

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡

INTRODUCTION

Principe du séquençage : À partir de la molécule d'ADN, retrouver la séquence de paires de bases constituant le message génétique.

Les séquenceurs automatiques ont permis la naissance de la génomique

De nombreuses applications, en santé, en biologie, en paléontologie, et même dans le milieu médico-légal...

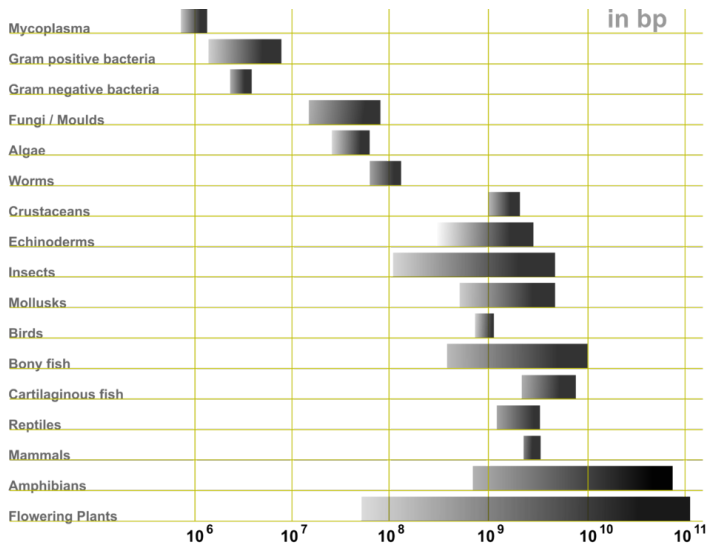
INTRODUCTION

Taille des génomes : variable de quelques millions de paires de bases chez les bactéries à plusieurs milliards de paires de bases chez la plupart des animaux et des plantes.

Mais les séquenceurs automatiques ne peuvent proposer que des séquences de quelques centaines de paires de base à la fois...

Problème : On se retrouve avec un gigantesque PUZZLE à assembler.

INTRODUCTION



- 1973 : Première publication d'une séquence de 24 bp

3'--ACCTTAACTTCGCTATTGTTAA--5'

- 1973 : Première publication d'une séquence de 24 bp

5'--TGGAA TTGTGAGCGGATAA CAATT--3'

3'--ACCTTAACACTCGCCTATTGTTAA--5'
- 1977-80 : Méthode de séquençage Sanger – Prix Nobel pour Wally Gilbert et Fred Sanger

- 1973 : Prem

- 1977-80 : Méthode de séquençage Sanger – Prix Nobel pour Wally Gilbert et Fred Sanger
- 1983 : Développement de la technique de PCR (Réaction en chaîne par polymérase)

- 1973 : Prem

- 1977-80 : Méthode de séquençage Sanger – Prix Nobel pour Wally Gilbert et Fred Sanger
- 1983 : Développement de la technique de PCR (Réaction en chaîne par polymérase)
- 1987 : 1er séquenceur automatique : Applied Biosystems Prism 373

- 1973 : Première publication d'une séquence de 24 bp

3'--ACCTTAACACTCGCCATTGTTAA--5'

- 1977-80 : Méthode de séquençage Sanger – Prix Nobel pour Wally Gilbert et Fred Sanger
- 1983 : Développement de la technique de PCR (Réaction en chaîne par polymérase)
- 1987 : 1er séquenceur automatique : Applied Biosystems Prism 373



- 1998 : Génome de *C. elegans* séquencé

- 1973 : Prem

- 1987 : 1er séquenceur automatique : Applied Biosystems Prism 373

-

-

NEXT GENERATION SEQUENCING (NGS)

- 2005 : 1er système de NGS : 454 (Roche) Life Sciences
- 2006 : 1er séquenceur de chez Solexa (Illumina)
- 2007 : 1er séquenceur de chez Applied Biosystems : SOLiD
- 2011 : 1er séquenceur de chez Ion Torrent : PGM

