



UCO
UNIVERSITÉ
CATHOLIQUE DE L'OUEST

IMA
Institut de Mathématiques
Appliquées



Harmonisation des données

Pierre CHAUVET, Nassib ABDALLAH, Jean-Marie MARION

Institut de Mathématiques Appliquées & Laboratoire LARIS

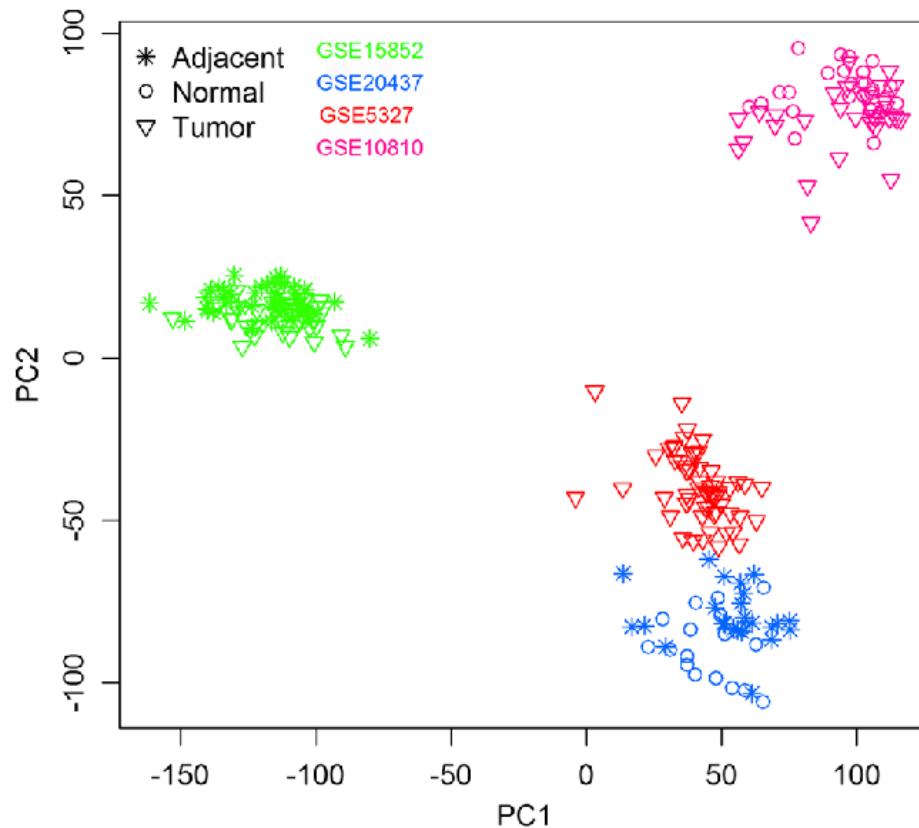
pierre.chauvet@uco.fr

Plan

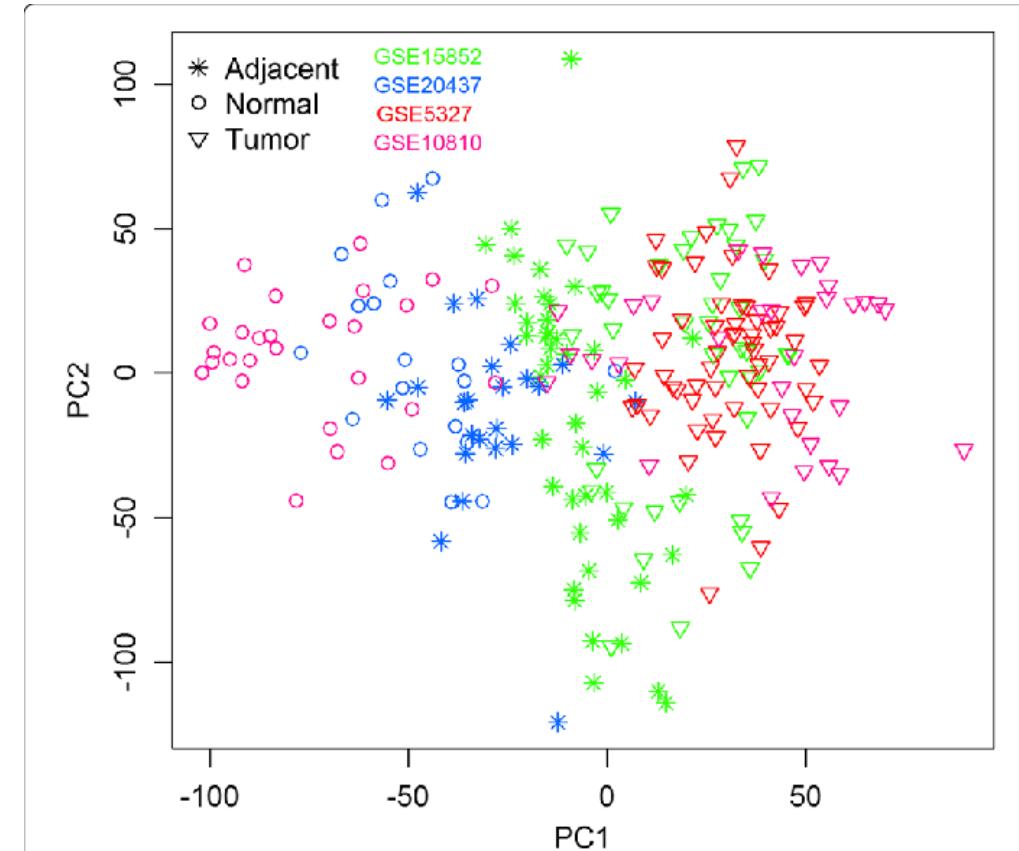
- Le « *Batch Effect* »
 - Problématique, état de l'art
- Ajustement position / échelle
- Le projet HARMONY
 - Premiers travaux : bases BreakHIS et IDC
 - Le challenge Hecktor 2021

Le « *Batch Effect* »

Quel problème ?



The PCA plot before batch effect removal. Three types of samples from 4 different datasets are shown on this figure; different colors indicate different datasets, while different symbols represent sample types (Normal, Tumor, or Adjacent Normal)



The PCA plot after ComBat batch effect removal. The same set of samples as in the previous figure, but after the ComBat batch effect removal procedure has been applied. Color and symbol schemes remain the same

Problématique (I)

- Le « *batch effect* » est un problème soulevé initialement par la biologie moléculaire (en particulier en génomique, transcriptomique, ...)

- Il s'observe en médecine, et est un frein à l'utilisation des données pour faire de la classification / prévision automatique.

Problématique (2)

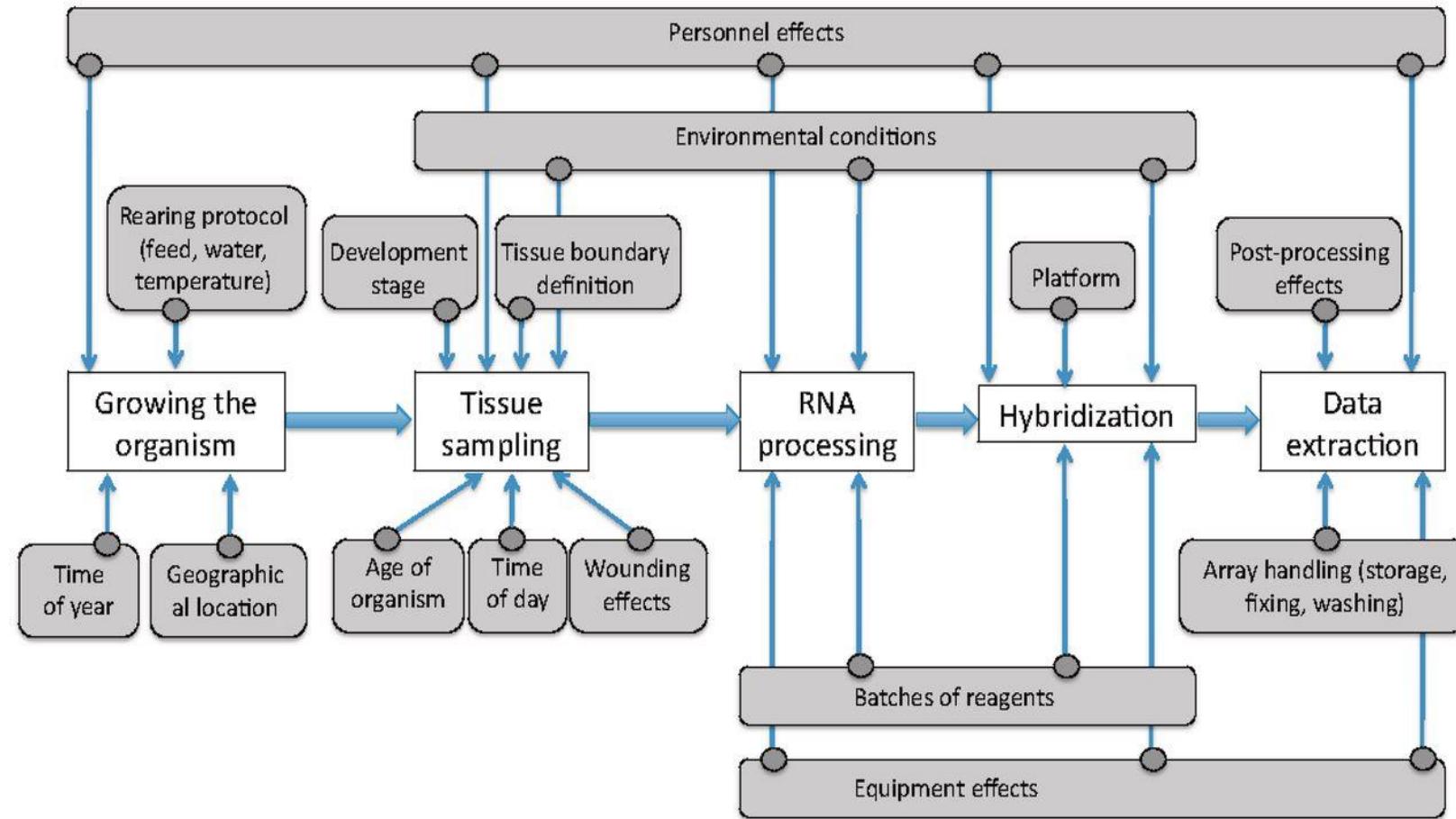
The most well-known source of latent variation in genomic experiments are batch effects—when samples are processed on different days, in different groups or by different people.

Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*. 2012;28(6):882-883. doi:10.1093/bioinformatics/bts034

In practical data analysis, the observations included in a dataset sometimes form distinct groups—denoted as “batches”; for example, measured at different times, under different conditions, by different persons or even in different labs.

