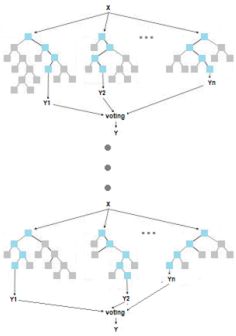**Dušan Popović**

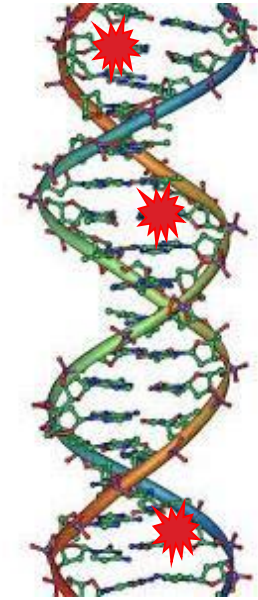# A computational framework for prioritization of disease-causing mutations
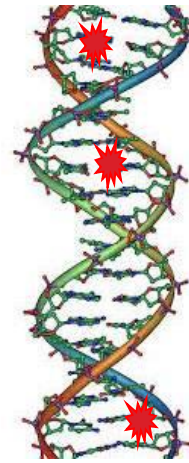
**PREVALENCE OF GENETIC DISORDERS**

**PREVELANCE OF GENETIC DISORDERS**

**CAUSED BY MUTATIONS**

**PREVELANCE OF GENETIC DISORDERS**

**CAUSED BY MUTATIONS**

**NEXT-GENERATION SEQUENCING**

**NOVEL GENETIC DISORDER**

**SAMPLE**

**SEQUENCING**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**NOVEL GENETIC DISORDER**

**SAMPLE**

**SEQUENCING**

**MUTATIONS**

**CONFIRMATORY EXPERIMENTS**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**SAMPLE**

**SEQUENCING**

**MUTATIONS**

**MUTATION PRIORITIZATION**

**NOVEL GENETIC DISORDER**

**PRIORITIZED VARIANTS**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**CONFIRMATORY EXPERIMENTS**

**NOVEL GENETIC DISORDER**

**SAMPLE**

**SEQUENCING**

**MUTATIONS**

**MUTATION PRIORITIZATION**

**PRIORITIZED VARIANTS**

**CONFIRMATORY EXPERIMENTS**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**MUTATIONS**

PolyPhen-2

Mutation Taster

SIFT

PRIORITIZED VARIANTS

**MUTATION PRIORITIZATION METHODS**

**NOVEL GENETIC DISORDER**

**SAMPLE**

**SEQUENCING**

**MUTATIONS**

eXtasy
Variant Prioritization

**MUTATION PRIORITIZATION**

**PRIORITIZED VARIANTS**

**CONFIRMATORY EXPERIMENTS**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**SAMPLE**

**SEQUENCING**

**MUTATIONS**

**DISEASE PHENOTYPES**

**MUTATION PRIORITIZATION**

**NOVEL GENETIC DISORDER**

**PRIORITIZED VARIANTS**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**CONFIRMATORY EXPERIMENTS**

**EXTASY**          **OTHER METHODS**

# HOW DOES EXTASY WORK?

**SAMPLE**

**SEQUENCING**

**MUTATIONS**

**DISEASE PHENOTYPES**

eXtasy
Variant Prioritization

**MUTATION PRIORITIZATION**

**NOVEL GENETIC DISORDER**

**PRIORITIZED VARIANTS**

**DISCOVERY OF MUTATION CAUSING THE DISEASE**

**CONFIRMATORY EXPERIMENTS**

MUTATIONS

DISEASE PHENOTYPES

PolyPhen-2

Mutation Taster

SIFT

ANNOTATION

Other scores

Endeavour

Phenotype scores

RANDOM FOREST CLASSIFIER

PRIORITIZED VARIANTS

**MUTATION PRIORITIZATION BY EXTASY**

**MUTATIONS**

**DISEASE PHENOTYPES**

PolyPhen-2

Mutation Taster

SIFT

**ANNOTATION**

Endeavour

**Other scores**

**Phenotype scores**

**RANDOM FOREST CLASSIFIER**

**MUTATION PRIORITIZATION BY EXTASY**

**PRIORITIZED VARIANTS**

**MUTATIONS**

**DISEASE PHENOTYPES**

PolyPhen-2

Mutation Taster

SIFT

**Other scores**

**ANNOTATION**

Endeavour

**Phenotype scores**

**RANDOM FOREST CLASSIFIER**

**PRIORITIZED VARIANTS**

**MUTATION PRIORITIZATION BY EXTASY**

# Endeavour

?

Endeavour

EW01X

EW02X

EW03X

EW04X

EW05X

EW06X

?

CHUBKK1

M01NK-Y

M02NK-Y

M03NK-Y

# Endeavour

Endeavour

?

**MUTATIONS**

**DISEASE PHENOTYPES**

PolyPhen-2

Mutation Taster

SIFT

**ANNOTATION**

**Other scores**

Endeavour

**Phenotype scores**

**RANDOM FOREST CLASSIFIER**

**PRIORITIZED VARIANTS**

**MUTATION PRIORITIZATION BY EXTASY**

# RANDOM FOREST CLASSIFIER

**MUTATIONS**

**DISEASE PHENOTYPES**

PolyPhen-2

Mutation Taster

SIFT

ANNOTATION

Endeavour

Other scores

Phenotype scores

**RANDOM FOREST CLASSIFIER**

PRIORITIZED VARIANTS

**MUTATION PRIORITIZATION BY EXTASY**

PHENOTYPE SCORES AGGREGATION BY EXTASY

PHENOTYPE SCORES AGGREGATION BY EXTASY

**MUTATION SCORE 1**

$t(X_i)$

$A(X)$

$t(X_i)$

**MUTATION SCORE n**

**PRIORITIZED VARIANT**

1. MAXIMUM :

$t(X_i) = X_i$

$A(X) = max(X)$

2. PARAMETRIC MODELING :

$t(X_i) = P(X_i)$

$A(X) = F(X)$

P – p-value of Gamma distribution

F – Fisher's omnibus statistics

PolyPhen-2

Mutation Taster

SIFT

**Other scores**

**ANNOTATION**

Endeavour   **Phenotype scores**

**RANDOM FOREST CLASSIFIER**

**MUTATIONS**

**DISEASE PHENOTYPES**

**MUTATION PRIORITIZATION BY EXTASY**

**PRIORITIZED VARIANTS**

*eXtasy: Variant Prioritization by Genomic Data Fusion*

**What is eXtasy?**
eXtasy is a pipeline for ranking nonsynonymous single nucleotide variants given a specific phenotype. It takes into account the putative deleteriousness of the variant, haploinsufficiency predictions of the underlying gene and the similarity of the given gene to known genes in the given phenotype.

