# Omics in human diseases
## Index

- Omics data and Biological databases
- NGS methods
- NGS data analysis
- **Prediction and interpretation of pathogenic variants**
- Protein-protein interaction networks

**Course organization 2022/2023**
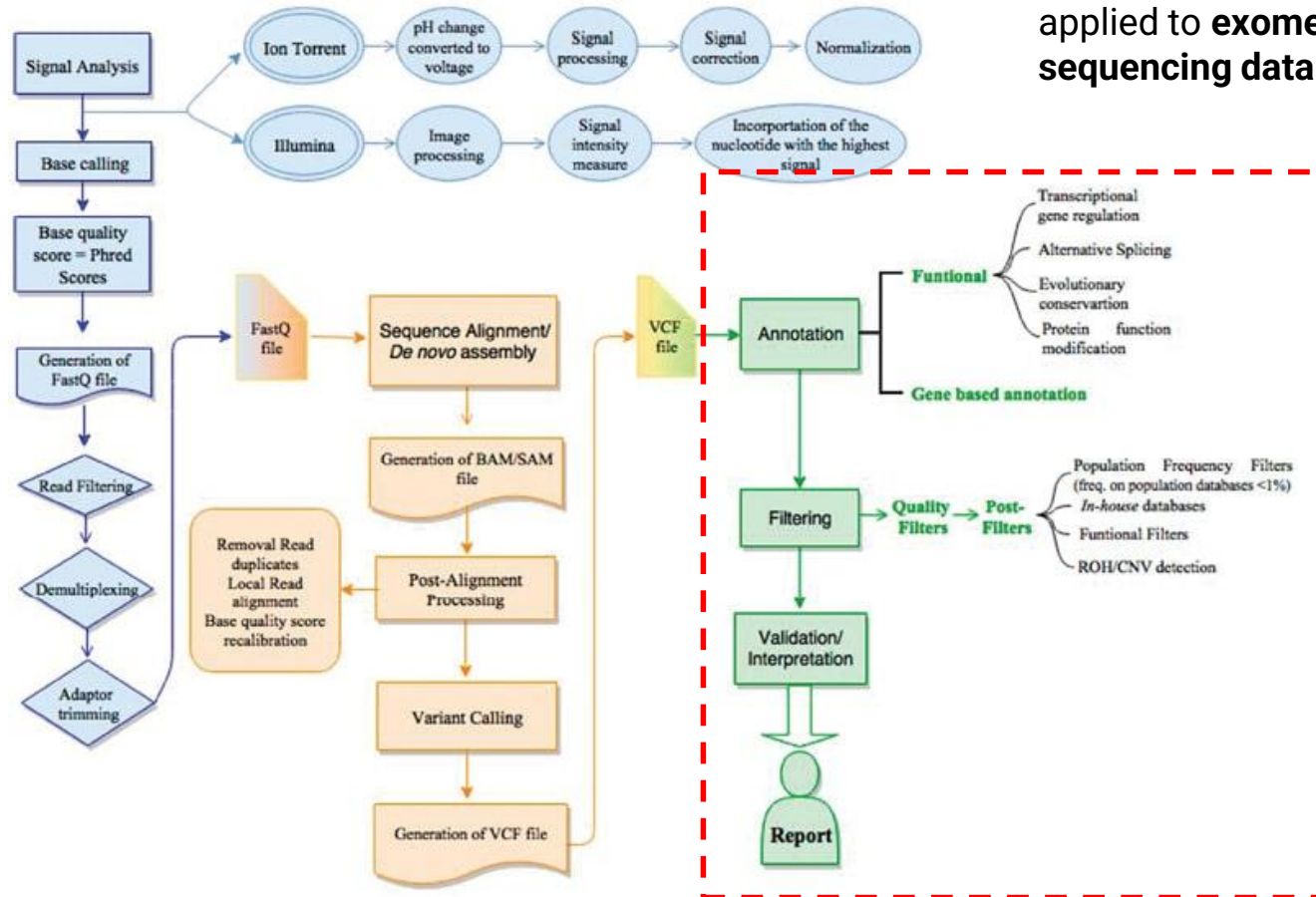
Frontal lecture/ guided practical activity

How to pass the exam: multiple choice quiz (50%) + results from practical activities (50%) + Bonus points, e.g. summary of previous lecture (up to 10%)

Mail: **emanuela.leonardi@unipd.it**

BioComputing

# NGS analysis workflow

The same data analysis tools used for WGS can be applied to **exome-sequencing data**

