



Introduction à Cassandra

Nicolas Romanetti
@nromanetti

www.jaxio.com

mercredi 18 avril 2012



University Talk

- Tout public
- 3h pour vraiment comprendre Cassandra
- Pédagogique



University Talk

- Démarrre doucement
- Continue doucement
- Mais en 3h nous allons couvrir beaucoup de sujets
- Devrait vous chatouiller un peu le cerveau
- Conseil: prenez des notes





Cassandra LAN Party

(Hands On Cassandra)

Cet après midi
de 13h30 à 16h30



A Propos

- Nicolas Romanetti
- @nromanetti
- www.jaxio.com
- Légitimité?



Cassandre

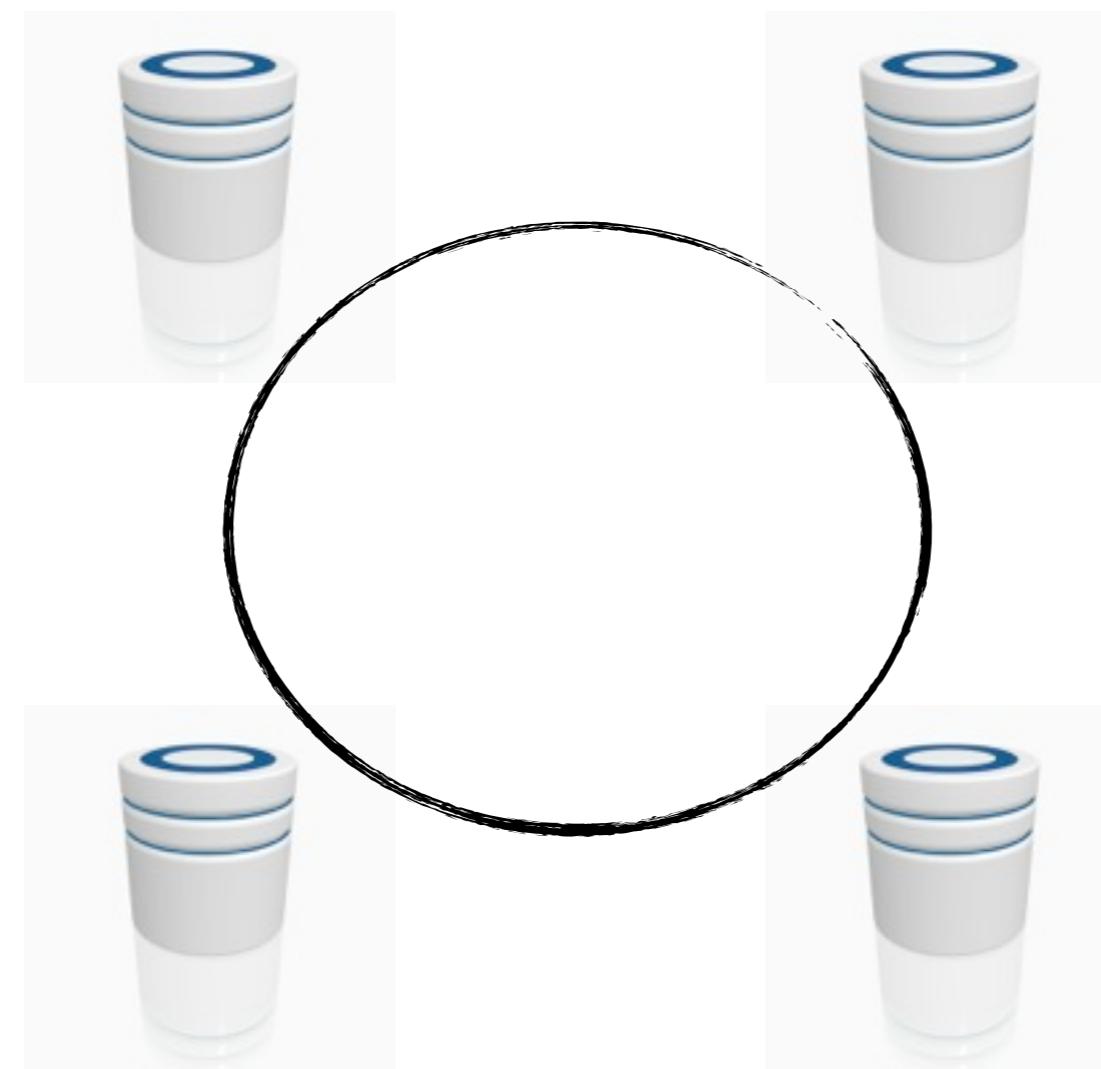
- Très belle
- Libre
- Avait raison



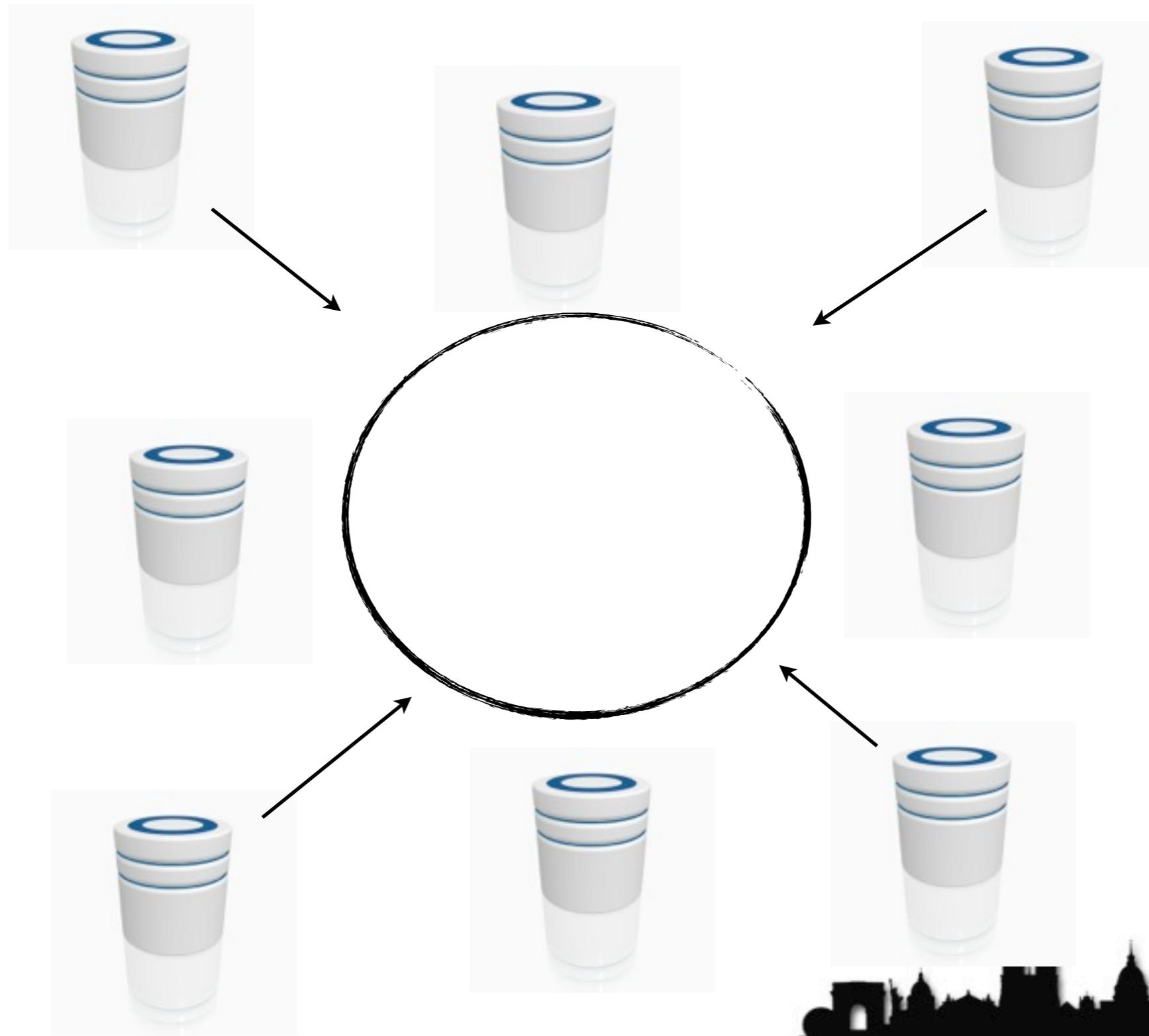
Cassandra c'est...



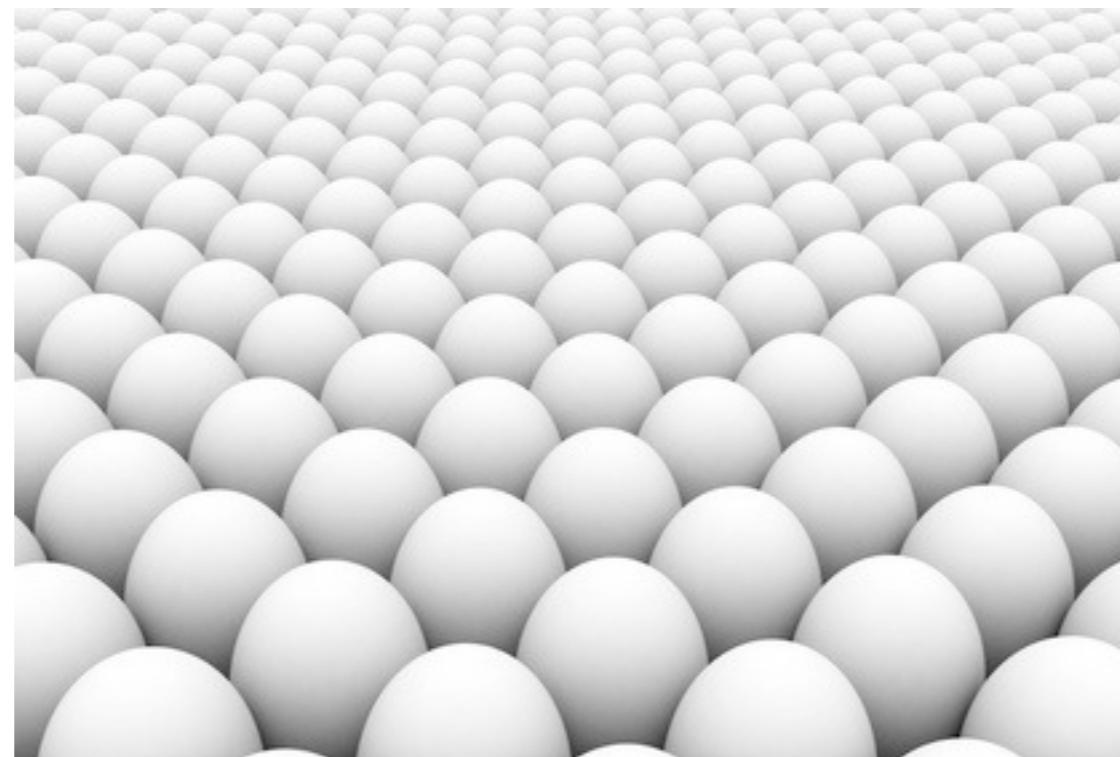
Cassandra c'est...



Cassandra c'est...



Symétrique



Conçue pour les pannes



RDBM

réPLICATION

INDEXES

CACHE

JOIN

SCALABILITÉ

DÉNORMALISATION

TRANSACTION



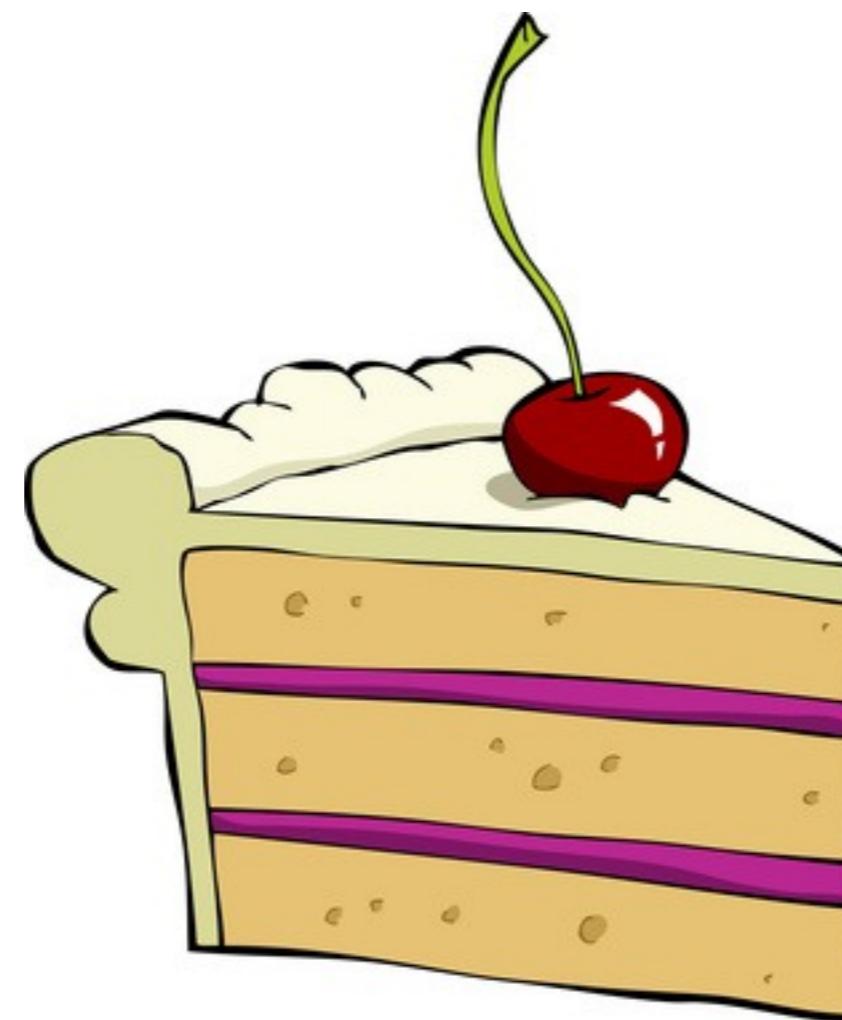
PERFORMANCE



Très important



Installation de Cassandra

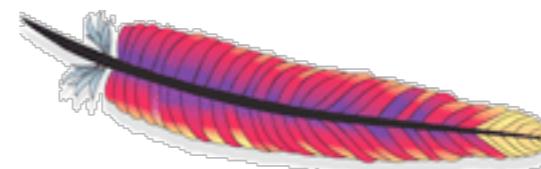


Historique

amazon
Dynamo

Google
Big Table

facebook



DATASTAX

NETFLIX

twitter



Data Model Cassandra

- Attention au vocabulaire



Version courte

- Une Map distribuée de Map



Column

name

value

timestamp

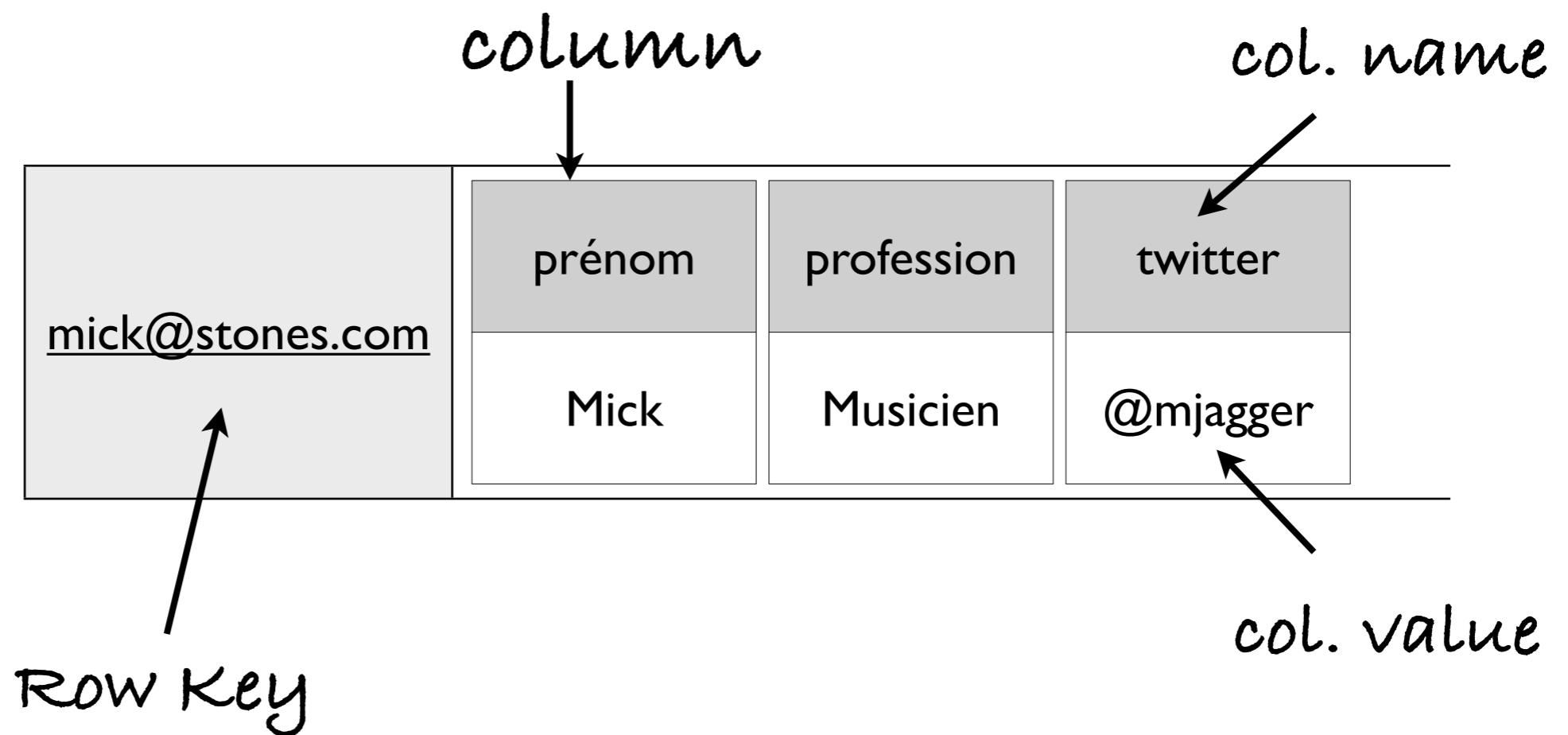


Row

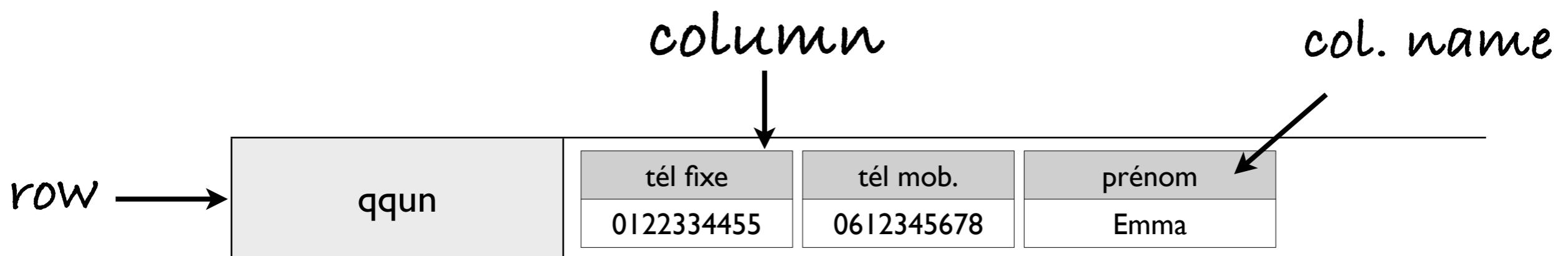
Key	col. name1	col. name2	col. name3	...	col. nameN
	col. value1	col. value2	col. value3		col. valueN



Exemple de Row



Libre...



qqunautre	fax 0102020202	naissance 01/12/2000	surnom Nico	tél fixe 0145060063
-----------	-------------------	-------------------------	----------------	------------------------

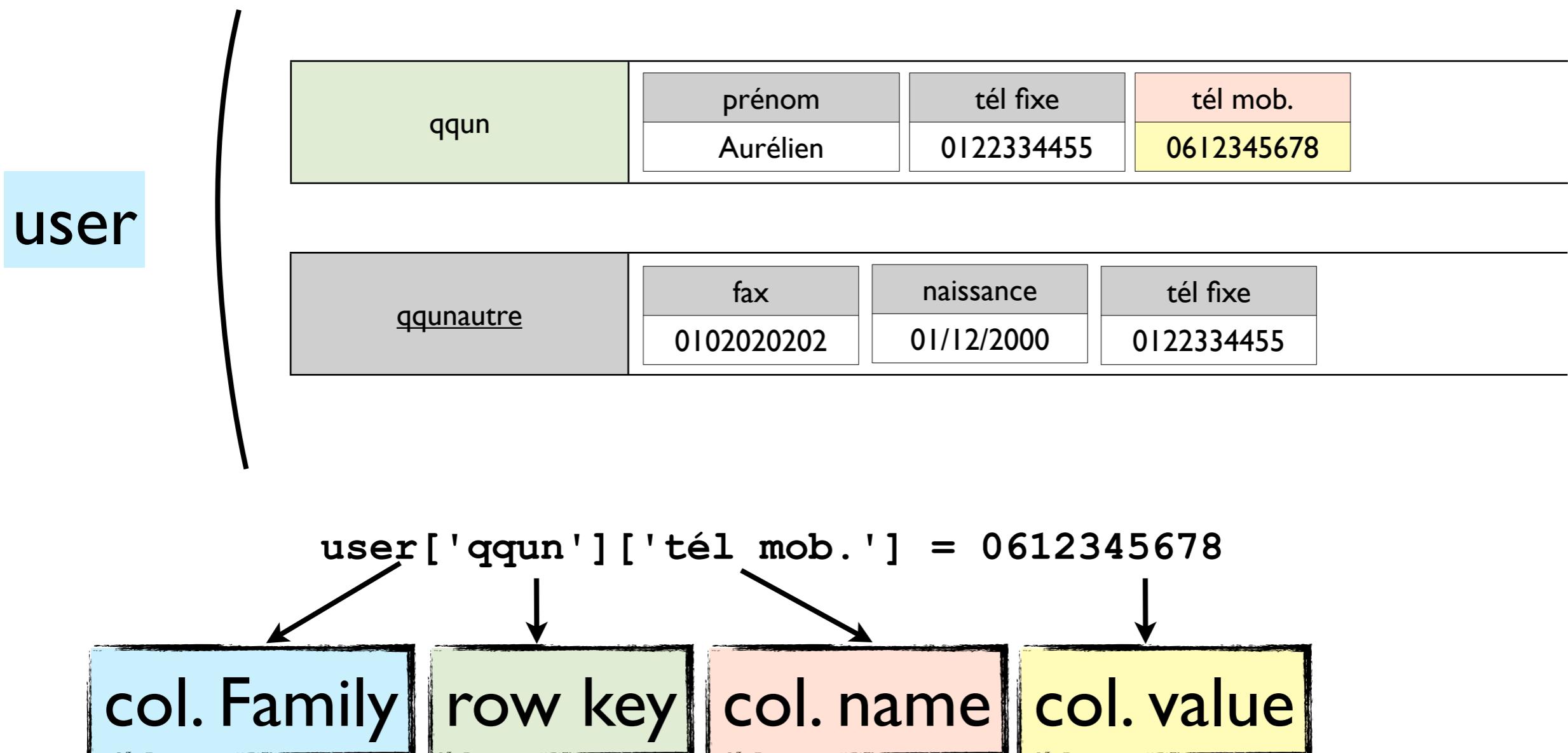
Key

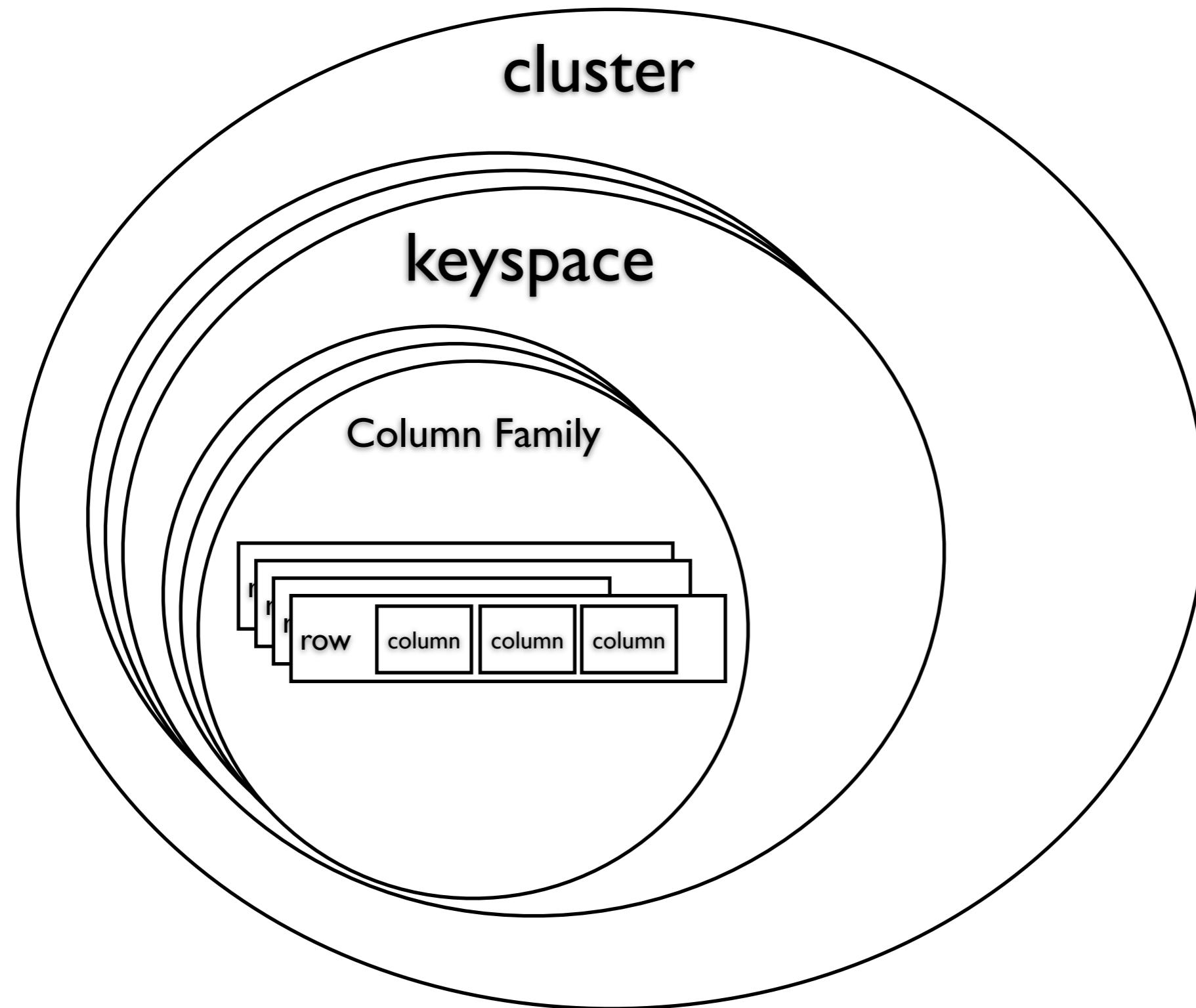
col. value

comparator/validator



Column Family





Le nom est une valeur!

‘vol’

XK_4501_2012 _01_06	6:55	7:00	7:20	etc...
	décollage face à la mer	survol Capo di Feno	survol Cannes	



Liens

‘BlogEntries’

jaxio	2012/04/18 10:52 rowkeyA	2012/04/12 10:12 rowkeyB	2012/04/01 10:11 rowkeyC	etc...
-------	------------------------------------	-----------------------------	-----------------------------	--------

‘BlogEntry’

rowkeyA	Content Plusieurs talks se ...	date 2012/04/18 10:52	Title DevoxxFR rox	etc...
----------------	-----------------------------------	--------------------------	-----------------------	--------



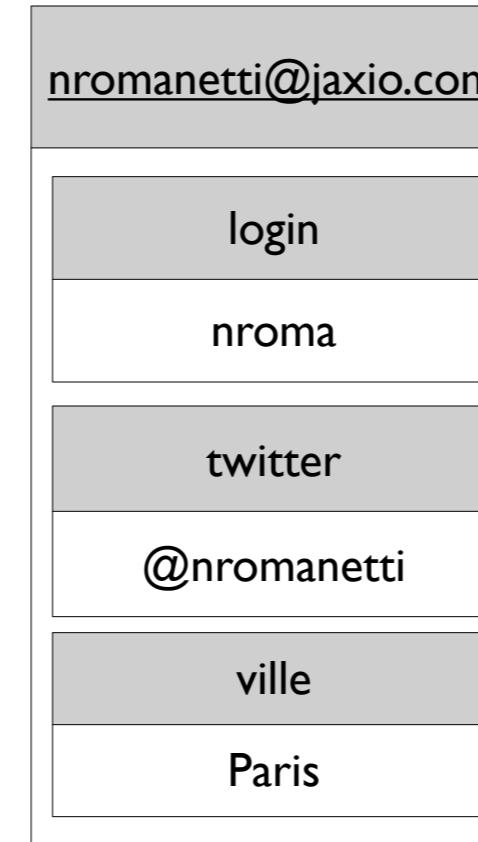
Composite name, valueless

‘TopScores’

2012/01/02	093:toto	105:titi	106:tutu	etc...
------------	----------	----------	----------	--------

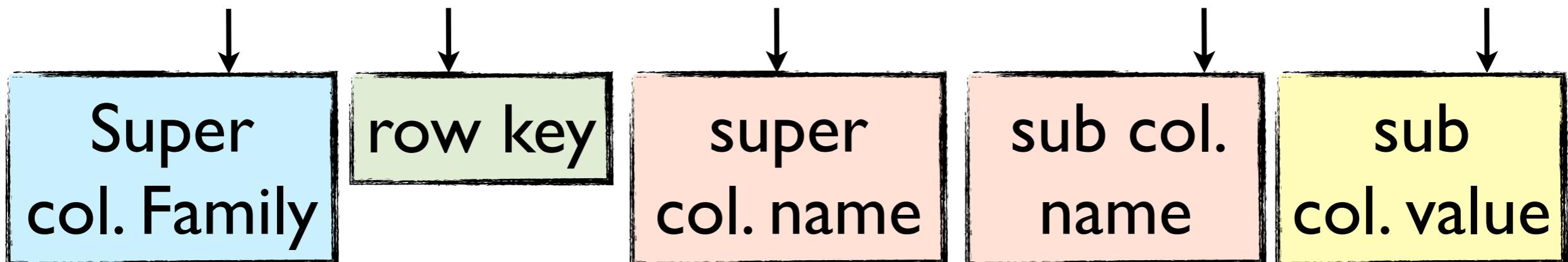


Super column



Deprecated
(va disparaître)

```
users['france']['nromanetti@jaxio.com']['ville'] = 'Paris'
```



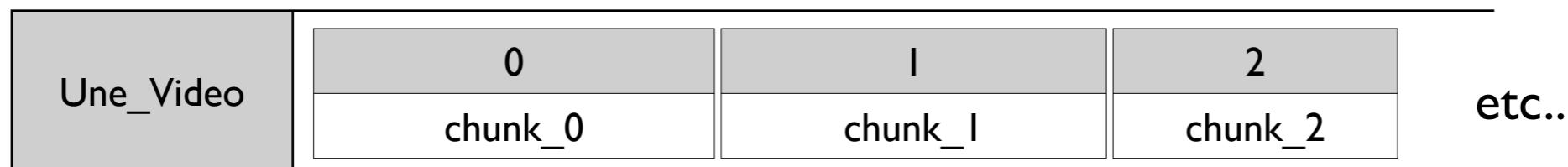
Fourre tout

contactInfo

tel:04...,
email: nr@...,
city:Ajaccio,
etc...



Gros contenu



Column expirante



Compteur Distribué

nb de like

57832



Tombstone



nom



Questions?



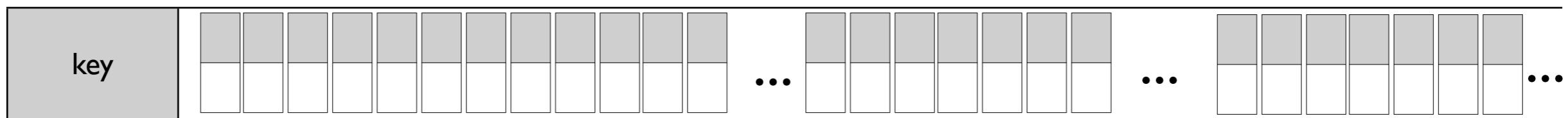
Query

- Attention au vocabulaire

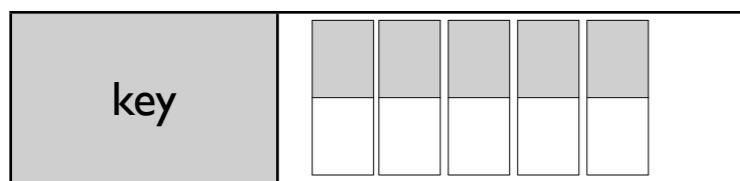


2 types de row

Wide row



Skinny row

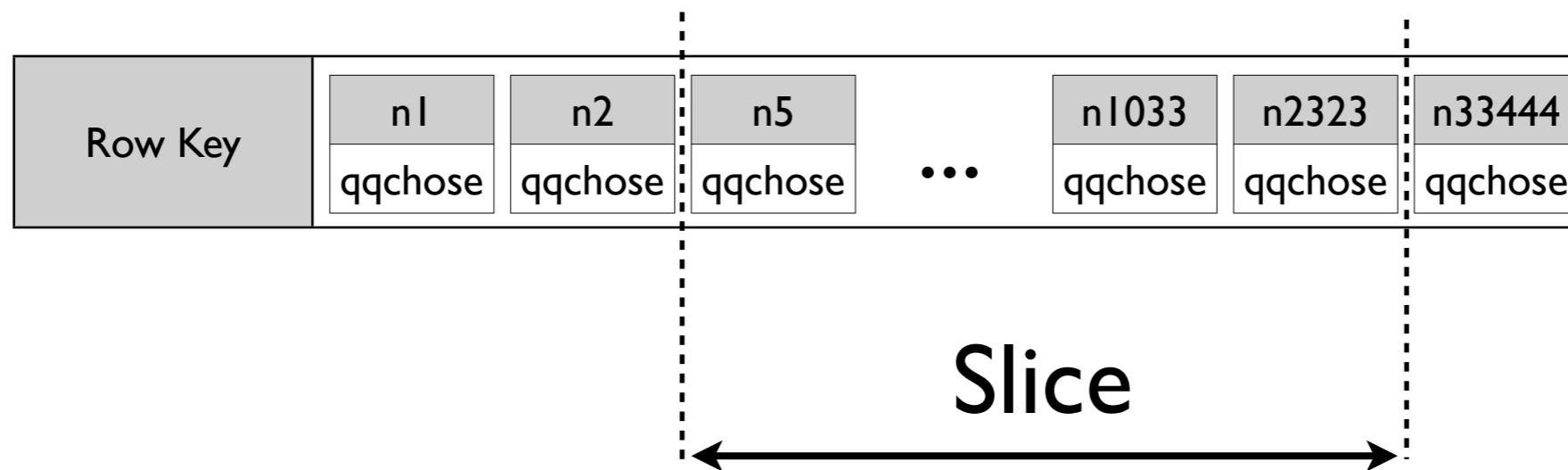


get

- `get user['toto']['firstname'];`
- `get user['toto']['state'];`



Column Slice



ordre..



