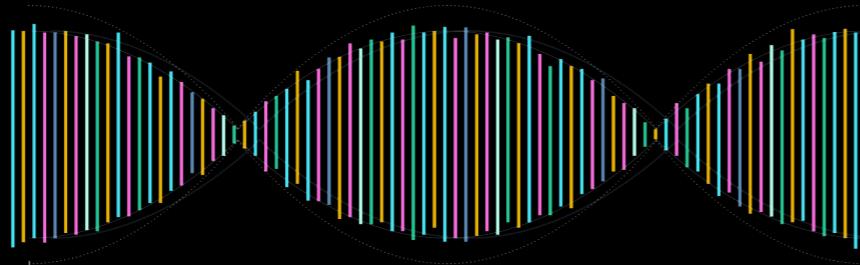




Практическая биоинформатика

Занятие 1

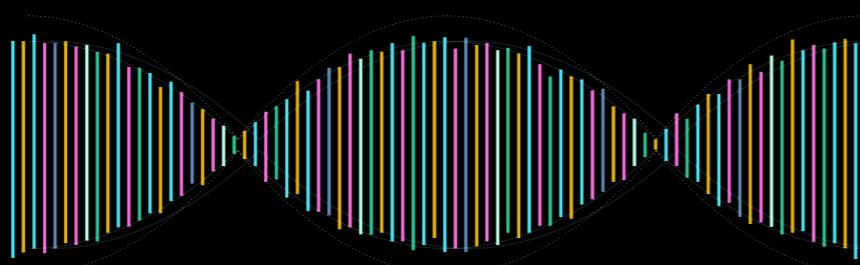
Александр Ракитко



...AATACAGTGC...
...AATGCAGTGC...



...AATGCAGTGC...
...AATGCAGTC...



...AATGCAGTGC...
...AATAACAGTGC...

Технологии анализа генома

ПЦР
(1 мутация)

Микрочиповое
генотипирование
(650 тыс. мутаций)

- Отдельные «горячие» мутации

Экзомное
секвенирование
(12-30 млн. мутаций)

- Риски заболеваний
- Наследственные заболевания
- Фармакогенетика
- Питание
- Происхождение

Полногеномное
секвенирование
(3 млрд. мутаций)

- Наследственные заболевания
- Фармакогенетика

- Риски заболеваний
- Наследственные заболевания
- Фармакогенетика
- Питание
- Происхождение

Форматы данных

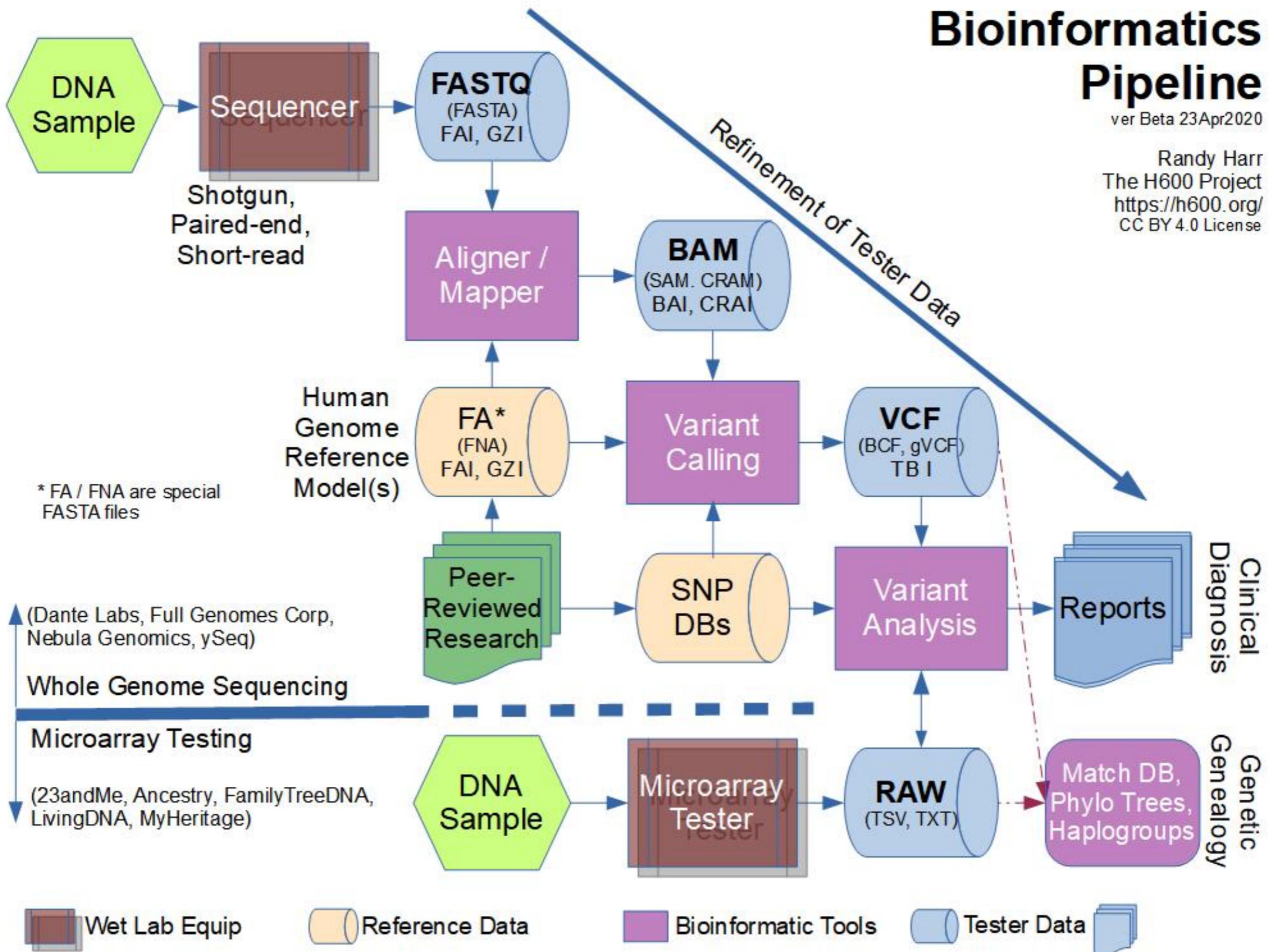
- Данные разного типа хранятся в разных форматах, универсального не существует
- При возможности стоит использовать стандартные и распространенные форматы данных
- XLSX – последнее, что стоит использовать
- Описание большинства форматов можно найти здесь:
<https://genome.ucsc.edu/FAQ/FAQformat.html>