SHOTGUN

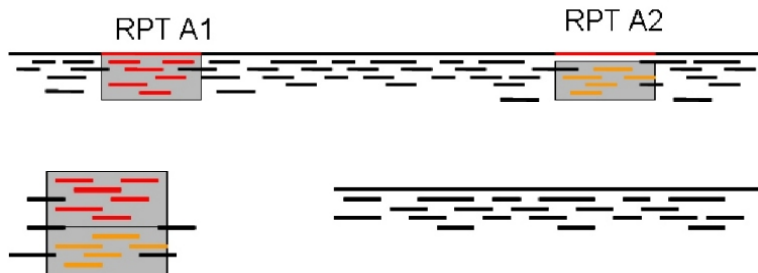
On casse l'ADN en millions de fragments aléatoires qui sont "lus" par le séquenceur, après amplification (PCR = Polymérase Chain Reaction)

Ensuite, la séquence totale est reconstituée grâce aux chevauchements recherchés par des méthodes informatiques.

SHOTGUN : LIMITATIONS

- Erreurs de lecture (recours à des méthodes inexactes, mais moins de précision sur le génome final)
- Impossibilité de séquencer correctement certaines zones du génome
- Problèmes dus à la présence de répétitions \Rightarrow fausse la recherche de chevauchements optimaux et produit des erreurs à l'assemblage.

Une solution : Masquer les répétitions, mais coûteux en ressources...



- Production de génome

- Production de génome
- Profil d'expression des transcrits

APPLICATIONS

- Production de génome
- Profil d'expression des transcrits
- Relation entre les facteurs de transcription

APPLICATIONS

- Production de génome
- Profil d'expression des transcrits
- Relation entre les facteurs de transcription
- Variations structurelles

APPLICATIONS

- Production de génome
- Profil d'expression des transcrits
- Relation entre les facteurs de transcription
- Variations structurelles
- Métagénomique
- ...

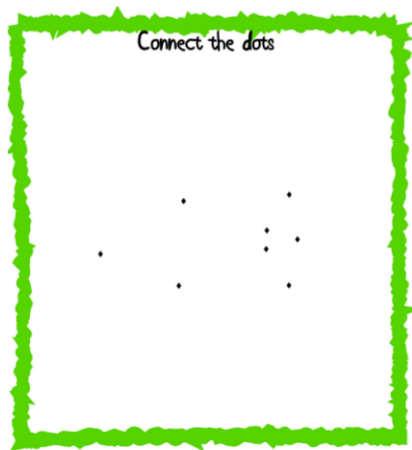
TYPES DE DONNÉES

	Capacity	Speed	Read Length	Homopolymers	Cost/run	Amplification
454 Roche	35-700 Mb	10-23 hours	400-700 bp	-	5.000 €	Yes
SOLiD	90-180 Gb	7-12 days	75 bp	+	5.000 €	Yes
Illumina	6-600 Gb	2-14 days	100-250 bp	+	10.000-20.000 €	Yes
Ion Torrent	20 Mb- 1Gb	4,5 hours	200 bp	-	1.000-2.000 €	Yes
Helicos	35 Gb	8 days	35 bp	+	20.000 €	No
PacBio	1Gb	30 minutes	3000 bp	+	600-800 €	No

d'après Andy Vierstraete

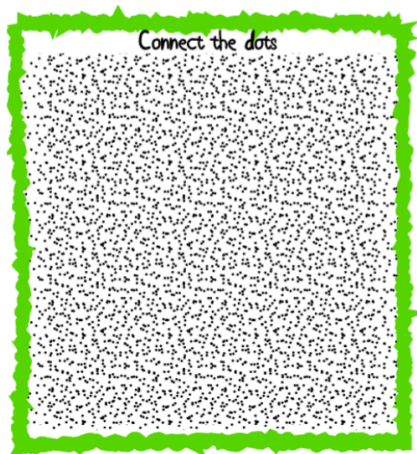
- le petites pièces
milliers de grosses
onne.
- le génome de

