

MON AMI B-TREE

Epitech Summit Lille - 4 février 2025

AGENDA

Les enjeux de la performance

HEAP : le mécanisme de stockage avec PostgreSQL

BTREE : structure d'indexation de la donnée

Conclusion

QUI SOMMES-NOUS ?



Florent
JARDIN

- Consultants en base de données
- Support, formations, conseils
- Animateurs du Meetup PostgreSQL User Group Lille
- Anciens collègues chez Claranet
- Anciens diplômés SUPINFO



Yoann
LA CANCELLERA

PATIENTEZ, S'IL VOUS PLAÎT

10 MS

Temps d'exécution moyen d'une requête SQL
dans des conditions optimales

100 MS

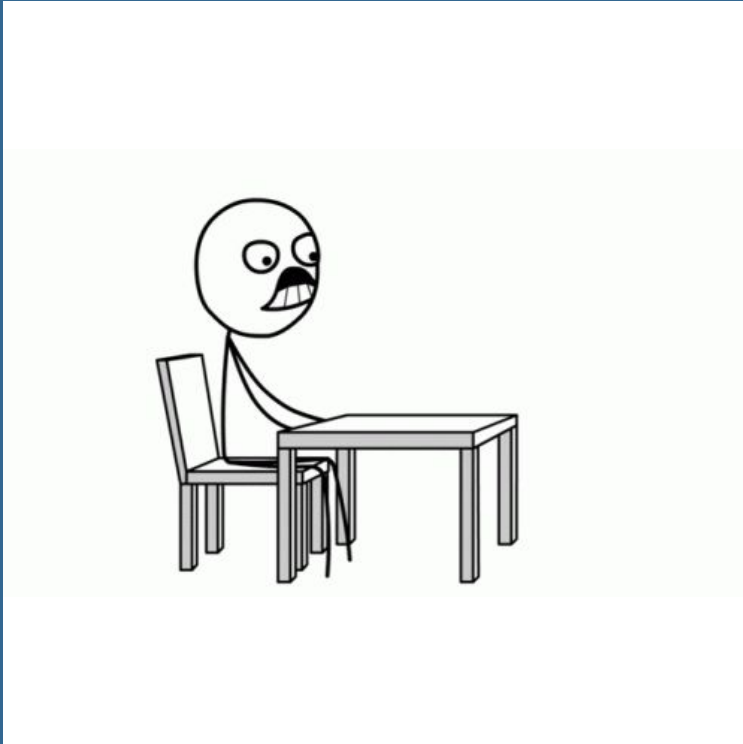
Délai cognitif moyen

pour une expérience utilisateur optimale



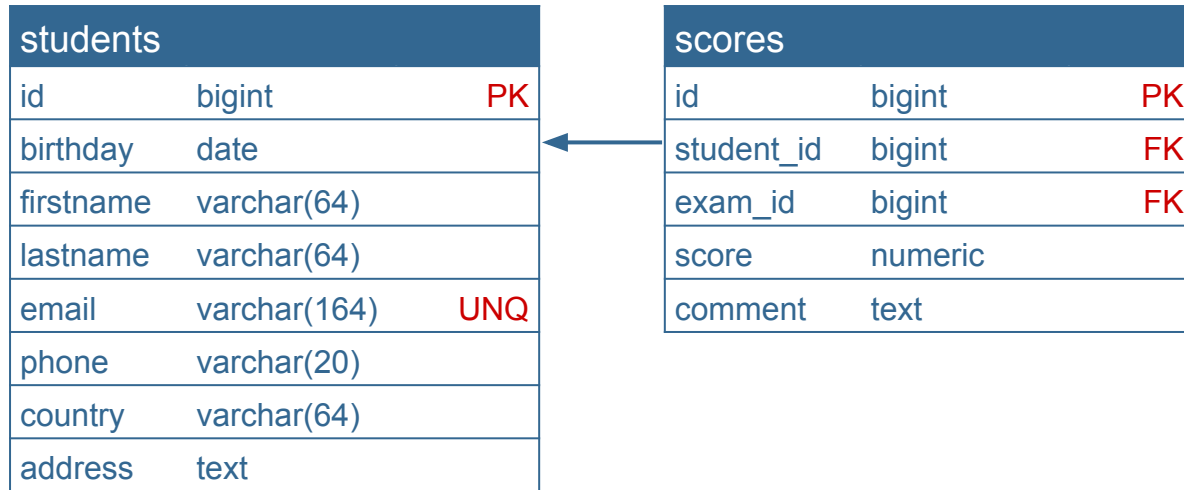
3 SEC

Selon Google, lorsque le temps de chargement
d'une page passe d'une à trois secondes,
le taux de rebond augmente de 32 %.

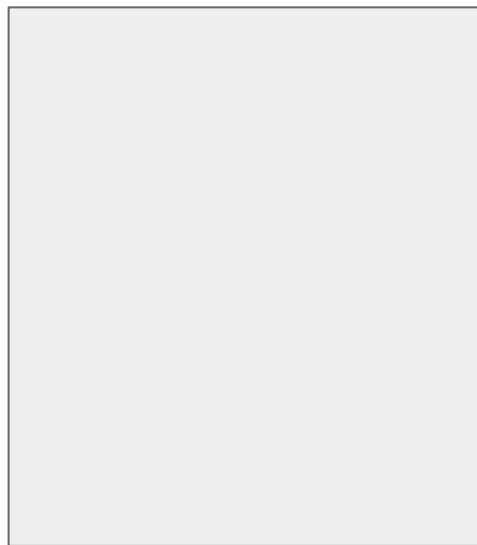


HEAP HEAP HEAP

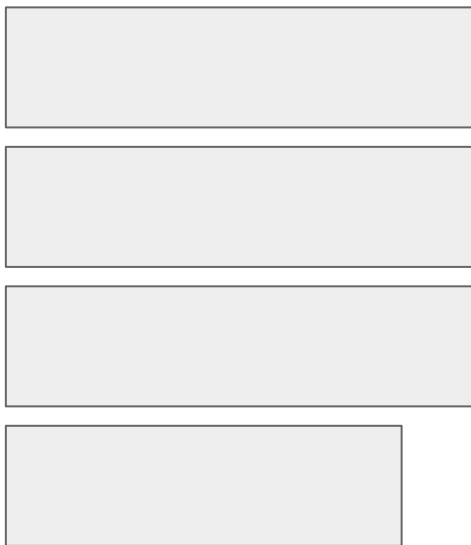
MODÉLISATION DU FIL ROUGE



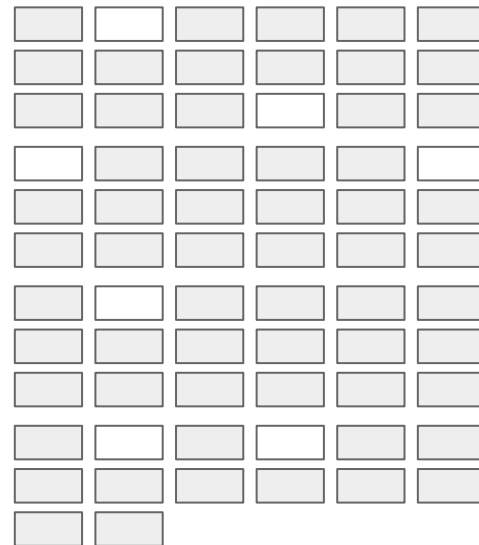
REPRÉSENTATION PHYSIQUE



HEAP Table



Fichiers (1 Go)

