

Méthodes de mapping de reads avec indexation des reads

Pierre Morisse, Thierry Lecroq, Arnaud Lefebvre

LITIS
Université de Rouen
76000 Rouen, France

Introduction

Depuis le milieu des années 2000 et le développement des séquenceurs à très haut débit (*Next Generation Sequencing*), la biologie doit faire face au traitement d'énormes quantités de données, formées par des millions de très courtes séquences appelées *reads*. Dans un papier relativement récent, Philippe et al. ont souligné l'importance de l'indexation de ces *reads* afin de résoudre des problèmes de correction ou de *mapping*, et ont développé un index supportant les 7 requêtes suivantes, pour une séquence f de longueur k donnée :

- Dans quels *reads* f apparaît ?
- Dans combien de *reads* f apparaît ?
- Quelles sont les occurrences de f ?
- Quel est le nombre d'occurrences de f ?
- Dans quels *reads* f n'apparaît qu'une fois ?
- Dans combien de *reads* f n'apparaît qu'une fois ?
- Quelles sont les occurrences de f dans les *reads* où f n'apparaît qu'une fois ?

Nous présentons ici l'état de l'art concernant les technologies de séquençage et les problèmes susmentionnés, ainsi que le problème sur lequel nous nous sommes principalement penchés durant le déroulement de ce stage.

État de l'art

Les tableaux présentés ci-dessous résument brièvement l'état de l'art concernant les technologies de séquençage et les outils existants, utilisant une structure d'index sur les *reads*, et permettant de résoudre des problèmes de correction, de *mapping*, ou de traitement des 7 requêtes précédentes.

Technologie	Technique de séquençage	Plateforme	Nombre de <i>reads</i>	Longueur	Précision	Temps	Débit	Coût	Erreurs
Illumina	Synthèse basé sur polymères	HiSeq 2500/1500	3 milliards	36 - 100	99	2 - 11 jours	600	740 000	Substitutions
		MiSeq	17 millions	25 - 250	~99	4 - 27 heures	125 000	8,5	
Roche	Polyséquençage	454 GS FLX+	1 million	700	99,997	23 heures	0,7	450 000	
		454 GS Junior	1 million	400	~99	10 heures	0,4	100 000	Indels.
ABI Life Technologies	Ligatures	2500xl SOLiD	2,8 millions	75	99,99	7 jours	180	595 000	
	Détection de ponts	Ion PGM Sequencer	60 - 80 millions	jusqu'à 200	~99	2 heures	10 - 100	245 000	Indels.
Pacific Biosciences	Simple molécule en temps réel	PacBio RS	50 000	3 000 en moyenne	85	2 heures	13	750 000	Indels.
Oxford Nanopore	Exonuclease par Nanopore	GridION	4 - 10	plusieurs milliers	96	variable	quelques dizaines	variable	
		MinION	70 000	1 065	70	48 heures	0,132	1 000	Indels.

TABLE 1: Récapitulatif des différentes technologies de séquençage. La précision est donnée en %, le débit en Gb, et le coût en \$.

Outil	Structure de données	Erreurs corrigées	Nombre de <i>reads</i> (longueur)	Espace mémoire (en Mo)	Temps (en min)	reads corrigés (en %)
SHREC	Arbre des suffixes	subs.	1 000 946 (70)	1 300	183	88,56
HybridSHREC	Arbre des suffixes	subs. + indels	977 971 (178)	15 000	28	98,39
HTREC	Table des suffixes	subs.	1 000 946 (70)	757	28	94,43
Flomo	Table des suffixes échantillonnée	subs. + indels	977 971 (178)	2 000	15	66,76
Coral	Table de hachage	subs. + indels	2 464 690 (142)	3 000	32	92,88
RACER	Table de hachage	subs.	977 971 (178)	6 000	5	92,88
BLESS	Filtres de Bloom	subs. + indels	2 119 404 (75)	1 437	23	76,65
			1 096 140 (101)	11	6	84,38
LoRDEC	Graphes de De Bruijn	subs. + indels	33 360 reads longs (2 938) et 2 513 613 <i>reads</i> courts (100)	960	10	85,78

TABLE 2: Récapitulatif des différentes méthodes de correction des *reads*. L'espace mémoire est donnée en Mo. Le temps est donné en minutes. Les *reads* corrigés sont donnés en %, et la valeur indiquée est une moyenne.

