

- ▶ `http://elle.epfl.ch/wiki`
- ▶ `Logiciels-Libres@groupes.epfl.ch`
- ▶ `epfl.comp.logiciels-libres`

Ce que vous ne voulez pas savoir sur les systèmes de fichiers

Vittoria Rezzonico, SB-IT

7 octobre 2009



Nagios®



Séminaires passés

Extensions Firefox	Michela	Tout public	3 novembre 2008
IPv6	Laurent	Avancé	28 novembre 2008
Tequila	Claude	Tout public	19 mars 2009
PXE	Vicky et Ale	Avancé	10 juin 2009

Bientôt ?

- ▶ GPG
- ▶ Extensions Thunderbird
- ▶ Extensions Firefox, II
- ▶ PAM
- ▶ vi
- ▶ Nagios

.config - Linux Kernel v2.6.30.7 Configuration**File systems**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

<> Second extended fs support

- <> Ext3 journalling file system support
- <> The Extended 4 (ext4) filesystem
- <> Reiserfs support
- <> JFS filesystem support
- <> XFS filesystem support
- <> GFS2 file system support
- <> OCFS2 file system support
- <> Btrfs filesystem (EXPERIMENTAL) Unstable disk format
- [*] Inotify support

v(+)

Select

< Exit >

< Help >

.config - Linux Kernel v2.6.30.7 Configuration**File systems**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

^(-)

```
<M> FUSE (Filesystem in Userspace) support
  Caches --->
  CD-ROM/DVD Filesystems --->
  DOS/FAT/NT Filesystems --->
  Pseudo filesystems --->
  [*] Miscellaneous filesystems --->
  [*] Network File Systems --->
  Partition Types --->
  {M} Native language support --->
  <M> Distributed Lock Manager (DLM) --->
```

Select

< Exit >

< Help >

.config - Linux Kernel v2.6.30.7 Configuration**Miscellaneous filesystems**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

+- Miscellaneous filesystems

- <> **A**DFS file system support (EXPERIMENTAL)
- <> **A**miga FFS file system support (EXPERIMENTAL)
- <> **e**Crypt filesystem layer support (EXPERIMENTAL)
- <> **A**pple Macintosh file system support (EXPERIMENTAL)
- <> **A**pple Extended HFS file system support
- <> **B**eOS file system (BeFS) support (read only) (EXPERIMENTAL)
- <> **H**FS file system support (EXPERIMENTAL)
- <> **E**FS file system support (read only) (EXPERIMENTAL)
- <> **J**ournalling Flash File System v2 (JFFS2) support

v(+)

Select>

< Exit >

< Help >

.config - Linux Kernel v2.6.30.7 Configuration**Miscellaneous filesystems**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

^(-)

- <> UBIFS file system support
- <M> Compressed ROM file system support (cramfs)
- <> SquashFS 4.0 - Squashed file system support
- <> FreeVxFS file system support (VERITAS VxFS(TM) compatible)
- <> Minix file system support
- <> SonicBlue Optimized MPEG File System support
- <> OS/2 HPFS file system support
- <> ONX4 file system support (read only)
- <M> ROM file system support**
 - RomFS backing stores (Block device-backed ROM file system)

v(+)

<Select>

< Exit >

< Help >

.config - Linux Kernel v2.6.30.7 Configuration**Miscellaneous filesystems**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

^(-)

- <> FreeVxFS file system support (VERITAS VxFS(TM) compatible)
- <> Minix file system support
- <> SonicBlue Optimized MPEG File System support
- <> OS/2 HPFS file system support
- <> QNX4 file system support (read only)
- <M> ROM file system support
 - RomFS backing stores (Block device-backed ROM file system)
- <> System V/Xenix/V7/Coherent file system support
- <> UFS file system support (read only)
- <> NILFS2 file system support (EXPERIMENTAL)**

<Select>

< Exit >

< Help >

.config - Linux Kernel v2.6.30.7 Configuration**Network File Systems**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

```
--- Network File Systems
<> NFS client support
<> NFS server support
<> SMB file system support (OBSOLETE, please use CIFS)
<> CIFS support (advanced network filesystem, SMBFS successor)
<> NCP file system support (to mount NetWare volumes)
<> Coda file system support (advanced network fs)
<*> Andrew File System support (AFS) (EXPERIMENTAL)
[ ] AFS dynamic debugging
```

Select < Exit > < Help >

.config - Linux Kernel v2.6.30.7 Configuration**Multiple devices driver support (RAID and LVM)**

Arrow keys navigate the menu. <Enter> selects submenus --->.
Highlighted letters are hotkeys. Pressing <Y> includes, <N> excludes,
<M> modularizes features. Press <Esc><Esc> to exit, <?> for Help, </>
for Search. Legend: [*] built-in [] excluded <M> module <>

Multiple devices driver support (RAID and LVM)

- <> RAID support
- <> Device mapper support

<Select>

< Exit >

< Help >



Partition de 50GB :

