

Indexer des fragments d'ADN

En collaboration avec:
N. Philippe, T. Commes, É. Rivals
M. Léonard, T. Lecroq

Journées au vert

21 juin 2011



Introduction

Contexte

- ▶ Séquenceurs haut-débit
- ▶ Des giga-octets de données par jour

Questions

- ▶ Comment tirer parti de tous les fragments produits ?
- ▶ Comment les organiser de façon intelligente ?
- ▶ Comment ne pas saturer nos machines (portables, ordinateurs de bureau, serveurs, ...)?

Des points communs...



Des points communs...



Des points communs...



Des points communs...



Un annuaire de k -mers

Fragments :

	0	1	2	3	4	5		0	1	2	3	4	5		0	1	2	3	4	5	
	A	A	G	A	A	G		G	C	A	A	G	A		A	G	A	A	G	T	
							0							1							2

Retrouver l'adresse des homonymes

Fragments : 0 1 2 3 4 5 0 1 2 3 4 5 0 1 2 3 4 5
 AAGAAG GCAAGA AGAAGT
 0 1 2

Annuaire (ou index...)

		F	P	
0	AAG	0	0	
1	AAG	0	3	
2	AAG	1	2	← 4
3	AAG	2	2	
4	AGT	2	3	1 1
5	AGA	0	1	
6	AGA	1	3	← 3
7	AGA	2	0	
8	CAA	1	1	3 1
9	GAA	0	2	
10	GAA	2	1	4 2
11	GCA	1	0	5 1

Annuaire inverse

F	P	
0	0	0
0	1	2
0	2	4
0	3	0
1	0	5
1	1	3
1	2	0
1	3	2
2	0	2
2	1	4
2	2	0
2	3	1

Des requêtes intéressantes...

Fragments : 0 1 2 3 4 5 0 1 2 3 4 5 0 1 2 3 4 5
 AAGAAG GCAAGA AGAAGT
 0 1 2

Annuaire (ou index...)

	F	P		F	P
0	0	0		0	0
1	0	3		0	1
2	1	2	0 4	0 2	4
3	2	2		0	3
4	2	3	1 1	1	0
5	0	1		1	1
6	1	3	2 3	1	2
7	2	0		1	3
8	1	1	3 1	2	0
9	0	2		2	1
10	2	1	4 2	2	2
11	1	0	5 1	2	3

Requêtes

1. Combien de fois voit-on le même k -mer qu'en 0 2 ?

Des requêtes intéressantes...

Fragments : 0 1 2 3 4 5 0 1 2 3 4 5 0 1 2 3 4 5
 A A G A A G G C A A G A A G A A G T
 0 1 2

Annuaire (ou index...)

	F	P		F	P
0	0	0		0	0
1	0	3		0	1
2	1	2	0 4	0	2
3	2	2		0	3
4	2	3	<u>1 1</u>	1	0
5	0	1		1	1
6	1	3	2 3	1	2
7	2	0		1	3
8	1	1	<u>3 1</u>	2	0
9	0	2		2	1
10	2	1	4 2	2	2
11	1	0	<u>5 1</u>	2	3

Requêtes

1. Combien de fois voit-on le même k -mer qu'en 0 2? $O(1)$

Des requêtes intéressantes...

Fragments : 0 1 2 3 4 5 0 1 2 3 4 5 0 1 2 3 4 5
 AAGAAG GCAAGA AGAAGT
 0 1 2

Annuaire (ou index...)

	F	P		F	P
0	0	0		0	0
1	0	3		0	1
2	1	2	0 4	0	2
3	2	2		0	3
4	2	3	<u>1 1</u>	1	0
5	0	1		1	1
6	1	3	2 3	1	2
7	2	0		2	0
8	1	1	<u>3 1</u>	2	1
9	0	2	← 4 2	2	2
10	2	1		2	1
11	1	0	<u>5 1</u>	2	3

Requêtes

1. Combien de fois voit-on le même k -mer qu'en 0 2? $O(1)$
2. À quels endroits trouve-t-on ce même k -mer?

Des requêtes intéressantes...

Fragments : 0 1 2 3 4 5 0 1 2 3 4 5 0 1 2 3 4 5
 A A G A A G G C A A G A A G A A G T
 0 1 2

Annuaire (ou index...)

