# Package: scpdata (via r-universe)

August 31, 2024

Type Package

Title Single-Cell Proteomics Data Package

**Version** 1.13.0

**Description** The package disseminates mass spectrometry (MS)-based single-cell proteomics (SCP) datasets. The data were collected from published work and formatted using the `scp` data structure. The data sets contain quantitative information at spectrum, peptide and/or protein level for single cells or minute sample amounts.

**Depends** R (>= 4.2.0), QFeatures, ExperimentHub

Imports utils, AnnotationHub, SingleCellExperiment, S4Vectors

**Suggests** scp, magrittr, dplyr, knitr, BiocStyle, BiocCheck, rmarkdown, testthat

**biocViews** ExperimentData, ExpressionData, ExperimentHub, ReproducibleResearch, MassSpectrometryData, Proteome, SingleCellData, PackageTypeData

License GPL-2

**Encoding** UTF-8

LazyData false

VignetteBuilder knitr

**Roxygen** list(markdown = TRUE)

RoxygenNote 7.3.1

Repository https://uclouvain-cbio.r-universe.dev

RemoteUrl https://github.com/uclouvain-cbio/scpdata

RemoteRef HEAD

**RemoteSha** a2c8b038c1284274a7ebf4be80a59a05b354113f

brunner2022

# **Contents**

brunner2022	 2
cong2020AC	 4
derks2022	 6
dou2019_boosting	 8
dou2019_lysates	 10
dou2019_mouse	 12
gregoire2023_mixCTRL	 14
guise2024	 16
khan2023	 18
leduc2022	 20
leduc2022_plexDIA	 20
leduc2022_pSCoPE	 22
liang2020_hela	 25
petrosius2023_AstralAML	
petrosius2023_mES	 28
schoof2021	
scpdata	 32
specht2019v2	
specht2019v3	
williams2020_lfq	
williams2020_tmt	
woo2022_lung	
woo2022_macrophage	
zhu2018MCP	
zhu2018NC_hela	
zhu2018NC_islets	
zhu2018NC_lysates	
zhu2019EL	 50
	53
	53

brunner2022

Brunner et al. 2022 (Mol. Syst. Biol.): cell cycle state study

# Description

Index

Single cell proteomics data acquired by the Mann Lab using a newly designed timsTOF instrument, referred to as timsTOF-SCP. The dataset contains quantitative information from single-cells blocked at 4 cell cycle stages: G1, G1-S, G2, G2-M. The data was acquired using a label-free sample preparation protocole combined to a data independent (DIA) acquisition mode.

# Usage

brunner2022

brunner2022 3

#### **Format**

A QFeatures object with 435 assays, each assay being a SingleCellExperiment object.

 Assay 1-434: DIA-NN main output report table split for each acquisition run. Since each run acquires 1 single cell, each assay contains a single column. It contains the results of the spectrum identification and quantification.

• protein: DIA-NN protein group matrix, containing normalised quantities for 2476 protein groups in 434 single cells. Proteins are filtered at 1% FDR, using global q-values for protein groups and both global and run-specific q-values for precursors.

The colData(brunner2022()) contains cell type annotations and batch annotations. The description of the rowData fields for the different assays can be found in the DIA-NN documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- **Cell isolation**: cells were detached with trypsin treatment, followed by strong pipetting, and isolate using FACS.
- Sample preparation: cell lysis by freeze-heat followed by sonication, overnight protein digestion with trypsin/lysC mix and desalting using EvoTips trap column (EvoSep)
- Separation: online EvoSep One LC system using a 5 cm x 75 μm ID column with 1.9μm C18 beads (EvoSep) at 100nL/min flow rate.
- **Ionization**: 10µm ID zero dead volume electrospray emitter (Bruker Daltonik) + nanoelectrospray ion source (Captive spray, Bruker Daltonik)
- Mass spectrometry: DIA PASEF mode. Correlation between IM and m/z was used to synchronize the elution of precursors from each IM scan with the quadrupole isolation window. Five consecutive diaPASEF cycles. The collision energy was ramped linearly as a function of the IM from 59 eV at 1/K0=1.6 Vs cm<sup>2</sup> to 20 eV at 1/K0=0.6 Vs cm<sup>2</sup>.
- Data analysis: DIA-NN (1.8).

#### **Data collection**

The data were collected from the PRIDE repository in the DIANN1.8\_SingleCells\_CellCycle.zip file.

We loaded the DIA-NN main report table and generated a sample annotation table based on the MS file names. We next combined the sample annotation and the DIANN tables into a QFeatures object following the scp data structure. We loaded the proteins group matrix as a SingleCellExperiment object, fixed ambiguous protein group names, and added the protein data as a new assay and link the precursors to proteins using the Protein. Group variable from the rowData.

#### Source

The data were downloaded from PRIDE repository with accession ID PXD024043.

4 cong2020AC

#### References

Brunner, Andreas-David, Marvin Thielert, Catherine Vasilopoulou, Constantin Ammar, Fabian Coscia, Andreas Mund, Ole B. Hoerning, et al. 2022. "Ultra-High Sensitivity Mass Spectrometry Quantifies Single-Cell Proteome Changes upon Perturbation." Molecular Systems Biology 18 (3): e10798. Link to article

# **Examples**

brunner2022()

cong2020AC

Cong et al. 2020 (Ana. Chem.): HeLa single cells

# Description

Single-cell proteomics using the nanoPOTS sample processing device in combination with ultranarrow-bore (20um i.d.) packed-column LC separations and the Orbitrap Eclipse Tribrid MS. The dataset contains label-free quantitative information at PSM, peptide and protein level. The samples are single Hela cells. Bulk samples (100 and 20 cells) were also included in the experiment to increase the idendification rate thanks to between-run matching (cf MaxQuant).

# Usage

cong2020AC

#### **Format**

A QFeatures object with 9 assays, each assay being a SingleCellExperiment object:

- 100/20 HeLa cells: 2 assays containing PSM data for a bulk of 100 or 20 HeLa cells, respectively.
- Blank: assay containing the PSM data for a blank sample
- Single cell X: 4 assays containing PSM data for a single cell. The X indicates the replicate number.
- peptides: quantitative data for 12590 peptides in 7 samples (all runs combined).
- proteins: quantitative data for 1801 proteins in 7 samples (all runs combined).

