

# Standardised and reproducible analysis of mass spectrometry-based single-cell proteomics data

Replication of the SCoPE2 analysis by Specht et al. 2019

Christophe Vanderaa, Laurent Gatto

Computational Biology Unit (CBIO)  
de Duve Institute  
UCLouvain  
Belgium

18 August 2020

# Outline

## Introduction

scp package

scp showcase

Replication results

Conclusion

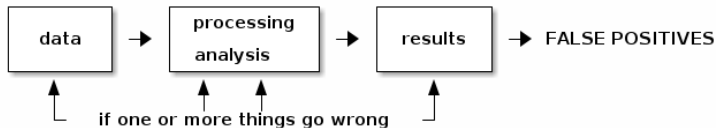
## ► Expectation



## ► Expectation



## ► Reality



- ▶ **Replication-based development** agreement between the developer, the data producer and the user.

Replication is the first step to define sound data infrastructure and principled analysis.

- ▶ SCoPE2 quantifies thousands of proteins x thousands single-cells
- ▶ Full protocole available
- ▶ Full analysis script and data available

- ▶ Contribute a standardised and principled data and analysis that is broadly applicable.
- ▶ Open, transparent and reproducible computational infrastructure to further improve data analysis and interpretation.
- ▶ R and Bioconductor (Huber et al. (2015)) offer an ideal environment to attain these goals.

Implemented in the `scp` package.

# Outline

Introduction

**scp package**

scp showcase

Replication results

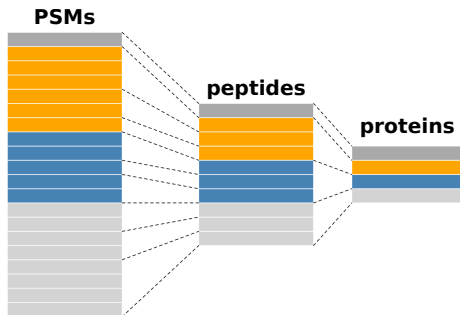
Conclusion



# Data infrastructure: QFeatures<sup>1</sup>

scp package

**QFeatures**: data framework dedicated to manipulate and process MS-based quantitative data.

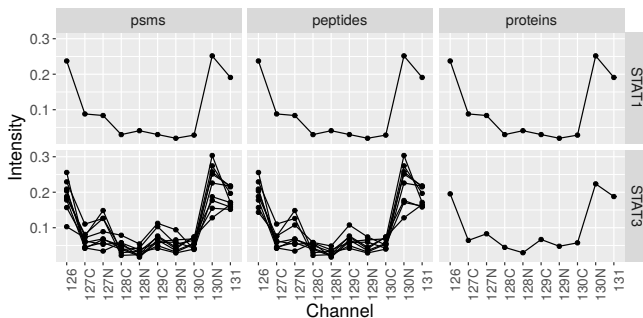
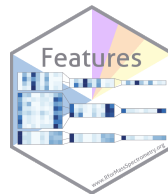


<sup>1</sup>Gatto (2020)

# Data infrastructure: QFeatures<sup>1</sup>

scp package

**QFeatures**: data framework dedicated to manipulate and process MS-based quantitative data.



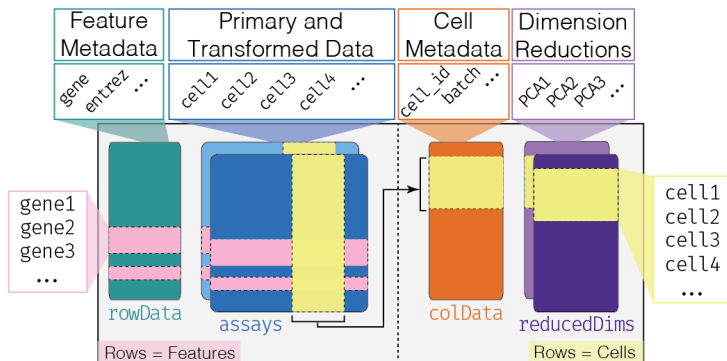
<sup>1</sup>Gatto (2020)

# Data infrastructure: SingleCellExperiment<sup>2,3</sup>

scp package



`SingleCellExperiment`: provides dedicated framework for single-cell data analysis.

