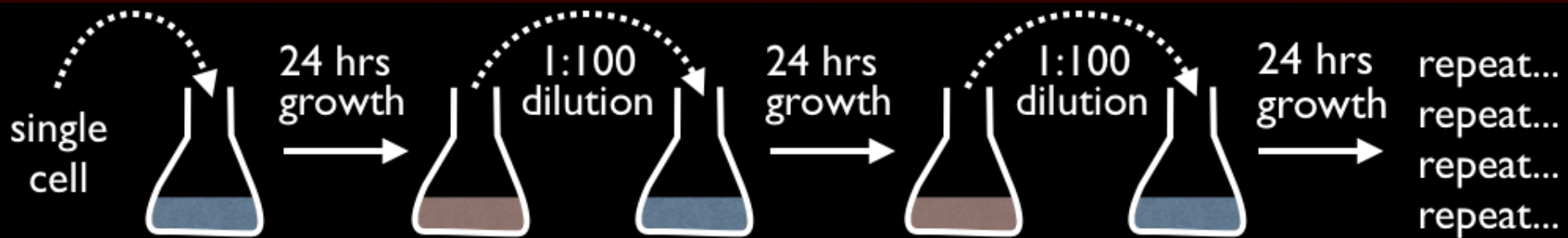# Advanced Variant Calling

# Types of Genome Sequence Variants

1. **Single Nucleotide Variants (SNVs)** ∗
   - Single base changes, e.g., A→T.

2. **Insertions-Deletions (Indels; DIPs)** ∗▶
   - Consisting of one or a few bases, e.g., +ATGA, ΔT.

3. **Structural Variants (SVs)** ▶
   - Everything else: large deletions, insertions, duplications, inversions, translocations, mobile element insertions, horizontal gene transfer

*Different sequencing information and different algorithms are used to predict each kind of variant.*

# Long-term *E. coli* evolution experiment

single cell → 24 hrs growth → 1:100 dilution → 24 hrs growth → 1:100 dilution → 24 hrs growth → repeat... repeat... repeat... repeat...

❖ 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 Frozen "fossil record" has been archived.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.
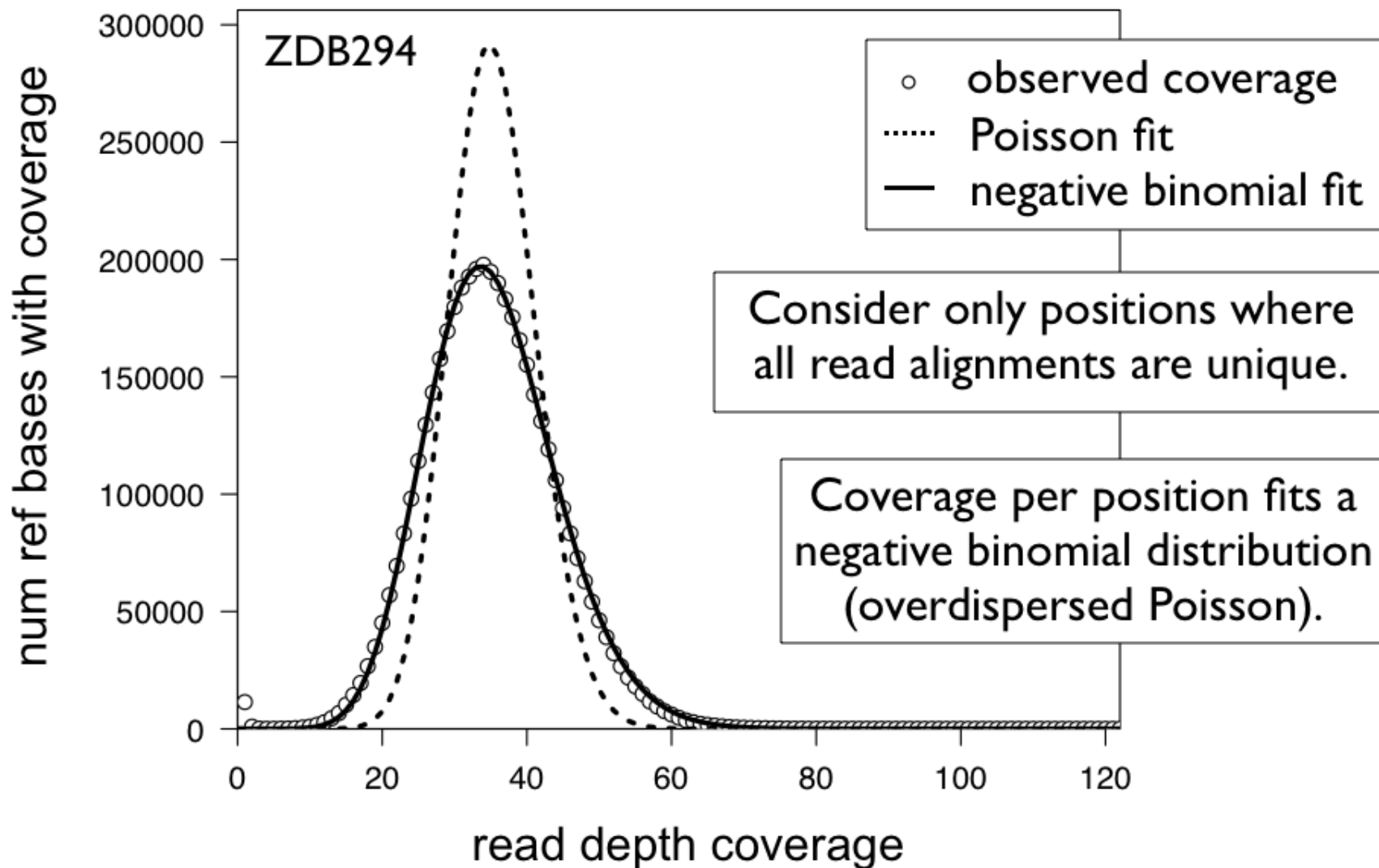
# Knowing what you don't know