Sample annotation is stored in colData(cong2020AC()).

cong2020AC 5

### **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: The HeLa cells were diluted and aspired using a microcapillary with a pulled tip.
- Sample preparation performed using the nanoPOTs device. Protein extraction using RapiGest (+ DTT) + alkylation (IAA) + Lys-C digestion + cleave RapiGest (formic acid)
- **Separation**: UltiMate 3000 RSLCnano pump with a home-packed nanoLC column (60cm x 20um i.d.; approx. 20 nL/min)
- **Ionization**: ESI (2,000V; Nanospray Flex)
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Eclipse. MS1 settings: accumulation time = 246ms; resolution = 120,000; AGC = 1E6. MS/MS settings depend on quantity. All: AGC = 1E5. 20-100 cels: accumulation time = 246ms; resolution = 120,000. Single cells: accumulation time = 500ms; resolution = 240,000.
- Data analysis: MaxQuant (v1.6.3.3) + Excel

#### **Data collection**

The PSM, peptide and protein data were collected from the PRIDE repository (accession ID: PXD016921). We downloaded the evidence.txt file containing the PSM identification and quantification results. The sample annotation was inferred from the samples names. The data were then converted to a QFeatures object using the scp::readSCP() function.

The peptide data were processed similarly from the peptides.txt file. The quantitative column names were adpated to match the PSM data. The peptide data were added to QFeatures object and link between the features were stored.

The protein data were similarly processed from the proteinGroups.txt file. The quantitative column names were adapted to match the PSM data. The peptide data were added to QFeatures object and link between the features were stored.

#### **Source**

All files can be downloaded from the PRIDE repository PXD016921. The source link is: ftp://ftp.pride.ebi.ac.uk/pride/data/ar

### References

Cong, Yongzheng, Yiran Liang, Khatereh Motamedchaboki, Romain Huguet, Thy Truong, Rui Zhao, Yufeng Shen, Daniel Lopez-Ferrer, Ying Zhu, and Ryan T. Kelly. 2020. "Improved Single-Cell Proteome Coverage Using Narrow-Bore Packed NanoLC Columns and Ultrasensitive Mass Spectrometry." Analytical Chemistry, January. (link to article).

### **Examples**

cong2020AC()

6 derks2022

derks2022	Derks et al. 2022 - plexDIA (Nat. Biotechnol.): PDAC vs melanoma cells vs monocytes

### **Description**

Single cell proteomics data acquired by the Slavov Lab using the plexDIA protocol. It contains quantitative information from pancreatic ductal acinar cells (PDAC; HPAF-II), melanoma cells (WM989-A6-G3) and monocytes (U-937) at precursor and protein level. The each run acquired 3 samples thanks to mTRAQ multiplexing.

# Usage

derks2022

#### **Format**

A QFeatures object with 66 assays, each assay being a SingleCellExperiment object. The assays either hold the DIA-NN main output report table or the DIA-NN MS1 extracted signal table. The DIA-NN main output report table contains the results of the spectrum identification and quantification. The DIA-NN MS1 extracted signal table contains quantification for all mTRAQ channels if its precursors was identified in at least one of the channels, regardless of whether there is sufficient evidence in those channels at 1% FDR.

The data is composed of three datasets

- 1. **Bulk**: dataset containing bulk (100-cell) data acquired using a Q-Exactive mass spectrometer. Assays 1-3 contain data from the DIA-NN main output report; assay 4 is the DIA-NN MS1 extracted signal.
- 2. **tims**: dataset containing single-cell data acquired using a timsTOF-SCP mass spectrometer. Assays 5-15 contain data from the DIA-NN main output report; assay 16 is the DIA-NN MS1 extracted signal.
- 3. **qe**: dataset containing single-cell data acquired using a Q-Exactive mass spectrometer. Assays 17-64 contain data from the DIA-NN main output report; assay 65 is the DIA-NN MS1 extracted signal.

The last assay proteins contains the processed protein data table generated by the authors.

The colData(derks2022()) contains cell type annotations and batch annotations. The description of the rowData fields for the different assays can be found in the DIA-NN documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

• Cell isolation: CellenONE cell sorting.

derks2022 7

• Sample preparation performed using the improved SCoPE2 protocol using the CellenONE liquid handling system. nPOP cell lysis (DMSO) + trypsin digestion + mTRAQ (3plex) labelling and pooling. A target library was generated as well to perform prioritized DDA (Huffman et al. 2022) using MaxQuant.Live (2.0.3).

- Separation: bulk online nLC (Dionex UltiMate 3000 UHPLC) with a 25 cm × 75 μm IonOpticks Aurora Series UHPLC column (AUR2-25075C18A), 200nL/min. qe online nLC (Dionex UltiMate 3000 UHPLC) with a 15 cm × 75 μm IonOpticks Aurora Series UHPLC column (AUR2-15075C18A), 200nL/min. tims nanoElute liquid chromatography system (Bruker Daltonics) using a 25 cm × 75 μm, 1.6-μm C18 (AUR2-25075C18A-CSI, IonOpticks).
- Ionization: ESI.
- Mass spectrometry: cf article.
- Data analysis: DIA-NN (1.8.1 beta 16).

#### **Data collection**

The data were collected from a shared Google Drive folder that is accessible from the SlavovLab website (see Source section).

For each dataset separately, we combined the sample annotation and the DIANN tables in a QFeatures object following the scp data structure. We then combined the three datasets in a single QFeatures object. We load the proteins table processed by the authors as a SingleCellExperiment object and adapted the sample names to match those in the QFeatures object. We added the protein data as a new assay and link the precursors to proteins using the Protein. Group variable from the rowData.

#### Source

The data were downloaded from the Slavov Lab website. The raw data and the quantification data can also be found in the massIVE repository MSV000089093.