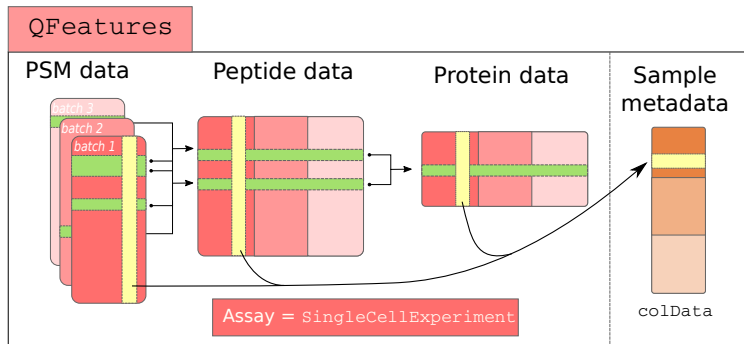


SingleCellExperiment

<sup>2</sup>Lun and Risso (2020)

<sup>3</sup>Amezquita et al. (2019)

`scp` = `SingleCellExperiment` + `QFeatures`



**AND** utility functions dedicated to managing and analyzing SCP data

Load the SCoPE2 dataset called `specht2019v2`<sup>4</sup>

```
1 library(scpdata)
2 data("specht2019v2")
```

## Dataset overview

```
1 show(specht2019v2)
```

```
An instance of class QFeatures containing 179 assays:
[1] 190222S_LCA9_X_FP94AA: SingleCellExperiment with 2823 rows and 11 col...
[2] 190222S_LCA9_X_FP94AB: SingleCellExperiment with 4297 rows and 11 col...
[3] 190222S_LCA9_X_FP94AC: SingleCellExperiment with 4956 rows and 11 col...
...
[177] 191110S_LCB7_X_APNOV16plex2_Set_9: SingleCellExperiment with 4626 r...
[178] peptides: SingleCellExperiment with 9208 rows and 1018 columns
[179] proteins: SingleCellExperiment with 2772 rows and 1018 columns
```

Tabular data (generated by MaxQuant, ProteomeDiscoverer, ...) can be converted to `QFeatures` using the `readSCP()` function.

---

<sup>4</sup>Specht et al. (2019)

Sample metadata common to **all assays** are stored in a single table, the `colData`

Set	Channel	SampleType	lcbatch	sortday	digest
190222S_LCA9_X_FP94AA	RI1	Carrier	LCA9	s8	N
190222S_LCA9_X_FP94AA	RI2	Reference	LCA9	s8	N
190222S_LCA9_X_FP94AA	RI3	Unused	LCA9	s8	N
190222S_LCA9_X_FP94AA	RI4	Macrophage	LCA9	s8	N
190222S_LCA9_X_FP94AA	RI5	Macrophage	LCA9	s8	N
190222S_LCA9_X_FP94AA	RI6	Macrophage	LCA9	s8	N
...	...	...	...	...	...

## 1. Load data

### PSM data

```
[1] 190222S_LCA9_X_FP94AA: SingleCellExperiment with 2823 rows and 11 columns
[2] 190222S_LCA9_X_FP94AB: SingleCellExperiment with 4297 rows and 11 columns
[3] 190222S_LCA9_X_FP94AC: SingleCellExperiment with 4956 rows and 11 columns
...
[177] 191110S_LCB7_X_APNOV16plex2_Set_9: SingleCellExperiment with 4626 rows and 16 columns
```

## 1. Load data

### PSM data

```
[1] 190222S_LCA9_X_FP94AA: SingleCellExperiment with 2823 rows and 11 columns  
[2] 190222S_LCA9_X_FP94AB: SingleCellExperiment with 4297 rows and 11 columns  
[3] 190222S_LCA9_X_FP94AC: SingleCellExperiment with 4956 rows and 11 columns  
...  
[177] 191110S_LCB7_X_APNOV16plex2_Set_9: SingleCellExperiment with 4626 rows and 16 columns
```

2. PSM filtering
3. Expression channel by reference channel division
4. PSM to peptides aggregation
5. Joining sets
6. Single cells filtering based on median CV
7. Normalization
8. Removal of highly missing peptides
9. Log-transformation

### Peptide data

```
[178] peptides: SingleCellExperiment with 9208 rows and 1018 columns
```



## 1. Load data

### PSM data

```
[1] 190222S_LCA9_X_FP94AA: SingleCellExperiment with 2823 rows and 11 columns  
[2] 190222S_LCA9_X_FP94AB: SingleCellExperiment with 4297 rows and 11 columns  
[3] 190222S_LCA9_X_FP94AC: SingleCellExperiment with 4956 rows and 11 columns  
...  
[177] 191110S_LCB7_X_APNOV16plex2_Set_9: SingleCellExperiment with 4626 rows and 16 columns
```

