# Laboratory Context

- Diagnosis of Breast Cancer

- Breast Cancer treatment response

- Design of signatures

# From Omics to Clinics

Omics

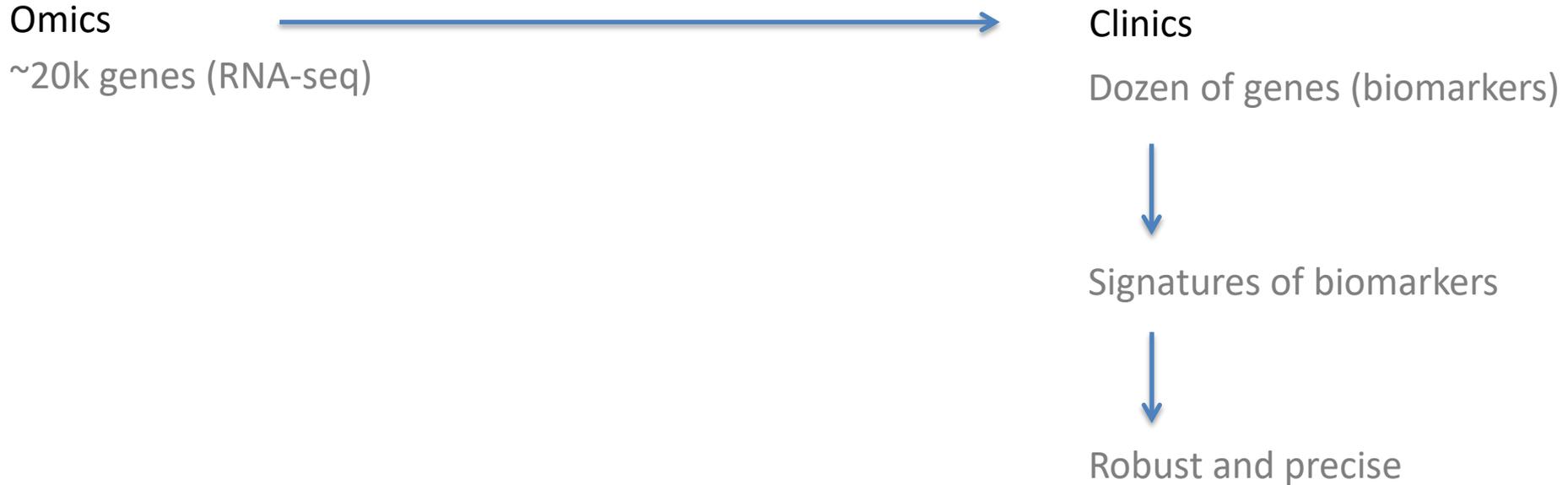~20k genes (RNA-seq)

Clinics

Dozen of genes (biomarkers)

# From Omics to Clinics

Omics

~20k genes (RNA-seq)

Clinics

Dozen of genes (biomarkers)

Signatures of biomarkers

Cancer Diagnosis

Cancer Prognosis

# From Omics to Clinics

Omics → Clinics

~20k genes (RNA-seq)

Dozen of genes (biomarkers)

↓

Signatures of biomarkers

↓

Robust and precise

# From Omics to Clinics

Omics

~20k genes (RNA-seq)

Clinics

Dozen of genes (biomarkers)

Deep Learning
SVM
KNN Genetic Algorithm
Unsupervised Clustering
Boosting Neural Networks
**Random Forest**

Signatures of biomarkers

Robust and precise
Models and signatures

# Random Forest

Omics
~20k genes (RNA-seq)

Clinics
Dozen of genes (biomarkers)

Signatures of biomarkers

Robust and precise
Models and signatures



**Image source:** redbubble.com

# Random Forest

Omics
~20k genes (RNA-seq)

Training = modeling

Clinics

Dozen of genes (biomarkers)

Signatures of biomarkers

Robust and precise
Models and signatures



**Image source:** redbubble.com

8

# Random Forest

Omics
~20k genes (RNA-seq)

Clinics
Dozen of genes (biomarkers)

Signatures of biomarkers

Robust and precise
Models and signatures

**I Entered
A Random Forest**

**Now I see
The Future** EncodedShirts

Training = modeling

Prediction

**Image source:** redbubble.com

Omics
~20k genes (RNA-seq)

Clinics
Dozen of genes (biomarkers)

A Random Forest is a set of

**random decision trees**

Signatures of biomarkers

Robust and precise
Models and signatures

# A decision helps to stratify the data

Omics
~20k genes (RNA-seq)                                                es (biomarkers)

Decision Threshold

< 0.5          ● (circle)          ≥ 0.5

SVM
Deep Learning
KNN Clustering Viscor KNN
Boosting    Random Fc                                          markers

Models and signatures

# A decision helps to stratify the data



samples

Decision Threshold

< 0.5          ≥ 0.5

Normal (0)          Tumor (1)

# A decision helps to stratify the data



Set of genes

samples

Labels: Tumor
Normal

Decision Threshold

< 0.5          ≥ 0.5

Normal (0)          Tumor (1)

# A decision helps to stratify the data



Set of genes

samples

Labels: Tumor
Normal

Decision Threshold

< 0.5    A    ≥ 0.5

Normal (0)    Tumor (1)

