# Variants prioritization: annotation and filtration steps



**2nd Variant Effect Prediction Training**
**Prague, Czech Republic, 2017**
Prof. Christophe Béroud,
Aix-Marseille University
INSERM UMR_S910 Medical Genetics and Functional Genomics

# RD Connect

| WP1: Coordination | WP2: Patient registries | WP3: Biobanks | WP4: Bioinformatics |
|---|---|---|---|
| **Hanns Lochmüller** Newcastle and TREAT-NMD | **Domenica Taruscio** ISS and EPIRARE | **Lucia Monaco** Fondaz. Telethon & EuroBioBank | **Christophe Béroud** AMU & INSERM Marseille |

| WP5: Unified platform | WP6 Ethical/legal/social | WP7: Impact and innovation |
|---|---|---|
| **Ivo Gut** CNAG Barcelona | **Mats Hansson** Uppsala | **Kate Bushby** Newcastle and EUCERD/ EJARD |

Aix Marseille université   Inserm Institut national de la santé et de la recherche médicale
GENETICS & BIOINFORMATICS TEAM

# The "genetics and bioinformatics" team

Medical school of la Timone



Marseille - France

Aix-Marseille Université and INSERM institute dedicated to Medical Genetics and functional Genomics
*"Translational research in Rare Diseases"*

# Our team

## Genetics and Bioinformatics team – Multidisciplinary group

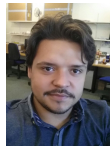**Team Leader**

Christophe Béroud
*Professor*

**Bioinformatics**

David Salgado
*PhD*

Céline Guien
*Engineer*

Jean-Pierre Desvignes
*Engineer*

Marc Garibal
*Engineer*

Coralie Grattepanche
Master *Student*

Mélanie Corcuff
Master *Student*

**Genetics of Dystonia**

Gwenaëlle Collod-Béroud
*Researcher*

Arnaud Blanchard
*PhD*

Kahina Medjber
*PhD*

Constance Renault
Master Student

# Our main projects

**Locus Specific Mutation databases - UMD (Universal Mutation databases)**

Gather >208,000 manually expert curated mutations in more than 60 databases:

- Genes involved in cancers (*APC, BRCA1, BRCA2, TP53, RB1, MEN1, SUR1, VHL, WT1*…)
- Genes involved in genetic disorders (*FBN1, LDLR, DMD, VLCAD, MCAD, LMNA, EMD, FKRP, SGCG, SGCA, ATP7B, TREAT-NMD_DMD, TREAT-NMD_SMA*…)

**Patients registries – Population databases**

Global TREAT-NMD DMD & SMA

International Dysferlinopathy Registry

French Database for Marfan and related Syndromes

National databases CNVs (BANCCO) + SNVs (RDVD Rare disease variant database)

**Pathogenicity prediction systems**

UMD-Predictor (http://umd-predictor.eu)

Human Splicing Finder (http://umd.be/HSF3/)

**Next generation sequencing**

NGS Data analysis

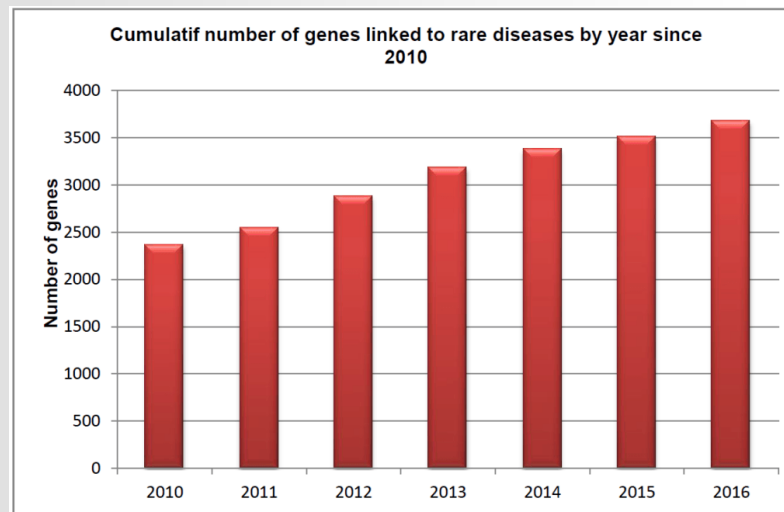Variant Annotation and Filtration tool (VarAFT)

**Clinical tools**

Skip-E for Antisense oligonucleotides to induce exon skipping

NR-Analyzer for nonsense mutations eligible for non-sense read-through

Crawfish for trans-splicing

# Context

- Next generation Sequencing has facilitated the discovery of new genes and genetic variants in a multitude of human disorders

- 1st Whole-Exome Sequencing (WES) done by Ng et al 2009 in Miller Syndrome

- In 2013, >150 Mendelian disorders were studied by WES (Rabbani *et al*, J Hum Genet, 2014)

- IRDiRC and Orphanet ➤ 3,700 genes involved in RD, >1,300 identified between 2010 and 2016 (IRDiRC website)



**Cumulatif number of genes linked to rare diseases by year since 2010**

IRDiRC
INTERNATIONAL
RARE DISEASES RESEARCH
CONSORTIUM

Aix Marseille université    Inserm
Institut national
de la santé et de la recherche médicale

