

# Genome Assembly Tutorial

Deb Triant

Based on the tutorial by Mike Schatz

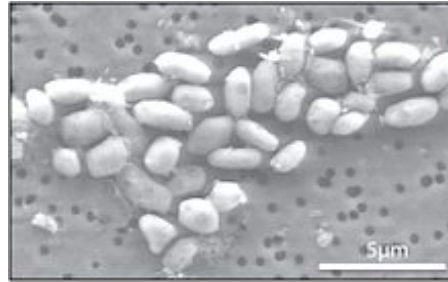
18 October, 2015  
Programming for Biology  
Cold Spring Harbor Labs



## Outline

1. Genome assembly with ALLPATHS-LG
2. Genome alignment with MUMmer

## Sample Data: *Halomonas* sp. GFAJ-I



### Library 1: Fragment

Avg Read length: 100bp

Insert length: 180bp

### Library 2: Short jump

Avg Read length: 50bp

Insert length: 2000bp

### **A Bacterium That Can Grow by Using Arsenic Instead of Phosphorus**

Wolfe-Simon et al (2010) *Science*. 332(6034)1163-1166.

3

## Files needed for exercises:

- **File location:** /home/pfb2015/data/GenomeAssembly\_data/asm
- **FASTQ files:**
  - Paired-end:** Illumina 2x100 reads from 180 +/- 20 bp fragments
    1. frag180.1.fq
    2. frag180.2.fq
  - Mate-pair/jumping:** Illumina 2x50 reads from 2000 +/- 200 bp fragments
    1. jump2k.1.fq
    2. jump2k.2.fq
- **Reference file**
  - ref.fa
- **ALLPATHS input files**
  - in\_groups.csv - name and location of files\*
  - in\_libs.csv - library details\*
  - \*names in both files must match

4

## Logging in for graphical display

- `ssh -i Student-PFB2015.pem -X deb@compute.programmingforbiology.org`
- Exit
- login again with:  
`ssh -i Student-PFB2015.pem -AX deb@compute.programmingforbiology.org`
- Type “xeyes” at terminal to test. Ctrl-C to end xeyes

5

## Sample fastq file

## Paired-end


[illegible]

## Mate-pair/Jumping

```
@ /1/
CGTACGCCATCATGATGCCCATACCCAGCGACAGCGTGGAAGAAAGCCTGG
+
555555666666666677777788888899999.....<=====>?
@2/1
CCACGGGTCAAGAAAGAACAGCTGCAGCGTAGTCAGCGAGCCCTCTGGCA
+
555555666666666677777788888899999.....<=====>?
```

6

**Running ALLPATHS-LG**  
**Iain MacCallum**



### How to use ALLPATHS-LG

---

1. **Data requirements (\*\* most critical \*\*)**
2. Computational requirements & Installation
3. Preparing your data
4. Assembling
5. What is an ALLPATHS-LG assembly?

### ALLPATHS-LG sequencing model

Libraries (insert types)	Fragment size (bp)	Read length (bases)	Sequence coverage (x)	Required
Fragment	180*	$\geq 100$	45	yes
Short jump	3,000	$\geq 100$ preferable	45	yes
Long jump	6,000	$\geq 100$ preferable	5	no**
Fosmid jump	40,000	$\geq 26$	1	no**

\*See next slide.

\*\*For best results. Normally not used for small genomes.  
However essential to assemble long repeats or duplications.

Cutting coverage in half still works, with some reduction in quality of results.

All: protocols are either available, or in progress.

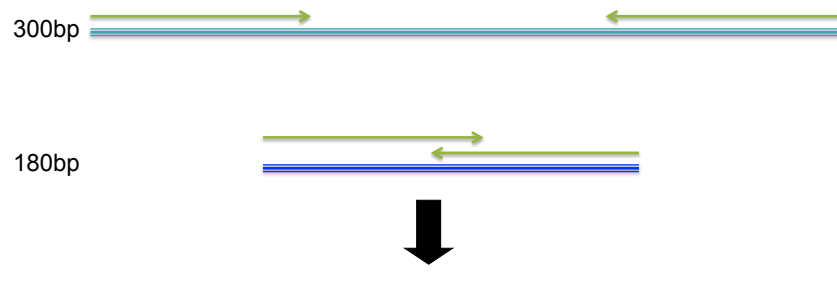
9

### Paired-end libraries from 180 bp fragments

Requires short insert separation – less than twice the read length so that reads may overlap.

