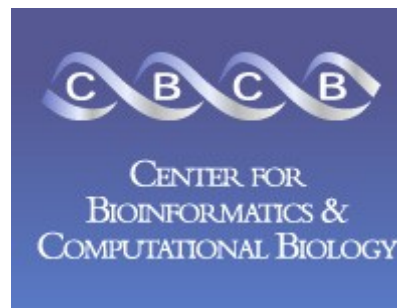


# First steps towards automated metagenomic assembly

Mihai Pop  
Sergey Koren  
Dan Sommer

University of Maryland, College Park

This presentation is licensed under the Creative Commons Attribution 3.0 Unported License available at <http://creativecommons.org/licenses/by/3.0/>

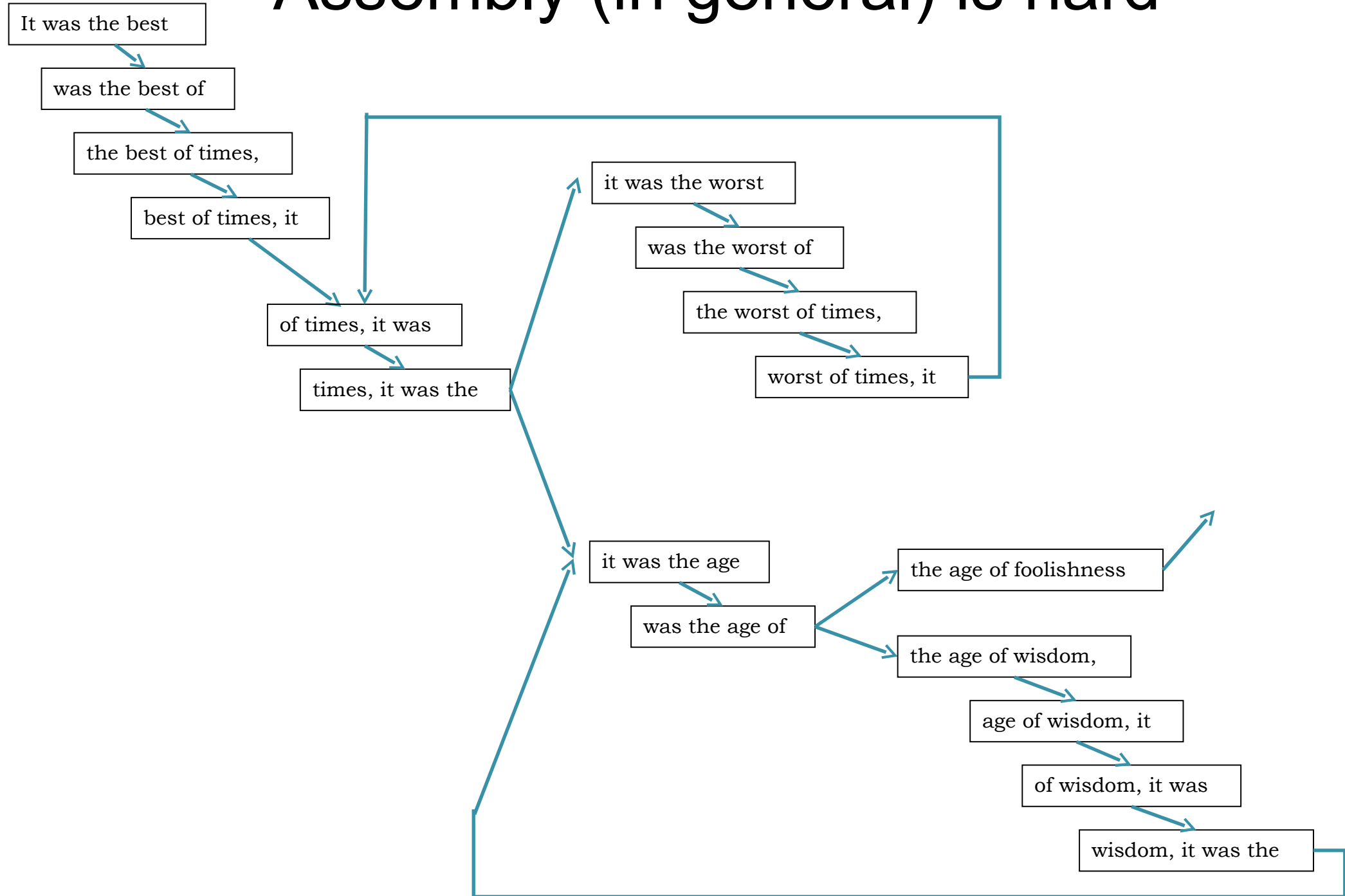


