# ABEILLE & VIOLA: novel tools to improve the diagnosis of mitochondrial diseases using omics and multi-omics data

15/05/2024

Justine LABORY, PhD student

# Medical context

Mitochondrial diseases

Medical context

Mitochondrial diseases

Rare diseases

# Medical context

**Less than 1**
person out of
**2000**

Mitochondrial diseases

Rare diseases

# Medical context

**Less than 1** person out of **2000**
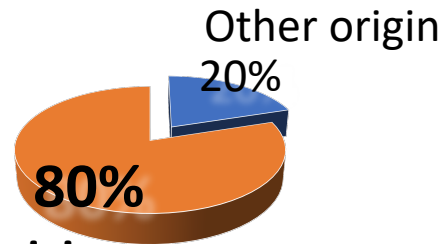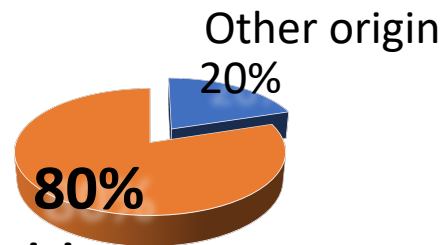
**30 millions of people**

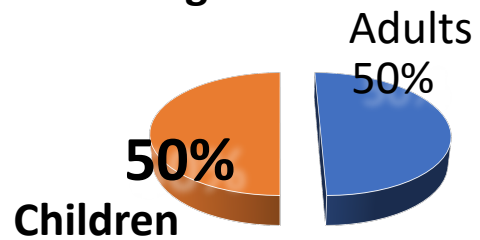Mitochondrial diseases

Rare diseases

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

Other origin 20%

**80%**

**Genetic origin**

Mitochondrial diseases

Rare diseases

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases

Mitochondrial diseases

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases

Mitochondria

Mitochondrial diseases

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

responsible for a wide variety of biochemical processes

Mitochondria

## Mitochondrial diseases

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

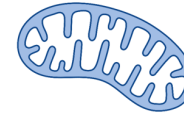**Children**

Rare diseases

responsible for a wide variety of biochemical processes

Mitochondria

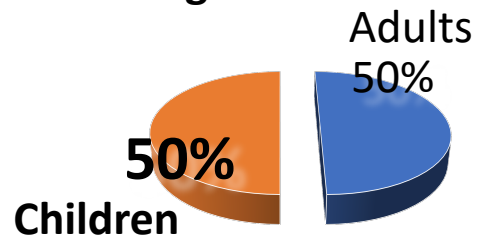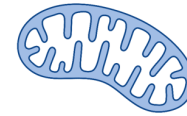under the double control of mtDNA and nDNA

Mitochondrial diseases

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases

responsible for a wide variety of biochemical processes

Mitochondria

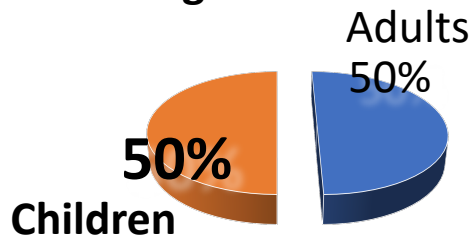under the double control of mtDNA and nDNA

Mitochondrial diseases

Disease

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

responsible for a wide variety of biochemical processes

Mitochondria

under the double control of mtDNA and nDNA

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases

deficiency of the mitochondrial respiratory chain

Mitochondrial diseases

Disease

# Medical context

Less than 1 person out of 2000

30 millions of people

Other origin 20%
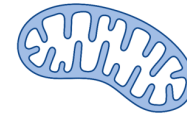
**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases
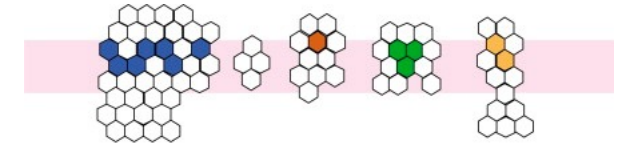
responsible for a wide variety of biochemical processes

Mitochondria

under the double control of mtDNA and nDNA
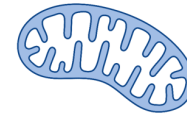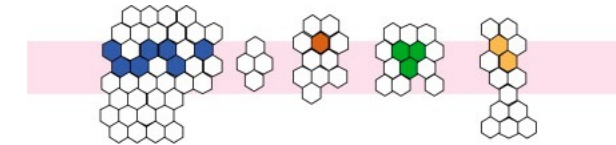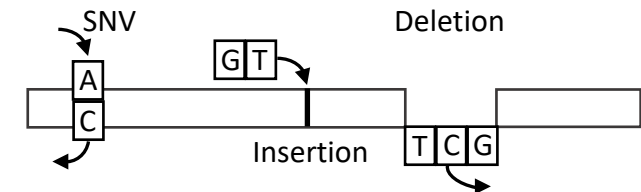
## Mitochondrial diseases

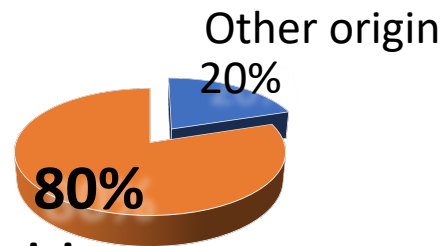deficiency of the mitochondrial respiratory chain

SNV

Deletion

A

C

G T

Insertion

T C G

due to rare hereditary or spontaneous variants of mtDNA or nDNA

Disease

# Medical context

Less than 1 person out of 2000

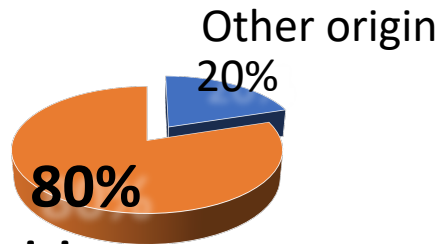30 millions of people

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

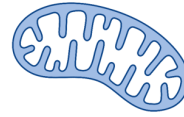**Children**

Rare diseases

responsible for a wide variety of biochemical processes
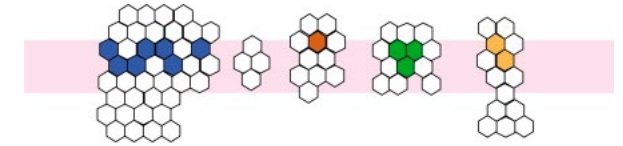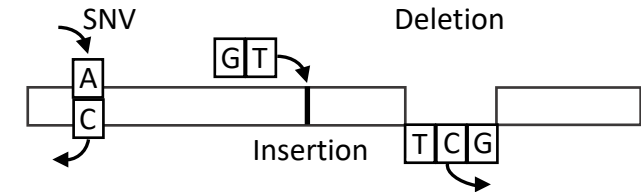
Mitochondria

under the double control of mtDNA and nDNA

Mitochondrial diseases

Diagnosis

deficiency of the mitochondrial respiratory chain

SNV          Deletion

A
C          G  T

Insertion          T  C  G

due to rare hereditary or spontaneous variants of mtDNA or nDNA

Disease

# Medical context

**Less than 1** person out of **2000**

**30 millions of people**

Other origin 20%

**80%**

**Genetic origin**

Adults 50%

**50%**

**Children**

Rare diseases

responsible for a wide variety of biochemical processes

Mitochondria

under the double control of mtDNA and nDNA

## Mitochondrial diseases

Responsible gene bearing the pathogenic variant

Diagnosis

patient    gene      patient    gene

Mitochondrial diseases

Other diseases

deficiency of the mitochondrial respiratory chain

SNV

A
C

Deletion

G T

Insertion

T C G

due to rare hereditary or spontaneous variants of mtDNA or nDNA

Disease

# Medical context

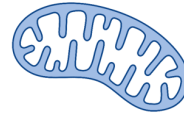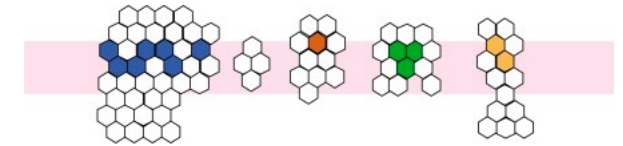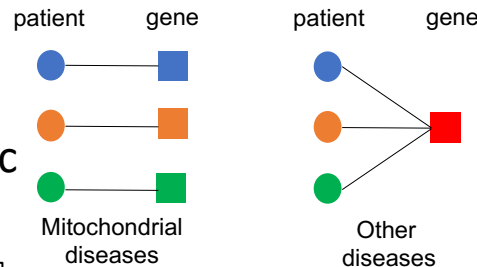**Less than 1** person out of **2000**

**30 millions of people**

responsible for a wide variety of biochemical processes
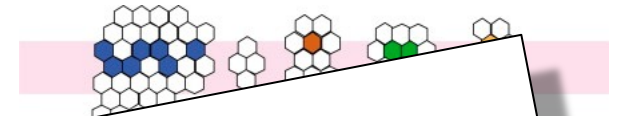
Mitochondria
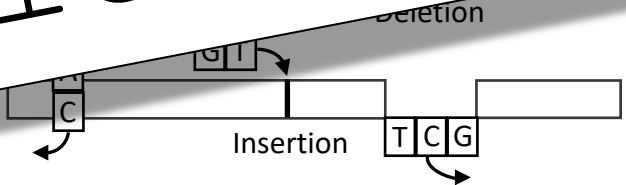
under the

1 patient out of 2 is in diagnostic stalemate

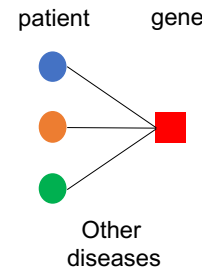due to rare hereditary or spontaneous variants of mtDNA or nDNA

Deletion

Insertion

Disease

Geneti

50%

**50%**
**Children**

Responsible gene bearing the pathogenic variant

patient   gene     patient   gene

Mitochondrial diseases

Other diseases

Rare diseases

Diagnosis

# Diagnosis of Mitochondrial Disease (MD)

Diagnostic = identification of the genetic variant responsible for the disease

# Diagnosis of Mitochondrial Disease (MD)

Diagnostic = identification of the genetic variant responsible for the disease

DNA
extraction

Sequencing

# Diagnosis of Mitochondrial Disease (MD)

Diagnostic = identification of the genetic variant responsible for the disease



DNA
extraction

Sequencing

**mtDNA**

Circular mitochondrial
genome (16Kb) independent
of nuclear genome

# Diagnosis of Mitochondrial Disease (MD)

Diagnostic = identification of the genetic variant responsible for the disease

DNA
extraction

Sequencing

**mtDNA**

Circular mitochondrial
genome (16Kb) independent
of nuclear genome
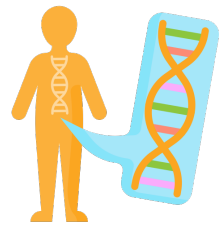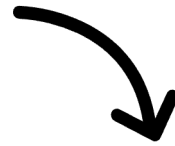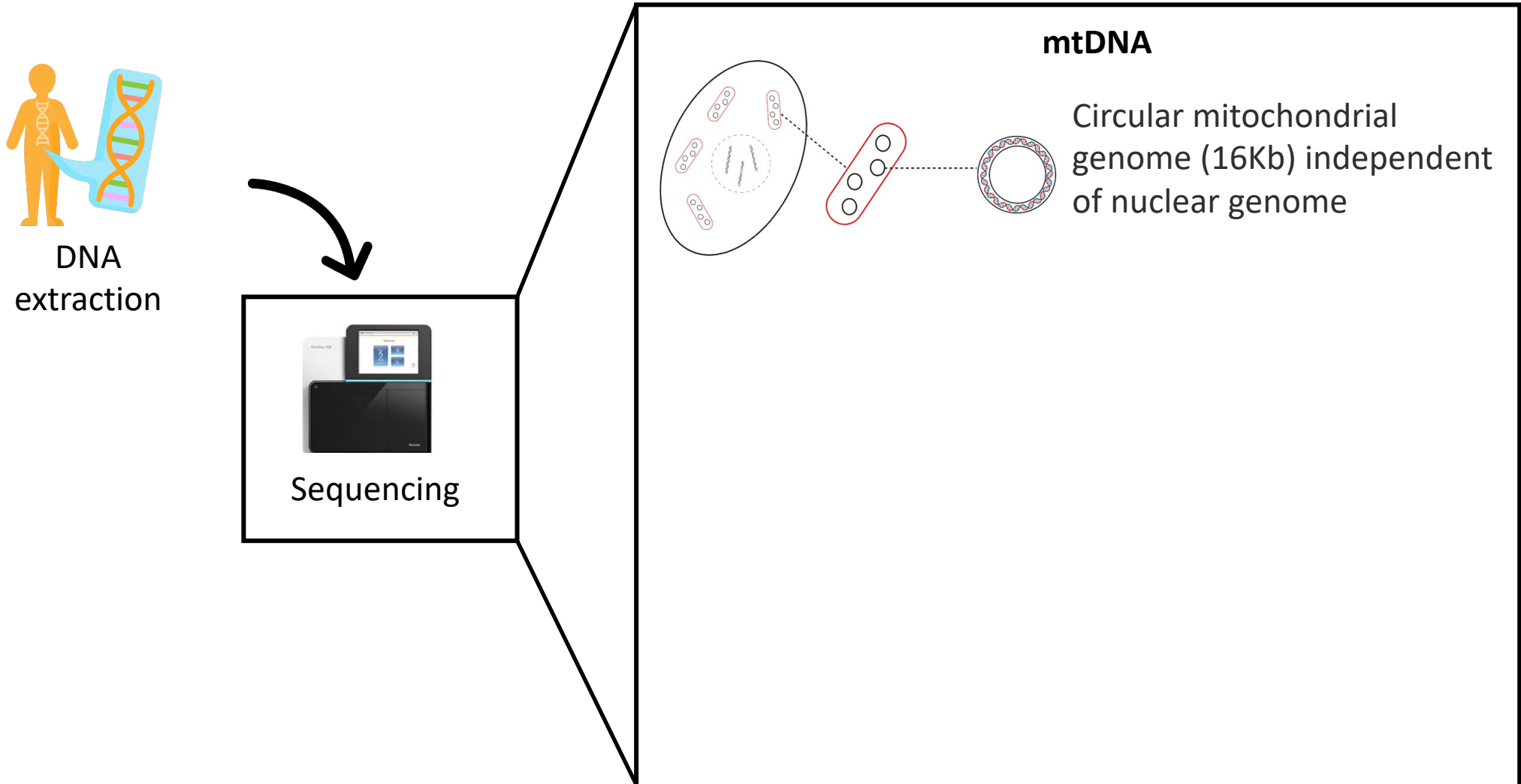
**Clinical exome**

