



Don't throw away *chimeras*!

(comics edition)

blogs.browardpalmbeach.com

Petar Ivanov

Mentors:

Sergey Nurk

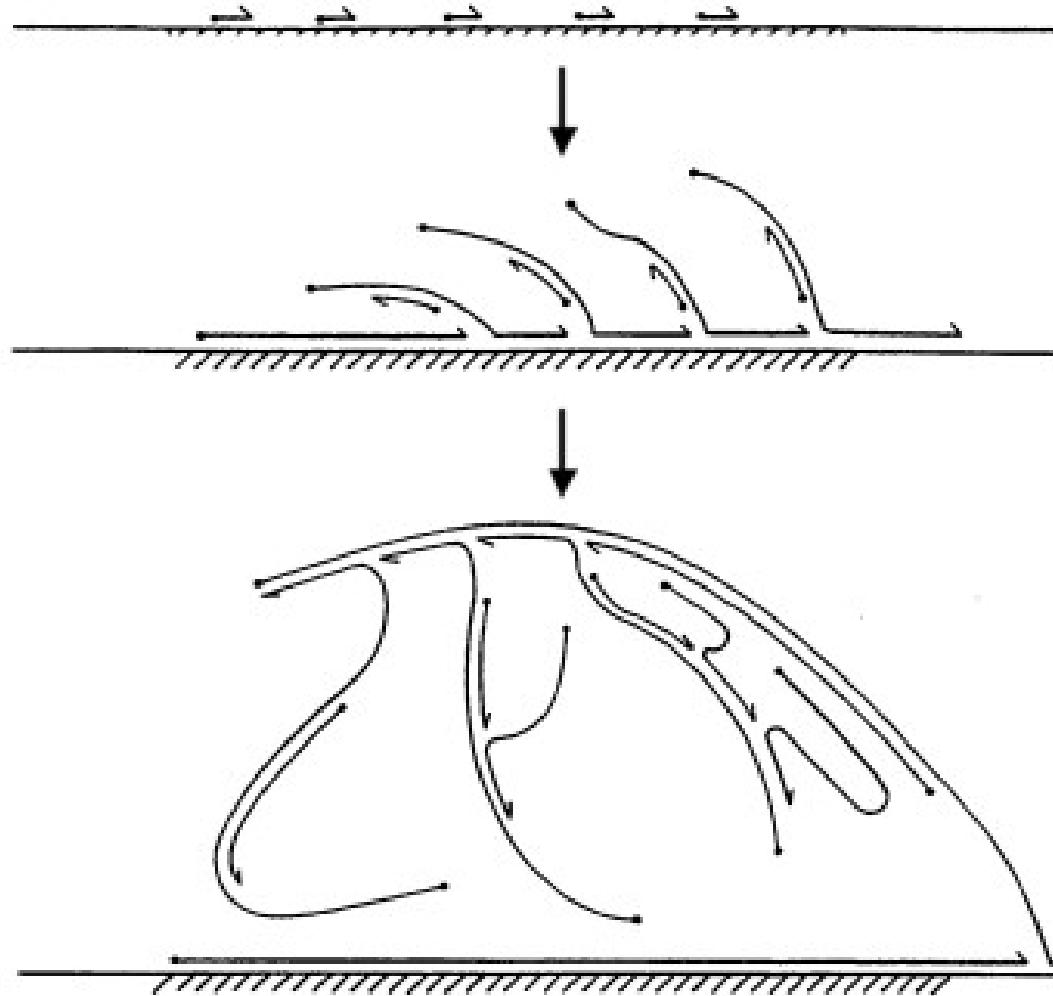
Anton Bankevich

Jul–Aug 2013

Algorithmic Biology Lab

St. Petersburg Academic University

Background: MDA



Background: MDA



- *Multiple Displacement Amplification* (MDA) is the currently used method for single cell DNA amplification
- MDA generates chimeric (i.e. not really existing in the genome) DNA rearrangements in the amplified DNA
- All the genome assemblers try to reduce and eliminate the chimeric reads

...but we can do better

Our dream



Make use of the proximity (span) expectation between both ends of the chimera:

- in path extending **(to be presented)**
- in scaffolding **(in the future)**

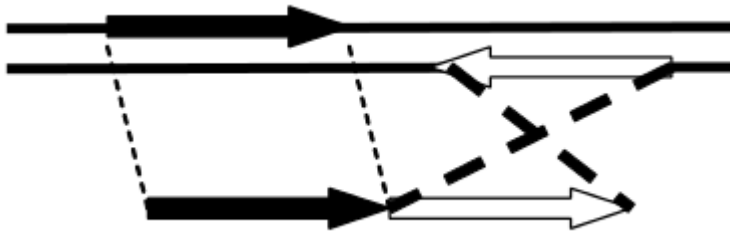


- (done before) MDA lane + SPAdes → simplified de Bruijn Graph → identify chimeras
- (current work) chimeras + reference → statistics
- (current work) de Bruijn Graph + statistics → chimeric path-chooser
- (future plan) integrate the chimera-flavored path-chooser in production SPAdes
- (future plan) chimera-based scaffolding

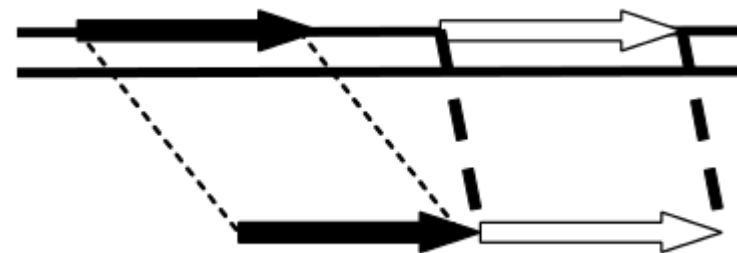
Chimeras: inverted vs direct



Reads 85% vs 15% (Lasken on *E. coli*)
Chimeras 71% vs 29% (Ours on *E. coli*)



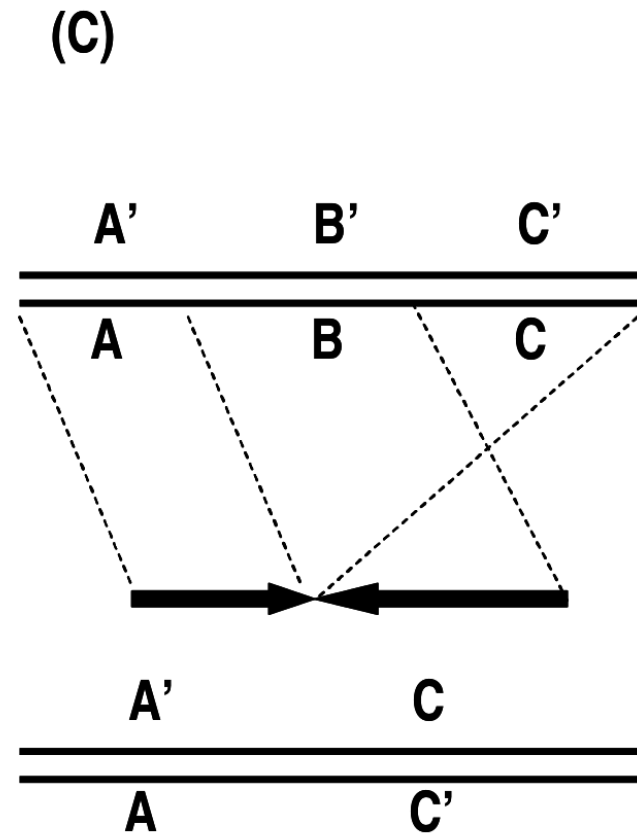
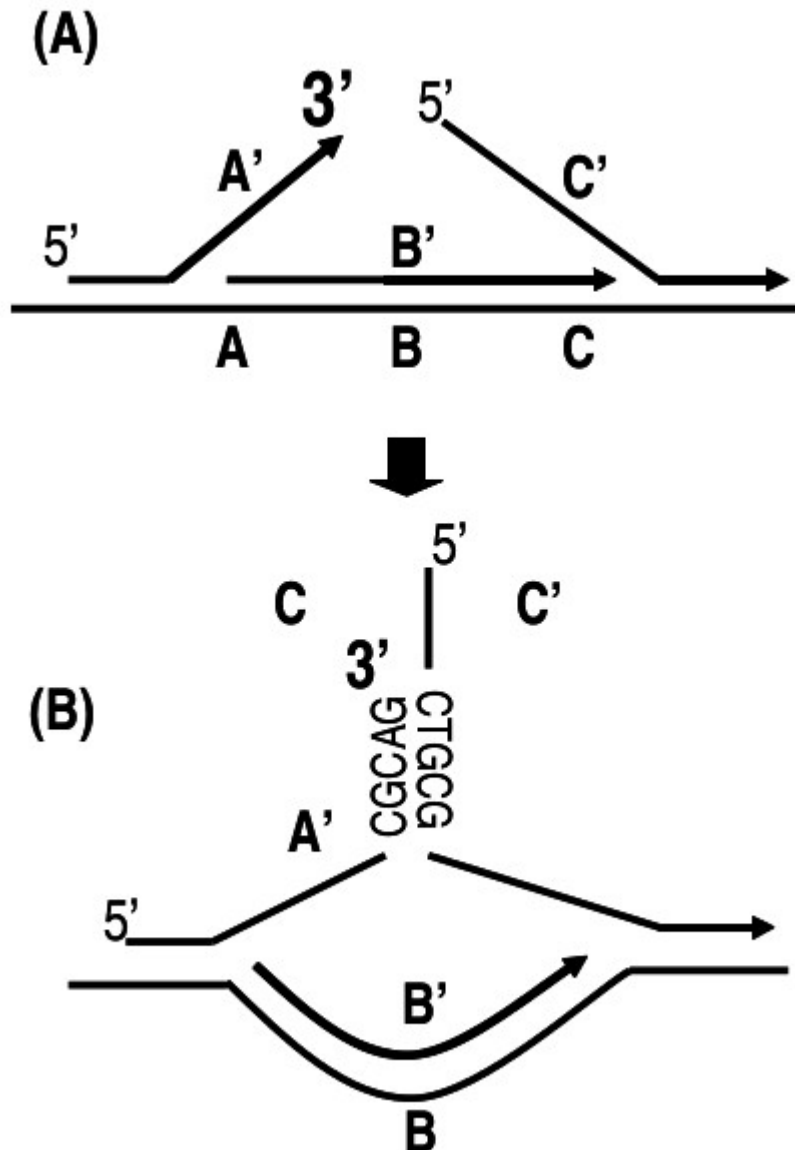
inverted



