

NGS FILE FORMATS

Peter FitzGerald, PhD

Head Genome Analysis Unit, CCR, NCI



DIFFERENT FILE FORMATS

Sequence Data

- FASTA
- FASTQ
- SRA

Alignment Data

- SAM
- BAM
- CRAM

Variant Data

- VCF

Annotation Data

- BED
- bigBED
- GFF
- GTF

Graphing Data

- bedGraph
- WIG
- bigWIG



DIFFERENT FILE FORMATS

COMPRESSED VARIANTS

- *.gz. - gzip compression
- *.zip - zip compression and or archive
- *.tar - archive of files
- *.tar.gz. - gzip compressed archive



SEQUENCE FILE FORMATS

FASTA FORMAT

FASTA

Standard text based format for storing simple sequence data.

Each entry consists of a header line that begins with a “>” followed by one or more lines of sequence data.

The format allows for multiple sequences in a single file.



SEQUENCE FILE FORMATS

FASTA FORMAT

FASTA

Single sequence example:

```
>HWI-ST398_0092:1:1:5372:2486#0/1
TTTTCGTTCTTTCATGTACCGCTTTGTTGGTTAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGAT
ACGTAGCAGCAGCATCAGTACGACTACGACGACTAGCACATGCGACGATCGATGCTAGCTGACTATCGATG
```

Multiple sequence example:

```
>Sequence Name 1
TTTTCGTTCTTTCATGTACCGCTTTGTTGGTTAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGAT
ACGTAGCAGCAGCATCAGTACGACTACGACGACTAGCACATGCGACGATCGATGCTAGCTGACTATCGATG
>Sequence Name 2
ACGTAGACACGACTAGCATCAGCTACGCATCGATCAGCATCGACTACGATCACACATCGATCAGCATCACGACTAGCAT
AGCATCGACTACACTACGACTACGATCCACGTACGACTAGCATGCTAGCGTAGCTAGCTAGTCGATCGATGAGT
AGCTAGCTAGCTAGC
>Sequence Name 3
ACTCAGCATGCATCAGCATCGACTACGACTACGACATCGACTAGCATCAGCAT
```



SEQUENCE FILE FORMATS

FASTQ FORMAT

FASTQ

Text based format for storing sequence data and corresponding quality scores for each base.

To enable a one-one correspondence between the base sequence and the quality score the score is stored as a single one letter/number code using an offset of the standard ASCII code.

Quality scores range from 0 to 40 and represent a \log^{10} score for the probability of being wrong.

E.g. score of 30 => 1:1000 chance of error



SEQUENCE FILE FORMATS

FASTQ FORMAT

FASTQ

Each fastq file contain multiple entries and each entry consists of 4 lines:

1. header line beginning with “@” and sequence name
2. sequence line
3. header line beginning with “+” which can have the name but rarely does
4. quality score line



SEQUENCE FILE FORMATS

FASTQ FORMAT

FASTQ

```
@HWI-ST398_0092:6:73:5372:2486#0/1
TTTTTCGTTCTTTCATGTACCGCTTTGTTGGTAGATCGGAAGAGCGGTTCAGCAGGAATGCCGAGACCGAT
+HWI-ST398_0092:1:1:5372:2486#0/1
fffffeedfcedfffffeffdefff_ffffffdccfdZdeeadefecZedaecdbRdTY^ZYT``_T`_`bc_Wceaa[
```

6 - Flowcell lane

73 - Tile number

5372:2486 - 'x','y'-coordinates of the cluster within the tile

#0 - index number for a multiplexed sample (0 for no indexing)

/1 - the member of a pair, /1 or /2 (paired-end or mate-pair reads only)

For paired end reads fastq files come in pairs, typically labelled R1 and R2 (reads are in same order in both files...header often does not distinguish between read1 and read2



SEQUENCE FILE FORMATS

SRA FORMAT

SRA (**S**equence **R**ead **A**rchive) is a binary archive, used by NCBI for distributing data from its SRA database.