Sequencing of genes
known to be involved
in rare diseases

# Diagnosis of Mitochondrial Disease (MD)

Diagnostic = identification of the genetic variant responsible for the disease
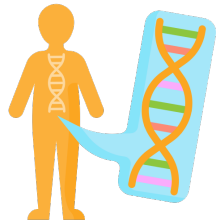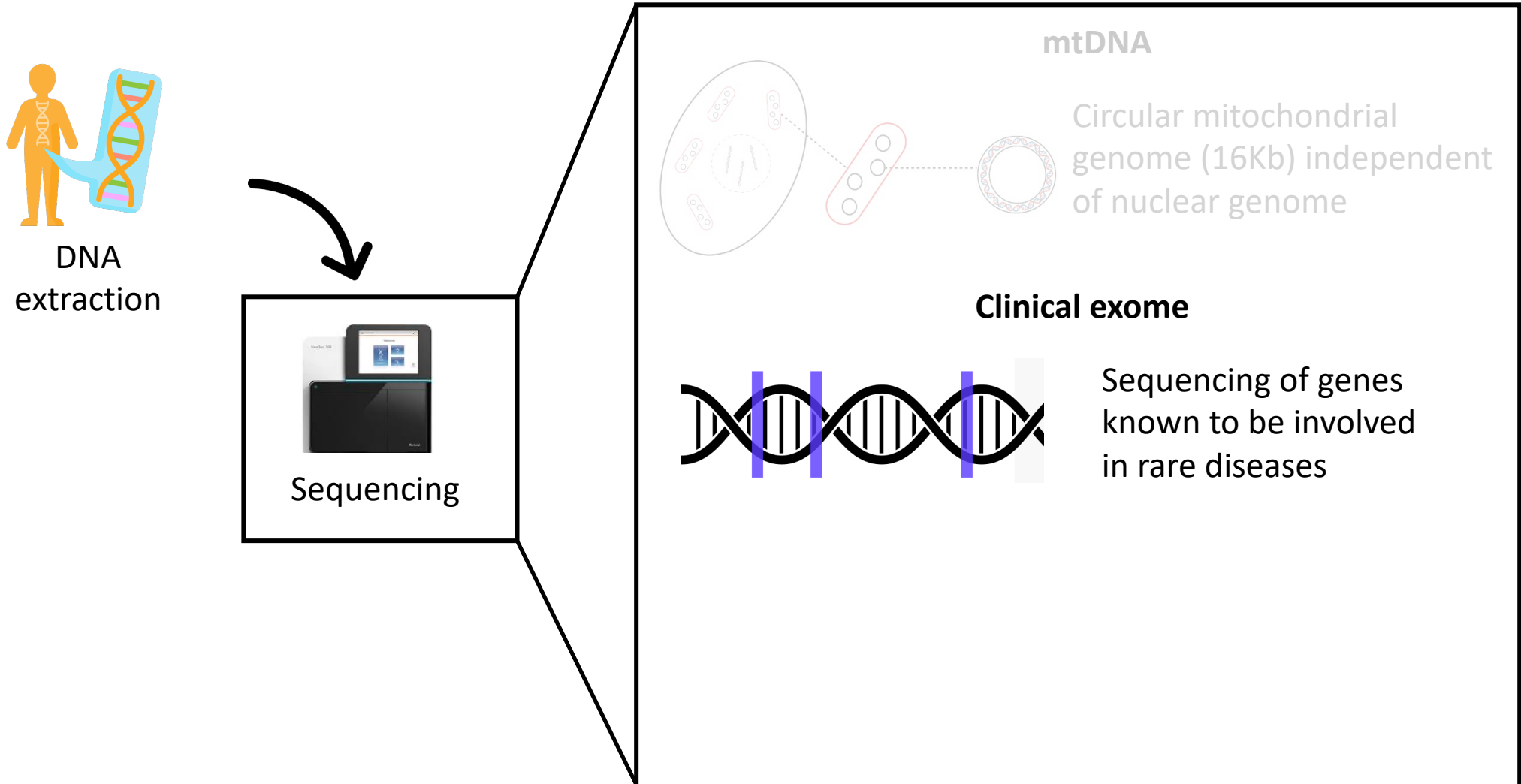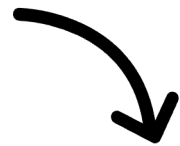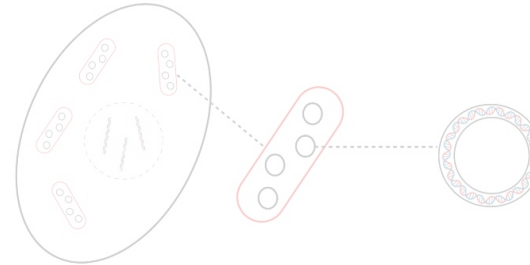
DNA extraction

Sequencing

**mtDNA**

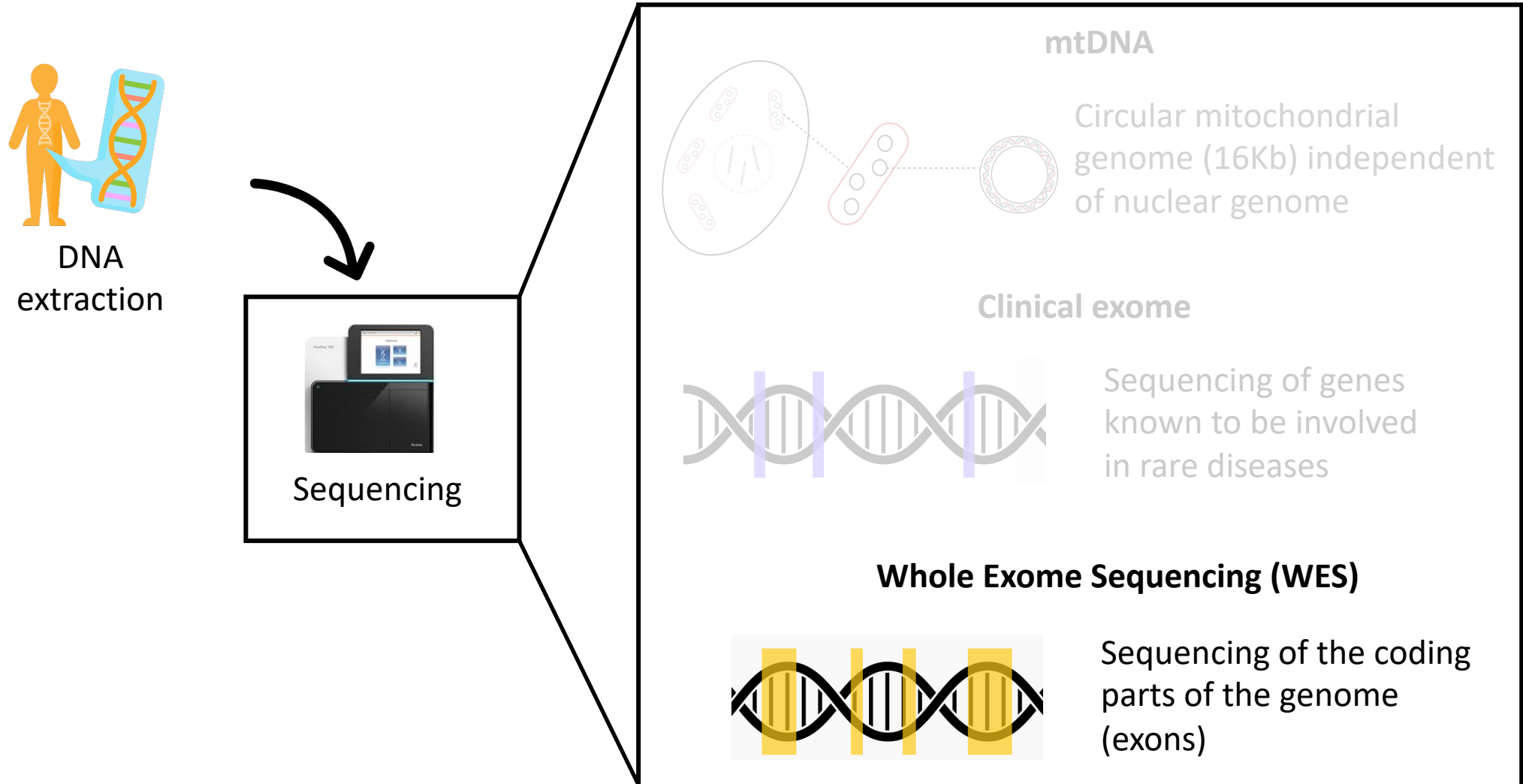Circular mitochondrial genome (16Kb) independent of nuclear genome

**Clinical exome**

Sequencing of genes known to be involved in rare diseases

**Whole Exome Sequencing (WES)**

Sequencing of the coding parts of the genome (exons)

# Diagnosis of Mitochondrial Disease (MD)

Diagnostic = identification of the genetic variant responsible for the disease

DNA extraction

Sequencing

Identification of genetic variants

Bioinformatic pipeline

**mtDNA**

Circular mitochondrial genome (16Kb) independent of nuclear genome

**Clinical exome**

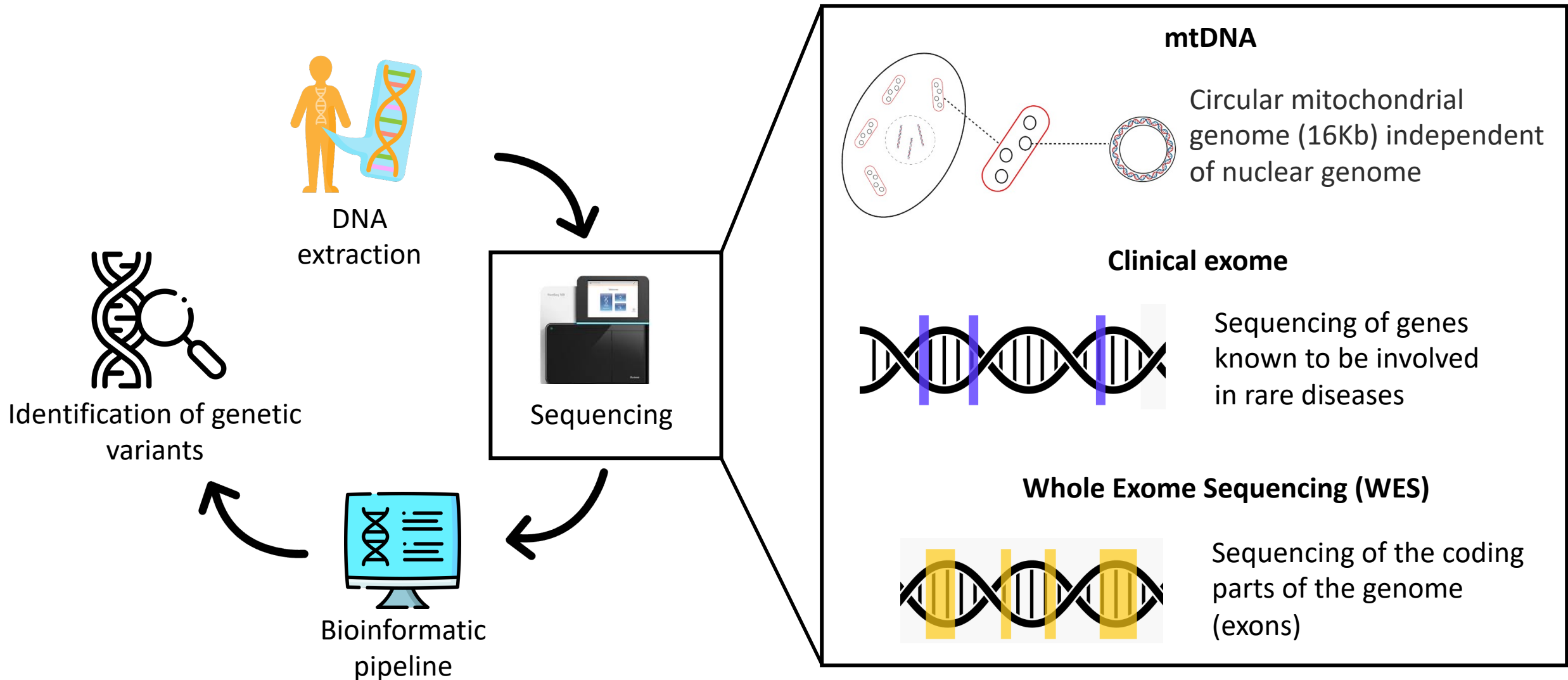Sequencing of genes known to be involved in rare diseases