**Who develops eXtasy?**
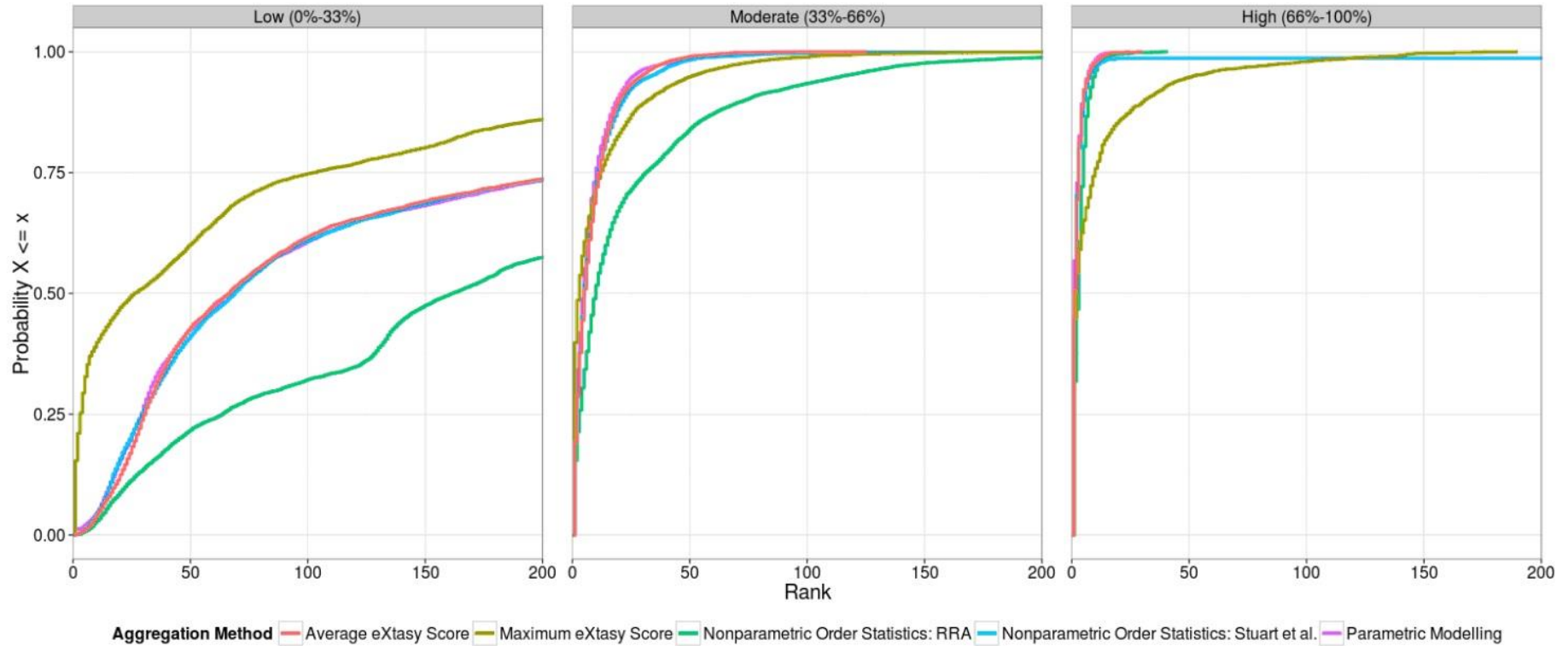eXtasy was developed in the Bioinformatics group at the Department of Electrical Engineering of the University of Leuven (part of the iMinds Future Health Department). It was implemented by Alejandro Sifrim and Dusan Popovic under the supervision of Prof. Jan Aerts, Prof. Bart de Moor and Prof. Yves Moreau.

**What is the input of eXtasy?**
One can run eXtasy on any VCF file mapped to hg19/Gchr37. As a second input the user can choose any of the precomputed gene prioritization files for a given HPO term (downloadable here). In the near future we will provide the user the possibility of creating custom gene prioritizations given a set of phenotype-associated genes.

**Run eXtasy online:**

HPO term(s): _____ (Comma-separated)

VCF file: _____ Browse...

Email: _____

Output file name: _____

Submit

Example Data: miller.vcf,schinzel_giedion.vcf

We provide two example vcf files which were generated by adding published disease causing variants for Miller syndrome (causative gene: DHODH, Ng et al., 2010, Nature Genetics) or Schinzel-Giedion syndrome (causative gene: SETBP1, Hoischen et al., 2010, Nature Genetics) to a publicly available VCF file of the exome of a healthy individual (obtained from here). These files can be prioritized against any of

**Speed considerations:**

Average job completion time over all jobs submitted to the eXtasy webtool: 00:25:48

Currently running eXtasy takes about 5-10 minutes for a single exome (~ 40 000 variants) on a standard single core (currently we don't perform any parallelization within one job). It uses only small amounts of RAM memory, allowing it to be run on almost any computer. Most of the time is spent annotating the variants. Significant increases in speed can be achieved by performing this step only once (using the -k and -r options) when prioritizing against mulitple phenotypes.

**News:**

- *July 5, 2015:* Major changes to the webtool and

*eXtasy: Variant Prioritization by Genomic Data Fusion*

**What is eXtasy?**
eXtasy is a pipeline for ranking nonsynonymous single nucleotide variants given a specific phenotype. It takes into account the putative deleteriousness of the variant, haploinsufficiency predictions of the underlying gene and the similarity of the given gene to known genes in the given phenotype.

**Who develops eXtasy?**
eXtasy was developed in the Bioinformatics group at the Department of Electrical Engineering of the University of Leuven (part of the iMinds Future Health Department). It was implemented by Alejandro Sifrim and Dusan Popovic under the supervision of Prof. Jan Aerts, Prof. Bart de Moor and Prof. Yves Moreau.

**What is the input of eXtasy?**
One can run eXtasy on any VCF file mapped to hg19/Gchr37. As a second input the user can choose any of the precomputed gene prioritization files for a given HPO term (downloadable here). In the near future we will provide the user the possibility of creating custom gene prioritizations given a set of phenotype-associated genes.

**Run eXtasy online:**

HPO term(s): | hair | (Comma-separated)

Abnormal hair whorl
HP:0010721

Abnormality of hair texture
HP:0010719

Abnormality of secondary sexual hair
HP:0009888

Abnormality of the frontal hairline
HP:0000599

Abnormality of the hair
HP:0001595

Abnormality of the hairline
HP:0009553

VCF file:

Email:

Output file name:

Submit

Example Data: miller.vcf,schinze...