GENETICS & BIOINFORMATICS TEAM
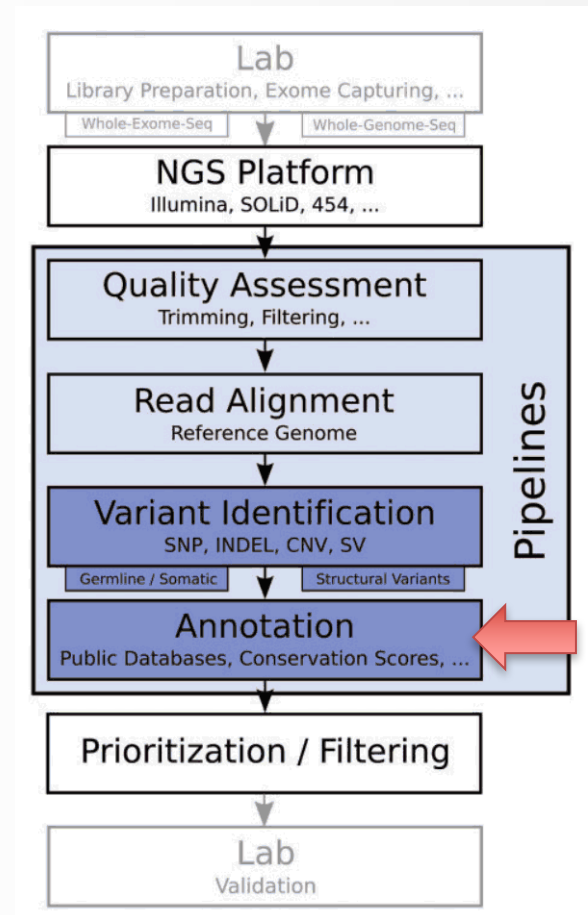
# …But…

- Despite all these encouraging figures

  - Only 23-26% of WES are successful (higher rate if several individuals from the same family are sequenced 34-37% for a trio) [Farwell et al. 2015, Sawyer et al. 2016]

    - Technical factors (homopolymers, GC reach regions, poor quality at read ends …)

    - Type of disease causing mutations (not captured, triplet repeat expansions, CNVs, pseudogenes …)

    - Bioinformatics pipeline to generate VCF (same sequencing technology, not same VCF)

    - Wrong annotations/filtrations

# …But…

- Despite all these encouraging figures

  - Only 23-26% of WES are successful (higher rate if several individuals from the same family are sequenced 34-37% for a trio) [Farwell et al. 2015, Sawyer et al. 2016]

    - Technical factors (homopolymers, GC reach regions, poor quality at read ends …)

    - Type of disease causing mutations (not captured, triplet repeat expansions, CNVs, pseudogenes …)

    - Bioinformatics pipeline to generate VCF (same sequencing technology, not same VCF)

    - **Wrong annotations/filtrations**

**Christophe Béroud– 2nd VEP Training**
**6th - 8th November 2017**

Aix Marseille universite    Inserm
Institut national de la santé et de la recherche médicale
GENETICS & BIOINFORMATICS TEAM

8

# Variant annotation

- Part of the data analysis process

    - Mandatory for prioritization and filtration of variants

- Two objectives

    - Help to refine our estimate of how likely a variant is to be true, genotype, quality …

    - Provide functional annotations to determine the links between a genetic variation and a disease



*Adapted from Pabinger et al. Briefings in Bioinformatics 2014*

# Variant annotation

- ■ It is performed at various levels

# Variant annotation



Variant level

**Quality data**
Quality score, quality filter, genotype, read depth, …

**Mutation localization**
exonic, intronic, splice site, 3'UTR, 5'UTR, … mutation type and HGVS nomenclature (g., c., p.)

**Mutation frequency**
EVS, dbSNP, 1000G, ExAC, …

**Pathogenicity predictions**
SIFT, Polyphen, CADD, Mutation taster, UMD-Predictor …

# Variant annotation



**Expression profile**
GTEx, Expression Atlas, RNA-seq data,…

**Gene Ontology**
Molecular function, cellular component, biological process

**Biological Pathways**
KEGG, BioCarta, WikiPathways, Reactome

Gene level

Variant level

GENETICS & BIOINFORMATICS TEAM

# Variant annotation



Phenotype level

Gene level

Variant level

**Involvement in Human diseases**
OMIM, UNIPROT, LSDB (LOVD, UMD,...), Clinvar, COSMIC

**Animal models (mutants, knockdowns, transgenics)**
MGI, Zfin, FlyBase, ...

# Variant annotation

# Variant annotation systems

| | Annovar | SNPeff | Ensembl VEP | SeattleSeq | AnnTools | Oncotator | Vanno | Variant Annotation Tools |
|---|---|---|---|---|---|---|---|---|
| Availability | Command line | Command line | Command line Webservices Web | web | command line | Command line Web | Web | Command line |
| Variant quality | Yes | Yes | Yes | - | Yes | - | Yes | Yes |
| Variant localization | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Gene/transcript annotation | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Genotype | Yes | Yes | Yes | Yes | Yes | - | Yes | Yes |
| Population frequency | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Impact at the RNA level | Yes | Yes | Yes | - | - | - | - | - |
| Impact at the protein level | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Conservation | Yes | Yes | Yes | Yes | - | Yes | Yes | Yes |
| Reported impact | - | - | Yes | Yes | - | Yes | Yes | Yes |
| Predicted pathogenicity | Yes | Yes | Yes | Yes | - | Yes | Yes | Yes |
| Gene ontology | Yes | - | - | - | - | Yes | Yes | - |
| Pathways | - | - | - | Yes | - | - | Yes | Yes |
| Tissue expression | - | - | - | - | - | - | - | - |

Various types of annotation software are available (command line/web)
No system is providing annotations at all levels ➤ need to be combined