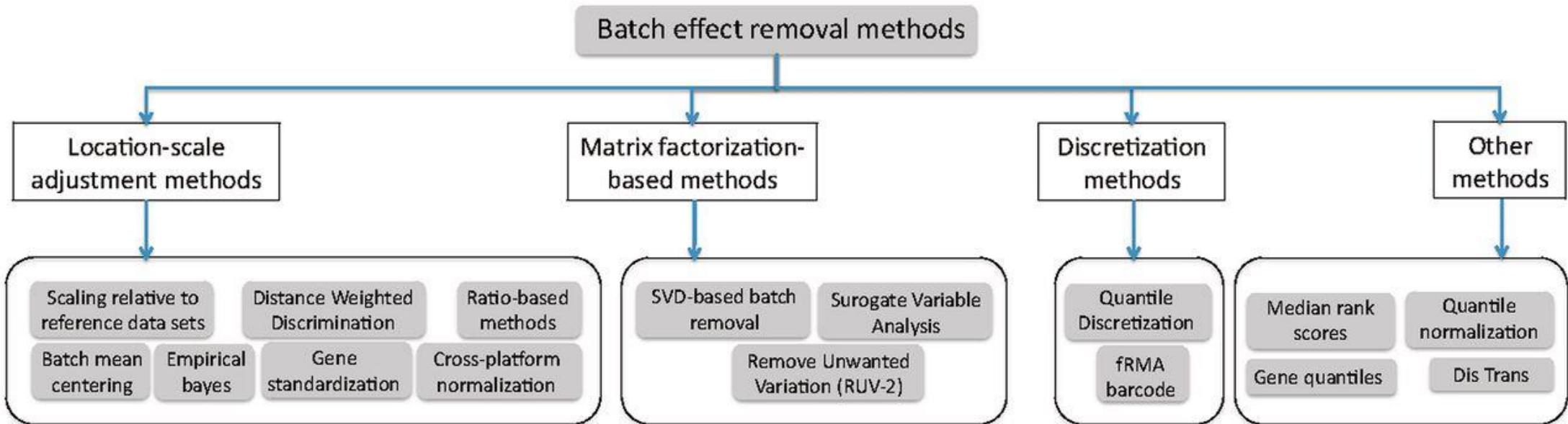
Hornung, R., Boulesteix, AL. & Causeur, D. Combining location-and-scale batch effect adjustment with data cleaning by latent factor adjustment. *BMC Bioinformatics* 17, 27 (2016). <https://doi.org/10.1186/s12859-015-0870-z>

Origines du *batch effect* en génomique



Cosmin Lazar, Stijn Meganck, Jonatan Taminau, David Steenhoff, Alain Coletta, Colin Molter, David Y. Weiss-Solis, Robin Duque, Hugues Bersini, Ann Nowé, *Batch effect removal methods for microarray gene expression data integration: a survey*, *Briefings in Bioinformatics*, Volume 14, Issue 4, July 2013, Pages 469–490, <https://doi.org/10.1093/bib/bbs037>

Taxonomie des méthodes de suppression du *batch effect*



*Cosmin Lazar, Stijn Meganck, Jonatan Taminau, David Steenhoff, Alain Coletta, Colin Molter, David Y. Weiss-Solís, Robin Duque, Hugues Bersini, Ann Nowé, Batch effect removal methods for microarray gene expression data integration: a survey, *Briefings in Bioinformatics*, Volume 14, Issue 4, July 2013, Pages 469–490, <https://doi.org/10.1093/bib/bbs037>*

Quelques références

Méthode Gene-wise standardization (2001)

Li C, Wong WH. Model-based analysis of oligonucleotide arrays: expression index computation and outlier detection. Proc Natl Acad Sci USA 2001;98(1):31–36.

Méthode Empirical Bayes ComBat (2007)

W. Evan Johnson, Cheng Li, Ariel Rabinovic, Adjusting batch effects in microarray expression data using empirical Bayes methods, Biostatistics, Volume 8, Issue 1, January 2007, Pages 118–127, <https://doi.org/10.1093/biostatistics/kxj037>

Méthode BMC (batch mean-centering - 2008)

Sims A, Smethurst G, Hey Y, et al. The removal of multiplicative, systematic bias allows integration of breast cancer gene expression datasets—improving meta-analysis and prediction of prognosis. BMC Medical Genom 2008;1(1):42–56.

Méthodes ratio-A et ratio-G (2010)

Luo J, Schumacher M, Scherer A, et al. A comparison of batch effect removal methods for enhancement of prediction performance using MAQC-II microarray gene expression data. PharmacogenomicsJ 2010;10(4):278–91

Méthode SVA (2012)

Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. Bioinformatics. 2012;28(6):882-883. <https://doi.org/10.1093/bioinformatics/bts034>

Méthode FAbatch (2016)

Hornung, R., Boulesteix, AL. & Causeur, D. Combining location-and-scale batch effect adjustment with data cleaning by latent factor adjustment. BMC Bioinformatics 17, 27 (2016). <https://doi.org/10.1186/s12859-015-0870-z>

Voir survey : Lazar C, Meganck S, Taminau J, Steenhoff D, Coletta A, Molter C, et al., Batch effect removal methods for microarray gene expression data integration: a survey. Brief Bioinformatics. 2013;14(4):469–90.

Ajustement position / échelle

- Pas de *batch* de référence



Centrer et / ou réduire

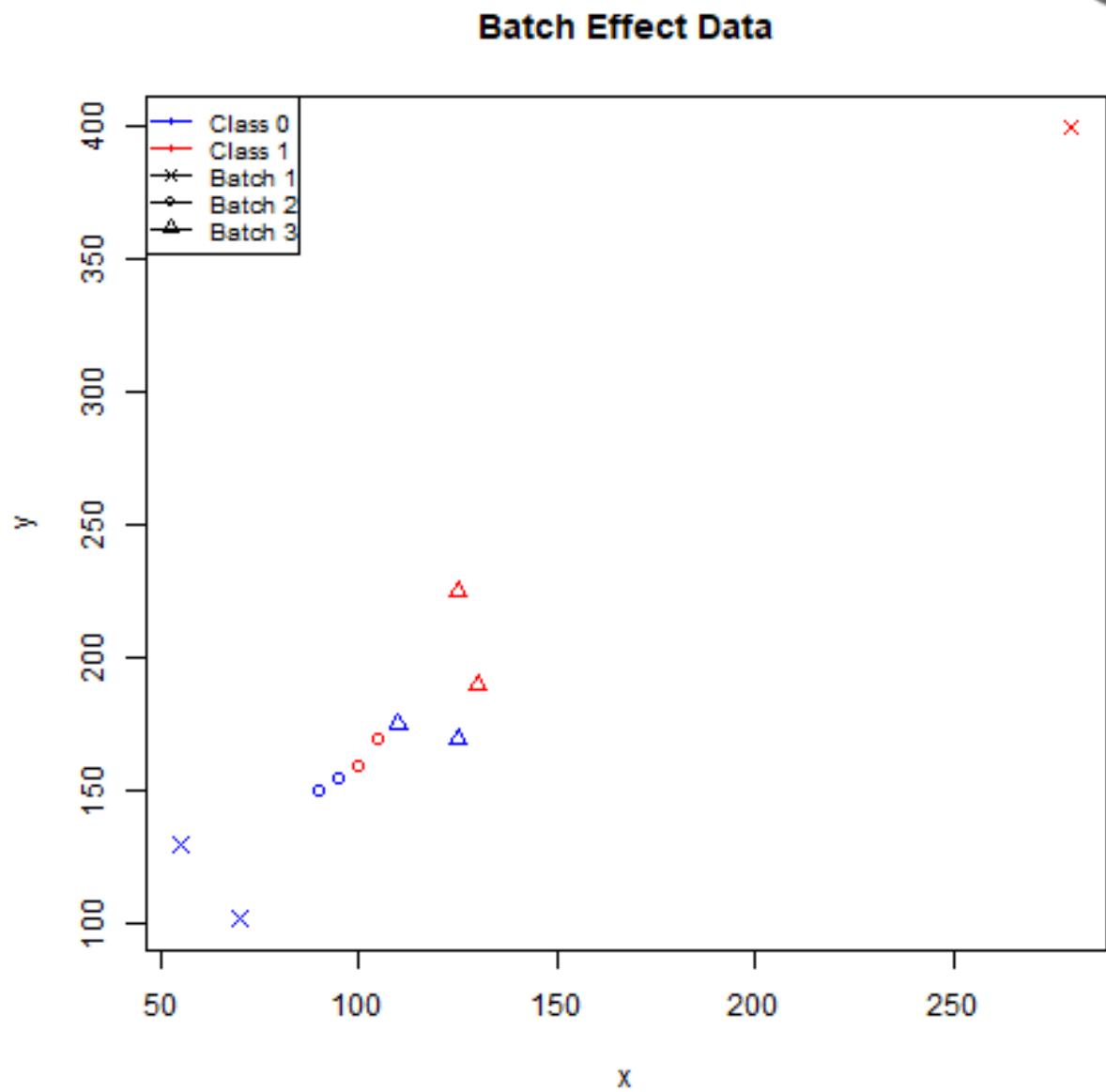
Calcul moyennes ou médianes

Calcul écarts-types ou quartiles

Exemple données

donnees.csv

x	y	class	batch
70,00	102,00	0	1
55,00	130,00	0	1
280,00	400,00	1	1
90,00	150,00	0	2
95,00	155,00	0	2
100,00	160,00	1	2
105,00	170,00	1	2
110,00	175,00	0	3
125,00	170,00	0	3
125,00	225,00	1	3
130,00	190,00	1	3



Centrer-réduire les données

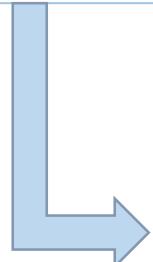
Calculer la moyenne et la variance de l'ensemble des données:

- C = centre (moyenne x, moyenne y) pour toutes les données
- S = écart-type (x, y) pour toutes les données

Transformer les données (x,y):

- Pour chaque échantillon j :

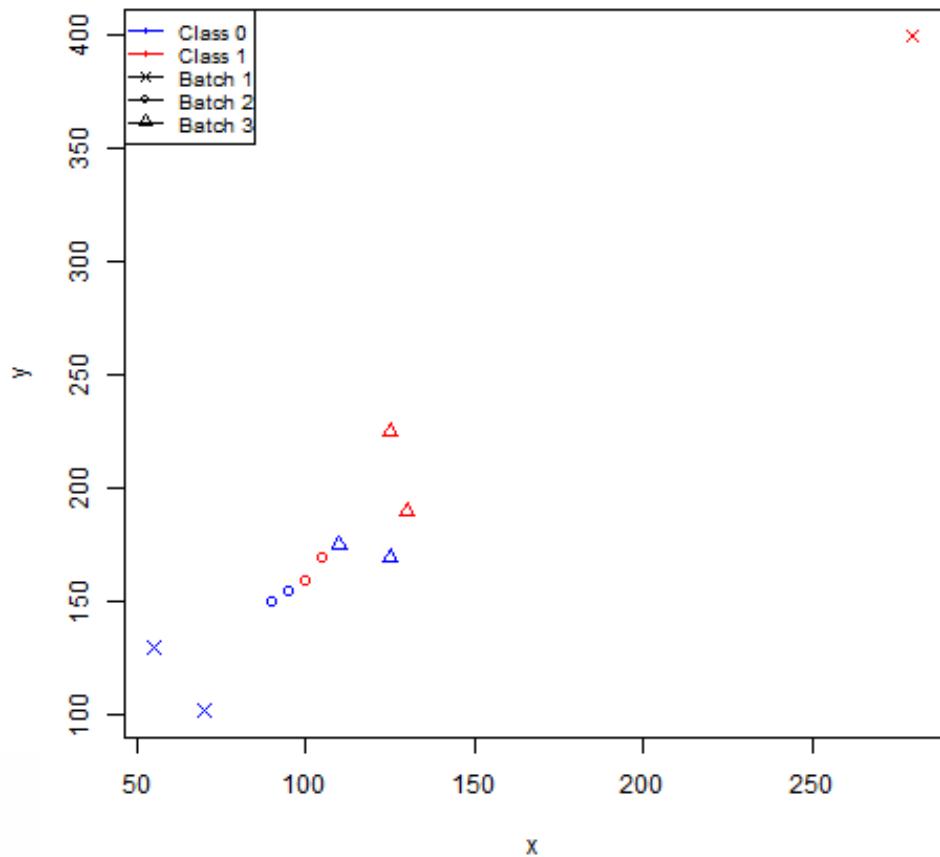
$$\begin{aligned}x_j &<- (x_j - C.x) / S.x \\y_j &<- (y_j - C.y) / S.y\end{aligned}$$



