# Statistical challenges in analyzing 16S microbiome data

*An application to the identification of microbe-regulated pathways in allergy and auto-immunity*

Marine Jeanmougin

*Institut Curie, U932 - Immunity and cancer*

Journées MAS, August 28th, 2014

## Goal

⤳ Unravel the **inflammatory pathways** during the **host-pathogen interactions** which may trigger allergic or autoimmune inflammation



## Clinical impact

⤳ Identify key microbes and molecular targets to develop **novel intervention strategies**
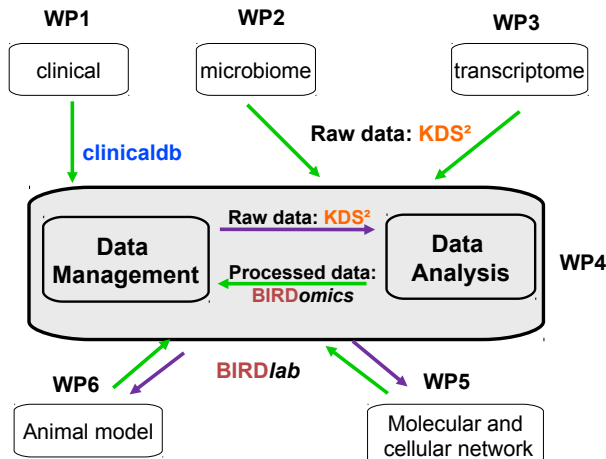
# WP4: data management and analysis



**WP1**
clinical

**WP2**
microbiome

**WP3**
transcriptome

**clinicaldb**

Raw data: **KDS²**

**Data Management**

Raw data: **KDS²**

Processed data: **BIRD*omics***

**Data Analysis**

**WP4**

**BIRD*lab***

**WP6**
Animal model

**WP5**
Molecular and cellular network

*King's college*
- Sophia Tsoka
- Gareth Muirhead

*FIOH*
- Dario Greco

*Karolinska*
- Juha Kere
- Shintaro Katayama

*Fios Genomics*
- Varrie Ogilvie
- Sarah Lynagh
- Max Bylesjo

*Institut Curie*
- Vassili Soumelis
- Philippe Hupé
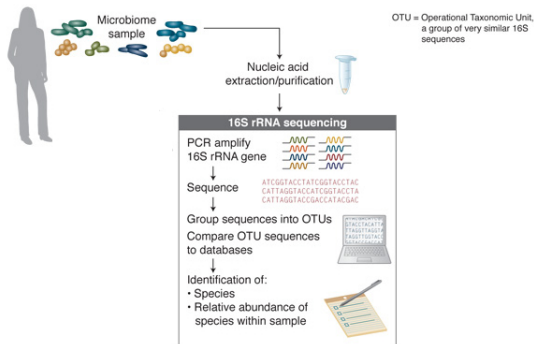- Gerome Jules-Clément
- Marine Jeanmougin

**K**nowledge and **D**ata **S**haring **S**ystem
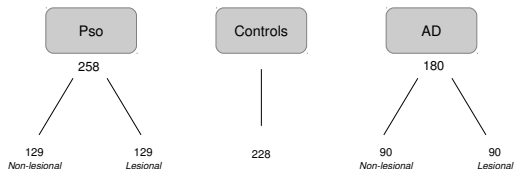**Clinical D**ata**B**ase
**BI**ological **R**esult **D**atabase

## The skin microbiome

- Ecosystem of microbes that live on the skin

- Culture independent microbiome research:

  - total microbiome DNA sequencing
  - 16S rRNA sequencing

institut**Curie**

- **Discrete counts** of sequence reads: number of time each OTU was found in a sample

- **Large-scale** data: ~ 17000 OTUs × 666 samples



- **Heterogeneneous** data due to:
  - ▸ biological phenomena: some species are found in only a small % of samples
  - ▸ technical reasons: others are not detected (insufficient seq depth)

→ **Library size** (total reads per sample) vary by orders of magnitude

→ **Sparsity**: *i.e.* most OTUs are rare (98% of sparsity in raw data)

→ **Overdispersion**: variance grows faster than the mean

Comparison across samples with different library sizes may induce biases in the downstream analysis