- ▶ mkfs.ext3 : 23.998s
- ▶ mkfs.reiserfs : 2.759s



Partition de 50GB :

- ▶ mkfs.ext3 : 23.998s
- ▶ mkfs.reiserfs : 2.759s
- ▶ mkfs.xfs : 0.522s

- ▶ ZFS
- ▶ ext4
- ▶ reiser4
- ▶ btrfs

```
01b0  00 00 00 00 00 00 00 00 18 7a 07 00 00 00 00 01 |.....z.....|
01c0  01 00 83 fe 3f 81 3f 00 00 00 c3 dd 1f 00 00 00 |....?.?.....|
01d0  00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 |.....|
*
01f0  00 00 00 00 00 00 00 00 00 00 00 00 00 55 aa |.....U.|
```

- ▶ La table des partitions commence à 0x1be (446)

Just kidding !

Je veux mes fichiers maintenant



J'accepte d'accéder à des données peut-être
corrompues dans un temps irraisonnable



- ▶ Fiabilité
- ▶ Performances

Définitions

Fiabilité

fsck

journal

checksums

checkpoints

Performances

Fiabilité et performances



Définitions

Fiabilité

fsck

journal

checksums

checkpoints

Performances

Fiabilité et performances

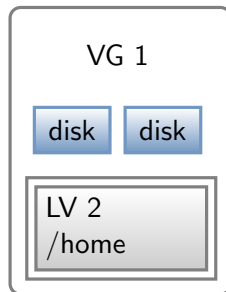
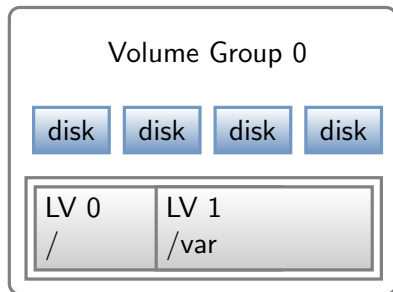
- ▶ façon d'organiser les données
 - ▶ sur les disques
 - ▶ sur le réseau
 - ▶ dans la RAM

Gestion par volumes logiques

- ▶ pas tous les systèmes de fichiers l'incluent
- ▶ on peut s'appuyer sur LVM

Le système de fichiers au sense strict

- ▶ sur un volume





Définitions

Fiabilité

fsck

journal

checksums

checkpoints

Performances

Fiabilité et performances

- ▶ Contrôle d'intégrité
 - ▶ fsck
 - ▶ comparaisons avec un log (journal)
 - ▶ checksums

- ▶ Contrôle d'intégrité
 - ▶ fsck
 - ▶ comparaisons avec un log (journal)
 - ▶ checksums
- ▶ Consistance des données à tout moment
 - ▶ checkpoints

- ▶ pour vérifier la cohérence d'un système de fichiers
- ▶ et corriger les incohérences
- ▶ lent, en général exécuté au démarrage après un arrêt brutal
- ▶ en général exécuté offline

- ▶ pour vérifier la cohérence d'un système de fichiers
- ▶ et corriger les incohérences
- ▶ lent, en général exécuté au démarrage après un arrêt brutal
- ▶ en général exécuté offline

Temps irraisonnable

- ▶ les opérations d'écriture sont enregistrées dans un log avant d'être exécutées
- ▶ on peut logger uniquement les métadonnées
- ▶ les données sont écrites deux fois !
- ▶ log-structured filesystems
- ▶ optimisations : regrouper les écritures

- ▶ les opérations d'écriture sont enregistrées dans un log avant d'être exécutées
- ▶ on peut logger uniquement les métadonnées
- ▶ les données sont écrites deux fois !
- ▶ log-structured filesystems
- ▶ optimisations : regrouper les écritures

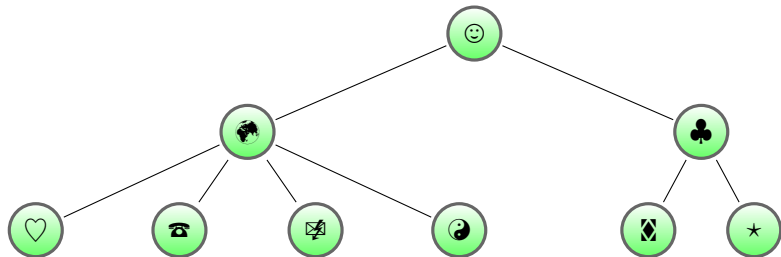
Est-ce adapté aux systèmes de fichiers de grande taille ?

