

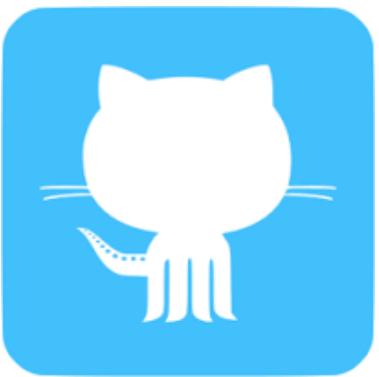
Stockage et analyse temps réel d'événements avec Riak chez Booking.com

Damien Krotkine

Damien Krotkine



- Ingénieur chez Booking.com



- github.com/dams



- [@damsieboy](https://twitter.com/@damsieboy)



- [dkrotkine](https://www.linkedin.com/in/dkrotkine)

Booking.com

Planet Earth's #1 Accommodation Site.

- 800,000 chambres réservées par jour

INTRODUCTION

B Booking.com: 606,351 hot x

www.booking.com

V GR GL monitor events Web iPhoto | open s Mon Modem

Booking.com

Recently Seen My Lists Sign In Manage booking

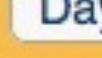
Make your next booking on the go! Download Booking.com's free apps  

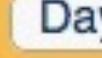
Find the Best Deals

606,000+ hotels, apartments, villas and more...

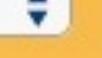
Destination/Hotel Name:
E.g. Best Western New York

Traveling for: Business Leisure

Check-in Date
 Day  Month 

Check-out Date
 Day  Month 

I don't have specific dates yet

Guests 2 adults, 0 children 

Additional search options

Search

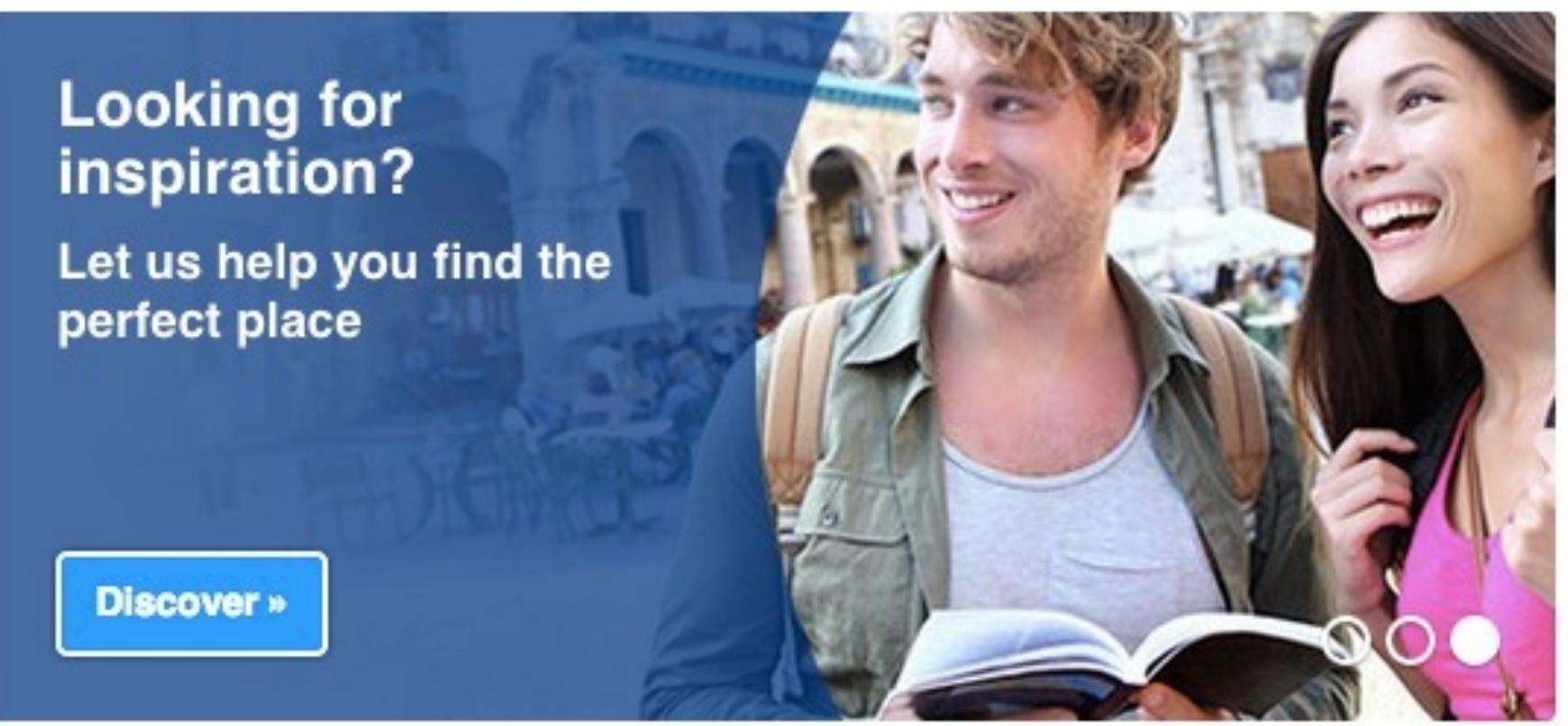
Subscribe for a 10% discount
Unlock Member Deals and customized inspiration

FREE cancellation on most rooms!

Looking for inspiration?

Let us help you find the perfect place

Discover »



Lyon 
313 places to stay

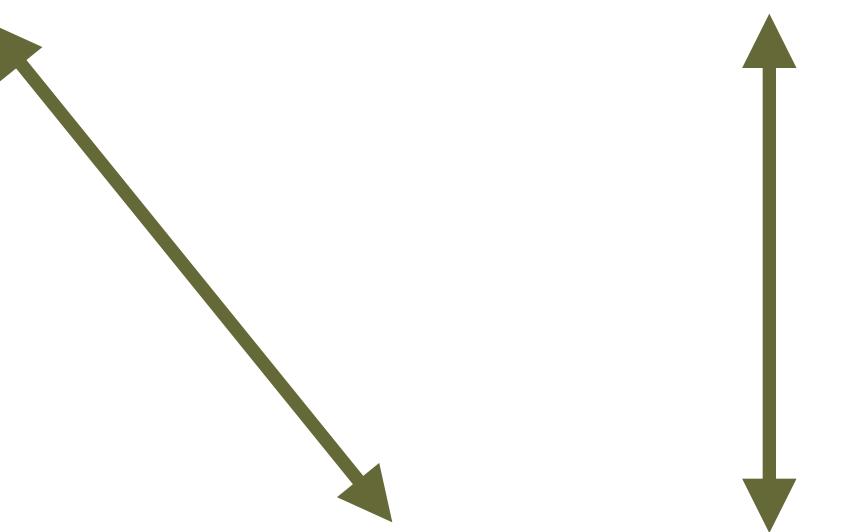
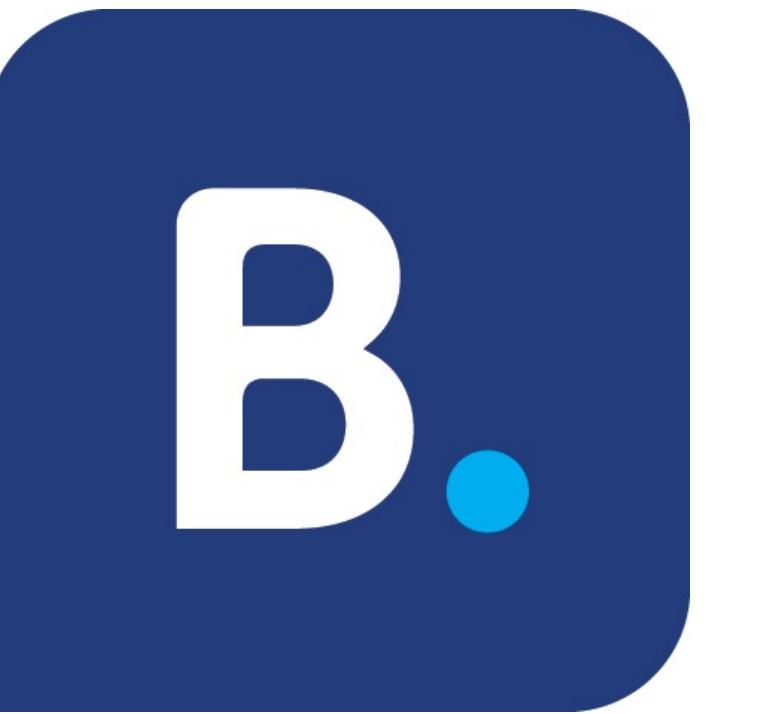
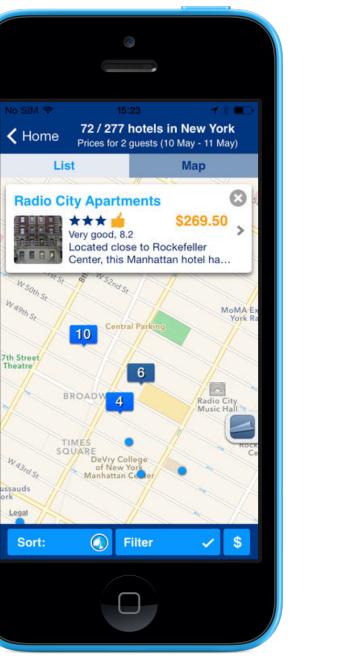


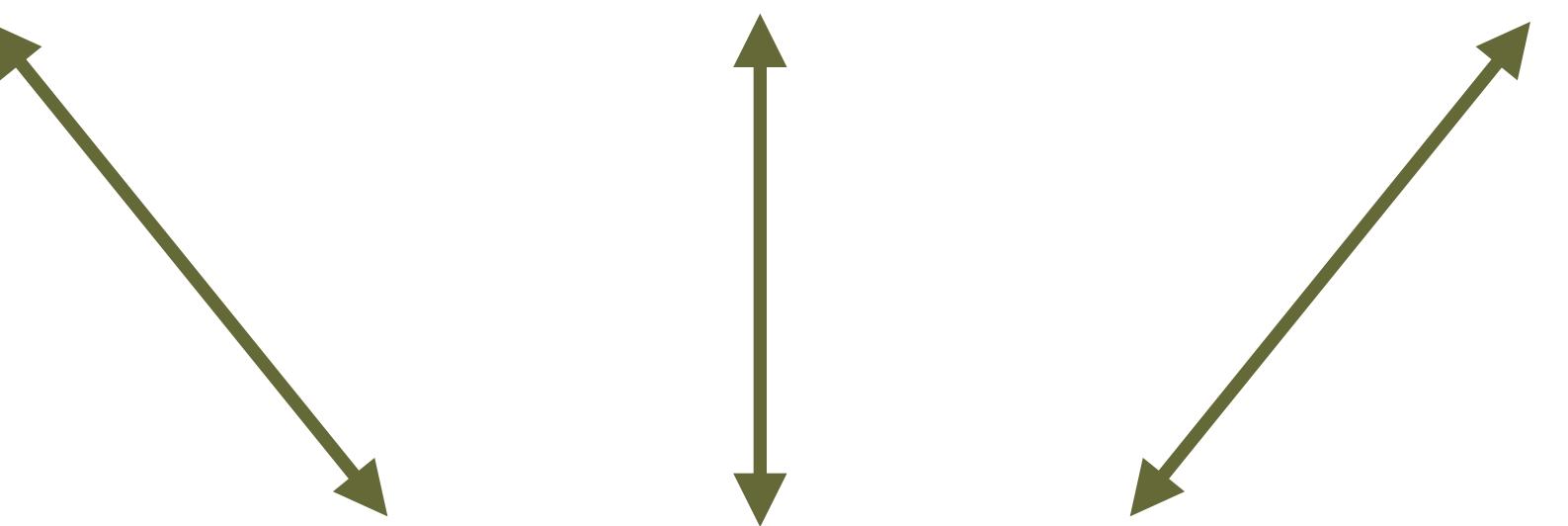
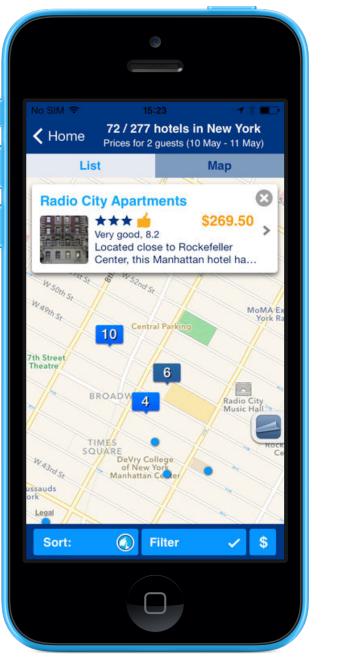
32 people looking at this destination

Find somewhere to stay in Lyon »



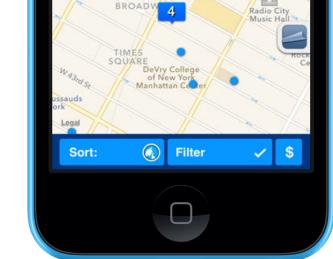
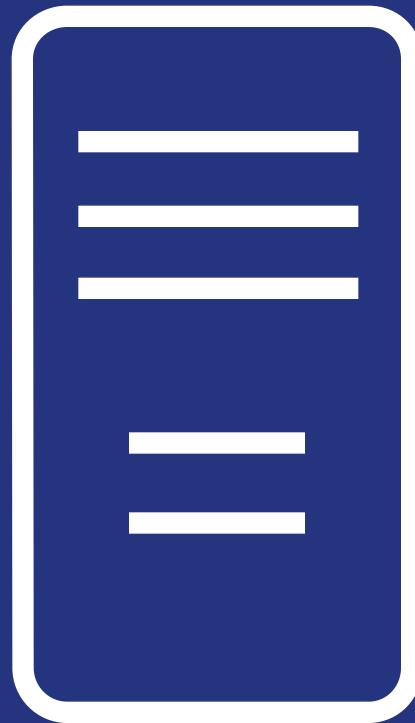
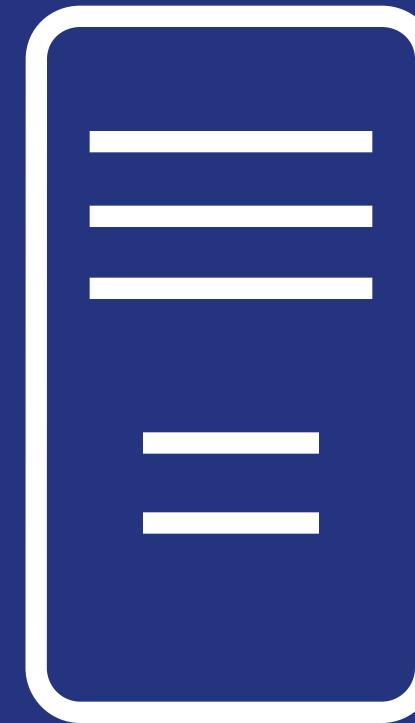
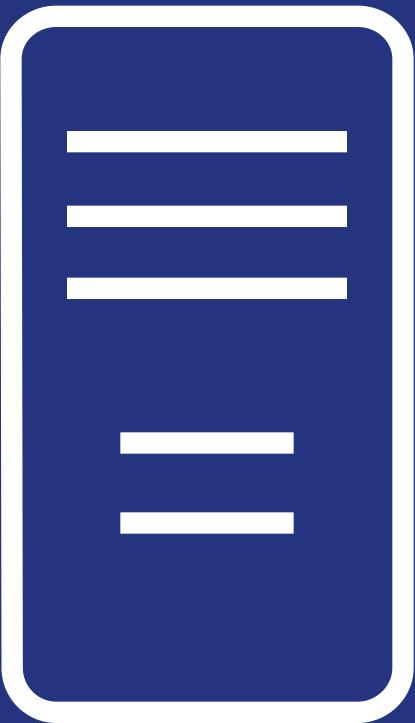






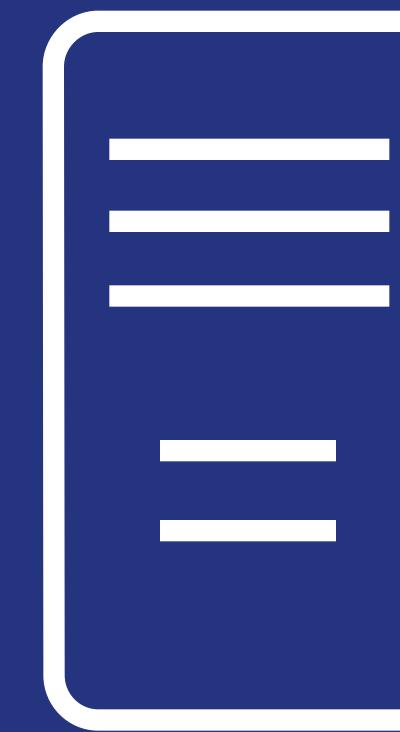
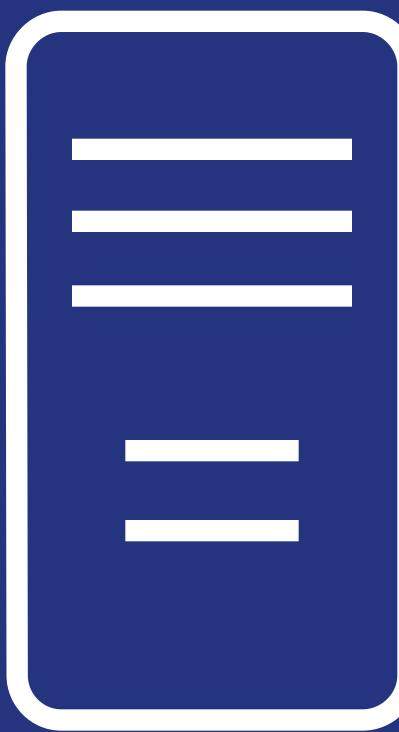
frontend

www mobi API



frontend

www mobi API

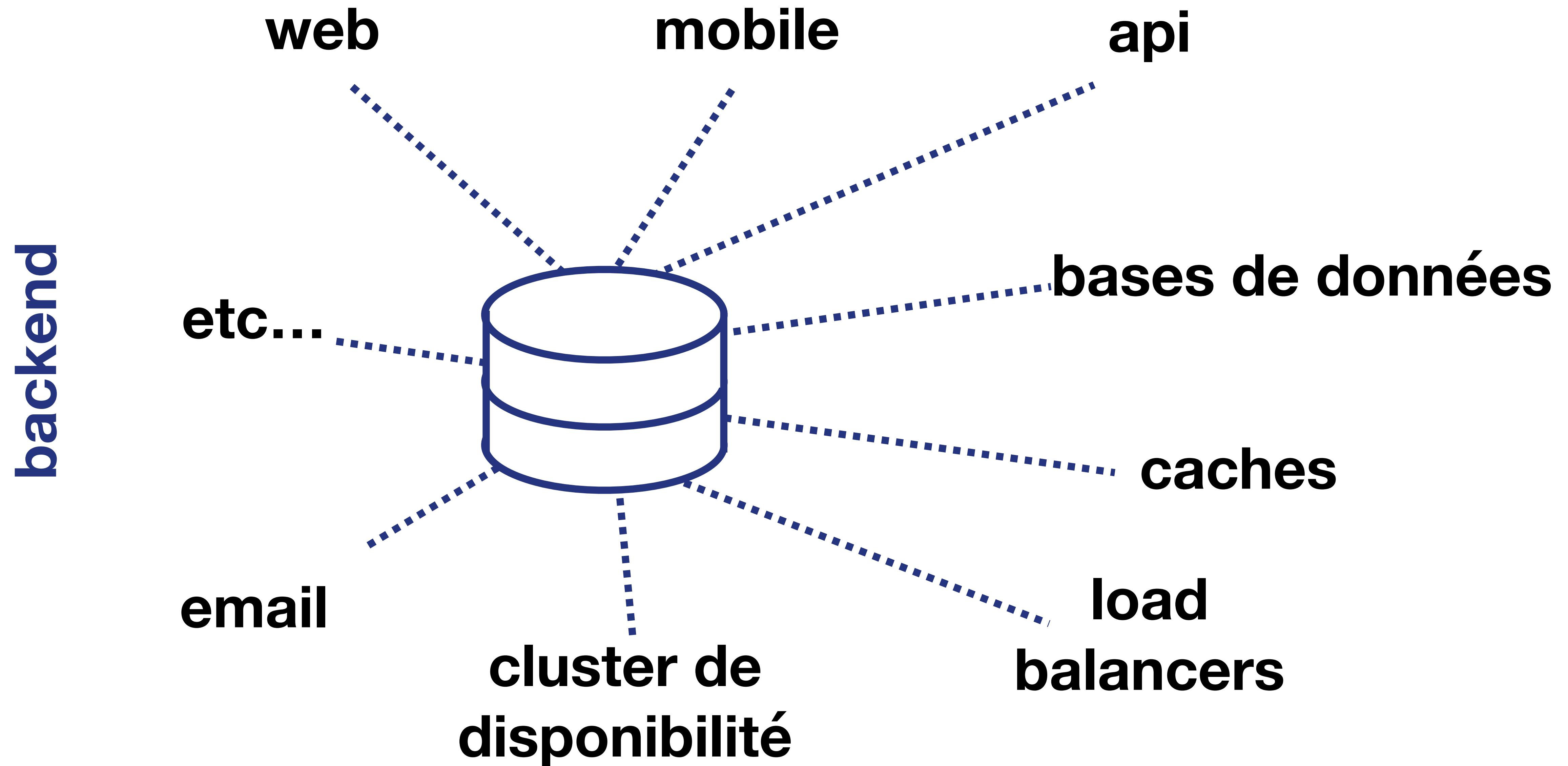


backend



stockage des événements

**événements: info sur
l'état des systèmes**



QU'EST-CE QU'UN EVENEMENT ?