ALLPATHS tries to close gap in overlapping reads to generate longer reads that are twice as long.

\* 300bp with inward orientation – 100bp reads

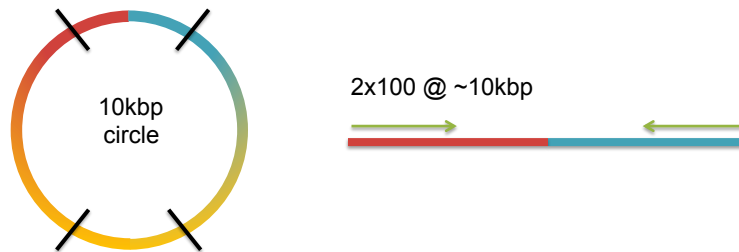


For longer reads, fragment size would be increased proportionally.

10

**Mate-pair sequencing**

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



11

**How to use ALLPATHS-LG**

1. Data requirements
- 2. Computational requirements & installation**
3. Preparing your data
4. Assembling
5. What is an ALLPATHS-LG assembly?

12

### Computational requirements

---

- **64-bit Linux**
- **runs multi-threaded on a single machine**
- **memory requirements**
  - about 160 bytes per genome base, implying
    - need 512 GB for mammal (Dell R315, 48 processors, \$39,000)
    - need 1 GB for bacterium (theoretically)
  - if coverage different than recommended, adjust...
  - potential for reducing usage
- **wall clock time to complete run**
  - 5 Mb genome → 1 hour (8 processors)
  - 2500 Mb genome → 500 hours (48 processors)

13

### Installing ALLPATHS-LG

---

Web page:

<http://www.broadinstitute.org/software/allpaths-lg/blog/>

General instructions:

<http://www.broadinstitute.org/science/programs/genome-biology/computational-rd/general-instructions-building-our-software>

**USE THE MANUAL!!!**

14

### Getting the ALLPATHS-LG source

---

**\*\*Do not need to do for course exercise.....just for reference\*\***

Our current system is to release code daily if it passes a test consisting of several small assemblies:

Download the latest build from:

<ftp://ftp.broadinstitute.org/pub/crd/ALLPATHS/Release-LG/>

Unpack it:

```
% tar xzf allpathslg-39099.tar.gz
```

(substitute the latest revision id for 39099)

This creates a source code directory `allpathslg-39099`:

```
% cd allpathslg-39099
```

15

### Building ALLPATHS-LG

---

**Step one:** `./configure`

Options:

```
-prefix=<prefix path>  
  put binaries in <prefix path>/bin, else ./bin
```

**Step two:** `make` and `make install`

Options:

```
-j<n>  
  compile with n parallel threads
```

**Step three:** add bin directory to your path

16



### How to use ALLPATHS-LG

---

1. Data requirements
2. Computational requirements & Installation
3. **Preparing your data** - `PrepareAllPathsInputs.pl`
4. Assembling
5. What is an ALLPATHS-LG assembly?

17

### Preparing data for ALLPATHS-LG

---

Before assembling, prepare and import your read data.

ALLPATHS-LG expects reads from:

- At least one fragment library.
  - One should come from fragments of size ~180 bp.
  - This is NOT checked but otherwise results will be bad.
- At least one long-insert library.

**IMPORTANT:** use all the reads, including those that fail the Illumina purity filter. These low quality reads may cover 'difficult' parts of the genome.

Trimming reads?

18



### Input files – in\_groups.csv

---

Each line in **in\_groups.csv** comma separated value file, corresponds to a BAM or FASTQ file you wish to import for assembly.

The library name must match the names in **in\_libs.csv**.