Bioinformatics Pipeline

ver Beta 23Apr2020

Randy Harr
The H600 Project
<https://h600.org/>
CC BY 4.0 License



FASTA file

```
>NG_008679.1:5001-38170 Homo sapiens paired box 6 (PAX6)
ACCCTTTTCTTATGACATTAAACTCTGGGGCAGGTCTCGCGTAGAACGCAGCTGTCA
GCCACTTCCCCTGCCGAGCGCGGTGAGAAGTGTGGAACCGCGCTGCCAGGCTCAC
CTGCCTCCCCGC
CCTCCGCTCCCAGGTAACCGCCCAGGCTCCGGCCCCGGCCGGCTCGGGGCCGCGGG
CCTCTCCGCTG
CCAGCGACTGCTGTCCCCAAATCAAAGCCGCCCAAGTGGCCCCGGGCTTGATT
TTGCTTTAAAAG
GAGGCATAAAAGATGGAAGCGAGTTACTGAGGGAGGGATAGGAAGGGGGTGGAGGAG
GGACTTGTCT
TGCCGAGTGTGCTCTGCAAAAGTAGCAAAATGTTCCACTCCTAAGAGTGGACTTCC
AGTCCGGCCCT
GAGCTGGGAGTAGGGGGCGGGAGTCTGCTGCTGTCTGCTAAAGCCACTCGCGACCGC
GAAAAATGCA
GGAGGTGGGGACGCACTTGCATCCAGACCTCCTCTGCATCGCAGTCACGACATCC
ACGCTTGGAAAG
TCCGTACCCGCGCTGGAGCGCTTAAAGACACCCCTGCCCGGGTCGGCGAGGTGCAG
CAGAAGTTCCC
GCGGTTGCAAAGTGCAGATGGCTGGACCGCAACAAAGTCTAGAGATGGGTTCTCAG
AAAGACGC
```

FIGURE 7.1: An example fasta file showing the first part of the PAX6 gene.

<https://compgenomr.github.io/book/fasta-and-fastq-formats.html>

<https://www.ncbi.nlm.nih.gov/nuccore/1798174254>

FASTQ file

Identifier ————— @HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Sequence ————— TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTNNNNNNNNNTAGTTCTTGAGA
+ sign & identifier ————— +HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
Quality scores ————— efcfffffffceefffcfffffffddf`feed] `] _Ba_ ^ __ [YBBBBBBBBBRTT\]]] [] dddd`