# A decision helps to stratify the data

**Set of genes**

A B C D E F

On ~2 (seq)

Labels: Tumor
Normal

samples

Decision Threshold

**Ability to discriminate groups of samples ➔ Importance**

< 0.5    A    ≥ 0.5

Normal (0)          Tumor (1)

Random Forest

# A decision helps to stratify the data



**Set of genes**

Labels: Tumor
Normal
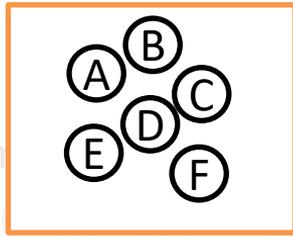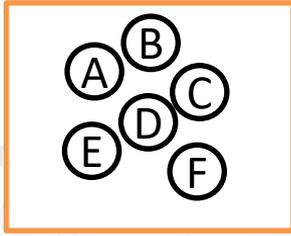
Decision Threshold

< 0.5   A   ≥ 0.5

Normal (0)          Tumor (1)

**Ability to discriminate groups of samples ➜ Importance**

discriminative power ↗ ➜ importance ↗

# A decision helps to stratify the data



**Set of genes**

**Ability to discriminate groups of samples ➔ Importance**

discriminative power ↗ ➔ importance ↗

Decision Threshold

< 0.5    (A)    ≥ 0.5

Normal (0)    Tumor (1)

samples

Labels: Tumor    Normal

Learning (training) process

# A decision helps to stratify the data



**Set of genes**

**Ability to discriminate groups of samples ➔ Importance**

discriminative power ↗ ➔ importance ↗

samples

Labels: Tumor
Normal

Clinics

Dozen of genes (biomarkers)

Decision Threshold

< 0.5    A    ≥ 0.5

Learning (training) process

Normal (0)    Tumor (1)

Signatures of biomarkers

Robust and precise
Models and signatures

1 tree = Multiple decisions

# Decision tree (real example)



Healthy: 0
Tumor:   1

1 tree = Multiple decisions

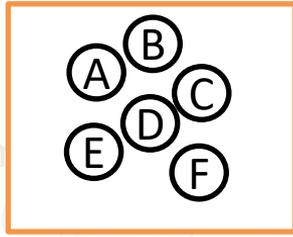# Random Forest is a set of random trees



$$
\begin{array}{c}
\quad miRNA1 \quad miRNA2 \quad \ldots \quad miRNA_n \\
\begin{array}{c}
sample_1 \\
sample_2 \\
. \\
. \\
. \\
sample_m
\end{array}
\left[
\begin{array}{cccc}
x_{11} & x_{12} & \ldots & x_{1n} \\
x_{21} & x_{22} & \ldots & x_{2n} \\
. & . & & . \\
. & . & \ldots & . \\
. & . & & . \\
x_{m1} & x_{m2} & \ldots & x_{mn}
\end{array}
\right]
\end{array}
\quad \text{and} \quad
\left[
\begin{array}{c}
Group \\
Healthy \\
Tumor \\
. \\
. \\
. \\
Healthy
\end{array}
\right]
$$

Omics
~20k genes (RI

**Learning set** / Validation set

Clinics
Dozen of genes (biomarkers)

Training set / OOB set

Random Forest

Signatures of biomarkers

Models and signatures

**1 Forest = Multiple trees**

21

# Random Forest is a set of random trees



$$\begin{bmatrix} & miRNA1 & miRNA2 & ... & miRNA_n \\ sample_1 & x_{11} & x_{12} & ... & x_{1n} \\ sample_2 & x_{21} & x_{22} & ... & x_{2n} \\ & . & . & ... & . \\ & . & . & ... & . \\ & . & . & ... & . \\ sample_m & x_{m1} & x_{m2} & ... & x_{mn} \end{bmatrix} \text{ and } \begin{bmatrix} Group \\ Healthy \\ Tumor \\ . \\ . \\ . \\ Healthy \end{bmatrix}$$

Omics
~20k genes (RI

Clinics
Dozen of genes (biomarkers)

**Learning set** / Validation set

Training set / OOB set       Training set / OOB set       Training set / OOB set

Signatures of biomarkers

Random Forest

Robust and precise
Models and signatures

1 Forest = Multiple trees

22

# Prediction procedure

# Prediction procedure



Learning set / Validation set

Omics
~20k genes (RNA-seq)

Clinics
Dozen of genes (biomarkers)

Unknown Sample

Healthy          Tumor          Healthy

Signatures of biomarkers

Robust and precise
Models and signatures

**Multiple Trees prediction procedure**

# Prediction procedure



Learning set / Validation set

Omics
~20k genes (RNA-seq)

Clinics
Dozen of genes (biomarkers)

Unknown Sample

Healthy          Tumor          Healthy

Majority voting = "Healthy"

**1 Forest = 1 Model**

# From Omics to Clinics

Omics

~20k genes (RNA-seq)

Clinics

Dozen of genes (biomarkers)

Deep Learning
SVM
KNN Genetic Algorithm
Unsupervised Clustering
Boosting Neural Networks

**Random Forest**

**R** implementation**s** of
Random Forest algorithm
proposed by *[Breiman et al. 2001]*

Signatures of biomarkers

Robust and precise
Models and signatures

# From Omics to Clinics

Omics

~20k genes (RNA-seq)

Clinics

Dozen of genes (biomarkers)

Signatures of biomarkers

**Random Forest**

**R** implementation**s** of
Random Forest algorithm
proposed by *[Breiman et al. 2001]*

Robust and precise
Models and signatures

# Objectives

**Toward a robust RF method for the Biological question asked**

Which method is suitable for which dataset (platform/technology) ?