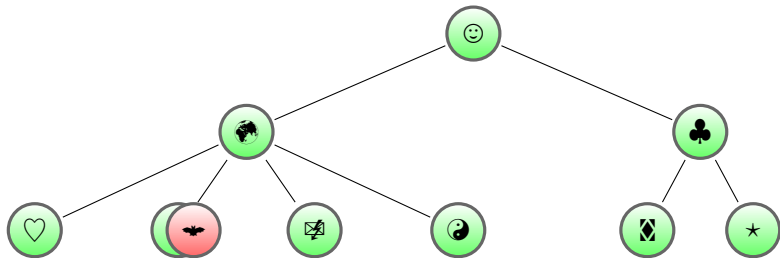
- ▶ comment savoir qu'il n'y a pas eu corruption des données
- ▶ checksum du fichier, des blocks,...
- ▶ stocké dans le "bloc parent"

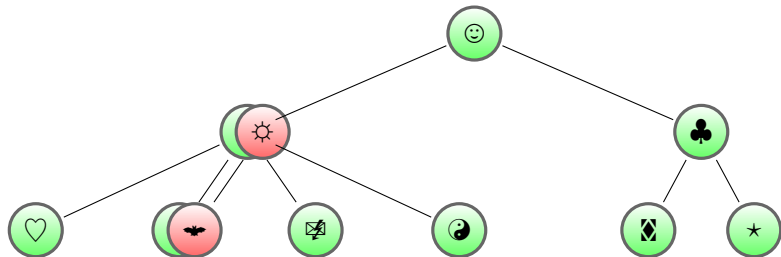
- ▶ éliminent la nécessité de journal et de fsck
- ▶ car le système de fichiers est toujours consistant

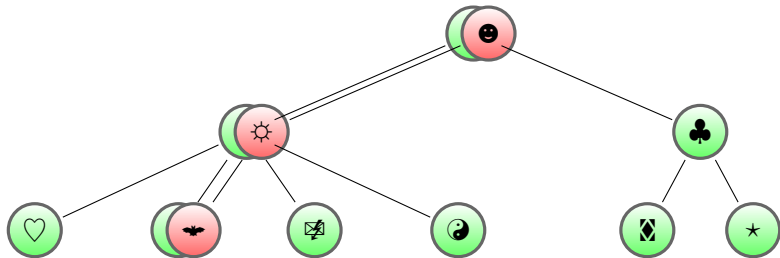
- ▶ éliminent la nécessité de journal et de fsck
- ▶ car le système de fichiers est toujours consistant

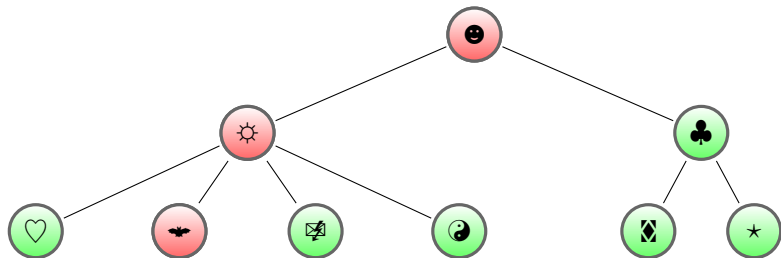
La solution à tous nos problèmes

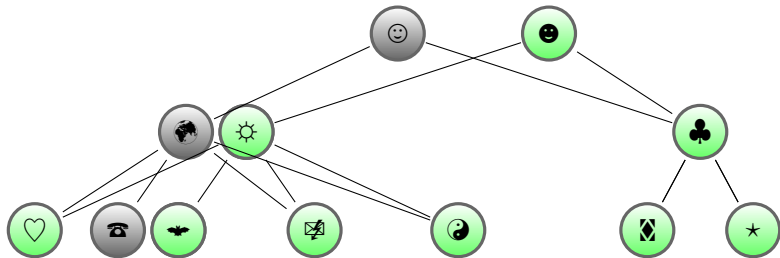












- ▶ snapshots (gratuits)
- ▶ clones (gratuits)