# Variant filtration



Adapted from Pabinger et al. Briefings in Bioinformatics 2014

Number of variants

| Number of variants | | Filters |
|---|---|---|
| 60000 | | |
| 20000 | Raw variants – PASS Filter | |
| 2000 | Mode of inheritance (Combine genotypes) Autosomal Recessive disease - Trio | |
| 1800 | Localization Exonic-Splicing | |
| 50 | Filter based on population frequency data (dbsnp 137) | |
| 4 | Pathogenicity prediction software | |
| 2 | Other evidences | |
| | Validation | |

Identify disease causing/private mutations

# Filtration tools

- **Automatic systems**
  - **Disease causing genes based on pedigree and phenotypic data**

| Software name | Availability | Mode of inheritance | Custom analysis | Mutation localization | Mutation type | Mutation frequency | Pathogenicity predictions | Functional evidences | Clinical report | Prioritization score |
|---|---|---|---|---|---|---|---|---|---|---|
| **ExomeWalker** | Web App | Yes | - | - | - | Yes | No but provided | No but provided | Yes | Yes |
| **Exomiser** | Command line | Yes | - | - | - | Yes | Yes | - | Yes | Yes |
| **eXtasy** | Command line Web App | - | - | - | - | - | No but provided | No but provided | - | Yes |
| **MirTRIOS** | Web App | Yes | - | Yes | Yes | Yes | Yes | - | No - But provided | Yes |
| **OMIM Explorer** | Web App | Yes* | - | - | - | - | - | - | Yes | Yes |
| **OVA** | Web App | Yes | - | Yes | Yes | Yes (2) | - | Yes | Yes | Yes |
| **wKGGSeq** | Web App | Yes | - | Yes | Yes | Yes | Yes | Yes | Yes | Yes |

- Fast and accurate method for known genes/diseases
- Automatically gather additional information
- Work only with known genes/diseases with annotations
- Limited flexibility

\* Does not combine multiple samples

# Filtration tools

- Semi automatic/manual systems
  - Users can select candidate mutations by applying various set of filters

| Software name | Availability | Mode of inheritance | Custom analysis | Mutation localization | Mutation type | Mutation frequency | Pathogenicity predictions | Functional evidences | Clinical report |
|---|---|---|---|---|---|---|---|---|---|
| ANNOVAR | Command line | - | - | - | - | Yes | Yes | - | Yes |
| BIERapp | Web App | Yes | Yes | Yes | Yes | Yes | - | - | - |
| FILTUS | GUI | Yes | Yes | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided |
| FMFilter | GUI | Yes | - | Yes - if provided | Yes - if provided | Yes - if provided | - | - | - |
| Gemini | Command line Web App | Yes | Yes | Yes | Yes | Yes | Yes | Yes | Yes |
| Vanno | Web App | - | - | Yes | Yes | Yes | No but provided | Yes | Yes |
| VarAFT | GUI | Yes | Yes | Yes | Yes | Yes | Yes | Yes | No - But provided |
| VarSifter | Command line Web App | Yes | Yes | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided |
| VCF-MINER | Local Web App | Yes | Yes | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided | Yes - if provided |

- Provide a good flexibility and traceability
- Can be tedious and difficult
- No filtration at all annotation levels

# Challenges

Variant annotation and filtration **are not a simple** and **solved problem**

- **Variant Annotation**

    - Require the management of multiple and large data sources (local)

    - Incorrect or incomplete annotations can cause researchers to overlook and dilute interesting variants in a pool of false positives

    - From Davis McCarthy (Genome Medecine 2014)

        - Choice of transcript set and software can have a large effect on the variant annotation

            - Matching annotation for ANNOVAR on Ensembl or RefSeq genesets is only 83% for all exonic variants (lof, missenses, synonymous …)

            - Comparison between ANNOVAR and VEP on Ensembl transcripts 87% of all exonic variants were in agreement

            - Even more discrepancies with splicing variants

        - Where a specific tissue of interest is known, restrict annotation to transcripts known to be expressed in that tissue (GENCODE)

**Christophe Béroud– 2nd VEP Training**
**6th - 8th November 2017**

Aix Marseille université    Inserm
Institut national de la santé et de la recherche médicale
GENETICS & BIOINFORMATICS TEAM

19

# Challenges

- **Variant filtration/prioritization**

  - Good phenotypic description of patients

  - Mode of inheritance should be known (identify the right model, hypothesis on the penetrance …)

  - No Gold standard but frequently used filters

    - Frequency in the population

    - Genotype

    - Mutation type / Pathogenicity

    - Need to be done interactively (add/remove filters)

  - Discrepancies between pathogenicity predictors (may introduce false +/-)

  - Population frequency (not all databases gather healthy individuals, ethnicity) e.g. presence in dbSNP does not mean "polymorphism"

  - Privacy issue may arise when using "online" system

# Our systems

## To improve and facilitate disease causing mutation identification



Pathogenicity prediction systems
variant annotation

Variant annotation and
filtration

# UMD-Predictor
## http://umd-predictor.eu

- Pathogenicity prediction system for any Human cDNA substitutions

  - Precomputed all possible substitutions for all nucleotides of any human transcripts (280,315,899 substitutions)

  - Combined multiple features in a unique score (0-100)

    - AA change – substitution and biochemical matrices (BLOSUM/Yu)

    - Exonic splicing signal (HSF - Acceptors and Donors splice sites)

    - Protein key residues (UNIPROT HCD)

    - Conserved and functional domains -100 species protein alignments (Phastcons) + Grantham

    - Allele frequency

- Available through a web application and webservices

INFORMATICS                                              Human Mutation

**UMD-Predictor: A High-Throughput Sequencing Compliant System for Pathogenicity Prediction of any Human cDNA Substitution**

