

MSBI 32400 – LAB 6 LARRY HELSETH, PHD AND JASON EDELSTEIN

July 26 2017

Making bams & calling variants

2

- Today we'll go from FASTQ → BAM → VCF
- Using samtools, bwa, bcftools
- Whole genome alignment requires hg19.fa (3.1 GB) + bwa index files for hg19.fa (~7 GB)
- Not enough space on VM!
 - Will search FASTQ for one gene region against one chromosome

Setup Lab6 folders then extract FASTQ

3

- ❑ Make /data/lab6/bin, /data/lab6/data, /data/lab6/doc, /data/lab6/results & /data/lab6/src
- ❑ Go to /data/lab6/data
- ❑ Run samtools fastq to extract reads from Vince Buffalo's sample BAM in /data/bds-files/chapter-11-alignment/NA12891_CEU_sample.bam

MSBI 32400 Lab 6 7/26/2017

Syntax

4

```

student@MSBI32400Lab1:/data/lab6/data
File Edit View Search Terminal Help
[student@MSBI32400Lab1 data]$ samtools fastq
Usage: samtools fastq [options...] <in.bam>
Options:
  -0 FILE      write paired reads flagged both or neither READ1 and READ2 to FILE
  -1 FILE      write paired reads flagged READ1 to FILE
  -2 FILE      write paired reads flagged READ2 to FILE
  -f INT       only include reads with all bits set in INT set in FLAG [0]
  -F INT       only include reads with none of the bits set in INT set in FLAG [0]
  -n           don't append /1 and /2 to the read name
  -o           output quality in the OQ tag if present
  -s FILE      write singleton reads to FILE (assume single-end)
  -t           copy RG, BC and QT tags to the FASTQ header line
  -v INT       default quality score if not given in file [1]
  --input-fmt-option OPT[=VAL]
                Specify a single input file format option in the form
                of OPTION or OPTION=VALUE
  --reference FILE
                Reference sequence FASTA FILE [null]
[student@MSBI32400Lab1 data]$ time samtools fastq -t /data/bds-files/chapter-11-alignment/NA12891_CEU_sample.bam > NA12891_CEU_sample.fastq
[M::bam2fq_mainloop] processed 636207 reads

real    0m2.116s
user    0m1.990s
sys     0m0.116s
[student@MSBI32400Lab1 data]$

```

NB-All times shown were on an 8GB laptop with 4GB allocated to the VM. Please see Larry or Jason if you need help increasing your VM's memory (stop VM, adjust memory under Settings/System then restart VM)

samtools fastq -t (path to BAM) > NA12891_CEU_sample.fastq

MSBI 32400 Lab 6 7/26/2017

Download a reference genome

5

- Buffalo's chapter 11 README.md shows USH2A gene coordinates from chromosome 1:

```
## 'NA12891.CEU_sample.bam' Sample BAM File

The 'NA12891.CEU_sample.bam' sample BAM file is from region
chr1:215,622,894-216,423,396, which is gene
[USH2A]([http://uswest.ensembl.org/Homo_sapiens/Gene/Summary?db=core;g=ENSG00000042781;r=1:215622894-216423396]).
The alignment data comes from the 1000 Genomes
Project([http://www.1000genomes.org]), and the file was created with:

$ samtools view -hb ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/pilot2_high_cov_GRCh37_bams/data/NA12891.alignment/NA12891.chrom1.ILLUMINA.bwa.CEU.high_coverage.20100
517.bam \
1:215622894-216423396 > NA12891.CEU_sample.bam

Note that this illustrates that 'samtools view' can work with (sorted and indexed) BAM files over networks.

## USH2A Region

I chose this region because it's of significant [medical
importance]([http://en.wikipedia.org/wiki/Usher_syndrome]) and has interesting
biology. The mismatches I discuss (positions 215,906,547 and 215,906,548) in
this chapter were chosen for the sake of a technical example to illustrate how
useful visual inspection of SNPs is). These mismatches are likely false
positive variant calls due to common technical issues in base calling and
alignment.
/data/bds/files/chapter-11-alignment/README.md
```

MSBI 32400 Lab 6 7/26/2017

Download chr1 from UCSC

6

Index of /goldenPath/hg19/chromosomes - Mozilla Firefox

hgdownload.soe.ucsc.edu/goldenPath/hg19/chromosomes/

All the files in this directory are freely available for public use.

