

# Comment trouver une aiguille dans une botte de foin

Amélie Gheerbrant  
[amelie@irif.fr](mailto:amelie@irif.fr)

IRIF, Université Paris Diderot  
Semaine d'accueil des LI  
10 septembre 2018





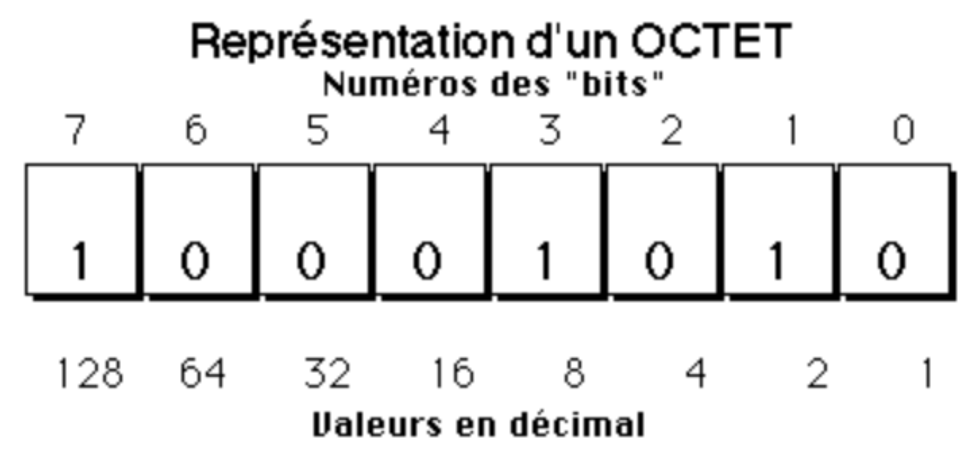


101010010110100110011001010101010101  
 101010010110100110011001010101010101  
 010001010 10011 01010101110010101110110110  
 1100101110101010010111001001010100101  
 1010100101101 00100101010101010101010  
 00110010111011001  
 0010010101010101010101010101010

# Big Data



- Le Big Data c'est quoi ?
- D'abord, des 0 et des 1 (bits)
- Unité de base : l'octet (suite de 8 bits)
- Ordres de grandeur : ko, Mo, Go, etc



138

décimal

10001010

binaire

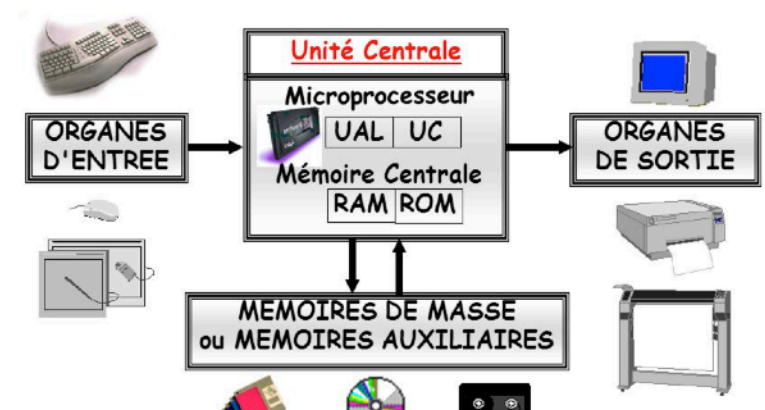



Schéma simplifié d'un micro-ordinateur


# Ordres de grandeur

Nom	Symbole	Valeur
kilooctet	ko	$10^3$
mégaoctet	Mo	$10^6$
gigaoctet	Go	$10^9$
téraoctet	To	$10^{12}$
pétaoctet	Po	$10^{15}$
exaoctet	Eo	$10^{18}$
zettaoctet	Zo	$10^{21}$
yottaoctet	Yo	$10^{24}$

 **temps\_algorithmes\_planches\_2.pdf** 760 Ko  
Modifié : aujourd'hui 00:17


Ajouter des tags...

▼ Général :  
Type : Adobe PDF document  
Taille : 759 843 octets (762 Ko sur disque)

 **Database System Concepts 6e By Abrah...** 15,2 Mo  
Modifié : mardi 18 août 2015 19:35


Ajouter des tags...

▼ Général :  
Type : Adobe PDF document  
Taille : 15 169 704 octets (15,2 Mo sur disque)

 **Enter The Wu-tang (36 Chambers)** 71,4 Mo  
Modifié : mardi 20 mai 2014 23:31

Ajouter des tags...

▼ Général :  
Type : Dossier  
Taille : 71 385 856 octets (71,4 Mo sur disque) pour 13 éléments

 **Bach** 2,54 Go  
Modifié : jeudi 25 décembre 2014 03:04

Ajouter des tags...

▼ Général :  
Type : Dossier  
Taille : 2 543 538 711 octets (2,54 Go sur disque) pour 90 éléments



# Ordres de grandeur

Nom	Symbole	Valeur
kilooctet	ko	$10^3$
mégaoctet	Mo	$10^6$
gigaoctet	Go	$10^9$
téraoctet	To	$10^{12}$
pétaoctet	Po	$10^{15}$
exaoctet	Eo	$10^{18}$
zettaoctet	Zo	$10^{21}$
yottaoctet	Yo	$10^{24}$

Tous les livres jamais écrits :  
quelques centaines de  
téraoctets (en texte brut)



# Ordres de grandeur

Nom	Symbole	Valeur
kilooctet	ko	$10^3$
mégaoctet	Mo	$10^6$
gigaoctet	Go	$10^9$
téraoctet	To	$10^{12}$
pétaoctet	Po	$10^{15}$
exaoctet	Eo	$10^{18}$
zettaoctet	Zo	$10^{21}$
yottaoctet	Yo	$10^{24}$

100 téraoctet :  
volume quotidien des  
contenus mis en ligne  
sur Facebook en 2013

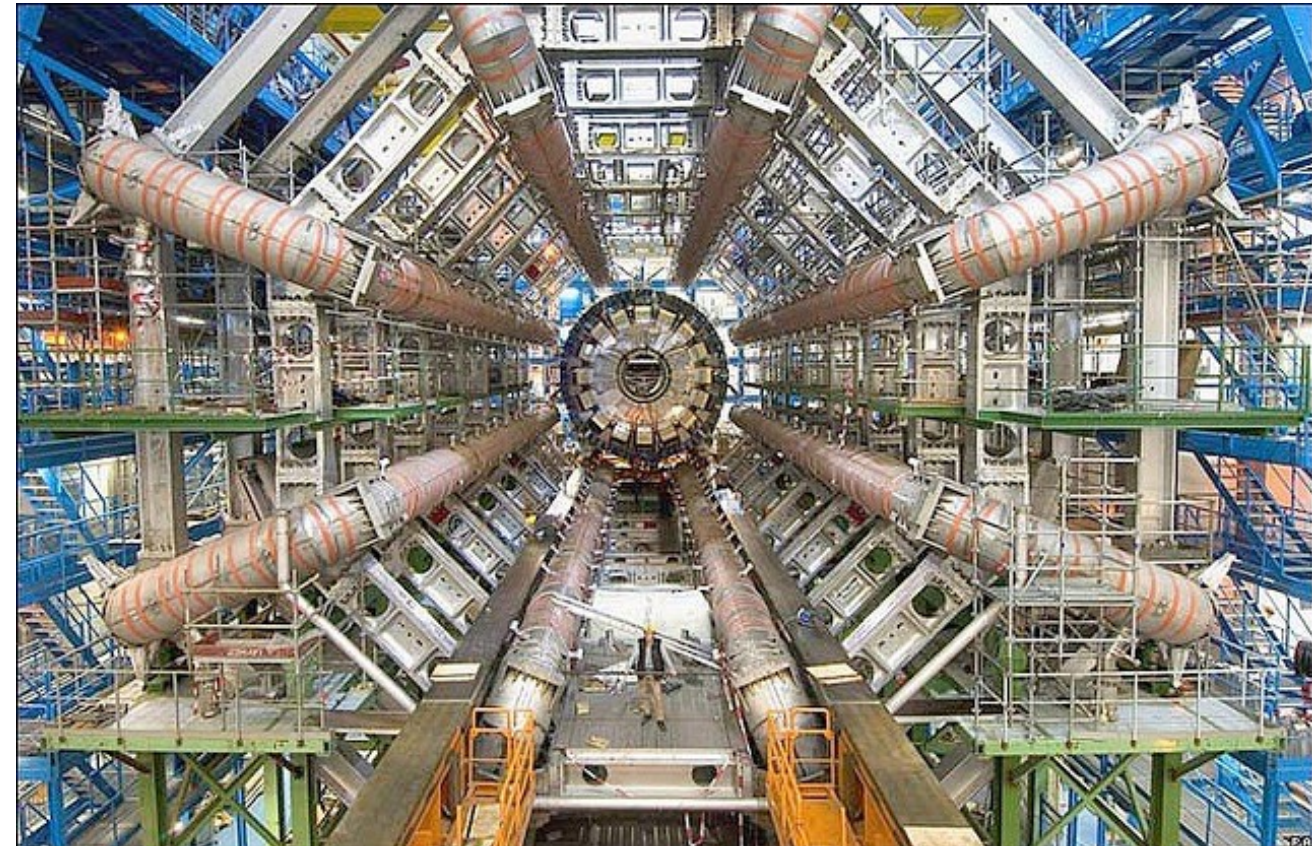
460 téraoctets :  
base de données  
clients de Wal-Mart,  
chaîne américaine de  
supermarché



# Ordres de grandeur

La quantité de données produite par l'accélérateur de particules du CERN en un jour : une centaine de pétaoctets

Nom	Symbole	Valeur
kilooctet	ko	$10^3$
mégaoctet	Mo	$10^6$
gigaoctet	Go	$10^9$
téraoctet	To	$10^{12}$
pétaoctet	Po	$10^{15}$
exaoctet	Eo	$10^{18}$
zettaoctet	Zo	$10^{21}$
yottaoctet	Yo	$10^{24}$

