

On the Linearization of Scaffolds sharing Repeated Contigs

Mathias Weller, **Annie Chateau**, Tom Davot, Rodolphe
Giroudeau

17 novembre 2017

DISCLAIMER

Be careful : An NP-complete problem could hide another one.
This talk is about the **third** one you'll meet.



INTRODUCTION

Sequencing is a technology used to infer genomic information out of DNA material.

- it produces short words, called *reads*, which have to be *assembled* to (try to) reconstruct the whole genomic sequence.
- Assembly can be modeled as an NP-complete problem (Shortest Superstring)¹

1. Do you follow? This is the first problem.

INTRODUCTION

Strategies for assembly :

- Greedy
- Overlap-Layout-Consensus
- De Bruijn graphs

INTRODUCTION

Results : sets of *contigs* of various sizes, disconnected

Far from the whole sequence...

Assembly doesn't take into account some pairing information on the reads.

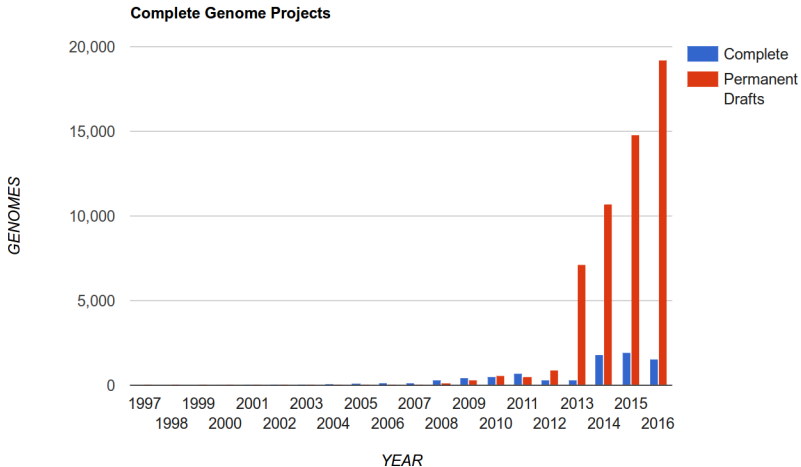
⇒ we need an additional step to use these information and (try to) produce chromosome-long sequences.

INTRODUCTION

Why?

- Observe genome-scale phenomena
- Improve reference genomes quality
- Lot of genomes have a "draft" status in databases

INTRODUCTION



<https://gold.jgi.doe.gov/statistics>

THE SCAFFOLDING PROBLEM²

To determine relative order and orientation of contigs, we need :

- Informations between contigs
 - ▶ pairing data (common, easy, cheap)
 - ▶ phylogenetic data (needs well assembled close species)
 - ▶ long reads (full of errors, expensive)
 - ▶ ...
- A weight on these information
 - ▶ number of pairs of reads
 - ▶ probabilistic measure
 - ▶ coverage depth
 - ▶ ...

2. yes, this is the second one, be patient

THE SCAFFOLDING PROBLEM

Data are modeled as a graph $G = (V, E)$:

- **Vertices** : contigs extremities
- **Edges** :
 - ▶ between both extremities of a given contig (contig edge)
 - ▶ between extremities of distinct contigs (inter-contigs edge)

Weight function : $w : E \rightarrow \mathbb{R}$.

WHAT ABOUT "SHARED CONTIGS" ?

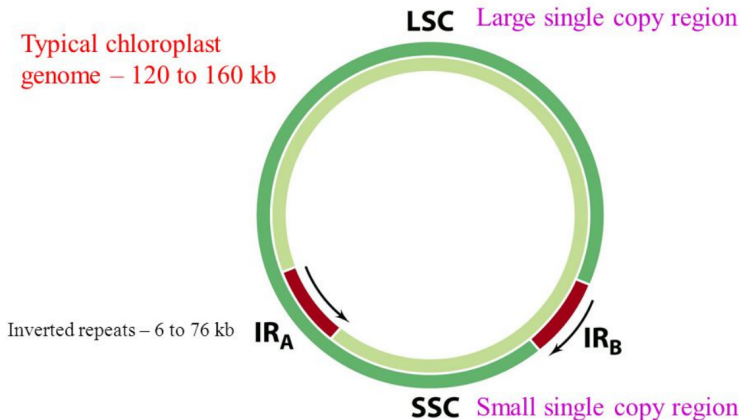
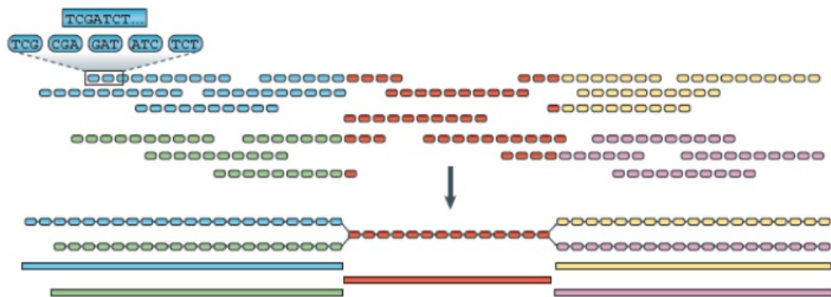


