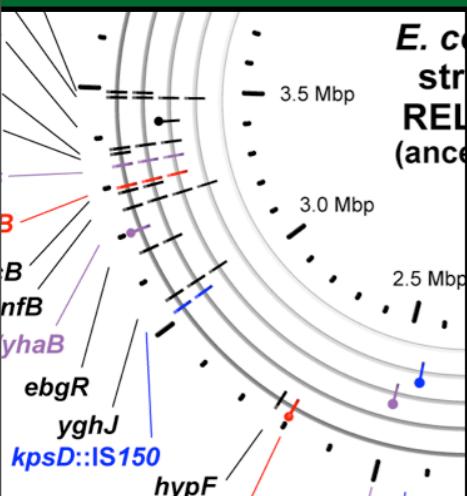
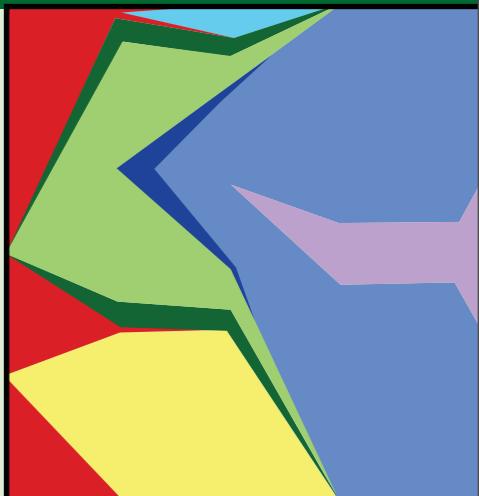


Re-sequencing hundreds of evolved *E. coli* genomes:



Finding non-SNP mutations,
analyzing mixed populations, and
knowing what you don't know

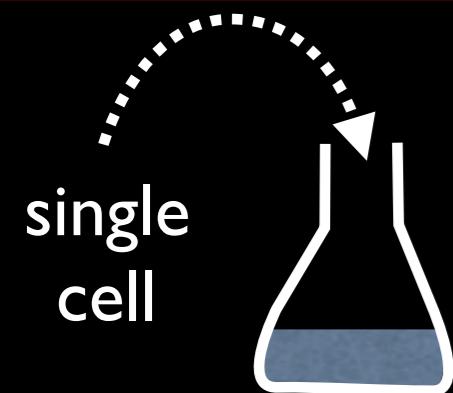


Jeff Barrick

Laboratory of Rich Lenski
Microbiology and Molecular Genetics

MICHIGAN STATE
UNIVERSITY

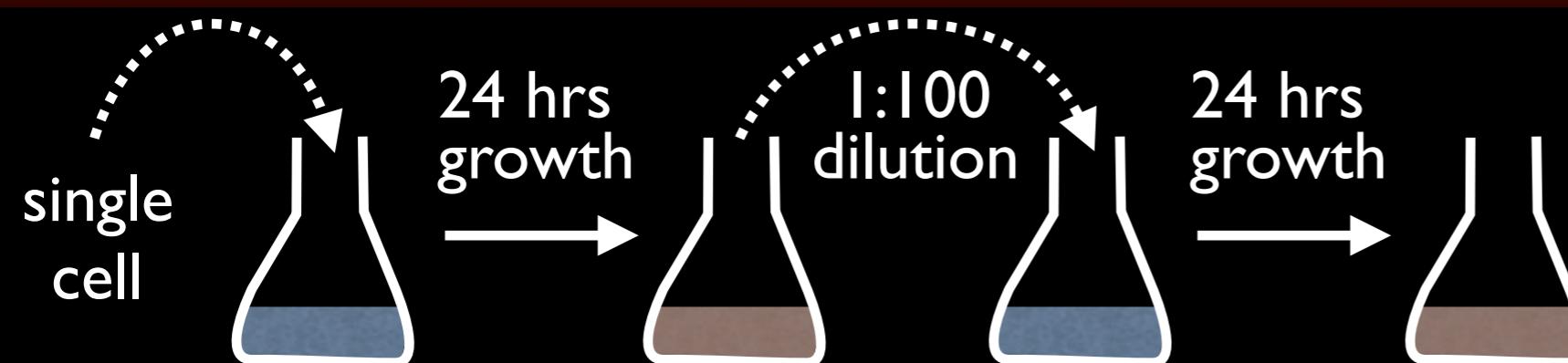
Lenski long-term evolution experiment



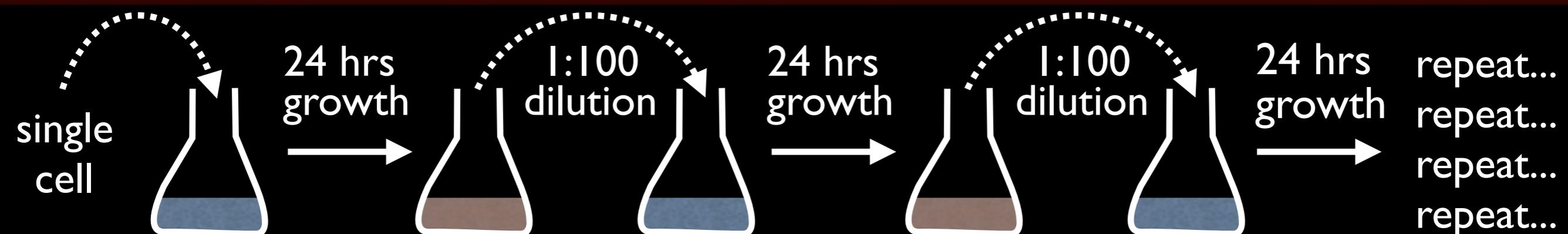
Lenski long-term evolution experiment



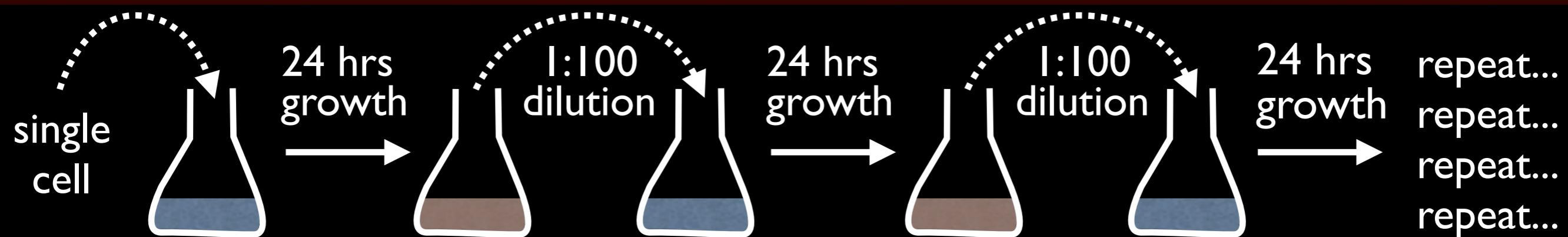
Lenski long-term evolution experiment



Lenski long-term evolution experiment



Lenski long-term evolution experiment



- ❖ 12 independent populations evolved >20 yrs.
Frozen “fossil record” has been archived.
- ❖ **How many and what mutations?**
- ❖ Compare rates of genomic change and fitness increase, monitor diversity in the population, understand molecular basis of adaptation.

Overview

- Application: **Re-sequencing** (*microbes*)
- Main Platform: **Illumina Genome Analyzer**

1. Overview of strategies and sequencing data
2. **breseq** (“brēz-sēk”) bacterial re-sequencing pipeline
 - characteristics of a typical data set
 - identifying different kinds of mutations
 - analysis of SNPs in a mixed population
3. Asking evolutionary questions

Before I forget...

Thanks to the Lenski lab, particularly
Brian Baer, Neerja Hajela,
Zachary Blount, and Justin Meyer.

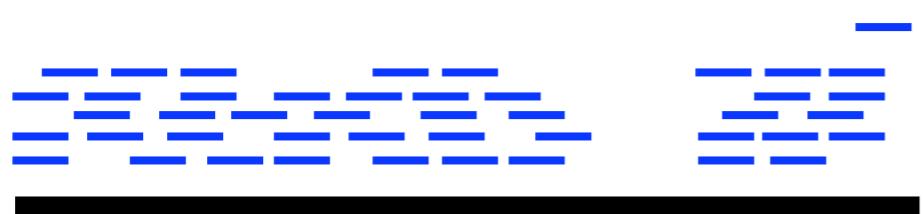
Thanks for coding discussions/contributions
Dave Knoester

Thanks to collaborators
Jihyun Kim et al. (KRIBB)
Dom Schneider et al. (Grenoble)
Genoscope

Strategies for finding mutations

re-sequencing

- ❖ map reads to known reference genome (ssaha2, maq, bowtie, ...)



- ❖ infer mutations

de novo assembly

- ❖ assemble reads by overlap (velvet, ABySS...)



- ❖ map contigs to reference genome



- ❖ infer mutations

Strategies for library preparation

single-end



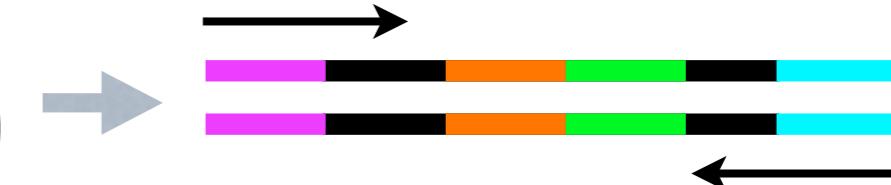
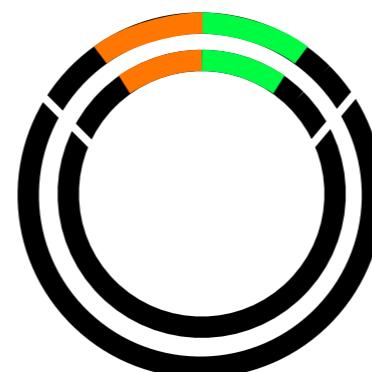
independent reads

paired-end



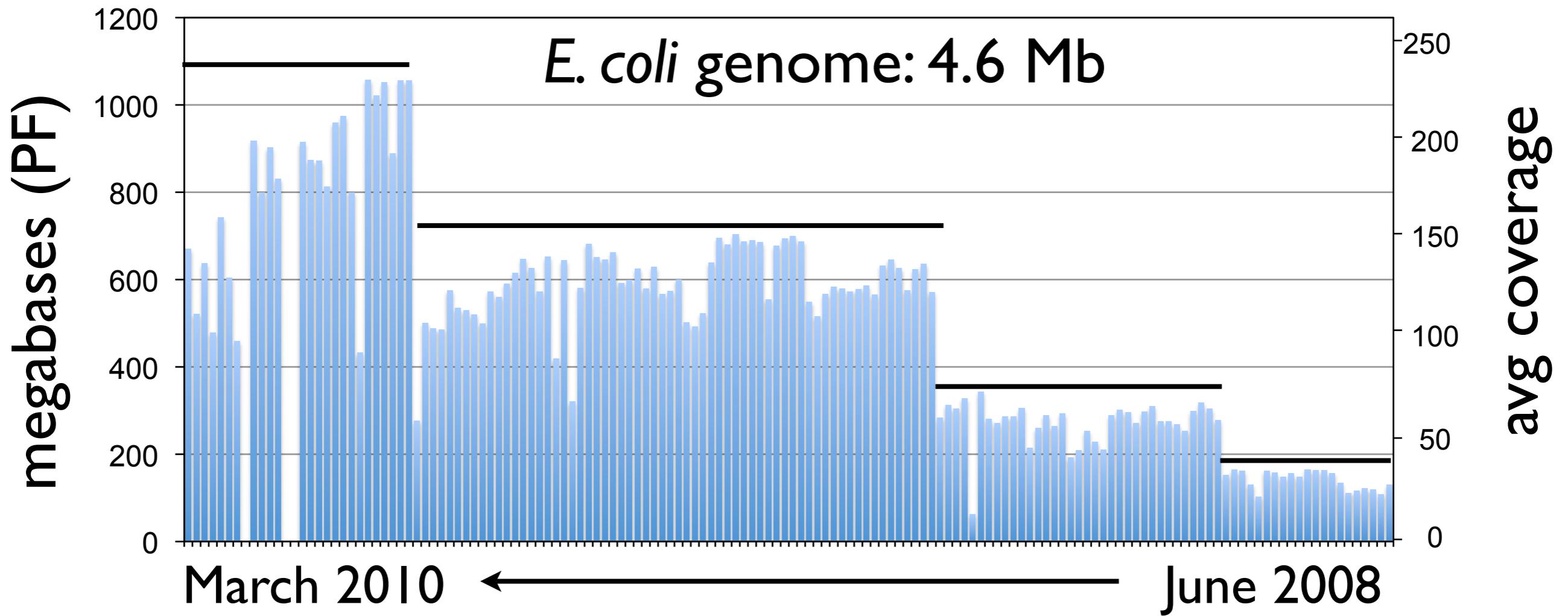
two inwardly oriented
reads separated by ~200 nt