LE PROBLÈME DE L'ASSEMBLAGE

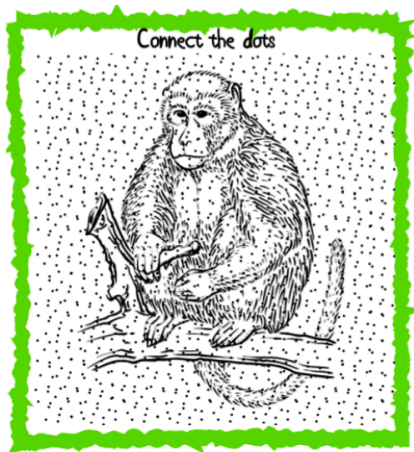


d'après Andy Vierstraete

LE PROBLÈME DE L'ASSEMBLAGE



LE PROBLÈME DE L'ASSEMBLAGE



LE PROBLÈME DE L'ASSEMBLAGE



un génome.

Comments

me de 5Mb donne

e reads chevauchants,

Économie de boucher

ale?

et que $K = \frac{PL}{n}$ est

UN BRIN DE THÉORIE

Alors la théorie de Landier et Waterman indique que l'on peut modéliser l'apparition d'une base dans un nombre donné de reads par une loi de Poisson.

La probabilité d'une base d'apparaître dans x reads est alors

$$\frac{P^x e^{-P}}{x!}$$

En particulier, la probabilité d'une base de n'apparaître dans aucun read est e^{-P} .

UN BRIN DE THÉORIE

$$\text{Proportion du génome couverte} = 1 - e^{-P}$$

$$\text{Longueur totale des trous} = Le^{-P}$$

$$\text{Nombre total de trous} = Ke^{-P}$$

$$\text{Taille moyenne des trous} = \frac{n}{p}$$

Taille moyenne des contigs = $\frac{n}{P}e^P$

LE PROBLÈME SSP

Pour étudier le problème de l'assemblage de façon algorithmique : on le formalise comme un problème combinatoire sur des *mots*.

Formellement, l'assemblage se résume au problème suivant : étant donné un ensemble de mots $F = \{F_1, \dots, F_n\}$ (les fragments), trouver une superséquence S de longueur minimale tel que chaque mot de F est un sous-mot de la superséquence S .

C'est le problème SSP (*Shortest Superstring Problem*)

LE PROBLÈME SSP

Exemple

F1 : AATGCC

F2 : GCCTTACAC

F3 : ACACTG

F4 : ACTGAAGG

F5 : GAAGGTTTA

B : AATGCCTTACACTGAAGGTTTA

LE PROBLÈME SSP

Problème !

SSP est un problème NP-complet pour plus de trois mots et un alphabet de plus de deux lettres...

On va devoir donc étudier des heuristiques pour résoudre le problème de façon polynomiale...

LE PROBLÈME SSP

Preuve de NP-complétude : En TD...

MÉTHODE GLOUTONNE

L'heuristique la plus naturelle : la méthode gloutonne.

1. Calculer les scores de chevauchement pour chaque paire de fragment (alignement semi-global)
2. Sélectionner une paire avec le meilleur chevauchement, fusionner les deux fragments, et recommencer

Complexité ?

