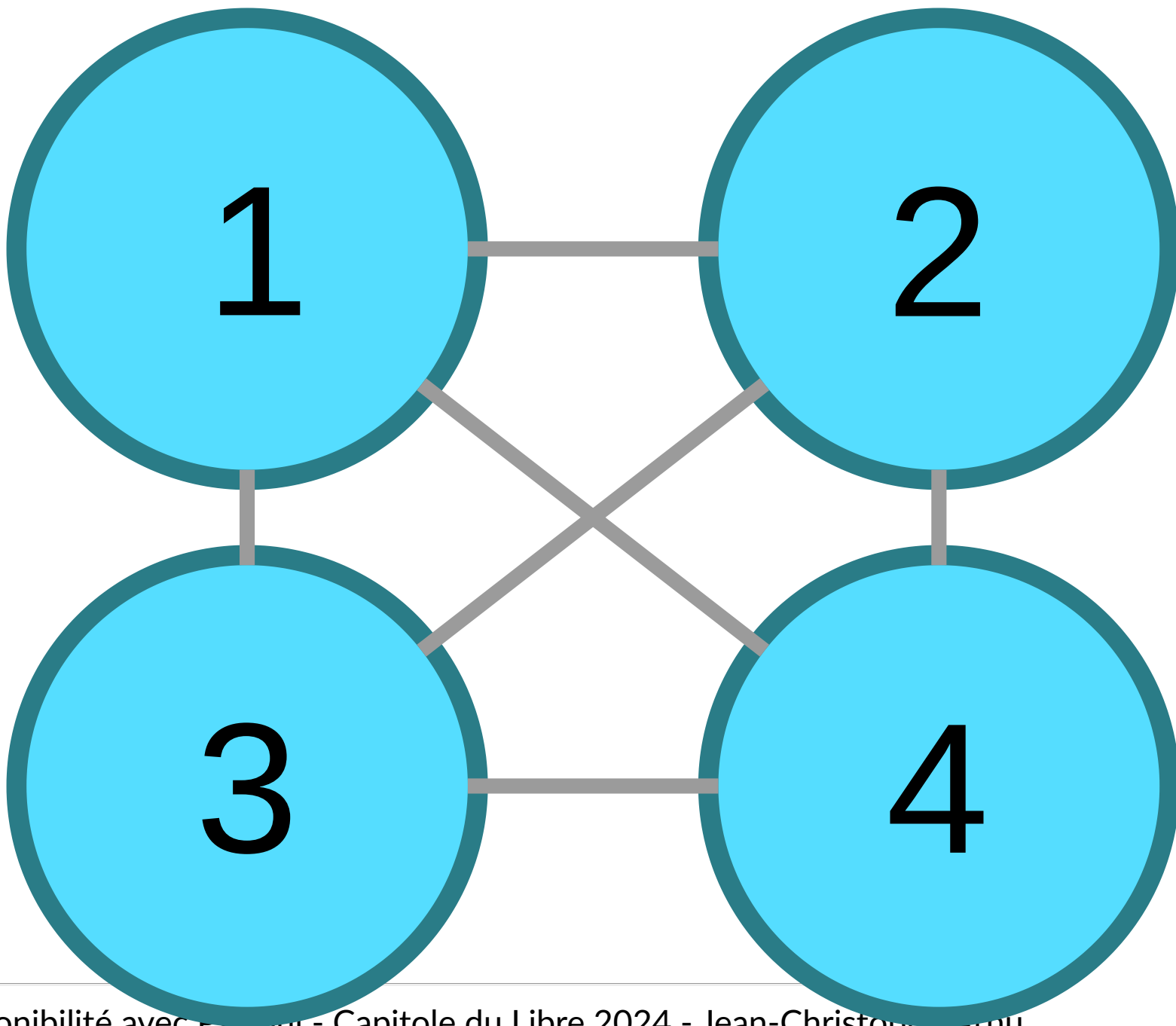
## UN SYSTÈME DISTRIBUÉ

ON VA REDONDER LE SYSTÈME SUR LE RÉSEAU, MAIS

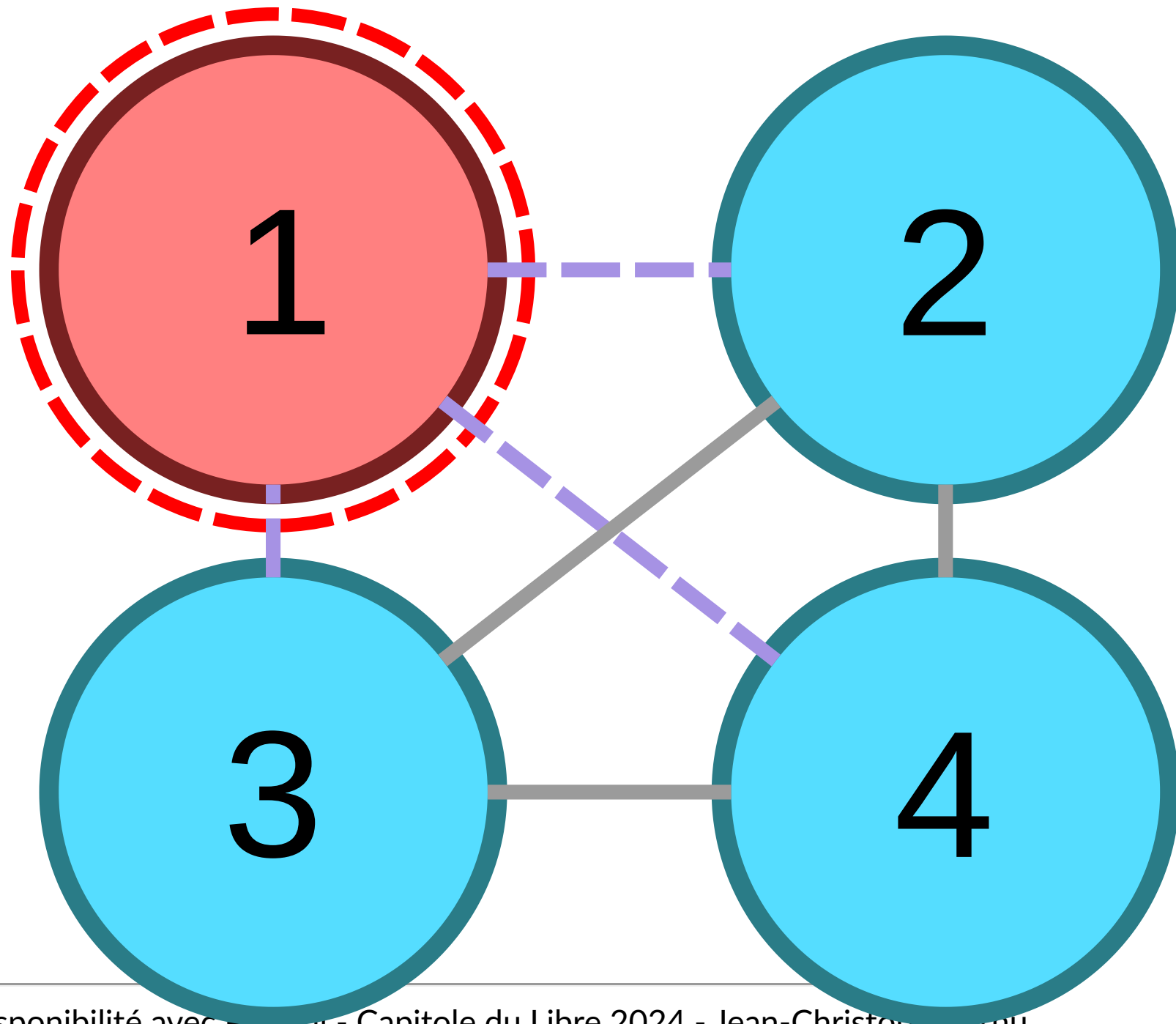
- On veut garantir la redondance, mais un nœud peut:
  - tomber en panne
  - se retrouver isolé
  - se retrouver indisponible