# Objectives

- Empirical comparison of random forest based methods

- Differences/Similarities of RF methods ➔ groups of methods

- Designing a high stability score to rank RF methods

**Toward a robust RF method for the Biological question asked**

Which method is suitable for which dataset (platform/technology) ?

# Materials and Methods

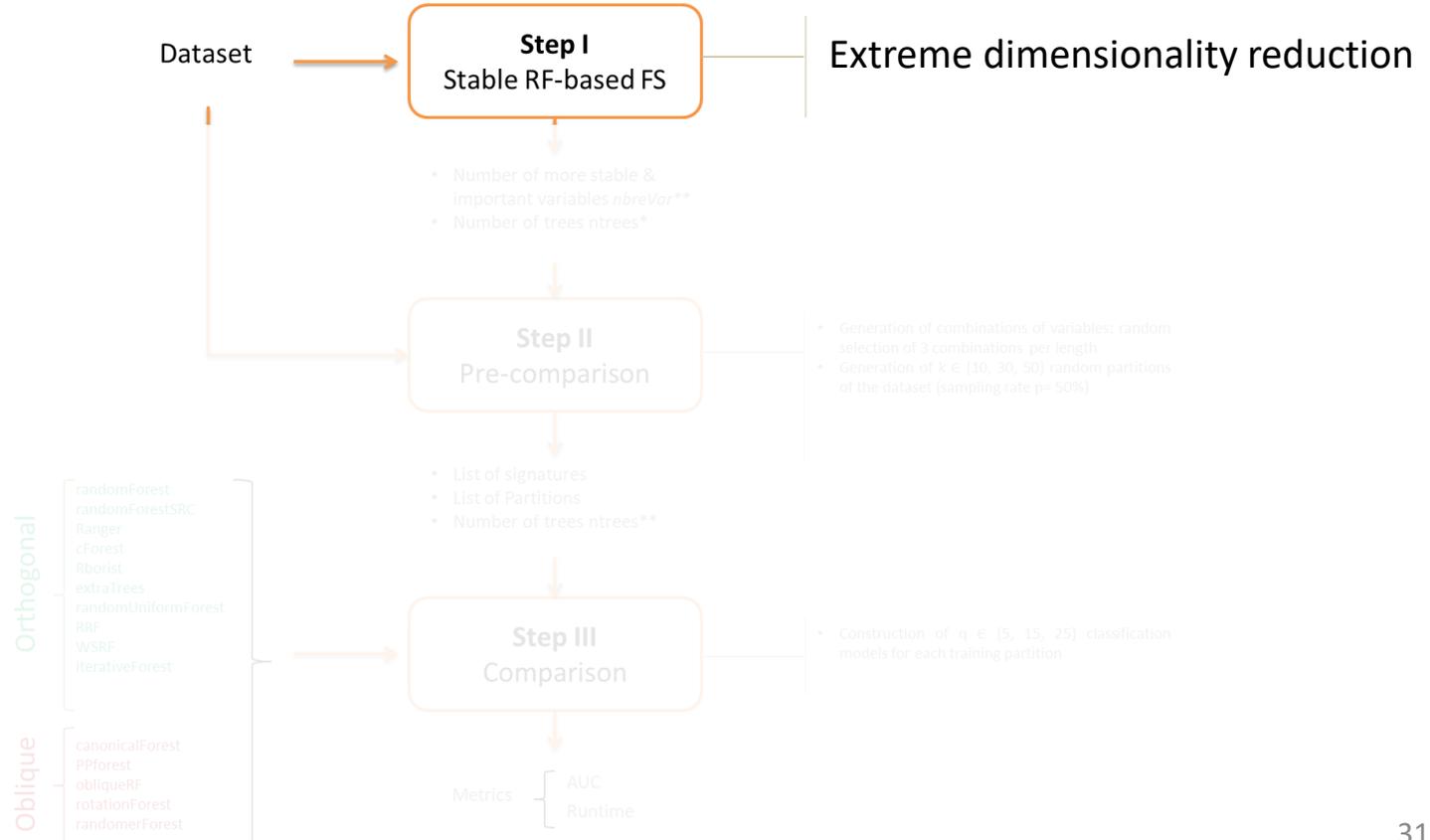- **Datasets (Perfectly balanced)**