**group\_name** - a unique nickname for this file  
**library\_name** - library to which the file belongs  
**file\_name** - the absolute path to the file  
 (should end in .bam or .fastq)  
 (use wildcards '?', '\*' for paired fastqs)

Example:

```
group_name, library_name, file_name
302GJ, Solexa-11541, /seq/Solexa-11541/302GJABXX.bam
303GJ, Solexa-11623, /seq/Solexa-11623/303GJABXX.?.fastq
```

### Input files – in\_groups.csv

---

Each line in **in\_groups.csv** comma separated value file, corresponds to a BAM or FASTQ file you wish to import for assembly.

The library name must match the names in **in\_libs.csv**.

**group\_name** - a unique nickname for this file  
**library\_name** - library to which the file belongs  
**file\_name** - the absolute path to the file  
 (should end in .bam or .fastq)  
 (use wildcards '?', '\*' for paired fastqs)

Example:

```
group_name, library_name, file_name
302GJ, Solexa-11541, /seq/Solexa-11541/302GJABXX.bam
303GJ, Solexa-11623, /seq/Solexa-11623/303GJABXX.?.fastq
```

## Libraries - in\_libs.csv

### EXAMPLE

library_name	type	paired	frag_size	frag_stddev	insert_size	insert_stddev	read_orientation	genomic_start	genomic_end
Solexa-11541	fragment	1	180	10			inward		
Solexa-11623	jumping	1			3000	500	outward	0	25

Pay attention to commas and empty fields

23

## Libraries - in\_libs.csv

### For fragment libraries only

**frag\_size** - estimated mean fragment size  
**frag\_stddev** - estimated fragment size std dev

### For jumping libraries only

**insert\_size** - estimated jumping mean insert size  
**insert\_stddev** - estimated jumping insert size std dev

These values determine how a library is used. If **insert\_size** is  $\geq 20,000$ , the library is assumed to be a Fosmid jumping library.

**paired** - always 1 (only supports paired reads)  
**read\_orientation** - inward or outward.

Paired reads can either point towards each other, or away from each other. Currently fragment reads must be **inward**, jumping reads **outward**, and Fosmid jumping reads **inward**.

24

## Libraries – in\_libs.csv

Reads can be trimmed to remove non-genomic bases produced by the library construction method:

**genomic\_start** - inclusive zero-based range of read bases  
**genomic\_end** to be kept; if blank or 0 keep all bases

Reads are trimmed in their original orientation.

Extra optional fields (descriptive only – ignored by ALLPATHS)

**project\_name** - a string naming the project.  
**organism\_name** - the organism name.  
**type** - fragment, jumping, EcoP15I, etc.

### EXAMPLE

library_name	type	paired	frag_size	frag_stddev	insert_size	insert_stddev	read_orientation	genomic_start	genomic_end
Solexa-11541	fragment	1	180	10			inward		24
Solexa-11623	jumping	1			3000	500	outward	0	25

## How to import assembly data files

### PrepareAllPathsInputs.pl

```
IN_GROUPS_CSV=<in groups file>
IN_LIBS_CSV=<in libs file>
DATA_DIR=<full path of data directory> *Required
PLOIDY=<ploidy, either 1 or 2>
PICARD_TOOLS_DIR=<picard tools directory> *in your
path - do not need to edit
```

- **IN\_GROUPS\_CSV** and **IN\_LIBS\_CSV**:  
/absolute\_path/in\_groups.csv and /absolute\_path/in\_libs.csv. These arguments determine where the data are found.
- **DATA\_DIR**: imported data will be placed here. Want to name your DATA\_DIR with the same pathway as will use in **RunAllPathsLG** for assembly.
- **PLOIDY**: either 1 (for a haploid or inbred organism), or 2 (for a diploid organism).
- **PICARD\_TOOLS\_DIR**: path to Picard tools, for data conversion from BAM. Do not need to specify if in path.

26

### Putting it all together

---

1. Create your input files:

**in\_libs.csv** - data file to describe your libraries and a  
**in\_groups.csv** - data file to describe your data files.