On peut se retrouver en situation de partitionnement réseau

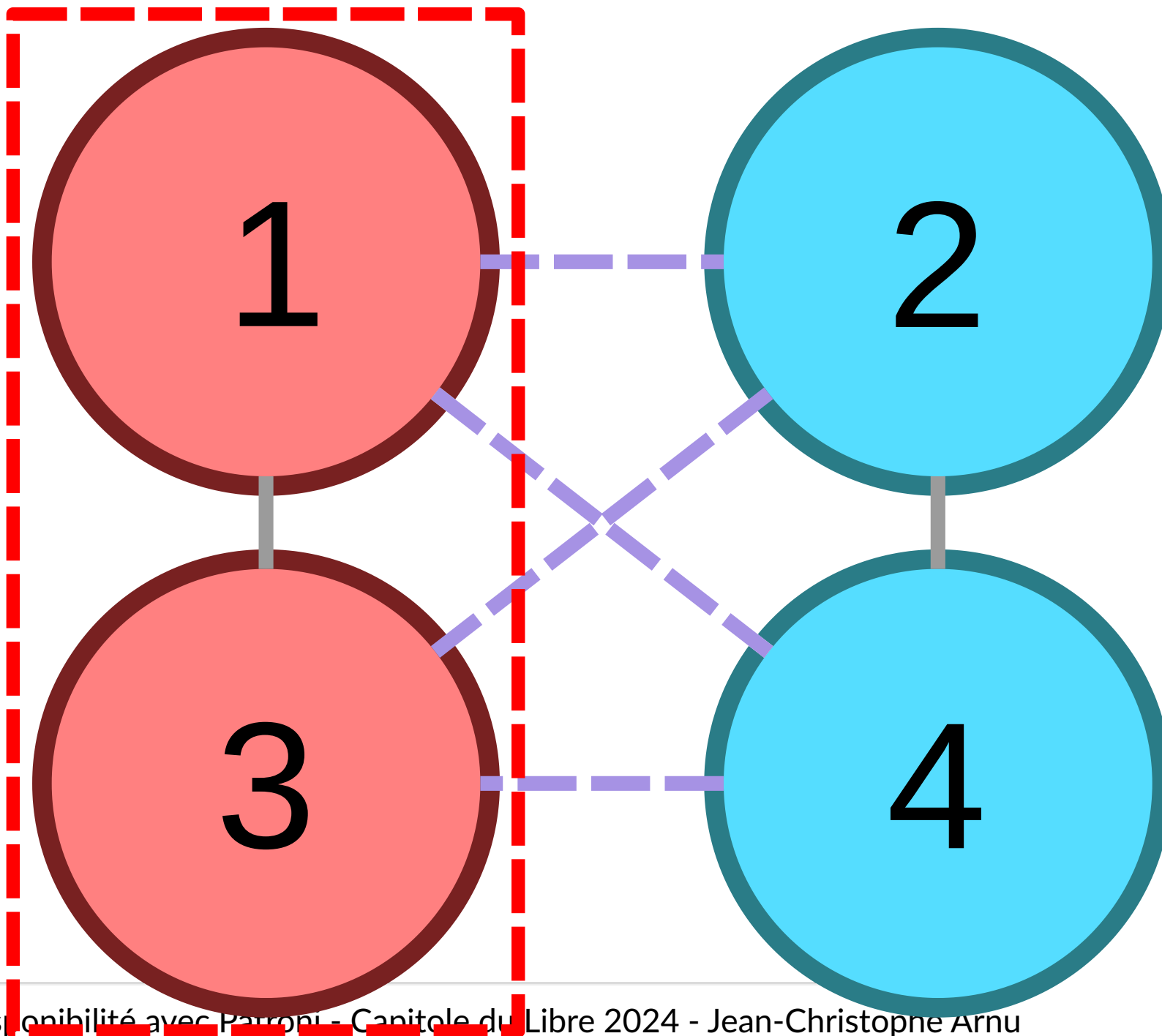
# LE PARTITIONNEMENT RÉSEAU



# LE PARTITIONNEMENT RÉSEAU



# LE PARTITIONNEMENT RÉSEAU



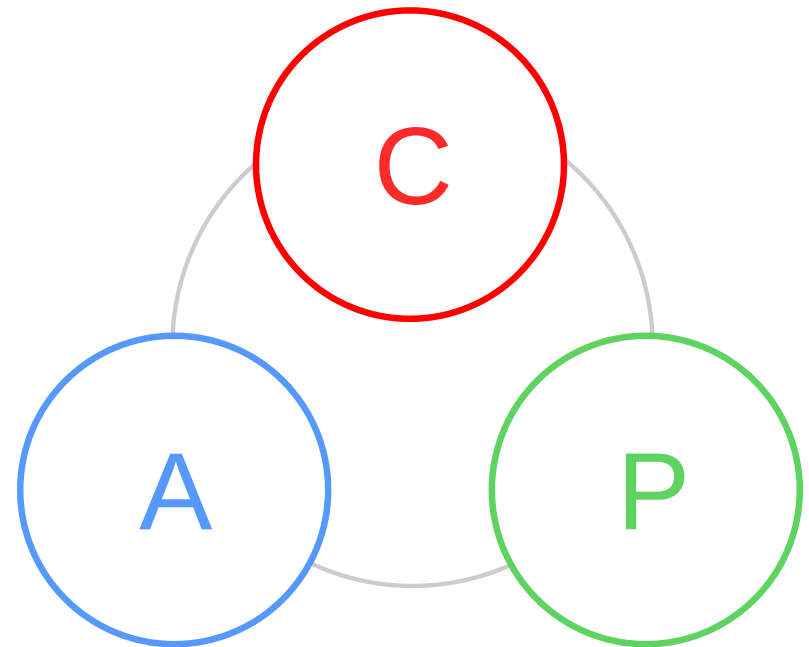
# CREUSONS

## THÉORÈME CAP

Eric Brewer

- C: Consistency : les données sont cohérentes entre tous les nœuds
- A: Availability : le service est disponible
- P: Partitioning : le service est résistant au partitionnement réseau

**ON NE PEUT EN AVOIR QUE DEUX PROPRIÉTÉS EN MÊME TEMPS.**

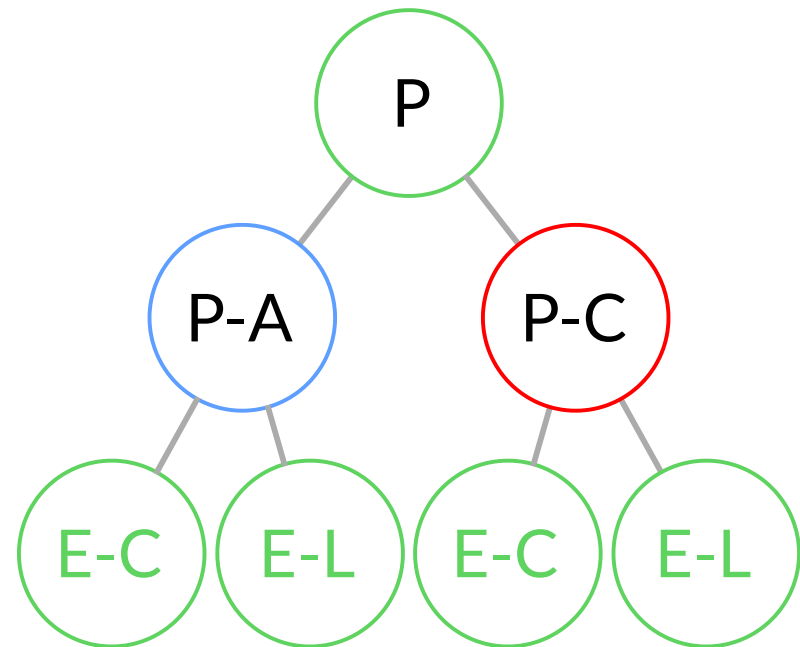


# IL FAUT DONC FAIRE DES CHOIX!

## THÉORÈME PACELC

Daniel Abadi

- Bascule:
  - Manuelle (P-C)
  - Automatique (P-A)
- Réplication:
  - Asynchrone (E-L)
  - Synchrone (E-C)





# ET POSTGRESQL DANS TOUT ÇA ?

## RÉPLICATION PAR FLUX NATIVE

- Un standby se connecte à un primaire
- Le standby reçoit le flux temps réel des modifications (WAL)
- Peut être asynchrone PC/EL ou synchrone PC/EC

# QUELQUES MOTS SUR LA RÉPLICATION SYNCHRONE

- Permet d'implémenter RPO=0
- Mais tout COMMIT est bloqué tant qu'au moins un standby n'a pas acquitté
  - Latence
  - Blocage en cas d'indispo d'un standby

## SAUVEGARDE

- Native (pg\_basebackup + archive command) [NIH]
- Restauration sans impact sur le primaire:
  - Barman
  - PgBackrest
- Archivage des WALs

## FONCTIONS ?

### POURQUOI PAS DE HAUTE DISPONIBILITÉ DANS POSTGRESQL ?

- PostgreSQL fournit une **infrastructure** pour l'implanter
  - Gestion de la réplication (lecture, contrôle)
  - Gestion de la promotion/« demotion » en primaire
- Permet de résoudre une grande variété de cas d'usage
- Sans contraindre l'utilisateur à une seule vision

# CONDITIONS DE LA HAUTE DISPONIBILITÉ

(AVEC POSTGRESQL)