Chaque colonne (x,y) est de moyenne nulle et de variance égale à 1

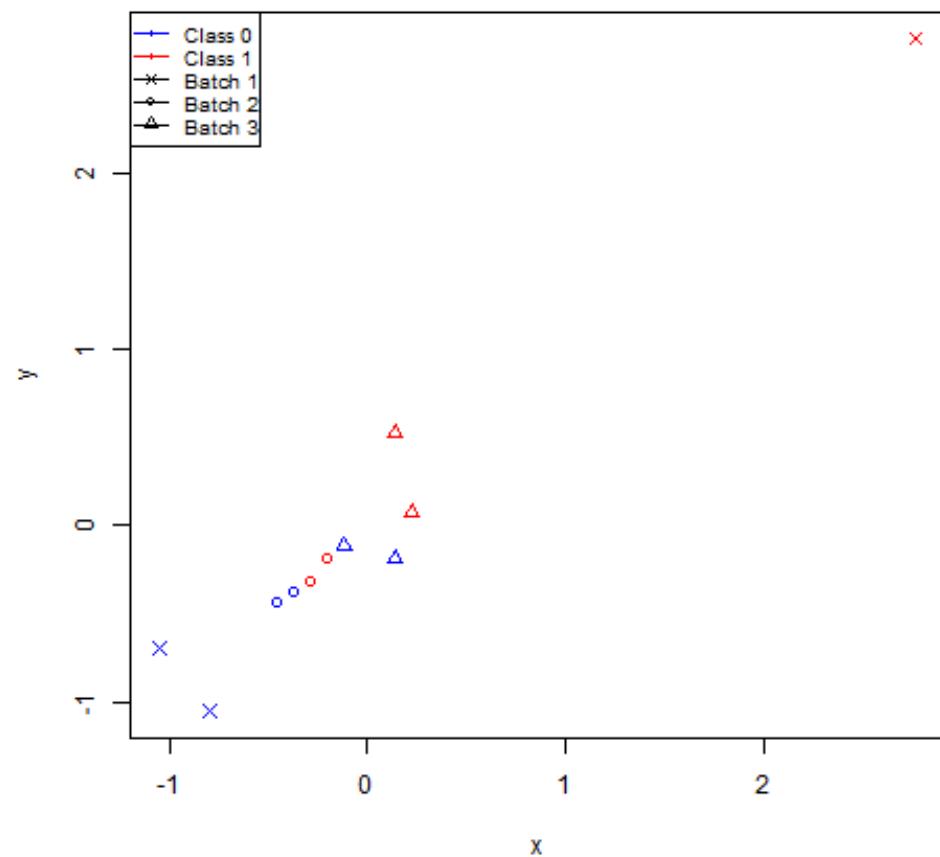
Centrer-réduire les données

Batch Effect Data



Avant...

Basic Scale



Après

Recentrage sur la moyenne / médiane globale

Calculer les centres moyens :

- $C = \text{centre}(\text{moyenne } x, \text{moyenne } y)$ pour toutes les données
- $M_i = \text{centre}(\text{moy } x, \text{moy } y)$ sur le batch i

Transformer les données (x,y) :

- Pour chaque batch i , translater chaque point selon $\overrightarrow{M_i C}$

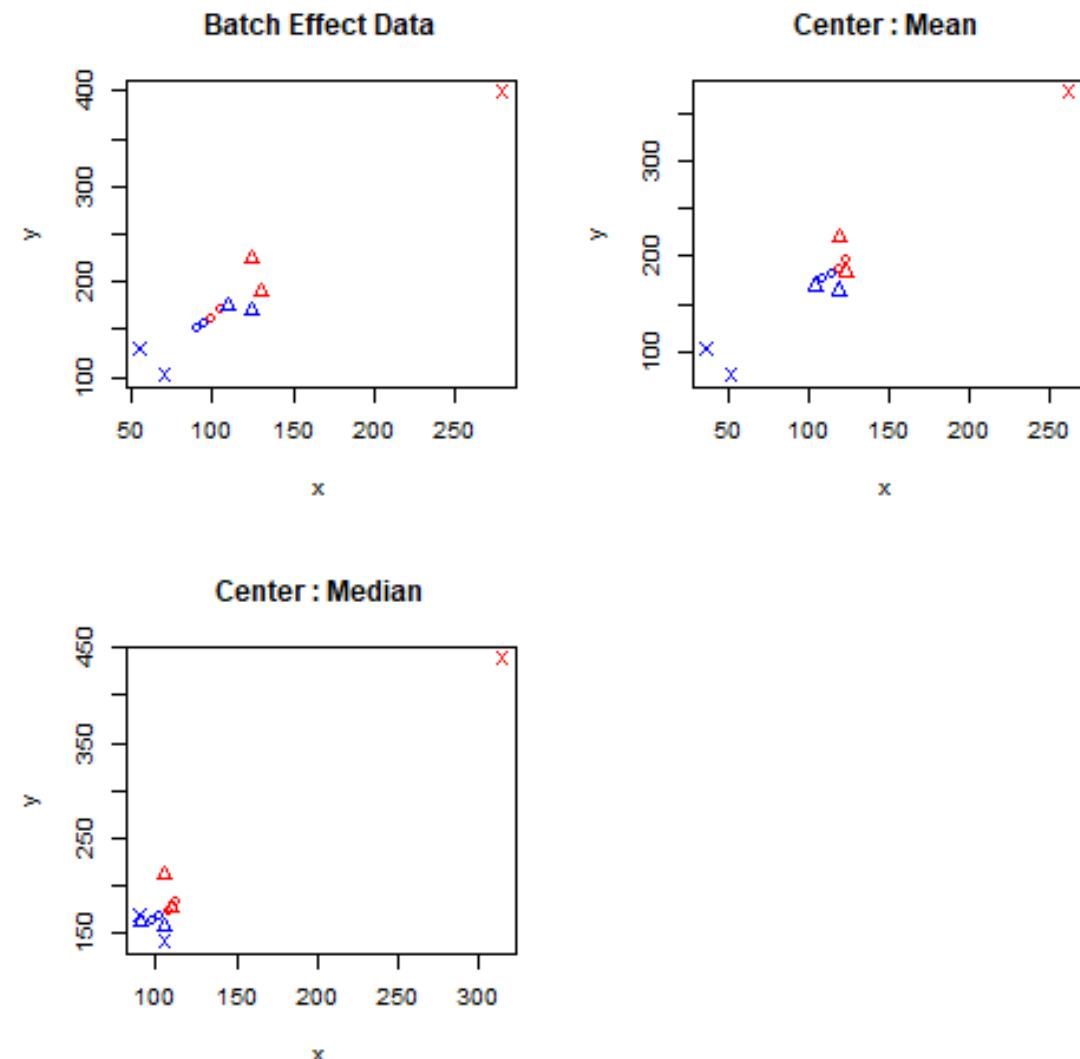
Calculer les centres médians :

- $C = \text{centre}(\text{médiane } x, \text{médiane } y)$ pour toutes les données
- $M_i = \text{centre}(\text{méd } x, \text{méd } y)$ sur le batch i

Transformer les données (x,y) :

- Pour chaque batch i , translater chaque point selon $\overrightarrow{M_i C}$

Recentrage sur la moyenne / médiane globale



Centrer chaque batch

Batch mean-centering

Calculer la moyenne de chaque batch i :

- $M_i = \text{centre}(\text{moy } x, \text{moy } y)$ pour les données du batch i

Transformer les données (x, y) par batch :

- Pour chaque échantillon j de chaque batch i :

$$x_{ij} \leftarrow (x_{ij} - M_i \cdot x)$$

$$y_{ij} \leftarrow (y_{ij} - M_i \cdot y)$$

Centrer - réduire chaque batch

Gene-wise standardization

Calculer la moyenne et la variance de chaque batch i :

- M_i = centre (moyenne x, moyenne y) pour toutes les données
- S_i = écart-type (x, y) sur le batch i

Transformer les données (x,y):

- Pour chaque échantillon j de chaque batch i :

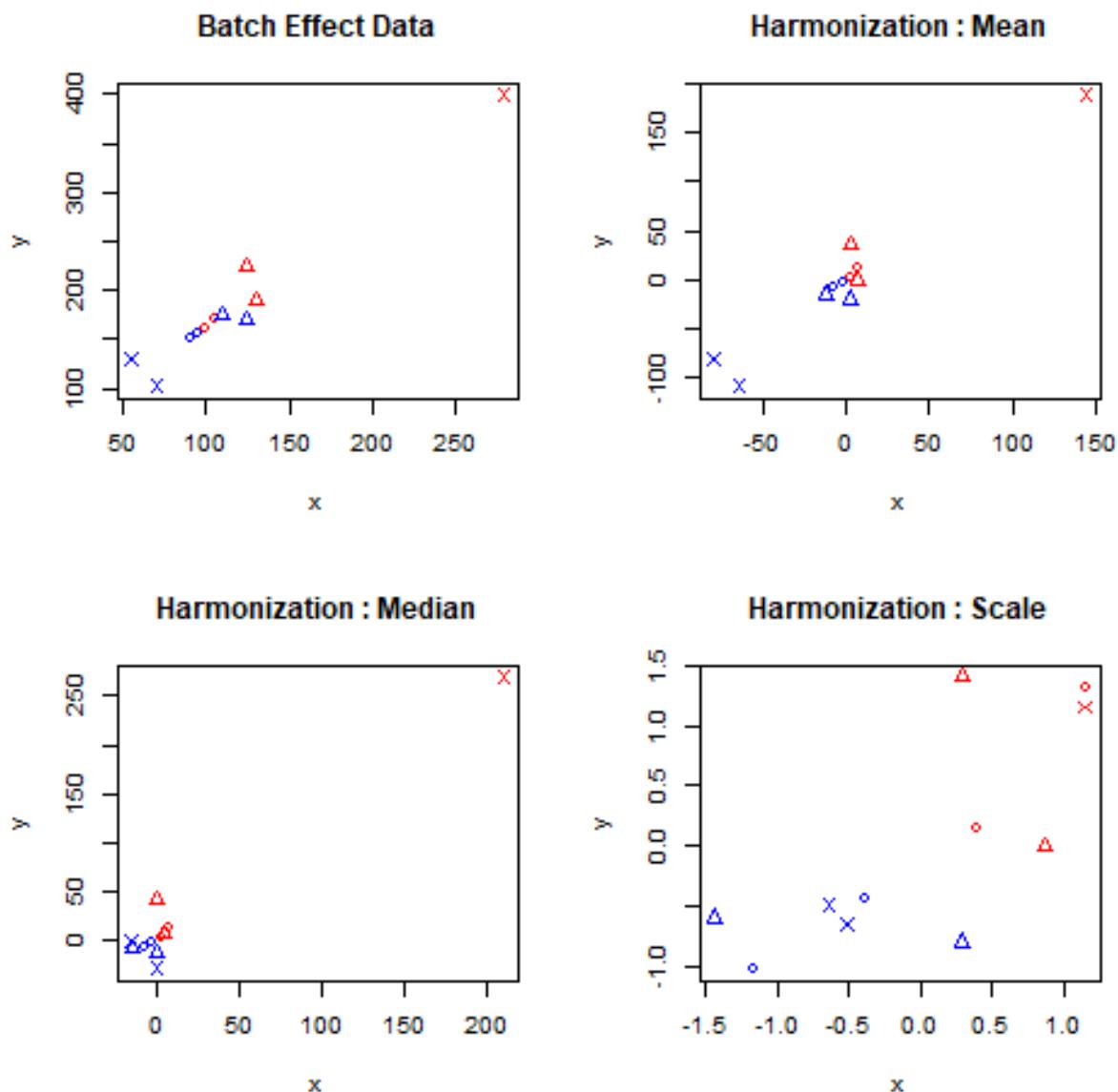
$$x_{ij} \leftarrow (x_{ij} - M_i \cdot x) / S_i \cdot x$$

$$y_{ij} \leftarrow (y_{ij} - M_i \cdot y) / S_i \cdot y$$



Chaque colonne (x,y) est de moyenne nulle
et de variance égale à 1

Comparaison *batch-mean*, *batch-median*, *batch-scale*



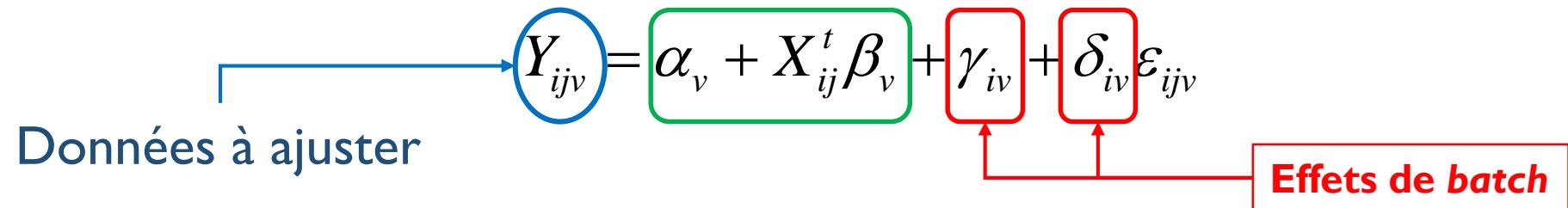
Méthode ComBat

Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics*. 2007 Jan;8(1):118-27. [doi:10.1093/biostatistics/kxj037](https://doi.org/10.1093/biostatistics/kxj037). Epub 2006 Apr 21. PMID: 16632515.