2. PSM filtering
3. Expression channel by reference channel division
4. PSM to peptides aggregation
5. Joining sets
6. Single cells filtering based on median CV
7. Normalization
8. Removal of highly missing peptides
9. Log-transformation

### Peptide data

```
[178] peptides: SingleCellExperiment with 9208 rows and 1018 columns
```

9. Peptides to proteins aggregation
10. Normalization
11. Imputation
12. Batch correction

### Protein data

```
[179] proteins: SingleCellExperiment with 2772 rows and 1018 columns
```

## 1. Load data

### PSM data

```
[1] 190222S_LCA9_X_FP94AA: SingleCellExperiment with 2823 rows and 11 columns  
[2] 190222S_LCA9_X_FP94AB: SingleCellExperiment with 4297 rows and 11 columns  
[3] 190222S_LCA9_X_FP94AC: SingleCellExperiment with 4956 rows and 11 columns  
...  
[i77] 191110S_LCB7_X_APNOV16plex2_Set_9: SingleCellExperiment with 4626 rows and 16 columns
```

2. **PSM filtering**
3. Expression channel by reference channel division
4. PSM to peptides aggregation
5. Joining sets
6. **Single cells filtering based on median CV**
7. Normalization
8. **Removal of highly missing peptides**
9. **Log-transformation**

### Peptide data

```
[i78] peptides: SingleCellExperiment with 9208 rows and 1018 columns
```

9. **Peptides to proteins aggregation**
10. Normalization
11. Imputation
12. **Batch correction**

### Protein data

```
[i79] proteins: SingleCellExperiment with 2772 rows and 1018 columns
```

# Outline

Introduction

scp package

scp showcase

Replication results

Conclusion

Filter out features based on the feature metadata

Example: filter out reverse hits. The filter is applied to the `Reverse` field in the feature metadata

```
1 filterFeatures(specht2019v2,  
2               ~ Reverse != "+")
```

Source code in `QFeatures`

Some QC metrics are not computed by MaxQuant:

- ▶ Sample to carrier ratio: discard samples with intensities higher than expected
- ▶ Peptide FDR: expected proportion of features wrongly assigned to a given peptide
- ▶ Cell median coefficient of variation: reliability of the protein quantification in a cell.

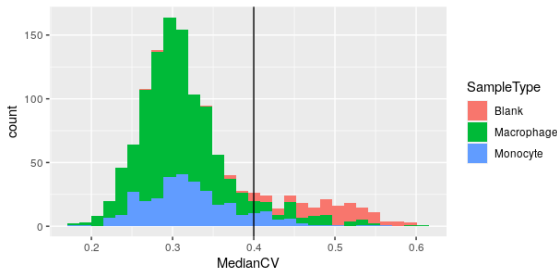
Example:

```
1 computeMedianCV(specht2019v2,
2                 i = "peptides",
3                 proteinCol = "protein",
4                 peptideCol = "peptide",
5                 batchCol = "Set")
```

Source code in scp

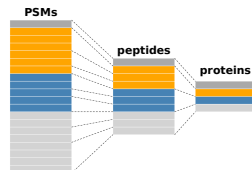
QC metrics are stored in the data set for plotting or subsetting

```
1 library(tidyverse)
2 specht2019v2[["peptides"]] %>%
3   colData %>%
4   data.frame %>%
5   ggplot(aes(x = MedianCV,
6             fill = SampleType)) +
7   geom_histogram() +
8   geom_vline(xintercept = 0.4)
```



Feature aggregation = combine features into a higher-level structure.

Example: aggregate peptides to proteins



- ▶ Combine the quantitative data from multiple peptides to a single protein
- ▶ Store the relationship between the protein and the aggregated peptides

```
1 aggregateFeatures(specht2019v2 ,  
2                   i = "peptides",  
3                   name = "proteins",  
4                   fcol = "protein",  
5                   fun = colMedians, na.rm = TRUE)
```

Source code in `QFeatures`

0's can be either **biological** or **technical** zero. They are better related by NA's.

```
1 zeroIsNA(specht2019v2,  
2         i = "peptides")
```

Remove highly-missing features (e.g.  $\geq 99\%$ )

```
1 filterNA(specht2019v2,  
2         i = "peptides",  
3         pNA = 0.99)
```

Impute missing data

```
1 impute(specht2019v2,  
2        i = "proteins",  
3        method = "knn",  
4        k = 3)
```

Source code in `QFeatures`



Common data transformation can easily be applied such as **log-transformation** or **normalization**.