1. Archive format that can hold many different types of data (reads and/or alignments etc)
2. Requires use of one or more of the programs in the SRA toolkit to extract usable data.
3. When used with NGS data the most useful tool is probably **fastq-dump**
4. It's challenging to know what data is in the archive



ALIGNMENT FILE FORMATS

SAM FORMAT

```
8_100_10000_12419    163  chrVII 271183 255  40M   =  271294 151
TGGTGTATTATACGCTACCGTGCAGTGCCGGGGCAACCG
bbbabbabbbbabbbbabbbbabbbbabbbbabbbbabbb  XA:i:0 MD:Z:40 NM:i:0
```

The **SAM Format** (**S**equence **A**lignment/**M**ap) is a text format for storing sequence alignment data in a series of tab delimited ASCII columns.

The file has two parts:

1. **Header** - Each line starts with a “@”.
@HD, @SQ, @RG, @PG
2. **Alignments** - One line for each entry.



ALIGNMENT FILE FORMATS

SAM FORMAT

Example of SAM Header

```
@HD VN:1.0      SO:unsorted
@SQ SN:chr1      LN:195471971
@SQ SN:chr2      LN:182113224
@SQ SN:chr3      LN:160039680
@SQ SN:chr4      LN:156508116
@SQ SN:chr5      LN:151834684
@SQ SN:chr6      LN:149736546
@SQ SN:chr7      LN:145441459
@SQ SN:chr8      LN:129401213
@SQ SN:chr9      LN:124595110
@SQ SN:chr10     LN:130694993
@SQ SN:chr11     LN:122082543
@SQ SN:chr12     LN:120129022
@SQ SN:chr13     LN:120421639
@SQ SN:chr14     LN:124902244
@SQ SN:chr15     LN:104043685
@SQ SN:chr16     LN:98207768
@SQ SN:chr17     LN:94987271
@SQ SN:chr18     LN:90702639
@SQ SN:chr19     LN:61431566
@SQ SN:chrX      LN:171031299
@SQ SN:chrY      LN:91744698
@SQ SN:chrM      LN:16299
@PG ID:bowtie2   PN:bowtie2 VN:2.2.9   CL:"/usr/local/apps/bowtie/2-2.2.9/bowtie2-align-s --wrapper basic-0 -x /fdb/bowtie2.DELETE/mm10 -q jun_minus_dex_rep1a -S jun_minus_dex_rep1a_mm10.sam -p8"
```



ALIGNMENT FILE FORMATS

SAM FORMAT

8_100_10000_12419 163 chrVII 271183 255 40M = 271294 151

TGGTGTATTATACGCTACCGTGC GG TGCCGGGGCAACCG

bbbabbbaaaaaaaaaaaaaaaaabbcbbaaaaaaaaaaaaaa XA:i:0 MD:Z:40 NM:i:0

8_100_10000_12	163	chr7	271183	255	40M	=	271294	151	TGGTGTATTATACG	bbbabbbaaaaaaaaaaaaaaaaabbcbbaaaaaaaaaaaaaa	XA:i:0	MD:Z:40
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	MRNM	MPOS	TLEN	SEQ	QUAL	OPT	

Col	Field	Description
1	QNAME	Query template/pair NAME
2	FLAG	bitwise FLAG
3	RNAME	Reference sequence NAME
4	POS	1-based leftmost POSition/coordinate of clipped sequence
5	MAPQ	MAPping Quality (Phred-scaled)
6	CIGAR	extended CIGAR string
7	MRNM	Mate Reference sequence NaMe ('=' if same as RNAME)
8	MPOS	1-based Mate POSition
9	TLEN	inferred Template LENGTH (insert size)
10	SEQ	query SEQuence on the same strand as the reference
11	QUAL	query QUALity (ASCII-33 gives the Phred base quality)
12+	OPT	variable OPTional fields in the format TAG:VTYPE:VALUE