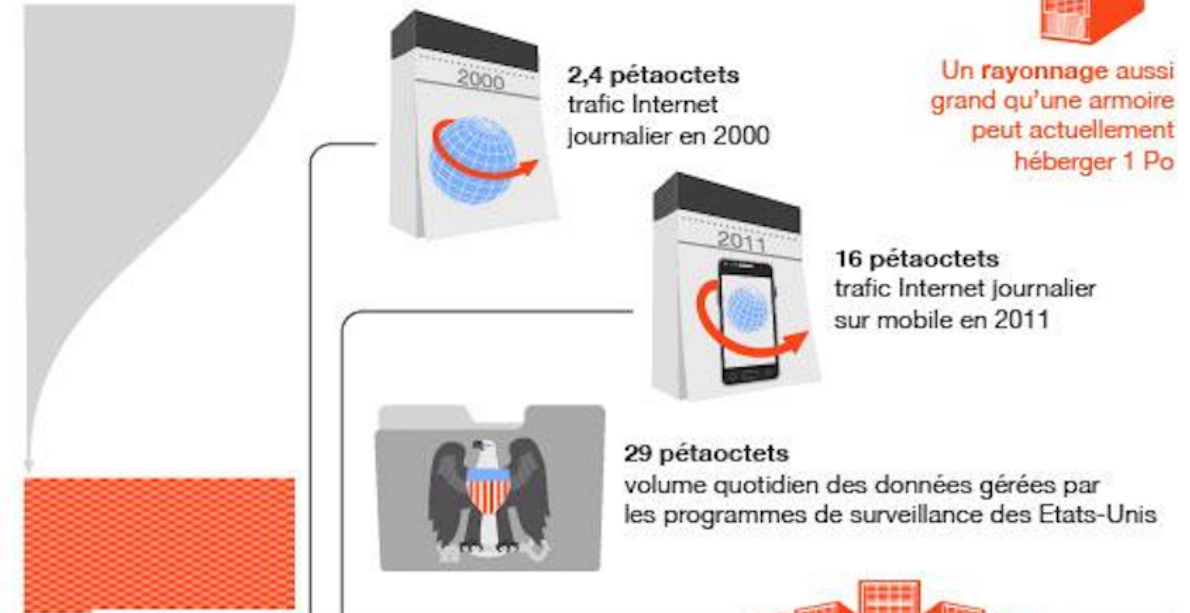




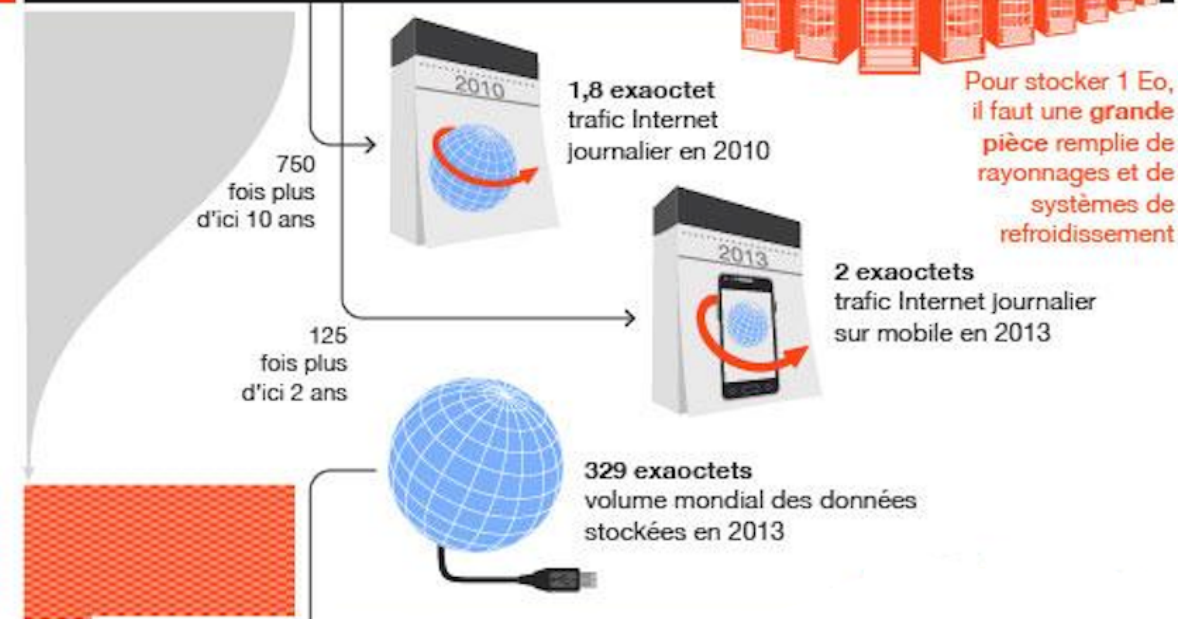
# Ordres de grandeur

Nom	Symbole	Valeur
kilooctet	ko	$10^3$
mégaoctet	Mo	$10^6$
gigaoctet	Go	$10^9$
téraoctet	To	$10^{12}$
pétaoctet	Po	$10^{15}$
exaoctet	Eo	$10^{18}$
zettaoctet	Zo	$10^{21}$
yottaoctet	Yo	$10^{24}$

**PO** 1 pétaoctet = 1000 téraoctets



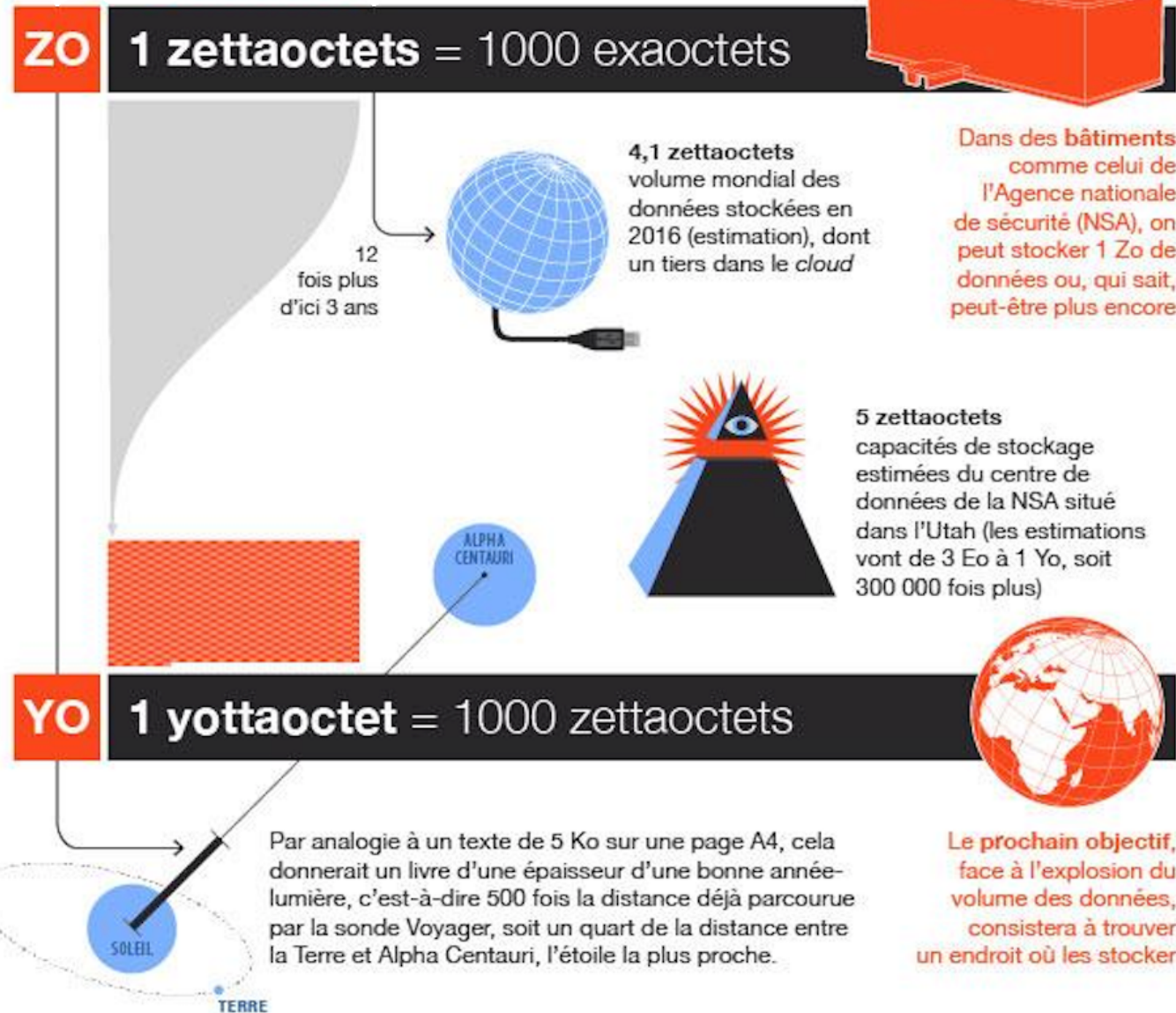
**EO** 1 exaoctet = 1000 pétaoctets





# Ordres de grandeur

Nom	Symbole	Valeur
kilooctet	Ko	$10^3$
mégaoctet	Mo	$10^6$
gigaoctet	Go	$10^9$
téraoctet	To	$10^{12}$
pétaoctet	Po	$10^{15}$
exaoctet	Eo	$10^{18}$
zettaoctet	Zo	$10^{21}$
yottaoctet	Yo	$10^{24}$



Counting Beyond a Yottabyte, or how SPARQL 1.1 Property Paths will Prevent Adoption of the Standard

Marcelo Arenas  
Department of Computer Science  
PUC Chile  
marenas@ing.puc.cl

Sebastián Conca  
Department of Computer Science  
PUC Chile  
saconca@puc.cl

Jorge Pérez  
Department of Computer Science  
Universidad de Chile  
jperez@dcc.uchile.cl

# Bases de données

- Dans Big Data il y a « data » : données
- Organisation sous forme de **bases de données** : Facebook, Amazon, Twitter, Instagram, Discogs, etc
- Gestion via un SGBD (Système de Gestion de Bases de Données) : système très puissant, capable de trouver une aiguille dans une botte de foin en une fraction de seconde