**Whole Exome Sequencing (WES)**
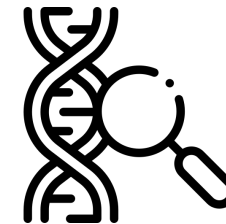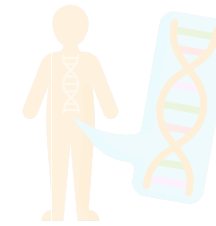
Sequencing of the coding parts of the genome (exons)

# Diagnosis of Mitochondrial Disease (MD)

**Challenges** :

The variant can be anywhere in genetic sequence (intron, exon, regulatory sequence)

DNA extraction

Sequencing

Identification of genetic variants
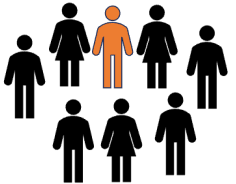
Bioinformatic pipeline
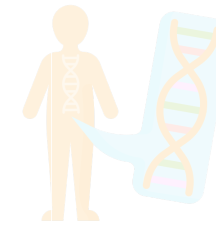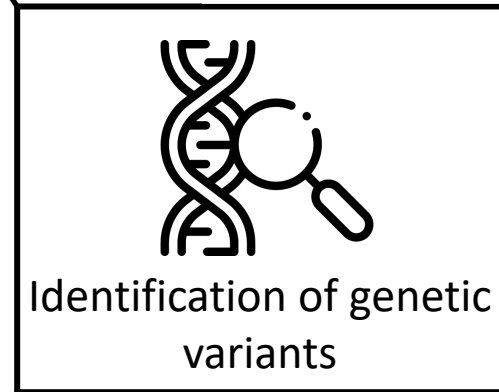
# Diagnosis of Mitochondrial Disease (MD)



**Challenges** :

The variant can be anywhere in genetic sequence (intron, exon, regulatory sequence)

The variant is rare, i.e is present in less than 1% of the population.
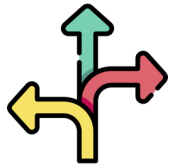
DNA extraction

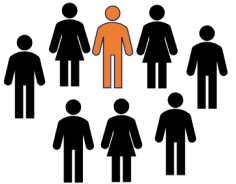Sequencing

Bioinformatic pipeline

Identification of genetic variants

# Diagnosis of Mitochondrial Disease (MD)

**Challenges** :

 The variant can be anywhere in genetic sequence (intron, exon, regulatory sequence)
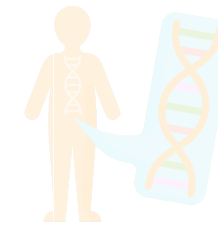
 The variant is rare, i.e is present in less than 1% of the population.

 Variant poorly characterized in databases

Identification of genetic variants

DNA extraction
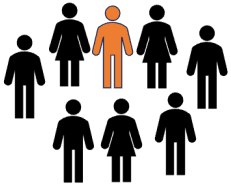
Sequencing

Bioinformatic pipeline

# Diagnosis of Mitochondrial Disease (MD)
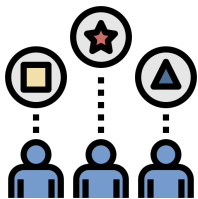
**Challenges** :

The variant can be anywhere in genetic sequence (intron, exon, regulatory sequence)

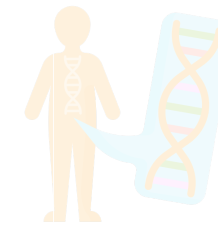The variant is rare, i.e is present in less than 1% of the population.

Variant poorly characterized in databases

1 different responsible variant for each patient

DNA extraction

Sequencing

Identification of genetic variants

Bioinformatic pipeline

# Diagnosis of Mitochondrial Disease (MD)

**Challenges** :

The variant can be anywhere in genetic
sequer
sequer

The va
than 1

Identification of too many variants

Identification of genetic
variants

Variant poorly characterized in databases

1 different responsible variant for each
patient

Bioinformatic
pipeline

# The diagnostic power of MD

# RNA-sequencing to improve MD diagnosis



Article | Published: 14 October 2019

**Diagnostic utility of transcriptome sequencing for rare Mendelian diseases**

RESEARCH ARTICLE | GENETIC DIAGNOSIS

**Improving genetic diagnosis in Mendelian disease with transcriptome sequencing**

## RNA-sequencing

- Aberrant expression
- Aberrant splicing
- Monoallelic expression

**Clinical implementation of RNA sequencing for Mendelian disease diagnostics**

**Identification of rare-disease genes using blood transcriptome sequencing and large control cohorts**

**ARTICLE**

Expanding the Boundaries of RNA Sequencing as a Diagnostic Tool for Rare Mendelian Disease

**Integration of proteomics with genomics and transcriptomics increases the diagnostic rate of Mendelian disorders**

# The diagnostic power of MD

# The diagnostic power of MD

# Objectives

## How to improve the diagnosis of mitochondrial diseases ?



From Labory et al. *Front Mol Biosci*.2020

# Table of contents



**1**   ABEILLE (ABerrant Expression Identification empLoying machine LEarning) to find candidate Aberrant Gene expression (AGEs)

**2**   VIOLA (Variant prIOritization using LAtent space) to find candidate pathogenic genetic variants

# Table of contents

**1**   ABEILLE (ABerrant Expression Identification empLoying machine LEarning) to find candidate Aberrant Gene expression (AGEs)

Gene expression

**ABEILLE: a novel method for ABerrant Expression Identification empLoying machine LEarning from RNA-sequencing data**

Justine Labory[1,2,†], Gwendal Le Bideau[2,†], David Pratella[1], Jean-Elisée Yao[1], Samira Ait-El-Mkadem Saadi[2], Sylvie Bannwarth[2], Loubna El-Hami[1,2], Véronique Paquis-Fluckinger[2,‡] and Silvia Bottini [1,*,‡]

ABEILLE

https://github.com/UCA-MSI/ABEILLE

# Context : RNA-seq questions

## 2 approaches

# Context : RNA-seq questions
## 2 approaches

*Which genes are differentially expressed between 2 groups ?*

**Differential Expression (DE)**

| | Control 1 | Control 2 | Control 3 | Patient 1 | Patient 2 | Patient 3 |
|---|---|---|---|---|---|---|
| Gene A | | | | | | |
| Gene B | | | | | | |
| Gene C | | | | | | |
| | | | | | | |

⟶ Gene A = **DE**

⟶ Gene B = Normal gene

⟶ Gene C = **DE**

🔍 Control group

🔍 Replicates

Tool: DESeq2[1]

[1]Love et al. *Genome Biology.*2014

# Context : RNA-seq questions
## 2 approaches

*Which genes are differentially expressed between 2 groups ?*

*Which genes are AGEs for each patient ?*

**Differential Expression (DE)**

| | Control 1 | Control 2 | Control 3 | Patient 1 | Patient 2 | Patient 3 |
|---|---|---|---|---|---|---|
| Gene A | | | | | | |
| Gene B | | | | | | |
| Gene C | | | | | | |
| | | | | | | |

→ Gene A = **DE**

→ Gene B = Normal gene

→ Gene C = **DE**

🔍 Control group

🔍 Replicates

Tool: DESeq2[1]

**Aberrant Gene Expression (AGE)**

| | Patient 1 | Patient 2 | Patient 3 |
|---|---|---|---|
| Gene A | | | |
| Gene B | | | |
| Gene C | | | |
| | | | |

→ [Gene A;Patient 2] = **AGE**

→ Gene B = Normal gene

→ [Gene C;Patient 1] = **AGE**

No control group

No replicates

[1]Love et al. *Genome Biology.*2014

# Methods to identify AGEs

OUTRIDER[1]

Autoencoder

**+**

Statistical test

[1]Brechtmann et al. *Am. J. Hum. Genet.* 2018

# Methods to identify AGEs

OUTRIDER[1]

Autoencoder



**+**

Statistical test
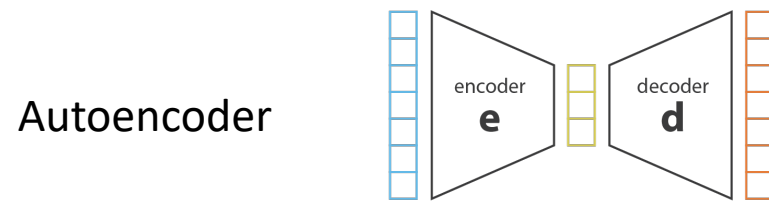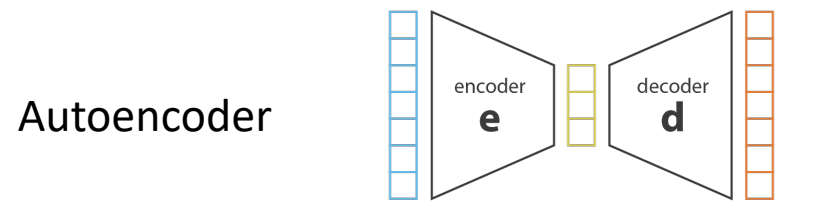


➔ OUTRIDER does not work on small data

[1]Brechtmann et al. *Am. J. Hum. Genet.* 2018

# Methods to identify AGEs

## OUTRIDER[1]

Autoencoder

+

Statistical test

➜ OUTRIDER does not work on small data

## OutPyR[2]

Bayesian model

+

Statistical test

[1]Brechtmann et al. *Am. J. Hum. Genet.* 2018

[2]Salkovic et al. *Journal of Computational Science,.*2020

# Methods to identify AGEs

| OUTRIDER[1]  | OutPyR[2] |
|---|---|

Autoencoder 

**+**

Statistical test 

➔ OUTRIDER does not work on small data

Bayesian model 

**+**

Statistical test 

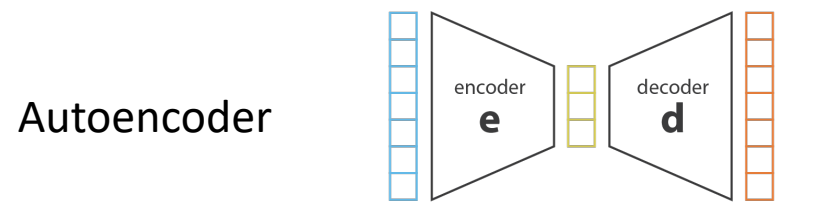➔ Tested only on a subset of real data

[1]Brechtmann et al. *Am. J. Hum. Genet.* 2018                    [2]Salkovic et al. *Journal of Computational Science,.*2020