# Metagenomic assembly is impossible

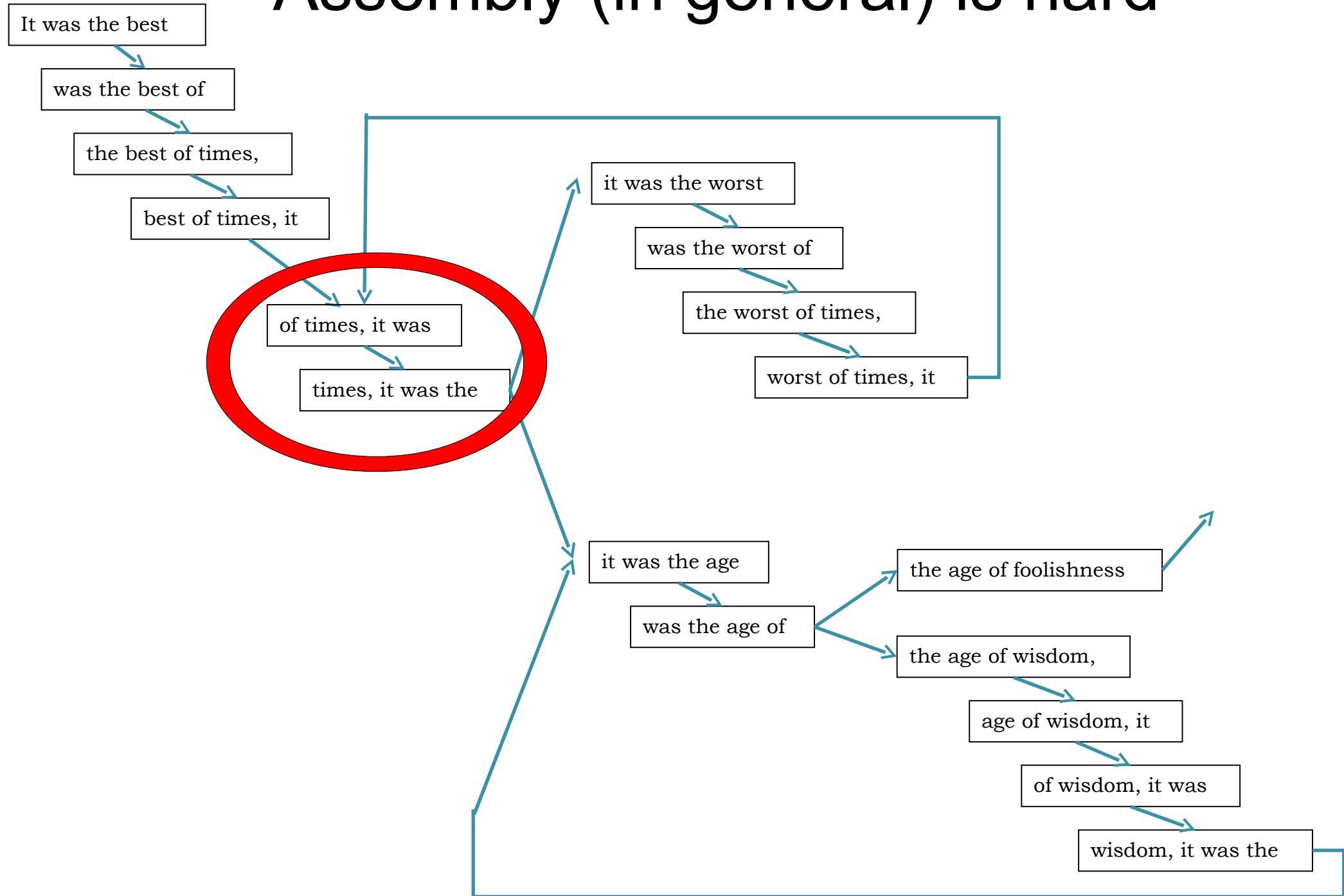
- Two competing goals:
  - assemble similar sequences from related genomes together
  - do not assemble similar sequences from unrelated genomes

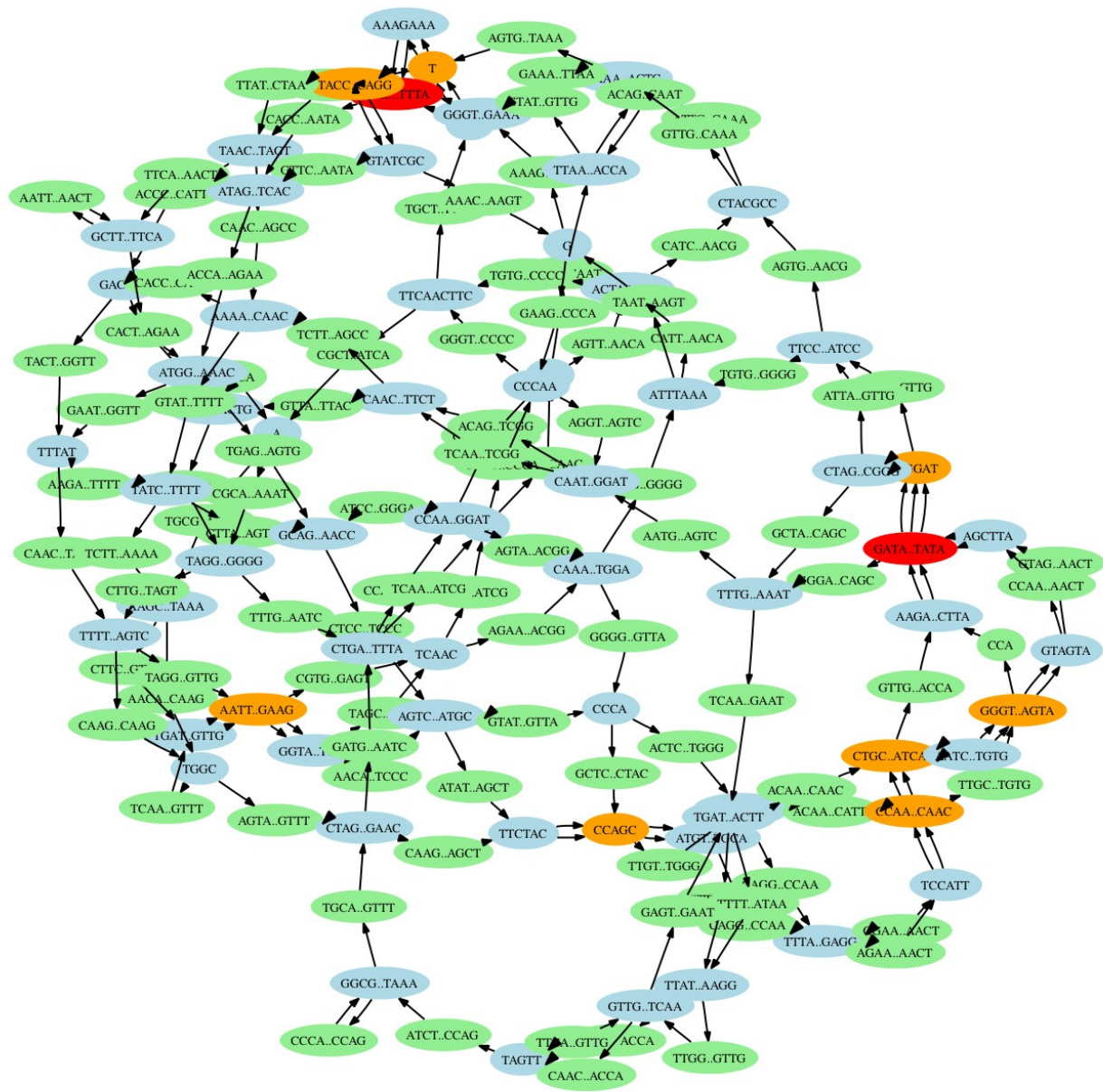
```
GCCTCCCGTAGGAGTTTGGACCGTGTCTCAGTTCCAATGTGGGGGACCTT
CATGCTGCCTCCCGTAGGAGTTTGGACCGTGTCTCAGTTCCAATGTG
TCCCGTAGGAGTCTGGTCCCGTGTCTCAGTACCAGTGTGGGGGACCTTCCTC
```