STRUCTURE D'UN EVENEMENT

- **Information** sur un **système**
- Données
- HashMap profond
- Timestamp
- Data Center + Type + Sous-type
- Le reste: données spécifiques
- **Sans Schema**

```
{ timestamp => 12345,
  type => 'WEB',
  subtype => 'app',
  dc => 1,
  action => { is_normal_user => 1,
               pageview_id => '188a362744c301c2',
               # ...
             },
  tuning => { the_request => 'GET /display/...'
               bytes_body => 35,
               wallclock => 111,
               nr_warnings => 0,
               # ...
             },
  # ...
}
```

```
{ type => 'FAV' ,  
subtype => 'fav' ,  
timestamp => 1401262979 ,  
dc => 1 ,  
tuning => {  
    flatav => {  
        cluster => '205' ,  
        sum_latencies => 21 ,  
        role => 'fav' ,  
        num_queries => 7  
    }  
}  
}
```

PROPRIETES DU FLUX D'EVENEMENTS

- Lecture seule
- Sans schema
- Flux continu, ordonné, temporel
- 15 K événements par sec
- **1.25 Milliards** événements par jour
- pics à 70 MB/s, min 25MB/s
- 100 GB par heure

UTILISATION

DEFINIR LES BESOINS

- Avant d'aborder le stockage
- Penser à l'utilisation

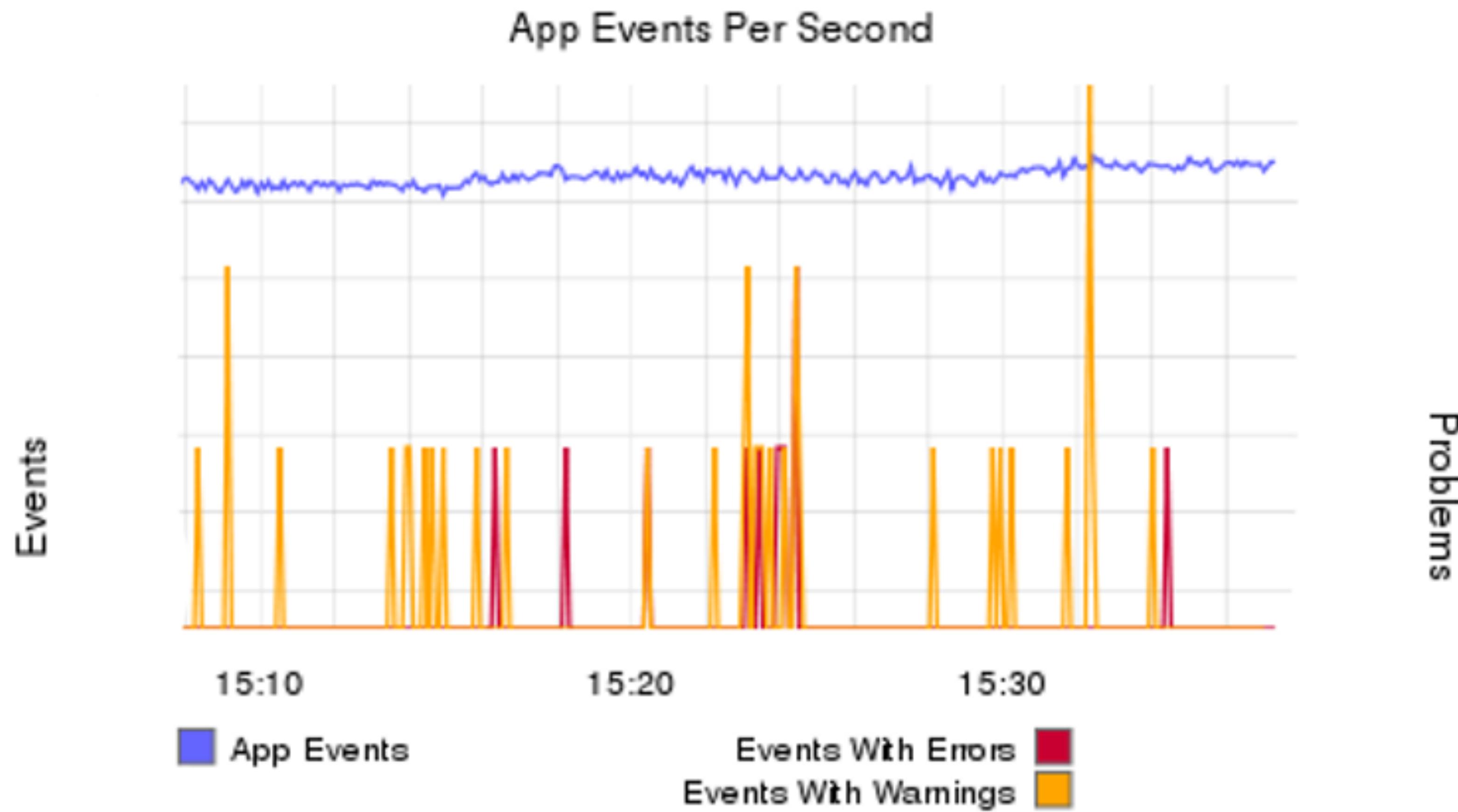
UTILISATION

1. GRAPHS
2. PRISE DE DECISION
3. ANALYSE COURT TERME
4. A/B TESTING

GRAPHS

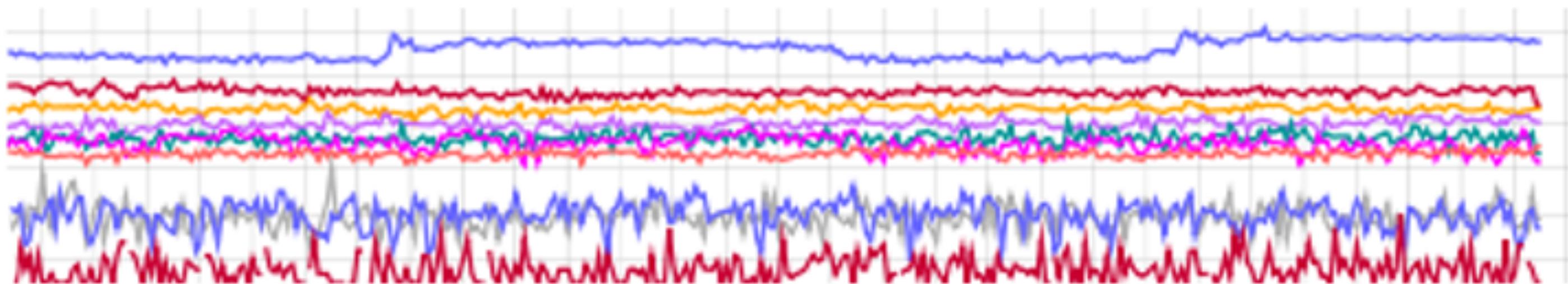
- Graph en temps réel (lag = quelques secondes)
- Graph pour le plus de systèmes possibles
- Etat de santé général de la plateforme

GRAPHS



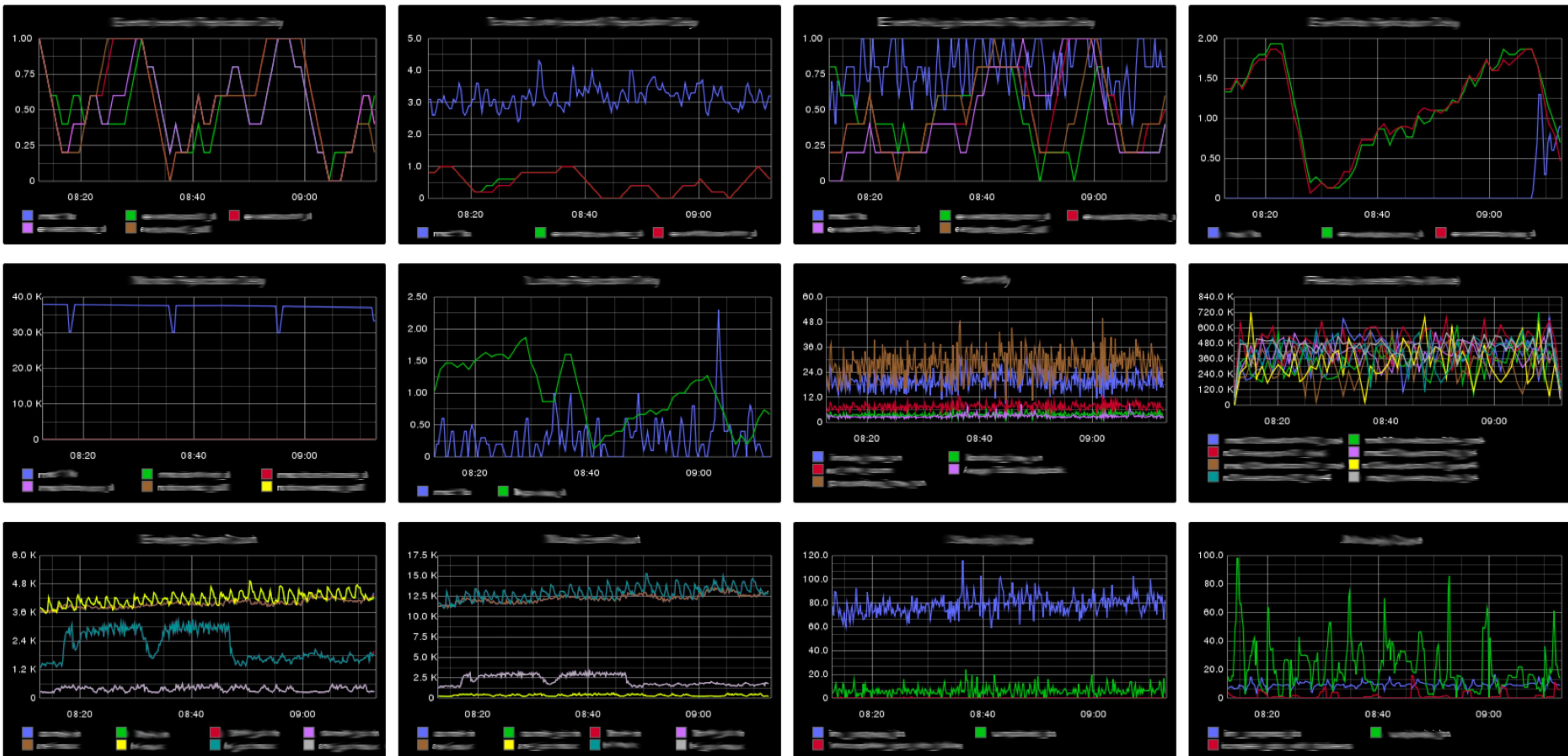
GRAPHS

App Top 10 Abnormal User Type Reasons Per Second Log10

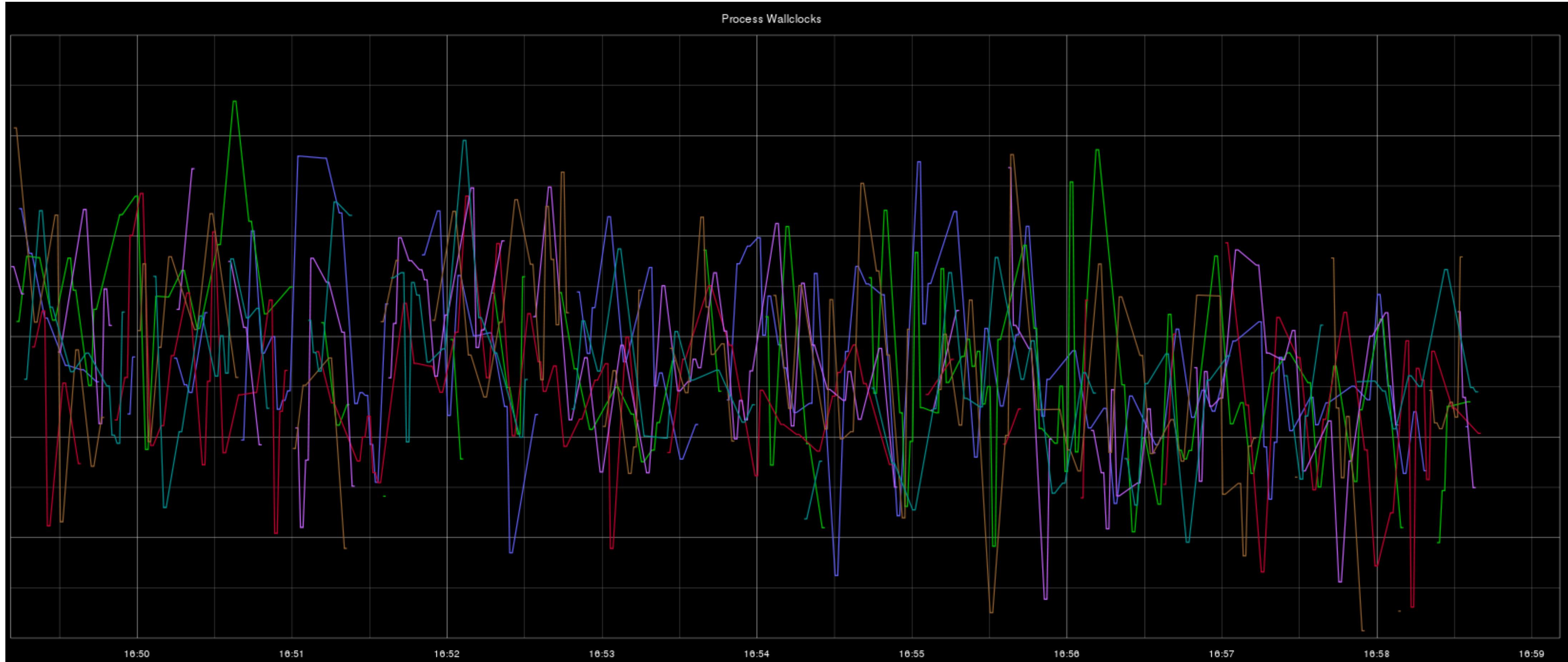


- abnormal_suspect_isp
- abnormal_ip_status_suspect
- abnormal_extranet_user
- abnormal_old_original_cookie
- abnormal_legacy_static_redirect_images
- abnormal_no_reason
- abnormal_legacy_static_redirect_static
- abnormal_too_many_visits_no_booking
- abnormal_internal_ip_internal
- abnormal_legacy_static_redirect_data_sp_aff

DASHBOARDS



META GRAPHS



USAGE

1. GRAPHS
2. PRISE DE DECISION
3. ANALYSE COURT TERME
4. A/B TESTING

DECISION MAKING

- Decisions Stratégiques (philosophie “use facts”)
- Court terme / long terme
- Décisions techniques / non techniques

USAGE

1. GRAPHS
2. PRISE DE DECISION
3. ANALYSE COURT TERME
4. A/B TESTING

SHORT TERM ANALYSIS

- 10 dernières secondes -> il y a 8 jours
- Déploiement de code et rollback
- Détecteur d'Anomalies

USAGE

1. GRAPHS
2. PRISE DE DECISION
3. ANALYSE COURT TERME
4. A/B TESTING

A/B TESTING

- Philosophie “use facts”
- Implique A/B testing
- Concept d’”experiments”: Experiments
- Evénements procurent la matière première

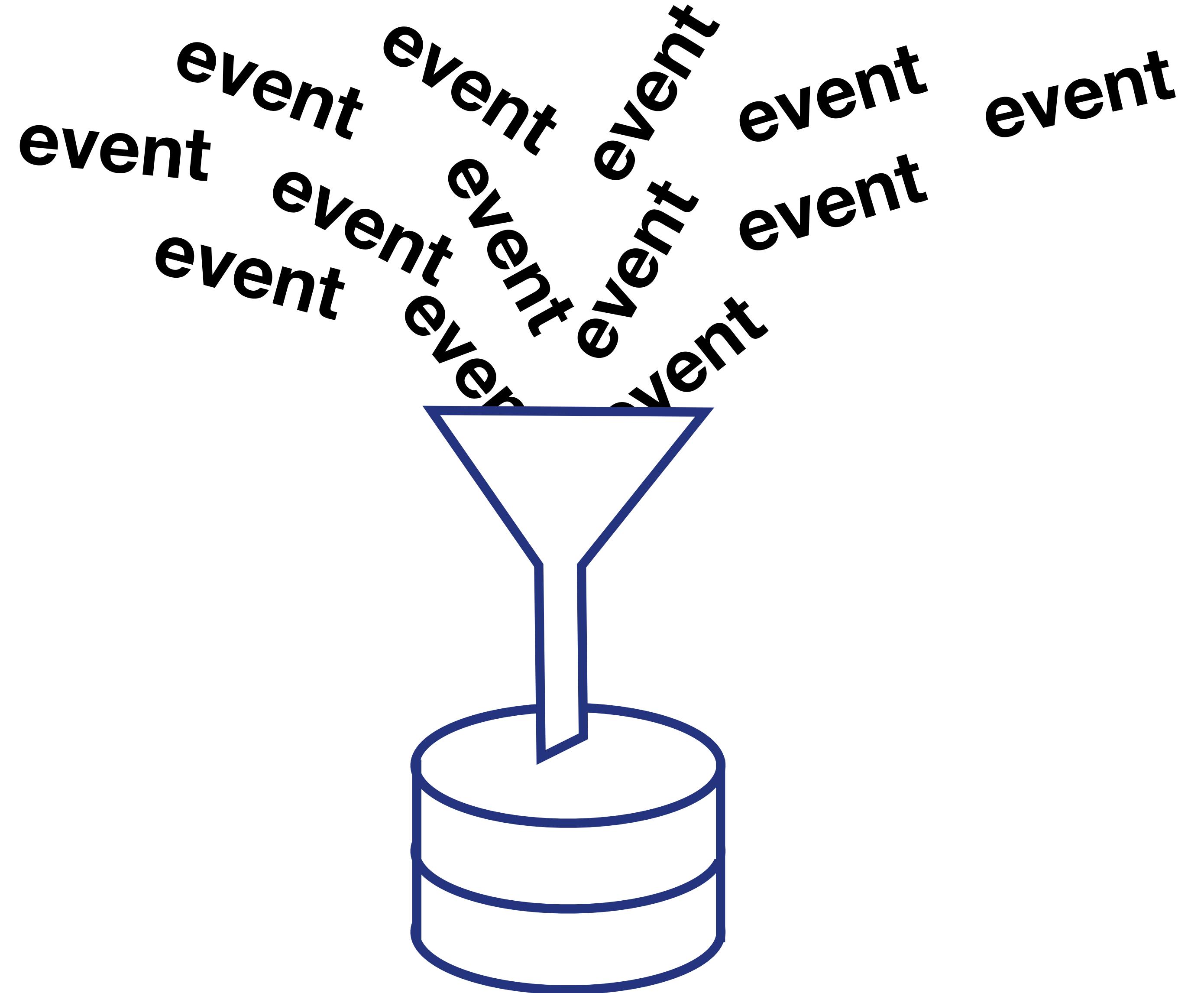
AGGREGATION D'EVENEMENTS

AGGREGATION D'EVENTEMENTS

- Regrouper les événements
- Granularité nécessaire: seconde

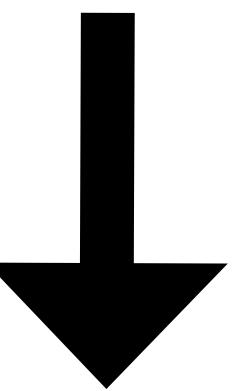
SERIALISATION

- JSON ne nous allait pas (lent, gros, manque de fonct.)
- Création de Sereal en 2012
- « *Sereal, a new, **binary** data serialization format that provides high-**performance**, **schema-less** serialization* »
- <https://github.com/Sereal/Sereal>