#### References

Derks, Jason, Andrew Leduc, Georg Wallmann, R. Gray Huffman, Matthew Willetts, Saad Khan, Harrison Specht, Markus Ralser, Vadim Demichev, and Nikolai Slavov. 2022. "Increasing the Throughput of Sensitive Proteomics by plexDIA." Nature Biotechnology, July. Link to article

# **Examples**

derks2022()

8 dou2019\_boosting

dou2019\_boosting

Dou et al. 2019 (Anal. Chem.): testing boosting ratios

# **Description**

Single-cell proteomics using nanoPOTS combined with TMT isobaric labeling. It contains quantitative information at PSM and protein level. The cell types are either "Raw" (macrophage cells), "C10" (epihelial cells), or "SVEC" (endothelial cells). Each cell is replicated 2 or 3 times. Each cell type was run using 3 levels of boosting: 0 ng (no boosting), 5 ng or 50 ng. When boosting was applied, 1 reference well and 1 boosting well were added, otherwise 1 empty well was added. Each boosting setting (0ng, 5ng, 50ng) was run in duplicate.

# Usage

dou2019\_boosting

#### **Format**

A QFeatures object with 7 assays, each assay being a SingleCellExperiment object:

- Boosting\_X\_run\_Y: PSM data with 10 columns corresponding to the TMT-10plex channels. The X indicates the boosting amount (0ng, 5ng or 50ng) and Y indicates the run number (1 or 2).
- peptides: peptide data containing quantitative data for 13,462 peptides in 60 samples (run 1 and run 2 combined).
- proteins: protein data containing quantitative data for 1436 proteins and 60 samples (all runs combined).

Sample annotation is stored in colData(dou2019\_boosting()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- **Cell isolation**: single-cells from the three murine cell lines were isolated using FACS (BD Influx II cell sorter). Boosting sample were prepared (presumably in bulk) from 1:1:1 mix of the three cell lines.
- Sample preparation performed using the nanoPOTs device. Protein extraction (DMM + TCEAP) + alkylation (IAA) + Lys-C digestion + trypsin digestion + TMT-10plex labeling and pooling.
- Separation: nanoLC (Dionex UltiMate with an in-house packed 50cm x 30um LC columns; 50nL/min)
- Ionization: ESI (2,000V)
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid (MS1 accumulation time = 50ms; MS1 resolution = 120,000; MS1 AGC = 1E6; MS2 accumulation time = 246ms; MS2 resolution = 60,000; MS2 AGC = 1E5)
- Data analysis: MS-GF+ + MASIC (v3.0.7111) + RomicsProcessor (custom R package)

dou2019\_boosting

#### **Data collection**

The PSM data were collected from the MassIVE repository MSV000084110 (see Source section). The downloaded files are:

- Boosting\_\*ng\_run\_\*\_msgfplus.mzid: the MS-GF+ identification result files.
- Boosting\_\*ng\_run\_\*\_ReporterIons.txt: the MASIC quantification result files.

For each batch, the quantification and identification data were combined based on the scan number (common to both data sets). The combined datasets for the different runs were then concatenated feature-wise. To avoid data duplication due to ambiguous matching of spectra to peptides or ambiguous mapping of peptides to proteins, we combined ambiguous peptides to peptides groups and proteins to protein groups. Feature annotations that are not common within a peptide or protein group are are separated by a;. The sample annotation table was manually created based on the available information provided in the article. The data were then converted to a QFeatures object using the scp::readSCP() function.

We generated the peptide data. First, we removed PSM matched to contaminants or decoy peptides and ensured a 1% FDR. We aggregated the PSM to peptides based on the peptide (or peptide group) sequence(s) using the median PSM instenity. The peptide data for the different runs were then joined in a single assay (see QFeatures::joinAssays), again based on the peptide sequence(s). We then removed the peptide groups. Links between the peptide and the PSM data were created using QFeatures::addAssayLink. Note that links between PSM and peptide groups are not stored.

The protein data were downloaded from Supporting information section from the publisher's website (see Sources). The data is supplied as an Excel file ac9b03349\_si\_004.xlsx. The file contains 7 sheets from which we took the 2nd, 4th and 6th sheets (named 01 - No Boost raw data, 03 - 5ng boost raw data, 05 - 50ng boost raw data, respectively). The sheets contain the combined protein data for the duplicate runs given the boosting amount. We joined the data for all boosting ration based on the protein name and converted the data to a SingleCellExperiment object. We then added the object as a new assay in the QFeatures dataset (containing the PSM data). Links between the proteins and the corresponding PSM were created. Note that links to protein groups are not stored.

### **Source**

The PSM data can be downloaded from the massIVE repository MSV000084110. FTP link: ftp://massive.ucsd.edu/MSV0000

The protein data can be downloaded from the ACS Publications website (Supporting information section).

#### References

Dou, Maowei, Geremy Clair, Chia-Feng Tsai, Kerui Xu, William B. Chrisler, Ryan L. Sontag, Rui Zhao, et al. 2019. "High-Throughput Single Cell Proteomics Enabled by Multiplex Isobaric Labeling in a Nanodroplet Sample Preparation Platform." Analytical Chemistry, September (link to article).

#### See Also

dou2019\_lysates, dou2019\_mouse

10 dou2019\_lysates

# **Examples**

dou2019\_boosting()

dou2019\_lysates

Dou et al. 2019 (Anal. Chem.): HeLa lysates

# Description

Single-cell proteomics using nanoPOTS combined with TMT multiplexing. It contains quantitative information at PSM and protein level. The samples are commercial Hela lysates diluted to single-cell amounts (0.2 ng). The boosting wells contain the same digest but at higher amount (10 ng).

# Usage

dou2019\_lysates

#### **Format**

A QFeatures object with 3 assays, each assay being a SingleCellExperiment object:

- Hela\_run\_1: PSM data with 10 columns corresponding to the TMT-10plex channels. Columns hold quantitative information for HeLa lysate samples (either 0, 0.2 or 10ng). This is the data for run 1.
- Hela\_run\_1: PSM data with 10 columns corresponding to the TMT-10plex channels. Columns hold quantitative information for HeLa lysate samples (either 0, 0.2 or 10ng). This is the data for run 2.
- peptides: peptide data containing quantitative data for 13,934 peptides in 20 samples (run 1 and run 2 combined).
- proteins: protein data containing quantitative data for 1641 proteins in 20 samples (run 1 and run 2 combined).