MÉTHODE GLOUTONNE

Exercice : Appliquer cette méthode à la famille de mots suivantes et retrouver la séquence initiale :

$$F_1 = \text{ACCTGA}$$

$$F_2 = \text{TGATTGC}$$

$$F_3 = \text{GCAGC}$$

$$F_4 = \text{AGCAA}$$

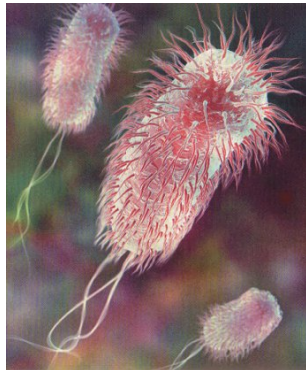
$$F_5 = \text{CAATG}$$

MÉTHODE GLOUTONNE



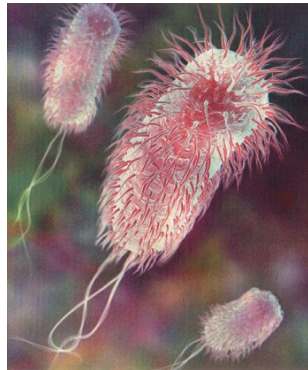
MÉTHODE GLOUTONNE

- Facile à implémenter



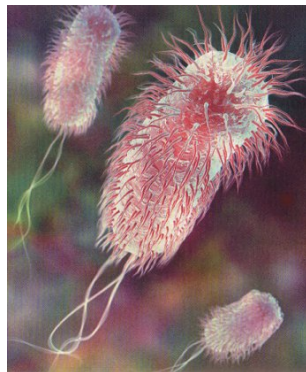
MÉTHODE GLOUTONNE

- Facile à implémenter
- Résultat assez satisfaisant mais...



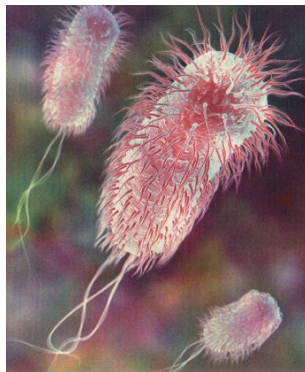
MÉTHODE GLOUTONNE

- Facile à implémenter
- Résultat assez satisfaisant mais...
- Prend énormément de mémoire : applicable seulement à des petits génomes



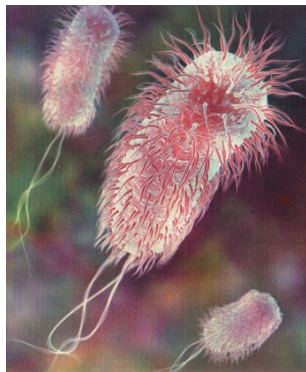
MÉTHODE GLOUTONNE

- Facile à implémenter
- Résultat assez satisfaisant mais...
- Prend énormément de mémoire : applicable seulement à des petits génomes
- Critère local pas adapté pour tenir compte des répétitions



MÉTHODE GLOUTONNE

- Facile à implémenter
- Résultat assez satisfaisant mais...
- Prend énormément de mémoire : applicable seulement à des petits génomes
- Critère local pas adapté pour tenir compte des répétitions
- \Rightarrow approches basées sur des graphes



OVERLAP-LAYOUT-CONSENSUS

Dans cette méthode, on se sert du fait que le calcul des chevauchements optimaux permet de construire un graphe sur lequel on peut travailler.

Overlap : trouver chaque paire de lectures qui se chevauchent ;

Layout : organiser ces lectures chevauchantes en une séquence contiguë (chemin hamiltonien) ;

Consensus : corriger les erreurs et générer une séquence consensus.

OVERLAP-LAYOUT-CONSENSUS

Definition

On définit $overlap(s_i, s_j)$ la longueur du plus long préfixe de s_j qui correspond à un suffixe de s_i .

OVERLAP-LAYOUT-CONSENSUS

On construit un graphe avec n sommets qui représentent les n chaînes s_1, s_2, \dots, s_n .

On insère les arcs portant un poids $overlap(s_i, s_j)$ entre les sommets s_i et s_j .

Problème : Trouver le chemin de poids maximal qui visite chaque sommet exactement une fois.

C'est une variante du problème du voyageur de commerce, qui est NP-complet !