Name	Last modified	Size	Description
Parent Directory		-	
chr1.fa.gz	20-Mar-2009 08:58	70M	
chr1_g1000191_random.fa.gz	20-Mar-2009 09:02	33K	
chr1_g1000192_random.fa.gz	20-Mar-2009 09:02	178K	
chr2.fa.gz	20-Mar-2009 08:58	75M	
chr3.fa.gz	20-Mar-2009 08:58	61M	
chr4.fa.gz	20-Mar-2009 08:59	59M	
chr4_ctg9_hap1.fa.gz	20-Mar-2009 09:02	190K	
chr4_g1000193_random.fa.gz	20-Mar-2009 09:02	57K	
chr4_g1000194_random.fa.gz	20-Mar-2009 09:02	61K	
chr5.fa.gz	20-Mar-2009 08:59	56M	
chr6.fa.gz	20-Mar-2009 08:59	52M	
chr6_apd_hap1.fa.gz	20-Mar-2009 09:02	768K	
chr6_cox_hap2.fa.gz	20-Mar-2009 09:02	1.5M	
chr6_dbb_hap3.fa.gz	20-Mar-2009 09:02	1.3M	
chr6_mann_hap4.fa.gz	20-Mar-2009 09:02	1.3M	
chr6_mcf_hap5.fa.gz	20-Mar-2009 09:02	1.2M	

Move chr1.fa.gz from your ~/Downloads to /data/lab6/data
and extract using gunzip

MSBI 32400 Lab 6 7/26/2017

bwa mem syntax:

7

```
[student@MSBI32400Lab1 ~]$ bwa mem

Usage: bwa mem [options] <idxbase> <in1.fq> [in2.fq]

Algorithm options:

-t INT      number of threads [1]
-k INT      minimum seed length [19]
-w INT      band width for banded alignment [100]
-d INT      off-diagonal X-dropoff [100]
-r FLOAT    look for internal seeds inside a seed longer than {-k} * FLOAT [1.5]
-y INT      seed occurrence for the 3rd round seeding [20]
-c INT      skip seeds with more than INT occurrences [500]
-D FLOAT    drop chains shorter than FLOAT fraction of the longest overlapping chain [0.50]
-W INT      discard a chain if seeded bases shorter than INT [0]
-m INT      perform at most INT rounds of mate rescues for each read [50]
-S          skip mate rescue
-P          skip pairing; mate rescue performed unless -S also in use

Scoring options:

-A INT      score for a sequence match, which scales options -TdBOELU unless overridden [1]
-B INT      penalty for a mismatch [4]
-O INT[,INT] gap open penalties for deletions and insertions [6,6]
-E INT[,INT] gap extension penalty; a gap of size k cost '{-O} + {-E}*k' [1,1]
-L INT[,INT] penalty for 5'- and 3'-end clipping [5,5]
-U INT      penalty for an unpaired read pair [17]

-x STR      read type. Setting -x changes multiple parameters unless overridden [null]
             pacbio: -k17 -W40 -r10 -A1 -B1 -O1 -E1 -L0 (PacBio reads to ref)
             ont2d: -k14 -W20 -r10 -A1 -B1 -O1 -E1 -L0 (Oxford Nanopore 2D-reads to ref)
             intractg: -B9 -O16 -L5 (intra-species contigs to ref)

Input/output options:

-p          smart pairing (ignoring in2.fq)
-R STR      read group header line such as '@RG\tID:foo\tSM:bar' [null]
-H STR/FILE insert STR to header if it starts with @; or insert lines in FILE [null]
-j          treat ALT contigs as part of the primary assembly (i.e. ignore <idxbase>.alt file)

-v INT      verbose level: 1=error, 2=warning, 3=message, 4=debugging [3]
-T INT      minimum score to output [30]
-h INT[,INT] if there are <INT> hits with score >80% of the max score, output all in XA [5,200]
-a          output all alignments for SE or unpaired PE
-C          append FASTA/FASTQ comment to SAM output
-V          output the reference FASTA header in the XR tag
```