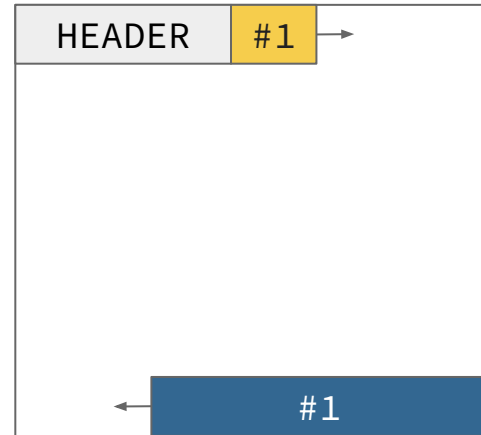


Blocs (8 Ko)

AJOUT D'UNE LIGNE (TUPLE)

```
INSERT INTO scores  
VALUES (1, 200, 10, .87, NULL);
```

Tuple #1
<code>\x010000000000000000c80000000000000000a 000000000000000000b00805700</code>

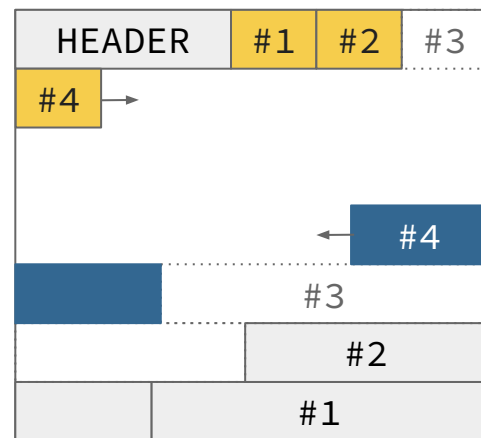


Bloc (8 Ko)

MISE À JOUR D'UNE LIGNE

```
UPDATE scores SET score = 1.0
WHERE student_id = 200
AND exam_id = 30
```

- Copy on Write (CoW)
- Dead Tuple
- Heap Only Tuple (HOT)

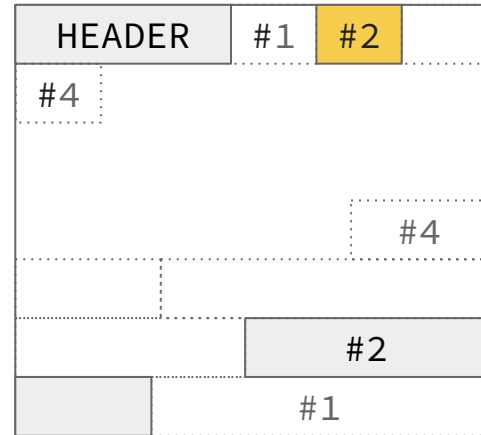


Bloc (8 Ko)

SUPPRESSION D'UNE LIGNE

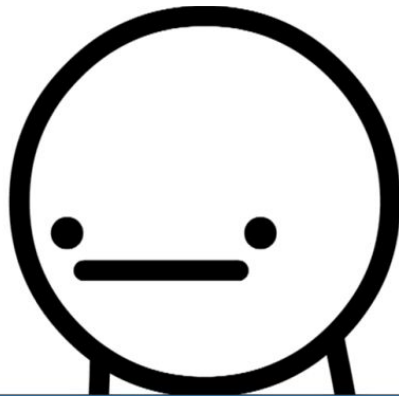
```
DELETE FROM scores  
WHERE student_id = 200
```

- Fragmentation (Bloat)
- VACUUM
- VACUUM FULL



Bloc (8 Ko)

```
SELECT * FROM scores  
WHERE student_id = 200
```



STRUCTURED QUERY LANGUAGE (SQL)

→ Langage déclaratif

Requête SQL

```
SELECT ...  
FROM ...  
JOIN ...  
WHERE ...
```

Résultat

1	@email.com
2	@email.com
3	@email.com
4	@email.com
5	@email.com

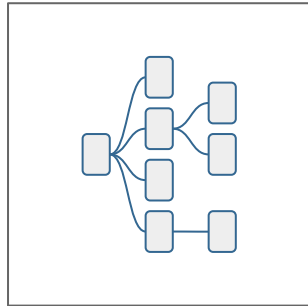
STRUCTURED QUERY LANGUAGE (SQL)

→ Langage déclaratif

Requête SQL

```
SELECT ...  
FROM ...  
JOIN ...  
WHERE ...
```

Arbre



Résultat

1	@email.com
2	@email.com
3	@email.com
4	@email.com
5	@email.com

Parser

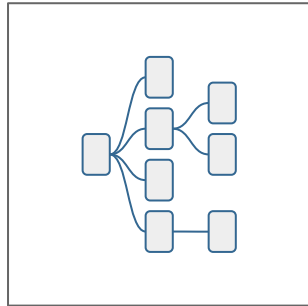
STRUCTURED QUERY LANGUAGE (SQL)

→ Langage déclaratif

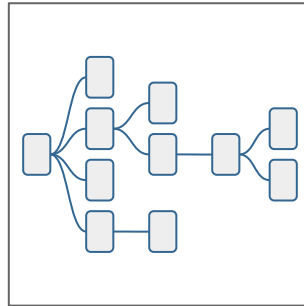
Requête SQL

```
SELECT ...  
FROM ...  
JOIN ...  
WHERE ...
```

Arbre



Arbre étendu



Résultat

1	@email.com
2	@email.com
3	@email.com
4	@email.com
5	@email.com

Parser

Rewriter

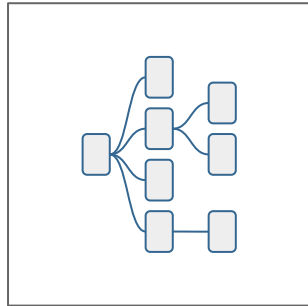
STRUCTURED QUERY LANGUAGE (SQL)

→ Langage déclaratif

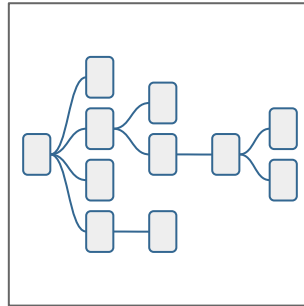
Requête SQL

```
SELECT ...  
FROM ...  
JOIN ...  
WHERE ...
```

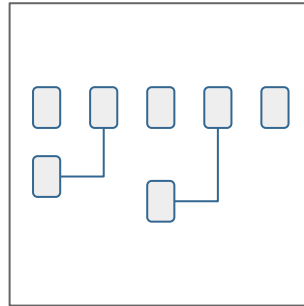
Arbre



Arbre étendu



Plan
d'exécution



Résultat

1	@email.com
2	@email.com
3	@email.com
4	@email.com
5	@email.com

Parser

Rewriter

Planner

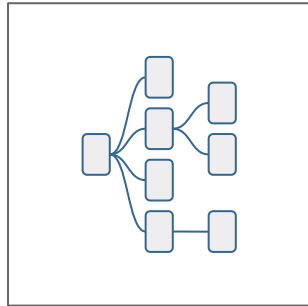
STRUCTURED QUERY LANGUAGE (SQL)

→ Langage déclaratif

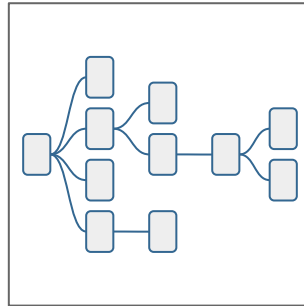
Requête SQL

```
SELECT ...  
FROM ...  
JOIN ...  
WHERE ...
```