Définitions

Fiabilité

fsck

journal

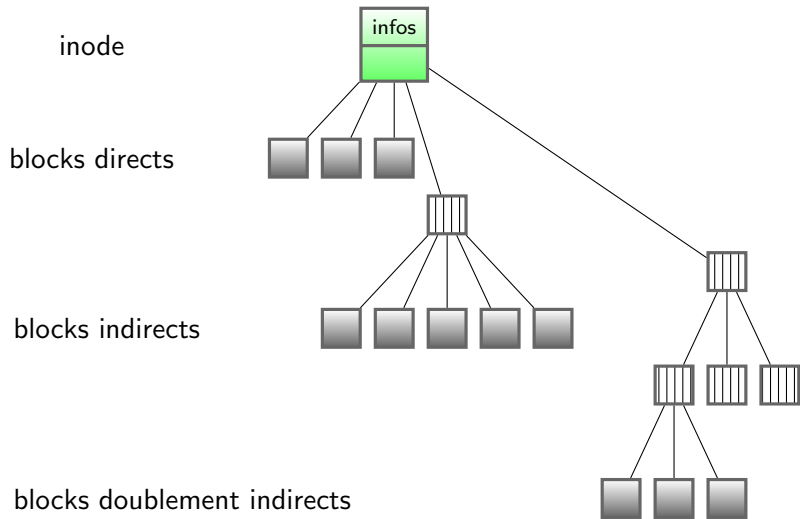
checksums

checkpoints

Performances

Fiabilité et performances

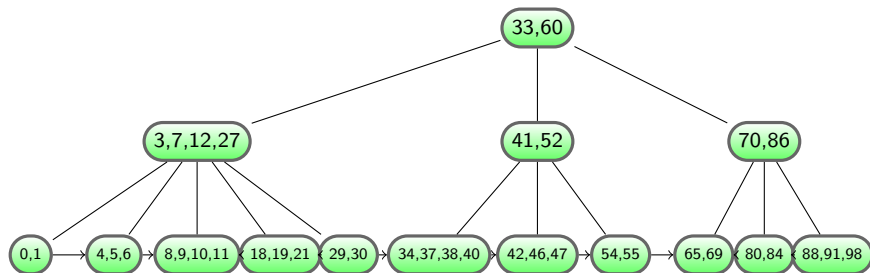
- ▶ Le block : séquence de bytes ou bits
- ▶ Block-size : taille de cette séquence
- ▶ niveau d'abstraction entre le hardware et les datas
- ▶ cluster : groupe de blocks
- ▶ Superblock (ou überblock) : block à la racine de tout, contient plus d'information que les blocks standard.
- ▶ inode : conteneur de métadonnées, généralement lié à des blocks de data.



- ▶ superblock
- ▶ inodes
- ▶ blocks avec uniquement datas
- ▶ blocks avec datas et pointeurs

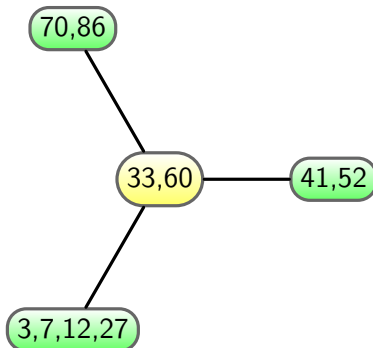
- ▶ le superblock se trouve physiquement après le boot record et ses données
- ▶ les inodes sont stockés dans un tableau
- ▶ les datas aussi, dans une autre zone du disque
- ▶ bitmaps pour les blocks des inodes et des datas
- ▶ tableau pour le nombre de références (clones, snapshots)

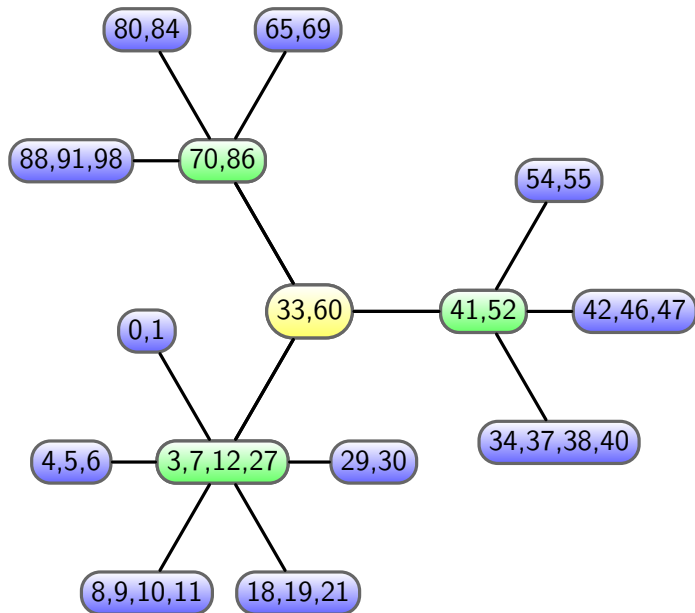
- ▶ *Balanced Tree*
- ▶ B-Tree d'ordre d : chaque noeud (sauf la racine) a entre d et $2d$ feuilles
- ▶ L-U B-Tree : chaque noeud a entre L et U feuilles.
- ▶ Clef de tri : index, p. ex ici **52**