Outil	Structure de données	Erreurs prises en compte	Nombre de <i>reads</i> (longueur)	Espace mémoire (en Mo)	Temps (en min)	reads mappés (en %)
MAQ	Table de hachage	subs. + indels	1 000 000 (44)	1 200	331	92,53
MinFAST	Table de hachage	subs.	1 000 000 (100)	20 000	169	90,70
MinFAST-Ultra	Table de hachage	subs.	2 000 000 (100)	2 000	57	91,41

TABLE 3: Récapitulatif des différentes méthodes de *mapping* de *reads*. L'espace mémoire est donnée en Mo. Le temps est donné en minutes. Les *reads* mappés sont donnés en %, et la valeur indiquée est une moyenne.

Peu d'outils sont présentés ici, mais de nombreuses méthodes de *mapping* de *reads*, n'utilisant pas de structure d'index sur les *reads*, existent et produisent de très bons résultats, aussi bien en espace et en temps, qu'en qualité de *mapping*.

Outil	Structure de données	Nombre de <i>reads</i> (longueur)	Espace mémoire (en Go)	Temps R1 (en ms)	Temps R2 (en ms)	Temps R3 (en ms)	Temps R4 (en ms)
GKA	Table des suffixes modifiés interne	42 400 000 (75)	20	16	25	25	0,1
	Table associant à une occurrence d'occurrences						
CGKA	Table des suffixes échantillonnée	42 400 000 (75)	3 - 7	1203	28	1278	28
	3 vecteurs de bits						
PgSA	Table des suffixes échantillonnée	42 400 000 (75)	1 - 4	70	58	70	58
	Table auxiliaire d'information sur les <i>reads</i> et <i>k</i> -mers						

TABLE 4: Récapitulatif des différentes méthodes permettant de traiter les 7 requêtes. L'espace mémoire est donné en Go. Le temps est donné en millisecondes. Les requêtes 5-7 sont exclues du comparatif, car non implémentées dans GkA et CGkA lors des tests réalisés dans le papier introduisant PgSA.

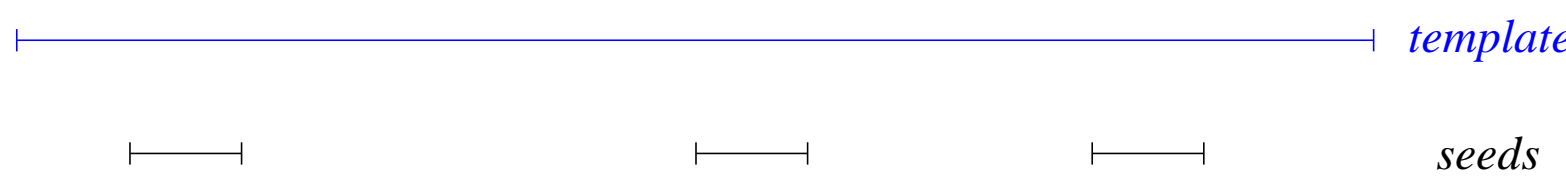
Correction de *reads* longs : Les *reads* NaS (Nanopore Synthetic-long)

Les nouvelles technologies de séquençage permettent de séquencer des *reads* de plus en plus longs, mais ceux-ci disposent d'un important taux d'erreur, avoisinant notamment les 30% pour les *reads* séquencés par la plateforme MinION. Comme le montrent les tableaux précédents, les méthodes de correction classiques ne sont pas adaptées à de tels *reads*, et sont donc très peu efficaces. Une solution alternative pour résoudre ce problème est la génération de *reads* dits synthétiques. Ces derniers sont générés via une approche hybride, utilisant des *reads* longs comme *templates* et des *reads* courts disposant d'un plus faible taux d'erreur. Les *reads* ainsi synthétisés, appelés *reads* NaS, car synthétisés à partir de *reads* de la technologie Nanopore, peuvent atteindre une longueur de 60 000, et s'aligner intégralement et sans erreurs. Nous présentons ici une première méthode permettant de synthétiser de tels *reads*, et la méthode que nous avons développée.

