

COMPLEXITY AND POLYNOMIAL-TIME APPROXIMATION ALGORITHMS AROUND THE SCAFFOLDING PROBLEM



Annie Chateau – Rodolphe Giroudeau
LIRMM, Montpellier

AICoB 2014 – Tarragona, July 1-3



Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

SOMMAIRE

The scaffolding problem

What is scaffolding?

The scaffold graph

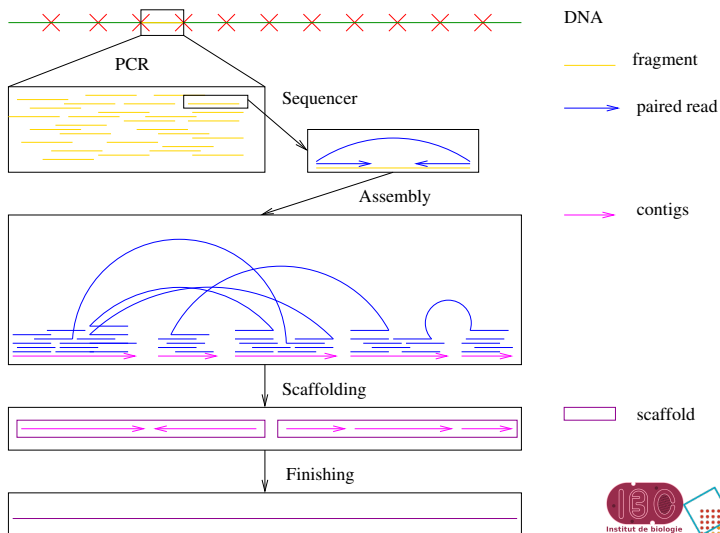
The (σ_p, σ_c) -Scaffold Problems

Complexity results

Approximability results

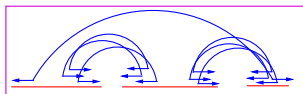
Conclusion and perspectives

FROM THE MOLECULE TO GENOMIC SEQUENCE



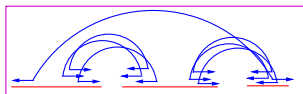
GENERATION OF THE SCAFFOLD GRAPH

- Paired reads are often used for scaffolding



GENERATION OF THE SCAFFOLD GRAPH

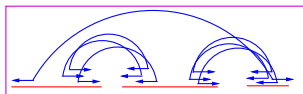
- Paired reads are often used for scaffolding



- **Mapping** of the reads on the contigs and compute the number of pairs of reads supporting each hypothesis
⇒ the scaffold graph

GENERATION OF THE SCAFFOLD GRAPH

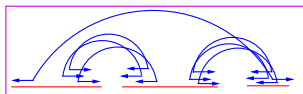
- Paired reads are often used for scaffolding



- **Mapping** of the reads on the contigs and compute the number of pairs of reads supporting each hypothesis
⇒ the scaffold graph
- **Nodes** = contigs extremities

GENERATION OF THE SCAFFOLD GRAPH

- Paired reads are often used for scaffolding



- **Mapping** of the reads on the contigs and compute the number of pairs of reads supporting each hypothesis
⇒ the scaffold graph

- **Nodes** = contigs extremities

- **Edges** = $\begin{cases} \text{contigs} \\ \text{stories} \end{cases}$ (weight = number of pairs of reads)

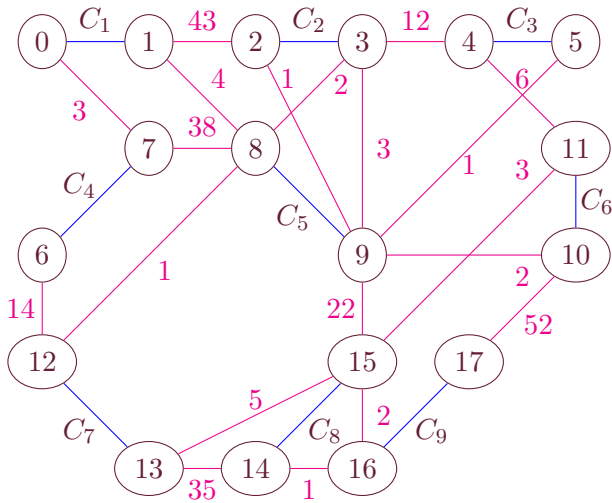
CONSTRAINED GENOMIC STRUCTURE

We consider a given genomic structure $(\sigma_p, \sigma_c) \in \mathbb{N} \times \mathbb{N} \setminus \{(0, 0)\}$, where :