Cassandra Query Language

- select user where country = 'FR';

Secondary Index



«Modélisation»



- Avoir un don de prophétie



Disponibilité ➔ RéPLICATION



Où stocker une row?

- Tous les serveurs sont identiques (symétrie)
- Postulats
 - 1 row doit tenir sur 1 seul serveur
 - On veut répartir les rows uniformément



Consistent Hashing

bytes
e.g. String

Consistent
Hash

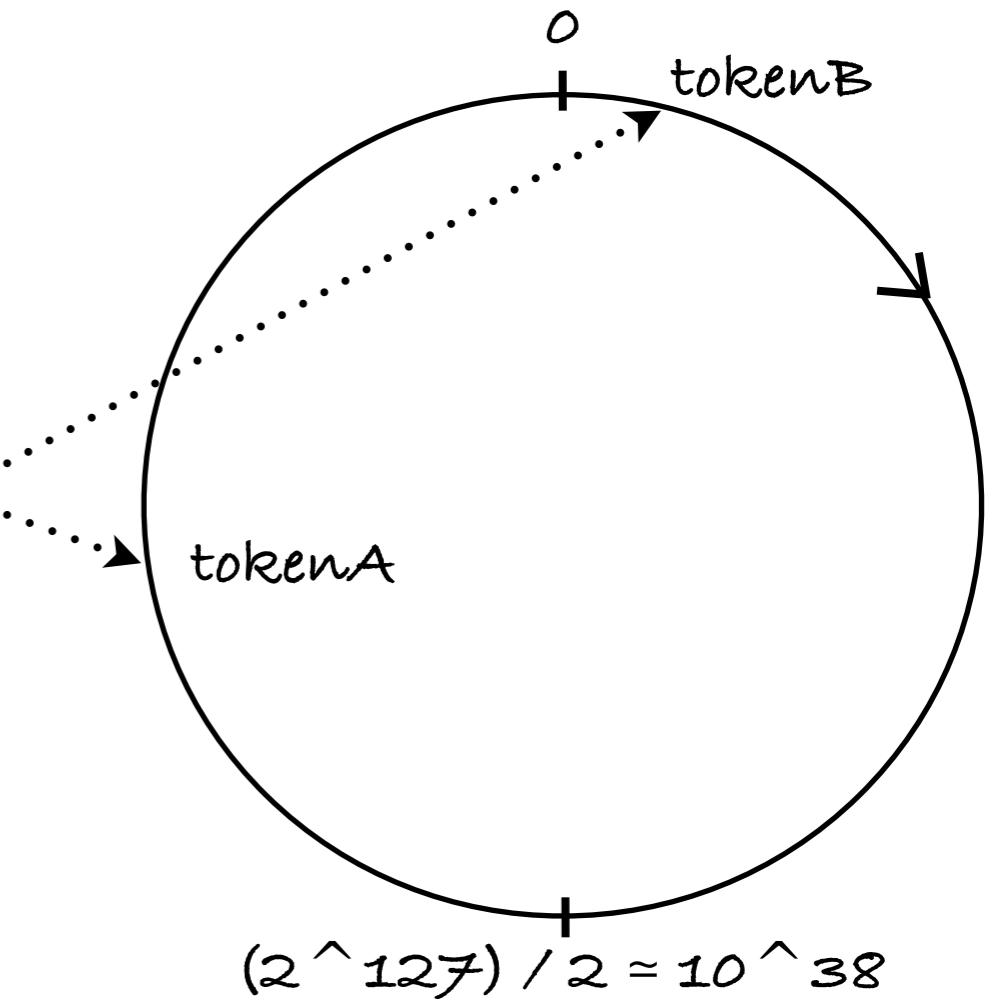
Un 'token' entre
0 et 2^{127}

row keyA

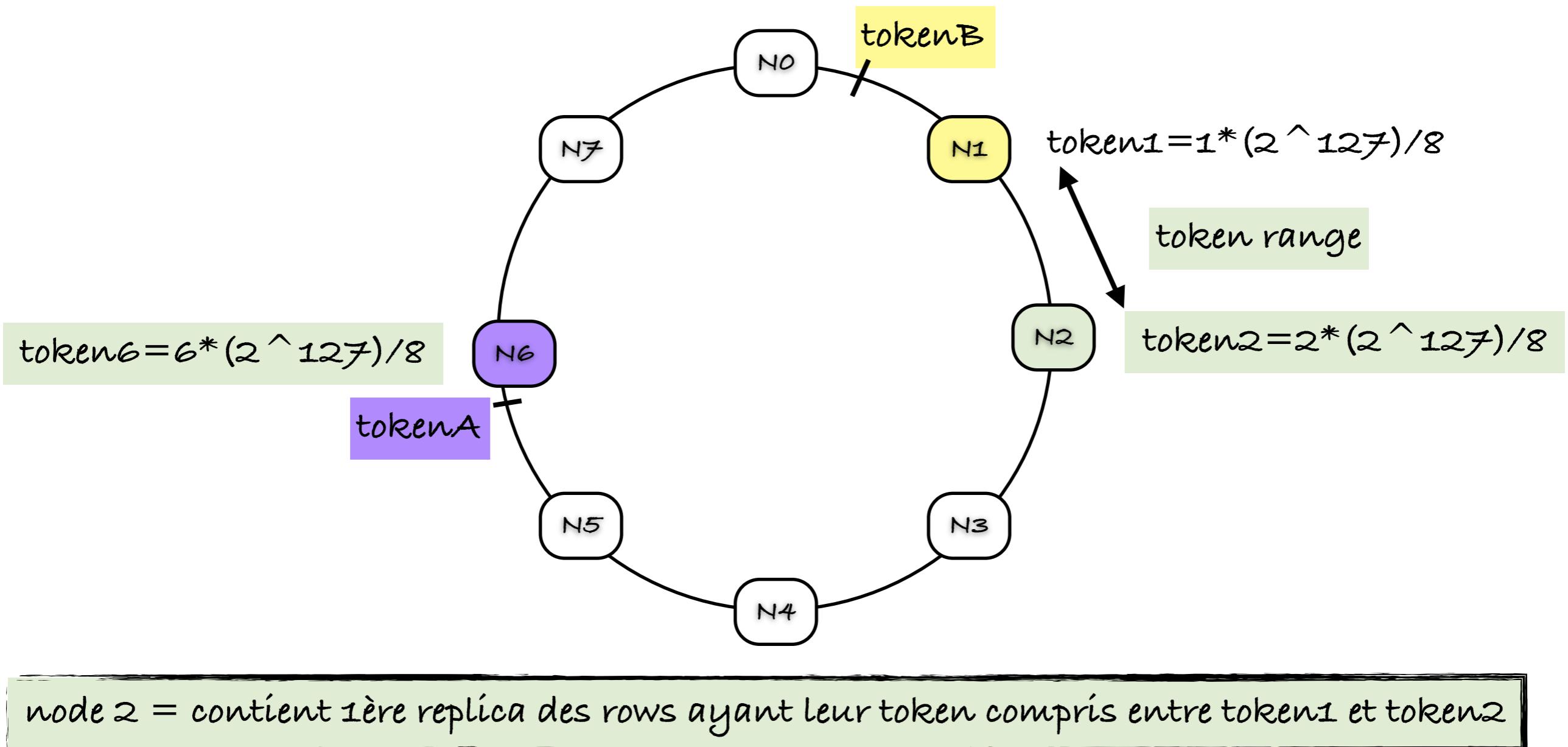
$f(\text{keyA})$

row keyB

$f(\text{keyB})$



Consistent Hashing



nodetool ring

```
MacBook-Pro-de-Nicolas-Romanetti:dsc-cassandra-1.0.8 nnromanetti$ bin/nodetool -h 10.1.1.1 ring
Address   DC   Rack   Status  State   Load          Owns      Token
10.1.1.1  1    1     Up      Normal  34,34 KB  3,33%    16448073285383832588838221182322480353
10.2.1.1  2    1     Up      Normal  34,34 KB  0,00%    1
10.3.1.1  3    1     Up      Normal  34,34 KB  0,00%    2
10.1.1.2  1    1     Up      Normal  43,33 KB  16,67%   28356863910078205288614550619314017621
10.2.1.2  2    1     Up      Normal  43,32 KB  0,00%    28356863910078205288614550619314017622
10.3.1.2  3    1     Up      Normal  43,33 KB  0,00%    28356863910078205288614550619314017623
10.1.1.3  1    1     Up      Normal  43,33 KB  16,67%   56713727820156410577229101238628035242
10.2.1.3  2    1     Up      Normal  38,72 KB  0,00%    56713727820156410577229101238628035243
10.3.1.3  3    1     Up      Normal  43,33 KB  0,00%    56713727820156410577229101238628035244
10.3.1.4  3    1     Up      Normal  43,33 KB  16,67%   85070591730234615865843651857942052866
10.1.1.5  1    1     Up      Normal  43,33 KB  16,67%   113427455640312821154458202477256070485
10.3.1.5  3    1     Up      Normal  43,33 KB  0,00%    113427455640312821154458202477256070487
10.2.1.5  2    1     Up      Normal  98,79 KB  3,25%    118959467553234022797339489649815746085
10.1.1.6  1    1     Up      Normal  145,25 KB 13,06%   141181442192599215569169119445432872465
10.2.1.6  2    1     Up      Normal  43,33 KB  0,35%    141784319550391026443072753096570088107
10.3.1.6  3    1     Up      Normal  43,33 KB  0,00%    141784319550391026443072753096570088108
10.1.1.4  1    1     Up      Normal  47,71 KB  3,03%   146946219305601782975974550082043979626
10.2.1.4  2    1     Up      Normal  47,71 KB  10,31%   16448073285383832588838221182322480353
MacBook-Pro-de-Nicolas-Romanetti:dsc-cassandra-1.0.8 nnromanetti$
```



Vocabulaire

Node

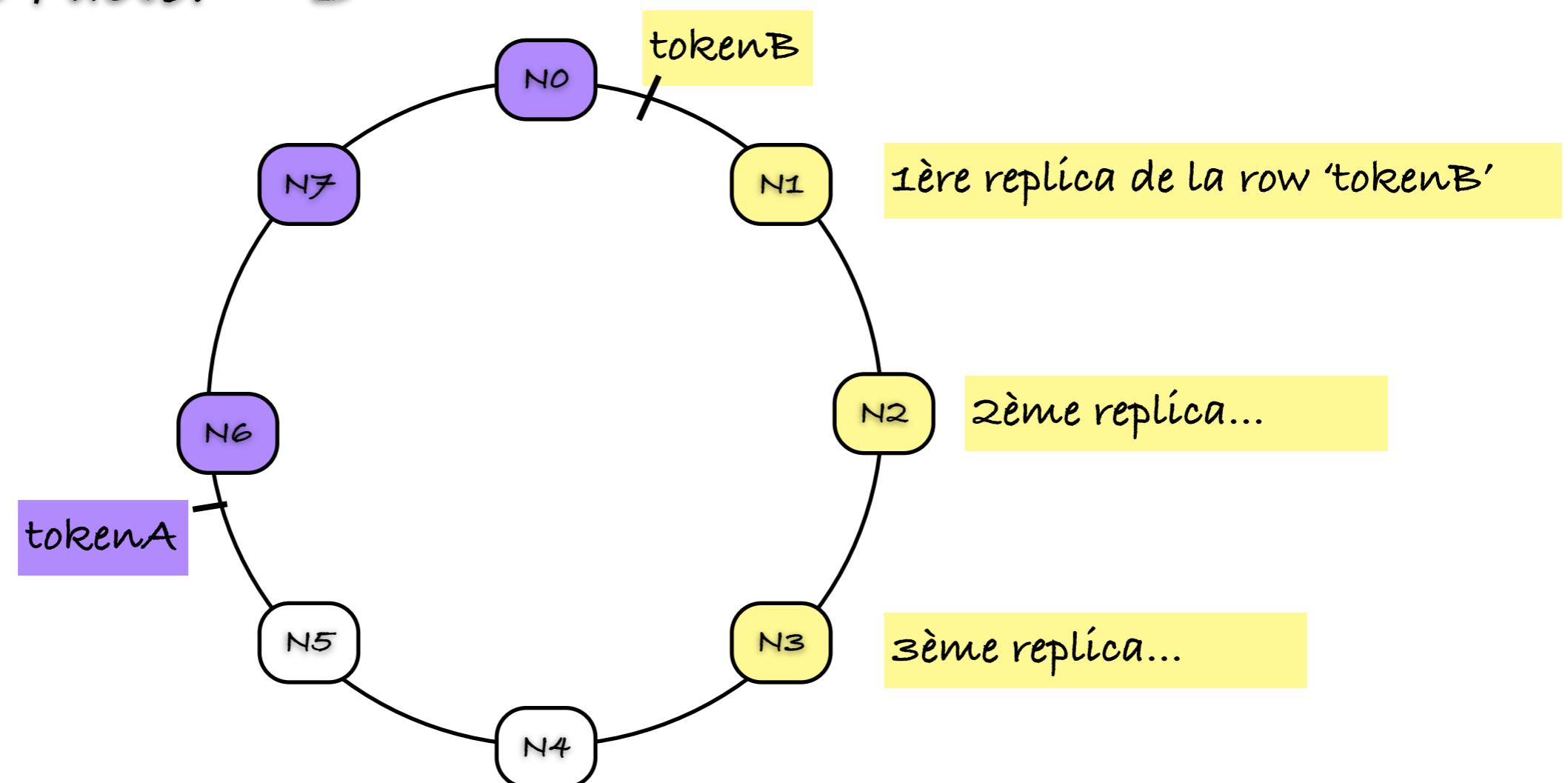
Replication Factor

Replica



Stratégie de placement

Replication Factor = 3



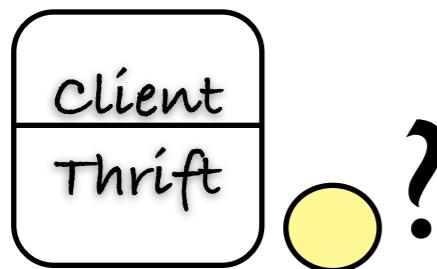
API client

- Hector... Frère de Cassandre
- Astyanax (créé par Netflix)... Fils d'Hector
- Driver CQL...

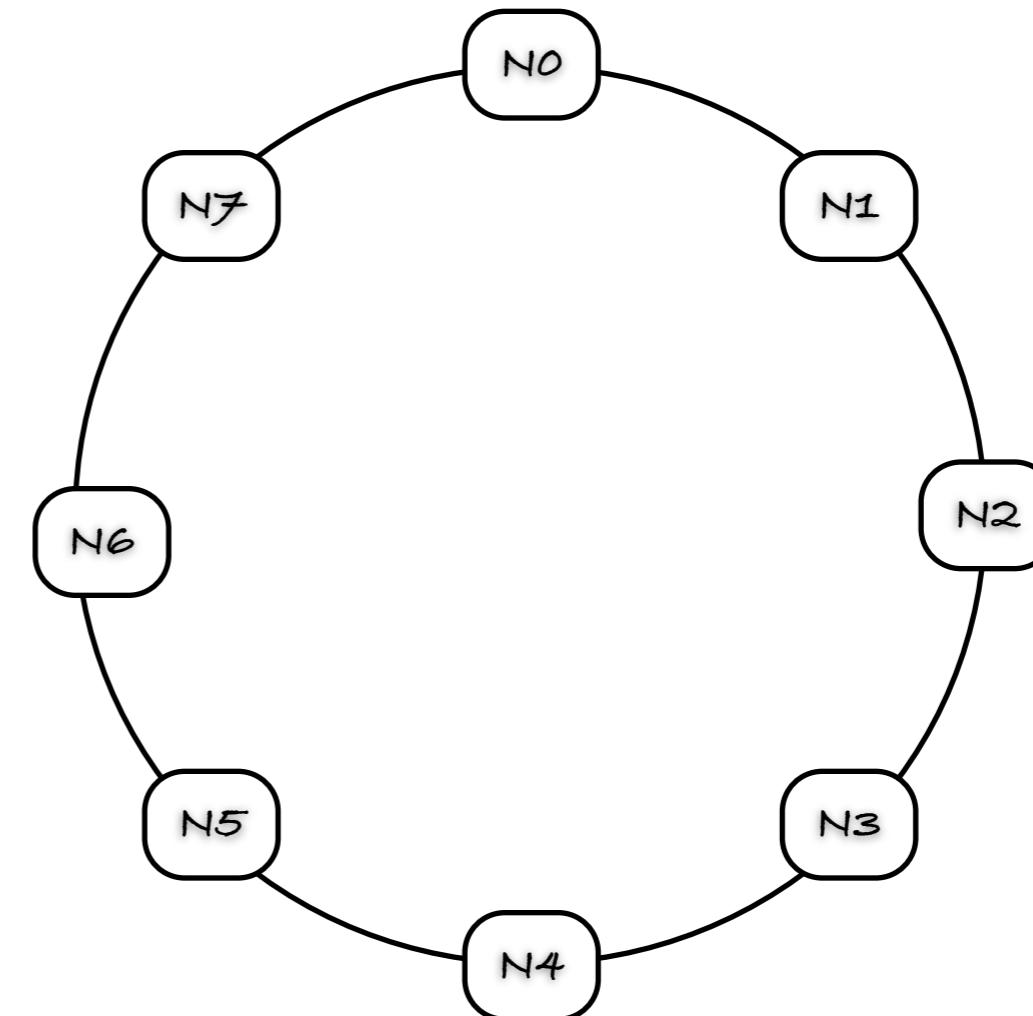
Thrift



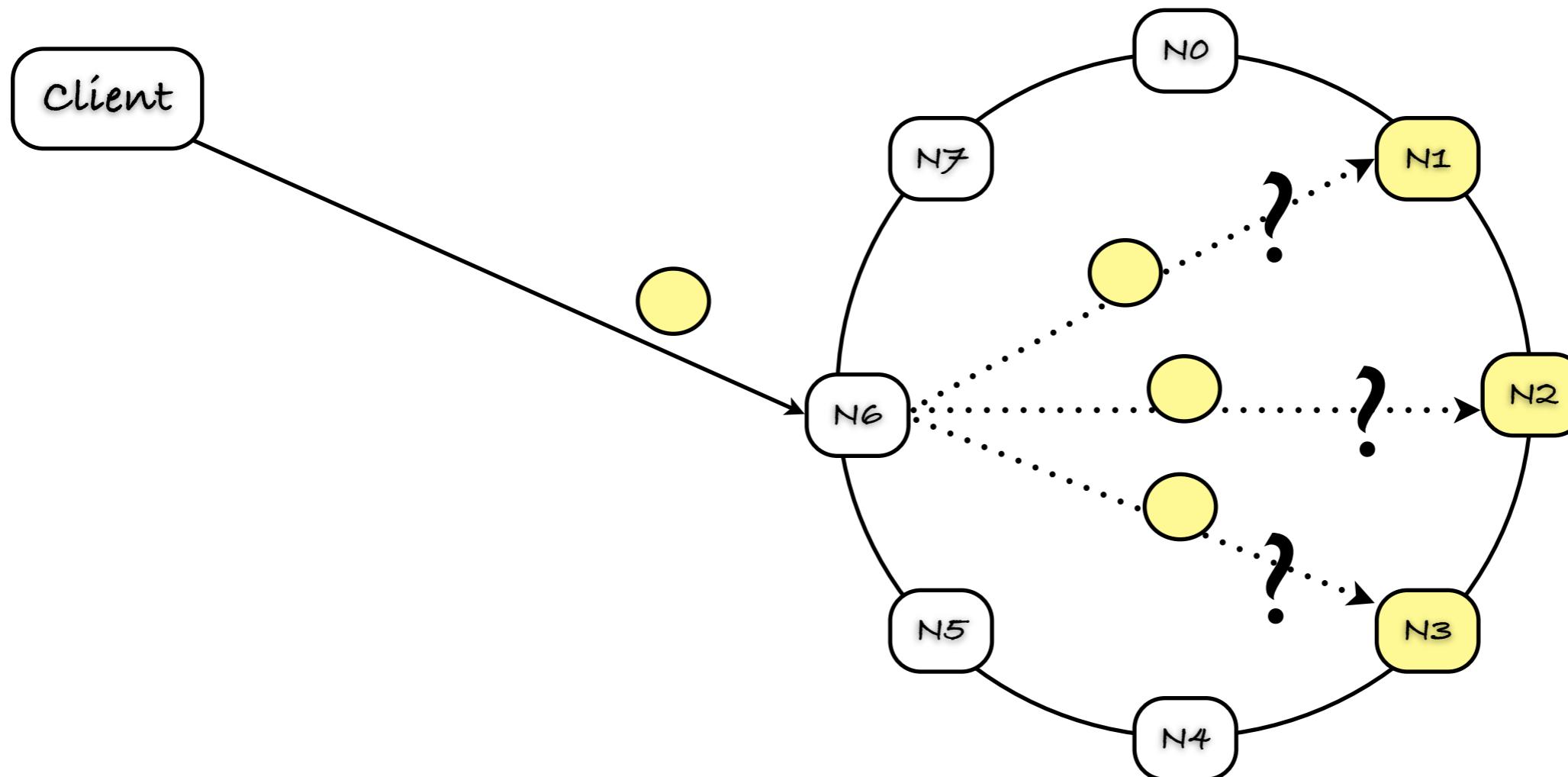
Ecriture: qui contacter?



Une replica?



Attendre les réponses?



solution:Tunable Consistency



CONSISTENCY LEVEL

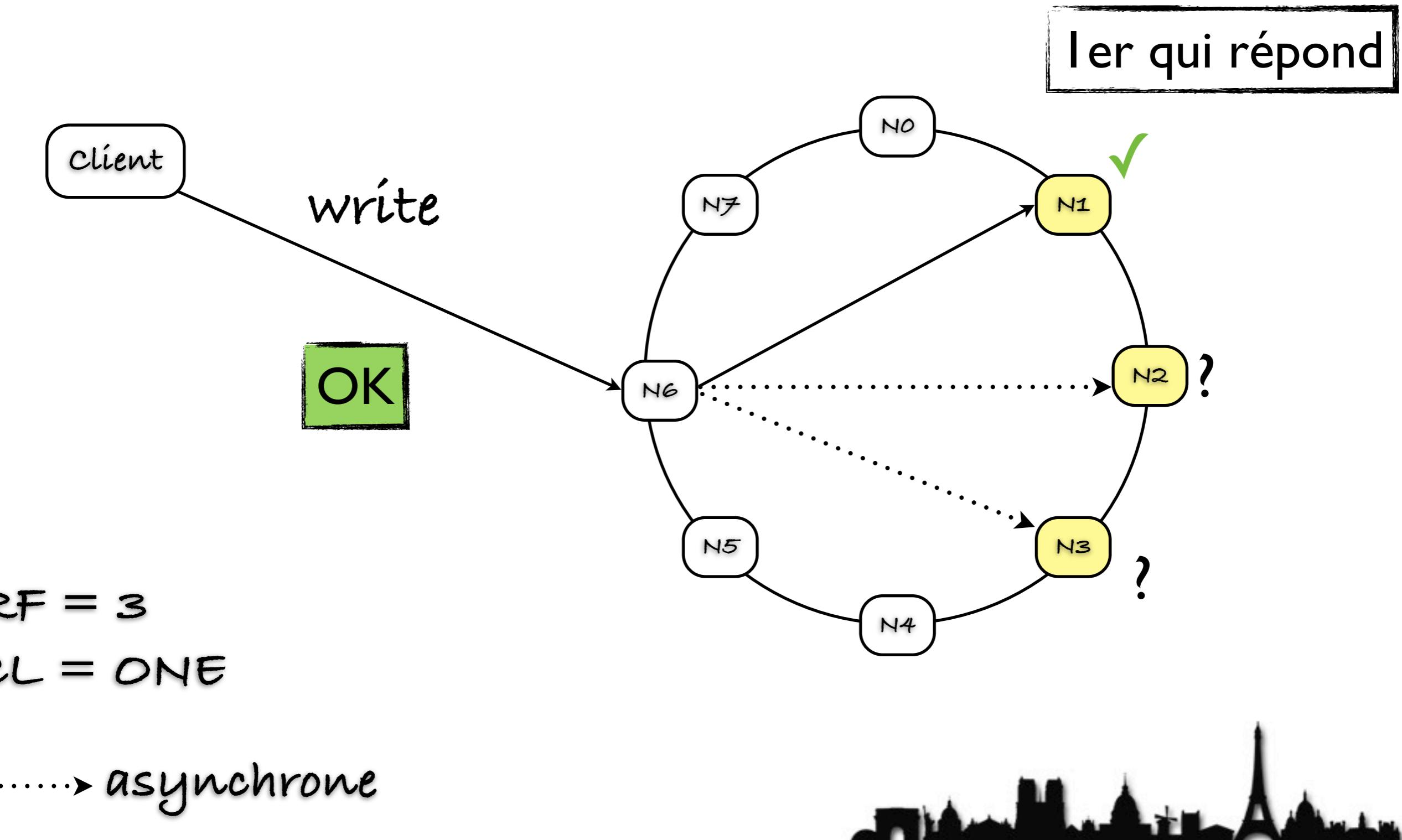


CL - ONE

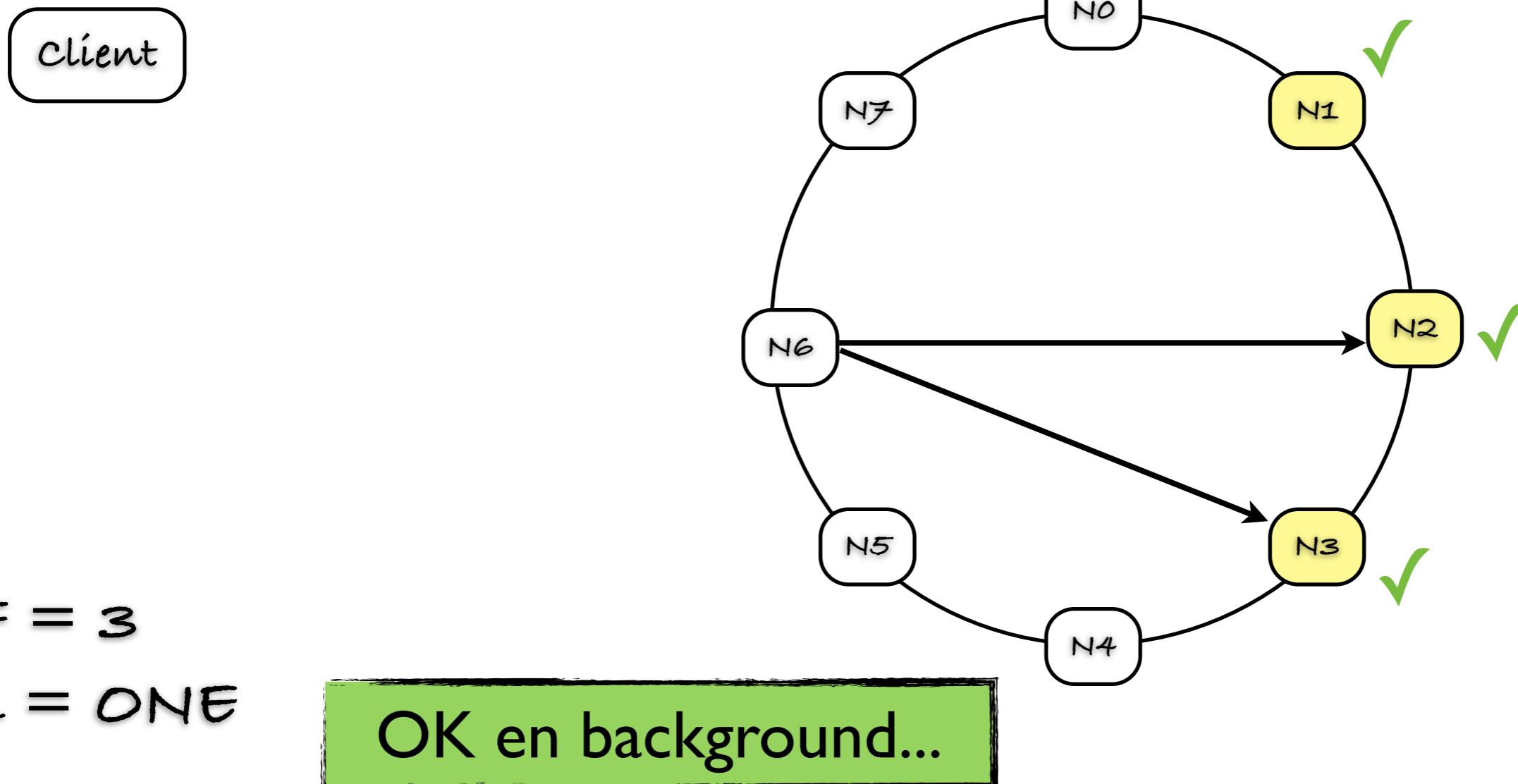
écriture



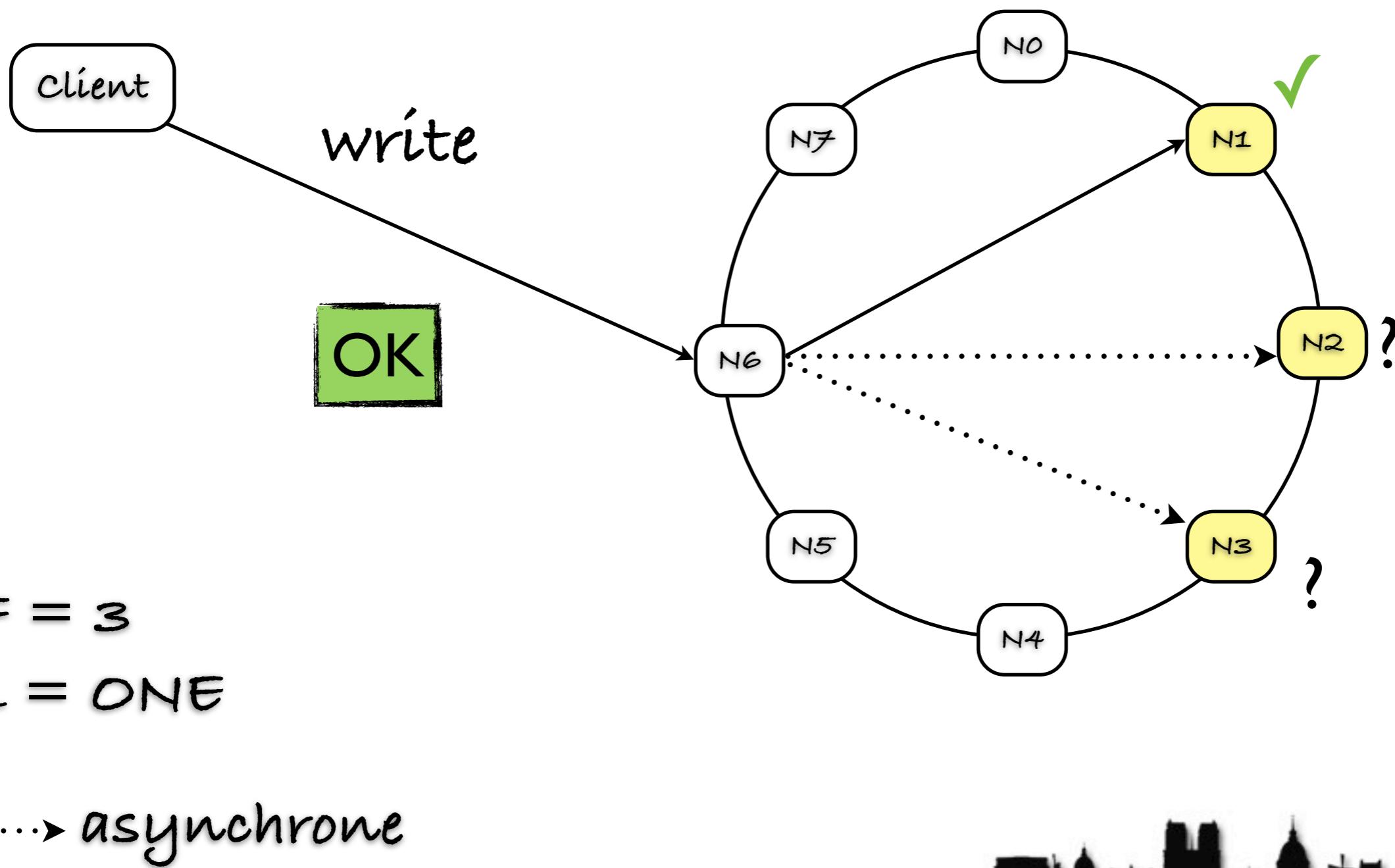
Exemple I



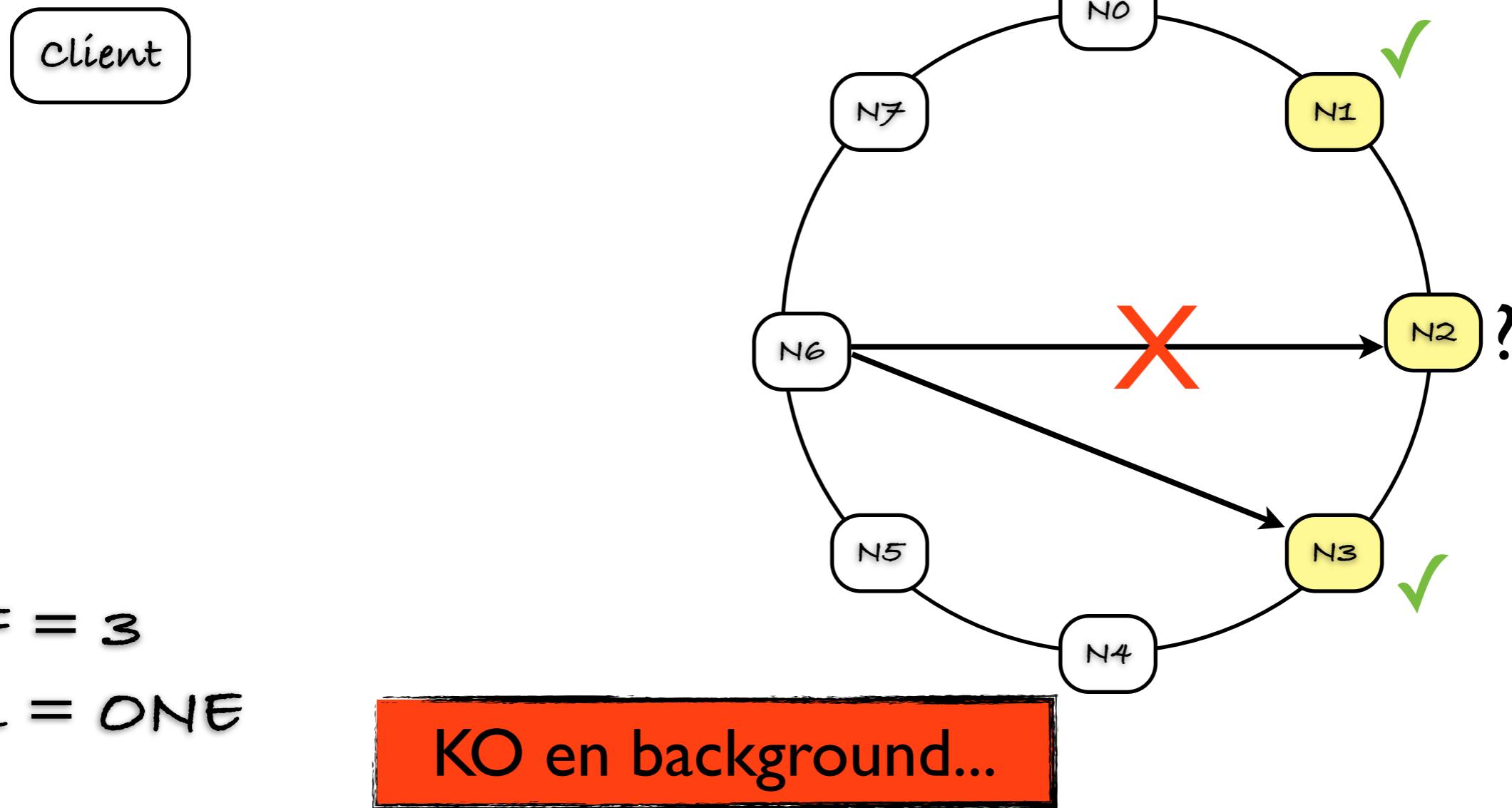
Exemple I (suite)