ALIGNMENT FILE FORMATS

SAM FORMAT

8_100_10000_12419 **163** chrVII 271183 255 **40M** = 271294 151

TGGTGTATTATACGCTACCGTGC GG TGCCGGGGCAACCG

bbbabbabbbbbbcbcbcbcbcbcbcbbbb XA:i:0 MD:Z:40 NM:i:0

Understanding Flag codes

<http://broadinstitute.github.io/picard/explain-flags.html>

1	read paired
2	read mapped in proper pair
4	read unmapped
8	mate unmapped
16	read reverse strand
32	mate reverse strand
64	first in pair
128	second in pair
256	not primary alignment
512	read fails platform/vendor quality checks
1024	read is PCR or optical duplicate
2048	supplementary alignment



ALIGNMENT FILE FORMATS

BAM/CRAM FORMAT

BAM (*.bam) is the compressed binary version of the Sequence Alignment/Map (SAM) format, a compact and index-able representation of nucleotide sequence alignments. **BAM** is compressed in the **BGZF** format that supports random access through the BAM file index (*.bam.bai).

HINT: Filename.bam and filename.bai always go together

CRAM (*.cram) - newer implementation of BAM like binary data.

1. Significantly better lossless compression than BAM
2. Full compatibility with BAM
3. Effortless transition to CRAM from using BAM files
4. Support for controlled loss of BAM data



ANNOTATION FILE FORMATS

BED FORMAT

1. **chrom** - name of the chromosome
2. **chromStart** - Start of feature (0-based)
3. **chromEnd** - End of the feature (not included in display)
+ 9 optional columns - most common are:
4. **name** - a label for the feature
5. **score** - a score (0-1000)
6. **strand** - which strand the feature on (+/-)

chr1	15000	20000	gene1	50	+
chr2	106000	108000	gene2	400	-



ANNOTATION FILE FORMATS

BED FORMAT

7. **thickStart** - The starting position at which the feature is drawn thickly (for example, the start codon in gene displays). When there is no thick part, thickStart and thickEnd are usually set to the chromStart position.
8. **thickEnd** - The ending position at which the feature is drawn thickly (for example the stop codon in gene displays).
9. **itemRgb** - An RGB value of the form R,G,B (e.g. 255,0,0). If the track line itemRgb attribute is set to "On", this RBG value will determine the display color of the data contained in this BED line. NOTE: It is recommended that a simple color scheme (eight colors or less) be used with this attribute to avoid overwhelming the color resources of the Genome Browser and your Internet browser.
10. **blockCount** - The number of blocks (exons) in the BED line.
11. **blockSizes** - A comma-separated list of the block sizes. The number of items in this list should correspond to blockCount.
12. **blockStarts** - A comma-separated list of block starts. All of the blockStart positions should be calculated relative to chromStart. The number of items in this list should correspond to blockCount.



ANNOTATION FILE FORMATS

GFF FORMAT

GFF (General Feature Format) GFF lines have nine required fields that *must* be tab-separated [GFF2 - UCSC & GFF3 - EMBL]

1. **squid** - The name of the chromosome or scaffold.
2. **source** - The program that generated this feature.
3. **feature** - The name of this type of feature. Some examples of standard feature types are "CDS" "start_codon" "stop_codon" and "exon"
4. **start** - The starting position of the feature in the sequence. The first base is numbered 1.
5. **end** - The ending position of the feature (inclusive).
6. **score** - floating point value
7. **strand** - Valid entries include "+", "-", or "." (for don't know/don't care).
8. **phase** - If the feature is a coding exon, frame should be a number between 0-2 that represents the reading frame of the first base. If the feature is not a coding exon, the value should be ":".
9. **attributes**- A list of feature attributes in the format tag=value pairs separated by ";"

GFF2 <http://genome.ucsc.edu/FAQ/FAQformat.html#format3>

GFF3 <https://github.com/The-Sequence-Ontology/Specifications/blob/master/gff3.md>
<http://useast.ensembl.org/info/website/upload/gff3.html>