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

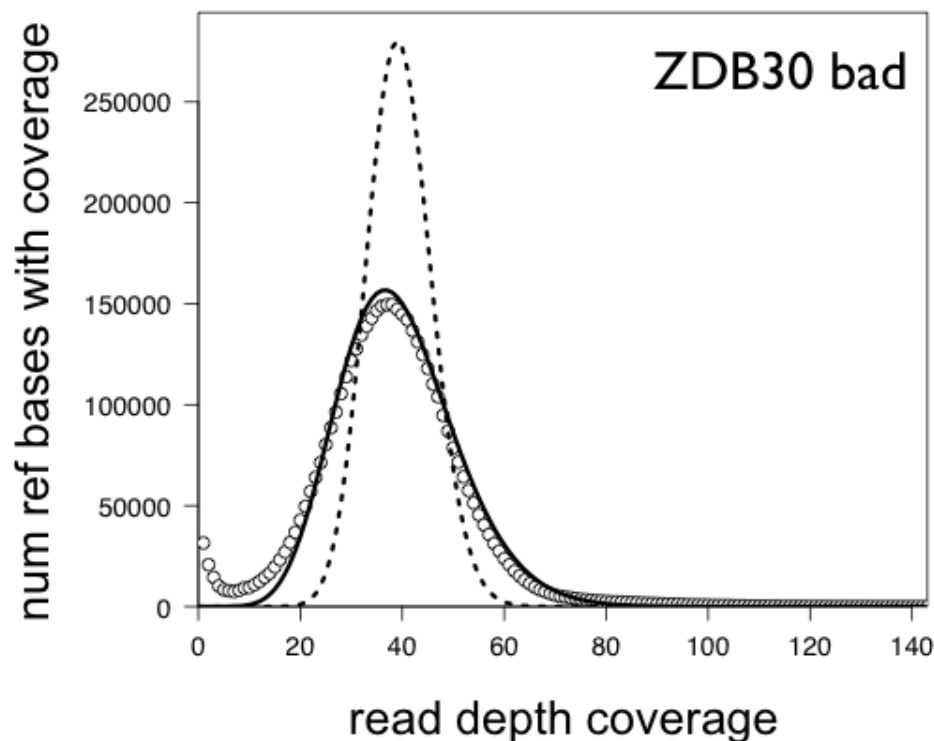|  | single-end | paired-end | mate-paired |
|---|:---:|:---:|:---:|
| IS insertions | * | * | * |
| duplications | * | * | * |
| inversions across IS | – | – | * |
| SNPs in repeats | – | – | * |
| insertion of new seq | – | – | – |

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.*
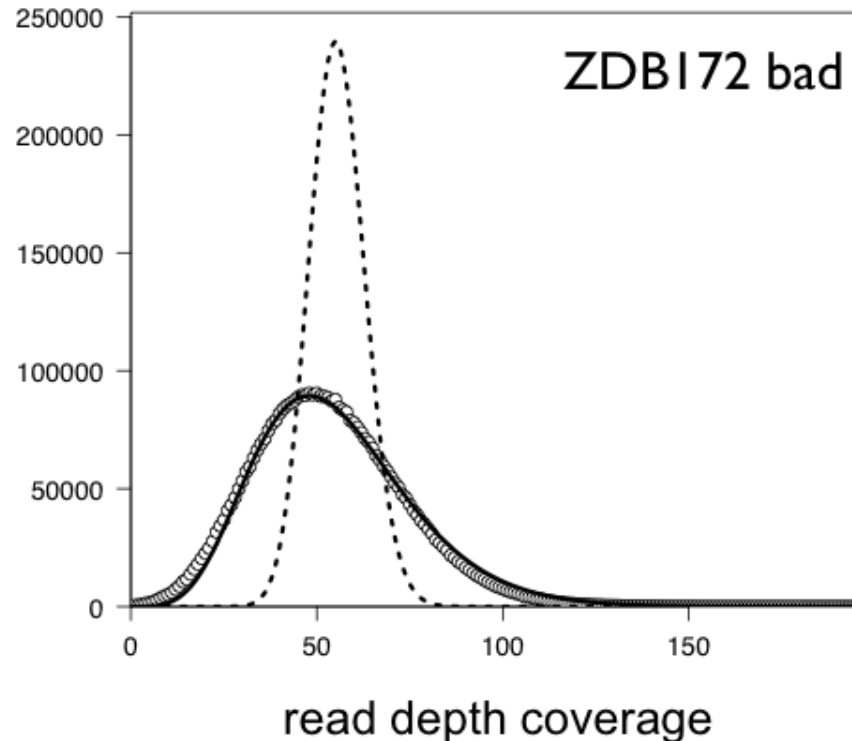
# Typical Coverage Distribution

