

# Improving Genome Assemblies With Syntenic Alignments

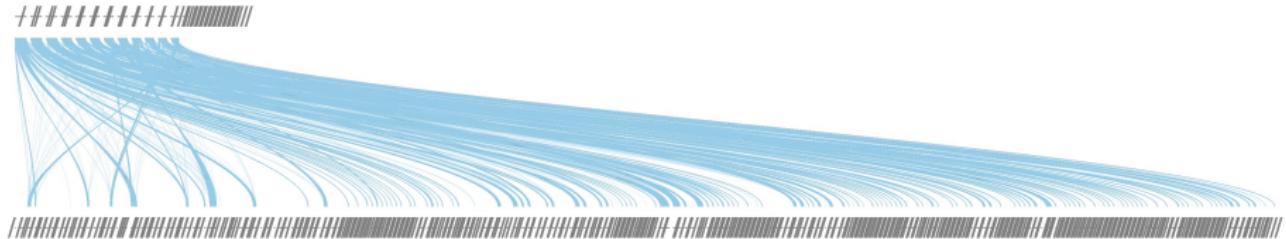
Karl Kaether

Abteilung Stadler, Leipzig

Bled 2025

## Context

- Synteny alignments provide signal for reference-guided contig scaffolding



## Context

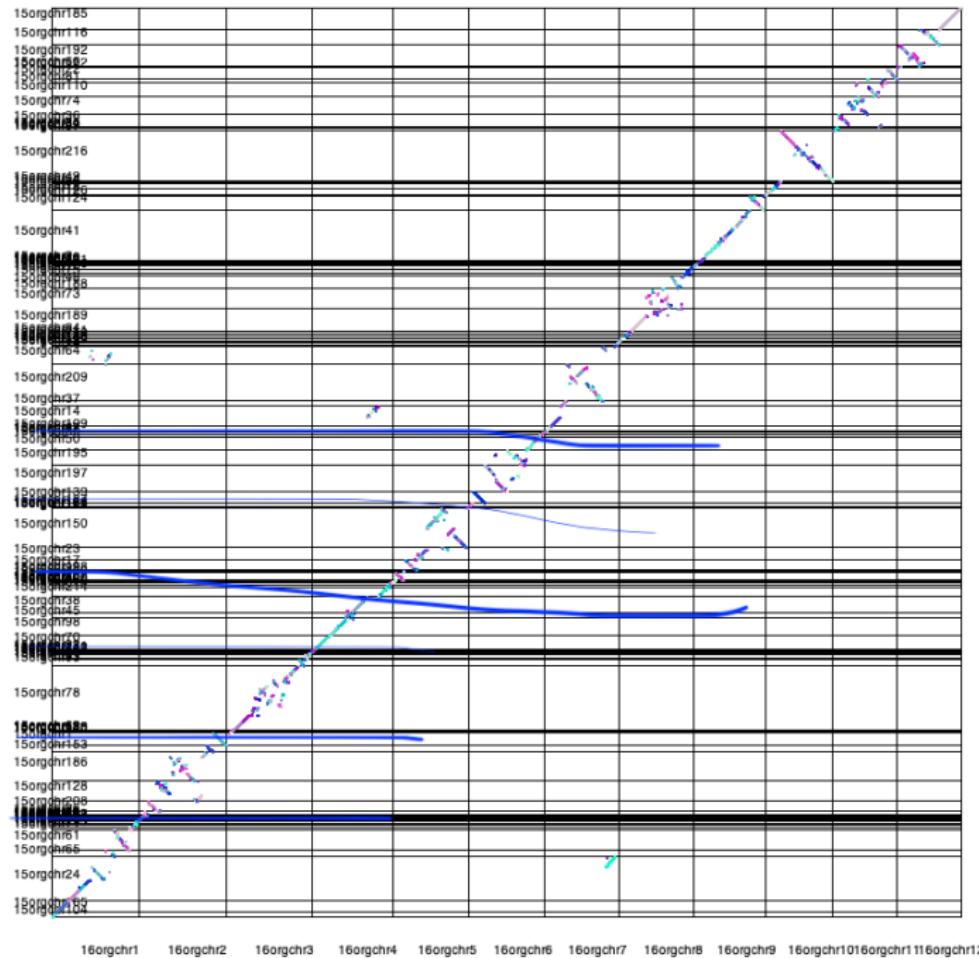
- Synteny alignments provide signal for reference-guided contig scaffolding
- note: not classical reference guided assembly but contig reordering
- tools exist: e.g. Ragout2 , MeDuSa, (Multi-)CSAR
- all use some kind of alignments + optimization to reduce number of rearrangements