  - TCGA-BRCA (RPKM): 182 samples x 9560 genes

  - TCGA-LUSC (RPKM): 96 samples x 9262 genes

  - TCGA-THCA (RPKM): 98 samples x 9353 genes

- **Main classification question**

  The difference between paired **Tumor / Normal** samples will be used as a **strong classification** parameter, allowing for **strong** modeling only
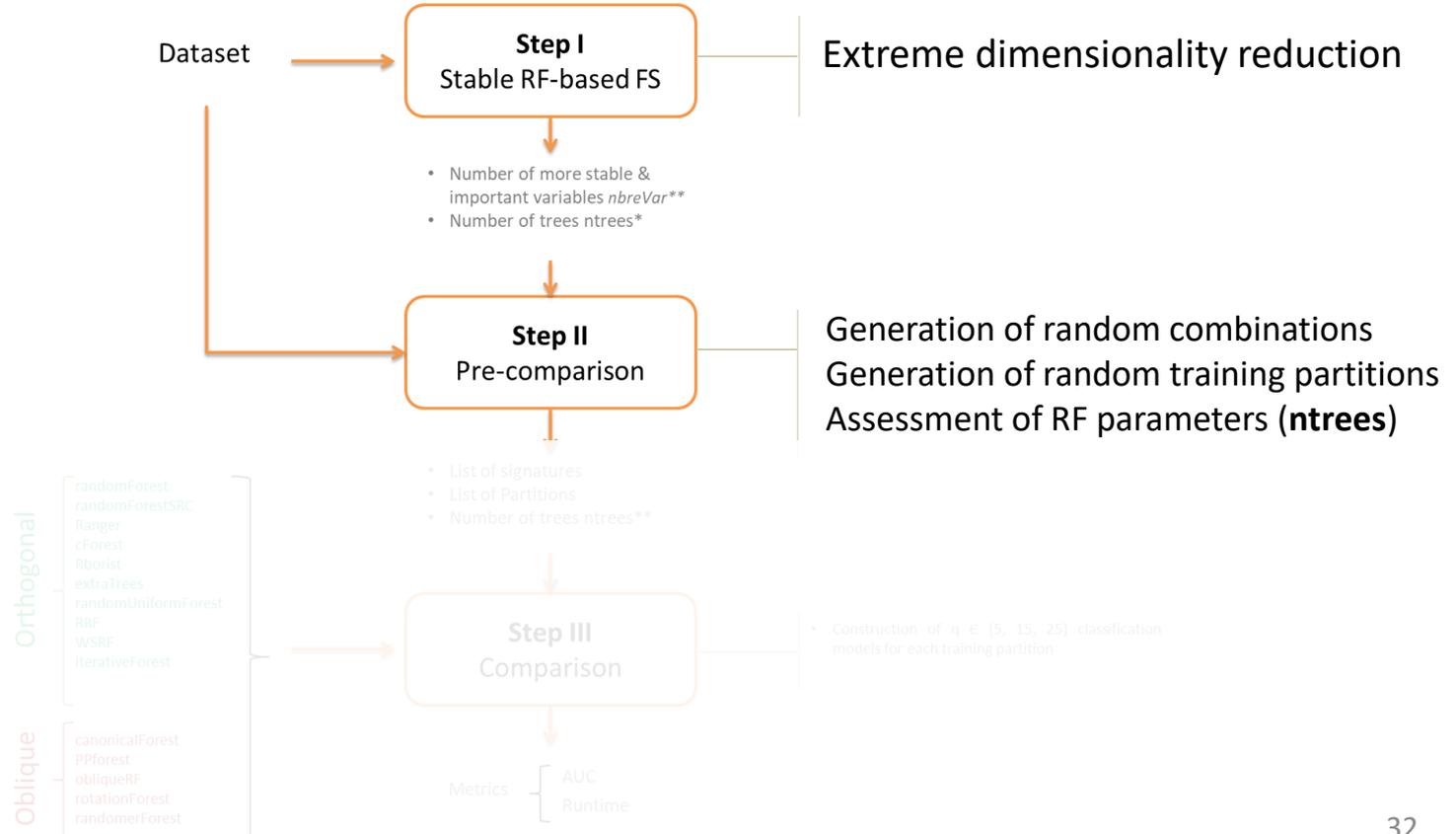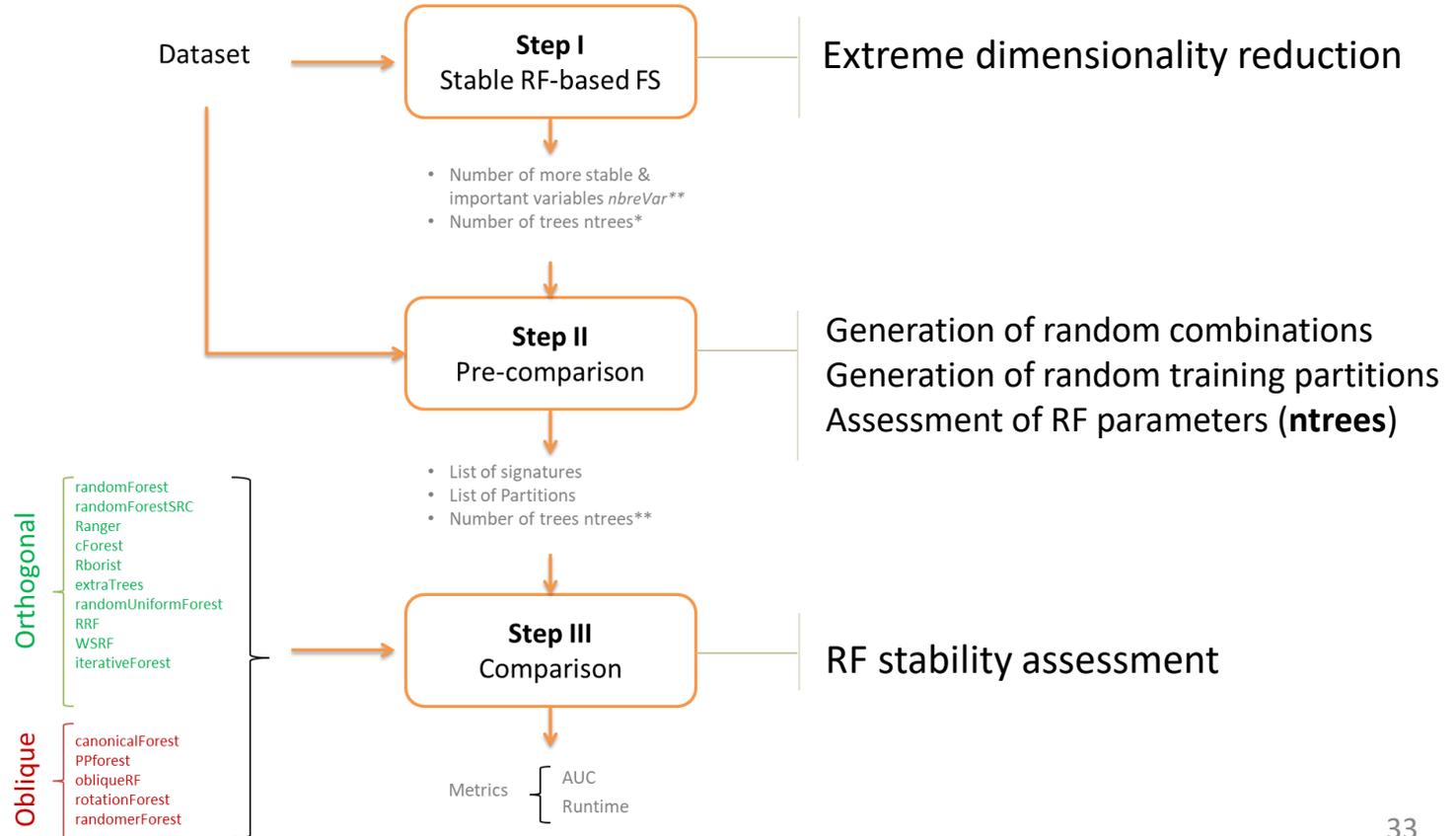
# Overview of the method