- σ_p is the number of linear chromosomes (paths)
- σ_c is the number of circular chromosomes (cycles)

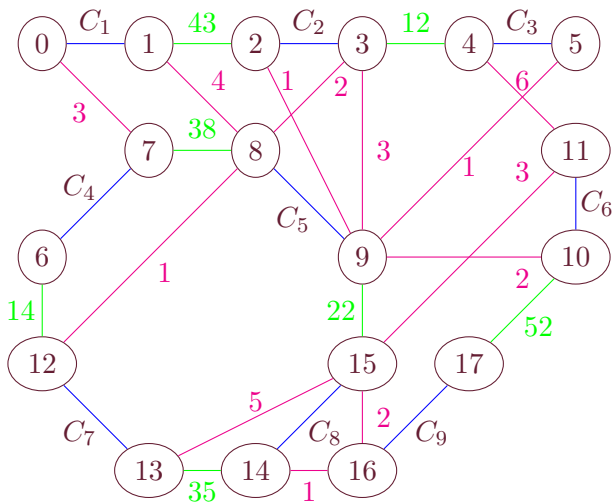
What is the best possible way to scaffold according to this constrained structure ?

THE SCAFFOLD GRAPH



$n = 9$
 $\sigma_p = 2$
 $\sigma_c = 1$
 18 vertices
 9 contig
 edges
 20 weighted
 edges

THE SCAFFOLD GRAPH



$$n = 9$$

$$\sigma_p = 2$$

$$\sigma_c = 1$$

18 vertices

9 contig

edges

20 weighted

edges

THE (σ_p, σ_c) -SCAFFOLD PROBLEMS

(Decision) (σ_p, σ_c) -SCAFFOLD PROBLEM (SP) :

Instance : $G = (V, E)$ graph with $2n$ vertices. Let M^* be a perfect matching of G , and $(\sigma_p, \sigma_c) \in \mathbb{N} \times \mathbb{N} \setminus \{(0, 0)\}$.

Question : Does it exist a vertex disjoint collection of exactly σ_p alternating-paths and σ_c alternating-cycles, covering the vertices of G ?

(Optimization) MIN/MAX- (σ_p, σ_c) -SCAFFOLD PROBLEM :

Instance : $G = (V, E, w)$ graph with $2n$ vertices. Let M^* be a perfect matching of G , and $(\sigma_p, \sigma_c) \in \mathbb{N} \times \mathbb{N} \setminus \{(0, 0)\}$.

Question : Find a vertex disjoint collection of exactly σ_p alternating-paths and σ_c alternating-cycles, covering the vertices of G , and of minimal (resp. **maximal**) total weight.

SOMMAIRE

The scaffolding problem

Complexity results

Polynomial case

NP-completeness

Approximability results

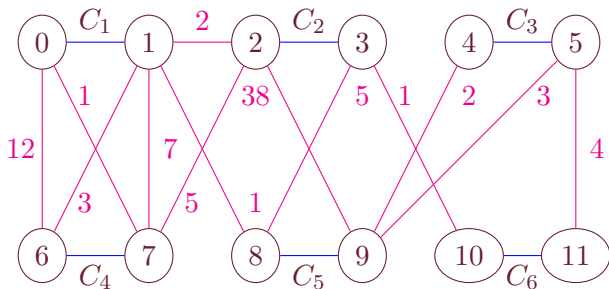
Conclusion and perspectives



POLYNOMIAL CASE

Theorem

When the size n of the matching M^* satisfies $n = \sigma_p + 2\sigma_c$, the problem (σ_p, σ_c) -SP (resp. MIN/MAX- (σ_p, σ_c) -SP) is polynomial.