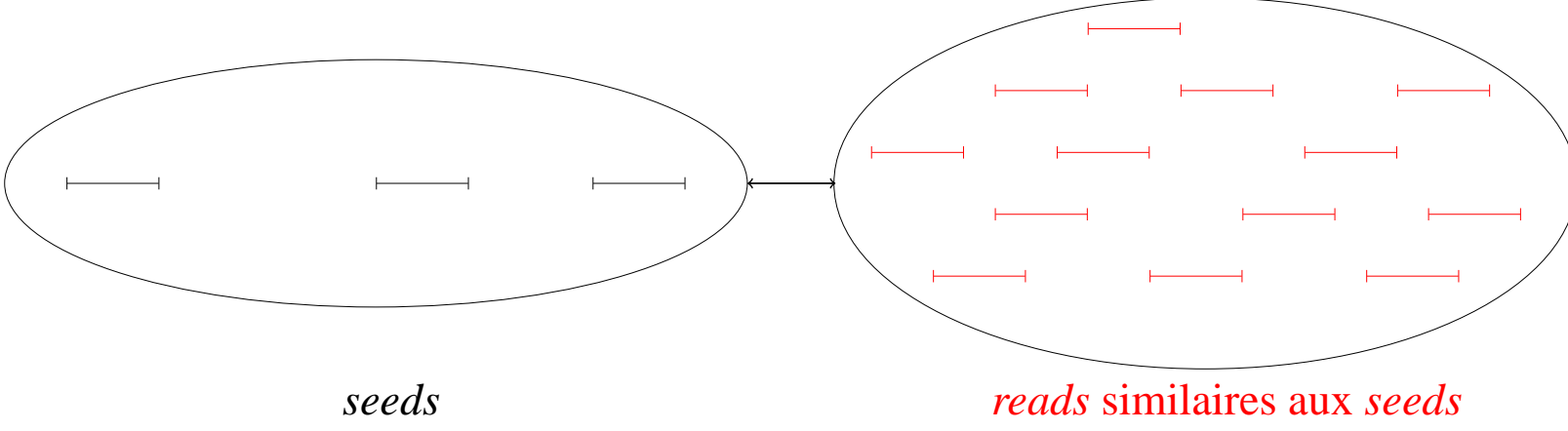
Première méthode

La première méthode de synthèse des *NaS* repose sur les étapes suivantes :

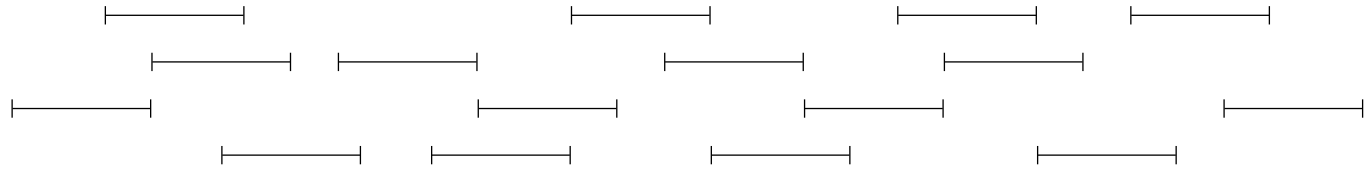
1. Alignement des *reads* courts sur le *read* long *template*, afin de trouver les *seeds*



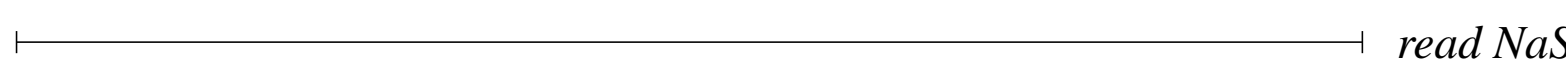
2. Recrutement de nouveaux *reads* courts, en recherchant des *reads* similaires aux *seeds*