Hornung, R., Boulesteix, AL. & Causeur, D. Combining location-and-scale batch effect adjustment with data cleaning by latent factor adjustment. *BMC Bioinformatics* 17, 27 (2016). <https://doi.org/10.1186/s12859-015-0870-z>

Jean-Philippe Fortin et al., Harmonization of cortical thickness measurements across scanners and sites, *NeuroImage*, Volume 167, 2018, Pages 104-120, ISSN 1053-8119,
<https://doi.org/10.1016/j.neuroimage.2017.11.024>.

Formulation Johnson 2007 (I)



**Design
matrix**

- Y_{ijv} = valeur de la variable à expliquer v , du *batch* i et de l'échantillon j
- X_{ij}^t = vecteur des variables explicatives, du *batch* i et de l'échantillon j
- α_v et β_v = coefficients régression liant Y à X
- $\varepsilon_{ijv} \sim N(0, \sigma_v^2)$ le terme d'erreur (hypothèse loi normale)
- γ_{iv} et δ_{iv} respectivement l'effet additif et l'effet multiplicatif du *batch* i sur la variable v

Formulation Johnson 2007 (2)

$$Y_{ijv}^* = \hat{\alpha}_v + X_{ij}^t \hat{\beta}_v + \left(Y_{ijv} - \hat{\alpha}_v - X_{ij}^t \hat{\beta}_v - \gamma_{iv}^* \right) / \delta_{iv}^*$$

Données ajustées

1. On calcule $(\hat{\alpha}_v, \hat{\beta}_v, \hat{\gamma}_{iv})$ qui minimise :

$$\sum_{i,j} \left(Y_{ijv} - \alpha_v - X_{ij}^t \beta_v - \gamma_{iv} \right)^2 \text{ sous la condition } \sum_i n_i \gamma_{iv} = 0.$$

2. On émet l'hypothèse de distributions a priori (approche Bayésienne) :

$\gamma_{iv} \sim N(\gamma_i, \tau_i^2)$ et $\delta_{iv}^2 \sim \text{InverseGamma}(\lambda_i, \theta_i)$ où γ_i , τ_i^2 , λ_i et θ_i sont estimés par la méthode des moments.

3. On en déduit γ_{iv}^* et δ_{iv}^* comme des moyennes conditionnelles :

$$\gamma_{iv}^* = \frac{n_i \tau_i^2 \hat{\gamma}_{iv} + \delta_{iv}^{2*} \bar{\gamma}_i}{n_i \tau_i^2 + \delta_{iv}^{2*}} \text{ et } \delta_{iv}^{2*} = \frac{\bar{\theta}_i + \frac{1}{2} \sum_j (Z_{ijv} - \gamma_{iv}^*)^2}{\frac{n_i}{2} + \bar{\lambda}_i - 1} \text{ avec } Z_{ijv} = \frac{Y_{ijv} - \hat{\alpha}_v - X_{ij}^t \hat{\beta}_v}{\hat{\sigma}_v}$$

Reformulation Fortin 2017

$$Y_{ijv} = \alpha_v + X_{ij}^t \beta_v + Z_{ij}^t \theta_v + \delta_{iv} \varepsilon_{ijv}$$

Z_{ij}^t θ_v
δ_{iv} ε_{ijv}
Effets de batch

- $\gamma_{iv} = Z_{ij}^t \theta_v$, avec Z_{ij} des variables indicatrices (co-variables),
- et θ_v les coefficients de régression à estimer.

Exemple R : packages sva et bladderbatch

```
# Etape 1: installer BiocManager
# Voir : https://rdrr.io/bioc/sva/
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")

# Etape 2: installer sva et dépendances
BiocManager::install("sva")

# Etape 3: installer bladderbatch et dépendances
BiocManager::install("bladderbatch")
```

Données « bladderbatch »

Jeu composé de 2 tableaux :

- Données génomiques : 22283 lignes et 57 colonnes
 - ligne = gène (ex. nom : « 1007_s_at »)
 - colonne = échantillon (ex. nom : « GSM71019.CEL »)
- Données cliniques (phénotype) : 57 lignes et 4 colonnes
 - ligne = échantillon (ex. nom : « GSM71019.CEL »)
 - colonnes : *sample; outcome; batch; cancer*

```
> pdata
    sample  outcome batch cancer
GSM71019.CEL     1   Normal    3  Normal
GSM71020.CEL     2   Normal    2  Normal
GSM71021.CEL     3   Normal    2  Normal
GSM71022.CEL     4   Normal    3  Normal
GSM71023.CEL     5   Normal    3  Normal
GSM71024.CEL     6   Normal    3  Normal
GSM71025.CEL     7   Normal    2  Normal
GSM71026.CEL     8   Normal    2  Normal
GSM71028.CEL     9 sTCC+CIS    5  Cancer
GSM71029.CEL    10 sTCC-CIS    2  Cancer
GSM71030.CEL    11 sTCC-CIS    5  Cancer
```

```
> head(edata)
          GSM71019.CEL  GSM71020.CEL  GSM71021.CEL  GSM71022.CEL  GSM71023.CEL
1007_s_at    10.115170    8.628044    8.779235    9.248569   10.256841
1053_at      5.345168    5.063598    5.113116    5.179410    5.181383
117_at       6.348024    6.663625    6.465892    6.116422    5.980457
121_at       8.901739    9.439977    9.540738    9.254368    8.798086
1255_g_at    3.967672    4.466027    4.144885    4.189338    4.078509
1294_at      7.775183    7.110154    7.248430    7.017220    7.896419
          GSM71024.CEL  GSM71025.CEL  GSM71026.CEL  GSM71028.CEL  GSM71029.CEL
1007_s_at    10.023133    9.108034    8.735616    9.803271   10.168602
1053_at      5.248418    5.252312    5.220931    5.595771    5.025180
117_at       5.796155    6.414849    6.846798    5.841478    6.352257
121_at       8.002870    9.093704    9.263386    7.789240    9.834564
1255_g_at    3.919740    4.402590    4.173666    3.590649    4.338196
1294 --      7.044676    7.466777    7.201225    7.267214    7.225725
```

Utilisation de ComBat.R

```
library(sva)
library(bladderbatch)
# Load data
data(bladderdata)
# Get the expression data : it is the data matrix
edata = exprs(bladderEset[1:400,])
# Get the pheno data: it is the design matrix
pdata = pData(bladderEset[1:400,])
# Get the batch column
batch = pdata$batch
# ComBat parametric adjustment
combat_edata1 = ComBat(dat=edata, batch=batch, mod=NULL, par.prior=TRUE, prior.plots=FALSE)
```

Résultats (I)

Utilisation ACP sur les données avant et après ajustement par ComBat :

```
> pc_edata <- prcomp(t(edata), center = TRUE, scale = TRUE)
> summary(pc_edata)
Importance of components:
```

	PC1	PC2	PC3	PC4	PC5	PC6
Standard deviation	89.2706	49.6448	37.00883	27.65554	25.04918	22.4930
Proportion of Variance	0.3576	0.1106	0.06147	0.03432	0.02816	0.0227
Cumulative Proportion	0.3576	0.4682	0.52971	0.56403	0.59219	0.6149

	PC7	PC8	PC9	PC10	PC11	PC12
Standard deviation	21.52591	21.07911	19.5214	19.31445	18.30804	17.52918
Proportion of Variance	0.02079	0.01994	0.0171	0.01674	0.01504	0.01379
Cumulative Proportion	0.63569	0.65563	0.6727	0.68947	0.70452	0.71831

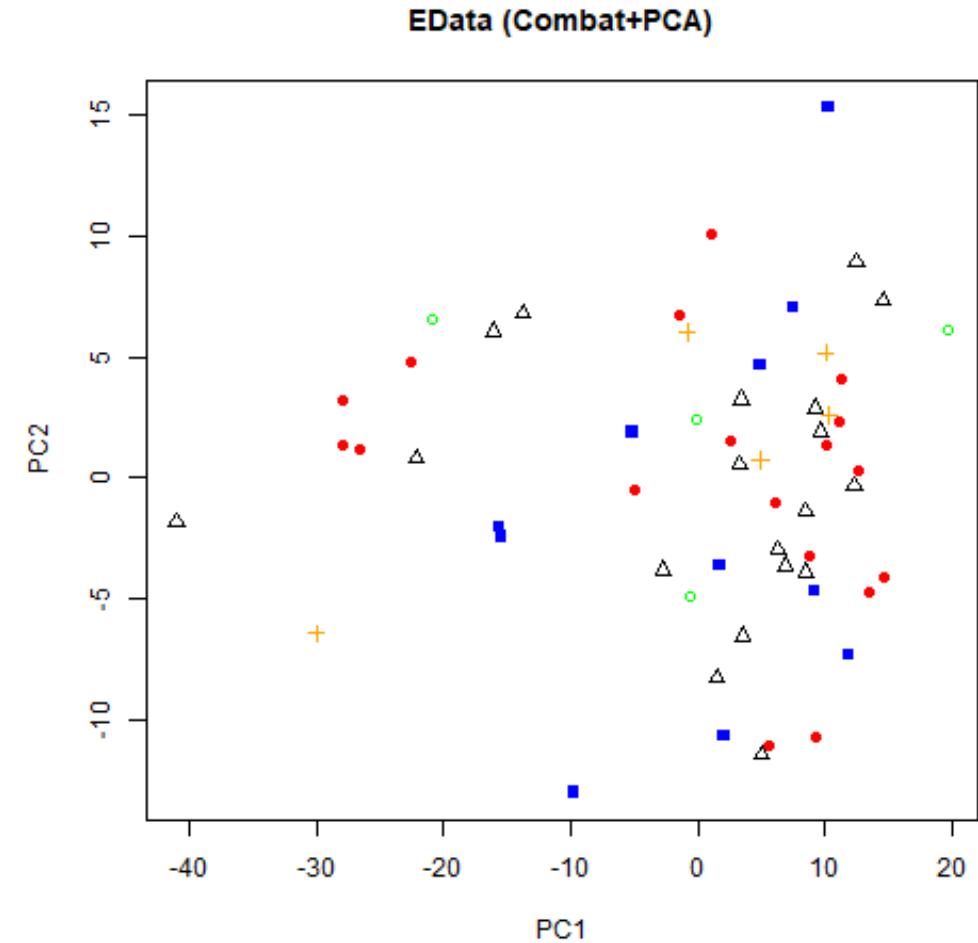
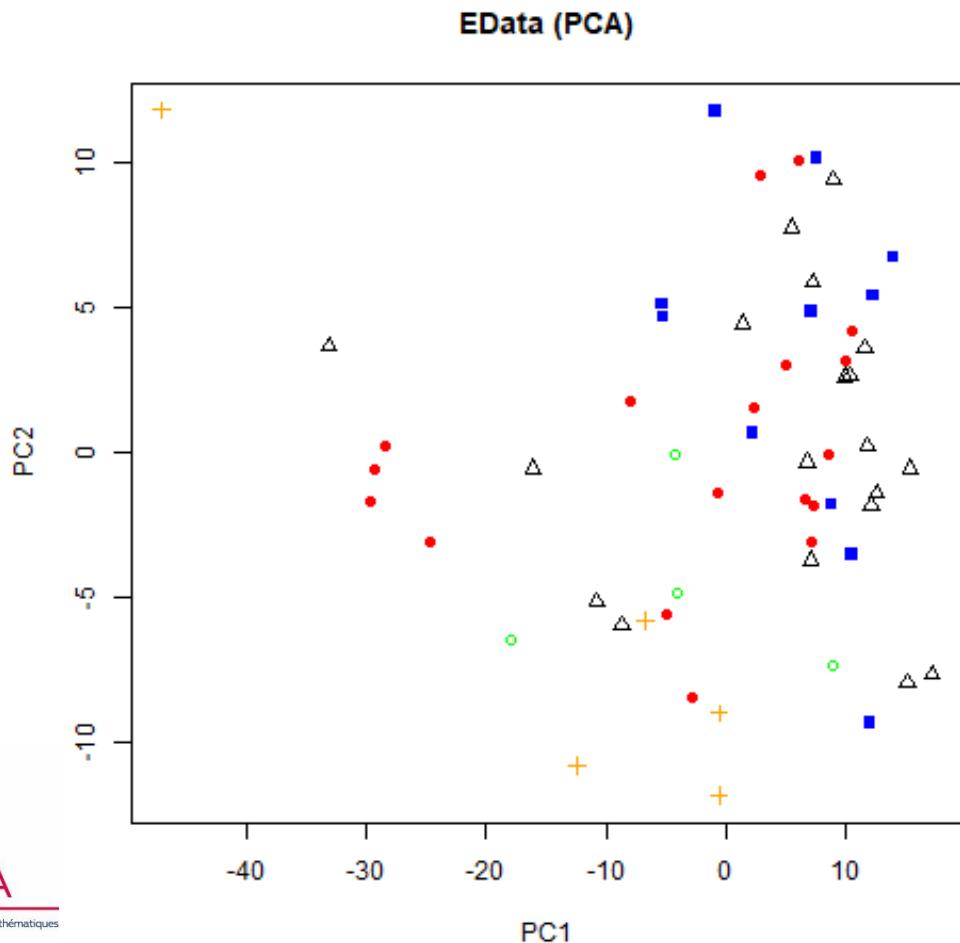

```
> pc_combatdata <- prcomp(t(combat_edatal), center = TRUE, scale = TRUE)
> summary(pc_combatdata)
Importance of components:
```