# LES DCS

## Distributed Consensus Store

- Consensus entre les nœuds du DCS
  - Un leader, N followers
  - Si k nœuds dans le cluster  $\Rightarrow$  quorum =  $(k/2)+1$
- Configuration du cluster PostgreSQL décentralisée
- Élection du primaire
  - Choix du primaire (le plus en avance)
  - Nouvelle topologie
  - Lock config
  - Primaire  $\rightarrow$  Update config

## EXEMPLE DE DCS

- zookeeper (praxos)
- etcd (raft)
- consul (raft)
- K8s API (raft)



# ACCÉDER AUX NOEUDS

## Entrypoints

- VIP (IP virtuelle)
- HAProxy
- PgBouncer
- N'oublions pas les drivers (libpq + java + DotNet)
- 12 factors directement auprès du DCS

# COMMENT RECONSTRUIRE DES RÉPLICAS ?

## Avec des outils de sauvegarde

- Basique et natif:
  - pg\_basebackup
- Plus évolués:
  - Barman
  - PgBackRest

# PATRONI

INITIÉ PAR ZALANDO (INTÉGRATION DOCKER : SPILO)

LICENCE MIT

# PATRONI EST UN TEMPLATE

## S'interface avec

- PostgreSQL
- Un DCS
- Un outil de backup
- Un seul fichier de configuration

# FONCTIONNALITÉS

## Pour le cluster

- Suivi de la réplication synchrone ou asynchrone (configuration)
- Réplication simple ou en cascade
- Failover
- Switchover
- Empêche l'état splitbrain  $\Rightarrow$  demote primary  $\neq$  Quorum
- Sait prendre en compte l'extension Citus (sharding)

# PRINCIPE

## Systeme modulaire

- Utilise un DCS pour
  - Propager la configuration du cluster
  - Stocker l'état du cluster
  - Propager certains éléments de configuration de PostgreSQL
- Utilise un système de backup pour construire des nœuds
- Configure PostgreSQL, pilote son fonctionnement et sa réplication

# QUIZZ

Quelle politique de haute-disponibilité

**PATRONI**

met-il en œuvre vis à vis:

- Du théorème CAP ?
- Du théorème PACELC ?

