

Annotation-free Identification of Synteny Anchors

Karl Kaether

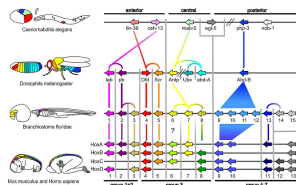
Abteilung Stadler, Leipzig

Bled 2024

Orthology Inference

- mostly based on sequence similarity
- e.g. (reciprocal) blast
- some problems^{1,2,3}
 - ▶ sequence divergence
 - ▶ genome rearrangements

- **synteny** can help in such situations



<https://en.wikipedia.org/wiki/Synteny>

- in principle solved by global alignments
- realistic approaches use genome annotations

¹Altenhoff et al., "The Quest for Orthologs benchmark service and consensus calls in 2020".

²Moyers and Zhang, "Further Simulations and Analyses Demonstrate Open Problems of Phylostratigraphy".

³Vakirlis, Carvunis, and McLysaght, "Synteny-based analyses indicate that sequence divergence is not the main source of orphan genes".

A



B

C

?



A

B

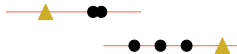
C

<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	98.63	73	1	0	1	73	47815972	47816044	3e-29	130
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	98.63	73	1	0	1	73	47818702	47818630	3e-29	130
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	94.52	73	4	0	1	73	233854997	233854925	3e-24	113
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	94.52	73	4	0	1	73	233859867	233859139	3e-24	113
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	91.78	73	6	0	1	73	233863496	233863568	7e-21	102
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	85.14	74	9	2	1	73	303312418	303312490	1e-12	75.0
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	93.75	48	3	0	1	48	303310365	303310412	5e-12	73.1
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	91.67	48	4	0	1	48	303308851	303308898	2e-10	67.6
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371225.1	92.31	26	2	0	1	26	303311081	303311106	0.19	38.1
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	98.63	73	1	0	1	73	363661361	363661433	3e-29	130
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	97.26	73	2	0	1	73	319778696	319778624	1e-27	124
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	97.26	73	2	0	1	73	372231794	372231866	1e-27	124
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	94.52	73	4	0	1	73	320281395	320281467	3e-24	113
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	94.52	73	4	0	1	73	372219431	372219503	3e-24	113
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	94.52	73	4	0	1	73	372231012	372230940	3e-24	113
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	94.52	73	4	0	1	73	372274162	372274090	3e-24	113
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	93.15	73	5	0	1	73	306015974	306015902	1e-22	108
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	93.06	72	4	1	2	73	306016192	306016262	2e-21	104
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	87.67	73	8	1	1	73	61873916187462	2e-15	84.2	
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	86.30	73	9	1	1	73	618287861828807	1e-13	78.7	
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	86.30	73	9	1	1	73	178173075	178173146	1e-13	78.7
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	86.30	73	9	1	1	73	198903251	198903322	1e-13	78.7
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	86.30	73	9	1	1	73	372199550	372199479	1e-13	78.7
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	86.30	73	9	1	1	73	372231203	372231274	1e-13	78.7
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	84.93	73	10	1	1	73	281774135	281774064	5e-12	73.1
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	83.78	74	9	3	1	73	81146068114677	2e-10	67.6	
<i>Drosophila_melanogaster_tRNA-Ala-AGC-1-1</i>	0X371224.1	87.10	31	4	0	8	38	319396726	319396696	0.67	36.2