We provide two example vcf files
Miller syndrome (causative gene
syndrome (causative gene: SETB
file of the exome of a healthy indi

se causing variants for
inzel-Giedion
ublicly available VCF
oritized against any of

**Speed considerations:**

**Average job completion time over all jobs submitted to the eXtasy webtool: 00:25:48**

Currently running eXtasy takes about 5-10 minutes for a single exome (~ 40 000 variants) on a standard single core (currently we don't perform any parallelization within one job). It uses only small amounts of RAM memory, allowing it to be run on almost any computer. Most of the time is spent annotating the variants. Significant increases in speed can be achieved by performing this step only once (using the -k and -r options) when prioritizing against mulitple phenotypes.

**News:**

• July 5, 2013: Major changes to the webtool and

eXtasy: Variant Prioritization by Genomic Data Fusion

**What is eXtasy?**
eXtasy is a pipeline for ranking nonsynonymous single nucleotide variants given a specific phenotype. It takes into account the putative deleteriousness of the variant, haploinsufficiency predictions of the underlying gene and the similarity of the given gene to known genes in the given phenotype.

**Who develops eXtasy?**
eXtasy was developed in the Bioinformatics group at the Department of Electrical Engineering of the University of Leuven (part of the iMinds Future Health Department). It was implemented by Alejandro Sifrim and Dusan Popovic under the supervision of Prof. Jan Aerts, Prof. Bart de Moor and Prof. Yves Moreau.

**What is the input of eXtasy?**
One can run eXtasy on any VCF file mapped to hg19/Gchr37. As a second input the user can choose any of the precomputed gene prioritization files for a given HPO term (downloadable here). In the near future we will provide the user the possibility of creating custom gene prioritizations given a set of phenotype-associated genes.

**Run eXtasy online:**

HPO term(s): HP:0001595, nose (Comma-separated)

Abnormality of the nose
HP:0000366

Aplasia/Hypoplasia involving the nose
HP:0009924

Bulbous nose
HP:0000443

Bulbous nose
HP:0000414

Flat nose
HP:0000457

Long nose
HP:0003189

Midline defect of the nose

VCF file:

Email:

Output file name:

Submit

Example Data: miller.vcf, schinze[...]

We provide two example vcf files [...] e causing variants for Miller syndrome (causative gene: [...] inzel-Giedion syndrome (causative gene: SETB[...] ublicly available VCF file of the exome of a healthy indi[...] ritized against any of

**Speed considerations:**

Average job completion time over all jobs submitted to the eXtasy webtool: 00:25:48

Currently running eXtasy takes about 5-10 minutes for a single exome (~ 40 000 variants) on a standard single core (currently we don't perform any parallelization within one job). It uses only small amounts of RAM memory, allowing it to be run on almost any computer. Most of the time is spent annotating the variants. Significant increases in speed can be achieved by performing this step only once (using the -k and -r options) when prioritizing against mulitple phenotypes.

**News:**

• July 5, 2013: **Major changes to the webtool and**

eXtasy: Variant Prioritization by Genomic Data Fusion

Your eXtasy job has been succesfully submitted!

If you provided an e-mail you will receive a notification when the job is completed. Depending on the number of variants to be prioritized and the number of scheduled jobs this can take anywhere from a couple of minutes to several hours. If you experience any difficulties please mail alejandro.sifrim@esat.kuleuven.be.
Return to the eXtasy Homepage.

Also You are able to see you results with the following Link.