	F	P		F	P
0	0	0		0	0
1	0	3		0	1
2	1	2	0 4	0	2
3	2	2		0	3
4	2	3	1 1	1	0
5	0	1		1	1
6	1	3	2 3	1	2
7	2	0		2	0
8	1	1	3 1	2	1
9	0	2	4 2	2	2
10	2	1	5 1	2	1
11	1	0		2	3

Diagram illustrating the mapping of fragments to positions. The table shows the frequency of each fragment (F) at each position (P). A yellow highlight is placed on the value '0 2' in the row for position 2, column for fragment '0 2'. Arrows indicate the mapping from the highlighted cell to the corresponding fragment and position in the original data.

Requêtes

1. Combien de fois voit-on le même k -mer qu'en 0 2? $O(1)$
2. À quels endroits trouve-t-on ce même k -mer? $O(\text{occ})$

Variantes

1. Combien de fragments partagent un même k -mer? $O(\text{occ})$
2. Quels fragments contiennent ce k -mer? $O(\text{occ})$

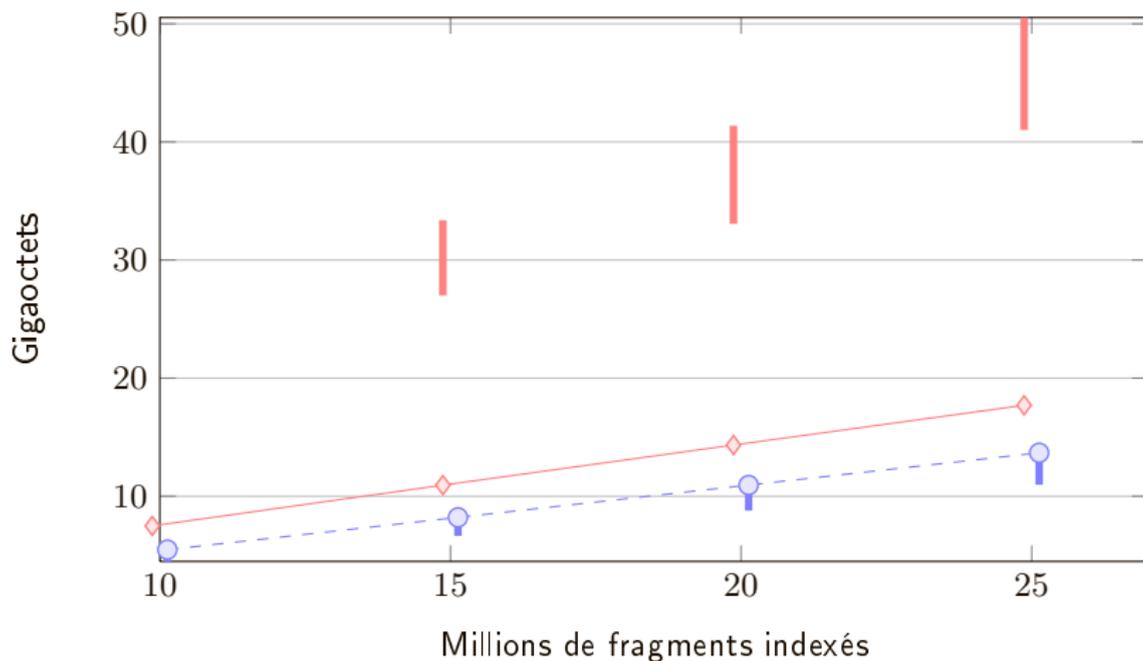
Quelques expériences

Ce qui nous intéresse

- ▶ Combien ça pèse ?
- ▶ Combien de temps ça prend de faire une requête ?
- ▶ Est-ce mieux qu'une solution alternative (table de hachage) ?

Quel espace?

Indexation de fragments de longueur 75, pour k allant de 10 à 30



$k = 10$

$k = 15 \dots 30$

--- (dashed blue line)

| (solid blue line)

Gk arrays

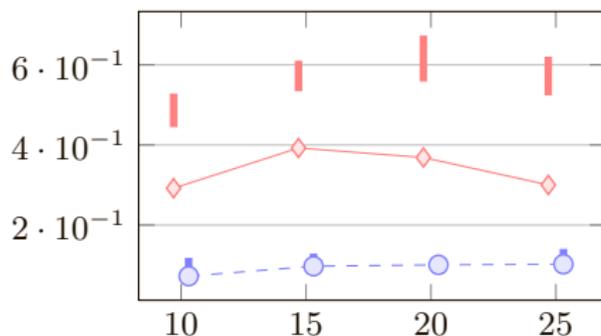
— (solid red line)

| (solid red line)