ORACLE®  
DATABASE

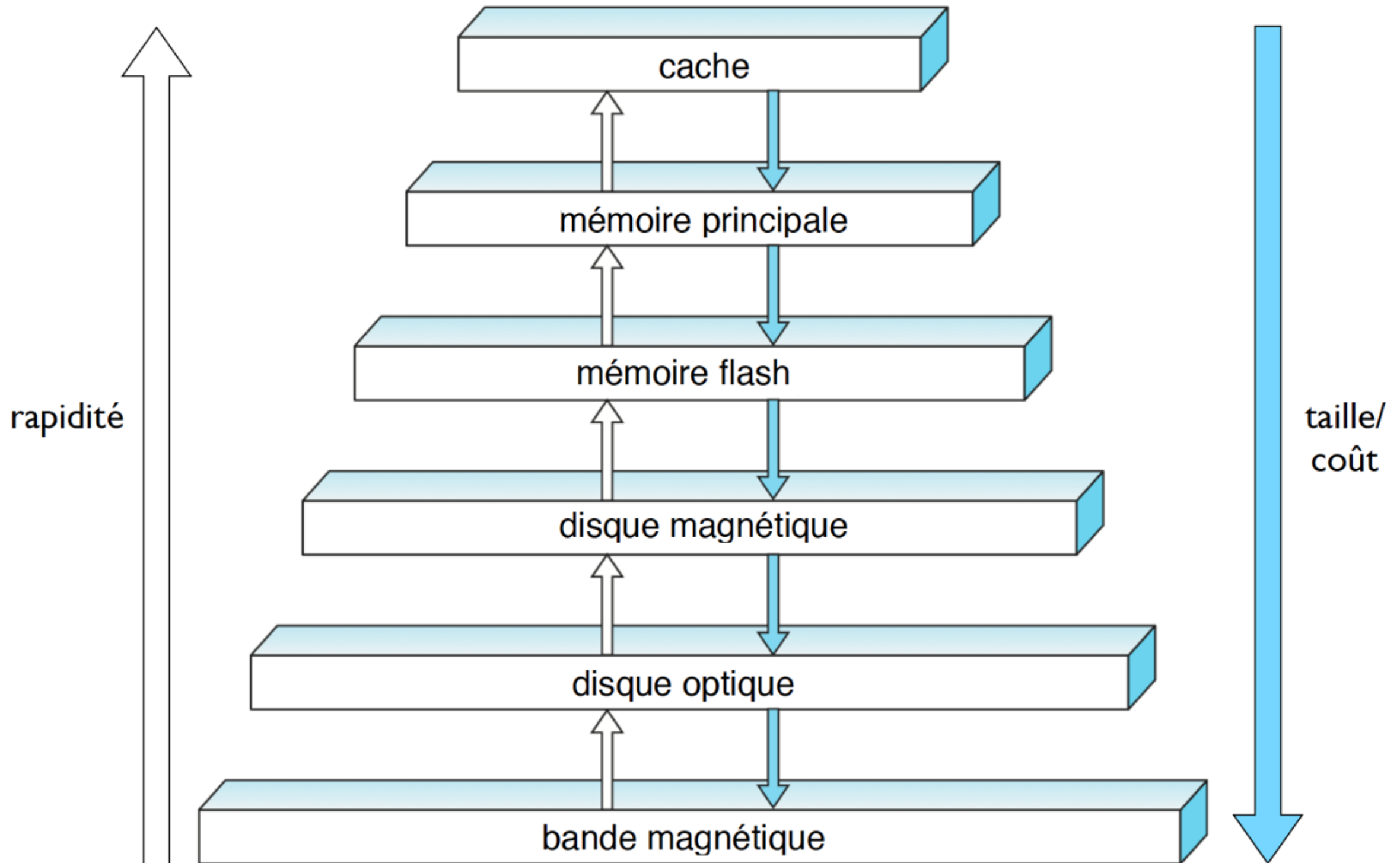




# Stockage des données

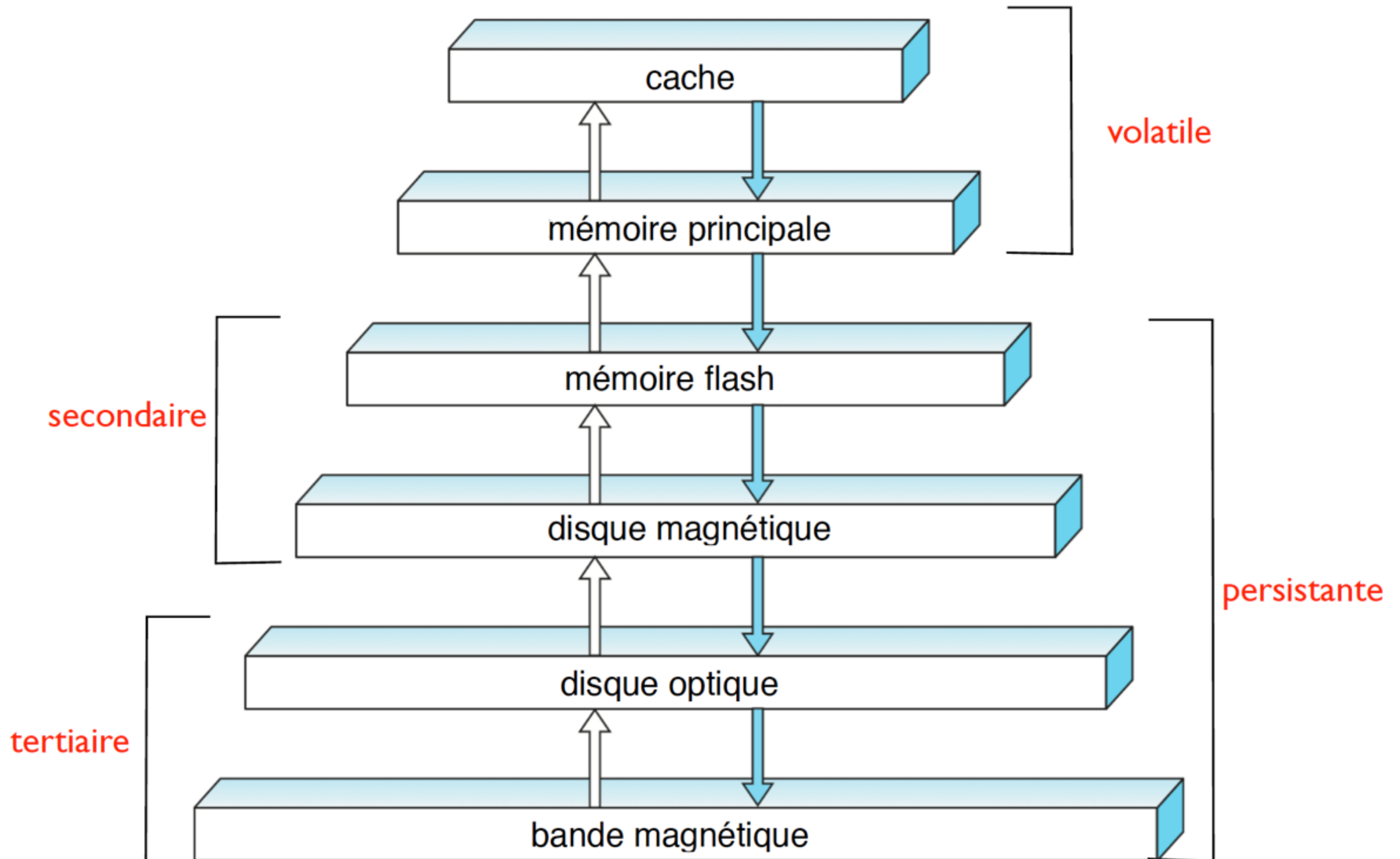
- Le SGBD stocke les données dans une mémoire **persistante**
  - ▶ **mémoire persistante**: les contenus sont préservés même en l'absence d'alimentation.
  - ▶ **mémoire volatile** : perd son contenu quand l'alimentation s'arrête
- Le SGBD interagit avec plusieurs types de mémoires pour gérer/traiter les données

# Hiérarchie de mémoire





# Hiérarchie de mémoire



# Hiérarchie de mémoire

- **Cache** – la plus rapide et la plus coûteuse; volatile  
Maintient les infos les plus récemment utilisées
- **Mémoire principale (RAM):**
  - ▶ volatile
  - ▶ accès rapide (10s à 100s de nanosecondes)
  - ▶ en générale trop petite (ou trop coûteuse) pour stocker la BD entière
    - en train de changer - on commence à parler de mémoires principales suffisamment grandes...
    - la BD doit de toute façon être également stockée en mémoire persistante
  - ▶ taille jusqu'à quelques Gigaoctets aujourd'hui
    - la taille augmente (et le coût diminue) d'un facteur 2 tous les deux/trois ans à peu près





# Hiérarchie de mémoire

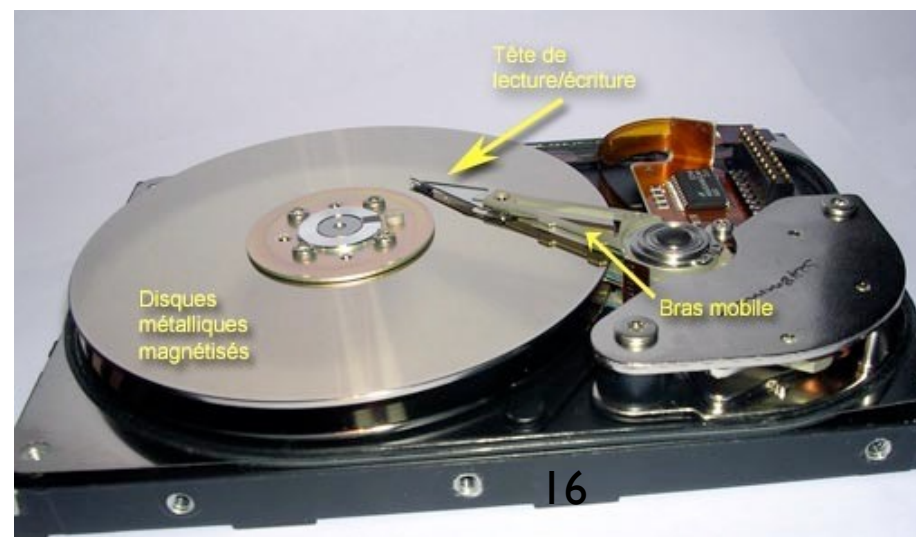
- **Mémoire flash** (*Electrically erasable programmable read-only memory*)
  - ▶ persistante
  - ▶ lectures presque aussi rapides qu'en mémoire principale
  - ▶ écritures plus lentes (quelques microsecondes)
  - ▶ largement utilisée dans les systèmes embarqués tels que cameras numériques, téléphones, clefs USB
  - ▶ remplace parfois aussi le disque



# Hiérarchie de mémoire

- **Disque magnétique**

- ▶ persistant, disque à écriture/ lecture magnétique
  - des pannes de disque peuvent détruire les données mais c'est rare
- ▶ support principal de stockage long-terme : stocke typiquement la BD entière
- ▶ accès beaucoup plus lent que la mémoire principale (millisecondes) : 10 000 à 100 000 fois plus lent que la RAM !!
- ▶ accès direct – il est possible de lire/écrire du disque dans n'importe quel ordre (à la différence de la bande magnétique)
- ▶ taille typiquement jusqu'à quelques teraoctets
  - augmente d'un facteur de 2 à 3 tous les deux ans environ





# Hiérarchie de mémoire

- **Disque optique**

- ▶ persistent, lectures/écritures optiques au laser
- ▶ CD-ROM (640 MB) et DVD (4.7 to 17 GB) les plus populaires
- ▶ disques Blu-ray : 27 GB à 54 GB
- ▶ les types *Write-one, read-many* (WORM) sont utilisés pour archivage (CD-R, DVD-R, DVD+R)
- ▶ Des versions à plusieurs écritures sont aussi disponibles (CD-RW, DVD-RW, DVD+RW, et DVD-RAM)
- ▶ lectures et écritures plus lente qu'avec le disque magnétique
- ▶ systèmes "Juke-box" pour le stockage de grande quantités de données



# Hiérarchie de mémoire

- **Bande magnétique**

- ▶ persistant, utilisé principalement pour *backup* et pour archivage
- ▶ accès séquentiel – beaucoup plus lent que le disque
- ▶ très grande taille par unité (bande de 40 à 300 GB disponibles)
- ▶ beaucoup moins chère que le disque
- ▶ “Juke-boxes” de bandes pour stocker des quantités massives de données
  - de centaines de terabytes (1 terabyte =  $10^{12}$  bytes)  
à plusieurs petabytes (1 petabyte =  $10^{15}$  bytes)