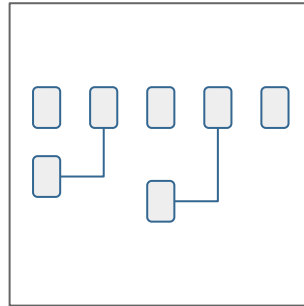
Arbre



Arbre étendu



Plan
d'exécution



Résultat

1	@email.com
2	@email.com
3	@email.com
4	@email.com
5	@email.com

Parser

Rewriter

Planner

Executor

PLANIFICATEUR & PLANS D'EXÉCUTION

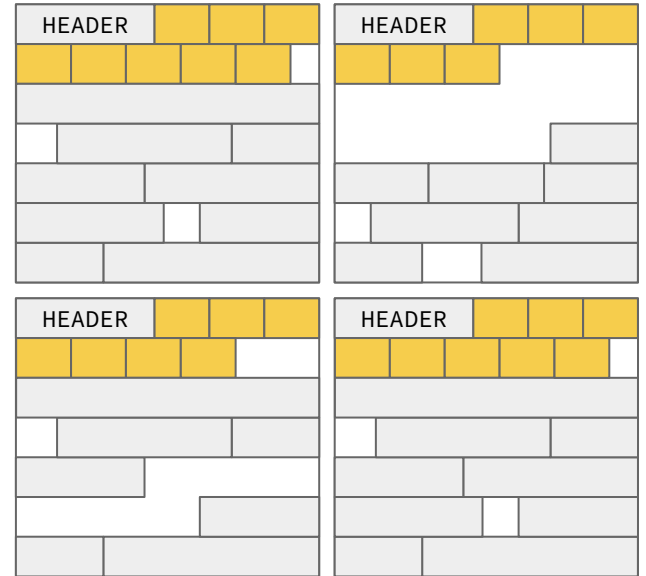
- Identifie le chemin le moins coûteux
- Statistiques de données
- Modélisation (les contraintes, les index, etc.)
- Paramètres d'instance



LECTURE DE LIGNES

```
SELECT * FROM scores  
WHERE student_id = 200
```

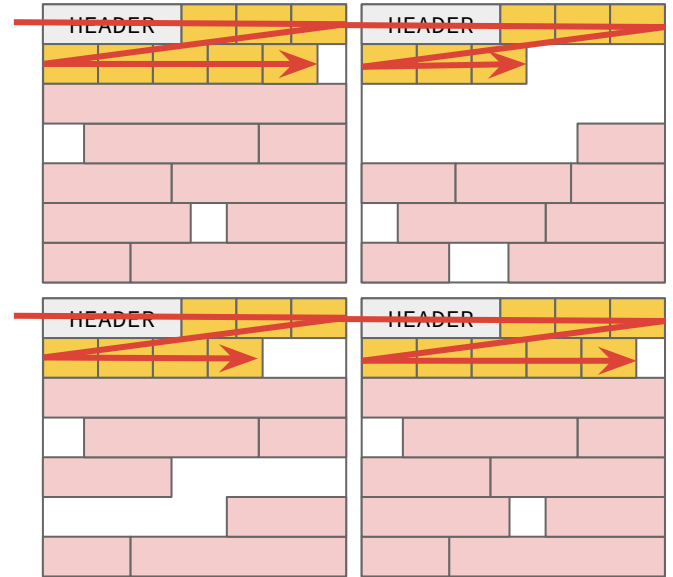
→ Quel plan d'exécution ?



LECTURE DE LIGNES

```
SELECT * FROM scores  
WHERE student_id = 200
```

- Sequential Scan
- Mise en cache des blocs
- Filtre en mémoire



EXPLAIN

```
scores=# EXPLAIN (analyze, buffers, costs off)
        SELECT * FROM scores WHERE student_id = 200;
```

QUERY PLAN

Seq Scan on scores (actual time=0.136..345.266 rows=27 loops=1)

Filter: (student_id = 200)

Rows Removed by Filter: 5810606

Buffers: shared hit=386 read=36625

Planning Time: 0.120 ms

Execution Time: 345.301 ms

EXPLAIN

```
scores=# EXPLAIN (analyze, buffers, costs off)
        SELECT * FROM scores WHERE student_id = 200;
```

QUERY PLAN

Seq Scan on scores (actual time=0.136..345.266 rows=27 loops=1)

Filter: (student_id = 200)

Rows Removed by Filter: 5810606

Buffers: shared hit=386 read=36625

Planning Time: 0.120 ms

Execution Time: 345.301 ms

EXPLAIN

```
scores=# EXPLAIN (analyze, buffers, costs off)
        SELECT * FROM scores WHERE student_id = 200;
```

QUERY PLAN

Seq Scan on scores (actual time=0.136..345.266 rows=27 loops=1)

Filter: (student_id = 200)

Rows Removed by Filter: 5810606

Buffers: **shared hit=386 read=36625 = 289 Mo**

Planning Time: 0.120 ms

Execution Time: 345.301 ms

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
        SELECT * FROM students st
           JOIN scores sc ON st.id = sc.student_id
           WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Hash Join (actual time=0.176..809.010 rows=27 loops=1)
  Hash Cond: (sc.student_id = st.id)
  Buffers: shared hit=358 read=36657
  -> Seq Scan on scores sc
        (actual time=0.055..334.512 rows=5810633 loops=1)
        Buffers: shared hit=354 read=36657
  -> Hash (actual time=0.033..0.034 rows=1 loops=1)
        Buckets: 1024 Batches: 1 Memory Usage: 9kB
        Buffers: shared hit=4
        -> Index Scan using students_email_key on students st
              (actual time=0.027..0.028 rows=1 loops=1)
              Index Cond: ((email)::text = 'johndoe@email.com'::text)
              Buffers: shared hit=4
```

Planning:

Buffers: shared hit=8

Planning Time: 0.310 ms

Execution Time: 809.064 ms

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
SELECT * FROM students st
JOIN scores sc ON st.id = sc.student_id
WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Hash Join (actual time=0.176..809.010 rows=27 loops=1)
  Hash Cond: (sc.student_id = st.id)
  Buffers: shared hit=358 read=36657
  -> Seq Scan on scores sc
      (actual time=0.055..334.512 rows=5810633 loops=1)
      Buffers: shared hit=354 read=36657
  -> Hash (actual time=0.033..0.034 rows=1 loops=1)
      Buckets: 1024 Batches: 1 Memory Usage: 9kB
      Buffers: shared hit=4
      -> Index Scan using students_email_key on students st
          (actual time=0.027..0.028 rows=1 loops=1)
          Index Cond: ((email)::text = 'johndoe@email.com'::text)
          Buffers: shared hit=4
```

Planning:

Buffers: shared hit=8

Planning Time: 0.310 ms

Execution Time: 809.064 ms

^ Hash Join #1
on sc.student_id = st.id

Hash Join Node joins two record sets by hashing one of them (using a Hash Scan).

General IO & Buffers Output Workers Misc

🕒 Timing: 434ms | 58%
☰ Rows: 27 (Planned: 15)
💰 Cost: 15,300 (Total: 110,000)

^ Seq Scan #2
on public.scores as sc

Seq Scan Node finds relevant records by sequentially scanning the input record set. When reading from a table, Seq Scans (unlike Index Scans) perform a single read operation (only the table is read).