mate-paired



two outwardly oriented reads separated by ~3000 nt

Sequence Data Sets

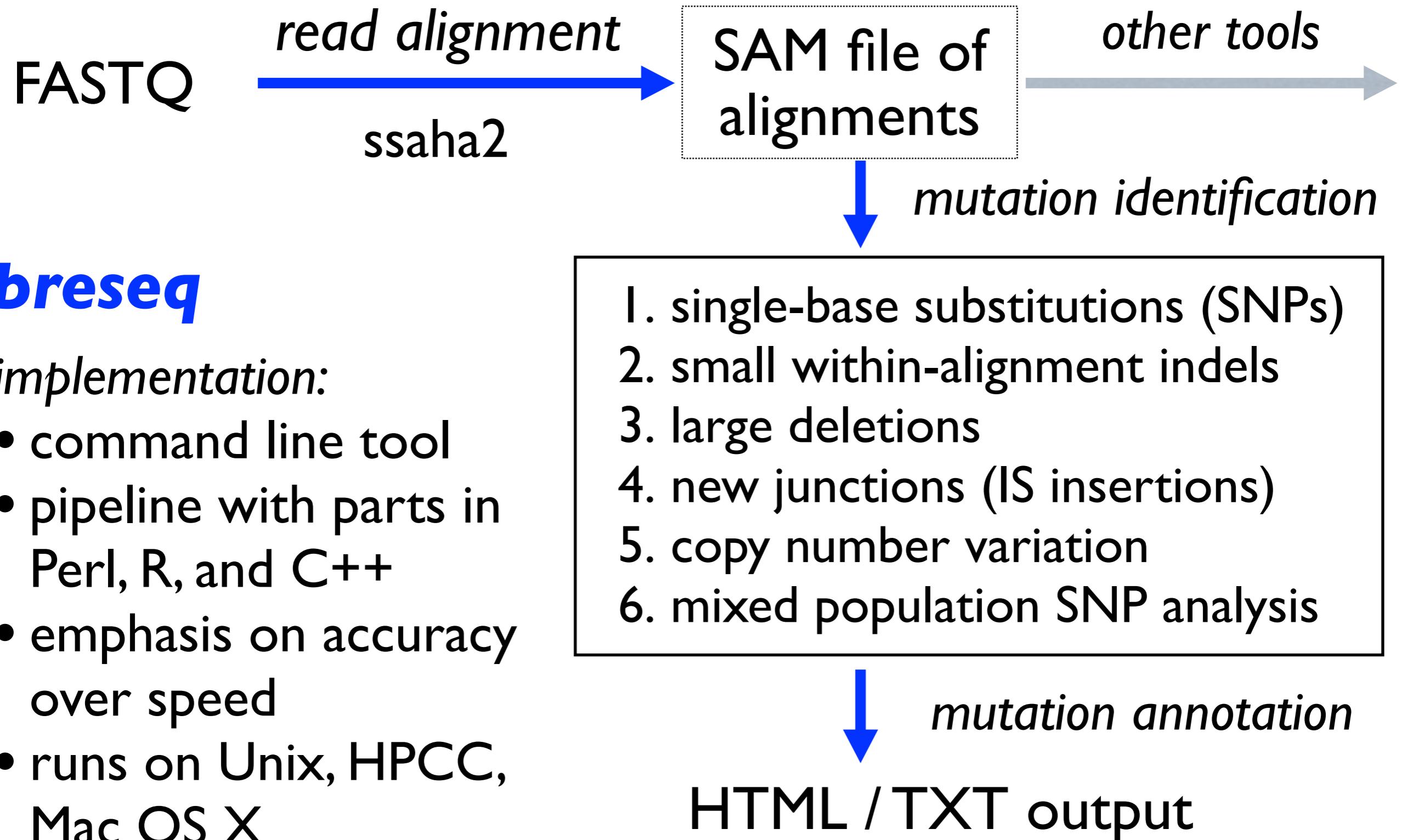


| 48 samples at RTSF (shown)

| 20 samples at Genoscope with Dom Schneider et al.

mostly 1 lane per genome, 36-bp single-end reads

breseq re-sequencing pipeline



SAM: Sequence Alignment/Map format

Text (SAM) and binary (BAM) files organized for quick retrieval of reads aligned to a certain position.

Created to support 1000 Genomes project by a team at the Sanger Center.

SAMtools (<http://samtools.sourceforge.net/>)

- C library with bindings to Java, Perl, Python, Lisp, etc.
- Command-line tools for manipulation, consensus/indel calling, viewing alignments as text, ...
- Many read aligners output in SAM format

Knowing what you don't know

1. Theoretical limits: Read length and pair distance.
2. Practical limits: Base quality and coverage evenness.

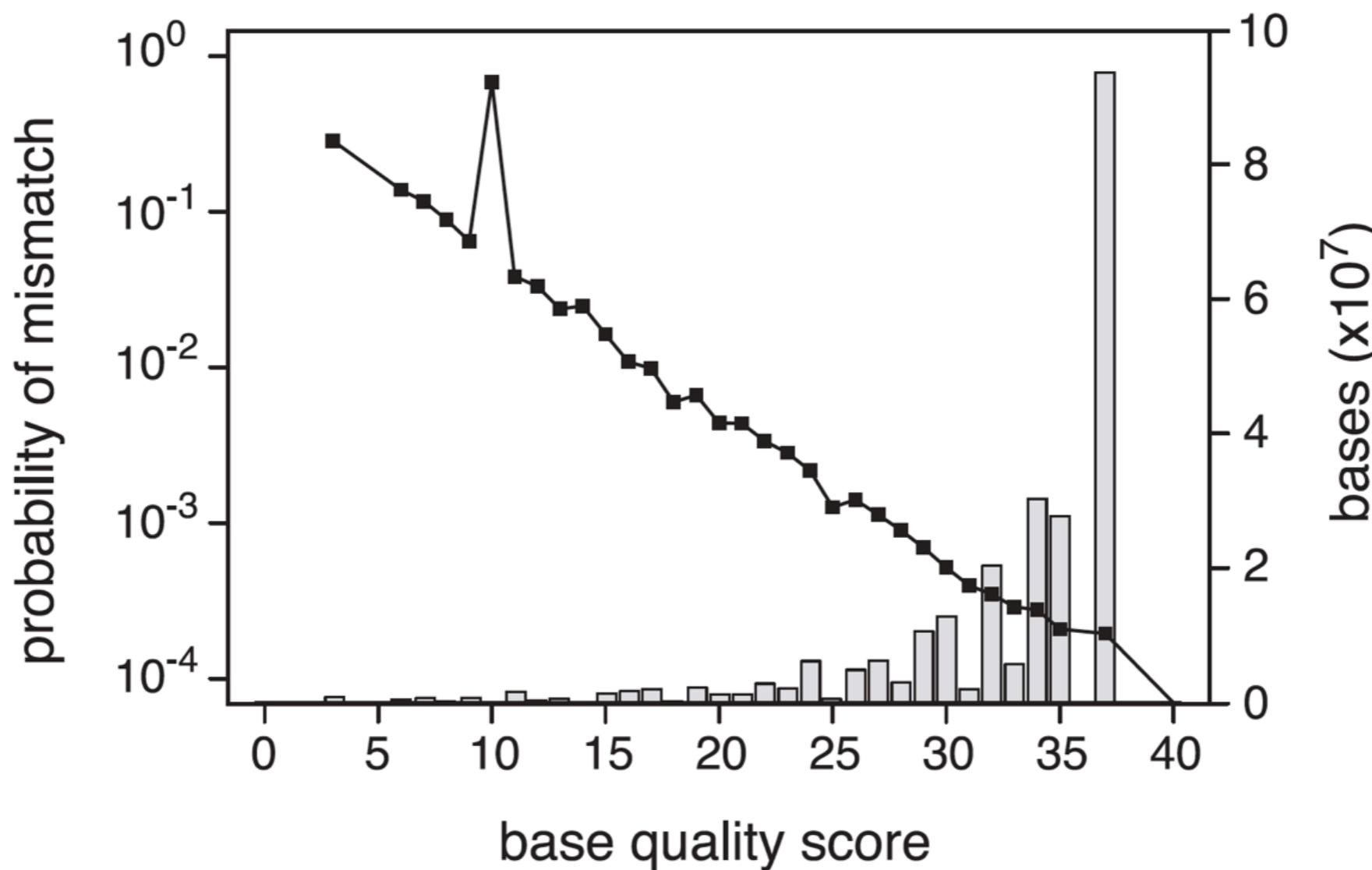
	single-end	paired-end	mate-paired
IS insertions	*	*	*
duplications	*	*	*
inversions across IS	—	—	*
SNPs in repeats	—	—	*
insertion of new seq	—	—	—

IS = bacterial mobile elements 0.8-1.5 kb in length.

Need standardized metrics to describe completeness of re-sequencing data on a per-base per-genome basis.