OVERLAP

Phase Overlap

1. Calcul des chevauchements optimaux sur les reads :
alignement semi-global
2. Construction progressive du graphe : on ajoute les sommets (reads) et les arêtes (chevauchements optimaux de score acceptable), pondérées par les scores.
3. On enlève du graphe (ou plutôt on ne met pas dans le graphe) les reads entièrement contenus

Remarque : Le graphe est orienté.

LAYOUT

TSP

Solution exactes

Algorithme naïf : tester les $(n - 1)!$ tours possibles.

Programmation dynamique : Plus rapide mais complexité en espace $\mathcal{O}((n - 1)2^{n-2}) \dots$

Branch and bound

LAYOUT

TSP

Quelques heuristiques :

- Approche gloutonne
- Le plus proche voisin
- 2-interchange
- k -opt
- Lin-Kernighan
- Programmation linéaire (Concorde) : jusqu'à 65000 sommets en quelques heures...

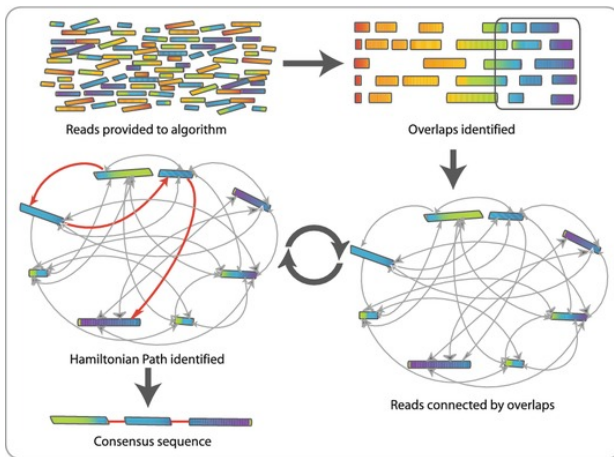
Assemblage : Cas asymétrique, on se ramène au cas symétrique en dupliquant les sommets.

CONSENSUS

Phase Consensus

On réalise cette phase en réalisant un alignement multiple par profil.

OVERLAP LAYOUT CONSENSUS



d'après J.Commins et al. *Computational Biology Methods and Their Application to the Comparative Genomics of Endocellular Symbiotic Bacteria of Insects* Bio.Proc.Online 2009

OVERLAP LAYOUT CONSENSUS

- L'information complète sur chaque read en entrée n'est nécessaire que pour les phases overlap et consensus.

OVERLAP LAYOUT CONSENSUS

- L'information complète sur chaque read en entrée n'est nécessaire que pour les phases overlap et consensus.
- La phase de raffinement du graphe permet de ne stocker qu'une quantité limitée d'information sur chaque chevauchement (coordonnées, longueur) \Rightarrow implémentation efficace en mémoire.

OVERLAP LAYOUT CONSENSUS

- L'information complète sur chaque read en entrée n'est nécessaire que pour les phases overlap et consensus.
- La phase de raffinement du graphe permet de ne stocker qu'une quantité limitée d'information sur chaque chevauchement (coordonnées, longueur) \Rightarrow implémentation efficace en mémoire.
- Mais trouver un chemin Hamiltonien est tout de même un problème NP-complet...

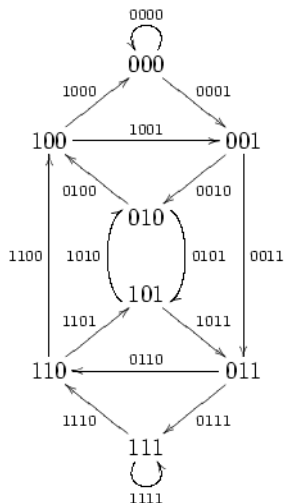
CHEMIN EULÉRIEN

On utilise dans cette approche un autre type de graphe lié aux reads : le graphe des k -mers, qui est un cas particulier de graphe de De Bruijn.

