# Genomic data

Libor Mořkovský, Václav Janoušek

# Genomic data

- Genome from the bioinformatic perspective

- Where does the genomic data come from?

- Common genomic data formats

- Specialized tools for genomic data

# Genome from the bioinformatic perspective

- sequence

```
AGTGGGCGAGGCGCGGAGGTCTGGCCTATAAAGTAGTCGCGGAGACGGGG
TGCTGGTTTGCGTCGTAGTCTCCTGCAGCGTCTGGGGTTTCCGTTGCAGT
CCTCGGAACCAGGACCTCGGCGTGGCCTAGCGAGTTATGGCGACGAAGGC
CGTGTGCGTGCTGAAGGGCGACGGCCCAGTGCAGGGCATCATCAATTTCG
AGCAGAAGGCAAGGGCTGGGACGGAGGCTTGTTTGCGAGGCCGCTCCCAC
CCGCTCGTCCCCCCGCGCACCTTTGCTAGGAGCGGGTCGCCCGCCAGGCC
TCGGGGCCGCCCTGGTCCAGCGCCCGGTCCCGGCCCGTGCCGCCCGGTCG
GTGCCTTCGCCCCCAGCGGTGCGGTGCCCAAGTGCTGAGTCACCGGGCGG
GCCCGGGCGCGGGGCGTGGGACCGAGGCCGCCGCGGGGCTGGGCCTGCGC
GTGGCGGGAGCGCGGGGAGGGATTGCCGCGGGCCGGGGAGGGGCGGGGGC
GGGCGTGCTGCCCTCTGTGGTCCTTGGGCCGCCGCCGCGGGTCTGTCGTG
GTGCCTGGAGCGGCTGTGCTCGTCCCTTGCTTGGCCGTGTTCTCGTTCCT
GAGGGTCCCGCGGACACCGAGTGGCGCAGTGCCAGGCCCAGCCCGGGGAT
GGCGACTGCGCCTGGGCCCGCCTGGTGTCTTCGCATCCCTCTCCGCTTTC
CGGCTTCAGCGCTCTAGGTCAGGGAGTCTTCGCTTTTGTACAGCTCTAAG
GCTAGGAATGGTTTTTATATTTTTAAAAGGCTTTGGAAAACAAAAATACG
CAACAGAGACCGTTTGTGTGACACTTTGCAGGGAAGTTTGCTGGCCTCTG
TTCTAGGTCATGATTGGGCTGCAAGGGCAGAGAAGGTAGCCTTGAACAGA
GGTCCTTTTCCTCCTCCTAAGCTCCGGGAGCCAGAGGTTTAACTGACCCT
```

# Genome from the bioinformatic perspective

- physical map

AGTGGGCGAGGCGCGGAGGTCTGGCCTATAAAGTAGTCGCGGAGACGGGGTGCTGGTTTGCGTCGTAGTCTCCTGCAGCGTCTGGGGTTTCCGTTGCAGT
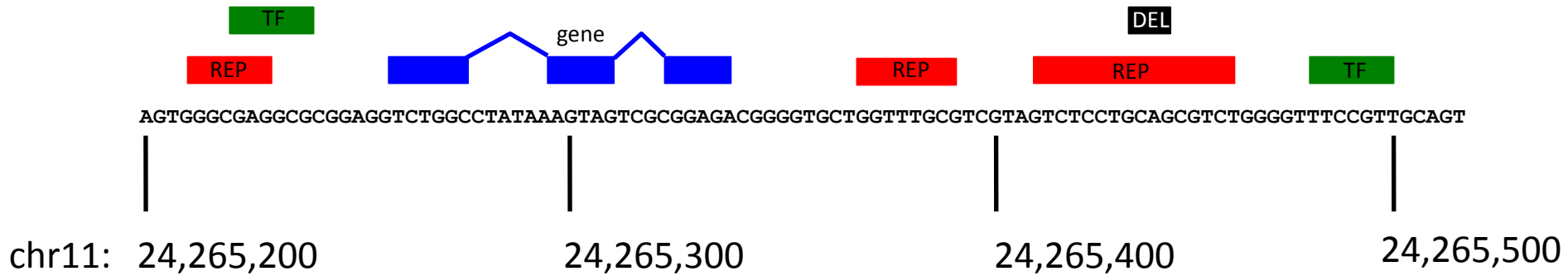
| | | | |