Sample annotation is stored in colData(dou2019\_lysates()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: commercially available HeLa protein digest (Thermo Scientific).
- Sample preparation performed using the nanoPOTs device. Protein extraction (DMM + TCEAP) + alkylation (IAA) + Lys-C digestion + trypsin digestion + TMT-10plex labeling and pooling.
- **Separation**: nanoLC (Dionex UltiMate with an in-house packed 50cm x 30um LC columns; 50nL/min)

dou2019\_lysates

- **Ionization**: ESI (2,000V)
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid (MS1 accumulation time = 50ms; MS1 resolution = 120,000; MS1 AGC = 1E6; MS2 accumulation time = 246ms; MS2 resolution = 60,000; MS2 AGC = 1E5)

• Data analysis: MS-GF+ + MASIC (v3.0.7111) + RomicsProcessor (custom R package)

#### **Data collection**

The PSM data were collected from the MassIVE repository MSV000084110 (see Source section). The downloaded files are:

- Hela\_run\_\*\_msgfplus.mzid: the MS-GF+ identification result files
- Hela\_run\_\*\_ReporterIons.txt: the MASIC quantification result files

For each batch, the quantification and identification data were combined based on the scan number (common to both data sets). The combined datasets for the different runs were then concatenated feature-wise. To avoid data duplication due to ambiguous matching of spectra to peptides or ambiguous mapping of peptides to proteins, we combined ambiguous peptides to peptides groups and proteins to protein groups. Feature annotations that are not common within a peptide or protein group are are separated by a; The sample annotation table was manually created based on the available information provided in the article. The data were then converted to a QFeatures object using the scp::readSCP() function.

We generated the peptide data. First, we removed PSM matched to contaminants or decoy peptides and ensured a 1% FDR. We aggregated the PSM to peptides based on the peptide (or peptide group) sequence(s) using the median PSM instenity. The peptide data for the different runs were then joined in a single assay (see QFeatures::joinAssays), again based on the peptide sequence(s). We then removed the peptide groups. Links between the peptide and the PSM data were created using QFeatures::addAssayLink. Note that links between PSM and peptide groups are not stored.

The protein data were downloaded from Supporting information section from the publisher's website (see Sources). The data is supplied as an Excel file ac9b03349\_si\_003.xlsx. The file contains 7 sheets from which we only took the sheet 6 (named 5 - Run 1 and 2 raw data) with the combined protein data for the two runs. We converted the data to a SingleCellExperiment object and added the object as a new assay in the QFeatures dataset (containing the PSM data). Links between the proteins and the peptides were created. Note that links to protein groups are not stored.

#### Source

The PSM data can be downloaded from the massIVE repository MSV000084110. FTP link: ftp://massive.ucsd.edu/MSV00007 The protein data can be downloaded from the ACS Publications website (Supporting information section).

#### References

Dou, Maowei, Geremy Clair, Chia-Feng Tsai, Kerui Xu, William B. Chrisler, Ryan L. Sontag, Rui Zhao, et al. 2019. "High-Throughput Single Cell Proteomics Enabled by Multiplex Isobaric Labeling in a Nanodroplet Sample Preparation Platform." Analytical Chemistry, September (link to article).

dou2019\_mouse

# See Also

dou2019\_mouse, dou2019\_boosting

# **Examples**

dou2019\_lysates()

dou2019\_mouse

Dou et al. 2019 (Anal. Chem.): murine cell lines

#### **Description**

Single-cell proteomics using nanoPOTS combined with TMT isobaric labeling. It contains quantitative information at PSM and protein level. The cell types are either "Raw" (macrophage cells), "C10" (epihelial cells), or "SVEC" (endothelial cells). Out of the 132 wells, 72 contain single cells, corresponding to 24 C10 cells, 24 RAW cells, and 24 SVEC. The other wells are either boosting channels (12), empty channels (36) or reference channels (12). Boosting and reference channels are balanced (1:1:1) mixes of C10, SVEC, and RAW samples at 5 ng and 0.2 ng, respectively. The different cell types where evenly distributed across 4 nanoPOTS chips. Samples were 11-plexed with TMT labeling.

# Usage

dou2019\_mouse

# Format

A QFeatures object with 13 assays, each assay being a SingleCellExperiment object:

- Single\_Cell\_Chip\_X\_Y: PSM data with 11 columns corresponding to the TMT channels (see Notes). The X indicates the chip number (from 1 to 4) and Y indicates the row name on the chip (from A to C).
- peptides: peptide data containing quantitative data for 15,492 peptides in 132 samples (run 1 and run 2 combined).
- proteins: protein data containing quantitative data for 2331 proteins in 132 samples (all runs combined).

Sample annotation is stored in colData(dou2019\_mouse()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

• **Cell isolation**: single-cells from the three murine cell lines were isolated using FACS (BD Influx II cell sorter).

dou2019\_mouse 13

• Sample preparation performed using the nanoPOTs device. Protein extraction (DMM + TCEAP) + alkylation (IAA) + Lys-C digestion + trypsin digestion + TMT-10plex labeling and pooling.

- **Separation**: nanoLC (Dionex UltiMate with an in-house packed 50cm x 30um LC columns; 50nL/min)
- **Ionization**: ESI (2,000V)
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid (MS1 accumulation time = 50ms; MS1 resolution = 120,000; MS1 AGC = 1E6; MS2 accumulation time = 246ms; MS2 resolution = 60,000; MS2 AGC = 1E5)
- Data analysis: MS-GF+ + MASIC (v3.0.7111) + RomicsProcessor (custom R package)

#### **Data collection**

The PSM data were collected from the MassIVE repository MSV000084110 (see Source section). The downloaded files are:

- Single\_Cell\_Chip\_\*\_\*\_msgfplus.mzid: the MS-GF+ identification result files.
- Single\_Cell\_Chip\_\*\_\*\_ReporterIons.txt: the MASIC quantification result files.