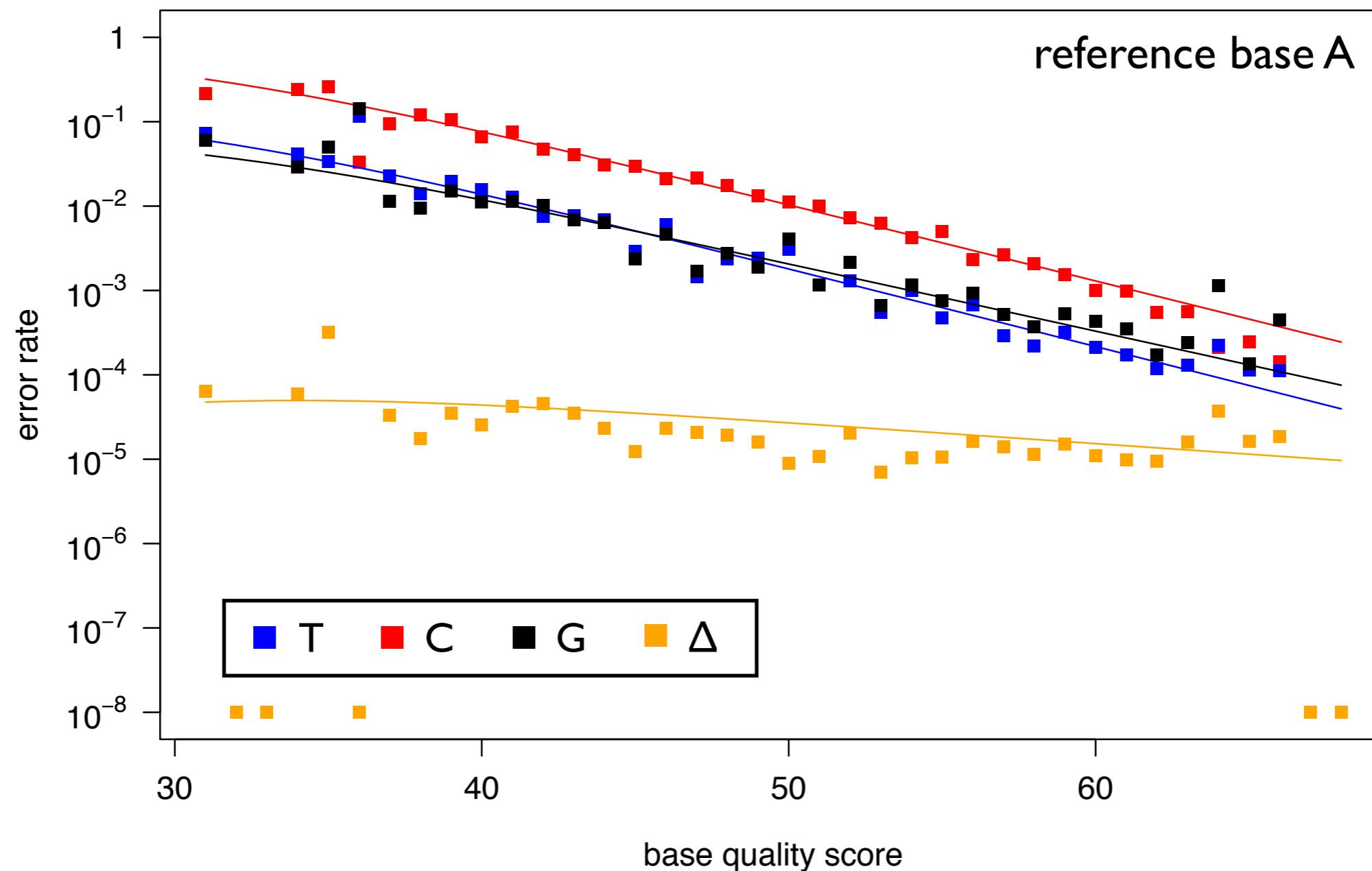
Typical Base Error Rates

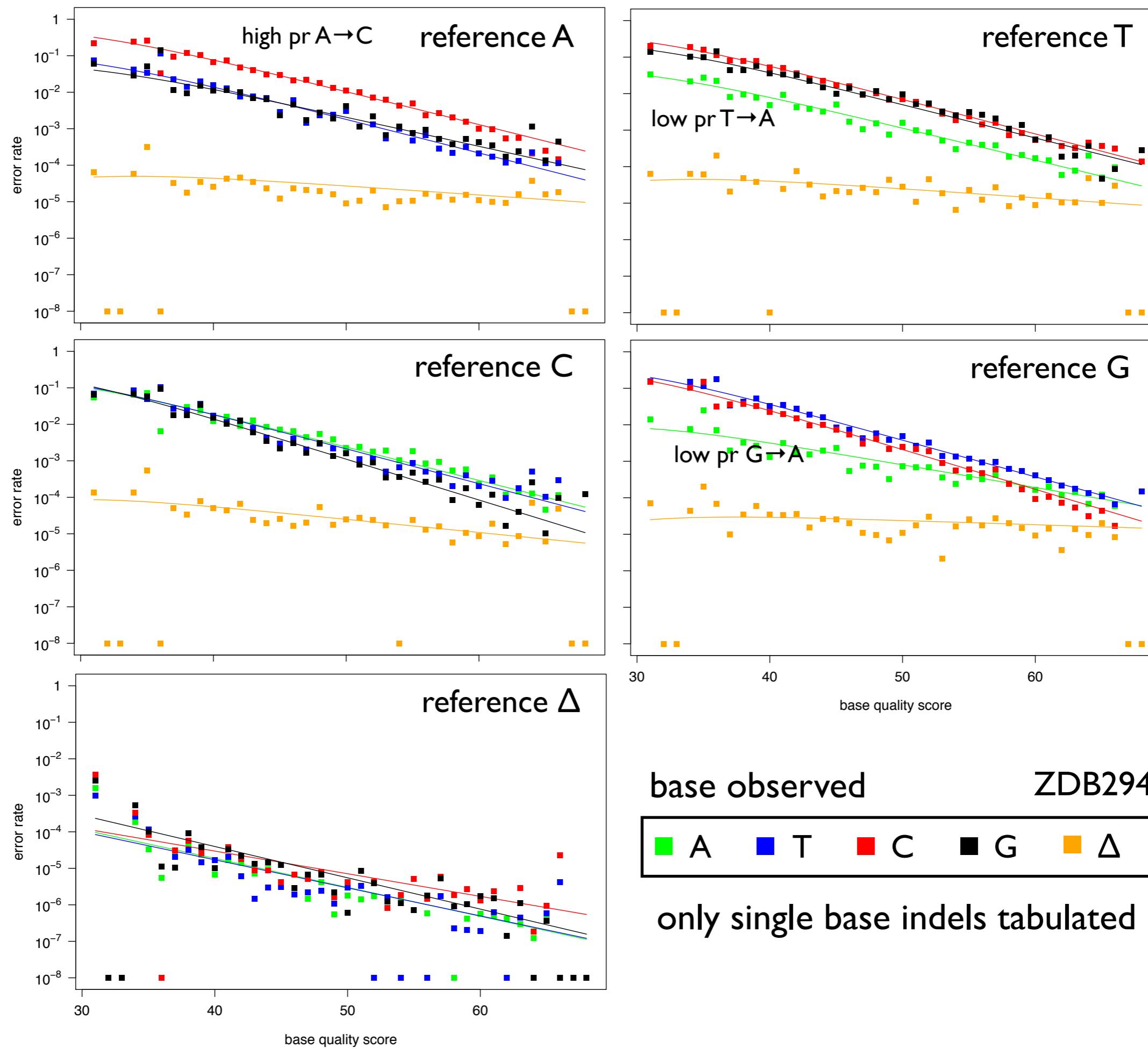
- Most bases in a run have error frequencies between 10^{-4} and 10^{-3} . Overall error rates agree well with Phred quality scores [$E=10^{-(Q/10)}$].



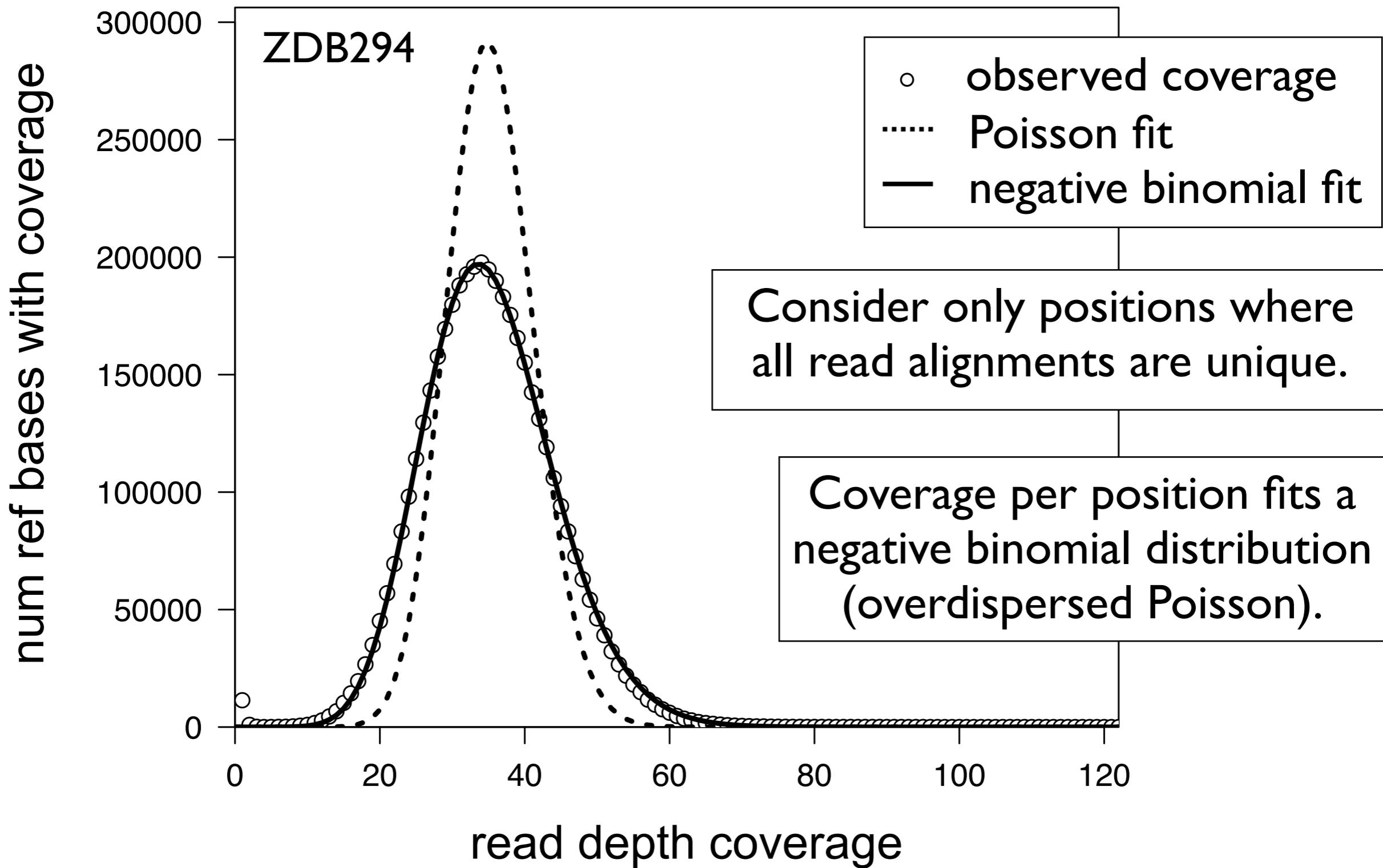
Typical Base Error Spectrum

- There is variation in the frequency at which different base errors occur at a given quality score.