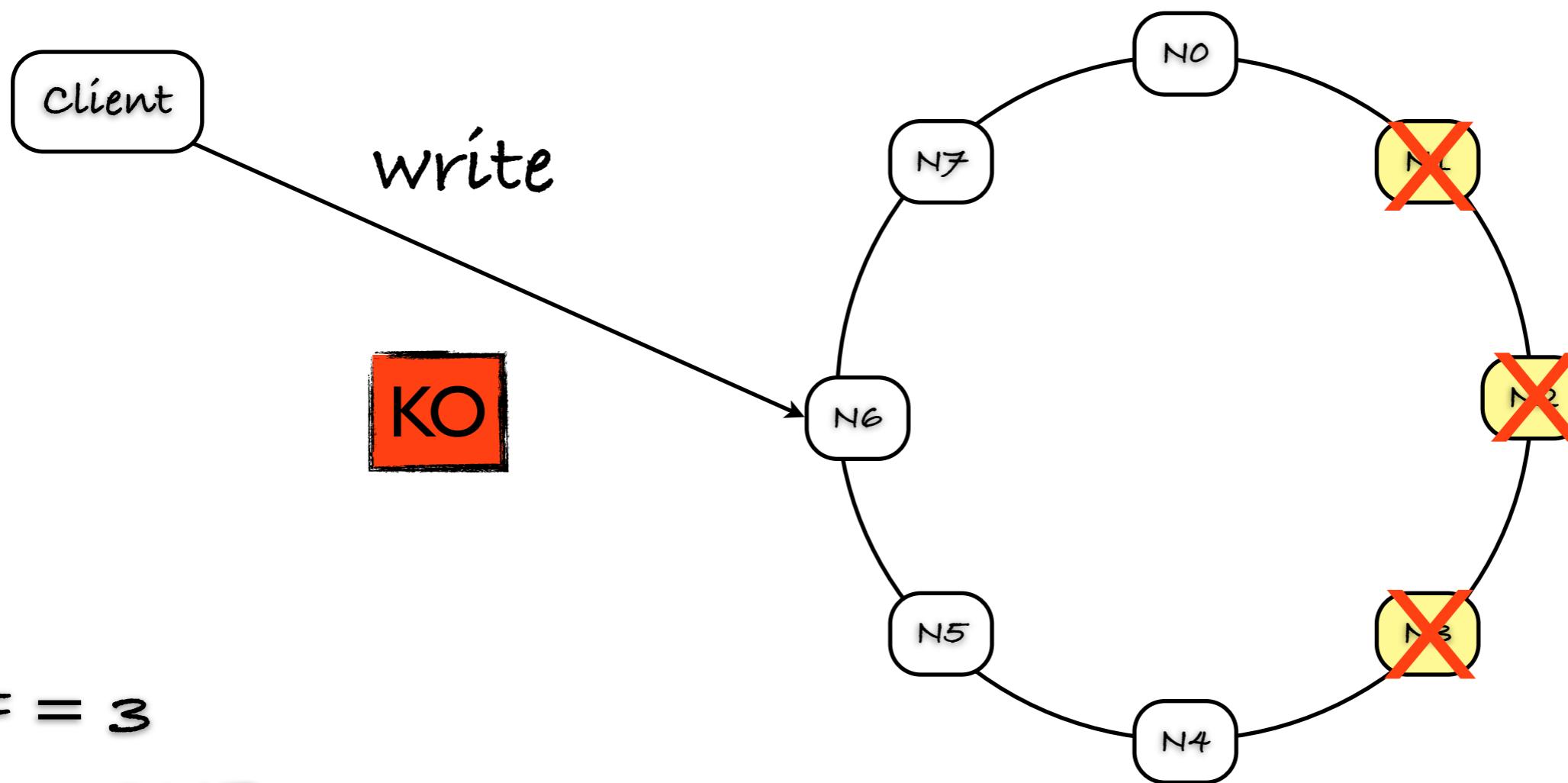
Exemple 2



Exemple 2 (suite)



Exemple 3

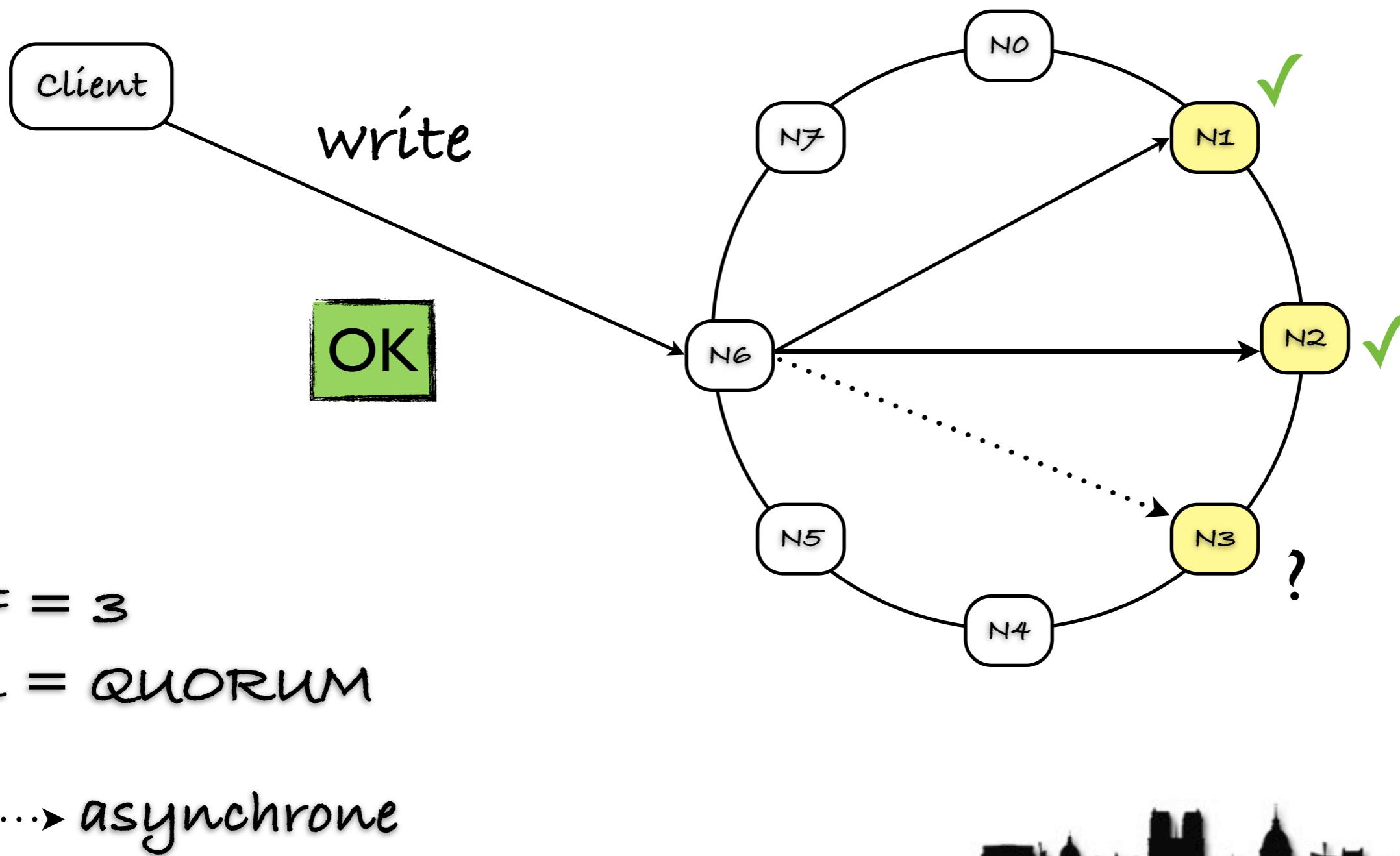


CL - QUORUM

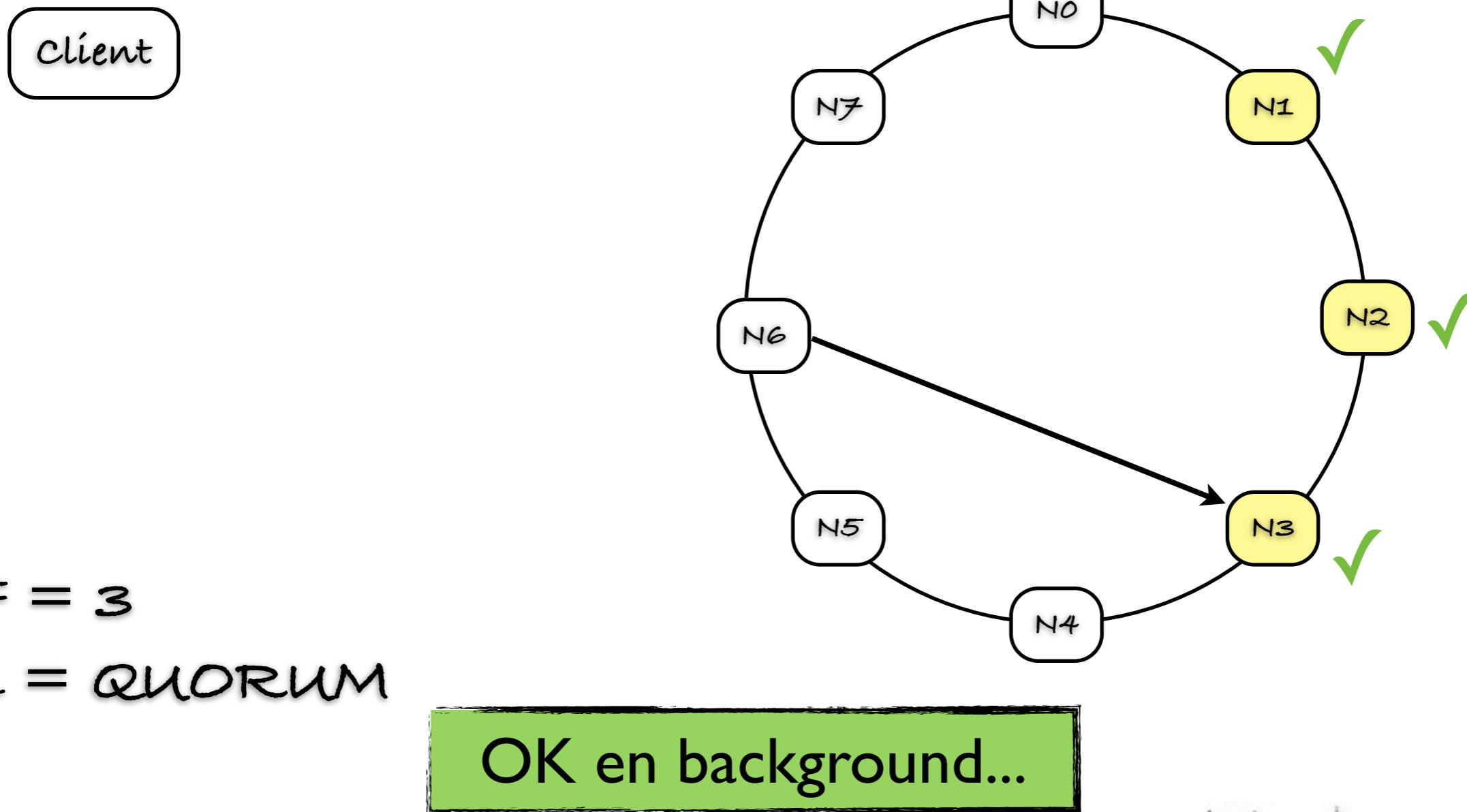
écriture



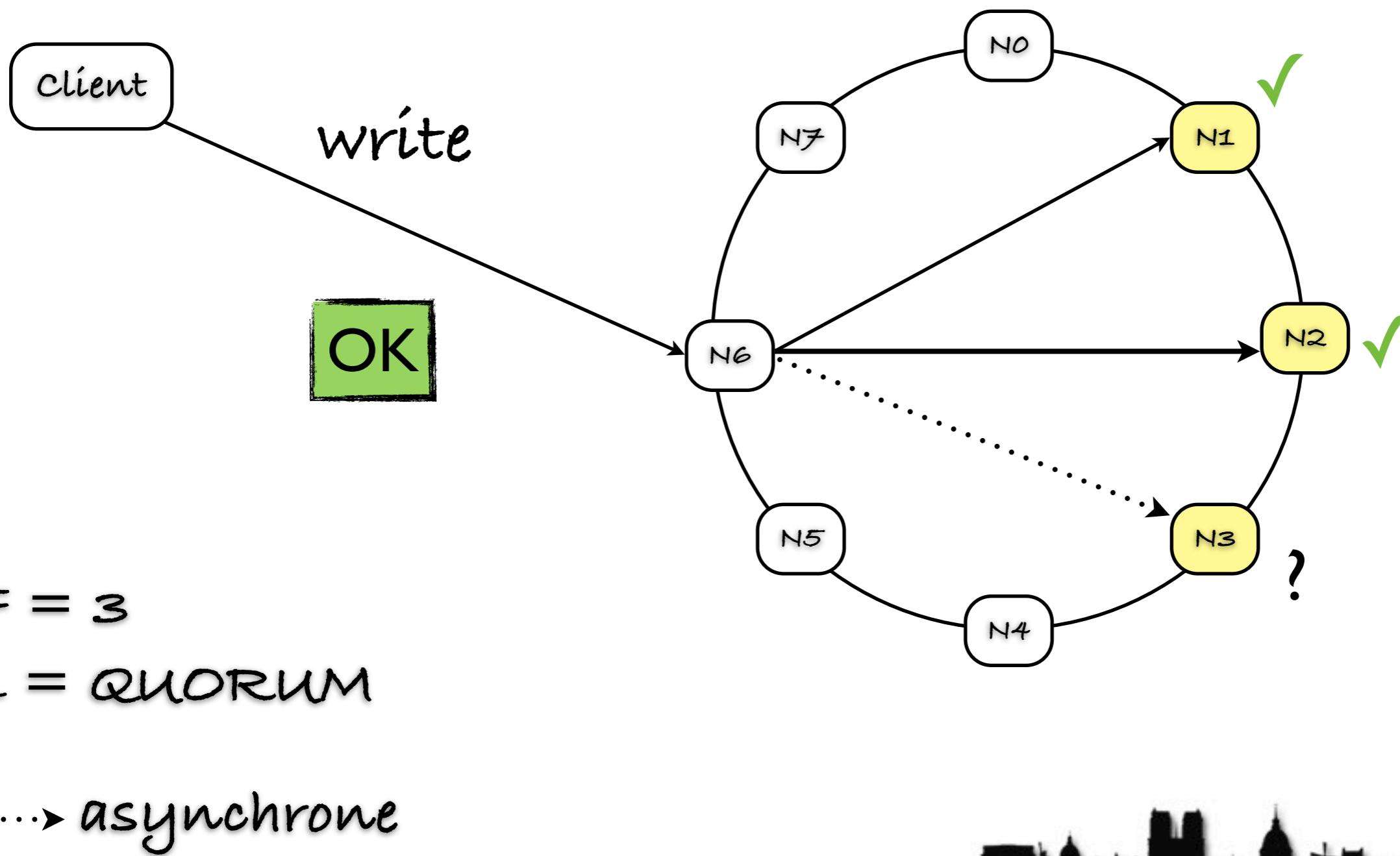
Exemple I



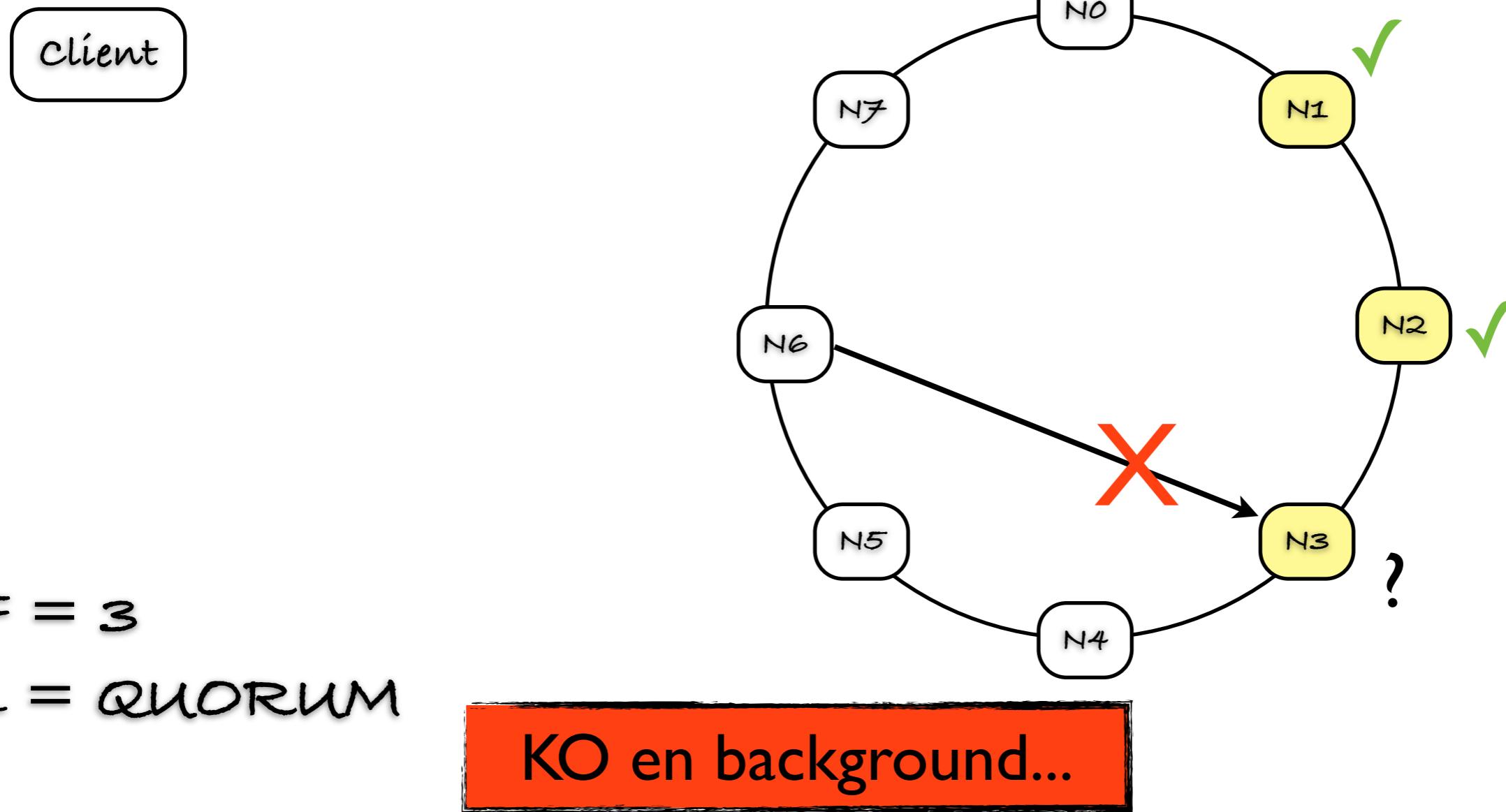
Exemple I (suite)



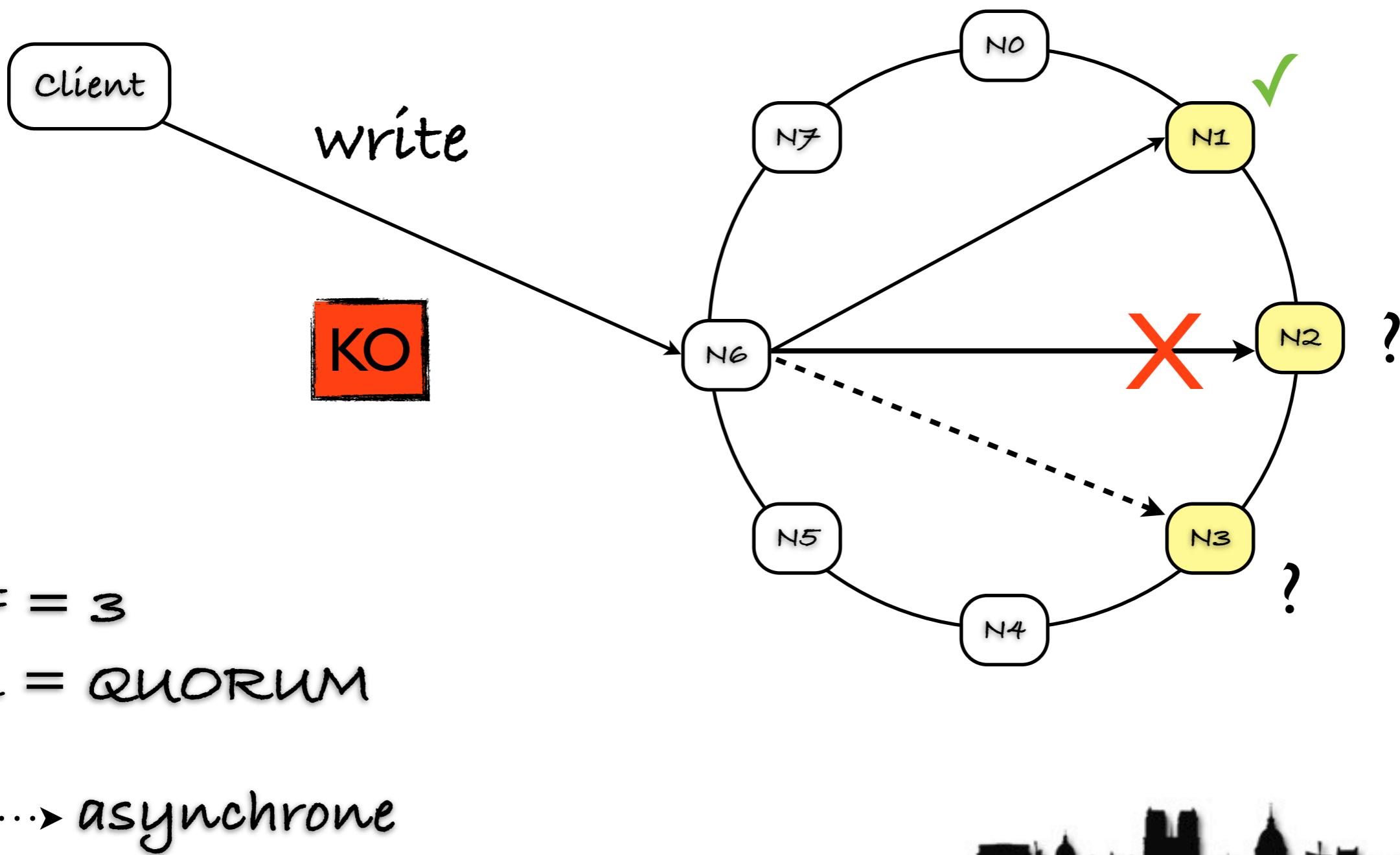
Exemple 2



Exemple 2 (suite)



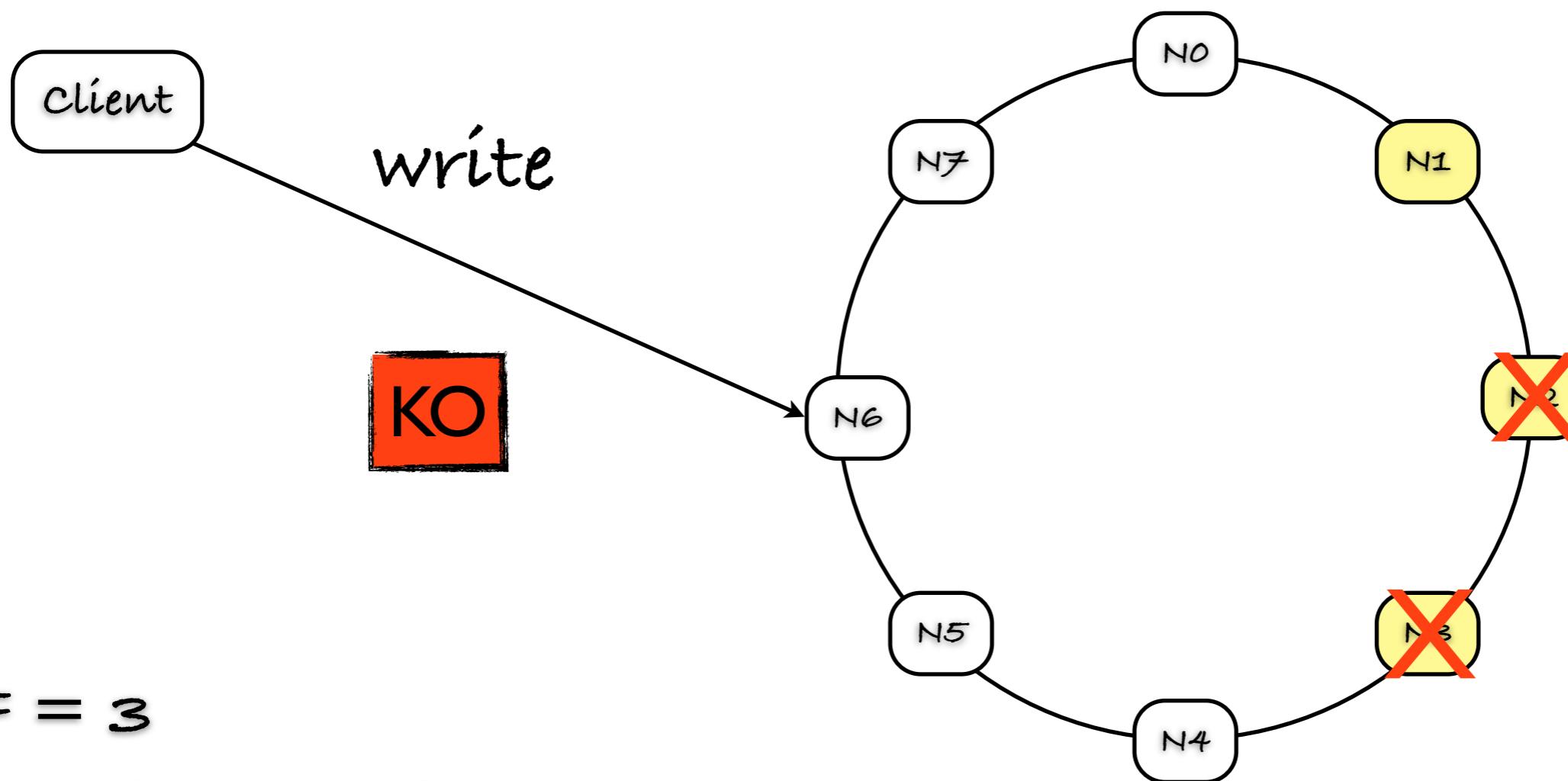
Exemple 3



.....→ asynchrone



Exemple 4



RF = 3

CL = QUORUM

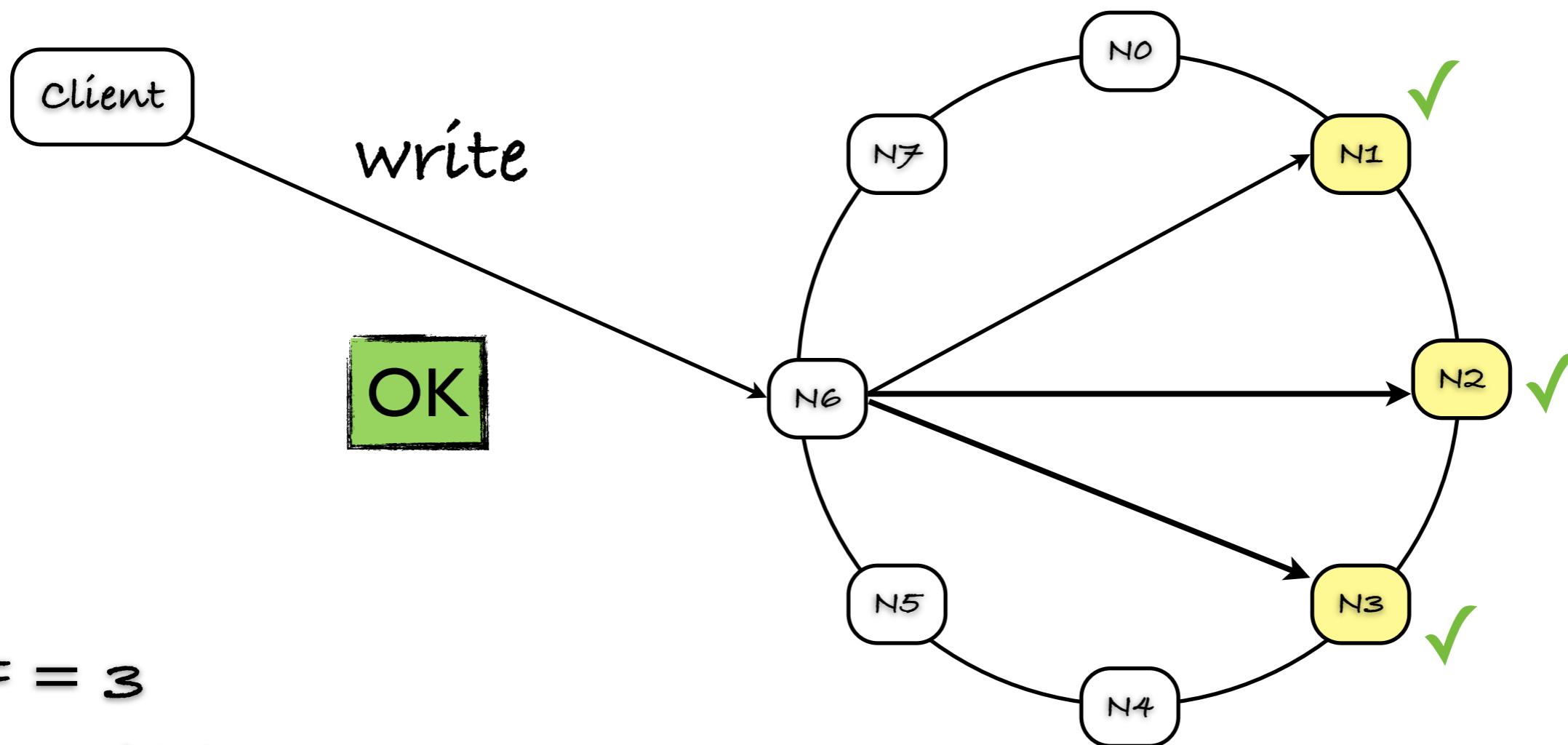


CL - ALL

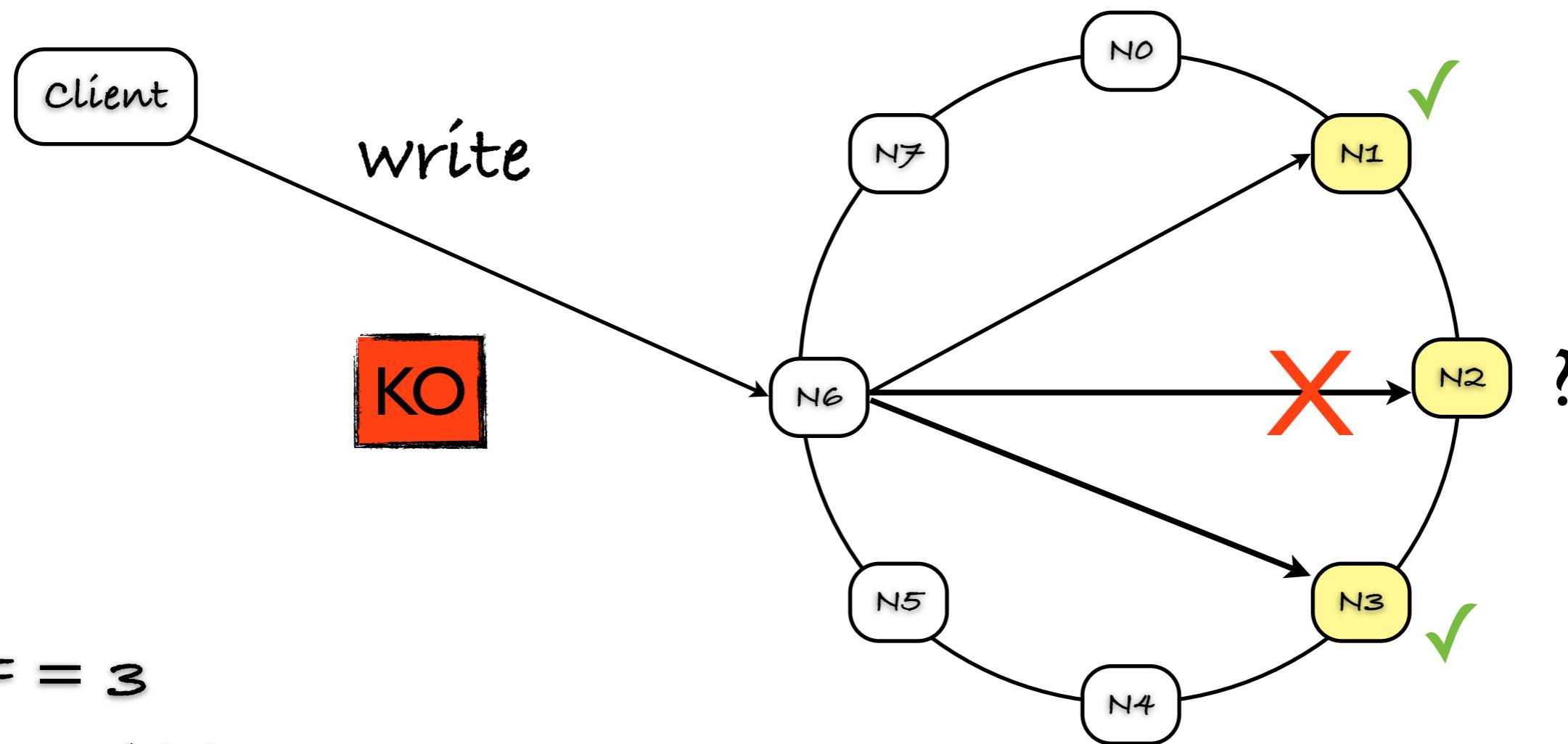
écriture



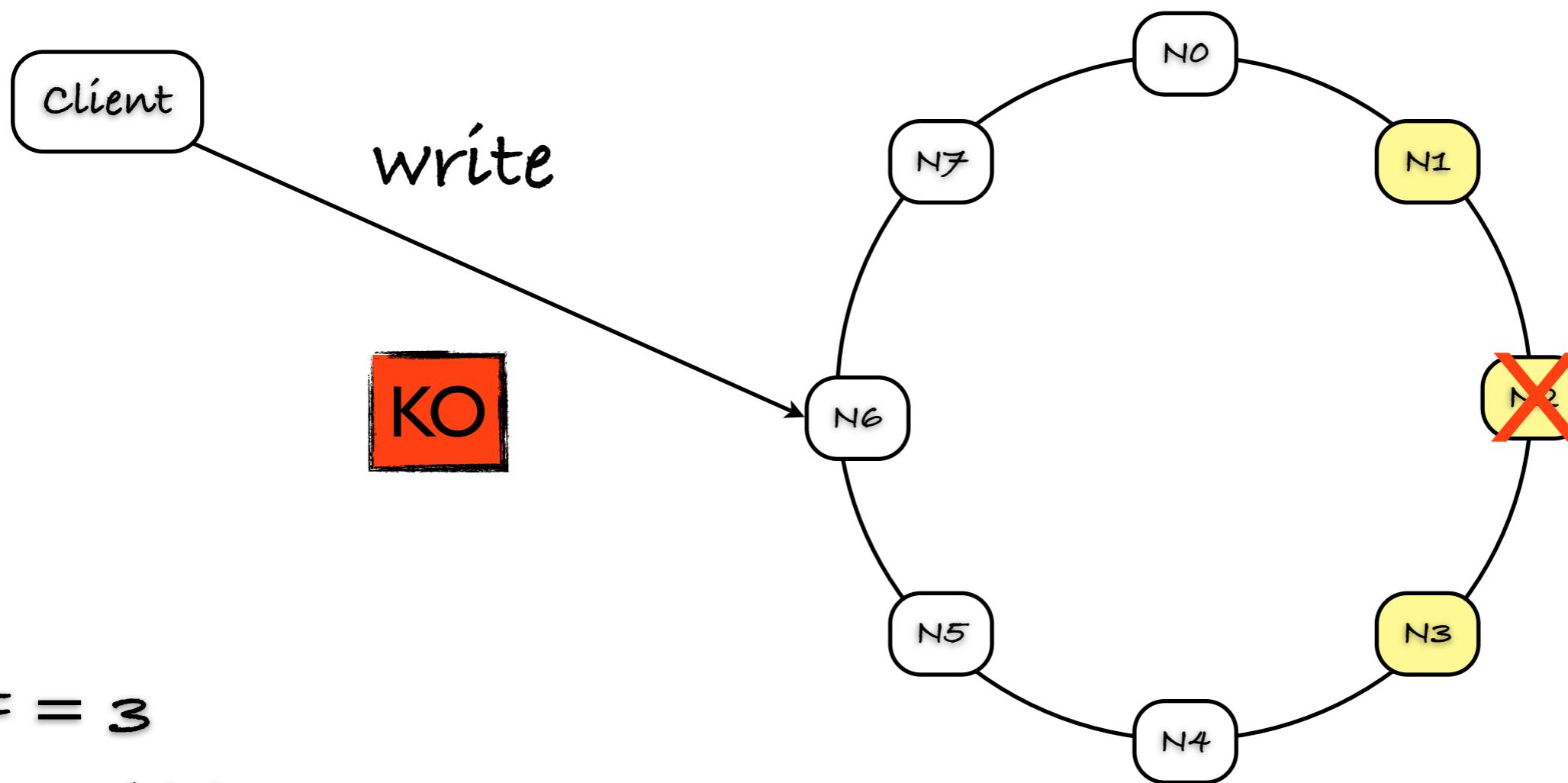
Exemple I



Exemple 2



Exemple 3



Sondage

Rollback ou pas?