# Problem Coverage Distributions



- Contamination with another sample?

- Large variance, missing coverage.

Both apparently from problems with library prep.

# Identifying within-alignment indels

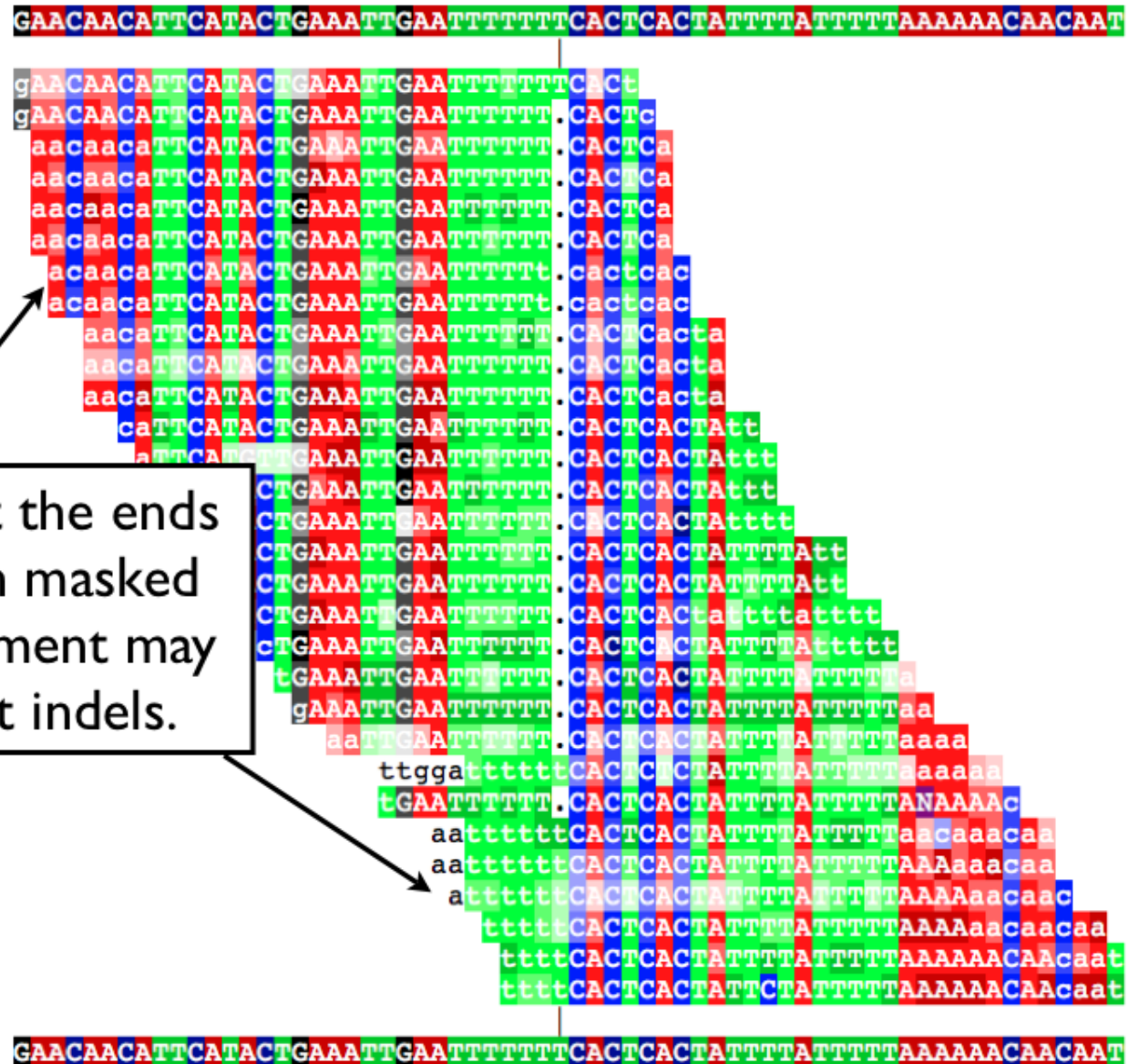- Need to be careful in repetitive sequences and at the edges of short reads...

```
TATATTAATGCGCGCGCTAGGCTAGCT
TATATTAAT--GCGCGCTAGGCTAGCT <
TATATTAATGCGCGC--TAGGCTAGCT >
TATATTAATGCGCGC............ >
...........GCGCGCTAGGCTAGCT <
```