$$n = 6$$

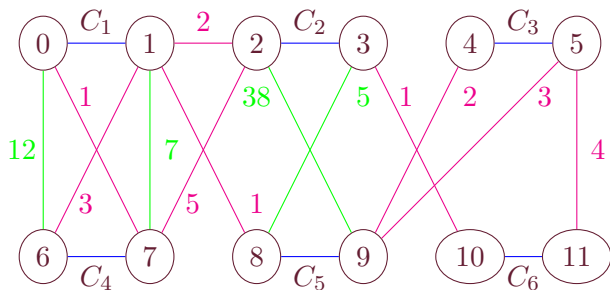
$$\sigma_p = 2$$

$$\sigma_c = 2$$

POLYNOMIAL CASE

Theorem

When the size n of the matching M^* satisfies $n = \sigma_p + 2\sigma_c$, the problem (σ_p, σ_c) -SP (resp. MIN/MAX- (σ_p, σ_c) -SP) is polynomial.



$$n = 6$$

$$\sigma_p = 2$$

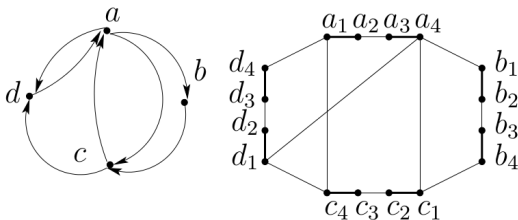
$$\sigma_c = 2$$

NP-COMPLETENESS

Theorem

The problem (σ_p, σ_c) -SP is NP-complete, even if the graph is bipartite.

Idea : reduction from DIRECTED HAMILTONIAN CYCLE



SOMMAIRE

The scaffolding problem

Complexity results

Approximability results

Approximability

Minimization

Maximization

Conclusion and perspectives



APPROXIMABILITY

Approximation algorithm \mathcal{A} with guaranteed ratio $\rho > 1$:

- $S_{\mathcal{A}}$ the solution given by \mathcal{A}



APPROXIMABILITY

Approximation algorithm \mathcal{A} with guaranteed ratio $\rho > 1$:

- $S_{\mathcal{A}}$ the solution given by \mathcal{A}
- S_{opt} the optimal solution

APPROXIMABILITY

Approximation algorithm \mathcal{A} with guaranteed ratio $\rho > 1$:

- $S_{\mathcal{A}}$ the solution given by \mathcal{A}
- S_{opt} the optimal solution
- $w(S_{\mathcal{A}}) \leq w(S_{opt}) \leq w(S_{\mathcal{A}}) \times \rho$ (maximization)

NON-APPROXIMABILITY RESULTS FOR THE MIN PROBLEM

Theorem

The problem $\text{MIN}-(\sigma_p, \sigma_c)\text{-SP}$ is non-approximable, unless $\mathcal{P} = \mathcal{NP}$, even if the graph is bipartite.

POLYNOMIAL-TIME APPROXIMATION ALGORITHMS FOR THE MAX PROBLEM

- Greedy algorithm : ratio 3
- Perfect matching based algorithm : ratio 3 in the general case, ratio 2 in the case $(0, 1)$

GREEDY ALGORITHM

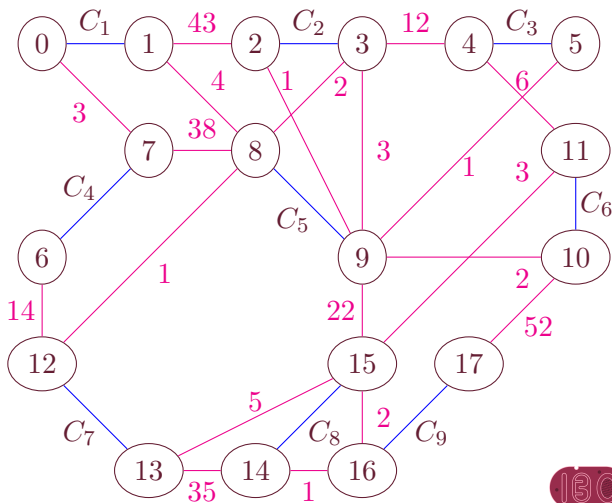
- Sort the edges by decreasing weight
- Add the edges in that order, if it is compatible with the desired structure
- When an edge is chosen, eliminate the contiguous edges
- Until we obtain exactly σ_p paths and σ_c cycles