# Variant filtering: criteria

- Genes list
- Sequencing parameters (filter artifactuals)
- Variant types (exonic – intronic)
- Variant class (frameshift - synonymous)
- Population frequency
- **Pathogenicity prediction score**
- **Conservation score**
- Clinical interpretation of genetic variants by the ACMG/AMP 2015 guideline
- Presence in mutation databases

*(PMID: 28118812;*
*https://www.nature.com/articles/s41525-021-00227-3)*

BioComputing

# Variant filtering: pathogenicity scores

<u>Computational methods to infer variant pathogenicity</u>

- performance of current methods: 90% of pathogenicity predictions are correct, such methods **identified only 10%–20% of pathogenic variants**

- Current guidelines for clinical variant interpretation recommend that all computational methods be (at best) treated as "**weak evidence**".

- MaveDB: "variant effect maps" , fewer than 1% of the 4,000 human disease-associated proteins
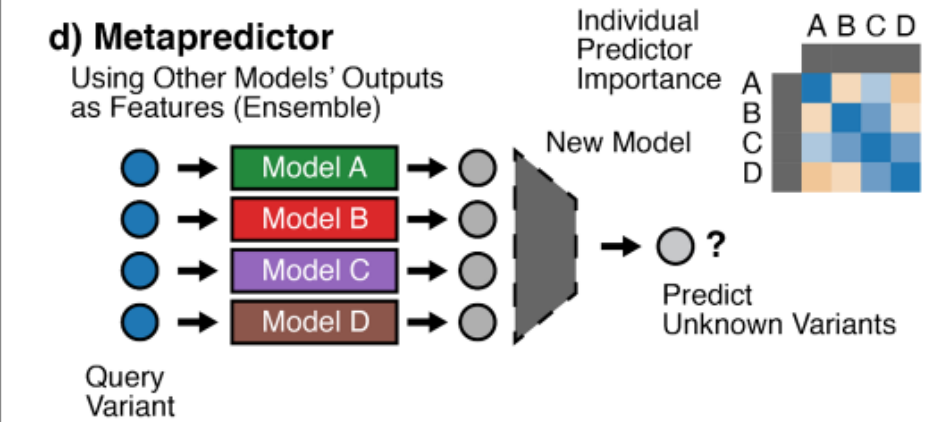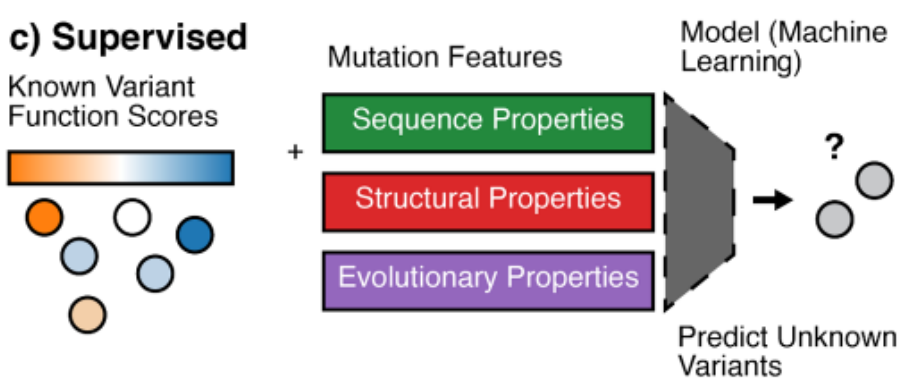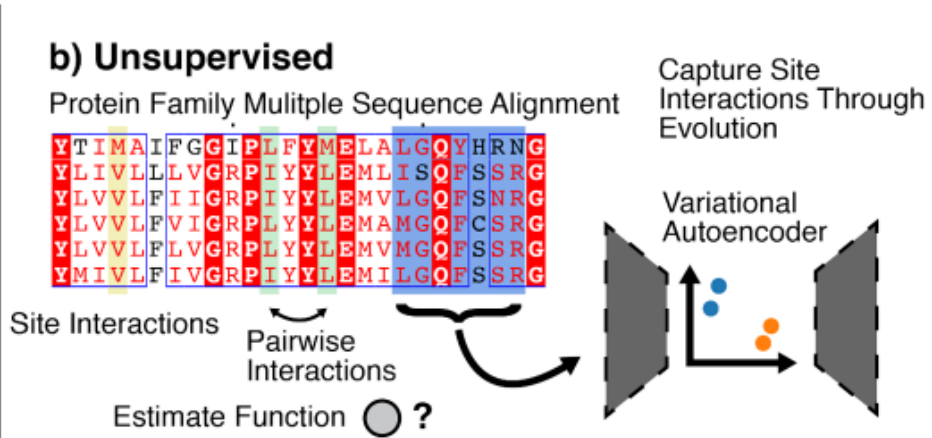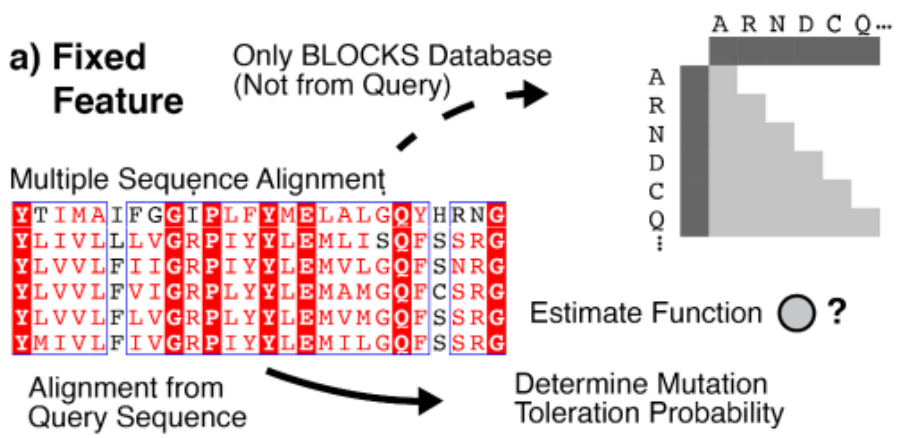
# Variant filtering: Annovar pathogenicity scores