- **Differential analysis**: the higher sequencing depth, the higher counts

- **Diversity/richness estimation**: rarefaction phenomenon
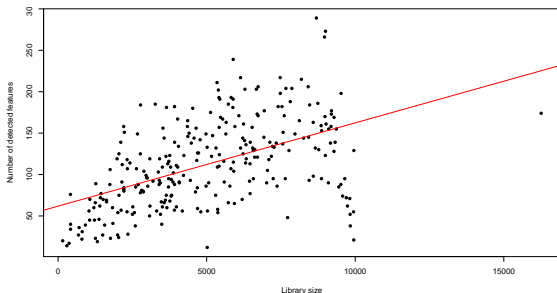  "The number of taxonomic features detected in a sample depends on the amount of sequencing performed"



Figure: Illustration of the rarefaction phenomenon on MAARS data

# Normalisation: current practices

### Rarefying

Random subsampling of each sample to a common depth:

- Omission of available data: **add artificial uncertainty**
- Inflate the variance and induce a **loss of power** in differential analysis

### Total-sum scaling (TSS): proportional abundance of species

Divide read counts by the total number of reads in each sample:

$$\widetilde{c}_{ij} = \frac{c_{ij}}{s_j}$$

where:

- $c_{ij}$ is the number of times taxonomic feature $i$ was observed in sample $j$
- $s_j = \sum_i c_{ij}$, sum of counts for sample $i$

In practice...

- Does not account for **heteroscedasticity**
- Dillies et al. demonstrated biais in RNA-seq data: undue influence of **high-count genes** on normalized counts
  - ↗ FPR when **differences in library composition**

Methods derived from the field of RNA-seq data analysis:

**1** **Quantile** (Q): Quantiles of the count distributions are matched between samples

**2** **Upper-Quartile** (UQ): scale factors are calculated from the 75% quantile of the counts for each library

**3** **Relative Log Expression** (RLE) - DESeq (Anders & Huber 2010):

$$\hat{s}_j = median_i \left( \frac{c_{ij}}{(\pi_{v=1}^n c_{iv})^{1/n}} \right)$$

where $n$ is the sample size.

**4** **Trimmed Mean of M-values** (TMM) - EdgeR (Robinson et al. 2010)
Trim data by log-fold-changes $M_i$ and absolute intensity $A_i$:

$$M_i = \log_2 \frac{c_{ij}/s_j}{c_{ij'}/s_{j'}}; \qquad A_i = \frac{1}{2} \log_2(c_{ij}/s_j \times c_{ij'}/s_{j'});$$

▷ Scaling factor: trimmed mean of the log-abundance ratios

**5** **Voom** (Law et al. 2014)
Log-counts per million (log-cpm) value:

$$y_{ij} = \log_2 \left( \frac{c_{ij} + 0.5}{s_j + 1} \times 10^6 \right)$$

The library size is offset by 1 to ensure that $0 < \frac{c_{ij}+0.5}{s_j+1} < 1$

Methods derived from the field of RNA-seq data analysis:

1. **Quantile** (Q): Quantiles of the count distributions are matched between samples
2. **Upper-Quartile** (UQ): scale factors are calculated from the 75% quantile of the counts for each library
3. **Relative Log Expression** (RLE) - DESeq (Anders & Huber 2010):

$$\hat{s}_j = median_i \left( \frac{c_{ij}}{\left( \pi_{v=1}^n c_{iv} \right)^{1/n}} \right)$$

where $n$ is the sample size.

4. **Trimmed Mean of M-values** (TMM) - EdgeR (Robinson et al. 2010)
   Trim data by log-fold-changes $M_i$ and absolute intensity $A_i$:

$$M_i = \log_2 \frac{c_{ij}/s_j}{c_{ij'}/s_{j'}}; \qquad A_i = \frac{1}{2} \log_2 (c_{ij}/s_j \times c_{ij'}/s_{j'});$$

   ▷ Scaling factor: trimmed mean of the log-abundance ratios

5. **Voom** (Law et al. 2014)
   Log-counts per million (log-cpm) value:

$$y_{ij} = \log_2 \left( \frac{c_{ij} + 0.5}{s_j + 1} \times 10^6 \right)$$

The library size is offset by 1 to ensure that $0 < \frac{c_{ij}+0.5}{s_j+1} < 1$

Methods derived from the field of RNA-seq data analysis:

1. **Quantile** (Q): Quantiles of the count distributions are matched between samples
2. **Upper-Quartile** (UQ): scale factors are calculated from the 75% quantile of the counts for each library
3. **Relative Log Expression** (RLE) - DESeq (Anders & Huber 2010):

$$\hat{s}_j = median_i \left( \frac{c_{ij}}{(\pi_{v=1}^n c_{iv})^{1/n}} \right)$$

where $n$ is the sample size.