Figure 15-10 Brock Biology of Microorganisms 11/e
© 2006 Pearson Prentice Hall, Inc.

WHAT ABOUT "SHARED CONTIGS" ?



THE SCAFFOLDING PROBLEM WITH MULTIPLICITIES

Input : $G = (V, E)$, $w : E \rightarrow \mathbb{N}$, M^* perfect matching,
 $\sigma_p, \sigma_c, k \in \mathbb{N}$, $m : E \rightarrow \mathbb{N}$

Query : Does it exist a set S of σ_p alternating open walks and σ_c alternating closed walks covering G such that $w(S) \geq k$ and satisfying the maximal multiplicity constraint?

THE SCAFFOLDING PROBLEM WITH MULTIPLICITIES

Guess what?³

3. If you read footnotes, you already know

THE SCAFFOLDING PROBLEM WITH MULTIPLICITIES

Guess what?³

It is NP-complete!

3. If you read footnotes, you already know

THAT'S ANOTHER STORY BUT...

We have algorithms and exact methods to solve this problem efficiently on real instances.

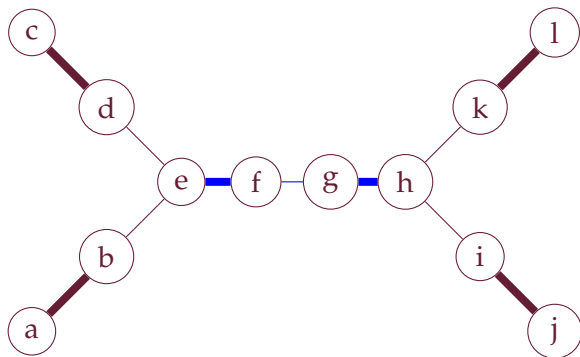
It scales : 1h30 to scaffold a mosquito genome...

Take repeats into account...

...everything seems to be perfect, but ...

SOLUTION GRAPH

Juste a little problem : the solution appears as a graph, not as a collection of cycles and paths. Multiples edges are not "attributed" to a particular path



LINEARIZATION OF SOLUTION GRAPH

Biologists like linear (or circular) genomes.

First solution : Convince biologists that a solution graph is cool.

It may take a while and lots of efforts...

Second solution : Transform the graph into sequences, without creating chimera

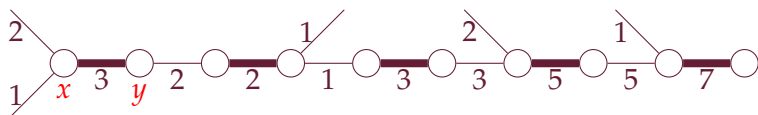
NOTATIONS

Definition

Let p be an alternating u - v -path in a solution graph. If all edges of p have the same multiplicity μ (that is, $m(e) = \mu$ for all $e \in p$), then p is called μ -uniform (or simply *uniform* if μ is unknown). Further, if p is μ -uniform and each of u and v is incident with a non-matching edge of multiplicity strictly less than μ , then p is called “ambiguous”.

NOTATIONS

Ambiguous path :



THE IDEAL CASE

Theorem

Let G be a solution graph. Then, G is made up of a unique multiset of alternating walks if and only if G does not contain ambiguous paths.

STRATEGIES TO GET RID OF AMBIGUOUS PATHS ?

Ignore. Chose an arbitrary multiset of walks making up G .
Weight is preserved, but risk to produce a chimeric sequence.

STRATEGIES TO GET RID OF AMBIGUOUS PATHS ?

- Ignore. Chose an arbitrary multiset of walks making up G .
Weight is preserved, but risk to produce a chimeric sequence.
- Clever. Chose walks that optimize some global criterion (i.e. N50). Again, risk to produce chimeric sequences.