General IO & Buffers Output Workers Misc

🕒 Timing: 316ms | 42%
☰ Rows: 5,810,633 (Planned: 5,810,611)
💰 Cost: 95,100 (Total: 95,100)

^ Hash #3

Hash Node generates a hash table from the records in the input recordset. Hash is used by Hash Join.

General IO & Buffers Output Workers Misc

🕒 Timing: 0.01ms | 0%
☰ Rows: 1 (Planned: 1)

^ Index Scan #4
on public.students as st
using students_email_key

Index Scan Node finds relevant records based on an **Index**. Index Scans perform 2 read operations: one to read the index and another to read the actual value from the table.

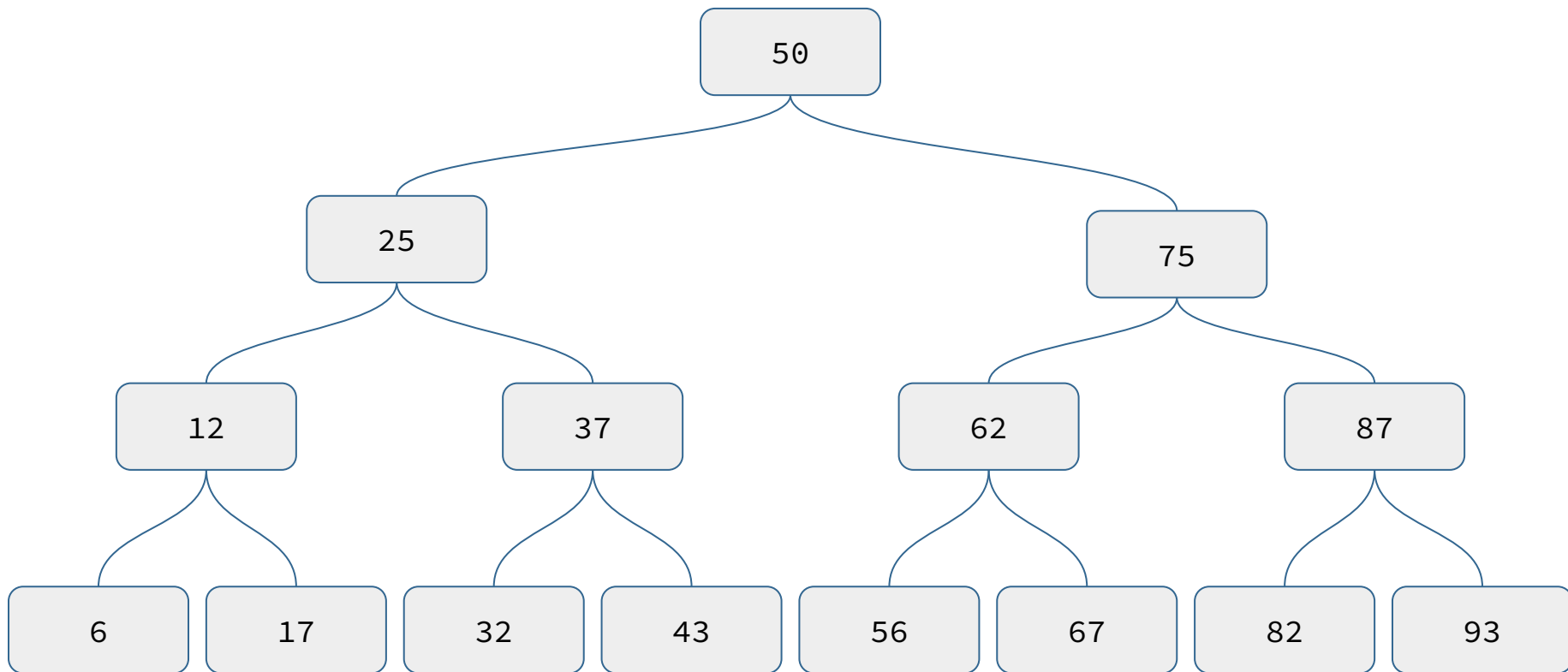
General IO & Buffers Output Workers Misc

🕒 Timing: 0.07ms | 0%
☰ Rows: 1 (Planned: 1)
💰 Cost: 8.44 (Total: 8.44)

<https://explain.dalibo.com/>

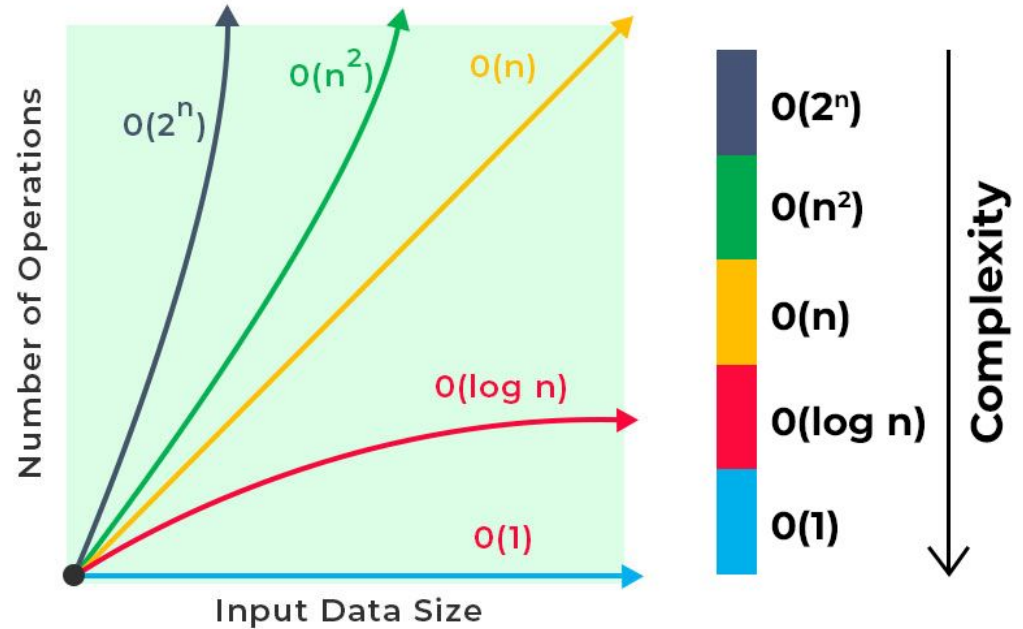
B(RO)-TREE

ARBRE BINAIRE BASIQUE



BTREE

- $\log_2(100) = 6.64$
- $\log_2(1000) = 9.96$
- **$\log_2(10^9) = 29.9$**
- $\log_2(10^{12}) = 39.86$