4. **Trimmed Mean of M-values** (TMM) - EdgeR (Robinson et al. 2010)
Trim data by log-fold-changes $M_i$ and absolute intensity $A_i$:

$$M_i = \log_2 \frac{c_{ij}/s_j}{c_{ij'}/s_{j'}}; \qquad A_i = \frac{1}{2} \log_2(c_{ij}/s_j \times c_{ij'}/s_{j'});$$

▷ Scaling factor: trimmed mean of the log-abundance ratios

5. **Voom** (Law et al. 2014)
Log-counts per million (log-cpm) value:

$$y_{ij} = \log_2 \left( \frac{c_{ij} + 0.5}{s_j + 1} \times 10^6 \right)$$

The library size is offset by 1 to ensure that $0 < \frac{c_{ij}+0.5}{s_j+1} < 1$

# Normalisation: alternative approaches

Methods derived from the field of RNA-seq data analysis:

1. **Quantile** (Q): Quantiles of the count distributions are matched between samples
2. **Upper-Quartile** (UQ): scale factors are calculated from the 75% quantile of the counts for each library
3. **Relative Log Expression** (RLE) - DESeq (Anders & Huber 2010):

$$\hat{s}_j = median_i\left(\frac{c_{ij}}{(\pi_{v=1}^n c_{iv})^{1/n}}\right)$$

where $n$ is the sample size.

4. **Trimmed Mean of M-values** (TMM) - EdgeR (Robinson et al. 2010)
Trim data by log-fold-changes $M_i$ and absolute intensity $A_i$:

$$M_i = \log_2 \frac{c_{ij}/s_j}{c_{ij'}/s_{j'}}; \qquad A_i = \frac{1}{2}\log_2(c_{ij}/s_j \times c_{ij'}/s_{j'});$$

▷ Scaling factor: trimmed mean of the log-abundance ratios

5. **Voom** (Law et al. 2014)
Log-counts per million (log-cpm) value:

$$y_{ij} = \log_2\left(\frac{c_{ij} + 0.5}{s_j + 1} \times 10^6\right)$$

The library size is offset by 1 to ensure that $0 < \frac{c_{ij}+0.5}{s_j+1} < 1$

institut**Curie**

Methods derived from the field of RNA-seq data analysis:

1. **Quantile** (Q): Quantiles of the count distributions are matched between samples

2. **Upper-Quartile** (UQ): scale factors are calculated from the 75% quantile of the counts for each library

3. **Relative Log Expression** (RLE) - DESeq (Anders & Huber 2010):

$$\hat{s}_j = median_i \left( \frac{c_{ij}}{(\pi_{v=1}^n c_{iv})^{1/n}} \right)$$

where $n$ is the sample size.

4. **Trimmed Mean of M-values** (TMM) - EdgeR (Robinson et al. 2010)
Trim data by log-fold-changes $M_i$ and absolute intensity $A_i$:

$$M_i = \log_2 \frac{c_{ij}/s_j}{c_{ij'}/s_{j'}}; \qquad A_i = \frac{1}{2} \log_2 (c_{ij}/s_j \times c_{ij'}/s_{j'});$$

$\triangleright$ Scaling factor: trimmed mean of the log-abundance ratios

5. **Voom** (Law et al. 2014)
Log-counts per million (log-cpm) value:

$$y_{ij} = \log_2 \left( \frac{c_{ij} + 0.5}{s_j + 1} \times 10^6 \right)$$

The library size is offset by 1 to ensure that $0 < \frac{c_{ij}+0.5}{s_j+1} < 1$

## CSS strategy

Paulson, J. *et al.* (2013), <u>Nature Methods</u>

- $q_j^l$: $l$th quantile of sample $j$

- $s_j^l = \sum_{i | c_{ij} \leq q_j^l} c_{ij}$

- $N$: normalization constant (ex: the $\text{med}_j(s_j^l)$)

$$\widetilde{c}_{ij} = \frac{c_{ij}}{s_j^l} N$$

▸ avoid placing undue influence on high-count features

## Selection of the appropriate quantile

- $\bar{q}^l = \text{med}_j(q_j^l)$, median $l$th quantile across samples

- $d_l = med_j|q_j^l - \bar{q}^l|$, **median absolute deviation** of sample-specific quantiles