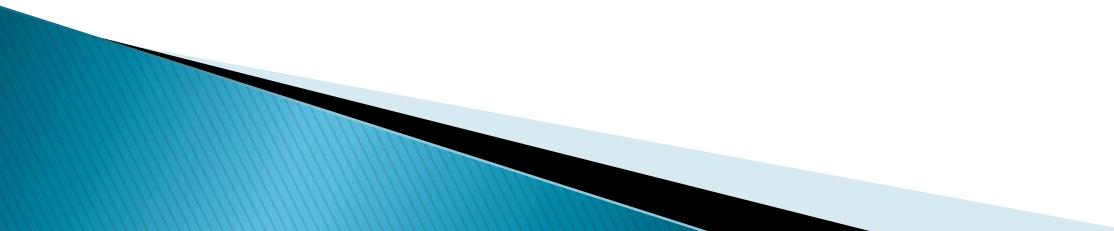| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | chromosome | refbase | altbase | position | genename | carol_scoi | ext | HP_0001595_fgs.extasy | HP_0000457_fgs.extasy | extasy_combined_max | extasy_combined_order_statistics |
| 2 | 15 | C | T | 91326099 | BLM | 0.99925~ | 007079724 | 0.736 | 0.496 | 0.736 | 3.37E-06 |
| 3 | 5 | T | C | 174156168 | MSX2 | 0.84 | 40715008 | 0.568 | 0.82 | 0.82 | 3.89E-06 |
| 4 | 2 | G | T | 121746956 | GLI2 | 0.479753 | 07079724 | 0.45 | 0.862 | 0.862 | 5.51E-06 |
| 5 | 12 | G | A | 121416797 | HNF1A | 0.99940£ | 270568577 | 0.586 | 0.76 | 0.76 | 6.22E-06 |
| 6 | 1 | G | A | 103379918 | COL11A1 | 0.8 | 07079724 | 0.716 | 0.384 | 0.716 | 1.40E-05 |
| 7 | 2 | C | T | 179643775 | TTN | 0.96 | 2092696 | 0.574 | 0.462 | 0.574 | 1.46E-05 |
| 8 | 12 | C | T | 48367976 | COL2A1 | | 007079724 | 0.636 | 0.442 | 0.636 | 1.46E-05 |
| 9 | 11 | G | A | 47470345 | RAPSN | 0.999918 | 07079724 | 0.632 | 0.434 | 0.632 | 1.94E-05 |
| 10 | 2 | G | A | 121747406 | GLI2 | 7.22E-08 | 007079724 | 0.428 | 0.784 | 0.784 | 2.00E-05 |
| 11 | 5 | G | A | 127873094 | FBN2 | 0.99 | 0.007079724 | 0.648 | 0.418 | 0.648 | 2.07E-05 |
| 12 | 16 | G | A | 14029033 | ERCC4 | 0.99997~ | 0.14827844 | 0.552 | 0.51 | 0.552 | 2.11E-05 |
| 13 | 5 | A | C | 42719239 | GHR | 0.994548 | 07079724 | 0.722 | 0.26 | 0.722 | 3.01E-05 |
| 14 | 14 | A | G | 75472653 | EIF2B2 | 99997 | 42900 11 | 0.442 | 0.686 | 0.686 | 3.42E-05 |
| 15 | 16 | T | C | 16295863 | ABCC6 | 0.999~ | 0.162481861 | 0.462 | 0.58 | 0.58 | 3.84E-05 |
| 16 | 12 | G | A | 121435427 | HNF1A | 0.968235 | 270568577 | 0.436 | 0.592 | 0.592 | 4.31E-05 |
| 17 | 5 | G | A | 112178795 | APC | 0.999~ | 2483136 | 0.382 | 0.72 | 0.72 | 4.70E-05 |
| 18 | 2 | G | A | 179650408 | TTN | 0. | 0.12092696 | 0.49 | 0.446 | 0.49 | 4.77E-05 |
| 19 | 11 | G | A | 86663296 | FZD4 | 0.91804, | 0.43283057 | 0.446 | 0.532 | 0.532 | 5.18E-05 |
| 20 | 12 | A | T | 56494998 | ERBB3 | 0.998191 | 07079724 | 0.764 | 0.09 | 0.764 | 5.98E-05 |
| 21 | 21 | G | A | 47545768 | COL6A2 | 0.999609 | 155138296 | 0.526 | 0.374 | 0.526 | 5.99E-05 |
| 22 | 16 | A | G | 14042077 | ERCC4 | 0.9818 | 0.14827844 | 0.594 | 0.338 | 0.594 | 6.32E-05 |
| 23 | 20 | C | T | 44579206 | ZNF335 | 0.999994 | 007079724 | 0.476 | 0.43 | 0.476 | 6.94E-05 |
| 24 | 13 | G | C | 103515085 | ERCC5 | 0.767353 | 007079724 | 0.706 | 0.206 | 0.706 | 7.94E-05 |
| 25 | 22 | A | G | 41548008 | EP300 | 0. ~ | 0.007079724 | 0.472 | 0.39 | 0.472 | 8.35E-05 |
| 26 | 12 | A | C | 121416650 | HNF1A | 0.9469. | 0.270568577 | 0.408 | 0.516 | 0.516 | 8.56E-05 |
| 27 | 22 | C | T | 18905964 | PRODH | 0.995252 | 4074579 | 0.444 | 0.44 | 0.444 | 8.63E-05 |
| 28 | 8 | C | G | 90990479 | NBN | 0.159482 | 07079724 | 0.502 | 0.352 | 0.502 | 8.83E-05 |
| 29 | 8 | T | C | 31024654 | WRN | 0.999845 | 44984829 | 0.56 | 0.31 | 0.56 | 9.57E-05 |
| 30 | X | G | C | 135956462 | RBMX | 0.999382 | 0.16195418 | 0.464 | 0.378 | 0.464 | 9.84E-05 |
| 31 | X | A | G | 135956408 | RBMX | 0.99. | 0.16195418 | 0.51 | 0.342 | 0.51 | 0.000103246 |
| 32 | X | C | G | 135956506 | RBMX | 0.98645 | 16195418 | 0.512 | 0.298 | 0.512 | 0.000131245 |
| 33 | 8 | C | T | 41566438 | ANK1 | 0.95651 | 007079724 | 0.616 | 0.232 | 0.616 | 0.000133319 |

