

# Introduction: Identifying Mutations and Studying Microbial Genome Evolution with *breseq*



Jeffrey E. Barrick

Department of Molecular Biosciences

THE UNIVERSITY OF  
**TEXAS**  
AT AUSTIN

<http://barricklab.org>



@barricklab

# Workshop Introduction

- When is *breseq* the right tool?
  - Installation
  - Basic usage
  - Input: references and reads
  - Output: HTML, GenomeDiff, etc.
- Analysis examples: Lenski LTEE
- Using *breseq* in research and education:  
**The other speakers in this workshop!**
- Online tutorials and workshop survey



# When is *breseq* the right tool?



Deatherage, D. E., Barrick, J. E. (2014) Identification of mutations in laboratory-evolved microbes from next-generation sequencing data using *breseq*.

*Methods Mol. Biol.* **1151**: 165–188. [https://doi.org/10.1007/978-1-4939-0554-6\\_12](https://doi.org/10.1007/978-1-4939-0554-6_12)

<https://barricklab.org/breseq>

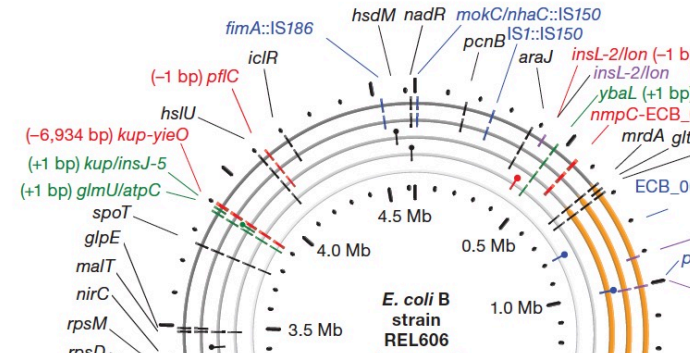
<https://github.com/barricklab/breseq>

- You have short-read NGS resequencing data.
- Your reference genome is ***haploid***.
  - Bacteria, Archaea, Phages, Plasmids, Haploid yeast
- You expect few genetic differences from the reference (a few to <1,000) in each sample.
- It's important that you identify all mutations.
- You are comfortable with using the terminal a little.
  - Changing directories, copying files, running a command

# Some uses of *breseq*

## Genetics

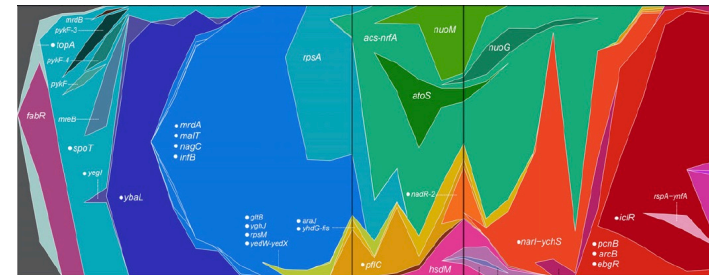
- Mechanisms of antibiotic resistance
- Mapping suppressor mutations



Barrick et al. (2009) *Nature*

## Experimental evolution

- Rates/nature of genome evolution
- Genetic diversity in populations



Maddamsetti et al. (2015) *Genetics*

## Biotechnology

- Verifying engineered plasmids/genomes
- Understanding beneficial mutations that arise during adaptive laboratory evolution

# Installing *breseq*

breseq 0.35.4 documentation » **breseq** Manual



Next topic

[Introduction](#)

This Page

[Show Source](#)

Quick search

Go

## breseq Manual

Contents:

- Introduction
  - Citing **breseq**
- Installation
  - Install external dependencies
  - Method 1. Binary download
  - Method 2. Source code download
    - Installing in a system-wide location
    - Installing in the source directory
    - Installing in a custom location
  - Method 3. GitHub source code
  - Installing on Cygwin (Windows)
  - Installing on Galaxy
  - Troubleshooting installation
- Usage Summary
  - **breseq**
  - Command: bam2aln
  - Command: bam2cov
- Test Drive
  - 1. Download data files
    - Reference sequence
    - Read files
  - 2. Run **breseq**
  - 3. Open **breseq** output

Latest release

v0.35.6

c7cf8df

Compare

## breseq v0.35.6

Edit

jeffreybarrick released this 25 days ago

- Fixed compatibility with GenBank reference files produced by Prokka and NCBI PGAP, and with GFF3 files produced by PGAP.

Assets 5

 <a href="#">breseq-0.35.6-Linux-x86_64.tar.gz</a>	13.7 MB
 <a href="#">breseq-0.35.6-MacOSX-10.9+.tar.gz</a>	13.9 MB
 <a href="#">breseq-0.35.6-Source.tar.gz</a>	12.4 MB
 <a href="#">Source code (zip)</a>	
 <a href="#">Source code (tar.gz)</a>	

<https://github.com/barricklab/breseq/releases>


Can be used on Linux, Mac OSX, and Windows machines; and in the Galaxy web platform.

Options to download and install by compiling from source code or using precompiled binaries.

Requires R and bowtie2.

# Installing *breseq*

breseq 0.35.4 documentation » **breseq** Manual



Next topic  
Introduction

This Page  
Show Source

Quick search  
 Go

## breseq Manual

Contents:

- Introduction
  - Citing **breseq**
- Installation
  - Install external dependencies
  - Method 1. **Binary download**
  - Method 2.
    - Install
    - Install
    - Install
  - Method 3.
  - Installing
  - Installing
  - Troubleshooting
- Usage Summary
  - **breseq**
  - Command
  - Command
- Test Drive
  - 1. Download
    - Reference
    - Read
  - 2. Run **breseq**
  - 3. Open b

Recommended method

**BIOCONDA**<sup>®</sup>

<http://bioconda.github.io/index.html>

**BIOCONDA**<sup>®</sup>

Navigation

User Docs

Contributing to Bioconda

Developer Docs

Tutorials

Bioconda @ Github

Package Index

chat on gitter

Quick search

Go

*recipe* **breseq**

A computational pipeline for finding mutations relative to a reference sequence in short-read DNA re-sequencing data.

Homepage: <https://github.com/barricklab/breseq>

License: GPL / GPL-3.0

Recipe: </breseq/meta.yaml>

*package* **breseq**

downloads 18k container none

Versions: ▶ 0.35.6-0, 0.35.5-1, 0.35.5-0, 0.35.4-0, 0.35.3-0, 0.35.2-0, 0.35.1-0, 0.35.0-0, 0.34.1-0, ...

Depends: [bowtie2](#) >=2.0.0, !=2.0.3, !=2.0.4, !=2.3.1, [libgcc-ng](#) >=9.3.0, [libstdcxx-ng](#) >=9.3.0, [r-base](#), [zlib](#) >=1.2.11, <1.3.0a0

Required By:

Installation

With an activated Bioconda channel (see [2. Set up channels](#)), install with:

```
conda install breseq
```

# Basic *breseq* usage

breseq 0.35.4 documentation » **breseq** Manual

## Basic *breseq* command

```
$ breseq -r reference.gbk reads_1.fastq reads_2.fastq
```

References can be in GenBank, GFF3, or FASTA format.

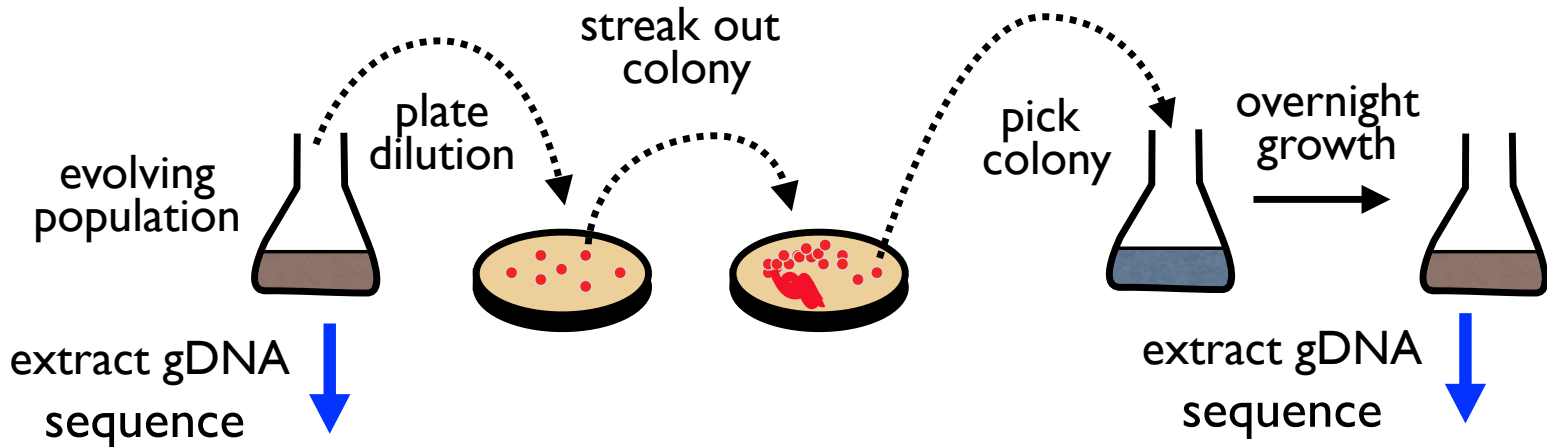
Multiple reference files can be used: `-r genome.fasta -r plasmid.gff3`

Read files can be gzipped: `reads_1.fastq.gz`

Speed up execution by using multiple cores: `-j 8`

- Troubleshooting installation
- Usage Summary
  - **breseq**
  - Command: `bam2aln`
  - Command: `bam2cov`
- Test Drive
  - 1. Download data files
    - Reference sequence
    - Read files
  - 2. Run **breseq**
  - 3. Open **breseq** output

# Two main types of samples

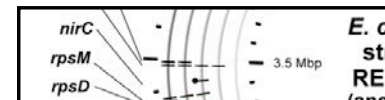


Every read could be from any individual.

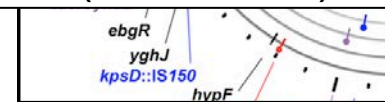
All reads are from a single clone.

Population or Polymorphism mode

```
$ breseq -p ...
```





Consensus mode (the default)



```
$ breseq ...
```

# Reference file considerations

- **Microbes (<20Mb):** download GenBank or GFF3 files with both DNA sequence and features.
- **Important:** having transposable elements annotated leads to better predictions!
- **What do I do if there is no reference?**
  - *de novo* assemble and annotate your own
  - **Recommendation:**  Unicycler **PROKKA** 
  - If you are using an assembly that has multiple contigs use `-c` instead of `-r` for specifying the contig reference:

```
$ breseq -c contigs.gbk reads_1.fastq reads_2.fastq
```
  - You may need to iteratively improve the assembly and annotation to get the best results. See **gdttools** **APPLY**.

# Downloading a reference from NCBI

⚠ Be sure you download a GenBank file that has both features and the sequence!

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide Search

Advanced Help

GenBank Send to: Change region shown

**Escherichia coli B str. REL606, complete sequence**

NCBI Reference Sequence: NC\_012967.1

[FASTA](#) [Graphics](#)

Go to: [v]

LOCUS NC\_012967 4629812 bp DNA circular CON 13-MAY-2021

DEFINITION Escherichia coli B str. REL606, complete sequence.

