

# Replicate Specht et al. 2019 mass spectrometry-based single-cell proteomics analysis

Laurent Gatto, Christophe Vanderaa

CBIO, de Duve Institute, UCLouvain

18 August 2020

# Outline

Introduction

Data framework

Standardized workflow

Replicating SCoPE2

Conclusion

MS-SCP: Mass spectrometry-based single-cell proteomics

MS-SCP consist of shotgun proteomics at single-cell level

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

**BUT**

Lack of standardized analysis software

Provide a suite of software package dedicated to MS-SCP that fulfill:

- ▶ User-friendly
- ▶ Computationally efficient
- ▶ Modularity: integrate other software packages
- ▶ Promote reproducibility
- ▶ Platform-independent
- ▶ Free of charge

R/Bioconductor is an ideal environment

# Outline

Introduction

Data framework

Standardized workflow

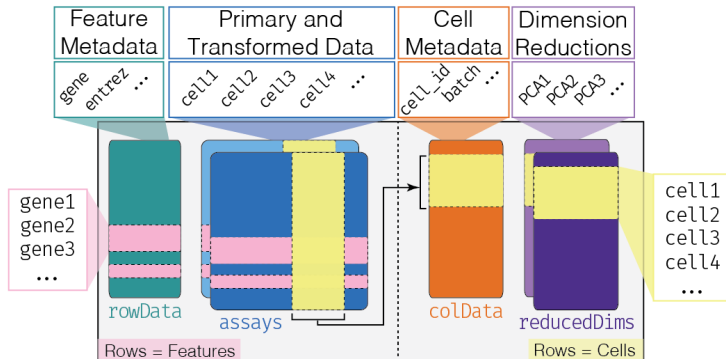
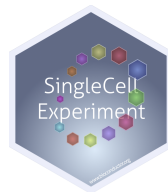
Replicating SCoPE2

Conclusion

# SingleCellExperiment

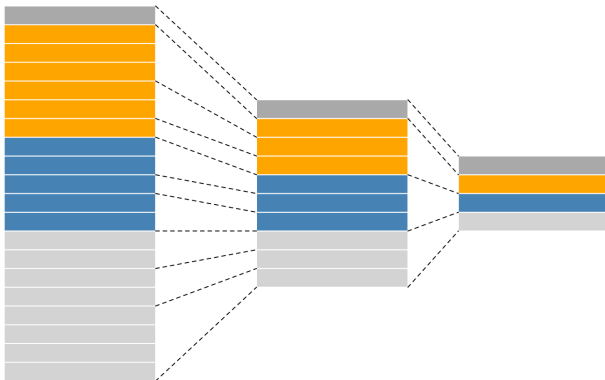
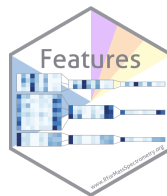
Data framework

SingleCellExperiment: provides dedicated framework for single-cell data analysis.  
Available on Bioconductor.

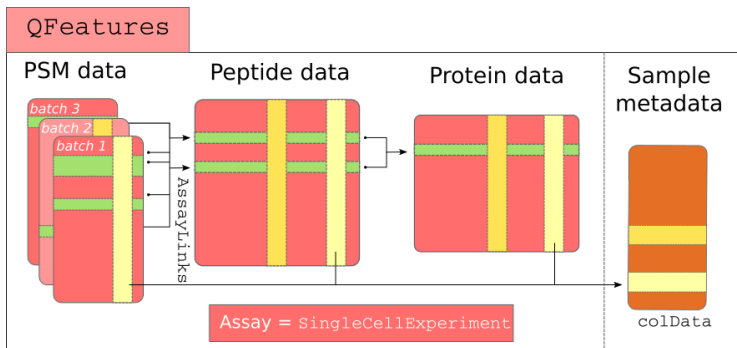


SingleCellExperiment

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



MS-SCP data framework = SingleCellExperiment + QFeatures





scpdata: distributes published MS-SCP datasets (e.g. SCoPE2 dataset)

scp: provides functionality for manipulating the MS-SCP data structure

# Outline

Introduction

Data framework

Standardized workflow

Replicating SCoPE2

Conclusion

Load the SCoPE2 dataset called `specht2019v2`

```
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
```

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`

Interesting metrics for MS-SCP quality control:

- ▶ Sample to carrier ratio: ratio of the carrier channel intensity signal over the sample channel intensity
- ▶ Peptide FDR<sup>1</sup>: expected rate of wrongly assigned features to a given peptide
- ▶ Cell median CV<sup>2</sup>: reliability of the protein quantification summarized over each cell.