Definition

Soit un alphabet $\mathfrak{A} = \{\alpha_1, \dots, \alpha_n\}$ de taille n . On note $G_{n,L} = (S, A)$ le graphe de de Bruijn d'ordre L construit sur l'alphabet \mathfrak{A} . Les sommets S de $G_{n,L}$ représentent les L^n mots de longueur L construits à partir des lettres de l'alphabet \mathfrak{A} . Il existe un arc $A_{i,j}$ allant du sommet s_i au sommet s_j du graphe $G_{n,L}$ si le suffixe de longueur $L - 1$ du mot u et le préfixe de longueur $L - 1$ du mot v correspondent.

CHEMIN EULÉRIEN



CHEMIN EULÉRIEN

Théorème

Les graphes de De Bruijn sont eulériens.

Rappel : un graphe est eulérien si on peut parcourir le graphe en utilisant exactement une fois chaque arc et en revenant au point de départ.

Théorème (Théorème d'Euler (version orientée))

Un graphe orienté fortement connexe est eulérien si et seulement pour chaque sommet le degré entrant est égal au degré sortant.

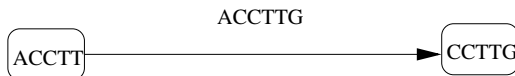
Séquence d'origine (read)

6-mers

- est petit pour capturer
 et pour éviter les
 que de détection des
 les k -mers dans une
 reurs par analyse de

CHEMIN EULÉRIEN

- Formalisation simple de la structure de graphe de De Bruijn



CHEMIN EULÉRIEN

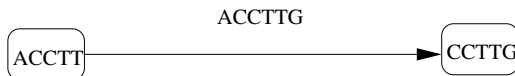
- Formalisation simple de la structure de graphe de De Bruijn



- Le problème du chemin eulérien est polynomial.

CHEMIN EULÉRIEN

- Formalisation simple de la structure de graphe de De Bruijn

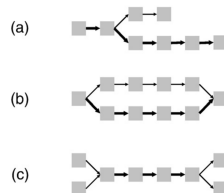


- Le problème du chemin eulérien est polynomial.
- Mais la taille de l'entrée peut tout de même être énorme, dépendamment de k et des reads \Rightarrow efforts à faire sur le côté "stockage" (algo du texte plus avancée...)

NETTOYER L'ASSEMBLAGE

La plupart des assembleurs ont des post-traitements pour réaliser les tâches suivantes :

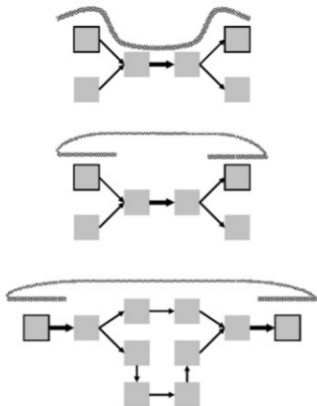
- Corriger les erreurs de séquençage (a)
- Détecter les polymorphismes (b)
- Gérer les répétitions (c)
- Détecter les régions mal assemblées



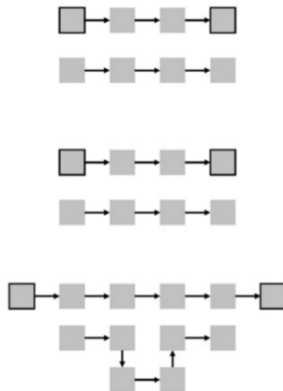
J.R.Miller et al., *Assembly algorithms for next-generation sequencing data*, Genomics, 2010

NETTOYER L'ASSEMBLAGE

(before)



(after)



FINITION

Il reste toujours des trous, des incohérences, des ambiguïtés à lever.

On réalise des expériences supplémentaires en laboratoire, on nettoie "à la main", afin de valider la correction de l'assemblage final.

C'est une opération très coûteuse, qui peut prendre des mois...

⇒ le temps et l'effort ne sont justifiés que pour les génomes "à haute priorité"

...les autres restent à l'état de brouillons (*draft genomes*)

FINITION