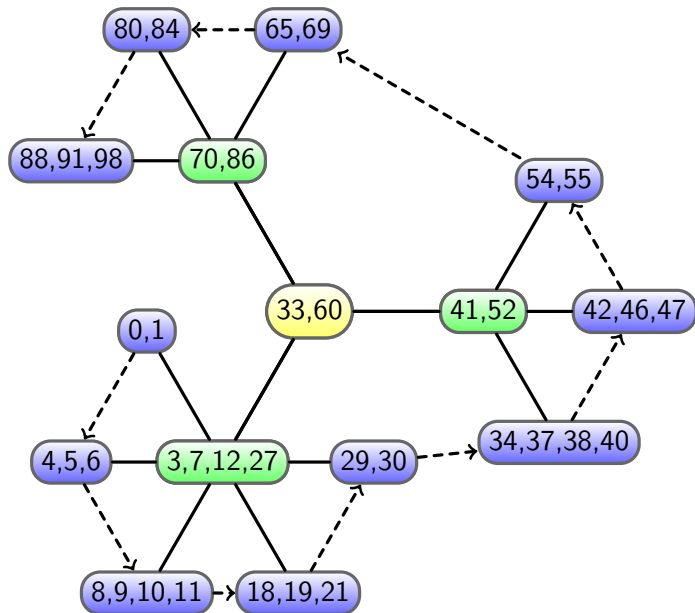


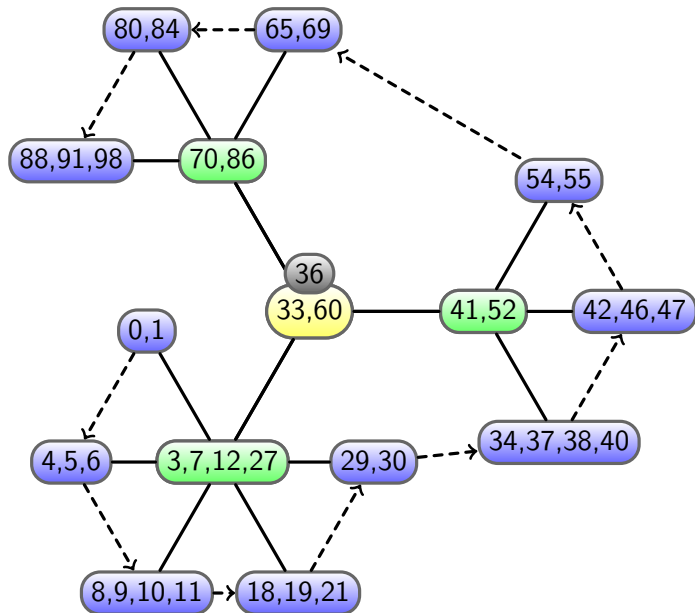
Exemples avec $n = 1M$ et $p = 2$

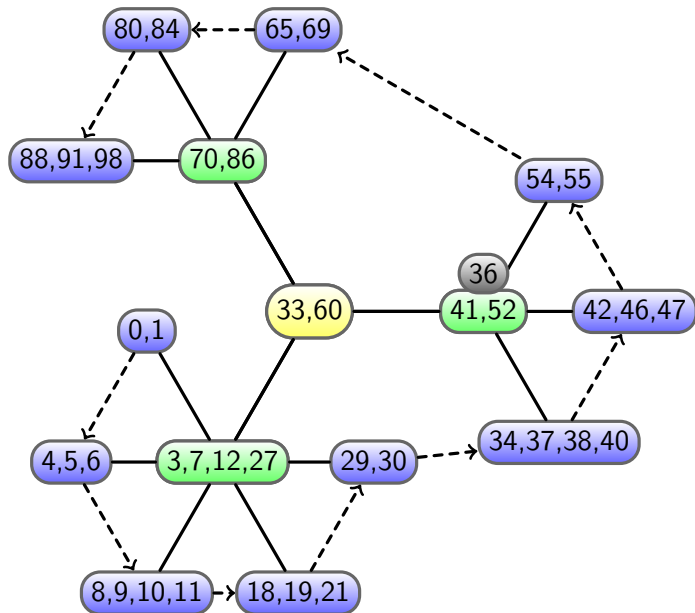
- ▶ Recherche $O(\log_p(n))$ 19.93
- ▶ Insertion $O(p \log_p(n))$ 20.93
- ▶ Remove $O(p \log_p(n))$ 20.93
- ▶ n clefs peuvent être casées dans au plus $\log_t \frac{n+1}{2}$ niveaux 18
- ▶ h niveaux peuvent contenir jusqu'à $2p \sum_{\ell=0}^h (2p+1)^\ell$ clefs 10^{13}