MSBI 32400 Lab 6 7/26/2017

Need to index for bwa

8

```
student@MSBI32400Lab1:~/data/lab6/data
File Edit View Search Terminal Help
[student@MSBI32400Lab1 data]$ time bwa index -a bwtsv chr1.fa
[bwa_index] Pack FASTA... 4.27 sec
[bwa_index] Construct BWT for the packed sequence...
[BWTInCreate] textLength=498501242, availableWord=47075968
[BWTInConstructFromPacked] 10 iterations done. 76384698 characters processed.
[BWTInConstructFromPacked] 20 iterations done. 142334426 characters processed.
[BWTInConstructFromPacked] 30 iterations done. 200946714 characters processed.
[BWTInConstructFromPacked] 40 iterations done. 253637450 characters processed.
[BWTInConstructFromPacked] 50 iterations done. 299331834 characters processed.
[BWTInConstructFromPacked] 60 iterations done. 340474362 characters processed.
[BWTInConstructFromPacked] 70 iterations done. 377037946 characters processed.
[BWTInConstructFromPacked] 80 iterations done. 409531690 characters processed.
[BWTInConstructFromPacked] 90 iterations done. 438408170 characters processed.
[BWTInConstructFromPacked] 100 iterations done. 464069642 characters processed.
[BWTInConstructFromPacked] 110 iterations done. 486873562 characters processed.
[bwt_gen] Finished constructing BWT in 116 iterations.
[bwa_index] 388.60 seconds elapsed.
[bwa_index] Update BWT... 7.48 sec
[bwa_index] Pack forward-only FASTA... 5.31 sec
[bwa_index] Construct SA from BWT and Occ... 115.07 sec
[main] Version: 0.7.15-r1140
[main] CMD: bwa index -a bwtsv chr1.fa
[main] Real time: 527.783 sec; CPU: 528.738 sec

real    8m47.705s
user    8m31.921s
sys     0m8.818s
[student@MSBI32400Lab1 data]$ ls -la chr1*
-rw-rw-r-- 1 student student 254235640 Feb 4 12:57 chr1.fa
-rw-rw-r-- 1 student student 707 Feb 4 13:37 chr1.fa.amb
-rw-rw-r-- 1 student student 44 Feb 4 13:37 chr1.fa.ann
-rw-rw-r-- 1 student student 249250696 Feb 4 13:37 chr1.fa.bwt
-rw-rw-r-- 1 student student 62312657 Feb 4 13:37 chr1.fa.pac
-rw-rw-r-- 1 student student 124625360 Feb 4 13:39 chr1.fa.sa
[student@MSBI32400Lab1 data]$
```

Use:

bwa index -a bwtsv chr1.fa

MSBI 32400 Lab 6 7/26/2017

Need samtools index of chr1.fa

11

- samtools faidx builds a .fai file

```

student@MSBI32400Lab1:/data/lab6/data
File Edit View Search Terminal Help
[student@MSBI32400Lab1 data]$ time samtools faidx chr1.fa

real    0m1.957s
user    0m1.903s
sys     0m0.038s
[student@MSBI32400Lab1 data]$

```

MSBI 32400 Lab 6 7/26/2017

Convert SAM to BAM

12

- samtools view -bt chr1.fa.fai
NA12891_CEU_sample.sam >
NA12891_CEU_sample.bam
- samtools sort -o
NA12891_CEU_sample_sorted.bam
NA12891_CEU_sample.bam
- samtools index NA12891_CEU_sample_sorted.bam

MSBI 32400 Lab 6 7/26/2017

View header of new sorted BAM

13

- The -R '@RG' syntax put our new ID and sample name in header along with the @PG (program) info for how we generated the alignment

```
[student@MSBI32400Lab1 data]$ samtools view -H NA12891_CEU_sample_sorted.bam
@HD VN:1.3 SO:coordinate
@SQ SN:chr1 LN:249250621
@RG ID:MSBI32400 test SM:NA12891_CEU_sample
@PG ID:bwa PN:bwa VN:0.7.15-r1140 CL:bwa mem -R @RG\tID:MSBI32400_test\tSM:NA12891_CEU_sample chr1.fa NA12891_CEU_sample.fastq
[student@MSBI32400Lab1 data]$
```

MSBI 32400 Lab 6 7/26/2017

Check samtools man page

14

- man samtools then search for mpileup (use '/mpileup')

o Call SNPs and short INDELS:

```
samtools mpileup -uf ref.fa aln.bam | bcftools call -mv > var.raw.vcf
bcftools filter -s LowQual -e '%QUAL<20 || DP>100' var.raw.vcf > var.flt.vcf
```

The **bcftools filter** command marks low quality sites and sites with the read depth exceeding a limit, which should be adjusted to about twice the average read depth (bigger read depths usually indicate problematic regions which are often enriched for artefacts). One may consider to add **-CS0** to **mpileup** if mapping quality is overestimated for reads containing excessive mismatches. Applying this option usually helps **BWA-short** but may not other mappers.

