

## Next Generation Sequencing Data Techniques: Reference-Based Mapping and de Novo Assembly

Matt Gitzendanner  
[magitz@ufl.edu](mailto:magitz@ufl.edu)


UF Information Technology www.it.ufl.edu

### Galaxy: Data intensive biology for everyone

- ▶ Accessible, reproducible, transparent computational biology
- ▶ galaxy.hpc.ufl.edu
  - Local instance of Galaxy
    - Faster access to storage, easier upload
    - Local compute resources
    - Local control

UF Information Technology www.it.ufl.edu

UNIVERSITY OF FLORIDA | High-Performance Computing



## HiPerGator

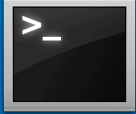
The University of Florida Supercomputer for Research

UF Information Technology www.it.ufl.edu

### Cluster basics

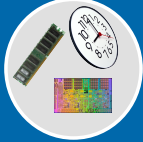
User interaction

Galaxy




Login node  
(Head node)

Scheduler



Tell the scheduler what you want to do

Compute resources




Your job runs on the cluster

UF Information Technology www.it.ufl.edu

### Reference-based mapping

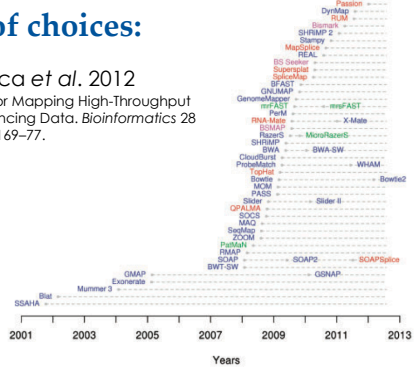
- ▶ Map NGS reads onto a reference genome
  - Identify SNPs
  - RNA-seq
  - CHIP-seq
  - Etc.



UF Information Technology www.it.ufl.edu

### Lots of choices:

- ▶ Fonseca et al. 2012
  - Tools for Mapping High-Throughput Sequencing Data. *Bioinformatics* 28 (24): 3169–77.



UF Information Technology www.it.ufl.edu



### Bowtie (Langmead *et al.* 2009)

- ▶ Keeping unmapped/mapped reads
  - --un <filename> unmapped reads
  - --al <filename> mapped reads
  - Can be helpful for downstream analyses
- ▶ Use -S for SAM output
  - Most likely will process output using SAM anyway
- ▶ -p: Bowtie is threaded, can run using multiple cores on **one** node
  - E.g.: nodes=1;ppn=8

UF | Information Technology | www.it.ufl.edu

### Bowtie2 (Langmead & Salzberg 2012)

- ▶ Adds gapped read alignment (indels)
- ▶ Faster than Bowtie for reads longer than 50bp
- ▶ Supports local alignment
  - Can trim ends that don't map
- ▶ Can map reads over Ns in reference
- ▶ No colorspace option

UF | Information Technology | www.it.ufl.edu

### Bowtie2 (Langmead & Salzberg 2012)

- ▶ Presets for both global and local
  - --very-fast(-local)
  - --fast(-local)
  - **--sensitive(-local) Defaults**
  - --very-sensitive(-local)

UF | Information Technology | www.it.ufl.edu

### SOLiD data

2nd base: A (blue), C (green), G (red), T (yellow)

Template Sequence: TA AC AA BA, GG CA CC TC, TA TT TT AT

**SNP site indicated by 2 adjacent color changes**

Reference in base space: A A C T A G G T G  
Reference in color space: ●●●●●●●●  
Read in color space: A A C G T A G G T G  
Read in base space: A A C G A T C C A C  
SNP

**Single color change is typically a measurement error**

Reference in base space: A A C T A G G T G  
Reference in color space: ●●●●●●●●  
Read in color space: A A C T A G G T G  
Read in base space: A A C T A G G T G  
Error

**1 Base Deletion**

Reference in base space: A A C T A G G T G  
Reference in color space: ●●●●●●●●  
Read in color space: A A C T A G G T G  
Read in base space: A A C T A G G T G  
Deletion

**Insertion**

Reference in base space: A A C T A G G T G  
Reference in color space: ●●●●●●●●  
Read in color space: A A C T A G G T G  
Read in base space: A A C T A G G T G  
Insertion

Use colorspace where possible

UF | Information Technology | www.it.ufl.edu

### Other mapping applications

- ▶ BWA
- ▶ Lastz
- ▶ Maq
  - Bowtie is generally faster
- ▶ Mosaik
  - Handles gapped alignments relative to reference
- ▶ PerM
- ▶ SRMA

UF | Information Technology | www.it.ufl.edu

### de Novo Assembly

- ▶ No reference genome
- ▶ Assemble contigs from reads
  - Assemble scaffolds using paired-end data
- ▶ Most short-read assemblers are de Buijn graph-based

Original sequence: GTAGTATAGTCAATCA

Sequence reads: GTAGA TAGAT AGTATA, GTAGG TAGAG, ATAGG TAGCA AGTCAG, GTCAG TAGTA, CAATG AGTATC GTATCA

k-mers (2-mers): GT TA AG AT TC CA

Consensus overlap assembly: GTAGA TAGAT AGTATA, TAGAT AGTATA, GTAGG TAGAG, ATAGG TAGCA AGTCAG, GTCAG TAGTA, CAATG AGTATC GTATCA

de Bruijn graph: A graph with nodes representing k-mers and edges representing overlaps between them.

Nature Reviews | Microbiology

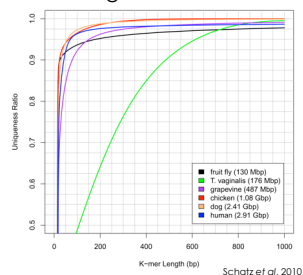
UF | Information Technology | www.it.ufl.edu

## kmers

- ▶ A kmer is a sequence of length k

- Longer kmer
  - More unique
  - Fewer reads/kmer
- Shorter kmer
  - Less unique
  - More reads/kmer

- ▶ The kmer you use does matter!
- Try different kmers



## Velvet (Zerbino & Birney 2008)

- ▶ Two stages
  - velvet
    - Creates the hash table of kmers
  - velvetg
    - Uses the de Bruijn graph to create contigs & scaffolds
- ▶ kmer is critical
  - 11-31: Default for Velvet, most memory efficient
  - Up to 249 available.

## Velvet (Zerbino & Birney 2008)

- ▶ Can use multiple types of sequencing inputs
  - Short, long
  - Paired, single
  - Different insert sizes
  - Reference
- ▶ A mix of library types is typically needed for de novo genome assembly
- ▶ Many helpful scripts distributed with Velvet
  - VelvetOptimiser—helps pick best kmer

## Other de novo assembly applications

- ▶ Abyss
- ▶ ALLPATHS-LG
  - Has very specific requirements for library types and coverage
- ▶ Metavelvet
  - Modified version of Velvet for metagenomics
- ▶ Newbler
  - Provided by Roche (454), but can use Illumina data
- ▶ SOAPdenovo
- ▶ For RNA-seq
  - Oases (builds on after Velvet)
  - SOAPdenovo-TRANS
  - Trinity

## UF Research Computing

- ▶ Help and Support (Continued)
  - <http://wiki.hpc.ufl.edu>
    - Documents on hardware and software resources
    - Various user guides
    - Many sample submission scripts
  - <http://hpc.ufl.edu/support>
    - Frequently Asked Questions
    - Account set up and maintenance