	PC1	PC2	PC3	PC4	PC5	PC6	PC7
Standard deviation	4.9233	2.3581	1.83592	1.4883	1.2409	1.14864	1.12351
Proportion of Variance	0.4848	0.1112	0.06741	0.0443	0.0308	0.02639	0.02525
Cumulative Proportion	0.4848	0.5960	0.66339	0.7077	0.7385	0.76488	0.79012

	PC8	PC9	PC10	PC11	PC12	PC13	PC14
Standard deviation	1.03600	0.97563	0.95019	0.87893	0.84690	0.77702	0.73575
Proportion of Variance	0.02147	0.01904	0.01806	0.01545	0.01434	0.01208	0.01083
Cumulative Proportion	0.81159	0.83063	0.84868	0.86413	0.87848	0.89055	0.90138

Résultats (2)

Utilisation ACP sur les données avant et après ajustement par ComBat :



Le projet HARMONY

HARMONization methods for optimized therapy
Appel d'offres structurant « Numérique en Oncologie » CGO / Régions 2019

Objectifs

- Développer des méthodes d'harmonisation de données
 - Ingénierie des caractéristiques
 - Apprentissage Profond
- Généraliser sur plusieurs types de cancer (dans le cadre du consortium défini) – 4 cancers différents.

Pipeline réutilisable testé sur les différentes unités du consortium

Région Bretagne

Equipe 1 : LaTIM, INSERM, UMR 1101, Univ. Brest

Région Pays de La Loire

Equipe 2 : Equipe Imagerie, biomarqueurs et thérapie, INSERM UMR 1253 « Imagerie et cerveau », Univ. Tours

Région Centre Val de Loire

Equipe 3 : Equipe “Oncologie nucléaire”, CRCINA, UMR 1232 Inserm, Univ. Nantes,

Equipe 4 : Equipe ISISV, LARIS EA 7315, Univ. Angers & IMA

Premiers travaux : bases IDCDB et BreaKHIS

Classification coupes histologiques IDC / Non IDC
IDC = *Invasive Ductal Carcinoma* (cancer du sein)

Choix Bases de données

Choix du cancer :

Cancer féminin - (ovaire et sein) – Spécialité ICO Angers.

Type de données :

Les données à disposition étaient des coupes histopathologiques (biopsie).

Notre pipeline devait s'appliquer sur celles-ci puis être utilisable sur d'autres types de cancer (Nantes, Tours, Brest).

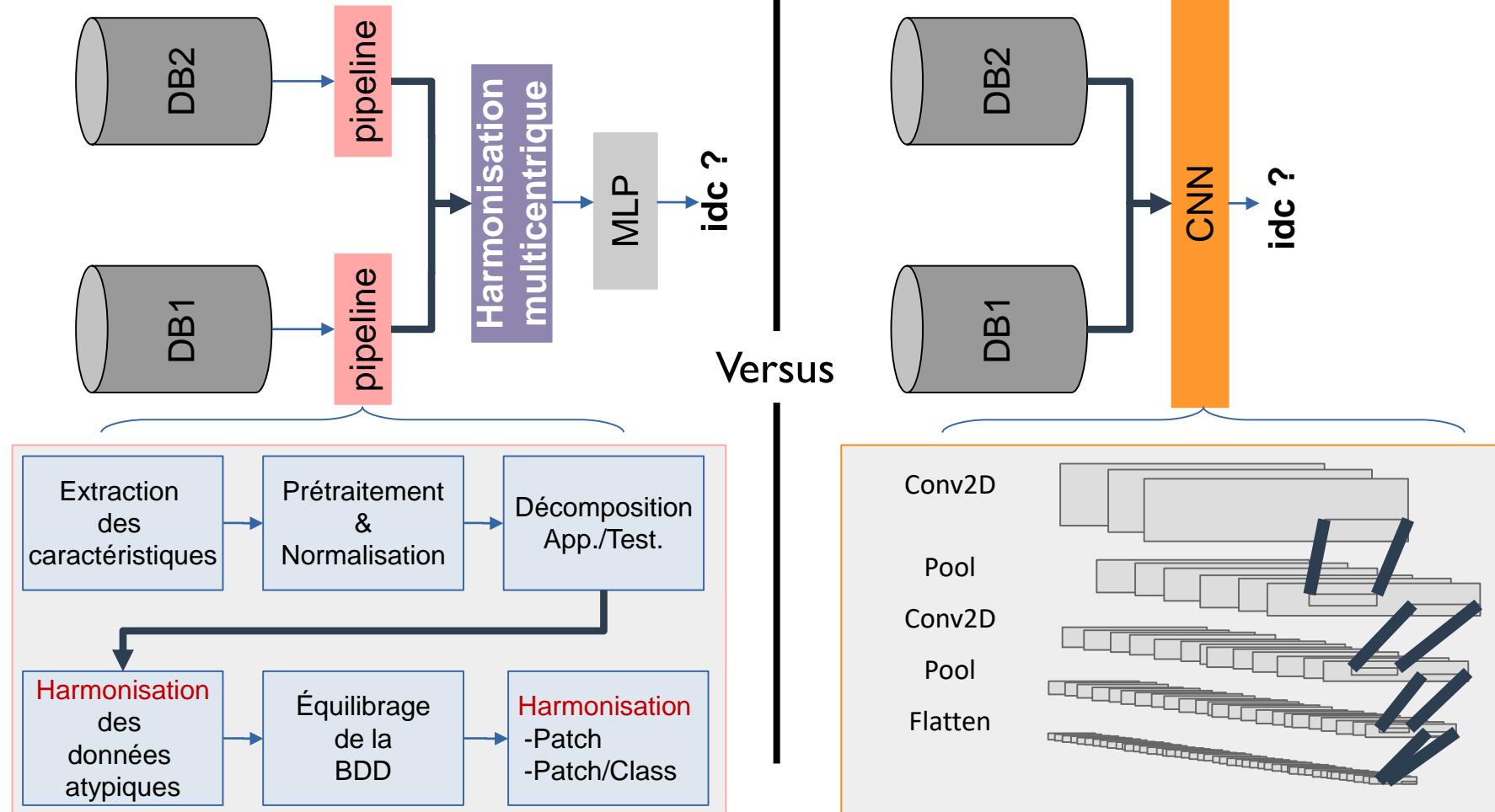
IDCDB – Database 1 (DB1):

- 162 images de lames entières de cancers du sein (Bca) numérisées à 40x [1] (162 patients).
- 277524 patches

BreakHIS – Database 2 (DB2):

- 9 109 images microscopiques de tissus tumoraux du sein prélevées sur 82 patients avec différents facteurs de grossissement (40X, 100X, 200X et 400X)[2].

Méthodologie



Extraction des caractéristiques

- Statistiques
32 caract.
- Texture R/G/B
50 caract. X Channel = 150 caract.
- Moments
8 caract.
- Entropies
3 caract.

Feature group	features
Histogram (R,G,B)	absolute energy, sum over the absolute value of consecutive changes, benford_correlation, count above mean, count_below_mean, first_location_of_maximum, first_location_of_minimum, cid_ce, minimum, maximum, median, kurtosis, longest_strike_above_mean, longest_strike_below_mean, mean_abs_change, mean_change, mean_second_derivative_central, variance, variance_coefficient percentage_of_reoccurring_datapoints_to_all_datapoints, percentage_of_reoccurring_values_to_all_values, skewness, ratio_value_number_to_series_length, standard_deviation, sum_of_reoccurring_data_points,sum_of_reoccurring_values, sum_values, variance
Texture	Correlation, Homogeneity, Energy, Contrast.
Entropy	Shanon Entropy, Simple Entropy, Sample entropy.
Moment	Moments, Hu Moments

Prétraitement des caractéristiques :

- 1) Quantification des features.
- 2) Suppression des Infs et des NaNs

Normalisation

StandardScaling

$$x' = \frac{x - \bar{x}}{\sigma}$$

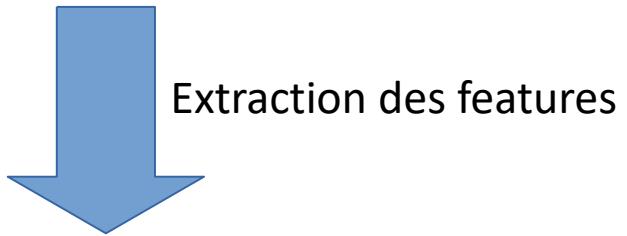
- Formule :
value = (value – mean) / std
- La présence de données aberrantes impacte la normalisation.
- Centre et réduit

RobustScaling

$$x' = \frac{x - median}{Q3 - Q1}$$

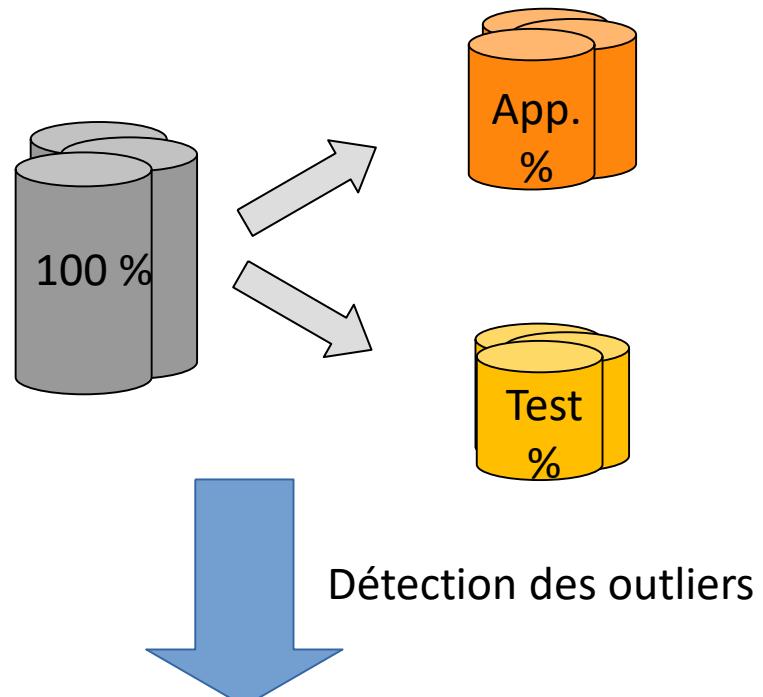
- Formule :
value = (value – median) / (p75 – p25)
- Non influencé par la présence de données aberrantes.
- Ne centre pas.