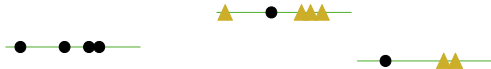
A



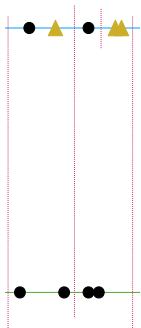
B



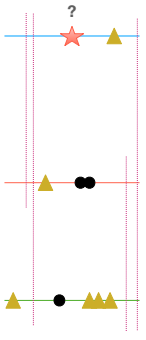
C



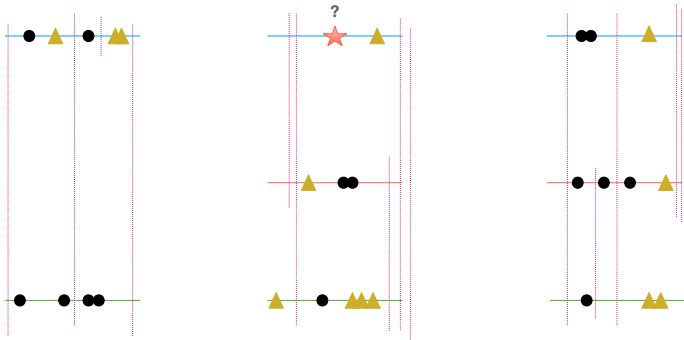
A



B



C



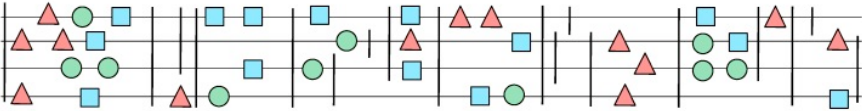
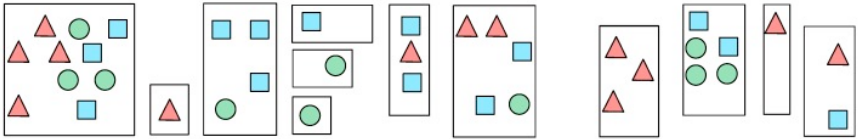
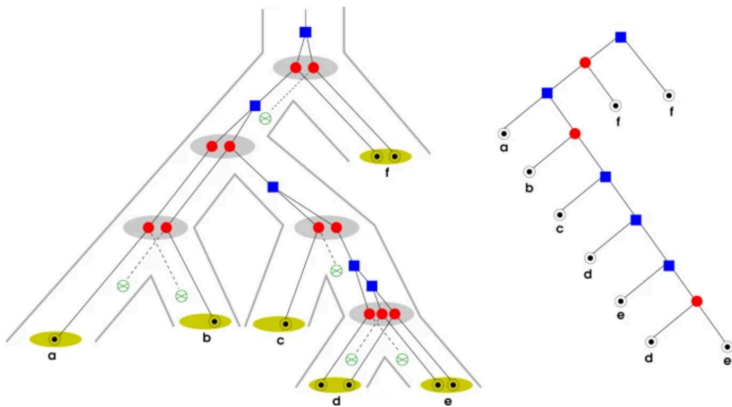
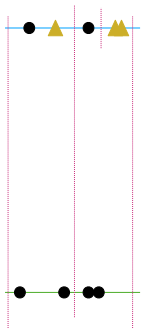
A**B**

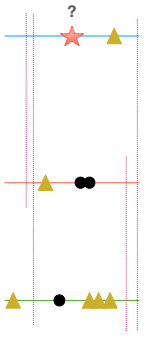
Figure 1



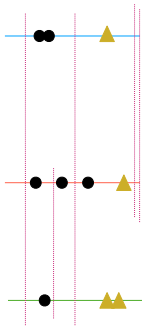
A



B



C



Finding Synteny Anchors

- usually based on genome annotations and possibly multiple sequence alignments^{4,5}
 - ▶ MSAs too expensive for many genomes
 - ▶ annotation might be unavailable
 - ▶ suffer from biases and errors⁶
 - ★ biased towards model species
 - ★ limited by assembly quality
 - ★ too sparse for synteny calculation
 - ★ contaminated by foreign species

⁴Wang et al., "MCSanX".

⁵Haas et al., "DAGchainer".

⁶Salzberg, "Next-generation genome annotation".