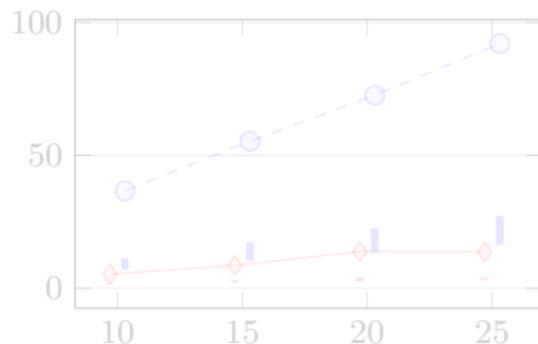
Table de hâchage

Quelle rapidité?

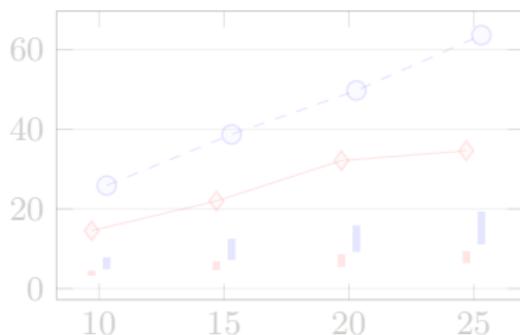
Combien d'homonymes?



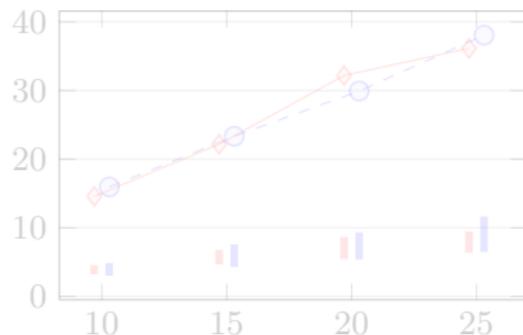
À quelles positions?



Dans combien de fragments?



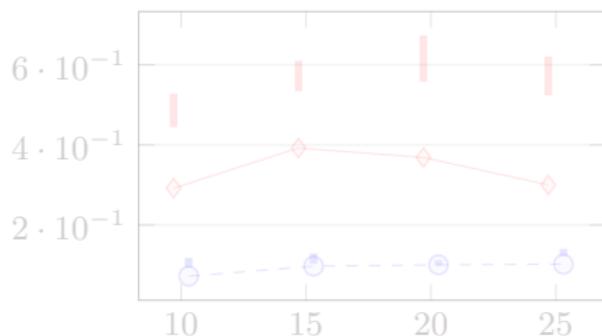
Dans quels fragments?



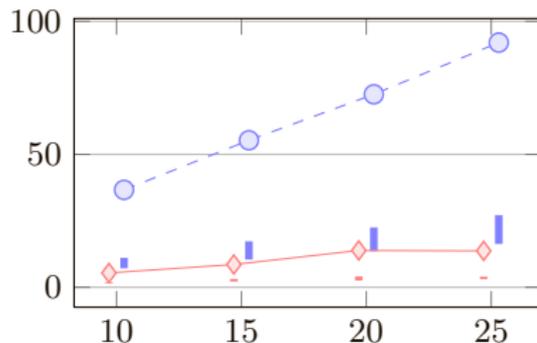
Millions de fragments indexés

Quelle rapidité?

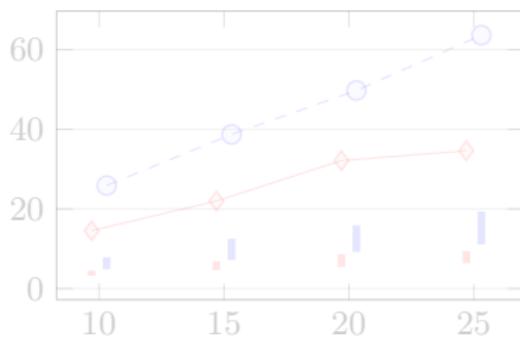
Combien d'homonymes?



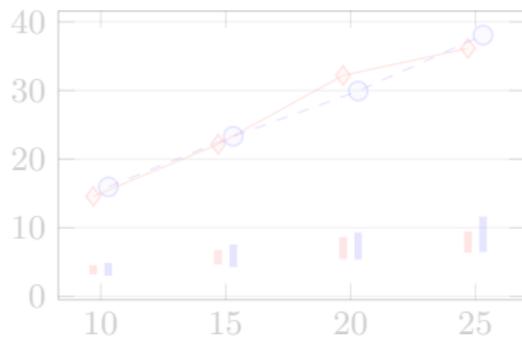
À quelles positions?