HGV$

OFFICIAL JOURNAL

HUMAN GENOME
VARIATION SOCIETY
www.hgvs.org

David Salgado,[1,2†] Jean-Pierre Desvignes,[1,2†] Ghadi Rai,[1,2] Arnaud Blanchard,[1,2] Morgane Miltgen,[1,2] Amélie Pinard,[1,2] Nicolas Lévy,[1,2,3] Gwenaëlle Collod-Béroud,[1,2] and Christophe Béroud[1,2,3*]
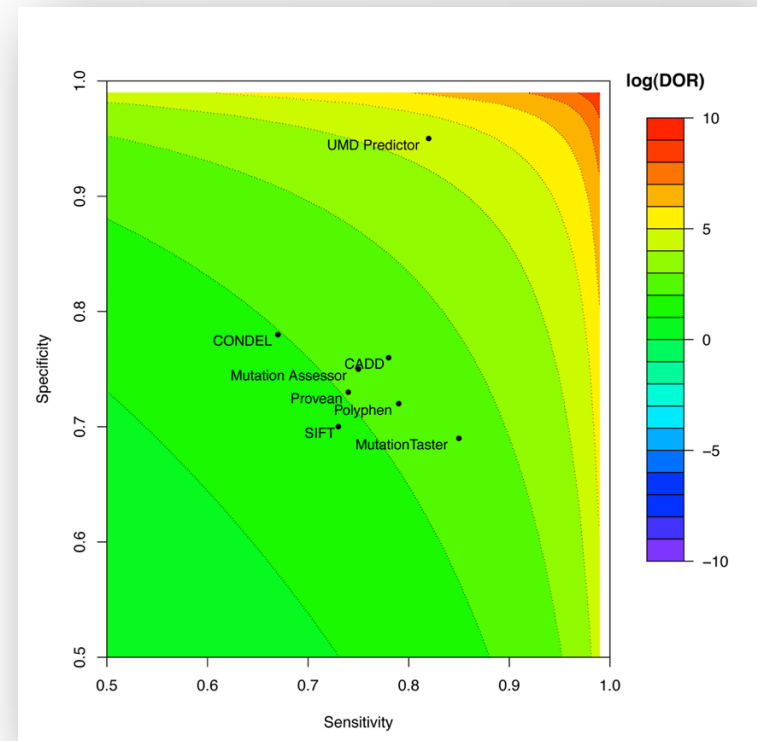
**Christophe Béroud– 2nd VEP Training**
**6th - 8th November 2017**

Aix Marseille université        Inserm
Institut national
de la santé et de la recherche médicale

22

GENETICS & BIOINFORMATICS TEAM

# UMD-Predictor
# Evaluation

- **4 datasets** (more than 140,000 mutations) - Varibench + dbSNP, Uniprot, ClinVar and PredictSNP

- **7 references pathogenicity predictors**

Varibench + dbSNP = 17,329 variations

| | SIFT | PPH2 | Provean | Mutation Assessor | CONDEL | Mutation Taster | CADD | UMD-Predictor |
|---|---|---|---|---|---|---|---|---|
| **TP** | 9596 | 10290 | 9638 | 9775 | 8797 | 11174 | 10182 | 10727 |
| **TN** | 2805 | 3045 | 3088 | 3162 | 3287 | 2937 | 3214 | 4024 |
| **FP** | 1229 | 1189 | 1147 | 1073 | 948 | 1298 | 1021 | 211 |
| **FN** | 3498 | 2803 | 3456 | 3319 | 4297 | 1920 | 2912 | 2367 |
| **Sensitivity** | 0.73 | 0.79 | 0.74 | 0.75 | 0.67 | 0.85 | 0.78 | 0.82 |
| **Specificity** | 0.70 | 0.72 | 0.73 | 0.75 | 0.78 | 0.69 | 0.76 | 0.95 |
| **DOR** | 6.3 | 9.7 | 7.7 | 9.0 | 7.2 | 12.6 | 11.2 | 86.6 |
| **log(DOR)** | 1.84 | 2.27 | 2.04 | 2.20 | 1.97 | 2.53 | 2.42 | 4.46 |



DOR : Measure the effectiveness of a diagnostic test
Trade-off between sensitivity and specificity

*Similar results with other datasets (Cf. Salgado, Desvignes, et al. Human Mutation 2016)*

# UMD-Predictor
# Real data evaluation

Comparison using 3 WES performed in a clinical diagnostic context

| | SIFT[A] | PPH2[A] | Provean[A] | Mutation Assessor[A,D] | CONDEL[A,B] | Mutation Taster | CADD[A] | UMD-Predictor |
|---|---|---|---|---|---|---|---|---|
| **NV1** | 1958 | 2881 | 1540 | 1339 | 1376 | 2677 | 3241 | 871 |
| **NV2** | 1341 | 2350 | 1332 | 1049 | 1111 | 2437 | 2555 | 540 |
| **NV3** | 1842 | 2781 | 1542 | 1350 | 1376 | 3401 | 3098 | 807 |

Shortest list of potential pathogenic mutations

Time required to process VCF files

| | SIFT[A] | PPH2[A] | Provean[A] | Mutation Assessor[A,D] | CONDEL[A,B] | Mutation Taster | CADD[A] | UMD-Predictor |
|---|---|---|---|---|---|---|---|---|
| **PT1 (s)** | 1200 | 420 | 3240 | 540 | 3000 | 2100 | 8700 | 93 |
| **PT2 (s)** | 240 | 420 | 8100 | 960 | 1500 | 2340 | 9360 | 206 |
| **PT3 (s)** | 540 | 420 | 4140 | 600 | 1500 | 2340 | 11160 | 240 |