# Assembly (in general) is hard



# Assembly (in general) is hard

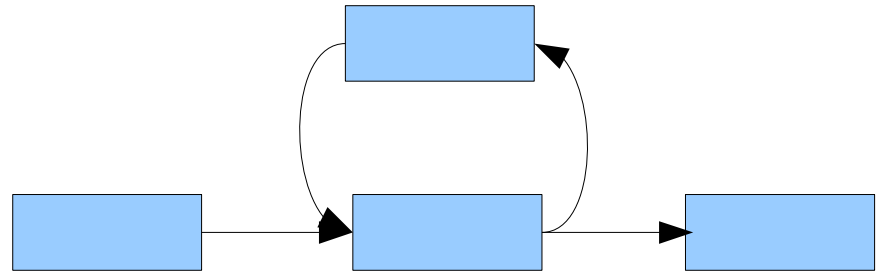




# Assembly (in general) is hard

- Repeats lead to ambiguity in reconstruction of genome  
(complexity exponential in # of repeats)

- Insufficient coverage:
  - gaps
  - obscures "true" assembly

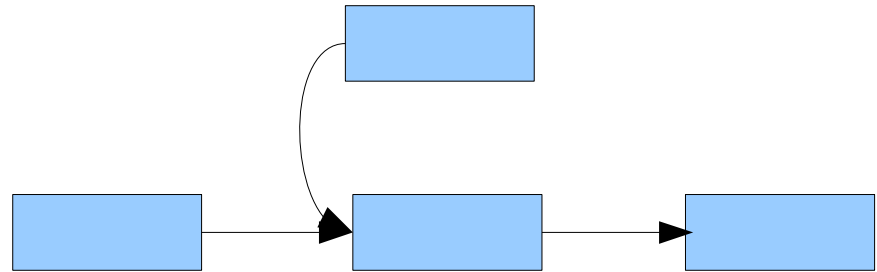


- Sequencing errors compound all other challenges

# Assembly (in general) is hard

- Repeats lead to ambiguity in reconstruction of genome  
(complexity exponential in # of repeats)

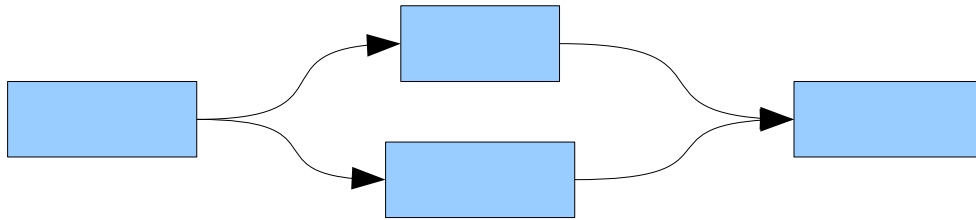
- Insufficient coverage:
  - gaps
  - obscures "true" assembly



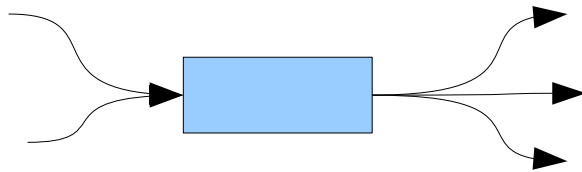
- Sequencing errors compound all other challenges

# Goals for a metagenomic assembler

- Must work well for clonal data
  - handle repeats
  - handle errors
  - deal with low coverage regions
- Must deal with polymorphisms
  - distinguish between errors and polymorphisms