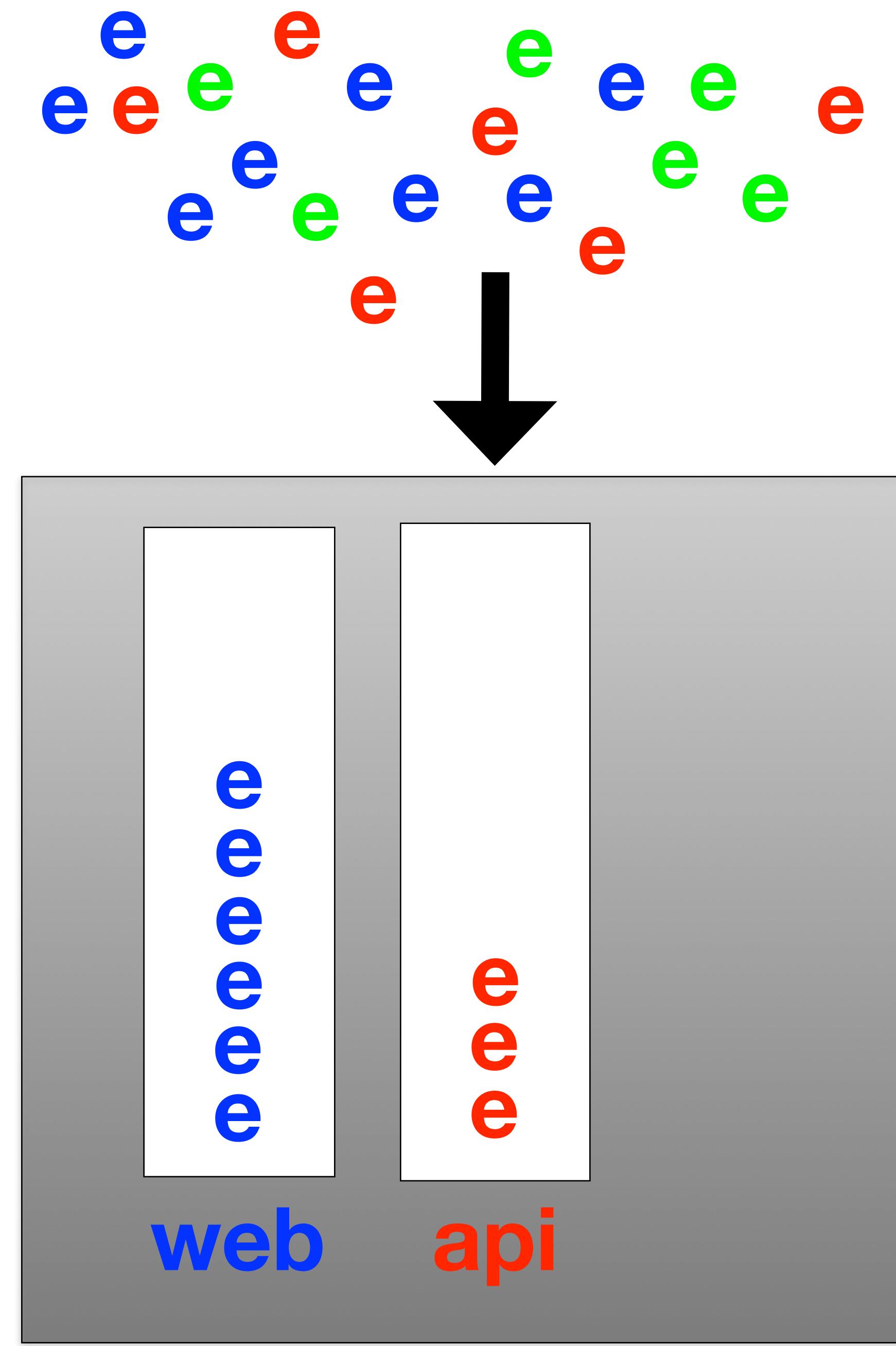


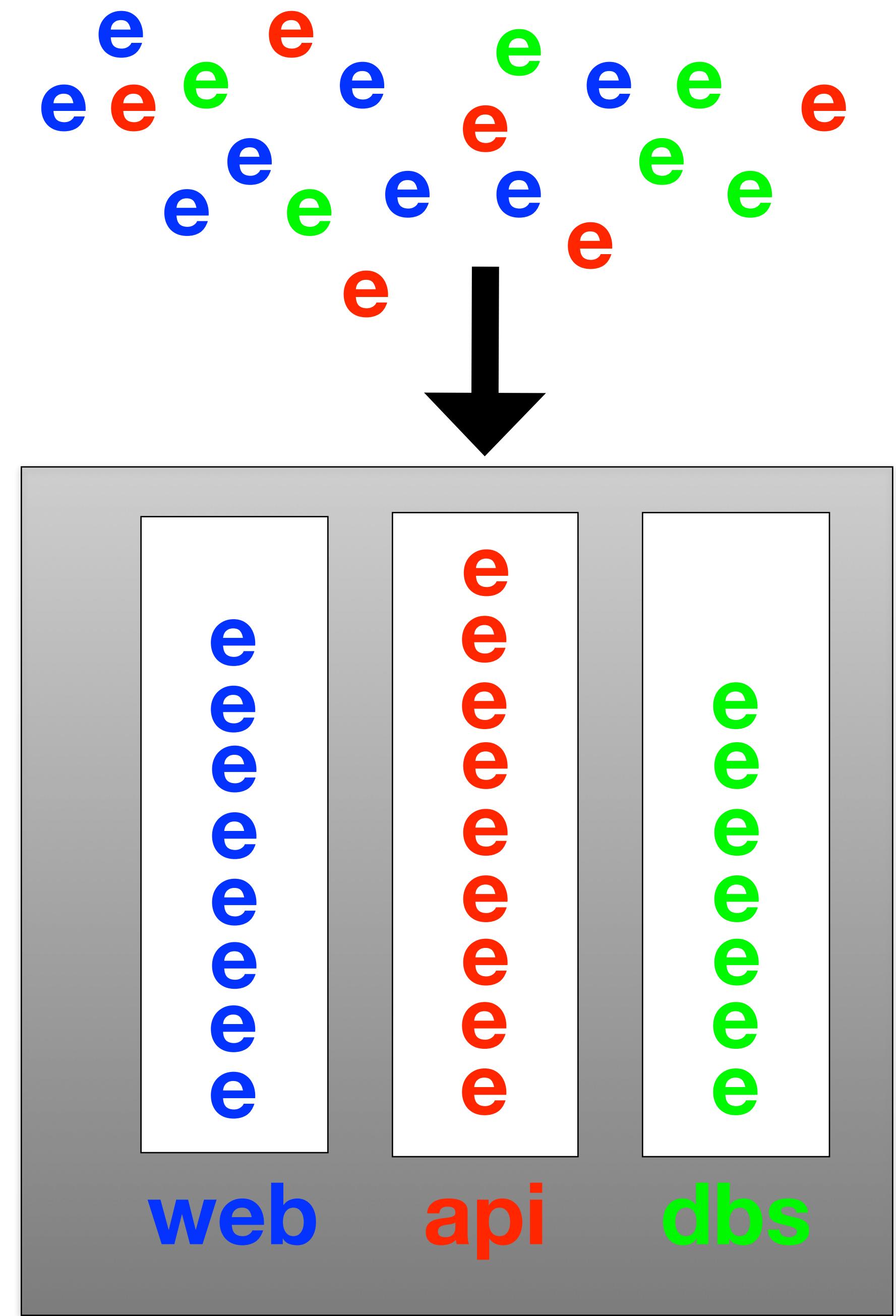
stockage des événements

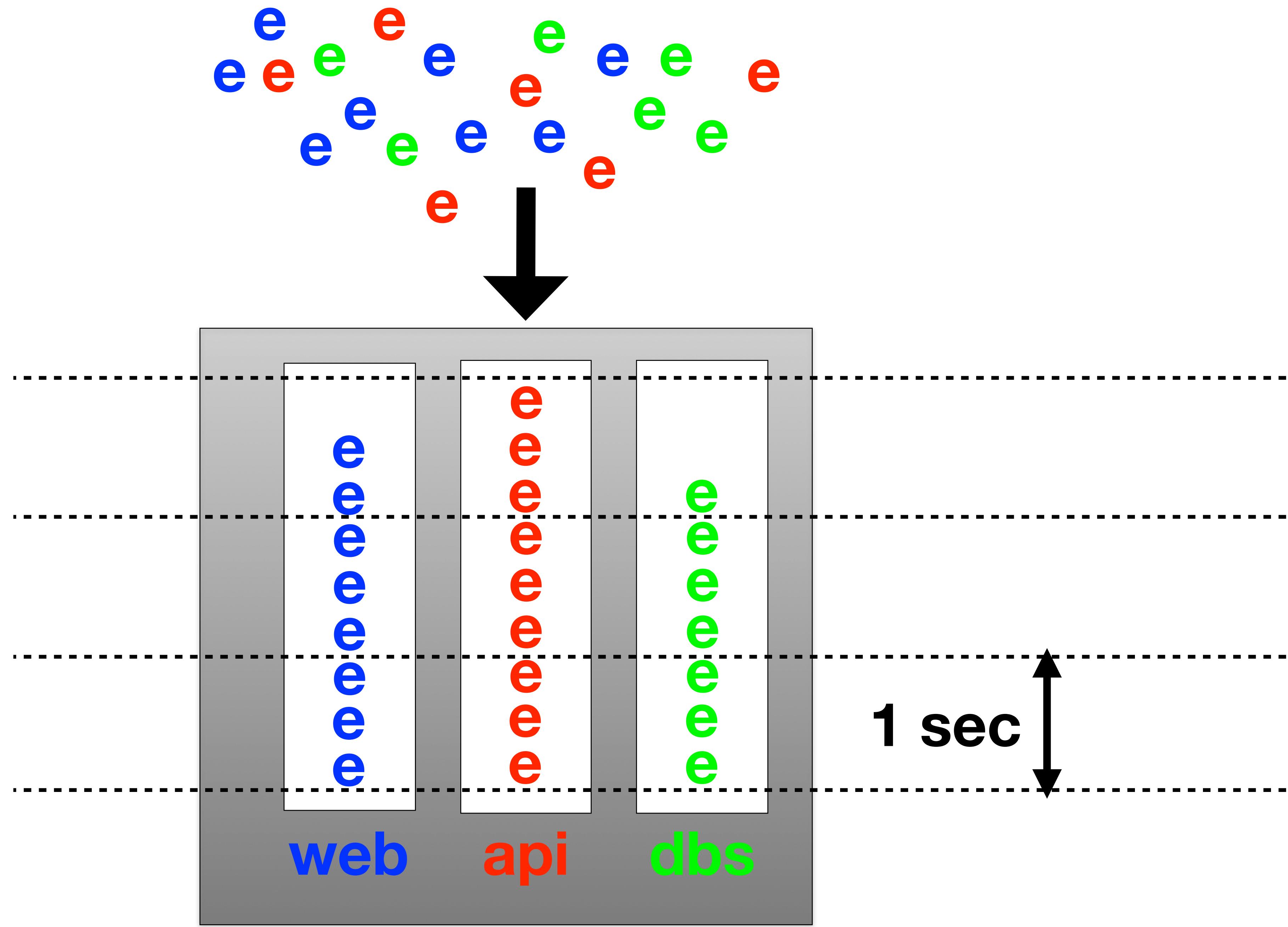
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e
e e e e e e e e e

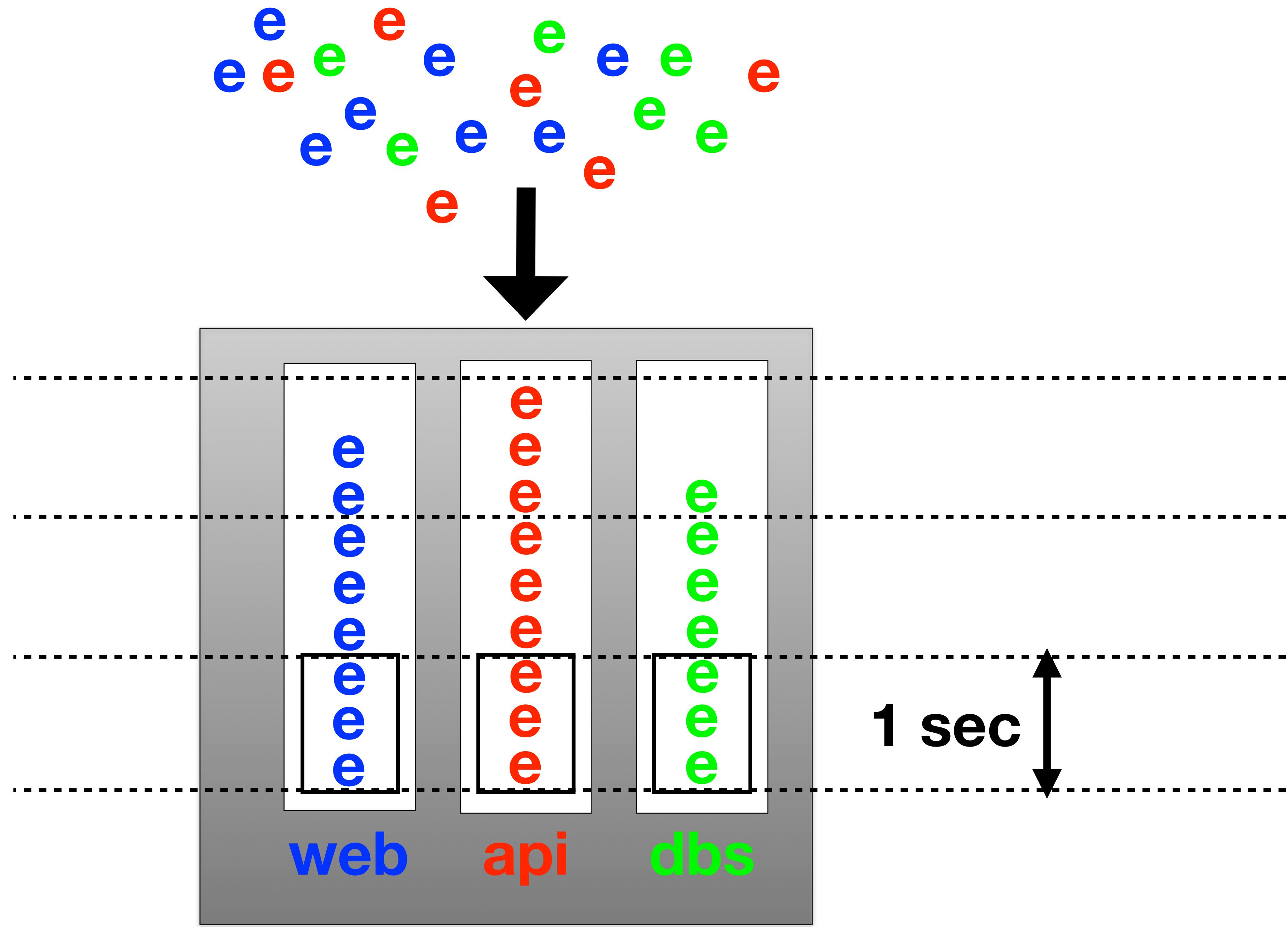


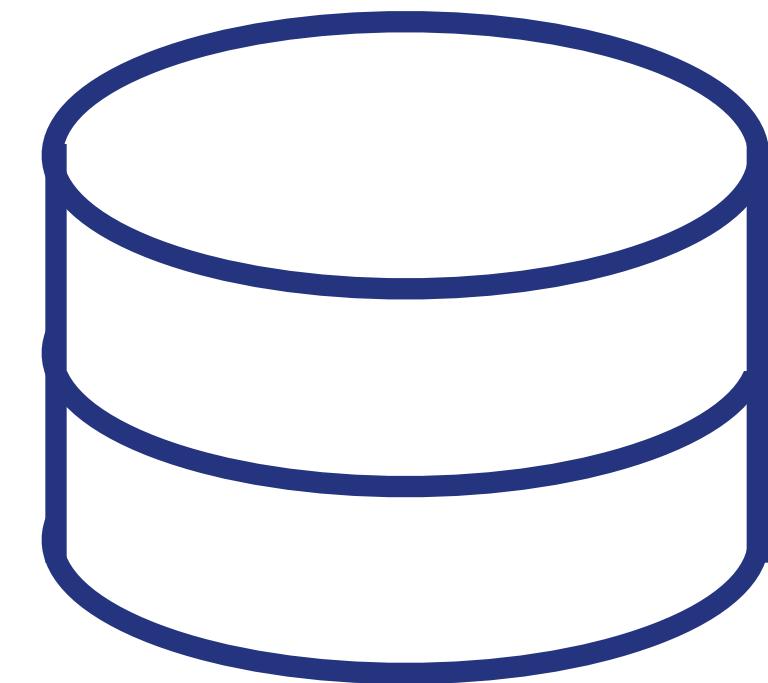
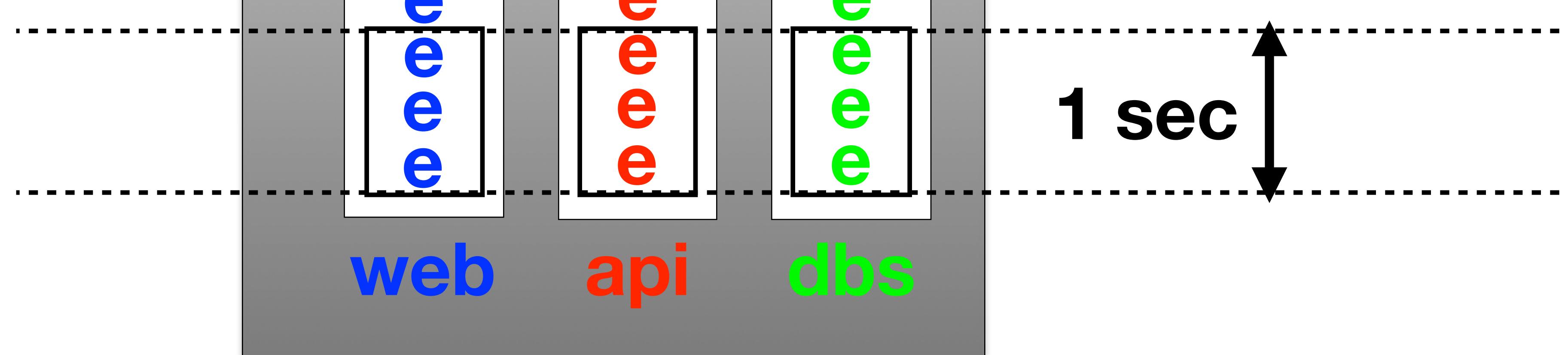
LOGGER