# WHY DOES EXTASY WORK SO WELL?

# 1. PROBLEM-TAILORED CHOICE OF THE TRAINING SET

# 1. PROBLEM-TAILORED CHOICE OF THE TRAINING SET

common polymorphisms

disease-causing variants

rare benign variants

# 1. PROBLEM-TAILORED CHOICE OF THE TRAINING SET

common polymorphisms
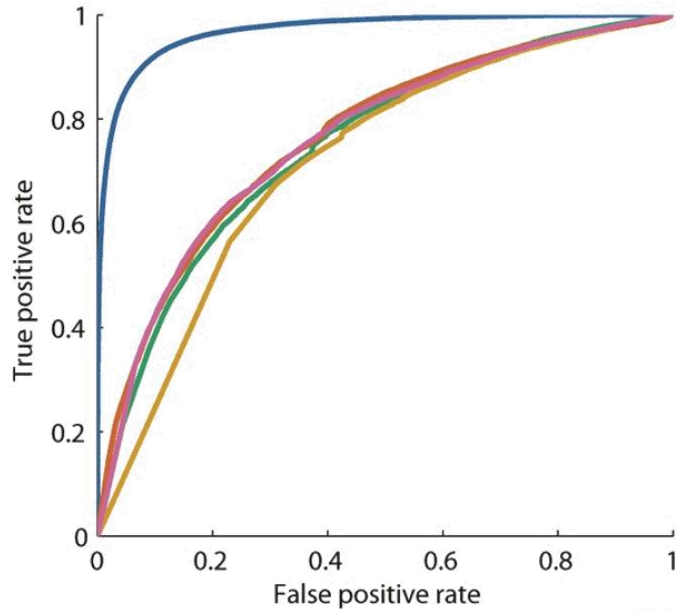
disease-causing variants
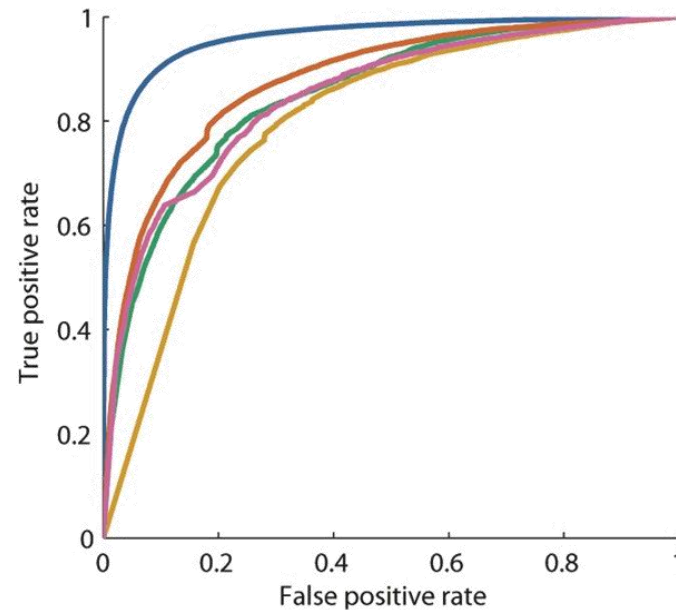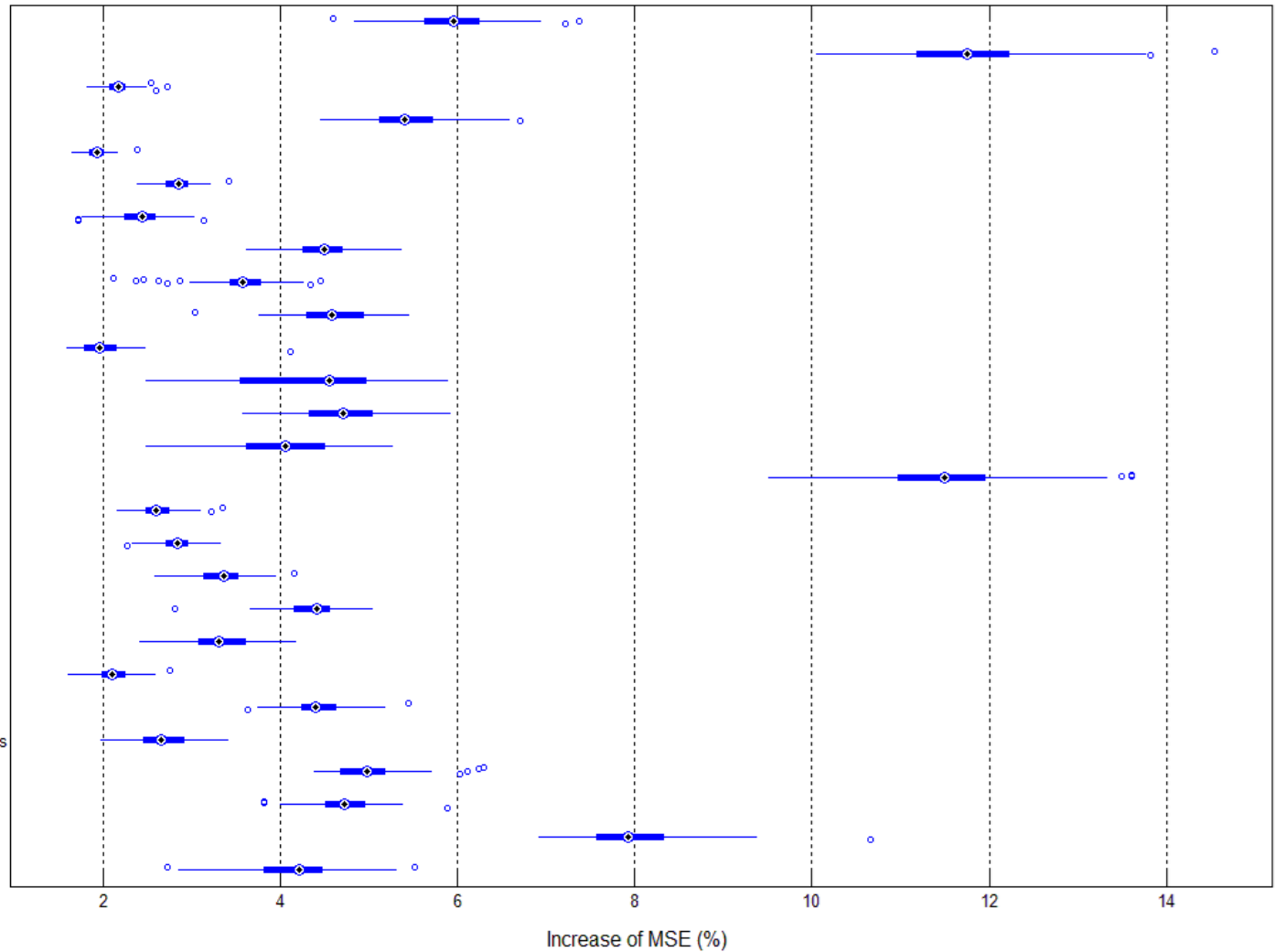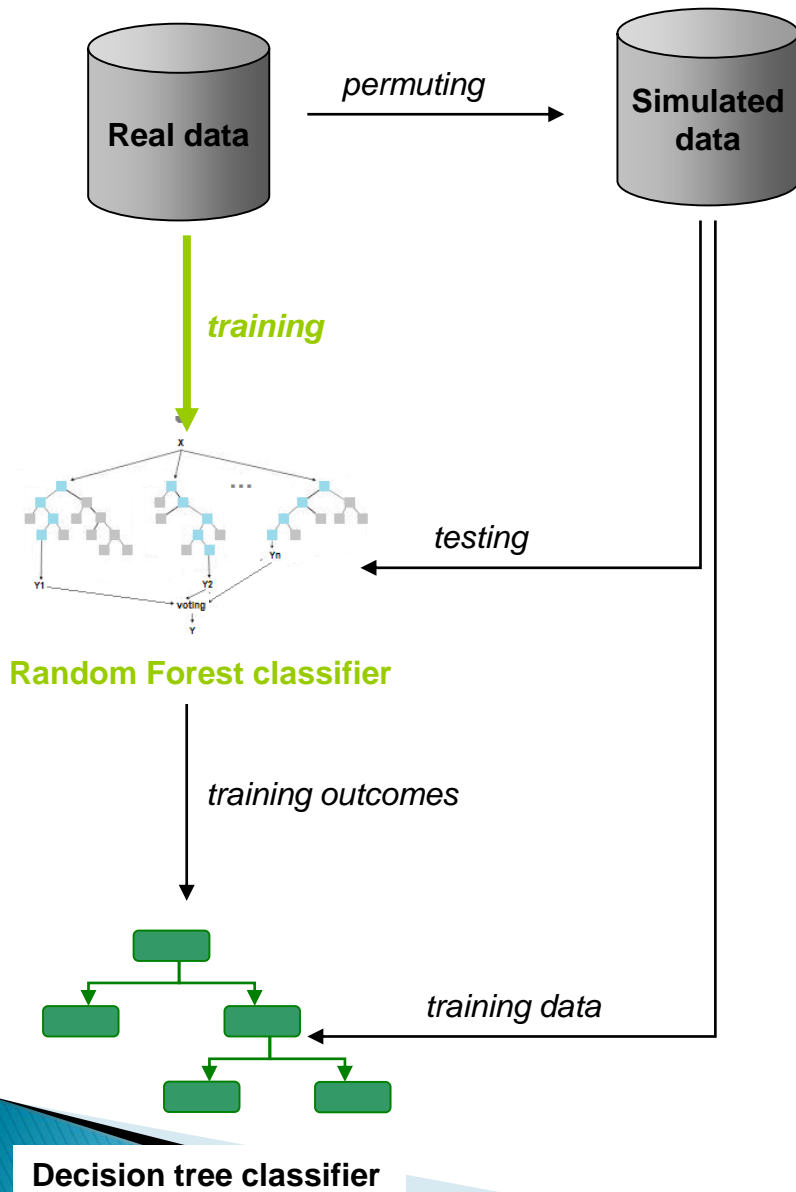
# 1. PROBLEM-TAILORED CHOICE OF THE TRAINING SET