2. Prepare input files

```
% cd /path/to/data/
% PrepareAllPathsInputs.pl \
  DATA_DIR=data_dir PLOIDY=1 >& prepare.log
```

Want to name your DATA\_DIR with the same pathway as will use in **RunAllPathsLG** for assembly: <PRE>/<REFERENCE>/<DATA>/<RUN>/ASSEMBLIES/test

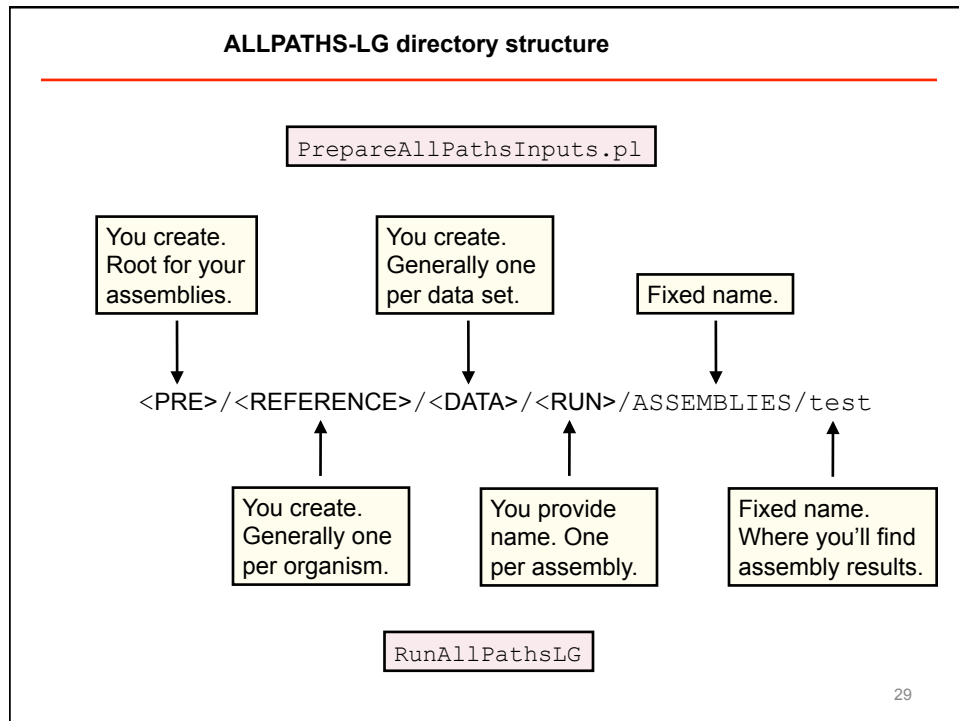

 CSHL/Halomonas/mydata/      #where to save your data

```
% mkdir -p CSHL/Halomonas/mydata # -p for making parent dirs
Our data:
% PrepareAllPathsInputs.pl \
  DATA_DIR=abs_path/CSHL/Halomonas/mydata/ PLOIDY=1 >&
  prepare.log #make sure you run in dir containing input files
```

### How to use ALLPATHS-LG

---

1. Data requirements
2. Computational requirements & installation
3. Preparing your data
4. **Assembling - RunAllPathsLG**
5. What is an ALLPATHS-LG assembly?



## How to assemble

### To run assembly

```
% RunAllPathsLG \
PRE=<prefix path> \
REFERENCE_NAME=<reference dir> \
DATA_SUBDIR=<data dir> \
RUN=<run dir>
```

**Automatic resumption.** If the pipeline crashes, fix the problem, then run the same `RunAllPathsLG` command again. Execution will resume where it left off. `% OVERWRITE=True`

What can go wrong: Not enough RAM, CPU time, artifacts or contamination in your data. Log files very informative!

Results. The assembly files are:

```
final.contigs.fasta      - fasta contigs
final.contigs.efasta     - efasta contigs
final.assembly.fasta     - scaffolded fasta
final.assembly.efasta    - scaffolded efasta
```

### Putting it all together

---

1. Collect the BAM or FASTQ files that you wish to assemble. Create a `in_libs.csv` metadata file to describe your libraries and a `in_groups.csv` metadata file to describe your data files.