# Methods to identify AGEs
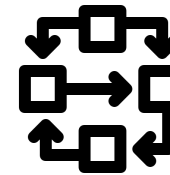
**How to identify AGEs for small cohorts ?**

# Methods to identify AGEs
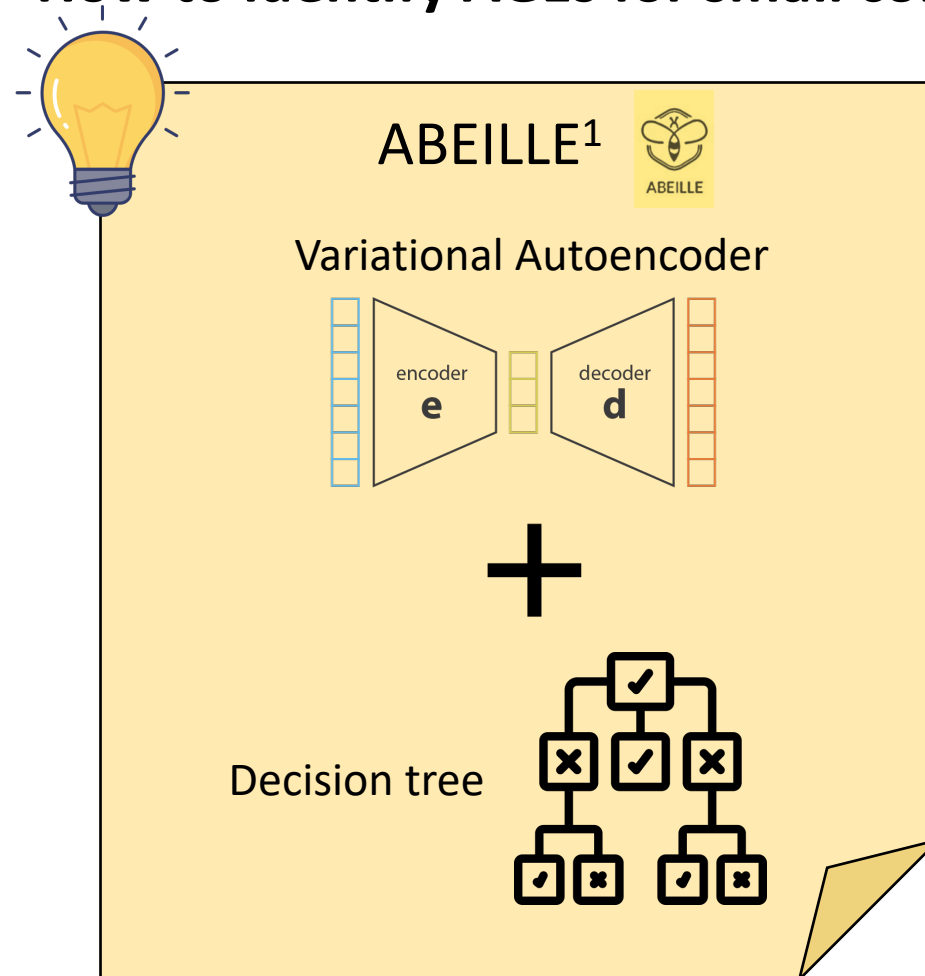
**How to identify AGEs for small cohorts ?**



ABEILLE[1]

Variational Autoencoder

encoder
e

decoder
d

+

Decision tree

[1]**Labory et al.** *Bioinformatics* 2022

# The autoencoder

the process that produce the "new features" representation

Latent space

the reverse process of an encoder

Initial data

**e** encoder

**d** decoder

Encoded – decoded data

encoded data

Dimensionality Reduction          Denoising          Sequence to sequence prediction

Compression          Feature Extraction          Generation          Recommendation system

# How to use AE to identify AGEs ?

**Input :**
Gene expression data from RNA-seq

**Output:**
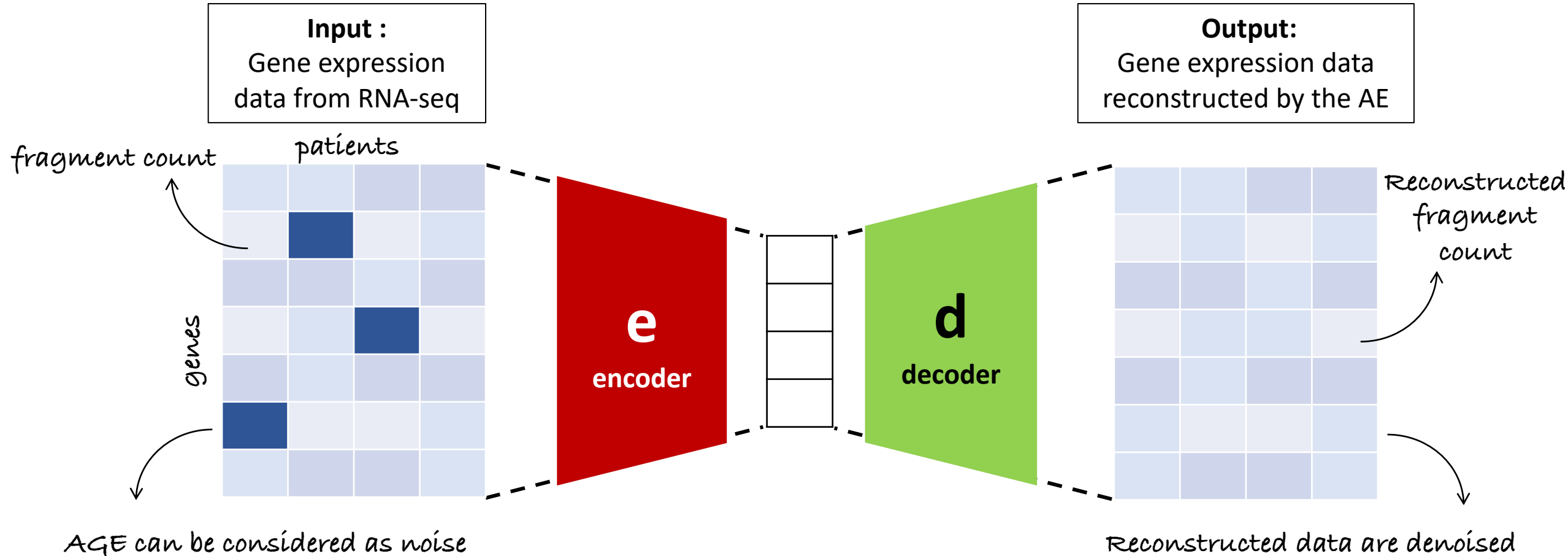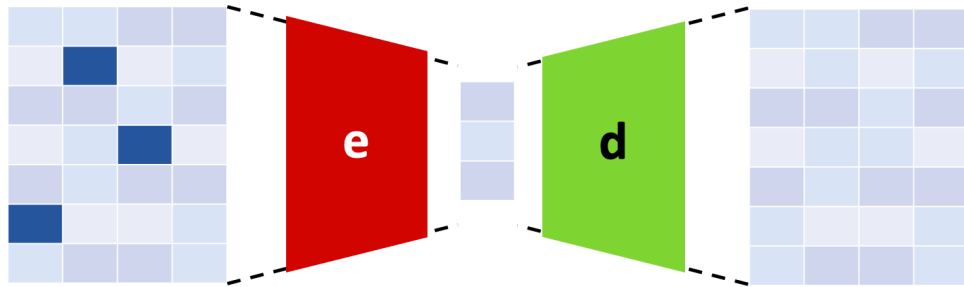Gene expression data reconstructed by the AE

fragment count

patients

genes

**e**
encoder

**d**
decoder

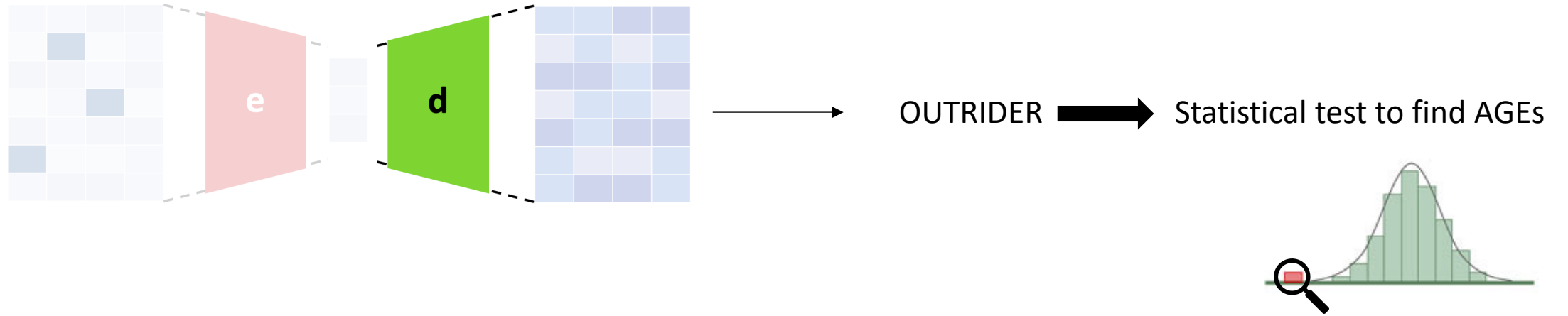Reconstructed fragment count

AGE can be considered as noise

Reconstructed data are denoised

# Difference between ABEILLE and OUTRIDER

# Difference between ABEILLE and OUTRIDER



OUTRIDER → Statistical test to find AGEs

# Difference between ABEILLE and OUTRIDER



OUTRIDER ⟹ Statistical test to find AGEs

ABEILLE

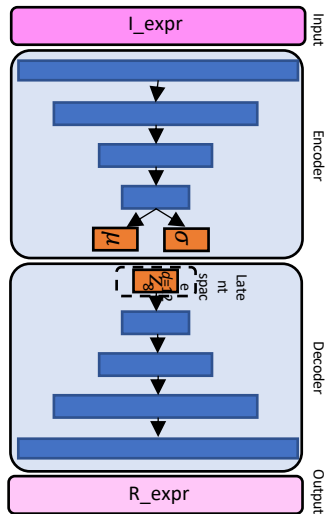To compare input and output of the VAE to find AGEs

# ABEILLE workflow

# ABEILLE workflow

**1**  **VAE**

To use VAE to generate reconstructed denoised counts

# ABEILLE workflow

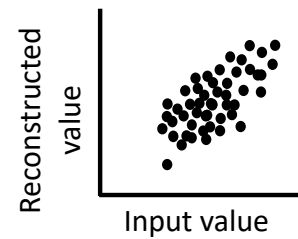# ABEILLE workflow

**1** **VAE**

To use VAE to generate reconstructed denoised counts



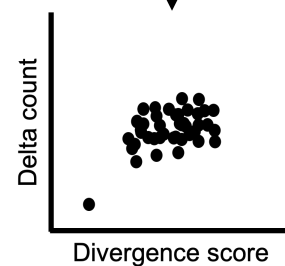**2** **Definition of metrics**

To compute metrics to assess the reconstruction fidelity
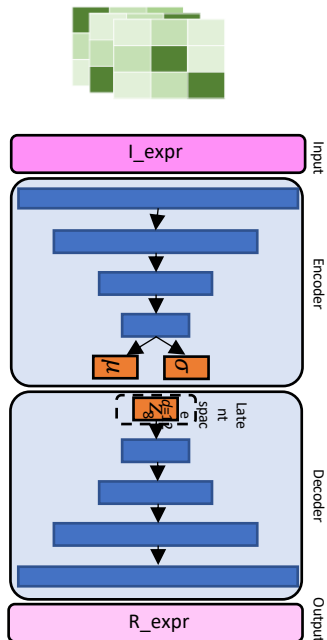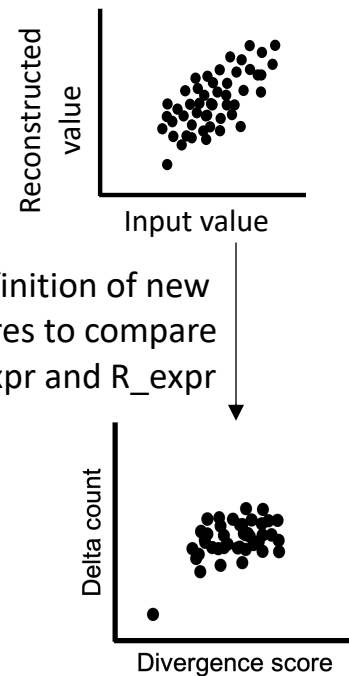
Definition of new scores to compare I_expr and R_expr