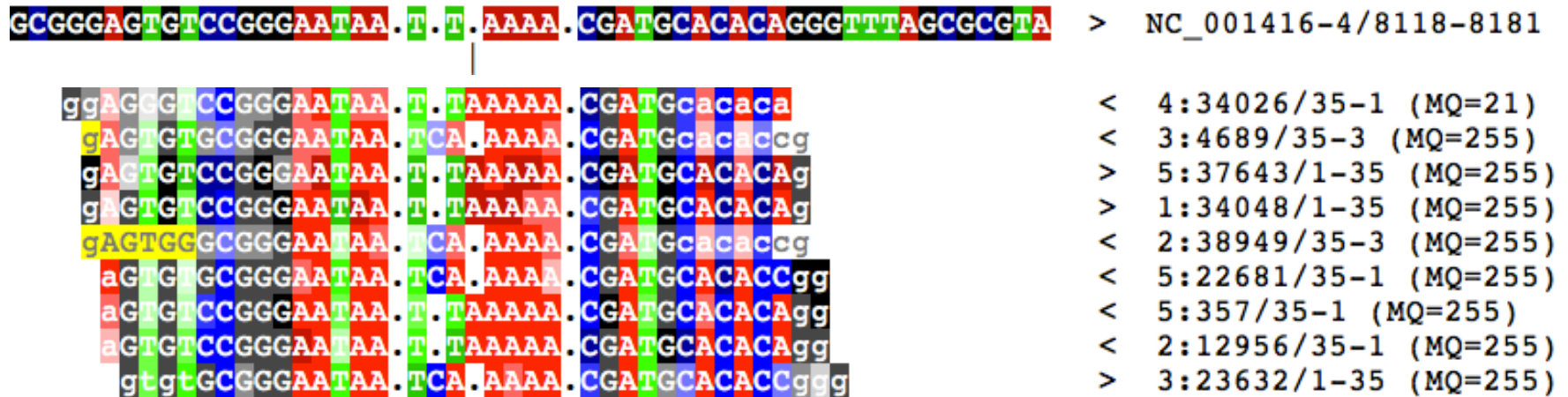...where reads aligned from different directions can be ambiguously aligned.

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

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

# Pitfalls of the column mindset



Requires local multiple sequence re-alignment to get it right!

Implemented in samtools mpileup and the
Genome Alignment Toolkit (GATK).
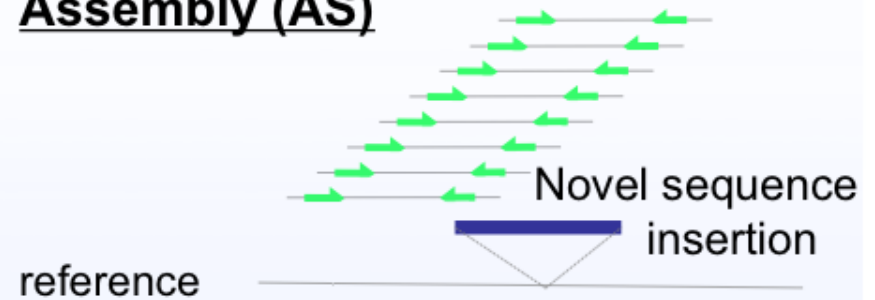
# Predicting structural variants



**Read Pairs (RP)**

No SV | Deletion | Mobile element (MEI) insertion | Tandem duplication

MEI

sample

reference

**Read Depth (RD)**

sample reads

Duplication

Deletion

reference

Unfortunately there is no program or pipeline that does **all** these things!!

