

Auteurs

GAO Guangze
JIAO Rusong

Encadrant

SUTRA Pierre

Partenaires



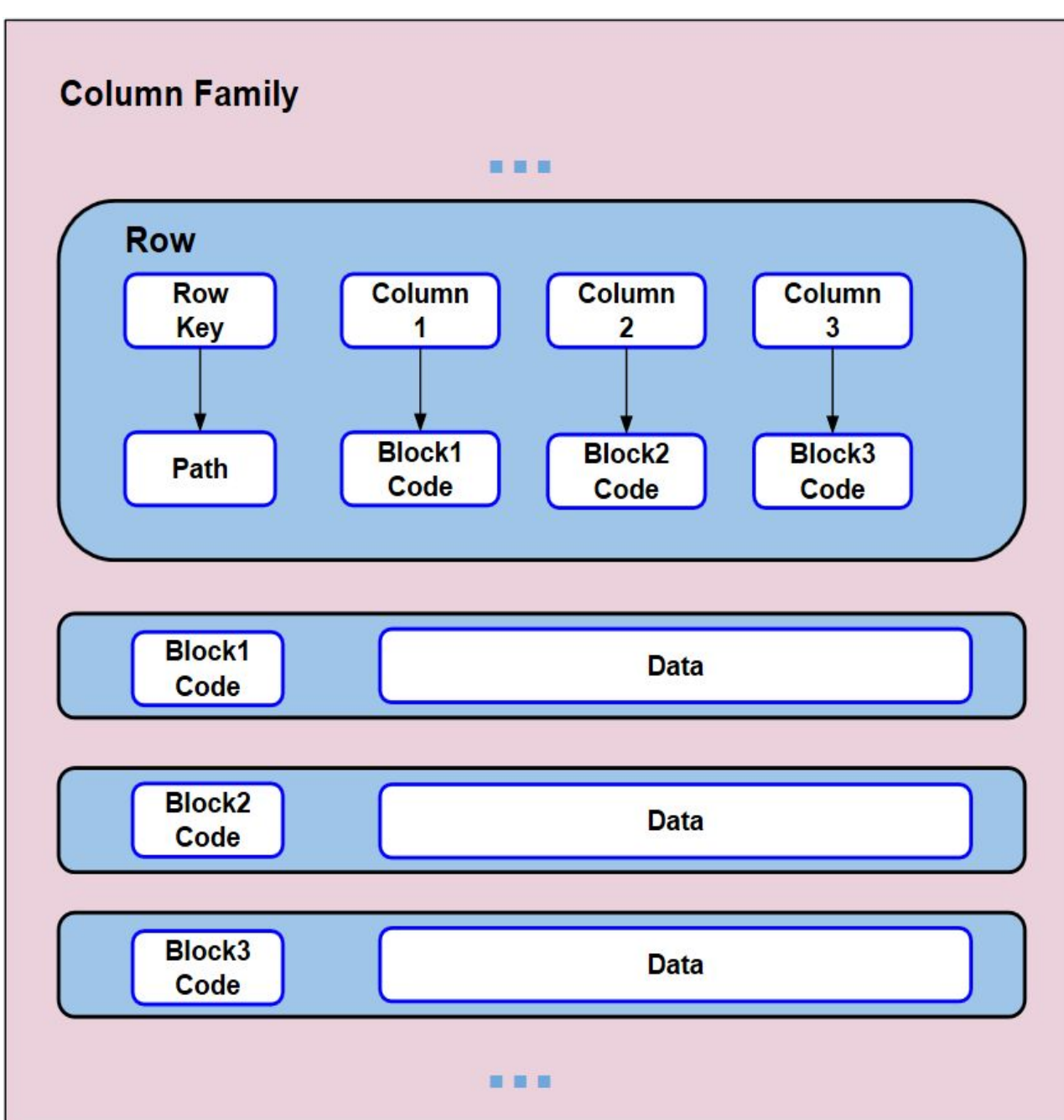
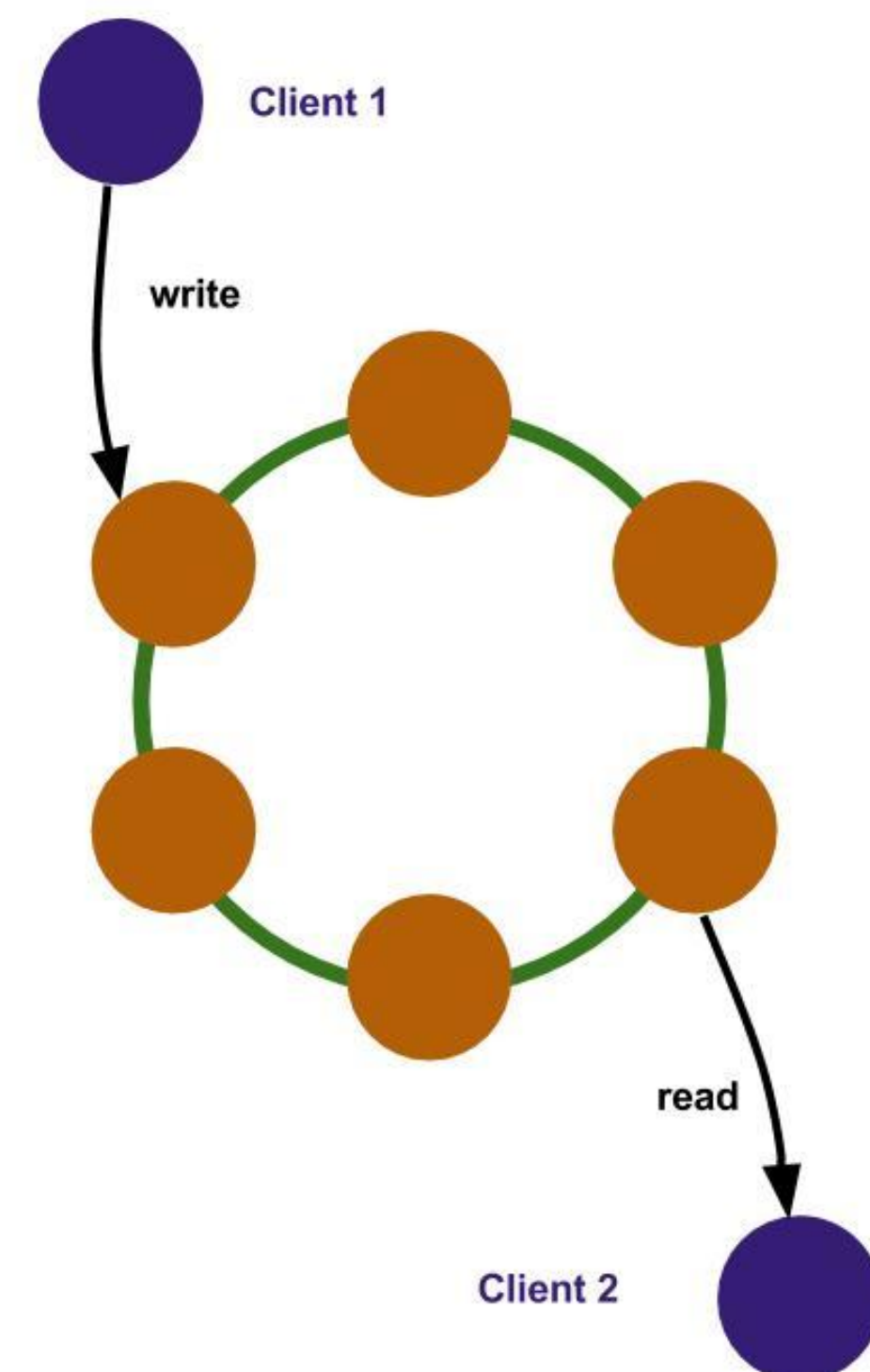
Technologies