How to do this differently

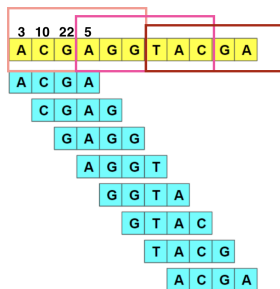
- 1 count k-mers

How to do this differently

- 1 count k-mers
- 2 chop genome into windows

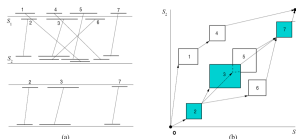
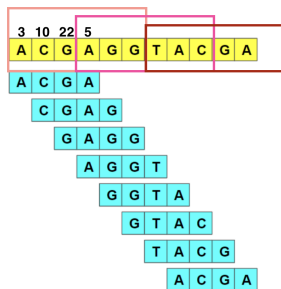
How to do this differently

- 1 count k-mers
- 2 chop genome into windows
- 3 sum up counts per window and further use x best %



How to do this differently

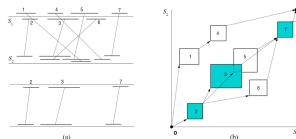
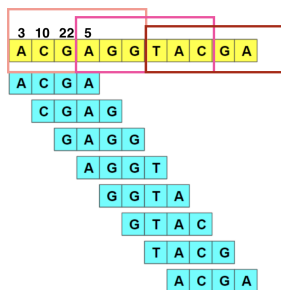
- 1 count k-mers
- 2 chop genome into windows
- 3 sum up counts per window and further use x best %



- 4 blast against own genome and chain hits $\rightarrow C(A, d^A)$ per genome

How to do this differently

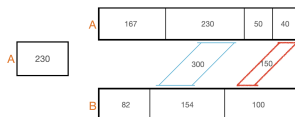
- 1 count k-mers
- 2 chop genome into windows
- 3 sum up counts per window and further use x best %

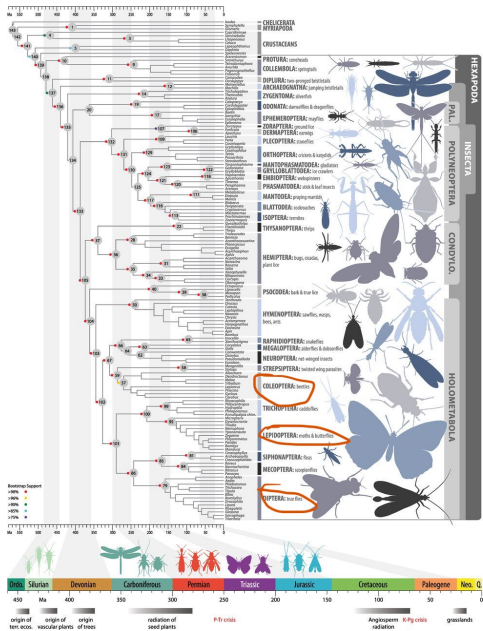


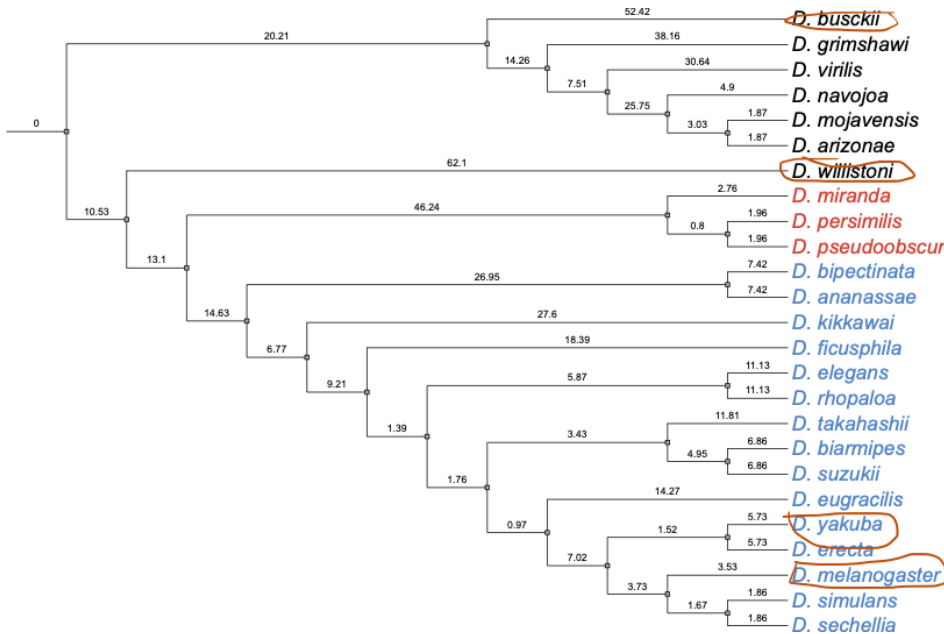
- 4 blast against own genome and chain hits $\rightarrow C(A, d^A)$ per genome
- 5 blast against $C(B, d^B)$ of other genomes and identify hits with distance d_1 satisfying $d_1 = \min(d_0^A, d_0^B) - tol \rightarrow$ syntenic anchors