**3** **Classification**

To classify gene expressions as AGEs or no AGEs

Supervised          Unsupervised

Semi-synthetics data          Real data

Decision tree

# ABEILLE workflow



**1** VAE

To use VAE to generate reconstructed denoised counts

I_expr          Input

Encoder

μ          σ

Late nt spac e          Decoder

R_expr          Output

**2** **Definition of metrics**

To compute metrics to assess the reconstruction fidelity

Reconstructed value

Input value

Definition of new scores to compare I_expr and R_expr

Delta count

Divergence score

**3** **Classification**

To classify gene expressions as AGEs or no AGEs

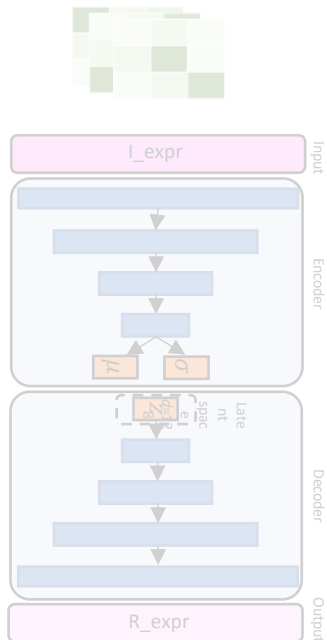Supervised          Unsupervised

Semi-synthetics data          Real data

Decision tree
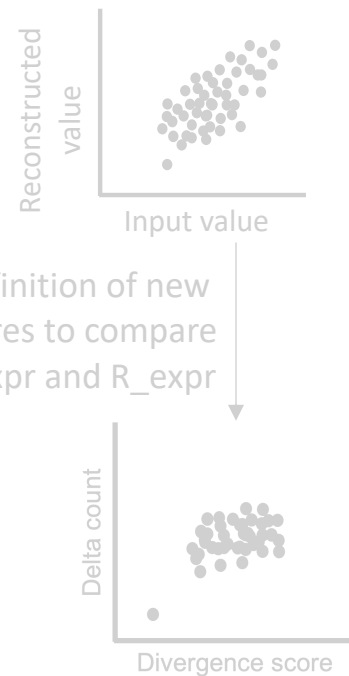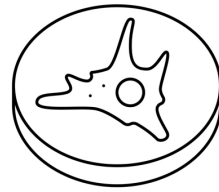
# Supervised phase – Creation of semi-synthetics datasets

54 tissue　　　　　1000 individuals　　　　56 200 transcripts

# Supervised phase – Creation of semi-synthetics datasets

**1** tissue          **504** individuals          56 200 transcripts

# Supervised phase – Creation of semi-synthetics datasets
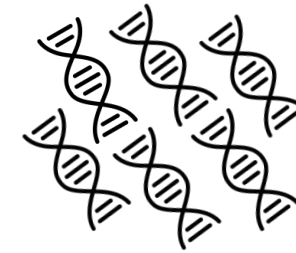
**1** tissue          **504** individuals          56 200 transcripts
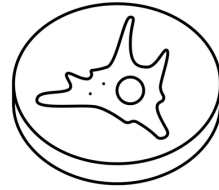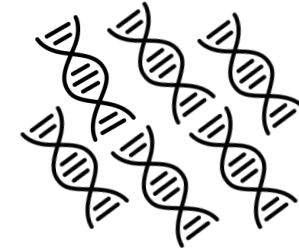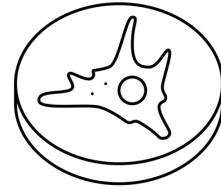
# Supervised phase – Creation of semi-synthetics datasets

**1** tissue      **504** individuals      56 200 transcripts



Generate computational AGEs
by replacing randomly 10 000
expression values

$$k_{ij}^O = \mathrm{round}\left(s_i 2^{\mu_j^u \pm \exp(N)\sigma_j^u}\right)$$

# Supervised phase – Creation of semi-synthetics datasets

**1** tissue                **504** individuals          56 200 transcripts
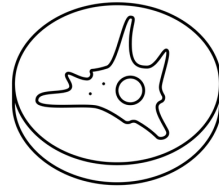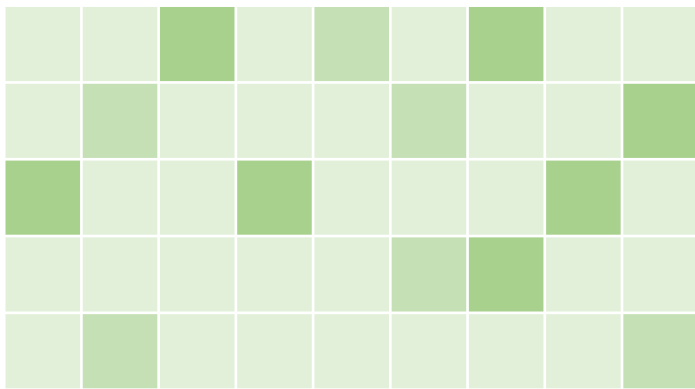


Generate computational AGEs by replacing randomly 10 000 expression values

$$k_{ij}^{O} = \text{round}\left(s_i 2^{\mu_j^u \pm \exp(N)\sigma_j^u}\right)$$

Repeat the process 20 times

# Supervised phase – To obtain the decision tree

①  **To use VAE to generate reconstructed denoised counts**



②  **To compute metrics to assess the reconstruction fidelity**



③  **To create a decision tree and identify thresholds for gene expression classification**

*Gene A*



● Injected AGE
● No AGE

Linear regression

Parameters calculated on each linear regression :
- Dfbetas
- Hat
- Type error

# Supervised phase – To obtain the decision tree

**1** **To use VAE to generate reconstructed denoised counts**



**2** **To compute metrics to assess the reconstruction fidelity**



**3** **To create a decision tree and identify thresholds for gene expression classification**



*Gene A*

Parameters calculated on each linear regression :
- Dfbetas
- Hat
- Type error

Linear regression

To be done for each gene

# Supervised phase – To obtain the decision tree

**①** **To use VAE to generate reconstructed denoised counts**



**②** **To compute metrics to assess the reconstruction fidelity**



**③** **To create a decision tree and identify thresholds for gene expression classification**
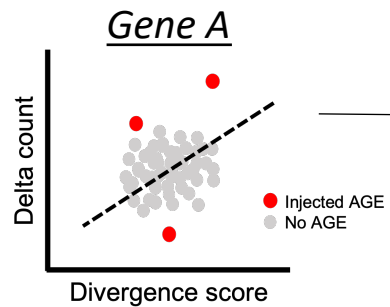


*Gene A*

Parameters calculated on each linear regression :
- Dfbetas
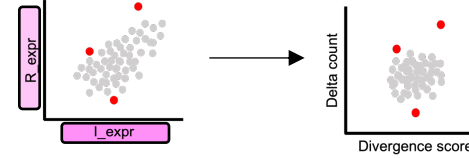- Hat
- Type error

Linear regression

To be done for each gene

Parameters calculated for all genes in all patients are used to feed a decision tree

This decision tree will be used in the unsupervised phase for gene expression classification

# ABEILLE workflow

**1** **VAE**

To use VAE to generate reconstructed denoised counts

**2** **Definition of metrics**

To compute metrics to assess the reconstruction fidelity

Definition of new scores to compare I_expr and R_expr

**3** **Classification**

To classify gene expressions as AGEs or no AGEs

Supervised

Unsupervised

Semi-synthetics data

Real data

Decision tree

# Unsupervised phase – gene expression classification

**1** To use VAE to generate reconstructed denoised counts

**2** To compute metrics to assess the reconstruction fidelity

**3** Classification of gene expressions as AGEs or no AGEs

# Unsupervised phase – gene expression classification

1 To use VAE to generate reconstructed denoised counts

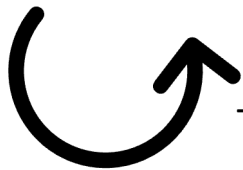2 To compute metrics to assess the reconstruction fidelity

3 Classification of gene expressions as AGEs or no AGEs
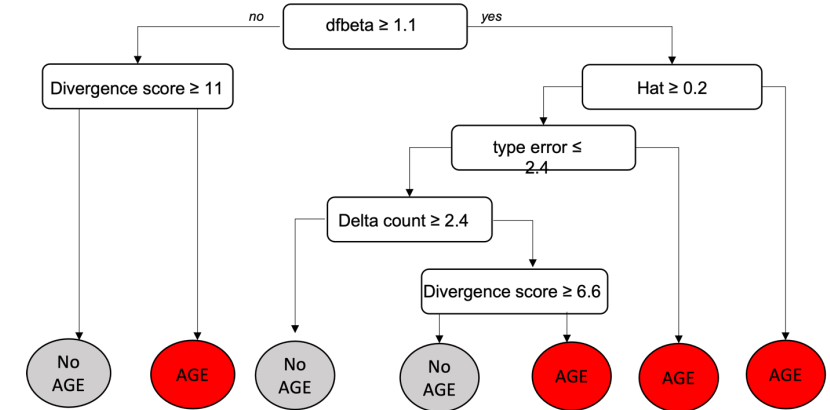


To be done for each gene

# Case study

119 patients with MD suspicion
(from Kremer et al. *Nat Comm* 2017)

RNA-seq

Validation of 5 candidate genes in 6 patients

# Case study



119 patients with MD suspicion
(from Kremer et al. *Nat Comm* 2017)

RNA-seq

Validation of 5 candidate genes in 6 patients

**Goal** : Compare ABEILLE to other methods

— ABEILLE
— OUTRIDER
— OutPyR
— DESeq2

# Performances of the four tools on real dataset



These observations rule out OutPyR as a tool for AGE identification in this context.

# Performances of ABEILLE and OUTRIDER



AGEs found by ABEILLE are more enriched in terms related to mitochondrial biology than the AGEs found by OUTRIDER.

# Validated pathogenic genes

| Validated pathogenic genes | Detected by | ABEILLE | OUTRIDER |
|---|---|---|---|
| *MGST1* | AGE | ✓ | ✓ |
| *TIMMDC1* | AGE | ✓ | ✓ |
| *MCOLN1* | AGE | ✓ | ✓ |

ABEILLE & OUTRIDER correctly classify the pathogenic genes as AGEs

# Validated pathogenic genes

| Validated pathogenic genes | Detected by | ABEILLE | OUTRIDER |
|---|---|---|---|
| *ALDH18A1* | MAE | X | ✔ |
| *CLPP* | AS | X | ✔ |

OUTRIDER classifies as AGEs two pathogenic genes that do not show aberrant expression (putative false positives)



ALDH18A1

MUC1404



CLPP

MUC1350

delta count

divergence score

# AGE detection on small dataset size

# AGE detection on small dataset size

| 110 | 90 | 60 | 30 | 20 | 10 | 110 | 90 | 60 | 30 | 20 | 10 |

Number of samples

# AGE detection on small dataset size



| 110 | 90 | 60 | 30 | 20 | 10 | 110 | 90 | 60 | 30 | 20 | 10 |

Number of samples

Number of datasets

0
2
4
6
8
10

# AGE detection on small dataset size

# AGE detection on small dataset size

# AGE detection on small dataset size



The performances of ABEILLE do not depend on the number of samples

# Conclusion of part 1

## ADVANTAGES

- ABEILLE identifies AGEs from RNA-seq data without the need of replicates

- ABEILLE showed good performances on small datasets

## LIMITATIONS

- The decision tree must be trained for each different type of data

- The choice of semi-synthetics datasets to feed the decision tree

## PERSPECTIVES

- Use a flexible model to work on any type of data

# Perspectives

We are developing a version 2 of ABEILLE :



DBSCAN : density based model

Multi-omics analyses are now possible

# Table of contents

**1** ABEILLE (ABerrant Expression Identification empLoying machine LEarning) to find candidate Aberrant Gene expression (AGEs)

**2** VIOLA (Variant prIOritization using LAtent space) to find candidate pathogenic genetic variants

Genomics

Phenomics

Transcriptomics

# Diagnosis of Mitochondrial Disease (MD)

**Challenges** :

The variant can be anywhere in genetic
sequence
sequence

The variant
than 1

Variant poorly characterized in databases

1 different responsible variant for each
patient

Identification of genetic
variants

### ⚠ Identification of too many variants

Bioinformatic
pipeline

# Diagnosis of Mitochondrial Disease (MD)

**Challenges** :

The variant can be anywhere in genetic sequence

The variant

than 1

Variant poorly characterized in databases

1 different responsible variant for each patient

Identification of genetic variants

Bioinformatic pipeline

💡 Variant prioritization

# Variant prioritization

Process of **selecting** and **ranking** genetic variants based on their potential **significance** or relevance to a specific **phenotype** or condition.

State-of-the-art tool :
(SOTAT)

EXOMISER

○ Exomiser ranks genetic variants according to a combination of criteria :
  • variant frequency
  • predicted pathogenicity
  • known disease associations
  • conservation
  • functional impact
  • phenotypic information

○ Drawbacks :
  • variants of a same gene have the same rank in Exomiser results
  • Exomiser is trained on large databases

# VIOLA's hypothesis

The disease-responsible variant(s) are **patient-specific** and **rare**.
→ unique combination of properties different from the rest of the patient variants.
The putative disease variants for MD are **outliers** of each patient variants' distribution.