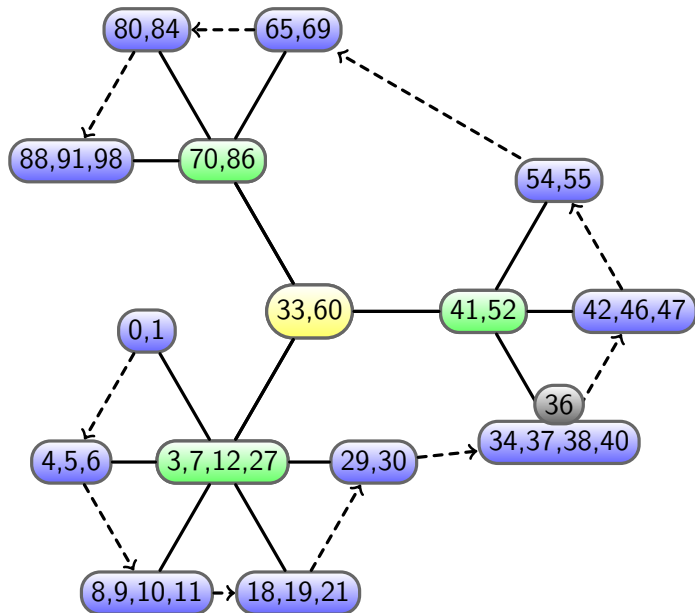


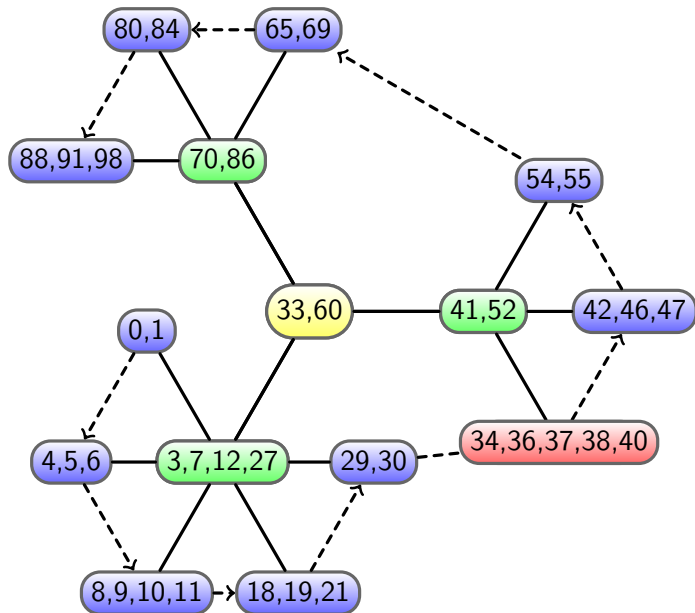


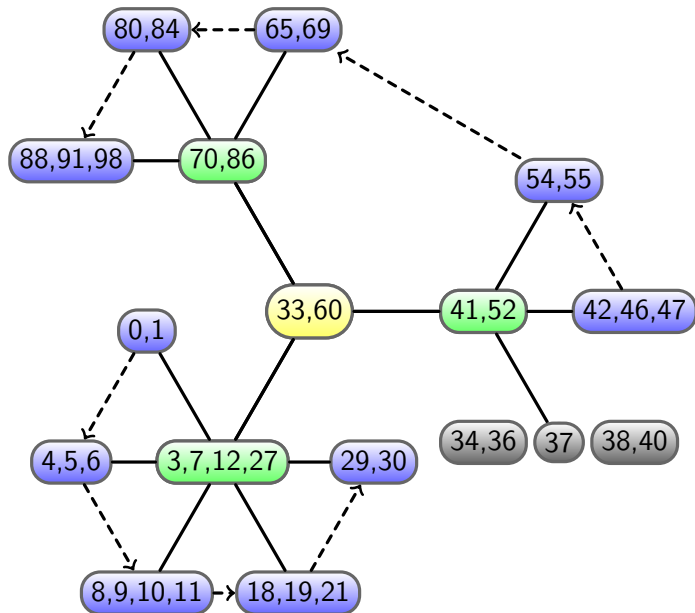


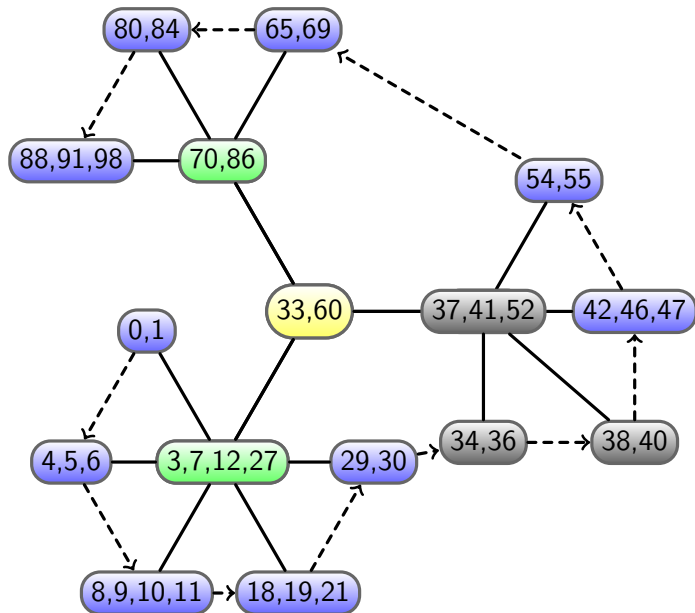


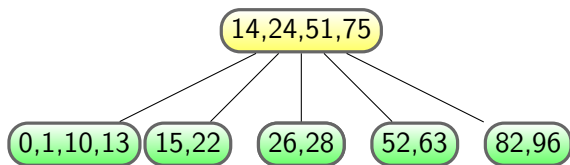


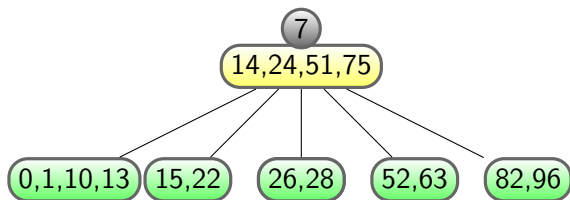


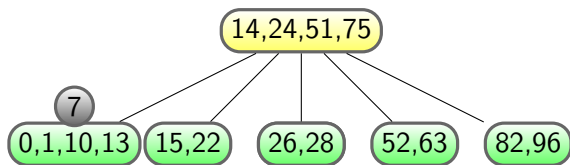


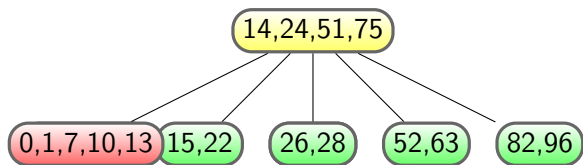


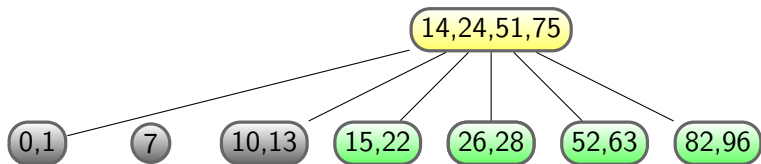


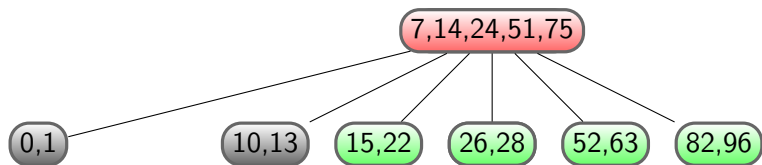


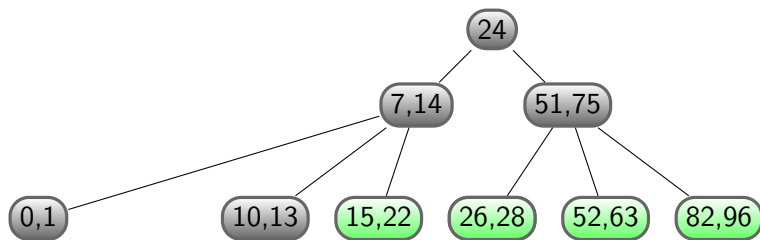








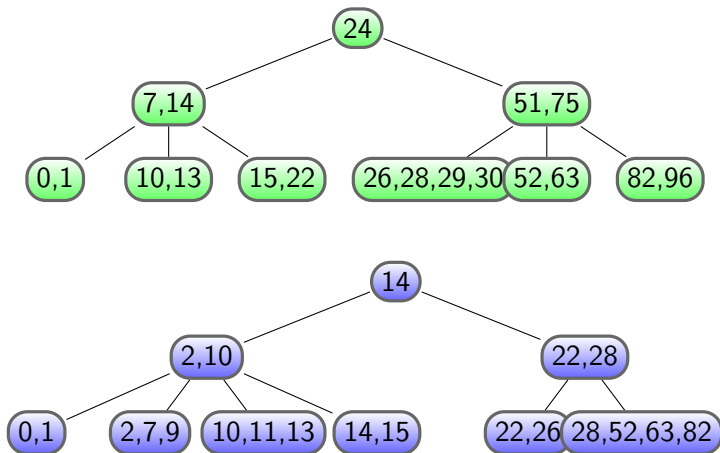


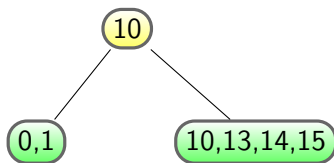


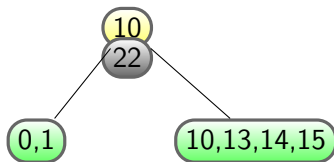
- ▶ Deux types de noeuds : les feuilles et les noeuds internes
- ▶ Les noeuds internes pour l'indexage
- ▶ Les feuilles contiennent les datas

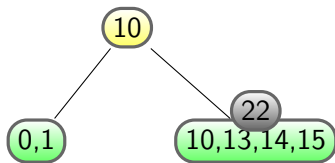
- ▶ Ordre d'un B+tree = nb max de descendants d'un noeud
- ▶ (B-tree d'ordre 2 \approx B+-tree d'ordre 5)
- ▶ In B-+Tree d'ordre p a : au maximum $n = p^h$ entrées

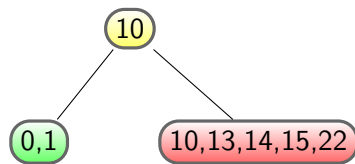
- ▶ Les feuilles sont souvent chaînées pour accélérer des requêtes multiples