Individuals are identified from the **SM** tags in the **@RG** header lines. Individuals can be pooled in one alignment file; one individual can also be separated into multiple files. The **-P** option specifies that indel candidates should be collected only from read groups with the **@RG-PL** tag set to **ILLUMINA**. Collecting indel candidates from reads sequenced by an indel-prone technology may affect the performance of indel calling.

- See also: http://proquestcombo.safaribooksonline.com.proxy.uchicago.edu/book/bioinformatics/9781449367480/visualizing-alignments-with-samtools-tview-and-the-integrated-genomics-viewer/idp33784528_html

MSBI 32400 Lab 6 7/26/2017

bcftools

15

From man page:

- BCFtools is a set of utilities that manipulate variant calls in the Variant Call Format (VCF) and its binary counterpart BCF. All commands work transparently with both VCFs and BCFs, both uncompressed and BGZF-compressed.
- Most commands accept VCF, bgzipped VCF and BCF with filetype detected automatically even when streaming from a pipe. Indexed VCF and BCF will work in all situations. Un-indexed VCF and BCF and streams will work in most, but not all situations. In general, whenever multiple VCFs are read simultaneously, they must be indexed and therefore also compressed.
- BCFtools is designed to work on a stream. It regards an input file "-" as the standard input (stdin) and outputs to the standard output (stdout). Several commands can thus be combined with Unix pipes.

MSBI 32400 Lab 6 7/26/2017

bcftools syntax:

16

```
[student@MSBI32400Lab1 data]$ bcftools
Program: bcftools (Tools for variant calling and manipulating VCFs and BCFs)
Version: 1.3.1 (using htslib 1.3.1)

Usage:  bcftools [--version|--version-only] [--help] <command> <argument>

Commands:

-- Indexing
  index          index VCF/BCF files

-- VCF/BCF manipulation
  annotate       annotate and edit VCF/BCF files
  concat         concatenate VCF/BCF files from the same set of samples
  convert        convert VCF/BCF files to different formats and back
  isec           intersections of VCF/BCF files
  merge         merge VCF/BCF files from non-overlapping sample sets
  norm          left-align and normalize indels
  plugin        user-defined plugins
  query         transform VCF/BCF into user-defined formats
  reheader      modify VCF/BCF header, change sample names
  view          VCF/BCF conversion, view, subset and filter VCF/BCF files

-- VCF/BCF analysis
  call          SNP/indel calling
  consensus     create consensus sequence by applying VCF variants
  cnv           HMM CNV calling
  filter        filter VCF/BCF files using fixed thresholds
  gtcheck       check sample concordance, detect sample swaps and contamination
  roh           identify runs of autozygosity (HMM)
  stats         produce VCF/BCF stats

Most commands accept VCF, bgzipped VCF, and BCF with the file type detected
automatically even when streaming from a pipe. Indexed VCF and BCF will work
in all situations. Un-indexed VCF and BCF and streams will work in most but
not all situations.
```

[student@MSBI32400Lab1 data]\$

MSBI 32400 Lab 6 7/26/2017

Generate mpileup & run bcftools

17

- ❑ `samtools mpileup -uf chr1.fa`
`NA12891_CEU_sample_sorted.bam | bcftools call`
`-mv > NA12891_CEU_sample_sorted_var.raw.vcf`
- ❑ `bcftools filter -s LowQual -e '%QUAL<20'`
`NA12891_CEU_sample_sorted_var.raw.vcf >`
`NA12891_CEU_sample_sorted_var.flt.vcf`
- ❑ How many variants are called in the final VCF?
 How many variants are called with "PASS"?
 - ❑ Include in your README for Jason

MSBI 32400 Lab 6 7/26/2017

bcftools call syntax

18

```
[student@MSBI32400Lab1 ~]$ bcftools call
```

About: SNP/indel variant calling from VCF/BCF. To be used in conjunction with samtools mpileup. This command replaces the former "bcftools view" caller. Some of the original functionality has been temporarily lost in the process of transition to htslib, but will be added back on popular demand. The original calling model can be invoked with the -c option.