2. Prepare input files

```
% PrepareAllPathsInputs.pl \
  DATA_DIR=abs_path/CSHL/Halomonas/mydata/ PLOIDY=1 >&
  prepare.log **This is where PrepareAllPathsInputs.pl
  saves our input files: CSHL/Halomonas/mydata/
```

3. Assemble.

```
% RunAllPathsLG \
  PRE= REFERENCE_NAME= \
  DATA_SUBDIR= RUN= THREADS=4 >& run.log
```

#### Our assembly example:

```
RunAllPathsLG \
  PRE=CSHL REFERENCE_NAME=Halomonas \
  DATA_SUBDIR=mydata RUN=run-1 THREADS=1 >& run-1.log
```

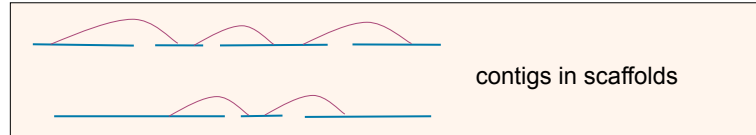
### How to use ALLPATHS-LG

---

1. Data requirements
2. Computational requirements & installation
3. Preparing your data
4. Assembling
5. What is an ALLPATHS-LG assembly?



## 1. Linear assemblies



**contig:** a contiguous sequence of bases....

```
CTGCCCCCTGTGCCAATGGGTTTGAGGCTCTCCCACTTCCTTTTCTATTAGATTCAATGTATCTGGTTTTATGTTGAGG
TCCTAGATCCACTTGGACTTGAGCTTTGTACAAGATGACATATATAGGCTGTGTTTTATCTTCTACATACAGACAGCCA
GTTATACCAGCACCATTATTGAAGACACTTCTTTATTCATTGTATATTTTTTACTTCCTTGTCAAAAATCAAGTGA
CCATGAGTATGTGGTTTCATTCTGGGTCTTCAATTGTATTCCATTAGTCAACATATCTGCTCTGTACCAATACCATGC
```

**scaffold:** a sequence of contigs, separated by gaps....

```
TCCTAGATCCACTTGGACTTGAGCTTTGTACAAGATGACATATATAGGCTGTGTTTTATCTTCTACATACAGACAGCCA
GTTATACCAGCACCATTATTGAAGACACTTCTTTATTCATTGTATATTTTTTACTTCCTTGTCAAAAATCAAGTGA
CCATGAGTATGTGGTTTCATTCTGGGTCTTCAATTGTATTCCATTAGTCAACATATCTGCTCTGTACCAATACCATGC
NNNNNNNN
AGTTTTTACCACAATTGCTCTATAGTAAAGCTTGAGGTCAGGGTTGGTGATCCCTCCAGCCATTCTTTCATTATTAAGAA
TTGTTTTCCCTAGTCTGGGTTTTTGCTTTTCCAGGCGAATTGAGAATTGCTCTTCCATGCTTTGAAGAATTGTGTT
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
GGGATTTTGATGGGGTTTGCAATCTGTAGATTGCTTTGGTAAGATGGTTAGTTTTACTATGTAAATCTGCCAAT
CCACAAGCATGGGAGCGCTCTCCATTTCTGAGATCTTCTCAATTTCTTCTTGAGAACTTGAAGTTATTGTCATACA
```

Number of Ns = predicted gap size, with error bars (can't be displayed in fasta format)