Anchor Candidate Mapping Across Genomes

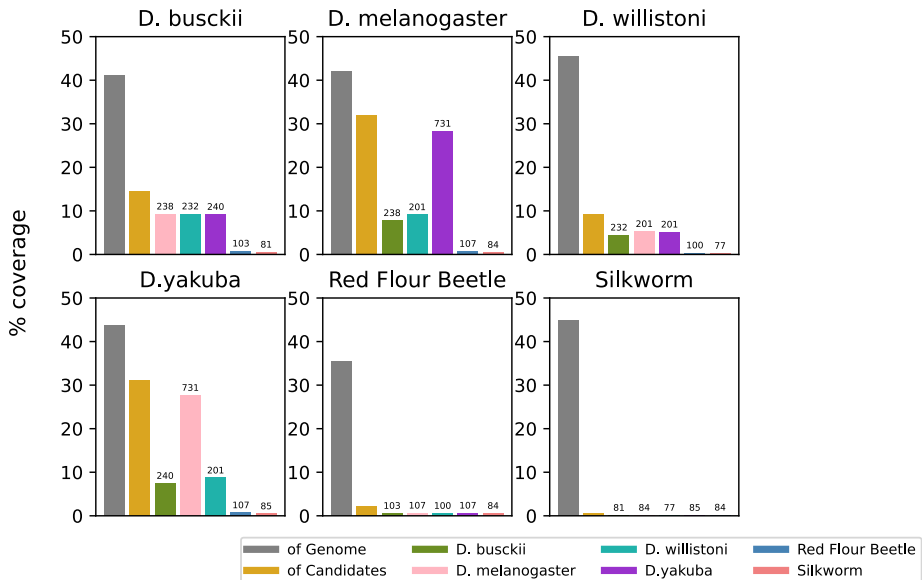
- triangle inequality implies hits between genomes with $d_1 = \min(d_0^A, d_0^B)/2$ are best hit in other genome
- in reality hits with score \geq score of region of candidate + tolerance taken