ACCESSION NC\_012967

VERSION NC\_012967.1

DBLINK BioProject: [PRJNA224116](#)  
BioSample: [SAMN02603421](#)  
Assembly: [GCF\\_000017985.1](#)

KEYWORDS RefSeq.

SOURCE Escherichia coli B str. REL606

ORGANISM [Escherichia coli B str. REL606](#)

REFERENCE 1 (bases 1 to 4629812)

AUTHORS Jeong,H., Barbe,V., Vallenet,D., Choi,S.-H., Lee,C.H., Lee,S.-W.,  
Vacherie,B., Yoon,S.H., Yu,D.-S., Cattolico,L., Hur,C.-G.,  
Park,H.-S., Segurens,B., Blot,M., Schneider,D., Studier,F.W.,  
Oh,T.K., Lenski,R.E., Daegelen,P. and Kim,J.F.

CONSRMT International E. coli B Consortium

TITLE Complete genome sequence of Escherichia coli (B) REL606

JOURNAL Unpublished

REFERENCE 2 (bases 1 to 4629812)

AUTHORS Daegelen,P., Vallenet,D., Barbe,V., Cattolico,L. and Segurens,B.

TITLE Direct Submission

JOURNAL Submitted (24-AUG-2007) UMR 8030 CNRS, Inserm, Genoscope

Customize view

☐ Abbreviated view

☒ Customize

Basic Features

☒ All features

☐ Gene, RNA, and CDS features only

Display options

☒ Show sequence

☐ Show reverse complement

☐ Show gap features

Update View

Analyze this sequence

Run BLAST

Pick Primers

Highlight Sequence Features

Related information

Assembly

BioProject

BioSample

Open

Select

Click

# Downloading a reference from NCBI

⚠ Be sure you download a GenBank file that has both features and the sequence!

```
gene
4629102..4629788
/feature="yjtd"
/locus_tag="ECB_RS22810"
/old_locus_tag="ECB_04279"
4629102..4629788
/feature="yjtd"
/locus_tag="ECB_RS22810"
/old_locus_tag="ECB_04279"
/EC_number="2.1.1.-"
/inference="COORDINATES: similar to AA
sequence:RefSeq:NP_710140.2"
/note="Derived by automated computational analysis using
gene prediction method: Protein Homology."
/codon_start=1
/transl_table=11
/product="tRNA/rRNA methyltransferase"
/protein_id="WP_001223167.1"
/translation="MRITIIIVAPARAENIGAAARAMKTMGFSELRIVDSPAHEPAT
RWVAHGSGDIIDNIKVFPPTLAESLHDVDFTVATTARSRAKYHYATPVELVPLLEES
SWMSHAALVFGREDSTNEELALADVLGTGPMVADYPSNLGQAVMYCYQLATLIQ
QPTKSDTTADQHQLQALRERVMALLTTLAVADDIKLVDWLQRLGLLEQDRTAMLRHL
LHDIEKNITK"

ORIGIN
1 agcttttcat tctgactgca acgggcaata tgtctctgtg tggattaaaa aaagagtgtc
61 tgatagcagc ttctgaactg gttacctgcc gtgagtaaat taaaatttta ttgacttagg
121 tcactaaata ctttaaccaa tataggcata gcgcacagac agataaaaat tacagagtac
181 acaacatcca tgaacacgat tagcaccacc attaccacca ccatcaccat taccacaggt
241 aacggtgcgg gctgacgcgt acagaaaaca cagaaaaaag cccgcacctg acagtgcggg
301 ctttttttcc gaccaaagggt aacgaggtta caacctgcg agtgttgaag ttcggcggtg
361 catcagtgcc aaatgcagaa cgttttctgc gggttgccga tattctggaa agcaatgcc
421 ggcaggggca ggtggccacc gtccctctctg cccccgccaa aatcaccaac cactgggtgg
481 cgatgattga aaaaaccatt agcggccagg atgctttacc caatatcagc gatgccgaac
541 gtatttttgc cgaacttttg acgggaactg cccgcgccca gccgggattc ccgctggcgc
601 aattgaaaaa tttcgtcgat caggaaattg cccaaataaa acatgtcctg catggcatta
661 gtttggttgg gcagtgcccg gatagcatca acgctgcgct gatttgcctg ggcgagaaaa
721 tgtcgcgcgc cattatggcc gccgtattag aagcgcgcgc tcacaacgct accgttatcg
781 atccgggtcg aaaaactgct gcagtggggc attacctcga atctaccgtc gatattgctg
841 agtccacccc ccgtattgcg gcaagtcgca ttccggctga tcacatgggt ctgatggcag
```

Scroll down until  
you see ORIGIN.



There should be a  
nucleotide  
sequence here!

# Read file considerations

## Sequencing technology

- Can work with any FASTQ
- Best results with short-read data (< 1000 bases)
- Not appropriate for **long-read** data (Nanopore, PacBio, etc.) In this case, you should *de novo* assemble and then compare assemblies.

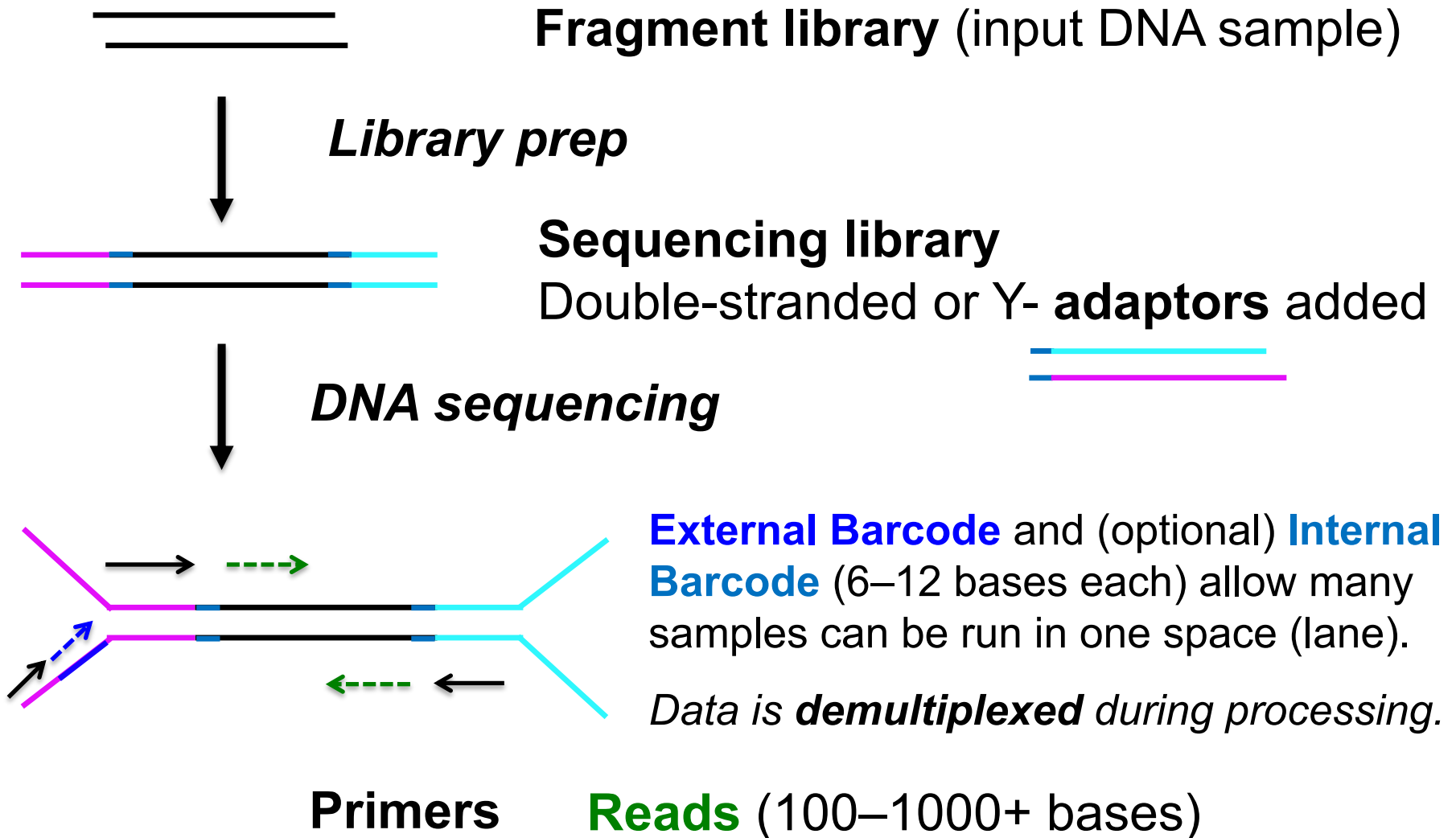
## Recommended depth of coverage

>40x for clonal samples

>120x for population samples

More coverage is unlikely to give improvements without error correction (ex: molecular barcodes).

# Read terminology



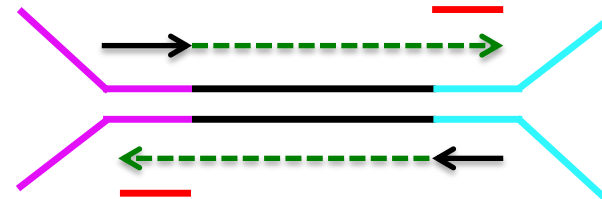
# FASTQ quality and trimming

Check the quality of your FASTQ data  FastQC.app

- Have internal barcodes been removed?
- Do I want to trim low-quality bases?

! Be careful to trim adaptors from your reads  
(*breseq* requires >90% of a read's length to map)

**Readthrough into adaptors** is  
especially common with  
new longer Illumina reads!



Programs that can help:

`fastp`, `trimmomatic`, `cutadapt`

```
$ breseq -j 8 -r REL606.gbk SRR030255_1.fastq.gz SRR030255_2.fastq.gz
```

36 minutes later... Open [output/index.html](#)

# HTML Output



breseq version 0.35.6 revision c7cf8df53bcd

[mutation predictions](#) | [marginal predictions](#) | [summary statistics](#) | [genome diff](#) | [command line log](#)