**Split Reads (SR)**

Deletion

reference

**Assembly (AS)**

Novel sequence insertion

reference

(http://www.genome.gov/Pages/Research/DER/1000GenomesProjectTutorials/StructuralVariants-JanKorbel.pdf)
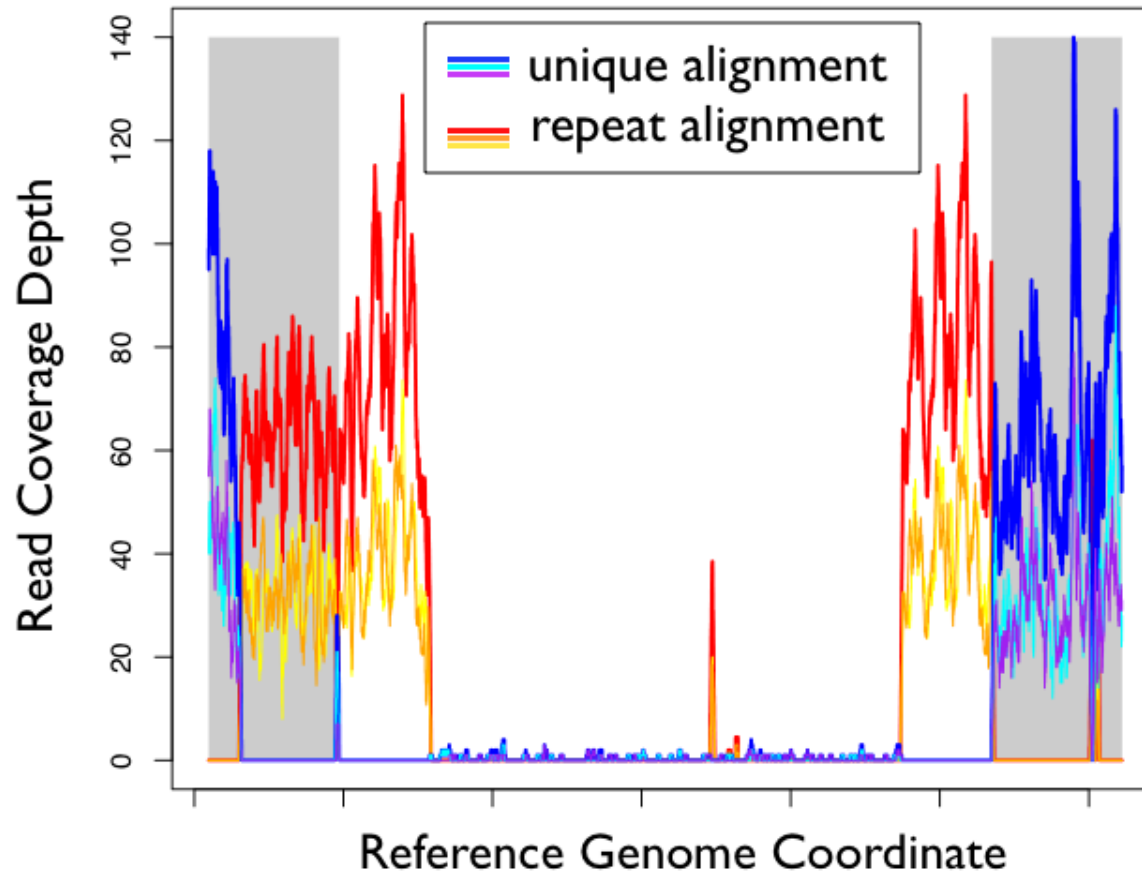
# Identifying large deletions

1. Seed deletions at positions with zero coverage.

2. Propagate boundaries outward until reaching a read-depth threshold based on the overall distribution.

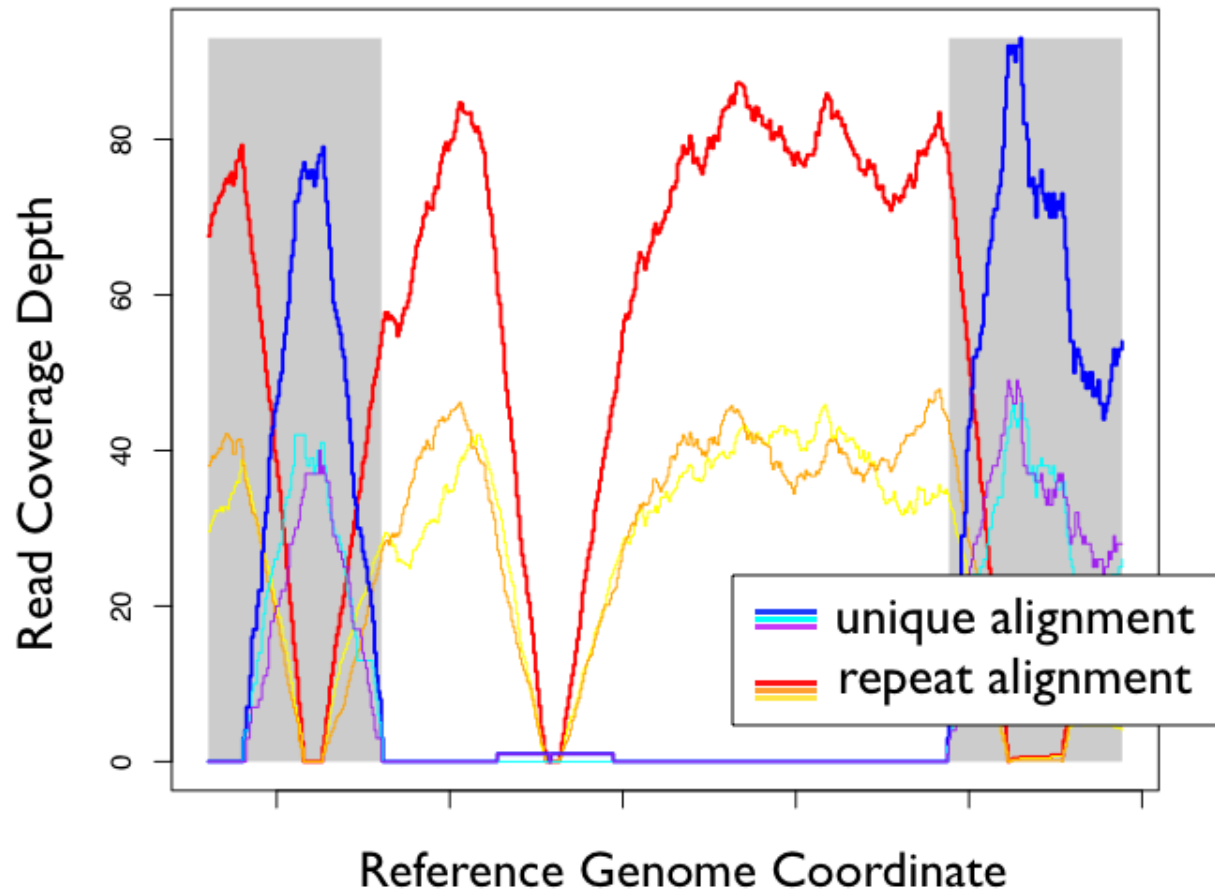3. Propagate through repeat regions, where a read aligns to multiple places in the genome.