# OUTILS/INTERFACES

## API REST

- Information sur la configuration
- État
- Santé
- ...



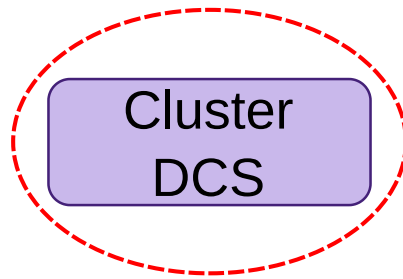
# OUTILS/INTERFACES

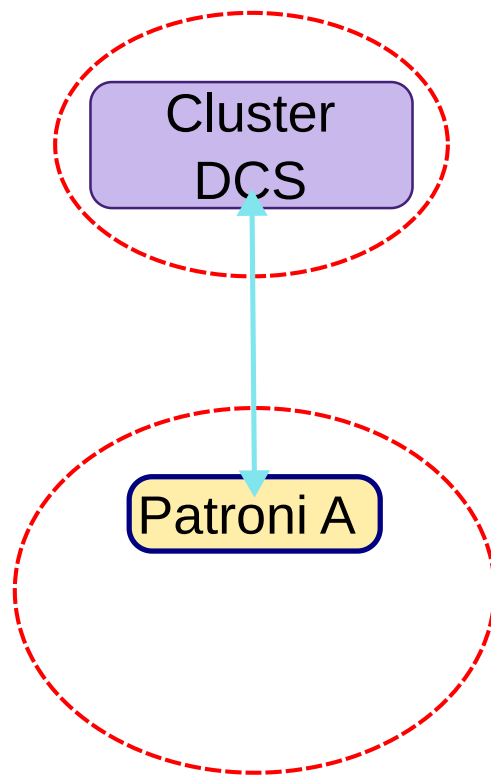
## Outil CLI pour contrôler le cluster

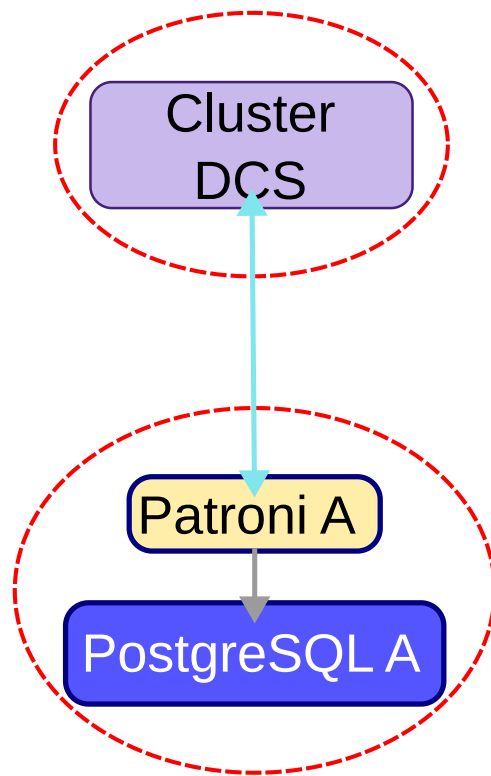
- Éditer/consulter la configuration
- Voir l'état, le lag
- Piloter le cluster (switchover)