Dans combien de fragments?



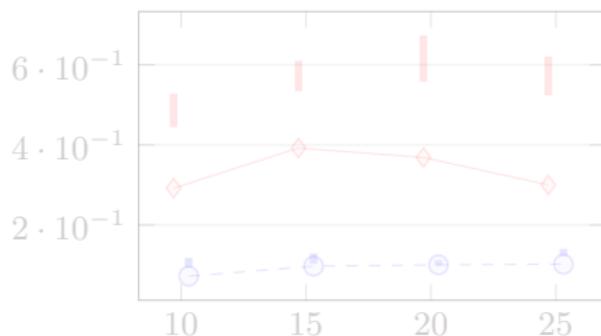
Dans quels fragments?



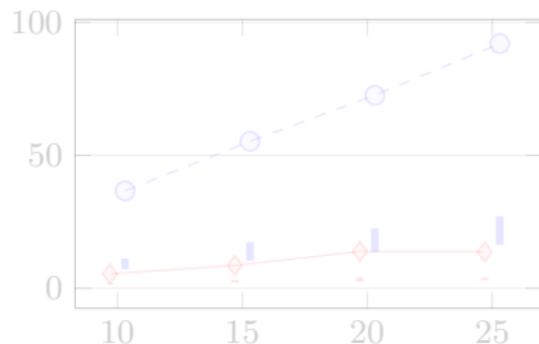
Millions de fragments indexés

Quelle rapidité?

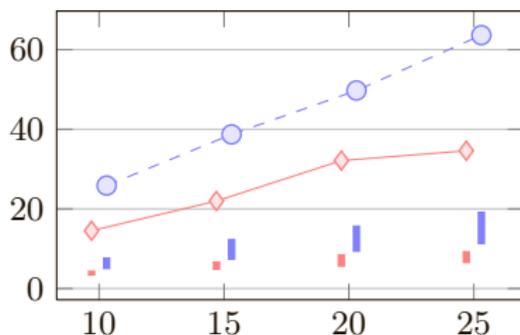
Combien d'homonymes?



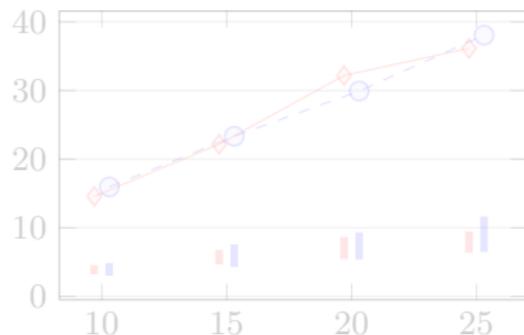
À quelles positions?



Dans combien de fragments?



Dans quels fragments?

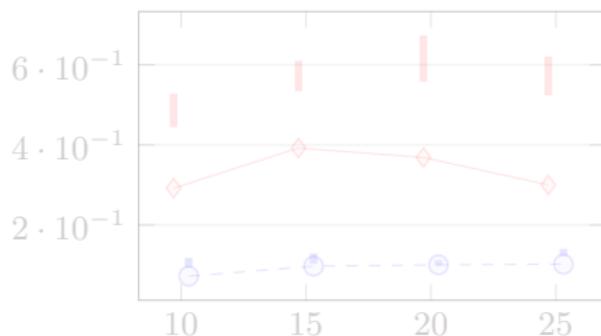


Millions de fragments indexés

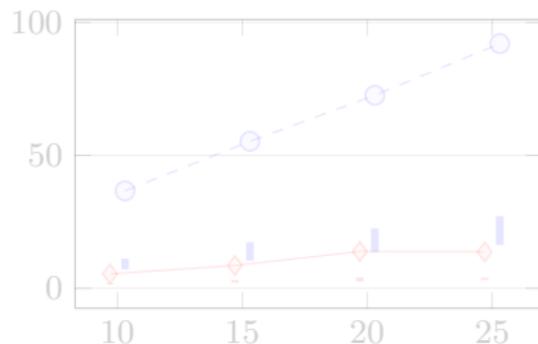
Durée en μs

Quelle rapidité?