- Sometimes the molecular event is obvious...



- Recombination between nearby IS3 copies.

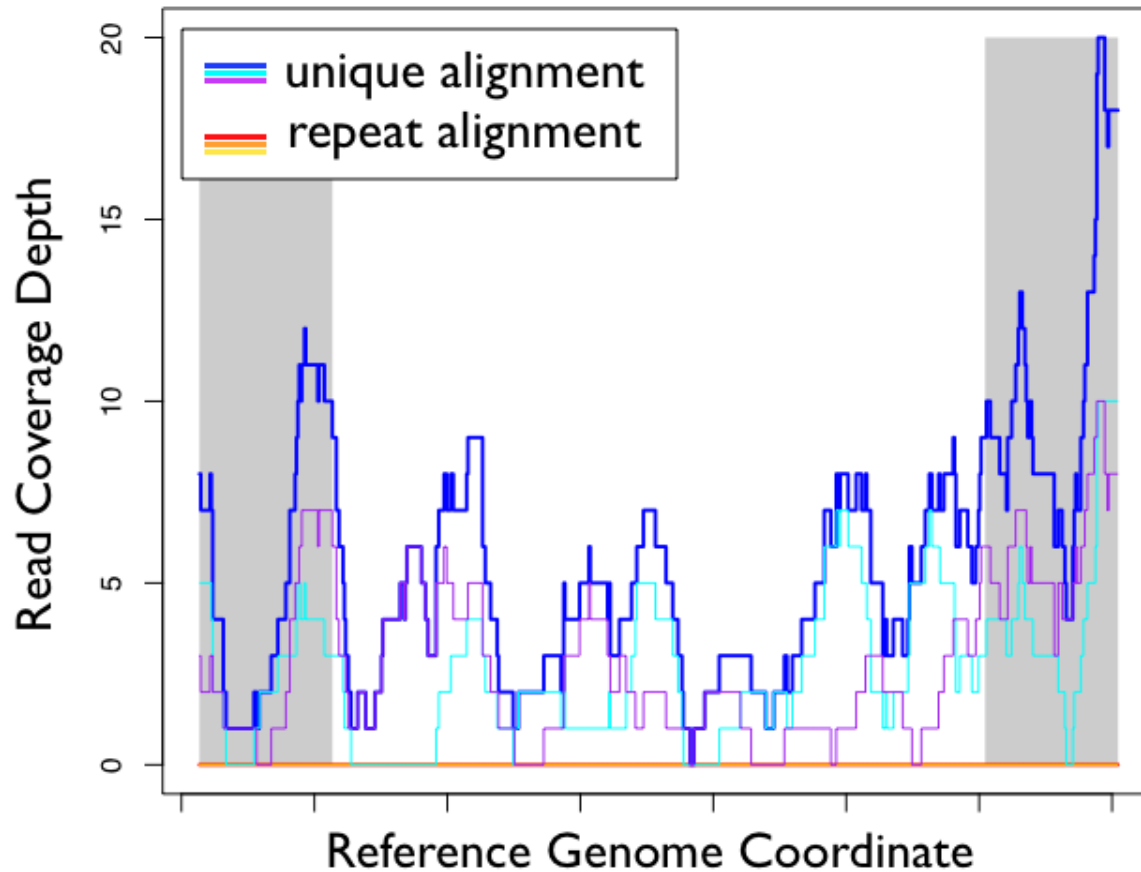# Example of a *breseq* prediction

- 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).
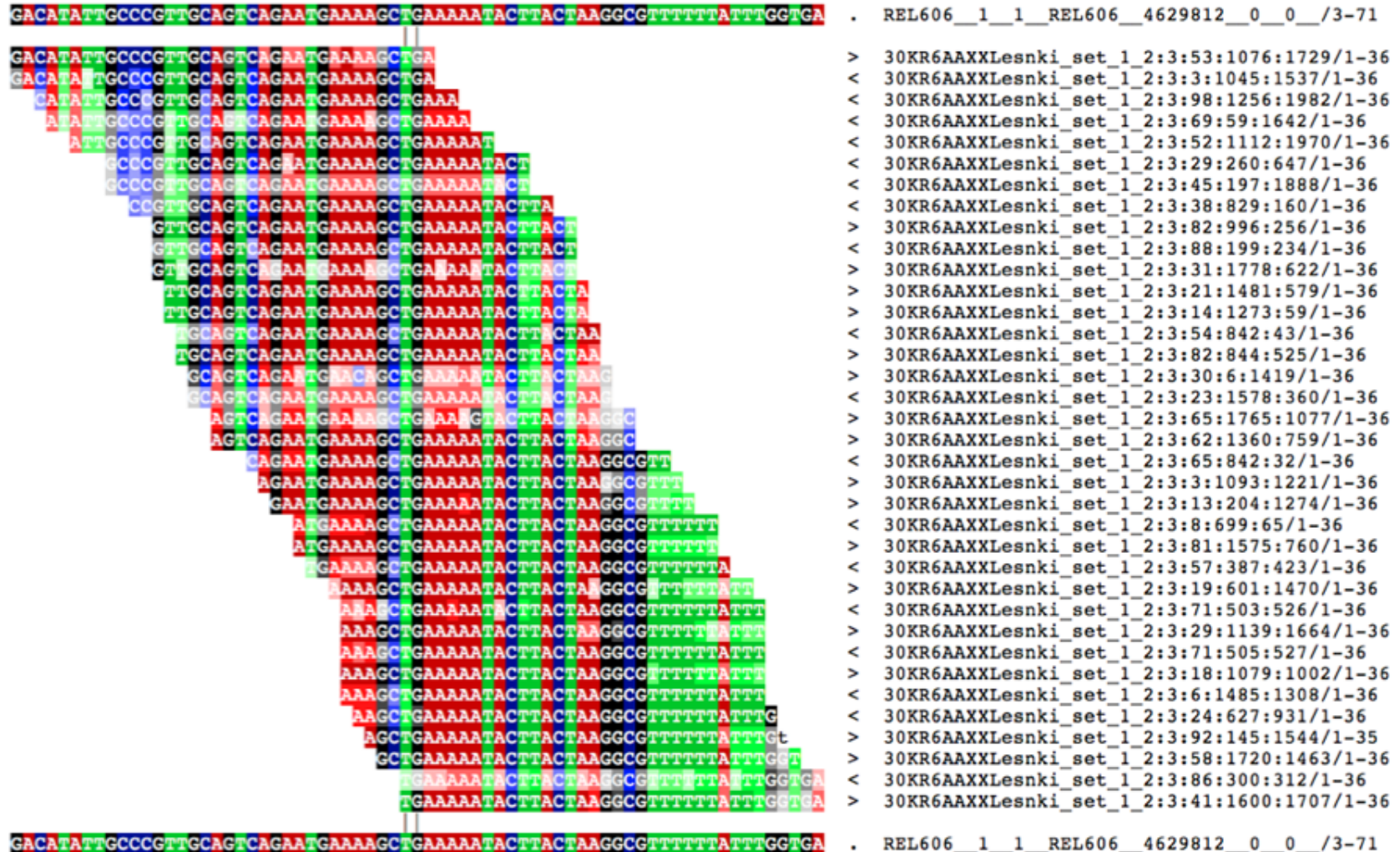


2. Create consensus list of possible new junctions.

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.

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

# Example of a bad junction

- Beware of reads ending in homopolymer runs!

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