STRATEGIES TO GET RID OF AMBIGUOUS PATHS ?

- Ignore. Chose an arbitrary multiset of walks making up G .
Weight is preserved, but risk to produce a chimeric sequence.
- Clever. Chose walks that optimize some global criterion (i.e. N50). Again, risk to produce chimeric sequences.
- Brutal. Isolate ambiguous paths by removing all non-matching edges incident to their extremities.
Weight could drastically be lowered.

STRATEGIES TO GET RID OF AMBIGUOUS PATHS ?

Ignore. Chose an arbitrary multiset of walks making up G .
Weight is preserved, but risk to produce a chimeric sequence.

Clever. Chose walks that optimize some global criterion (i.e. N50). Again, risk to produce chimeric sequences.

Brutal. Isolate ambiguous paths by removing all non-matching edges incident to their extremities.
Weight could drastically be lowered.

Semi-brutal. Choose a proper set of endpoints of ambiguous path and remove all non-matching edges incident to it. Optimize a criteria.

STRATEGIES TO GET RID OF AMBIGUOUS PATHS ?

- Ignore. Chose an arbitrary multiset of walks making up G . Weight is preserved, but risk to produce a chimeric sequence.
- Clever. Chose walks that optimize some global criterion (i.e. N50). Again, risk to produce chimeric sequences.
- Brutal. Isolate ambiguous paths by removing all non-matching edges incident to their extremities. Weight could drastically be lowered.
- Semi-brutal. **Choose a proper set of endpoints of ambiguous path and remove all non-matching edges incident to it. Optimize a criteria.**

SEMI-BRUTAL CUT⁴

Input : a solution graph (G, M^*, w, m) and some $k \in \mathbb{N}$

Query : Is there a set X of extremities of ambiguous paths in G such that removing all non-contig edges incident to vertices of X destroys all ambiguous paths and the score of X is at most k ?

4. Here it is!

SCORING FUNCTION FOR SBC

Cut score. Pay one per side of an ambiguous path that is cut :
 $\text{score}(X) := |X|.$

Path score. Pay one for each multiplicity that is cut :
 $\text{score}(X) := \sum \{m(uv) \mid uv \in E \setminus M^* \wedge uv \cap X \neq \emptyset\}.$

Weight score. Pay the total cost of edges that are cut :
 $\text{score}(X) := \sum \{m(uv) \cdot w(uv) \mid uv \in E \setminus M^* \wedge uv \cap X \neq \emptyset\}.$

COMPLEXITY OF SBC



COMPLEXITY

Theorem

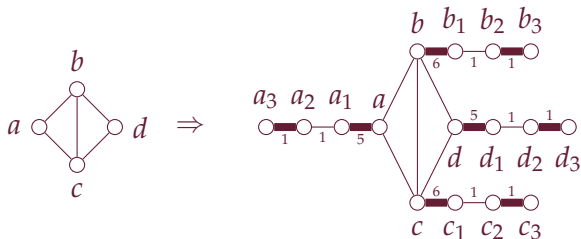
It is \mathcal{NP} -hard to decide whether all ambiguous paths in a solution graph can be destroyed by removing the non-matching edges incident to at most k endpoints.

Theorem

It is \mathcal{NP} -hard to decide whether a solution graph without ambiguous paths can be obtained by removing at most k non-matching edges.

(IDEA OF) PROOFS

For Cut-Score : reduction from Vertex Cover



Corollary

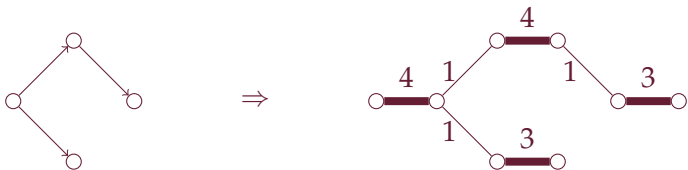
SBC with cut-score cannot be solved in $2^{o(n)}$ time unless ETH fails, and cannot be approximated within a ratio of 1.3606 (resp. better than factor 2) unless $\mathcal{P} = \mathcal{NP}$ (resp. UGC fails).

(IDEA OF) PROOFS

For Path-Score : reduction from Transitivity Deletion

Input : A triangle-free directed acyclic graph (V, A) and $k \geq 0$.

Question : Is there an $A' \subseteq A$ with $|A'| \leq k$ and $(V, A \setminus A')$ is transitive?