Dataset → **Step I** Stable RF-based FS — Extreme dimensionality reduction

- Number of more stable & important variables *nbreVar*\*\*
- Number of trees ntrees\*

**Step II** Pre-comparison

- Generation of combinations of variables: random selection of 3 combinations per length
- Generation of k ∈ {10, 30, 50} random partitions of the dataset (sampling rate p= 50%)

- List of signatures
- List of Partitions
- Number of trees ntrees\*\*

Orthogonal
randomForest
randomForestSRC
Ranger
cForest
Rborist
extraTrees
randomUniformForest
RRF
WSRF
iterativeForest

**Step III** Comparison

- Construction of q ∈ {5, 15, 25} classification models for each training partition

Oblique
canonicalForest
PPforest
obliqueRF
rotationForest
randomerForest

Metrics
AUC
Runtime

# Overview of the method



Dataset → **Step I** Stable RF-based FS — Extreme dimensionality reduction

- Number of more stable & important variables *nbreVar***
- Number of trees ntrees*

**Step II** Pre-comparison — Generation of random combinations
Generation of random training partitions
Assessment of RF parameters (**ntrees**)

- List of signatures
- List of Partitions
- Number of trees ntrees**

Orthogonal
randomForest
randomForestSRC
Ranger
cForest
Rborist
extraTrees
randomUniformForest
RRF
WSRF
iterativeForest

Oblique
canonicalForest
PPforest
obliqueRF
rotationForest
randomerForest

**Step III** Comparison — Construction of q ∈ {5, 15, 25} classification models for each training partition

Metrics — AUC, Runtime

32

# Overview of the method



Dataset →

**Step I**
Stable RF-based FS

— Extreme dimensionality reduction

- Number of more stable & important variables *nbreVar***
- Number of trees ntrees*

**Step II**
Pre-comparison

— Generation of random combinations
Generation of random training partitions
Assessment of RF parameters (**ntrees**)

- List of signatures
- List of Partitions
- Number of trees ntrees**

**Orthogonal**
randomForest
randomForestSRC
Ranger
cForest
Rborist
extraTrees
randomUniformForest
RRF
WSRF
iterativeForest

**Oblique**
canonicalForest
PPforest
obliqueRF
rotationForest
randomerForest

**Step III**
Comparison

— RF stability assessment

Metrics
AUC
Runtime

33

# Stable Feature Selection

Extreme Dimensionality reduction

# First pass Feature Selection results



TCGA-BRCA dataset

# First pass Feature Selection results



$ntrees^* = 2000$

$nVar^* = 200$

TCGA-BRCA dataset

# First pass Feature Selection results



$ntrees^* = 2000$

$nVar^* = 200$

~9000 to **200** variables (Genes)