Décomposition App./Test

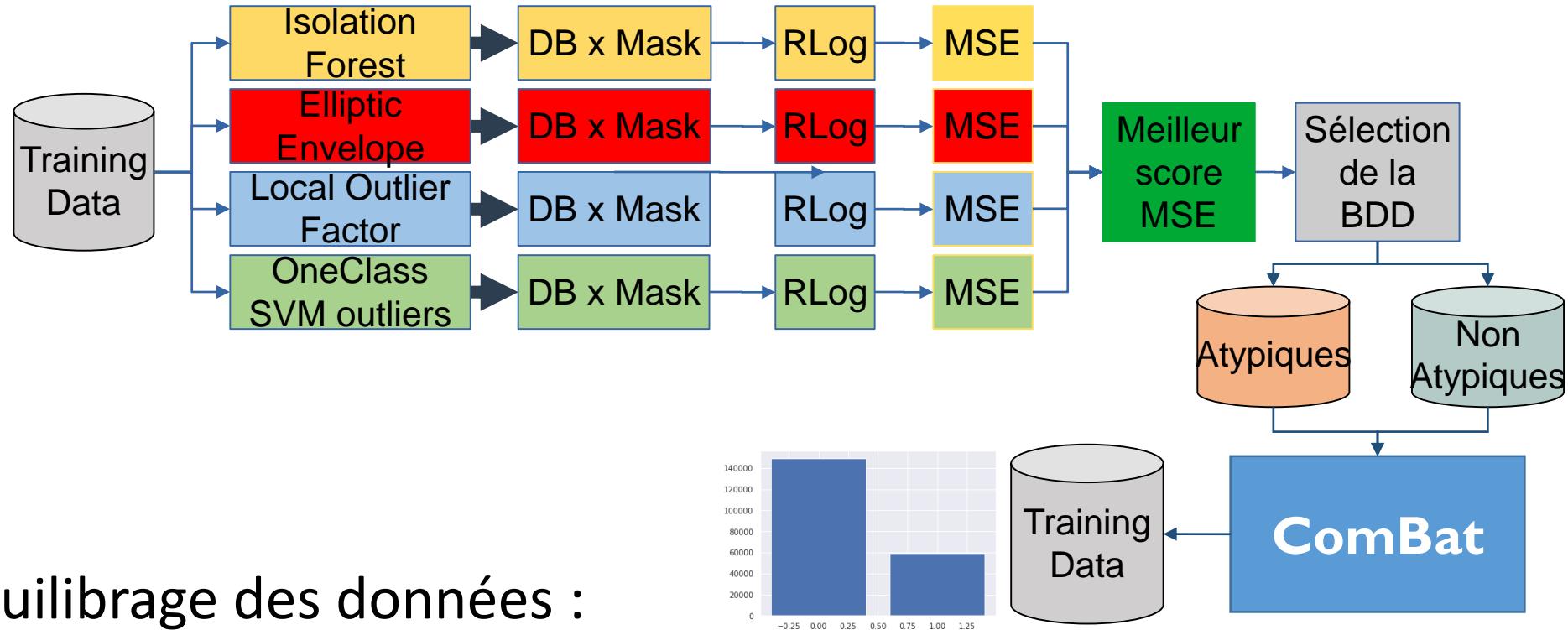


Split des données pour App. et Test (sur chacune des deux bases) :

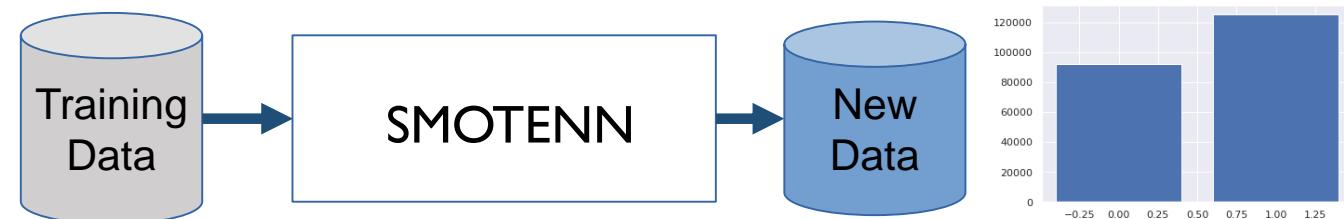
- **75/25**
- **70/30**
- **80/20**



Détections données aberrantes/atypiques

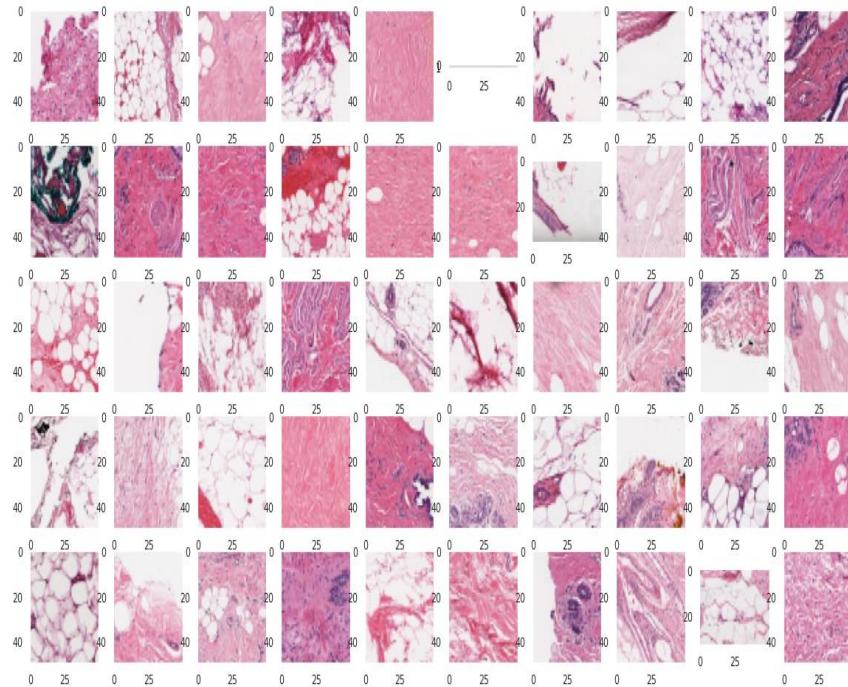


Équilibrage des données :

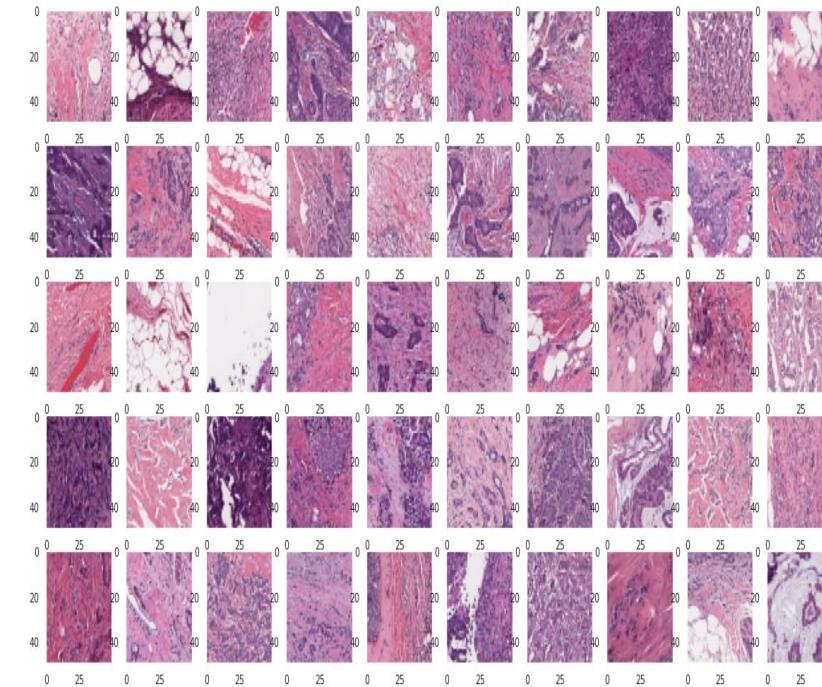


Harmonisation intra-base (I)

Pourquoi une harmonisation par patchs ?