For each batch, the quantification and identification data were combined based on the scan number (common to both data sets). The combined datasets for the different runs were then concatenated feature-wise. To avoid data duplication due to ambiguous matching of spectra to peptides or ambiguous mapping of peptides to proteins, we combined ambiguous peptides to peptides groups and proteins to protein groups. Feature annotations that are not common within a peptide or protein group are are separated by a;. The sample annotation table was manually created based on the available information provided in the article. The data were then converted to a QFeatures object using the scp::readSCP() function.

We generated the peptide data. First, we removed PSM matched to contaminants or decoy peptides and ensured a 1% FDR. We aggregated the PSM to peptides based on the peptide (or peptide group) sequence(s) using the median PSM instenity. The peptide data for the different runs were then joined in a single assay (see QFeatures::joinAssays), again based on the peptide sequence(s). We then removed the peptide groups. Links between the peptide and the PSM data were created using QFeatures::addAssayLink. Note that links between PSM and peptide groups are not stored.

The protein data were downloaded from Supporting information section from the publisher's website (see Sources). The data is supplied as an Excel file ac9b03349\_si\_005.xlsx. The file contains 7 sheets from which we only took the 2nd (named 01 - Raw sc protein data) with the combined protein data for the 12 runs. We converted the data to a SingleCellExperiment object and added the object as a new assay in the QFeatures dataset (containing the PSM data). Links between the proteins and the corresponding PSM were created. Note that links to protein groups are not stored.

# Note

Although a TMT-10plex labeling is reported in the article, the PSM data contained 11 channels for each run. Those 11th channel contain mostly missing data and are hence assumed to be empty channels.

#### Source

The PSM data can be downloaded from the massIVE repository MSV000084110. FTP link: ftp://massive.ucsd.edu/MSV0000 The protein data can be downloaded from the ACS Publications website (Supporting information section).

#### References

Dou, Maowei, Geremy Clair, Chia-Feng Tsai, Kerui Xu, William B. Chrisler, Ryan L. Sontag, Rui Zhao, et al. 2019. "High-Throughput Single Cell Proteomics Enabled by Multiplex Isobaric Labeling in a Nanodroplet Sample Preparation Platform." Analytical Chemistry, September (link to article).

#### See Also

dou2019\_lysates, dou2019\_boosting

### **Examples**

dou2019\_mouse()

gregoire2023\_mixCTRL Grégoire et al. 2023 - mixCTRL (arXiv): benchmark using monocytes/macrophages

# **Description**

Single cell proteomics data acquired using the SCoPE2 protocol. The dataset contains two monocytes cell lines (THP1 and U937) as well as controlled mixtures of both and macrophage-like cells produced upon PMA treatment. It contains quantitative information at PSM, peptide and protein levels. Data was acquired using Lumos Orbitrap (mainly) and timsTOF SCP mass spectrometers.

#### **Usage**

gregoire2023\_mixCTRL

#### **Format**

A QFeatures object with 119 assays, each assay being a SingleCellExperiment object:

- Assays 1-42: PSM data acquired with a TMT-16plex protocol, hence those assays contain 16 columns. Columns hold quantitative information from single-cell channels, carrier channels, blank (negative control) channels and unused channels.
- Assays 43-84: peptide data resulting from the PSM to peptide aggregation of the 42 PSM assays.

- Assays 85-91: peptide data for each of the 7 acquisition batches. Peptide data were joined based on their respective acquisition batches.
- Assays 92-98: normalised peptide data.
- Assays 99-105: normalised and log-transformed peptide data.
- Assays 106-112: protein data for each of the 7 acquisition batches. Normalised and log-transformed peptide data were agregated to protein.
- Assays 113-119: Batch corrected protein data. Normalised and log-transformed protein data were batch corrected to remove technical variability induced by runs and channels.

All the data has been filtered to keep high quality features and samples.

The colData(gregoire2023\_mixCTRL()) contains cell type annotation and batch annotation that are common to all assays. The description of the rowData fields for the PSM data can be found in the sage documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see *References*).

- Cell isolation: BD FACSAria III cell sorting.
- **Sample preparation** performed using the SCoPE2 protocol: mPOP cell lysis + trypsin digestion + TMT-16plex labeling and pooling.
- **Separation**: online nLC (Ultimate 3000 LC System or Vanquish Neo UHPLC System) with a BioZen Peptide Polar C18 250 x 0.0075mm column.
- Mass spectrometry: Orbitrap Fusion Lumos Tribrid (MS1 resolution = 70,000; MS2 accumulation time = 120ms; MS2 resolution = 70,000) and timsTOF SCP.
- Data preprocessing: Sage.

# Data collection

The PSM data were collected from a Zenodo archive (see Source section). The folder contains the following files of interest:

- results.sage.cbio.tsv: the sage identification output file for batches acquired on the Lumos MS.
- results.sage.giga.tsv: the sage identification output file for batches acquired on the timsTOF SCP MS.
- ullet quant.cbio.tsv: the sage quantification output file for batches acquired on the Lumos MS.
- quant.giga.tsv: the sage quantification output file for batches acquired on the timsTOF SCP MS.
- sampleAnnotation\_batch.csv: sample annotation for each acquisition batch. There are in total 8 different annotation files.

We combined the sample annotations in a single table. We also combined cbio and giga tables together and merged resulting identification and quantification tables. Both annotation and features tables are then combined in a single QFeatures object using the scp::readSCP() function.

16 guise 2024

The QFeatures object was processed as described in the author's manuscript (see source). Note that the imputed assays were used in the paper for illustrative purposes only and have not been reproduced here.

#### Source

The data were downloaded from the Zenodo repository. The raw data and the quantification data can also be found in the ProteomeXchange Consortium via the PRIDE partner repository, project PXD046211.