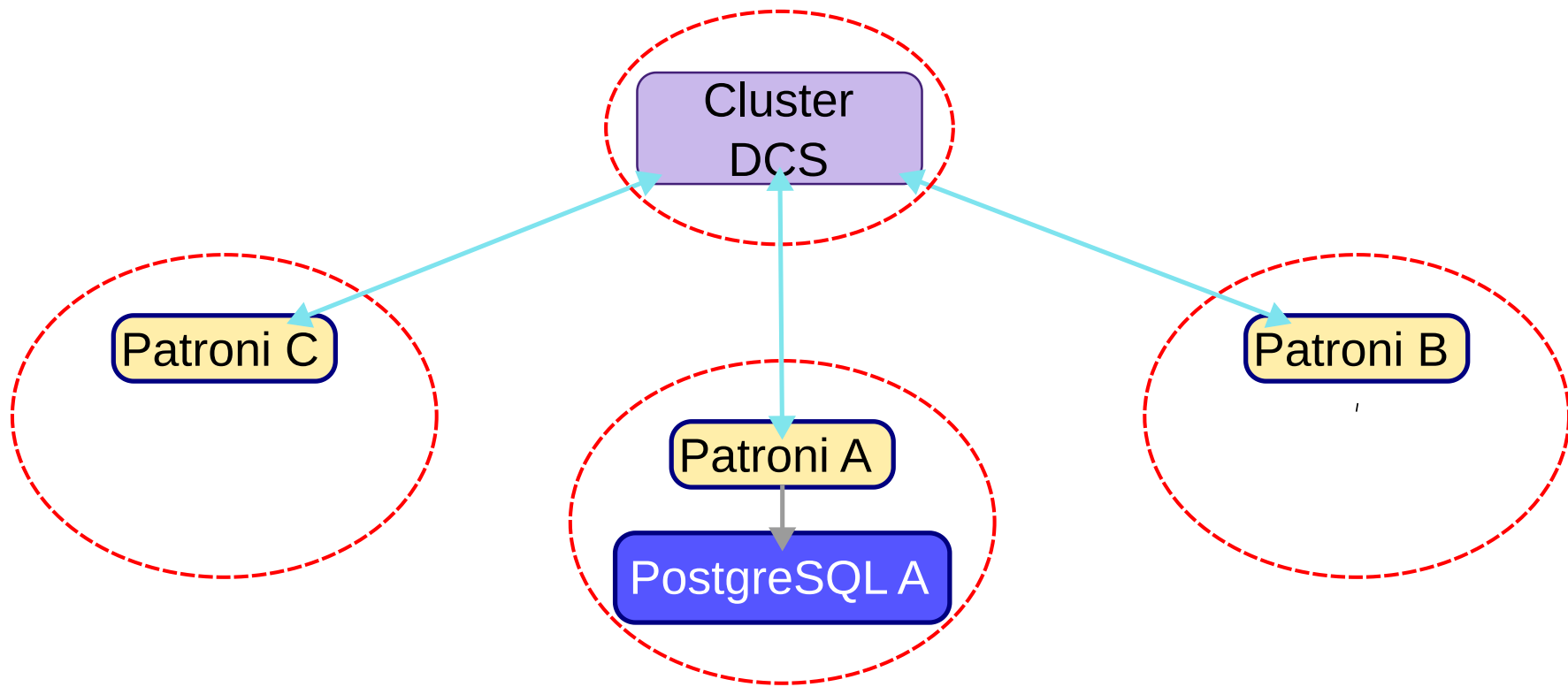
# FONCTIONNEMENT

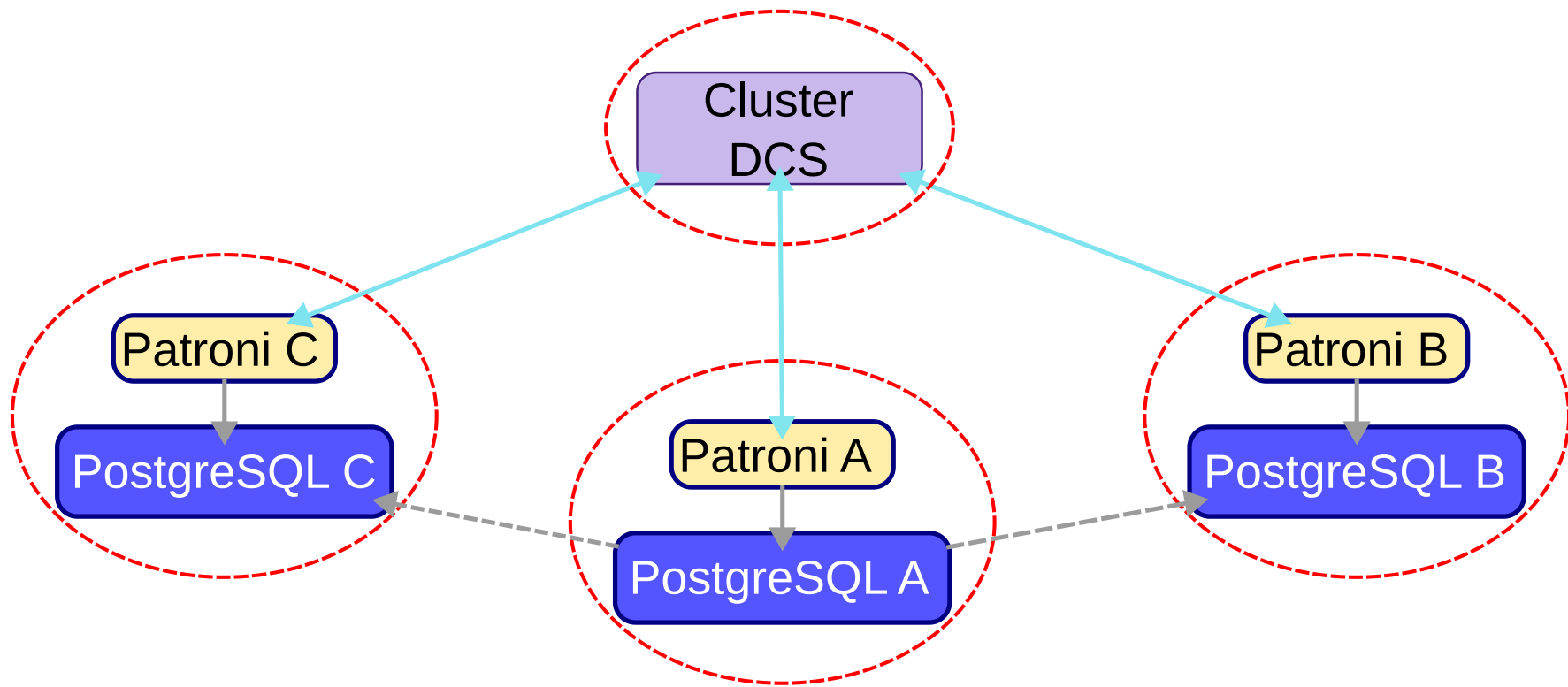
- **Patroni :**
  - Pilote PostgreSQL
  - Bootstrap du primaire
  - Création de réplicas
  - Gère les bascules
- **Mais on doit configurer indépendamment :**
  - Le cluster DCS
  - Le système de sauvegarde
- **Et les entrypoints ?**
  - Utiliser les données du DCS pour les autoconfigurer
  - Par exemple avec confd/remco