Even faster by using webservices

Fastest system to process variations from VCF files

*Salgado D, et al. Human Mutation 2016*

# UMD-Predictor (http://umd-predictor.eu)
# Real data evaluation

Data from yesterday presentation: *BRCA1* c.5207T>C

## ENST00000357654

More at Ensembl    More at NCBI

| Transcript position | Genomic position | Mutation (c.) | Mutant (p.) | Pathogenicity | Conclusion |
|---|---|---|---|---|---|
| | | 5207 | | | Select Filter |
| 5207 | 41209139 | c.5207T>A | p.Val1736Asp | 100 | Pathogenic |
| 5207 | 41209139 | c.5207T>C | p.Val1736Ala | 84 | Pathogenic |
| 5207 | 41209139 | c.5207T>G | p.Val1736Gly | 99 | Pathogenic |

Go to page: 1   Show rows: 25   1–3 of 3

Export to XML    Export to CSV    Export to TSV    Export to JSON

▶ **Show transcript statistical mutation distribution**

*Salgado D, et al. Human Mutation 2016*

# Human Splicing Finder
http://umd.be/HSF3/

- Pathogenicity prediction system for any mutations on splicing signals

- Reference system ( 828 citations, Web of Science; 1,125, Google Scholar)

- A "One stop-Shop" system
  - Splicing signals
  - Branch points
  - Auxiliary signals (ESE, ESS, ISE, ISS)

- Combine various predictive systems, matrices and specific algorithms

- Expert system for data interpretation (establishment of rules to provide a conclusion) e.g. *MSH2* c.274_276del



- Compliant with NGS technologies: webservices (available in the HSF3-Pro version)

# Human Splicing Finder
## http://umd.be/HSF3/

- Use a *BRCA1/2* mutations dataset from ClinVar
  - 5' ss (3 last exonic nt + 6 first intronic nt)
  - 3' ss (12 last intron... ... ... ...)

  - *BRCA1*
    - 135 pathogen...
    - 16 non-patho...
  - *BRCA2*
    - 88 pathogenic mu...
    - 15 non-pathogenic...

# Human Splicing Finder
## http://umd.be/HSF3/

| *BRCA1 & BRCA2* dataset | |
|---|---|
| True Positives | 223 |
| True Negatives | 28 |
| False Positives | 3 |
| False Negatives | 0 |
| Positive Predictive Value | 0.986 |
| Negative Predictive Value | 1.000 |
| Sensitivity | 1.000 |
| Specificity | 0.903 |
| Accuracy | 0.988 |
| Matthews Correlation Coefficient | 0.944 |

GENETICS & BIOINFORMATICS TEAM

# Human Splicing Finder
http://umd.be/HSF3/

- Impact of mutations on branch points



Pre-mRNA    5'SS    BPS    3'SS

5' exon  GU ——— A — AG  3' exon

intron

Harris et Senapathy, 1988

Desmet *et al*, 2009

# Human Splicing Finder
http://umd.be/HSF3/

- Impact of mutations on branch points

- Few BPS are described in the literature
    - BRAF c.1141-51 C>G (Pupo et al. 2017), BPS identification
    - IKBKG/NEMO c.519-23A>T (Jorgensen et al. 2016), alteration of the BPS
    - C21orf2 c.643-23A>T (Wang et al. 2016), alteration of the BPS
    - ITGB4 c.1762-25T>A (Masunaga et al. 2015), alteration of the BPS
    - PC c.1369-29A>G (Ostergaardet al. 2012), alteration of the BPS
    - KCNH2 c.IVS9-28A>G (Crotti et al. 2009), alteration of the BPS
    - …

# Human Splicing Finder
## http://umd.be/HSF3/

- Few BPS are described in the literature
  - BRAF c.1141-51 C>G (Pupo et al. 2017), BPS identification

# Human Splicing Finder
## http://umd.be/HSF3/

- Few BPS are described in the literature
  - ITGB4 c.1762-25T>A (Masunaga et al. 2015), alteration of the BPS

# Human Splicing Finder
http://umd.be/HSF3/

- Examples from Andreas Laner yesterday (impact on ESE/ESS)

# Human Splicing Finder
http://umd.be/HSF3/

- Examples from Andreas Laner yesterday (impact on ESE/ESS)

# Human Splicing Finder
http://umd.be/HSF3/

HSF system is highly accurate to predict the impact of mutations on 5' and 3' splice sites

Now contains specific matrices for non-canonical splice sites

**HSF** efficiently predicts the **BPS** and the impact of mutations on these sites

**HSF** predicts the impact of mutations on ESE/ESS

**Christophe Béroud– 2nd VEP Training**
**6th - 8th November 2017**

Aix Marseille université    Inserm
Institut national
de la santé et de la recherche médicale
GENETICS & BIOINFORMATICS TEAM

35

# Splicing predictors comparison

# Splicing predictors comparison

| Software | Global | 3'ss | 5'ss |
|---|---|---|---|
| GeneSplicer | 63.89 % | 82.61 % | 55.1 % |
| GenScan | 45.83 % | 82.61 % | 30.61 % |
| **Human Splicing Finder** | **100%** | **100%** | **100%** |
| MaxEntScan | **100%** | **100%** | **100%** |
| NNSplice | 97.22 % | **100%** | 95.92 % |
| SplicePort | 87.5 % | 82.61 % | 89.8 % |
| SplicePredictor | 87.5 % | 95.65 % | 83.67 % |
| SpliceView | **100%** | **100%** | **100%** |
| SROOGLE | **100%** | **100%** | **100%** |
| *Average* | *86.88 %* | *93.72 %* | *83.9 %* |