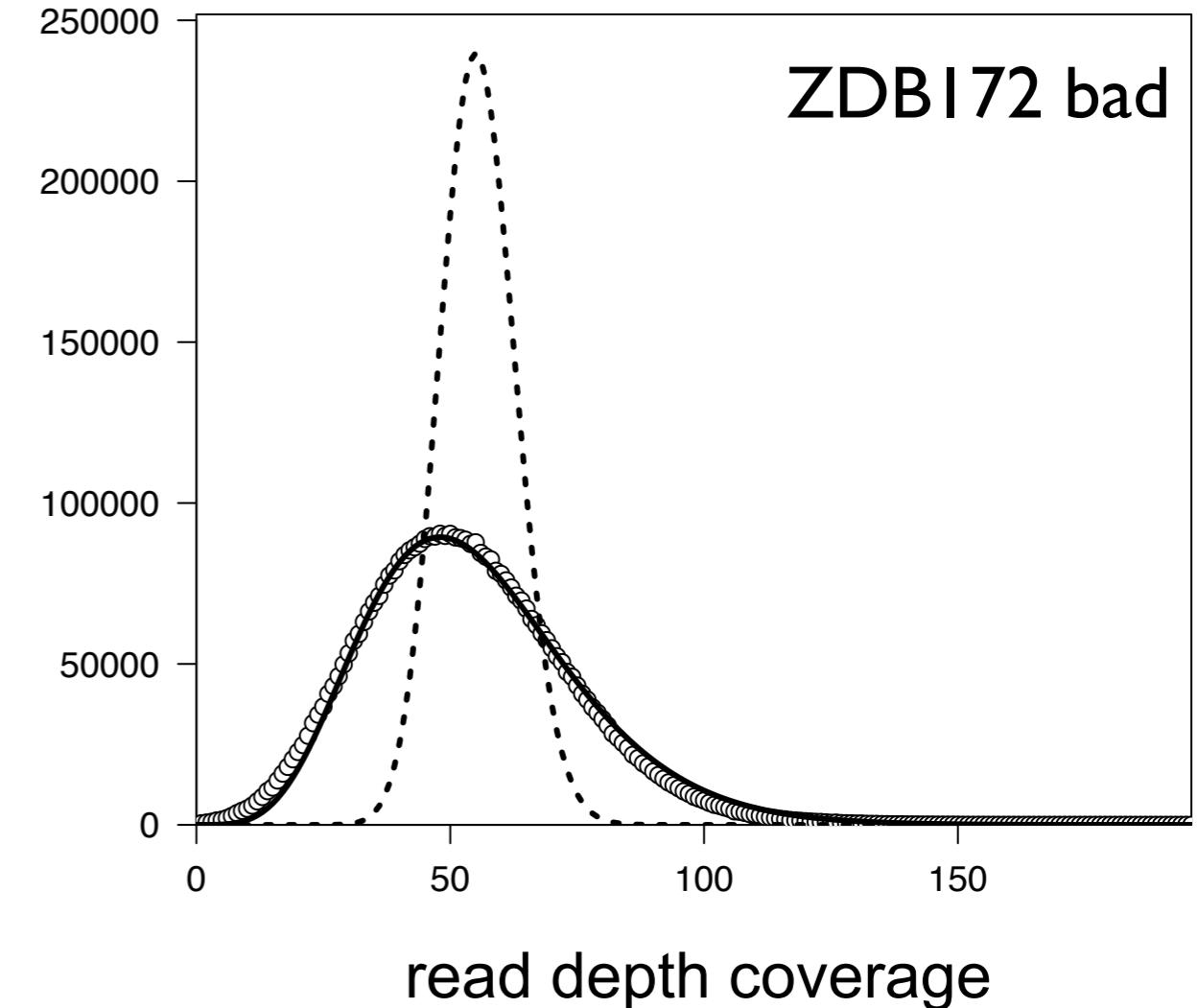
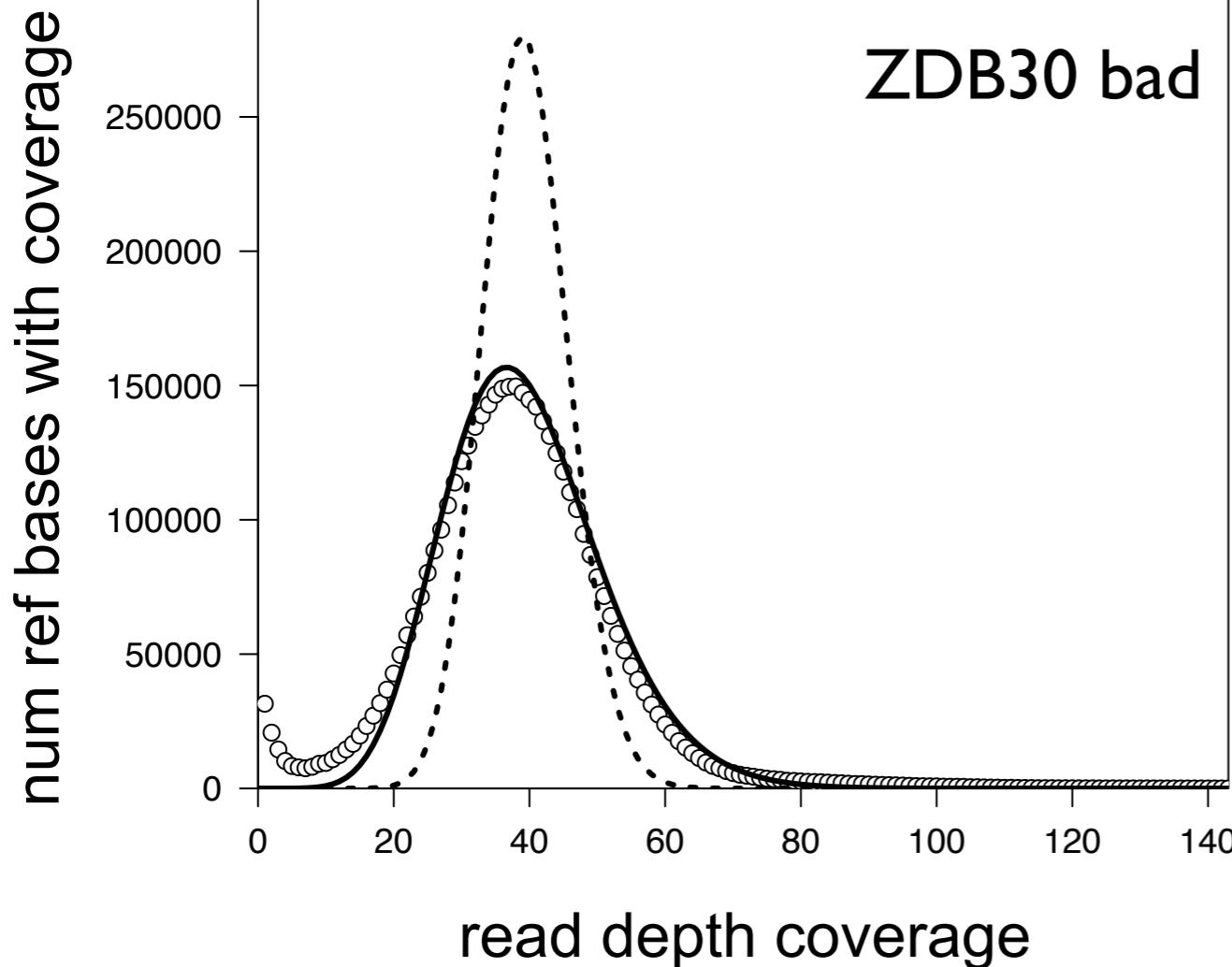




Typical Coverage Distribution



Problem Coverage Distributions



- Contamination with another sample?

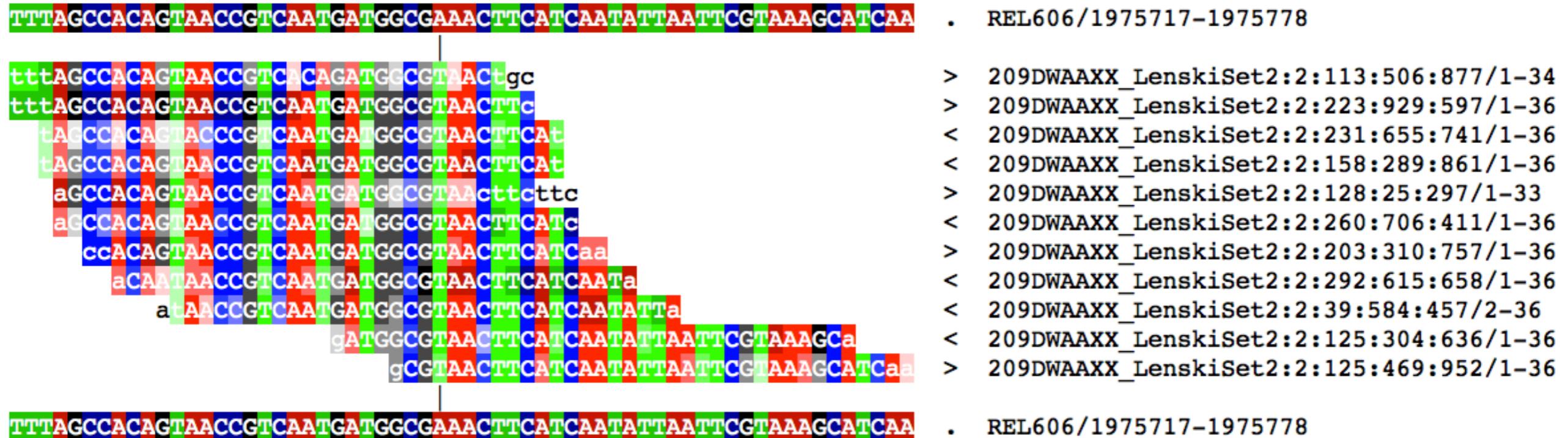
- Large variance, missing coverage.

Both apparently from problems with library prep.

Identifying single-base substitutions

- Calculate probability of best base versus other bases given observed bases and error model.
- Accept as consensus if E-value < 0.01.

start	end	ref	change	quality	cov	tot_cov	type	gene position	codon change	aa change	gene	product
1975746	1975746	A	T	22.6	6/5	6/5	substitution	406 (136)	TCG→ACG	S→T	yedV	predicted sensory kinase in two-component regulatory system with YedW



Legend: ATCG < 36 ≤ ATCG < 43 ≤ ATCG < 57 ≤ ATCG < 71 ≤ ATCG

Identifying within-alignment indels

- Need to be careful in repetitive sequences and at the edges of short reads...

TATATTAAAT**GCGCGCGC**TAGGCTAGCT

TATATTAAAT--**GCGCGC**TAGGCTAGCT <

TATATTAAAT**GCGCGC**--TAGGCTAGCT >

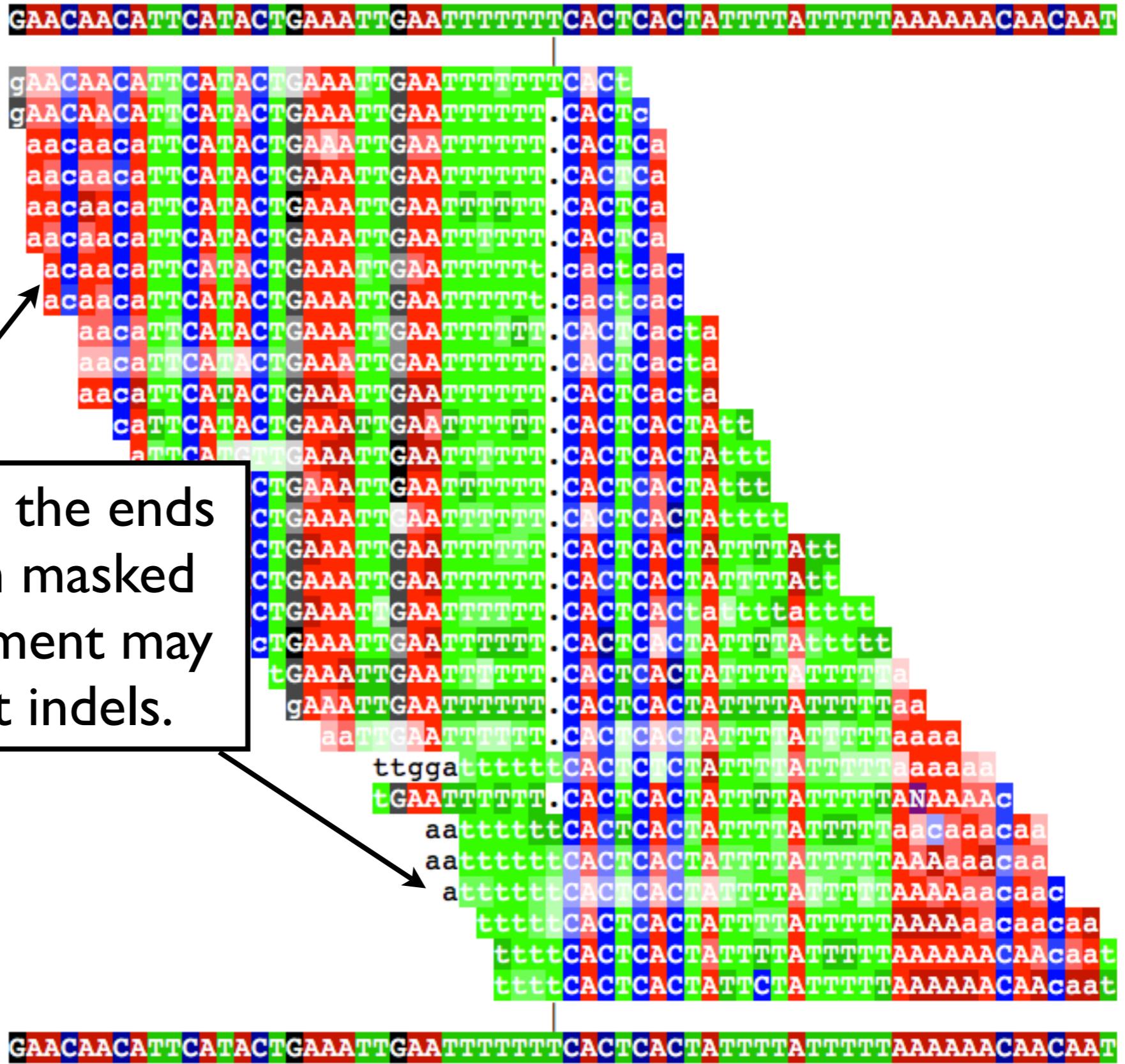
TATATTAAAT**GCGCGC**.....>

.....**GCGCGC**TAGGCTAGCT <

...where reads aligned from different directions can be ambiguously aligned.

...where reads from different directions that end in a simple sequence repeat may hide indels.