# VIOLA's workflow

89'000 variants

89'000 variants

30 scores

89'000 variants

16 scores

**30**'000 variants

16 scores

**3**'000 variants

16 scores

**500** variants

**1** score

500 variants

**2** scores

Variant annotation

Score transformation

Outlier detection

Phenotype integration

VIOLA
score (VS)

VCS

# Creation of the VIOLA combined score (VCS)

**Goal** : Incorporate knowledge of mitochondrial diseases into VIOLA score

VIOLA SCORE

Variant uniqueness

Transcriptomics data

MitoBook

The variant is not an artifact

Gene bearing the variant is already known to be involved in MD

*VCS = 0.5 ( VS + transcriptomics + uniqueness ) + 0.01 ( known gene + artifact )*

# Results on in-house cohort

# VIOLA results on the other patients



input

Patients

Variants in input
Variants selected by VIOLA

# VIOLA results on the other patients



input

output

Patients

Patients

☐ Variants in input
■ Variants selected by VIOLA

# VIOLA results on the other patients

input

output



Patients

Patients

☐ Variants in input
■ Variants selected by VIOLA

VIOLA selects 1% of input variants as potential candidates for MD

# Can VIOLA find the responsible variant for positive patients?

Ranking

1st - 10th

11th - 40th

41th and +

Scores

VIOLA Score ●

VIOLA Combined Score ★

State Of the Art Tool ▲

# Can VIOLA find the responsible variant for positive patients?

Patient A    Patient B

Ranking                    Scores

1st - 10th          VIOLA Score          ●

VIOLA Combined    ★
Score

11th - 40th        State Of the Art      ▲
Tool

41th and +

# Can VIOLA find the responsible variant for positive patients?

Patient A    Patient B

8th

7th

Ranking

1st - 10th

11th - 40th

41th and +

Scores

VIOLA Score    ●

VIOLA Combined
Score    ★

State Of the Art
Tool    ▲

# Can VIOLA find the responsible variant for positive patients?

Patient A     Patient B

1st
8th

1st
7th

Ranking

1st - 10th

11th - 40th

41th and +

Scores

VIOLA Score          ●

VIOLA Combined
Score          ★

State Of the Art
Tool          ▲
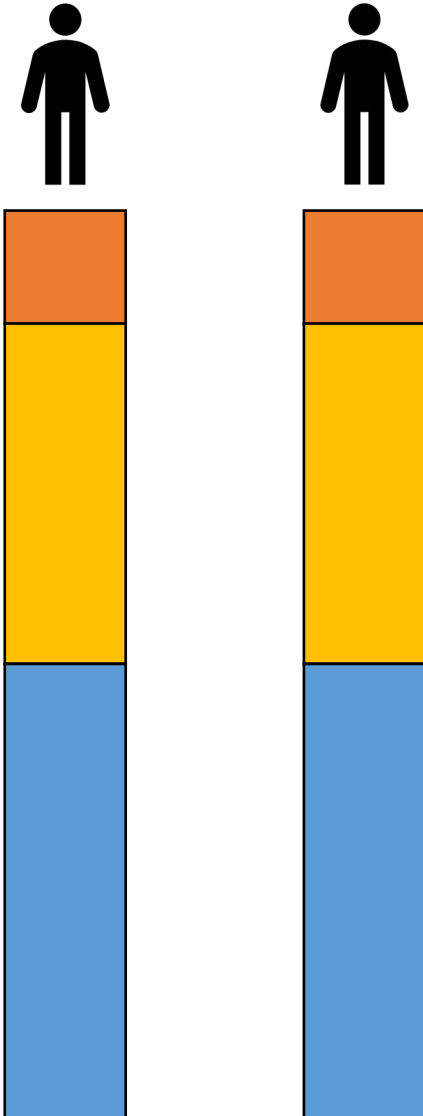
# Can VIOLA find the responsible variant for positive patients?

# Can VIOLA find the responsible variant for positive patients?
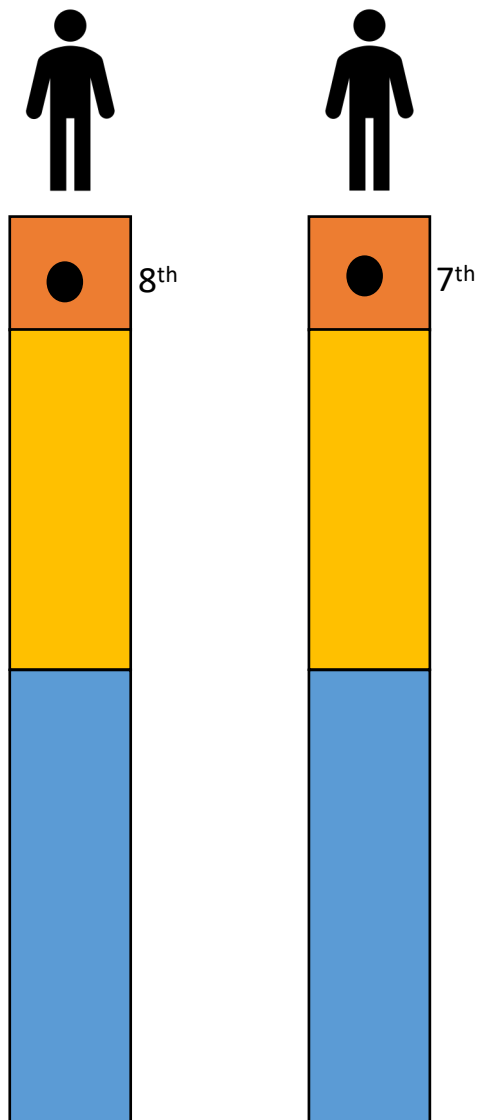
Patient A          Patient B          Patient C

Ranking

Scores

Patient A:
- ★ 1st
- ● 8th
- ▲ 9th

Patient B:
- ★ 1st
- ● 7th
- ▲ 14th

Ranking:
- 1st - 10th
- 11th - 40th
- 41th and +

Scores:
- VIOLA Score ●
- VIOLA Combined Score ★
- State Of the Art Tool ▲

# Can VIOLA find the responsible variant for positive patients?

# Can VIOLA find the responsible variant for positive patients?

# Can VIOLA find the responsible variant for positive patients?
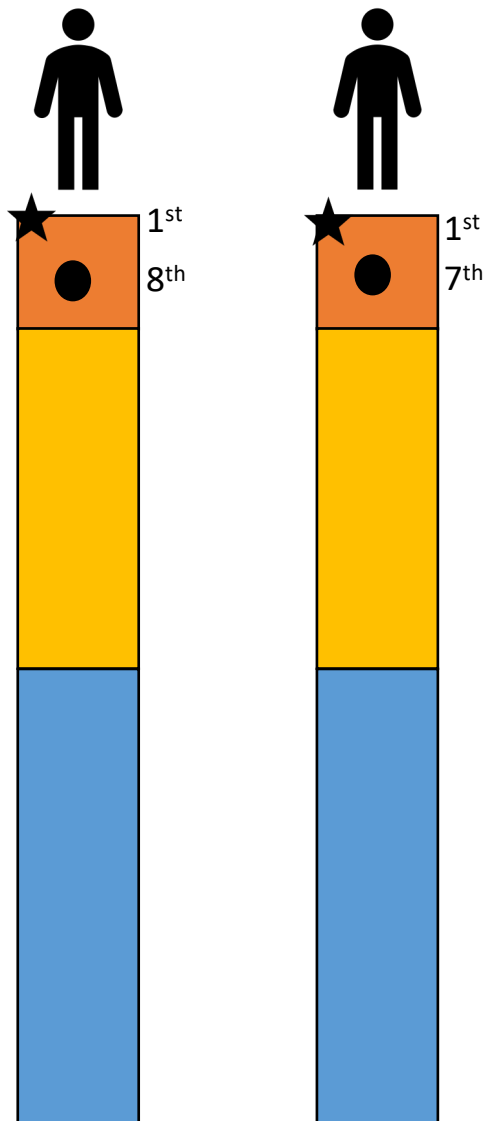
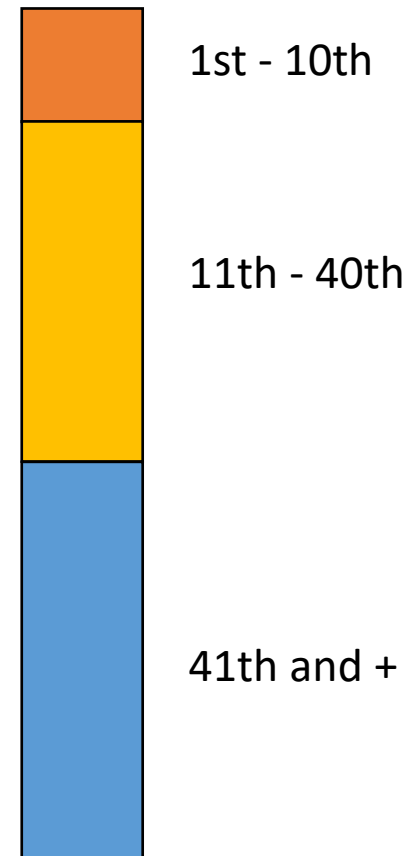Can VIOLA find the responsible variant for positive patients?