- distinguish between repeats and polymorphisms

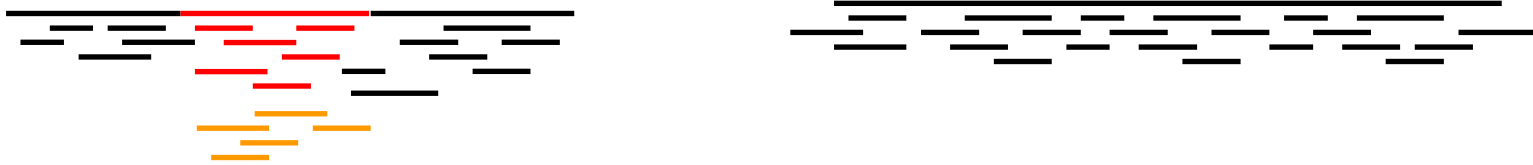
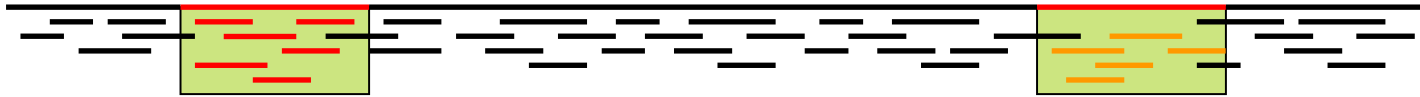


- enable discovery of polymorphisms/variation



# Repeat detection

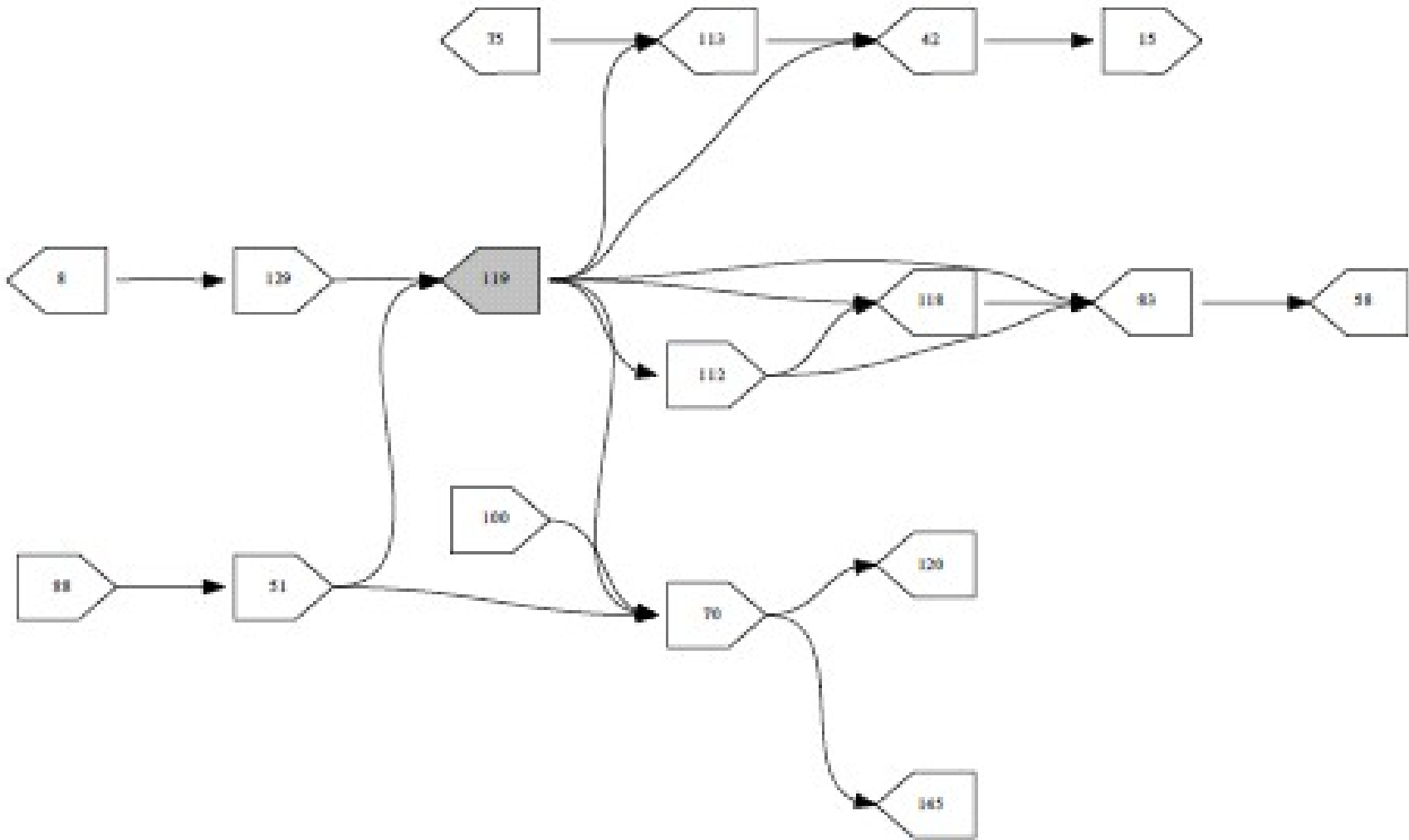
- Basic idea: deeply covered unitigs are repeats



- But.... in metagenomics the deeply covered contigs are often abundant organisms (which should be assembled)