TCGA-BRCA dataset

# Second pass Feature Selection results

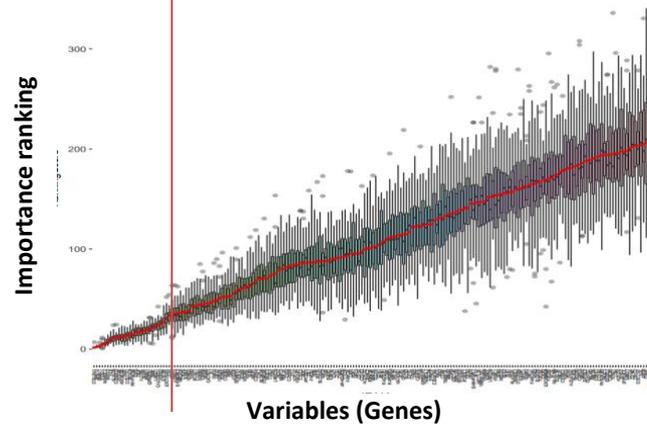

TCGA-BRCA dataset

# Second pass Feature Selection results



$$nVar^{**} = 30$$

TCGA-BRCA dataset

# Second pass Feature Selection results



$$nVar^{**} = 30$$

TCGA-BRCA dataset

# Second pass Feature Selection results



$$nVar^{**} = 30$$

~200 to **30** variables (Genes)

TCGA-BRCA dataset

# Pre-Comparison

Assessment of RF parameters & Generation of random combinations

# Pre-comparison

## I- Generation of random combinations (Cancer signatures)

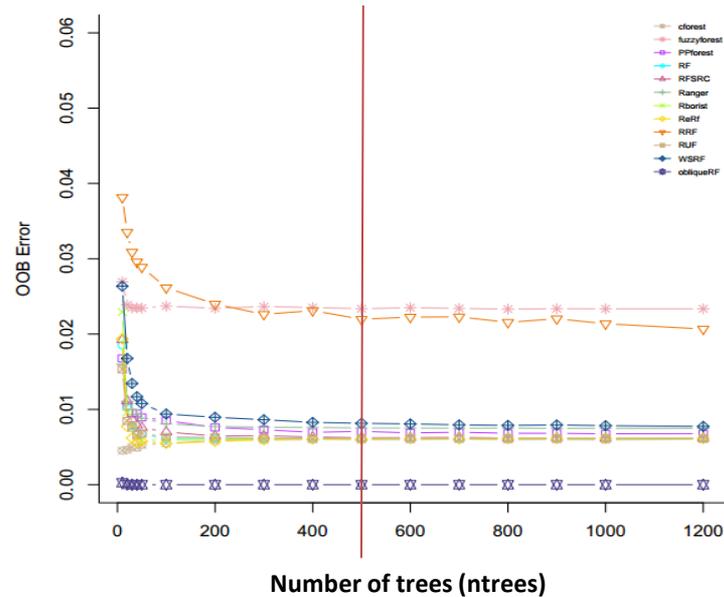- Multiple predictive models using combinations of different lengths

$nVar^{**}$ → $$\left(2^{nVar^{**}} - 1\right) \\ combinations$$ → Random selection of 3 signatures per length

## II- Generation of random training partitions

- 50 random training partitions    training partition = a set of samples used to construct a model

# Pre-comparison

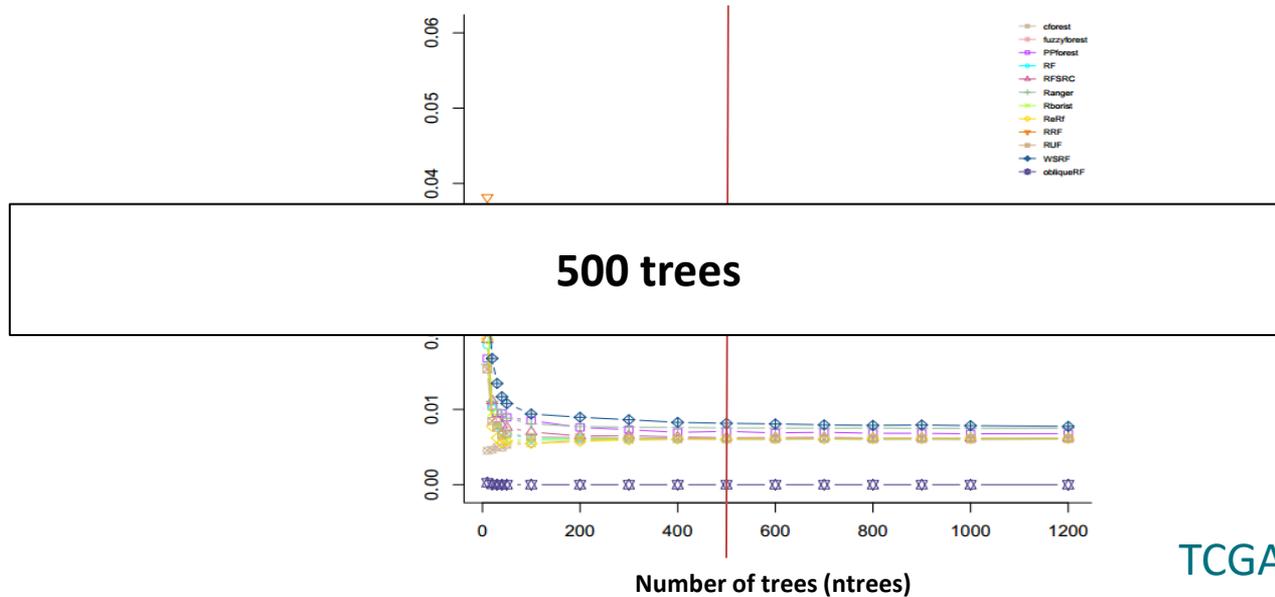III- Tuning the parameter $ntrees$ for each RF method