chr11: 22,341,400          22,341,500          22,341,600          22,341,700

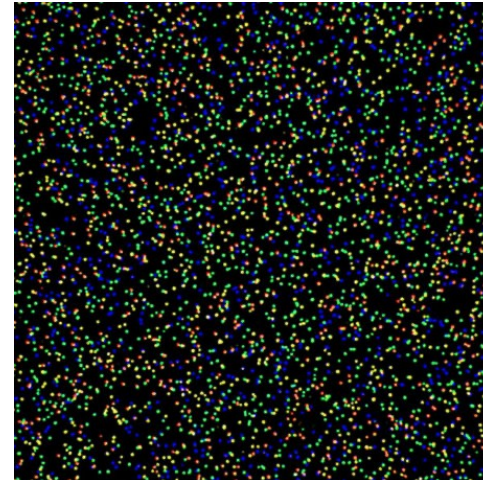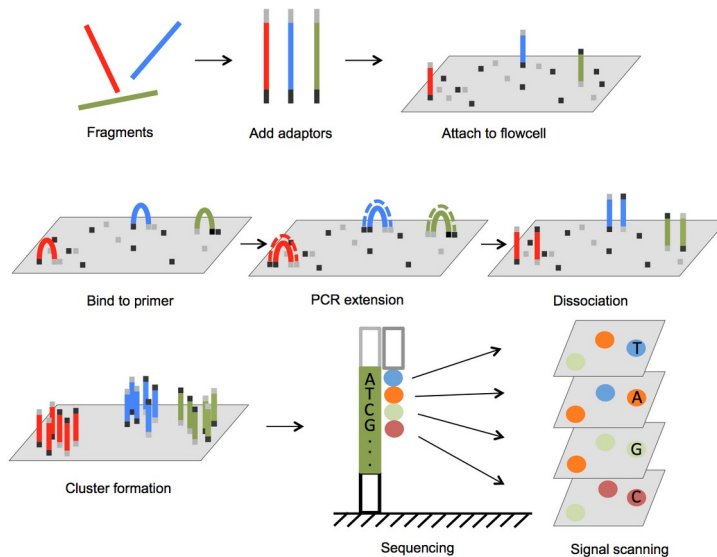# Genome from the bioinformatic perspective