Example of a *breseq* prediction



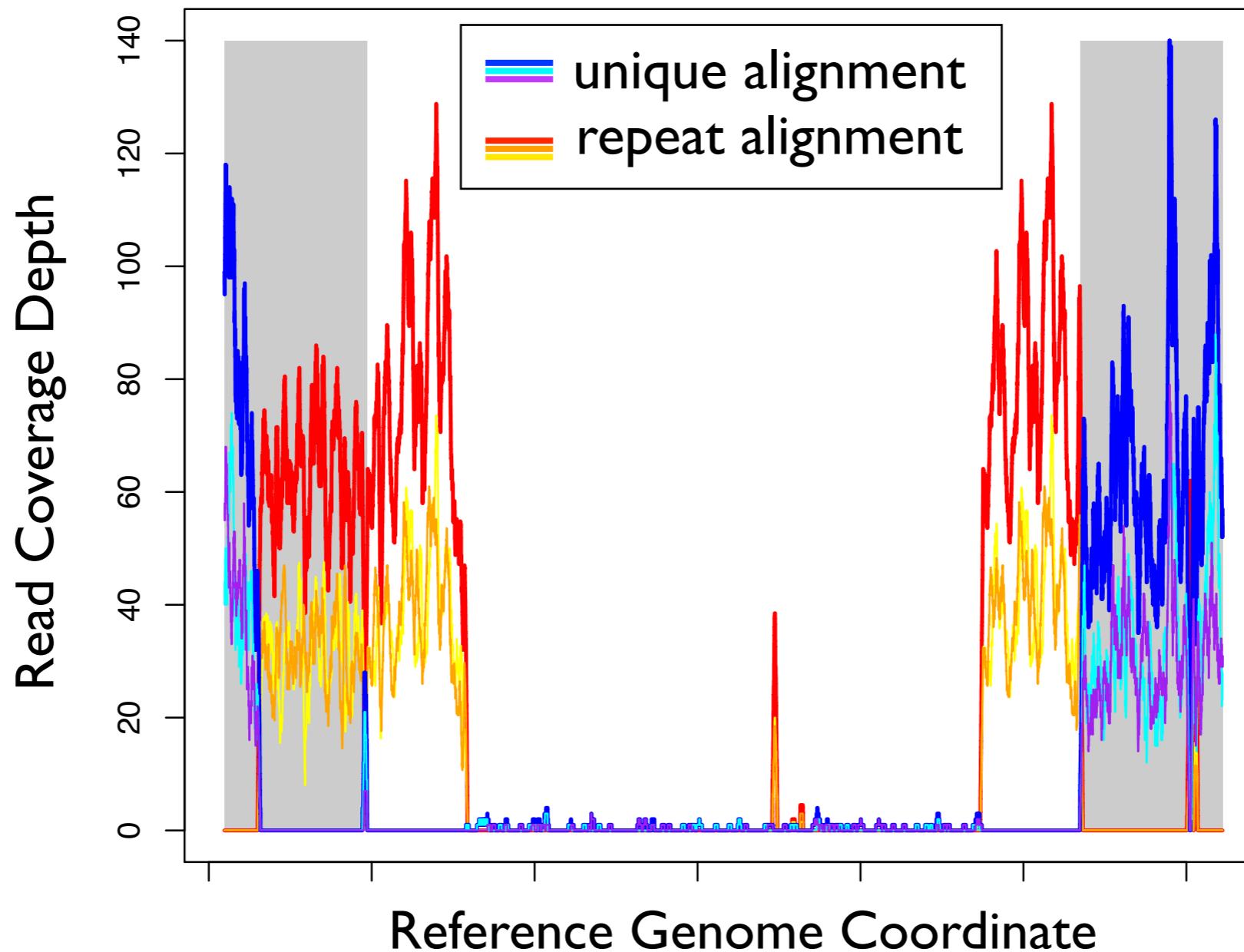
Lowercase bases at the ends of reads have been masked because their alignment may be ambiguous wrt indels.

Identifying large deletions

1. Seed deletions at positions with zero coverage.
2. Propagate boundaries outward until reaching a read-depth threshold based on the overall distribution.
3. Propagate through repeat regions, where a read aligns to multiple places in the genome.

Example of a *breseq* prediction

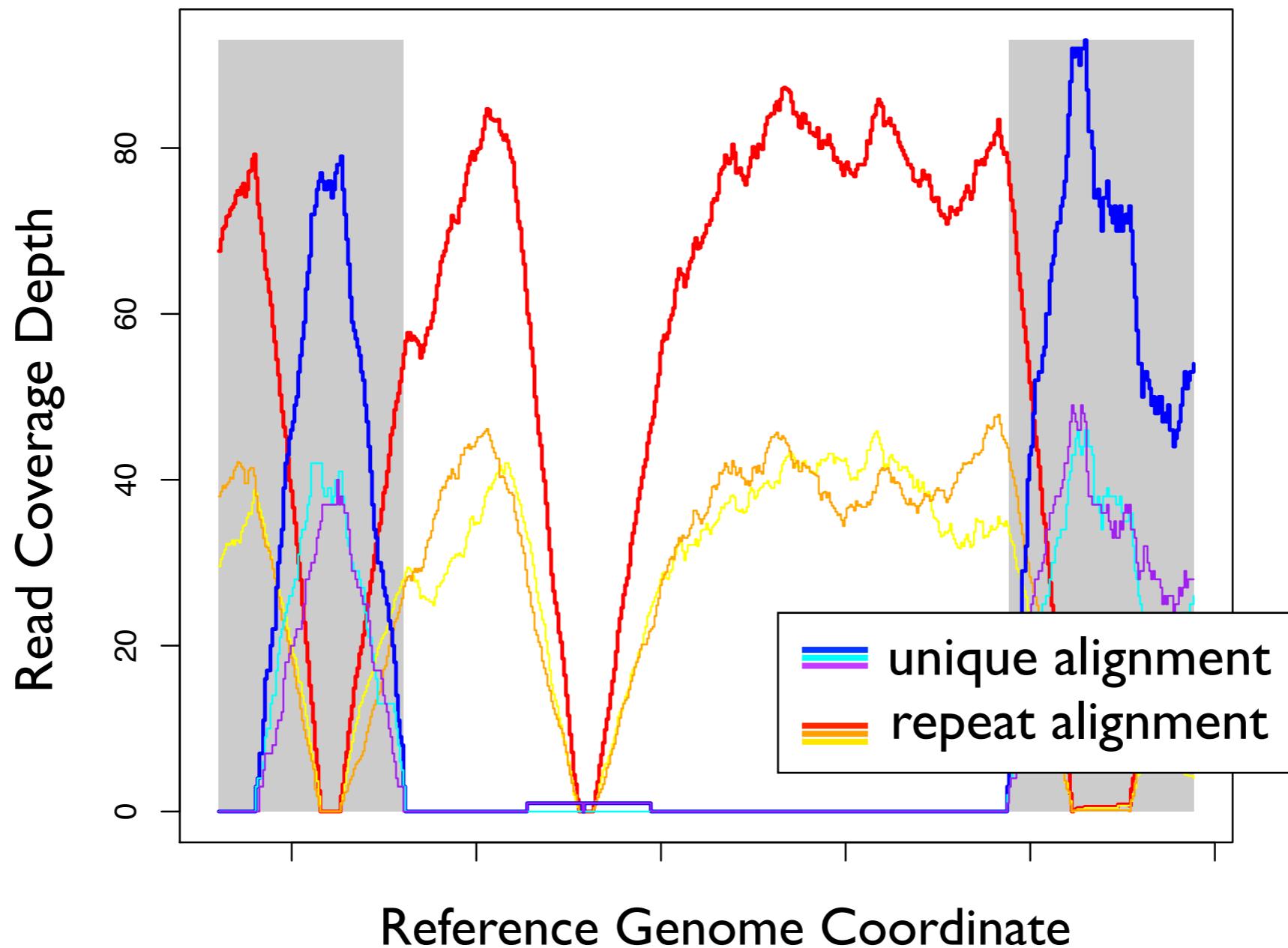
- Sometimes the molecular event is obvious...



- Recombination between nearby IS3 copies.

Example of a *breseq* prediction

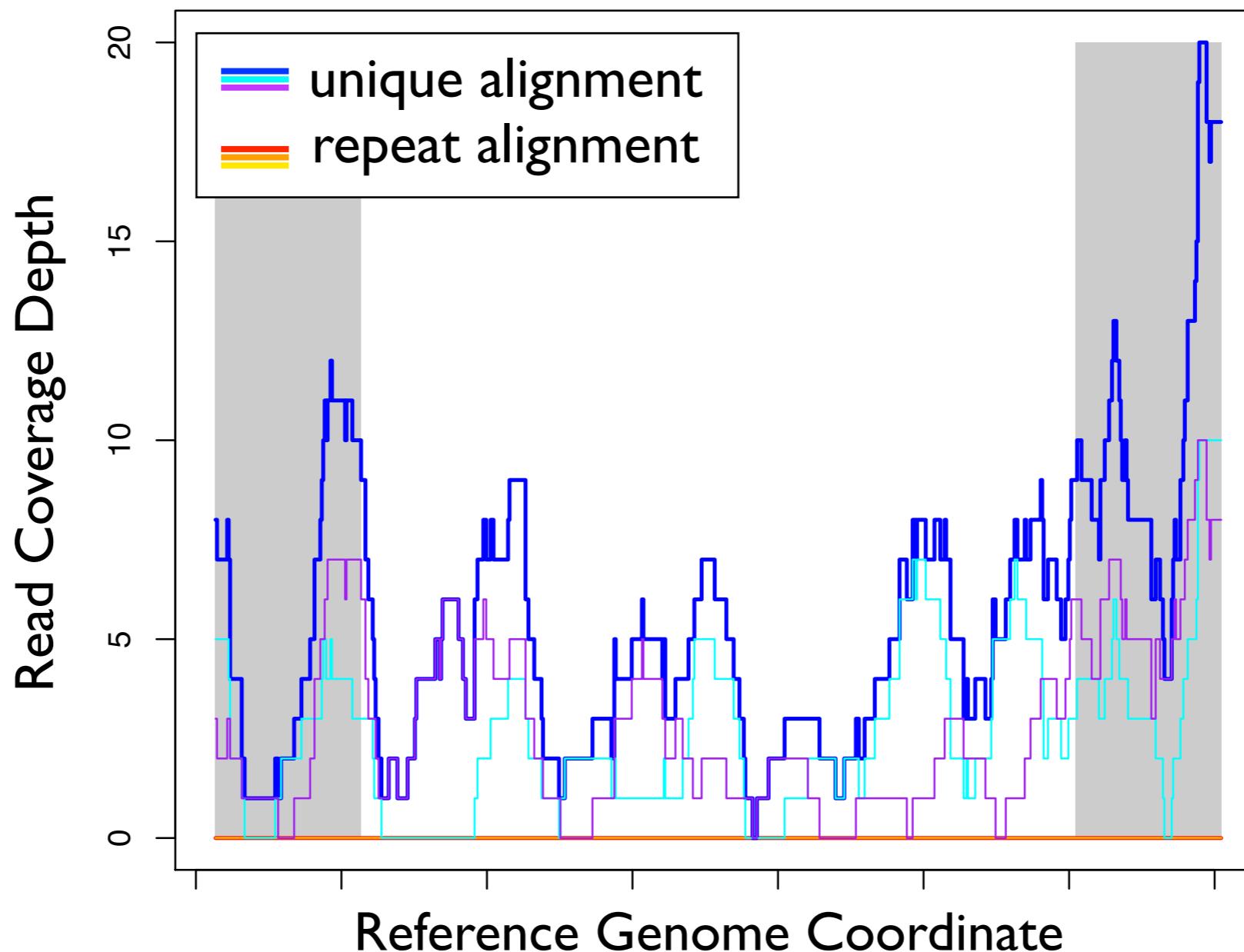
- Sometimes the mutation is not obvious...



- Gene conversion of 23S rRNA copy!!

Example of a *breseq* prediction

- Sometimes overall low or biased coverage leads to false predictions of deletions.



- Recognizable by sloped vs. steep edges.