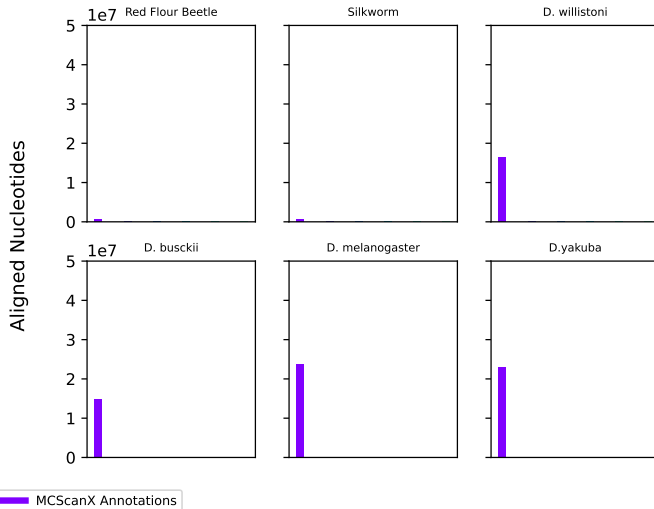




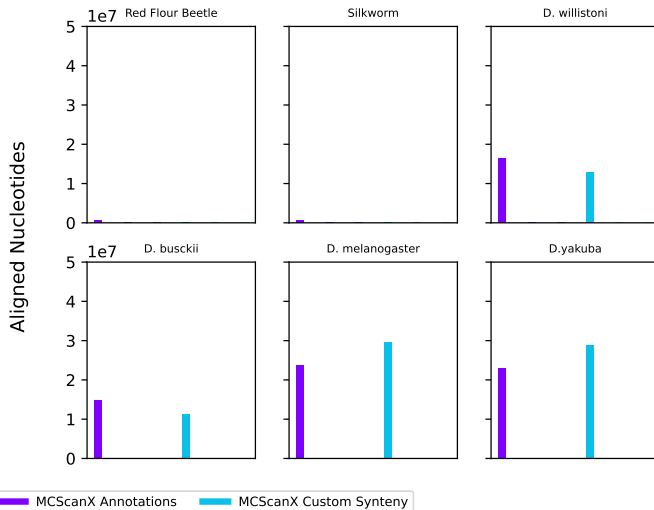
Pairwise Alignments of Anchor Candidates



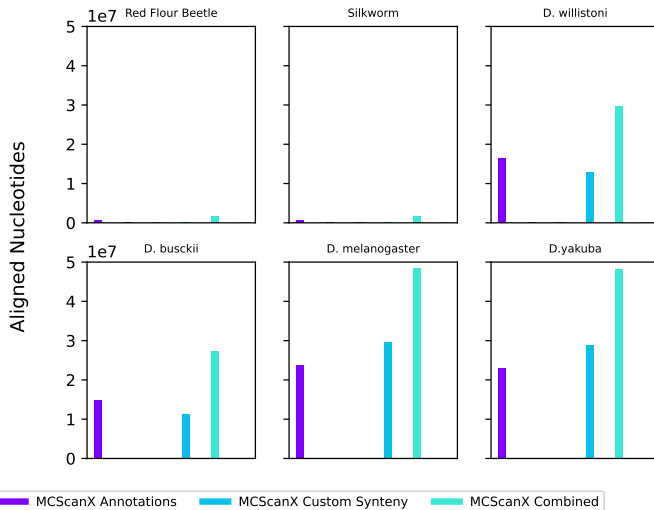
Performance of Different Synteny Calculation Strategies



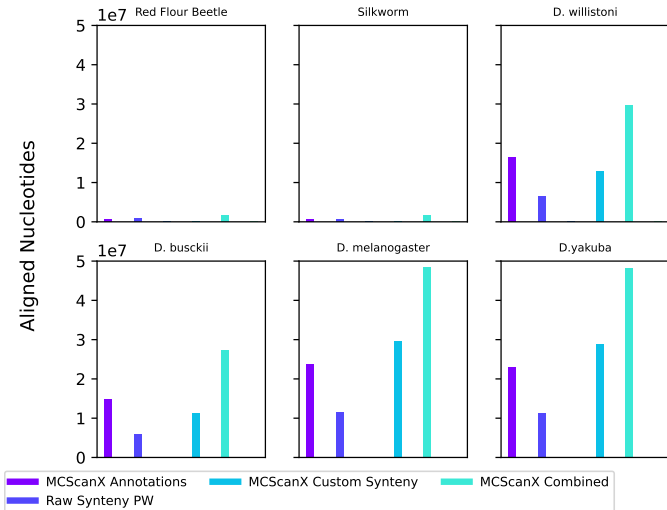
Performance of Different Synteny Calculation Strategies



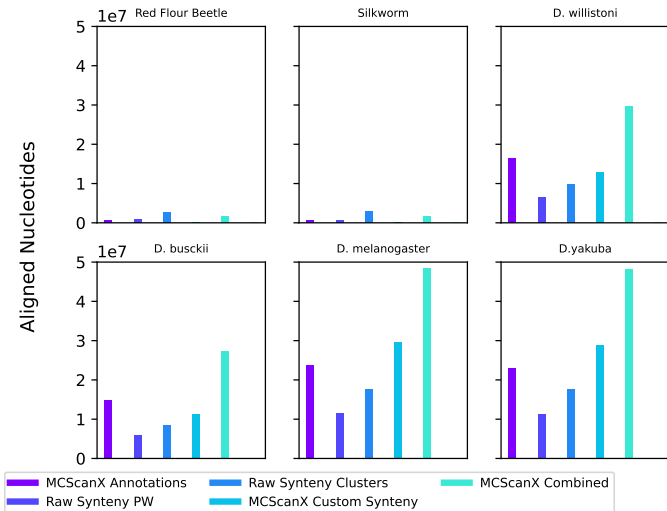
Performance of Different Synteny Calculation Strategies