Usage: `bcftools call [options] <in.vcf.gz>`

File format options:

<pre>--no-version -o, --output <file> -O, --output-type <b v z v> --ploidy <assembly>[?] --ploidy-file <file> -r, --regions <region> -R, --regions-file <file> -s, --samples <list> -S, --samples-file <file> -t, --targets <region> -T, --targets-file <file> --threads <int></pre>	<pre>do not append version and command line to the header write output to a file [standard output] output type: 'b' compressed BCF; 'u' uncompressed BCF; 'r' compressed VCF; 'v' uncompressed VCF [v] predefined ploidy, 'list' to print available settings, append '?' for details space/tab-delimited list of CHROM, FROM, TO, SEX, PLOIDY restrict to comma-separated list of regions restrict to regions listed in a file list of samples to include [all samples] PED file or a file with an optional column with sex (see man page for details) [all samples] similar to -r but streams rather than index-jumps similar to -R but streams rather than index-jumps number of extra output compression threads [0]</pre>
--	---

Input/output options:

<pre>-A, --keep-alts -f, --format-fields <list> -g, --gvcf <int>[,...] -i, --insert-missed -M, --keep-masked-ref -V, --skip-variants <type> -v, --variants-only</pre>	<pre>keep all possible alternate alleles at variant sites output format fields: GQ,GP (lowercase allowed) [1] group non-variant sites into gvcf blocks by minimum per-sample DP output also sites missed by mpileup but present in -T keep sites with masked reference allele (REF=N) skip indels/snps output variant sites only</pre>
---	--

Consensus/variant calling options:

<pre>-c, --consensus-caller -C, --constrain <str> -m, --multiallelic-caller -n, --novel-rate <float>[,...] -p, --pval-threshold <float> -P, --prior <float></pre>	<pre>the original calling method (conflicts with -m) one of: alleles, trio (see manual) alternative model for multiallelic and rare-variant calling (conflicts with -c) likelihood of novel mutation for constrained trio calling, see man page for details [1e-8,1e-9,1e-9] variant if P(ref D)<FLOAT with -c [0.5] mutation rate (use bigger for greater sensitivity) [1.1e-3]</pre>
---	---

```
[student@MSBI32400Lab1 ~]$
```

MSBI 32400 Lab 6 7/26/2017

bcftools filter syntax

19

```
[student@MSBI32400Lab1 data]$ bcftools filter
```

About: Apply fixed-threshold filters.

Usage: bcftools filter [options] <in.vcf.gz>

Options:

-e, --exclude <expr>	exclude sites for which the expression is true (see man page for details)
-g, --SnpGap <int>	filter SNPs within <int> base pairs of an indel
-G, --IndelGap <int>	filter clusters of indels separated by <int> or fewer base pairs allowing only one to pass
-i, --include <expr>	include only sites for which the expression is true (see man page for details)
-m, --mode [x]	"+": do not replace but add to existing FILTER; "x": reset filters at sites which pass
--no-version	do not append version and command line to the header
-o, --output <file>	write output to a file [standard output]
-O, --output-type <b u z v>	b: compressed BCF, u: uncompressed BCF, z: compressed VCF, v: uncompressed VCF [v]
-r, --regions <region>	restrict to comma-separated list of regions
-R, --regions-file <file>	restrict to regions listed in a file
-s, --soft-filter <string>	annotate FILTER column with <string> or unique filter name ("Filter%d") made up by the program ("+")
-S, --set-GTs <. 0>	set genotypes of failed samples to missing (.) or ref (0)
-t, --targets <region>	similar to -r but streams rather than index-jumps
-T, --targets-file <file>	similar to -R but streams rather than index-jumps
--threads <int>	number of extra output compression threads [0]

```
[student@MSBI32400Lab1 data]$
```