3. Micro-assemblage de l'ensemble de *reads* obtenu

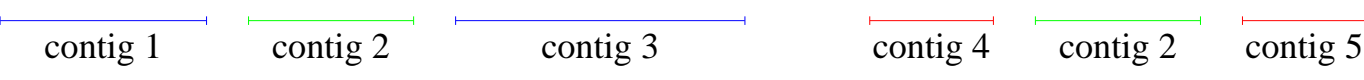


4. Obtention du *read* NaS

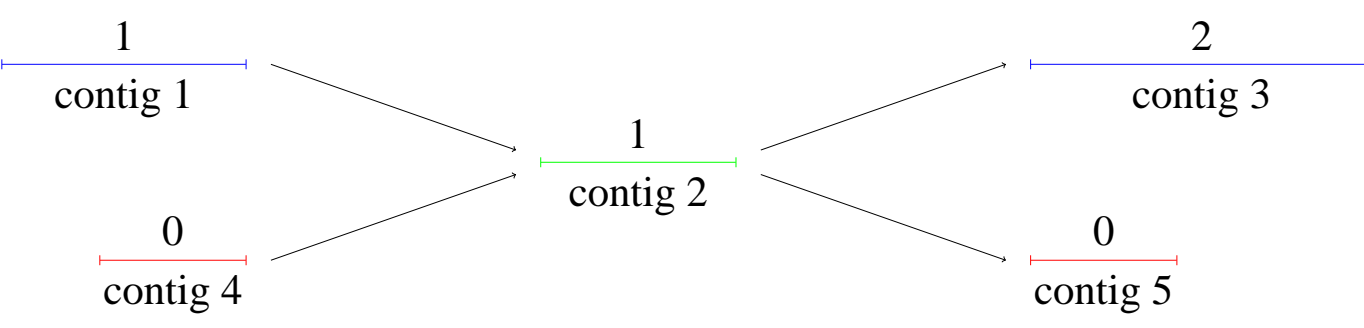


En général, un seul contig est produit par cette méthode, mais il est cependant possible que la phase de recrutement des *reads*, notamment dans les régions répétitives, recrute de mauvais *reads*, et que des contigs erronés, ne devant pas être associé au *template*, soient alors produits. Pour résoudre ce problème, et ne produire qu'un seul contig en sortie, il suffit d'employer la démarche suivante :

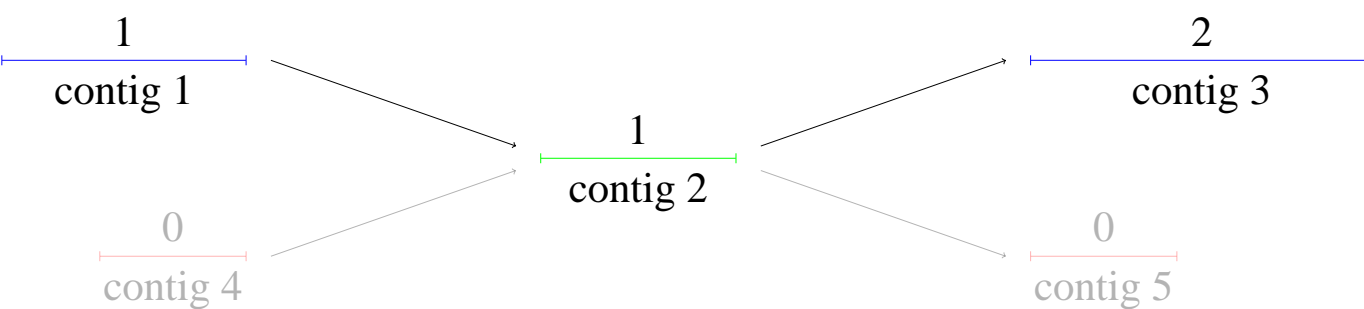
1. Obtention de plusieurs contigs



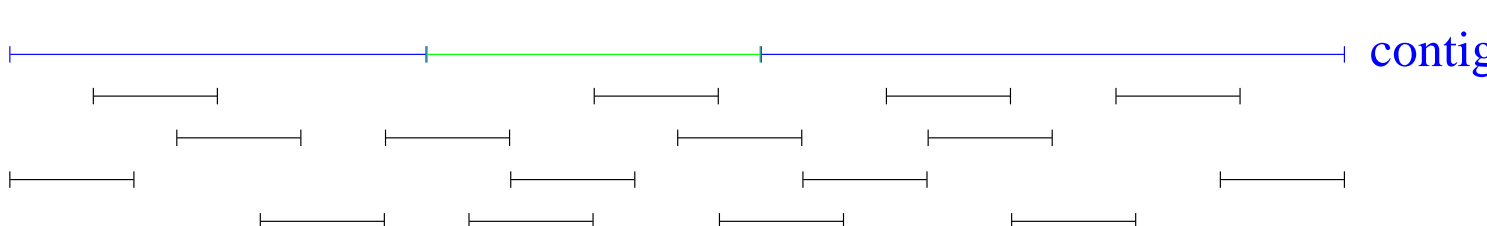
2. Construction du graphe des contigs, weighted par le degré de couverture des contigs par les *seeds*