disease-causing variants

rare benign variants

# 1. PROBLEM-TAILORED CHOICE OF THE TRAINING SET



rare benign vs. disease-causing

common polymorphisms vs. disease causing

Legend:
- eXtasy
- Polyphen score
- SIFT score
- Mutation Taster
- Carol score

1. **PROBLEM-TAILORED CHOICE OF THE TRAINING SET**
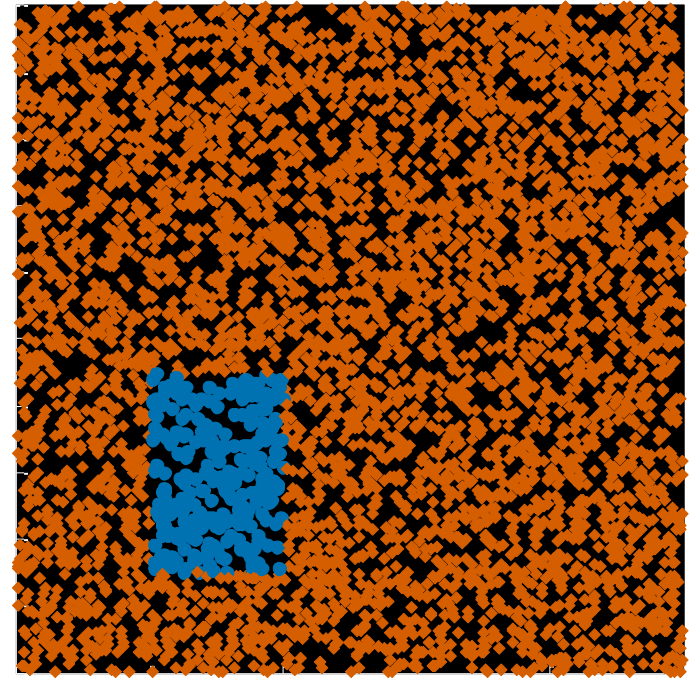
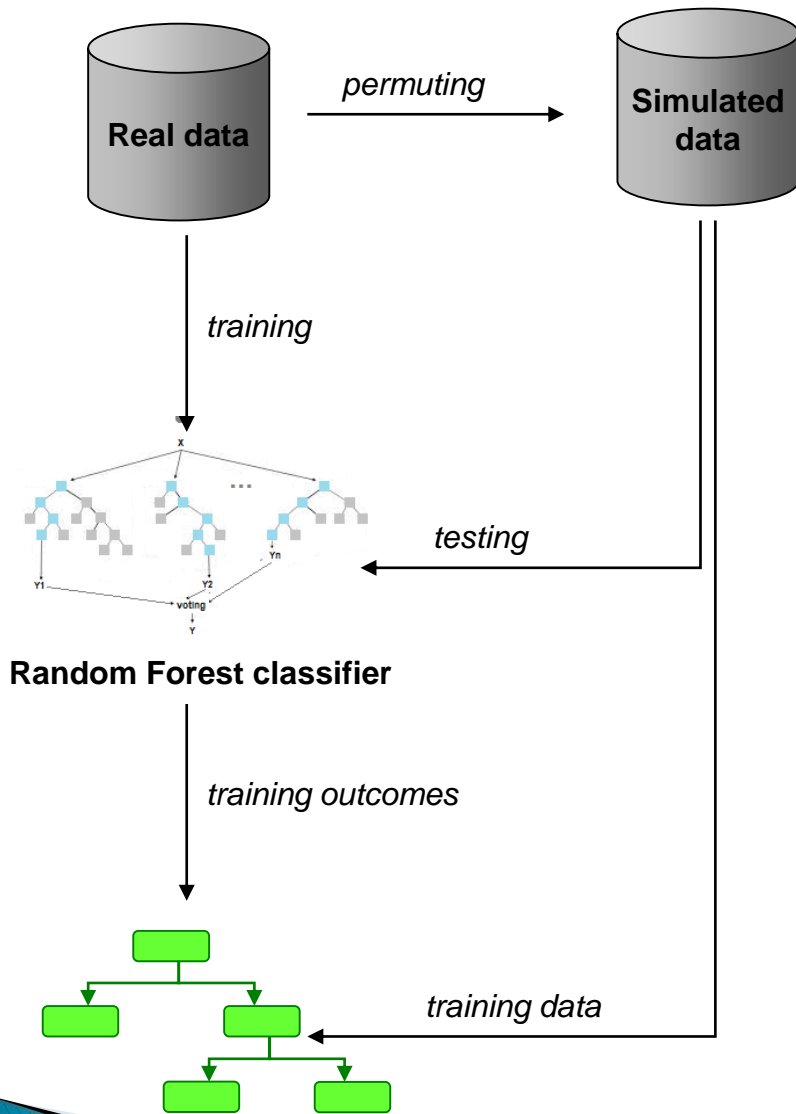2. **HETEROGENOUS DATA FUSION**