Identifying new junctions

1. Find “mosaic” reads that partially map to two locations in the genome (possibly with overlap).



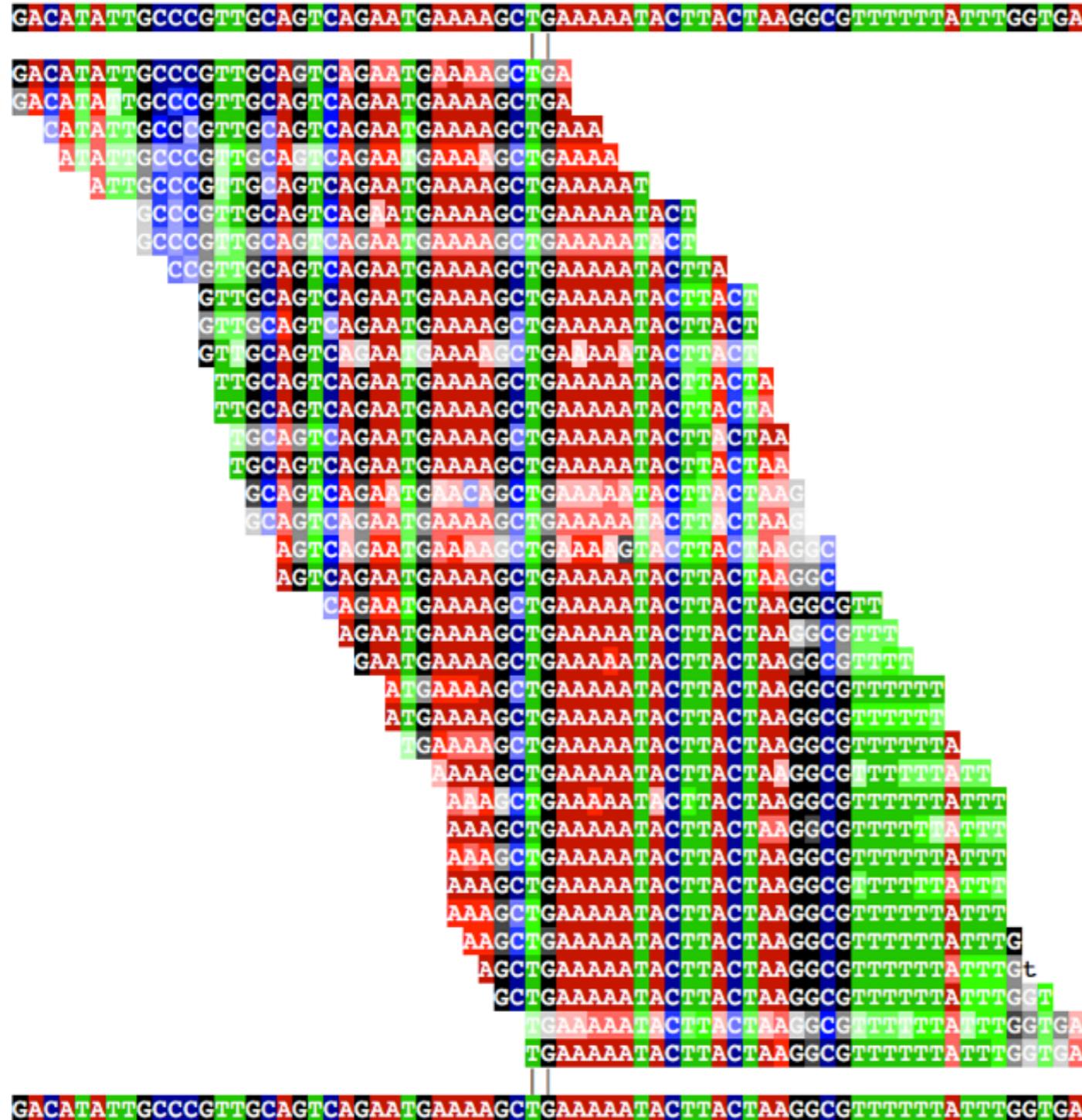
2. Create consensus list of possible new junctions.
3. Re-align all reads to candidate junctions.



4. Predict a new junction if reads map better to it than to the reference across its whole length.

Example of a good junction

position	overlap	reads	gene	coords	product
1 =	0	36	-/thrL	/189	-/thr operon leader peptide
= 4629812			lasT/-	4629789/	predicted rRNA methyltransferase/-



 . REL606_1_1_REL606_4629812_0_0_/3-71

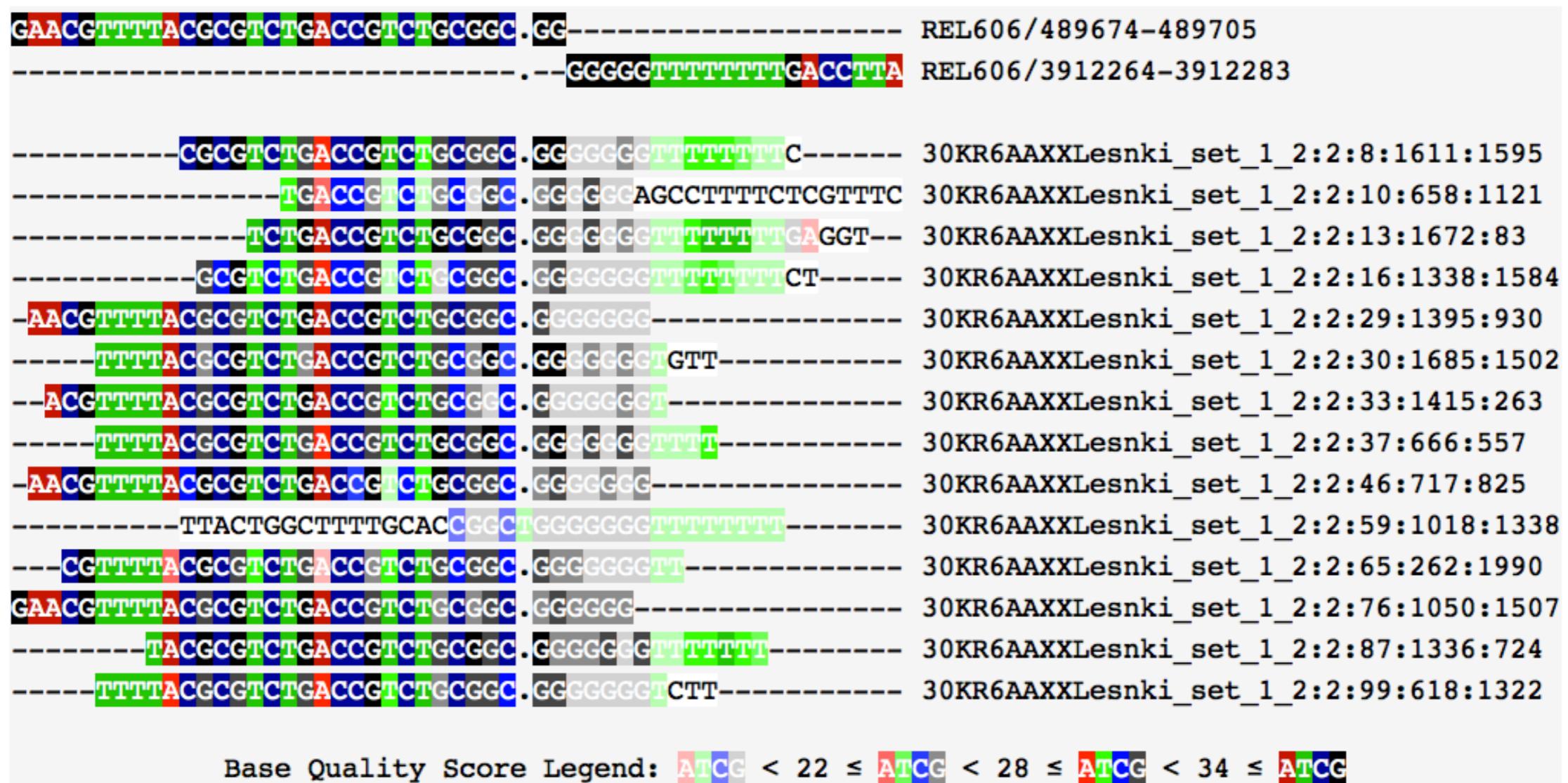
> 30KR6AAXXLesnki_set_1_2:3:53:1076:1729/1-36
 < 30KR6AAXXLesnki_set_1_2:3:3:1045:1537/1-36
 < 30KR6AAXXLesnki_set_1_2:3:98:1256:1982/1-36
 < 30KR6AAXXLesnki_set_1_2:3:69:59:1642/1-36
 < 30KR6AAXXLesnki_set_1_2:3:52:1112:1970/1-36
 < 30KR6AAXXLesnki_set_1_2:3:29:260:647/1-36
 < 30KR6AAXXLesnki_set_1_2:3:45:197:1888/1-36
 < 30KR6AAXXLesnki_set_1_2:3:38:829:160/1-36
 > 30KR6AAXXLesnki_set_1_2:3:82:996:256/1-36
 < 30KR6AAXXLesnki_set_1_2:3:88:199:234/1-36
 > 30KR6AAXXLesnki_set_1_2:3:31:1778:622/1-36
 > 30KR6AAXXLesnki_set_1_2:3:21:1481:579/1-36
 > 30KR6AAXXLesnki_set_1_2:3:14:1273:59/1-36
 < 30KR6AAXXLesnki_set_1_2:3:54:842:43/1-36
 > 30KR6AAXXLesnki_set_1_2:3:82:844:525/1-36
 > 30KR6AAXXLesnki_set_1_2:3:30:6:1419/1-36
 < 30KR6AAXXLesnki_set_1_2:3:23:1578:360/1-36
 > 30KR6AAXXLesnki_set_1_2:3:65:1765:1077/1-36
 > 30KR6AAXXLesnki_set_1_2:3:62:1360:759/1-36
 < 30KR6AAXXLesnki_set_1_2:3:65:842:32/1-36
 > 30KR6AAXXLesnki_set_1_2:3:3:1093:1221/1-36
 > 30KR6AAXXLesnki_set_1_2:3:13:204:1274/1-36
 < 30KR6AAXXLesnki_set_1_2:3:8:699:65/1-36
 > 30KR6AAXXLesnki_set_1_2:3:81:1575:760/1-36
 < 30KR6AAXXLesnki_set_1_2:3:57:387:423/1-36
 > 30KR6AAXXLesnki_set_1_2:3:19:601:1470/1-36
 < 30KR6AAXXLesnki_set_1_2:3:71:503:526/1-36
 > 30KR6AAXXLesnki_set_1_2:3:29:1139:1664/1-36
 < 30KR6AAXXLesnki_set_1_2:3:71:505:527/1-36
 > 30KR6AAXXLesnki_set_1_2:3:18:1079:1002/1-36
 < 30KR6AAXXLesnki_set_1_2:3:6:1485:1308/1-36
 < 30KR6AAXXLesnki_set_1_2:3:24:627:931/1-36
 > 30KR6AAXXLesnki_set_1_2:3:92:145:1544/1-35
 > 30KR6AAXXLesnki_set_1_2:3:58:1720:1463/1-36
 < 30KR6AAXXLesnki_set_1_2:3:86:300:312/1-36
 > 30KR6AAXXLesnki_set_1_2:3:41:1600:1707/1-36

. REL606_1_1_REL606_4629812_0_0_/3-71

Example of a bad junction

- Beware of reads ending in homopolymer runs!

position	overlap	full / total reads	gene	coords	product
= 489705	0	7 / 14	<i>ybbN</i>	490447-489593	predicted thioredoxin domain-containing protein
3912264 =			<i>ilvL/ilvG</i>	3912221/3912359	ilvG operon leader peptide/acetolactate synthase II, valine insensitive, large subunit

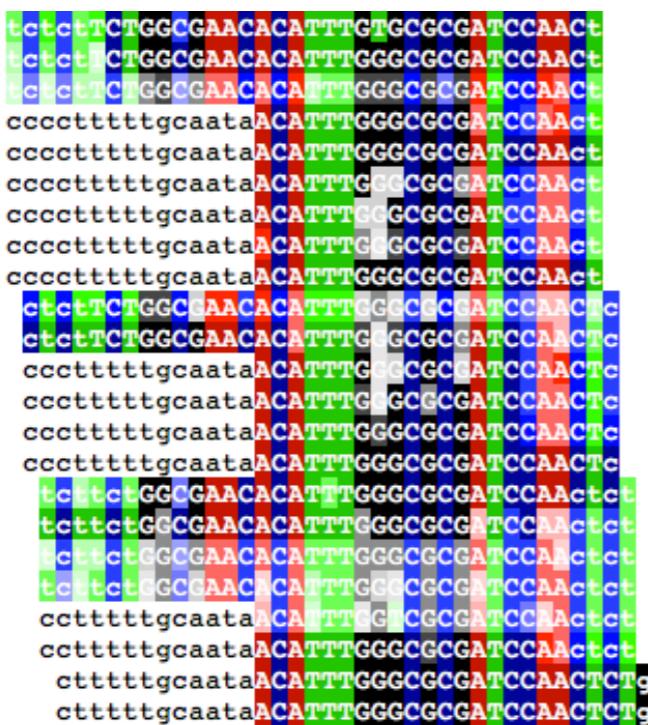


Example of a *breseq* prediction

- IS insertions create two new junctions...

	position	overlap	reads	gene	coords	product
	16989			IS150 (+)	+1443 (+3) bp	
* ?	16990 =	0	44	<i>mokC/nhaA</i>	16959/17487	regulatory protein for HokC, overlaps CDS of hokC/pH-dependent sodium/proton antiporter
? = 3652533				<i>IS150</i>	3651091-3652533	repeat region
* ?	= 16992	0	41	<i>mokC/nhaA</i>	16959/17487	regulatory protein for HokC, overlaps CDS of hokC/pH-dependent sodium/proton antiporter
? 3893554 =				<i>IS150</i>	3893554-3894996	repeat region

- Sometimes both new and old junctions exist...