- $\hat{l}$: smallest value for which high instability is detected

Figure: Distribution of library sizes across normalisation approaches
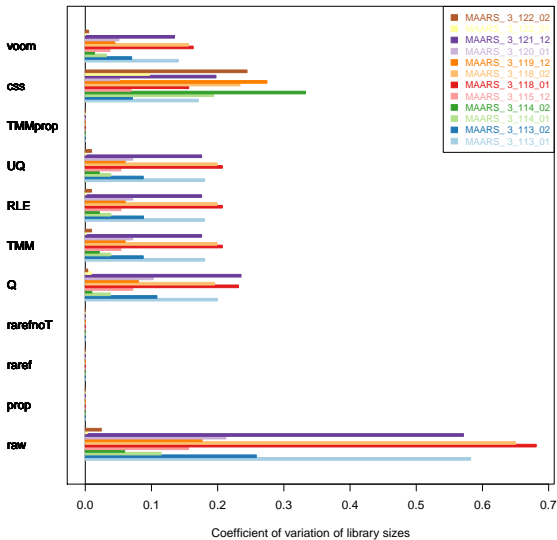
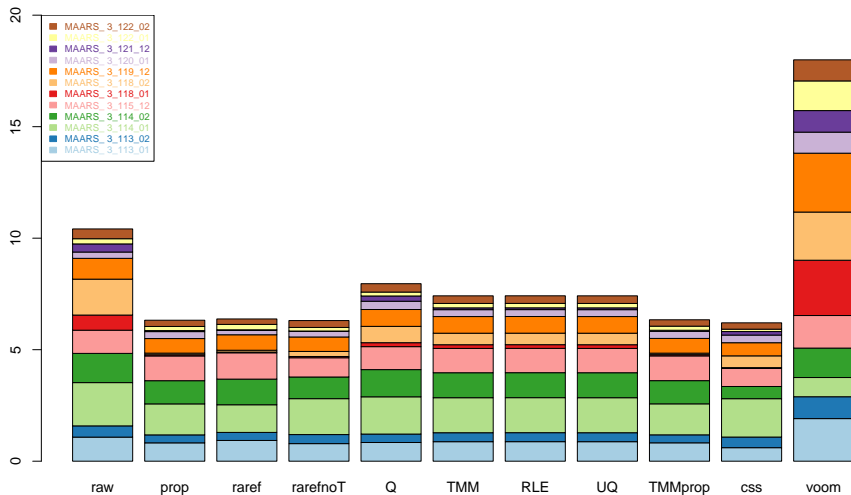Figure: Homogeneity of library sizes between technical replicates
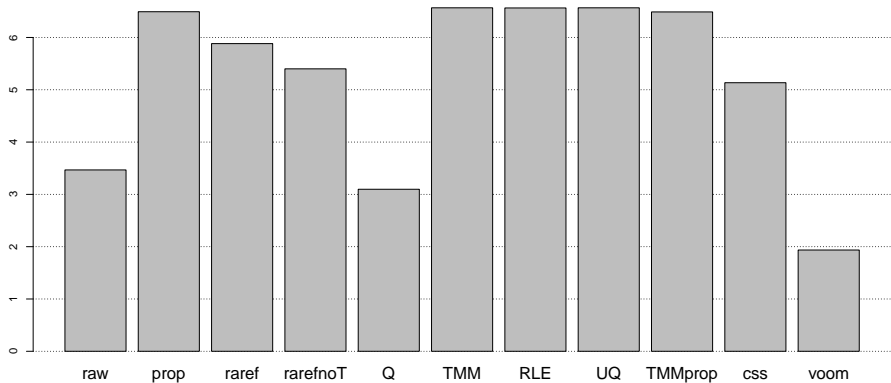
Figure: Distances between technical replicates