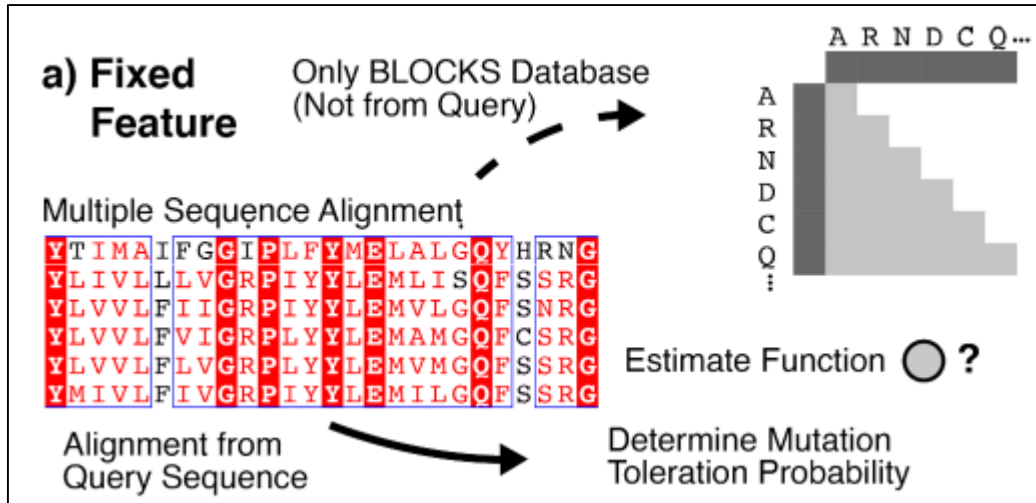| Score (dbtype) | Categorical Prediction |
|---|---|
| SIFT (sift) | D: Deleterious (sift<=0.05); T: tolerated (sift>0.05) |
| PolyPhen 2 HDIV (pp2 hdiv) | D: Probably damaging (>=0.957), P: possibly damaging (0.453<=pp2_hdiv<=0.956); B: benign (pp2_hdiv<=0.452) |
| PolyPhen 2 HVar (pp2 hvar) | D: Probably damaging (>=0.909), P: possibly damaging (0.447<=pp2_hdiv<=0.909); B: benign (pp2_hdiv<=0.446) |
| LRT (lrt) | D: Deleterious; N: Neutral; U: Unknown |
| MutationTaster (mt) | A" ("disease_causing_automatic"); "D" ("disease_causing"); "N" ("polymorphism"); "P" ("polymorphism_automatic" |
| MutationAssessor (ma) | H: high; M: medium; L: low; N: neutral. H/M means functional and L/N means non-functional |
| FATHMM (fathmm) | D: Deleterious; T: Tolerated |
| PROVEAN pred | |
| MetaSVM (metasvm) | D: Deleterious; T: Tolerated |
| MetaLR (metalr) | D: Deleterious; T: Tolerated |
| M-CAP pred | |
| CADD | D >25 |
| fathmm-MKL_coding_pred | D: Deleterious; N: Neutral |
| GERP++ (gerp++) | higher scores are more deleterious (>3) |
| PhyloP (phylop) | higher scores are more deleterious |
| SiPhy (siphy) | higher scores are more deleterious |
| DANN | |
| Eigen | |