direct

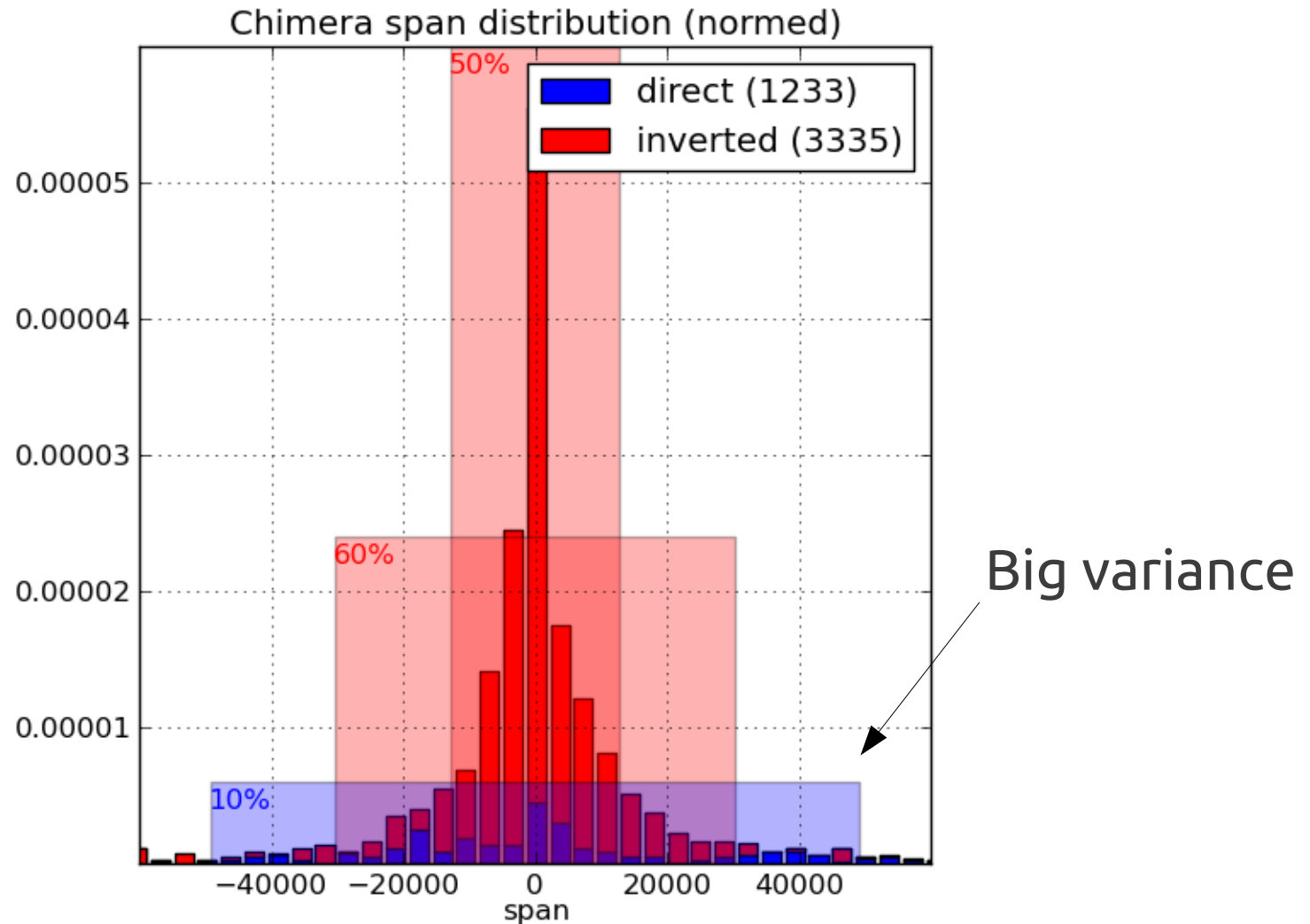
Lasken, Stockwell – “Mechanism of chimera formation during the Multiple Displacement Amplification reaction”, 2007

Inverted chimera formation



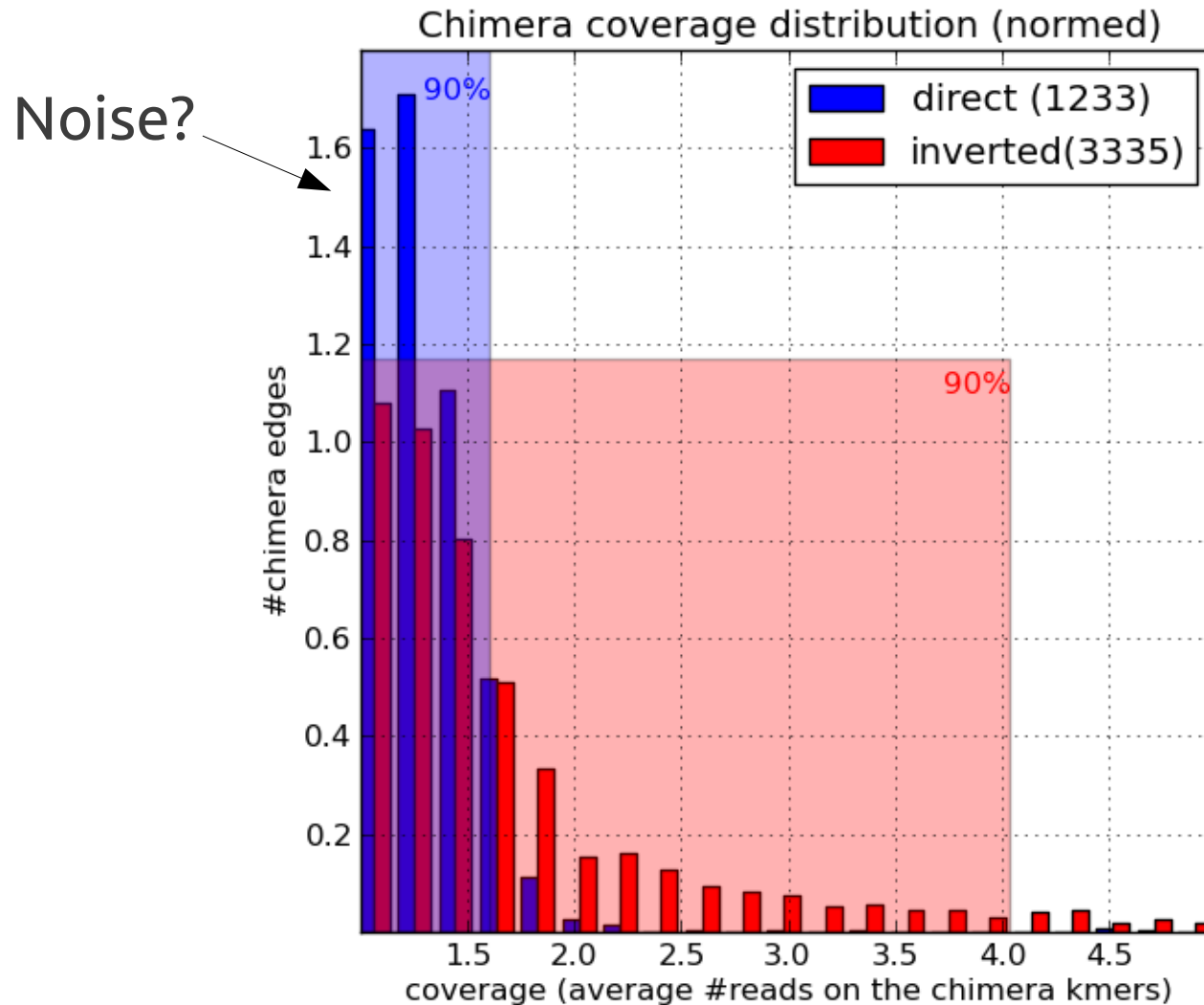
Lasken, Stockwell – “Mechanism of chimera formation during the Multiple Displacement Amplification reaction”, 2007