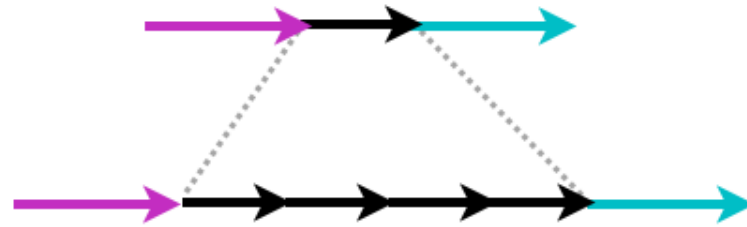Base Quality Score Legend: ATCG < 22 ≤ ATCG < 28 ≤ ATCG < 34 ≤ ATCG

# 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 | *mokC/nhaA* | 16959/17487 | regulatory protein for HokC, overlaps CDS of hokC/pH-dependent sodium/proton antiporter |
| | ? | = 3652533 | | | *IS150* | 3651091-3652533 | repeat region |
| * | ? | = 16992 | 0 | 41 | *mokC/nhaA* | 16959/17487 | regulatory protein for HokC, overlaps CDS of hokC/pH-dependent sodium/proton antiporter |
| | ? | 3893554 = | | | *IS150* | 3893554-3894996 | repeat region |

- Sometimes both new and old junctions exist...