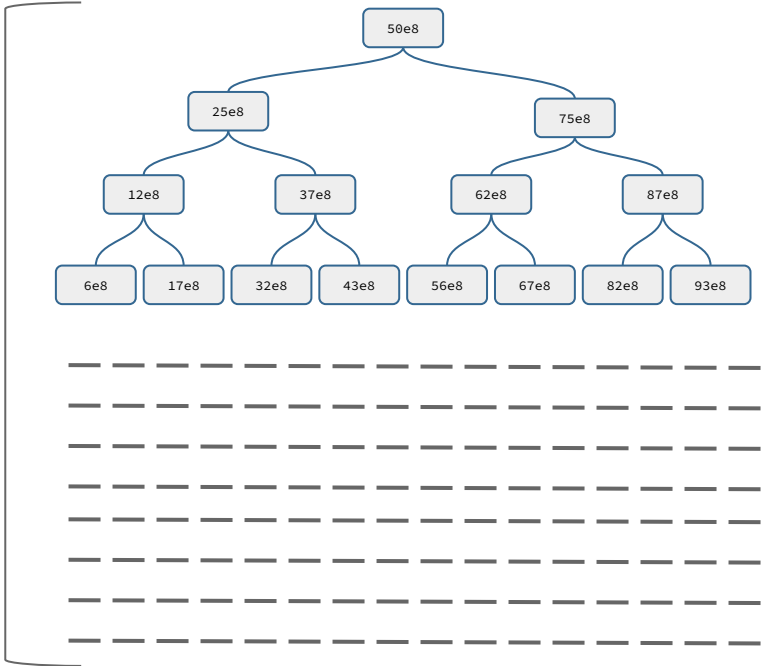
BINARY TREE != B+TREE

- 29.9 “nodes” internes
- 29.9 opérations

$$\log(10^9) =$$

29.9 niveaux

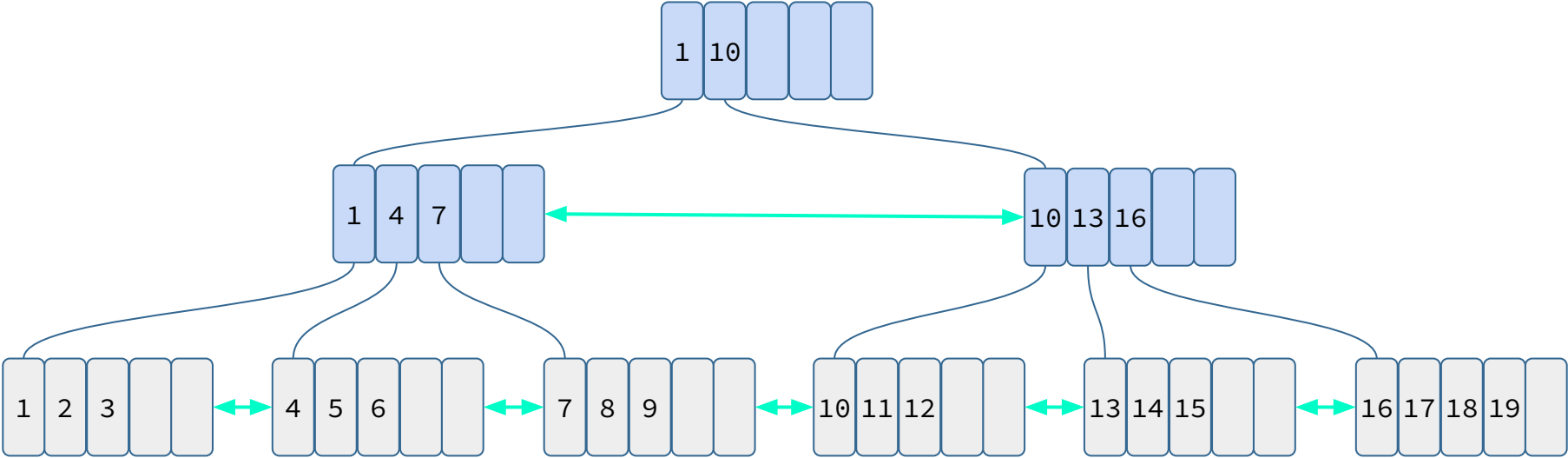
29.9 niveaux si BTREE de facteur 1



B+TREE

node "interne"

node "feuille"

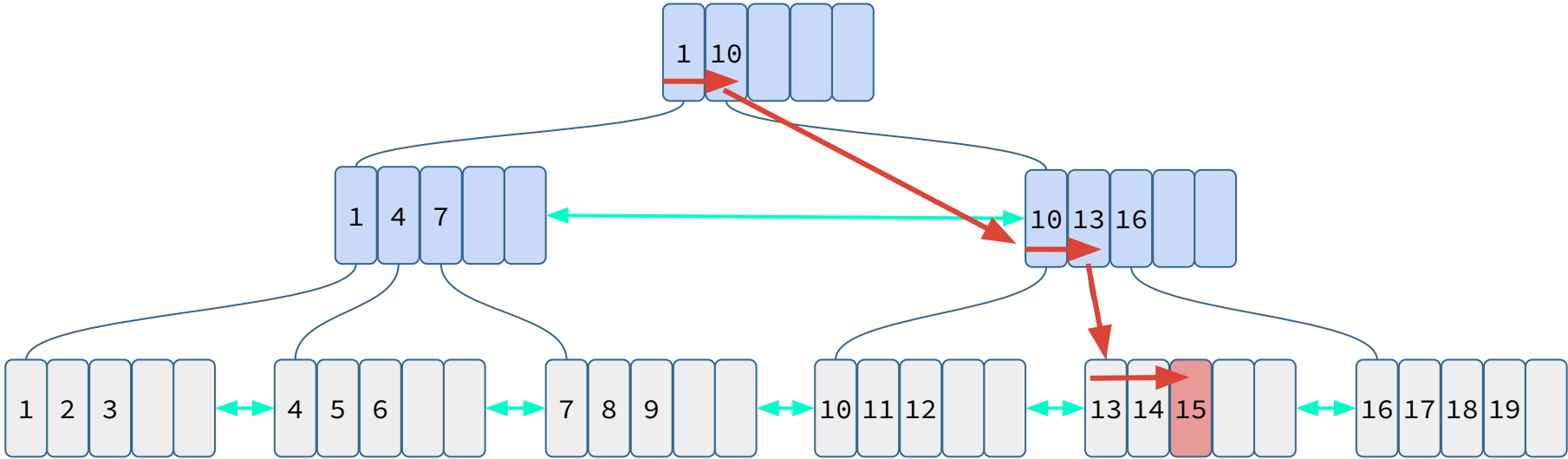


RECHERCHE

```
SELECT * FROM scores  
WHERE student_id = 15
```

node "interne"

node "feuille"

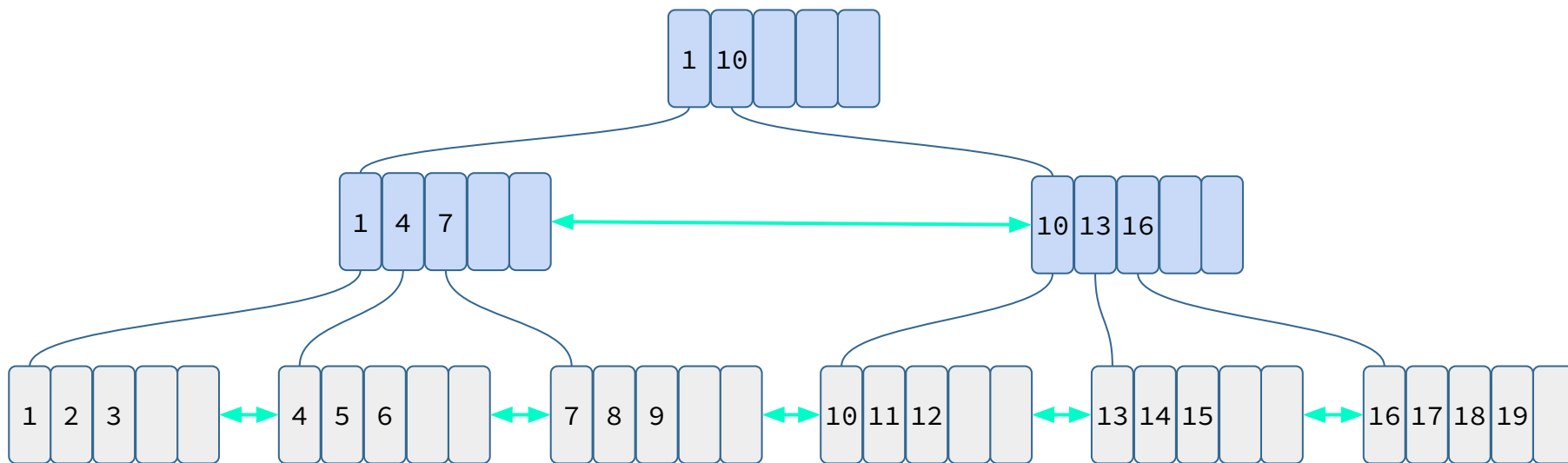


→ Nombre de valeurs par nodes pour des int8

8ko / 8 ~= 1024 valeurs par nodes

node "interne"

node "feuille"



$$\lceil \log_{1024}(10^9) \rceil =$$

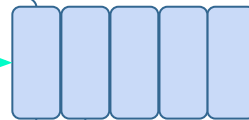
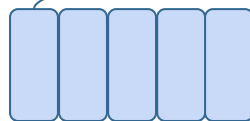
2.98