Example:  $\log_2$ -transformation:

```
1 logTransform(specht2019v2,  
2             i = "peptides",  
3             base = 2,  
4             name = "peptides_log")
```

Source code in `QFeatures`

**Custom function** can be applied to the data set, for example batch correction using 'ComBat'. Three-step procedure:

1. Extract the assay to process

```
1 x <- specht2019v2[["proteins"]]
```

2. Apply the custom function

```
2 assay(x) <- ComBat(assay(x), ...)
```

3. Add the processed assay in the dataset

```
3 addAssay(specht2019v2, x, name = "proteins_batch_corrected")
```

# Outline

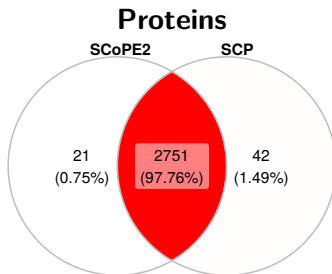
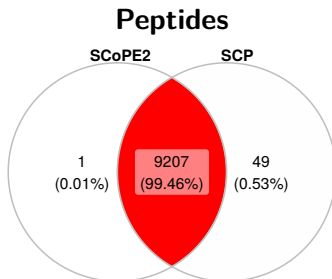
Introduction

scp package

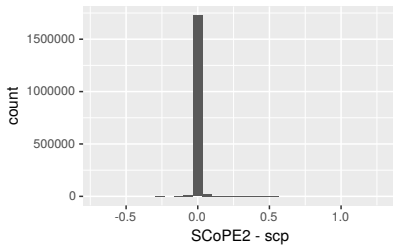
scp showcase

Replication results

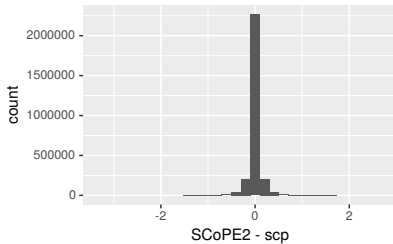
Conclusion



## Peptides



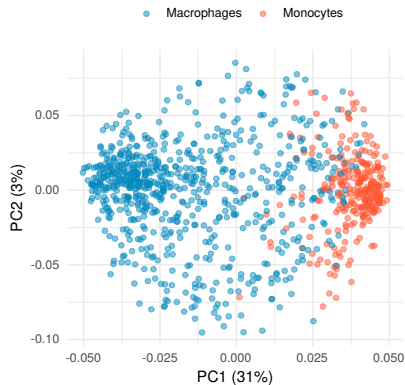
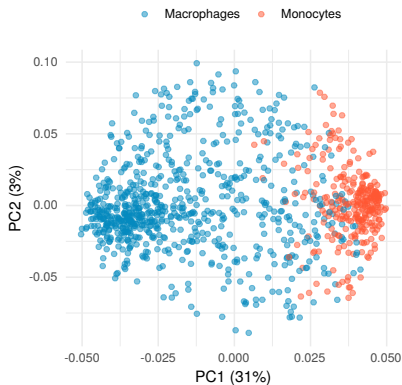
## Proteins



Weighted PCA on the **protein** data

SCoPE2 (Figure 3a in preprint)

scp

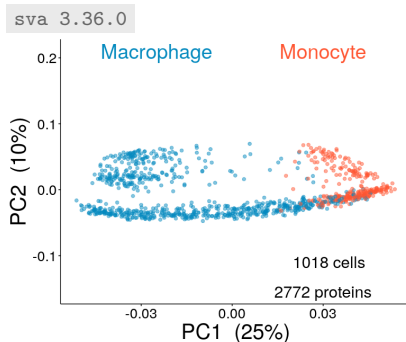
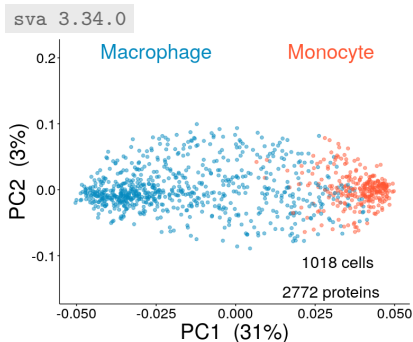


# Note about reproducibility

Replication results

Good software development includes continuous maintenance and improvement **BUT** might impact reproducibility

Example: batch correction using 2 versions of the **ComBat** algorithm  
(**sva** package)