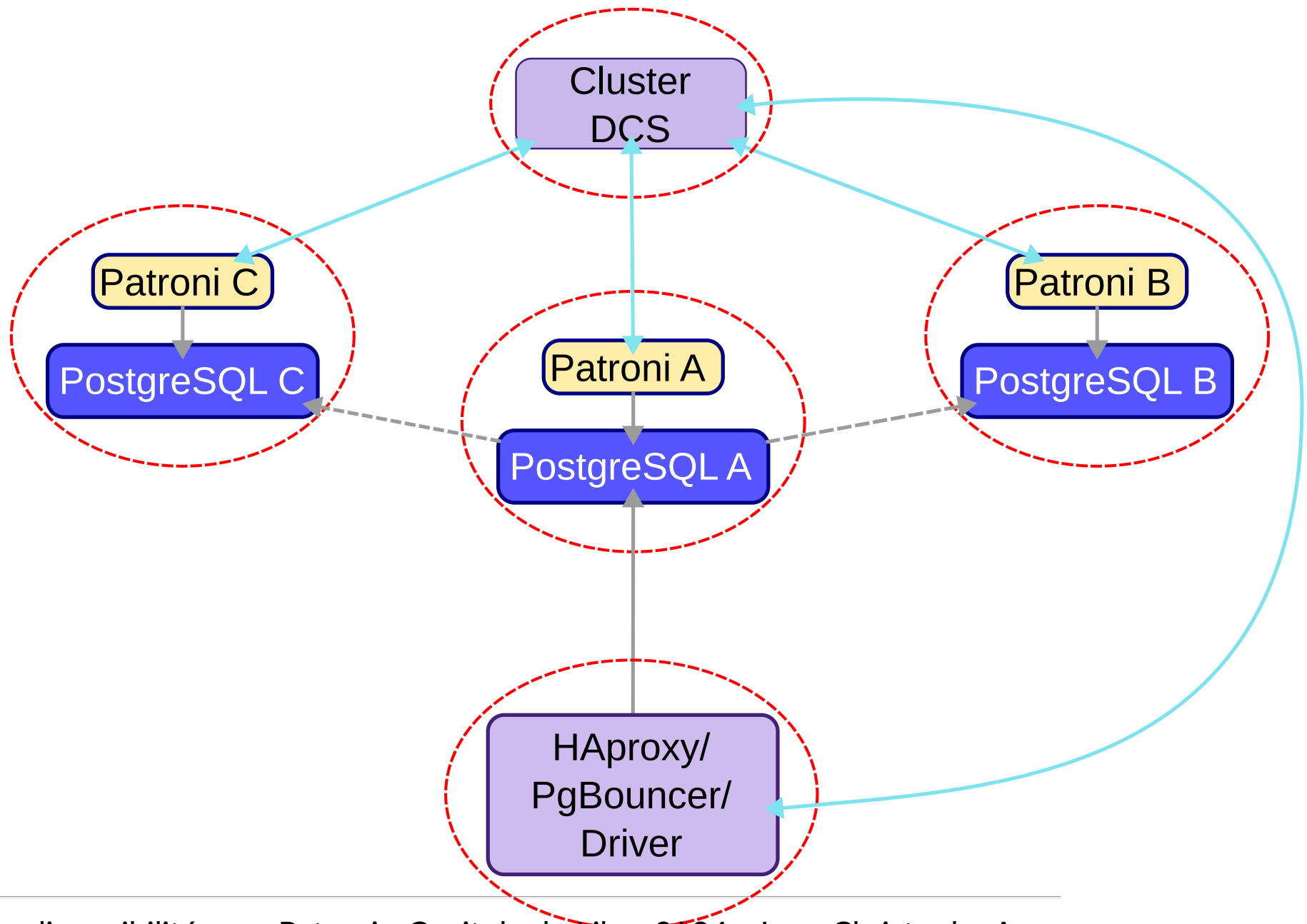




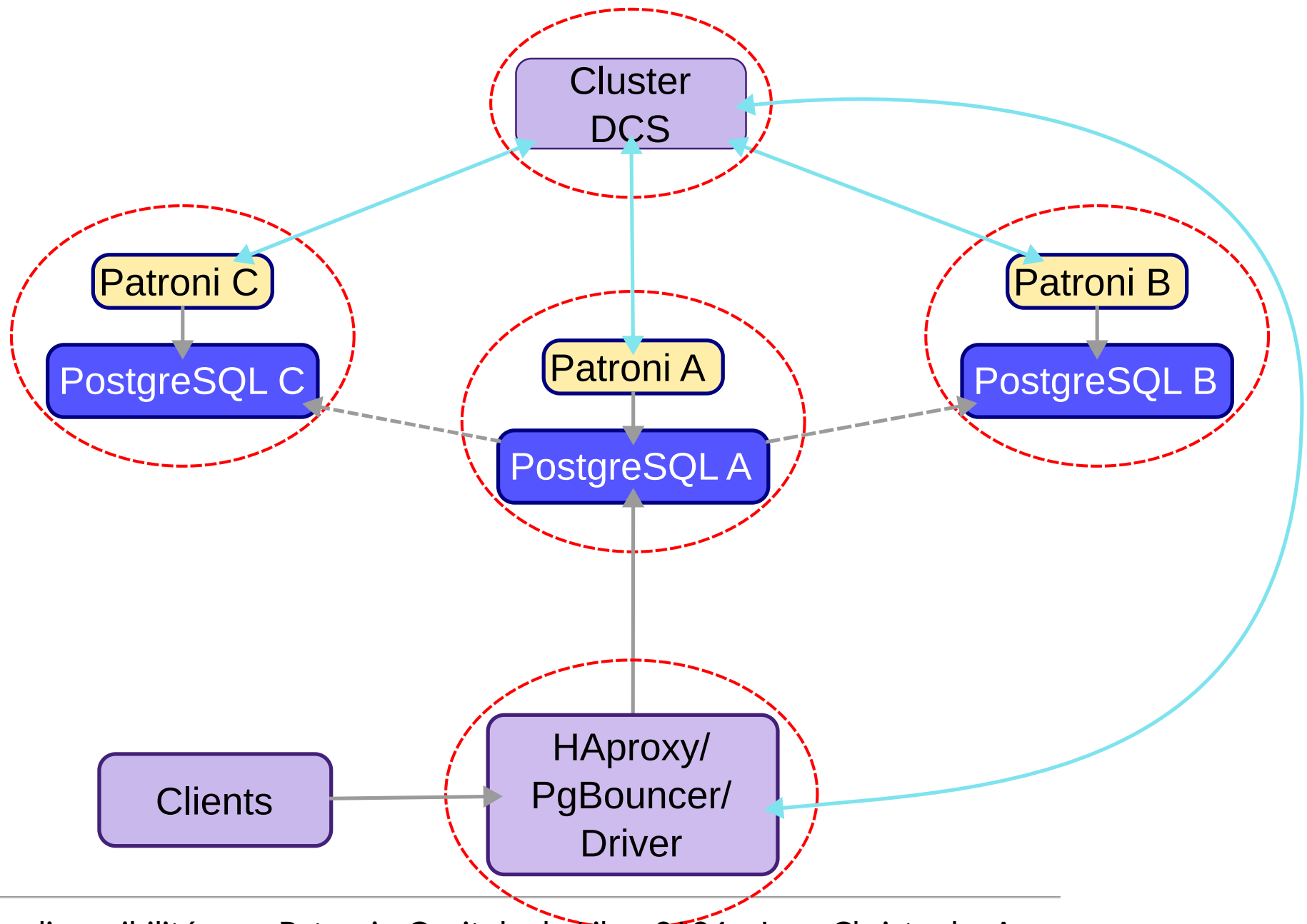












```
1 scope: loxodemo
2 name: pgdeb01
3
4 restapi:
5   listen: 10.200.0.11:8008
6   connect_address: 10.200.0.11:8008
7
8 etcd3:
9   hosts: # membres du cluster etcd.
10     - 10.200.0.11:2379
11     - 10.200.0.12:2379
12     - 10.200.0.13:2379
13
14 log:
15   level: INFO
16   dir: /var/log/patroni
17   file_size: 50000000
18   file_num: 10
19   format: '%(asctime)s %(levelname)s: %(message)s'
20   dateformat: '%Y-%m-%d %H:%M:%S'
21   loggers:
22     etcd.client: INFO
```

```
1 bootstrap: #DCS et initdb
2   dcs:      # Utiliser edit-config pour modifier après init
3     ttl: 30 #tps (s) validité d'une clé dans le DCS après mise à jour
4     loop_wait: 10 # tps(s) entre 2 update DCS
5     retry_timeout: 10 # tps (s) à attendre (2*) si DCS non joignable
6     maximum_lag_on_failover: 1048576
7     postgresql:
8       use_pg_rewind: false
9       use_slots: false
10      parameters:
11        archive_mode: "on"
12        archive_command: pgbackrest --stanza=loxodemo archive-push %p
13        recovery_conf:
14          restore_command: pgbackrest --stanza=loxodemo archive-get %f %p
15  initdb:
16    - encoding: UTF8
17    - data-checksums
18
19  pg_hba:
20    - host all all 10.200.0.0/24 trust
21    - host replication repli 10.200.0.0/24 trust
22  users:
```