Predicted mutations					
evidence	position	mutation	annotation	gene	description
<a href="#">RA</a>	380,188	A→C	F239L (TTT→TTG)	<i>araJ</i> ←	predicted transporter
<a href="#">RA</a>	475,292	+G	coding (14/1677 nt)	<i>ybaL</i> ←	predicted transporter with NAD(P)-binding Rossmann-fold domain
<a href="#">RA</a>	649,391	T→A	I471F (ATC→ITC)	<i>mrdA</i> ←	transpeptidase involved in peptidoglycan synthesis (penicillin-binding protein 2)
<a href="#">RA</a>	683,496	A→C	V65G (GTT→GGT)	<i>nagC</i> ←	DNA-binding transcriptional dual regulator, repressor of N-acetylglucosamine
<a href="#">JC JC</a>	969,836	IS 150 (+) +3 bp	coding (810-812/2283 nt)	<i>pflB</i> ←	pyruvate formate lyase I
<a href="#">RA</a>	1,329,516	C→T	H33Y (CAC→TAC)	<i>topA</i> →	DNA topoisomerase I
<a href="#">JC JC</a>	1,544,289	IS 150 (-) +3 bp	coding (150-152/1536 nt)	<i>xasA</i> ←	predicted glutamate:gamma-aminobutyric acid antiporter
<a href="#">JC JC</a>	1,733,647	IS 150 (-) +3 bp	coding (683-685/1413 nt)	<i>pykF</i> →	pyruvate kinase
<a href="#">RA</a>	1,976,879	T→G	intergenic (-57/-76)	<i>yedW</i> ← / → <i>yedX</i>	predicted DNA-binding response regulator in two-component system with YedV/hypothetical protein
<a href="#">RA</a>	2,082,685	G→A	A494V (GCT→GIT)	<i>yegI</i> ←	hypothetical protein
<a href="#">RA</a>	2,499,315	G→A	intergenic (-110/-179)	<i>maeB</i> ← / → <i>talA</i>	malic enzyme/transaldolase A
<a href="#">RA</a>	3,045,069	G→T	T312N (ACC→AAC)	<i>yghJ</i> ←	predicted inner membrane lipoprotein
<a href="#">RA</a>	3,248,957	A→T	D764E (GAT→GAA)	<i>infB</i> ←	translation initiation factor IF-2
<a href="#">MC JC</a>	3,289,962	Δ16 bp	coding (96-111/4554 nt)	<i>gltB</i> →	glutamate synthase, large subunit
<a href="#">RA</a>	3,339,158	A→C	intergenic (+22/-4)	<i>yhdG</i> → / → <i>fis</i>	tRNA-dihydrouridine synthase B/DNA-binding protein Fis
<a href="#">RA</a>	3,370,027	T→A	K117M (AAG→ATG)	<i>rpsM</i> ←	30S ribosomal protein S13
<a href="#">RA</a>	3,424,910	G→A	M1M (ATG→ATA) †	<i>nirC</i> →	nitrite transporter
<a href="#">RA</a>	3,483,047	C→A	R455S (CGC→AGC)	<i>malT</i> →	transcriptional regulator MalT
<a href="#">RA</a>	3,762,741	A→T	K662I (AAA→ATA)	<i>spoT</i> →	bifunctional (p)ppGpp synthetase II/ guanosine-3',5'-bis pyrophosphate 3'-pyrophosphohydrolase
<a href="#">RA</a>	3,875,632	(T)7→8	intergenic (-66/+287)	<i>glmU</i> ← / → <i>atpC</i>	bifunctional N-acetylglucosamine-1-phosphate uridyltransferase/glucosamine-1-phosphate acetyltransferase/F0F1 ATP synthase subunit epsilon
<a href="#">RA</a>	3,893,551	+G	intergenic (+6/-50)	<i>kup</i> → / → <i>insJ-5</i>	potassium transporter/IS150 hypothetical protein
<a href="#">MC JC</a>	3,894,997	Δ6,934 bp	IS150-mediated	<i>rbsD</i> –[ <i>yieO</i> ]	<i>rbsD</i> , <i>rbsA</i> , <i>rbsC</i> , <i>rbsB</i> , <i>rbsK</i> , <i>rbsR</i> , [ <i>yieO</i> ]
<a href="#">RA</a>	4,100,655	C→T	M192I (ATG→ATA)	<i>hslU</i> ←	ATP-dependent protease ATP-binding subunit
<a href="#">RA</a>	4,126,706	(T)8→7	coding (342/879 nt)	<i>pflC</i> →	pyruvate formate lyase II activase
<a href="#">RA</a>	4,560,632	T→C	Y131C (TAC→TGC)	<i>hsdM</i> ←	DNA methylase M

Unassigned missing coverage evidence									
	seq id	start	end	size	←reads	reads→	gene	description	
* - ±	REL606	546953–547700	555934–555877	8178–8982	20 [18]	[16] 19	[ <i>insB-6</i> ]–[ <i>ECB_00513</i> ]	[insB-6],insA-6,nmpC,ycbR,ycbS,ycbT,ycbU,ECB_00510,nohB,ECB_00512,[ECB_00513]	
* - ±	REL606	2031675–2031718	2054970–2054943	23226–23296	21 [17]	[18] 21	[ <i>manB</i> ]–[ <i>cpsG</i> ]	[manB],manC,insB-14,insA-14,wbbD,wbbC,wzy,wbbB,wbbA,vioB,vioA,wzx,rmlC,rfbA,rfbD,rfbB,galf,wcaM,wcaL,wcaK,wzc,wcaJ,[cpsG]	

Unassigned new junction evidence										
	seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotation	gene	product
* ?	REL606	= 547699	NA (NA)	80 (1.360)	37/70	0.2	NA	noncoding (1/768 nt)	IS1	repeat region
- ?	REL606	555924 =	NA (NA)					coding (1209/2346 nt)	ECB_00513	conserved hypothetical protein

Predicted mutations					
evidence	position	mutation	annotation	gene	
<a href="#">RA</a>	380,188	A→C	<a href="#">239L (TTT→TTG)</a>	<i>araJ</i> ←	predicted trans
<a href="#">RA</a>	475,292	+G	<a href="#">coding (14/1677 nt)</a>	<i>ybaI</i> ←	predicted trans
<a href="#">RA</a>	649,391	T→A	<a href="#">I471F (ATC→ITC)</a>	<i>mrdA</i> ←	transpeptidase
<a href="#">RA</a>	683,496	A→C	<a href="#">V65G (GTT→GGT)</a>	<i>nagC</i> ←	DNA-binding tr
<a href="#">JC JC</a>	969,836	IS 150 (+) +3 bp	coding (810-812/2283 nt)	<i>pflB</i> ←	pyruvate form
<a href="#">RA</a>	1,329,516	C→T	<a href="#">H33Y (CAC→TAC)</a>	<i>topA</i> →	DNA topoisom
<a href="#">JC JC</a>	1,544,289	IS 150 (-) +3 bp	coding (150-152/1536 nt)	<i>xasA</i> ←	predicted gluta
<a href="#">JC JC</a>	1,733,647	IS 150 (-) +3 bp	coding (683-685/1413 nt)	<i>pykF</i> →	pyruvate kinas
<a href="#">RA</a>	1,976,879	T→G	intergenic (-57/-76)	<i>yedW</i> ← / <i>yedX</i>	predicted DNA
<a href="#">RA</a>	2,082,685	G→A	<a href="#">A494V (GCT→GTT)</a>	<i>yegI</i> ←	hypothetical pr
<a href="#">RA</a>	2,499,315	G→A	intergenic (-110/-179)	<i>maeB</i> ← / <i>talA</i>	malic enzyme/ transcarboxylase
<a href="#">RA</a>	3,045,069	G→T	<a href="#">T312N (ACC→AAC)</a>	<i>yghJ</i> ←	predicted inner membrane lipoprotein
<a href="#">RA</a>	3,248,957	A→T	<a href="#">D764E (GAT→GAA)</a>	<i>infB</i> ←	translation initiation factor IF-2
<a href="#">MC JC</a>	3,289,962	Δ16 bp	coding (96-111/4554 nt)	<i>gltB</i> →	glutamate synthase, large subunit
<a href="#">RA</a>	3,339,158	A→C	intergenic (+22/-4)	<i>yhdG</i> ← / <i>fis</i>	tRNA-dihydrouridine synthase B/DNA-binding protein Fis
<a href="#">RA</a>	3,370,027	T→A	<a href="#">K117M (AAG→ATG)</a>	<i>rpsM</i> ←	30S ribosomal protein S13
<a href="#">RA</a>	3,424,910	G→A	<a href="#">M11M (ATG→ATA) †</a>	<i>nirC</i> →	nitrite transporte
<a href="#">RA</a>	3,483,047	C→A	<a href="#">R455S (CGC→AGC)</a>	<i>malT</i> →	transcriptional re
<a href="#">RA</a>	3,762,741	A→T	<a href="#">K662I (AAA→ATA)</a>	<i>spoT</i> →	bifunctional (p)pr
<a href="#">RA</a>	3,875,632	(T) <sub>7</sub> →8	intergenic (-66/+287)	<i>glmU</i> ← / <i>atpC</i>	bifunctional N-ac subunit epsilon
<a href="#">RA</a>	3,893,551	+G	intergenic (+6/-50)	<i>kup</i> → / <i>insJ-5</i>	potassium transp
<a href="#">MC JC</a>	3,894,997	Δ6,934 bp	IS 150-mediated	<i>rbsD-[yieO]</i>	<i>rbsD</i> , <i>rbsA</i> , <i>rbsC</i>
<a href="#">RA</a>	4,100,655	C→T	<a href="#">M192I (ATG→ATA)</a>	<i>hslH</i> ←	ATP-dependent
<a href="#">RA</a>	4,126,706	(T) <sub>8</sub> →7	coding (342/879 nt)	<i>pflC</i> →	pyruvate format
<a href="#">RA</a>	4,560,632	T→C	<a href="#">Y131C (TAC→TGC)</a>	<i>hsdM</i> ←	DNA methylase

## Mutations (fully predicted)

- Base substitutions
- Small indels
- IS element insertions
- Large deletions

Evidence for other genetic differences that can't be fully resolved


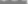
Unassigned missing coverage evidence									
seq id	start	end	size	←reads	reads→	gene	description		
* - ± REL606	546953-547700	555934-555877	8178-8982	20 [18]	[16] 19	<i>[insB-6]-[ECB_00513]</i>	[insB-6],insA-6,nmpC,ybcR,ybcS,ybcT,ybcU,ECB_00510,nohB,ECB_00512,[ECB_00513]		
* - ± REL606	2031675-2031718	2054970-2054943	23226-23296	21 [17]	[18] 21	<i>[manB]-[cpsG]</i>	[manB],manC,insB-14,insA-14,wbbD,wbbC,wzy,wbbB,wbbA,vioB,vioA,wzx,rmlC,rfaA,rfaD,rfaB,galF,wcaM,wcaL,wcaK,wzc,wcaJ,[cpsG]		

Unassigned new junction evidence									
seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotation	gene	product
* ? REL606	= 547699	NA (NA)	80 (1.360)	37/70	0.2	NA	noncoding (1/768 nt)	<i>IS1</i>	repeat region
- ? REL606	555924 =	NA (NA)					coding (1209/2346 nt)	<i>ECB_00513</i>	conserved hypothetical protein