Context du projet

DFS : Distributed file system

- Un système de fichiers distribués (en anglais, *distributed file system*) est un système de fichiers qui permet le partage de fichiers à plusieurs clients au travers du réseau informatique.
- De nombreuses applications pratiques ont entré dans notre vie quotidienne : Sun Microsystems' Network File System (*NFS*), Google File System (*GFS*) et Hadoop Distributed File System, etc.
- L'objectif du projet est de réaliser un système de fichiers distribués
- Apache Cassandra est utilisé pour le stockage de données distribuées; Filesystem in Userspace (*FUSE*) nous permet de créer un système de fichiers dans l'espace d'utilisateur.



Solution et implémentation

Deux étapes pour l'implémentation

- Nous avons utilisé Cassandra et FUSE pour créer ce système
- La solution naïve. Stocker toutes les données dans une seule colonne. Un fichier occupe une ligne.
- La solution améliorée. Découper des données en plusieurs blocs de données. Générer un hashcode pour chaque bloc. Stocker chaque pair (code, bloc) dans une ligne. Stocker des codes d'un fichier dans une ligne. Écriture et lecture des données sont parallèles.
- Nous avons choisi python comme le langage de programmation. Nous avons utilisé les bibliothèques de Python (*fusepy* et *pycassa*) pour créer ce système. Package *future* est utilisé pour la parallélisme.

Résultats du projet

Un prototype de système de fichiers distribués

- Un système pour stocker et récupérer des données à distance.
- Insuffisance de la solution naïve : refaire l'écriture de toutes les données d'un fichier même si juste changer une partie du fichier.
- La version améliorée a une meilleure performance que la solution naïve en découpant un fichier en plusieurs blocs pour éviter l'écriture de tout le fichier et en parallélisant les opérations.
- Une autre amélioration possible à continuer est d'ajouter la technologie de caching qui est une raison importante pour Network File System (*NFS*) d'avoir une meilleure performance.