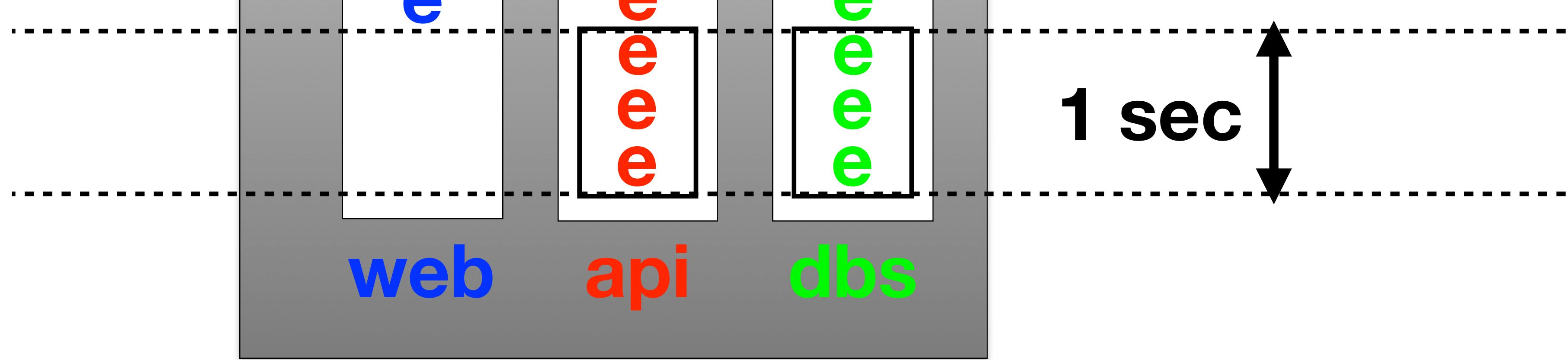




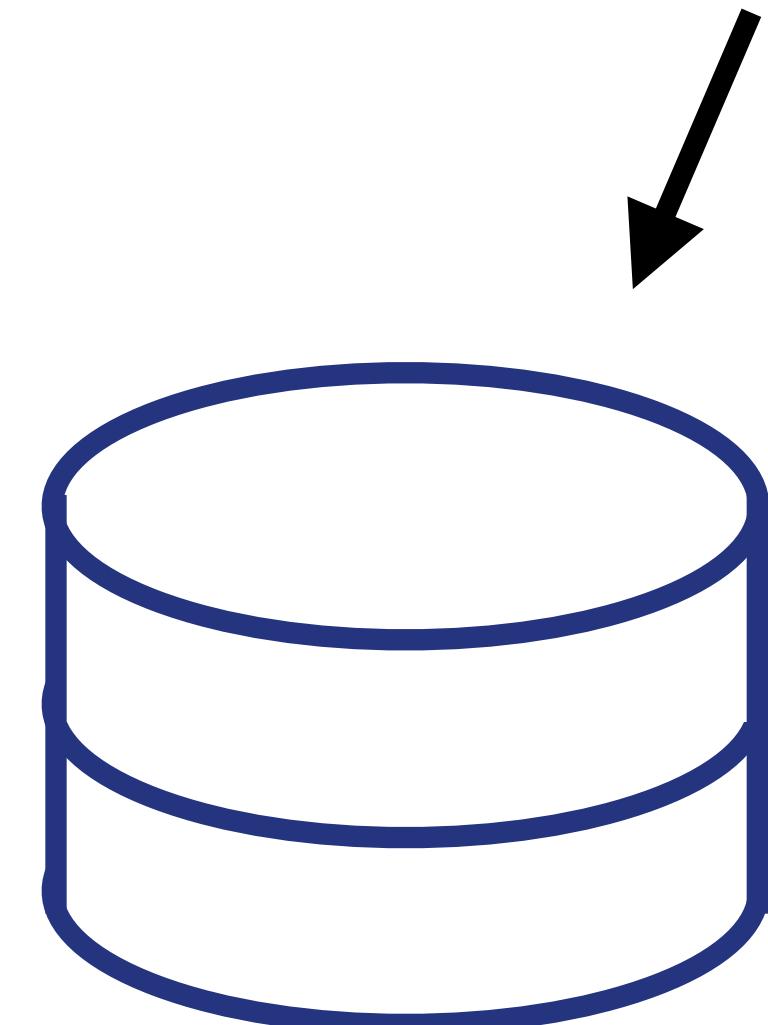




stockage des événements



eee → **re-sérialisation
+ compression**



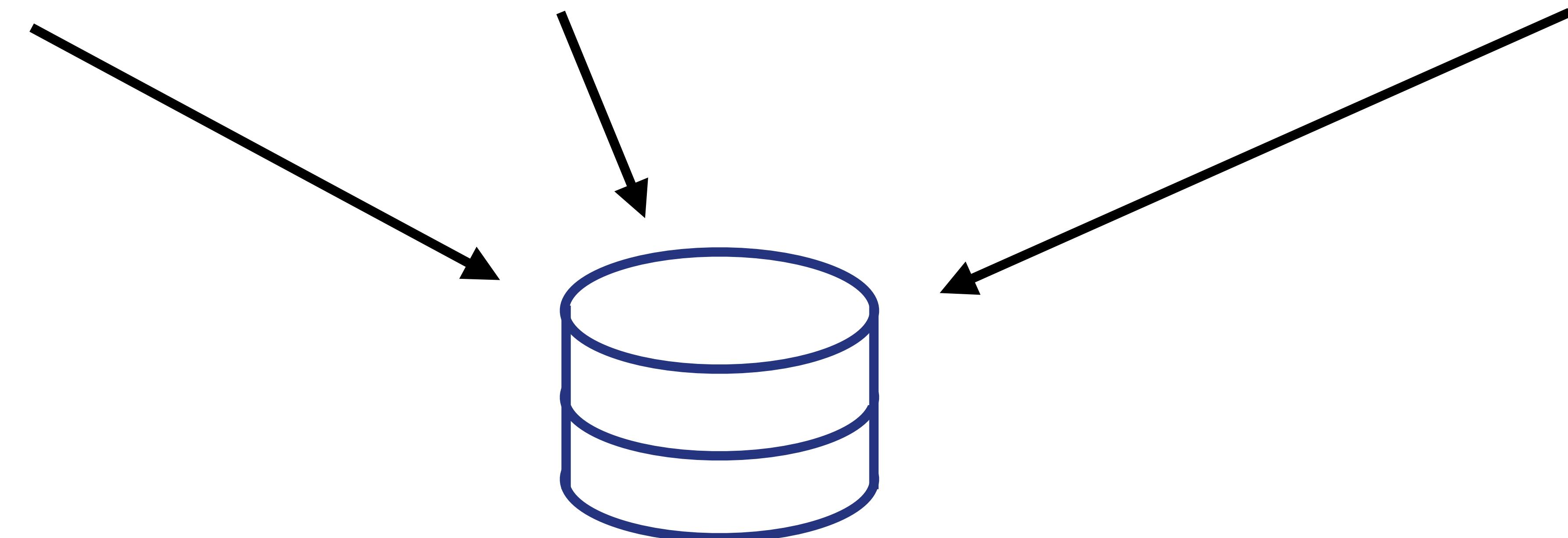
stockage des événements

LOGGER

LOGGER

...

LOGGER



stockage des événements

STOCKAGE

BESOINS

- Sécurité de stockage
- Performance écriture massive
- Performance lecture massive
- Administration Facile
- Très scalable

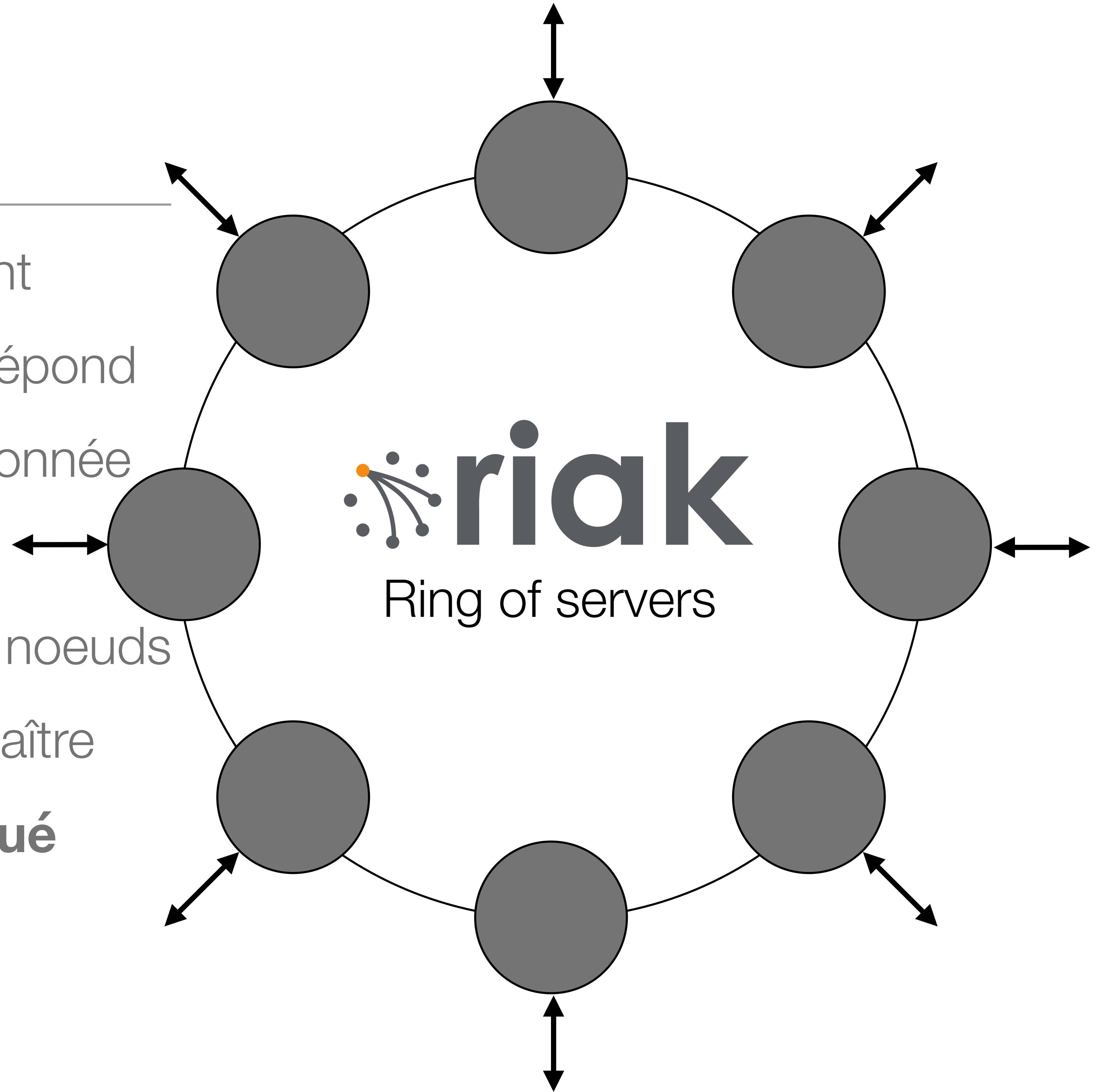
NOUS AVONS CHOISI RIAK

- Sécurité: cluster, distribué, replication, Ring très robuste
- Performances bonnes et prévisibles
- Très facile à administrer
- Fonctions avancées (MapReduce, triggers, 2i, CRDTs ...)
- Riak Search
- Réplication Multi-Datacenter



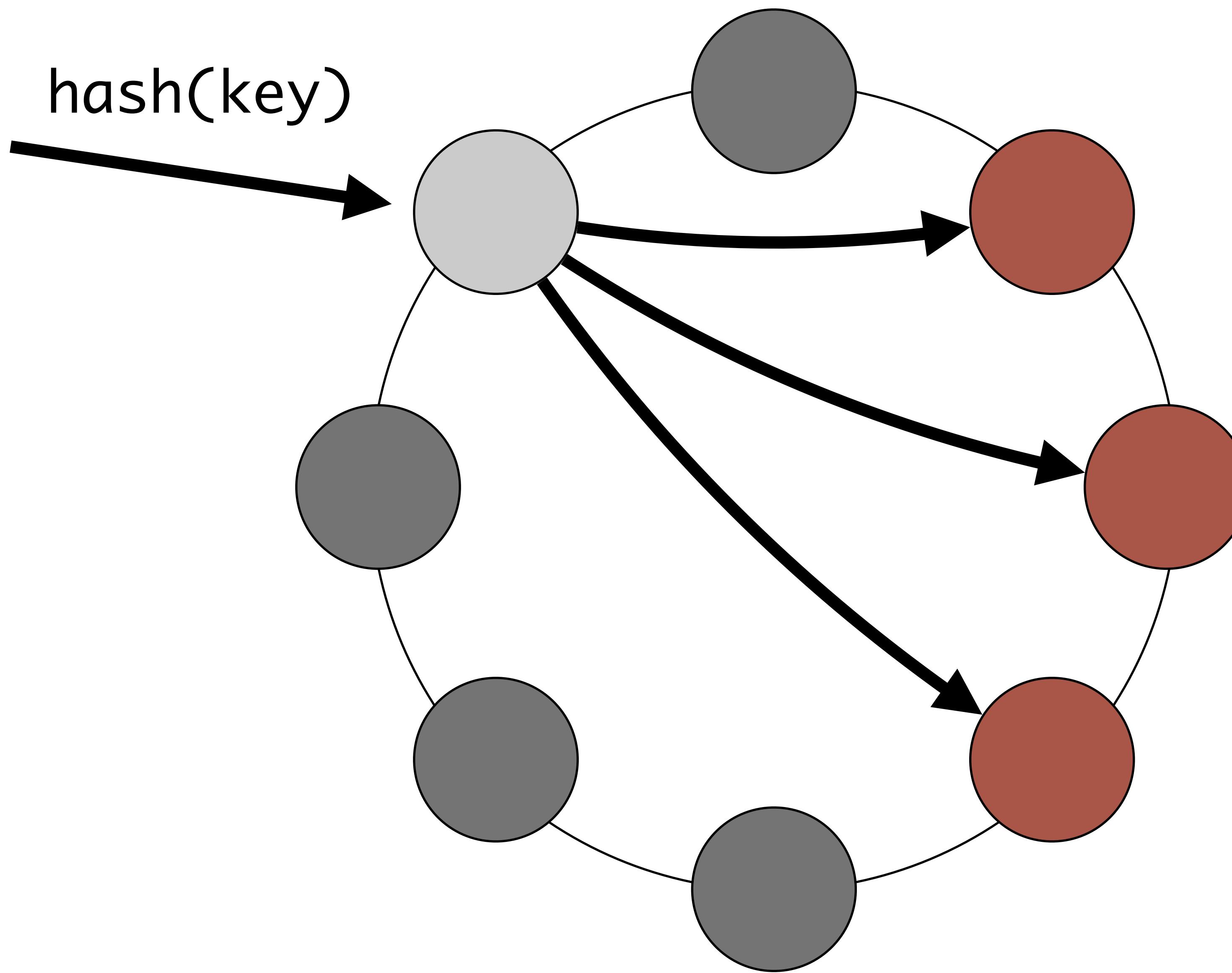
CLUSTER

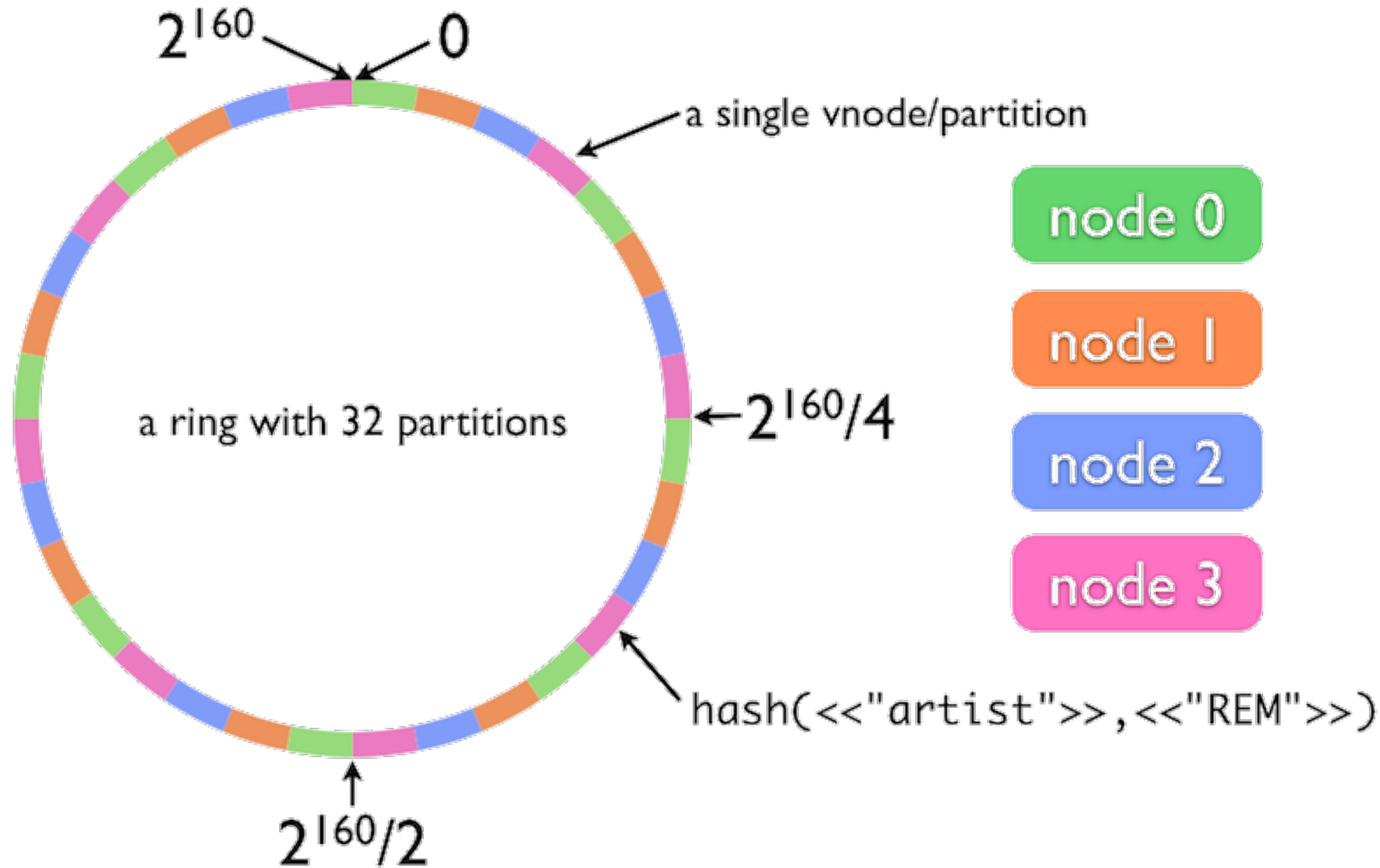
- Hardware courant
- Chaque noeud répond
- RéPLICATION de donnée
- Gossip entre les noeuds
- Pas de noeud maître
- Système **distribué**





hash(key)





STOCKAGE CLEF VALEUR

- Espaces de noms: **bucket**
- Valeurs: **opaque** ou CRDTs

RIAK: FONCTIONALITES AVANCEES

- MapReduce
- Index secondaires (2i)
- Riak Search
- Réplication Multi-DataCenter

MULTI-BACKEND DE STOCKAGE

- Bitcask
- Leveldb
- Memory

BACKEND: BITCASK

- Backend basé sur un système de logs
- AOF (Append-only files)
- Expiration de données fine
- Performance prévisibles (1 seek max)
- Parfait pour des données séquentielles

CONFIGURATION DU CLUSTER

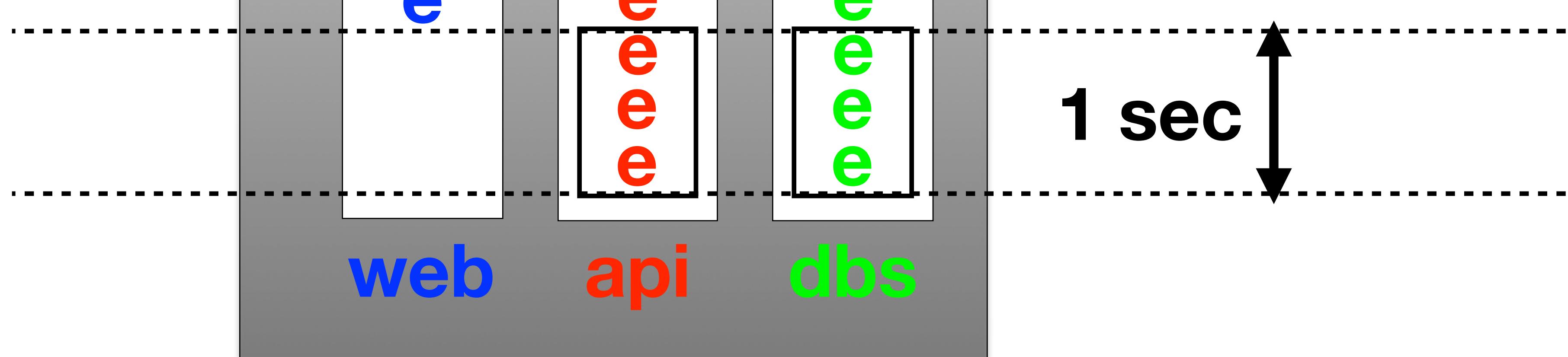
ESPACE DISQUE NECESSAIRE

- 8 jours
- 100 GB par heure
- RéPLICATION = 3
- $100 * 24 * 8 * 3$
- Besoin de **60 T**

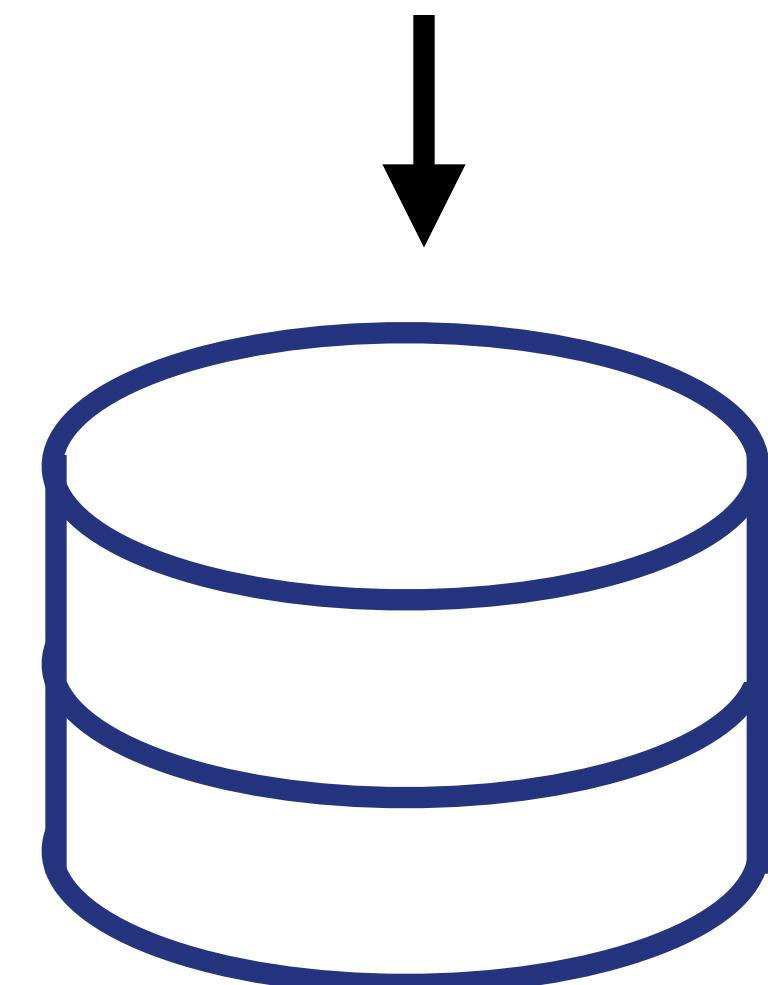
HARDWARE UTILISE

- 12 puis 16 nodes (bientôt 24)
- 12 CPU cores (Xeon 2.5Ghz) par noeud
- 192 GB RAM
- réseau 1 Gbit/s
- 8 TB (raid 6)
- Espace total: 128 TB