Base T
phred Quality] = 29

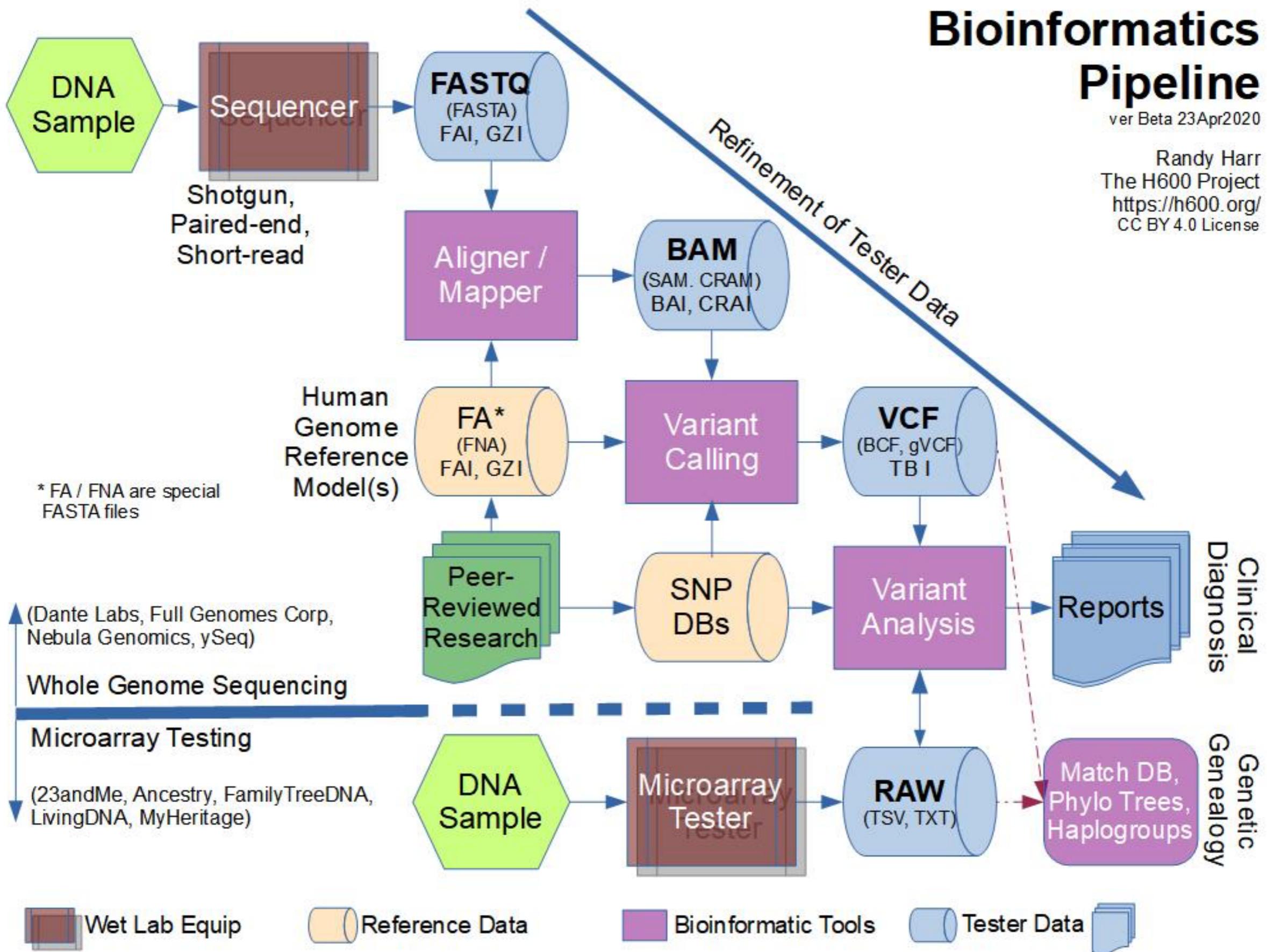
FIGURE 7.2: FASTQ format and a brief explanation of each line in the format.