**Number of trees (ntrees)**

TCGA-BRCA dataset

44

# Pre-comparison

III- Tuning the parameter $ntrees$ for each RF method



**500 trees**

Number of trees (ntrees)

TCGA-BRCA dataset

# Summary of step I + step II

| Dataset | $nVar$ | $nVar^*$ | $nVar^{**}$ | $ntrees^*$ | $ntrees^{**}$ | Nbre combinations |
|---------|--------|----------|-------------|------------|---------------|-------------------|
| **TCGA-BRCA** | 9560 | 200 | 30 | 2000 | 500 | 78 |
| **TCGA-LUSC** | 9262 | 200 | 10 | 2000 | 500 | 21 |
| **TCGA-THCA** | 9353 | 200 | 40 | 2000 | 500 | 108 |

# Comparison

Random Forest stability assessment

# Random Forest Method Comparison

- Comparison of RF methods under perfect conditions

- Using **same** random training partitions

- Assessing the **same** signatures

- On computational cores of **same** characteristics

# Random Forest Method Comparison

- For each signature, we'll focus on:
    - **50 resampling** to build the Training and the Validation set.
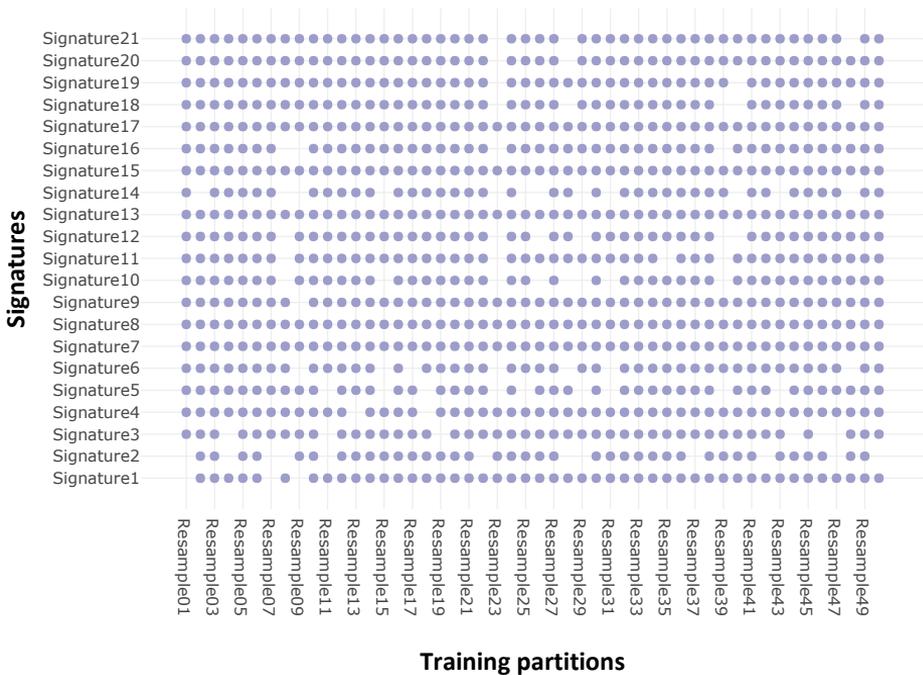    - **25 modeling and validations.**

- Analysis of:
    - **Coefficient of Variation of 1,250 models & AUCs**

**1,250 models & AUCs**

# Random Forest Method Comparison

- For each signature, we'll focus on:
  - **50 resampling** to build the Training and the Validation set.
  - **25 modeling and validations.**