Theorem

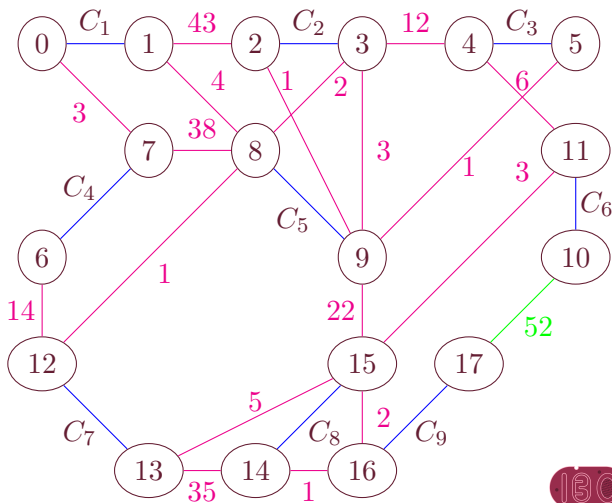
This algorithm gives a solution for the MAX- (σ_p, σ_c) -SCAFFOLD PROBLEM in complete graphs with non-negative weight, with a tight approximation rate of 3 when $n \geq 2(\sigma_p + 2\sigma_c)$, and a time complexity $\mathcal{O}(n^2 \log n)$.



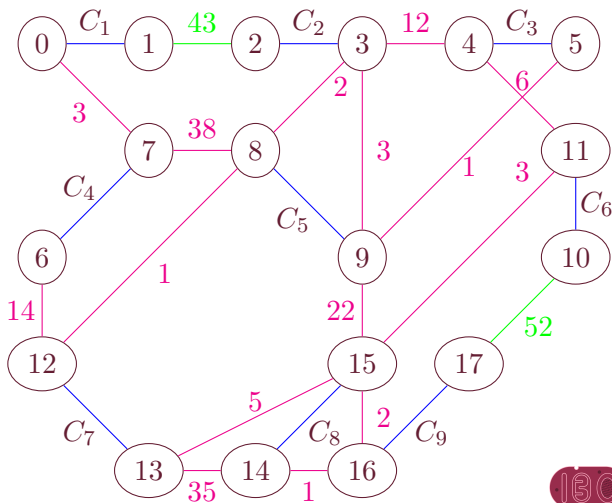
GREEDY ALGORITHM



GREEDY ALGORITHM

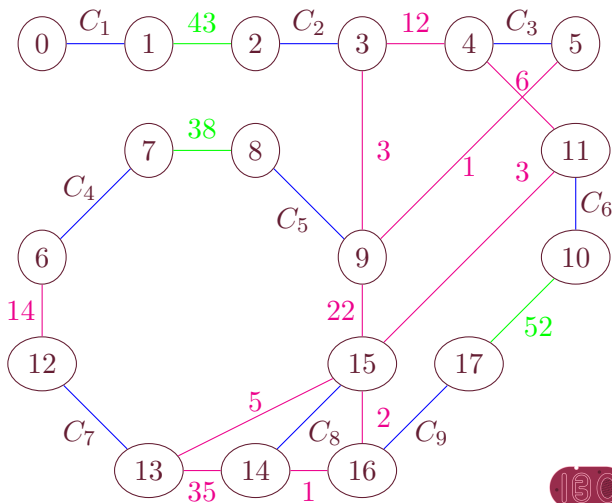


GREEDY ALGORITHM



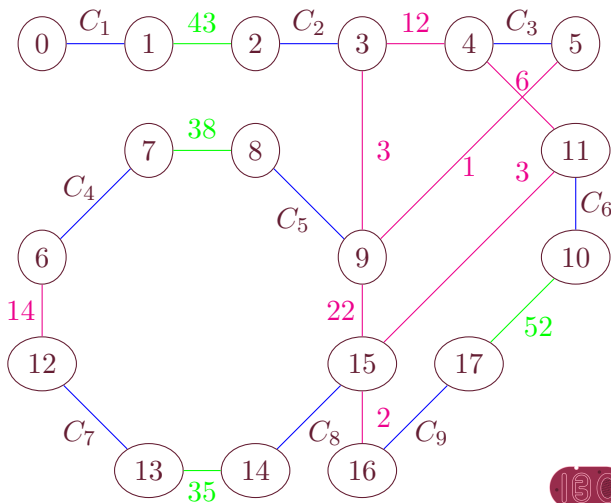
Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