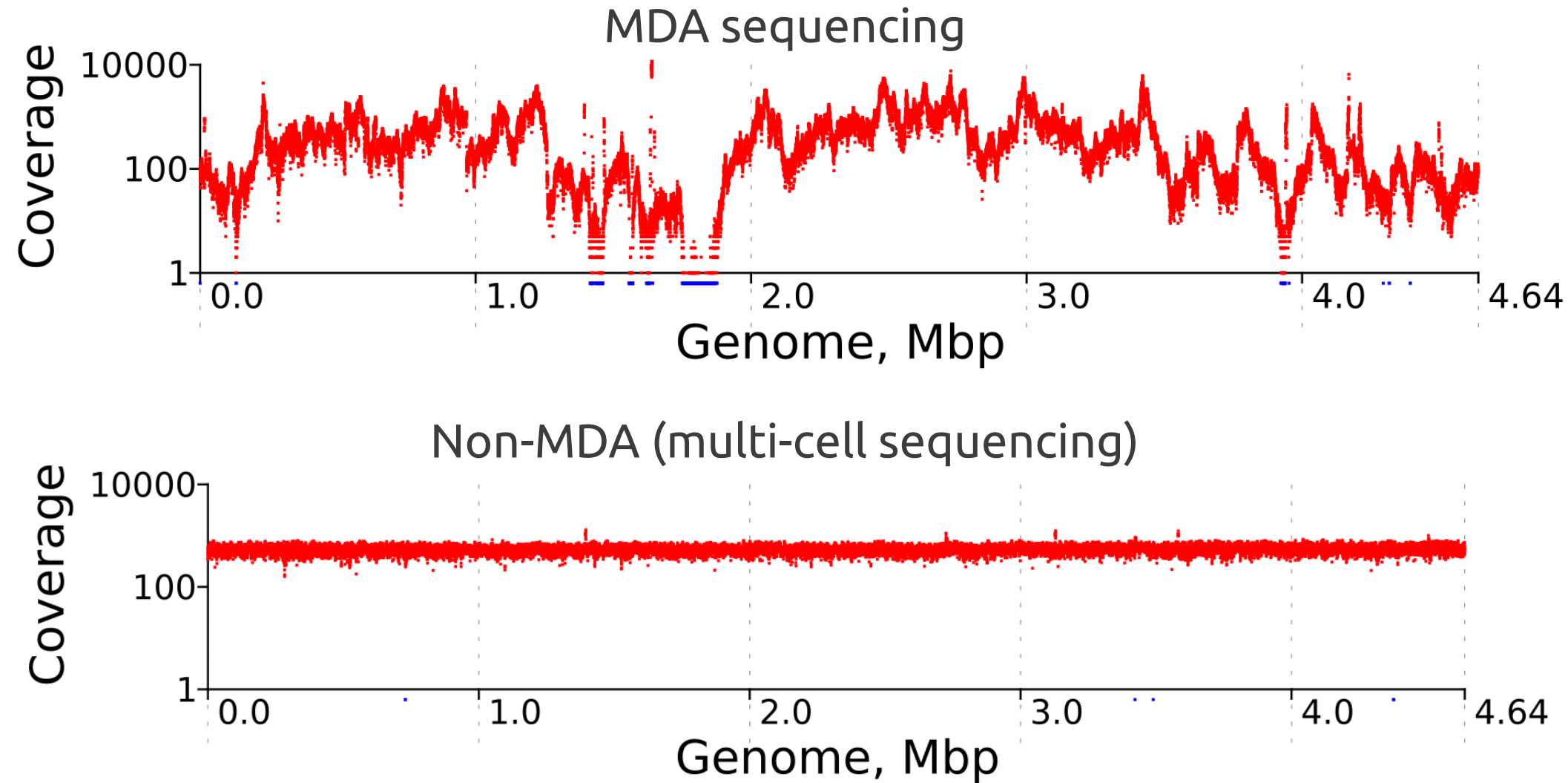
Symmetric & unimodal Nice!



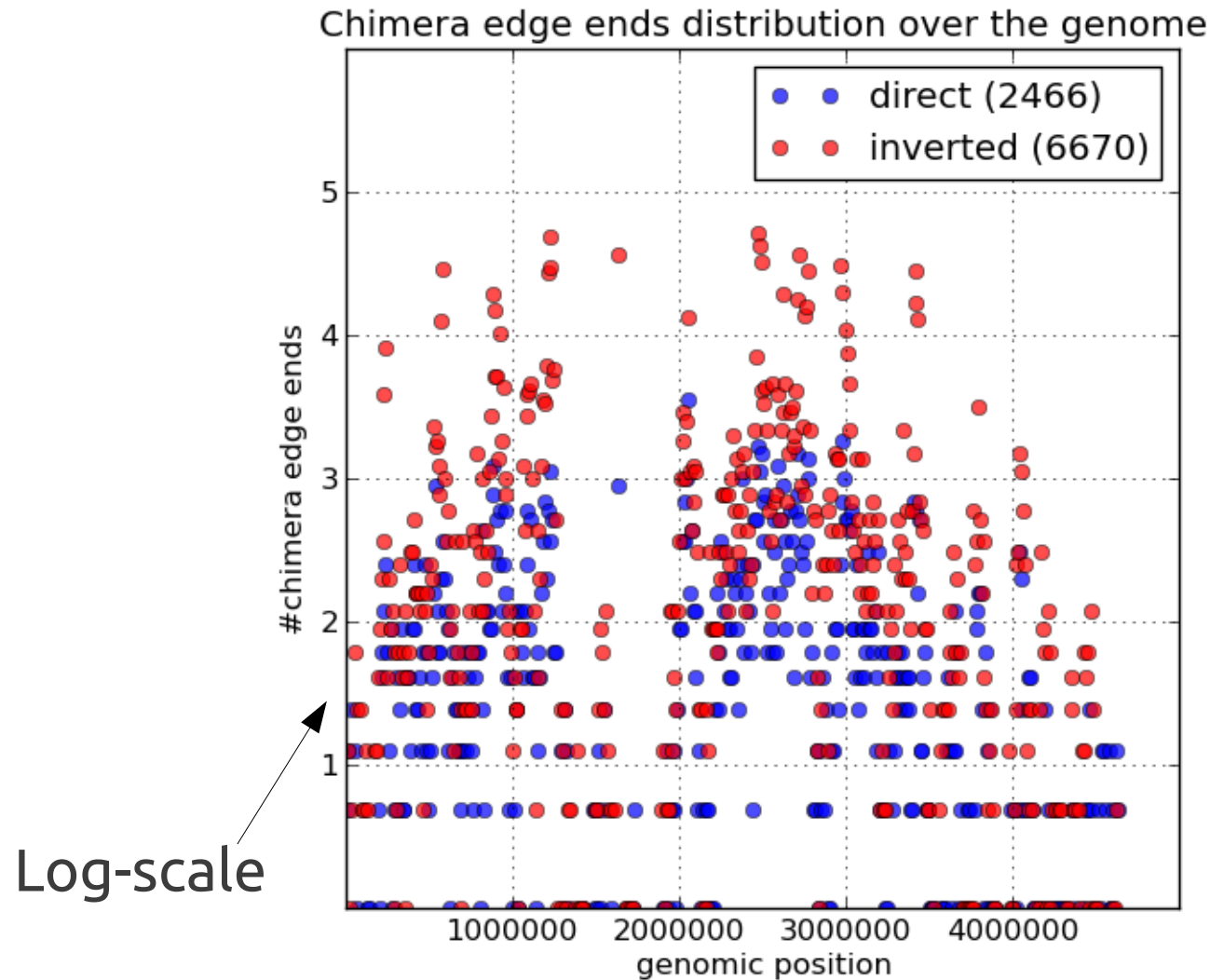


Inverted chimeras have greater coverage (?!)





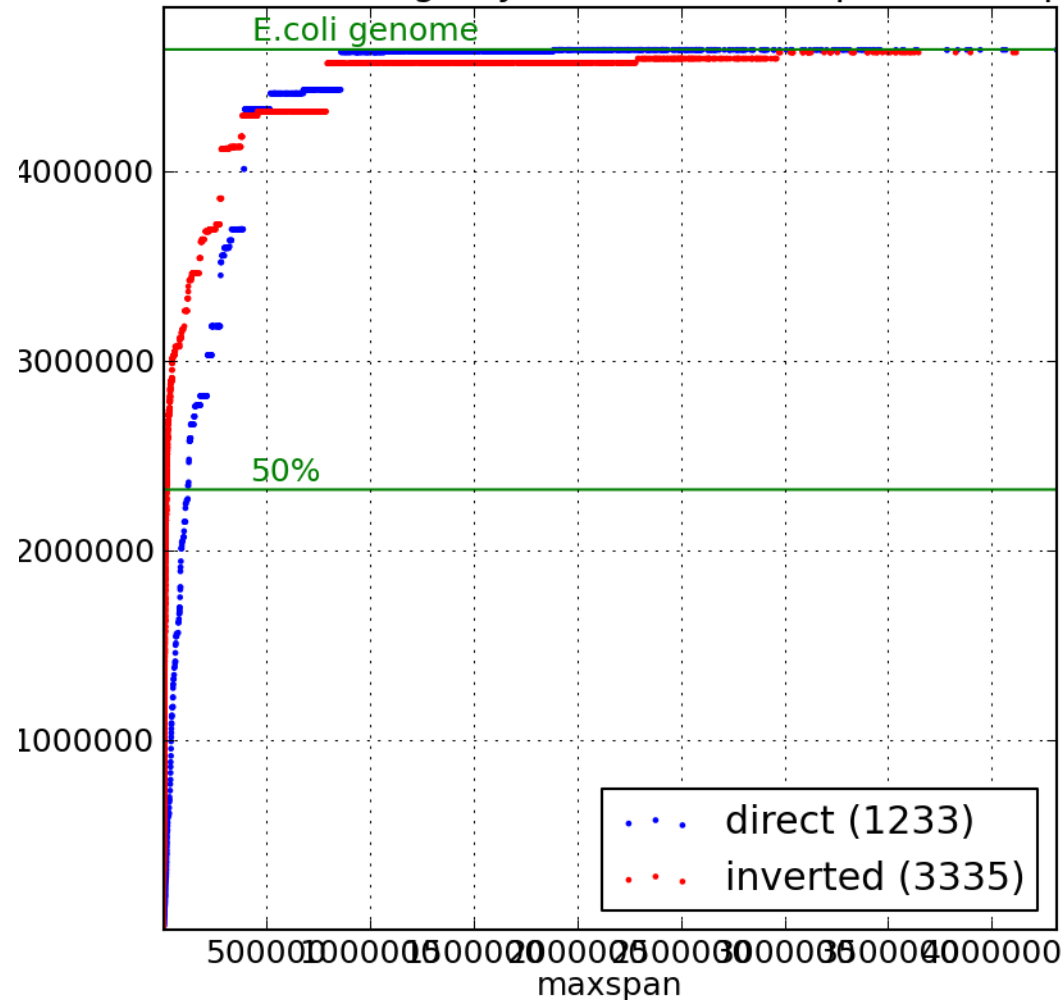
Chimeras on the genome are as fuzzy as reads