3. Sélection du chemin optimal, passant par les contigs ayant le plus haut degré de couverture par les *seeds*



4. Vérification du contig produit par alignement des *reads* courts, et acceptation si la couverture est suffisante



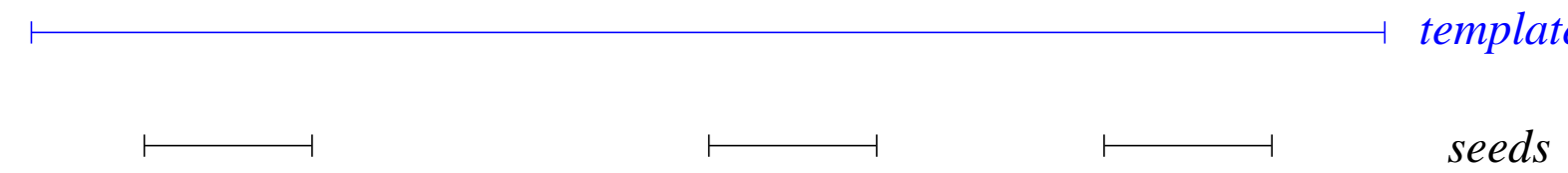
Le temps de traitement d'un *read* long, et donc la synthèse d'un *read* NaS, prend en moyenne moins d'une minute, la majorité de ce temps étant provoquée par la méthode peu efficace de recrutement de *reads* similaires. Celle-ci est réalisée en indexant l'ensemble des k -mers des *seeds* dans celui-ci. Un *seed* et un *read* non aligné sont alors considérés comme similaires s'ils partagent au moins t k -mers ne se chevauchant pas.

La synthèse de *reads* NaS par cette méthode a été testée sur un ensemble de 66 492 *reads* longs séquencés par MinION, et à l'aide de plusieurs sous-ensembles de *reads* courts de la technologie Illumina. 11 275 *reads* NaS, d'une longueur maximale de 59 863, ont ainsi été produits. Seulement 17% des *reads* longs ont donc produit un NaS, ce qui est du au fort taux d'erreurs des *reads* MinION. De plus, 97% des *reads* ainsi synthétisés ont pu être alignés sur le génome de référence sans aucune erreur, prouvant ainsi l'efficacité de cette méthode.

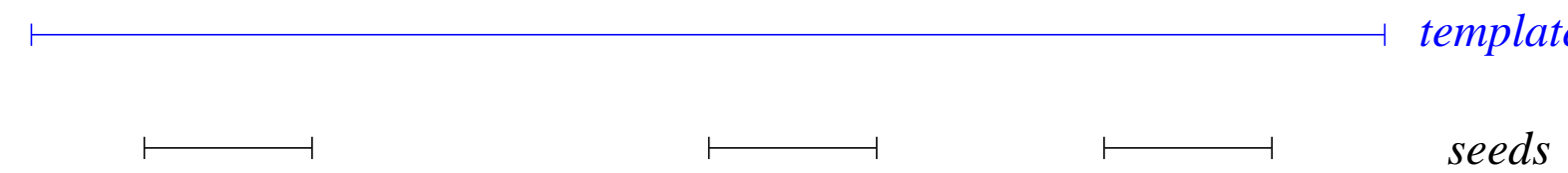
Notre méthode

Notre méthode repose sur le même principe que la méthode précédente, mais vise à diminuer le temps d'exécution en proposant une méthode différente pour le recrutement de *reads*. Nous alignons donc tout d'abord les *reads* courts sur les *reads* longs *templates*, en se fixant un seuil $lmin$, et récupérons les *reads* :

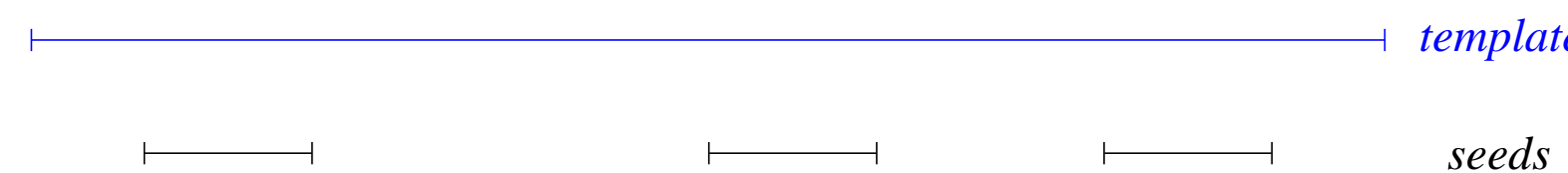
- Totalemment alignés



- Avec un préfixe de longueur $\geq lmin$ aligné



- Avec un suffixe de longueur $\geq lmin$ aligné



Ces différents ensembles de *reads* sont ajoutés à trois listes, triées en fonction des positions de début d'alignements préalablement calculées. Ces listes sont alors parcourues en parallèle, afin de recruter de nouveaux *reads* similaires aux *seeds*, en considérant qu'un *read* est similaire à un *seed* si son préfixe (respectivement son suffixe) correctement aligné chevauche ce *seed* sur une longueur supérieur ou égale au seuil fixé $lmin$.