# Variant filtering: pathogenicity scores

| Approach | Training set | Conservation analysis | Structural attributes | Annotations |
|---|---|---|---|---|
| Randome forest; Alignment score; Bayesan classification; SVM (Support Vector Machine); Machine learning | HGMD, Swiss Prot; Protein Mutant Database | PSIC, position-specific independent counts; PFAM; PSI-BLAST; Sequence environment, sequence profiles | Predicted attributes; Homologue mapping | Swiss-Prot; Pfam domain; GO |

BioComputing

# Variant filtering: Variant effect prediction (VEP) tools

# Variant effect prediction (VEP) tools



a) **Fixed Feature**

Only BLOCKS Database (Not from Query)

Multiple Sequence Alignment

Estimate Function ◯ ?

Alignment from Query Sequence

Determine Mutation Toleration Probability

A R N D C Q...


Sorting Intolerant From Tolerant

- Evaluate **mutation toleration** at a given site
- Direct statistics from given inputs and
- features
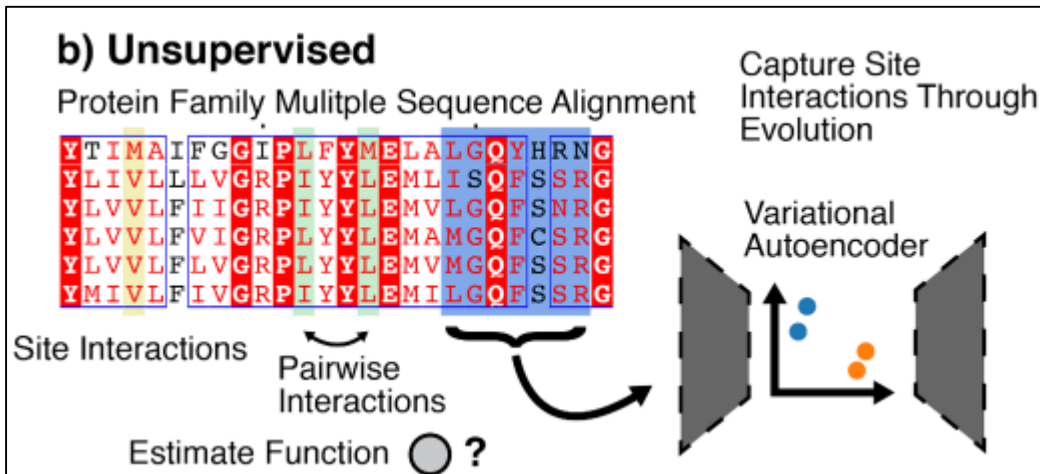- based on averaged quantities about amino acid frequency and amino acid properties (size, charge, and hydrophobicity)

**Statistical methods**:

**BLOSUM62**, a block substitution scoring Matrix, used as the default scoring matrix for many multiple sequence alignment (MSA) methods.
**SIFT (and SIFT 4G)** is the commonly used static feature method