<https://compgenomr.github.io/book/fasta-and-fastq-formats.html>

Bioinformatics Pipeline

ver Beta 23Apr2020

Randy Harr
The H600 Project
<https://h600.org/>
CC BY 4.0 License



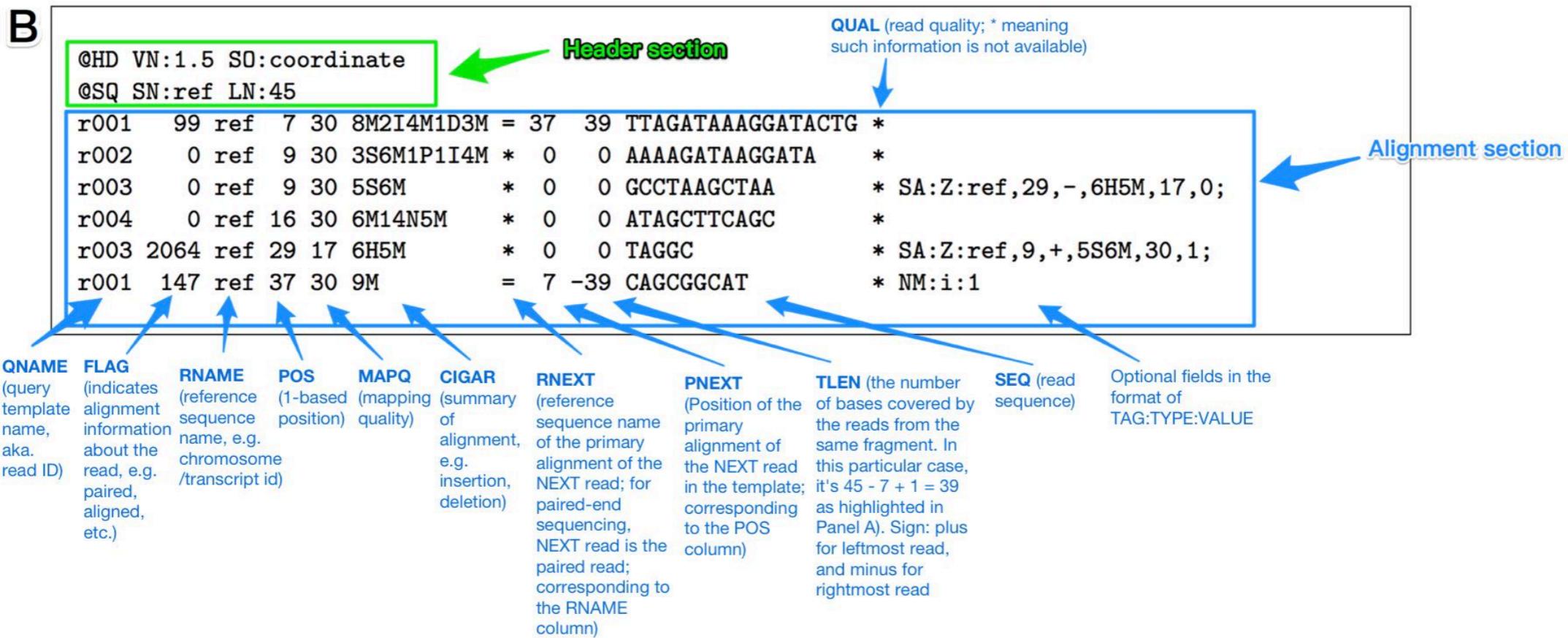
SAM file

<http://zyxue.github.io/2017/09/26/sam-format-example.html>

A

Coor	12345678901234	10	20	30	40
ref	AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT				
+r001/1	TTAGATAAAGGATA*CTG				
+r002	aaaAGATAA*GGATA				
+r003	gcctaAGCTAA				
+r004		ATAGCT.....TCAGC			
-r003		ttagctTAGGC			
-r001/2			CAGCGGCAT		