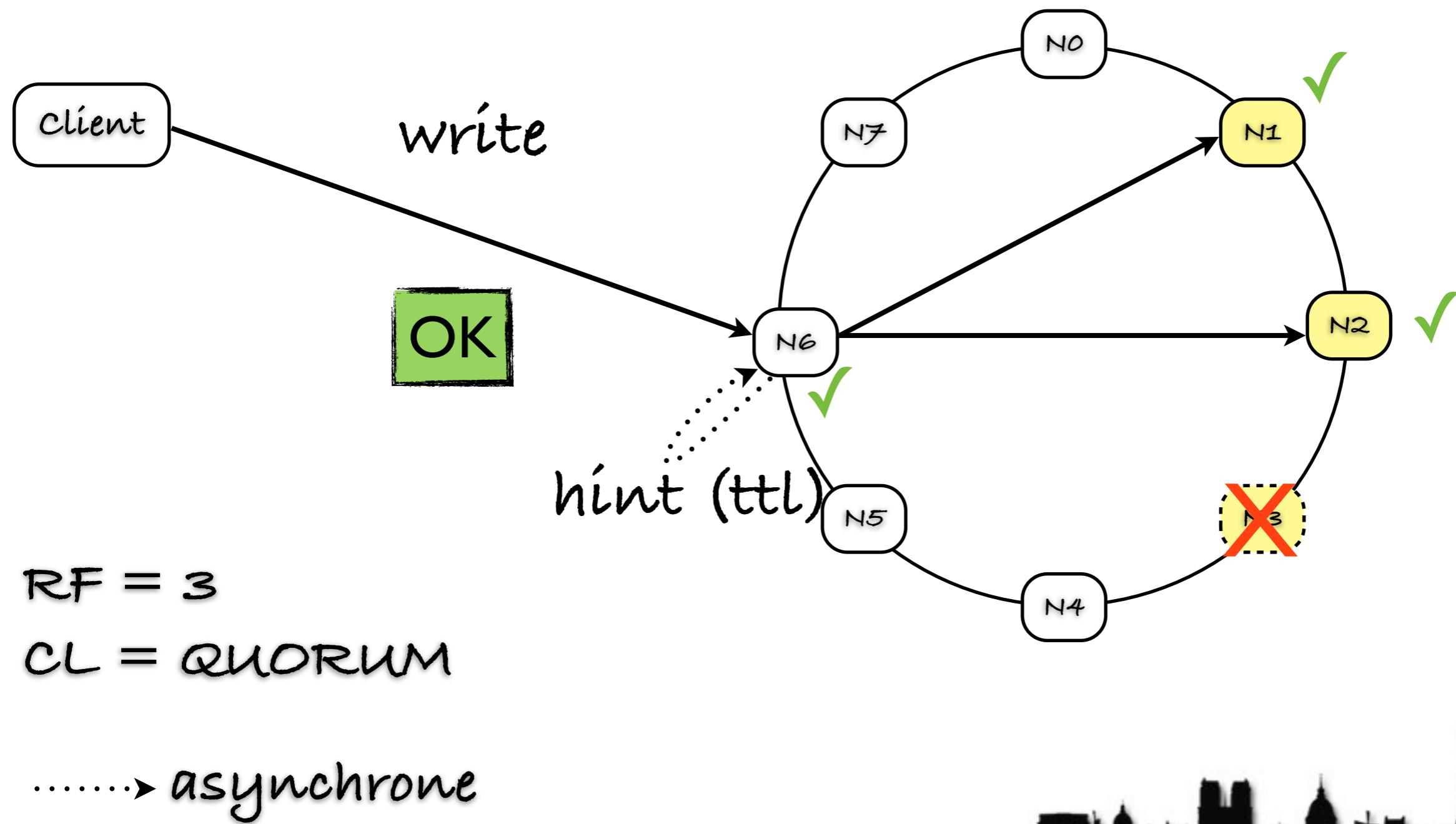


Pas de Rollback...

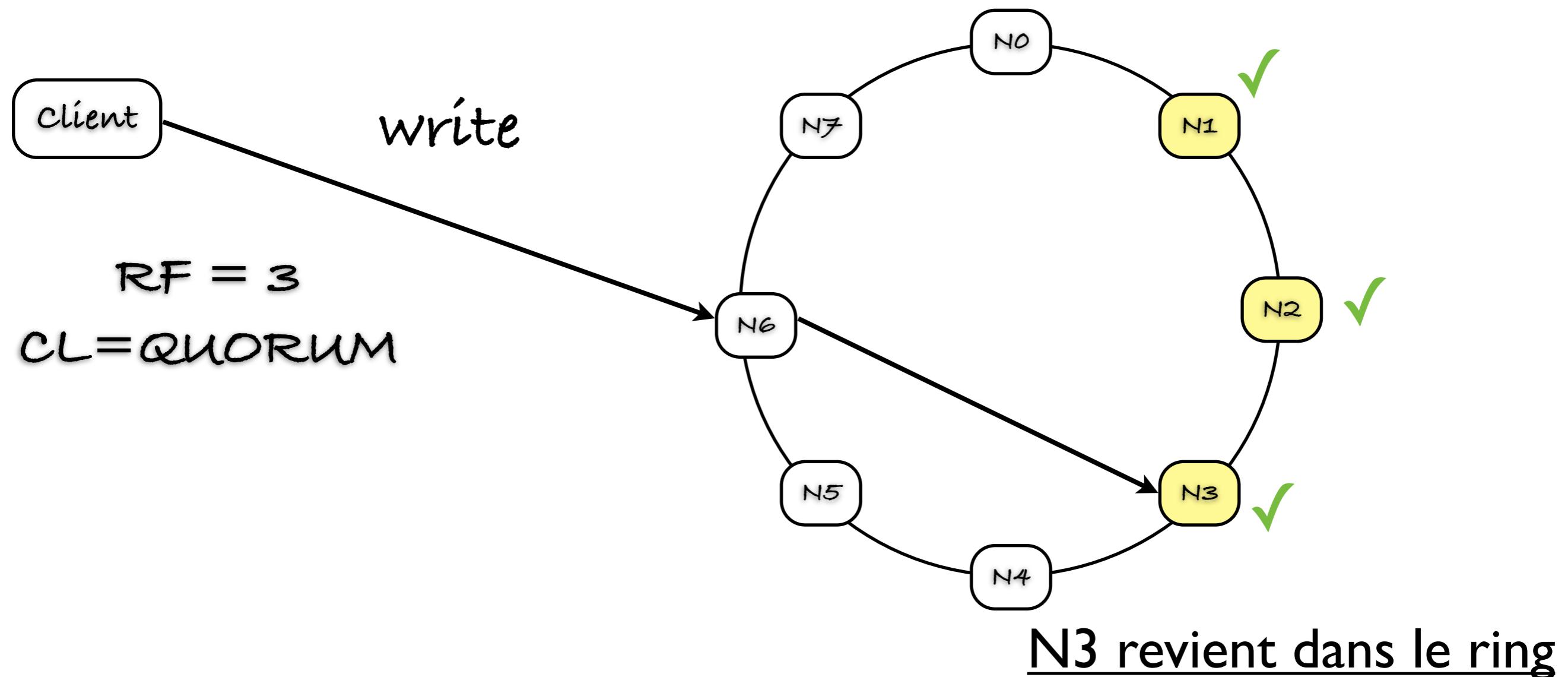
- La donnée va être propagée
- Mais comment?



Hinted Handoff



Hinted Handoff



Hinted Handoff

- 1er mécanisme de reprise sur erreur
- TTL du hint

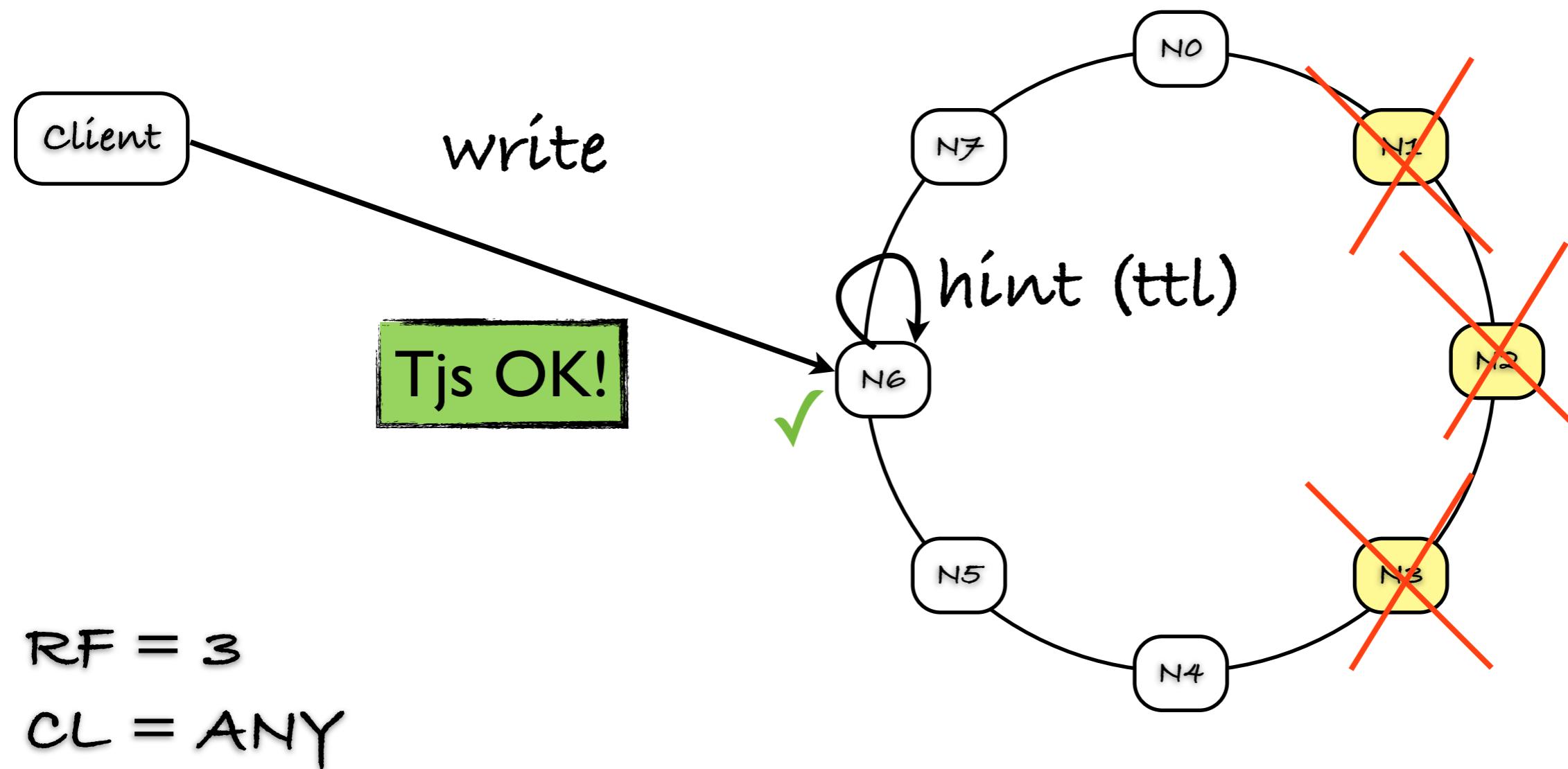


CL - ANY

écriture



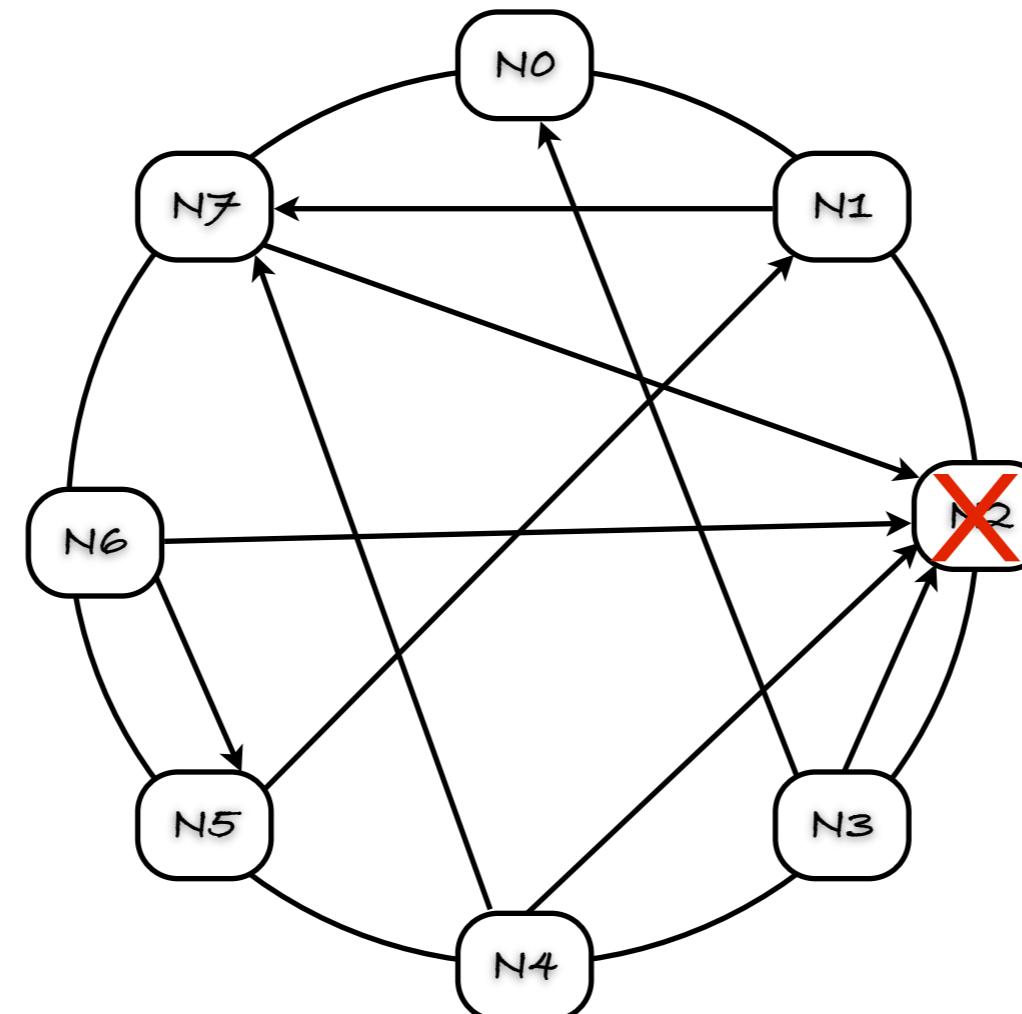
Exemple CL Any



Etat du ring?



Gossip Protocol



- Pas d'ilots
- Découverte du ring
- Echange d'info sur le système



$\sim\varphi$ Accrual Failure Detector

- Ni oui ni non
- Probabiliste

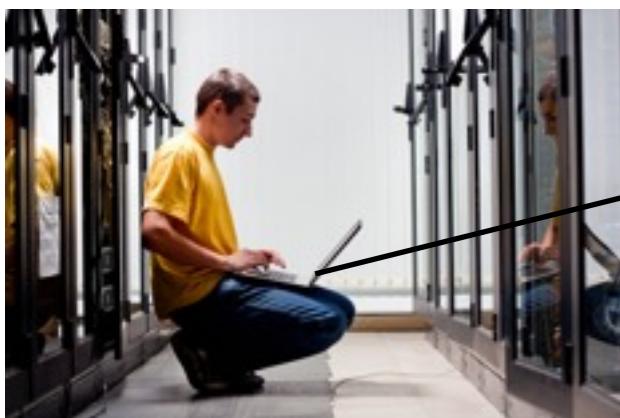
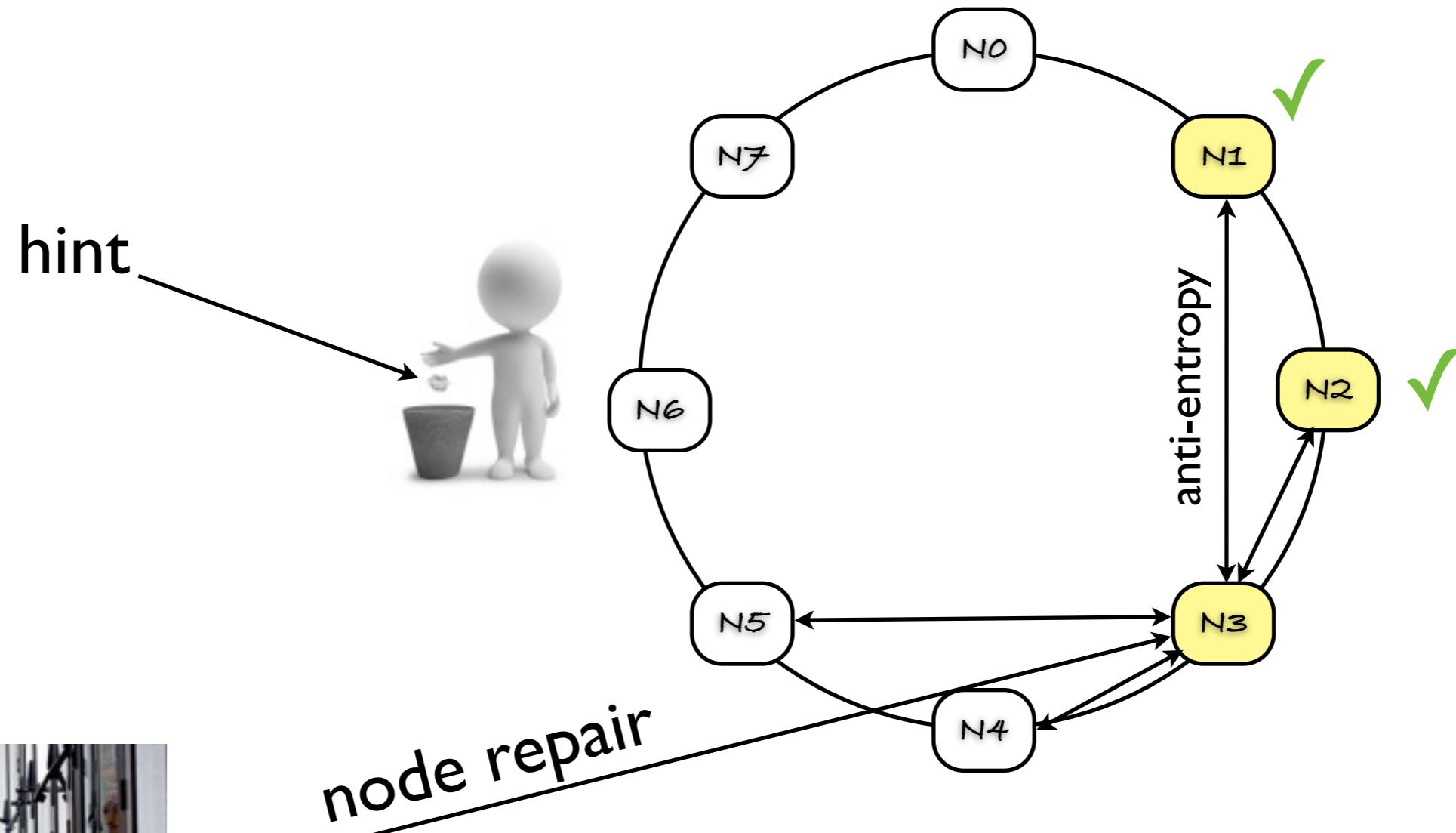


Expiration du hint

- Comment synchroniser les données lorsque le noeud revient dans le ring après expiration des hints?



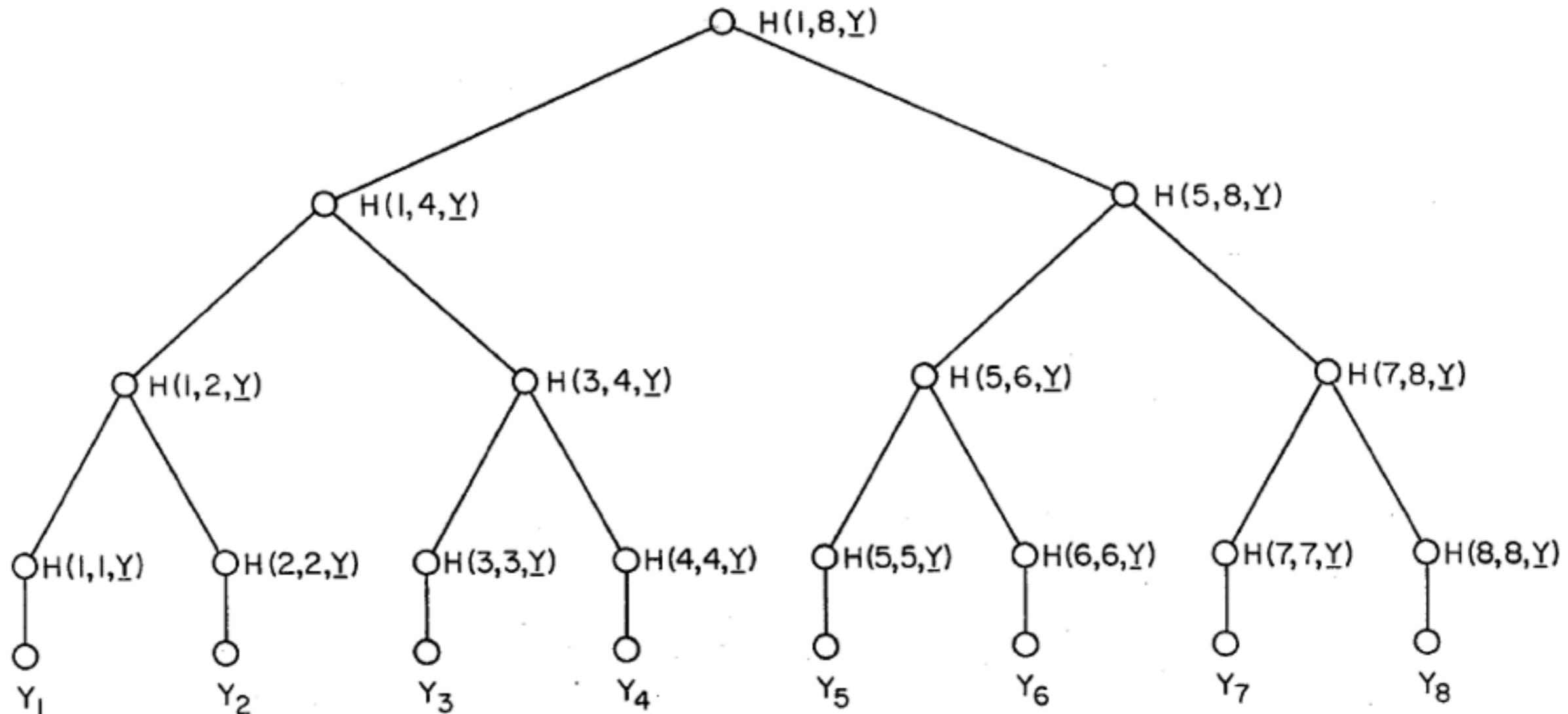
Anti-Entropy



2ème mécanisme de reprise sur erreur



Merkle Tree (1982)



United States Patent [19]
Merkle Best At

[54] METHOD OF PROVIDING DIGITAL
SIGNATURES

[75] Inventor: Ralph C. Merkle, Mountain View,
Calif.



En lecture

- Tunable consistency

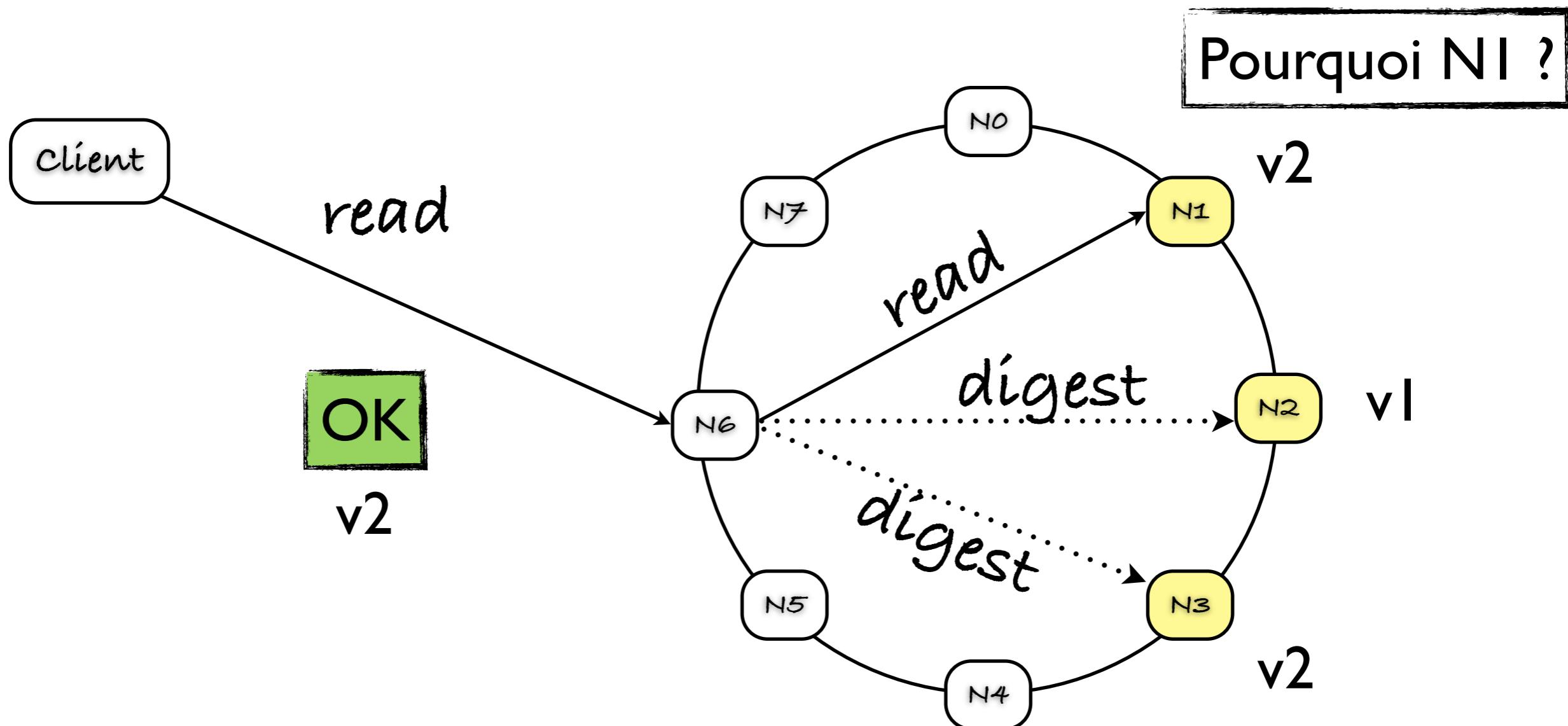


CL - ONE

lecture



Exemple I



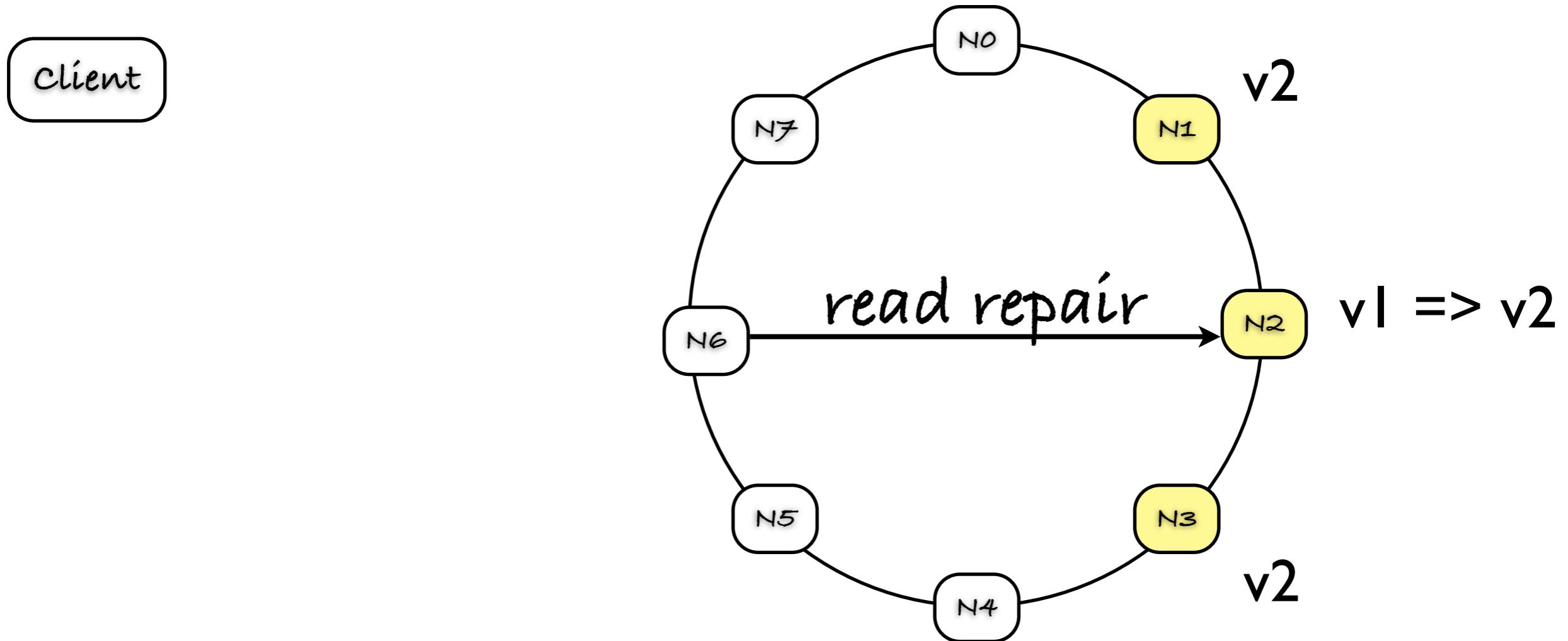
RF = 3

CL = ONE

Donnée la plus récente parmi ONE replica



Exemple I (suite)