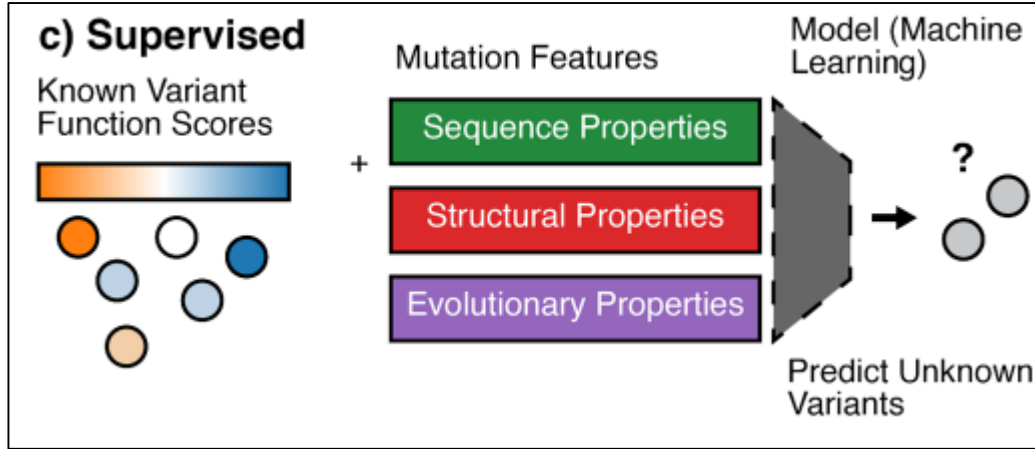BioComputing

# Variant effect prediction (VEP) tools



- **do not fit experimental data**

- capture the frequencies and dependencies among amino acid residues given **evolutionary pressure**

- dependencies among residues (**epistasis**) affect molecular function

**EVmutation** is not a deep learning-based (DL) method but rather a statistical model capturing dependencies across pairs of residues

**DeepSequence,** variational autoencoder framework to model the evolutionary fitness landscape, from a protein family's evolutionary history

While unsupervised methods readily generalize across protein space, they neglect to learn from the many labeled mutagenesis **datasets** appearing in the literature

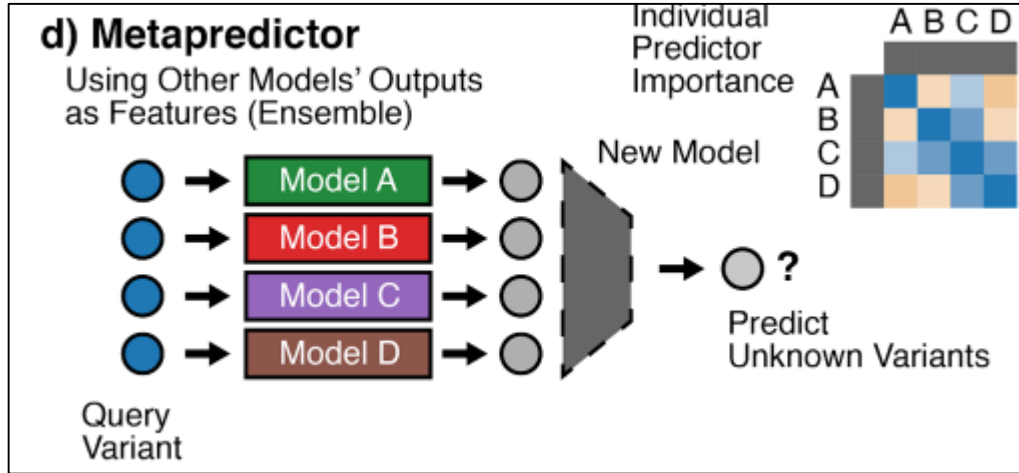# Variant effect prediction (VEP) tools



- Use large-scale datasets
- Rely on fitting models to experimental
- measurements of the fitness landscape often is derived from Deep mutational scan (DMS) experiments

**SNPs&GO** a support vector machine (SVM) ML model was used to categorize mutations as disease-causing or not

Envision a random forest (RF) regressor on DMS datasets of eight proteins

using datasets for assessing model performance, caution should be taken as model architectures, parameters, and hyperparameters may over-represent training set protein families, and real-world prediction accuracy may be overstated