# Stockages des données et types de mémoires

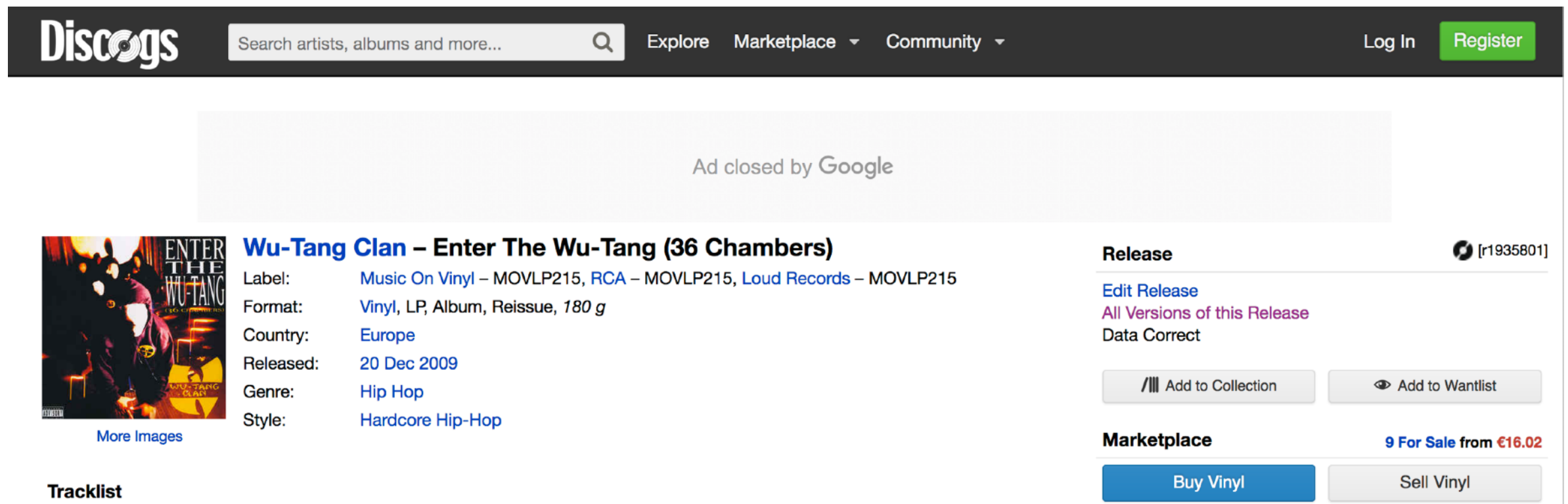
- Une BD est stockée en général en **mémoire secondaire** (disques magnétiques, flash)
  - ▶ raisons :
    - persistance, stockage de grande quantité à coût raisonnable
    - accès direct et en temps réel (par rapport à la mémoire tertiaire)
  - ▶ aujourd'hui BD très volumineuse  $\Rightarrow$  un seul disque ne suffit pas de plus : besoin de fiabilité des données
    - RAID: *Redundant Arrays of Independent Disks*
    - à plus grande échelle : BD distribuée sur le réseau (*Cloud*)
- Les données sont déplacées du disque en **mémoire principale** pour l'accès, puis réécrites sur disque pour le stockage
- Les données plus modifiées et rarement lues sont transférées en **mémoire tertiaire** (disques optiques, bandes magnétiques) pour l'archivage

# Stockages des données et types de mémoires

- Une BD est stockée en général en **mémoire secondaire** (disques magnétiques, flash)
  - ▶ raisons :
    - persistance, stockage de grande quantité à coût raisonnable
    - accès direct et en temps réel (par rapport à la mémoire tertiaire)
  - ▶ aujourd'hui BD très volumineuse ⇒ un seul disque ne suffit pas de plus : besoin de fiabilité des données
    - RAID: *Redundant Arrays of Independent Disks*
    - à plus grande échelle : BD distribuée sur le réseau (*Cloud*)
- Les données sont déplacées du disque en **mémoire principale** pour l'accès, puis réécrites sur disque pour le stockage
- Les données plus modifiées et rarement lues sont transférées en **mémoire tertiaire** (disques optiques, bandes magnétiques) pour l'archivage



# Discogs : une petite BD



The screenshot shows the Discogs website interface. At the top is a navigation bar with the Discogs logo, a search bar, and links for Explore, Marketplace, and Community. A 'Log In' link and a green 'Register' button are on the right. Below the navigation bar is a grey box with the text 'Ad closed by Google'. The main content area features the album 'Wu-Tang Clan – Enter The Wu-Tang (36 Chambers)'. On the left is the album cover. To the right of the cover are the album details: Label (Music On Vinyl, RCA, Loud Records), Format (Vinyl, LP, Album, Reissue, 180 g), Country (Europe), Released (20 Dec 2009), Genre (Hip Hop), and Style (Hardcore Hip-Hop). Below the cover is a 'More Images' link. To the right of the album details is a 'Release' section with a release ID [r1935801], links for 'Edit Release' and 'All Versions of this Release', and a 'Data Correct' status. Below this are two buttons: 'Add to Collection' and 'Add to Wantlist'. Further down is a 'Marketplace' section showing '9 For Sale from €16.02' and two buttons: 'Buy Vinyl' and 'Sell Vinyl'. At the bottom left of the screenshot is a 'Tracklist' section.

- Une petite base de données ~ 80Go (0,08To)
- Un catalogue de plus de 8 800 000 références, 5 000 000 d'artistes et 1 000 000 de labels
- Plus de 369 000 contributeurs

# Discogs

Discogs

Search artists, albums and more...



Explore

Marketplace

Community

Log In

Register

## Genre

- Rock 408,615
- Electronic 398,818
- Pop 210,781
- Folk, World, & Country 108,200
- Funk / Soul 98,565
- All

## Style

- House 67,548
- Pop Rock 60,747
- Punk 49,286
- Techno 45,613
- Experimental 40,931
- All

## Format

- Vinyl 740,855
- Album 395,449
- CD 348,999
- LP 302,468
- 7" 225,873
- All

All Release 8,889,969 Master Artist 5,077,218 Label 1,066,193

Search Marketplace

1 - 50 of 1,231,091 < Prev Next >

Sort Title, Z-A



Zzzzz...  
Skull Kontrol



Zzzzz.....  
Tor.Men.To



ZZZ  
Sun Bather



...Zzzipp!  
Tab Two



Zzzero  
Various



ZZZ EP  
Dubble D presents Mo...



Zzz  
sanmi



Zzz  
J. Mono



Zzyzx Rd.  
Stone Sour



Zzyzx  
Zeromancer





# Blocs et enregistrements

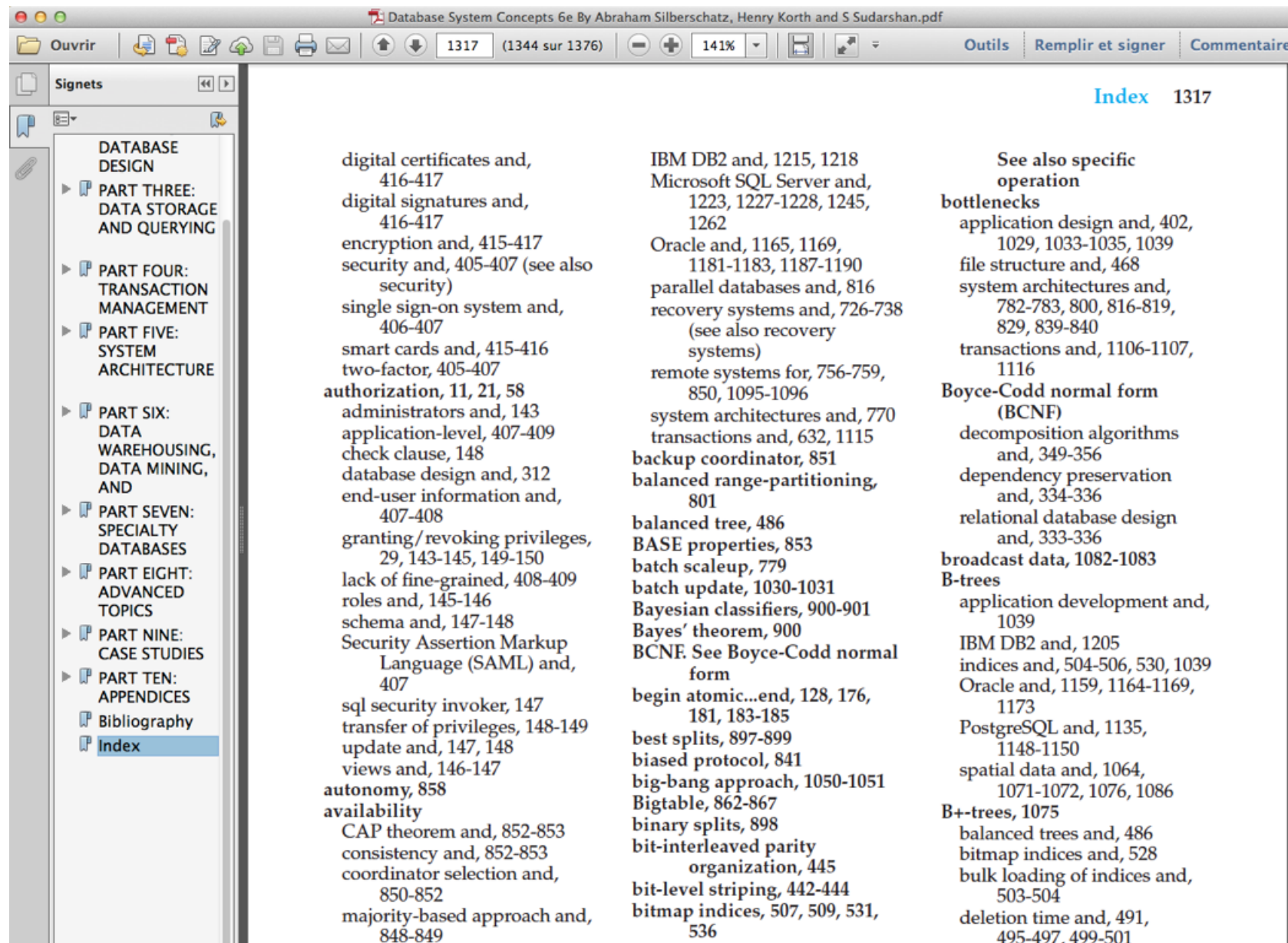
Une BD contient des tableaux dont on appelle chaque ligne «enregistrement», par exemple :

ALBUM	Num_Album	Nom_Artiste	Nom_Album	Genre	Année
	42	Wu-Tang Clan	Enter The Wu-Tang (36 Chambers)	Hip Hop	1993
	69	Mykki Blanco	Cosmic Angel : The Illuminati Prince/ss		2012
	555	Five Finger Posse	Hood Goth	Hip Hop	2016
	3	Mason vs. Princess Superstar	Perfect (Exceeder)	Electro	2006

- Ces enregistrements sont rangés dans des «pages» stockées dans des «blocs» du disque
- Chaque enregistrement a une **adresse physique** en mémoire
- Le bloc est l'**unité de transfert** avec la mémoire principale

# Indexes

- Afin de déplacer un enregistrement du disque à la mémoire principale, il nous faut son adresse physique.
- Les SGBD rangent ces adresses physiques grâce à des «**indexes**»





# Efficacité de la recherche

- **Objectif** : être capable de trouver efficacement les enregistrements d'un tableau avec une valeur particulière pour un certain champ
  - **Modèle de coût** pour évaluer l'efficacité :
    - ▶ **Mesure du coût** : nombre de lectures/écritures de blocs du disque
    - ▶ **Justification** :
      - coût d'un accès au disque (lecture/ écriture d'un bloc) : millisecondes
      - coût d'un accès en mémoire centrale : nanosecondes (milliardièmes de seconde)
- ⇒ toutes les opérations de traitement d'un bloc en mémoire centrale ont un coût négligeable

# Recherche séquentielle

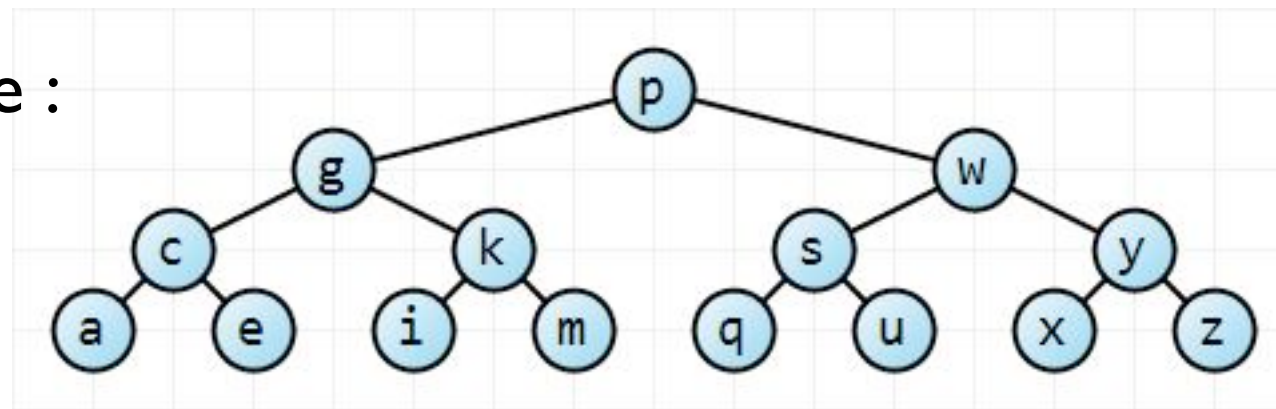
- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»
- **Approche naïve** : parcourir tous les blocs de Album triés par ordre alphabétique
  - ▶ **Coût** : environ  $\lceil n r / B \rceil$ 
    - $n$  : nombre d'enregistrements dans Album,
    - $r$  : taille d'un enregistrement
    - $B$  : taille d'un bloc
- Si  $n$  est très grand ça peut être prohibitif
- **Exemple** :  $n = 8\,000\,000$ ,  $r = 100$ ,  $B = 8000$ 
  - ▶ **Coût** : environ  $100\,000 * 10$  millisecondes =  $1\,000$  secondes =  $16,66$  mins !!!!
    - $(8\,000\,000 * 100) / 8000 = 100\,000$  et  $1$ seconde =  $1\,000$  millisecondes
    - $10$  millisecondes : temps typique d'accès au disque

# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple :



- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.