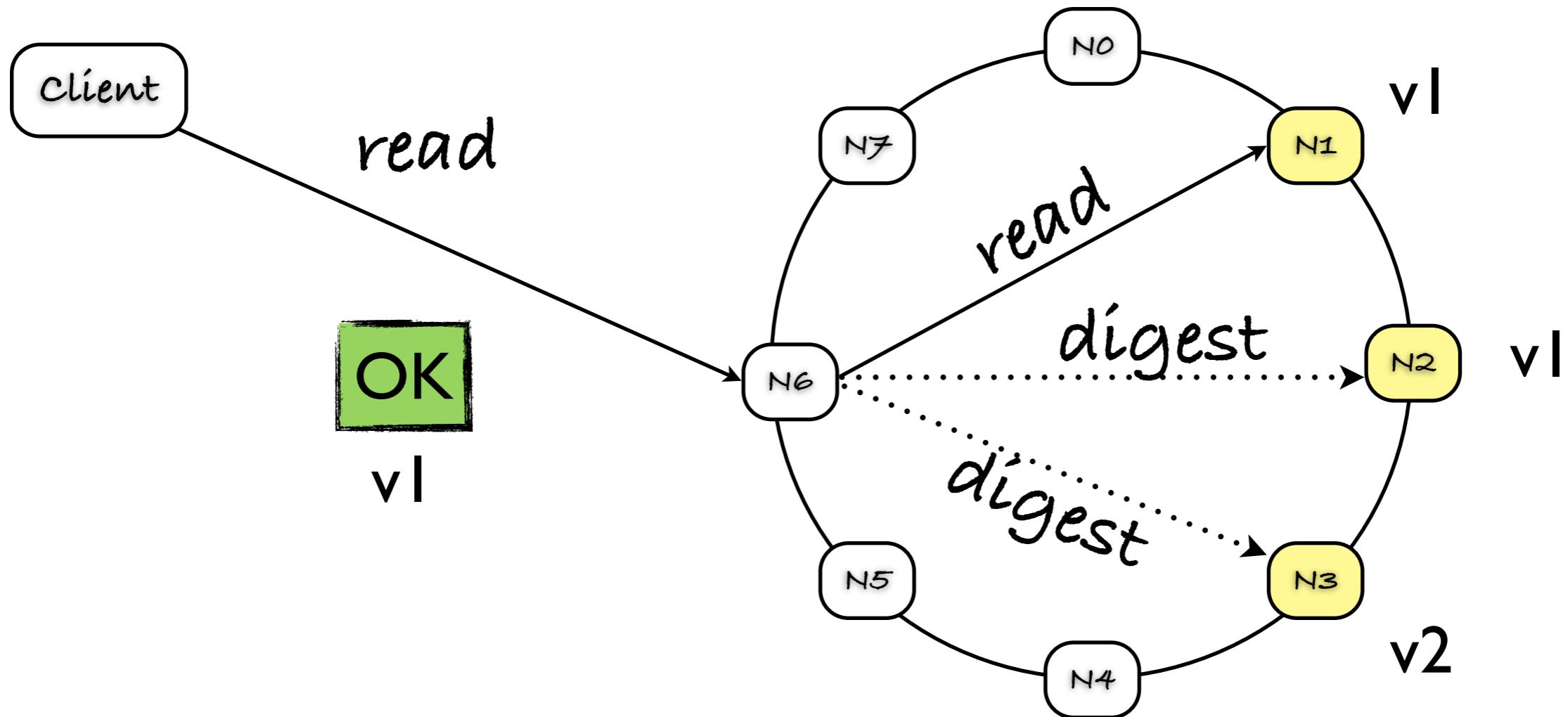


RF = 3

CL = ONE



Exemple 2



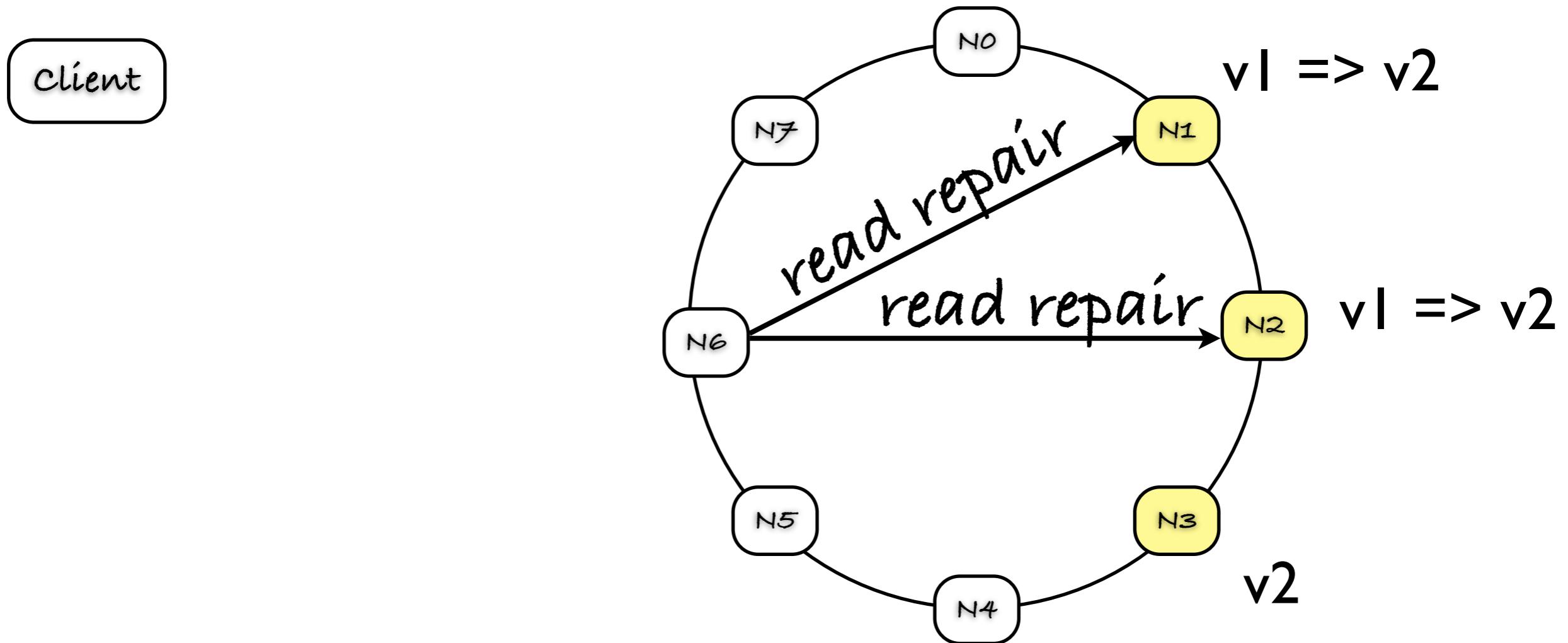
RF = 3

CL = ONE

Donnée la plus récente parmi ONE replica



Exemple 2 (suite)



RF = 3

CL = ONE



Read Repair

- 3ème mécanisme de reprise sur erreur
- Configurable (% de chance)

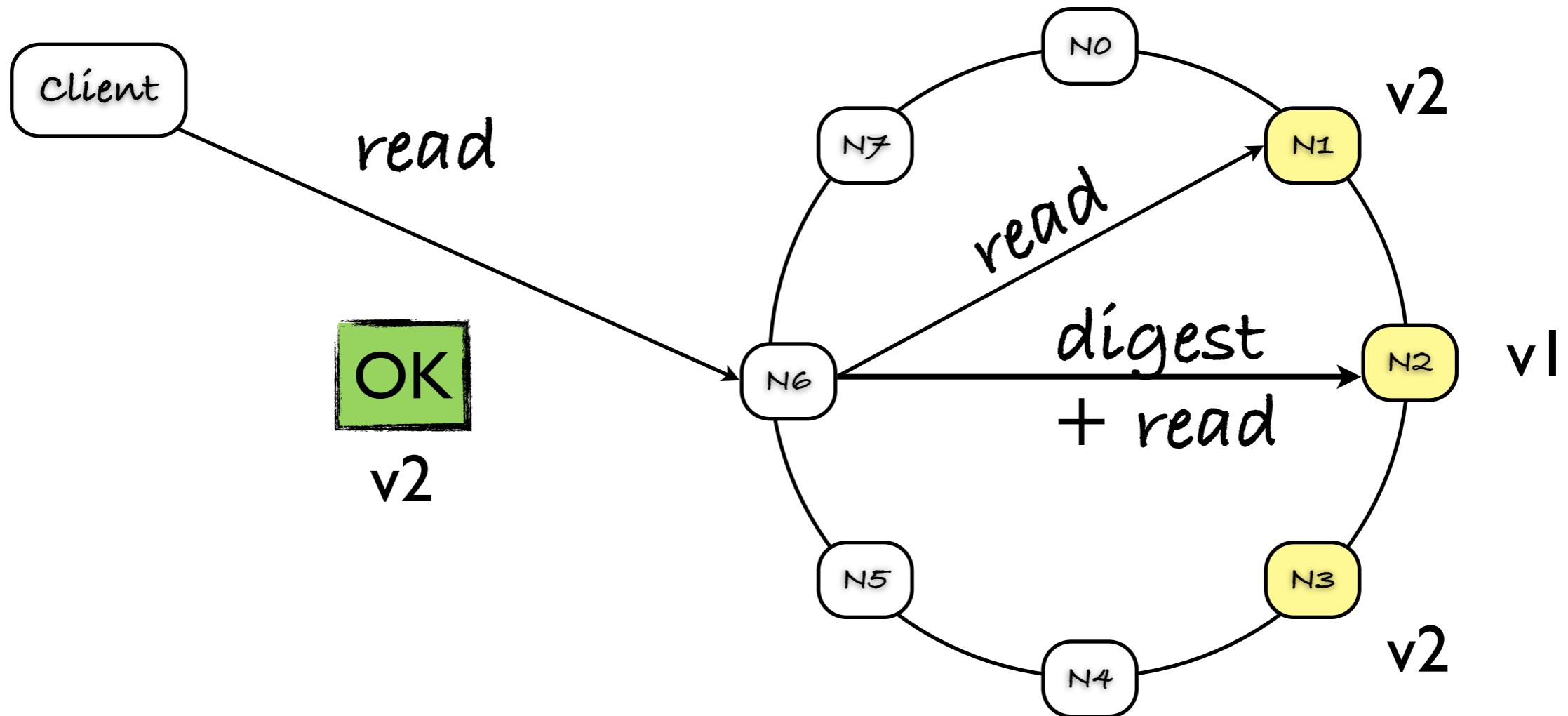


CL - QUORUM

lecture



Exemple I



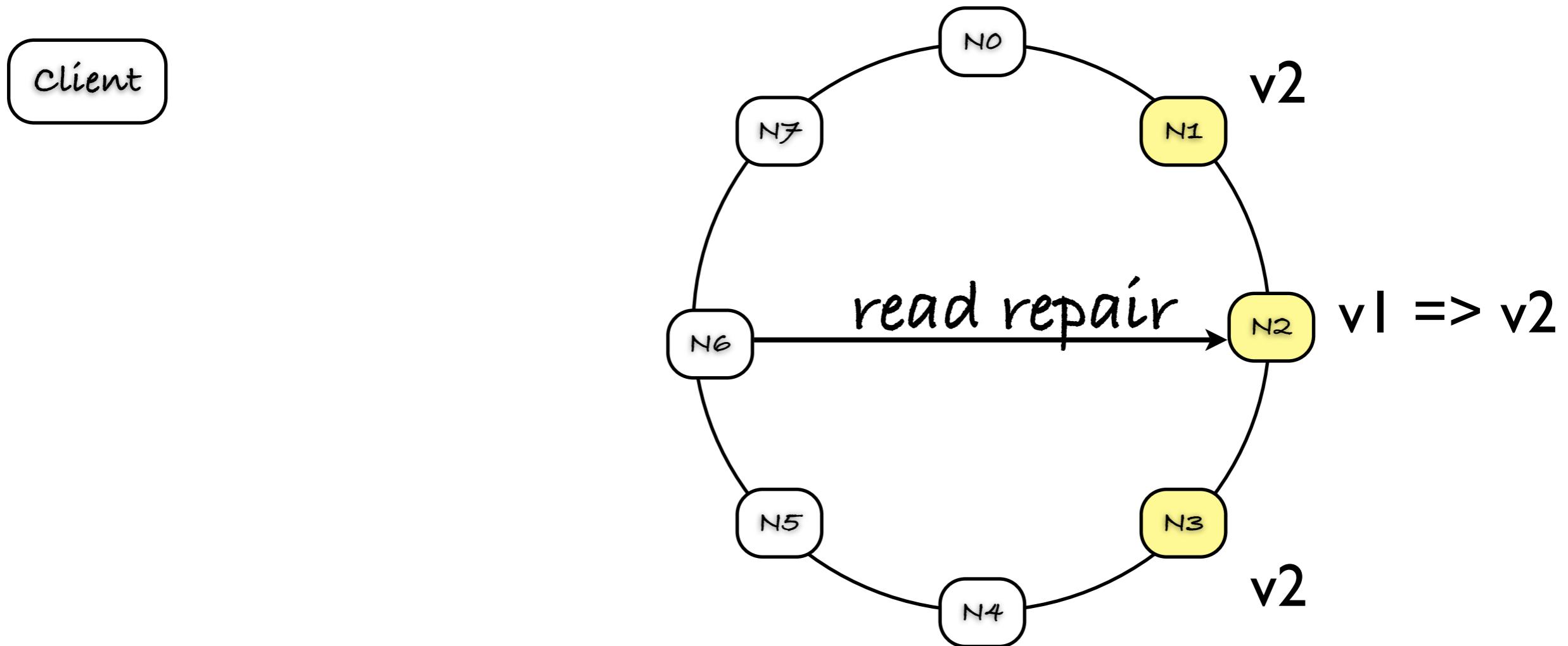
RF = 3

CL = QUORUM

Donnée la plus récente parmi 2 replicas



Exemple I (suite)

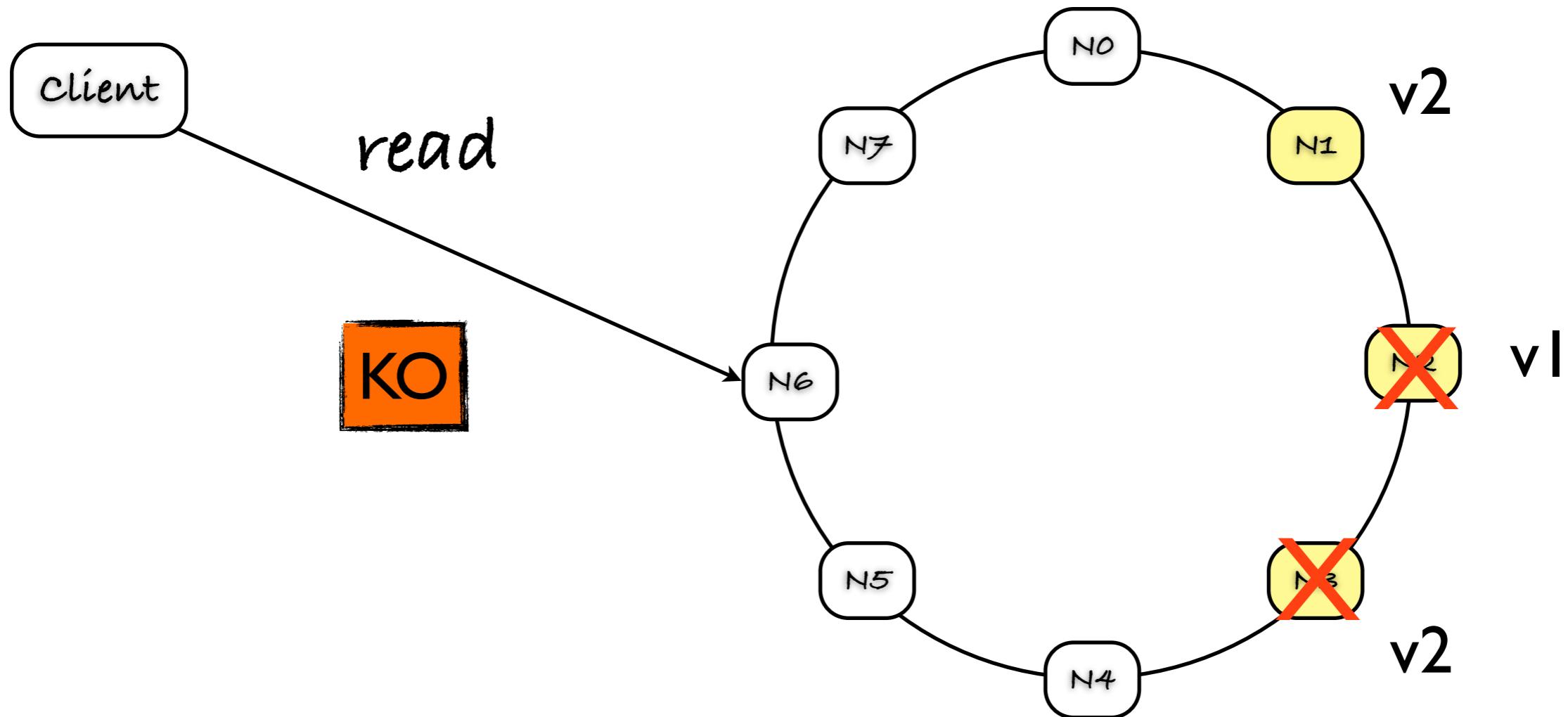


RF = 3

CL = QUORUM



Exemple 2



RF = 3

CL = QUORUM

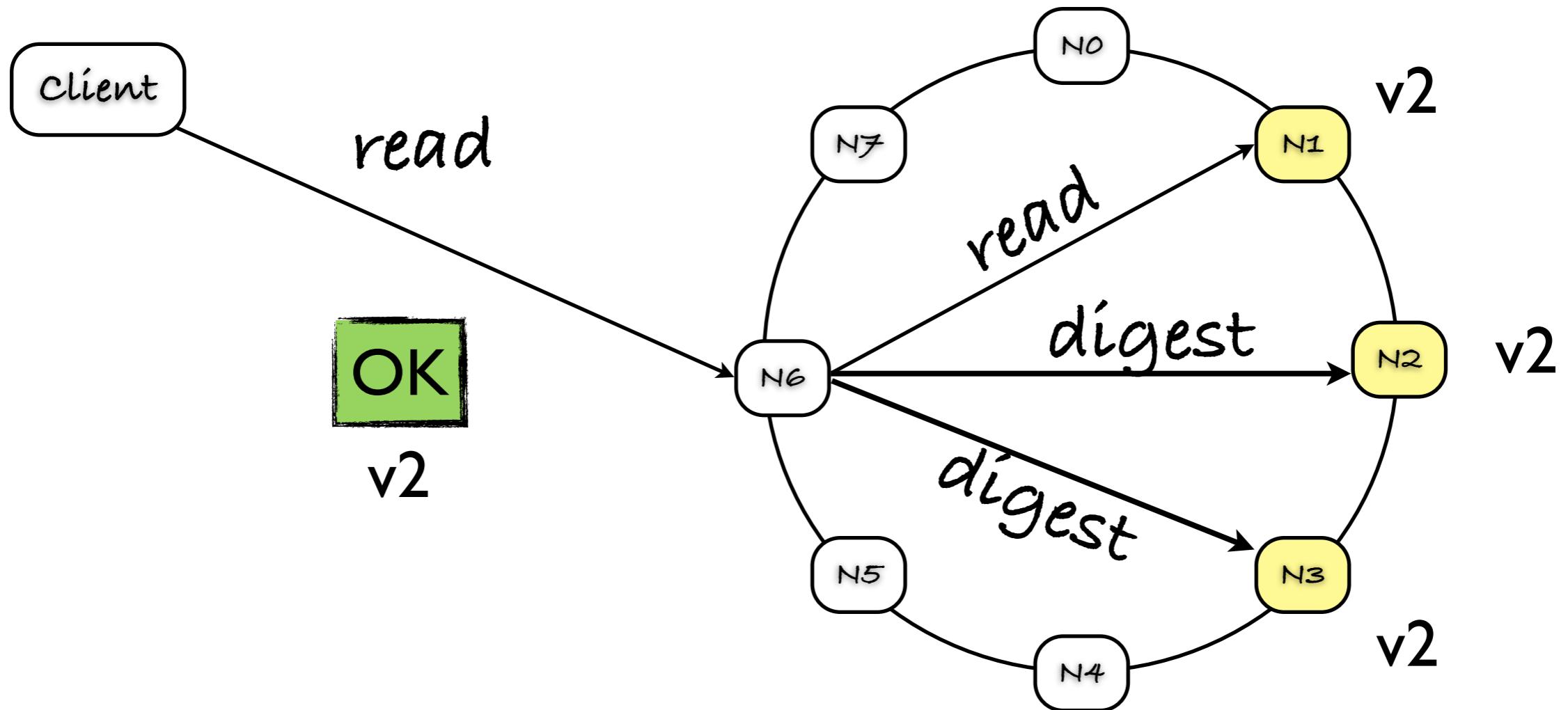


CL - ALL

lecture



Exemple I



RF = 3

CL = ALL

Donnée la plus récente parmi 3 replicas



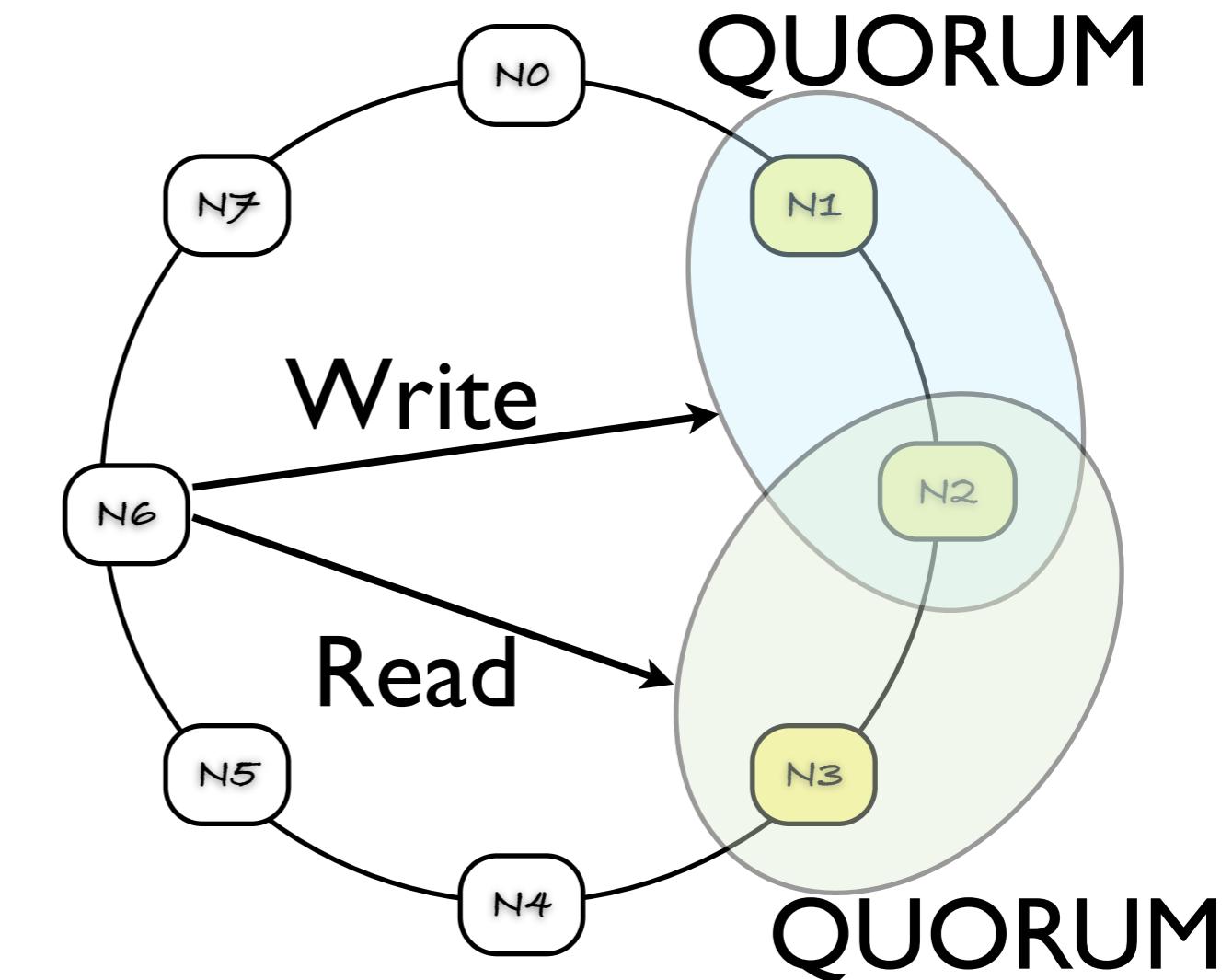
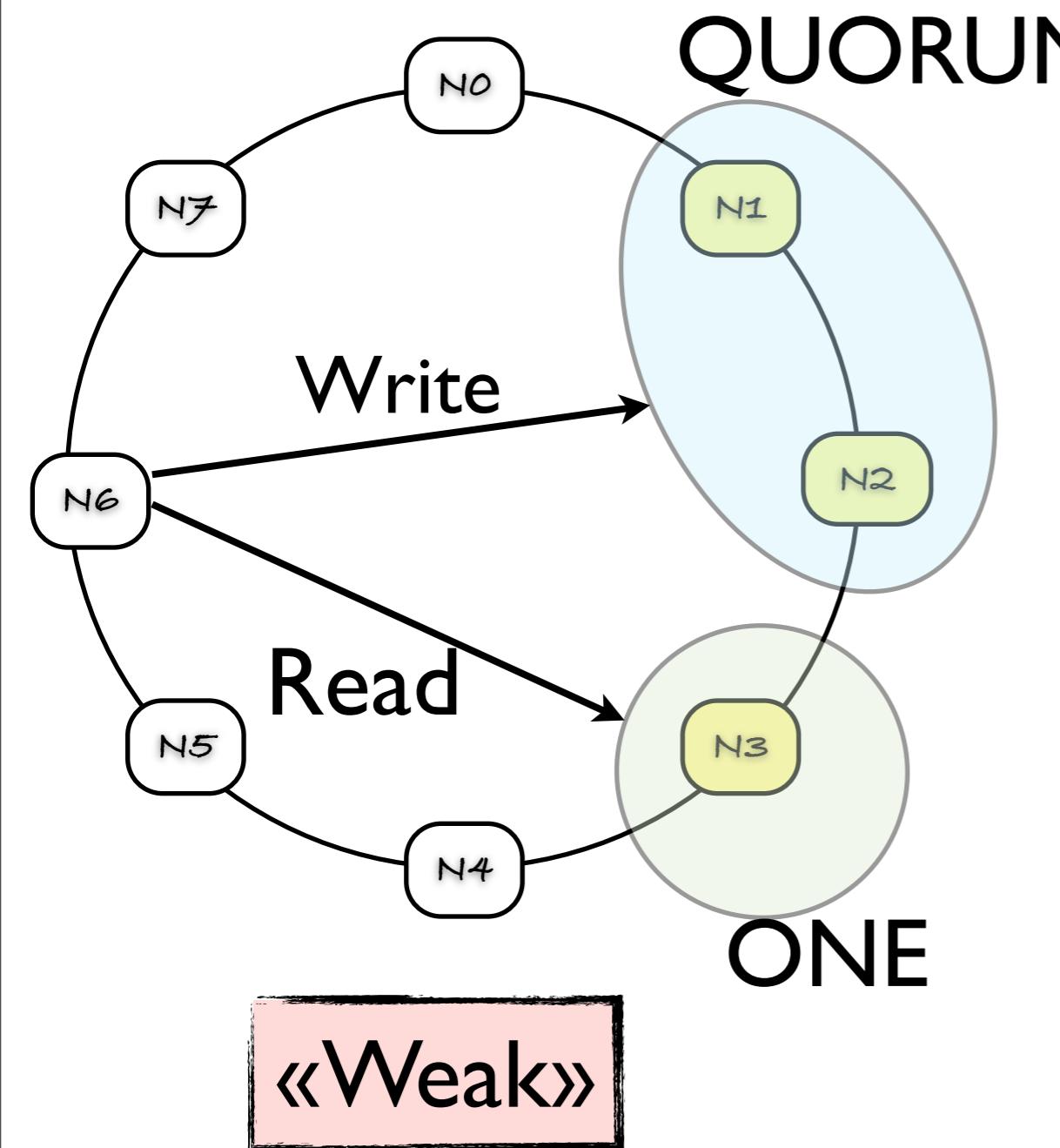
Consistance Forte

- Comment être certain que la donnée lue est la plus récente?
- formule $W + R > RF$
 - W: nb de replica contactées en écriture
 - R: nb de replica contactées en lecture
 - RF: replication factor
- Cas trivial: écriture en ALL et lecture en ONE



Consistency

$RF = 3$



«Strong»



Démo - calculer Token

```
import java.math.BigInteger;

public class SimpleTokenCalculator {

    public static void main(String args[]) {
        int nbNodes = Integer.parseInt(args[0]);

        for (int i = 0; i < nbNodes; i++) {
            BigInteger token = new BigInteger("2");
            token = token.pow(127);
            token = token.multiply(new BigInteger(" " + i));
            token = token.divide(new BigInteger(" " + nbNodes));
            System.out.println("token " + i + ": " + token);
        }
    }
}
```



Démo - création KS + CF

- create keyspace **ks** with placement_strategy = 'SimpleStrategy' and strategy_options = {replication_factor:2};
- use ks;
- create column family **user** with key_validation_class = 'AsciiType' and comparator = 'AsciiType' and default_validation_class = 'AsciiType';
- set user['antonio@devoxx.fr']['nom'] = 'Goncalves';
- get user['antonio@devoxx.fr'];



Questions?



Performance en écriture

- En écriture, Cassandra est extrêmement rapide, pourquoi?