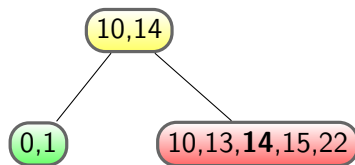


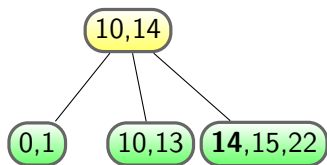


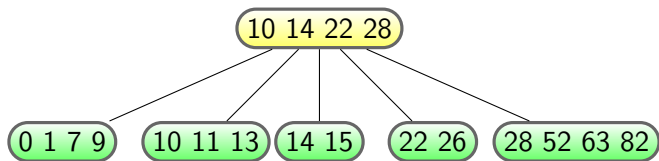


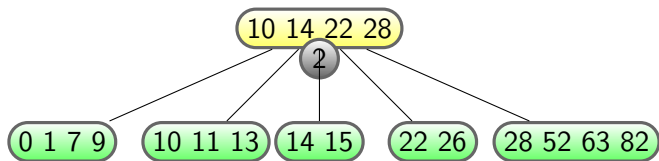


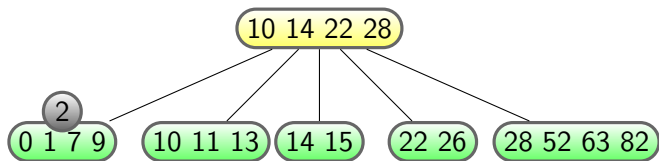


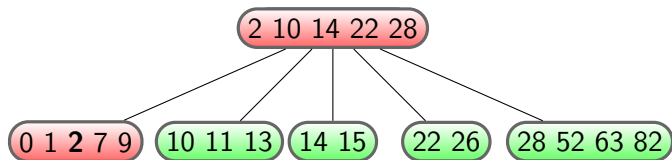


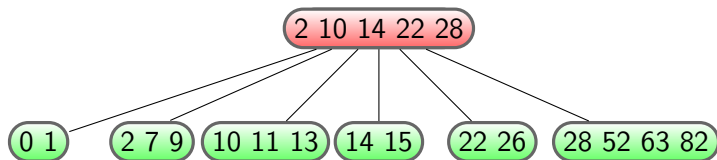


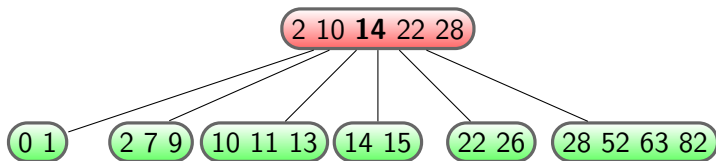


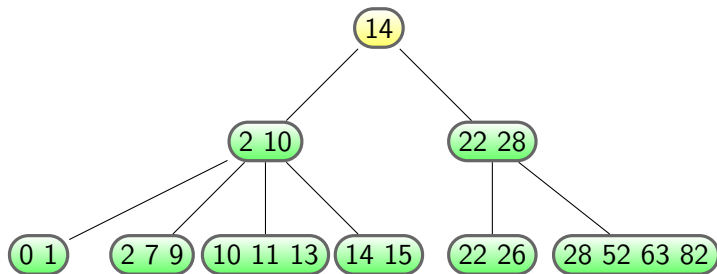




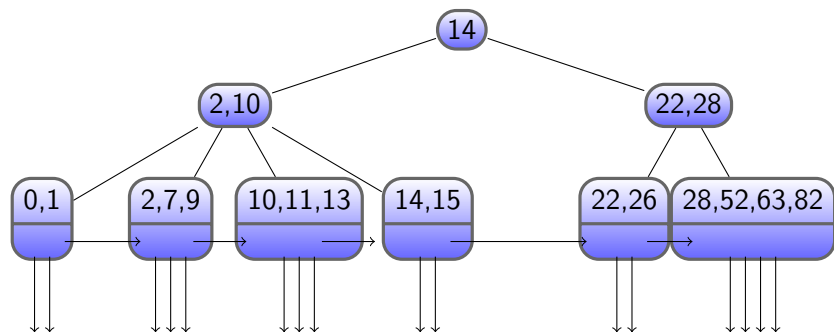








- ▶ les noeuds internes pointent vers leur noeuds fils
- ▶ les feuilles pointent vers les datas
- ▶ et vers la feuille successive





To tree

To tree

Liste exhaustive des fs qui utilisent les arbres balancés



Définitions

Fiabilité

fsck

journal

checksums

checkpoints

Performances

Fiabilité et performances

- ▶ les checkpoints sont incompatibles avec les arbres comme on les connaît

- ▶ les checkpoints sont incompatibles avec les arbres comme on les connaît
- ▶ il faut re-écrire (la moitié de) la théorie sur les B-Trees.

- ▶ les checkpoints sont incompatibles avec les arbres comme on les connaît
- ▶ il faut re-écrire (la moitié de) la théorie sur les B-Trees.
- ▶ ça a été fait **en 2007**

- ▶ principale incompatibilité entre arbres et checkpoints : le fait que les feuilles sont enchainées
⇒ on enlève ces liens, en renonçant entre autre à l'accès rapide à des blocs contigus
- ▶ comment compter les références ?

Questions ?

Trolls ?

- ▶ Dave Hitz et al. : File System Design for an NFS File Server Appliance
- ▶ Zachary N. J. Peterson : Ext3cow : A Time-Shifting File System for Regulatory Compliance
- ▶ Jeff Bonwick et al. : The Zettabyte File System
- ▶ Ohad Rodeh : B-trees, Shadowing, and Clones