Figure: Ratio of distances between clinical groups and technical replicates

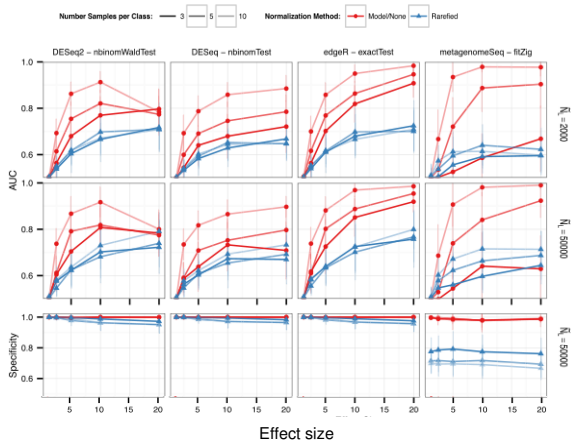McMurdie, P.J. and Holmes, S. (2014), <u>PLOS Comp. Biol.</u>
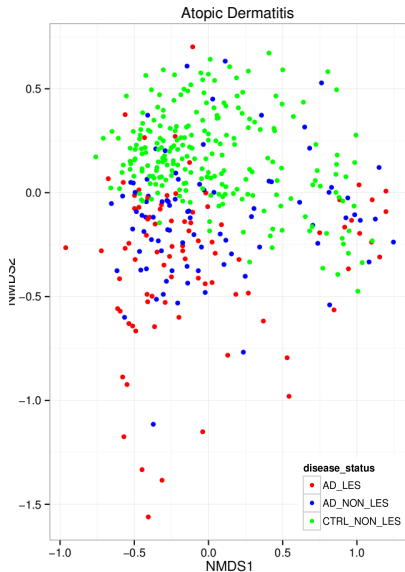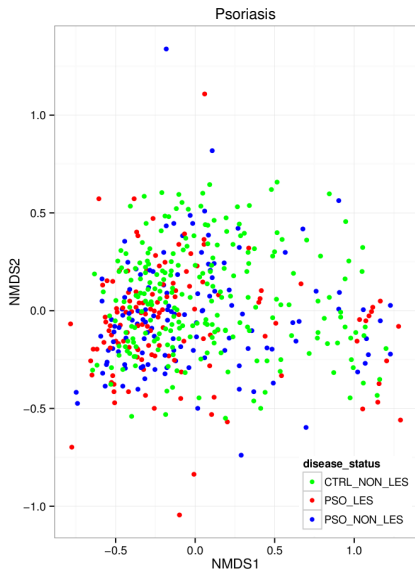


Figure: Performance of differential abundance detection on simulated data

Preliminary results on permuted data show that proportions and rarefying exhibit a FPR of 30%

- TMM and RLE are the best compromises :
  - show good results on simulated data (McMurdie 2014)
  - reduce the heterogeneity in library sizes
  - lower the distances between technical replicates
  - do not degrade the biological signal

- UQ performs well but need to be tested on simulated data

- Voom, Q and CSS normalisation approaches to be proscribed

- Perspectives for differential abundance testing: zero-inflated negative binomial model

# Outline

**Non-metric MultiDimensional Scaling**



Psoriasis

Atopic Dermatitis

disease_status
- CTRL_NON_LES
- PSO_LES
- PSO_NON_LES

disease_status
- AD_LES
- AD_NON_LES
- CTRL_NON_LES

# Integration of microbiome and transcriptome data

institut**Curie**

▷ Unravel the interdependencies between skin microbiome and transcriptome

## Univariate analysis

- Associate the presence of a given microbe with different transcriptome profiles

## Multivariate exploratory analysis

**Canonical Correlation Analysis**:

⤳ identify largest correlations between linear combinations of transcriptome and OTU profiles

Let us consider two matrices $X$ and $Y$ of order $n \times p$ and $n \times q$ respectively, with $p \leq q$.