GREEDY ALGORITHM

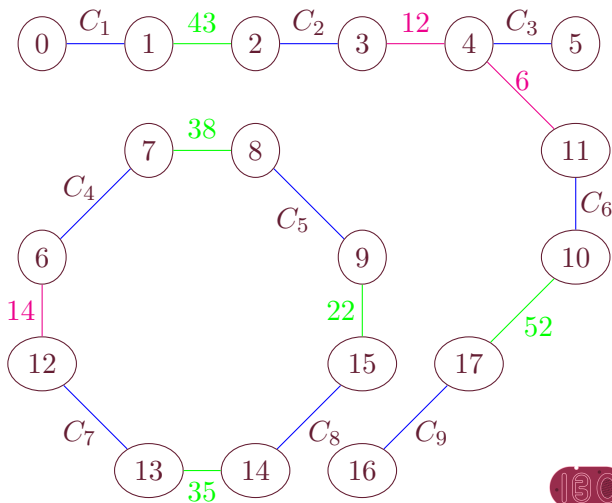


Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

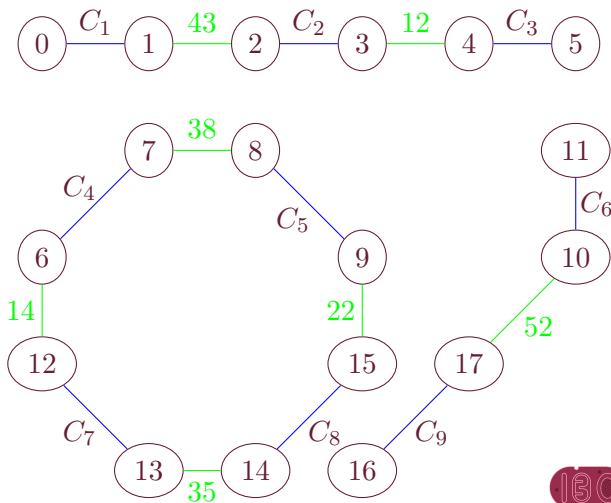
GREEDY ALGORITHM



GREEDY ALGORITHM



GREEDY ALGORITHM



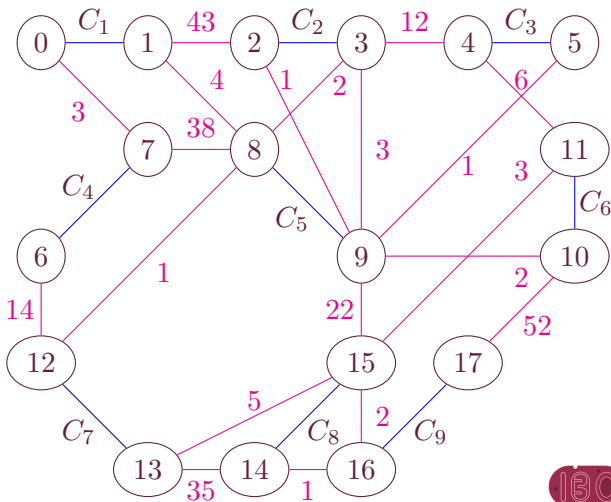
PERFECT MATCHING ALGORITHM

- First compute a perfect matching of maximal weight \rightarrow cycle cover of the graph
- Select a subset of "removable edges" that can be temporarily erased
- Reconnect the paths to form σ_p paths and σ_c cycles

Theorem

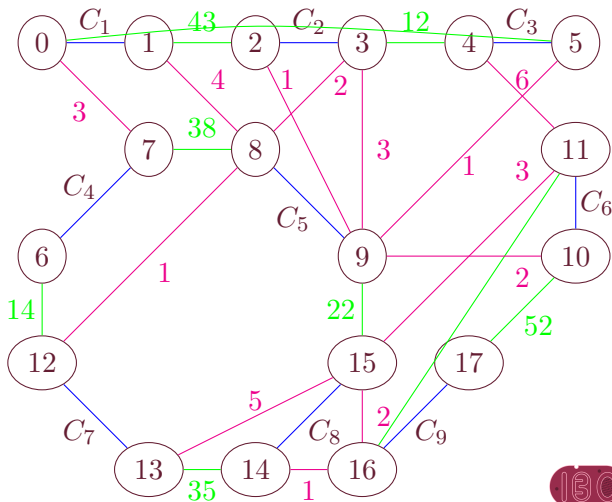
When $n \geq 2(\sigma_p + 2\sigma_c)$, this algorithm provides a solution for the MAX- (σ_p, σ_c) -SP in complete graphs with non-negative weights, with a tight approximation ratio of three and a time complexity $\mathcal{O}(n^3)$.