# Can VIOLA find the responsible variant for positive patients?



Patient A          Patient B          Patient C          Patient D

1st
8th
9th

1st
7th

14th

17th
18th
19th

19th

34th

51th

Ranking

1st - 10th

11th - 40th

41th and +

Scores

VIOLA Score          ●

VIOLA Combined Score          ★

State Of the Art Tool          ▲

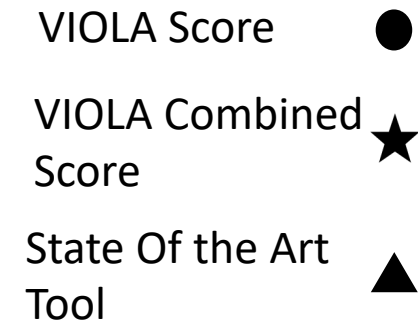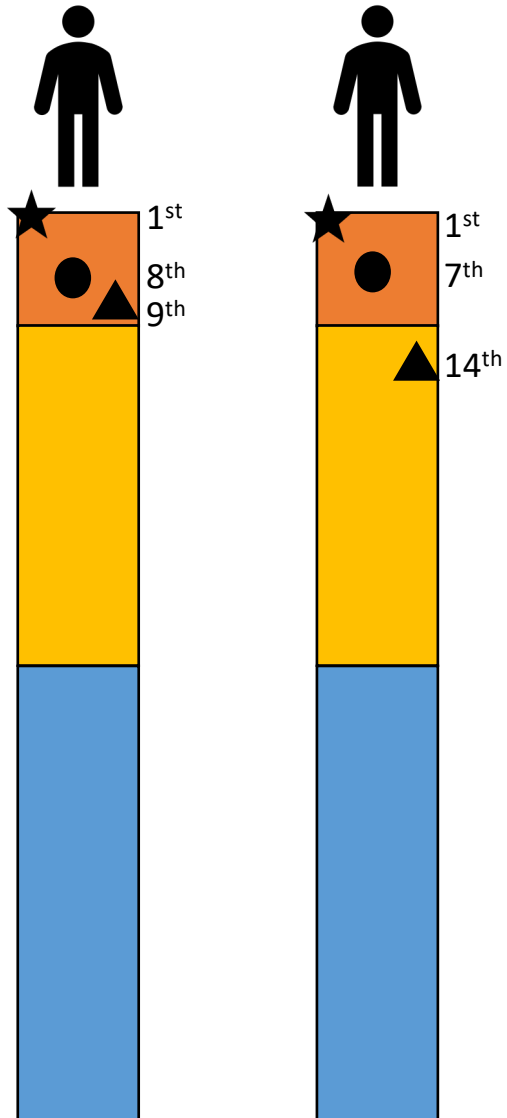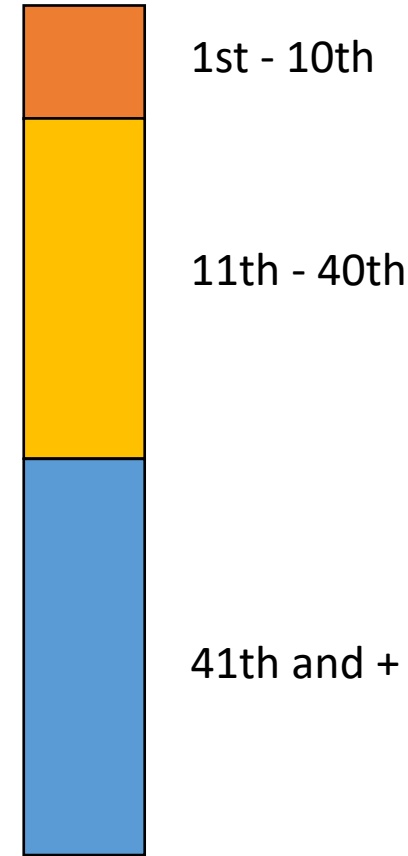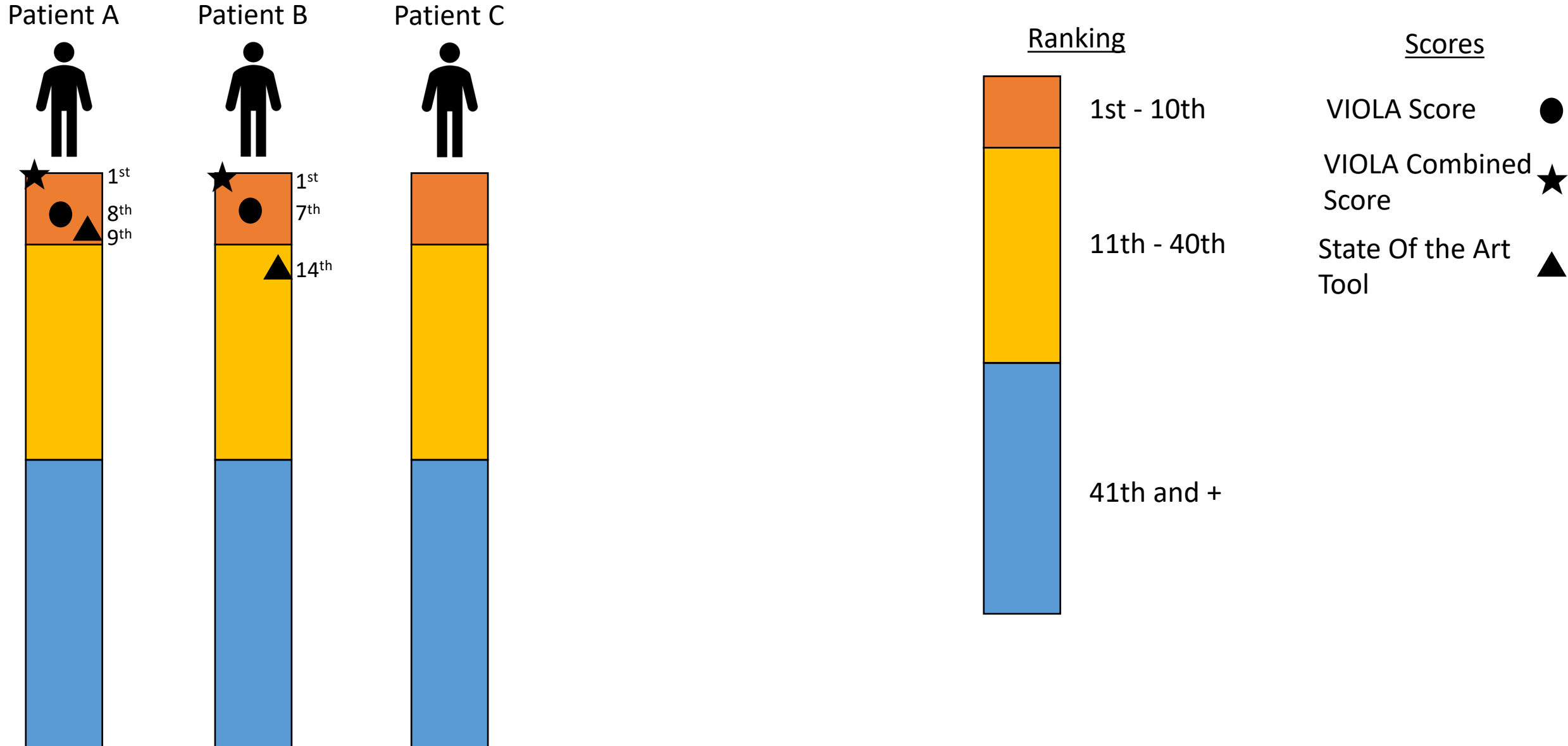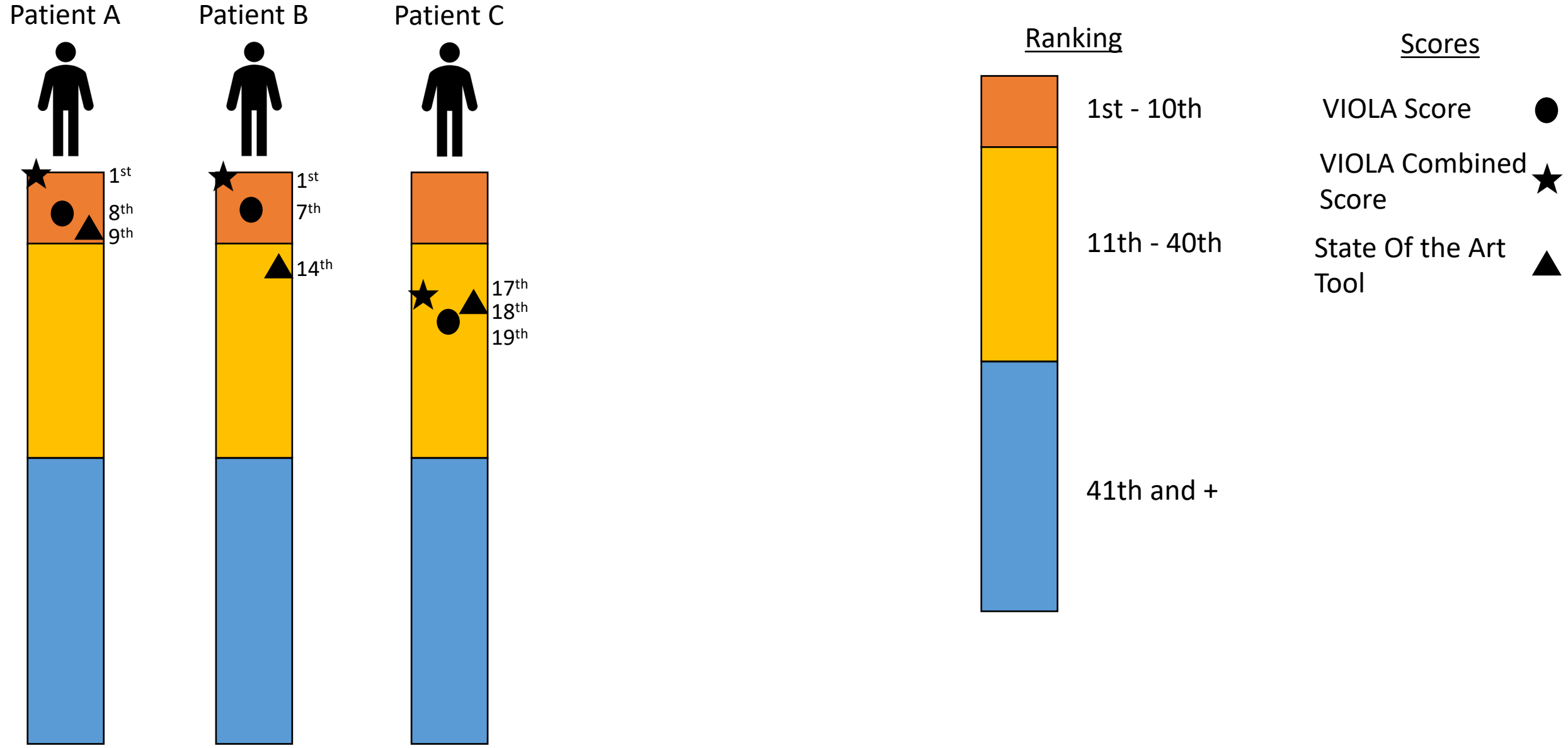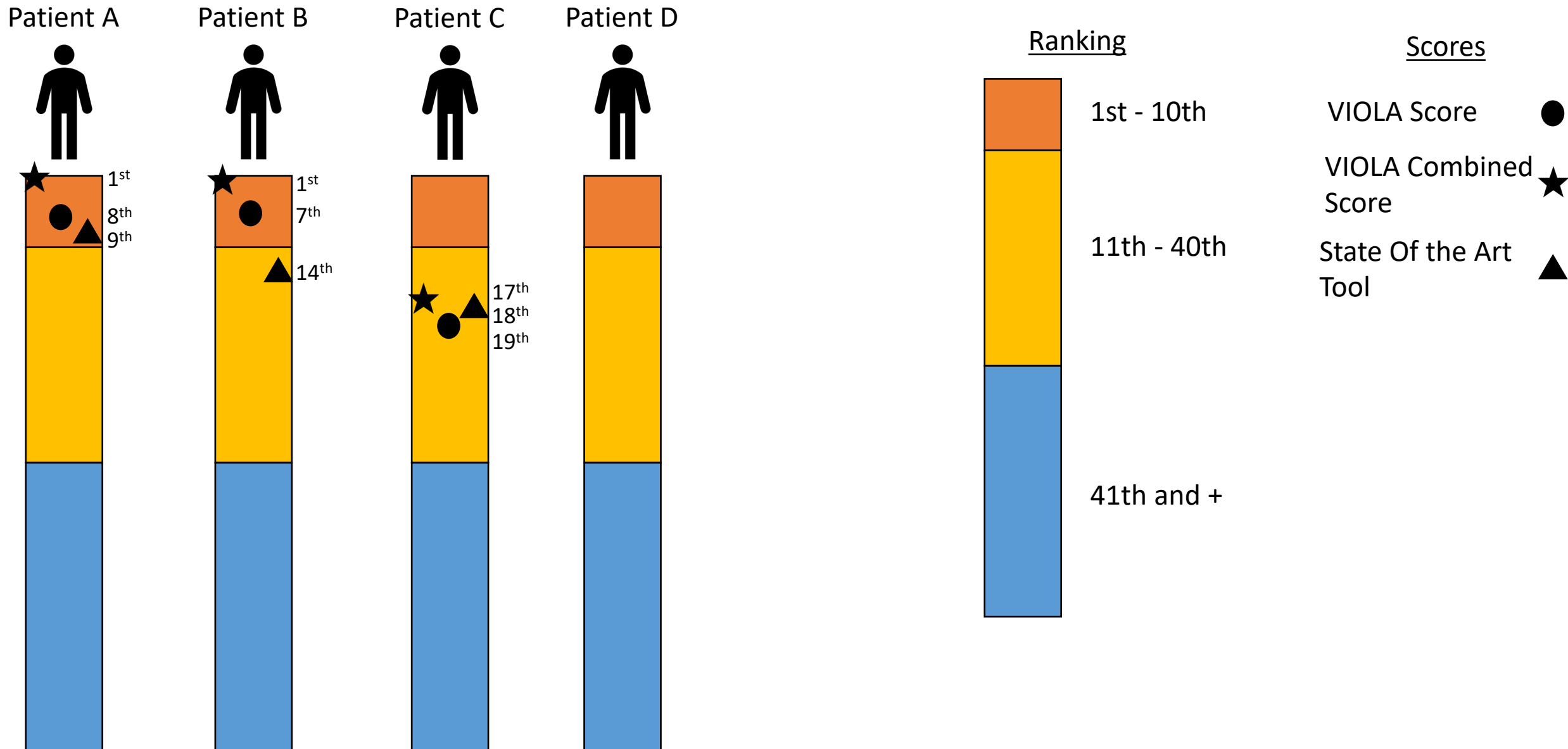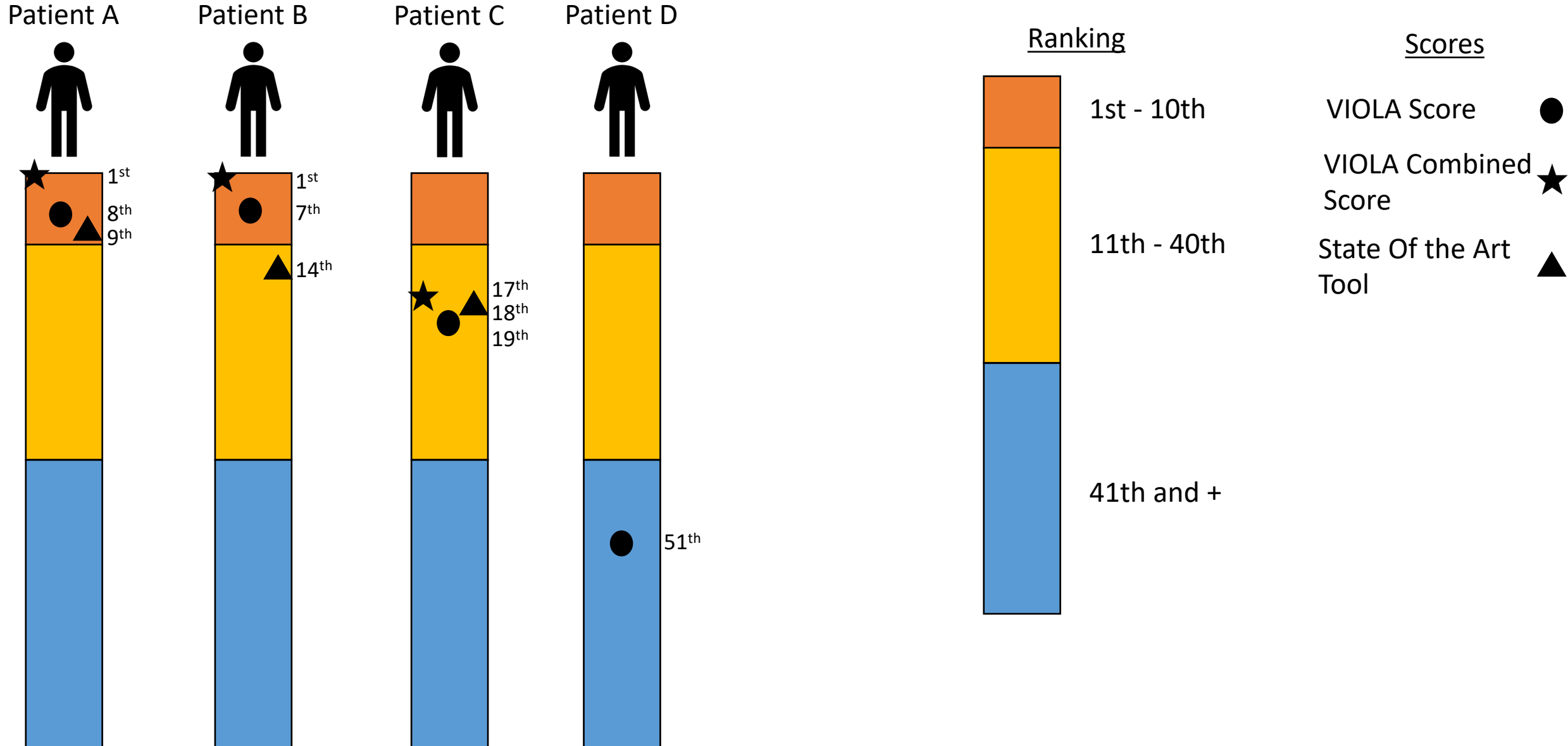# Can VIOLA find the responsible variant for positive patients?