Documenting software version is **essential** for reproducible work

# Outline

Introduction

scp package

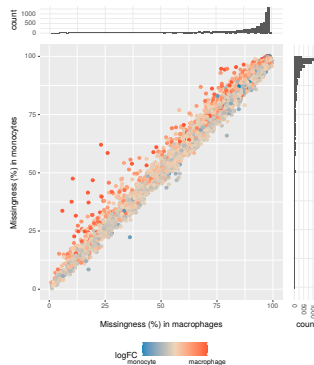
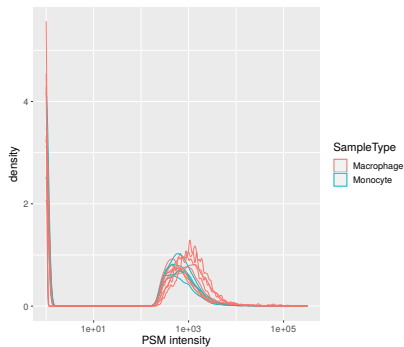
scp showcase

Replication results

Conclusion



- ▶ MS-based single cell proteomics: young field, with many challenges and great progress. `scp` to address the need for principled and reproducible data analysis.
- ▶ `scp` isn't specific to SCoPE2/TMT data, applicable to other LF protocols such as nanoPOTS (Williams et al. (2020); Cong et al. (2020)).
- ▶ `scp` and `SingleCellExperiment`: same infrastructure for single cell proteomics and RNA sequencing.
- ▶ Tool for novel computational developments.



## Resources

- ▶ `scp`: <http://UClouvain-CBIO.github.io/scp>
- ▶ `scpdata`: coming soon
- ▶ `QFeatures`: <http://rformassspectrometry.org>
- ▶ `SingleCellExperiment`: Bioconductor
- ▶ **Slides**: <http://bit.ly/2020SCP> under CC-BY SA

## Resources

- ▶ `scp`: <http://UClouvain-CBIO.github.io/scp>
- ▶ `scpdata`: coming soon
- ▶ `QFeatures`: <http://rformassspectrometry.org>
- ▶ `SingleCellExperiment`: Bioconductor
- ▶ **Slides**: <http://bit.ly/2020SCP> under CC-BY SA

## Acknowledgements

- ▶ Nikolai Slavov, Harrison Specht, Ed Emmott.
- ▶ Fonds National de la Recherche Scientifique (FNRS)
- ▶ **Thank you for your attention**

# References I

- Robert A Amezcua, Aaron T L Lun, Etienne Becht, Vince J Carey, Lindsay N Carpp, Ludwig Geistlinger, Federico Martini, Kevin Rue-Albrecht, Davide Risso, Charlotte Soneson, Levi Waldron, Hervé Pagès, Mike L Smith, Wolfgang Huber, Martin Morgan, Raphael Gottardo, and Stephanie C Hicks. Orchestrating single-cell analysis with bioconductor. *Nat. Methods*, pages 1–9, December 2019.
- Yongzheng Cong, Yiran Liang, Khatereh Motamedchaboki, Romain Huguet, Thy Truong, Rui Zhao, Yufeng Shen, Daniel Lopez-Ferrer, Ying Zhu, and Ryan T Kelly. Improved single cell proteome coverage using Narrow-Bore packed NanoLC columns and ultrasensitive mass spectrometry. *Anal. Chem.*, January 2020.
- Laurent Gatto. *QFeatures: Quantitative features for mass spectrometry data*, 2020. URL <https://github.com/RforMassSpectrometry/QFeatures>. R package version 0.7.0.
- W Huber, V J Carey, R Gentleman, S Anders, M Carlson, B S Carvalho, H C Bravo, S Davis, L Gatto, T Girke, R Gottardo, F Hahne, K D Hansen, R A Irizarry, M Lawrence, M I Love, J MacDonald, V Obenchain, A K Oleś, H Pagès, A Reyes, P Shannon, G K Smyth, D Tenenbaum, L Waldron, and M Morgan. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat. Methods*, 12(2): 115–121, January 2015.
- Aaron Lun and Davide Risso. *SingleCellExperiment: S4 Classes for Single Cell Data*, 2020. R package version 1.10.1.

# References II

- Harrison Specht, Edward Emmott, Toni Koller, and Nikolai Slavov. High-throughput single-cell proteomics quantifies the emergence of macrophage heterogeneity. June 2019.
- Sarah M Williams, Andrey V Liyu, Chia-Feng Tsai, Ronald J Moore, Daniel J Orton, William B Chrisler, Matthew J Gaffrey, Tao Liu, Richard D Smith, Ryan T Kelly, Ljiljana Paša-Tolić, and Ying Zhu. Automated coupling of nanodroplet sample preparation with liquid Chromatography-Mass spectrometry for High-Throughput Single-Cell proteomics. *Anal. Chem.*, July 2020.