Example:

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

Source code in `scp`

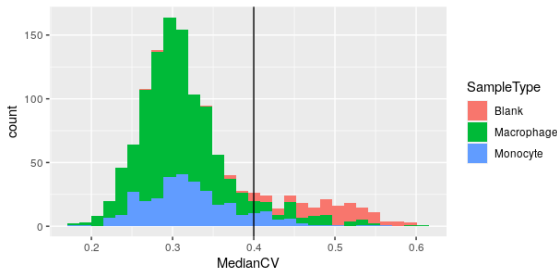
---

<sup>1</sup>false discovery rate

<sup>2</sup>coefficient of variation

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 includes 2 steps:

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

Example: aggregate peptides to proteins

```
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 relaced by NA's.

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

Features containing too many missing data (e.g.  $\geq 99\%$ ) should be removed

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

Source code in `QFeatures`



Common data transformation can easily be applied:

- ▶ Normalization
- ▶ Log-transformation
- ▶ Imputation

Example:  $\log_2$ -transformation:

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

Source code in `QFeatures`

Some custom function can be applied to the data set too.

Example: batch correction using `sva::ComBat`. First, extract the data to correct

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

Build the correction matrix and apply the ComBat algorithm

```
1 batch <- colData(sce)$Set
2 model <- model.matrix(~ SampleType, data = colData(sce))
3 assay(sce) <- ComBat(dat = assay(sce),
4                       batch = batch,
5                       mod = model)
```

Add the corrected protein to the dataset and keep feature relationships

```
1 addAssay(specht2019v2,
2          sce,
3          name = "proteins_batchC") %>%
4 addAssayLinkOneToOne(from = "proteins",
5                       to = "proteins_batchC")
```

# Outline

Introduction

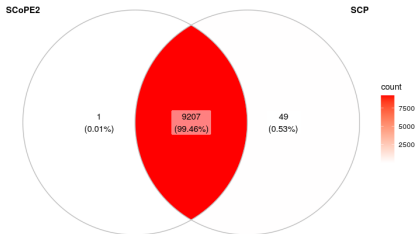
Data framework

Standardized workflow

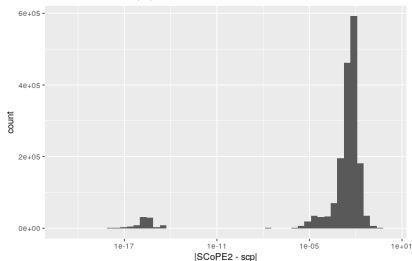
Replicating SCoPE2

Conclusion

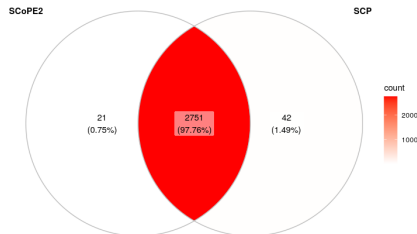
### Peptides



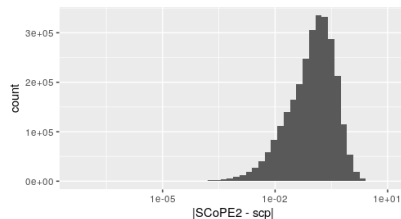
Benchmark of the peptide data



### Proteins



Benchmark of the protein data



# Replicate figures from SCoPE2 (1)

Replicating SCoPE2

# Replicate figures from SCoPE2 (2)

Replicating SCoPE2

# Replicate figures from SCoPE2 (3)

Replicating SCoPE2

# Outline

Introduction

Data framework

Standardized workflow

Replicating SCoPE2

Conclusion



- ▶ `scp` package suite provides a standardized environment for performing MS-SCP data analysis
- ▶ Flexibly reproduce existing analyses from different groups or protocols (multiplex vs label free)

## Advantages:

- ▶ Allow automation of the analysis
- ▶ Facilitate new computational developments
- ▶ Promotes reproducibility
- ▶ Increases field visibility

## Packages

- ▶ `scp`: GitHub repository `UClouvain-CBIO/scp`
- ▶ `scpdata`: GitHub repository `UClouvain-CBIO/scpdata`
- ▶ `QFeatures`: GitHub repository `rformassspectrometry/QFeatures`
- ▶ `SingleCellExperiment`: Bioconductor

## SCoPE2 reproduction vignette

Available at...

## Slides and source code

Available at...