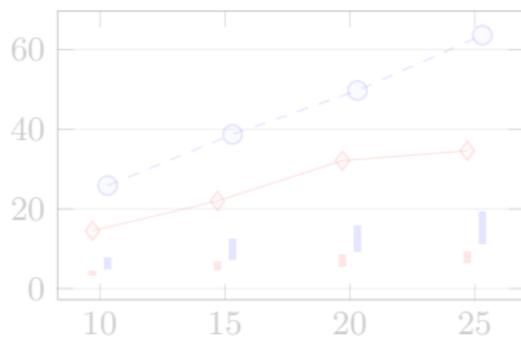
Combien d'homonymes?



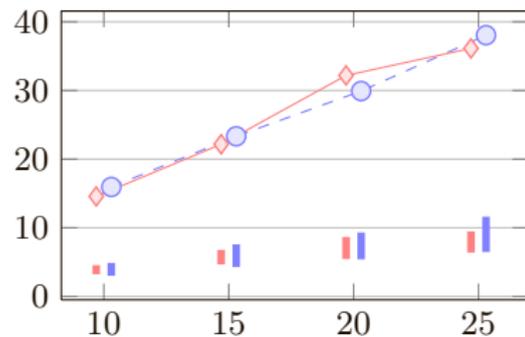
À quelles positions?



Dans combien de fragments?



Dans quels fragments?

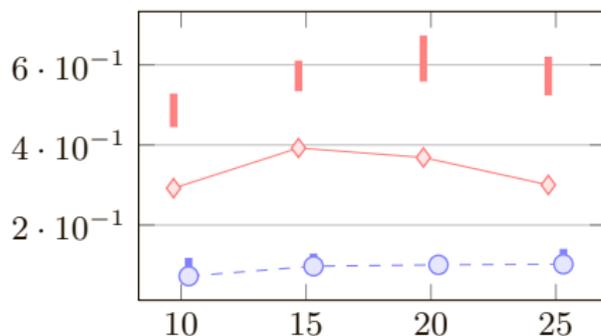


Millions de fragments indexés

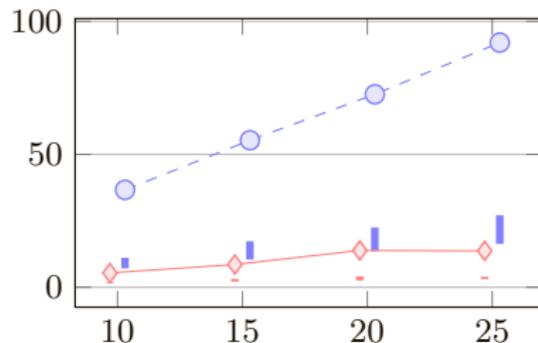
Durée en μs

Quelle rapidité?

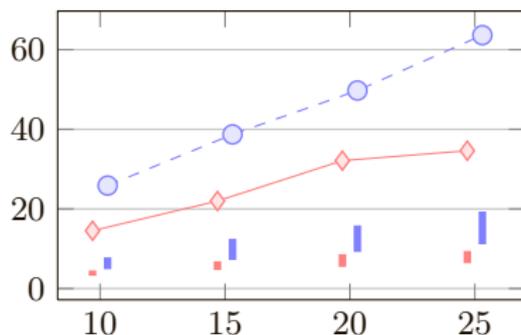
Combien d'homonymes?



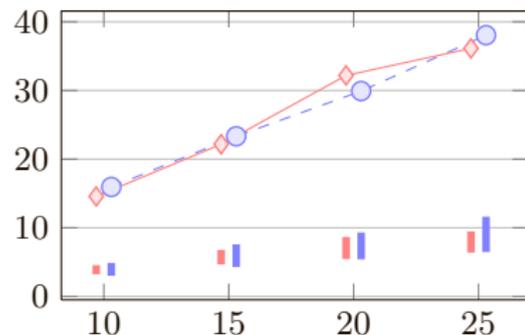
À quelles positions?



Dans combien de fragments?

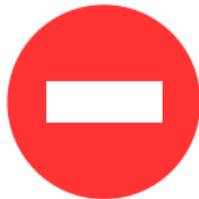


Dans quels fragments?



Millions de fragments indexés

Conclusions



- ▶ Mauvaises performances avec un grand nombre d'occurrences (k petit)
- ▶ Espace peu dépendant de k



- ▶ Prend moins de place que les tables de hâchage
- ▶ Espace peu dépendant de k
- ▶ Requêtes à peu près aussi rapides qu'une table de hâchage
- ▶ Peut gérer des reads de longueur variable