- annotations

# Genome from the bioinformatic perspective

- versioned reference

# Where does the genomic data come from?
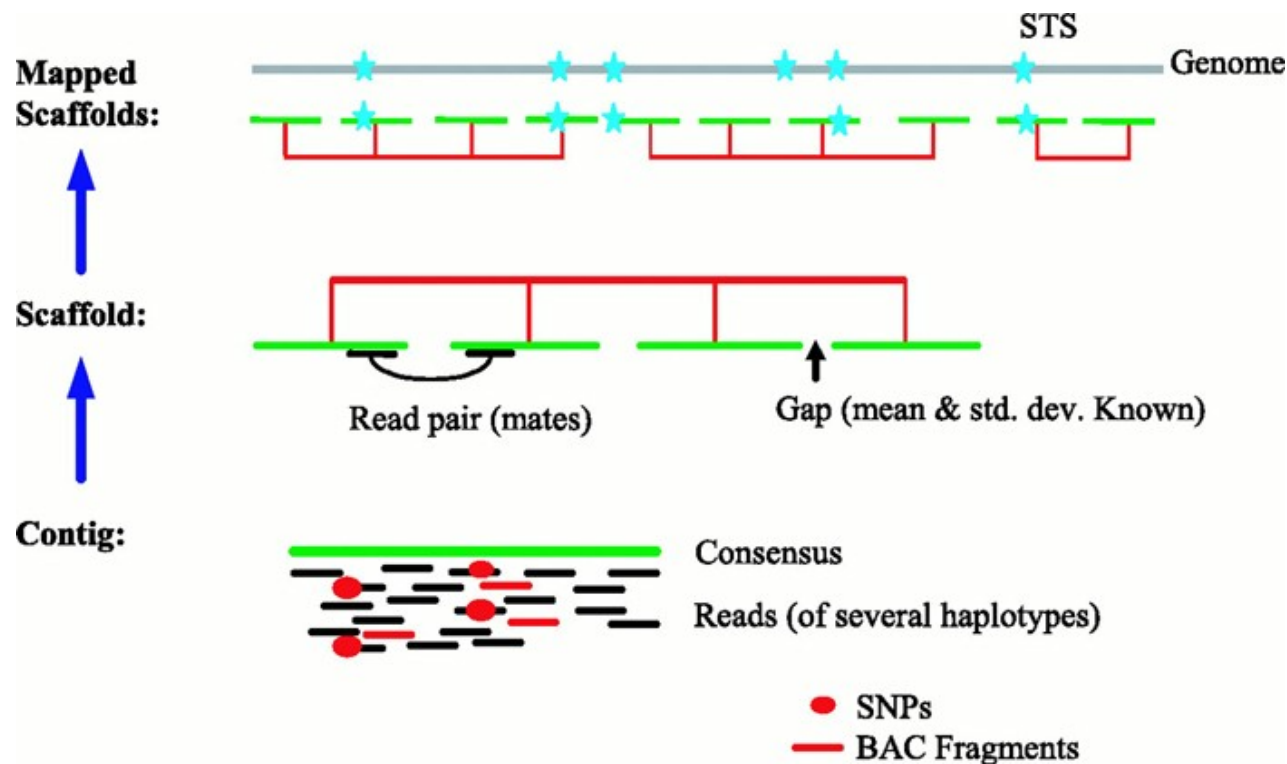
# Get a sequence

- Various methods:
  - NGS: Illumina, IonTorrent
  - TGS: PacBio, NanoPore
- They all produce short stretches of DNA (<u>reads</u>) of various length (100 bp - 100 kbp)
- Reads can form <u>pairs</u> (i.e. physical distance known between them) which is used for assembly
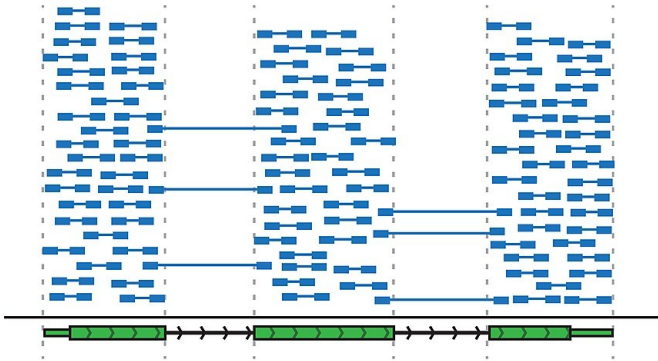
# Map the sequence