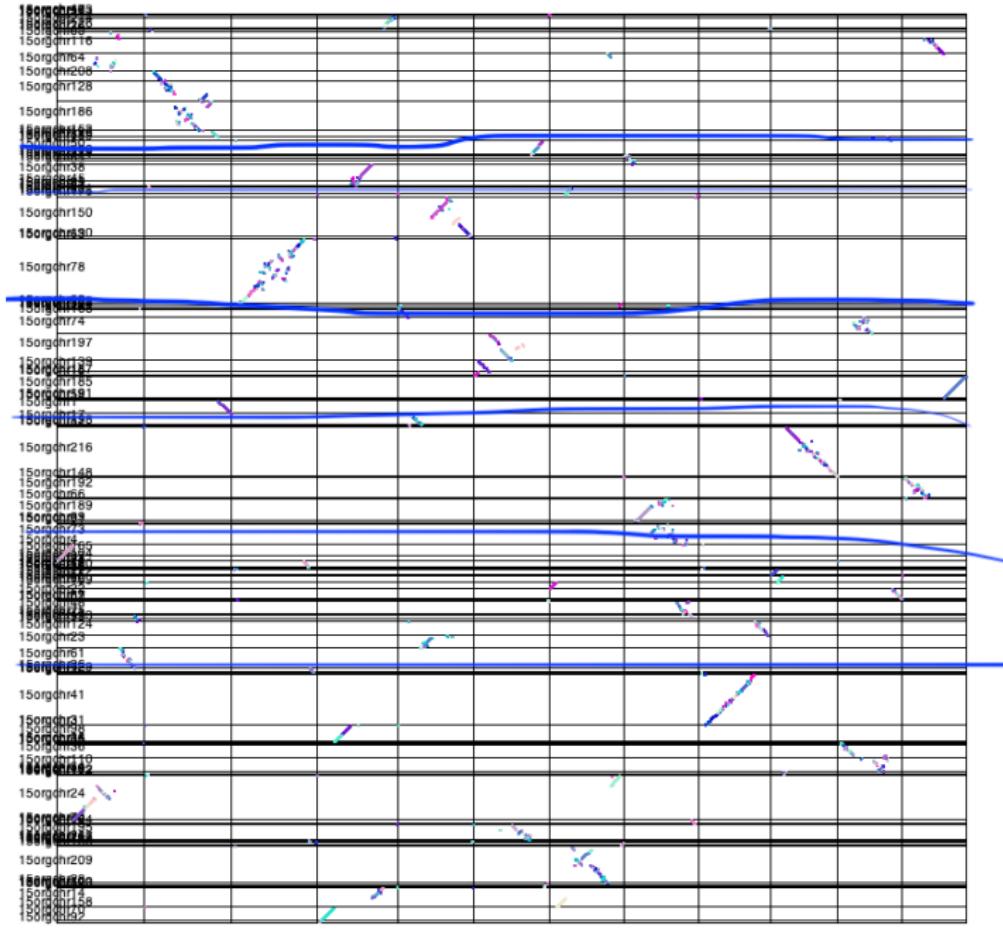
## Lets give it a try

- need some model for evolutionary events
- would be nice to be able to include multiple references
- example: *Solanum stellatiglandulosum* - kindly assembled by Thomas

## Lets give it a try

- need some model evolutionary events:
  - ▶ simple model - regard reference as fully assembled and contigs as reliable
  - ▶ enumerate where and which way anchors on target contigs align on reference chromosomes
  - ▶ orient contig according to the majority of alignments (and enumerate possible inversion events) - corresponds to minimizing inversion events
  - ▶ define the order of target contigs on reference chromosomes based on majority of alignments - corresponds to minimizing translocations