# Splicing predictors comparison

| Software | Intronic (<100 bp) | Intronic (> 100 bp) |
|---|---|---|
| GeneSplicer | 41.18% | 46.15% |
| GenScan | 11.76% | 0.00% |
| **Human Splicing Finder** | **70.59%** | **92.31%** |
| MaxEntScan | 23.53% | **92.31%** |
| NNSplice | 5.88% | 69.23% |
| SplicePort | 35.29% | 61.54% |
| SplicePredictor | 23.53% | 69.23% |
| SpliceView | 17.65% | 84.62% |
| Sroogle | 17.65% | 30.77% |
| *Average* | *27.45%* | *60.68%* |

# Splicing predictors comparison

| Software | Polymorphisms |
|---|---|
| GeneSplicer | 50% |
| GenScan | **0%** |
| **Human Splicing Finder** | **0%** |
| MaxEntScan | 50% |
| NNSplice | 25% |
| SplicePort | 100% |
| SplicePredictor | **0%** |
| SpliceView | **0%** |
| Sroogle | **0%** |
| *Average* | *10%* |

# Splicing predictors comparison

GENETICS & BIOINFORMATICS TEAM

# **VAR**iant **A**nnotation and **F**iltration **T**ool
## http://varaft.eu

- An user-friendly variant annotation and filtration tool

- Standalone multiplatform application

- Evaluation of the data coverage for WGS/WES/panel

- One click annotation (based on ANNOVAR) and other sources

- Provides UMD-Predictor and HSF annotations

- Interactive filtration at many levels

  - Combine multiple samples

  - Automatic selection of variants (mode of inheritance)

  - Genetic population studies

  - Cancers

- Standardize your variant analysis processes

  - save/reuse/share applied filters



**VarAFT**    About   Latest Changes   Download   Documentation   Tutorials   FAQ   Our tools   Contact   Citation

**VarAFT**

Variant Annotation and Filter Tool

**ABOUT VarAFT**

The identification of disease-causing mutations in human genetics remains challenging despite the **NGS** revolution as up to 70% of cases are still unsolved. To tackle this challenge, we developed the

VarAFT is a standalone and multiplatform (Windows, Mac, Unix) system which is easy to install and easy to use. It does not require skills in informatics or a powerful computer.

# **Var**iant **A**nnotation and **F**iltration **T**ool

# Conclusions

- No ideal annotation/filtration systems
  - Many of them are built to be used by bioinformaticians (command line systems)
- To maximize chances of identify causing mutations
  - be aware of challenges posed by each steps of the data analysis pipeline
  - use family members (when possible)
  - collect exact and complete phenotypic information
- Many challenges remain to be solved for both annotation and filtration systems - benchmarking initiatives – CAGI challenges
- Most of the current available systems are created for WES and need to be adapted to WGS
- Many more issues arise with WGS
  - Annotation for non-coding regions but also for non-protein coding genes
  - Need to develop and improve pathogenicity prediction system for non-exonic mutations
- Facilitated with data-sharing initiatives (shared variants and conclusion)

**Christophe Béroud– 2nd VEP Training**
**6th - 8th November 2017**

Aix Marseille université    Inserm
Institut national de la santé et de la recherche médicale
GENETICS & BIOINFORMATICS TEAM

43

# A simple use case

- A trio with an affected daughter and two healthy parents

- Mabry syndrome: intellectual disability, distinctive facial features, hyperphosphatasia, and other signs and symptoms.

| Father | Daughter | Mother |
|--------|----------|--------|
| 20,486 variants | 20,645 variants | 20,486 variants |
| 8,802 genes | 8,767 genes | 8,774 genes |

# A simple use case

- A trio with an affected daughter and two healthy parents

- Mabry syndrome: intellectual disability, distinctive facial features, hyperphosphatasia, and other signs and symptoms.

- 4 steps filtration workflow:

  - Mode of inheritance (Autosomal recessive)

    ➤ Homozygous or compound heterozygous

| Father | Daughter | Mother |
|---|---|---|
| 20,486 variants | 20,645 variants | 20,486 variants |
| 8,802 genes | 8,767 genes | 8,774 genes |

❶ Mode of inheritance

| |
|---|
| 1,936 variants |
| 520 genes |

**Christophe Béroud– 2nd VEP Training**
**6th - 8th November 2017**

Aix Marseille université    Inserm
Institut national de la santé et de la recherche médicale
GENETICS & BIOINFORMATICS TEAM