- *Reads* are *assembled* into continuous *contigs*
- *Paired-end reads* help to create a *scaffold* of contings
- Scaffolds are then mapped to *chromosomes*
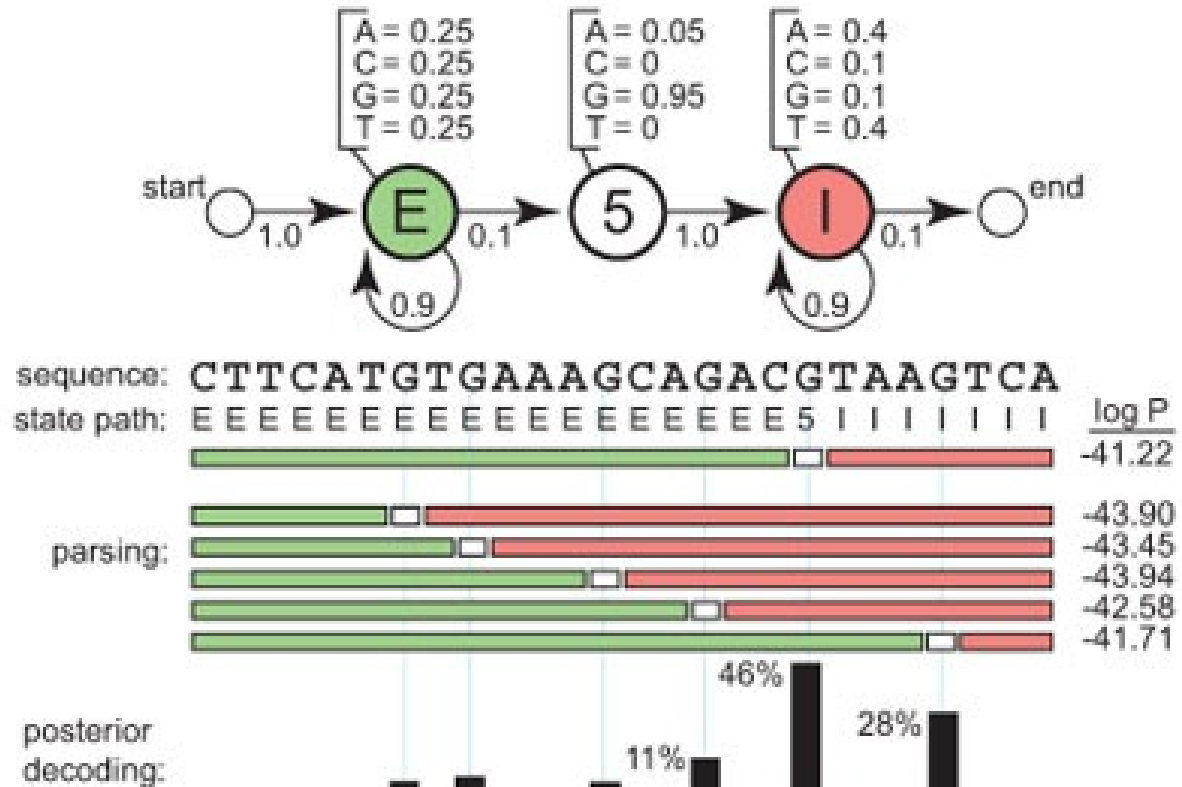
# Sequence annotation

- sequence similarity:
  - to known features (sequence similarity to ESTs, RNA-seq)
  - to homologous features in other organisms (homology – gene/protein families)

# Sequence annotation

- feature prediction using models:
  - using Hidden Markov Models to predict gene structure

# Sequence annotation

- Other non-coding functional elements
  - TF binding sites, etc.
  - interspecies sequence conservation
  - ChIP-seq (protein-DNA interaction)
  - DNAseI Hypersensitive Sites (open chromatin sites)

# Sequence annotation

- Other features
    - Variation data (SNPs, INDELS)
    - Structural variation data (CNVs)
    - Repeat data (RepeatMasker)
    - Epigenomic data (methylation, histone acetylation)
    - Functional data (Gene Ontology, KEGG, …)
    - Gene Expression

# Where are genomic data stored?

# Common genomic data formats

# Common genomic data formats

- Regular text files of a specific format
  - easy to open and explore
  - easy to work with
  - .fasta, .fastq, .bed, .gff, .gtf, .vcf, ...