[mutation predictions](#) | [marginal predictions](#) | [summary statistics](#) | [genome diff](#) | [command line log](#)

feature	position	mutation	annotation	gene
<b>RA</b>	380,188	A→C	<b>F239L</b> (TT <b>T</b> →TT <b>G</b> )	<i>araJ</i> ←
<b>RA</b>	475,232	+G	coding (14/1677 nt)	<i>ybaL</i> ←
<b>RA</b>	649,391	T→A	<b>I471F</b> (ATC→ <b>ITC</b> )	<i>mrdA</i> ←
<b>RA</b>	683,496	A→C	<b>V65G</b> (G <b>IT</b> →G <b>GT</b> )	<i>nagC</i> ←
<b>JC JC</b>	969,836	IS150 (+) +3 bp	coding (810-812/2283 nt)	<i>pflB</i> →
<b>RA</b>	1,329,516	C→T	<b>H33Y</b> (C <b>AC</b> →T <b>AC</b> )	<i>topA</i> →
<b>JC JC</b>	1,544,289	IS150 (-) +3 bp	coding (150-152/1556 nt)	<i>xasA</i> ←
<b>JC JC</b>	1,733,647	IS150 (-) +3 bp	coding (683-685/1413 nt)	<i>pykF</i> →
<b>RA</b>	1,976,879	T→G	intergenic (-57/-76)	<i>yedW</i> ← / → <i>yedX</i>
<b>RA</b>	2,082,685	G→A	<b>A494V</b> (G <b>CT</b> →G <b>TT</b> )	<i>yegl</i> ←
<b>RA</b>	2,499,315	G→A	intergenic (-110/-179)	<i>maeB</i> ← / → <i>talA</i>
<b>RA</b>	3,045,069	G→T	<b>T312N</b> (A <b>CC</b> →A <b>AC</b> )	<i>yghJ</i> ←
<b>RA</b>	3,248,957	A→T	<b>D764E</b> (G <b>AT</b> →G <b>AA</b> )	<i>infB</i> ←
<b>MC JC</b>	3,289,962	Δ16 bp	coding (96-111/4554 nt)	<i>gltB</i> →
<b>RA</b>	3,339,158	A→C	intergenic (+22/-4)	<i>yhdG</i> → / → <i>fis</i>
<b>RA</b>	3,370,027	T→A	<b>K117M</b> (A <b>AG</b> →A <b>IG</b> )	<i>rpsM</i> ←
<b>RA</b>	3,424,910	G→A	<b>M11M</b> (AT <b>G</b> →AT <b>A</b> ) †	<i>nirC</i> →
<b>RA</b>	3,483,047	C→A	<b>R455S</b> (C <b>GC</b> → <b>AGC</b> )	<i>malT</i> →
<b>RA</b>	3,762,741	A→T	<b>K662I</b> (A <b>AA</b> →A <b>TA</b> )	<i>spoT</i> →
<b>RA</b>	3,875,632	(T) <sub>7→8</sub>	intergenic (-66/+287)	<i>glmU</i> ← / ← <i>atpC</i>
<b>RA</b>	3,893,551	+G	intergenic (+6/-50)	<i>kup</i> → / → <i>insJ-5</i>
<b>MC JC</b>	3,894,997	Δ6,934 bp	IS150-mediated	<i>rbsD</i> -[ <i>yieO</i> ]
<b>RA</b>	4,100,655	C→T	<b>M192I</b> (AT <b>G</b> →AT <b>A</b> )	<i>hslU</i> ←
<b>RA</b>	4,126,706	(T) <sub>8→7</sub>	coding (342/879 nt)	<i>pflC</i> →
<b>RA</b>	4,560,632	T→C	<b>Y131C</b> (T <b>AC</b> →T <b>GC</b> )	<i>hsdM</i> ←

	seq id	start	end	size	←reads	reads→	gene	
	REL606	546953–547700	555934–555877	8178–8982	20 [18]	[16] 19	[ <i>insB-6</i> ]–[ <i>ECB_00513</i> ]	[ <i>ins</i> ]
	REL606	2031675–2031718	2054970–2054943	23226–23296	21 [17]	[18] 21	[ <i>manB</i> ]–[ <i>cpsG</i> ]	[ <i>ma</i> ]

	seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotat
+	REL606	= 547699	NA (NA)	80 (1.360)	37/70	0.2	NA	noncoding (1
-	REL606	555924 =	NA (NA)					coding (1209

evidence	seq id	position	mutation	annotation	gene	description
<a href="#">RA</a>	REL606	380,188	A→C	<a href="#">F239L</a> (TTT→TTG)	<i>araJ</i> ←	predicted transporter

seq id	position	ref	new	freq	score (cons/poly)	reads	annotation	genes	product
* REL606	380,188	0	A	C	100.0%	129.8 / NA	40	F239L (TTT→TTG)	<i>araJ</i> predicted transporter

Reads supporting (aligned to +/- strand): ref base A (0/0); new base C (18/22); total (18/22)

CAGCACCATCCCTAGCCCAACTAACATCATTAATAAAGGTCATCGCCGTTTCCGAAAAACCGGAAAT > REL606/380155-380220

```

CAGCACCA|CCC|AGCCCAACT|AATCAT|CATAAT|AAGGTC|CAT|CGCG|TT|CCG|AAAAAC|CGG|AAA| > REL606/380155-380220

caac|ACC|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|aa < 1:4049041/33-1 (MQ=38)
cAG|CACC|AT|CCC|AG|CCCA|AACT|AAT|CAT|ATAAT|aa > 2:3988152/1-36 (MQ=38)
aGC|CACC|AT|CCC|AG|CCCA|AACT|AAT|CAT|ATAAT|AAg > 2:2846392/1-36 (MQ=255)
aGC|CACC|AT|CCC|AG|CCCA|AACT|AAT|CAT|ATAAT|AAg > 1:1710600/1-36 (MQ=255)
aGC|CACC|AT|CCC|AG|CCCA|AACT|AAT|CAT|ATAAT|AAg < 2:478931/36-1 (MQ=255)
c|ACC|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGt < 2:3427822/36-1 (MQ=255)
c|ACC|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGt > 2:2868696/1-36 (MQ=255)
cc|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCa < 1:613685/36-1 (MQ=255)
cc|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCa < 1:3037406/36-1 (MQ=255)
cc|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCa < 1:2129644/36-1 (MQ=255)
c|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCat > 2:113100/1-36 (MQ=255)
c|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCat > 2:2149235/1-36 (MQ=255)
c|AT|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCat > 1:309277/1-36 (MQ=255)
a|T|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCatc > 1:1146898/1-36 (MQ=255)
a|T|CCCT|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCatc < 2:880957/36-1 (MQ=37)
ccc|TAG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCATCGc < 2:3039312/36-1 (MQ=255)
ccc|TAG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCATCGc < 1:1074806/36-1 (MQ=255)
ccc|TAG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCATCGc < 1:2828983/36-1 (MQ=37)
cc|TAG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCATCGcc > 2:904743/1-36 (MQ=37)
c|TAG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCg < 1:2956106/36-1 (MQ=255)
ta|aac|c|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCg < 2:2953907/33-1 (MQ=25)
t|AG|CCCA|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCg > 1:2693636/1-36 (MQ=37)
a|GCC|CA|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCgt > 2:1240813/36-1 (MQ=37)
ccc|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCGTTtc < 2:465024/36-1 (MQ=255)
cc|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCGTT|Tcc > 2:136753/1-36 (MQ=255)
c|AACT|AAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCg > 1:430752/1-36 (MQ=38)
aac|aac|c|T|CAAT|AAT|AAGGTCATCGCCGTTTCCGa < 1:566412/36-1 (MQ=37)
ac|T|CAAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa < 1:4078565/36-1 (MQ=255)
aa|c|T|CAAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa < 1:1319202/36-1 (MQ=255)
a|cT|AAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa > 1:2592025/1-36 (MQ=255)
a|cT|AAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa > 2:422670/1-36 (MQ=255)
cT|AAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa > 2:3759532/36-1 (MQ=255)
cT|AAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa < 1:2197499/36-1 (MQ=255)
t|AAT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa < 2:1187588/36-1 (MQ=255)
aa|cT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa < 2:243349/36-1 (MQ=21)
a|cT|CAT|ATAAT|AAGGTCATCGCCGTTTCCGa < 2:1690209/36-1 (MQ=255)
cat|cat|AAT|AAGGTCATCGCCGTTTCCGa < 1:2225790/36-1 (MQ=255)
cat|cat|AAT|AAGGTCATCGCCGTTTCCGa > 2:3173081/1-36 (MQ=255)
cat|cat|AAT|AAGGTCATCGCCGTTTCCGa < 1:1341065/36-1 (MQ=255)
at|cat|AAT|AAGGTCATCGCCGTTTCCGa > 2:270940/1-36 (MQ=25)
cat|at|AAT|AAGGTCATCGCCGTTTCCGa > 1:1008139/1-36 (MQ=255)

```

# MC = Missing Coverage Evidence



breseq version 0.35.6 revision c7cf8df53bcd

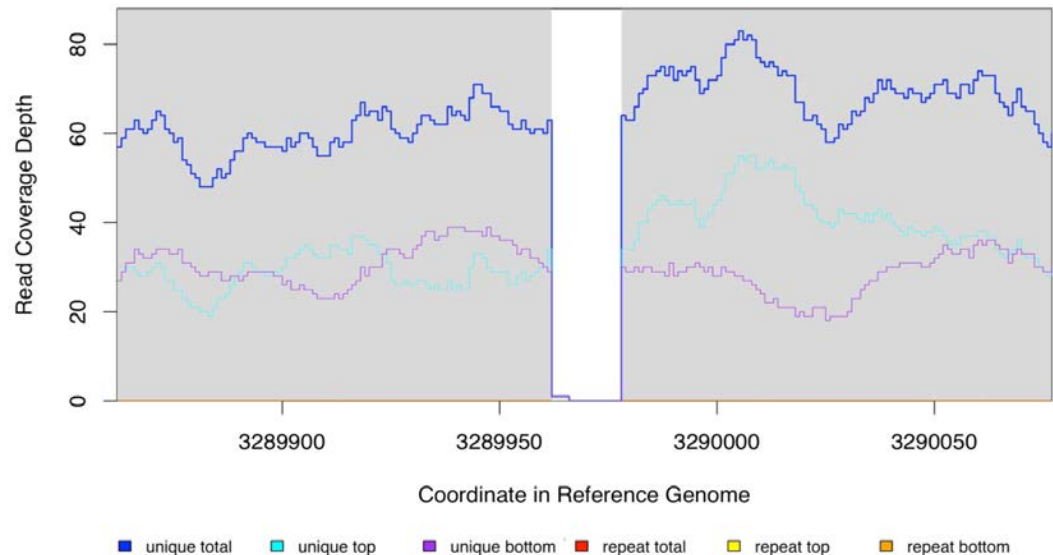
[mutation predictions](#) | [marginal predictions](#) | [summary statistics](#) | [genome diff](#) | [command line log](#)

Predicted mutations				
evidence	position	mutation	annotation	gene
<a href="#">RA</a>	380,188	A→C	F239L (TTI→TTG)	<i>araJ</i> ←
<a href="#">RA</a>	475,292	+G	coding (14/1677 nt)	<i>ybaL</i> ←
<a href="#">RA</a>	649,391	T→A	I471F (ATC→ITC)	<i>mrda</i> ←
<a href="#">RA</a>	683,496	A→C	V65G (GIT→GGT)	<i>nagC</i> ←
<a href="#">JC JC</a>	969,836	IS150 (+) +3 bp	coding (810-812/2283 nt)	<i>pflB</i> ←
<a href="#">RA</a>	1,329,516	C→T	H33Y (CAC→TAC)	<i>topA</i> →
<a href="#">JC JC</a>	1,544,289	IS150 (-) +3 bp	coding (150-152/1536 nt)	<i>xasA</i> ←
<a href="#">JC JC</a>	1,733,647	IS150 (-) +3 bp	coding (683-685/1413 nt)	<i>pykF</i> →
<a href="#">RA</a>	1,976,879	T→G	intergenic (-57/-76)	<i>yedW</i> ← / → <i>yed</i>
<a href="#">RA</a>	2,082,685	G→A	A494V (GCT→GTT)	<i>yegl</i> ←
<a href="#">RA</a>	2,499,315	G→A	intergenic (-110/-179)	<i>maeB</i> ← / → <i>talA</i>
<a href="#">RA</a>	3,045,069	G→T	T312N (ACC→AAC)	<i>yghJ</i> ←
<a href="#">RA</a>	3,248,957	A→T	D764E (GAT→GAA)	<i>infB</i> ←
<b>MC JC</b>	3,289,962	Δ16 bp	coding (96-111/4554 nt)	<i>gltB</i> →
<a href="#">RA</a>	3,339,158	A→C	intergenic (+22/-4)	<i>yhdG</i> → / → <i>fis</i>
<a href="#">RA</a>	3,370,027	T→A	K117M (AAG→ATG)	<i>rpsM</i> ←
<a href="#">RA</a>	3,424,910	G→A	M11M (ATG→ATA) †	<i>nirC</i> →
<a href="#">RA</a>	3,483,047	C→A	R455S (CGC→AGC)	<i>malT</i> →
<a href="#">RA</a>	3,762,741	A→T	K662I (AAA→ATA)	<i>spoI</i> →
<a href="#">RA</a>	3,875,632	(T) <sub>7</sub> →8	intergenic (-66/+287)	<i>glmU</i> ← / ← <i>atpC</i>
<a href="#">RA</a>	3,893,551	+G	intergenic (+6/-50)	<i>kup</i> → / → <i>insJ-3</i>
<b>MC JC</b>	3,894,997	Δ6,934 bp	IS150-mediated	<i>rbsD</i> -[ <i>yieO</i> ]
<a href="#">RA</a>	4,100,655	C→T	M192I (ATG→ATA)	<i>hslU</i> ←
<a href="#">RA</a>	4,126,706	(T) <sub>8</sub> →7	coding (342/879 nt)	<i>pflC</i> →
<a href="#">RA</a>	4,560,632	T→C	Y131C (TAC→TGC)	<i>hsdM</i> ←