Pour un index B+TREE de entier 8 octet, parcourir un milliard de valeurs ne prendra que 3 blocs d'index 8kb.

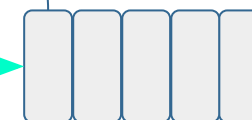
Racine



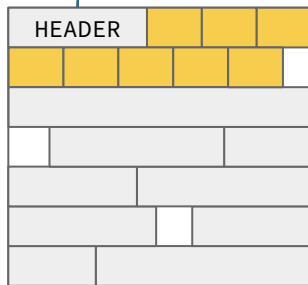
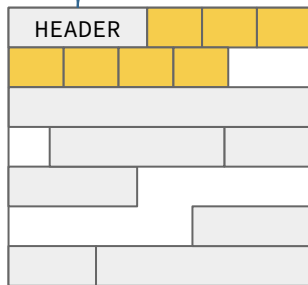
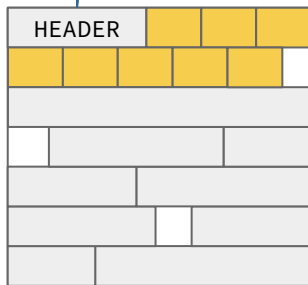
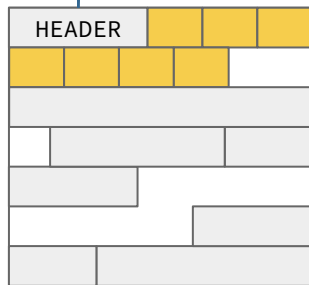
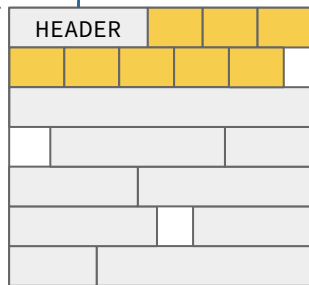
Noeuds internes



Feuilles



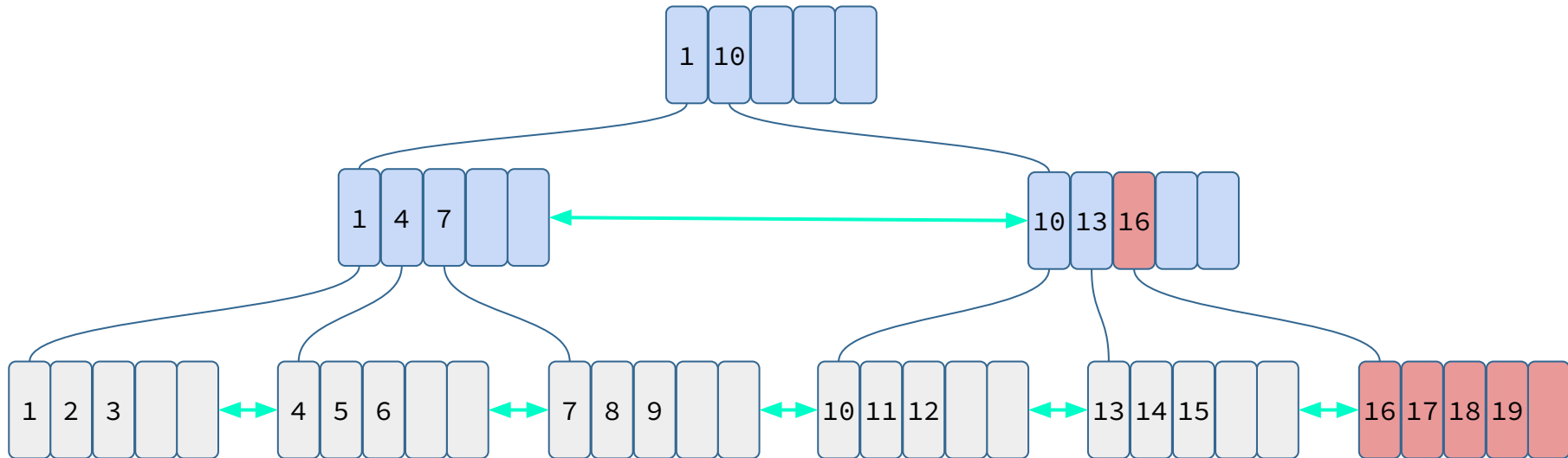
Mapping



HEAP

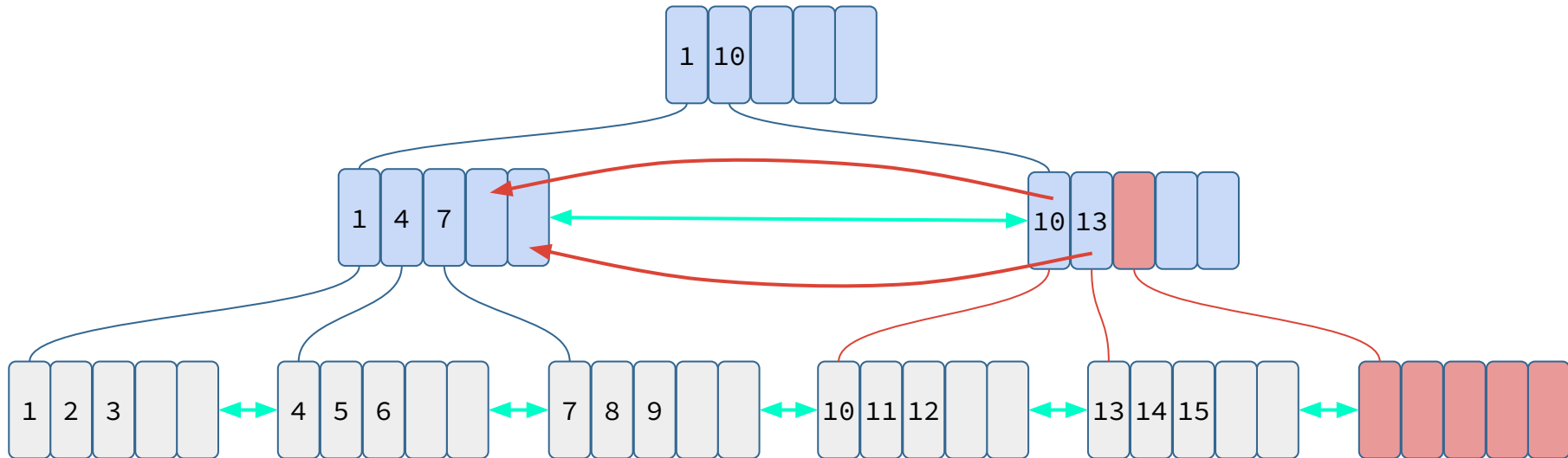
SUPPRESSION DE PLUSIEURS LIGNES

```
DELETE FROM scores  
WHERE student_id > 15
```