Performance en écriture

Mémoire

MemTable

put

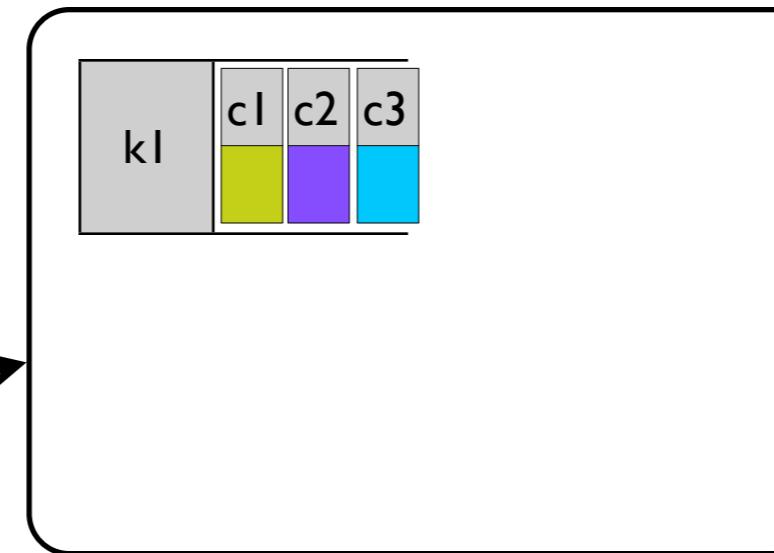
append

Commit
Log

Fichier

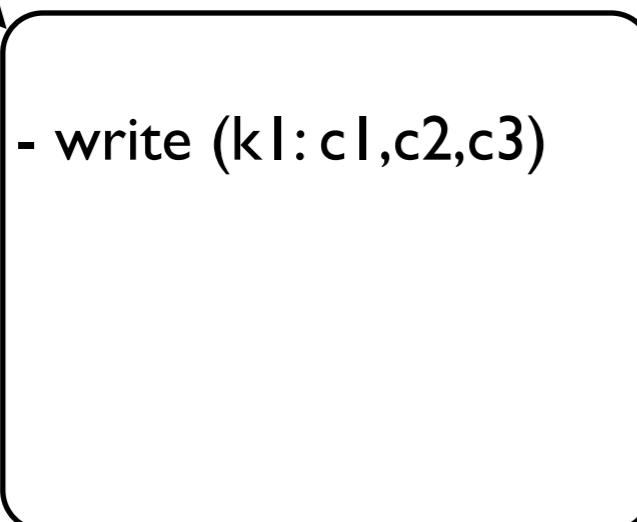


Mémoire



put

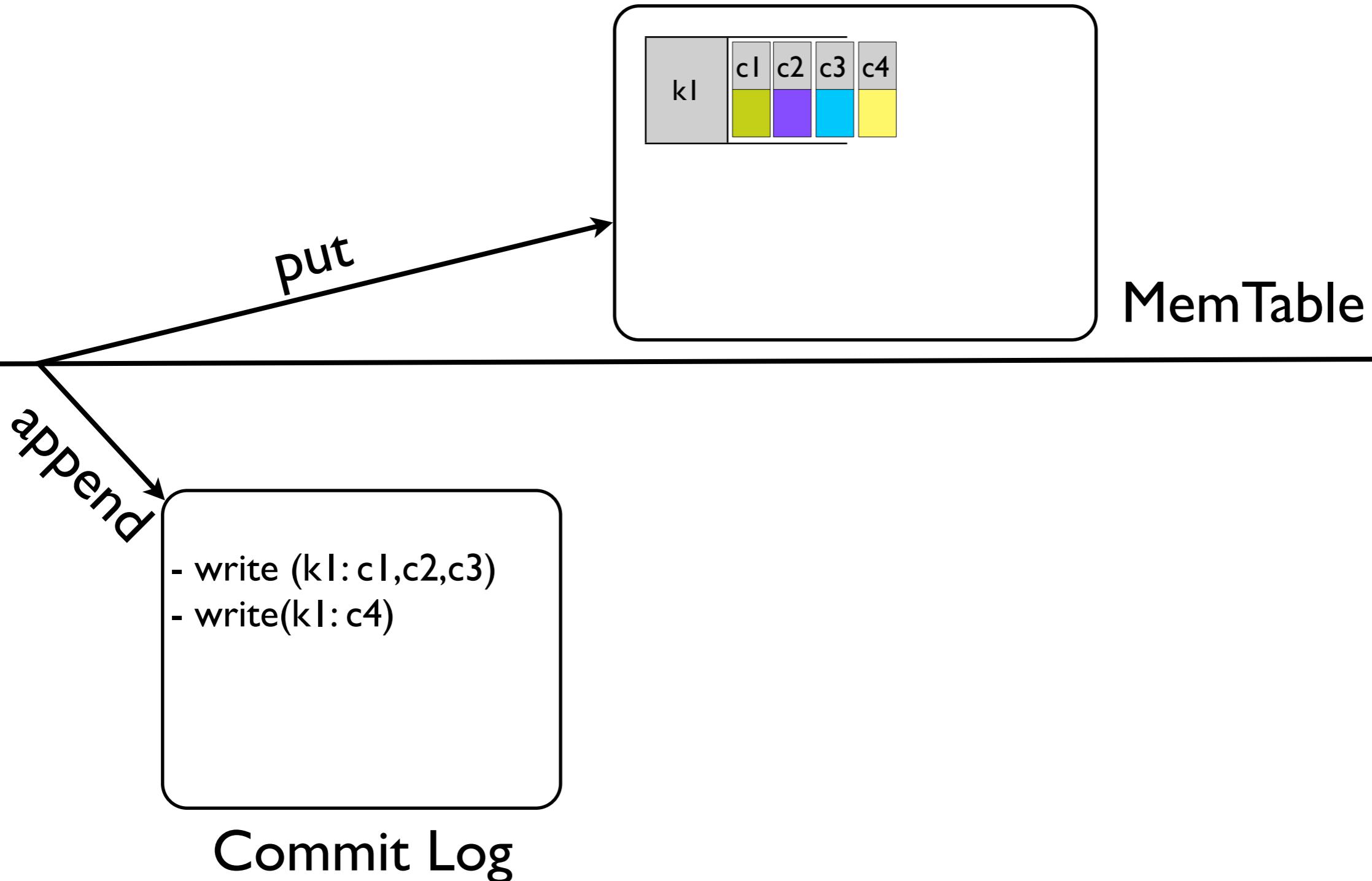
append

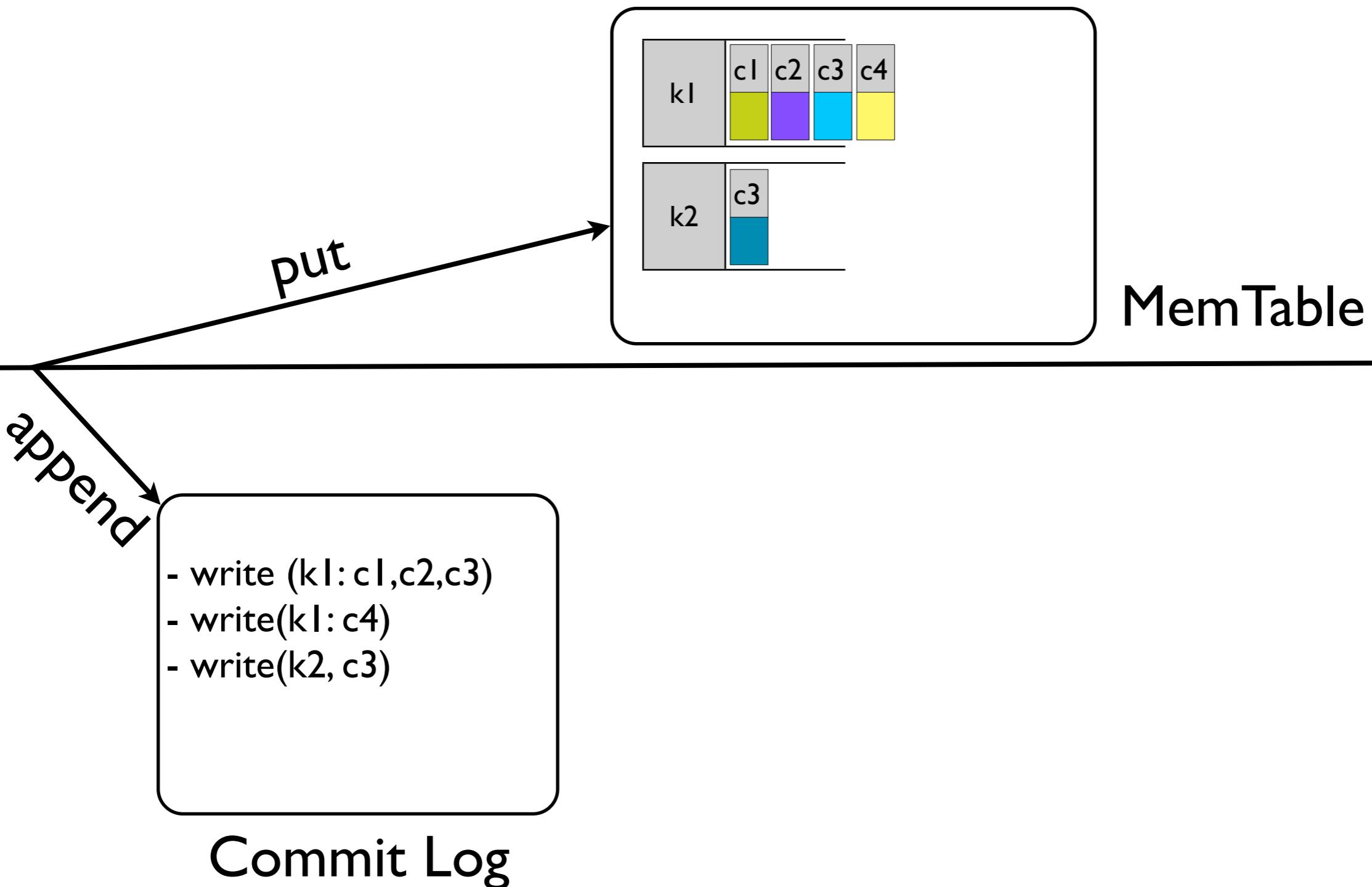


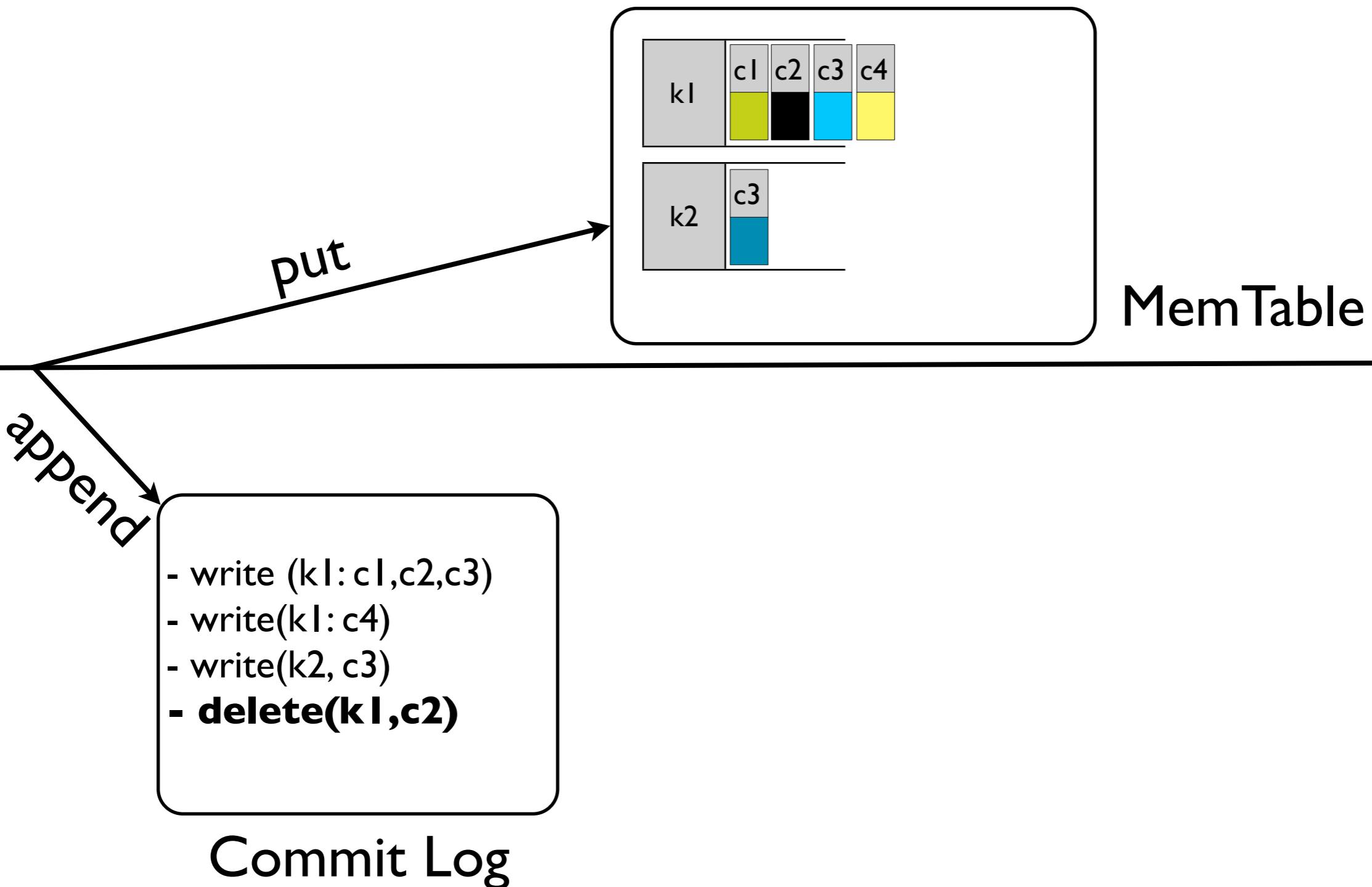
MemTable

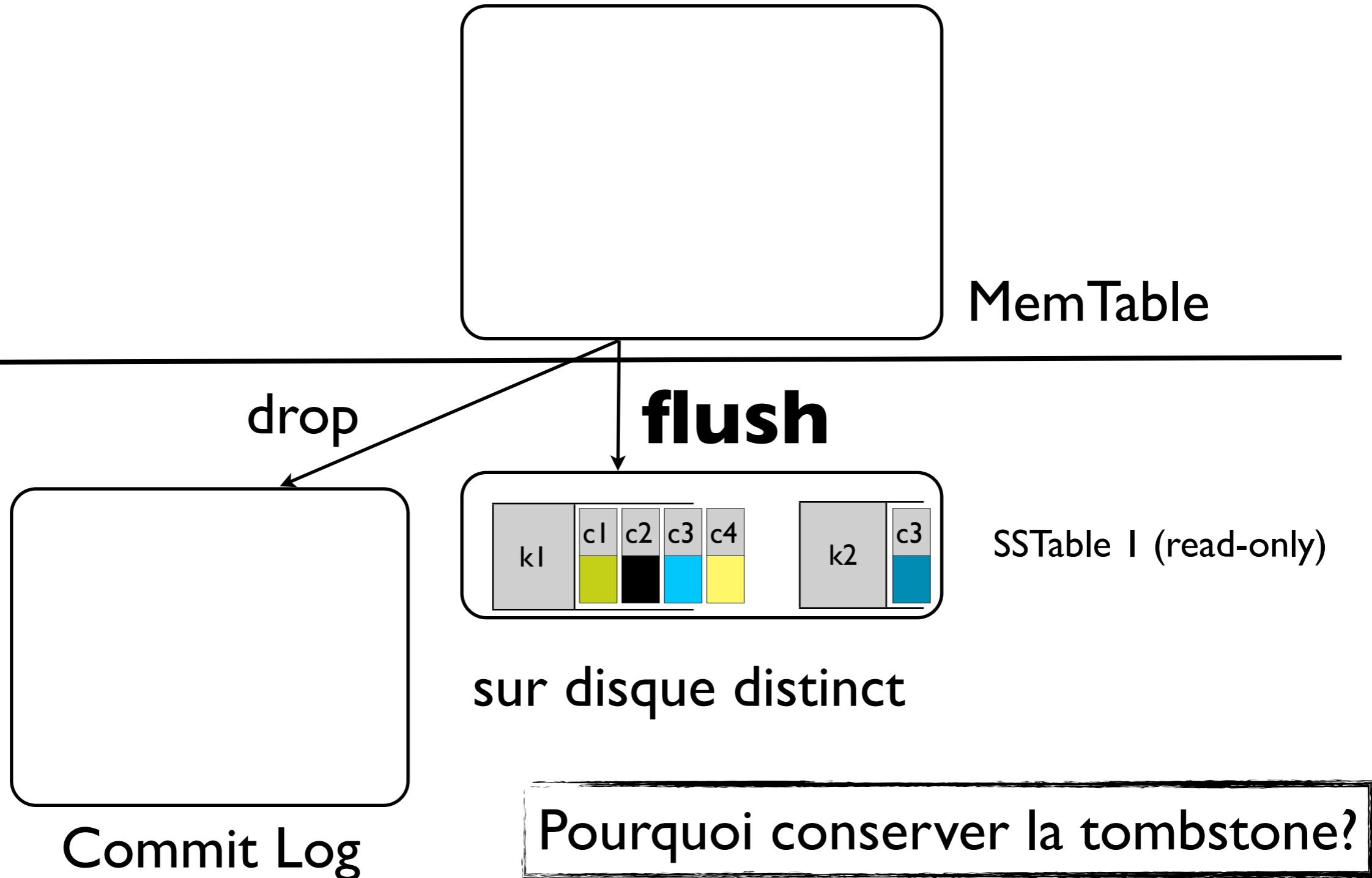
Commit Log







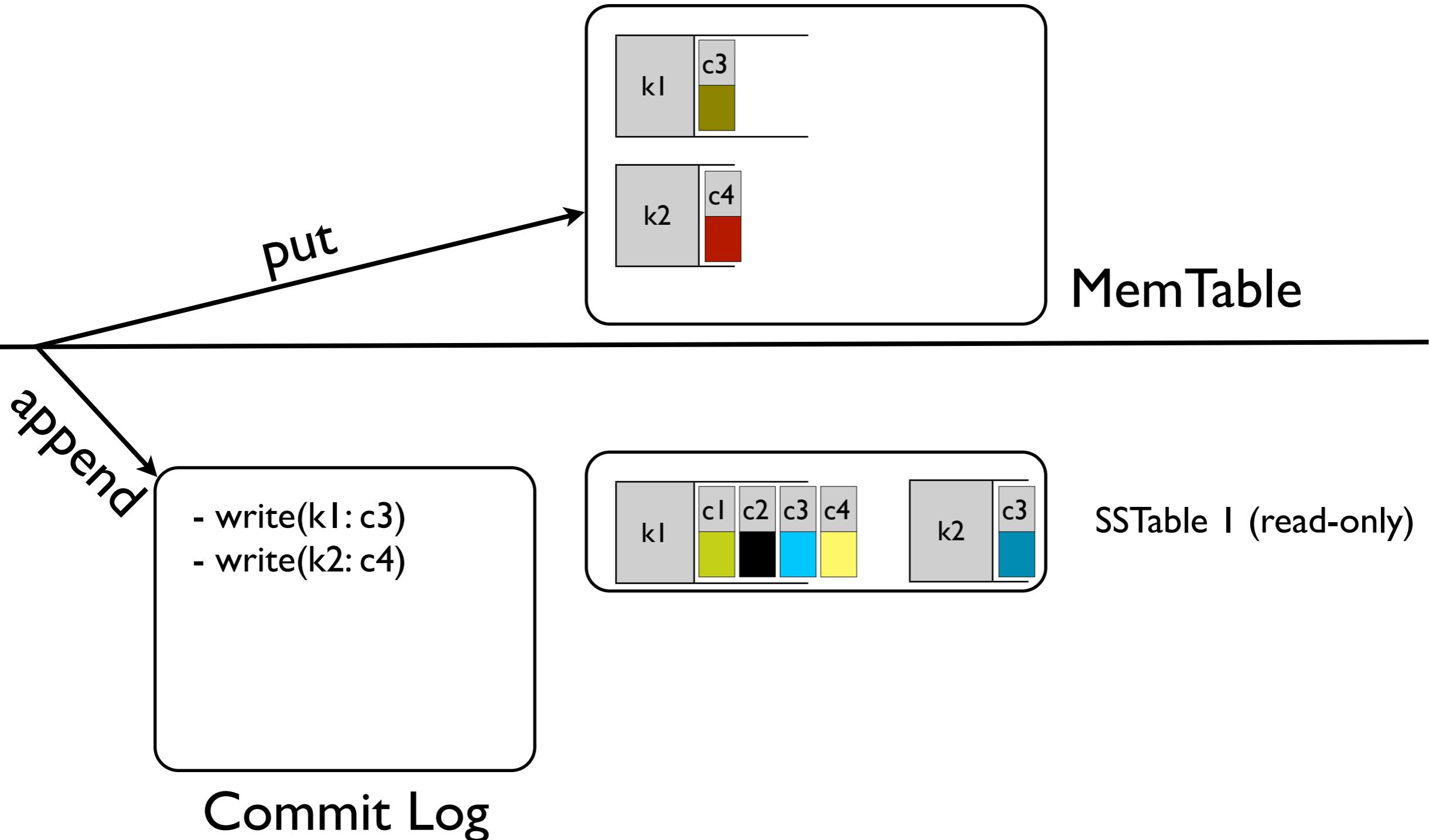




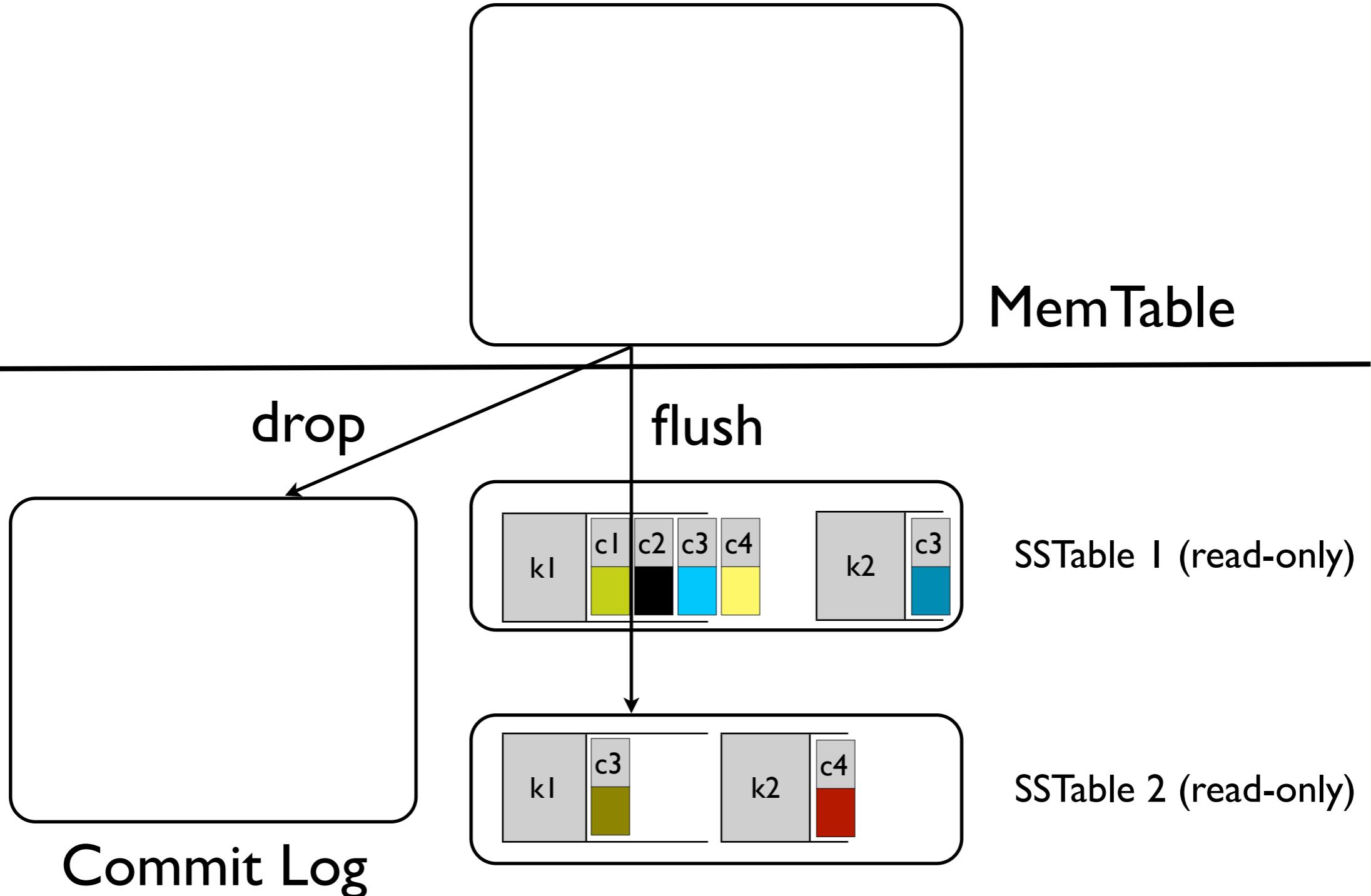
Delete distribué

- Grace period: 10 jours par défaut





Pas d'update au flush

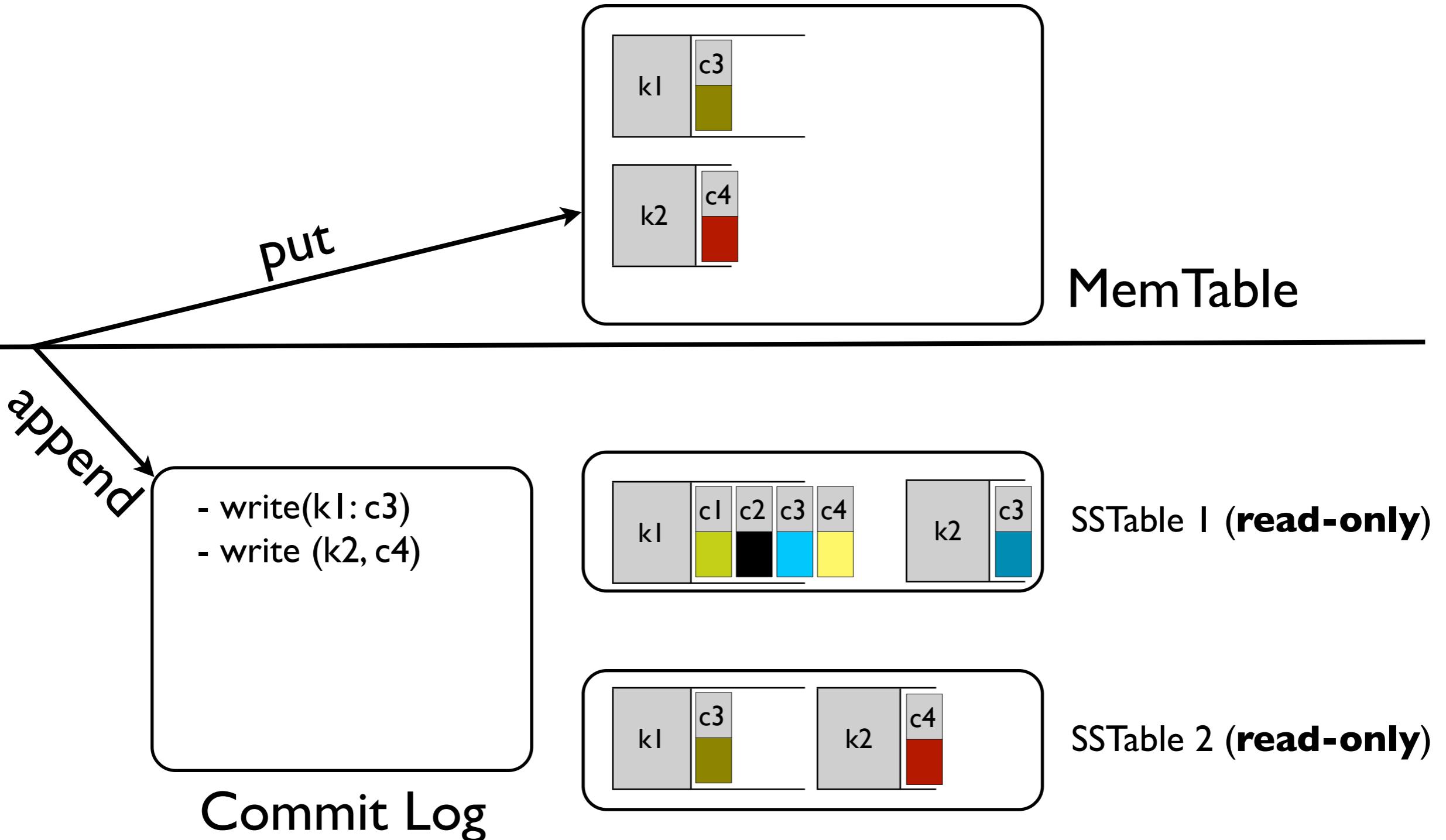


Backups?

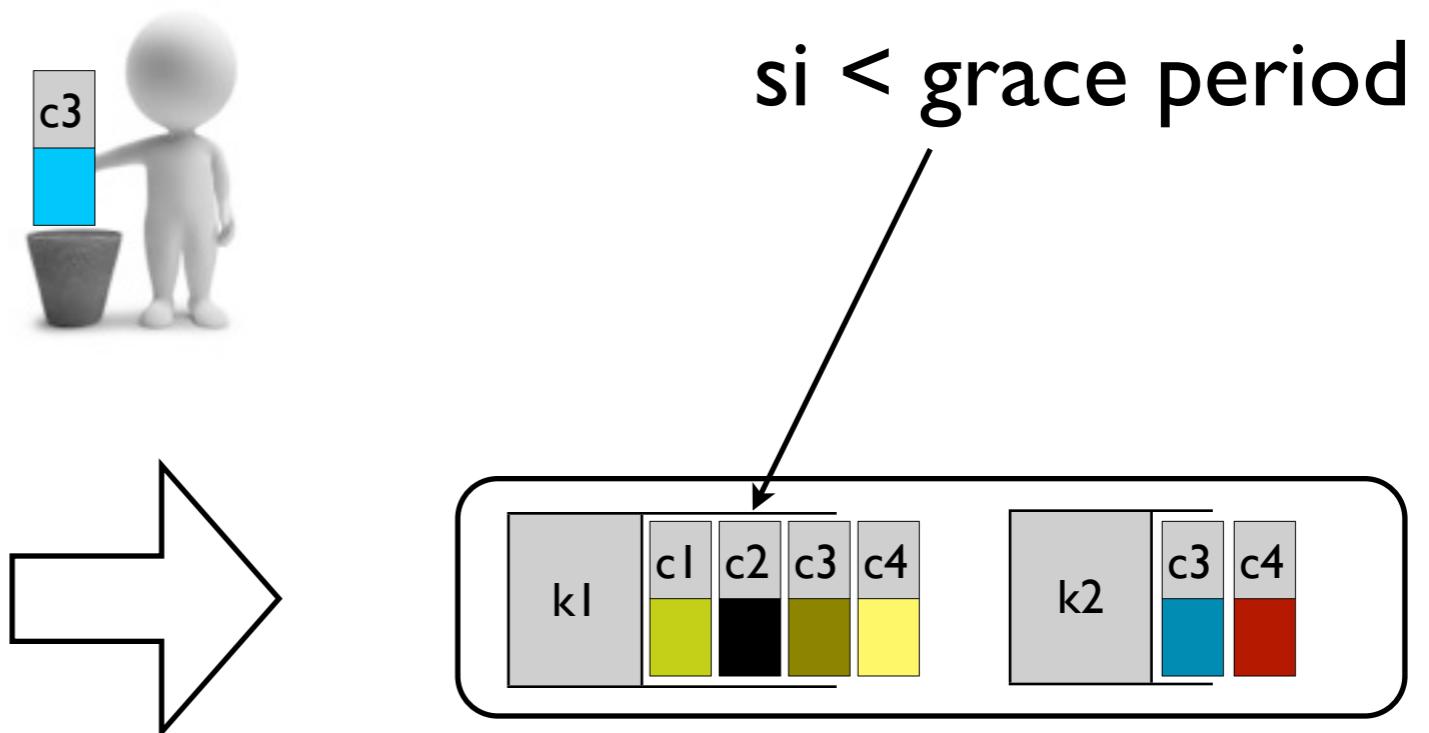
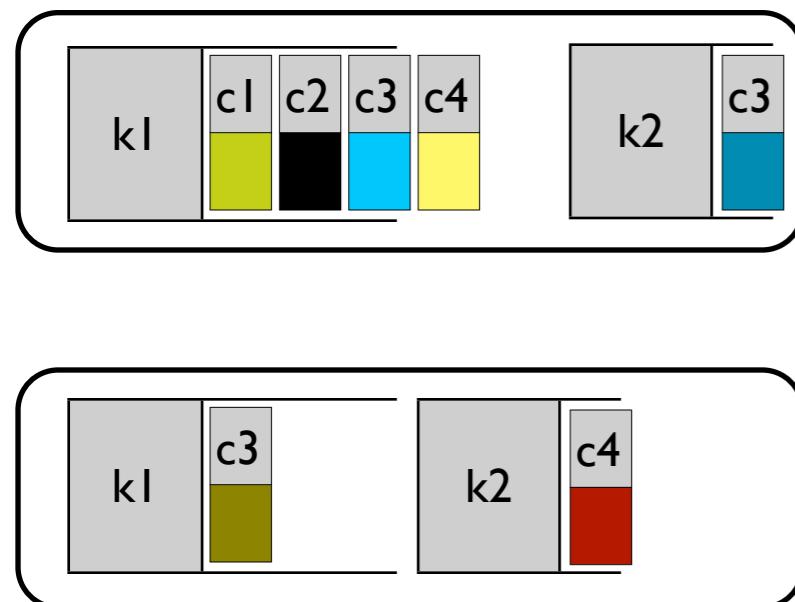
- SSTables sont readonly....



Quid de la lecture?



Compaction

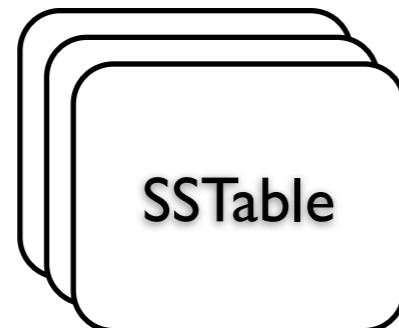


Lecture

MemTable

Mémoire

Fichiers



SSTable

6- Contient la valeur...



Lecture

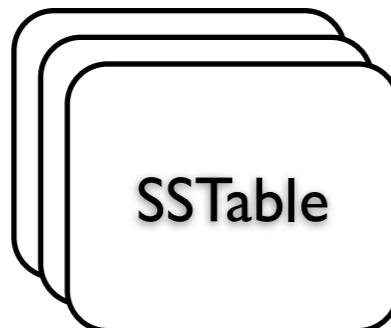
MemTable

I- Pertinent lorsque le client spécifie le nom de la colonne

Row Cache

Mémoire

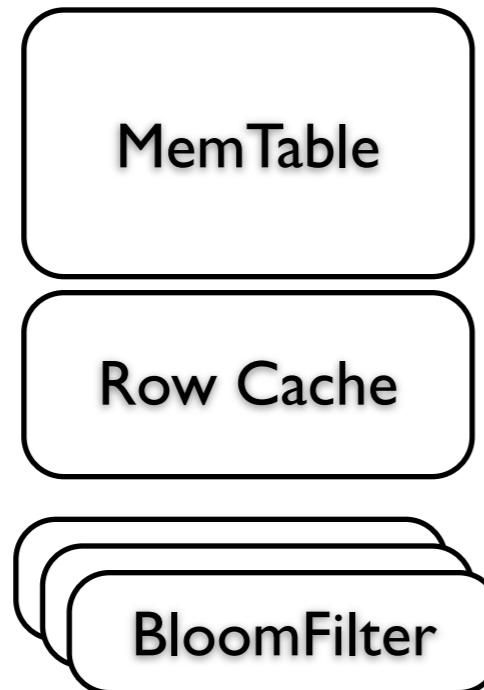
Fichiers



6- Contient la valeur...



Lecture

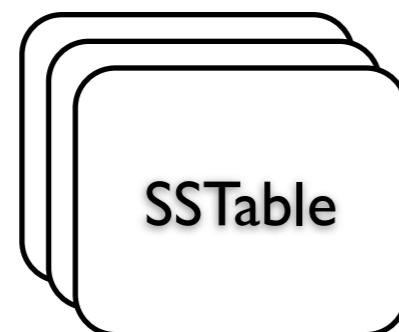


1- Pertinent lorsque le client spécifie le nom de la colonne

2- Stocke l'intégralité d'une row

Mémoire

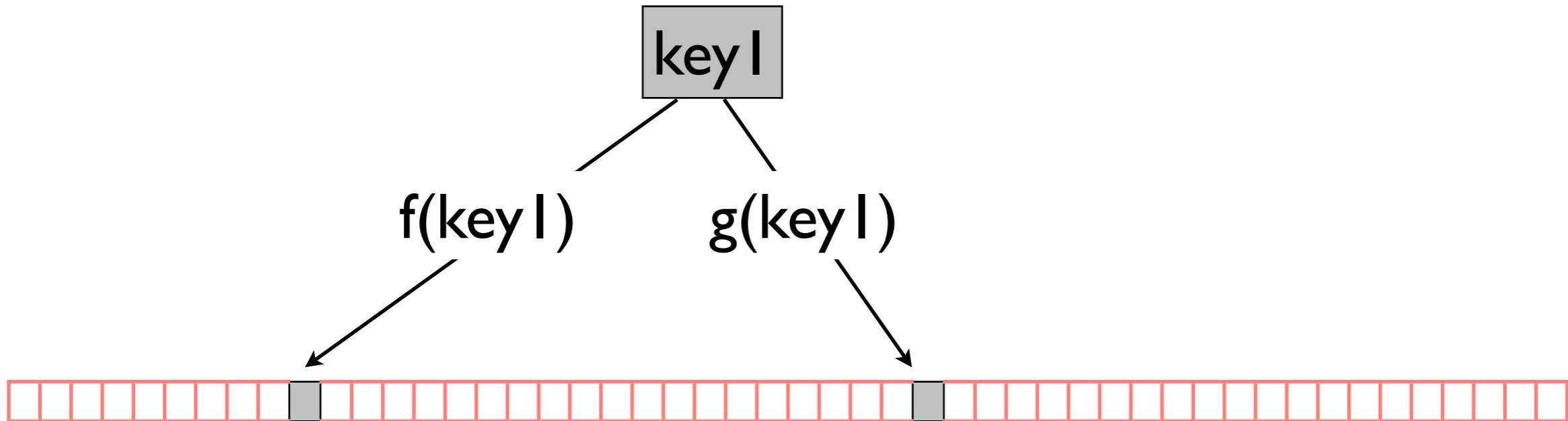
Fichiers



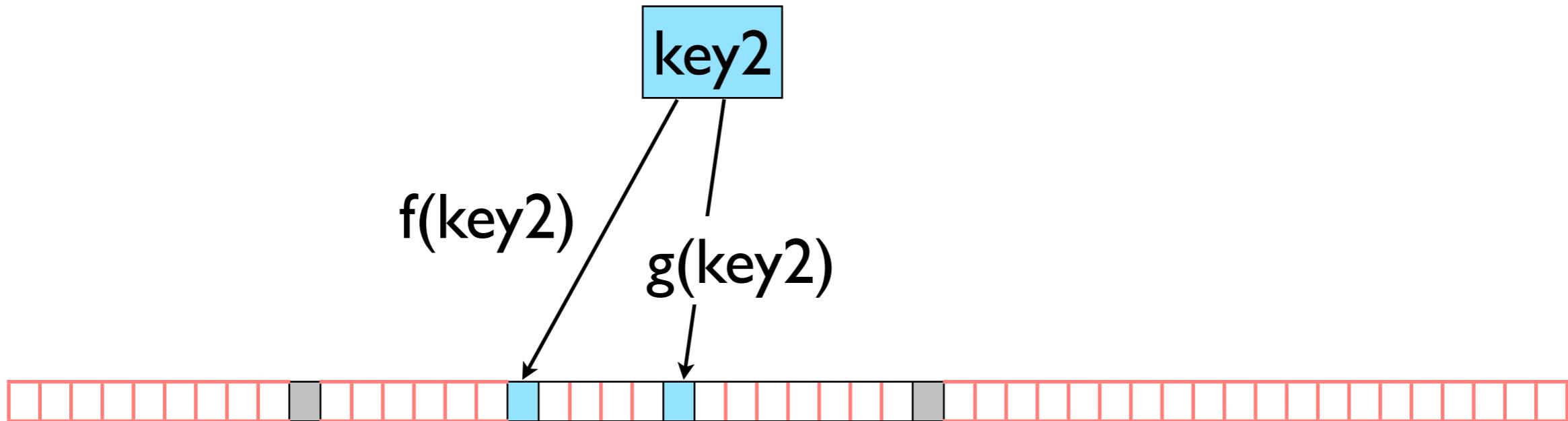
6- Contient la valeur...