Predicted mutation						
evidence	seq id	position	mutation	annotation	gene	description
<a href="#">MC JC</a>	REL606	3,289,962	Δ16 bp	coding (96-111/4554 nt)	<i>gltB</i> →	glutamate synthase, large subunit

Missing coverage evidence...										
	seq id	start	end	size	←reads	reads→	gene	description		
<a href="#">*</a>	<a href="#">±</a>	<a href="#">+</a>	REL606	3289962	3289977	16	63 [1]	[0] 64	<i>gltB</i>	glutamate synthase, large subunit

New junction evidence									
seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotation	gene	product
<a href="#">2</a>	REL606 = 3289961	1 (0.020)	62 (1.050)	39/70	0.1	99.2%	coding (95/4554 nt)	<i>gltB</i>	glutamate synthase, large subunit
<a href="#">2</a>	REL606 3289978 =	0 (0.000)					coding (112/4554 nt)	<i>gltB</i>	glutamate synthase, large subunit

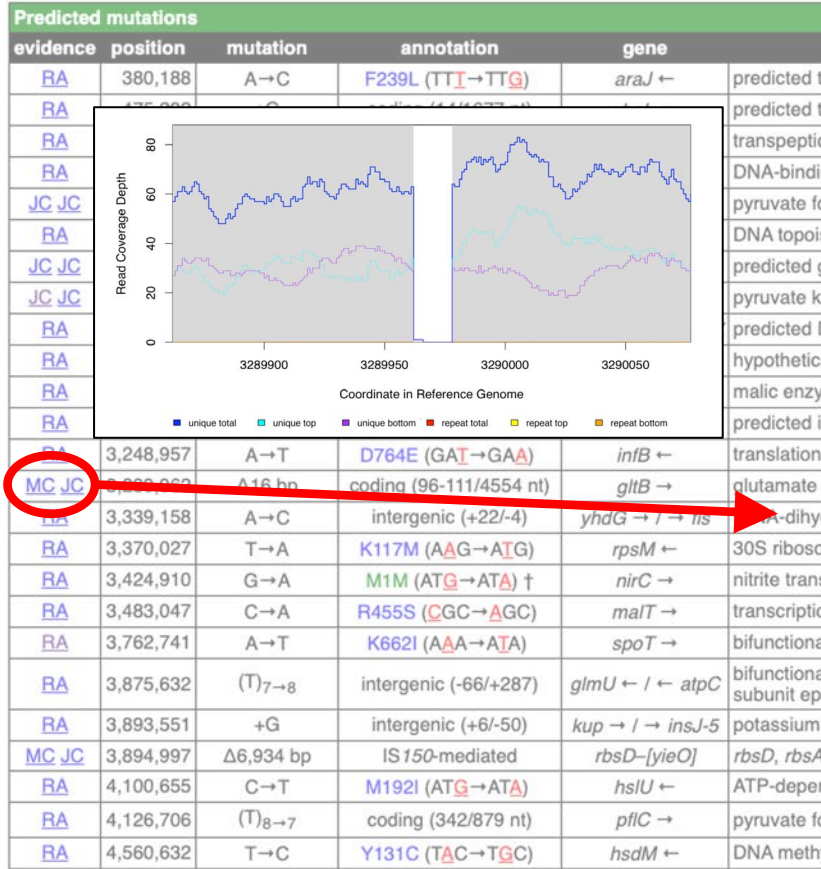


Unassigned missing coverage evidence								
seq id		start	end	size	←reads	reads→	gene	
<a href="#">*</a>	<a href="#">*</a>	REL606	546953–547700	555934–555877	8178–8982	20 [18]	[16] 19	<i>[insB-6]–[ECB_00513]</i>
<a href="#">*</a>	<a href="#">*</a>	REL606	2031675–2031718	2054970–2054943	23226–23296	21 [17]	[18] 21	<i>[manB]–[cpsG]</i>

Unassigned new junction evidence									
seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotation	gene	product
<a href="#">2</a>	REL606 = 547699	NA (NA)	80 (1.360)	37/70	0.2	NA	noncoding (1/768 nt)	<i>IS1</i>	repeat region
<a href="#">2</a>	REL606 555924 =	NA (NA)					coding (1209/2346 nt)	<i>ECB_00513</i>	conserved hypothetical protein

# JC = New Junction Evidence

breseq version 0.35.6 revision c7cf8df53bcd  
mutation predictions | marginal predictions | summary statistics | genome diff | commands



**Unassigned missing coverage evidence**

seq id	start	end	size	←reads	reads→	gene
* - ± REL606	546953-547700	555934-555877	8178-8982	20 [18]	[16] 19	[insB-6]-[ECB_00513] [insB-6],insA-6
* - ± REL606	2031675-2031718	2054970-2054943	23226-23296	21 [17]	[18] 21	[manB]-[cpsG] [manB],manC,i

**Unassigned new junction evidence**

seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotation
* ? REL606 = 547699	NA (NA)	80 (1.360)	37/70	0.2	NA	noncoding (1/768 nt)	
- ? REL606 555924 =	NA (NA)					coding (1209/2346 nt)	

**Predicted mutation**

evidence	seq id	position	mutation	annotation	gene	description
MC JC	REL606	3,289,962	Δ16 bp	coding (96-111/4554 nt)	glfB →	glutamate synthase, large subunit

**Missing coverage evidence...**

seq id	start	end	size	←reads	reads→	gene	description
* - ± REL606	3289962	3289977	16	63 [1]	[0] 64	glfB	glutamate synthase, large subunit

**New junction evidence**

seq id	position	reads (cov)	reads (cov)	score	skew	freq	annotation	gene	product
* ? REL606 = 3289961		1 (0.020)	62 (1.050)	39/70	0.1	99.2%	coding (95/4554 nt)	glfB	glutamate synthase, large subunit
* ? REL606 3289978 =		0 (0.000)					coding (112/4554 nt)	glfB	glutamate synthase, large subunit

GGGTCCCGCAGAGCCTGGGGAGGTTCCACGATATCTTTGAGAGGGATAACTGTGGTTCGGCCGATCGC

> REL606/3289927-3289961  
> REL606/3289978-3290012

< 2:3623510/36-1  
< 2:3380704/1-36  
< 2:3460319/1-36  
< 1:131137/1-36  
< 1:1353443/36-1  
< 1:3014326/36-1  
< 1:3346641/1-36  
< 1:689686/1-36  
< 2:657798/36-1  
< 1:2511982/36-1  
< 1:2025880/36-1  
< 2:3369418/36-1  
< 1:2431606/1-36  
< 2:2593722/1-36  
< 2:2252332/36-1  
< 2:892776/36-1  
< 1:2634884/1-36  
< 2:55338/36-1  
< 2:2034737/36-1  
< 2:16872/1-36  
< 2:223924/1-36  
< 1:187041/36-1  
< 1:1339127/36-1  
< 2:2339905/1-36  
< 2:2339574/1-36  
< 1:1082563/1-36  
< 1:2546169/1-36  
< 1:1973857/1-36  
< 2:1448864/1-36  
< 1:732312/1-36  
< 2:2707039/36-1  
< 1:2303795/1-36  
< 1:2120882/1-36  
< 1:372974/1-36  
< 2:2173522/36-1  
< 2:2610370/36-1  
< 2:2779388/36-1  
< 2:2438260/1-36

# Summary Statistics



breseq version 0.35.6 revision c7cf8df53bcd  
[mutation predictions](#) | [marginal prediction](#) | [summary statistics](#) | [genome diff](#) | [command line log](#)

## Read File Information

	read file	reads	bases	passed filters	average	longest	mapped
<a href="#">errors</a>	SRR030255_1	4,092,676	147,336,336	98.7%	36.0 bases	36 bases	95.3%
<a href="#">errors</a>	SRR030255_2	4,103,100	147,711,600	98.9%	36.0 bases	36 bases	93.9%
	total	8,195,776	295,047,936	98.8%	36.0 bases	36 bases	94.6%

## Reference Sequence Information

		seq id	length	fit mean	fit dispersion	% mapped reads	description
<a href="#">coverage</a>	<a href="#">distribution</a>	REL606	4,629,812	60.6	3.1	100.0%	Escherichia coli strain REL606.
		total	4,629,812			100.0%	

fit dispersion is the ratio of the variance to the mean for the negative binomial fit. It is =1 for Poisson and >1 for over-dispersed data.

## New Junction Evidence

### Junction Candidates Tested

option	limit	actual
Number of alignment pairs examined for constructing junction candidates	$\leq 100000$	100047
Coverage evenness (position-hash) score of junction candidates	$\geq 2$	$\geq 2$
Test this many junction candidates (n). May be smaller if not enough passed the coverage evenness threshold	$100 \leq n \leq 5000$	60
Total length of all junction candidates (factor times the reference genome length)	$\leq 0.1$	0.001

### Junction Skew Score Calculation

reference sequence	pr(no read start)
REL606	0.48689

# Reference Sequence Coverage



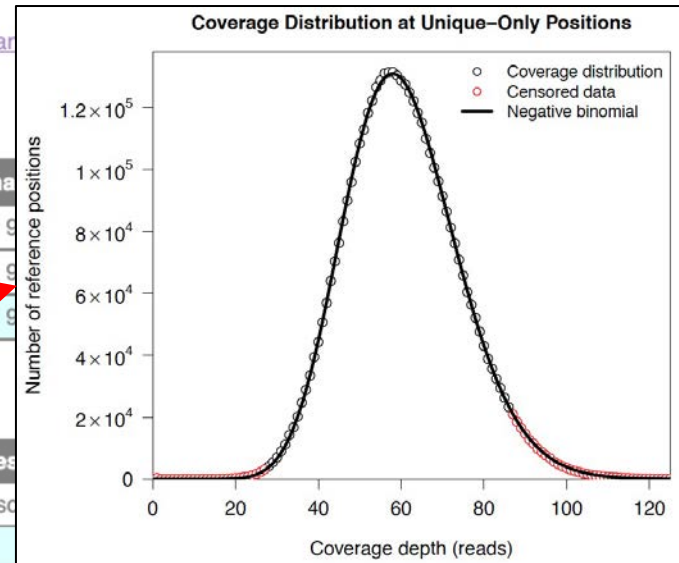
**breseq** version 0.35.6 revision c7cf8df53bcd  
[mutation predictions](#) | [marginal predictions](#) | [summary statistics](#) | [genome diff](#) | [command line](#)

## Read File Information

	read file	reads	bases	passed filters	average	longest	max
<a href="#">errors</a>	SRR030255_1	4,092,676	147,336,336	98.7%	36.0 bases	36 bases	9
<a href="#">errors</a>	SRR030255_2	4,103,100	147,711,600	98.9%	36.0 bases	36 bases	9
	<b>total</b>	<b>8,195,776</b>	<b>295,047,936</b>	<b>98.8%</b>	<b>36.0 bases</b>	<b>36 bases</b>	<b>9</b>

## Reference Sequence Information

	seq id	length	fit mean	fit dispersion	% mapped reads	des
<b>coverage</b>	REL606	4,629,812	60.6	3.1	100.0%	Esc
	<b>total</b>	<b>4,629,812</b>			<b>100.0%</b>	



fit dispersion is the ratio of the variance to the mean for the negative binomial fit. It is =1 for Poisson and >1 for over-dispersed data.