1. **PROBLEM-TAILORED CHOICE OF THE TRAINING SET**

2. **HETEROGENOUS DATA FUSION**

3. **PHENOTYPIC INFORMATION**

Real data

*permuting*

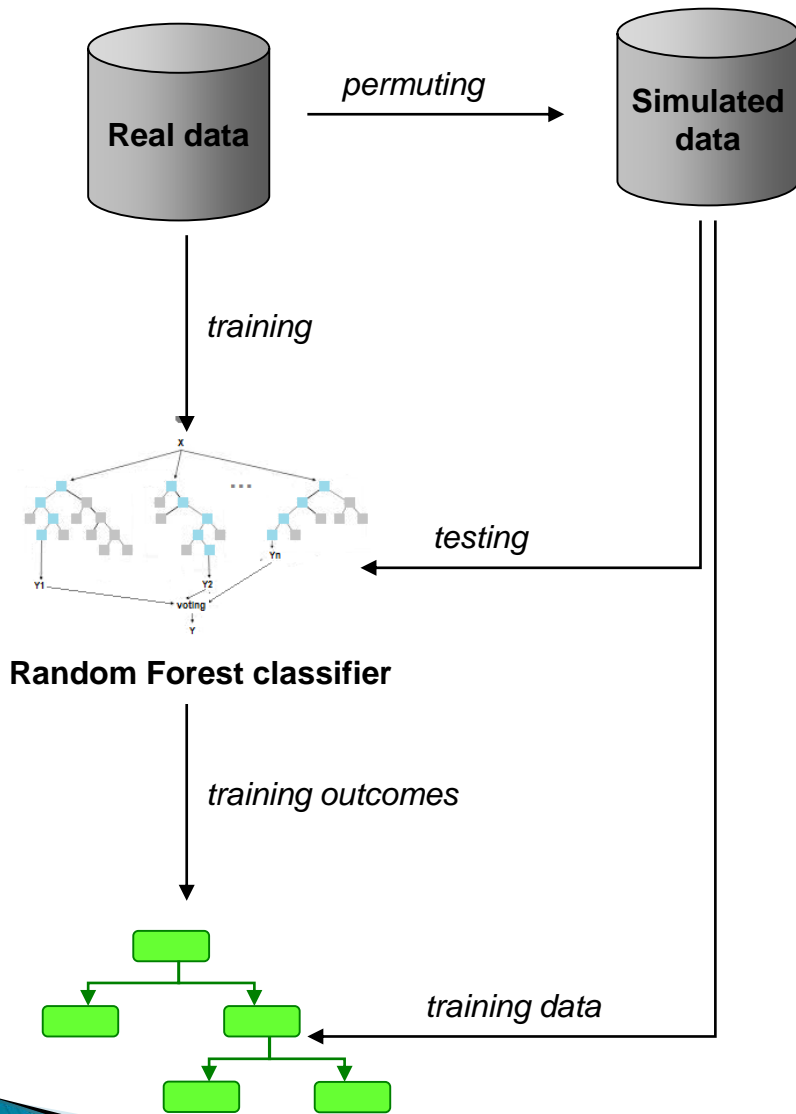Simulated data

*training*

*testing*

**Random Forest classifier**

*training outcomes*

*training data*

**Decision tree classifier**

Real data

*permuting* → Simulated data

*training*

**Random Forest classifier**

*testing*

*training outcomes*

*training data*

**Decision tree classifier**

permuting

Real data

Simulated data

training

testing

**Random Forest classifier**

**training outcomes**

training data

**Decision tree classifier**

**Real data**

*permuting*

**Simulated data**

*training*

x

...

Y1   Y2   Yn

voting

Y

*testing*

**Random Forest classifier**

*training outcomes*

*training data*

**Decision tree classifier**

permuting

Real data → Simulated data

training

**Random Forest classifier**

testing

training outcomes

training data

**Decision tree classifier**

1. **PROBLEM-TAILORED CHOICE OF THE TRAINING SET**

2. **HETEROGENOUS DATA FUSION**

3. **PHENOTYPIC INFORMATION**

# WHAT COMES NEXT?

**Data :**
- Appropriate control set
- Heterogenous sources
- Phenotypic information

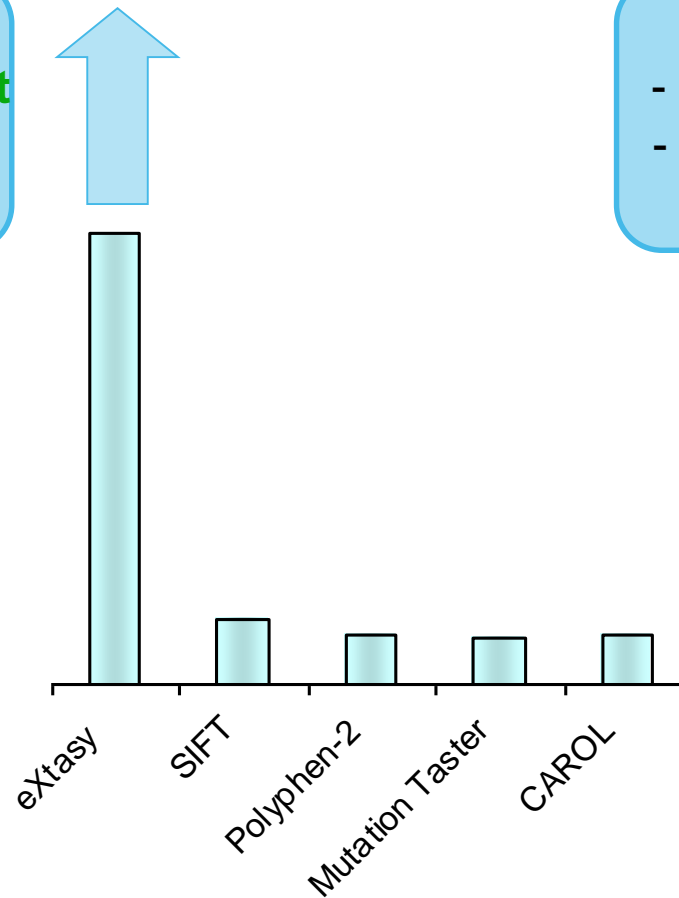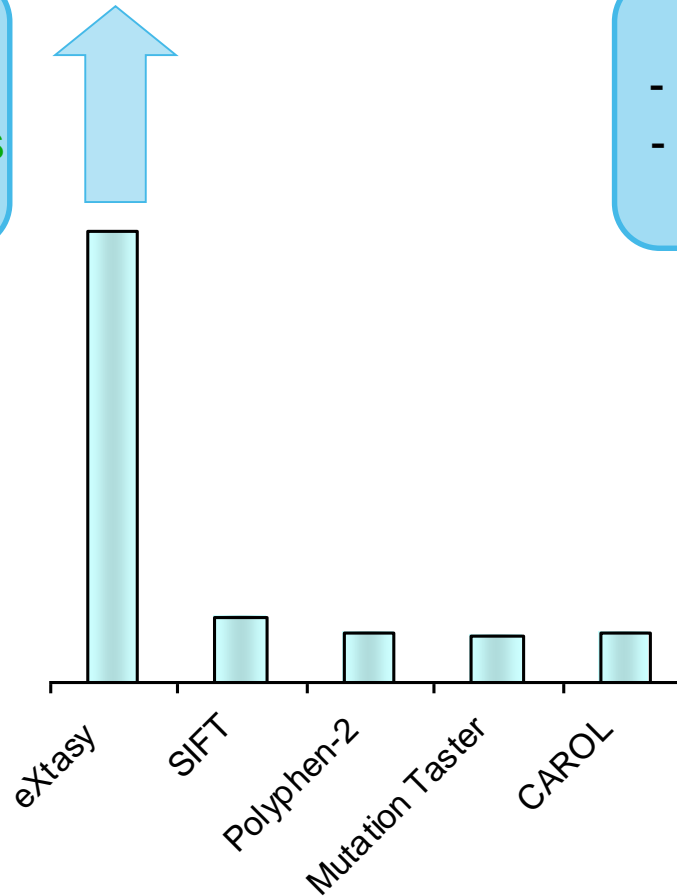**Algorithms :**
- Flexible learner
- Score aggregation