#### References

Samuel Grégoire, Christophe Vanderaa, Sébastien Pyr dit Ruys, Gabriel Mazzucchelli, Christopher Kune, Didier Vertommen and Laurent Gatto. 2023. Standardised workflow for mass spectrometry-based single-cell proteomics data processing and analysis using the scp package. arXiv. DOI:10.48550/arXiv.2310.13598

# **Examples**

```
gregoire2023_mixCTRL()
```

guise2024

Guise et al. 2020 (Cell Rep.): postmortem ALS spinal moto neurons

# **Description**

Single-cell proteomics data from postmortem human spinal moto neurons (MN) obtained from control donors or donors with amyotrophic lateral sclerosis (ALS). The data were generated following the NanoPOTS protocol. Cells were isolated from samples obtained by the university of Miami Brain Bank using laser capture microdissection (LCM). Additional information about the amount of TDP-43 intra-cellular levels has been assigned into levels 0 to 4.

# Usage

guise2024

#### **Format**

A QFeatures object with 102 assays, each assay being a SingleCellExperiment object:

- F\*: 100 assays containing PSM data.
- peptides: quantitative data for 34,315 peptides in 108 samples. All samples combined, along with 8 additional unannotated samples.
- proteins: quantitative data for 4,437 protein groups in 108 samples. All samples combined, along with 8 additional unannotated samples.

Sample annotation is stored in colData(guise2024()).

guise2024 17

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- **Cell isolation**: The MN were isolated from samples obtained by the university of Miami Brain Bank using LCM.
- **Sample preparation** performed using the nanoPOTS workflow. Cells are treated with 0.1% DDM (for lysis) added with DTT (protein reduction), then IAA (alkylation), then Lys-C and trypsin (protein digestion).
- **Separation**: Samples were injected on the column using an Ultimate 3000 RSLCnano pump. The in-line loading column is a home-packed SPE column (5cm x 75um) while the peptide separation is performed on a an in-house-packed analytical SPE column (50 cm x 30um), using a 20nL/min flow rate.
- **Ionization**: nanospray emmitter (2,000V)
- Mass spectrometry: Orbitrap Exploris 480. HCD fragmentation. MS1 settings: accumulation time = 200 ms; resolution = 120,000; AGC = 1E6. MS2 settings: exclusion duration = 90 s; accumulation time = 500 ms; resolution = 30,000; AGC = 1E5.
- **Data analysis**: Sequest HT in Proteome Discoverer (v2.5) and the search database is Swiss-Prot (July 2020).

#### **Data collection**

All data were collected from the MassIVE repository (accession ID: MSV000092119).

and in Groups.txt.

The PSM data were found in the Biogen TDP43 Pound? Peanalysis 10-13-2021 PSMs txt

The PSM data were found in the Biogen\_TDP43\_Round2\_Reanalysis\_10-13-2021\_PSMs.txt file. The data were converted to a QFeatures object using the scp::readSCP() function. We could not find sample annotations for MS run ID: F61, F34, F42, F88, F77, F8, F21, F5.

The peptide data were found in the Biogen\_TDP43\_Round2\_Reanalysis\_10-13-2021\_PeptideGroups.txt file. The column names holding the quantitative data were adapted to match the sample names in the QFeatures object. The data were then converted to a SingleCellExperiment object and then inserted in the QFeatures object.

A similar procedure was applied to the protein data. The data were found in the Biogen\_TDP43\_Round2\_Reanalysis\_10-13-file. The column names were adapted, the data were converted to a SingleCellExperiment object and then inserted in the QFeatures object.

The sample annotations were combined from the tables in Biogen\_TDP43\_Round2\_Reanalysis\_10-13-2021\_InputFiles.

# Source

All data can be downloaded from the MassIVE repository MSV000092119. The source link is: ftp://massive.ucsd.edu/v05/MSV000092119/

#### References

Guise, Amanda J., Santosh A. Misal, Richard Carson, Jen-Hwa Chu, Hannah Boekweg, Daisha Van Der Watt, Nora C. Welsh, et al. 2024. "TDP-43-Stratified Single-Cell Proteomics of Postmortem Human Spinal Motor Neurons Reveals Protein Dynamics in Amyotrophic Lateral Sclerosis." Cell Reports 43 (1): 113636. (link to article).

18 khan2023

### **Examples**

guise2024()

khan2023

Khan et al, 2023 (biorRxiv): Epithelial–Mesenchymal Transition

# **Description**

Single-cell samples were prepared using the nPOP sample preparation method. Proteomics data were acquired using the SCoPE2 protocol on a Thermo Scientific Q-Exactive mass spectrometer. The dataset contains quantitative information on 421 MCF-10A single cells undergoing epithelial—mesenchymal transition (EMT) triggered by TGF beta. The data are available at the PSM, and protein levels. The paper investigates the dynamics of correlation modules at the protein level.

# Usage

khan2023

#### **Format**

A QFeatures object with 47 assays, each assay being a SingleCellExperiment object:

- Assay 1-44: PSM data acquired with a TMTPro 16plex protocol, hence those assays contain 16 columns. Columns hold quantitative information from single-cell channels, carrier channels, reference channels, empty (negative control) channels and unused channels.
- peptides: peptide data containing quantitative data for 10055 peptides and 421 single-cells.
- proteins\_imputed: protein data containing quantitative data for 4096 proteins and 421 single-cells with k-nearest neighbors (KNN) imputation.
- proteins\_unimputed: protein data containing quantitative data for 4096 proteins and 421 single-cells without imputation.

The colData(khan2023()) contains cell type and batch annotations that are common to all assays. The description of the rowData fields for the PSM data can be found in the MaxQuant documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: CellenONE cell sorting.
- **Sample preparation** performed using the SCoPE2 protocol. nPOP cell lysis (DMSO) + trypsin digestion + TMTPro 16plex protocol.
- **Separation**: online nLC (DionexUltiMate 3000 UHPLC with a 25cm x 75um IonOpticks Odyssey Series column (ODY3-25075C18); 200nL/min).

khan2023

- **Ionization**: ESI (1,700 V).
- Mass spectrometry: Thermo Scientific Q-Exactive (MS1 resolution = 70,000; MS1 accumulation time = 300ms; MS2 resolution = 70,000).

• Data analysis: MaxQuant(2.4.13.0) + DART-ID.