# Better repeat detection

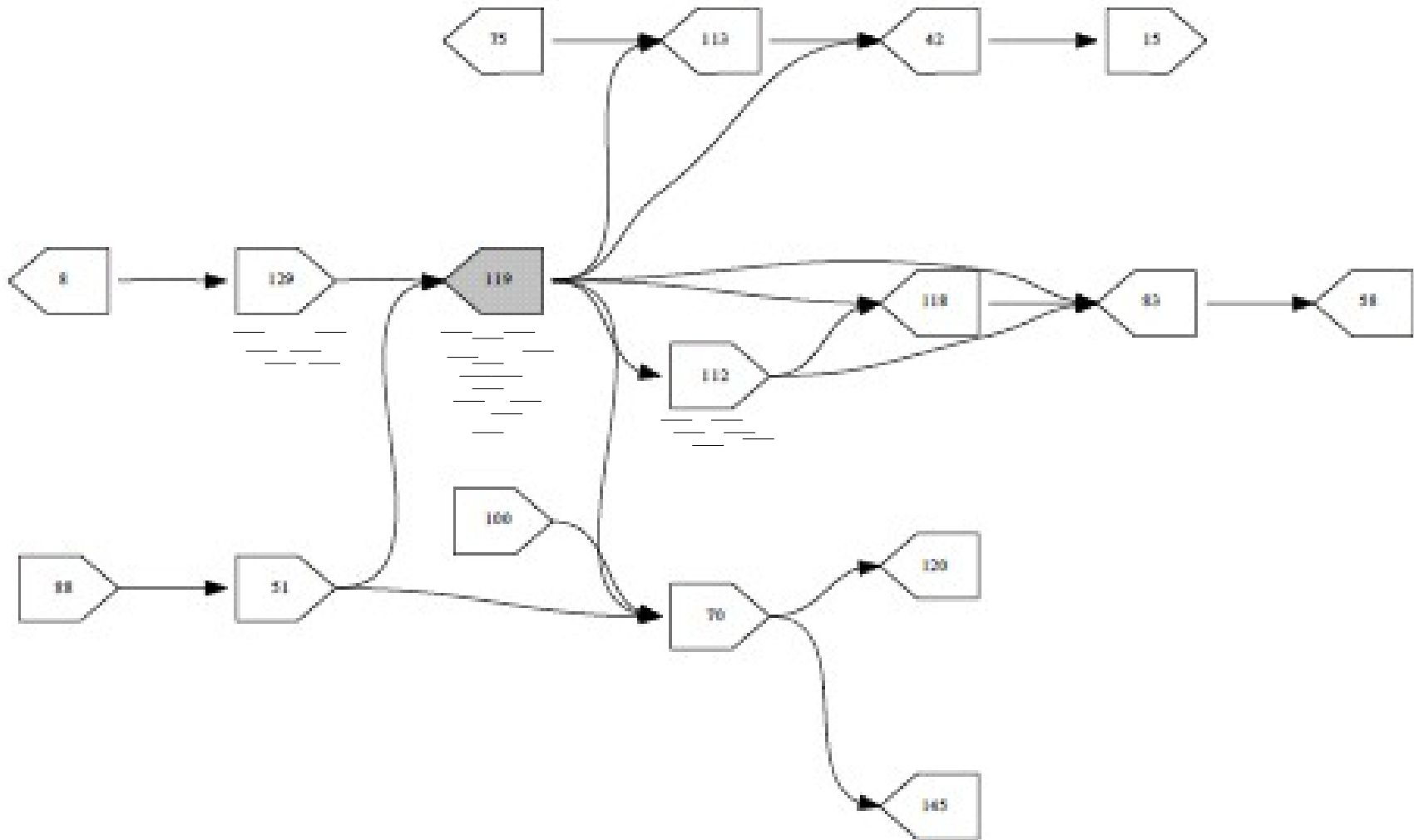
- Observation: repeats "tangle" the assembly graph



- Solution: find nodes "central" to the graph

# Better repeat detection...2

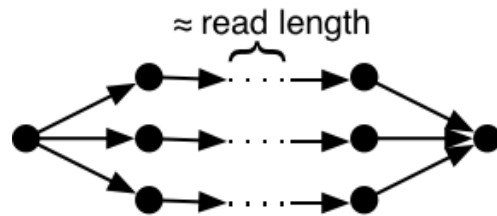
- Observation: depth of coverage is good local marker of repeats



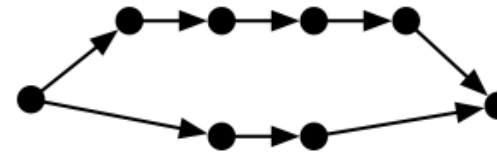
- Solution: run local depth-of-coverage statistics

# Detection of variation

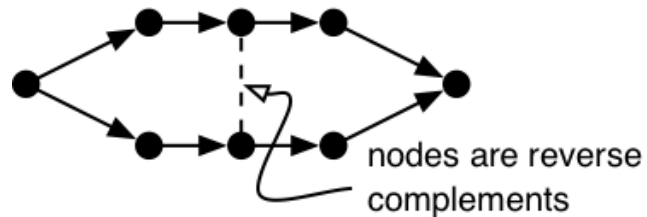
- Key idea: find assembly motifs that look like genome variation