- Binaries
  - more efficient for large datasets
  - fast retrieval by specific tools
  - .2bit, .gz, .bcf

# Storing sequences: FASTA

```
>ID_seq|specific_info
AGTGGGCGAGGCGCGGAGGTCTGGCCTATAAAGTAGTCGCGGAGACGGGG
TGCTGGTTTGCGTCGTAGTCTCCTGCAGCGTCTGGGGTTTCCGTTGCAGT
CCTCGGAACCAGGACCTCGGCGTGGCCTAGCGAGTTATGGCGACGAAGGC
CGTGTGCGTGCTGAAGGGCGACGGCCCAGTGCAGGGCATCATCAATTTCG
AGCAGAAGGCAAGGGCTGGGACGGAGGCTTGTTTGCGAGGCCGCTCCCAC
CCGCTCGTCCCCCGCGCACCTTTGCTAGGAGCGGGTCGCCCGCCAGGCC
TCGGGGCCGCCCTGGTCCAGCGCCCGGTCCCGGCCCGTGCCGCCCGGTCG
GTGCCTTCGCCCCCAGCGGTGCGGTGCCCAAGTGCTGAGTCACCGGGCGG
```

# Storing reads: FASTQ

```
@ID_seq1

AGTGGGCGAGGCGCGGAGGTCTGGCCTATAAAGTAGTCGCGGAGACGGGG

+          ASCII

!''*(((*(**+))%%%++)(%%%).1***-+*''))**55CCF>>>>>

@ID_seq2

CCTCGGAACCAGGACCTCGGCGTGGCCTAGCGAGTTATGGCGACGAAGGC

+

')%'* (***+)*''))*%%++5(
```

ASCII Table

ASCII = American Standard Code for Information Interchange

# FASTQ: ASCII to PHRED

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.................................................
.......................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX...........................
..........................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII...................
.............................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ...................
..LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL...............................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                               |     |        |                               |     |
33                              59    64       73                              104   126
 0.......................26...31.......40
                         -5....0........9............................40
                               0........9............................40
                               3.....9............................40
 0.2.......................26...31.......41
```

S - Sanger          Phred+33,  raw reads typically (0, 40)
X - Solexa          Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64,  raw reads typically (3, 40)
    with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
    (Note: See discussion above).
L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)

# PHRED: quality scores

| Phred Quality Score | Probability of incorrect base call | Base call accuracy |
|---|---|---|
| 10 | 1 in 10 | 90% |
| 20 | 1 in 100 | 99% |
| 30 | 1 in 1000 | 99.9% |
| 40 | 1 in 10,000 | 99.99% |

$$\text{Phred} = -10 \log_{10} P$$

# Storing annotations: GFF/GTF

- GFF
  - General Feature Format (any kind of annotation/feature)
- GTF
  - Gene Transfer Format (specific form of GFF used to store gene annotation)
- 9 TAB separated fields
- actual content of individual fields depends on the database and type of data

| seqname | source | feature | start | end | score | strand | frame | attribute |
|---|---|---|---|---|---|---|---|---|
| 2 | protein_coding | CDS | 2419108 | 2419128 | . | + | 0 | gene_id "ENSG00000223972"; |
| X | protein_coding | CDS | 1186934 | 1440976 | . | – | 0 | gene_id "ENSG00000123546"; |

gene_id "ENSG00000223972"; gene_name "DDX11L1"; gene_source "havana"; gene_biotype "protein_coding";

tag "value";

# Storing annotations: BED

- 3/4/6/12 columns
- used by UCSC Genome Browser to visualize various features

| chrom | chromStart | chromEnd | name | score | strand |
|-------|-----------|----------|------|-------|--------|
| 2 | 2419108 | 2419128 | ENSG00000223972 | . | + |
| X | 1186934 | 1440976 | ENSG00000123546 | . | – |

# Storing annotations: BED

- 0-based vs. 1-based coordinate system

| chr1 | T | A | C | G | T | C | A |
|------|---|---|---|---|---|---|---|
| 1-based | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| 0-based | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |

|  | 1-based | 0-based |
|---|---|---|
| Indicate a single nucleotide | chr1:4-4  G | chr1:3-4  G |
| Indicate a range of nucleotides | chr1:2-4  ACG | chr1:1-4  ACG |
| Indicate a single nucleotide variant | chr1:5-5  T/A | chr1:4-5  T/A |

# Storing variation data: VCF

- Variant Call Format