#### **Data collection**

The PSM data were collected from a shared Google Drive folder that is accessible from the SlavovLab website (see Source section). The folder ('/002-singleCellDataGeneration') contains the following files of interest:

- ev\_updated\_NS.DIA.txt: the MaxQuant/DART-ID output file
- annotation.csv: sample annotation
- batch.csv: batch annotation

We combined the sample annotation and the batch annotation in a single table. We also formatted the quantification table so that columns match with those of the annotation and filter only for single-cell runs. Both table are then combined in a single QFeatures object using the scp::readSCP() function.

The peptide data were taken from the same google drive folder (EpiToMesen.TGFB.nPoP\_trial1\_pepByCellMatrix\_NSThr The data were formatted to a SingleCellExperiment object and the sample metadata were matched to the column names (mapping is retrieved after running the SCoPE2 R script, EMTTGFB\_singleCellProcessing.R) and stored in the colData. The object is then added to the QFeatures object and the rows of the PSM data are linked to the rows of the peptide data based on the peptide sequence information through an AssayLink object.

The imputed protein data were taken from the same google drive folder (EpiToMesen.TGFB.nPoP\_trial1\_ProtByCellMatr The data were formatted to a SingleCellExperiment object and the sample metadata were matched to the column names (mapping is retrieved after running the SCoPE2 R script, EMTTGFB\_singleCellProcessing.R) and stored in the colData. The object is then added to the QFeatures object and the rows of the peptide data are linked to the rows of the protein data based on the protein sequence information through an AssayLink object.

The unimputed protein data were taken from the same google drive folder (EpiToMesen.TGFB.nPoP\_trial1\_ProtByCellMa The data were formatted and added exactly as imputed data.

### Source

The data were downloaded from the Slavov Lab website via a shared Google Drive folder. The raw data and the quantification data can also be found in the MassIVE repository MSV000092872: ftp://MSV000092872@massive.ucsd.edu/.

#### References

Saad Khan, Rachel Conover, Anand R. Asthagiri, Nikolai Slavov. 2023. "Dynamics of single-cell protein covariation during epithelial–mesenchymal transition." bioRxiv. (link to article).

20 leduc2022\_plexDIA

# **Examples**

khan2023()

leduc2022

Deprecated leduc2022 dataset

# **Description**

The leduc2022 dataset has been updated to include plexDIA and pSCoPE data. The new datasets names are leduc2022\_pSCoPE (previously leduc2022) and leduc2022\_plexDIA (new). See the respective documentation pages for more information.

### Usage

leduc2022()

#### Value

The leduc2022\_pSCoPE dataset.

leduc2022\_plexDIA

Leduc et al. 2022 - plexDIA (biorRxiv): melanoma cells

# **Description**

Single cell proteomics data acquired by the Slavov Lab. This is the dataset associated to the fourth version of the preprint (and the Genome Biology publication). It contains quantitative information of melanoma cells at precursor, peptide and protein level. This version of the data was acquired using the plexDIA MS acquisition protocol.

#### Usage

leduc2022\_plexDIA

# **Format**

A QFeatures object with 48 assays, each assay being a SingleCellExperiment object:

- Assay 1-45: precursor data acquired with a mTRAQ-3 protocol, hence those assays contain 3 columns. Columns hold quantitative information from single cells or negative control samples.
- Ms1Extracted: the DIA-NN MS1 extracted signal, it combines the information from assays 1-45.
- peptides: peptide data containing quantitative data for 3,608 peptides and 104 single cells. The data were filtered to 1% protein FDR.

leduc2022\_plexDIA 21

• proteins: protein data containing quantitative data for 508 proteins and 105 single cells. Note that the peptide and protein data provided by the authors differ by 3 samples. The precursor data were aggregated to protein intensity using maxLFQ. The protein data were further median normalized by column and by row, log2 transformed, impute using KNN (k = 3), again median normalized by column and by row, batch corrected using ComBat, and median normalized by column and by row once more.

The colData(leduc2022\_plexDIA()) contains cell type annotation and batch annotation that are common to all assays. The description of the rowData fields for the precursor data can be found in the DIA-NN documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: CellenONE cell sorting.
- Sample preparation performed using the improved SCoPE2 protocol using the CellenONE liquid handling system. nPOP cell lysis (DMSO) + trypsin digestion + mTRAQ-3 labeling and pooling.
- **Separation**: online nLC (DionexUltiMate 3000 UHPLC with a 25cm x 75um IonOpticks Aurora Series UHPLC column; 200nL/min).
- Ionization: ESI (1,800V).
- Mass spectrometry: Thermo Scientific Q-Exactive. The duty cycle = 1 MS1 + 4 DIA MS2 windows (120 Th, 120 Th, 200 Th and 580 Th, spanning 378-1,402 m/z). Each MS1 and MS2 scan was conducted at 70,000 resolving power, 3×10E6 AGC and 300ms maximum injection time.
- Data analysis: DIA-NN.

#### **Data collection**

The PSM data were collected from a shared Google Drive folder that is accessible from the SlavovLab website (see Source section). The folder contains the following files of interest:

- annotation\_plexDIA.csv: sample annotation
- report\_plexDIA\_mel\_nPOP. tsv: the DIA-NN output file with the precursor data
- report.pr\_matrix\_channels\_ms1\_extracted.tsv: the DIA-NN output file with the combined precursor data
- plexDIA\_peptide.csv: the processed data table containing the peptide data
- plexDIA\_protein\_imputed.csv: the processed data table containing the protein data

We removed the failed runs as identified by the authors. We also formatted the annotation and precuror quantification tables to facilitate matching between corresponding columns. Both annotation and quantification tables are then combined in a single QFeatures object using scp::readSCPfromDIANN().

The plexDIA\_peptide.csv and plexDIA\_protein\_imputed.csv files were loaded and formatted as SingleCellExperiment objects. The columns names were adapted to match those in the QFeatures object. The SingleCellExperiment objects were then added to the QFeatures object and the rows of the peptide data are linked to the rows of the precursor data based on the peptide sequence or the protein name through an AssayLink object.