Classe 1

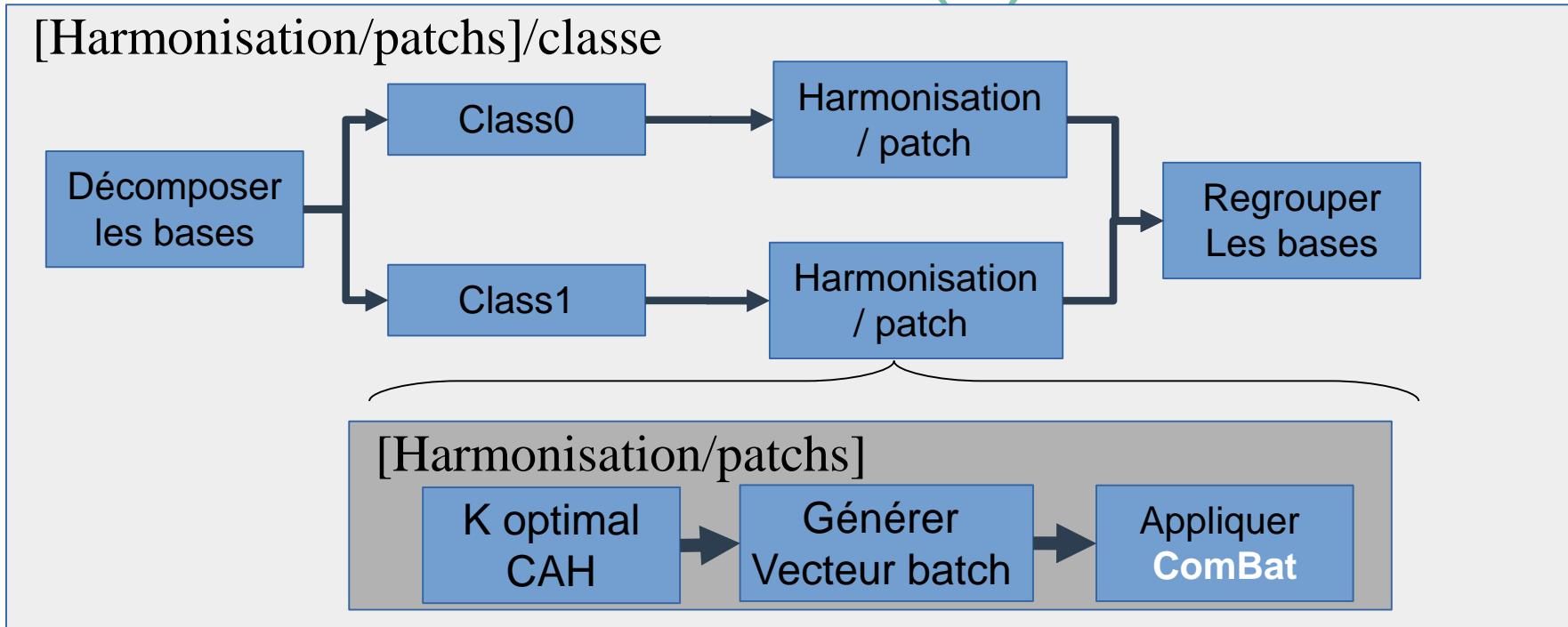


Classe 2

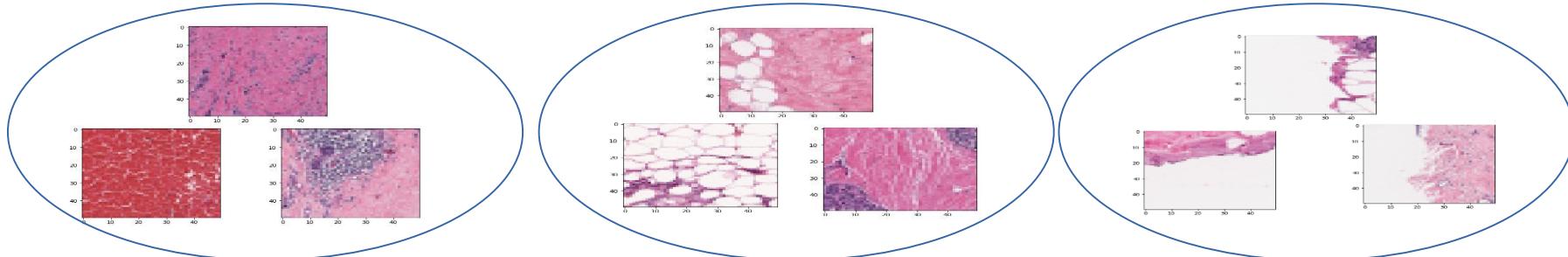
Deux variabilités identifiées :

- 1) Patchs de bords / patchs centraux
- 2) Texture, couleur

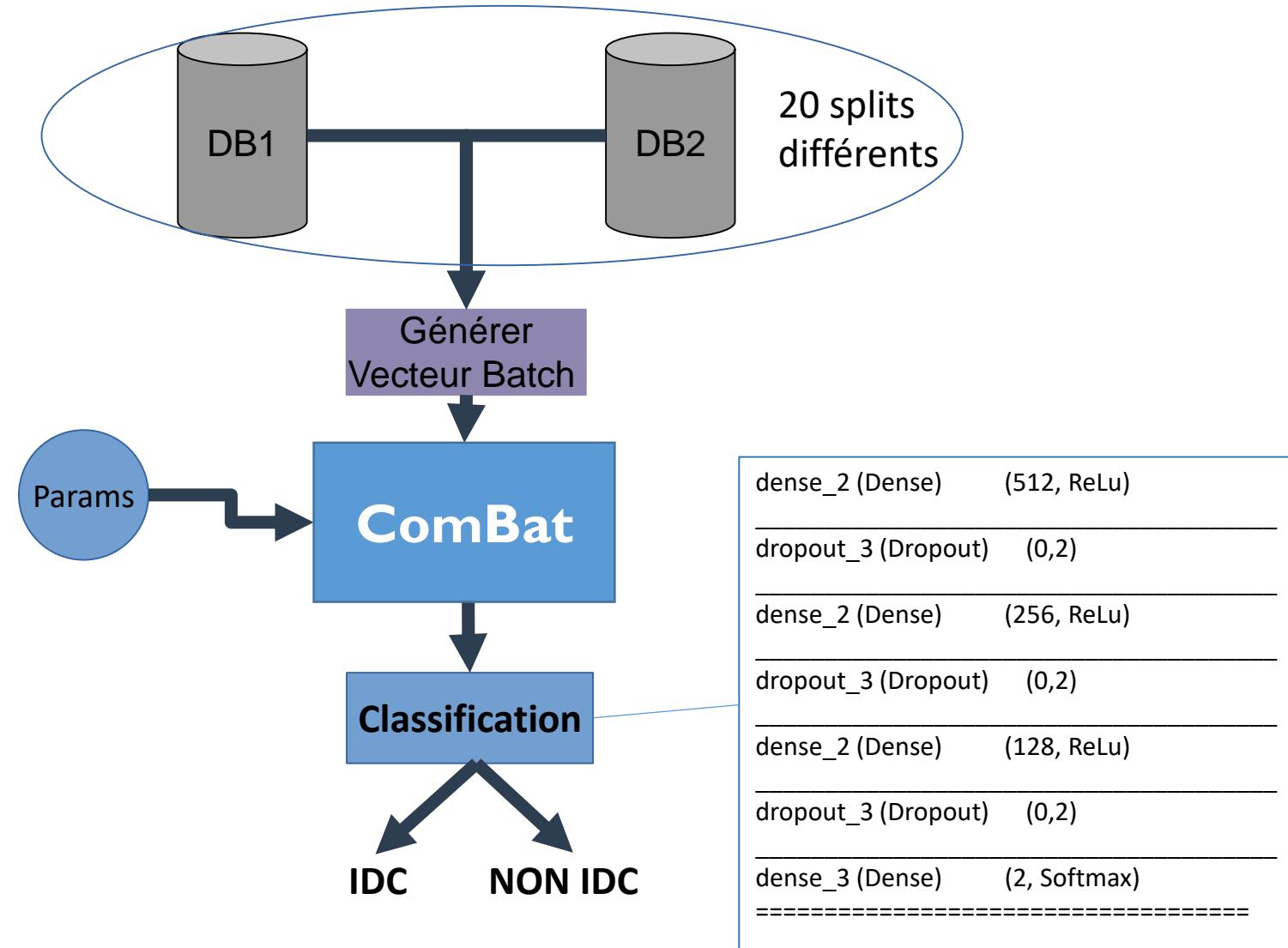
Harmonisation intra-base (2)



batch_0				batch_1				batch_2			
0	0	...	0	0	1	1	...	1	1	2	2



Classification



Résultats harmonisation intra-base

Modèle de Base

IDCDB :

Training Accuracy: 83.40%
 Validation Accuracy: 83.47%
 Testing Accuracy : **67.04 %**

BreKHis :

Training Accuracy: 98.27%
 Validation Accuracy: 88.95%
 Testing Accuracy : **89.61 %**

Avant Harmonisation

IDCDB_SMOTENN :

Training Accuracy: 92.74%
 Validation Accuracy: 92.32%
 Testing Accuracy : **79.16%**

BreKHis_SMOTENN :

Training Accuracy: 99.86%
 Validation Accuracy: 96.25%
 Testing Accuracy : **88.52%**

Après Harmonisation - Config3 : ComBatOutliers > SMOTENN

IDCDB_Config3 :

Training Accuracy: 95.47%
 Validation Accuracy: 95,03%
 Testing Accuracy : **83±1,7%**

BreKHis_Config3 :

Training Accuracy: 99.86%
 Validation Accuracy: 96.25%
 Testing Accuracy : **93,4±1,8 %**

Pour IDCDB	Test
Soumya et al. [6]	92,55 %
Choudhary et al. [7]	92,07 %
Notre meilleur modèle Config3>CombatByPatch	94,7 %

Pour BreakHIS	Test
Sanchez-Morillo et al. [4]	88,3 %
Boumaraf et al. [5]	87,7 %
Notre meilleur modèle	95,2 %

Résultats harmonisation inter-base (multicentrique)

MultiC	Apprentissage	Validation	Test sur IDC	Test sur BreakHis
Model IDC	89,46 %	87,87 %	87,55 %	58,16 %
Model BreakHis	99,27 %	92,34 %	47,2 %	93,28 %
Model_MultiC	95,67 %	95,04 %	80 %	81 %

Conclusion

L'harmonisation apporte une amélioration des résultats de classification :

- Intra-base : +15,5 % sur IDCDB et +6,68 % sur BreaKHis par rapport au modèle avec base équilibrée.
- Inter-base : l'harmonisation améliore la robustesse vis-à-vis des données provenant d'un centre extérieur.

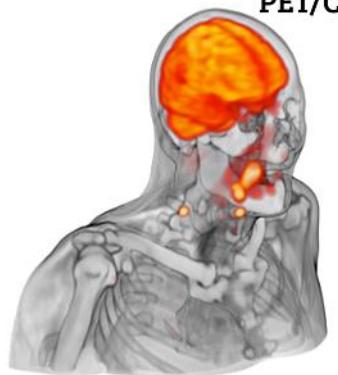
Limitation :

- Harmonisation sur l'ensemble d'App. + de Test.
- Solution : PCA ou Harmonisation par base de référence