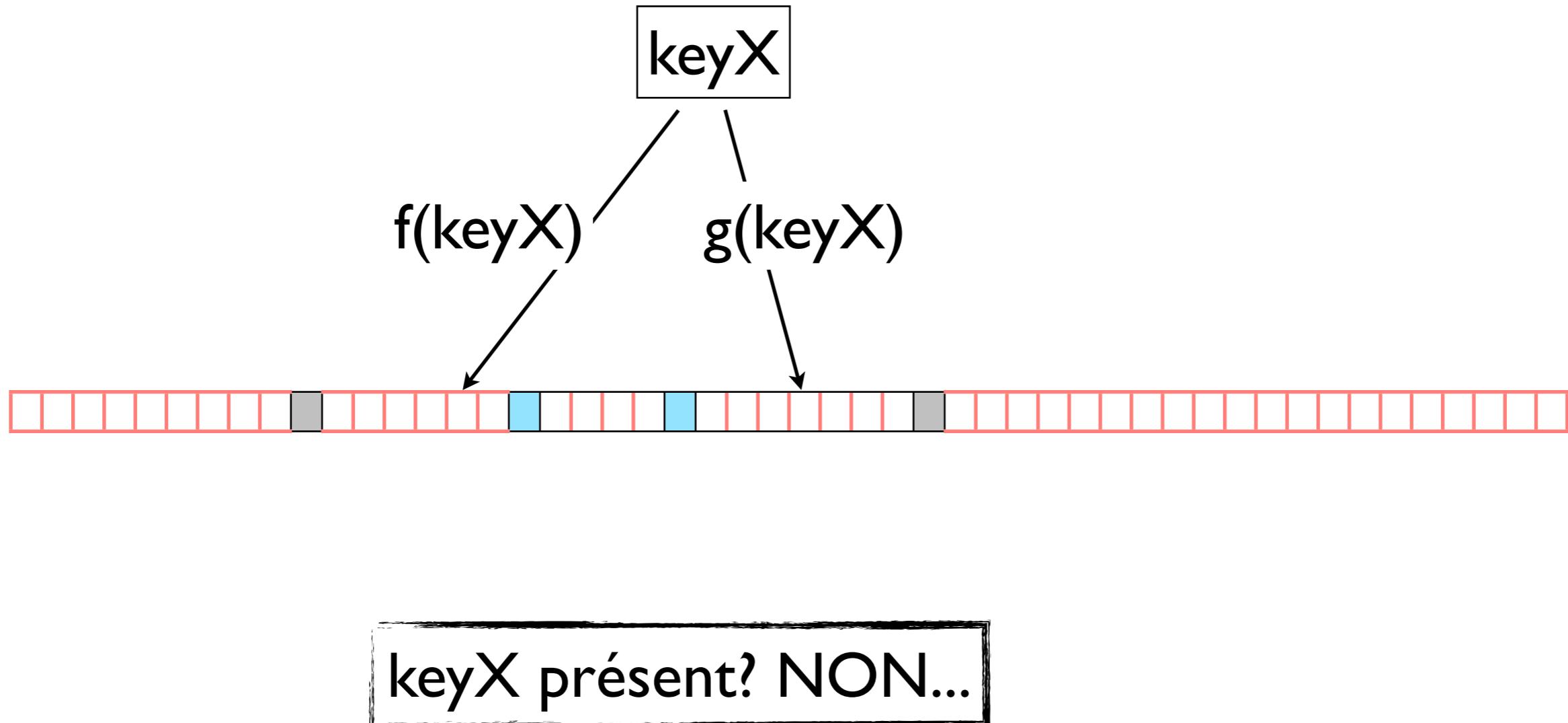
Bloom Filter (création)



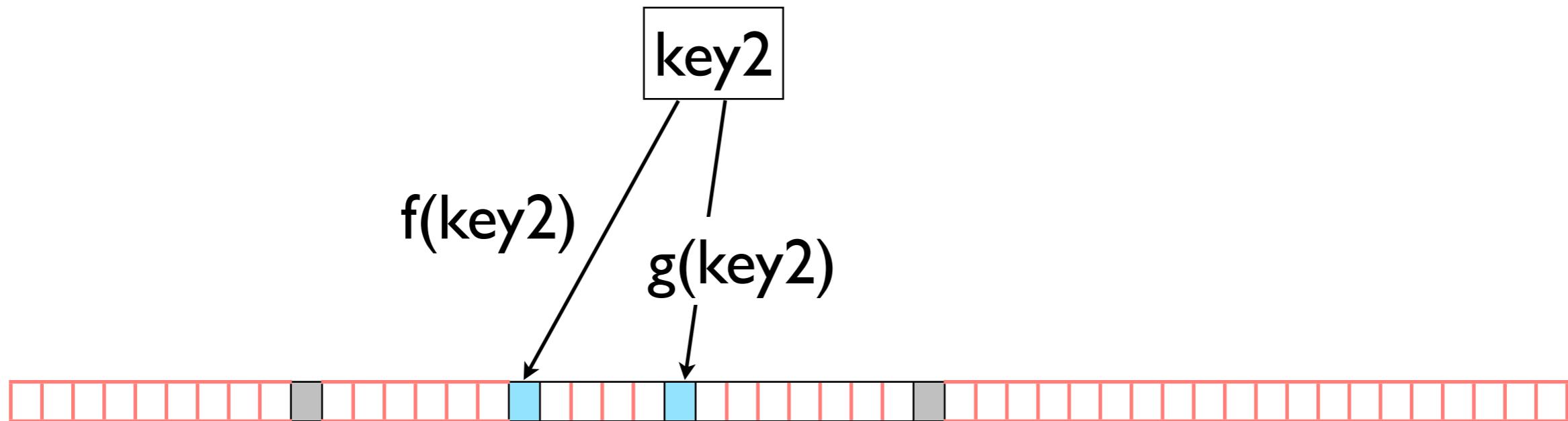
Bloom Filter (création)



Bloom Filter (lecture)



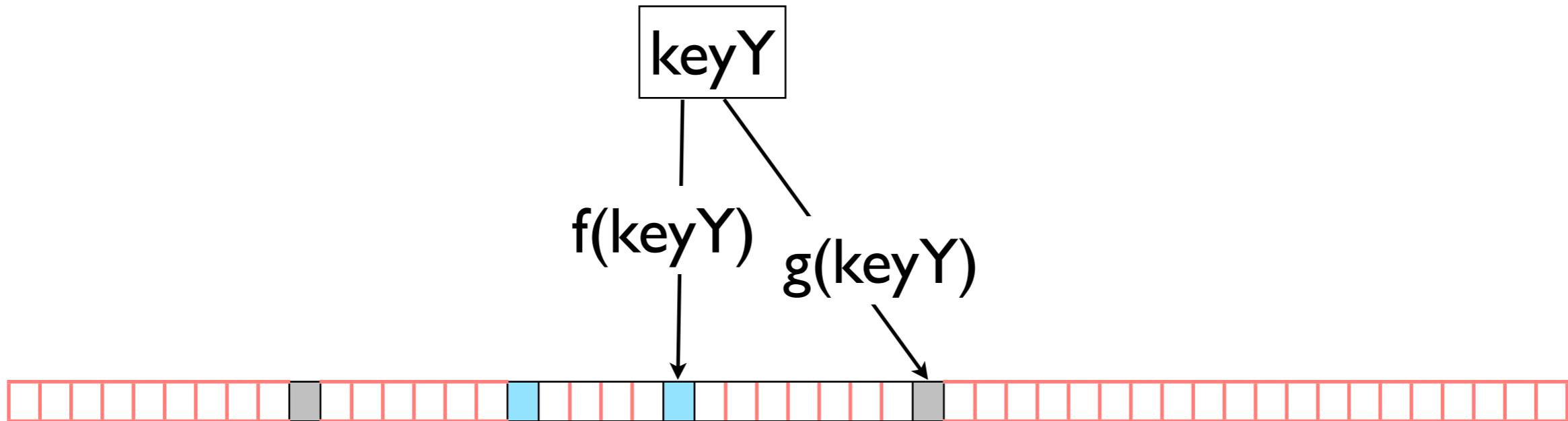
Bloom Filter (lecture)



key2 présent? OUI... vrai positif



Bloom Filter (lecture)



keyY présent? OUI... mais c'est un faux positif



Lecture

MemTable

1- Pertinent lorsque le client spécifie le nom de la colonne

Row Cache

2- Stocke l'intégralité d'une row

BloomFilter

3- répond à «La SSTable contient-elle un fragment de row ?»
Attention, faux positif possible.

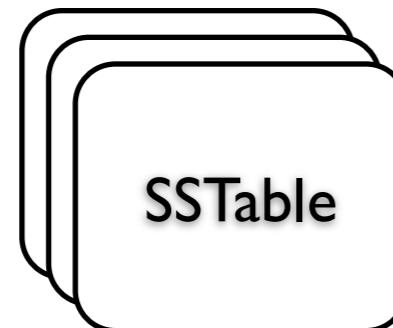
Key Cache

Mémoire

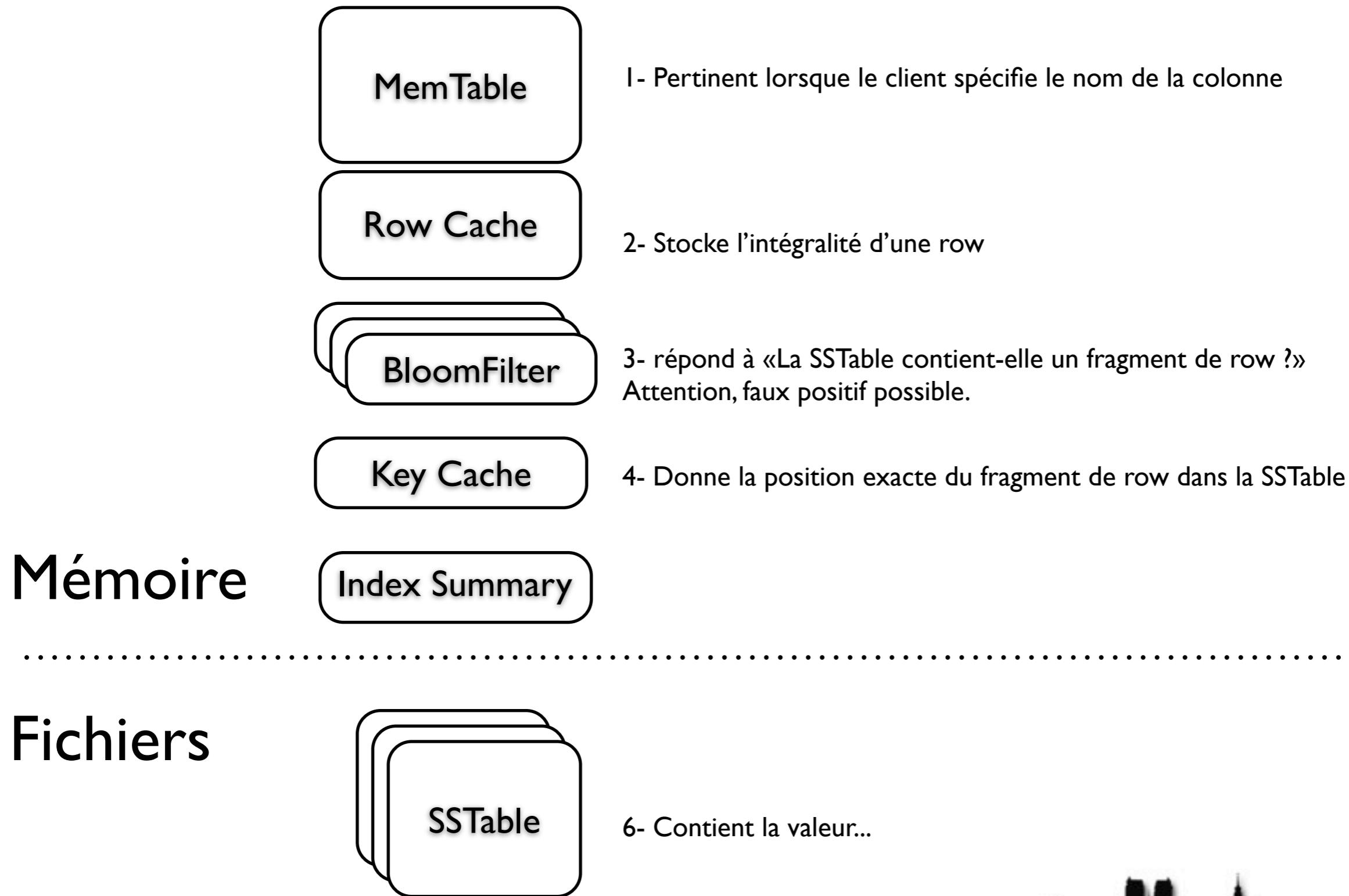
Fichiers

SSTable

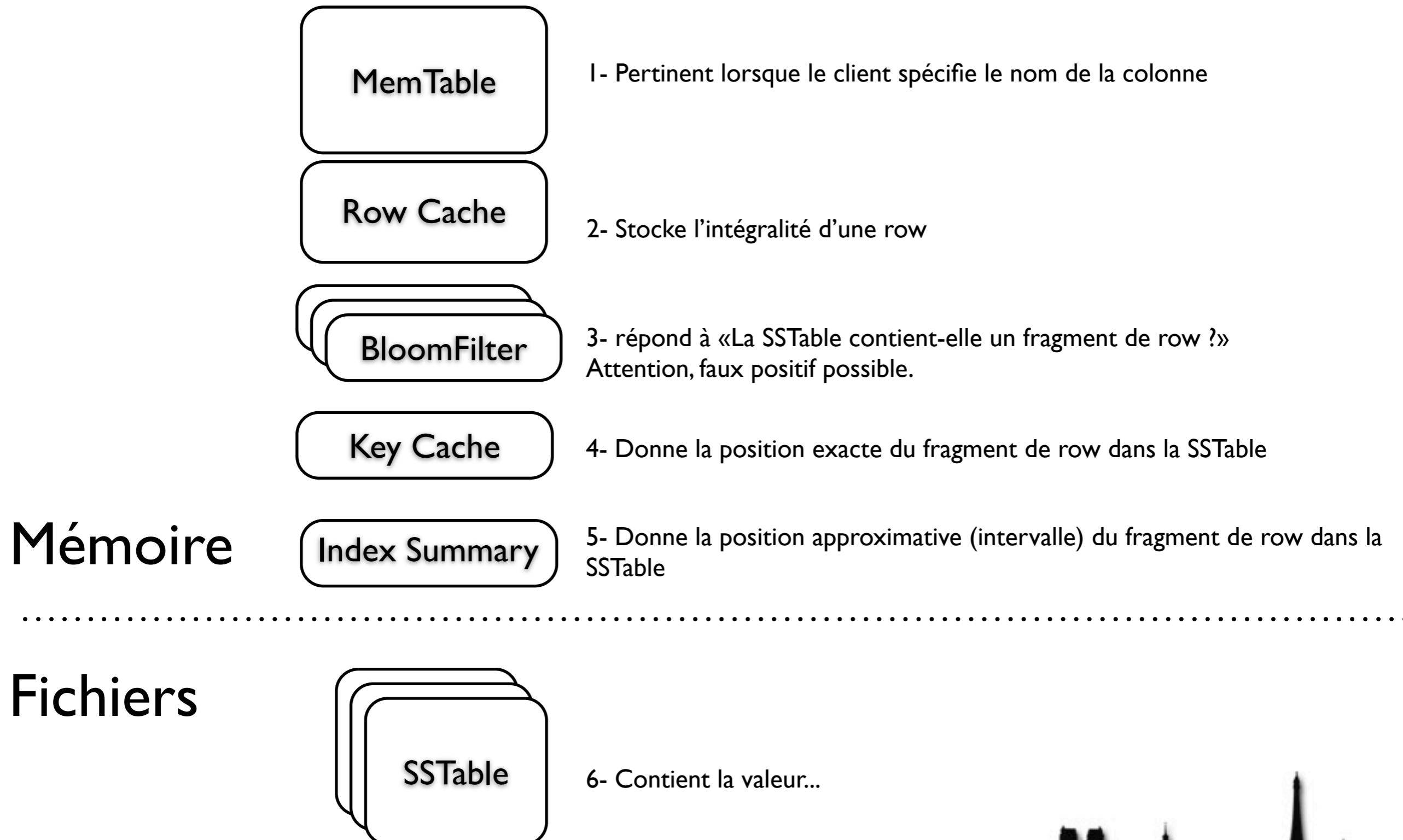
6- Contient la valeur...



Lecture



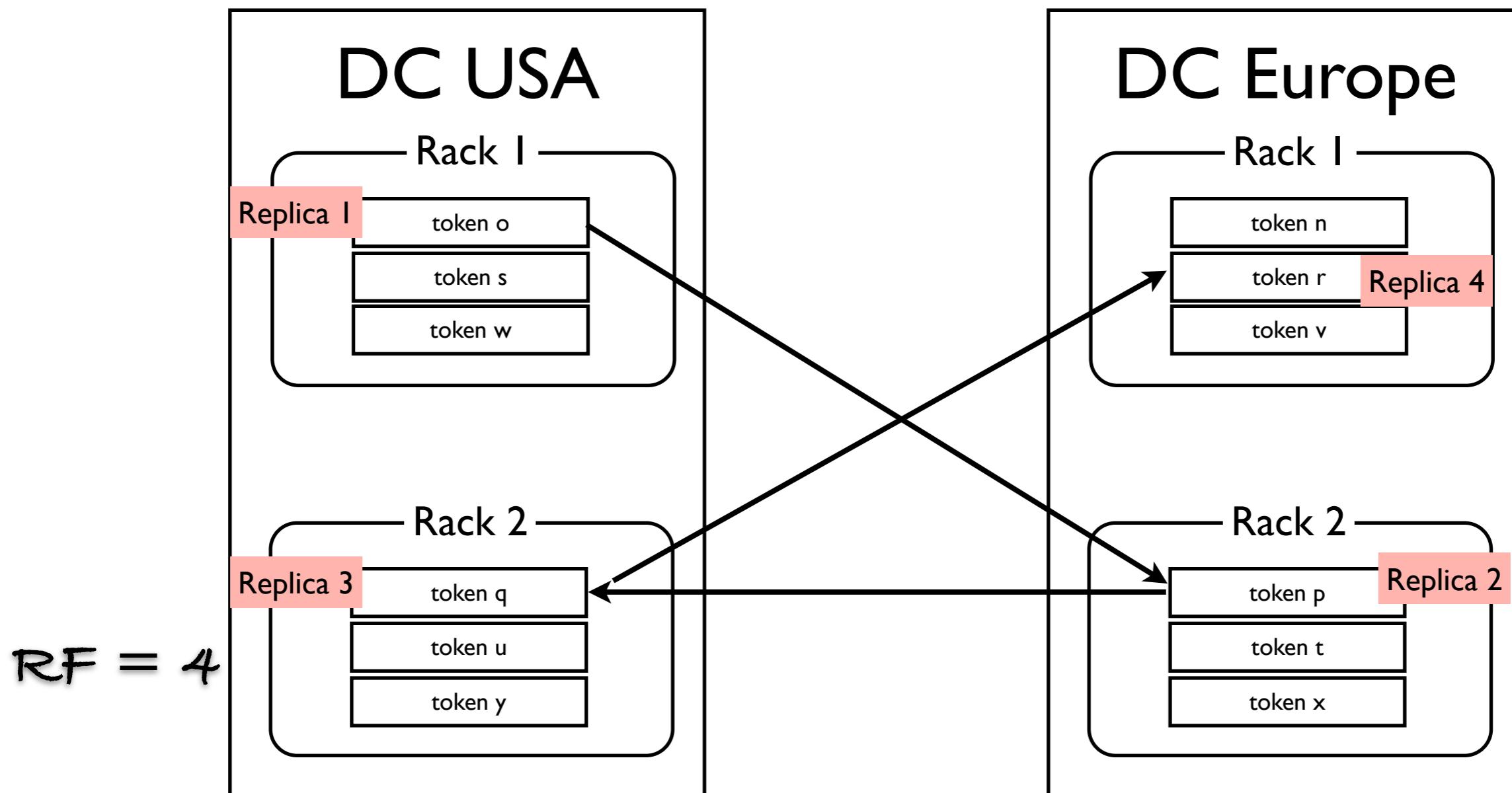
Lecture



Questions?



Multi Data Center



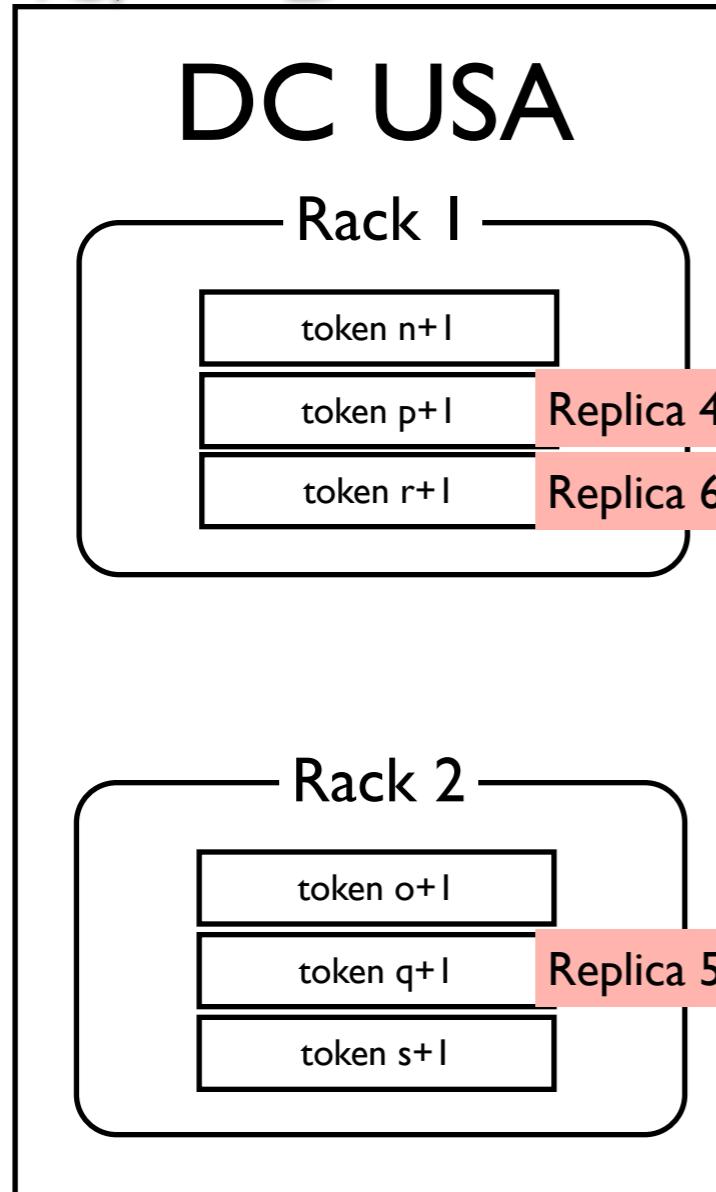
Pas de symétrie entre les DC
QUORUM implique nécessairement 2 DC => latence



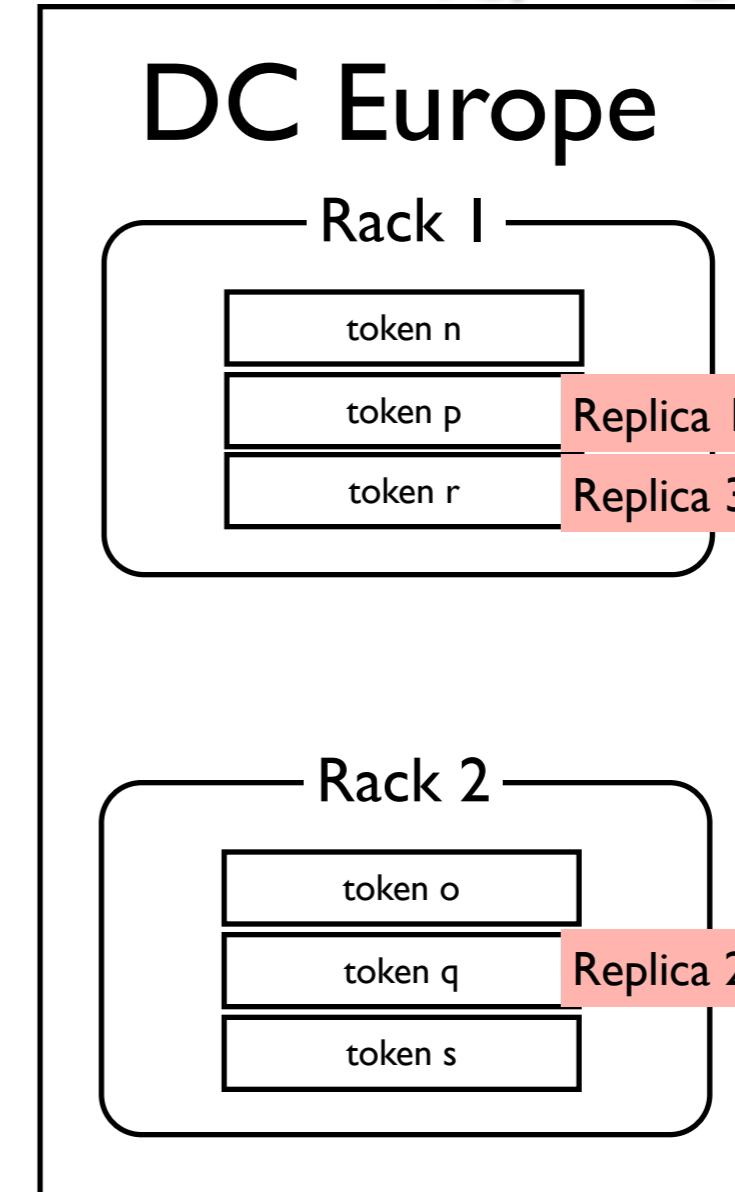
DataCenter «Aware»

NetworkTopologyStrategy + alternance tokens

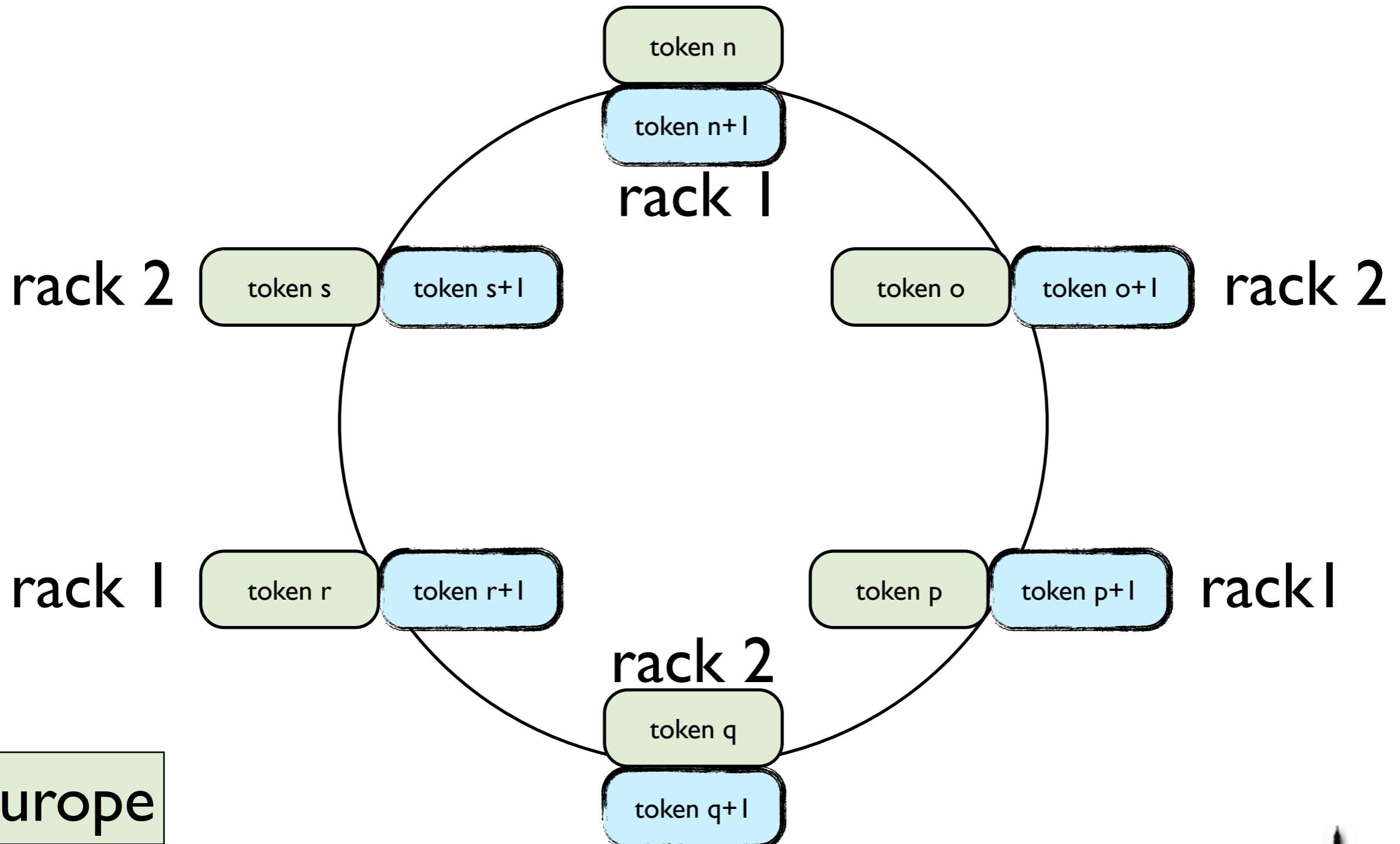
$RF = 3$



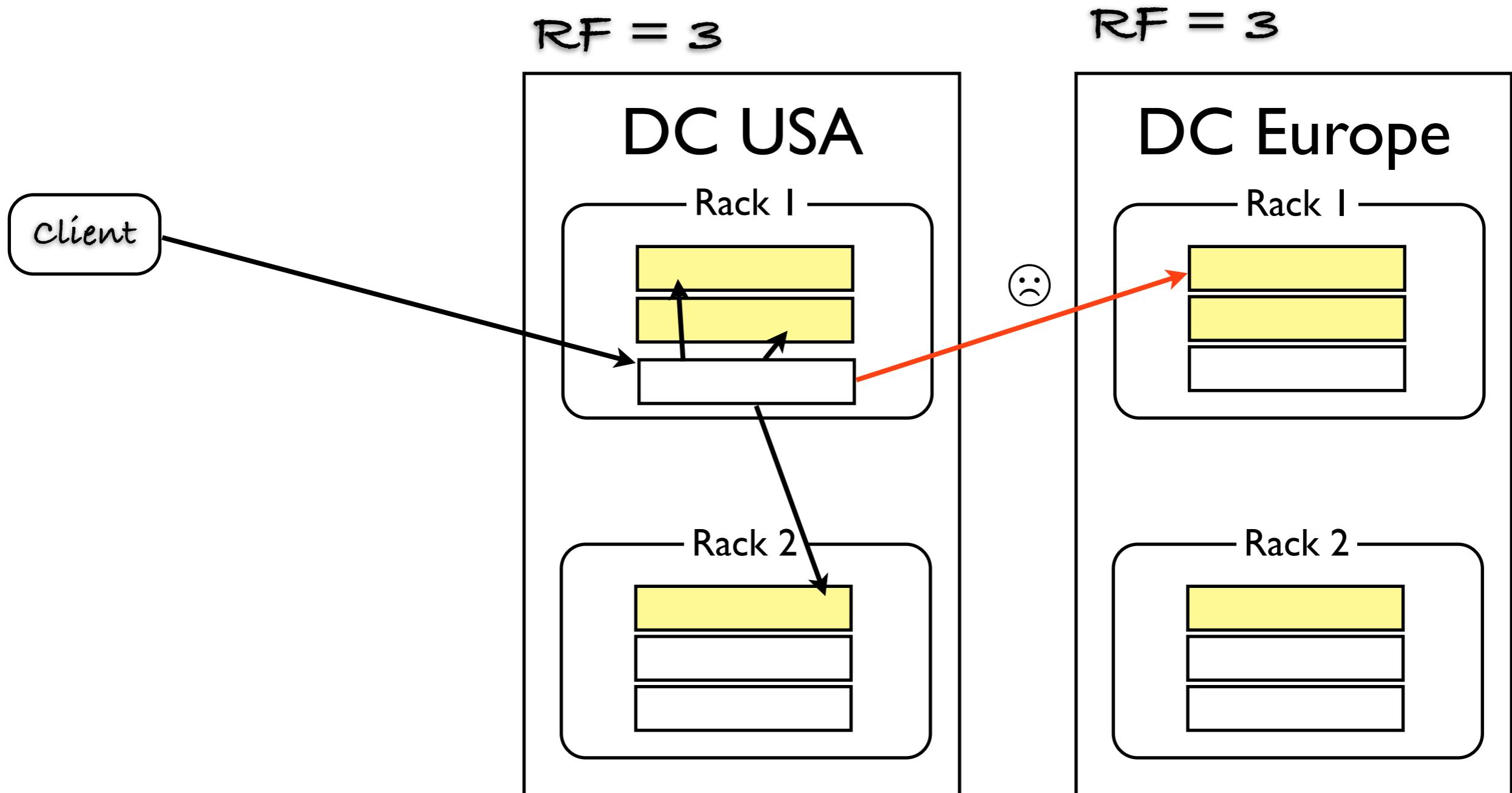
$RF = 3$



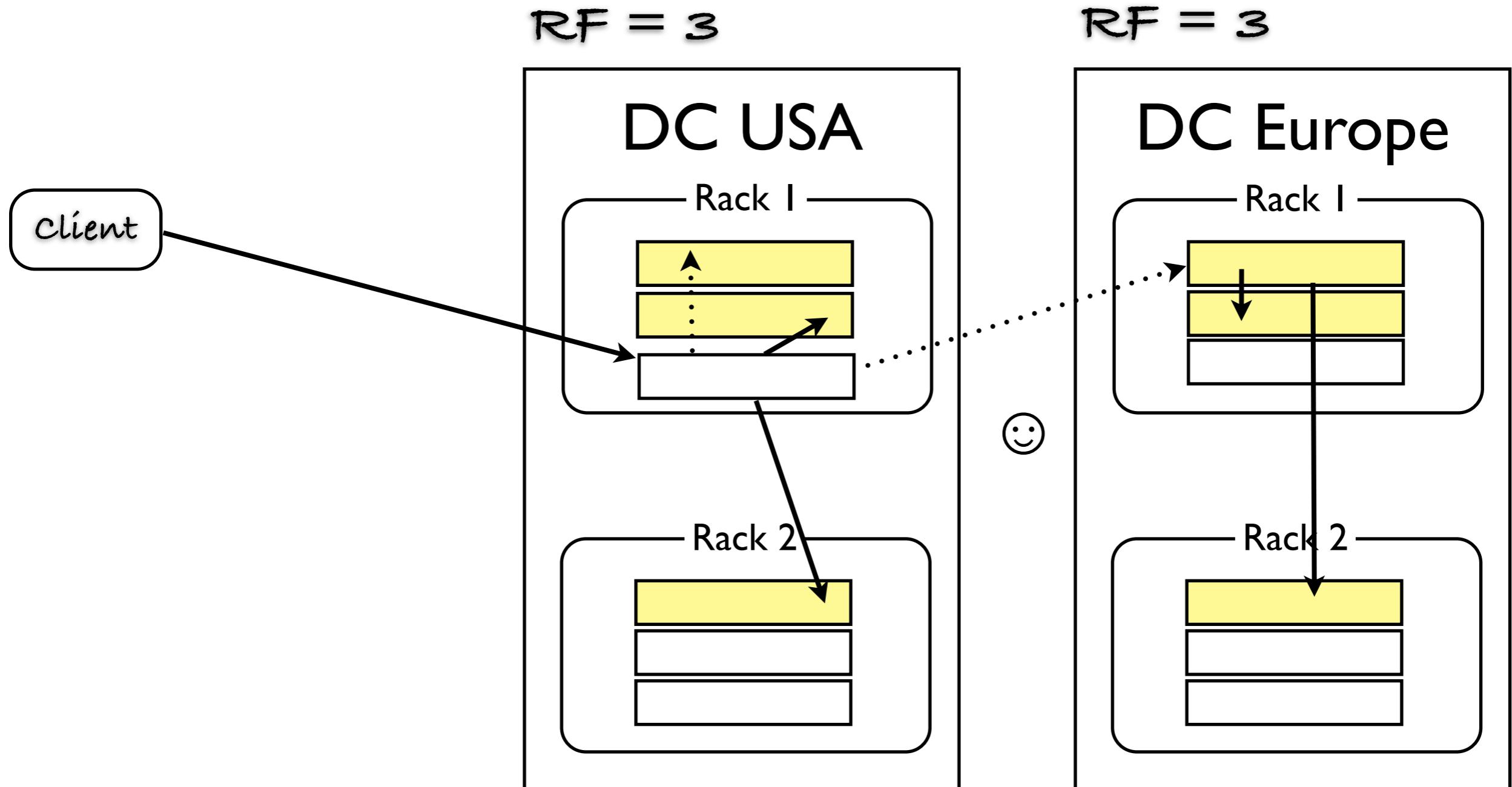
Ring 2 DataCenter (+I)