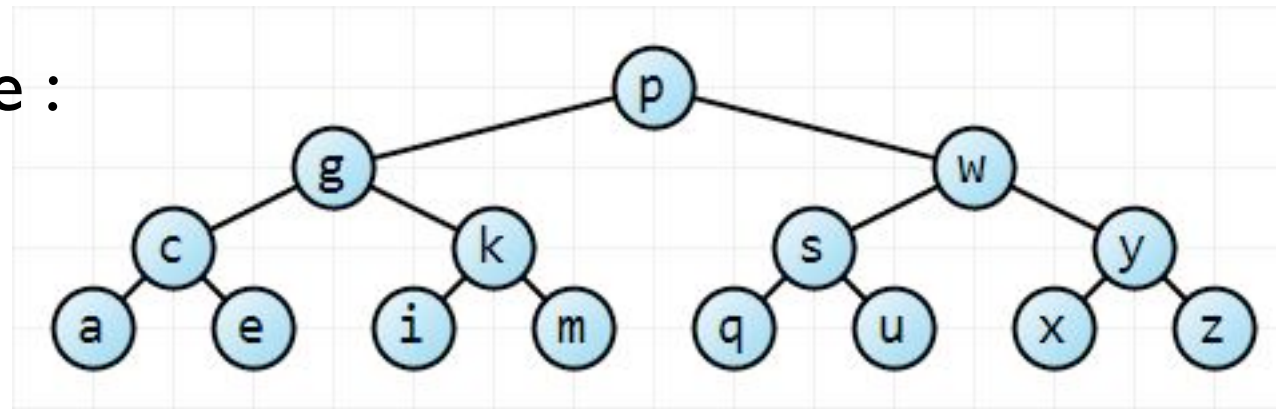


# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple :



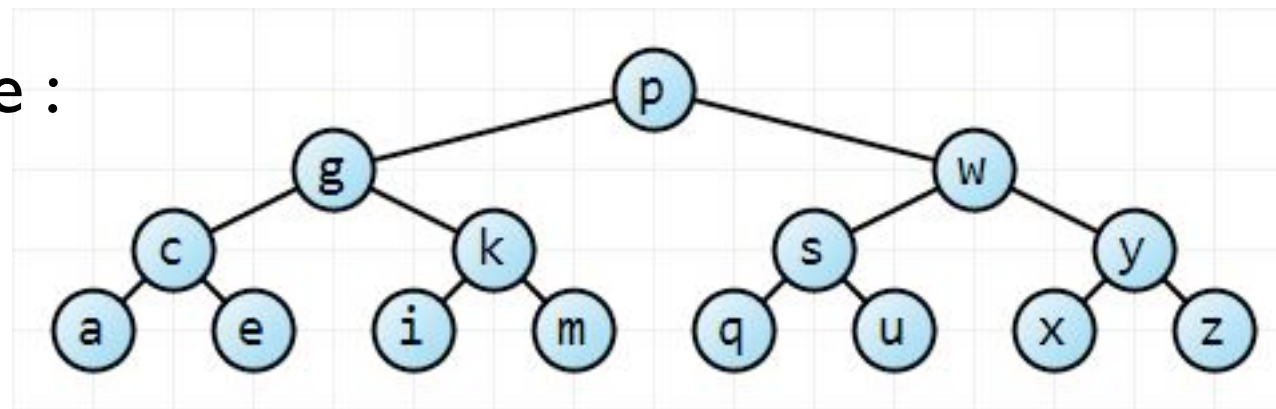
- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p

# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple :



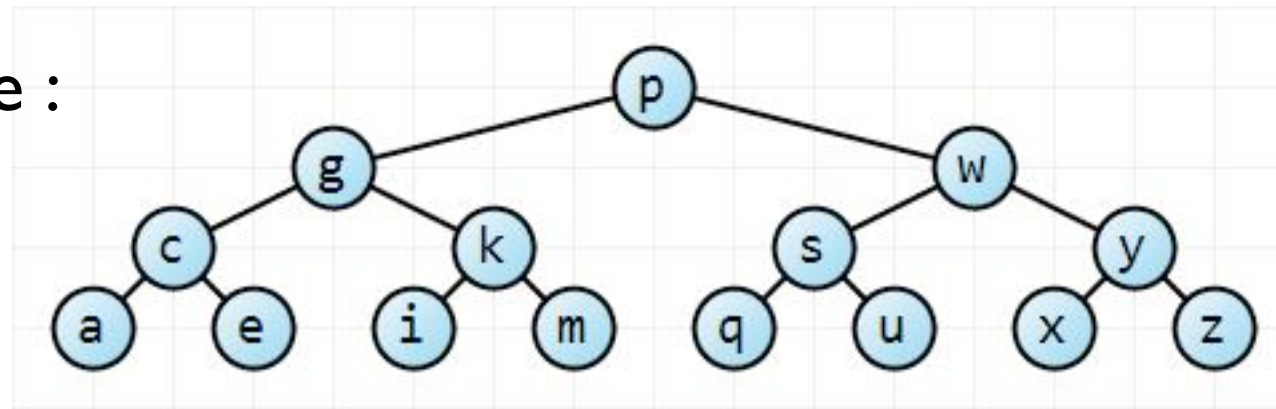
- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p
- Est-ce que x est égal à p, avant p, ou après p dans l'ordre alphabétique ? après

# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple :



- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p
- Est-ce que x est égal à p, avant p, ou après p dans l'ordre alphabétique ? après
- Toute la partie allant jusqu'à p ne nous intéresse donc plus

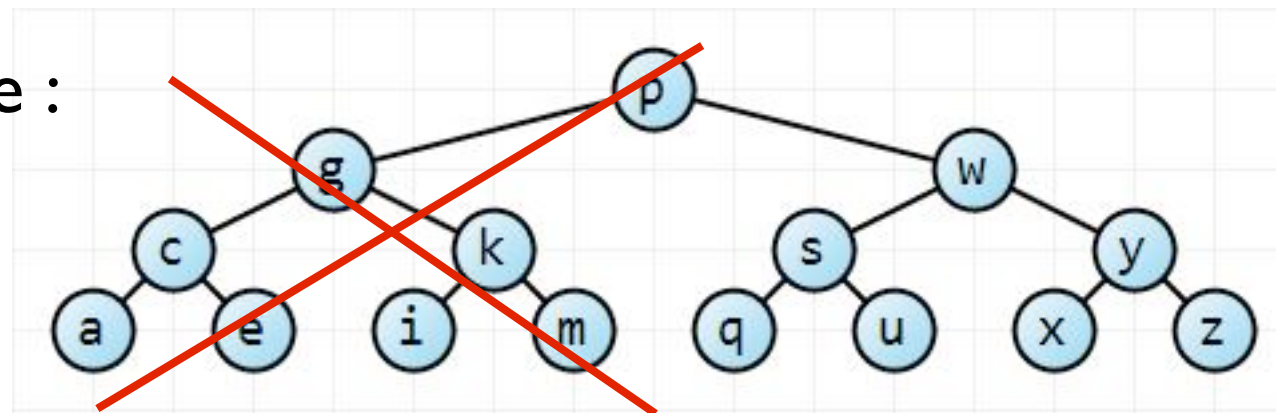


# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple :



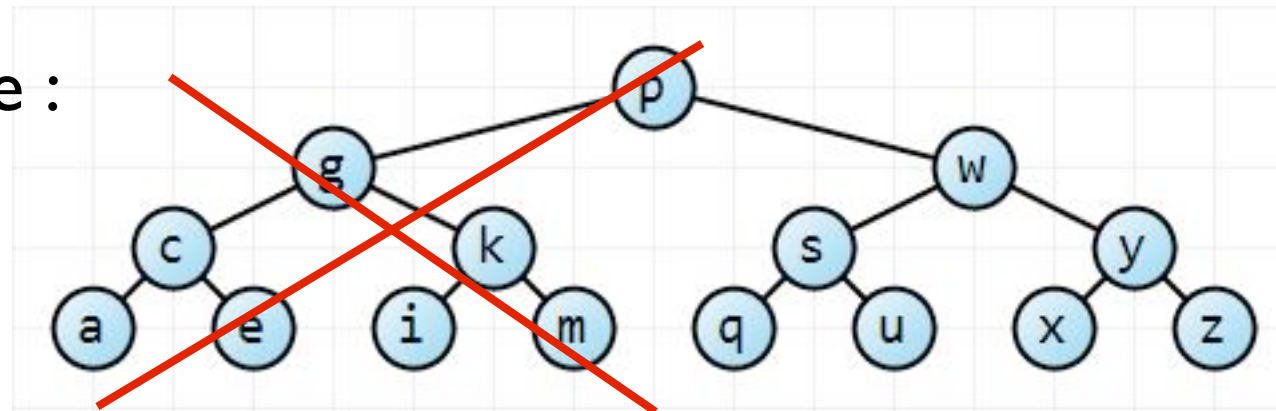
- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p
- Est-ce que x est égal à p, avant p, ou après p dans l'ordre alphabétique ? après
- Toute la partie allant jusqu'à p ne nous intéresse donc plus

# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple :



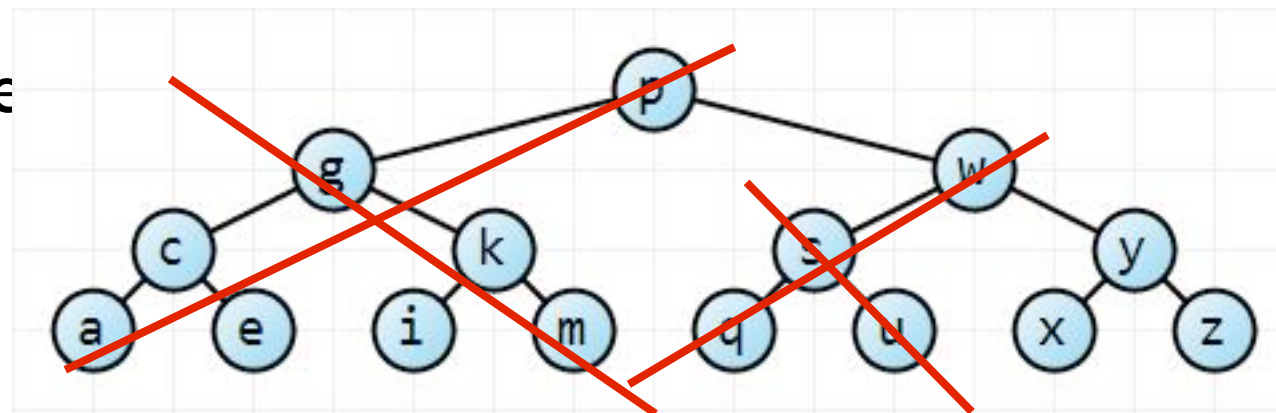
- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p
- Est-ce que x est égal à p, avant p, ou après p dans l'ordre alphabétique ? après
- Toute la partie allant jusqu'à p ne nous intéresse donc plus
- On répète alors le même raisonnement pour le fragment de l'alphabet allant de q à z.

# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

- Un exemple plus simple



- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p
- Est-ce que x est égal à p, avant p, ou après p dans l'ordre alphabétique ? après
- Toute la partie allant jusqu'à p ne nous intéresse donc plus
- On répète alors le même raisonnement pour le fragment de l'alphabet allant de q à z. Etc.

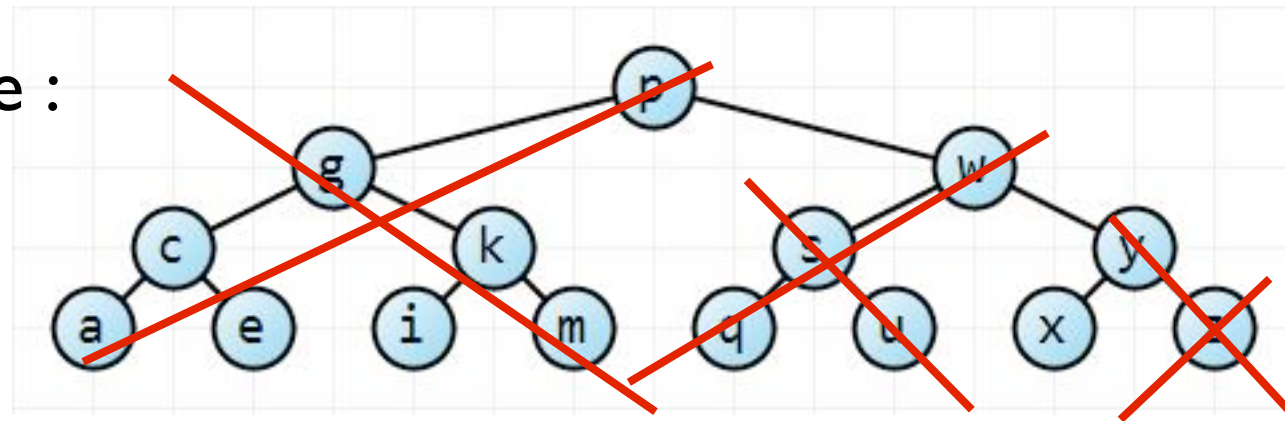


# Recherche dichotomique

- **Objectif** : trouver l'enregistrement correspondant à «Enter the Wu-Tang (36 Chambers)»

- **Approche naïve 2** : diviser pour régner

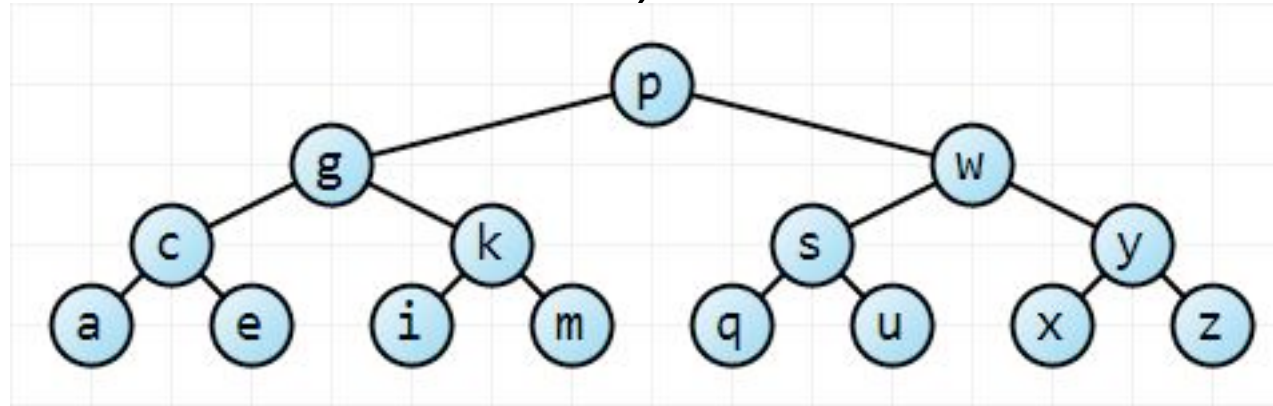
- Un exemple plus simple :



- On cherche la lettre x parmi 15 lettres ordonnées de manière alphabétique.
- Quelle lettre se trouve au milieu dans l'ordre alphabétique ? p
- Est-ce que x est égal à p, avant p, ou après p dans l'ordre alphabétique ? après
- Toute la partie allant jusqu'à p ne nous intéresse donc plus
- On répète alors le même raisonnement pour le fragment de l'alphabet allant de q à z. Etc.

# Recherche dichotomique

- Approche naïve 2 : diviser pour régner
- (Oublions provisoirement les blocs.)



- ▶ Coût : pas plus de  $\lceil \log_2 15 \rceil = 4$  étapes
  - 15 : nombre de lettres,
  - $\log_2$  : le logarithme en base 2 de 15 est la puissance à laquelle il faut élever 2 pour arriver à 15
  - $\lceil \log 15 \rceil$  la partie entière par excès (ou «plafond») de  $\log_2 15$
  - $\log_2 15 \sim 3,9$  ( $2*2*2 = 8$  et  $2*2*2*2=16$ , or  $8 < 15 < 16$ )
- Beaucoup mieux que l'approche 1, qui aurait pris 14 étapes !

# Recherche dichotomique

- Remarque: ce coût peut être élevé pour un grand nombre d'enregistrements
  - ▶ Exemple: 8 000 000 d'enregistrements
    - supposer par exemple 80 enregistrements par bloc
    - $\Rightarrow$  100 000 blocs
    - $\Rightarrow$  un recherche demande  $\lceil \log_2 100\,000 \rceil = 17$  accès aux blocs du disque (en effet  $65\,536 < 100\,000 < 131\,072$ )
    - temps typique d'accès à un bloc : 10 ms  $\Rightarrow$  170 ms = 0,17 s par recherche
    - Beaucoup beaucoup beaucoup mieux que l'approche 1 !!! (16,66 mins)
    - Mais quand même.... seulement 5 recherches par seconde...



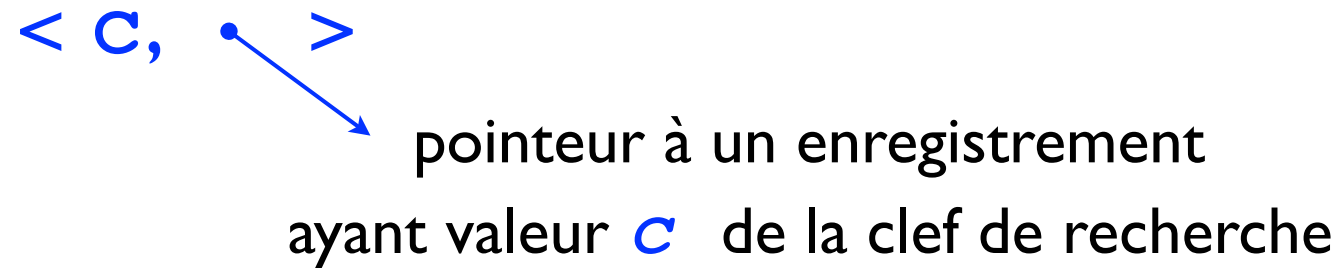
# Indexation

- **Approche 3** : Utiliser une structure de données auxiliaire
- On garde si possible cette structure en mémoire principale
- Permet d'accéder facilement à l'enregistrement cherché
  - Donne le bloc où se trouve l'enregistrement cherché
  - Donne l'adresse dans le bloc
- - Lecture uniquement du (ou des) bloc(s) utile(s)

**Index** : une structure physique auxiliaire construite sur une clef de recherche qui permet d'accéder "presque directement" aux enregistrements ayant une valeur particulière de la clef de recherche

# Index

- Un index sur une clef de recherche est une collection de couples



$\langle 10101, \bullet \rangle$	→	10101	Srinivasan	Comp. Sci.	65000
$\langle 12121, \bullet \rangle$	→	12121	Wu	Finance	90000
$\langle 15151, \bullet \rangle$	→	15151	Mozart	Music	40000
$\langle 22222, \bullet \rangle$	→	22222	Einstein	Physics	95000
$\langle 32343, \bullet \rangle$	→	32343	El Said	History	60000
$\langle 33456, \bullet \rangle$	→	33456	Gold	Physics	87000
$\langle 45565, \bullet \rangle$	→	45565	Katz	Comp. Sci.	75000
$\langle 58583, \bullet \rangle$	→	58583	Califieri	History	62000
$\langle 76543, \bullet \rangle$	→	76543	Singh	Finance	80000
$\langle 76766, \bullet \rangle$	→	76766	Crick	Biology	72000
$\langle 83821, \bullet \rangle$	→	83821	Brandt	Comp. Sci.	92000
$\langle 98345, \bullet \rangle$	→	98345	Kim	Elec. Eng.	80000

index

fichier de données

# Indexation : principe

- Rechercher des enregistrements ayant une valeur  $c$  de la clef de recherche :
  - ▶ parcourir l'index
  - ▶ si un couple  $\langle c, \cdot \rangle$  est trouvé suivre le pointeur
    - charger en mémoire principale le bloc contenant l'enregistrement
- Nombres d'entrées dans l'index : au plus le nombre  $n$  d'enregistrements dans le fichier de données
- Mais un couple  $\langle c, \cdot \rangle$  occupe beaucoup moins de place qu'un enregistrement!