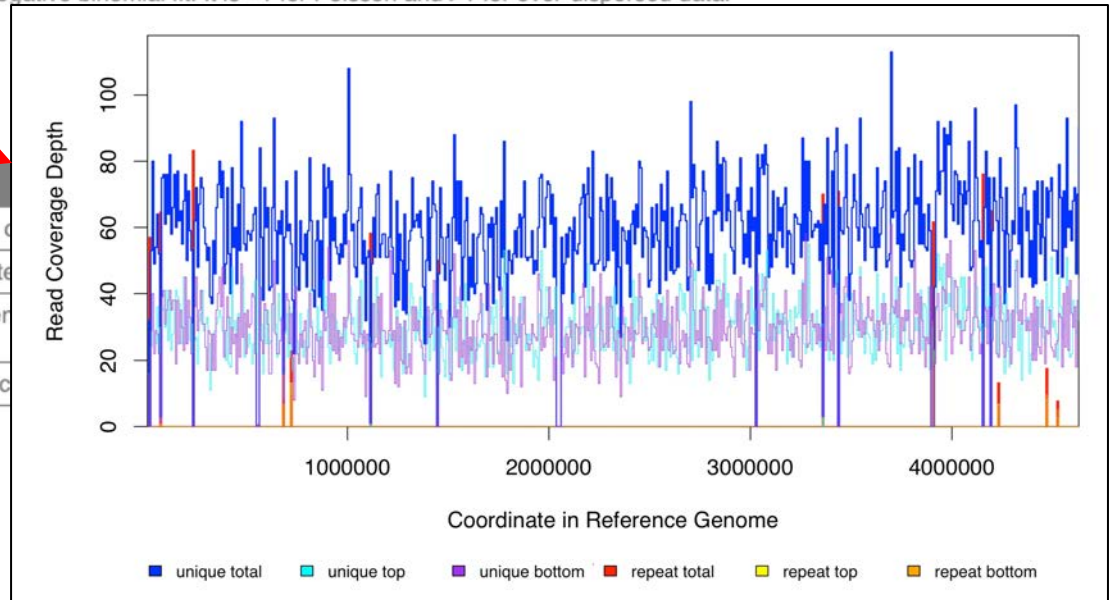
## New Junction Evidence

### Junction Candidates Tested

option
Number of alignment pairs examined for constructing junction candidates
Coverage evenness (position-hash) score of junction candidate
Test this many junction candidates (n). May be smaller if not enough threshold
Total length of all junction candidates (factor times the reference)

### Junction Skew Score Calculation

reference sequence	pr(no read start)
REL606	0.48689



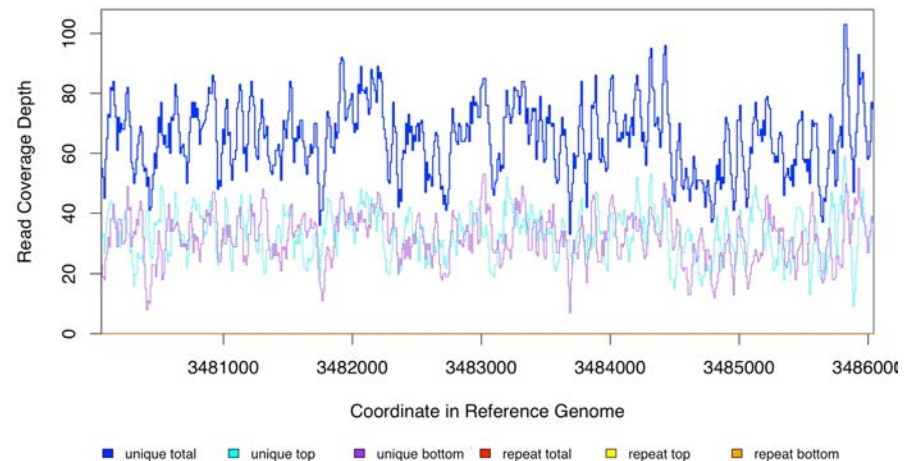
# Utilities to explore output

You can run utility subcommands from inside the main output directory of a *breseq* run. `$ breseq --help` to see others.

```
$ breseq BAM2ALN
-o alignment.html
REL606:3483047-3483047
```

```
AAGACACCATGCACGCAGAAATTAACGCTCTGCGCGCCCAAGTGCCGATTAAAGATGGTAATCCG > REL606/3483015-3483079
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg < 1:2369690/36-1 (MQ=255)
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg > 1:577628/1-36 (MQ=255)
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg > 2:1772887/1-36 (MQ=255)
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg < 1:130379/36-1 (MQ=255)
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg < 2:3079501/36-1 (MQ=255)
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg > 1:1820887/1-36 (MQ=255)
aagaCACCATGCACGCAGAAATTAACGCTCTGgpcg > 1:2369308/36-1 (MQ=255)
agaCACCATGCACGCAGAAATTAACGCTCTGgpcg > 2:3469595/1-36 (MQ=255)
agaCACCATGCACGCAGAAATTAACGCTCTGgpcg < 2:1489970/36-1 (MQ=255)
cacCATGCACGCAGAAATTAACGCTCTGgCGCCCa > 1:1927484/1-36 (MQ=255)
cacCATGCACGCAGAAATTAACGCTCTGgCGCCCa > 2:2734863/36-1 (MQ=255)
cacCATGCACGCAGAAATTAACGCTCTGgCGCCCa < 2:2587112/36-1 (MQ=255)
cacCATGCACGCAGAAATTAACGCTCTGgCGCCCa < 2:1926447/36-1 (MQ=255)
acCATGCACGCAGAAATTAACGCTCTGgCGCCCAg < 2:885743/36-1 (MQ=255)
ccATGCACGCAGAAATTAACGCTCTGgCGCCCAg < 2:2448233/1-36 (MQ=255)
ccATGCACGCAGAAATTAACGCTCTGgCGCCCAg < 1:3403951/36-1 (MQ=255)
ccATGCACGCAGAAATTAACGCTCTGgCGCCCAg > 2:3361806/1-36 (MQ=255)
cATGCACGCAGAAATTAACGCTCTGgCGCCCAg > 2:3230993/1-36 (MQ=255)
aTGACGCAGAAATTAACGCTCTGgCGCCCAg < 2:1743516/36-1 (MQ=255)
aTGACGCAGAAATTAACGCTCTGgCGCCCAg < 2:3672937/36-1 (MQ=255)
aTGACGCAGAAATTAACGCTCTGgCGCCCAg > 1:3325866/1-36 (MQ=255)
aTGACGCAGAAATTAACGCTCTGgCGCCCAg < 1:3348771/36-1 (MQ=255)
tGCACGCAGAAATTAACGCTCTGgCGCCCAg < 2:3403193/36-1 (MQ=255)
tGCACGCAGAAATTAACGCTCTGgCGCCCAg < 2:1611056/1-36 (MQ=255)
gCACCgTAATTAACGCTCTGgCGCCCAg < 1:2589008/1-36 (MQ=38)
taCGCAGAAATTAACGCTCTGgCGCCCAg < 1:2979881/35-1 (MQ=25)
```

```
$ breseq BAM2COV
-o coverage.png
REL606:3480047-3486047
```



These can help with identifying copy number changes (e.g., duplications) and understanding complex structural variation.

# Explore aligned reads using IGV

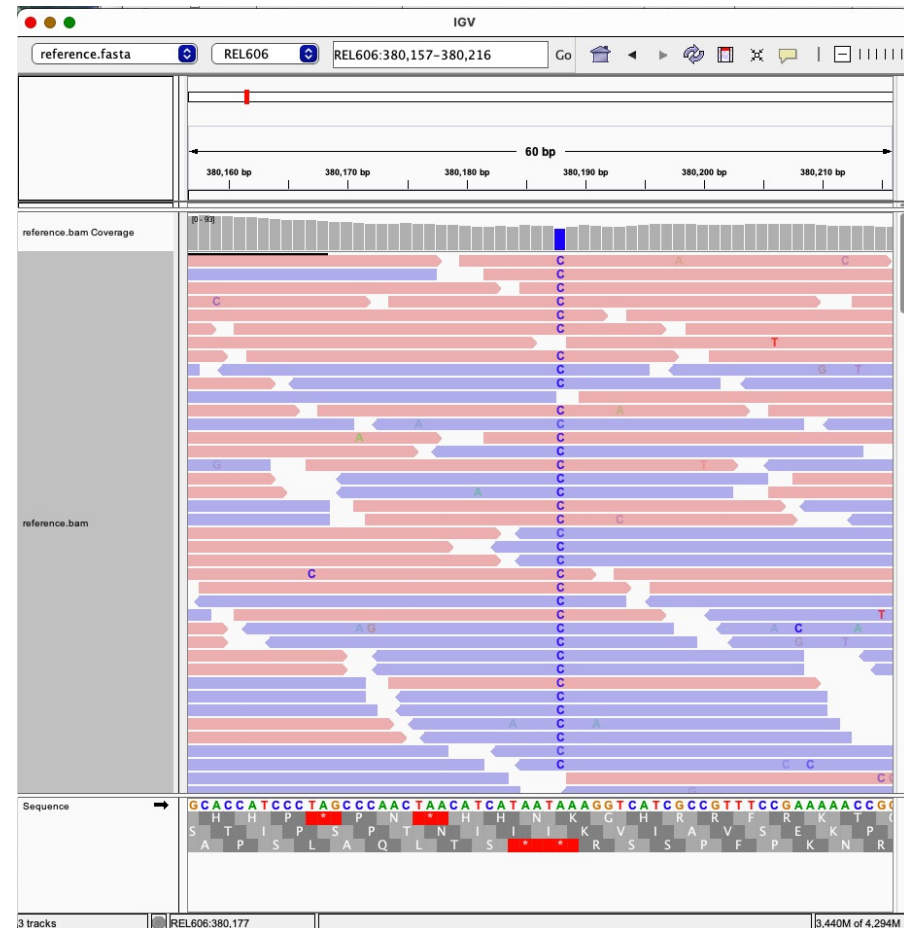


<https://software.broadinstitute.org/software/igv/>

## Viewing Output / Aligned Reads in the IGV

You can visualize the “raw data” (how **breseq** aligned reads to the reference genome) using the **Integrative Genomics Viewer (IGV)** and files located in the **data** folder created by **breseq**.