33

## 1. Linear assemblies

### Example of an assembly in fasta format

```
>scaffold_1
TCCTAGATCCACTTGGACTTGAGCTTTGTACAAGATGACATATATAGGCTGTGTTTTATCTTCTACATACAGACAGCCA
GTTATACCAGCACCATTATTGAAGACACTTCTTTATTCATTGTATATTTTTTACTTCCTTGTCAAAAATCAAGTGA
CCATGAGTATGTGGTTTCATTCTGGGTCTTCAATTGTATTCCATTAGTCAACATATCTGCTCTGTACCAATACCATGC
NNNNNNNN
AGTTTTTACCACAATTGCTCTATAGTAAAGCTTGAGGTCAGGGTTGGTGATCCCTCCAGCCATTCTTTCATTATTAAGAA
TTGTTTTCCCTAGTCTGGGTTTTTGCTTTTCCAGGCGAATTGAGAATTGCTCTTCCATGCTTTGAAGAATTGTGTT
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
GGGATTTTGATGGGGTTTGCAATCTGTAGATTGCTTTGGTAAGATGGTTAGTTTTACTATGTAAATCTGCCAAT
CCACAAGCATGGGAGCGCTCTCCATTTCTGAGATCTTCTCAATTTCTTCTTGAGAACTTGAAGTTATTGTCATACA
>scaffold_2
CTGAAGTTGTTTATCAGCTGGAGAAGTTCTCAGGTAGAATTTTGGGATTGCTTATGTATGCTATCATATCACTTGCAAA
TAGTGATACCTTGATTCTTTTACCAGTATGATATCCATTGATCTCTTCTGTTGCTTATTGTTCTAGCTAACACTT
CAAGTACTATATTGAATAGATATGGGAGAGTGGGAATCCTTGCTTGTCTCCGATTTCAAGTGGGATTGCTTCAAGTATG
```

34

### 3. Linearized graph assemblies

Efasta: extended fasta

...ACTGTTT{A,C}GAAAT...      A or C at site  
 ...CGCGTTTTTTTTT{T,TT}CAT...      0 or 1 or 2 Ts at site

#### Example of an assembly in efasta format

```
>scaffold_1
TCCTAGATCCACTTGGACTTGAGCTTTGTATATATATATATATATA{,TA}CAAGATGACATATATAGGAGACAGCCA
GTTATACCAGCACCATTATTTGAAGACACTTTCTTTATTCATTGTATATTTTTTTTACTTCCTTGTCAAAAATCAAGTGA
CCATGAGTATGTGGTTTCATTTCTGGGTCTTCAATTGTATTCATTAGTCAACATATCTGTCTGTACCAATACCATGC
NNNNNNNN
AGTTTTTACCACAATTGCTCTATAGTAAAGCTTGAGGTCAGGGTTGGTGATCCCTCCAGCCATTCTTTCATTATTAAGAA
TTGTTTCCCTAGTCTGGGTTTTTGCTTTCCAGGCGAATTGAGAATTGCTCTTCCATGTCTTTGAAGAATTGTGTT
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
GGGATTTTGATGGGTTTGCAATTGAATCTGTAGATTGCTTTGGTAAGATGGTTAGTTTACTATGTTAATTCTGCCAAT
CCACAAGCATGGGAGCGCTCTCCATTTCTGAGATCTTCTCAATTCTTCTTGAGAACTGAAGTTATTGTCATACA
>scaffold_2
CTGAAGTTGTTTATCAGCTGGAGAAGTTCTCAGGTAGAATTTTGGGATT{A,C,G}GCTTATGTATGCTATCTTGCAAA
TAGTGATACCTTGATTCTTTTTTACCAATATGTATCCCATTGATCTCTTCTGTTGCTTATGTTCTAGCTAACACTT
CAAGTACTATATTGAATAGATATGGGAGAGTGGGAATCCTTGCTTGCTCCGATTTCAAGTGGGATGCTTCAAGTATG
```

35

### Putting it all together

1. Collect the BAM or FASTQ files that you wish to assemble. Create a `in_libs.csv` metadata file to describe your libraries and a `in_groups.csv` metadata file to describe your data files.

2. Prepare input files

```
% PrepareAllPathsInputs.pl \
  DATA_DIR=`pwd` PLOIDY=1 >& prepare.log
```

3. Assemble.

```
% RunAllPathsLG \
  PRE=. REFERENCE_NAME=. \
  DATA_SUBDIR=. RUN=default THREADS=4 >& run.log
```

4. Get the results (four files).

```
% cd default/ASSEMBLIES/test/
% less final.{assembly,contigs}.{fasta,efasta}
```

Once input files are generated (**PrepareAllPathsInputs.pl**)  
 can run assemblies with different parameters.