MSBI 32400 Lab 6 7/26/2017

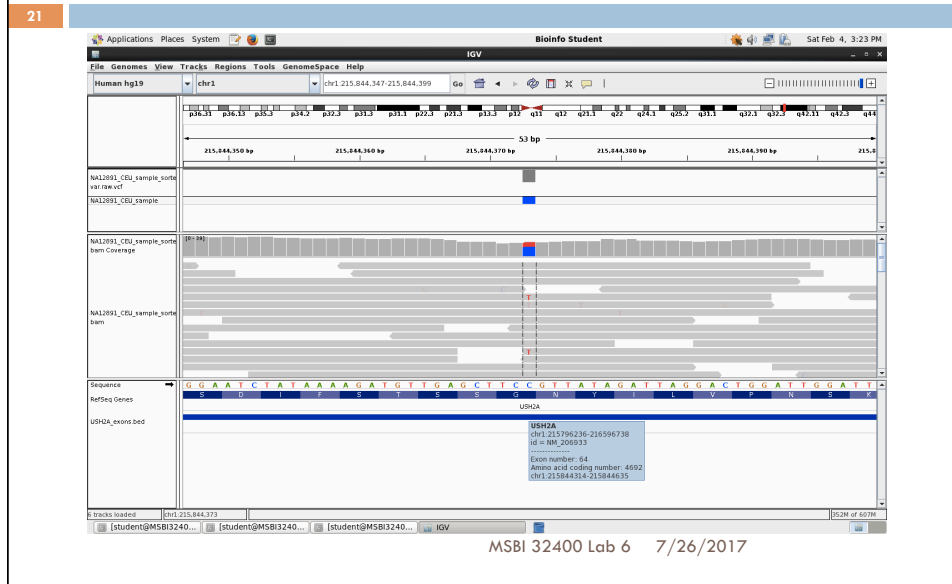
Open BAM & VCF in IGV

20

- View exon 64 and look for SNPs called in VCF
 - ▣ Most SNPs in introns, but a few in exons
 - ▣ Record the coordinates and Amino Acid # to send to Jason

MSBI 32400 Lab 6 7/26/2017

IGV view (BAM + bai + VCF + BED)



Another way

22

Vince Buffalo shows:

- samtools mpileup -v --no-BAQ --region 1:215906528-215906567 -fasta-ref...
 - His coordinates won't work for our BAM since it uses chr1
 - Also, his coordinates are only 39 bp!
 - If you try his notes, use the full sequence from his Chapter 11 README file (chr1:215622894-216423396)

MSBI 32400 Lab 6 7/26/2017

Vince's way

23

```
[student@MSB132400Lab1 data]$ time samtools mpileup -v --no-BAQ --region chr1:215622894-216423396 --fasta-ref chr1.fa NA12891_CEU_sample_sorted.bam > NA12891_CEU_sample_sorted_full_region.vcf.gz
[mpileup] 1 samples in 1 input files
[mpileup] Set max per-file depth to 8000

real    0m46.709s
user    0m45.198s
sys     0m0.625s
[student@MSB132400Lab1 data]$ time bcftools call -v -m NA12891_CEU_sample_sorted_full_region.vcf.gz > NA12891_CEU_sample_sorted_full_region_calls.vcf.gz
Note: Neither --ploidy nor --ploidy-file given, assuming all sites are diploid

real    0m4.141s
user    0m4.059s
sys     0m0.029s
[student@MSB132400Lab1 data]$
```

- His VCF is the very similar to the one generated before, though he outputs a vcf.gz which is not recognized by IGV or gunzip
- Solution: `bgzip -d`
`NA12891_CEU_sample_sorted_full_region.vcf.gz`
 then open in IGV or text editor for viewing

MSB1 32400 Lab 6 7/26/2017

samtools mpileup with BED file

24

- `samtools mpileup -B -C50 -f chr1.fa -l USH2A_exons.bed -o NA12891_CEU_sample_sorted.vcf -v -u NA12891_CEU_sample_sorted.bam`
- Check the samtools man page to see what `-B` and `-C50` mean for mpileup
 - ▣ Put that in your README

MSB1 32400 Lab 6 7/26/2017

Homework

25

- E-mail Jason (jasone@uchicago.edu) the README with the file information requested above before next class with “**Lab #6**” in the subject line

MSBI 32400 Lab 6 7/26/2017