45

# A simple use case

- A trio with an affected daughter and two healthy parents

- Mabry syndrome: intellectual disability, distinctive facial features, hyperphosphatasia, and other signs and symptoms.
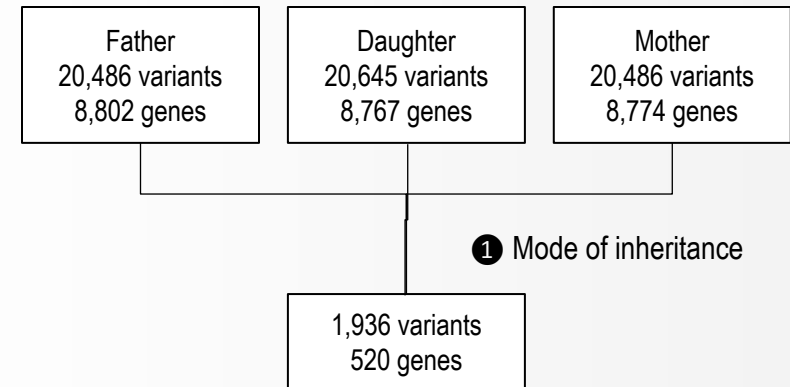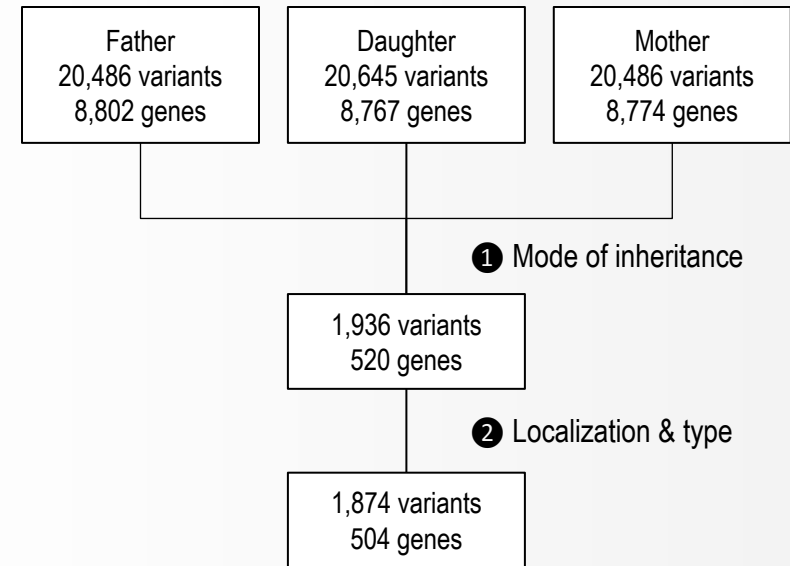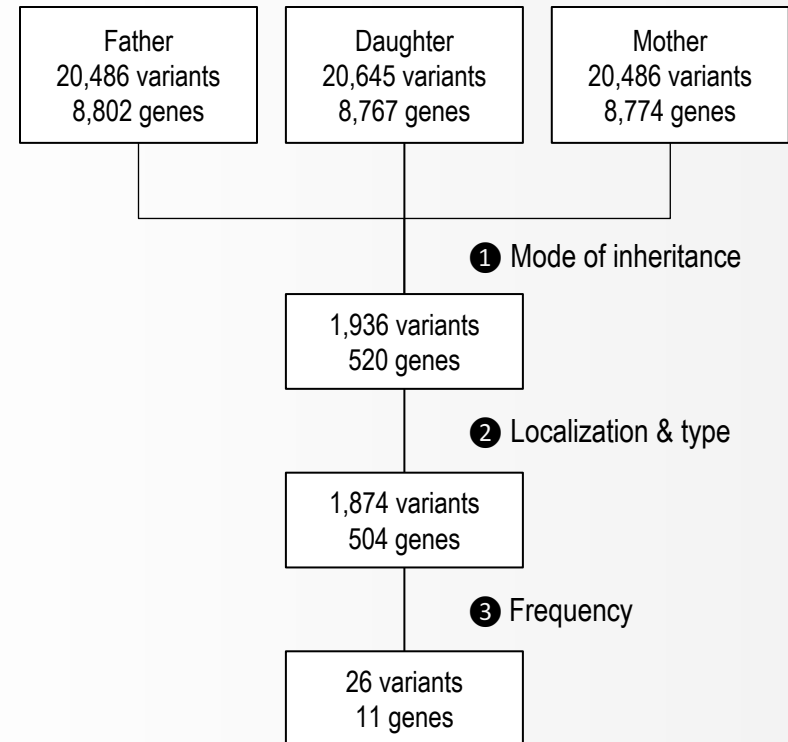
- 4 steps filtration workflow:

  - Mode of inheritance (Autosomal recessive)

    ➤ Homozygous or compound heterozygous

  - Mutation localization and type

| Father 20,486 variants 8,802 genes | Daughter 20,645 variants 8,767 genes | Mother 20,486 variants 8,774 genes |
|---|---|---|