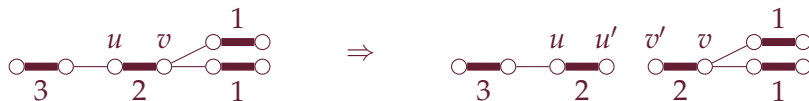


For Weight-Score : Path-Score is a special case of Weight-Score

POLYNOMIAL CASES

Rule

Let $uv \in M^*$ be a contig edge that does not occur in ambiguous paths and let u and v have degree at least two. Then, remove uv , add new vertices u' and v' and add the contig edges uv' and vu' with multiplicity $m(uv)$.



POLYNOMIAL CASES

Trees : dynamic programming

$c(x)$ = cost of a solution below x in which all non-contigs incident with x are cut

$\bar{c}(x)$ = cost of any other solution below x .

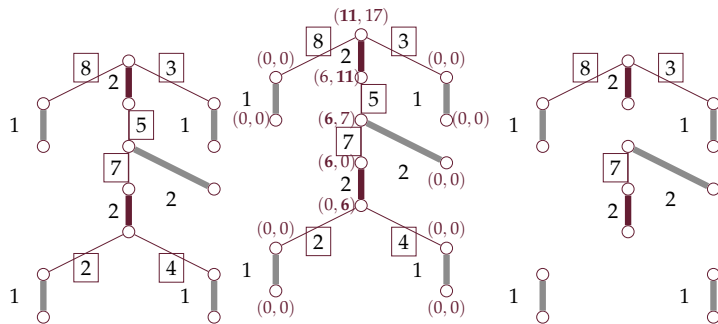
If x is a leaf of G , $c(x) = \bar{c}(x) = 0$.

For any non-leaf x , we set

$$c(x) = \sum_{y \in \text{Children}(x)} \min(\bar{c}(y), c(y)) + \sum_{y \in \text{Children}(x) \setminus \{M^*(x)\}} w_{xy}$$

$$\bar{c}(x) = \begin{cases} c(M^*(x)) & \text{if } M^*(x) \text{ is below } x \\ 0 & \text{otherwise} \end{cases} + \sum_{y \in \text{Children}(x) \setminus \{M^*(x)\}} \min(\bar{c}(y), c(y) + w_{xy})$$

POLYNOMIAL CASES



POLYNOMIAL CASES

Max degree two : collection of cycles and paths.

ILP formulation yields totally unimodular matrices

$$\begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 & 0 & 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & -1 & 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & -1 & 0 & 0 & 1 \end{bmatrix}$$

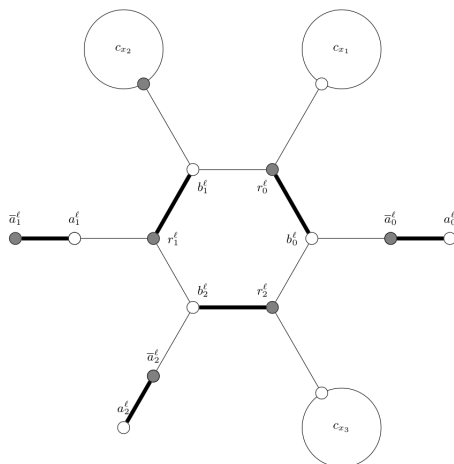
TOWARDS THE FRONTIER

Theorem

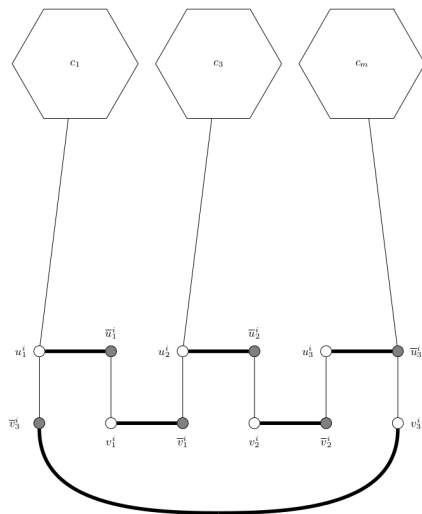
The problem SBC for Cut-Score is \mathcal{NP} -complete, even if the graph is bipartite, planar, has maximum degree three and has its multiplicities in the set $\{1, 2\}$.

TOWARDS THE FRONTIER

Idea of proof : reduction from 3-SAT



TOWARDS THE FRONTIER



CONCLUSION

There is a lot of work left :

- Approximation
- Define and test heuristics
- Test ILP
- Lower bounds of complexity