1. Install and open IGV
2. Import the reference genome sequence:
  - Click ‘File’, and then ‘Import Genome...’
  - Fill out the requested information: ‘ID’, ‘Name’
  - Choose the FASTA file: `data/reference.fasta`.
  - The other fields are optional.
3. Import the reference genome feature information:
  - Click ‘File’, and then ‘Load from File...’
  - Choose the GFF3 file: `data/reference.gff3`.
4. Import the read alignments to the reference genome:
  - Click ‘File’, and then ‘Load from File...’
  - Choose the BAM file: `data/reference.bam`.



# GenomeDiff output

## Machine-readable text files for further processing

```
#=GENOME_DIFF      1.0
#=CREATED      15:16:00 24 May 2021
#=PROGRAM      breseq 0.35.6 revision c7cf8df53bcd
#=COMMAND      breseq -j 8 -o tests/long_Ara-1_10000gen_4536A ...
#=REFSEQ       tests/long_Ara-1_10000gen_4536A/../../data/long_tests/REL606.gbk
#=READSEQ      tests/long_Ara-1_10000gen_4536A/../../data/long_tests/SRR030255_1.fastq.gz
#=READSEQ      tests/long_Ara-1_10000gen_4536A/../../data/long_tests/SRR030255_2.fastq.gz
#=CONVERTED-BASES 295047936
#=CONVERTED-READS 8195776
#=INPUT-BASES    298701576
#=INPUT-READS    8297266
#=MAPPED-BASES   277772336
#=MAPPED-READS   7750270
SNP   1      29      REL606      380188      C
INS   2      32      REL606      475292      G
SNP   3      36      REL606      649391      A
SNP   4      37      REL606      683496      C
MOB   5      101,102      REL606      969836      IS150 1      3
SNP   6      41      REL606      1329516      T
MOB   7      103,109      REL606      1544289      IS150 -1      3
MOB   8      110,111      REL606      1733647      IS150 -1      3
SNP   9      46      REL606      1976879      G
SNP  10      49      REL606      2082685      A
...
```

GenomeDiff format  
output/output.gd

Format specification provided in the *breseq* manual

# What can you do with a GenomeDiff?

Generate an HTML table comparing multiple clones/populations:

```
$ gdttools COMPARE -o compare.html -r reference.gbk input1.gd input2.gd ...
```

Convert to TSV, VCF or other formats for interchange with other programs:

```
$ gdttools ANNOTATE -o -f TSV -r reference.gbk input1.gd input2.gd ...
```

Count mutations and numbers of sites at risk for mutations:

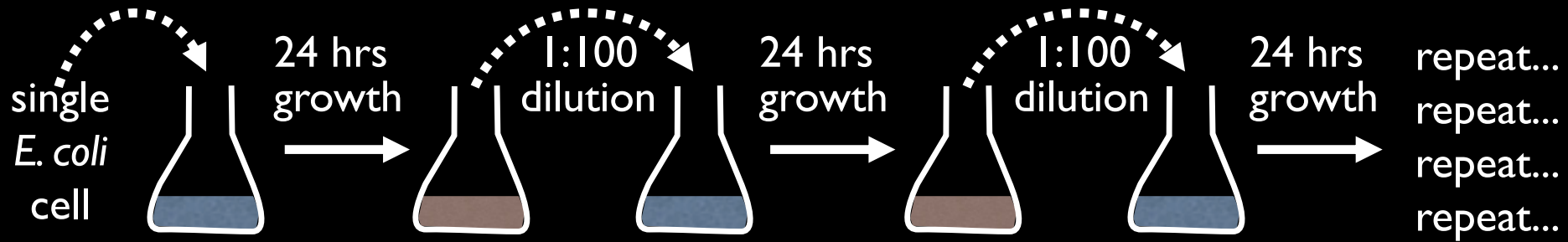
```
$ gdttools COUNT -o output.csv -r reference.gbk input1.gd input2.gd ...
```

Apply the mutations to generate an updated reference sequence:

```
$ gdttools APPLY -f GENBANK -o updated.gbk -r reference.gbk input.gd
```

And more... `$ gdttools --help`

# Lenski Long-Term Evolution Experiment



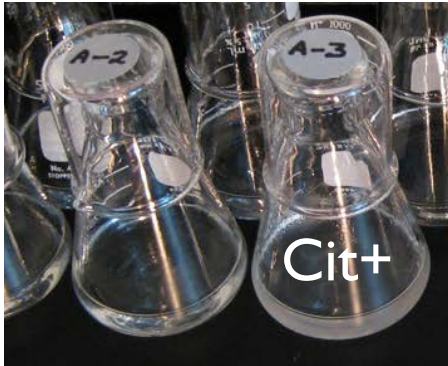
Richard Lenski  
Michigan State

- ❖ 12 independent populations
- ❖ Deep evolutionary history
- ❖ Viable frozen "fossil record"

>73,000 generations of  
*E. coli* growth (>30 years)!

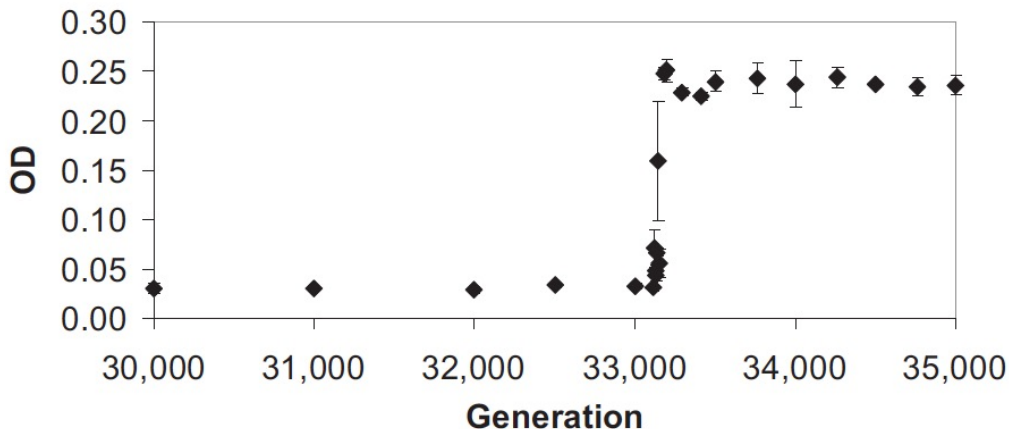


# Analysis: Causative Mutations

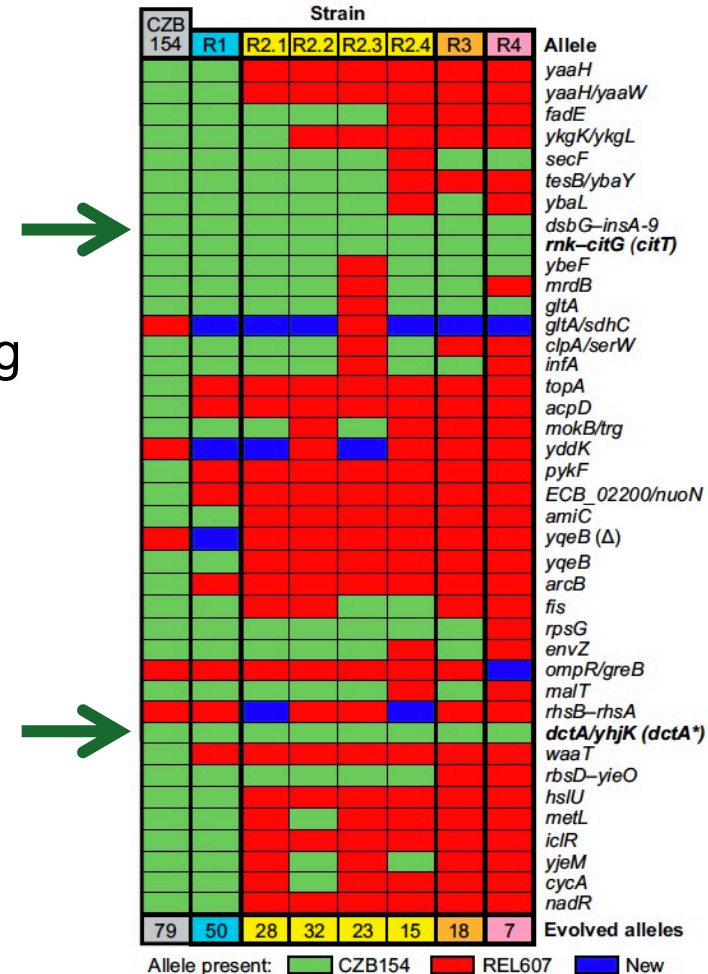


How?

Backcross and  
sequence: Only  
**two mutations**  
required for strong  
Cit<sup>+</sup> phenotype

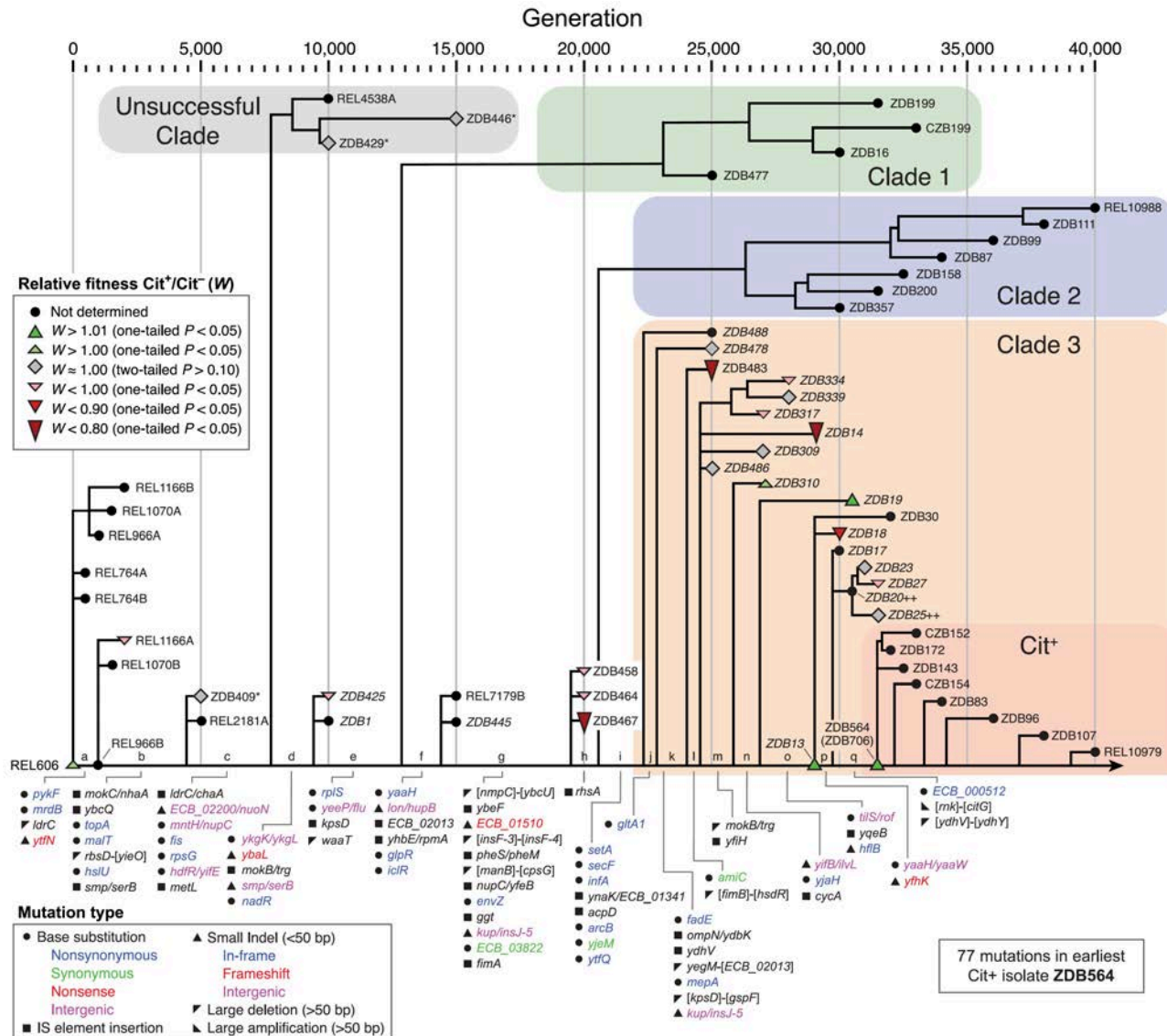


Citrate utilization evolved after 33,000 generations in one LTEE population



\$ gdttools COMPARE...