- Analysis of:
  - **Coefficient of Variation of 1,250 models & AUCs**

**1,250 models & AUCs**

**Clinics** ➔ **Hyper Stability : CV == 0**

# Hyper Stability discriminates RF methods



TCGA-LUSC dataset

# Hyper Stability discriminates RF methods



**Best Methods** 🙂

**Worst Methods** ☹️

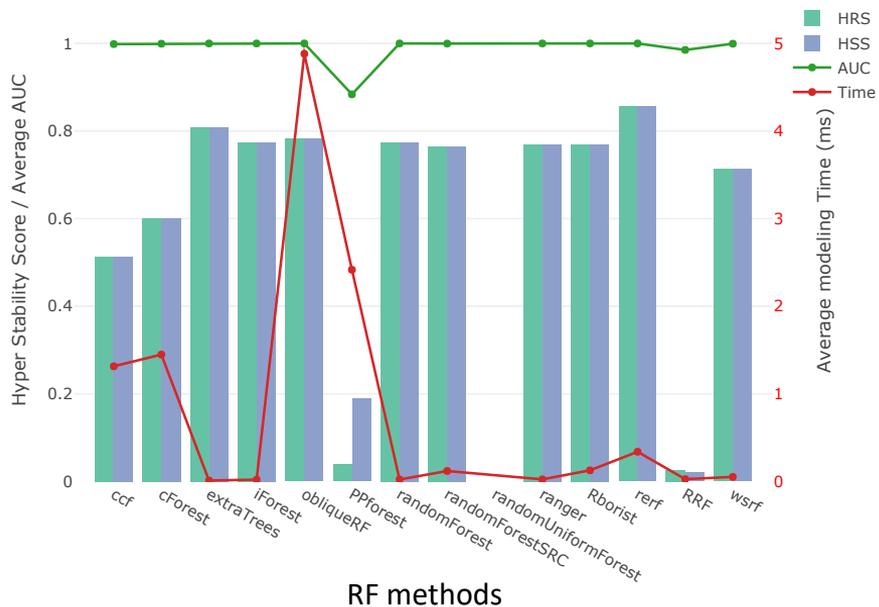# Hyper Stability Score helps finding the best method(s)
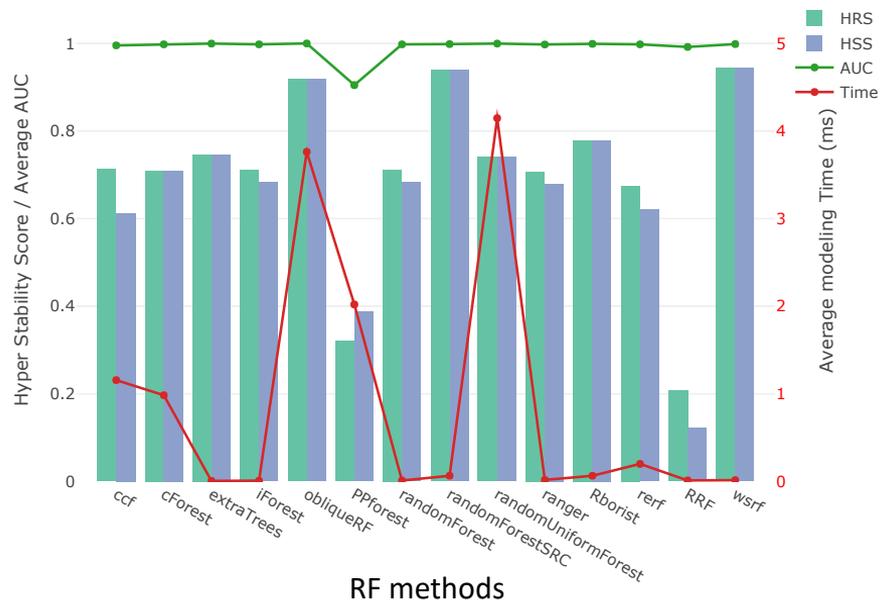


TCGA-BRCA

# Hyper Stability Score is dataset dependent



TCGA-BRCA

TCGA-LUSC

# Conclusions

- The AUC precision is dataset dependent
  - The Methods are dataset dependent.


- Trade-off:
  - AUC precision (hyper-stability)
  - Average AUC value
  - Modeling Time

Classification of classification methods

Towards robust signatures and predictions

# Acknowledgment

## Human Genetics (GIGA)

- Prof. Vincent Bours, MD, PhD
- Christophe Poulet, PhD
- Corinne Fasquelle, Ir

## Oncology (CHU-Liege)

- Prof. Guy Jerusalem, MD
- Dr. Claire Josse
- Aurelie Poncin, MD
- Jerome Thiry

## BIO3 Unit (GIGA)

- Prof. Kristel Van Steen

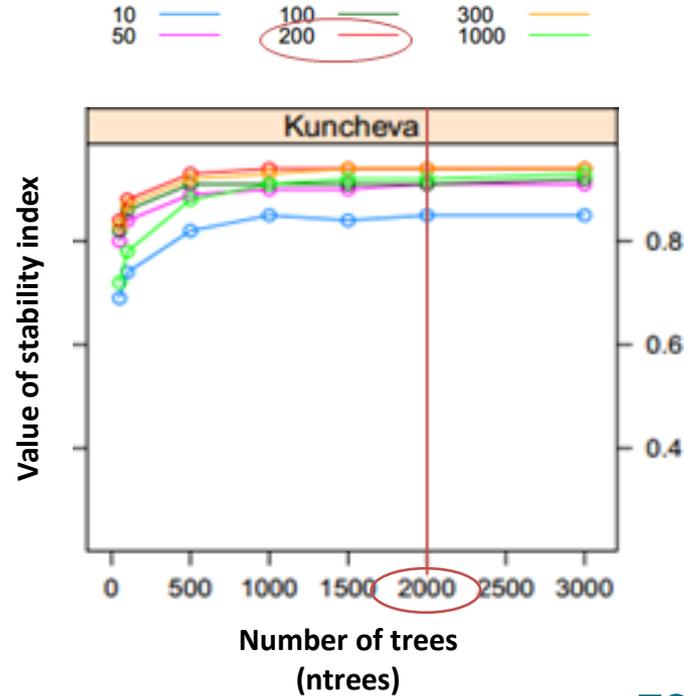## CBIO, Mines ParisTech, Institut Curie

- Chloe-Agathe Azencott, PhD

# First pass Feature Selection results



$ntrees^* = 2000$

$nVar^* = 200$

TCGA-BRCA dataset