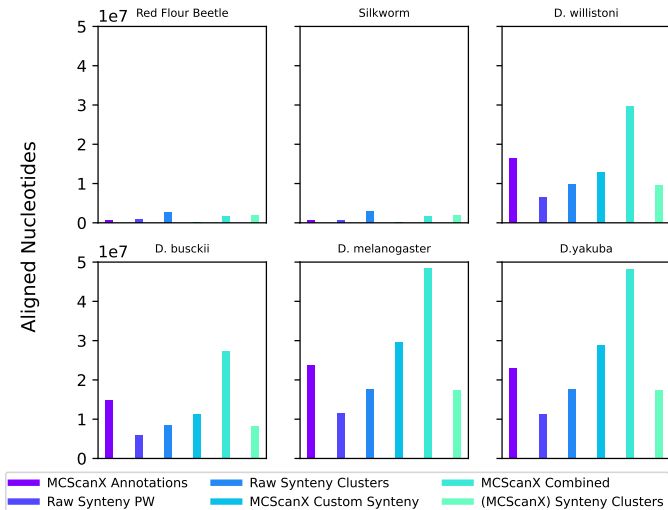
Performance of Different Synteny Calculation Strategies



Performance of Different Synteny Calculation Strategies



Performance of Different Synteny Calculation Strategies



Computational Resources

- set of 70 diptera (true flies) genomes
 - ▶ avg. size ~ 415 MB
 - ▶ from very poor to very good assemblies
 - ▶ I am told they cover the phylogeny of true flies well
- results
 - ▶ takes around 2 weeks on Leipzig Bioinf cluster
 - ▶ core results around 15 GB in python dict(s)
 - ▶ coverage of genome with candidates overall: $32.58 \% \pm 4.05$
 - ▶ length of $\sim 471 \pm 663$ and spacing in between of $\sim 918 \pm 190$
 - ▶ for candidates with ≥ 5 matches length is around 1000 and spacing 16000
 - ▶ how much of candidates are aligned somewhere (to any other species): $20.53 \% \pm 18.23$
 - ▶ how much of candidates are aligned counting the mean of all alignments per candidate: $3.20 \% \pm 2.29$



Altenhoff, Adrian M et al. "The Quest for Orthologs benchmark service and consensus calls in 2020". In: *Nucleic Acids Research* 48 (W1 July 2, 2020), W538–W545. ISSN: 0305-1048. DOI: [10.1093/nar/gkaa308](https://doi.org/10.1093/nar/gkaa308). URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7319555/> (visited on 09/08/2023).



Haas, Brian J. et al. "DAGchainer: a tool for mining segmental genome duplications and synteny". In: *Bioinformatics* 20.18 (Dec. 12, 2004), pp. 3643–3646. ISSN: 1367-4803. DOI: [10.1093/bioinformatics/bth397](https://doi.org/10.1093/bioinformatics/bth397). URL: <https://doi.org/10.1093/bioinformatics/bth397> (visited on 09/08/2023).



Moyers, Bryan A. and Jianzhi Zhang. "Further Simulations and Analyses Demonstrate Open Problems of Phylostratigraphy". In: *Genome Biology and Evolution* 9.6 (June 1, 2017), pp. 1519–1527. ISSN: 1759-6653. DOI: [10.1093/gbe/evx109](https://doi.org/10.1093/gbe/evx109). URL: <https://doi.org/10.1093/gbe/evx109> (visited on 09/08/2023).



Salzberg, Steven L. "Next-generation genome annotation: we still struggle to get it right". In: *Genome Biology* 20.1 (May 16, 2019)



Thanks for the attention :)