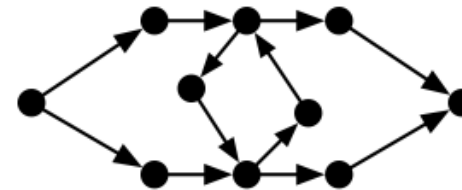
(a) SNPs



(b) Indels



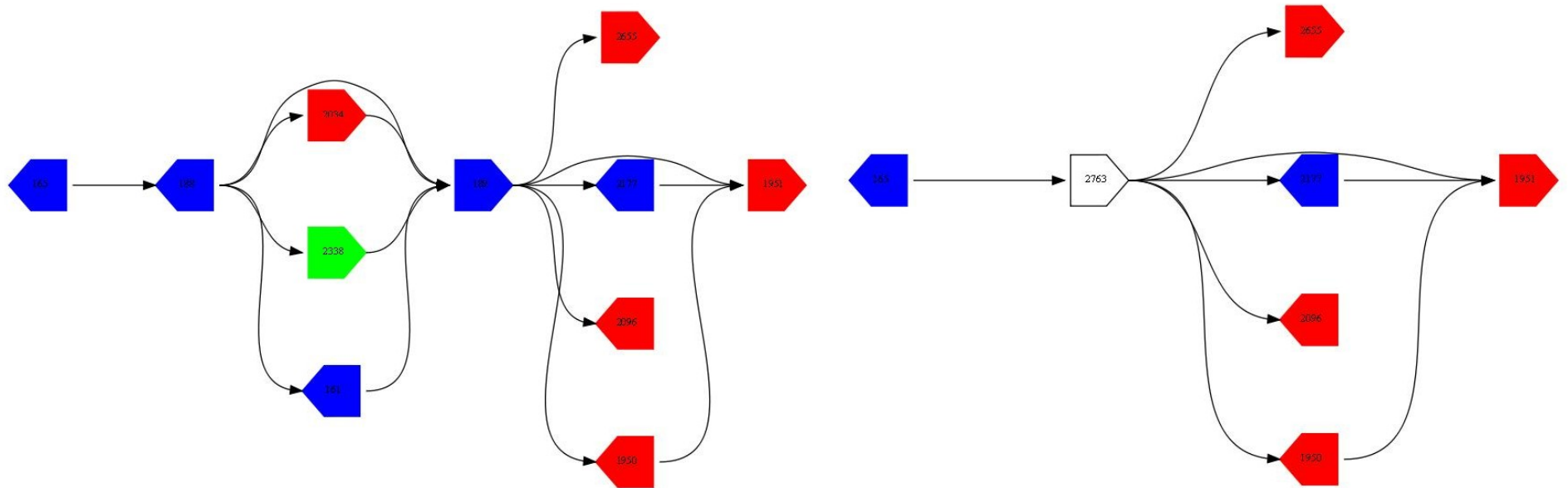
(c) Inversions



(d) Transpositions

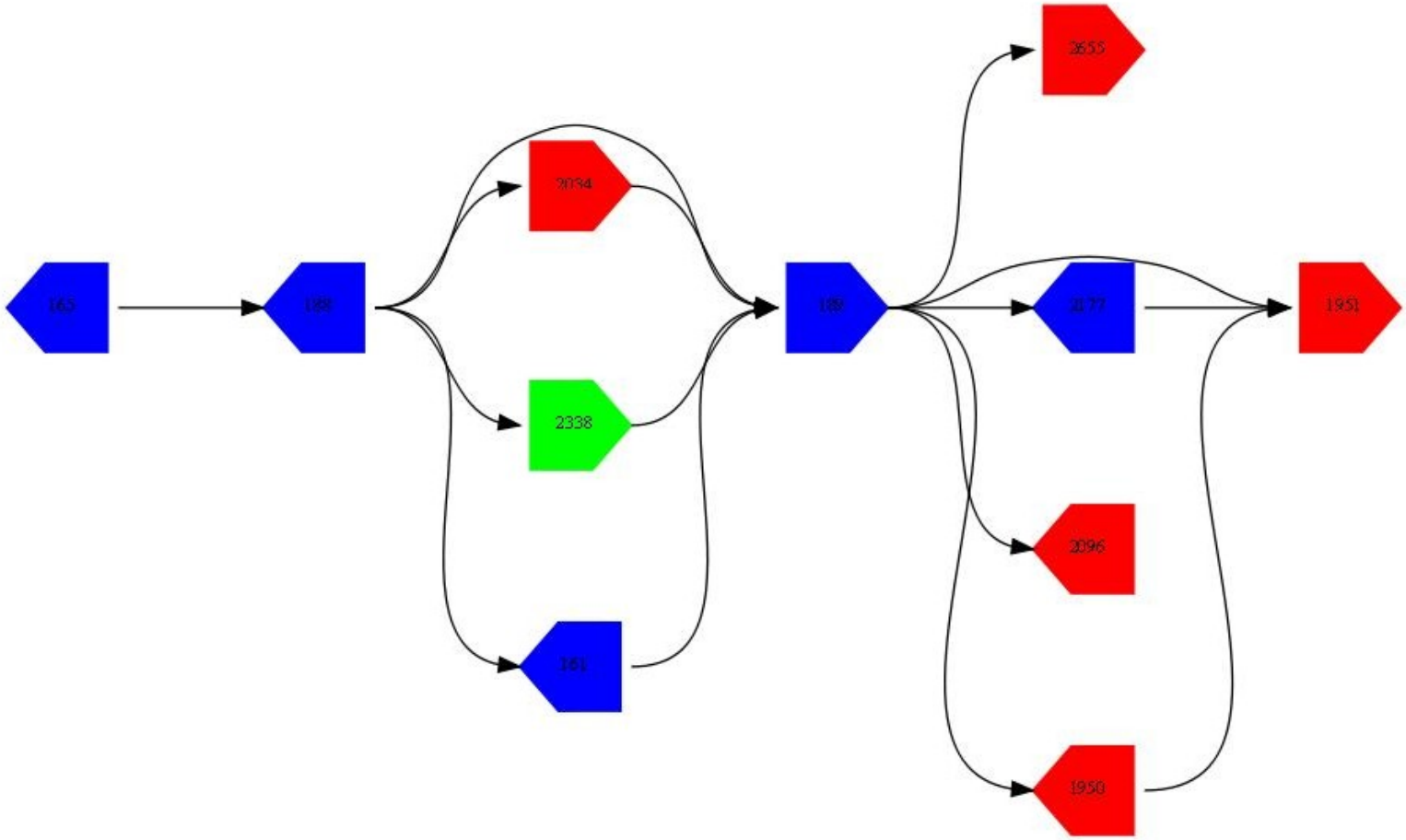
# Scaffolding through variation

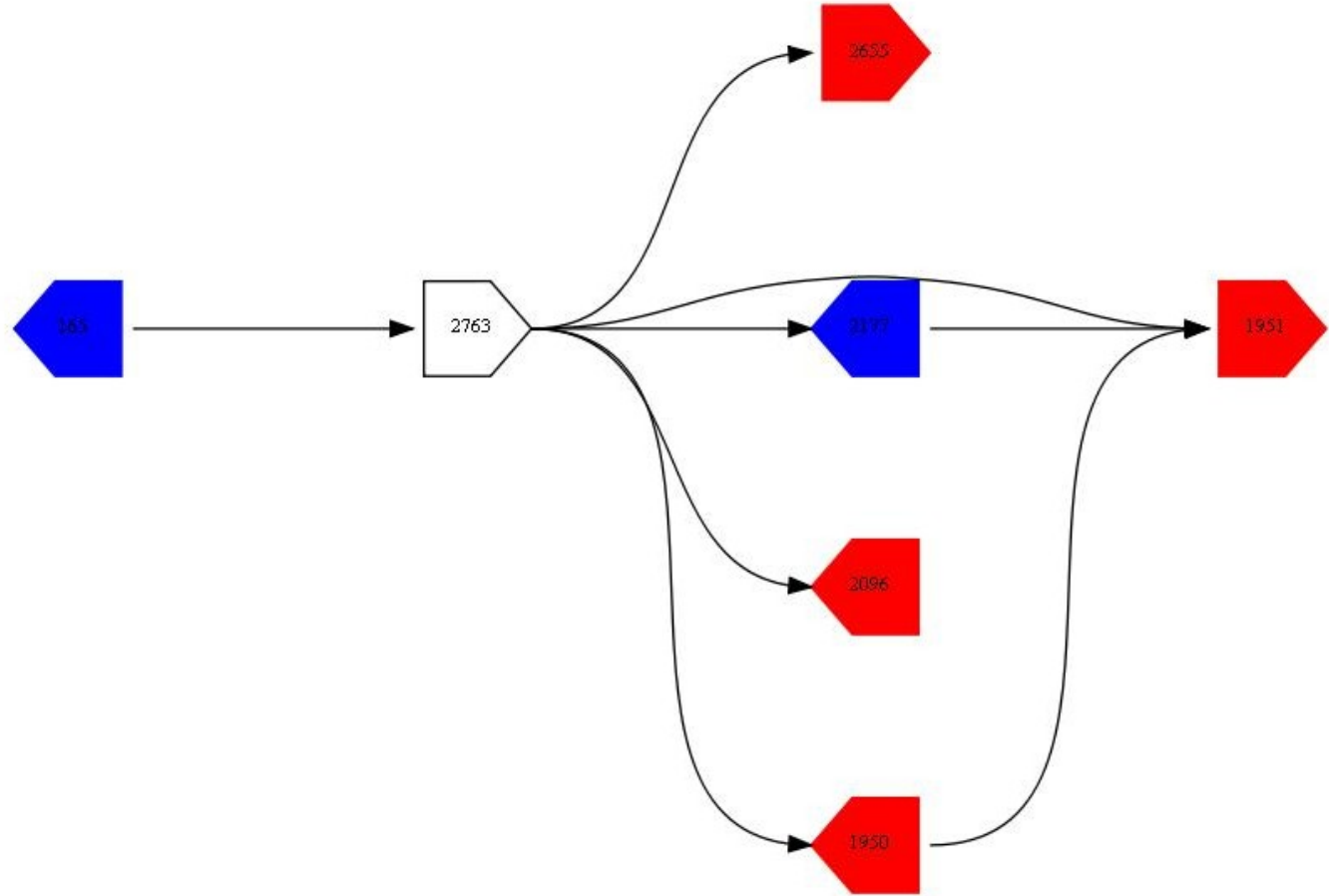
- Identify motifs and collapse them



# Detection of variation...cont

- Still work in progress
  - simple motifs can be found
  - looking for "enriched" motifs
  - still analyzing motifs found by our code
- Basic idea (find motif, collapse it) very computer graphic-y (that's how video games work...)
- Corrolary: This can lead to interactive visualizations of assembly graphs







# Bambus 2

- <http://www.cbcb.umd.edu/software/bambus>
- Can be used with output from most assemblers (tested with CA, Minimus, Newbler)
- Good repeat detection

Organism	ASM		# Repeats	# TP	# FP	# FN	SN	SP
<i>Brucella suis</i> 1330	Bambus 2	Component-Joining	6	6	0	5	54.54	100
		Local Coverage	14	5	9	6	45.45	94.61
		Total	20	11	9	0	100	94.61
	CA		30	8	22	0	100	79.43
Acid Mine	Bambus 2	Component-Joining	93	13	80	44	22.80	99.25
		Local Coverage	2,508	25	2483	32	43.85	76.77
		Total	2,601	38	2,563	19	66.66	76.03
	CA		4,749	49	4,700	8	85.96	51.88
	CA-met		1,126	24	1,102	28	46.15	82.47