B

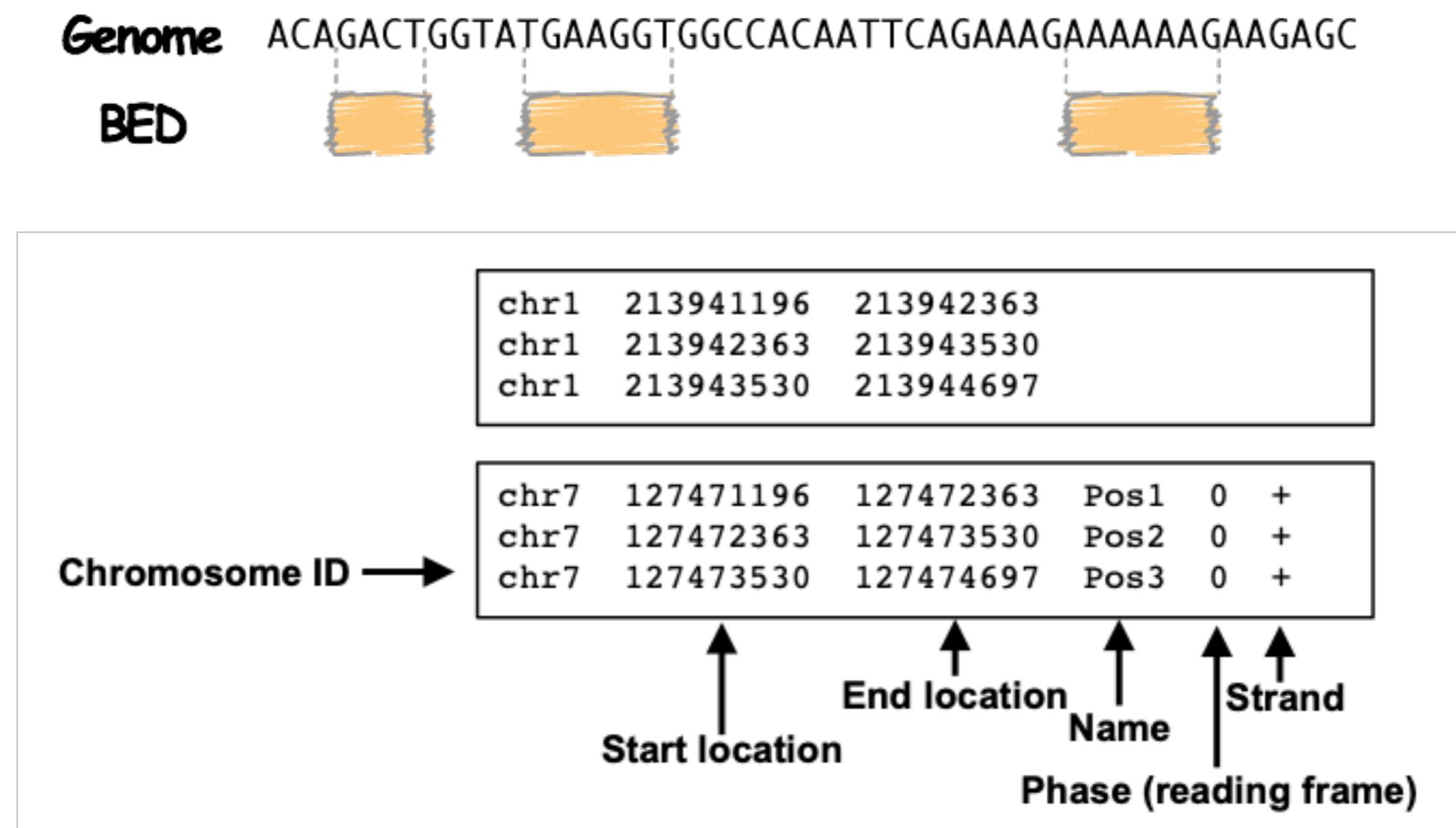


SAM / BAM / CRAM

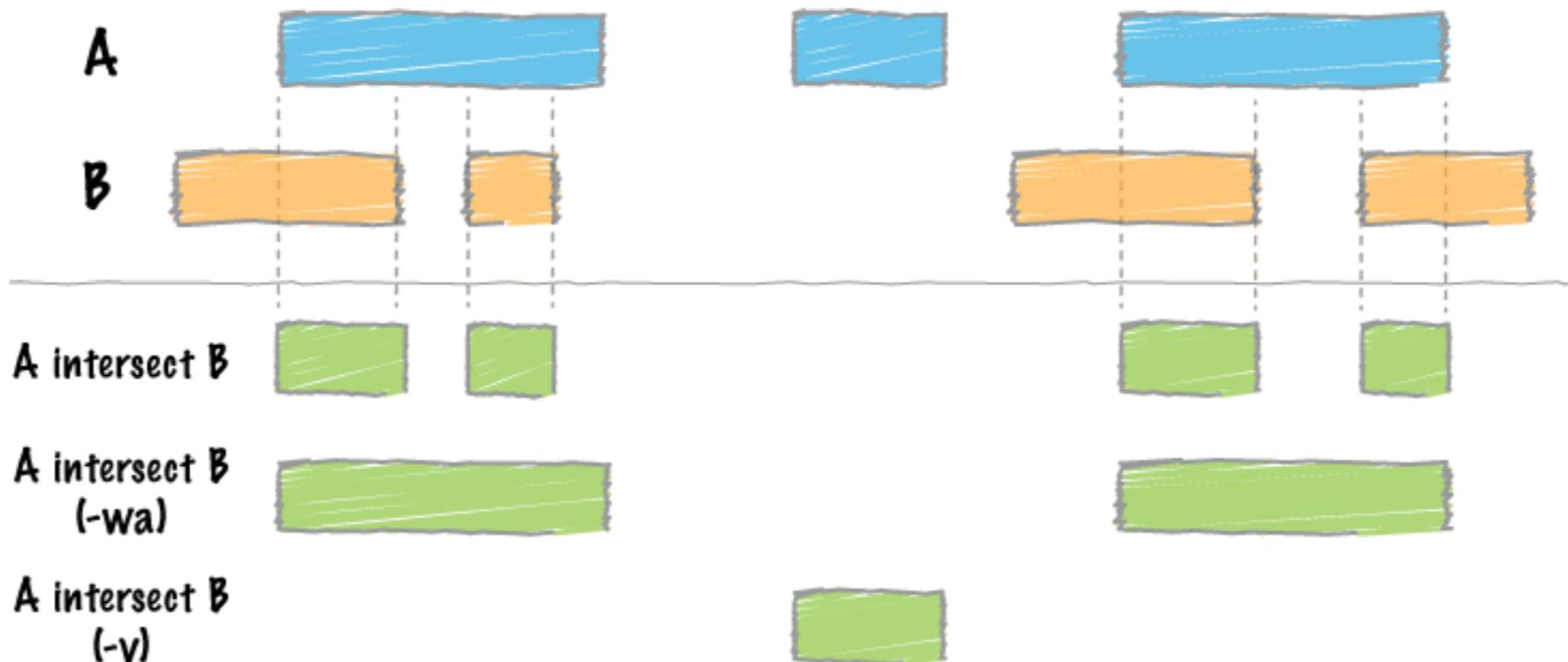
- Выравнивания хранятся в SAM формате. SAM можно сжать в BAM или CRAM
- Работать с ними можно с помощью **samtools** (просмотр, объединение, фильтрация)
- Всегда сортируем и индексируем
- Для визуализации можно использовать IGV
- BAM можно конвертировать в FASTA/FASTQ