DATA DESIGN



regroupé par EPOCH / DC / CELL / TYPE / SUBTYPE
500 KB max par morceaux



stockage des événements

DONNEES

- Bucket : “data”
- Clef: “12345:1:cell0:WEB:app:chunk0”
- Valeur: Liste sérialisée compressée des events HashMap
- 200 clefs par secondes

META-DONNEES

- Bucket : “metadata”
- Clef: <epoch>-<dc> “1428415043-2”
- Valeur: liste de clef données:
 - [“1428415043:1:cell0:WEB:app:chunk0”,
“1428415043:1:cell0:WEB:app:chunk1”
...
“1428415043:4:cell0:EMK::chunk3”]
- PSV (pipe separated value)

ECRIRE LES DONNEES

ECRIRE LES DONNEES

- Dans chaque DC, cellule, Loggers poussent vers Riak
- 2 protocoles: REST ou ProtoBuf
- Chaque seconde:
 - Ecrire les valeurs “data” vers Riak, en asynchrone
 - Attendre retour succès
 - Ecrire les meta-données “metadata”

JAVA

```
Bucket DataBucket = riakClient.fetchBucket("data").execute();
DataBucket.store("12345:1:cell0:WEB:app:chunk0", Data1).execute();
DataBucket.store("12345:1:cell0:WEB:app:chunk1", Data2).execute();
DataBucket.store("12345:1:cell0:WEB:app:chunk2", Data3).execute();

Bucket MetaDataBucket = riakClient.fetchBucket("metadata").execute();
MetaDataBucket.store("12345-1", metaData).execute();
riakClient.shutdown();
```

Perl

```
my $client = Riak::Client->new(...);

$client->put(data => '12345:1:cell0:WEB:app:chunk0', $data1);
$client->put(data => '12345:1:cell0:WEB:app:chunk1', $data2);
$client->put(data => '12345:1:cell0:WEB:app:chunk2', $data3);

$client->put(metadata => '12345-1', $metadata, 'text/plain' );
```

LIRE LES DONNEES

LIRE UNE SECONDE

- Pour **une seconde** (une epoch donnée)
- Lire les meta-données pour l’<epoch>-DC
- Parser la valeur
- **Filtrer les types / sous-types**
- Récupérer les données de “data”

Perl

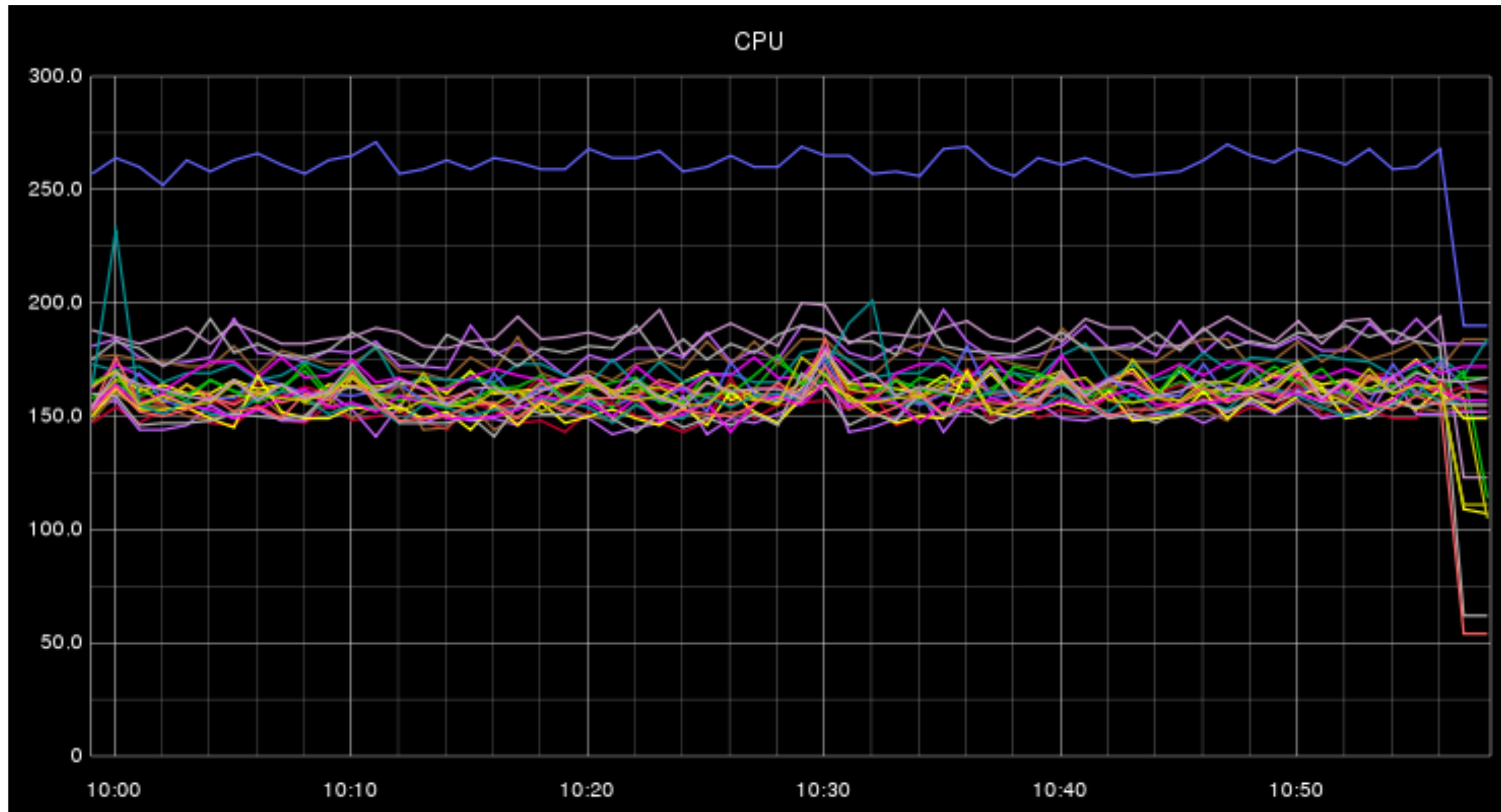
```
my $client = Riak::Client->new(...);
my @array = split '\|', $client->get(metadata => '1428415043-1');
@filtered_array = grep { /WEB/ } @array;
$client->get(data => $_) foreach @filtered_array;
```

LIRE UNE PERIODE

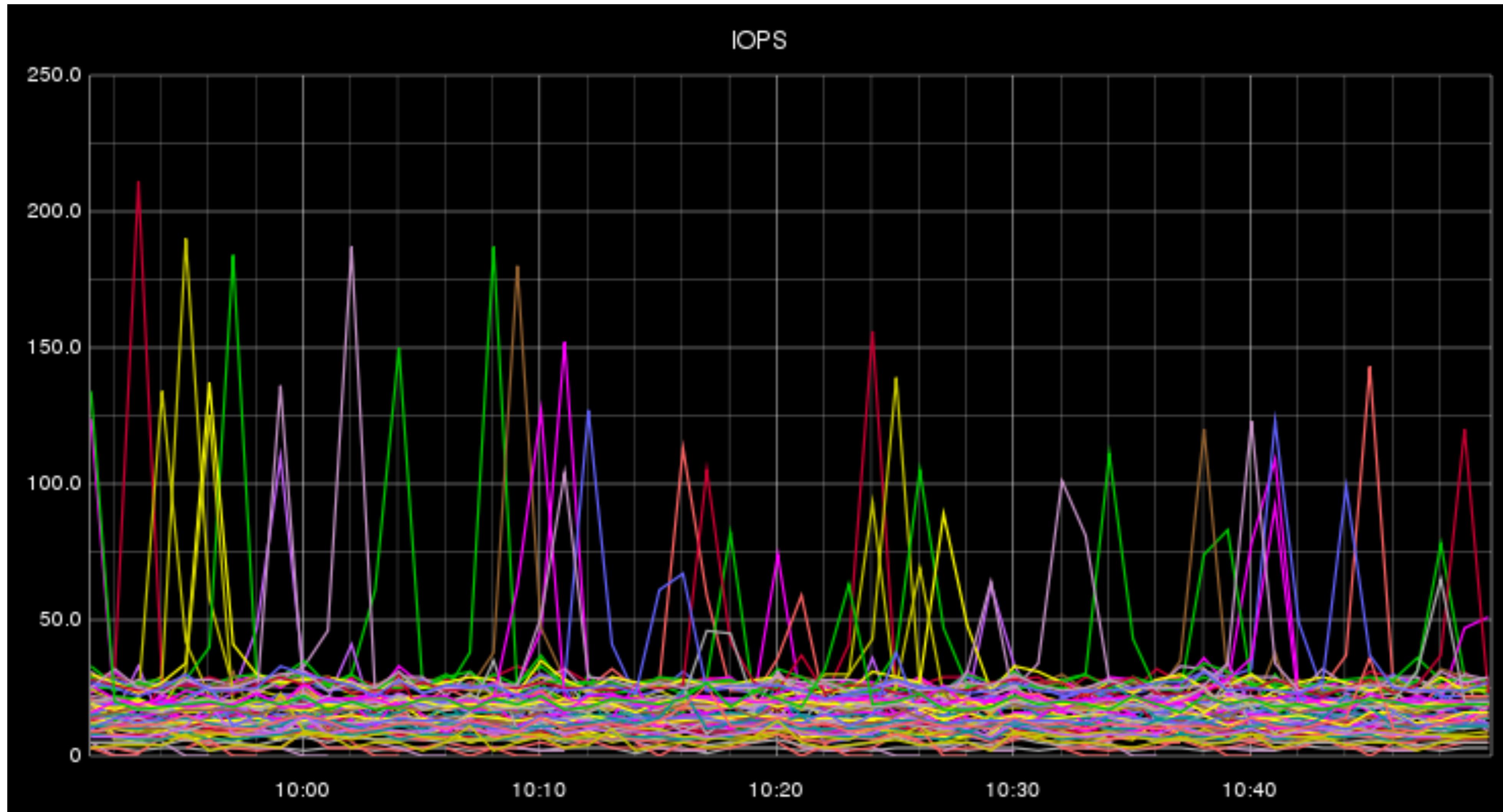
- Pour **un intervalle** epoch1 -> epoch2
- Générer la liste des epochs
- Récupérer en parallèle

RIAK CLUSTER HEALTH

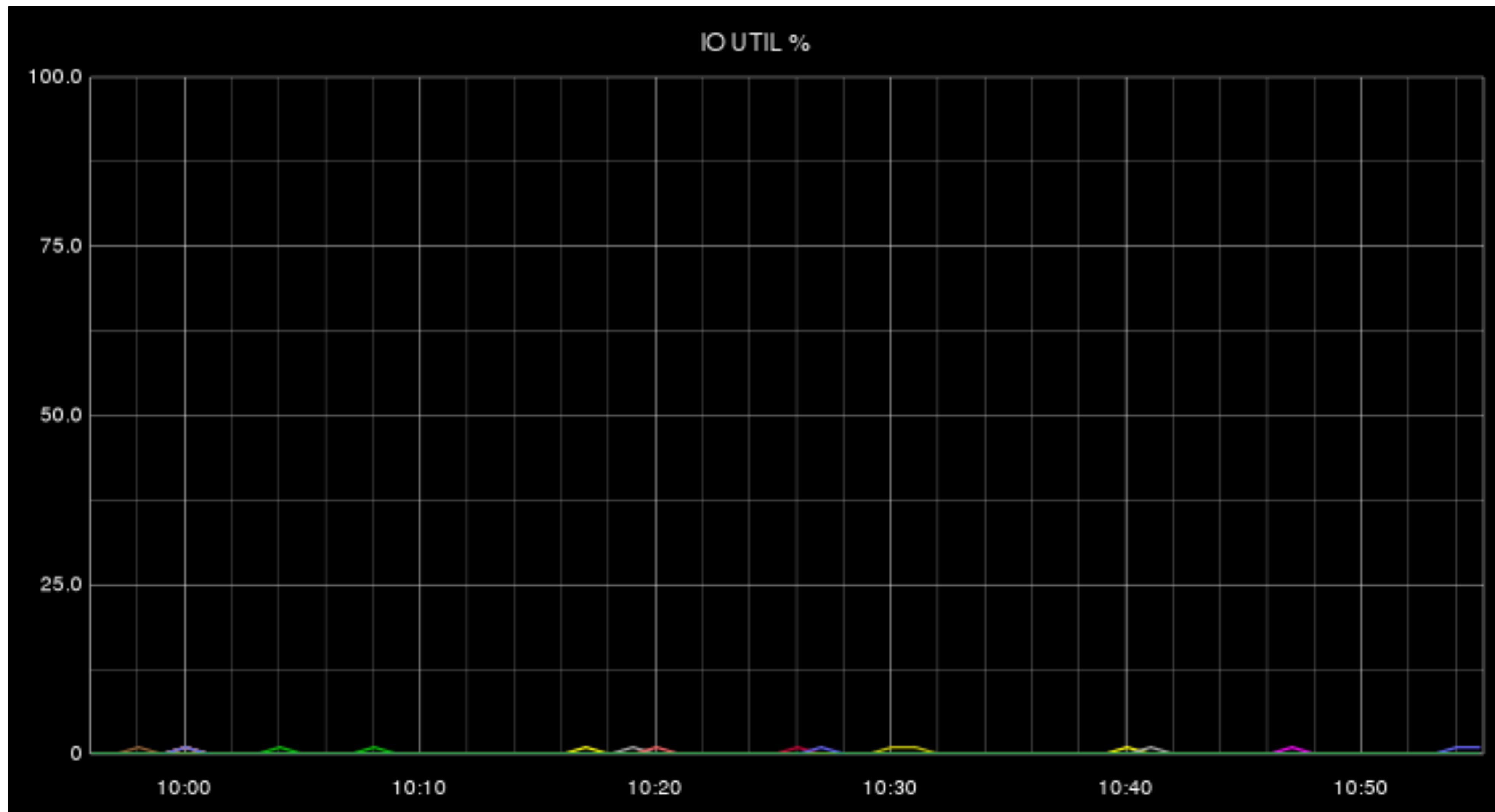
CPU



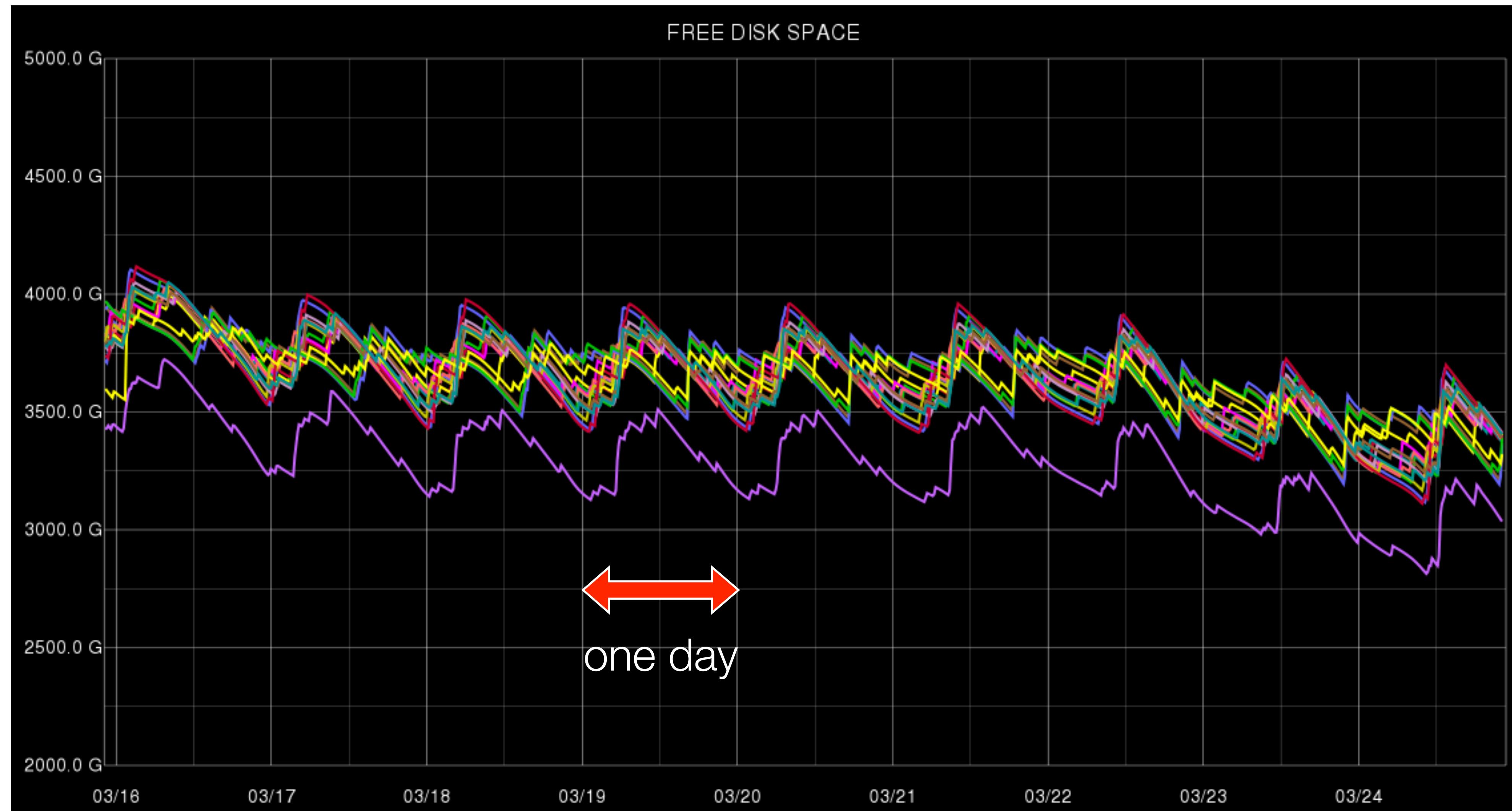
DISK IOPS



DISK IO %



EXPIRATION & ESPACE DISQUE

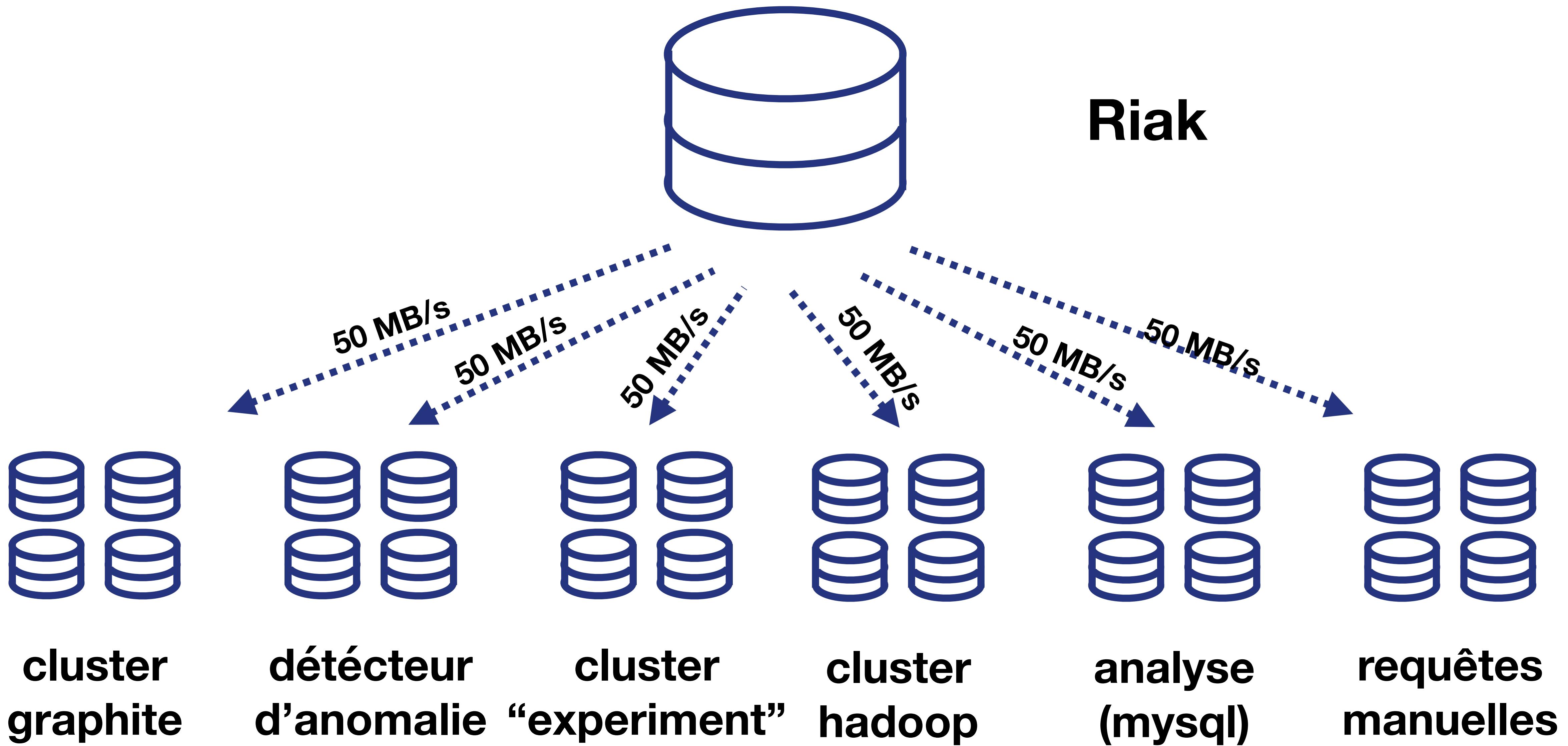


TRAITEMENT TEMPS REEL HORS RIAK

EXEMPLES

- Streaming => Cluster Graphite (chaque seconde)
- Streaming => DéTECTEUR d'Anomalie (agrégation 2 min.)
- Streaming => Cluster “Experiment” (agrégation 1 jour)
- Chaque minute => Cluster Hadoop
- Requête manuelle => tests, debug, investigation
- Requête Batch => analyse ad hoc
- => quantité **énorme** de requête en **lecture**