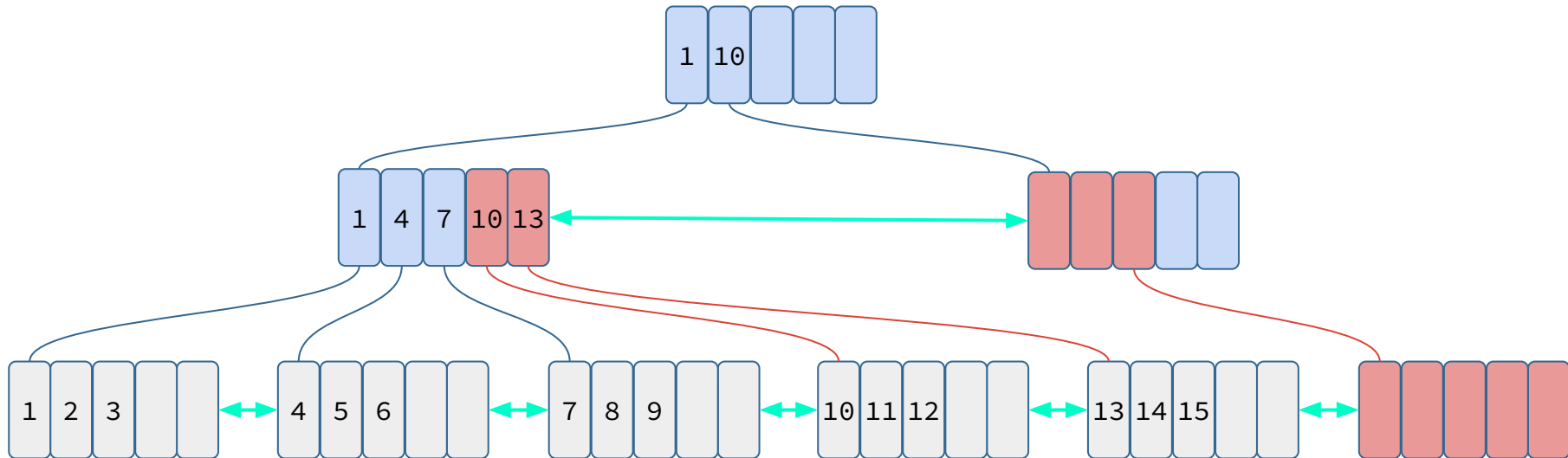
SUPPRESSION DE PLUSIEURS LIGNES

```
DELETE FROM scores  
WHERE student_id > 15
```



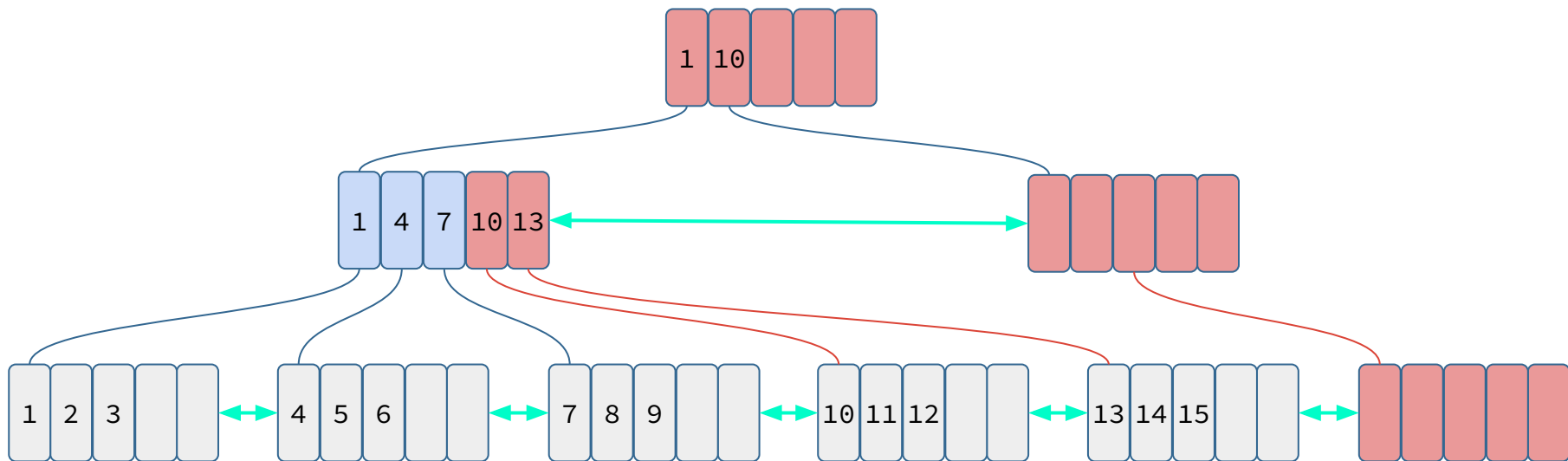
SUPPRESSION DE PLUSIEURS LIGNES

```
DELETE FROM scores  
WHERE student_id > 15
```



SUPPRESSION DE PLUSIEURS LIGNES

```
DELETE FROM scores  
WHERE student_id > 15
```



SUPPRESSION DE PLUSIEURS LIGNES

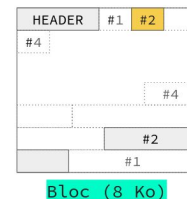
```
DELETE FROM scores  
WHERE student_id > 15
```



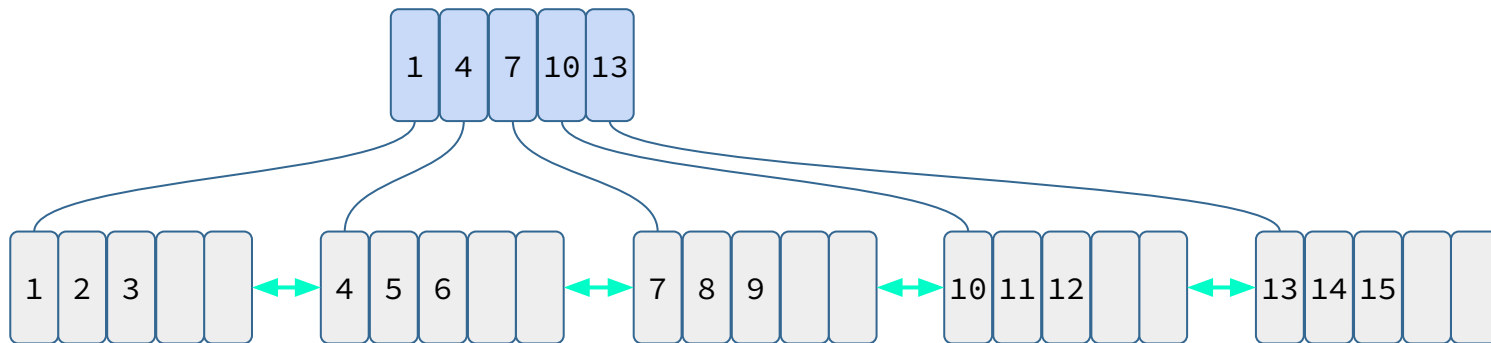
SUPPRESSION D'UNE LIGNE

```
DELETE FROM scores  
WHERE student_id = 200
```