ANNOTATION FILE FORMATS

GFF FORMAT

GFF example

```
0 ##gff-version 3.2.1
1 ##sequence-region ctg123 1 1497228
2 ctg123 . gene 1000 9000 . + . ID=gene0001;Name=EDEN
3 ctg123 . TF_binding_site 1000 1012 . + . ID=tfbs0001;Parent=gene0001
4 ctg123 . mRNA 1050 9000 . + . ID=mRNA0001;Parent=gene0001;Name=EDEN.1
5 ctg123 . mRNA 1050 9000 . + . ID=mRNA0002;Parent=gene0001;Name=EDEN.2
6 ctg123 . mRNA 1300 9000 . + . ID=mRNA0003;Parent=gene0001;Name=EDEN.3
7 ctg123 . exon 1300 1500 . + . ID=exon0001;Parent=mRNA0003
8 ctg123 . exon 1050 1500 . + . ID=exon0002;Parent=mRNA0001,mRNA0002
9 ctg123 . exon 3000 3902 . + . ID=exon0003;Parent=mRNA0001,mRNA0003
10 ctg123 . exon 5000 5500 . + . ID=exon0004;Parent=mRNA0001,mRNA0002,mRNA0003
11 ctg123 . exon 7000 9000 . + . ID=exon0005;Parent=mRNA0001,mRNA0002,mRNA0003
12 ctg123 . CDS 1201 1500 . + 0 ID=cds0001;Parent=mRNA0001;Name=edenprotein.1
13 ctg123 . CDS 3000 3902 . + 0 ID=cds0001;Parent=mRNA0001;Name=edenprotein.1
14 ctg123 . CDS 5000 5500 . + 0 ID=cds0001;Parent=mRNA0001;Name=edenprotein.1
15 ctg123 . CDS 7000 7600 . + 0 ID=cds0001;Parent=mRNA0001;Name=edenprotein.1
16 ctg123 . CDS 1201 1500 . + 0 ID=cds0002;Parent=mRNA0002;Name=edenprotein.2
17 ctg123 . CDS 5000 5500 . + 0 ID=cds0002;Parent=mRNA0002;Name=edenprotein.2
18 ctg123 . CDS 7000 7600 . + 0 ID=cds0002;Parent=mRNA0002;Name=edenprotein.2
19 ctg123 . CDS 3301 3902 . + 0 ID=cds0003;Parent=mRNA0003;Name=edenprotein.3
20 ctg123 . CDS 5000 5500 . + 1 ID=cds0003;Parent=mRNA0003;Name=edenprotein.3
21 ctg123 . CDS 7000 7600 . + 1 ID=cds0003;Parent=mRNA0003;Name=edenprotein.3
22 ctg123 . CDS 3391 3902 . + 0 ID=cds0004;Parent=mRNA0003;Name=edenprotein.4
23 ctg123 . CDS 5000 5500 . + 1 ID=cds0004;Parent=mRNA0003;Name=edenprotein.4
24 ctg123 . CDS 7000 7600 . + 1 ID=cds0004;Parent=mRNA0003;Name=edenprotein.4
```



ANNOTATION FILE FORMATS

GTF FORMAT

GTF (Gene Transfer Format) is a refined form of the GFF with group attributes - essentially the same as GFF2

1. **seqname** - The name of the sequence. Must be a chromosome or scaffold. (chr1 or 1)
2. **source** - The program that generated this feature.
3. **feature** - The name of this type of feature. Some examples of standard feature types are "CDS" "start_codon" "stop_codon" and "exon"
4. **start** - The starting position of the feature in the sequence. The first base is numbered 1.
5. **end** - The ending position of the feature (inclusive).
6. **score** - A score between 0 and 1000 (UCSC) OR floating point value
7. **strand** - Valid entries include "+", "-", or "." (for don't know / don't care).
8. **frame** - If the feature is a coding exon, frame should be a number between 0-2 that represents the reading frame of the first base. If the feature is not a coding exon, the value should be ":".
9. **attributes/group** - A list of feature attributes in the format tag=value pairs separated by ";"

GTF/GFF2 <http://useast.ensembl.org/info/website/upload/gff.html>



ANNOTATION FILE FORMATS

VCF/BCF FORMAT

VCF (Variant Call Format) is a flexible and **extendable** format for variation data such as single nucleotide variants, insertions/deletions, copy number variants and structural variants. The file may be compressed and indexed for faster access.