```
##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##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=.,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=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##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      Sample1
2      4370    rs6057    G    A       29   .      NS=2;DP=13;AF=0.5;DB;H2          GT:GQ:DP:HQ 0|0:48:1:52,51
2      7330    .         T    A       3    q10    NS=5;DP=12;AF=0.017              GT:GQ:DP:HQ 0|0:46:3:58,50
2      110696  rs6055    A    G,T     67   PASS   NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27
2      130237  .         T    .       47   .      NS=2;DP=16;AA=T                  GT:GQ:DP:HQ 0|0:54:7:56,60
2      134567  microsat1 GTCT G,GTACT 50   PASS   NS=2;DP=9;AA=G                   GT:GQ:DP    0/1:35:4
```

```
< /data-shared/vcf_examples/luscinia_vars_flags.vcf.gz zcat |  less -S
```

# Storing variation data: VCF

- Variant Call Format

```
##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##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=.,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=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##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     Sample1
2      4370   rs6057   G    A      29   .      NS=2;DP=13;AF=0.5;DB;H2           GT:GQ:DP:HQ 0|0:48:1:52,51
2      7330   .        T    A      3    q10    NS=5;DP=12;AF=0.017               GT:GQ:DP:HQ 0|0:46:3:58,50
2      110696 rs6055   A    G,T    67   PASS   NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT:GQ:DP:HQ 1|2:21:6:23,27
2      130237 .        T    .      47   .      NS=2;DP=16;AA=T                   GT:GQ:DP:HQ 0|0:54:7:56,60
2      134567 microsat1 GTCT G,GTACT 50  PASS   NS=2;DP=9;AA=G                    GT:GQ:DP    0/1:35:4
```

**Header part**
**(description of abbreviations used in the data part)**

**Data part**

# Storing variation data: VCF

- Variant Call Format