# Analysis: Phylogenetic trees



What mutations led to Cit+ evolution?

Generate an alignment of genomic changes

```
$ gdttools COMPARE
-f PHYLIP clone1.gd
clone2.gd ....
```

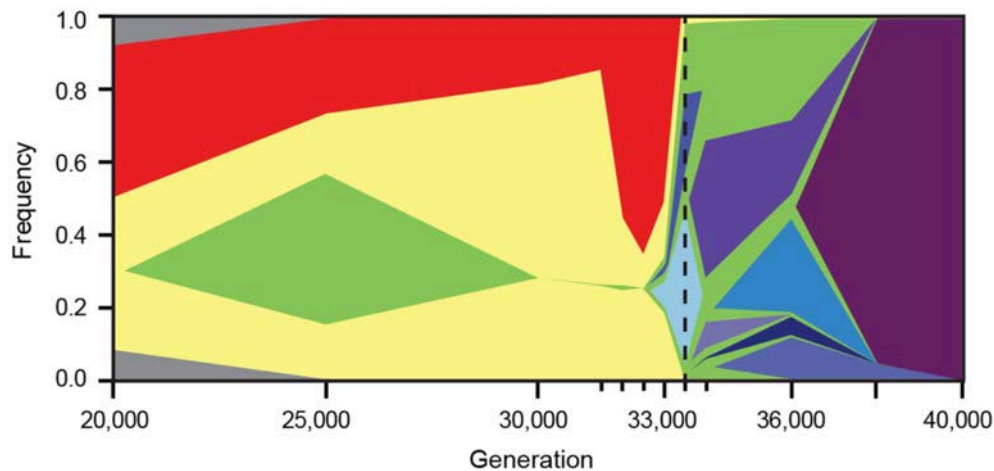
or

```
$ gdttools COMPARE
-f FASTA clone1.gd
clone2.gd ....
```

Build and visualize a maximum parsimony tree using PHYLIP, MEGAX, etc.

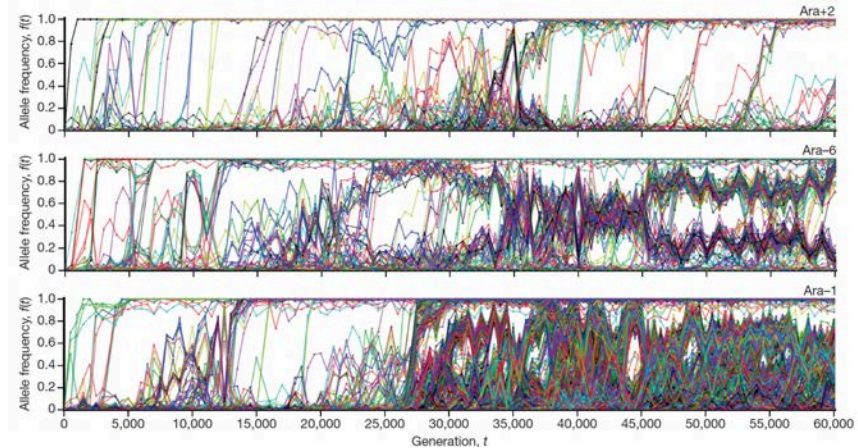
# Analysis: Allele/Genotype Frequencies

Muller Plot (Genotype Frequency)



Quandt *et al.* (2015) *eLife*

Allele Frequency



Good *et al.* (2017) *Nature*

For tracking how genetic diversity evolves within populations, visualizing dynamics, selective sweeps, and stable coexistence.

```
gdtools COMPARE -f CSV pop1.gd pop2.gd ....
```

Programs/packages that can help:

R, ggplot, ggMuller, EvoFreq, MullerPlot

# Workshop Presentations

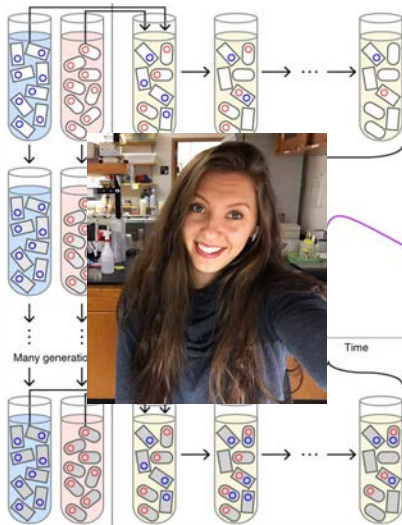
JLS	K0001	K0022	K0014	K0006	K0030	K0031	
Position	A1-1	S1-1	G1-1	E1-1	G1-1	G1-1	Mutation
56,273							G→A
62,682							+G
105,411							G→A
156,056							T→G
444,525							T→G
490,544							Δ6 bp
490,998							Δ1 bp
520,989							
556,778							
572,671							
666,783							
683,143							
755,210							
909,411							
991,031							
1,213,665							
1,218,024							
1,337,160							
1,349,606							



Antibiotic Resistance Reversal:  
*breseq* Analysis of Experimental  
Evolution, Compared with FACS  
Competition Assays of Relative Fitness

**Joan Slonczewski**

*Kenyon College*



Identifying Adaptive Paths in Host-  
Plasmid Coevolution Using *breseq*

**Olivia Kosterlitz**

*University of Washington*

# Workshop Presentations



EVOLVINGSTEM



Decoding Evolution-In-Action in  
Classroom Experiments That Simulate  
Infection Biology Using *breseq*

**Vaughn Cooper**

*University of Pittsburgh*



ALEdb: A Living High-Quality Database  
of Mutations from Adaptive Evolution  
Experiments Powered by *breseq*

**Adam Feist**

*University of California, San Diego*



## Table of Contents

### Tutorial: Population Samples (Polymorphism Mode)

- 1. Download data files
  - Reference sequence
  - Read files
- 2. Run **breseq** with default filters
- 3. Run **breseq** with no filters
- 4. Compare predictions of mutations
- 5. Examine allele frequency time courses

## Previous topic

Tutorial: Clonal Samples (Consensus Mode)

## Next topic

Tutorial: Ultra-rare variant detection using consensus reads and targeted sequencing

## This Page

Show Source

## Quick search

# Tutorial: Population Samples (Polymorphism Mode)

In this exercise, you will analyze two population (metagenomic) samples using **breseq** to track the frequencies of evolved alleles and changes in genetic diversity in population Ara-3 of the Lenski long-term evolution experiment (LTEE). As discussed in [Tutorial: Clonal Samples \(Consensus Mode\)](#) this population evolved citrate utilization after 31,500 generations.

breseq 0.35.4 documentation » Tutorial: Clonal Samples (Consensus Mode)

previous | next | index



## Table of Contents

### Tutorial: Clonal Samples (Consensus Mode)

- 1. Download data files
  - Reference sequence
  - Read files
- 2. Run **breseq**
- 3. Open **breseq** output
- 4. Resolving the Cit+ mutation
  - A. *mk-citG* junction
  - B. Zoomed-in coverage
  - C. Add the amplification to the *GenomeDiff* file
- 5. Generating a mutated reference sequence
- 6. Characterizing genetic diversity and genome evolution
  - Example 1. Compare mutations in different genomes
  - Example 2. Analyze rates and nature of genome evolution

# Tutorial: Clonal Samples (Consensus Mode)

This tutorial expands on the [Test Drive](#). You will analyze mutations in the genomes of multiple clones isolated from population Ara-3 of the Lenski long-term evolution experiment (LTEE). A complex mutation is present in these samples that was necessary for evolution of the aerobic citrate utilization trait (Cit+). In addition to some tips on **breseq** usage and examples of interpreting more complex mutations in the output, this tutorial also introduces functionality in the **gdtools** utility command that can be used to compare and analyze mutations in an entire set of evolved genomes.

**Note:** This tutorial was created for the EMBO Practical Course [Measuring intra-species diversity using high-throughput sequencing](#) held 27–31 July 2015 in Oeiras, Portugal.

**Warning:** If you encounter any **breseq** or **gdtools** errors or crashes in running this tutorial, please [report issues on GitHub](#).

## 1. Download data files

First, create a directory called `tutorial_clonal`:

```
$ mkdir tutorial_clones
$ cd tutorial_clones
```

## Reference sequence

**breseq** prefers the reference sequence in [Genbank](#) or [GFF3](#) format. In this example, the

# Let us know how we can help!

These slides can be downloaded at <http://barricklab.org/breseq>



## breseq Workshop Survey

We would like to plan one or more interactive virtual sessions to help you use breseq to analyze your data.

<https://forms.gle/qkvkjbqCXZAhY7GW6>

## Interactive Workshop

- Install on your system
- Use on your data
- Help interpret output
- Provide advice on further analysis

## Post bug reports and issues on GitHub

Please check that you are using the newest *breseq* version first!

A screenshot of the GitHub repository page for 'barricklab / breseq'. The page shows the repository name at the top, followed by statistics: 22 Unwatched, 75 Unstars, and 11 Forks. Below this is a navigation bar with tabs for Code, Issues (31), Pull requests (1), Actions, Projects, Wiki, Security, and Insights. The 'Issues' tab is selected, showing a list of 31 open issues. The first issue is titled 'Advice with annotating \*.gd file with deletions and SNPs' and was opened on Jan 29 by lthomp06. The second issue is 'How someone can concatenate the info of syn/non.syn mutations to the predicted'. The page also includes filters for 'is:issue is:open' and buttons for 'Labels' (19) and 'Milestones' (0). A 'New issue' button is visible in the top right corner of the issues list.

# Acknowledgments

## Breseq Developers



Dan Deatherage

David Knoester

Geoffrey Colburn

Matt Strand

Jordan Borges

Aaron Reba

## Funding

NIH K99/R00  
(GM087550)

NSF CAREER  
(CBET-1554179)

NSF BEACON Center  
(DBI-0939454)

Thanks to many *breseq* users and research collaborators who have given feedback over the past decade!

Including Richard Lenski, Dominique Schneider, Olivier Tenaillon, Vaughn Cooper, Michael Desai, Yousif Shamoo, Zachary Blount, Genoscope, the Gulbenkian Institute, and members of these and many other research groups and communities.