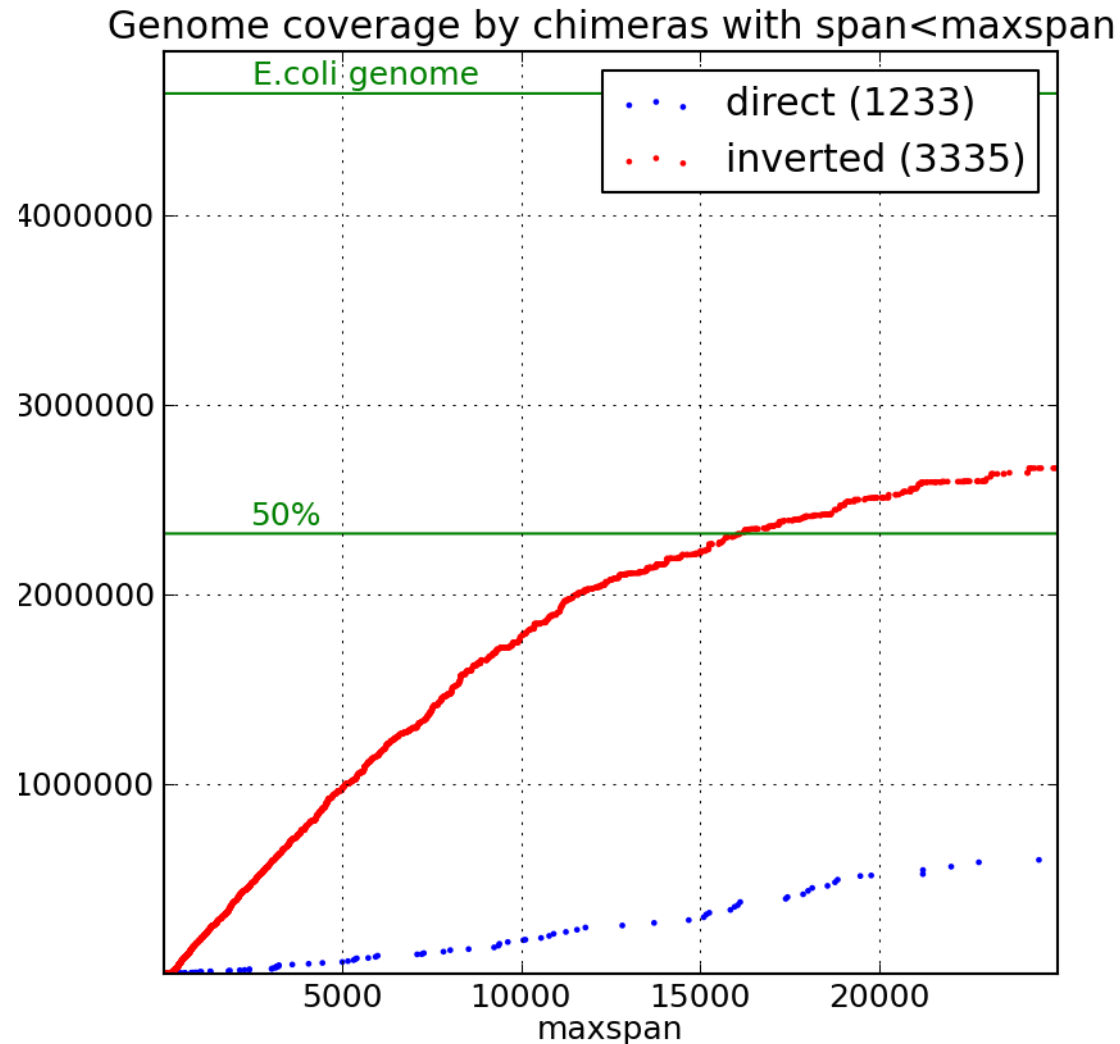
Upper bound of informativeness



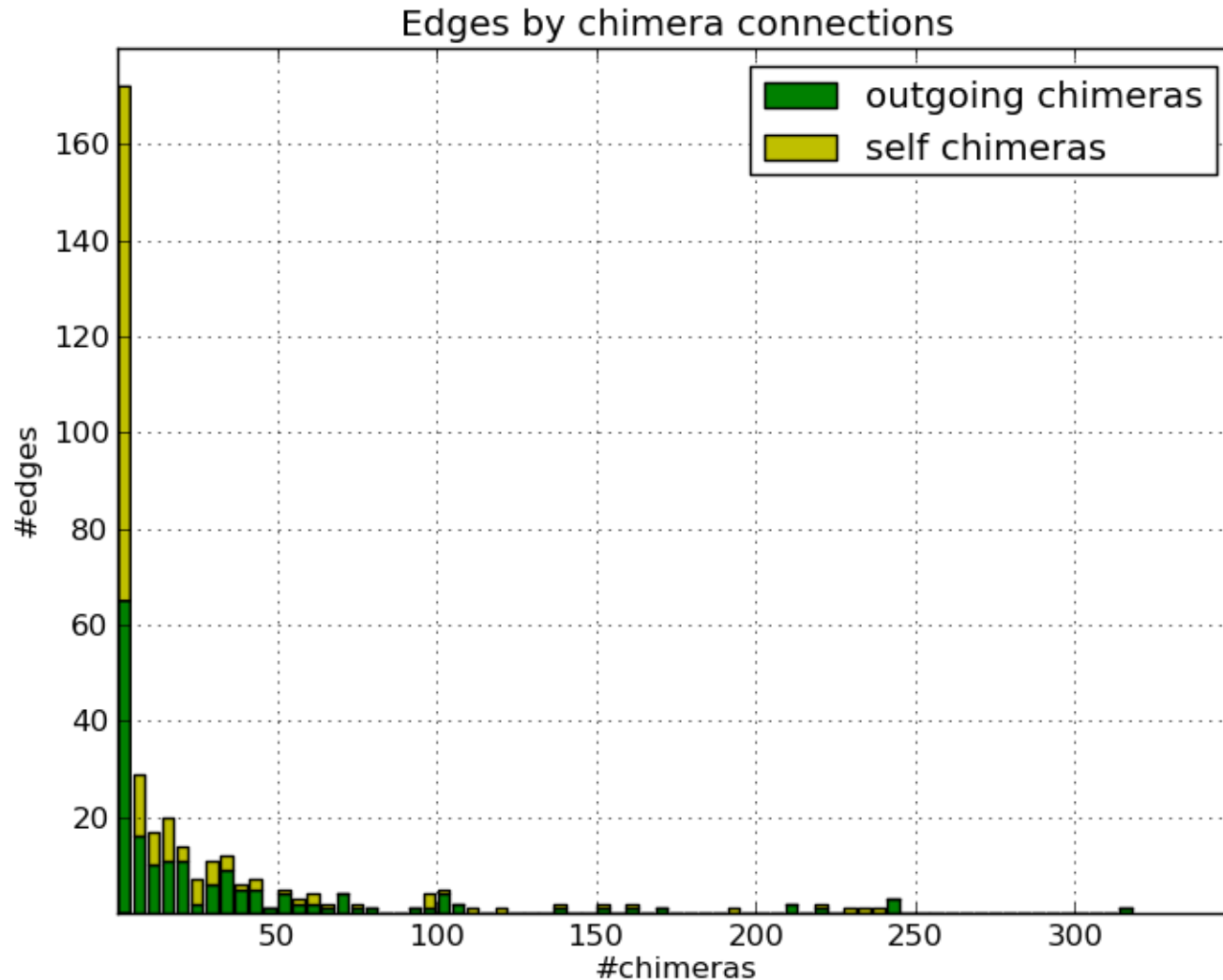
Genome coverage by chimeras with span < maxspan