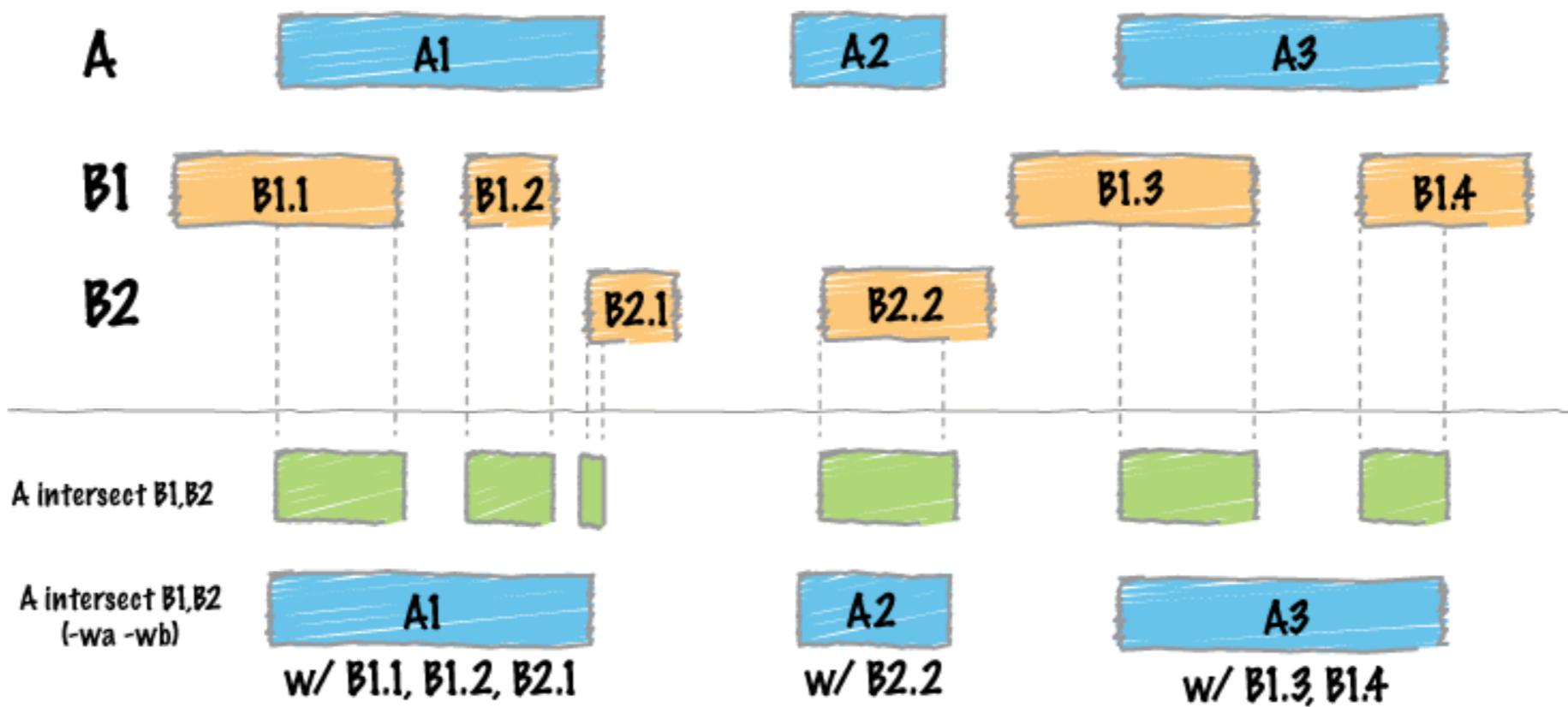
BED file



Intersect w/ 1 database



Intersect w/ 2 or more databases



BED file

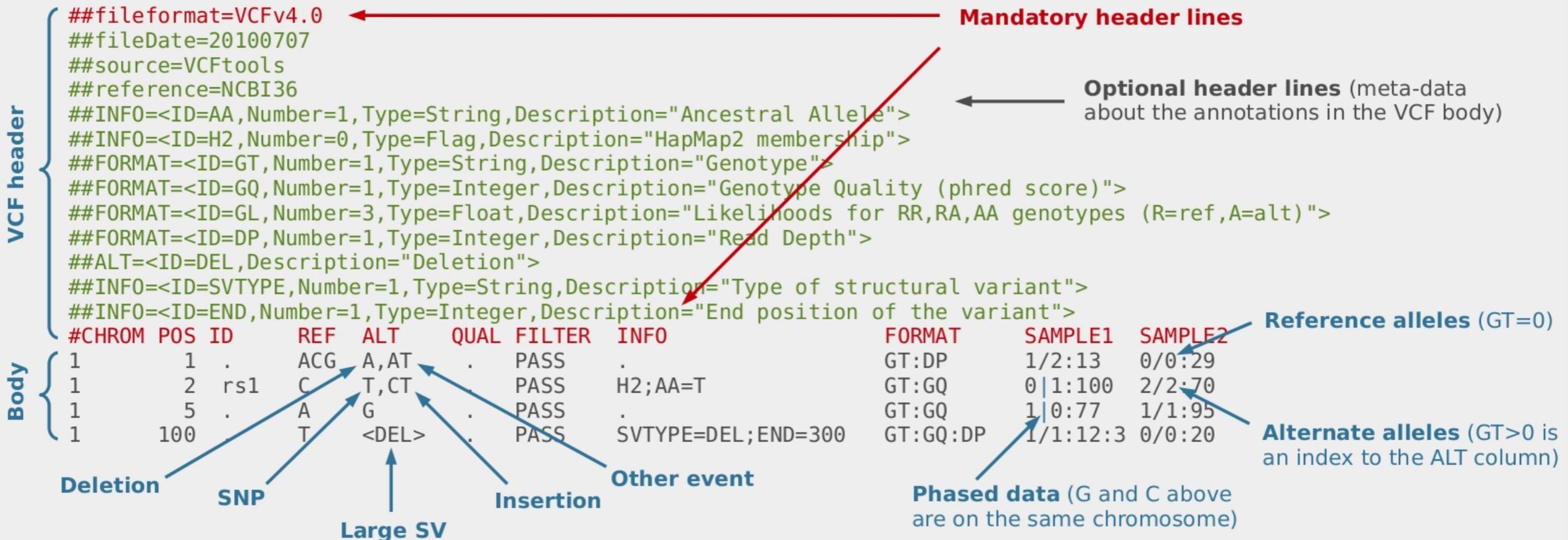
- Используется для хранения и работы с геномными интервалами
- Например, можно использовать для координат генов, экзонов и тд.
- Не содержит header
- Удобный инструмент для работы **bedtools**