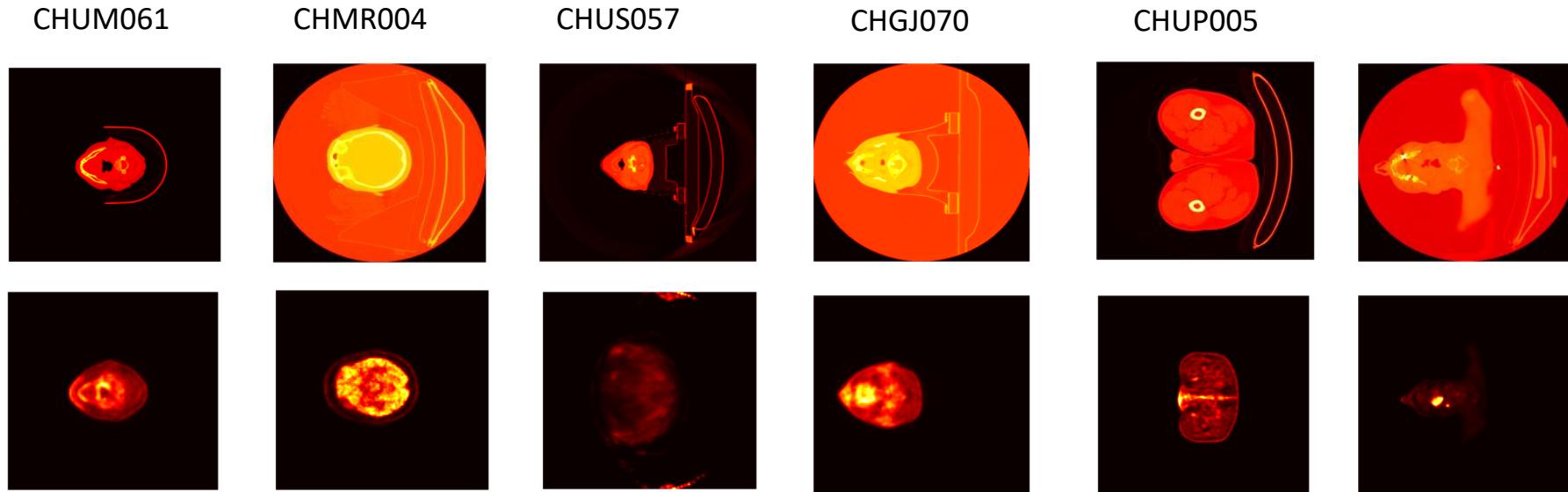
Le challenge Hecktor 2021

HECKTOR 2021

HEad and neCK TumOR segmentation and outcome prediction in
PET/CT images

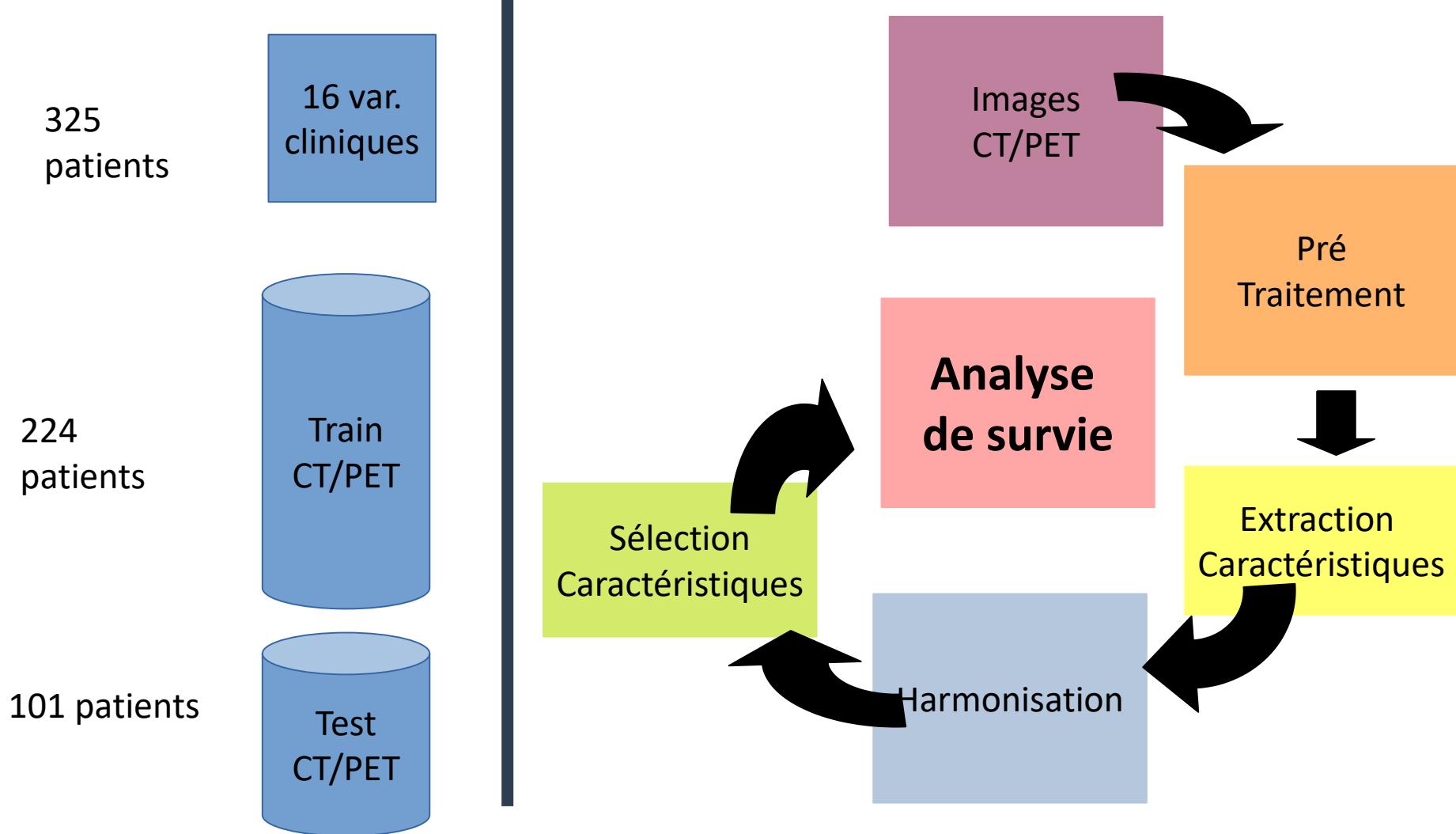


Le challenge

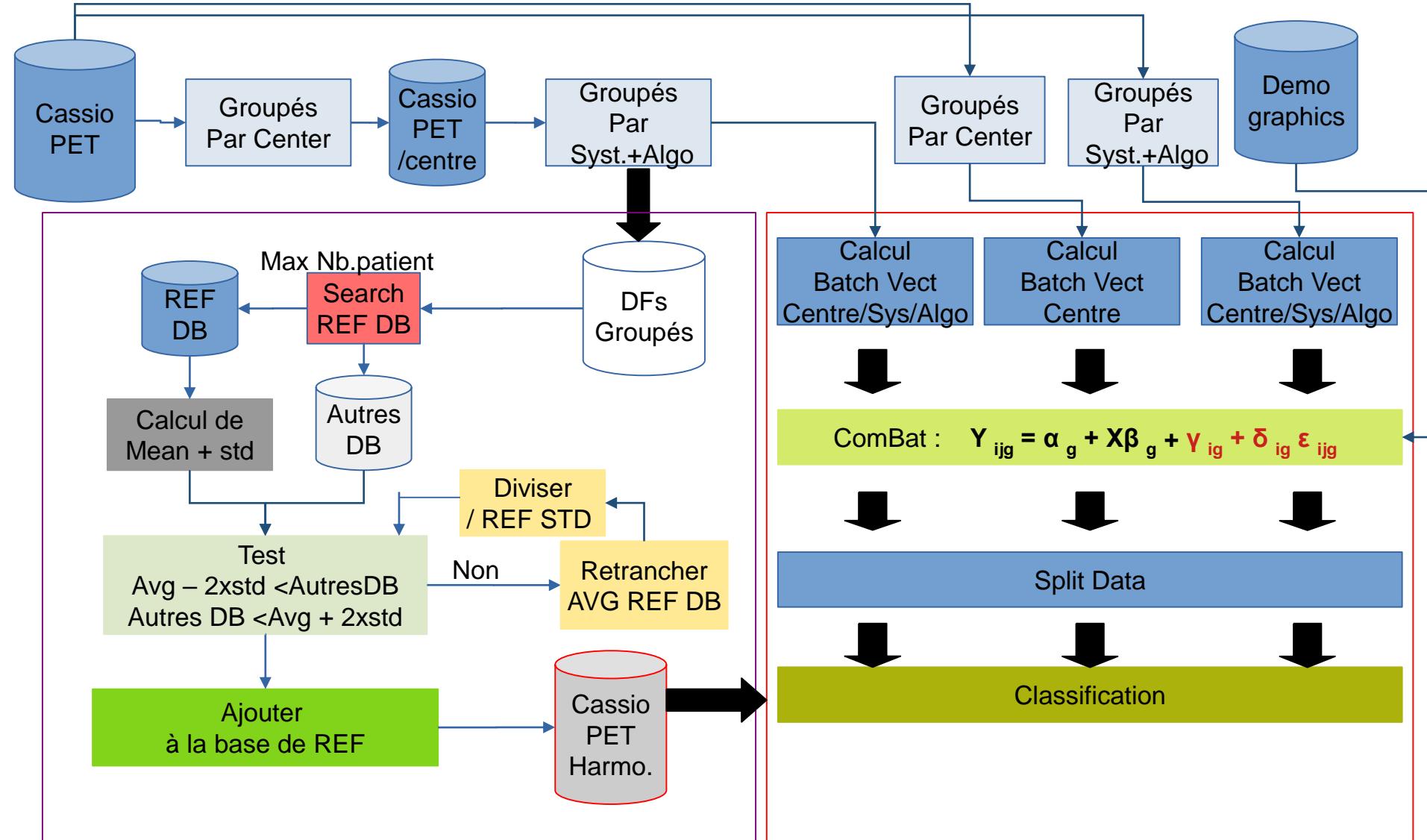


1. Segmentation de tumeurs primaires de la tête et du dos.
2. Prédiction de la survie sans progression à partir d'images PET/CT + données cliniques.
3. Tâche 2, sauf que les annotations (masques) sont fournies aux challengers.

Méthodologie



Harmonisation par base de référence



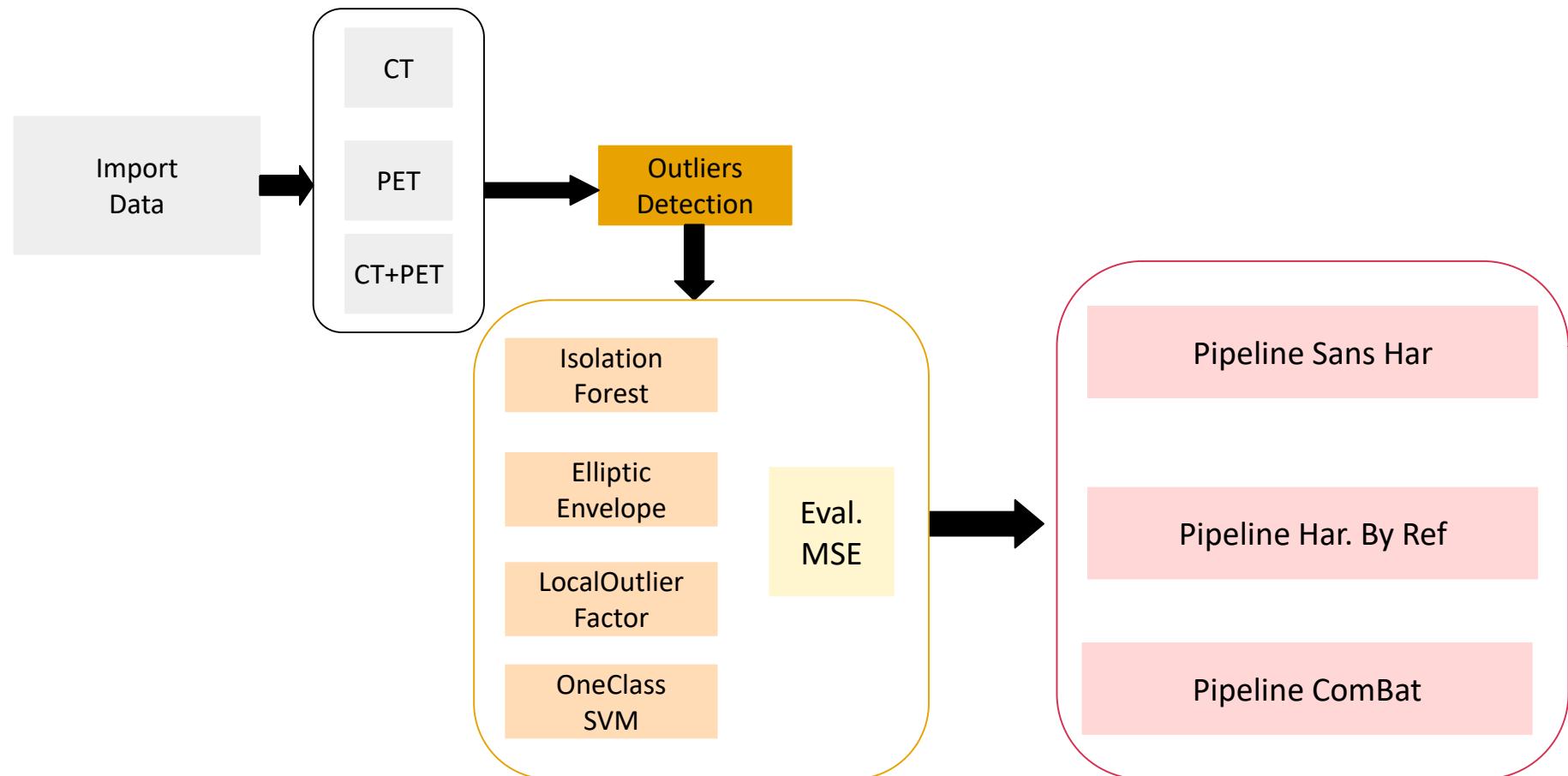
Pipeline / Benchmark

Appliquer sur

App +Test =

224+101 =

305 patients



QUESTIONS ?