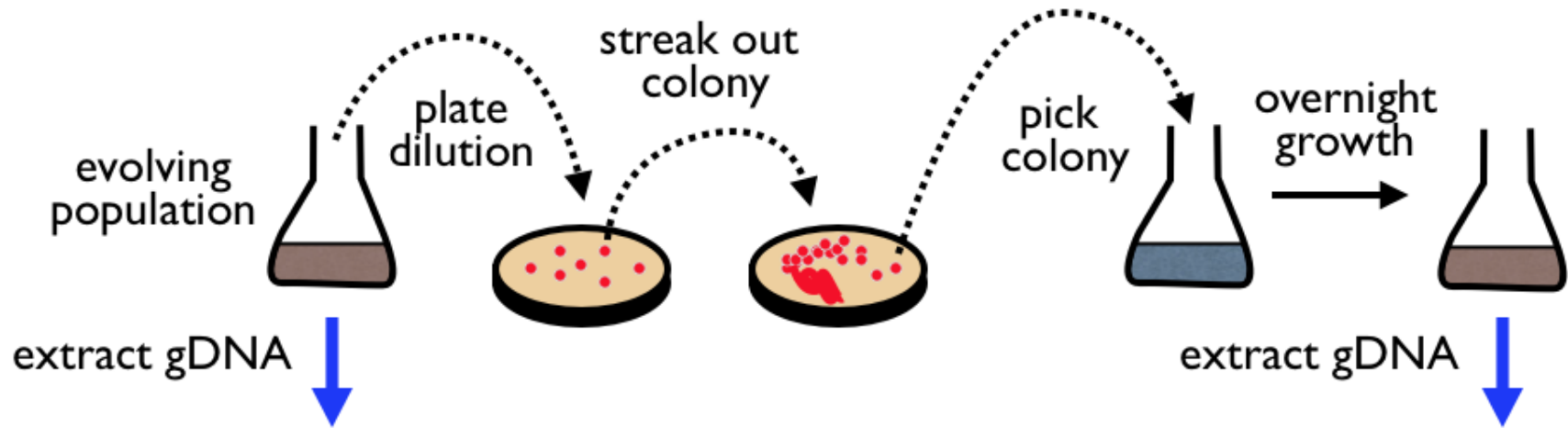tandem head-to-tail duplications

# Identifying copy number variation

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



- Tile into segments, train model on many genomes, look for deviation

# 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**

**Aligned reads**

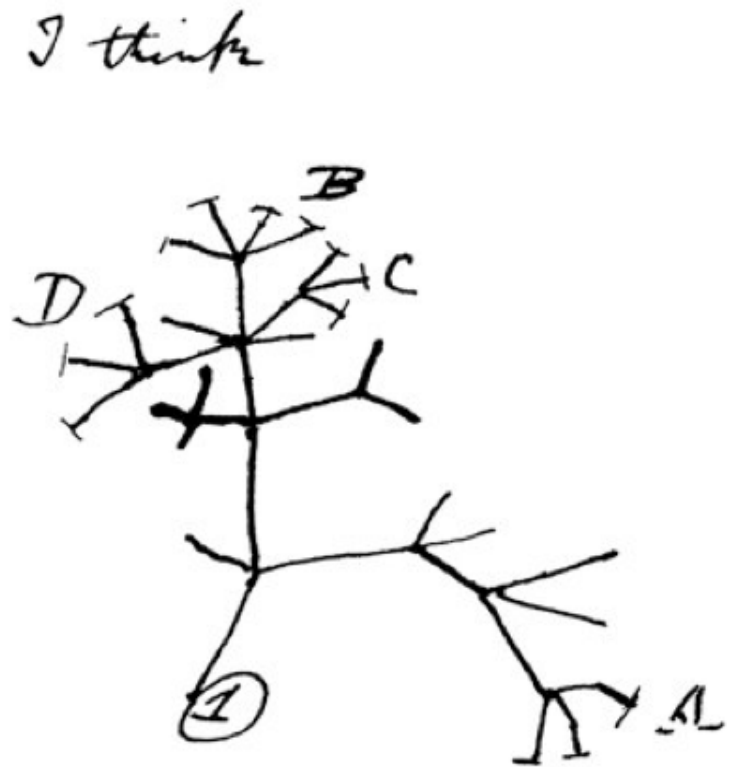- Map reads to ancestor genome. Only consider single-base substitutions.

- Log-likelihood test for polymorphism:

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

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

- Filter out predictions with other biases: strand bias, systematically low quality scores

- 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. SNV in region that is later deleted)

# 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 Comp Sci 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, ...