# Indexation : principe

- Coût de la recherche : dépend de l'implémentation de l'index mais au pire

environ  $\lceil ni/B \rceil + 1$  (  $i$  : taille d'une entrée de l'index )  
(  $n$  : nombre d'enregistrements )

$B$  : nombre de blocs occupés par l'index      un accès au fichier de données

$$\lceil ni/B \rceil \ll \lceil nr/B \rceil \quad \text{puisque } i \ll r$$

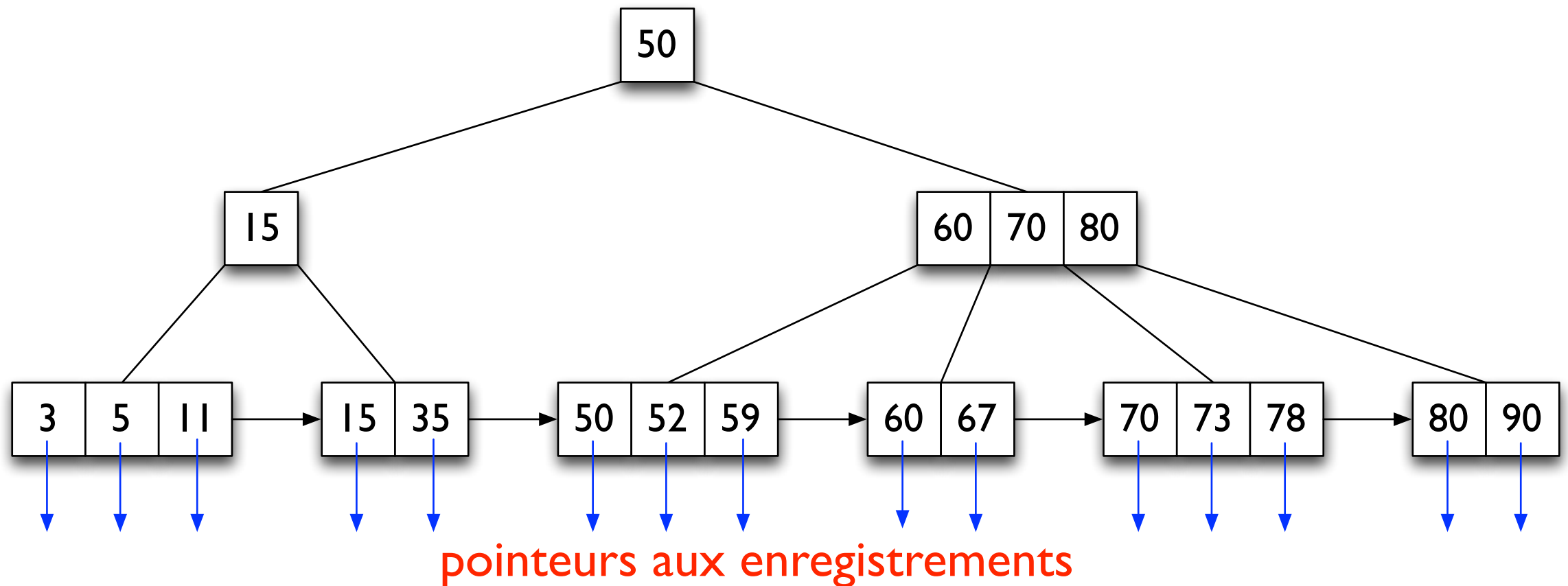
(  $r$  : taille d'un enregistrements )

- De plus il existe des implémentations efficaces des indexes :  
le coût de la recherche d'un couple  $\langle c, \cdot \rangle$  dans l'index peut passer de  $\lceil ni/B \rceil$  à "presque" constant !

# Implementation des indexes

- Différentes structures de données sont adaptées à l'implémentation d'indexes
- Les plus fréquentes :
  - ▶ arbres B+, arbres B
  - ▶ Hash
  - ▶ Bitmap
- Nous allons brièvement présenter les arbres B+.
- Vous verrez toutes ces structures en détail dans le cours de Bases de Données avancées de MI :-)

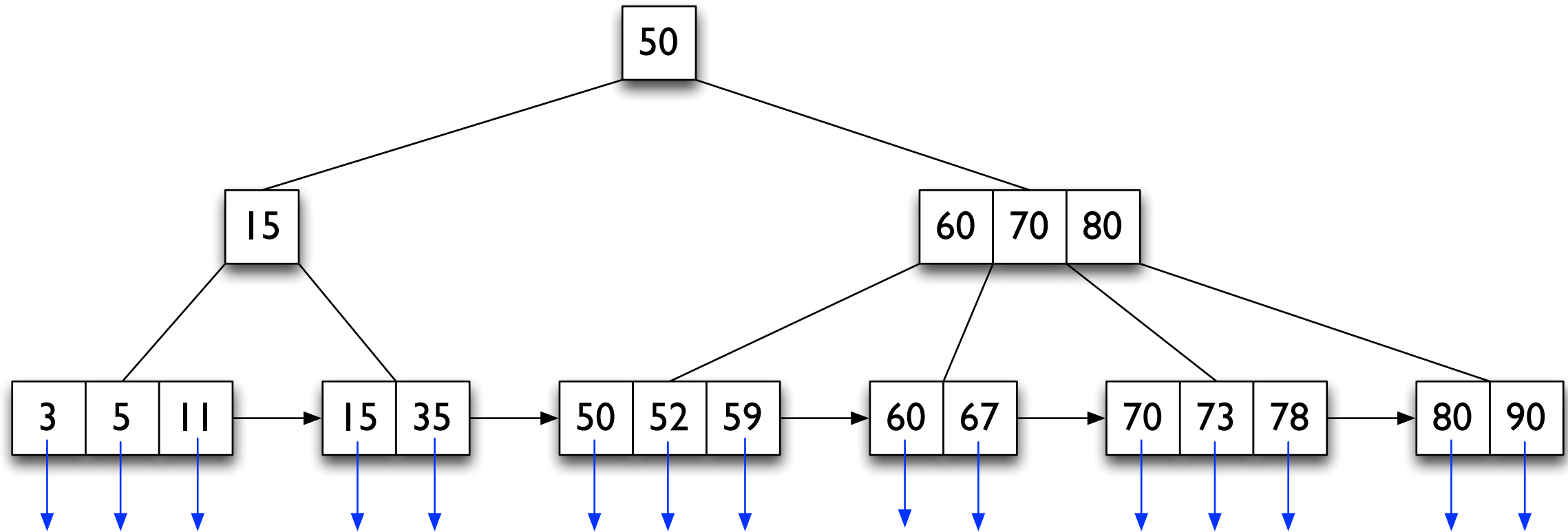
# Arbres B<sup>+</sup> : introduction



- Structure arborescente (index multi-niveau)
  - ▶ un noeud de l'arbre = un bloc du disque
  - ▶ les noeuds feuille contiennent les couples  $\langle \text{clef}, \text{pointeur} \rangle$  de l'index
  - ▶ les noeuds internes contiennent des doublons des clefs (appelés *balises*) pour orienter la recherche
  - ▶ tous les noeuds feuille sont liés en liste chaînée, de gauche à droite (raison : permettre un parcours linéaire de la table par clef de recherche)



# Arbres B<sup>+</sup> : introduction



- Si un noeud (interne ou feuille) contient  $k$  clefs, il contient  $k+1$  pointeurs
  - ▶ dans un noeud interne : pointeurs aux sous-arbres
  - ▶ dans un noeud feuille : pointeurs aux enregistrements et à la prochaine feuille

Soit  $n$  l'entier maximal tel que un bloc du disque contienne  $n$  clefs et  $n+1$  pointeurs

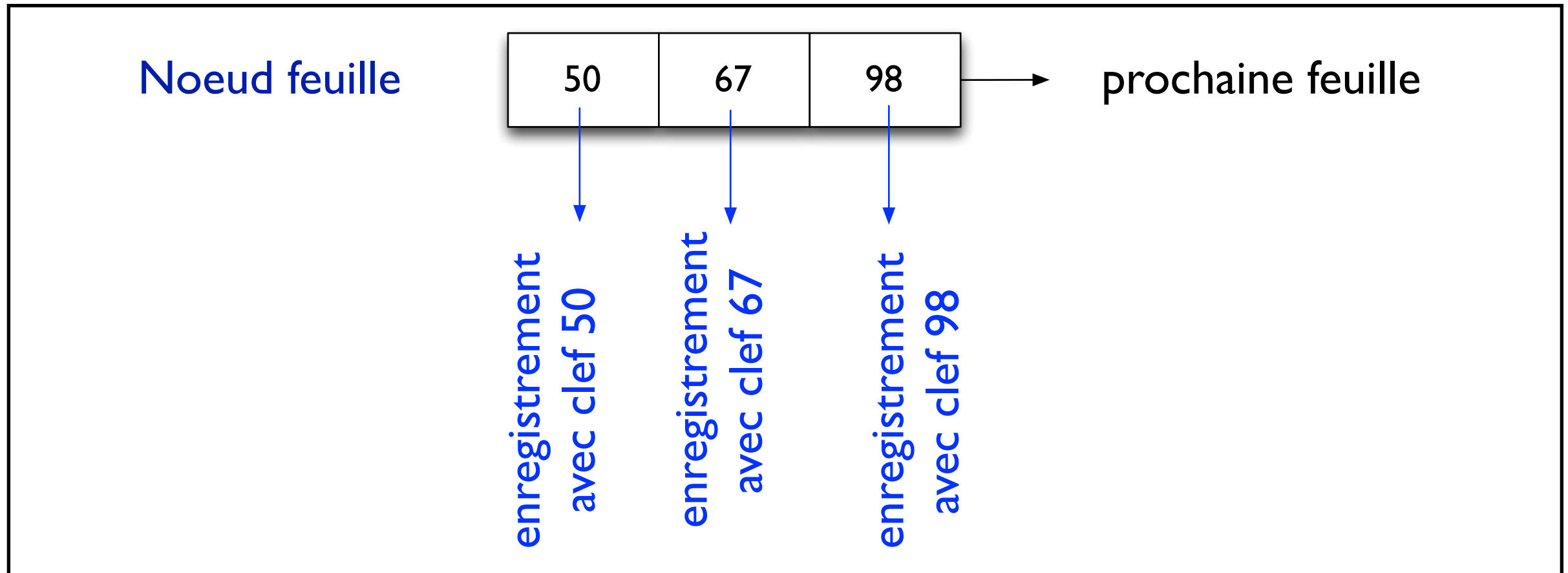
⇒  $n$  est le nombre maximal de clefs dans un noeud de l'arbre

# Arbres B<sup>+</sup> : définition

Un arbre B<sup>+</sup> possède les propriétés suivantes :

## I) les feuilles

- ont toutes la même profondeur,
- contiennent des couples  $\langle \text{clef}, \text{pointeur} \rangle$ , triés par clef, et un pointeur à la prochaine feuille



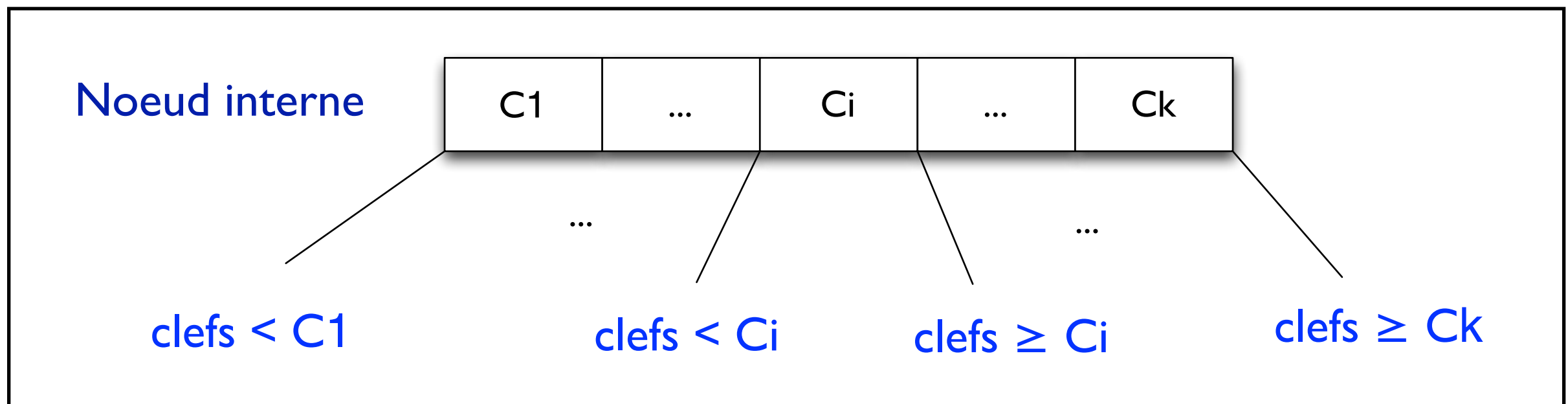
- contiennent entre  $\lceil n/2 \rceil$  et  $n$  clefs (n défini comme au slide précédent)