# CONFIGURATION

```
1 postgresql:
2   listen: 10.200.0.11:5432
3   connect_address: 10.200.0.11:5432
4   data_dir: /data/postgres/loxodemo
5   bin_dir: /usr/lib/postgresql/16/bin/
6   authentication:
7     replication:
8       username: repli
9       password: azerty
10    superuser:
11      username: postgres
12      password: azerty
13    parameters:
14      unix_socket_directories: '/tmp'
15      logging_collector: on
16    create_replica_methods:
17      - pgbackrest
18    pgbackrest:
19      command: pgbackrest --stanza=loxodemo --delta restore
20      keep_data: True
21      no_params: True
22
```

```
patronictl topology loxodemo
+ Cluster: loxodemo (7436044389321027980) -----+-----+-----+
| Member      | Host          | Role   | State   | TL | Lag in MB |
+-----+-----+-----+-----+-----+-----+
| pgdeb01     | 10.200.0.11  | Leader | running | 1  |           |
| + pgdeb02   | 10.200.0.12  | Replica| streaming| 1  |          0 |
| + pgdeb03   | 10.200.0.13  | Replica| streaming| 1  |          0 |
+-----+-----+-----+-----+-----+-----+
```

```
patronictl -c /etc/patroni/config.yml switchover --candidate=pgdeb02 loxodemo --force
Current cluster topology
+ Cluster: loxodemo (7436044389321027980) -----+-----+-----+
| Member      | Host          | Role   | State   | TL | Lag in MB |
+-----+-----+-----+-----+-----+-----+
| pgdeb01     | 10.200.0.11  | Leader | running | 1  |           |
| pgdeb02     | 10.200.0.12  | Replica| streaming| 1  |          0 |
| pgdeb03     | 10.200.0.13  | Replica| streaming| 1  |          0 |
+-----+-----+-----+-----+-----+-----+
2024-11-11 15:35:05.12331 Successfully switched over to "pgdeb02"
+ Cluster: loxodemo (7436044389321027980) -----+-----+-----+
| Member      | Host          | Role   | State   | TL | Lag in MB |
+-----+-----+-----+-----+-----+-----+
| pgdeb01     | 10.200.0.11  | Replica| stopped  |   | unknown   |
| pgdeb02     | 10.200.0.12  | Leader | running  | 1  |           |
| pgdeb03     | 10.200.0.13  | Replica| in archive recovery| 1  |          0 |
+-----+-----+-----+-----+-----+-----+
```

```
patronictl -c /etc/patroni/config.yml topology loxodemo
+ Cluster: loxodemo (7436044389321027980) -----+-----+-----+
| Member      | Host          | Role   | State   | TL | Lag in MB |
+-----+-----+-----+-----+-----+-----+
| pgdeb02     | 10.200.0.12  | Leader | running  | 2  |           |
| + pgdeb01   | 10.200.0.11  | Replica| streaming| 2  |          0 |
| + pgdeb03   | 10.200.0.13  | Replica| streaming| 2  |          0 |
+-----+-----+-----+-----+-----+-----+
```

```
patronictl -c /etc/patroni/config.yml list loxodemo --format=json|jq .
[
  {
    "Cluster": "loxodemo",
    "Member": "pgdeb01",
    "Host": "10.200.0.11",
    "Role": "Replica",
    "State": "streaming",
    "TL": 2,
    "Lag in MB": 0
  },
  {
    "Cluster": "loxodemo",
    "Member": "pgdeb02",
    "Host": "10.200.0.12",
    "Role": "Leader",
    "State": "running",
    "TL": 2
  },
  {
    "Cluster": "loxodemo",
    "Member": "pgdeb03",
    "Host": "10.200.0.13",
    "Role": "Replica",
    "State": "streaming",
    "TL": 2,
    "Lag in MB": 0
  }
]
```

## ICI AVEC ETCD

```
etcdctl get --from-key /service/loxodemo/  
/service/loxodemo/config  
{  
  "ttl":30,"loop_wait":10,"retry_timeout":10,"maximum_lag_on_failover":1048576,"postgresql":{"use_pg  
/service/loxodemo/failover  
{  
/service/loxodemo/history  
[[1,83886240,"no recovery target specified","2024-11-11T15:35:04.535410+00:00","pgdeb02"]]  
/service/loxodemo/initialize  
7436044389321027980  
/service/loxodemo/leader  
pgdeb02  
/service/loxodemo/members/pgdeb01  
{  
  "conn_url":"postgres://10.200.0.11:5432/postgres","api_url":"http://10.200.0.11:8008/patroni","sta  
/service/loxodemo/members/pgdeb02  
{  
  "conn_url":"postgres://10.200.0.12:5432/postgres","api_url":"http://10.200.0.12:8008/patroni","sta  
/service/loxodemo/members/pgdeb03  
{  
  "conn_url":"postgres://10.200.0.13:5432/postgres","api_url":"http://10.200.0.13:8008/patroni","sta  
/service/loxodemo/status  
{  
  "optime":83886696,"slots":{"pgdeb02":83886696},"retain_slots":["pgdeb01","pgdeb02"]  
}
```