# Variant effect prediction (VEP) tools



d) **Metapredictor**
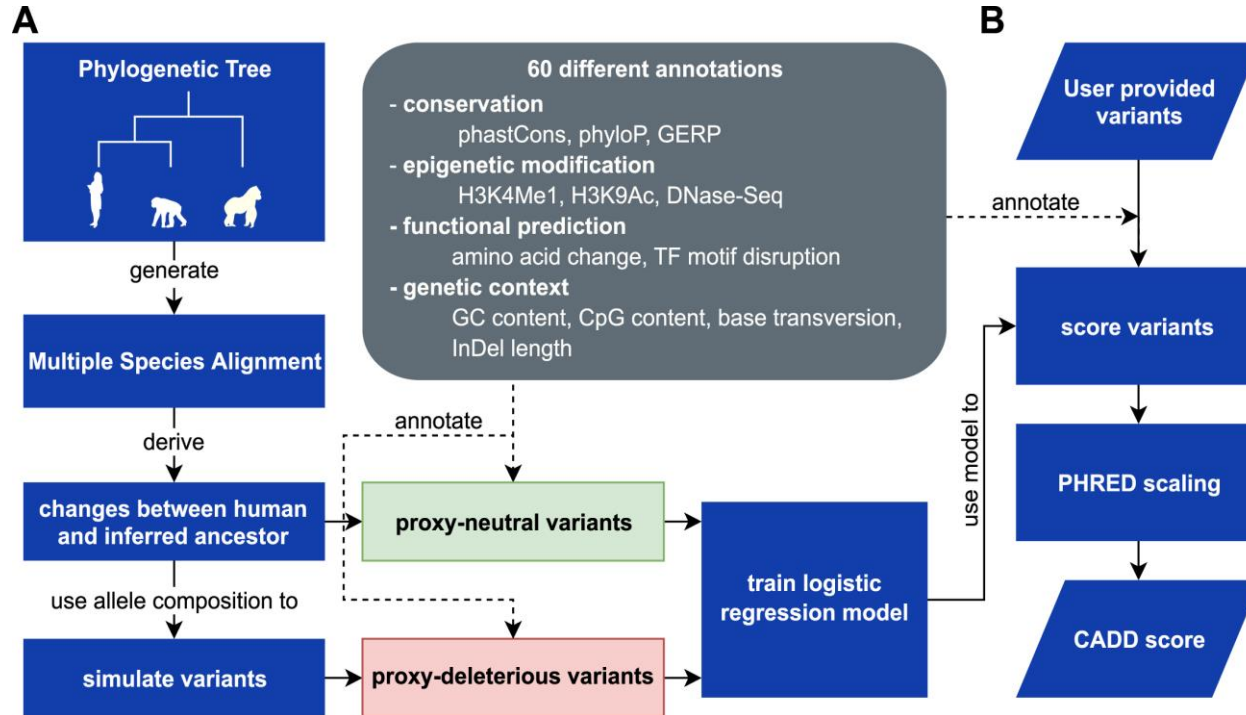Using Other Models' Outputs as Features (Ensemble)

- Leverage the predictive power of ensembling models together to improve performance
- The outputs of other VEP models are used as inputs and typically trained similarly to supervised learning methods

**REVEL** (Rare Exome Variant Ensemble Learner) combines 13 other VEP tools as features (MutPred, FATHMM v2.3, VEST 3.0, PolyPhen-2, SIFT, PROVEAN, MutationAssessor, MutationTaster, LRT, GERP++, SiPhy, phyloP, and phastCons)

- is an RF ML method trained on rare neutral and disease variants
- The features of greatest significance in the developed RF were the FATHMM and VEST models.
- FATHMM employs a hidden Markov modeling method to analyze MSAs, which is unique among the other techniques

# Combined Annotation-Dependent Depletion (CADD)



Using more than 60 diverse annotations, a machine learning model is trained to classify variants as *proxy-neutral* versus *proxy-deleterious*

# Combined Annotation-Dependent Depletion (CADD)

Score that ranks genetic variants, including single nucleotide variants (SNVs) and short inserts and deletions (InDels), throughout the human genome reference assembly

Reference genome SNVs at the 10th-% of CADD scores are assigned to CADD-10
top 1% to CADD-20,
top 0.1% to CADD-30