36

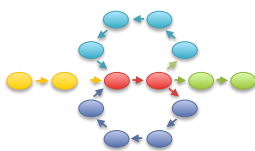
## Assembly parameters

- Parameters – all listed in ALLPATHS manual
- What about kmers?

“The user should **not** adjust the kmer size from the default value of  $K=96$ .

The relationship between kmer size  $K$  and read size is not a direct one in ALLPATHS-LG, unlike in many other assemblers. ALLPATHS-LG actually uses a number of different sizes of  $K$  internally, and because of this, it is not intended that users change the  $K$  values for an assembly”.

37



### Whole Genome Alignment with MUMmer

38

## Structural Variant Types

<p>Insertion into Reference</p> <p>R: AIB Q: AB</p>	<p>Insertion into Query</p> <p>R: AB Q: AIB</p>	<ul style="list-style-type: none"> <li>• Different structural variation types / misassemblies will be apparent by their pattern of breakpoints</li> <li>• Most breakpoints will be at or near repeats</li> <li>• Things quickly get complicated in real genomes</li> </ul> <p><a href="http://mummer.sf.net/manual/AlignmentTypes.pdf">http://mummer.sf.net/manual/AlignmentTypes.pdf</a></p>
<p>Collapse Query</p> <p>R: ARRB Q: ARB</p>	<p>Collapse Reference</p> <p>R: ARB Q: ARRB</p>	
<p>Collapse Query w/ insertion</p> <p>R: ARIRB Q: ARB</p> <p>Exact tandem alignment if I=R</p>	<p>Collapse Reference w/ insertion</p> <p>R: ARB Q: ARIRB</p> <p>Exact tandem alignment if I=R</p>	
<p>Collapse Query</p> <p>R: ARRRB Q: ARRB</p>	<p>Collapse Reference</p> <p>R: ARRB Q: ARRRB</p>	
<p>Inversion</p> <p>R: ABC Q: A<sup>b</sup>C</p>	<p>Rearrangement w/ Disagreement</p> <p>R: ABCDE Q: AFCBE</p>	

## Using MUMmer

- **mummer** – core program of package. Outputs matches between reference and query
- **nucmer** - all-vs-all comparison of nucleotide sequences contained in fasta files. Best for highly similar sequences with possible rearrangements
- **mummerplot** – generates dotplots and coverage plots
- **show-coords** – parses the output of nucmer and displays the coordinates

40

## Using MUMmer

```
mummer ref.fa CSHL/Halomonas/mydata/run-1/ASSEMBLIES/test/ \
final.contigs.fasta > mummer.out
```

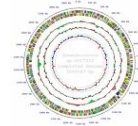
```
nucmer --maxmatch ref.fa \
CSHL/Halomonas/mydata/run-1/ASSEMBLIES/test/
final.contigs.fasta -p refctg
-maxmatch      Find maximal exact matches (MEMs) without repeat filtering
-p refctg      Set the output prefix for delta file
```

```
mummerplot --layout refctg.delta
-r Show the dotplot
```

```
show-coords -rclo refctg.delta
-r Sort alignments by reference position
-c Show percent coverage
-l Show sequence lengths
-o Annotate each alignment with BEGIN/END/CONTAINS
```

41

## Resources



- Assembly Competitions
  - Assemblathon: <http://assemblathon.org/>
  - GAGE: <http://gage.cbc.umd.edu/>
- Assembler Websites:
  - ALLPATHS-LG: <http://www.broadinstitute.org/software/allpaths-lg/blog/>
  - SOAPdenovo: <http://soap.genomics.org.cn/soapdenovo.html>
  - Celera Assembler: <http://wgs-assembler.sf.net>
- Tools:
  - MUMmer: <http://mummer.sourceforge.net/>
  - Quake: <http://www.cbc.umd.edu/software/quake/>
  - AMOS: <http://amos.sf.net>

42