- Fragmentation (Bloat)
- VACUUM
- VACUUM FULL



× 4 lignes



```
CREATE INDEX  
ON scores(student_id)
```



```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
        SELECT * FROM students st
           JOIN scores sc ON st.id = sc.student_id
           WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Nested Loop (actual time=0.067..0.251 rows=27 loops=1)
->  Index Scan using students_email_key on students st
    (actual time=0.019..0.020 rows=1 loops=1)
    Index Cond: ((email)::text = 'johndoe@email.com'::text)
    Buffers: shared hit=4
->  Bitmap Heap Scan on scores sc
    (actual time=0.043..0.217 rows=27 loops=1)
    Recheck Cond: (st.id = student_id)
    Heap Blocks: exact=27
    Buffers: shared hit=6 read=24
->  Bitmap Index Scan on scores_student_id_idx
    (actual time=0.029..0.029 rows=27 loops=1)
    Index Cond: (student_id = st.id)
    Buffers: shared hit=2 read=1
```

Planning:

Buffers: shared hit=43 read=6

Planning Time: 0.556 ms

Execution Time: 0.285 ms

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
        SELECT * FROM students st
           JOIN scores sc ON st.id = sc.student_id
           WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Nested Loop (actual time=0.067..0.251 rows=27 loops=1)
->  Index Scan using students_email_key on students st
    (actual time=0.019..0.020 rows=1 loops=1)
    Index Cond: ((email)::text = 'johndoe@email.com'::text)
    Buffers: shared hit=4
->  Bitmap Heap Scan on scores sc
    (actual time=0.043..0.217 rows=27 loops=1)
    Recheck Cond: (st.id = student_id)
    Heap Blocks: exact=27
    Buffers: shared hit=6 read=24
->   Bitmap Index Scan on scores_student_id_idx
      (actual time=0.029..0.029 rows=27 loops=1)
      Index Cond: (student_id = st.id)
      Buffers: shared hit=2 read=1
```

Planning:

Buffers: shared hit=43 read=6

Planning Time: 0.556 ms

Execution Time: 0.285 ms

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
        SELECT * FROM students st
           JOIN scores sc ON st.id = sc.student_id
           WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Nested Loop (actual time=0.067..0.251 rows=27 loops=1)
->  Index Scan using students_email_key on students st
    (actual time=0.019..0.020 rows=1 loops=1)
    Index Cond: ((email)::text = 'johndoe@email.com'::text)
    Buffers: shared hit=4
->  Bitmap Heap Scan on scores sc
    (actual time=0.043..0.217 rows=27 loops=1)
    Recheck Cond: (st.id = student_id)
    Heap Blocks: exact=27
    Buffers: shared hit=6 read=24
->   Bitmap Index Scan on scores_student_id_idx
    (actual time=0.029..0.029 rows=27 loops=1)
    Index Cond: (student_id = st.id)
    Buffers: shared hit=2 read=1
```

Planning:

Buffers: **shared hit=43 read=6 = 2 Mo**

Planning Time: 0.556 ms

Execution Time: 0.285 ms

CASCADE FOREIGN KEYS

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
        DELETE FROM students WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Delete on students (actual time=0.082..0.083 rows=0 loops=1)
```

```
  Buffers: shared hit=9 dirtied=2
```

```
   -> Index Scan using students_email_key on students
```

```
        (actual time=0.029..0.032 rows=1 loops=1)
```

```
        Index Cond: ((email)::text = 'johndoe@email.com'::text)
```

```
        Buffers: shared hit=4
```

```
Planning Time: 0.089 ms
```

```
Trigger for constraint scores_student_id_fkey: time=562.464 calls=1
```

```
Execution Time: 562.635 ms
```

```
(8 rows)
```


CASCADE FOREIGN KEYS

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
        DELETE FROM students WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Delete on students (actual time=0.082..0.083 rows=0 loops=1)
```

```
  Buffers: shared hit=9 dirtied=2
```

```
   -> Index Scan using students_email_key on students
```

```
        (actual time=0.029..0.032 rows=1 loops=1)
```

```
        Index Cond: ((email)::text = 'johndoe@email.com'::text)
```

```
        Buffers: shared hit=4
```

```
Planning Time: 0.089 ms
```

```
Trigger for constraint scores_student_id_fkey: time=562.464 calls=1
```

```
Execution Time: 562.635 ms
```

```
(8 rows)
```

CASCADE FOREIGN KEYS - AVEC INDEX

```
scores=# EXPLAIN (ANALYZE, BUFFERS, COSTS OFF)
         DELETE FROM students WHERE email = 'johndoe@email.com';
```

QUERY PLAN

```
Delete on students (actual time=0.076..0.077 rows=0 loops=1)
```

```
  Buffers: shared hit=6
```

```
   -> Index Scan using students_email_key on students
```

```
       (actual time=0.051..0.053 rows=1 loops=1)
```

```
      Index Cond: ((email)::text = 'johndoe@email.com'::text)
```

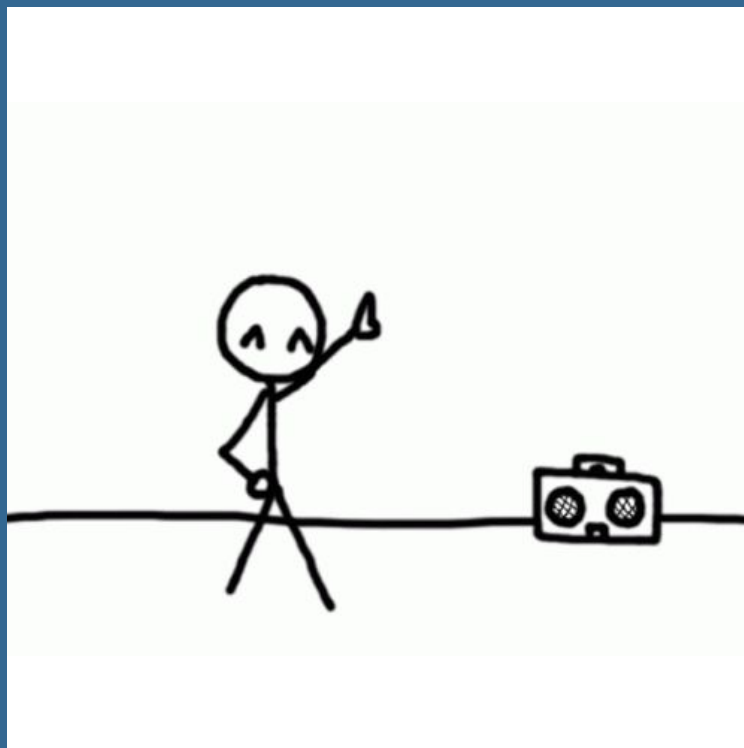
```
      Buffers: shared hit=4
```

```
Planning Time: 0.120 ms
```

```
Trigger for constraint scores_student_id_fkey: time=0.650 calls=1
```

```
Execution Time: 0.757 ms
```

```
(8 rows)
```



CONCLUSION

- N'indexez que ce qui est nécessaire
- Indexez vos colonnes de clés étrangères
- Surveillez vos plans d'exécution avec EXPLAIN

ENVIE DE PARTICIPER, DE PARTAGER VOTRE EXPÉRIENCE
OU DE SPONSORISER LES PROCHAINS ÉVÉNEMENTS ?

Contactez-nous sur les
pages Meetup et
Linkedin du groupe !