## The CADD score

Various prediction tools in color coding from damaging to tolerated

Gene Y

Cancer

Normal

Gene X

kernel

Cancer

Decision surface

Normal

Support vector machines use a transformation of "messy data" into a higher dimensional structure where data points in both groups can be separated easily by a "hyperplane"
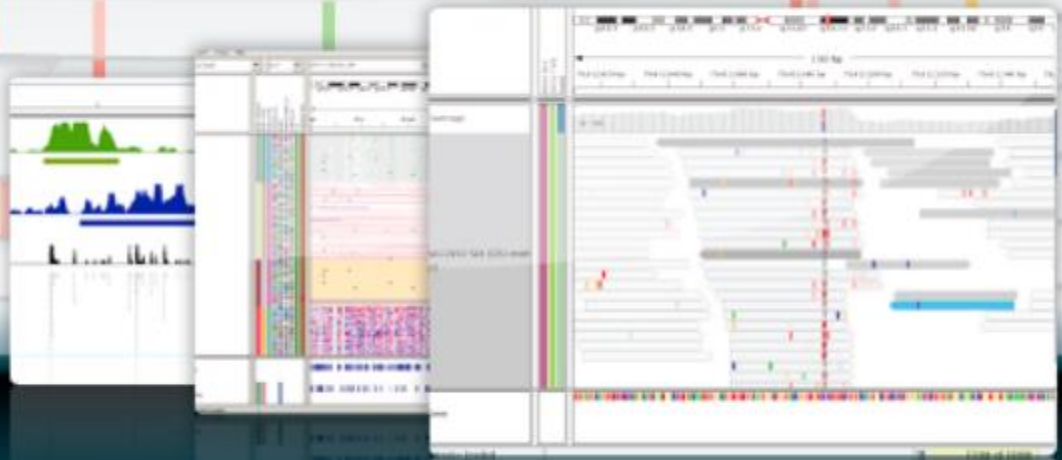
The Epi4K *de novo* mutations sorted by CADD score

BioComputing

# Variant filtering



Total Number of variants

Filtered by Internal Lab database
Filtered low quality variants

Filtered by gene list (Not for gene panel)

Removed intronic and synonimous variants

Filtered by population frequency

Filtered by pathogenicity score
**< 10 variants to check**

# Variant Filtering: visual inspection

# Variant Filtering: Integrative Genome Viewer (IGV)

Visual inspection can greatly increase the confidence in calls, reduce the risk of false positives, and help characterize complex events

A number of **IGV features** have been developed specifically to aid this manual review step:

- highlighting mismatched bases in individual reads in color to aid detection of unusual patterns and **mismapped alignment**

- highlighting ambiguously mapped reads (**mapping quality**), indicative of high reference sequence homology, as such regions are known to produce many false positives

- shading of mismatched bases by **read base quality**, as clusters of bases of low quality can be indicative of sequencing errors

- **sorting, grouping, and coloring alignments** by alignment, sequencing, and platform metadata, which can be useful for detecting **systematic errors** upstream of read alignment

BioComputing

# Variant Filtering: Integrative Genome Viewer (IGV)

- User Interface

- Download Reference Genome (Human hg19 for the gene panel sequencing)

- Viewing Sequencing Data with IGV

  - Alignment (BAM, SAM)

  - Variants (VCF)

https://youtu.be/E_G8z_2gTYM

BioComputing

# IGV: user interface
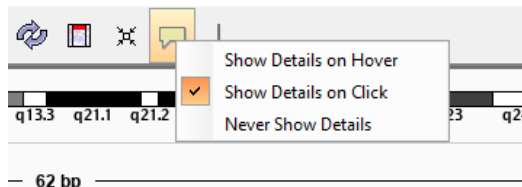


Search by gene or chromosome position

Select the reference genome

Genome data track

Sample data track

Reference sequence track

# IGV: viewing alignment

- Upload a **BAM** or SAM file

- The main data file must include the *.bam*

- The index file should have the same filename but with the *.bai*

- When loading by file, IGV automatically searches for the index file within the same directory as the data file

# IGV: viewing alignment

Coloring and sorting alignments
1. Use **Pop-up menù**

2. From menù
**View/Preferences**

# IGV: viewing variant

- Upload a **VCF** file

- VCF data files must be indexed for viewing in IGV

- Search for a variant by chr position

- Use pop-up menù or change preferences to change color code

# IGV: viewing variant (nucleotide substitution)

chr16:89352586:C:T

ANKRD11

# IGV: viewing variant (nt substitution)



chr16:89352586:C:T

ANKRD11

NM_013275:exon8

c.753G>A p.K251K

Synonymous variant

# IGV: viewing variant (nt deletion)

**chrX:53285126:AGGGGGGC:AGGGGG,AGGGGGC**

# IGV: viewing variant (nt deletion)

**chrX:53285126:AGGGGGGC:AGGGGG,AGGGGGC**

GQ: 557, DP: 461, AF: 0.426304,0

IQSEC2:NM015075:exon3:c.239delC:p.P80fs

IQSEC2:NM015075:exon3:c.233-234del:p.G78fs

Annowar annotation from dbsnp: rs782460038

NM_015075.2:c.239del

**Frequency!!!**
delG=0.00000 (0/44894, ExAC)

Is in a polymeric region

Is present in with other VCFs

**Probably error sequencing**

# Variant interpretation: criteria