<15k span chimeras cover 50% of *E. coli*



How many edges can we join together



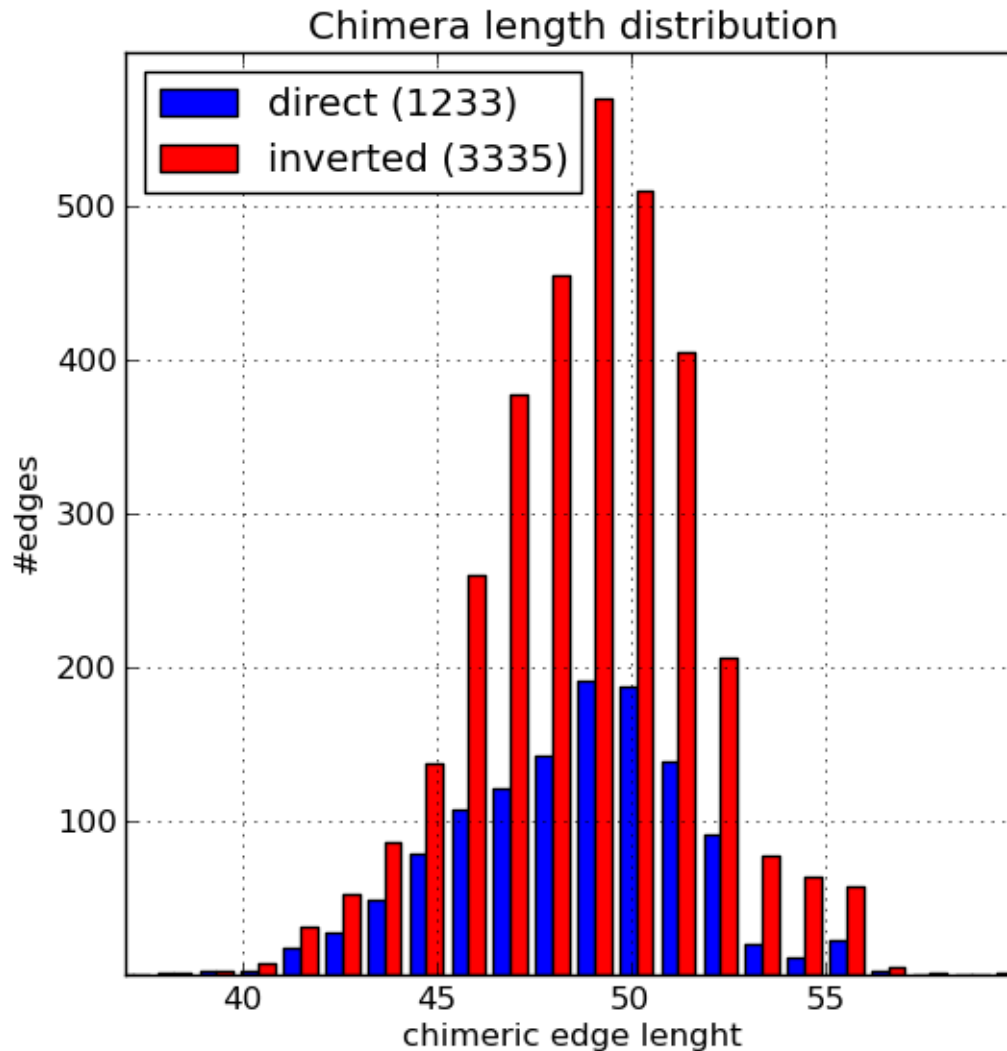
Chimera length



5' - CCAGTGAATTTCACTTCG**CCAACG** - 3'
 3' - **GGTTG**CCCGACCGGGTGTTCAAC - 5'
 5' - TGTACTAAAAGGGTAGTC**AGAAAAA** - 3'
 3' - **TCTTTTTT**GGTCCAGAGCTAAAAT - 5'
 5' - GAGGCAACATTTGATC**GTCAGTG** - 3'
 3' - **CAGTC**ATACTTTTCAGGCACCGTCG - 5'
 5' - CGCCAGGAAACATTGCAC**ACCACGC** - 3'
 3' - **GTGGGGCG**CTAGCGCTCCGTTTGG - 5'
 5' - CATTCCCGGAATTAC**ATATCTTT** - 3'
 3' - **TATAGAAAA**AGTAATCCGTCACCGGA - 5'
 5' - GCATATCTCCATCCT**GAGTGACGC** - 3'
 3' - **CTCATTGCG**AAAACCAACCCGCTCTT - 5'
 5' - TTTGAAATATCCACTATTAAGCTAG**TGTTTAACG** - 3'
 3' - **CACAAATTGCG**TCGGAA - 5'

Lasken, Stockwell – “Mechanism of chimera formation during the Multiple Displacement Amplification reaction”, 2007

Length can also vote for chimeras

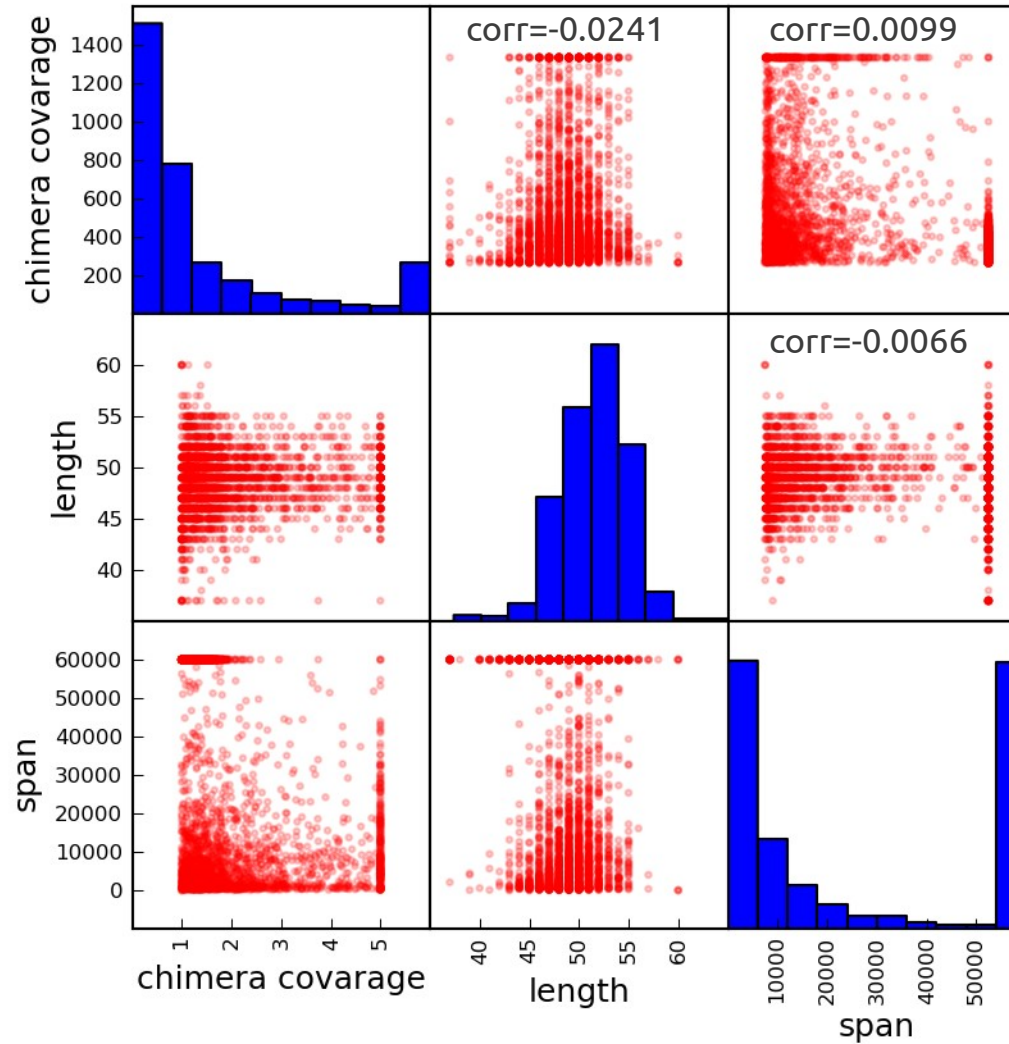


Mean of 49 for both
chimera types
(agreement with
Lasken & Stockwell)

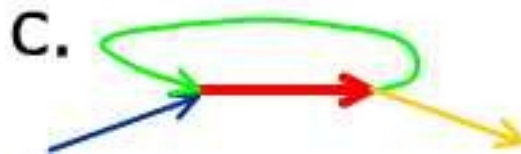
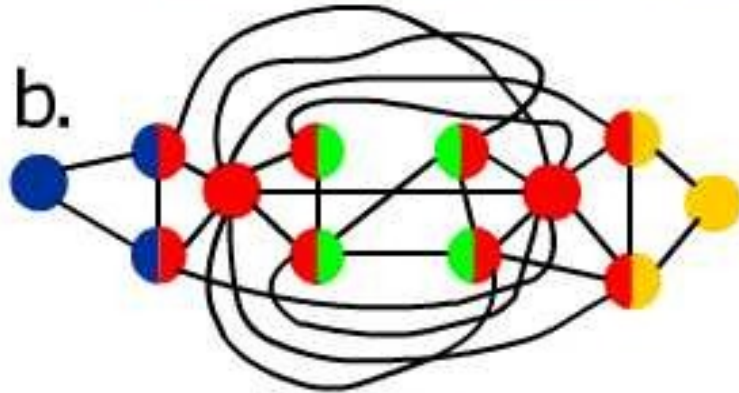
No correlation... :(



Relative distributions for inverted chimeras



de Bruijn Graph



Our hero...
the De Bruijn Graph



How chimera looks like

Path extender queries choosers



Given: **path** + set of **outgoing edges**

Return: the correct **extension edge** (if any)
or an empty set if not sure

Idea: lets make a path chooser on chimeras!

SPAdes has several different choosers
(by paired-end reads, mate-pair reads, long reads, etc.)

Naïve path chooser



Considers only chimeras with
 $span \leq maxspan$ (~15'000)