VCF format

HEADER	<pre>##fileformat=VCFv4.1 ##fileDate=20090805 ##tcgaversion=1.1 ##vcfProcessLog=<InputVCF=<file1.vcf>,InputVCFSource=<caller1>,InputVCFVer=<1.0>,InputVCFParam=<a1,b>,InputVCFgeneAnno=<anno1.gaf>> ##reference=ftp://ftp.ncbi.nih.gov/genbank/genomes/Eukaryotes/vertebrates_mammals/Homo_sapiens/GRCh37/special_requests/GRCh37-lite.fa ##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"></pre>	INFO meta-information																																																
	<pre>##FILTER=<ID=q10,Description="Quality below 10"> ##FILTER=<ID=s50,Description="Less than 50% of samples have data"></pre>	FILTER meta-information																																																
BODY	<pre>##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"></pre>	FORMAT meta-information																																																
	<pre>##SAMPLE=<ID=NORMAL,Individual=TCGA-01-1000,File=TCGA-01-1000-1.bam,Platform=Illumina,Source=dbGAP,Accession=1234> ##SAMPLE=<ID=TUMOR,Individual=TCGA-01-1000,File=TCGA-01-1000-2.bam,Platform=Illumina,Source=dbGAP,Accession=4567> ##PEDIGREE=<Name_0=TUMOR,Name_1=NORMAL></pre>	Optional: FORMAT field specifying data type + Per-sample genotype data																																																
Fixed fields																																																		
		<table border="1"><thead><tr><th>FORMAT</th><th>NORMAL</th><th>TUMOR</th></tr></thead><tbody><tr><td>GT:GQ:DP:HQ</td><td>0 0:48:1:51,51</td><td>1 0:48:8:51,51</td></tr><tr><td>GT:GQ:DP:HQ</td><td>0 0:49:3:58,50</td><td>0 1:3:5:65,3</td></tr><tr><td>GT:GQ:DP:HQ</td><td>1 2:21:6:23,27</td><td>2 1:2:0:18,2</td></tr><tr><td>GT:GQ:DP:HQ</td><td>0 0:54:7:56,60</td><td>0 0:48:4:51,51</td></tr><tr><td>GT:GQ:DP</td><td>0/1:35:4</td><td>0/2:17:2</td></tr></tbody></table>	FORMAT	NORMAL	TUMOR	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	GT:GQ:DP	0/1:35:4	0/2:17:2																														
FORMAT	NORMAL	TUMOR																																																
GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51																																																
GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3																																																
GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2																																																
GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51																																																
GT:GQ:DP	0/1:35:4	0/2:17:2																																																
<table border="1"><thead><tr><th>#CHROM</th><th>POS</th><th>ID</th><th>REF</th><th>ALT</th><th>QUAL</th><th>FILTER</th><th>INFO</th></tr></thead><tbody><tr><td>20</td><td>14370</td><td>rs6054257</td><td>G</td><td>A</td><td>29</td><td>PASS</td><td>NS=3;DP=14;AF=0.5;DB;H2</td></tr><tr><td>20</td><td>17330</td><td>.</td><td>T</td><td>A</td><td>3</td><td>q10</td><td>NS=3;DP=11;AF=0.017</td></tr><tr><td>20</td><td>1110696</td><td>rs6040355</td><td>A</td><td>G,T</td><td>67</td><td>PASS</td><td>NS=2;DP=10;AF=0.333,0.667;DB</td></tr><tr><td>20</td><td>1230237</td><td>.</td><td>T</td><td>.</td><td>47</td><td>PASS</td><td>NS=3;DP=13;AA=T</td></tr><tr><td>20</td><td>1234567</td><td>microsat1</td><td>GTC</td><td>G,GTCTC</td><td>50</td><td>PASS</td><td>NS=3;DP=9;AA=G</td></tr></tbody></table>		#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;DB	20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	20	1234567	microsat1	GTC	G,GTCTC	50	PASS	NS=3;DP=9;AA=G	
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO																																											
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2																																											
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017																																											
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;DB																																											
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T																																											
20	1234567	microsat1	GTC	G,GTCTC	50	PASS	NS=3;DP=9;AA=G																																											

VCF format



VCF file

- Содержит данные генотипирования
- Часто содержит только позиции с альтернативными аллелями
- Подходит для хранения больших референсных панелей, фазированных данных
- Удобный инструмент для работы **bcftools / vcftools**



PLINK format

*.ped

FID	IID	PID	MID	Sex	P	rs1	rs2	rs3
1	1	0	0	2	1	CT	AG	AA
2	2	0	0	1	0	CC	AA	AC
3	3	0	0	1	1	CC	AA	AC

*.map

Chr	SNP	GD	BPP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

*.fam

FID	IID	PID	MID	Sex	P
1	1	0	0	2	1
2	2	0	0	1	0
3	3	0	0	1	1

*.bed

Contains binary version of the SNP info of the *.ped file.
(not in a format readable for humans)

*.bim

Chr	SNP	GD	BPP	Allele 1	Allele 2
1	rs1	0	870000	C	T
1	rs2	0	880000	A	G
1	rs3	0	890000	A	C

Covariate file

FID	IID	C1	C2	C3
1	1	0.00812835	0.00606235	-0.000871105
2	2	-0.0600943	0.0318994	-0.0827743
3	3	-0.0431903	0.00133068	-0.000276131

Legend

FID	Family ID	rs{x}	Alleles per subject per SNP
IID	Individual ID	Chr	Chromosome
PID	Paternal ID	SNP	SNP name
MID	Maternal ID	GD	Genetic distance (morgans)
Sex	Sex of subject	BPP	Base-pair position (bp units)
P	Phenotype	C{x}	Covariates (e.g., Multidimensional Scaling (MDS) components)

PLINK format

- Содержит данные генотипирования
- Позволяет сохранять родственные связи и фенотипы
- Легче, чем VCF
- «Теряет» фазирование
- Удобный инструмент для работы **plink**



Форматы данных

Формат	Данные	Инструменты
FASTA	Последовательность нуклеотидов	
FASTQ	Последовательности ридов	
SAM/BAM/CRAM	Выровненные риды	samtools
BED	Регионы	bedtools
VCF	Генотипы	bcftools/vcftools
PLINK	Генотипы	plink