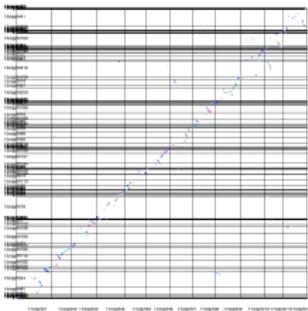




16orgchr1 16orgchr2 16orgchr3 16orgchr4 16orgchr5 16orgchr6 16orgchr7 16orgchr8 16orgchr9 16orgchr10 16orgchr11 16orgchr12

## Lets give it a try

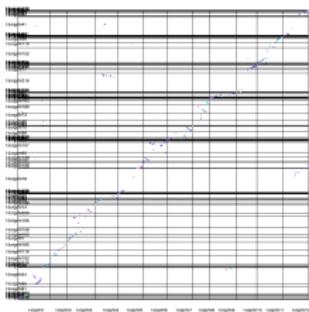
- would be nice to be able to include multiple references:
  - ▶ followed Multi-CSAR algorithm



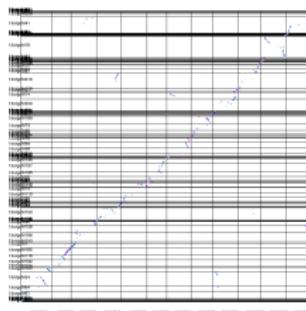
(a) *Solanum  
tuberousum*



(b) *Solanum  
lycopersicum*

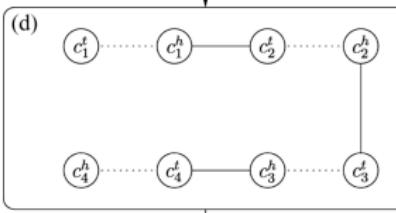
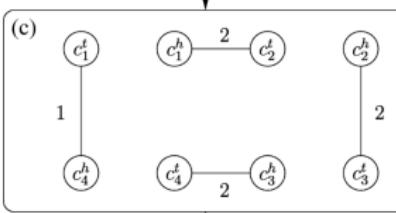
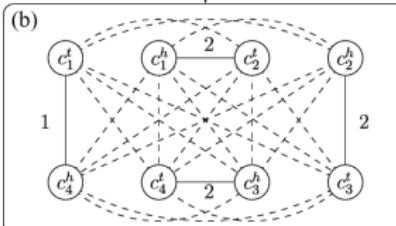