# Arbres B<sup>+</sup> : définition

Un arbre B<sup>+</sup> possède les propriétés suivantes :

2) les **noeuds internes** contiennent

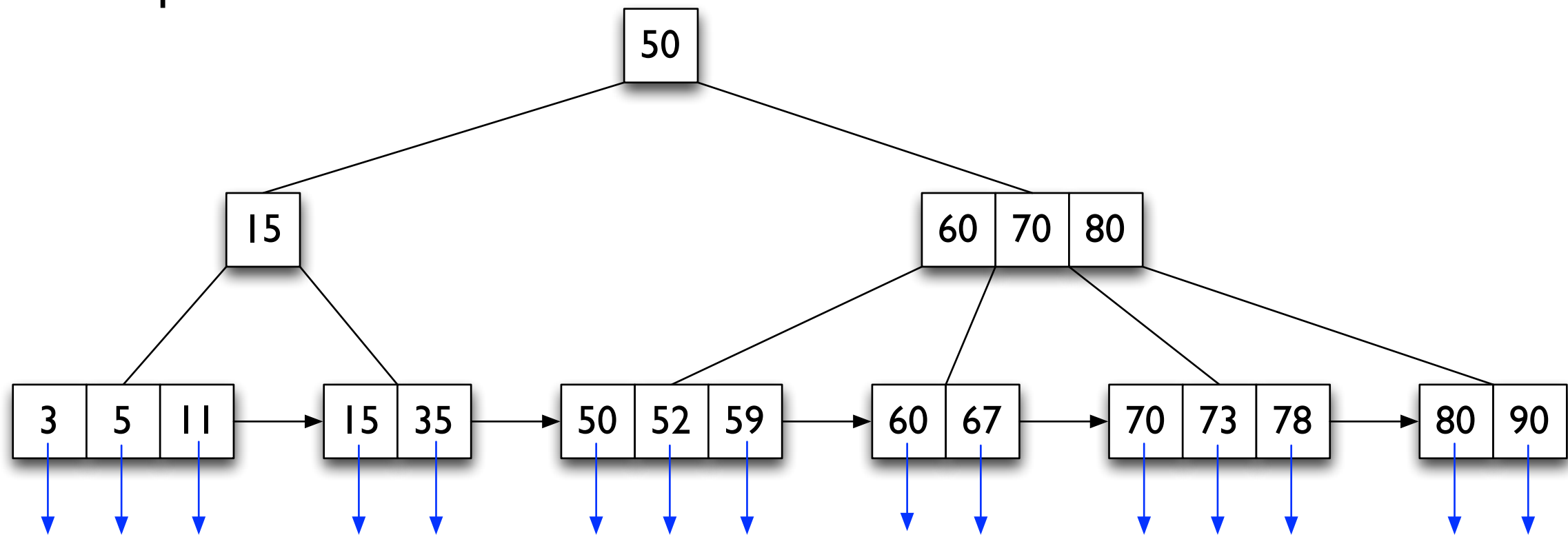
- une suite de clefs triées par ordre croissant
- un nombre de pointeurs à sous-arbres = nombre de clefs + 1, avec la propriété suivante pour chaque clef C<sub>i</sub>:



- tous noeuds internes sauf la racine : entre  $\lfloor n/2 \rfloor$  et  $n$  clefs
- la racine : entre 1 et  $n$  clefs

# Arbres B<sup>+</sup> : exemple

- Exemple avec n= 3

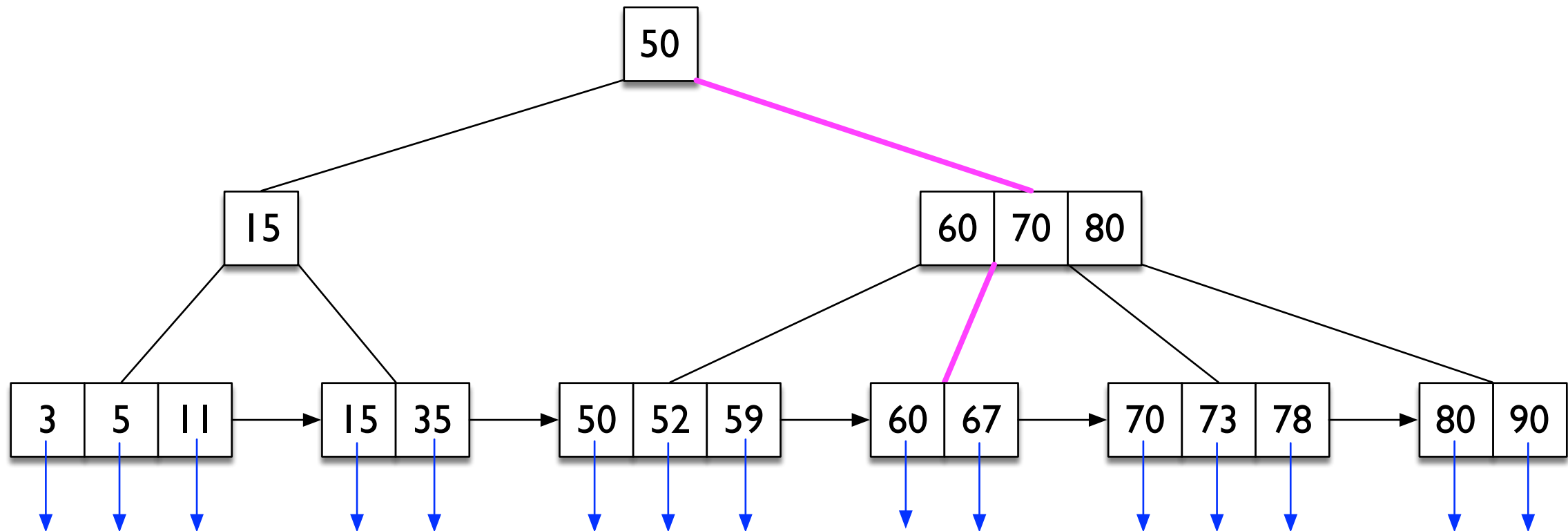


- noeuds feuilles : entre 2 et 3 clefs
- noeuds internes sauf la racine : entre 1 et 3 clefs
- la racine : entre 1 et 3 clefs



# Arbres B<sup>+</sup> : recherche

- Exemple  $c = 67$



- Ici le nombre de noeud maximal depuis la racine jusqu'à une feuille est 3 : on dit que cet arbre est de **hauteur 3**.
- Ici le nombre de noeuds depuis la racine jusqu'à une feuille est toujours le même : on dit que l'arbre est **équilibré**.
- Le nombre de lectures de noeuds pour arriver à un enregistrement est égal à la longueur d'une branche, c'est à dire à la hauteur de l'arbre.

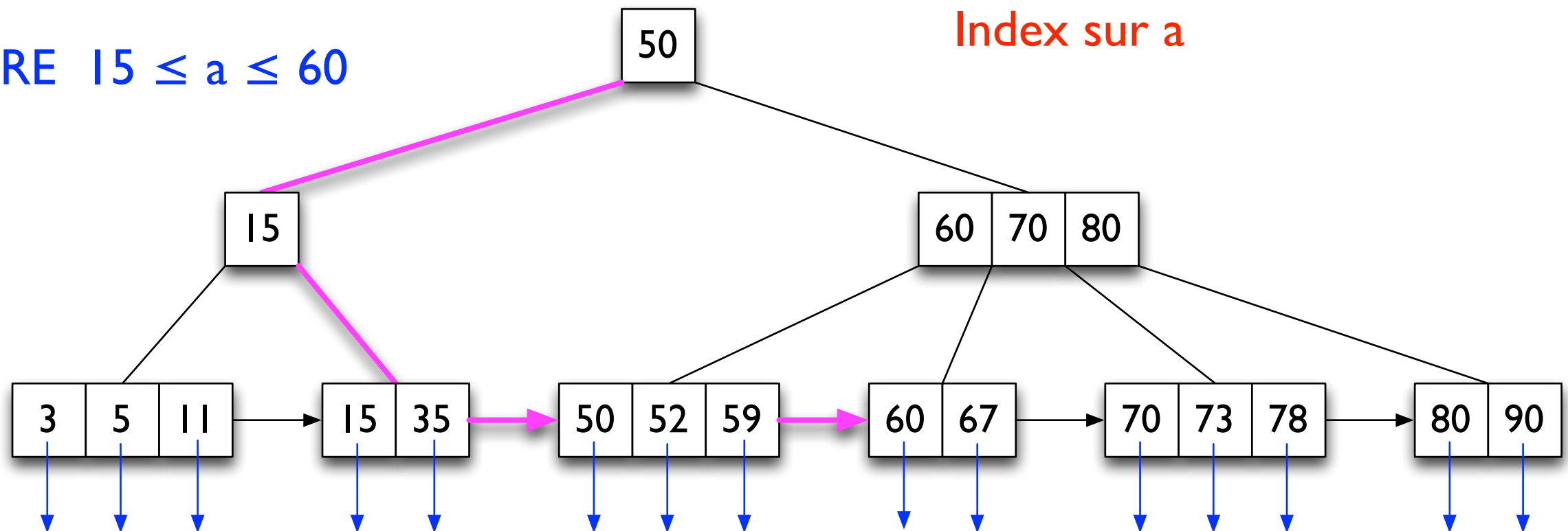
# Arbres B<sup>+</sup> : en pratique

- Très utilisés en pratique
  - ▶ coût de la recherche et du rééquilibrage pratiquement constant pour des index pas trop gros
- Très utilisés surtout pour les **recherches d'intervalle**, i.e.

SELECT \*

FROM Tableau

WHERE  $15 \leq a \leq 60$



# Arbres B<sup>+</sup> : bornes

- But des bornes sur le nombre de clefs:
  - ▶ borne supérieure : pour pouvoir stocker un noeud dans un bloc
  - ▶ borne inférieure :
    - pour la procédure de rééquilibrage de l'arbre après insertion / suppression
    - pour éviter des blocs trop vides (le blocs sont au moins à moitié pleins)

# Arbres B<sup>+</sup> : coût de la recherche

- **Exemple** : 8 000 000 enregistrements, 80 enregistrements par noeud, 41 = nombre minimal d'enregistrement par noeud interne
- **Hauteur de l'arbre** :  $\lceil \log_{41}(8\,000\,000/2 + 1) \rceil = 5$
- $41 * 41 = 1681$
- $1681 * 41 = 68921$
- $68921 * 41 = 2\,825\,761$
- $2\,825\,761 * 41 = 115\,856\,201$
- La recherche dans l'arbre B demande donc 5 accès disques seulement.
  - temps typique d'accès à un bloc : 10 ms  $\Rightarrow$  50 ms = 0,05 s par recherche
- La recherche dichotomique demandait  $\lceil \log_2 100\,000 \rceil = 17$  accès disques.
  - temps typique d'accès à un bloc : 10 ms  $\Rightarrow$  170 ms = 0,17 s par recherche
- Gain de performance très important ! (5 versus 1200 requêtes par seconde)



# Conclusions

## Dimensions du Big Data