Riak



VRAI TEMPS REEL

- 1 seconde de données
- Stockée en moins d'une seconde
- Disponible en moins d'une seconde
- Problème long terme : saturation réseau

TRAITEMENT TEMPS REEL DANS RIAK

L'IDEE

- Au lieu de
 - Lire la donnée
 - Traiter la donnée (ex: moyenne)
 - Produire un petit résultat
- Faire
 - Envoyer le code près des données
 - Traiter les données sur Riak
 - Récupérer le résultat

CE QUI PREND DU TEMPS

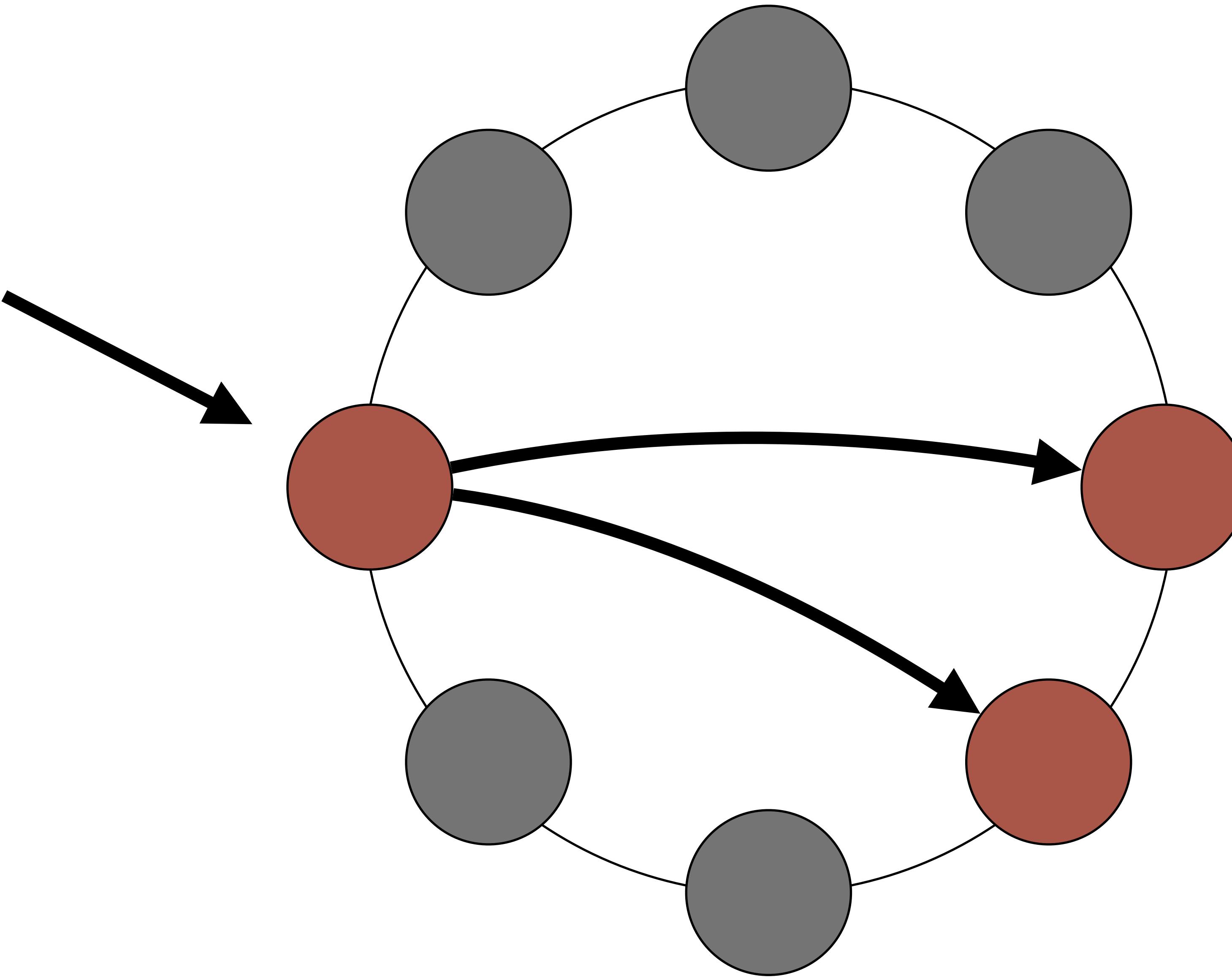
- Prend du temps et de la bande passante
 - Récupérer les données: **BP Réseau**
 - Décompresser/déserialiser: **Temps CPU**
- Ne prend pas de temps
 - Traiter les données

MAPREDUCE

- Entrée: epoch-dc
- Map1: clef meta donnée -> clefs données
- Map2: calcul sur la donnée
- Reduce: regrouper
- Temps réel: OK
- OK pour la **BP réseau**, pas pour le **Temps CPU**

HOOKS

- A chaque écriture des méta-données
- Post-Commit hook déclenché
- Traite les données sur les noeuds



NOEUD

Riak post-commit hook

service RIAK

clef

socket

service REST

clef

decompression
et traitement des
n taches

résultat renvoyé



HOOK CODE

```
metadata_stored_hook(RiakObject) ->
    Key = riak_object:key(RiakObject),
    Bucket = riak_object:bucket(RiakObject),
    [ Epoch, DC ] = binary:split(Key, <<"-">>),
    MetaData = riak_object:get_value(RiakObject),
    DataKeys = binary:split(MetaData, <<"|">>, [ global ]),
    send_to_REST(Epoch, Hostname, DataKeys),
    ok.
```

```
send_to_REST(Epoch, Hostname, DataKeys) ->
    Method = post,
    URL = "http://" ++ binary_to_list(Hostname)
        ++ ":5000?epoch=" ++ binary_to_list(Epoch),
    HTTPOptions = [ { timeout, 4000 } ],
    Options = [ { body_format, string },
                { sync, false },
                { receiver, fun(ReplyInfo) -> ok end }
            ],
    Body = iolist_to_binary(mochijson2:encode(DataKeys)),
    httpc:request(Method, {URL, [], "application/json", Body},
                  HTTPOptions, Options),
    ok.
```

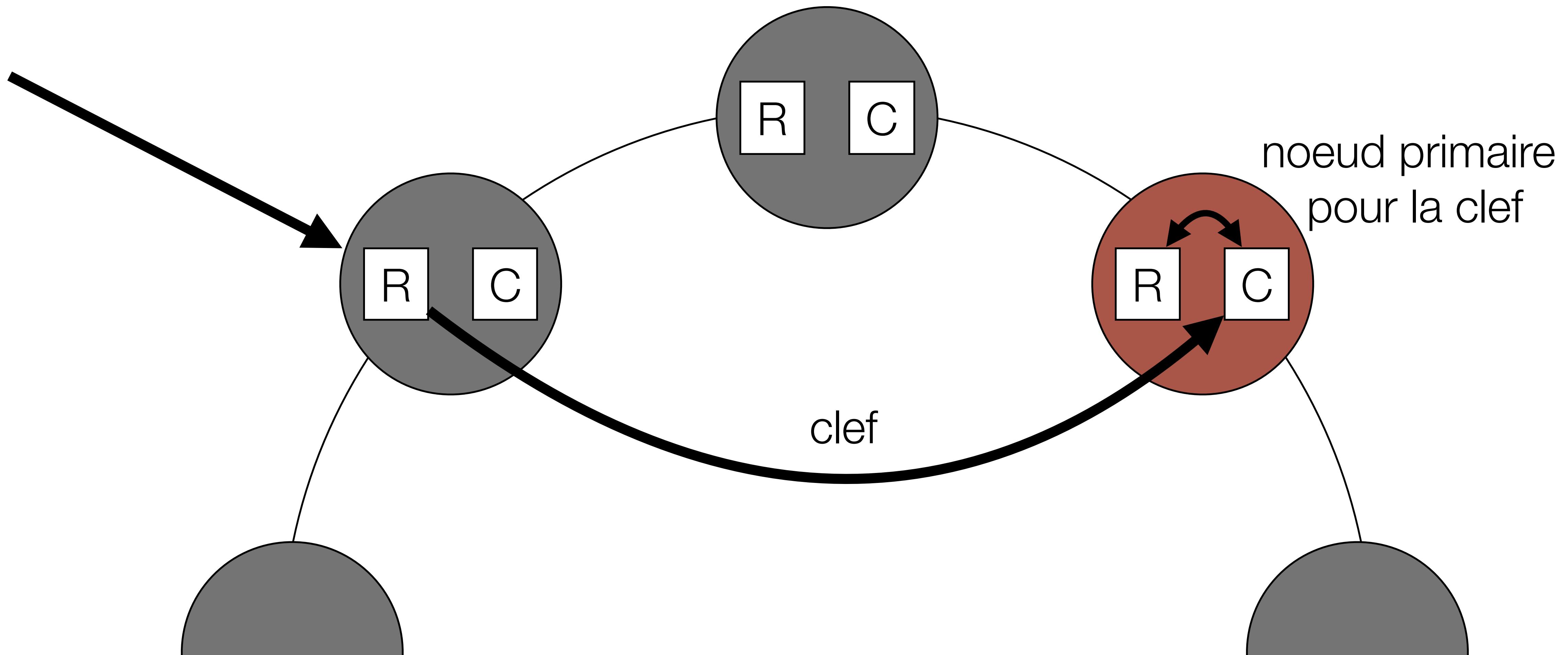
SERVICE REST

- En Perl, norme PSGI (WSGI), Starman, preforks
- Permet d'écrire les traitements en Perl
- Supporte chargement de code à la demande

SERVICE REST: SCALABLE

- **Scalable**
- 1 seconde = 200 clefs
- 16 noeuds, 10 CPUs (2 pour Riak)
- 1 clef doit être traitée en $16*10/200$ sec = 0.8 sec
- Inclus réseau + CPU + overhead
- => Amplement le temps

OPTIMISATION



AVANTAGES

- Utilisation CPU et tps d'execution peut être limité
- Donnée est **locale** au traitement
- Donnée décompressée **une seule fois**
- Tous les traitements sont fait sur une donnée
- Les deux systèmes sont découplés (Erlang surveille)
- Peut être écrit en n'importe quel langage

DESAVANTAGES

- Seulement pour données temps réel
- Traitement inter-epoch moins facile

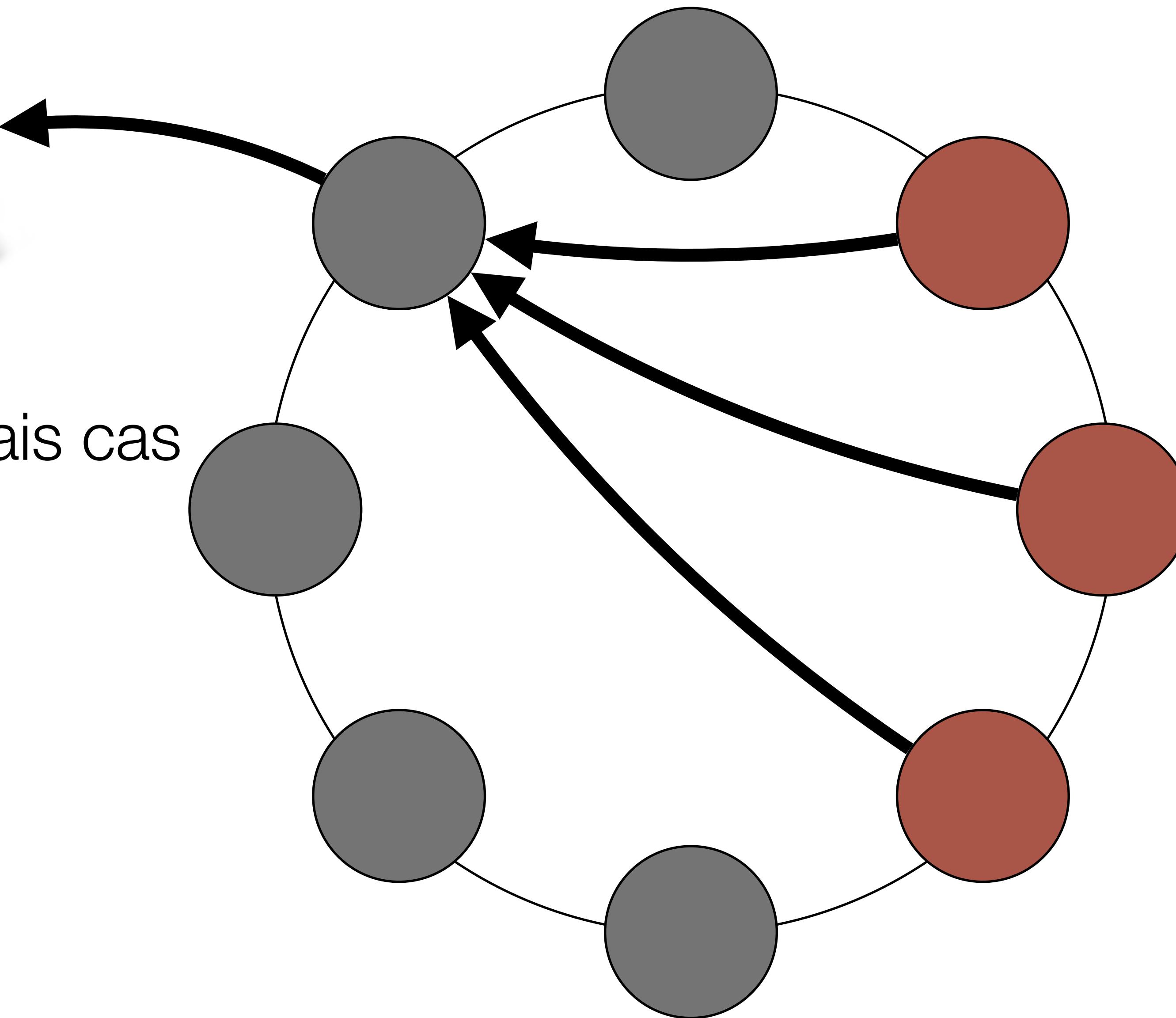
FUTUR

- Utiliser ce système pour réduire l'utilisation réseau
- Explorer la piste de Riak Search

BANDE PASSANTE RESEAU: PROBLEME

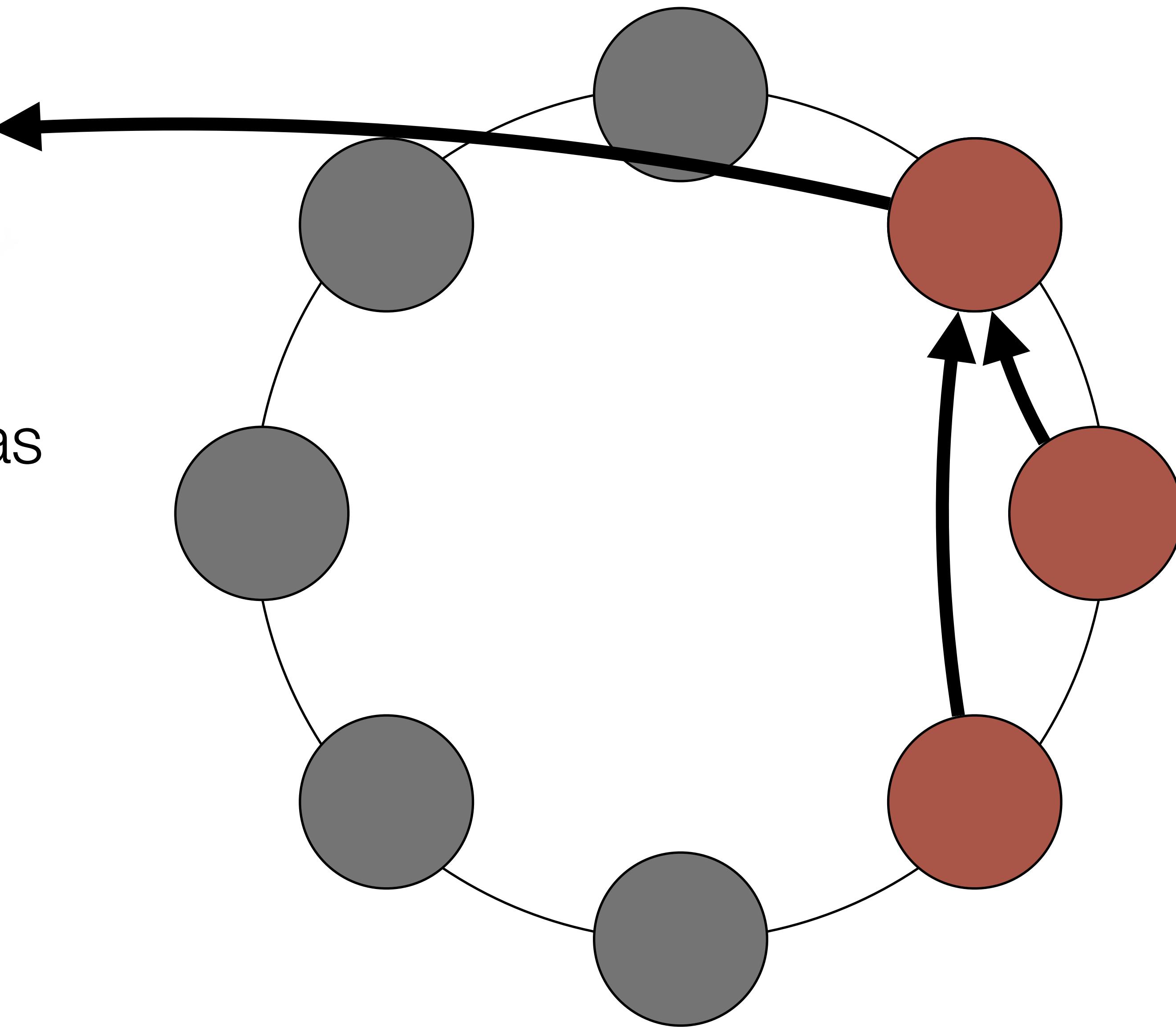


- **GET** - mauvais cas
- intérieur =
3 x extérieur





- **GET** - bon cas
- intérieur =
2 x extérieur



BANDE PASSANTE

- Ratio tend vers 3
- Cas général : OK
- Dans notre cas: problème
- Grosses valeurs, PUT constants, beaucoup de GET
- Limité à 1 Gbits/s