(c) *Capsicum annuum*

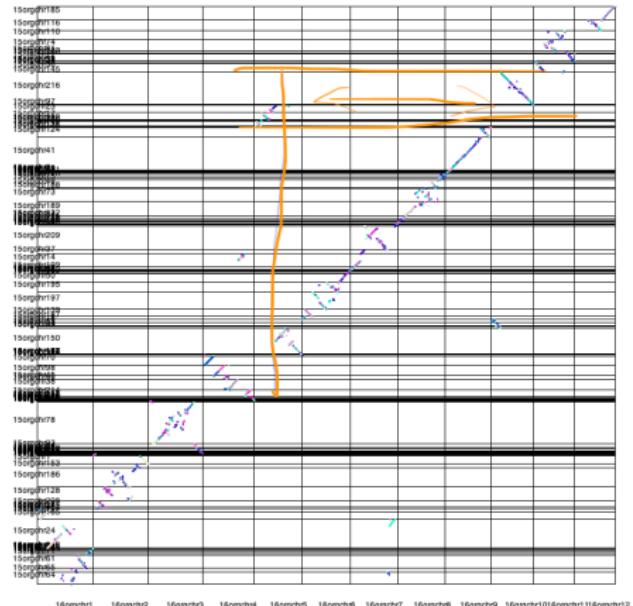


(d) *Solanum  
americanum*

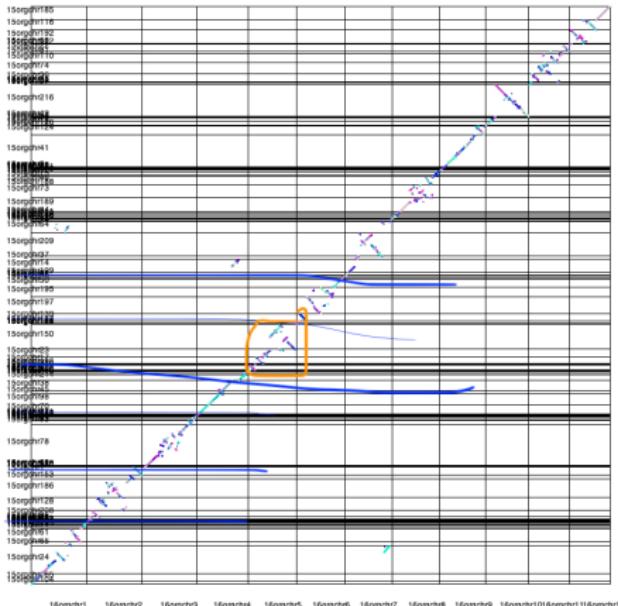
(a)  $T = \{c_1, c_2, c_3, c_4\}$   
 $S_1 = (+c_1, +c_2, +c_3)$   
 $S_2 = (+c_2, +c_3, +c_4)$   
 $S_3 = (-c_2, -c_1, -c_4, -c_3)$



(e) Scaffold of  $T = (+c_1, +c_2, +c_3, +c_4)$



(a) Consensus vs. *S.torvum*



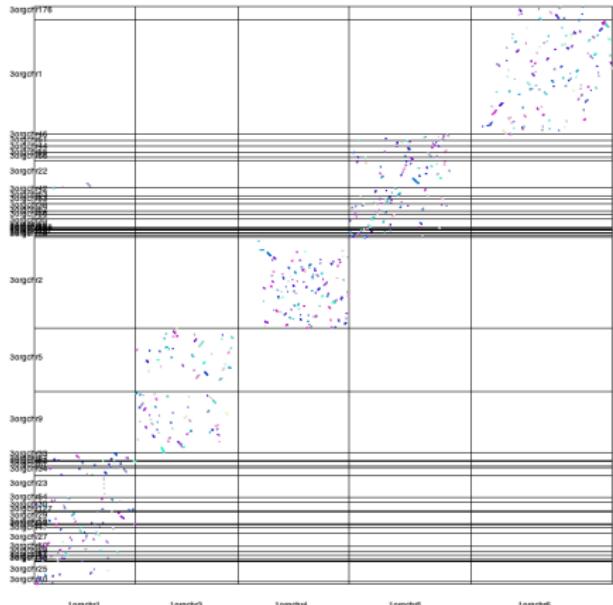
(b) *Solanum Torvum*

Table: Improved N50 from original 28 MB

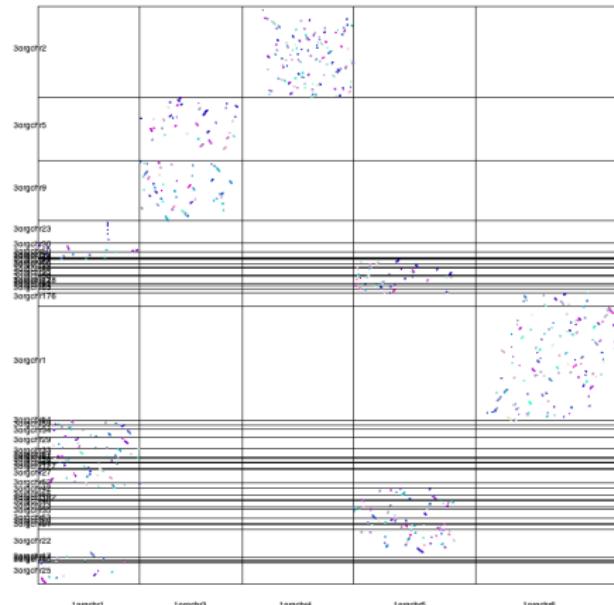
Reference Genome	scaffolded N50 all	scaffolded N50 limited
<i>C.annuum</i>	119 MB	92 MB
<i>S.americanum</i>	107 MB	95 MB
<i>S.candolleanum</i>	83 MB	57 MB
<i>S. dulcamara</i>	104 MB	87 MB
<i>S. galapagense</i>	107 MB	100 MB
<i>S. laciniatum</i>	51 MB	35 MB
<i>S. lycopersicum</i>	107 MB	107 MB
<i>S. lyratum</i>	108 MB	104 MB
<i>S. melongena</i>	107 MB	100 MB
<i>S. pimpinellifolium</i>	107 MB	100 MB
<i>S. pinnatisectum</i>	107 MB	95 MB
<i>S. pseudocapsicum</i>	79 MB	63 MB
<i>S. retroflexum</i>	77 MB	31 MB
<i>S. spirale</i>	73 MB	28 MB
<i>S. torvum</i>	108 MB	108 MB
<i>S. tuberosum</i>	107MB	100 MB
<i>S. viarum</i>	107 MB	99 MB
Consensus neutral	100 MB	102 MB
weights		
Consensus informed	101 MB	100 MB
weights		