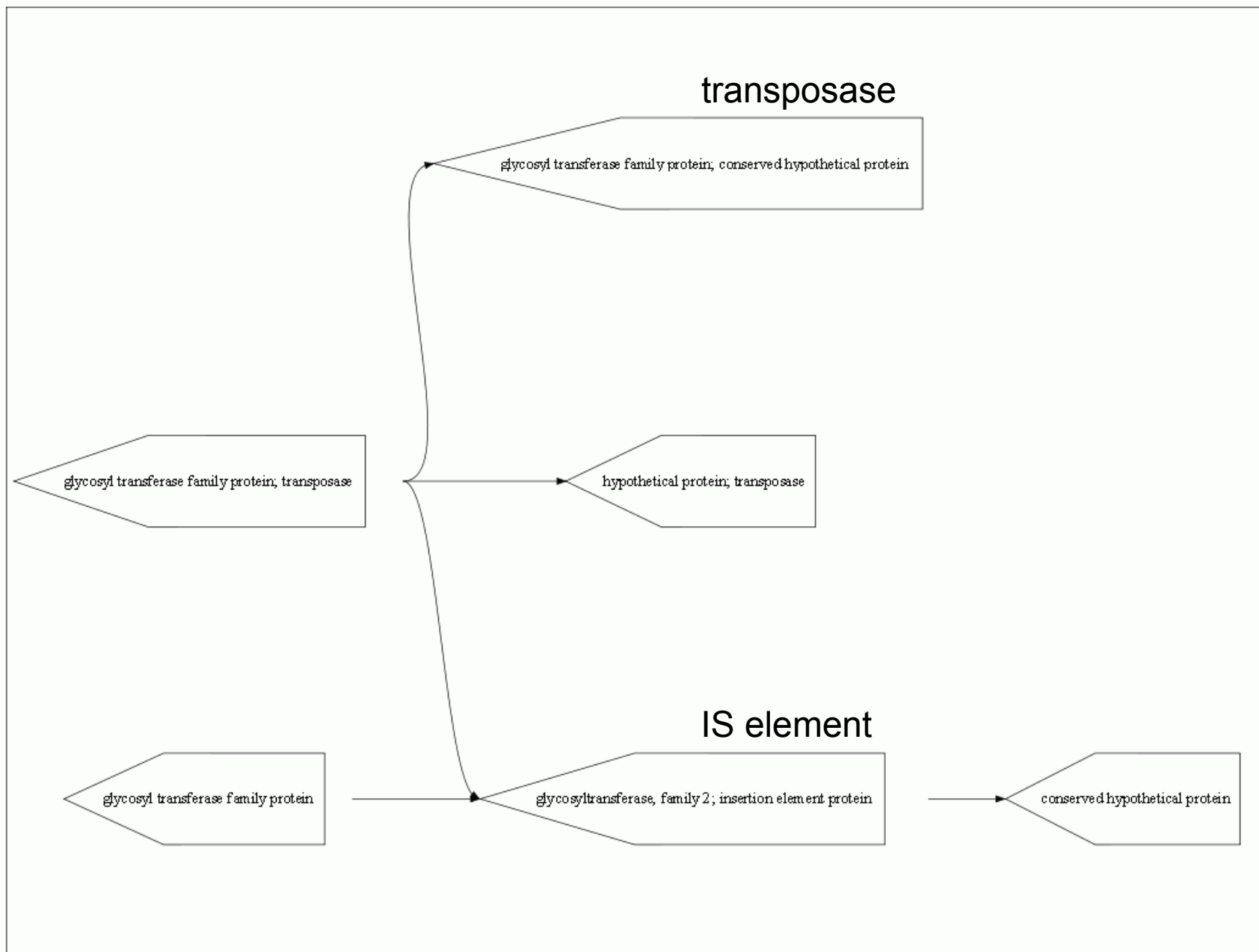
# Bambus 2

- Good genome reconstruction

ASM	Organism	# Scfs	% Genome
Published	<i>Leptospirillum sp</i> Group II '5-way CG'	70	81.86 %
	<i>Leptospirillum sp</i> Group III	474	120.49 %
	<i>Ferroplasma acidarmanus</i> Type I	170	112.66 %
	<i>Ferroplasma sp</i> Type II	59	138.54 %
	<i>Thermoplasmales archaeon</i> Gpl	410	121.44 %
CA	<i>Leptospirillum sp</i> Group II '5-way CG'	198	102.42 %
	<i>Leptospirillum sp</i> Group III	277	82.04 %
	<i>Ferroplasma acidarmanus</i> Type I	151	91.21 %
	<i>Ferroplasma sp</i> Type II	342	112.48 %
	<i>Thermoplasmales archaeon</i> Gpl	405	87.93 %
CA-met	<i>Leptospirillum sp</i> Group II '5-way CG'	101	97.14 %
	<i>Leptospirillum sp</i> Group III	234	81.69 %
	<i>Ferroplasma acidarmanus</i> Type I	62	90.15 %
	<i>Ferroplasma sp</i> Type II	90	99.46 %
	<i>Thermoplasmales archaeon</i> Gpl	179	83.60 %
Bambus 2	<i>Leptospirillum sp</i> Group II '5-way CG'	109	102.40 %
	<i>Leptospirillum sp</i> Group III	103	84.78 %
	<i>Ferroplasma acidarmanus</i> Type I	26	102.13 %
	<i>Ferroplasma sp</i> Type II	237	112.04 %
	<i>Thermoplasmales archaeon</i> Gpl	167	94.90 %

# Variation motif

## Glycosyl transferase hypervariable locus in *Leptospirillum*



# Future work

- Better documentation
- Better integration with other assemblers
- Tool for inspection of scaffolding data
- New types of variation
- Variation analysis toolkit

# Acknowledgments

- Steven Salzberg
- Carl Kingsford
- Sergey Koren
- Dan Sommer
- Bo Liu
- Mohammad Ghodsi
- Ted Gibbons
- Todd Treangen

NIH R01-HG-004885

NSF IIS-0812111