```
##fileformat=VCFv4.0
##fileDate=20110705
##reference=1000GenomesPilot-NCBI37
##phasing=partial
##INFO=<ID=NS,Num...
##INFO=<ID=DP,Num...
##INFO=<ID=AF,Num...
##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,...
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
```

Table: Variants (rows) vs. Samples (columns)

(description of abbreviations used in the data part)

**Samples + Genotypes**

Variation details (location, quality, type, etc.)

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | Sample1 |
|--------|-----|-----|-----|-----|------|--------|------|--------|---------|
| 2 | 4370 | rs6057 | G | A | 29 | . | NS=2;DP=13;AF=0.5;DB;H2 | GT:GQ:DP:HQ | 0|0:48:1:52,51 |
| 2 | 7330 | . | T | A | 3 | q10 | NS=5;DP=12;AF=0.017 | GT:GQ:DP:HQ | 0|0:46:3:58,50 |
| 2 | 110696 | rs6055 | A | G,T | 67 | PASS | NS=2;DP=10;AF=0.333,0.667;AA=T;DB | GT:GQ:DP:HQ | 1|2:21:6:23,27 |
| 2 | 130237 | . | T | . | 47 | . | NS=2;DP=16;AA=T | GT:GQ:DP:HQ | 0|0:54:7:56,60 |
| 2 | 134567 | microsat1 | GTCT | G,GTACT | 50 | PASS | NS=2;DP=9;AA=G | GT:GQ:DP | 0/1:35:4 |

**Data part**

# Specialized tools for genomic data

# samtools

- Working with SAM/BAM files (i.e read alignment data)
- Manipulation with SAM/BAM (sorting, merging, subsetting)
- Summary statistics (read depth by position)
- Viewing read alignment in command line:
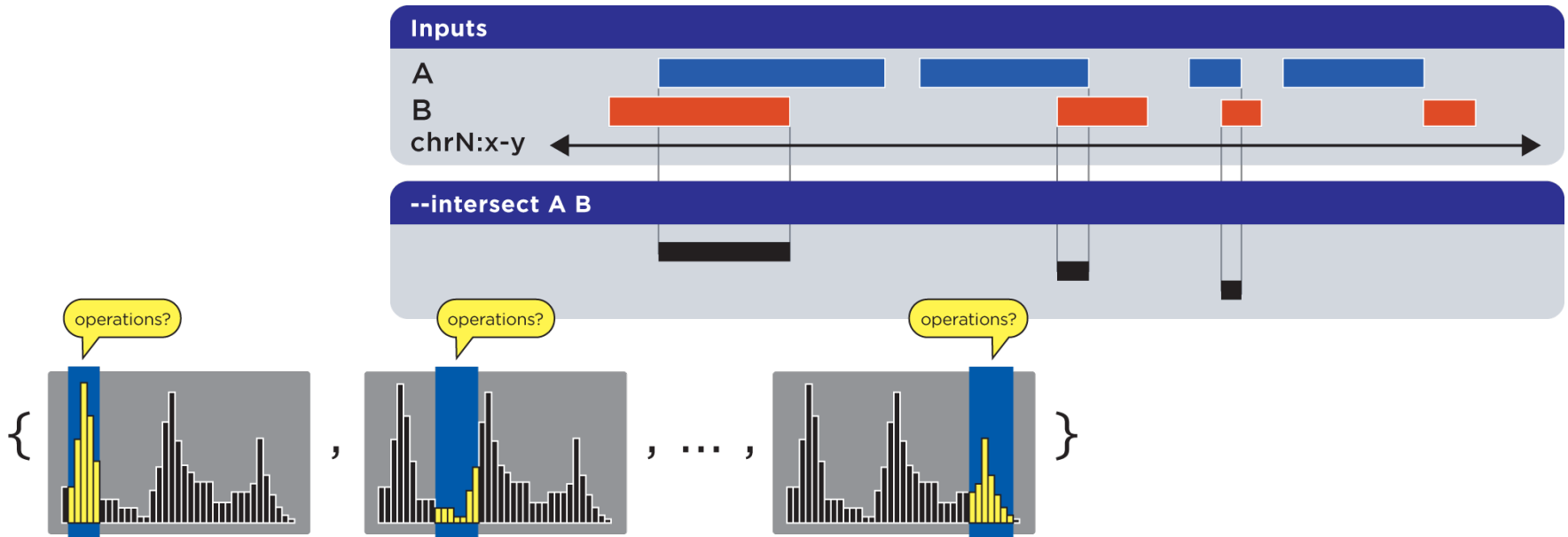


https://samtools.github.io

# bcftools/vcftools

- variant call files (vcf/bcf)
- bcftools:
  - annotation, concatenation, merging, converting to different formats, filtering based on various criteria, variant calling
- vcftools:
  - mainly filtering/creating subsets
  - population genetics (allele frequency, Hardy-Weinberg, Fst, Pi, Tajima, linkage disequilibrium,...)

https://vcftools.github.io/index.html
https://samtools.github.io/bcftools/

# bedtools/bedops

- Operations with genomic data based on their physical position in genome (chromosome, feature start, feature end, strand)
- Usually intersections, overlaps, summary by specific regions (e.g. coverage), sliding window analysis, randomization

# What did we learned?

- How does genome look from the bioinformatic perspective

- Where does the genomic data come from?

- Common genomic data formats

- Specialized tools for genomic data