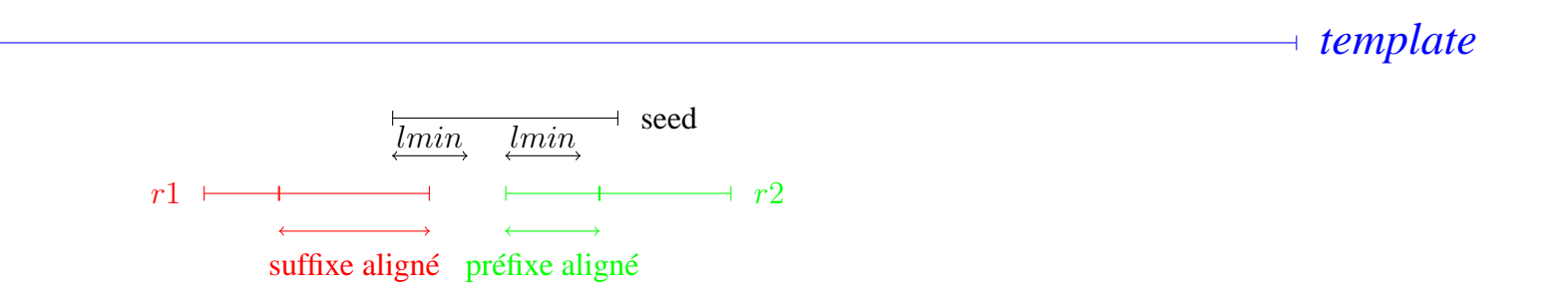


FIGURE 1: Illustration de la similarité entre *seeds* et *reads* courts partiellement alignés. $lmin$ représente le seuil permettant de déterminer cette similarité. Ici, on remarque que le *seed* est similaire à $r2$, mais pas à $r1$.

Lors du processus de recrutement, la liste des *seeds* est mise à jour, afin de prendre un compte les allongements des alignements provoqués par les recrutements, et ainsi couvrir d'avantage le *read* long *template* considéré.

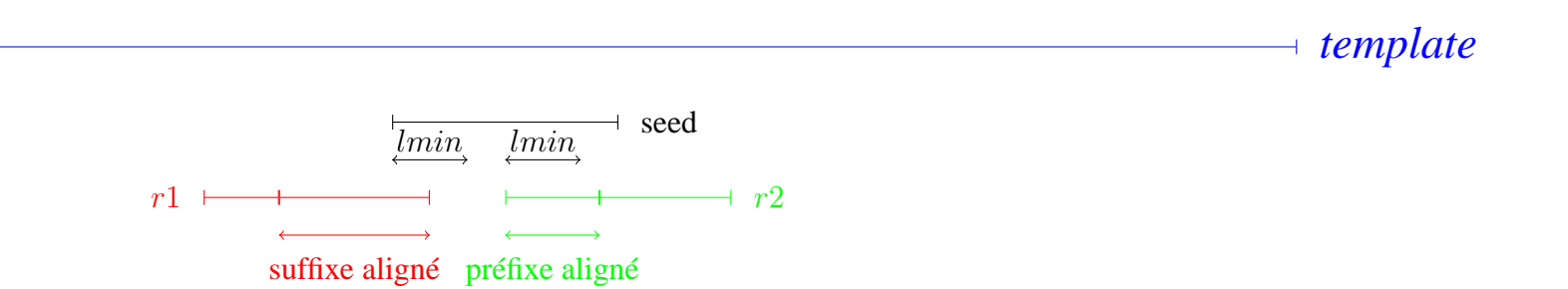


FIGURE 2: Illustration du processus de recrutement pour un *seed*. Le *seed* est mis à jour avec les parties ne le chevauchant pas des *reads* ayant un préfixe ou un suffixe aligné, afin d'allonger l'alignement.