The formal format is a moving target with several revisions and modifications.

Reference: <https://github.com/samtools/hts-specs>



ANNOTATION FILE FORMATS

VCF/BCF FORMAT

VCF Example

```
##fileformat=VCFv4.3
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens",taxonomy=x>
##phasing=partial##INFO=<ID=NS,Number=1>Type=Integer,Description="Number of Samples With Data"> ##INFO=<ID=DP,Number=1>Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129"> ##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10"> ##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT:GQ:DP:HQ0|0:48:1:51,511|0:48:8:51,511/1:43:5
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ0|0:49:3:58,500|1:3:5:65,30/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ1|2:21:6:23,272|1:2:0:18,22/2:35:4
20 1230237 . T . 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ0|0:54:7:56,600|0:48:4:51,510/0:61:2
20 1234567 microsat1 GTC G,GTCT 50 PASS NS=3;DP=9;AA=G GT:GQ:DP0/1:35:40/2:17:21/1:40:3
```



GRAPHING FILE FORMATS

WIG (BIGWIG) FORMAT

1) FixedStep

fixedStep	chrom=chr1 start=3001 step=1
24	
56	
100	

2) VariableStep

variableStep	chrom=chr1
3001	24
3002	56
3003	100

variableStep	chrom=chr1
3001	24
3003	56
3010	100



GRAPHING FILE FORMATS

BEDGRAPH FORMAT

1. **chrom** - name of the chromosome
2. **chromStart** - Start of feature (0-based)
3. **chromEnd** - End of the feature (not included in display)
4. **score** - a score (integer or real positive / negative number)

chr1	15000	20000	1
chr2	106000	108000	0.75



Format Conversion Utilities

- Galaxy (<http://galaxy.psu.edu/> - <http://galaxy.cit.nih.gov/>)
 - Galaxy is an open, web-based platform for data intensive biomedical research. Whether on the free public server or your own instance, you can perform, reproduce, and share complete analyses.
- Samtools (<http://samtools.sourceforge.net>)
 - SAM Tools provide various utilities for manipulating alignments in the SAM format, including sorting, merging, indexing and generating alignments in a per-position format. Also, note TABIX for indexing generic tab delimited files.
- Picard (<http://picard.sourceforge.net/>)
 - Picard comprises Java-based command-line utilities that manipulate SAM files, and a Java API (SAM-JDK) for creating new programs that read and write SAM files. Both SAM text format and SAM binary (BAM) format are supported.
- UCSC Utilities (<http://hgdownload.cse.ucsc.edu/admin/exe/>)



Format Conversion Utilities

- Bamtools -(<https://github.com/pezmaster31/bamtools>)
 - BamTools provides both a programmer's API and an end-user's toolkit for handling BAM files.
- Bedtools (<http://bedtools.readthedocs.io/en/latest/>)
 - Collectively, the bedtools utilities are a swiss-army knife of tools for a wide-range of genomics analysis tasks. The most widely-used tools enable genome arithmetic: that is, set theory on the genome. For example, bedtools allows one to intersect, merge, count, complement, and shuffle genomic intervals from multiple files in widely-used genomic file formats such as BAM, BED, GFF/GTF, VCF. While each individual tool is designed to do a relatively simple task (e.g., intersect two interval files), quite sophisticated analyses can be conducted by combining multiple bedtools operations on the UNIX command line.
- FASTX-Toolkit (http://hannonlab.cshl.edu/fastx_toolkit/)
 - The FASTX-Toolkit is a collection of command line tools for Short-Reads FASTA/FASTQ files preprocessing.
- SRA ToolKit (<https://github.com/ncbi/sra-tools>)
 - The SRA Toolkit and SDK from NCBI is a collection of tools and libraries for using data in the INSDC Sequence Read Archives.



Binary Formats & Indices

Indexed binary file formats are much more efficient.

Only the portions of the files needed for the region currently being processed or visualized are transferred and loaded as needed. Thus for large data sets they are considerably faster than regular files.
(e.g. bigBED, bigWIG, BAMindexed)



THE END