❶ Mode of inheritance

1,936 variants
520 genes

❷ Localization & type

1,874 variants
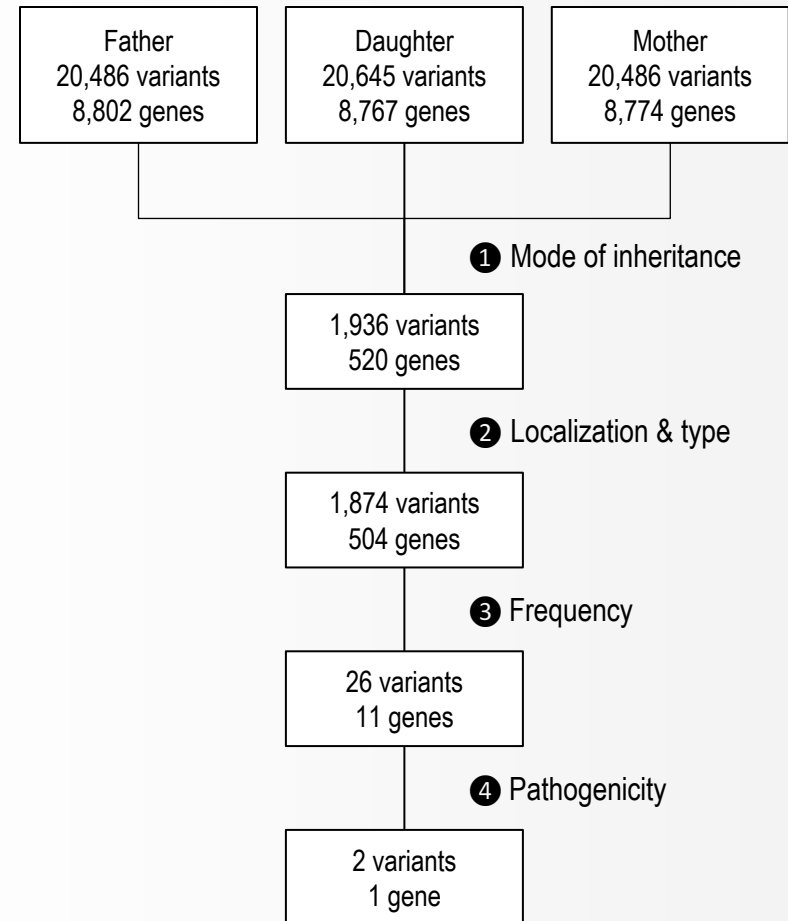504 genes

# A simple use case

- A trio with an affected daughter and two healthy parents

- Mabry syndrome: intellectual disability, distinctive facial features, hyperphosphatasia, and other signs and symptoms.

- 4 steps filtration workflow:

  - Mode of inheritance (Autosomal recessive)

    ➤ Homozygous or compound heterozygous

  - Mutation localization and type

  - Allele frequency

| Father 20,486 variants 8,802 genes | Daughter 20,645 variants 8,767 genes | Mother 20,486 variants 8,774 genes |
|---|---|---|

❶ Mode of inheritance

1,936 variants
520 genes

❷ Localization & type

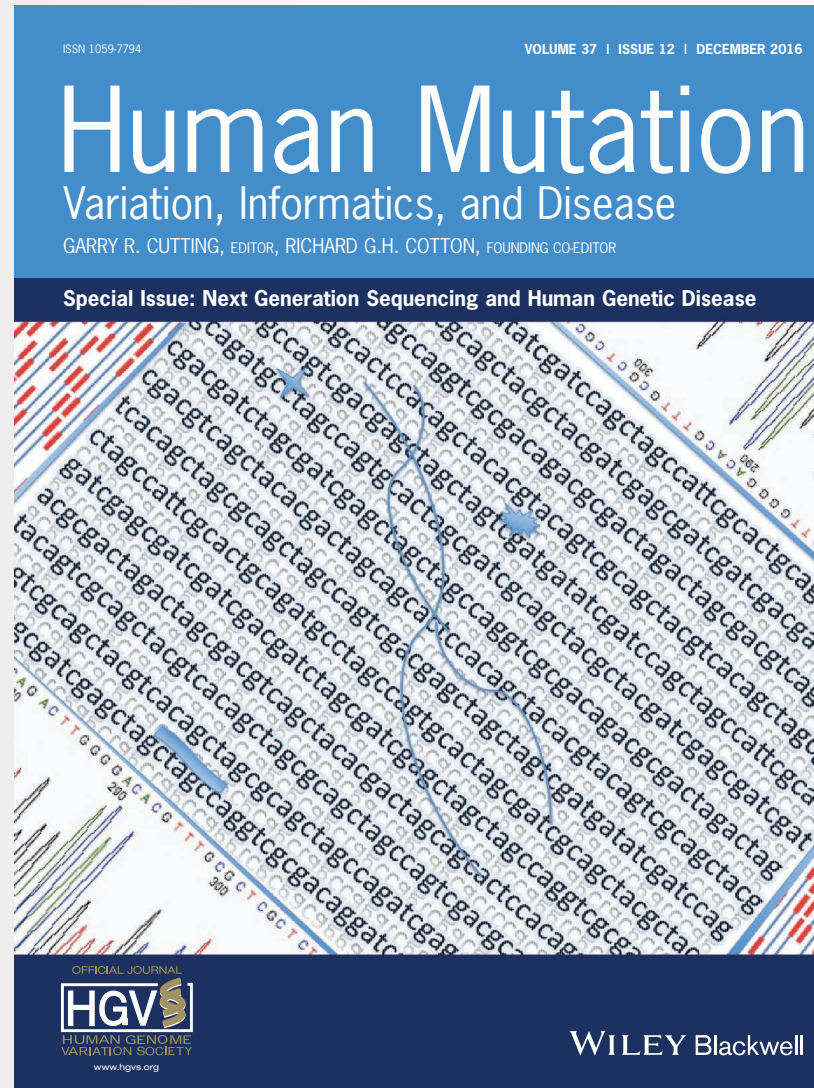1,874 variants
504 genes

❸ Frequency

26 variants
11 genes

# A simple use case

- A trio with an affected daughter and two healthy parents

- Mabry syndrome: intellectual disability, distinctive facial features, hyperphosphatasia, and other signs and symptoms.

- 4 steps filtration workflow:

  - Mode of inheritance (Autosomal recessive)

    - ➤ Homozygous or compound heterozygous

  - Mutation localization and type

  - Allele frequency

  - Pathogenicity predictions



| Father 20,486 variants 8,802 genes | Daughter 20,645 variants 8,767 genes | Mother 20,486 variants 8,774 genes |

❶ Mode of inheritance

1,936 variants
520 genes

❷ Localization & type

1,874 variants
504 genes

❸ Frequency

26 variants
11 genes

❹ Pathogenicity

2 variants
1 gene

# More details

# More details

# Acknowledgements


David Salgado


Jean-Pierre Desvignes

<u>**all members of the Genetics & Bioinformatics Team**</u>









**GenOmnis company (https://genomnis.com)**

UMD-Predictor and HSF professional version for commercial users