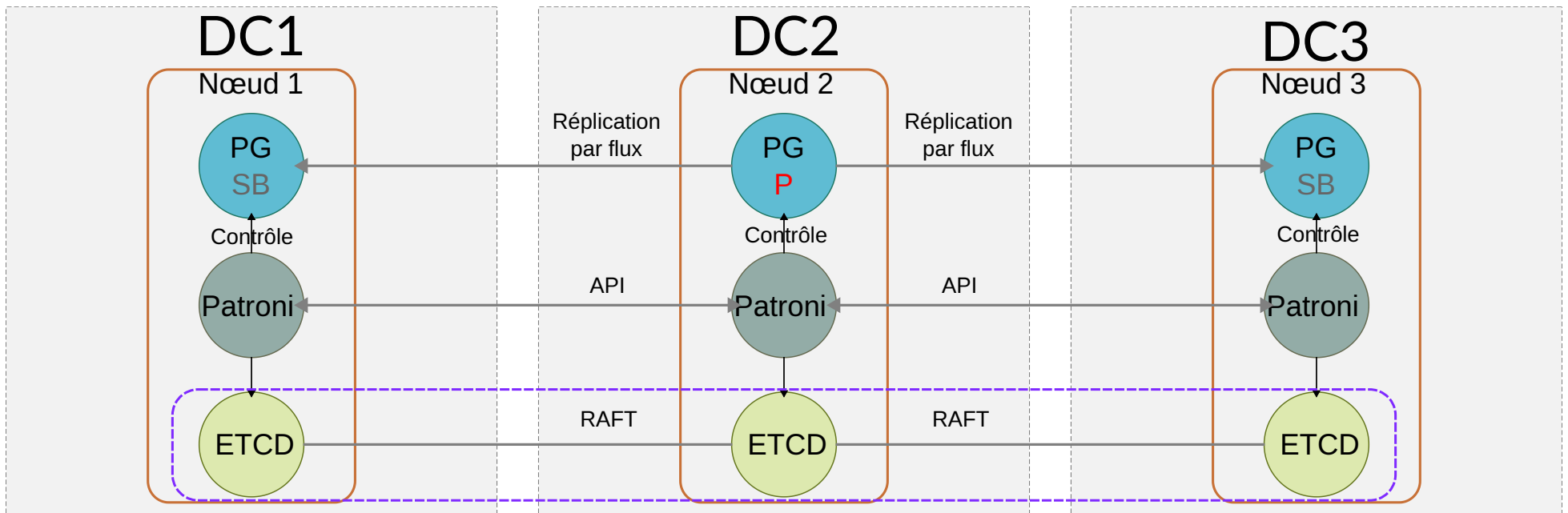
## PENSEZ BIEN

- Quorum pas atteint  $\Rightarrow$  tous les nœuds passent read-only (demote du primaire)
- Mode DCS failsafe depuis patroni 3.0
- Outil de sauvegarde = offload du primaire lors de la création de standbys
- Architectures:
  - Simple
  - Cascade
  - Réplication asynchrone/synchrone (configuration spécifique!!!)

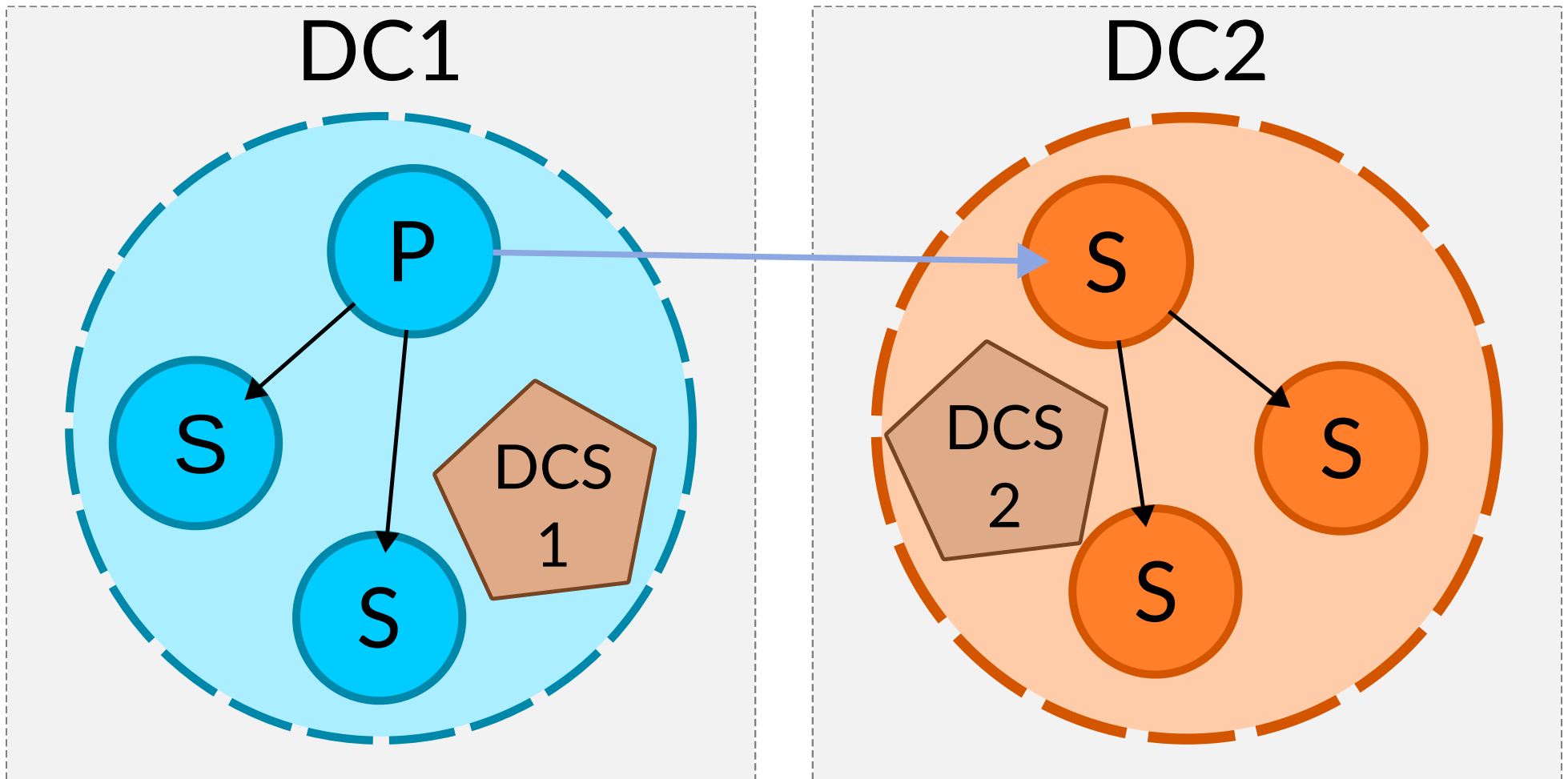


# CAS D'USAGE

1 DC est vu comme un standby de l'autre.



1 DC est vu comme un standby de l'autre.



MERCI!