## Take home

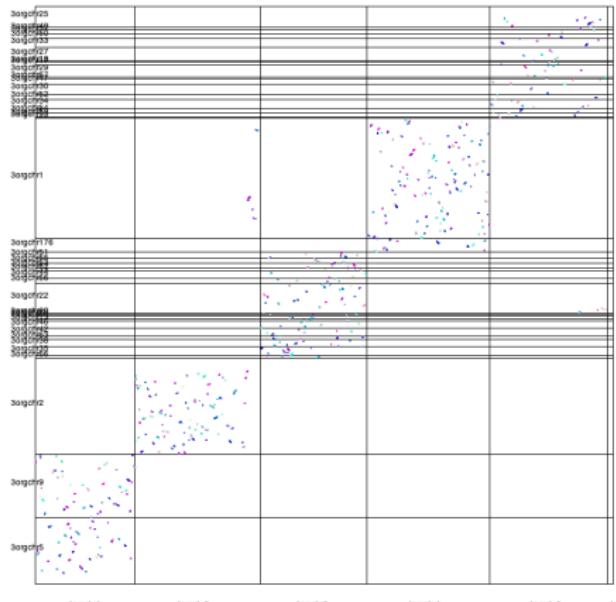
- reference-guided contig scaffolding possible without MSA or gene homology data
- way more benchmarking needed to see within what boundaries it can work
- potentially faster and more scalable
- really interesting bit are the potential breakpoints



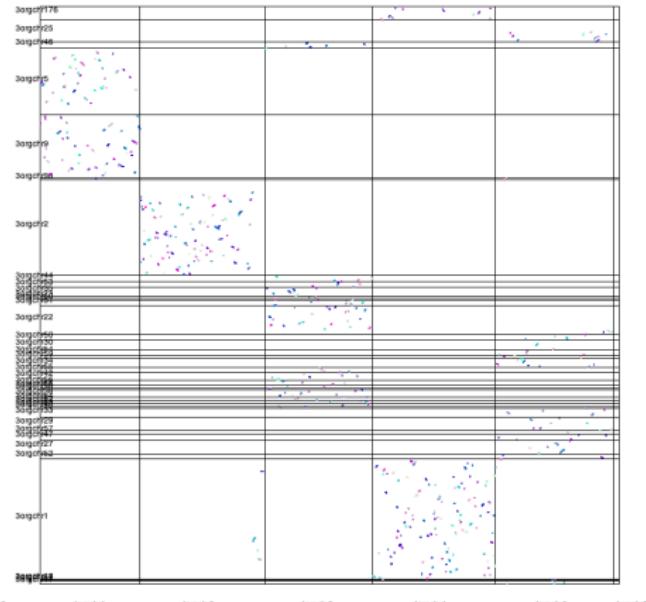
(a) D.mel



(b) D.mel CSAR



(a) *D. busckii*



(b) *D. busckii* CSAR