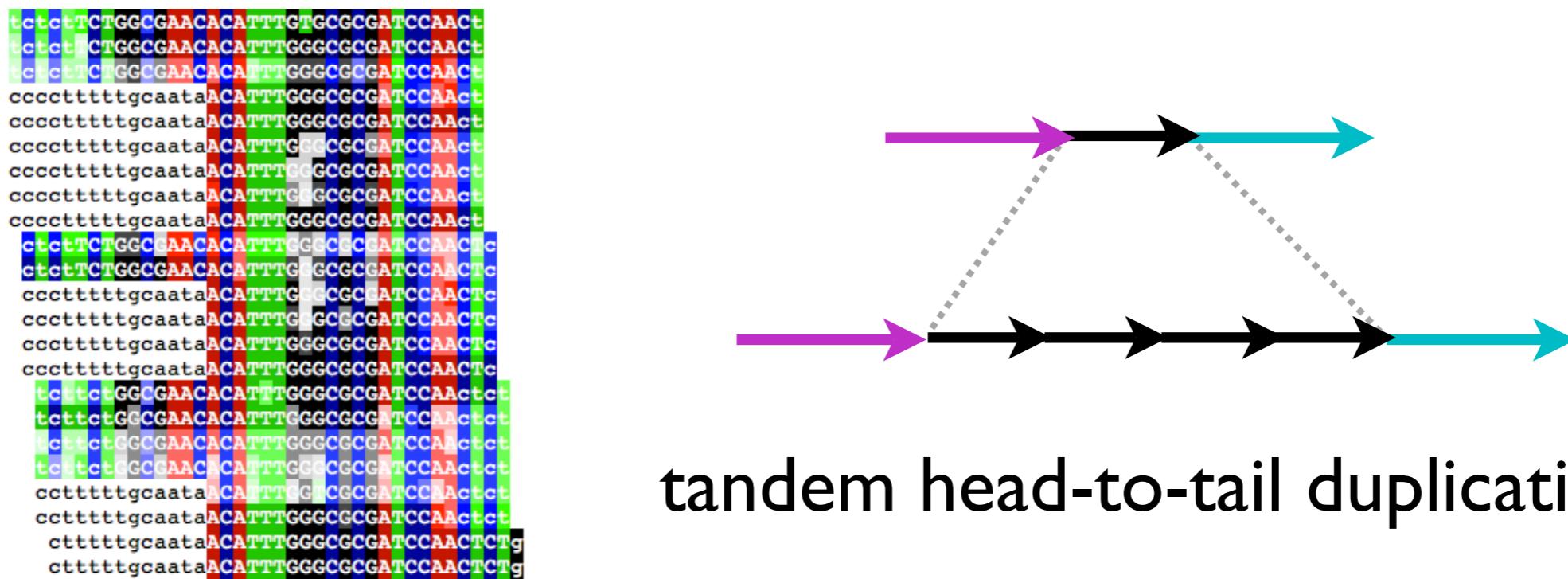


Example of a *breseq* prediction

- IS insertions create two new junctions...

	position	overlap	reads	gene	coords	product
	16989			IS150 (+)	+1443 (+3) bp	
* ?	16990 =	0	44	<i>mokC/nhaA</i>	16959/17487	regulatory protein for HokC, overlaps CDS of hokC/pH-dependent sodium/proton antiporter
? = 3652533				<i>IS150</i>	3651091-3652533	repeat region
* ?	= 16992	0	41	<i>mokC/nhaA</i>	16959/17487	regulatory protein for HokC, overlaps CDS of hokC/pH-dependent sodium/proton antiporter
? 3893554 =				<i>IS150</i>	3893554-3894996	repeat region

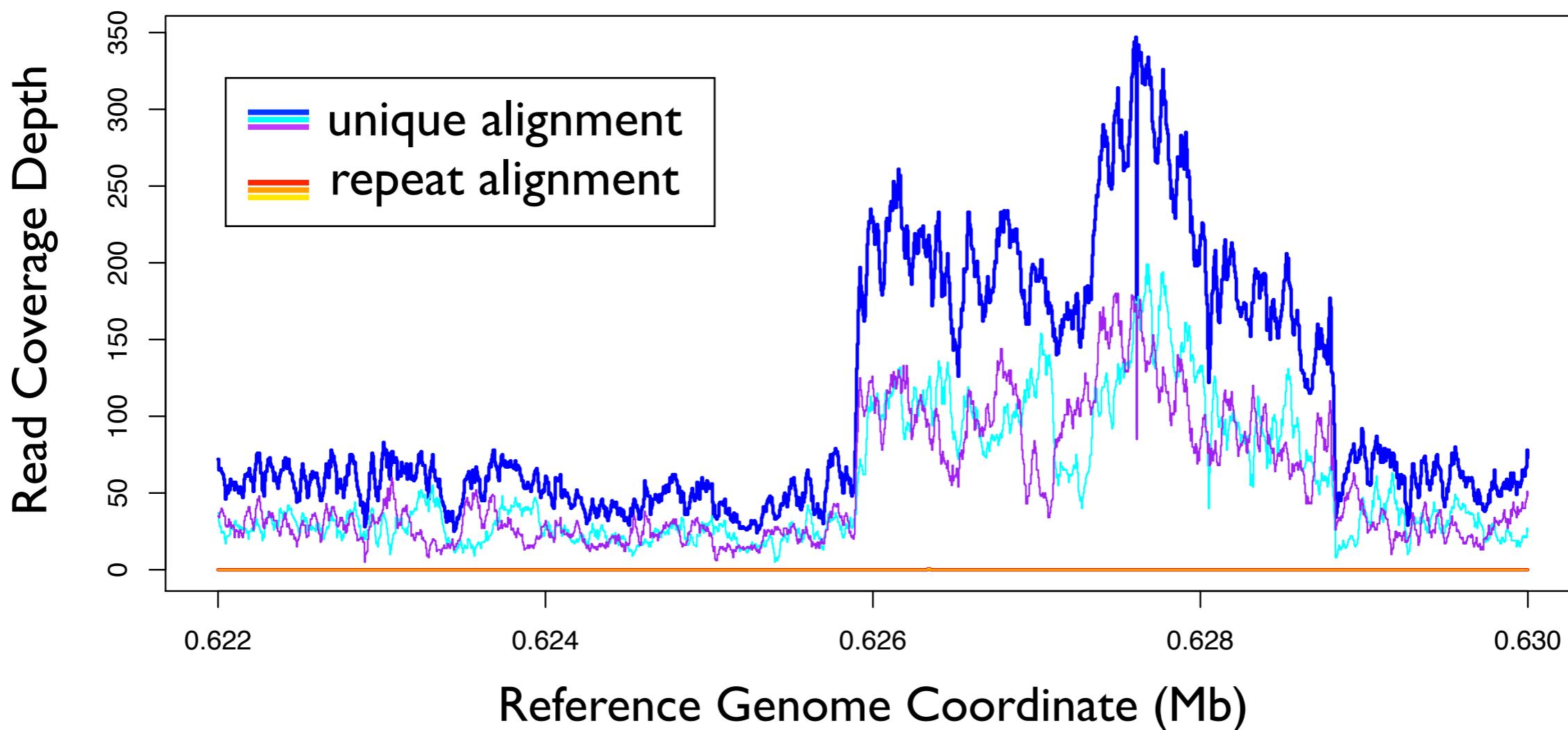
- Sometimes both new and old junctions exist...



tandem head-to-tail duplications

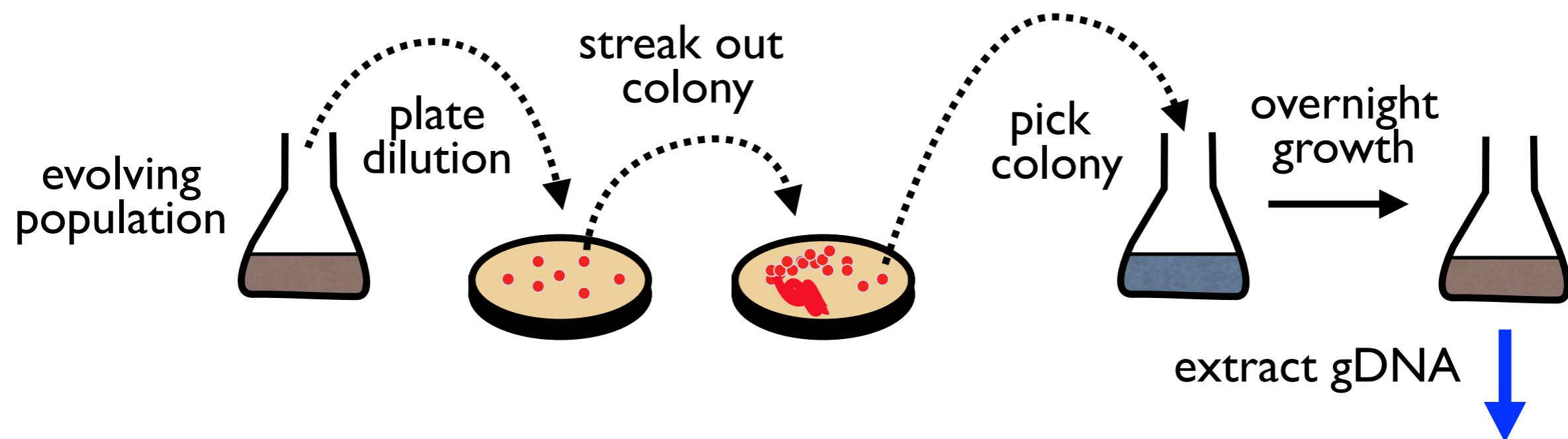
Identifying copy number variation

- Coverage is very noisy, but a fingerprint is (somewhat) consistent across runs.



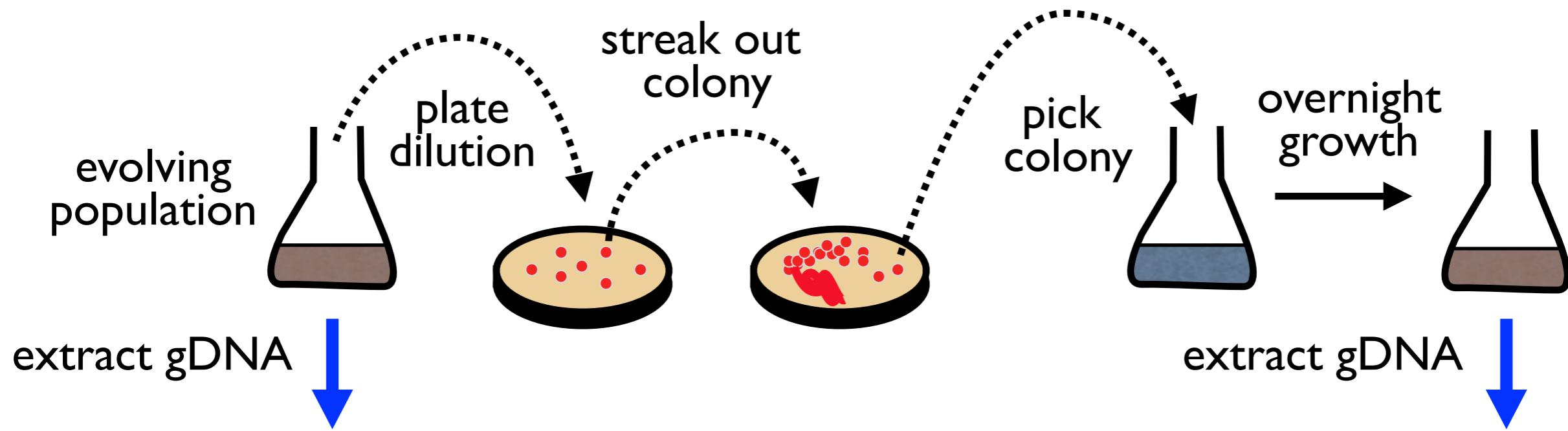
- Tile into 100 bp segments, train bg model on many genomes, look for deviation (in progress).