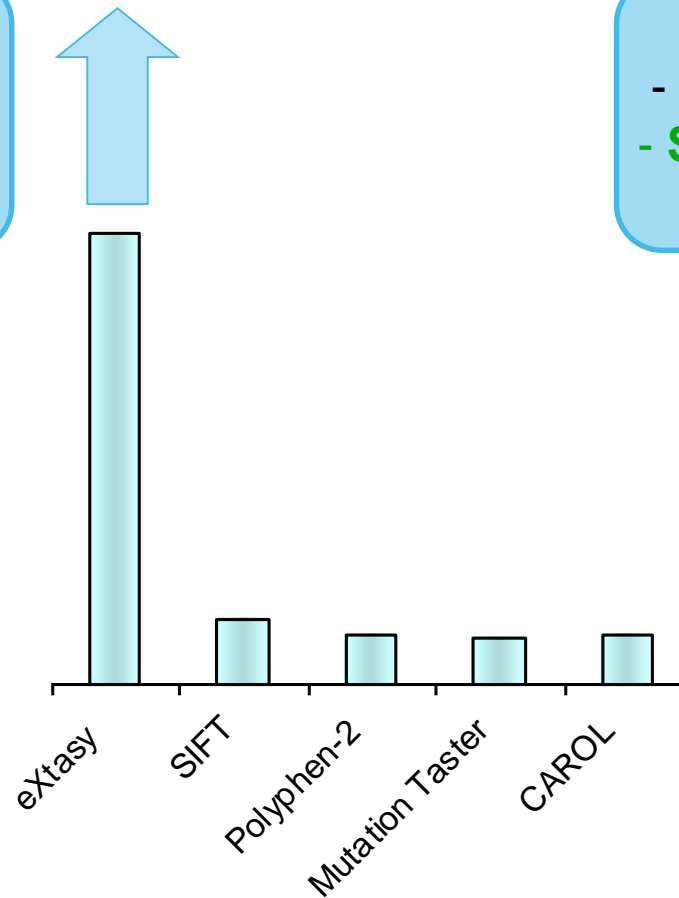eXtasy    SIFT    Polyphen-2    Mutation Taster    CAROL

**Data :**
- **Appropriate control set**
- Heterogenous sources
- Phenotypic information

**Algorithms :**
- Flexible learner
- Score aggregation



eXtasy    SIFT    Polyphen-2    Mutation Taster    CAROL

**Data :**
- Appropriate control set
- **Heterogenous sources**
- Phenotypic information

**Algorithms :**
- Flexible learner
- Score aggregation
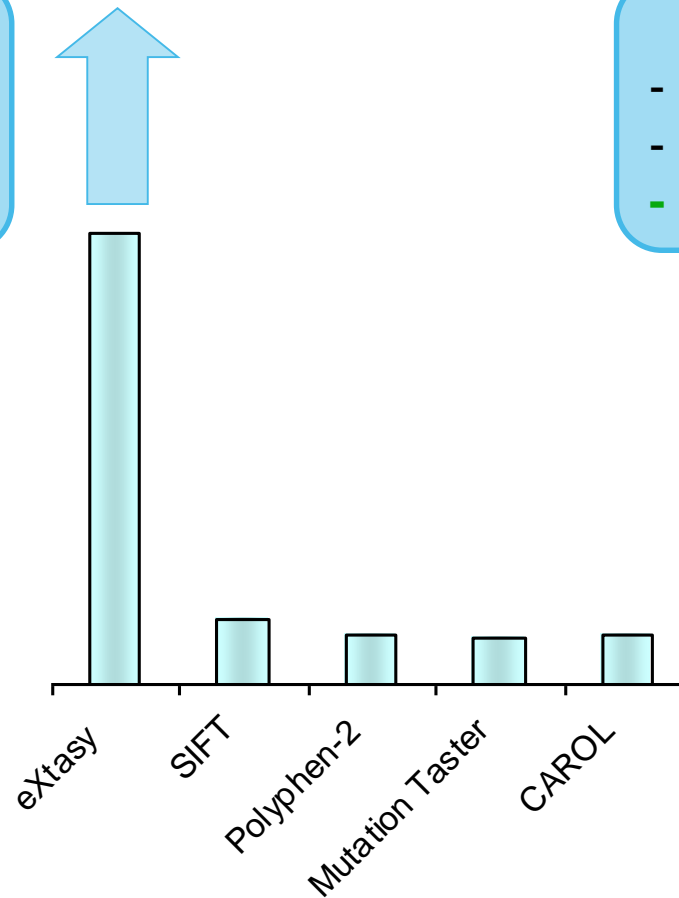


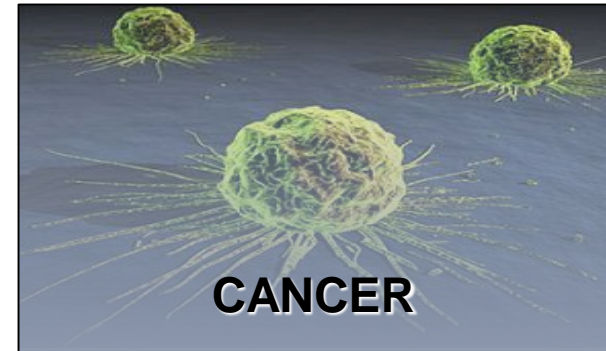eXtasy   SIFT   Polyphen-2   Mutation Taster   CAROL

**Data :**
- Appropriate control set
- Heterogenous sources
- Phenotypic information

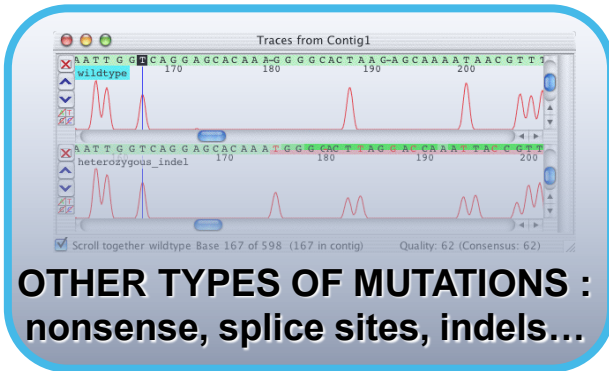**Algorithms :**
- Flexible learner
- **Score aggregation**

eXtasy   SIFT   Polyphen-2   Mutation Taster   CAROL

OTHER TYPES OF MUTATIONS :
nonsense, splice sites, indels…

CANCER

# THANK YOU!