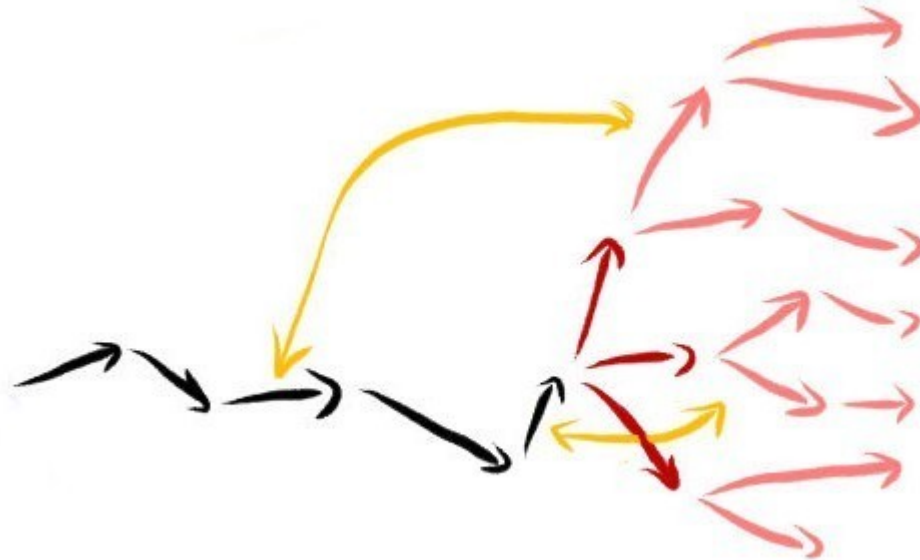


Choose the edge connected to the path
by a maximum number of chimeras

chooser invocations: 416
corrects: 3
incorrects: 0

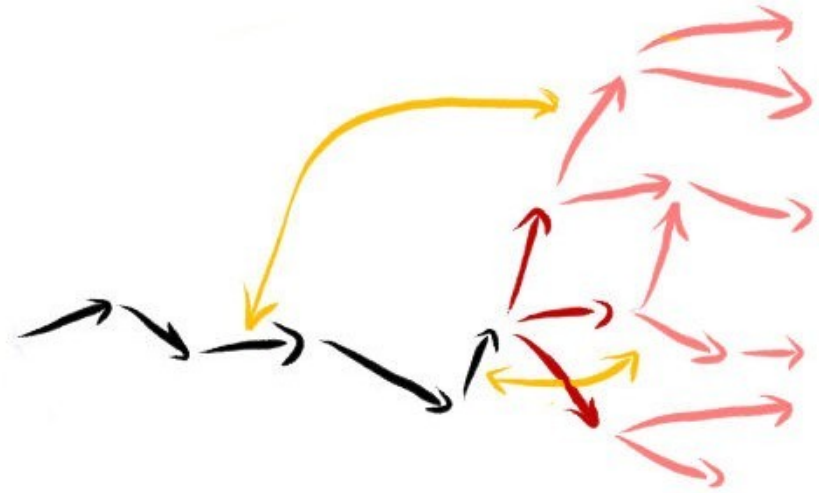
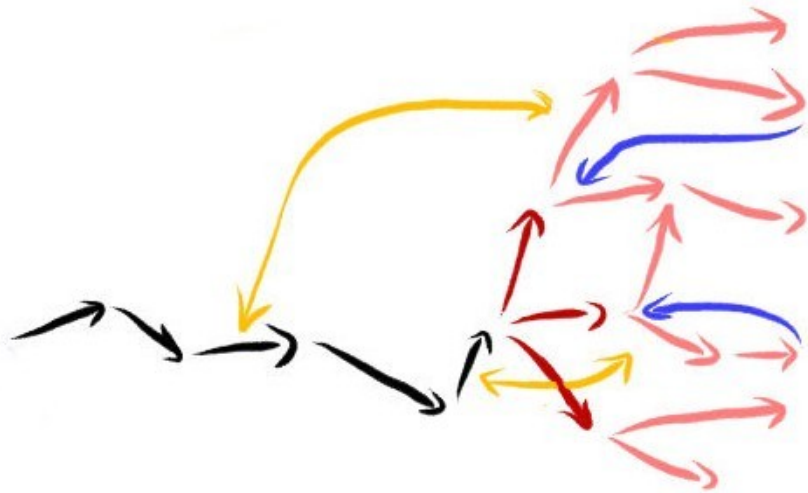
Concl: nice but the #chimeras is not enough

Not edges but paths



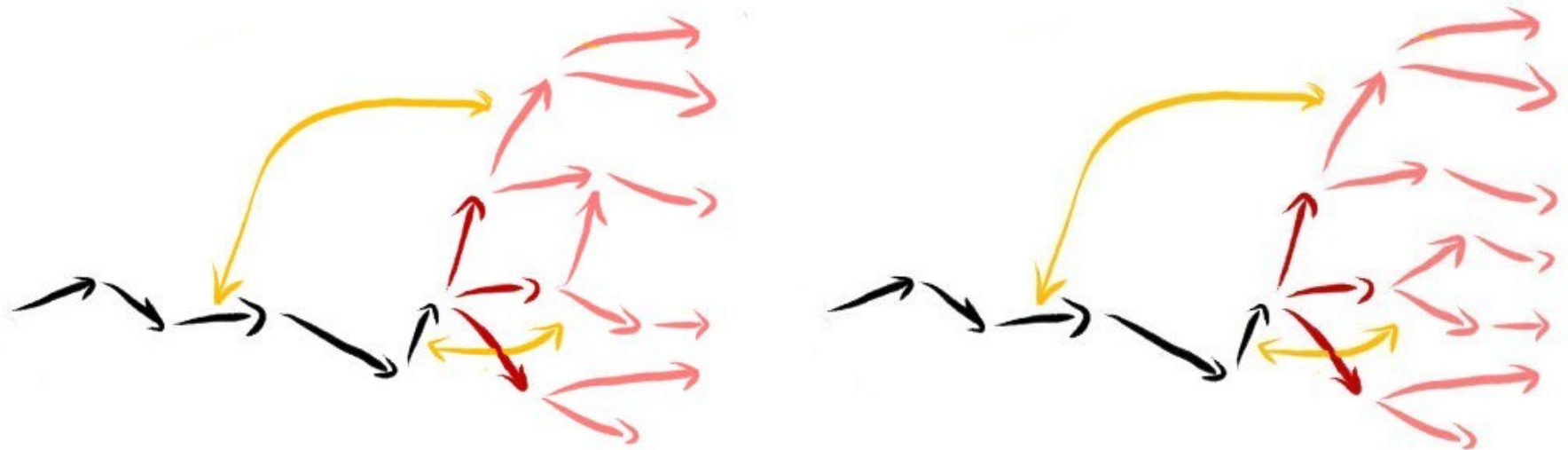
Consider not only the chimeras to edges but to whole path extensions — choose the one connected to the path by a maximum number of chimeras

Cycles?



Kill them somehow
(e.g. by constructing a **DFS** tree)

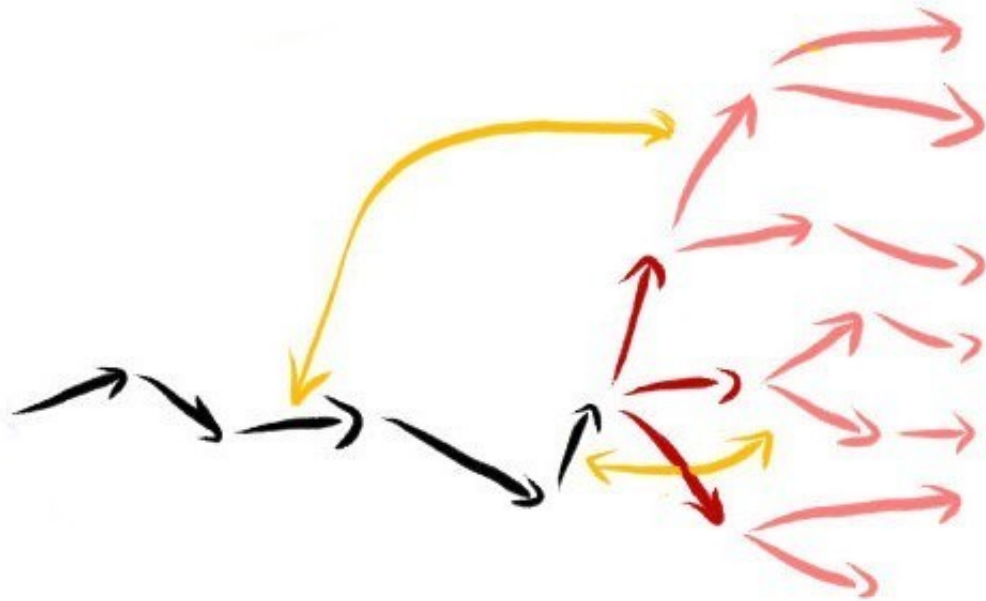
Not a tree?



In case of a **Directed Acyclic Graph**
duplicate the joint **path extensions**



A tree again



So the notion of “path” gets clear



Let $\{\text{result}^i\}$ be a subset of edges, s.t.

- starting at resultⁱ there is an extension path connected with at least *min#chimeras* ($\sim 2 \div 3$)

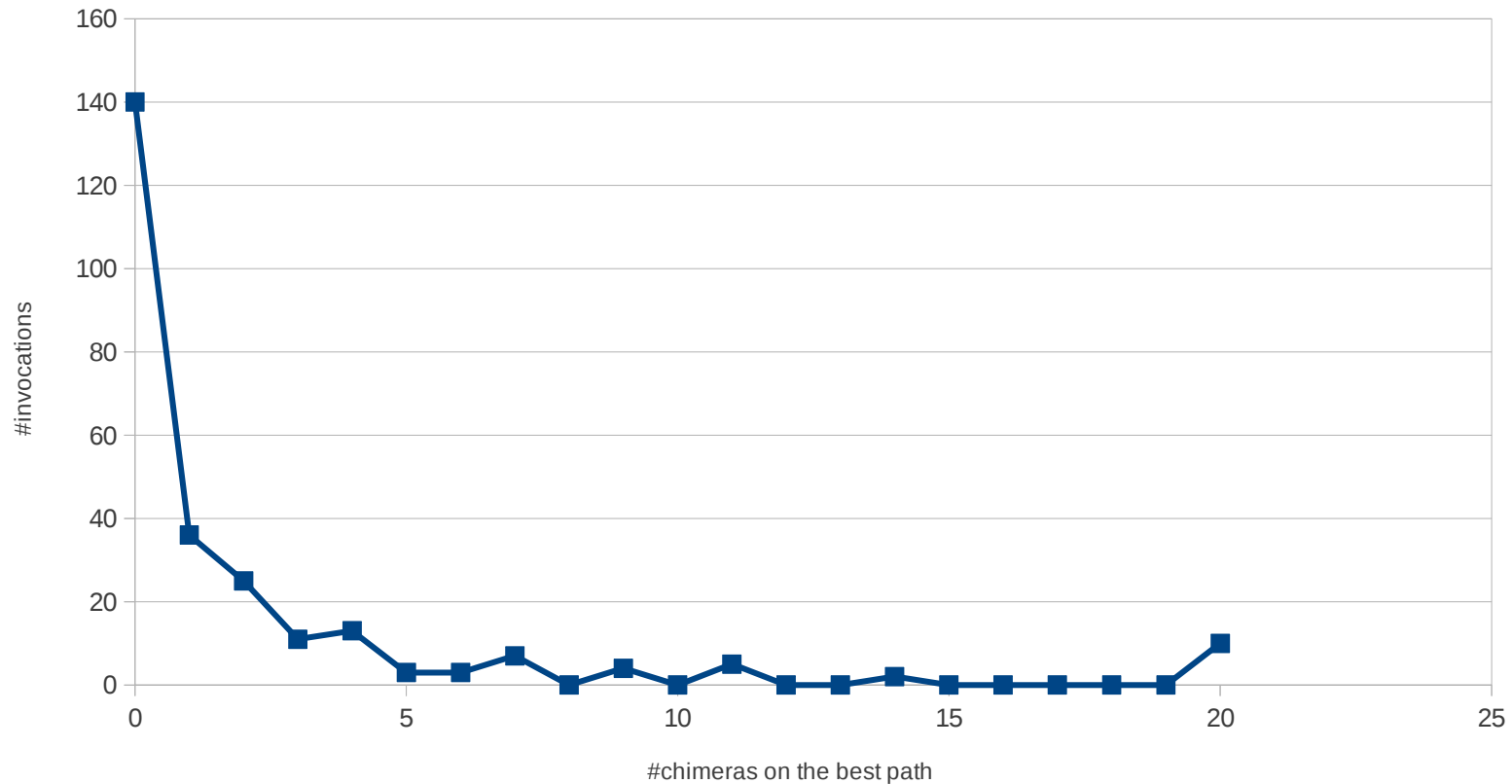
Return nothing in case of:

- multiple results; or
- another extension edge leading to a path with $\geq \text{leader_coef}$ ($\sim 1 \div 1.5$) times the #chimeras in the best path

Else:

- Return $\{\text{result}^i\}$ which has a single or no edges

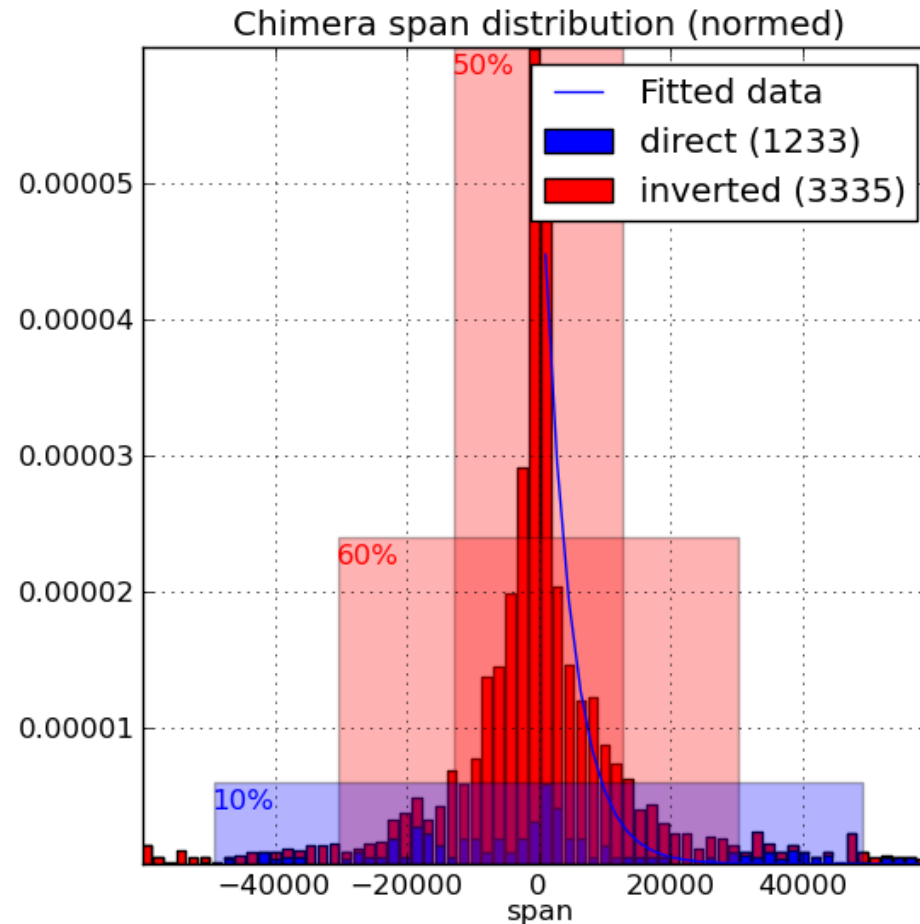
#invocations by #chimeras in the best path



We need more chimeras, folks...



Artificial chimeras



Blue curve:
 $\lambda \cdot \exp(-\lambda \cdot x)$
 $\lambda = 1/16'145$

The span distribution
of the inverted chimeras is almost exponential

E.coli results



Artificial inverted chimeras	Edges with chimeras	Corrects	Incorrects	#tmp
50000	676	51	8	213
20000	626	44	2	208
10000	570	35	3	210
5000	496	27	1	211
0	332	19	1	212

Edges: 1935
 Natural inverted chimeras: 3336
 Natural direct chimeras: 1331
 MaxSpan: 15000
 min#chimeras: 3
Invocations: ~430÷460

Not yet ready for production :")

Assembly	Scaffolds without any choosers	Current best results	+ Chimera chooser
N50	67332	109825	121369
#misassemblies	3	2	3

Some useless parameters we tried



- Give up if no long edges in the extension path
- Don't allow extension path to continue to the path: only few such situations
- Limit the number of forks on the extension path: no visible correlation with correctness
- Limit the number of edges the extension path



Pain in the ass: Correctness



It is tricky to certify the chooser correctness:

Sometimes the given **path** doesn't match
anywhere in the genome

Different heuristics are tried with no right
way to do it



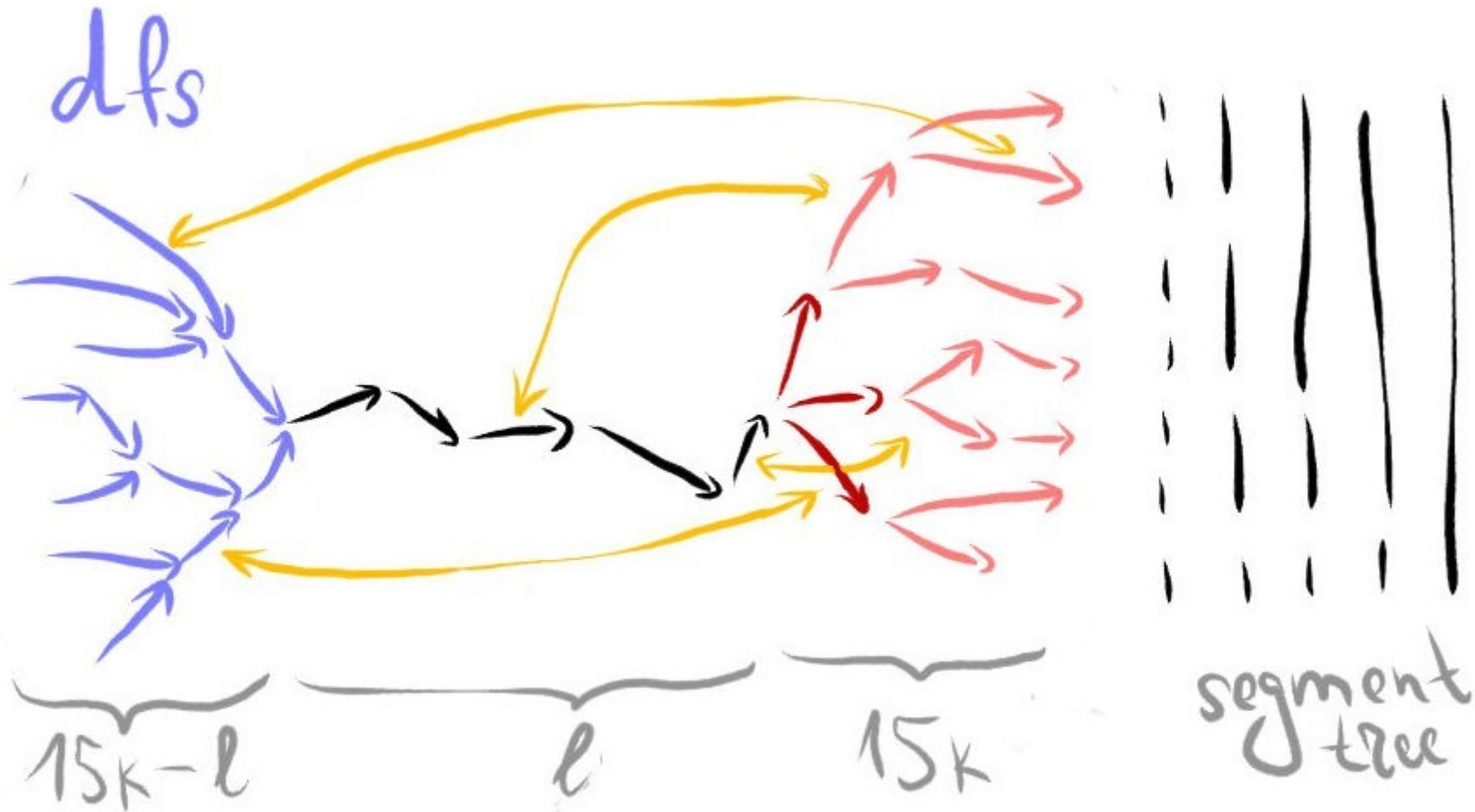
Fundamental problems



Too few chimeras

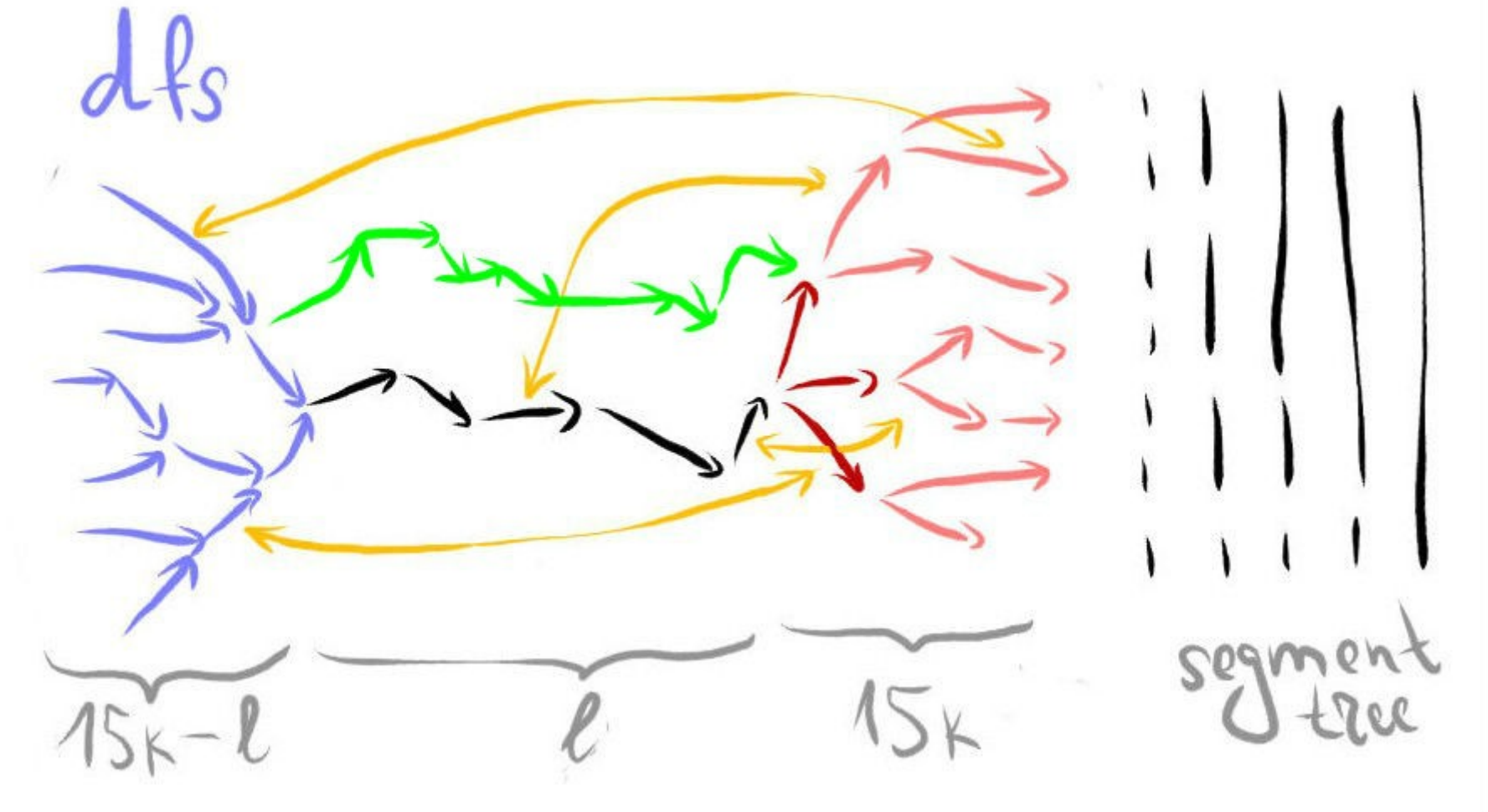
- Can we ask the biologists to produce more chimera?
- Let's try to use even more chimeras...

Joint optimization of the incoming and outgoing paths



Recurse the incoming tree while updating the #chimeras going to all the extension paths

Joint optimization of the incoming and outgoing paths



Recurse the incoming tree while updating the #chimeras going to all the extension paths

Joint optimization of the incoming and outgoing paths



$O(?)$: finds BypassingPaths from $v \in \text{incoming tree}$ or $v' \in \text{extension paths}$
 $\theta(N + C \cdot \log C)$: map the **chimeras** going to the **extension paths**
 $\theta(C \cdot \log C + N \cdot \log C)$: recurse the **incoming tree** maintaining a segment tree for the #chimeras in all **extension paths**

N – #vertices, C – #chimeras

```

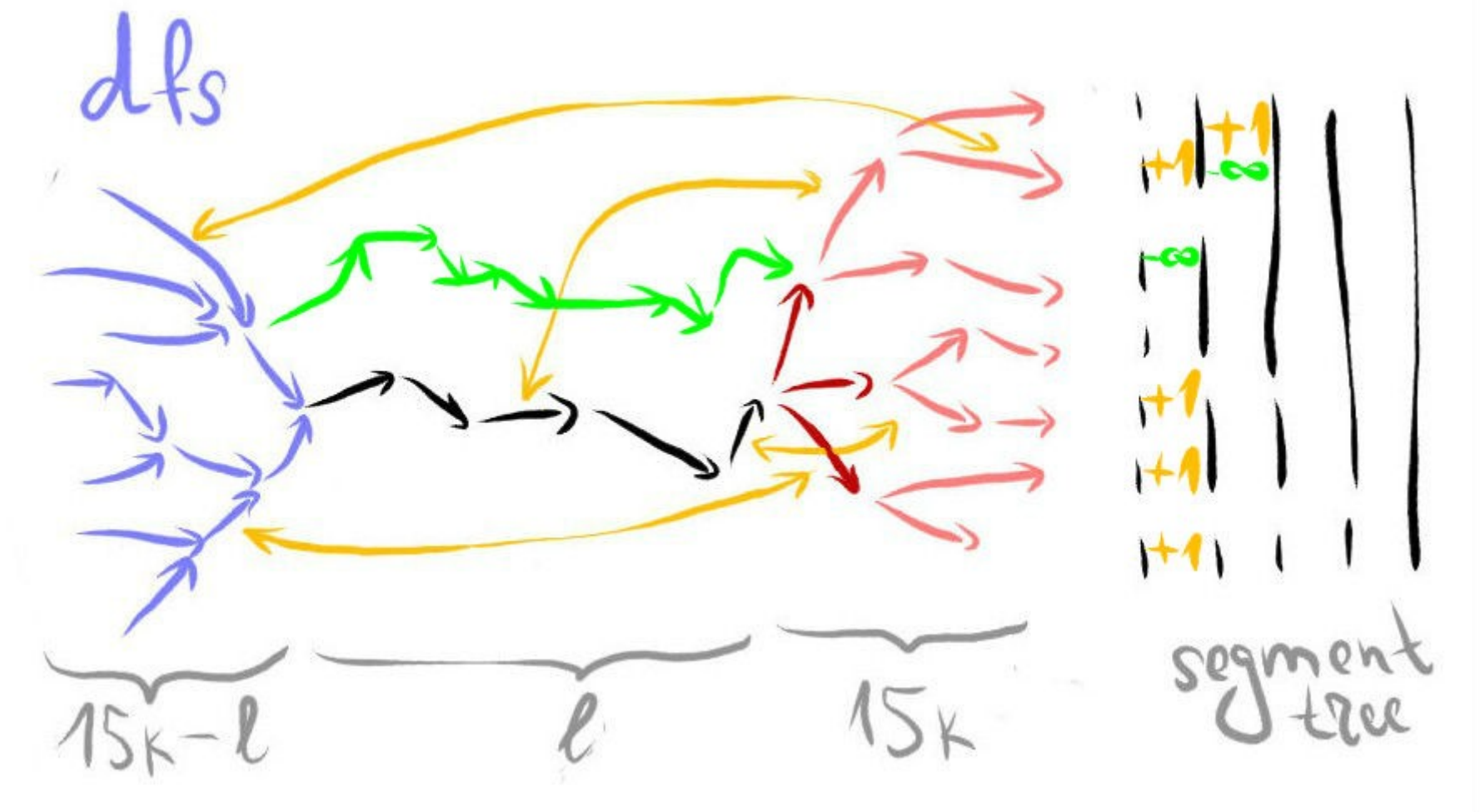
dfs(v, w):                                     // v, w ∈ incoming tree:
    UpdateBypassings(w, -inf)
    UpdateChimeras(v, w, +1)
    for edge(u, v) ∈ Edges, u ∈ incoming tree:
        dfs(u, v)
    UpdateChimeras(v, w, -1)
    UpdateBypassings(w, +inf)

UpdateBypassings(v, val):                     // v ∈ incoming tree:
    for edge(v, v') ∈ BypassingPath, v' ∈ extension:
        updateSegTree(v', val)

UpdateChimeras(v, w, val):                   // v, w ∈ incoming tree:
    for chimera from edge(v, w) to edge(v', w') ∈ extension:
        updateSegTree(w', val)

updateSegTree(v', val):                      // v' ∈ extension:
    chimeras := updateVal(v', val)
    updatePathRes(v', chimeras)
    
```

Joint optimization of the incoming and outgoing paths



Recurse the incoming tree while updating the #chimeras going to all the extension paths



- Get better graph visualizations for debugging
- Test on other datasets (S.aureus)
- Ask the technicians to increase #chimeras
- Chimeras → Scaffold
- Inverted vs Direct chimera classify by graph topology
- Define a probabilistic interpretation

Thank you for
discussing :")



*Let we use the MDA “bugs”(chimeras) to
assemble single cells better than multi-cells!*