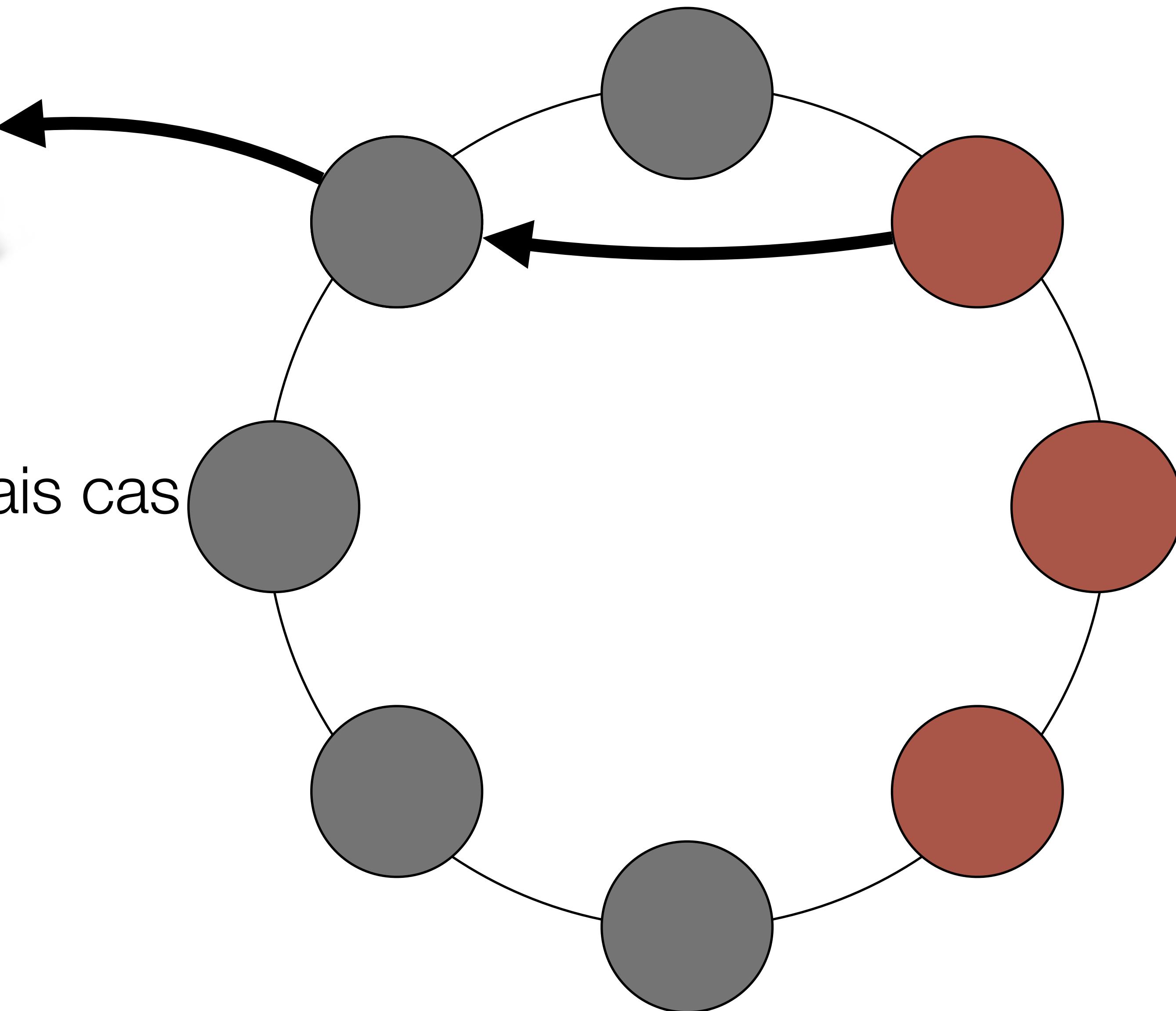
BANDE PASSANTE RESEAU: SOLUTIONS

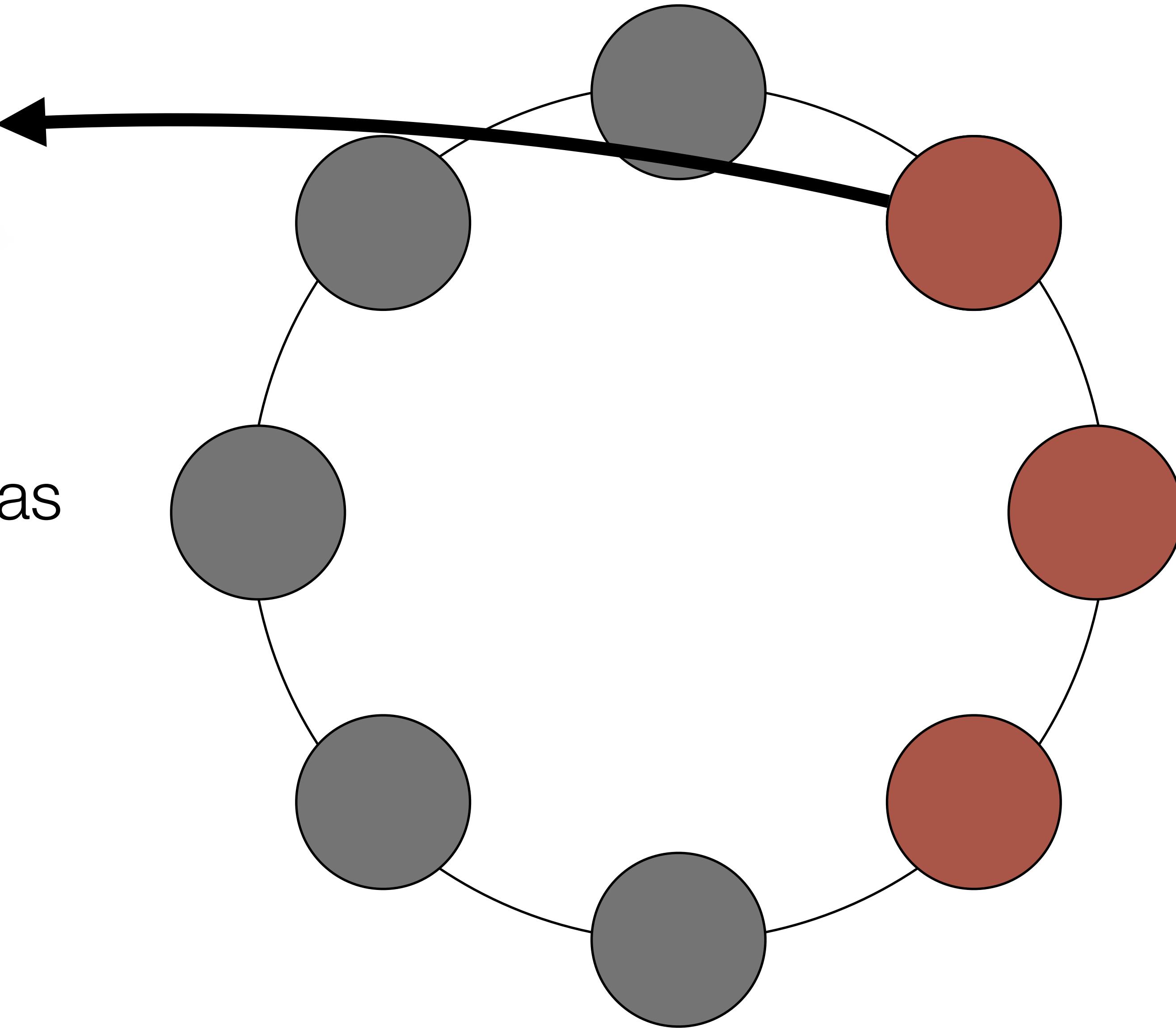
BANDE PASSANTE RESEAU: SOLUTIONS

1. GET à faible consommation réseau
2. Ne pas choisir un noeud au hasard



- **GET** - mauvais cas
- **n_val = 1**
- extérieur =
1 x intérieur





- **GET** - bon cas
 - **n_val = 1**
 - **exterieur =**
0 x intérieur

RESULTAT

- Utilisation réseau réduite par 2 !

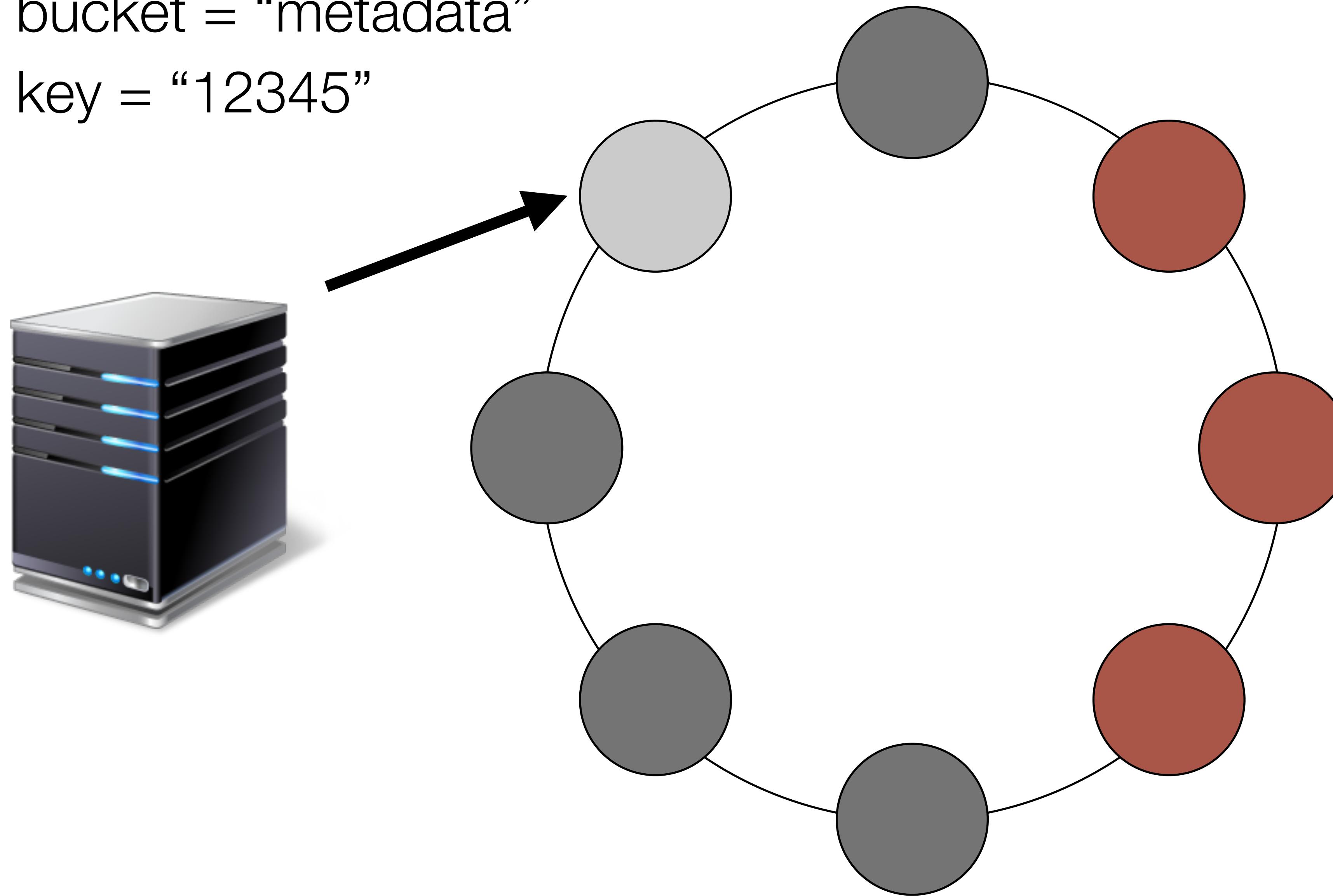
ATTENTION

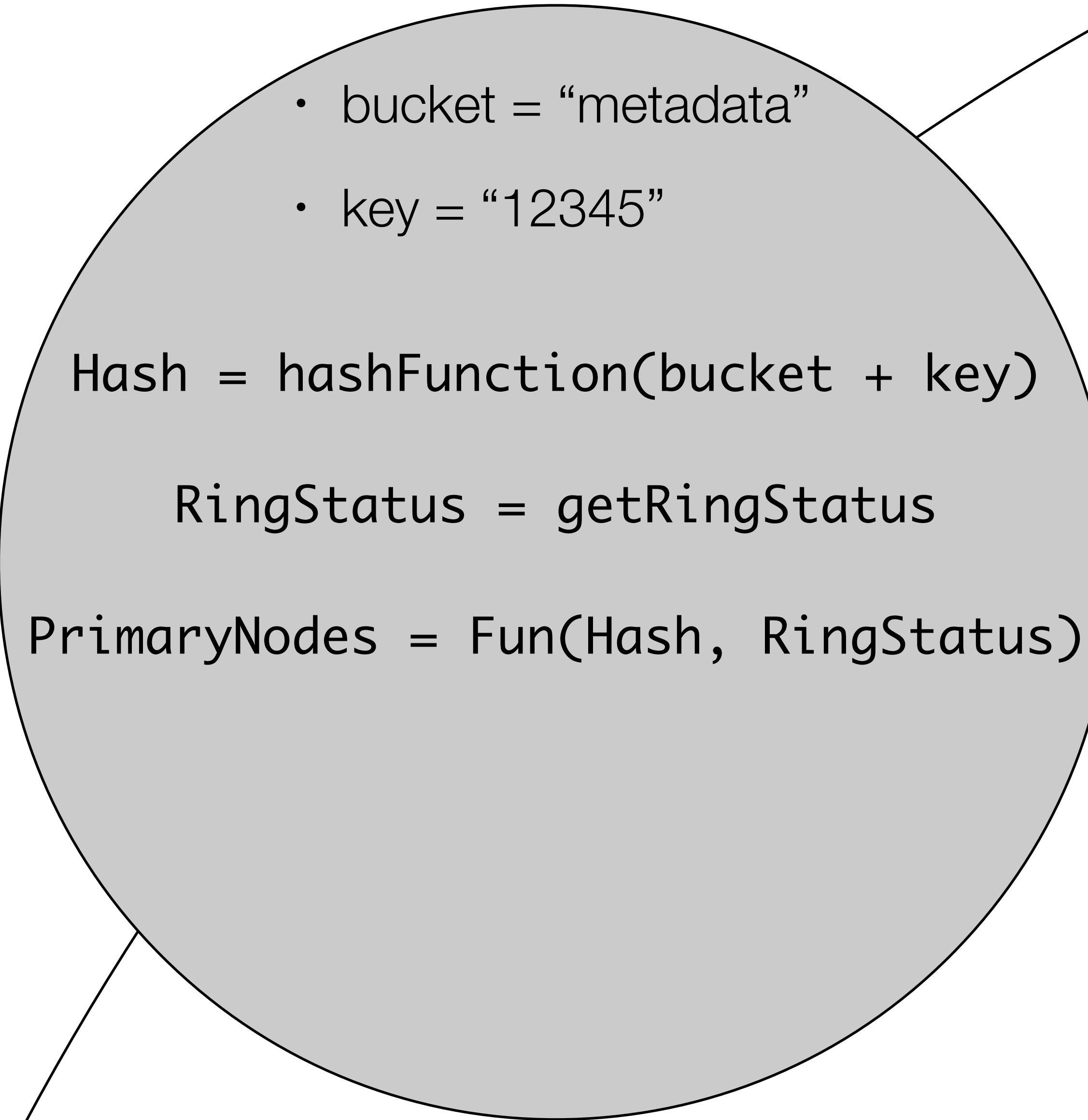
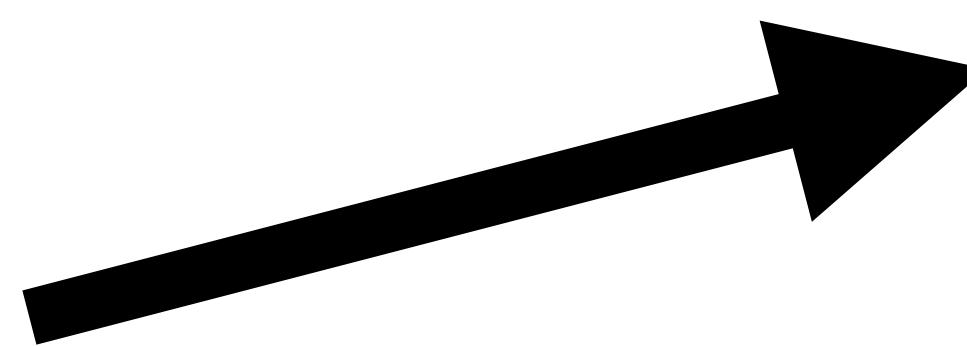
- Possible car les données sont en lecture seule
- Données ont un checksum interne
- Aucun conflit n'est possible
- Corruption de donnée auto-déTECTée
- Prévoir un fallback avec $n_val=3$

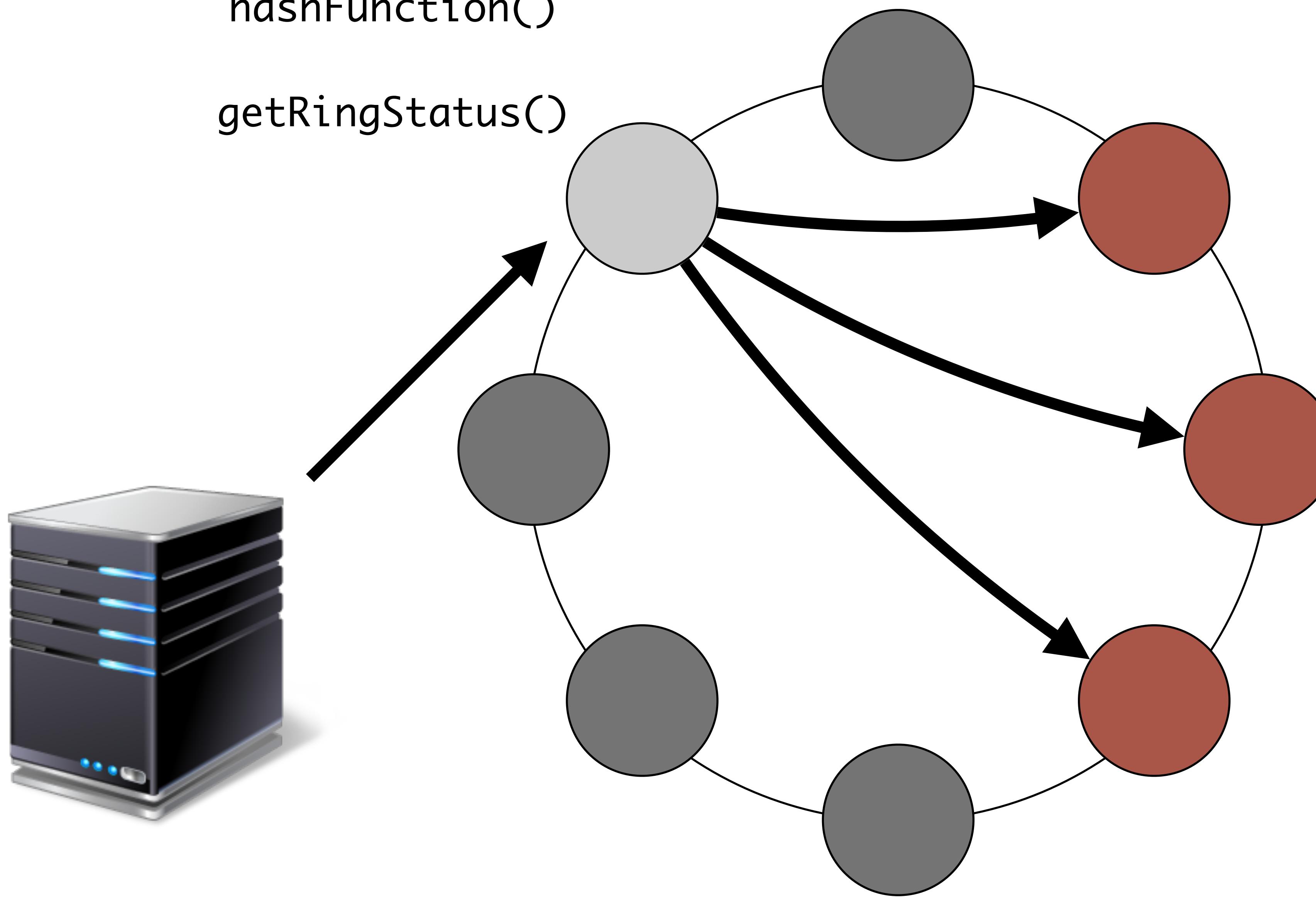
BANDE PASSANTE RESEAU: SOLUTIONS

1. GET à faible consommation réseau
2. Ne pas choisir un noeud au hasard

- bucket = “metadata”
- key = “12345”