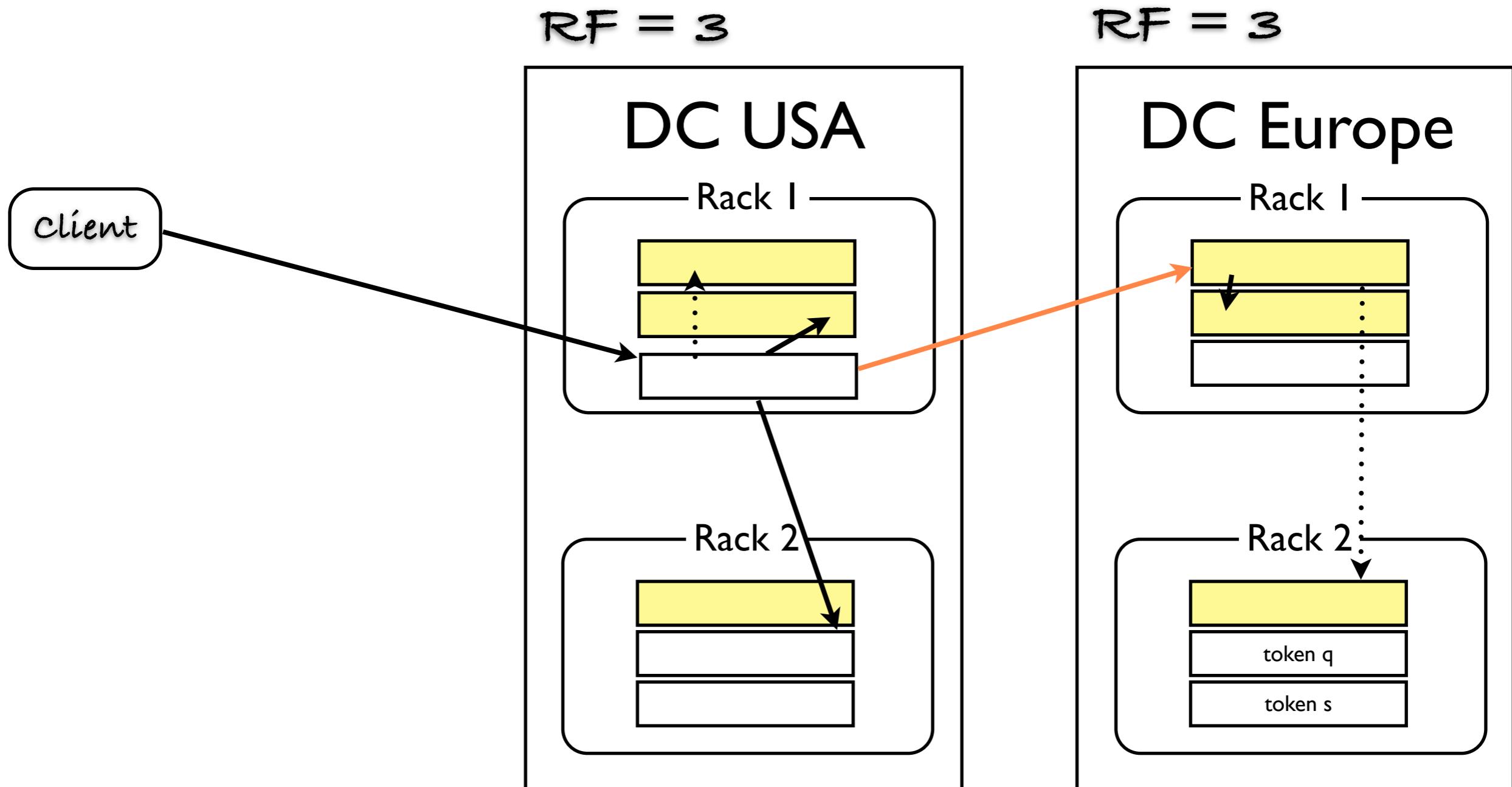
Latence du QUORUM



LOCAL_QUORUM



EACH_QUORUM

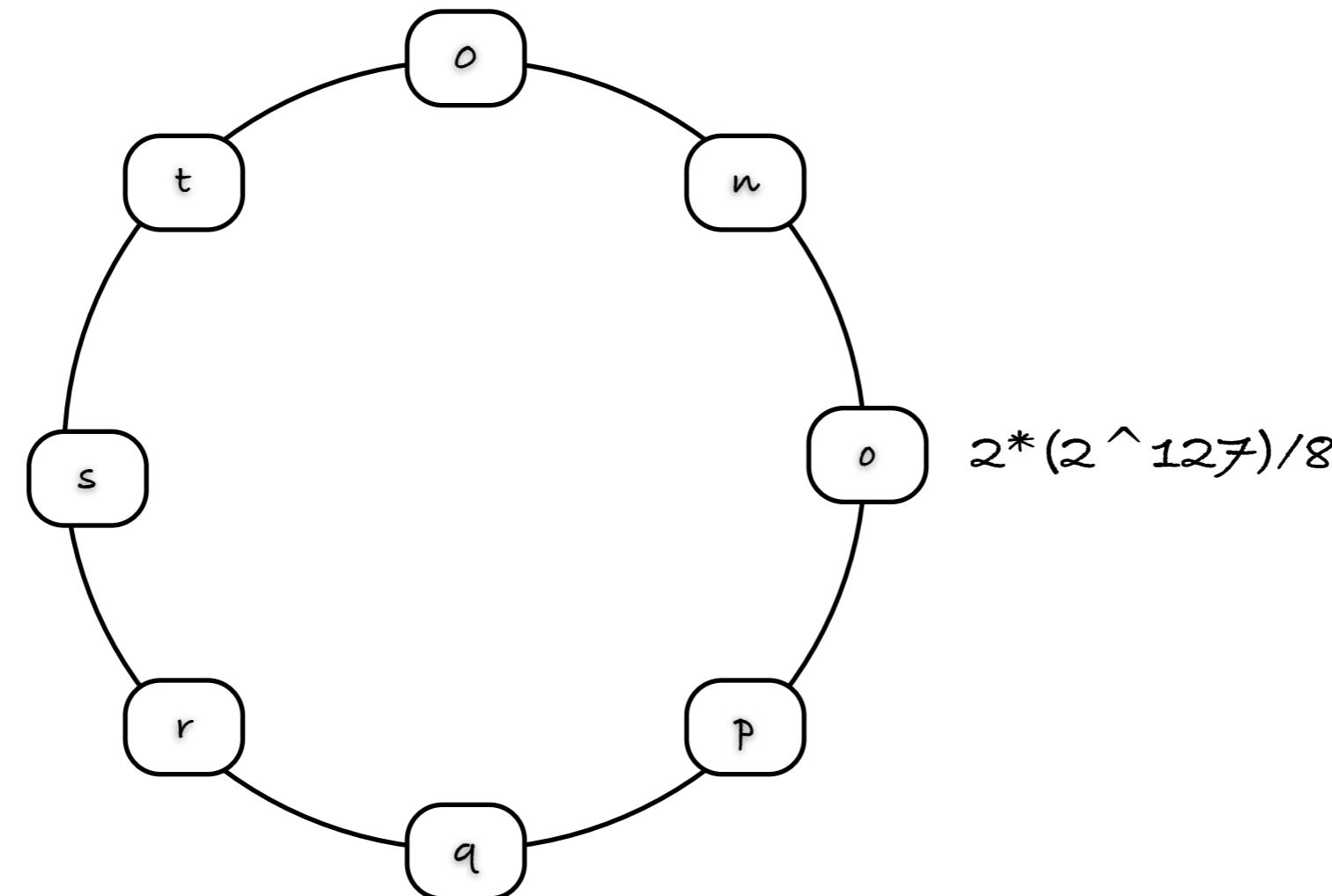


Autres usages Multi DC

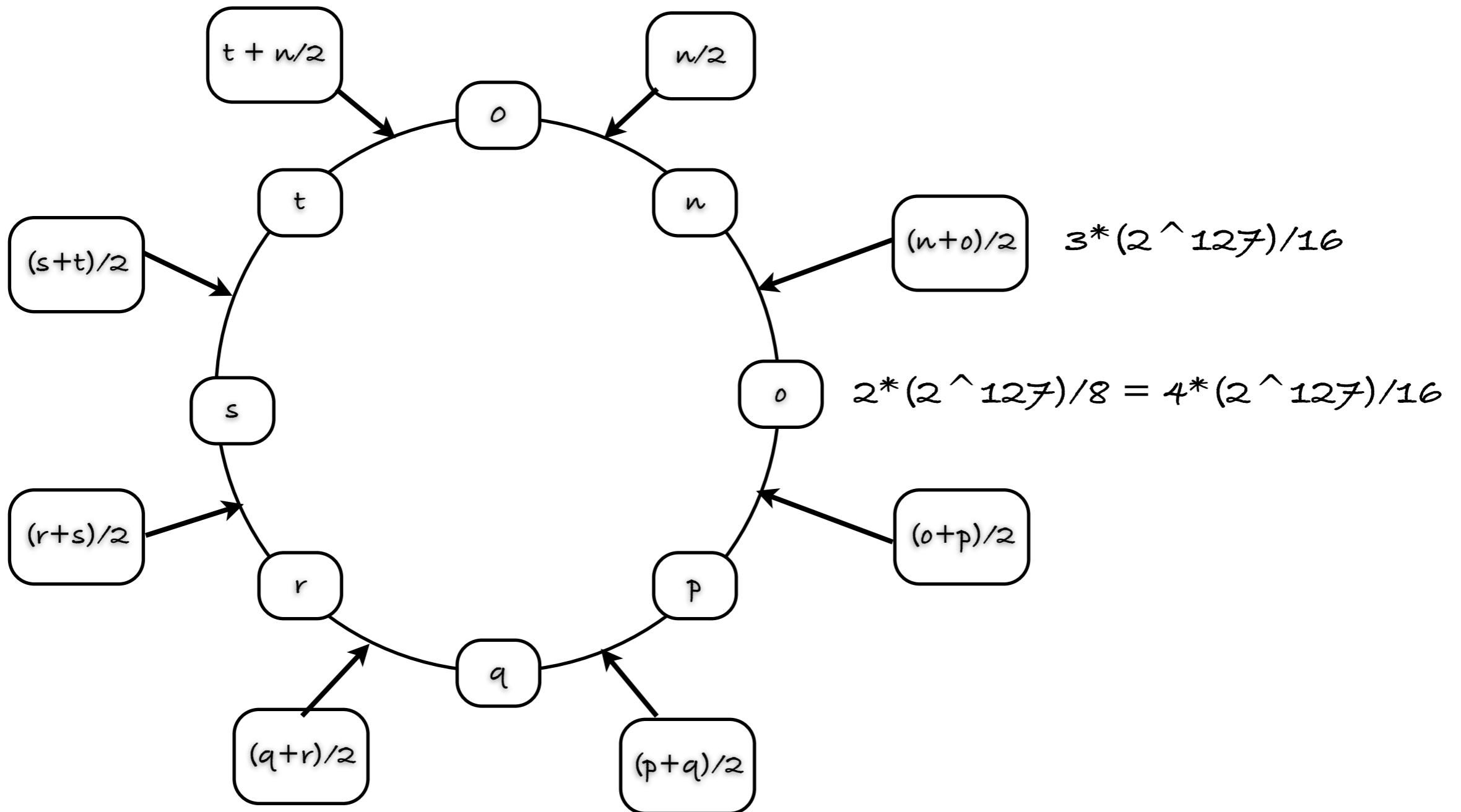
- 1 DC pour backup
- 1 DC pour traitement particulier



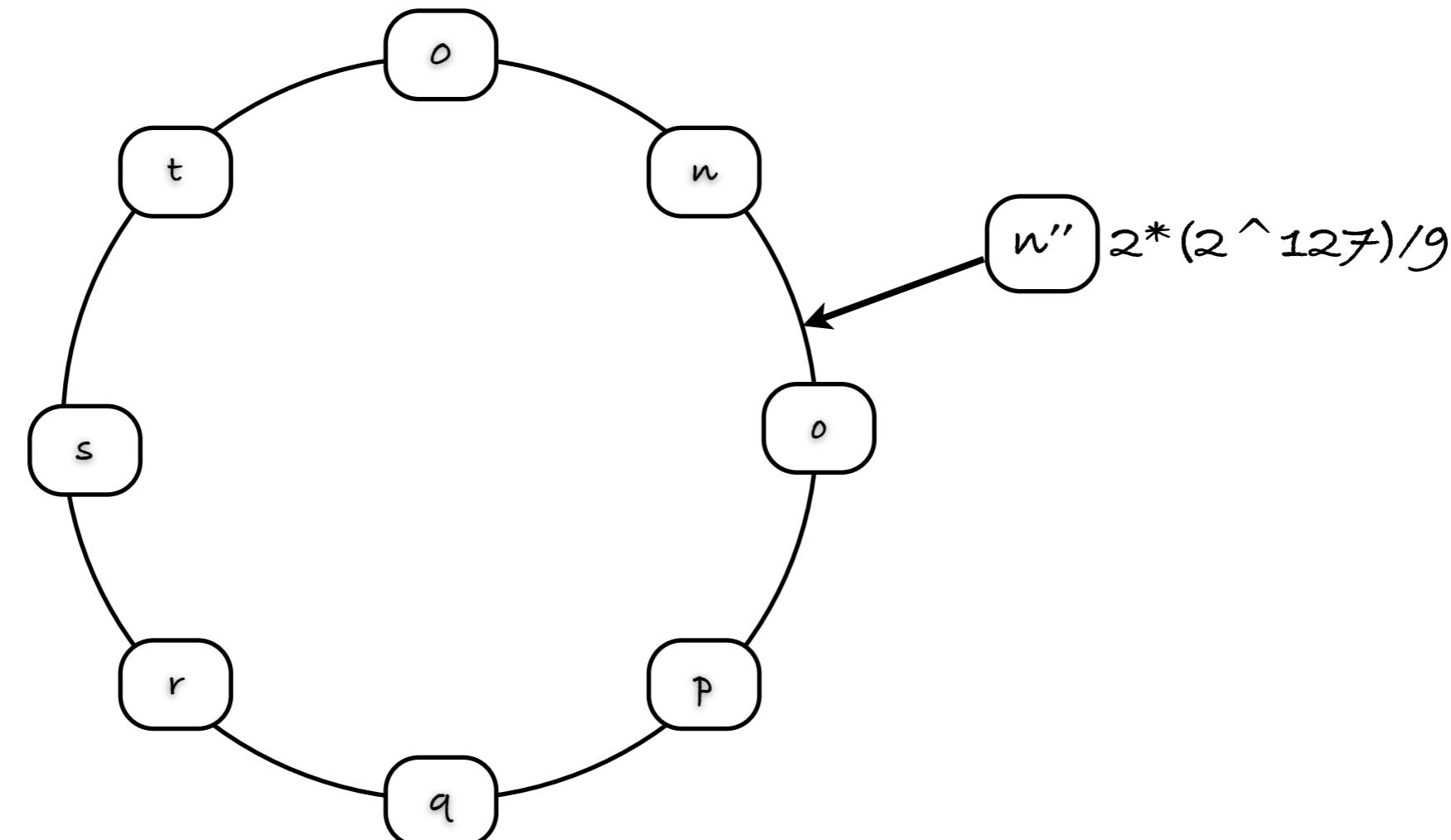
Doubler la taille du cluster



Doubler la taille du cluster



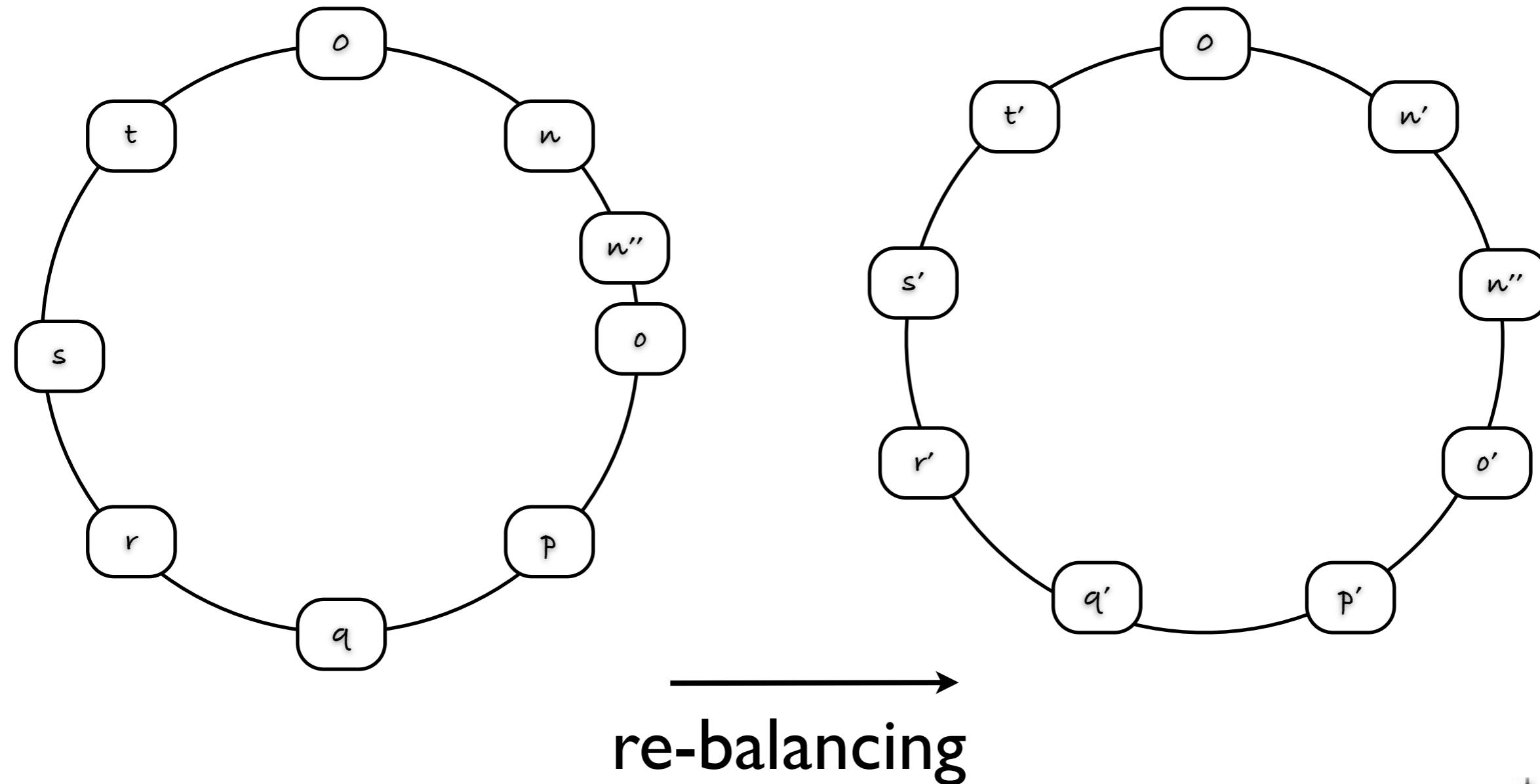
Ajouter 1 noeud



On crée un déséquilibre...



Ajouter 1 noeud

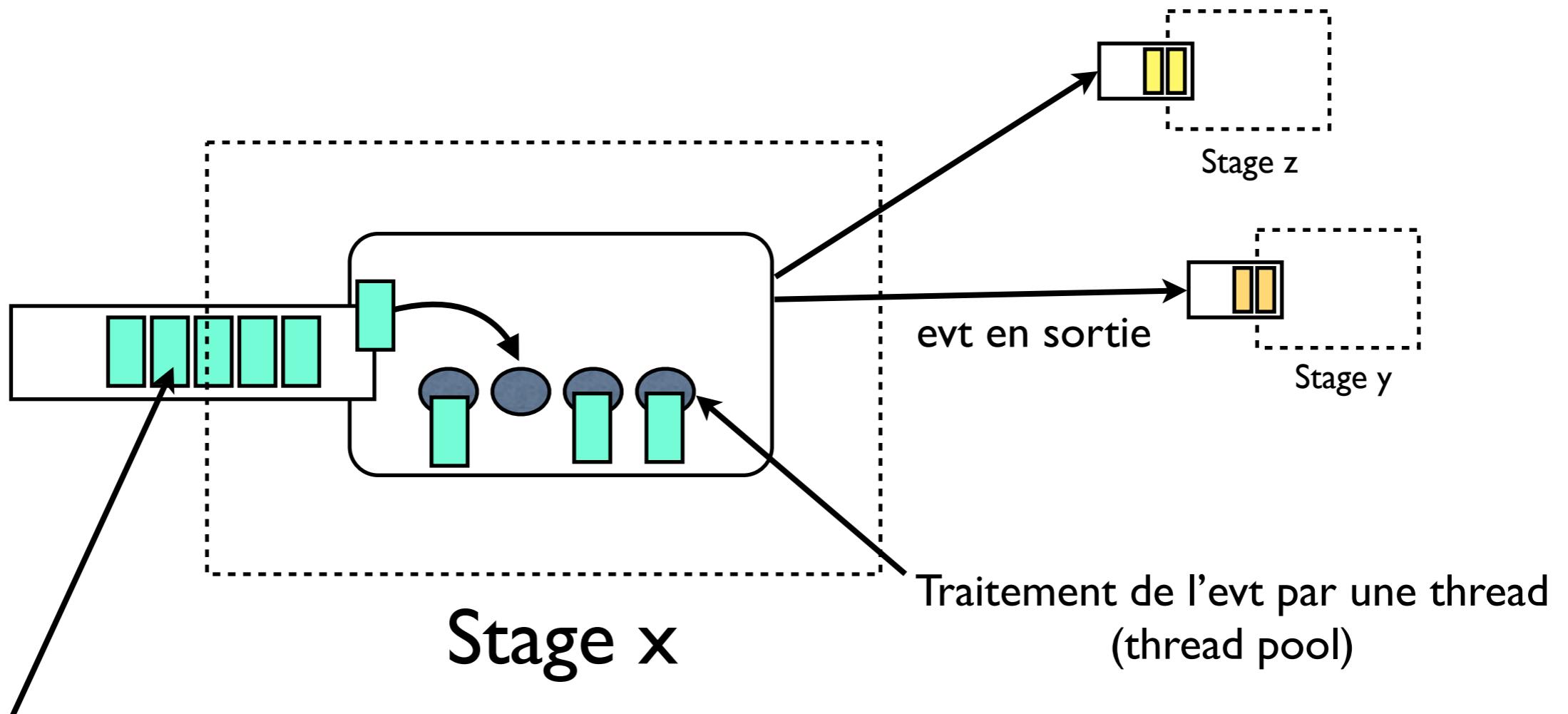


Ajout/Retrait/Upgrade

- Se font à chaud



Architecture Interne & Monitoring



SEDA: Staged Event-Driven Architecture...



Architecture Interne & Monitoring

- Monitoring des «stages» cassandra:
(1) => 1 thread seulement

Pool Name	Active	Pending	Completed	Blocked	All time blocked
ReadStage	0	0	1	0	0
RequestResponseStage	0	0	0	0	0
MutationStage	0	0	20	0	0
ReadRepairStage (1)	0	0	0	0	0
ReplicateOnWriteStage	0	0	0	0	0
GossipStage (1)	0	0	86919	0	0
AntiEntropyStage (1)	0	0	0	0	0
MigrationStage (1)	0	0	0	0	0
MemtablePostFlusher	0	0	4	0	0
StreamStage	0	0	0	0	0
FlushWriter	0	0	4	0	0
MiscStage (1)	0	0	0	0	0
InternalResponseStage	0	0	0	0	0
HintedHandoff	0	0	13	0	0
Message type	Dropped				
RANGE_SLICE	0				
READ_REPAIR	0				
BINARY	0				
READ	0				
MUTATION	0				
REQUEST_RESPONSE	0				

Messages expirés sont «droppés»

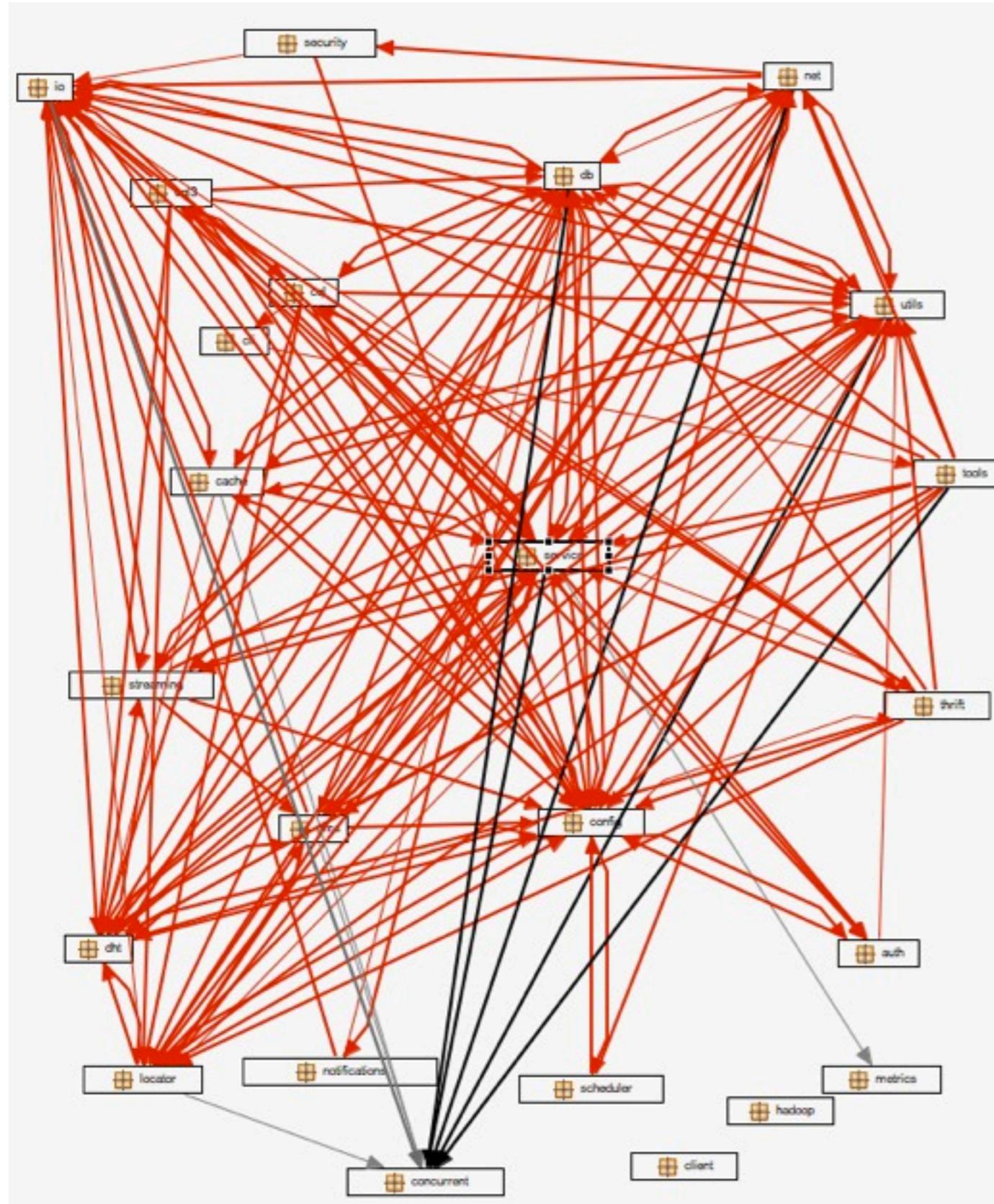


Une critique... si si

- Difficile de rentrer dans le code source



Une critique... si si



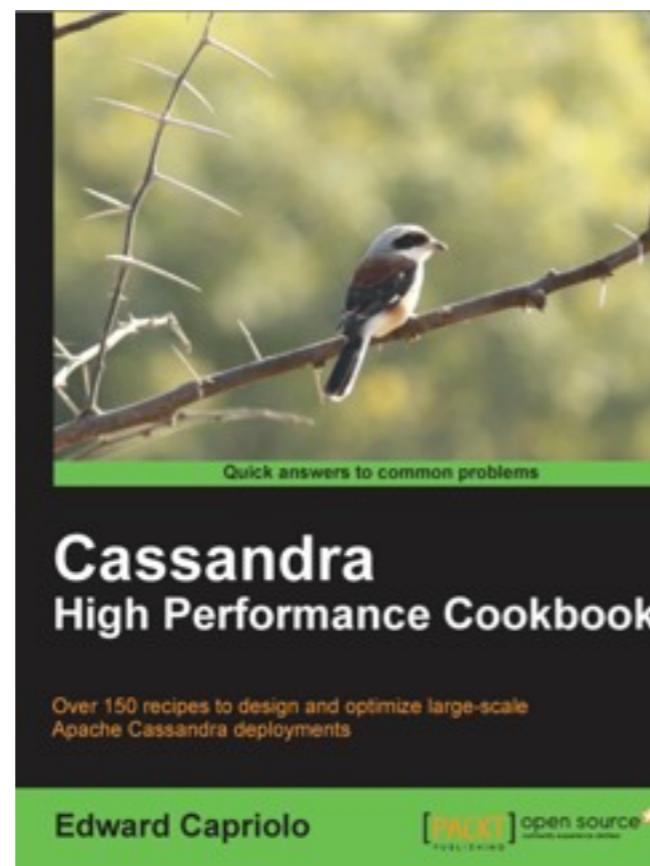
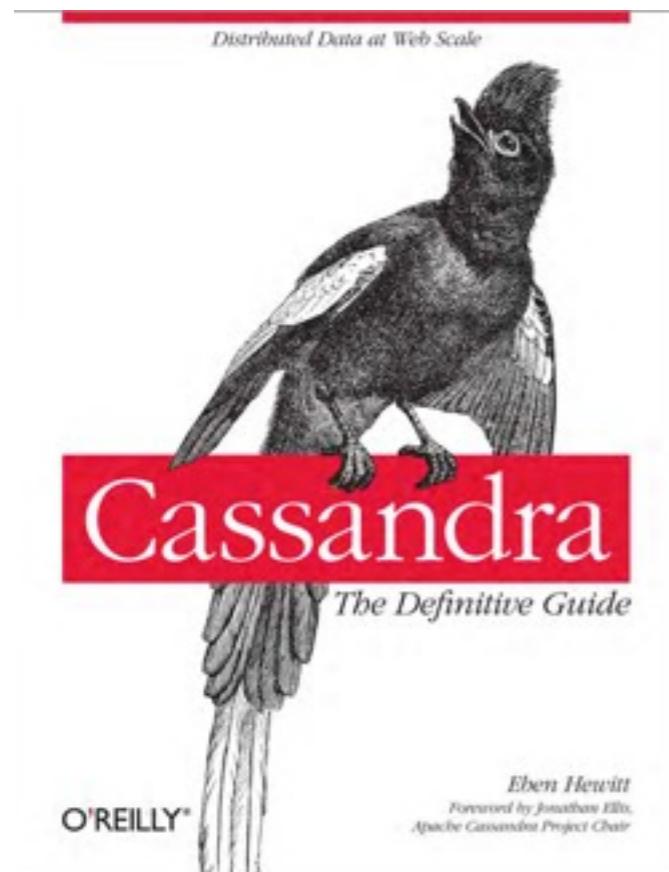
- Difficile de rentrer dans le code source



Questions?



Pour approfondir



Apache Cassandra in Action

By [Jonathan Ellis](#)

Publisher: O'Reilly Media

Released: February 2011

Run time: 3 hours 4 minutes



[Read 1 Review](#) | [Write a Review](#)

<http://www.datastax.com/docs/1.0/index>

<http://wiki.apache.org/cassandra/>

<http://ria101.wordpress.com/2010/02/24/hbase-vs-cassandra-why-we-moved/>