For $S = 1, ..., p$, find $\rho_1 \geq \rho_2 \geq ... \geq \rho_p$ such as:

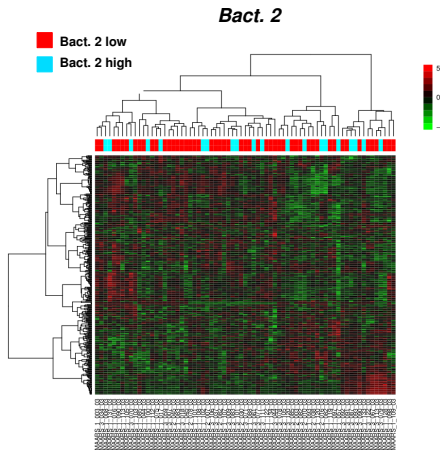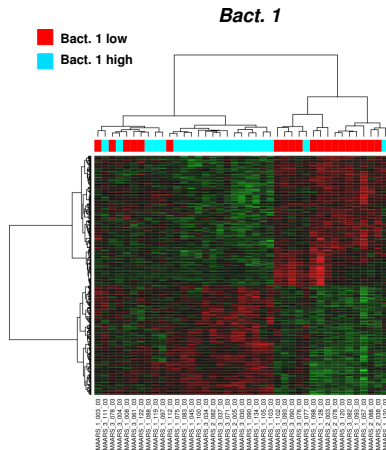$$\rho_s = \max_{a^S, b^S} cor(Xa^S, Yb^S) \tag{1}$$

$$= cor(U^S, V^S) \tag{2}$$
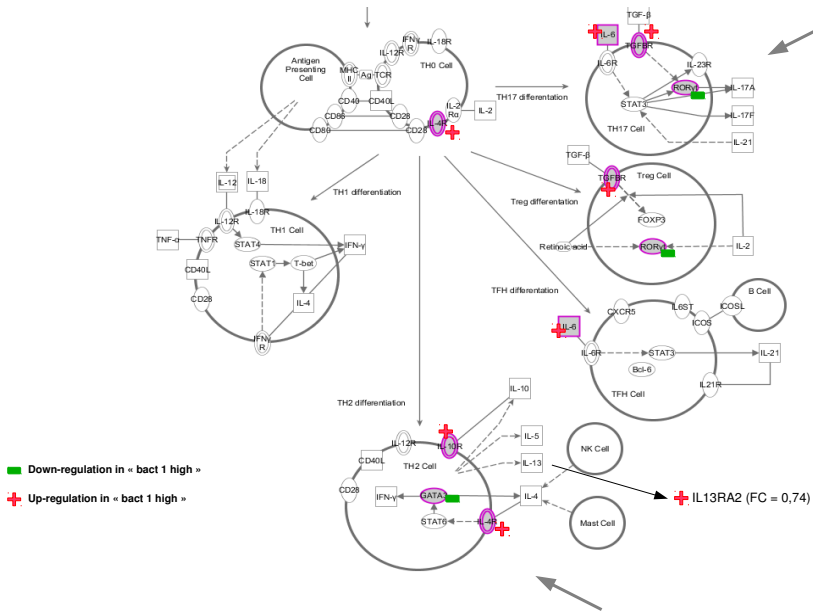
with $cor(U^S, U^K) = cor(V^S, V^K) = 0$ for $S \neq K$.

- $U^S$ and $V^S$: canonical variates
- $\rho_S$: canonical correlations

institut**Curie**
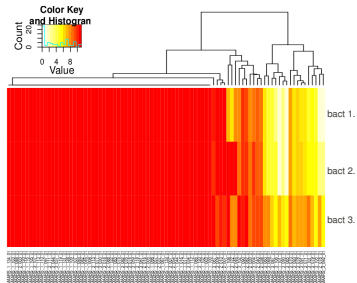*Ensemble, prenons le cancer de vitesse*

Abundance of given bacteria in AD is associated with different transcriptome profiles



*Bact. 1*

■ Bact. 1 low
■ Bact. 1 high

*Bact. 2*

■ Bact. 2 low
■ Bact. 2 high

**Heatmap of differentially expressed genes in bact. 1 low/high patients**

**Heatmap of the 400 « most significant » genes in *bact 2.* low/high patients**

Down-regulation in « bact 1 high »

Up-regulation in « bact 1 high »

IL13RA2 (FC = 0,74)

- 16S data: large-scale count data
  - similar features than RNA-seq data
  - BUT with a higher level of sparsity
- Normalization methods used in RNA-seq analysis
  - perform well on 16S data
  - should be transferred to microbiome research (instead of rarefying)
- No consensus for differential analysis
- Investigate co-occurences/co-exclusions of microbes

Alix, Mahé, Paula, Sol, Caro, Max, Gérôme, Maude, Phil, Lucia, Irit, Vassili, Salvo, Sofia, Anto, Colline, Aurore.

Thank you !