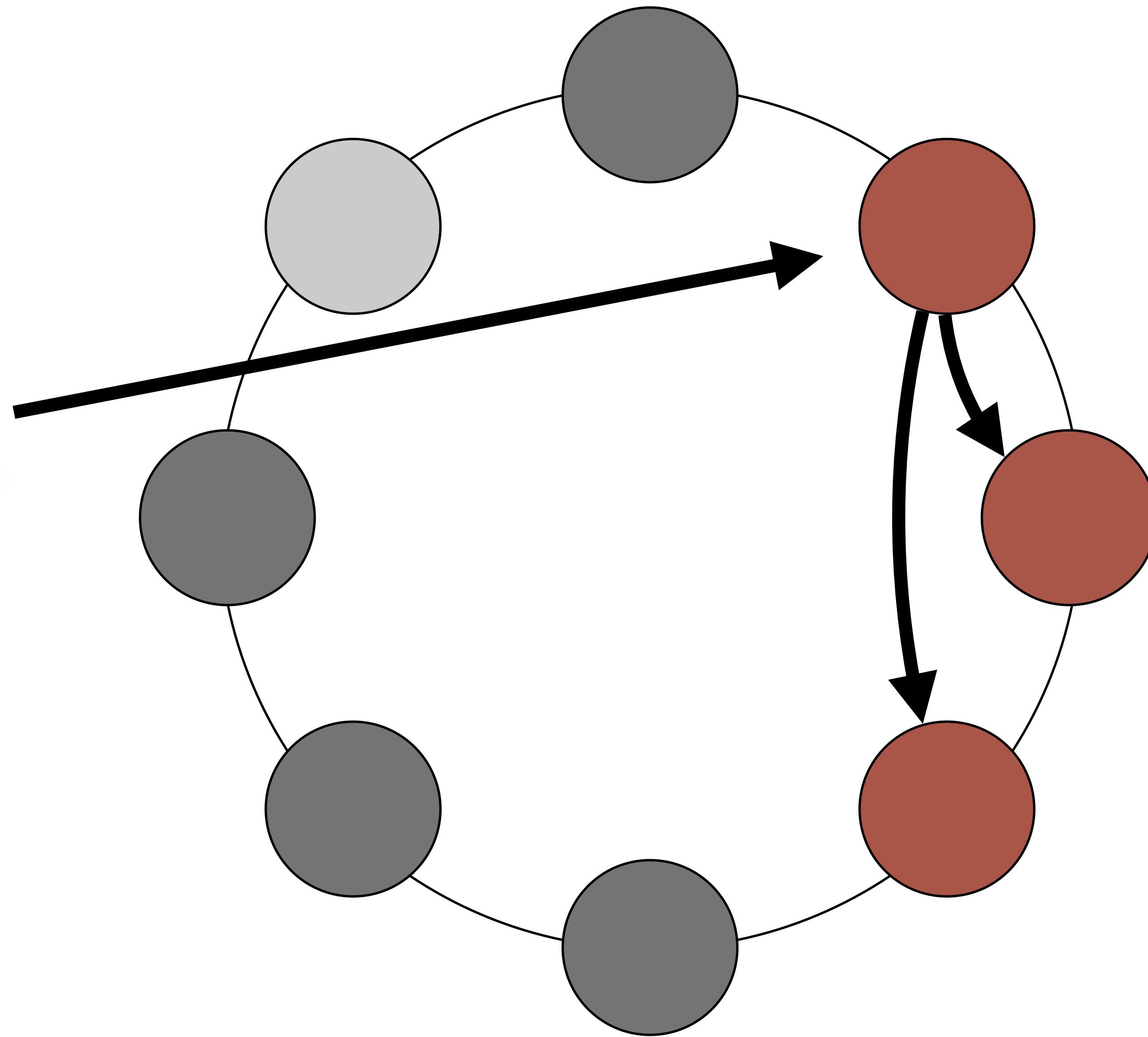








hashFunction()
getRingStatus()



LE PRINCIPE

- Faire le hashing **côté client**
- Fonction de hash par défaut:
 - “`chash_keyfun`”: {“`mod`”: “`riak_core_util`”,
“`fun`”: “`chash_std_keyfun`”},
- Allons voir le code sur github

LA FONCTION DE HASH

```
%% @spec chash_std_keyfun(BKey :: riak_object:bkey()) -> chash:index()  
%% @doc Default object/ring hashing fun, direct passthrough of bkey.  
chash_std_keyfun({Bucket, Key}) -> chash:key_of({Bucket, Key}).
```

```
%% @doc Given any term used to name an object, produce that object's key  
%%       into the ring. Two names with the same SHA-1 hash value are  
%%       considered the same name.  
-spec key_of(ObjectName :: term()) -> index().  
key_of(ObjectName) ->  
    sha(term_to_binary(ObjectName)).
```

- `riak_core/src/chash.erl`

LA FONCTION DE HASH

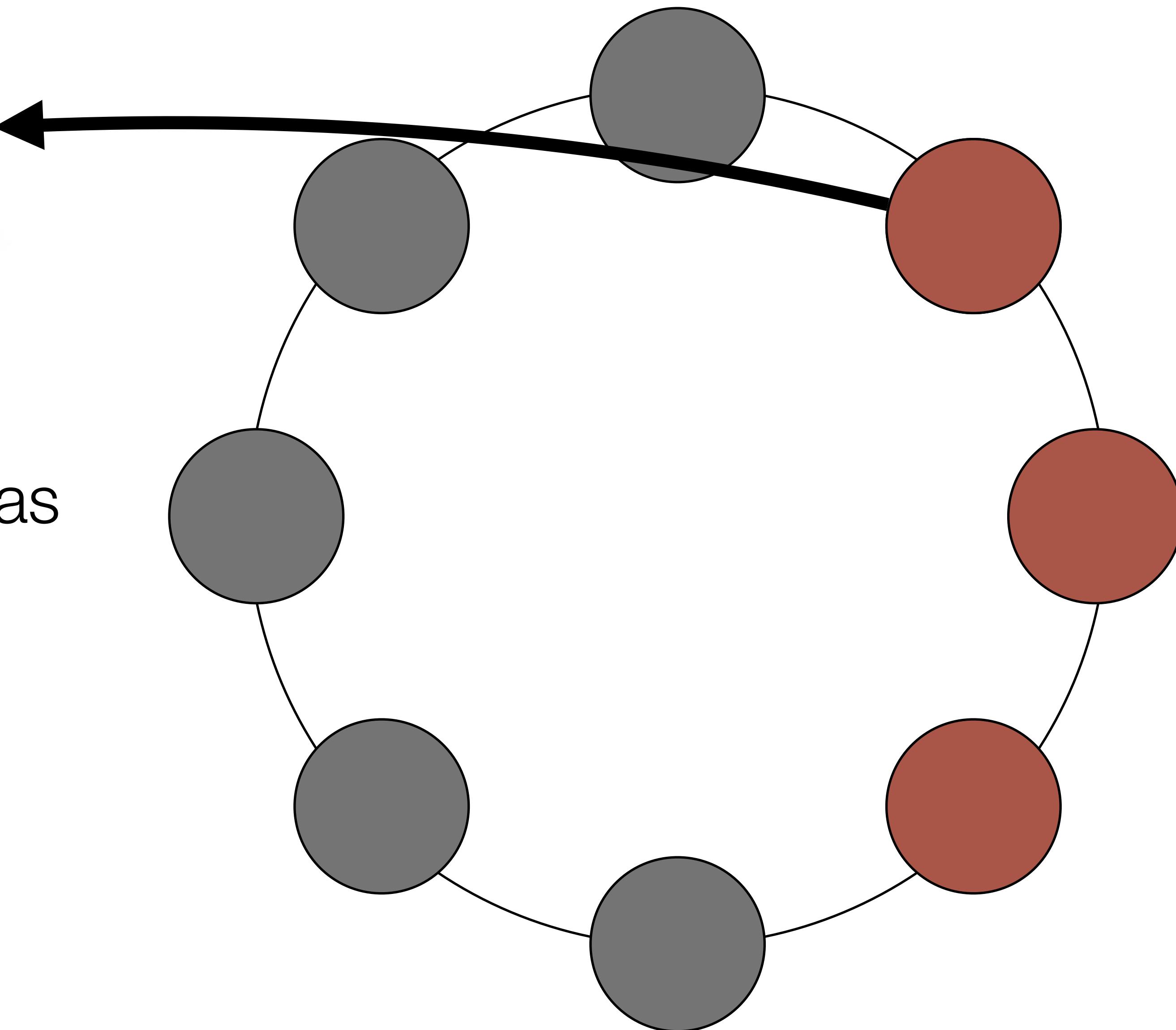
- Facile à réimplémenter côté client:
- sha1 représenté en BigInteger
- 4008990770823839524233143877797980545530986496
- entre 0 et 2^{160}

ETAT DU RING

```
$ curl -s -XGET http://$host:8098/riak_preflists/myring
{
  "0" : "riak@host-16"
  "570899077082383952423314387797980545530986496" : "riak@host-17"
  "11417981541647679048466287755595961091061972992" : "riak@host-18"
  "17126972312471518572699431633393941636592959488" : "riak@host-19"
  ...
  ...
}
```



- **GET** - bon cas
- **n_val = 1**
- exterieur =
0 x intérieur



RESULTAT

- Utilisation réseau encore diminuée
- Particulièrement pour les GETs

ATTENTION

- Possible seulement si
 - Noeuds sont monitorés (ici Zookeeper)
 - En cas d'échec, retour à noeud au hasard
 - Données lues de manière uniforme

CONCLUSION

CONCLUSION

- Version Open Source de Riak
- Pas de training, autodidacte, équipe très petite
- Riak est une bonne solution
 - Robuste, rapide, scalable, facile
 - Très flexible and hackable
- Nous permet de continuer à scaler chez Booking.com

Q&A



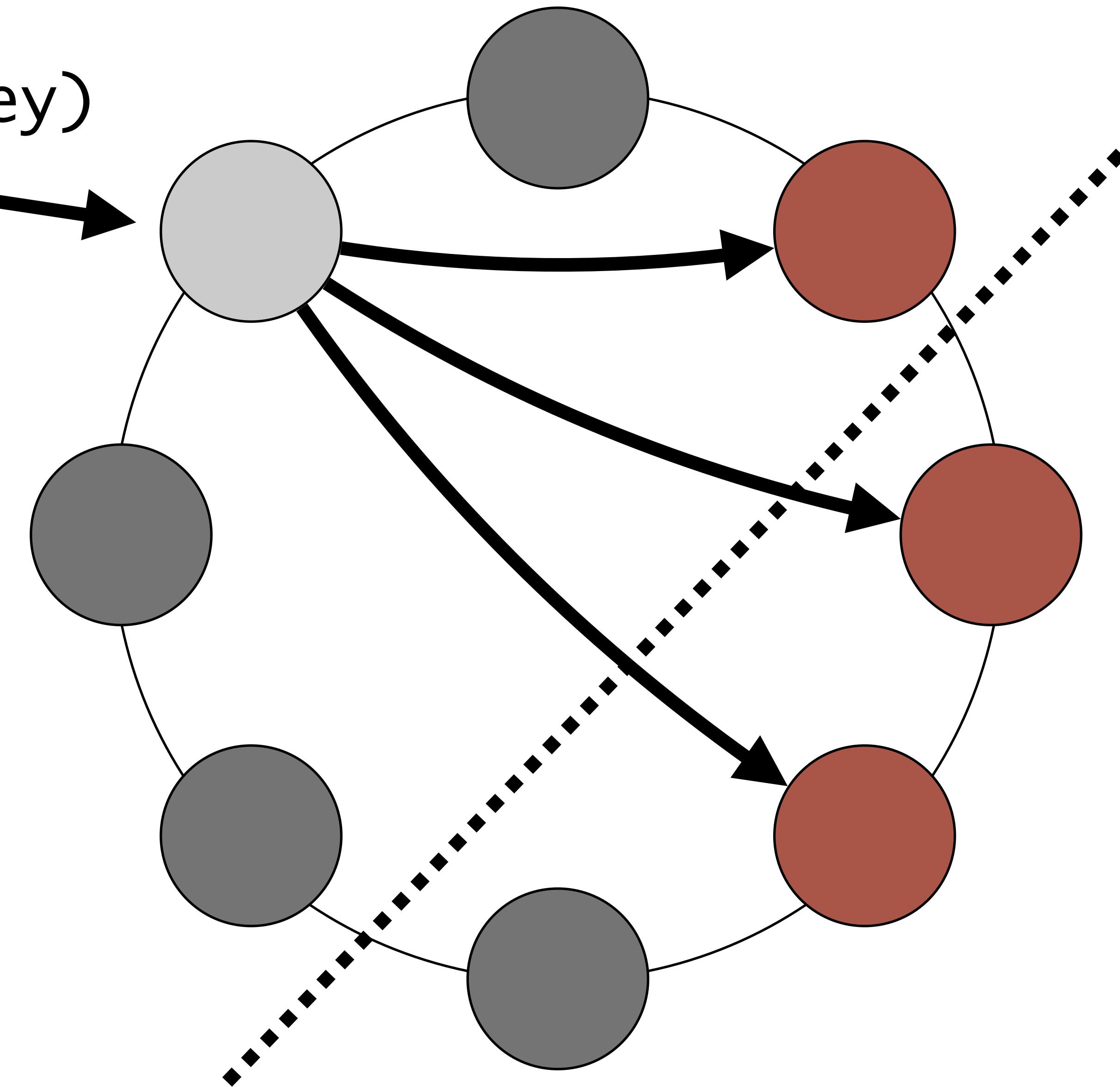
@damsieboy

EXTRA: RIAK ET LE “CAP THEOREM”



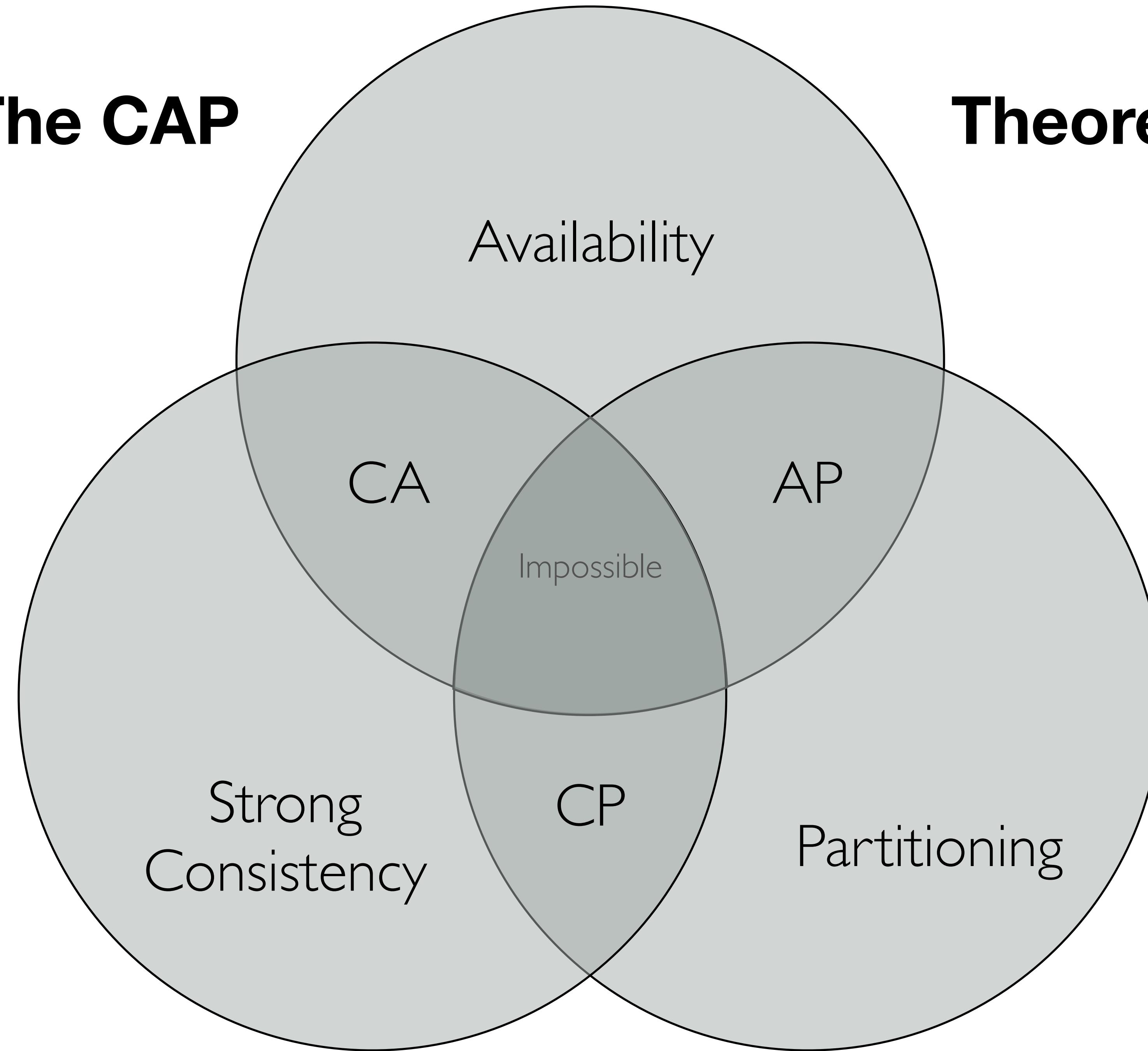
partition

hash(key)



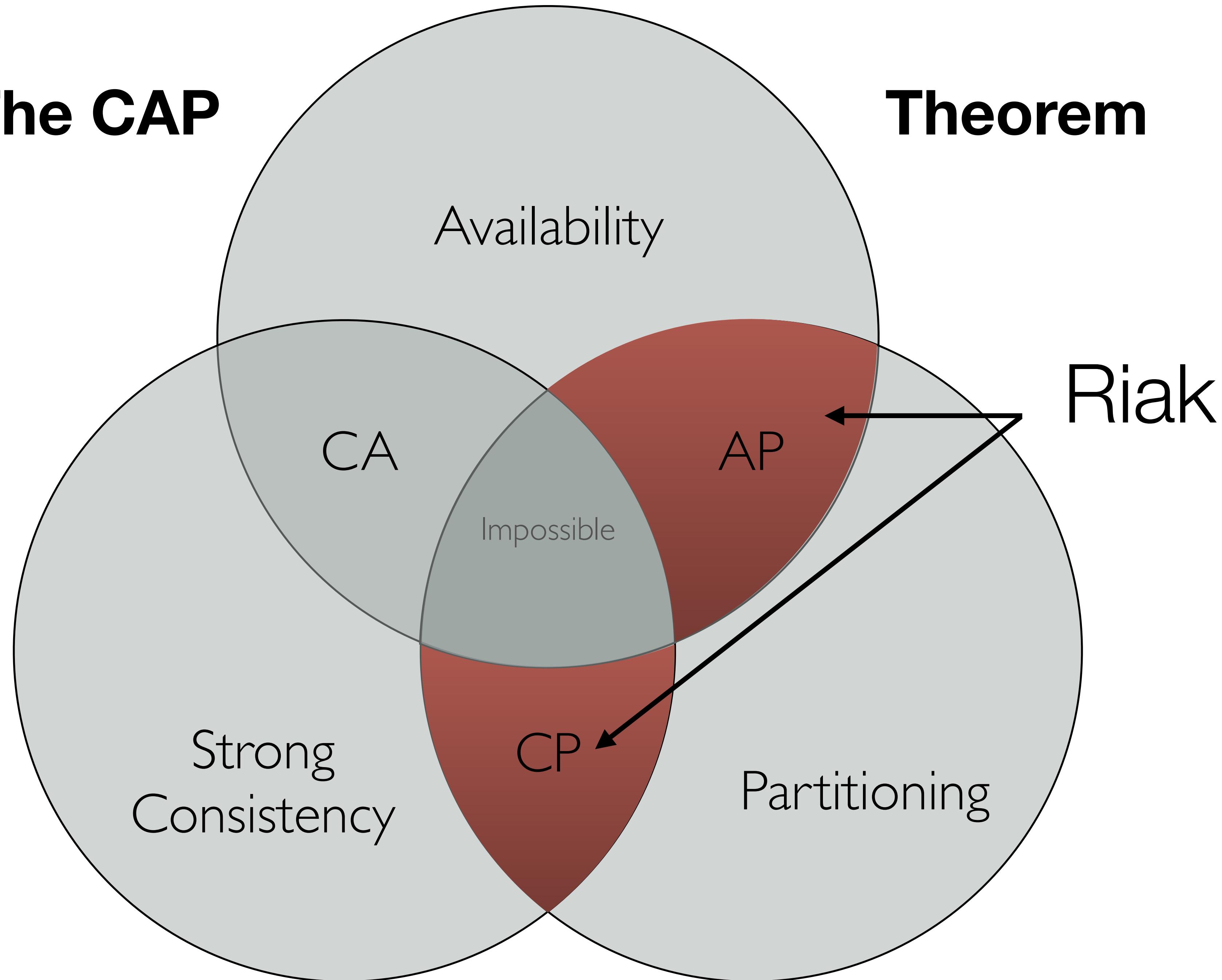
The CAP

Theorem



The CAP

Theorem



The CAP Theorem