Patient A          Patient B          Patient C          Patient D

Patient A: 1st, 8th, 9th
Patient B: 1st, 7th, 14th
Patient C: 17th, 18th, 19th
Patient D: 19th, 34th, 51th

Ranking

1st - 10th

11th - 40th

41th and +

Scores

VIOLA Score ●

VIOLA Combined Score ★

State Of the Art Tool ▲

- Ranks with VCS are better than those with VS.

# Can VIOLA find the responsible variant for positive patients?



- Ranks with VCS are better than those with VS.
- VIOLA outperformed SOTAT in 3 out of 4 patients.

# VIOLA results

- Enrichment in genes already known to be involved in Mitochondrial disease (MD)



Top variants (upper quartile)

Bottom variants (lower quartile)

Involved in MD
- TRUE
- FALSE

Genes bearing top variants are more enriched in genes already known to be involved in MD than genes bearing bottom variants

# VIOLA results

- Enrichment in MitoCarta genes

Top variants (upper quartile)

Bottom variants (lower quartile)



Genes bearing top variants are more enriched in MitoCarta genes than genes bearing bottom variants

# VIOLA results

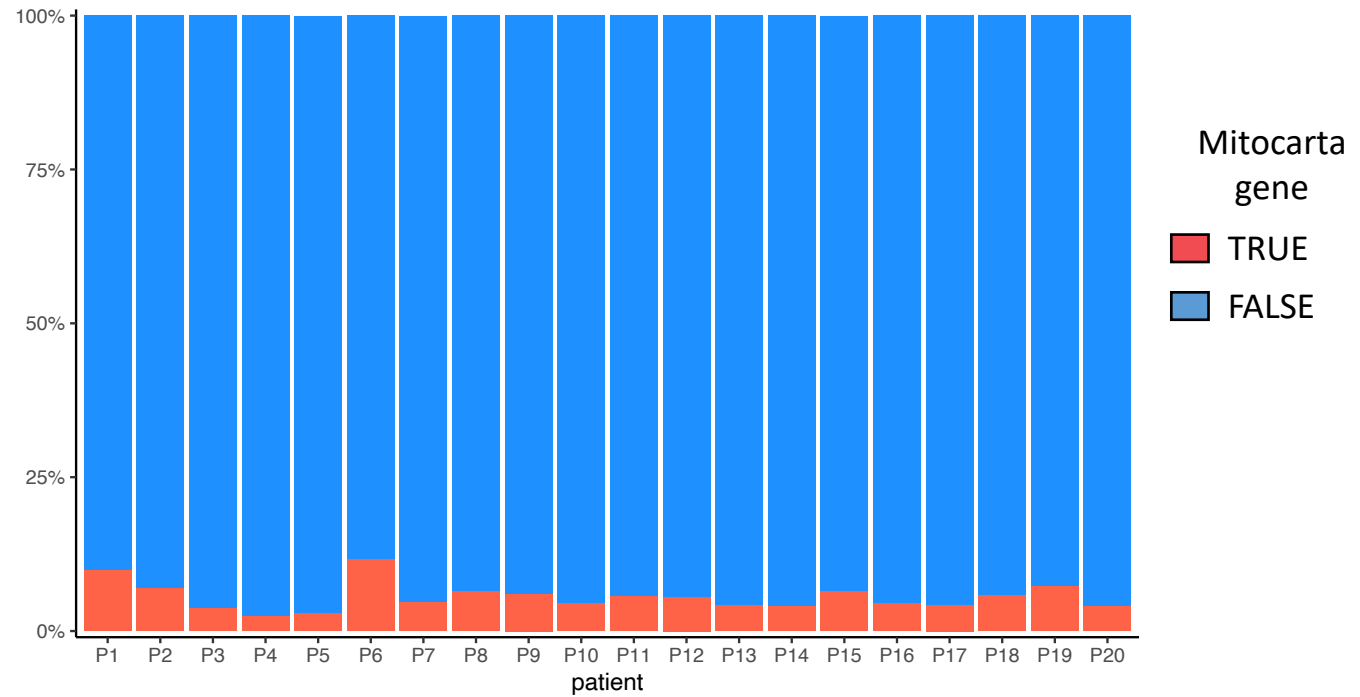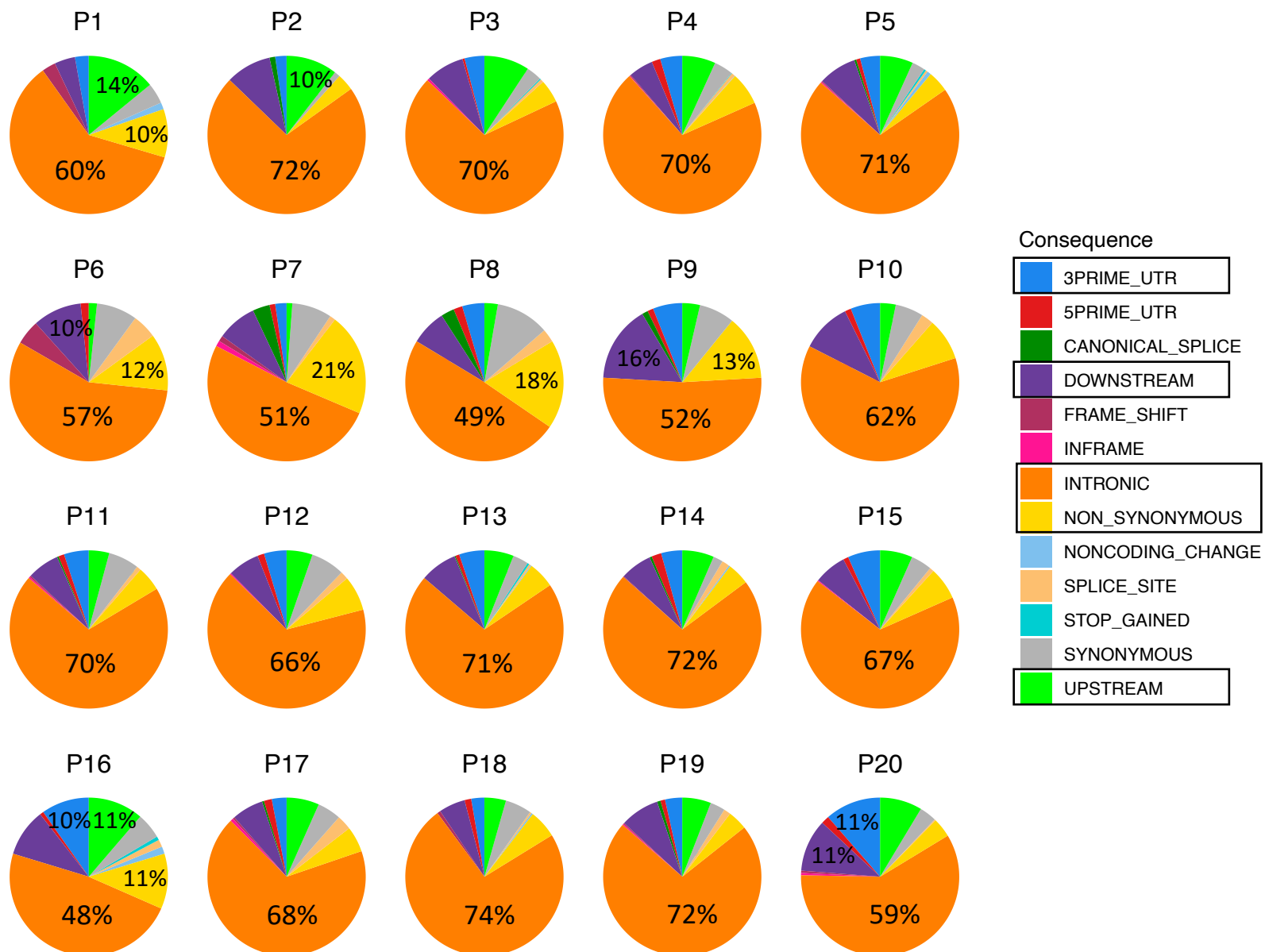- Consequences of **top** variants (upper quartile)

# VIOLA results

- Consequences of **bottom** variants (lower quartile)

# VIOLA find 2 potential candidates for 2 patients of the cohort

- **Case 1**:

**Clinic** ➕

**Variant** 🧬

**Gene** 🧬

- Male baby
- Died shortly after birth
- Dilated cardiomyopathy with elevated lactates

- Intronic SNV in the C1QBP gene
- Heterozygous and rare (not listed in databases)
- Only found for this patient
- Ranked 7th with the VCS

- C1QBP = Encodes a multifunctional protein found mainly in the mitochondrial matrix.
- Listed in MitoCarta and known to be involved in MD
- Similar symtoms for 2 other patients with a variant in C1QBP gene

# VIOLA find 2 potential candidates for 2 patients of the cohort

- **Case 2**:

### Clinic ⊕

- Male adult (24 years old)
- Cardiomyopathy
- Transplanted

### Variant

- Intronic SNV in the LAMA4 gene
- Heterozygous and rare (frequency of 0.000014 in GnomAD)
- Only found for this patient
- Ranked 3[th] with the VCS

### Gene

- LAMA4 = Encodes extracellular matrix glycoprotein
- Known to be involved in cardiomyopathy

# Conclusion of part 2

**Patient-specific** tool, very convenient in a diagnostic context

Model based on the **integration** of **genomics**, **transcriptomics** and **phenomics** data

**Devolpment** of a new model to prioritize genetic variants

For 3 out of 4 patients, VIOLA **outperfoms** Exomiser by ranking the responsible variant in the top 20

VIOLA found 2 potential **candidate** variants for 2 patients in the cohort
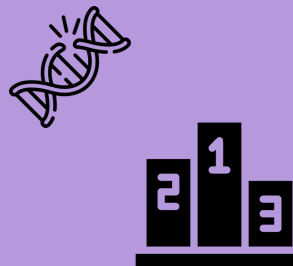
# Take Home Messages

**ABEILLE**

Identification of Aberrant Gene Expression from transcriptomics data for small cohorts

**VIOLA**

Prioritization of genetic variants potentially responsible for MD using latent space

**Personalized medicine**

The diagnosis of MD is complex, important to move as far as possible towards a personalized medicine approach

# Acknowledgements



Thesis supervisors :

Dr. Silvia Bottini

Dr. Sylvie Bannwarth

Genetic team (CHU de Nice):
Dr. Samira Ait-El-Mkadem Saadi

M2P2 team (ISA, Sophia Antipolis)

Master students:
Youssef Boulaimen
Gwendal Le Bideau
Jean Elisée Yao
Jasmine Kaur