PERFECT MATCHING ALGORITHM

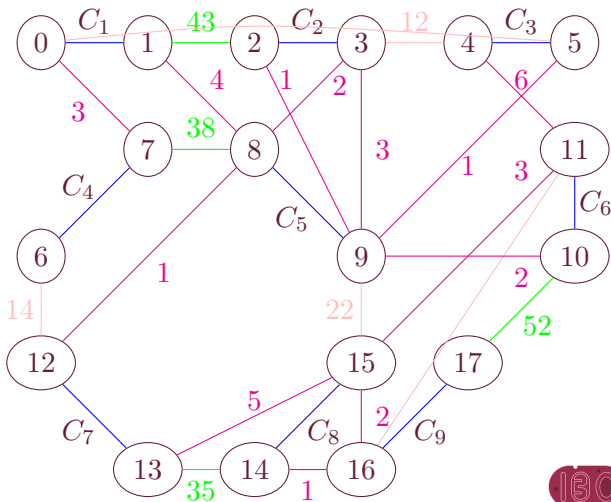


Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

PERFECT MATCHING ALGORITHM



PERFECT MATCHING ALGORITHM



Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

SOMMAIRE

The scaffolding problem

Complexity results

Approximability results

Conclusion and perspectives

Conclusion

Completing the model

More on approximability

Biological validation



CONCLUSION

Problems	Complexity	Approximability	
	Decision	Min	Max
(σ_p, σ_c) -SP	\mathcal{NP} -Complete	Non-approximable	Ratio 3
$(0, 1)$ -SP	\mathcal{NP} -Complete	Non-approximable	Ratio 2

COMPLETING THE MODEL

- **Multiplicities on the contigs** $m : \mathcal{C} \rightarrow \mathbb{N}$. Problem : find an optimal vertex cover (maximal weight) such that each contig edge $e \in \mathcal{C}$ is used exactly $m(e)$ times (or interval)

COMPLETING THE MODEL

- **Multiplicities on the contigs** $m : \mathcal{C} \rightarrow \mathbb{N}$. Problem : find an optimal vertex cover (maximal weight) such that each contig edge $e \in \mathcal{C}$ is used exactly $m(e)$ times (or interval)
- **Multicriteria** Average quality score on the stories \Rightarrow two criteria to optimize.

COMPLETING THE MODEL

- **Multiplicities on the contigs** $m : \mathcal{C} \rightarrow \mathbb{N}$. Problem : find an optimal vertex cover (maximal weight) such that each contig edge $e \in \mathcal{C}$ is used exactly $m(e)$ times (or interval)
- **Multicriteria** Average quality score on the stories \Rightarrow two criteria to optimize.
- **Integrating insert sizes** especially in the case where multiple sources are used (e.g. PacBio + Illumina reads)

APPROXIMABILITY

- Lower bounds ?



APPROXIMABILITY

- Lower bounds ?

- Other polynomial-time algorithms with a better ratio ?



APPROXIMABILITY

- Lower bounds ?
- Other polynomial-time algorithms with a better ratio ?
- Discuss the ratio according to n vs. $2\sigma_c + \sigma_p$



APPROXIMABILITY

- Lower bounds ?
- Other polynomial-time algorithms with a better ratio ?
- Discuss the ratio according to n vs. $2\sigma_c + \sigma_p$
- Randomized algorithms ?



BIOLOGICAL VALIDATION

- Need to compare to other methods (Implementation and experiments in progress...)

BIOLOGICAL VALIDATION

- Need to compare to other methods (Implementation and experiments in progress...)
- Difficulties to find a good benchmark and good measures

Hunt *et al. Genome Biology* 2014, **15**:R42
<http://genomebiology.com/2014/15/3/R42>



RESEARCH

Open Access

A comprehensive evaluation of assembly scaffolding tools

Martin Hunt^{1*}, Chris Newbold^{2,1}, Matthew Berriman¹ and Thomas D Otto¹



Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

THAT'S ALL FOLKS

- Questions ?



Institut de biologie
computationnelle

Laboratoire
Informatique
Robotique
Microélectronique
Montpellier

THAT'S ALL FOLKS

- Questions ?
- Answers ?