Mixed population analysis



All reads are from a
single clone.

Mixed population analysis



Every read could be
from any individual.

Frequencies of mutations
competing in population.

No linkage information.

All reads are from a
single clone.

Information about which
mutations occur together.

Sequencing error or polymorphism?

Ref

TAG

↓

TCG

TAG

TAG

TAG

TCG

TCG

TCG

TAG

TCG

TCG

TCG

TAG

TAG

Aligned reads

- Map reads to ancestor genome.
Only consider single-base substitutions.

Sequencing error or polymorphism?

Ref

TAG



TCG

TAG

TAG

TAG

TCG

TCG

TCG

TAG

TCG

TCG

TCG

TAG

TAG

Aligned reads

- Map reads to ancestor genome.
Only consider single-base substitutions.

- Log-likelihood test for polymorphism:

$$D = -2 \ln \frac{\Pr(\text{obs} \mid \text{no polymorphism, i.e. all error})}{\Pr(\text{obs} \mid \text{ML fraction new allele})}$$

Sequencing error or polymorphism?

Ref

TAG



TCG

TAG

TAG

TAG

TCG

TCG

TCG

TAG

TCG

TCG

TCG

TAG

TAG

Aligned reads

- Map reads to ancestor genome.
Only consider single-base substitutions.

- Log-likelihood test for polymorphism:

$$D = -2 \ln \frac{\Pr(\text{obs} \mid \text{no polymorphism, i.e. all error})}{\Pr(\text{obs} \mid \text{ML fraction new allele})}$$

- Clone sequence data serves as a negative control (all errors, no polymorphisms).

Sequencing error or polymorphism?

Ref

TAG



TCG >

TAG >

TAG <

TAG >

- Map reads to ancestor genome.
 - Only consider single-base substitutions.
- Log-likelihood test for polymorphism:
- $$D = -2 \ln \frac{\Pr(\text{obs} \mid \text{no polymorphism, i.e. all error})}{\Pr(\text{obs} \mid \text{ML fraction new allele})}$$

Aligned reads

TCG >

TCG >

TAG <

TCG >

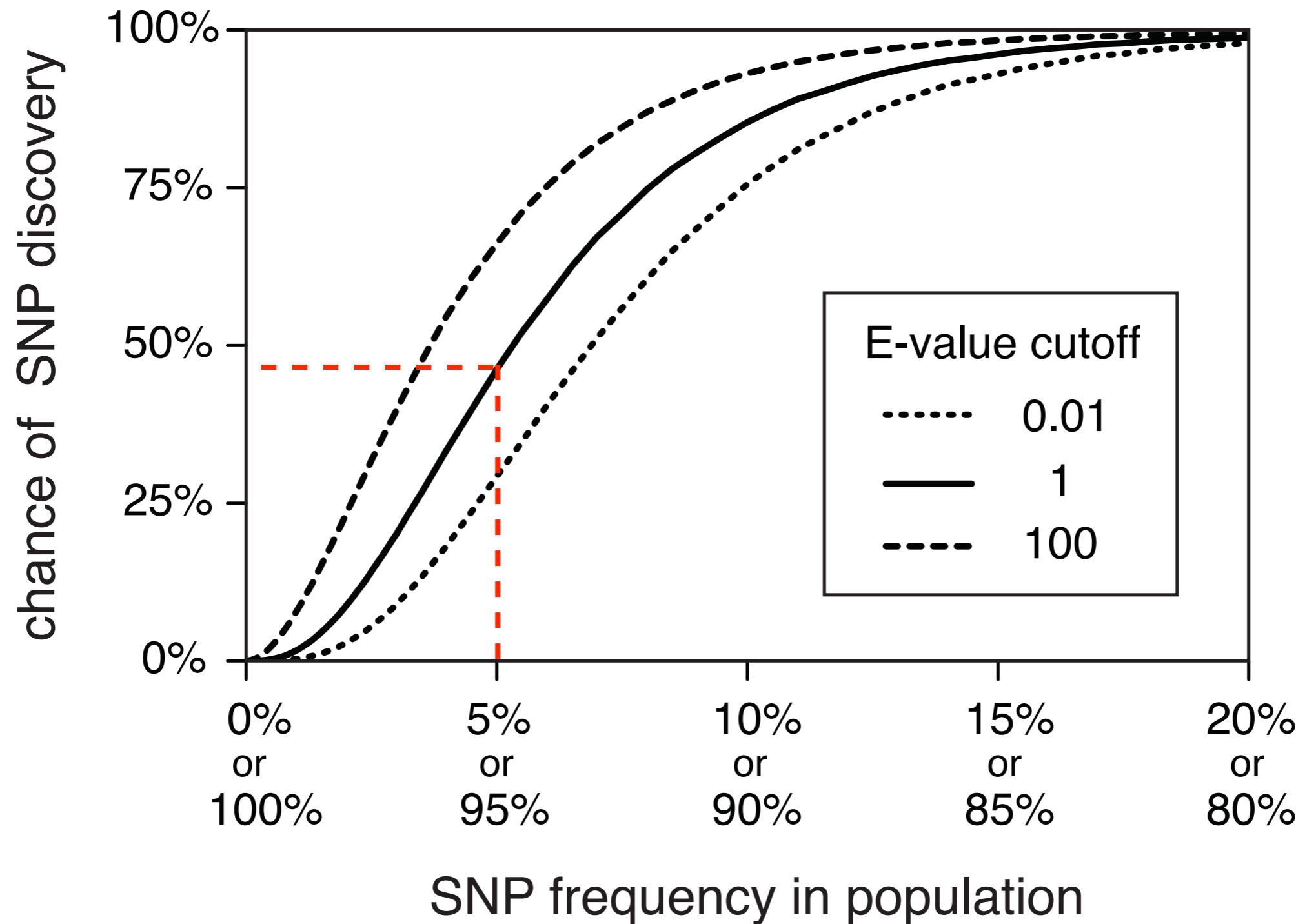
TCG >

TAG >

TAG <

- Clone sequence data serves as a negative control (all errors, no polymorphisms).
- Filter out predictions with other biases: strand bias, systematically low quality scores

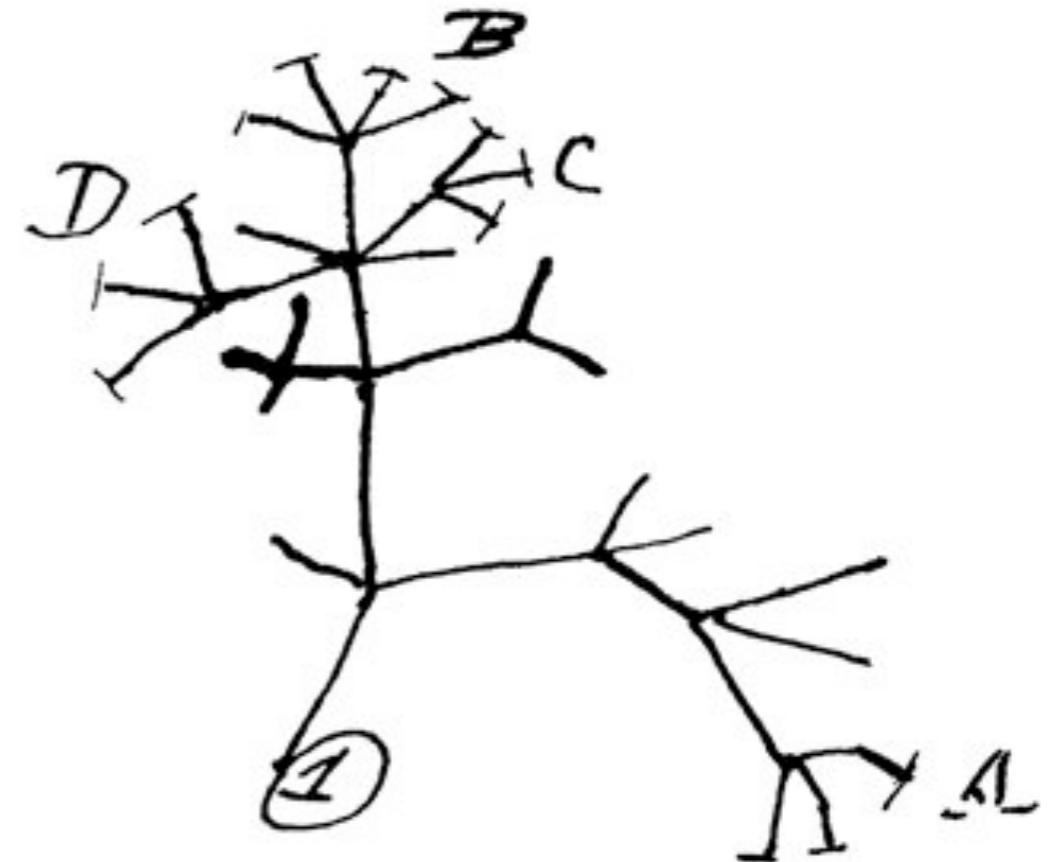
Sensitivity estimate



From evidence to mutations

- Genome sequencing data gives us evidence of changes: read alignments, missing coverage, new junctions, ...
- But we really want a list of biological mutations to study evolutionary history and molecular mechanisms.
- Complication: Later events may sometimes hide earlier events (e.g. SNP in region that is later deleted)

I think



Genome Diffs

- To submit a changed genome sequence to GenBank you must currently re-submit the entire genome – even *if it has only a one base difference*.
- Mutational events are essentially *genome differences*. (In a CS sense of applying “patches” to files)
- Supplementary tables are not a sustainable, standardized, or re-usable way to report this data.
- An ideal genome analysis also reports what is not known, frequency information for mixed population samples, quality metrics, ...

Genome Diff

- Include both evidence and mutations in a file.
- Evidence entries can support mutations.
- Can also record experiments (PCR, targeted resequencing, PFGE) that accumulate more evidence.
- It may require human intervention to interpret evidence. Thus, this file should be human-editable. (This is more fun than the Sunday crossword.)
- **breseq** generates and uses Genome Diffs, so you'll see them in action in the tutorials.

Where are things going

- Re-sequencing will be used to routinely check mutant constructs. (\$1000 human = \$1 *E. coli*)
- *De novo* assembly will become more common, as technologies with longer read lengths come online.
- Studies of within-host diversity of virus populations and genetic diversity of neoplastic tumors.
- Every strain in the Lenski freezer will be sequenced (~11,000 to go)...

Please let me <jbarrick@msu.edu> know if...

- ...you have any questions about **breseq** internals.
- ...you have any ideas for better analysis strategies.
- ...you find **breseq** useful or not on your data.

For more example data you can...

Search “Lenski” in the Short Read Archive (SRA)
<http://www.ncbi.nlm.nih.gov/sra>

For more information...

1. Barrick, J.E. et al. Genome evolution and adaptation in a long-term experiment with *Escherichia coli*. *Nature* **461**, 1243-1247 (2009).
2. Barrick, J.E. & Lenski, R.E. Genome-wide mutational diversity in an evolving population of *Escherichia coli*. *Cold Spring Harbor Symp. Quant. Biol.* **74**, ePub Sept. 23, 2009 (2009).