22 leduc2022\_pSCoPE

#### Source

The links to the data were found on the Slavov Lab website. The data were downloaded from the Google drive folder 1 and Google drive folder 2. The raw data and the quantification data can also be found in the massIVE repository MSV000089159: ftp://massive.ucsd.edu/MSV000089159.

#### References

Andrew Leduc, Gray Huffman, and Nikolai Slavov. 2022. "Droplet Sample Preparation for Single-Cell Proteomics Applied to the Cell Cycle." bioRxiv. Link to article

Andrew Leduc, Gray Huffman, Joshua Cantlon, Saad Khan, and Nikolai Slavov. 2022. "Exploring Functional Protein Covariation across Single Cells Using nPOP." Genome Biology 23 (1): 261. Link to article

Jason Derks, Andrew Leduc, Georg Wallmann, Gray Huffman, Matthew Willetts, Saad Khan, Harrison Specht, Markus Ralser, Vadim Demichev, and Nikolai Slavov. 2023. "Increasing the Throughput of Sensitive Proteomics by plexDIA." Nature Biotechnology 41 (1): 50–59. Link to article

#### See Also

leduc2022\_pSCoPE

# **Examples**

leduc2022\_plexDIA()

leduc2022\_pSCoPE

Leduc et al. 2022 - pSCoPE (biorRxiv): melanoma cells vs monocytes

# Description

Single cell proteomics data acquired by the Slavov Lab. This is the dataset associated to the third version of the preprint. It contains quantitative information of melanoma cells and monocytes at PSM, peptide and protein level. This version of the data was acquired using the pSCoPE MS acquisition approach.

# Usage

leduc2022\_pSCoPE

## **Format**

A QFeatures object with 138 assays, each assay being a SingleCellExperiment object:

• Assay 1-134: PSM data acquired with a TMT-18plex protocol, hence those assays contain 18 columns. Columns hold quantitative information from single-cell channels, carrier channels, reference channels, empty (negative control) channels and unused channels.

leduc2022\_pSCoPE 23

• peptides: peptide data containing quantitative data for 20,804 peptides and 1556 single-cells. These data have been filtered to keep high-quality PSMs, all batches have been normalized to the reference channel, PSMs were aggregated to peptides, and single-cells with low median coefficient of variation were kept.

- peptides\_log: peptide data containing quantitative data for 12,284 peptides and 1543 single-cells. The peptides data was further normalized, highly missing peptides were removed and the quantifications were log-transformed.
- proteins\_norm2: protein data containing quantitative data for 2844 proteins and 1543 single-cells. The peptides from peptides\_log were aggregated to proteins and normalized.
- proteins\_processed: protein data containing quantitative data for 2844 proteins and 1543 single-cells. The proteins\_norm2 data were imputed, batch corrected and normalized.

The colData(leduc2022\_pSCoPE()) contains cell type annotation, LC batch information, the TMT label, the MS run ID. We also added the sample prep annotations provided by the cellenONE dispensing device (only for single cells): time stamp of cell isolation by the device, the diameter and elongation of the cell, the ID of the sample glass side (4 slides in total), the field within the glass (each slide is divided in 4 field), the pooled well ID (each field contains 9 pools), the x and y coordinates of each cell dropped in a field and of each cell pool upon pickup. Finally, we also retrieved the melanoma subpopulation generated by the authors upon data analysis. The main population is encoded as A while the small population is encoded B. The description of the rowData fields for the PSM data can be found in the MaxQuant documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: CellenONE cell sorting.
- Sample preparation performed using the improved SCoPE2 protocol using the CellenONE liquid handling system. nPOP cell lysis (DMSO) + trypsin digestion + TMT-18plex labeling and pooling. A target library was generated as well to perform prioritized DDA (Huffman et al. 2022) using MaxQuant.Live (2.0.3).
- **Separation**: online nLC (DionexUltiMate 3000 UHPLC with a 25cm x 75um IonOpticks Aurora Series UHPLC column; 200nL/min).
- Ionization: ESI (1,800V).
- Mass spectrometry: Thermo Scientific Q-Exactive (MS1 resolution = 70,000; MS2 accumulation time = 300ms; MS2 resolution = 70,000). Prioritized data acquisition was performed using the pSCoPE protocol (Huffman et al. 2022)
- Data analysis: MaxQuant (1.6.17.0) + DART-ID

#### **Data collection**

The PSM data were collected from a shared Google Drive folder that is accessible from the SlavovLab website (see Source section). The folder contains the following files of interest:

- ev\_updated.txt: the MaxQuant/DART-ID output file
- annotation.csv: sample annotation

24 leduc2022\_pSCoPE

- batch.csv: batch annotation
- t0.csv: the processed data table containing the peptides data
- t3.csv: the processed data table containing the peptides\_log data
- t4b.csv: the processed data table containing the proteins\_norm2 data
- t6.csv: the processed data table containing the proteins\_processed data

We combined the sample annotation and the batch annotation in a single table. We also formatted the quantification table so that columns match with those of the annotations. Both annotation and quantification tables are then combined in a single QFeatures object using the scp::readSCP() function.

The 4 CSV files were loaded and formatted as SingleCellExperiment objects and the sample metadata were matched to the column names (mapping is retrieved after running the author's original R script) and stored in the colData. The object is then added to the QFeatures object (containing the PSM assays) and the rows of the peptide data are linked to the rows of the PSM data based on the peptide sequence information through an AssayLink object.

#### Source

The data were downloaded from the Slavov Lab website. The raw data and the quantification data can also be found in the massIVE repository MSV000089159: ftp://massive.ucsd.edu/MSV000089159.

### References

Andrew Leduc, Gray Huffman, and Nikolai Slavov. 2022. "Droplet Sample Preparation for Single-Cell Proteomics Applied to the Cell Cycle." bioRxiv. Link to article

Gray Huffman, Andrew Leduc, Christoph Wichmann, Marco di Gioia, Francesco Borriello, Harrison Specht, Jason Derks, et al. 2022. "Prioritized Single-Cell Proteomics Reveals Molecular and Functional Polarization across Primary Macrophages." bioRxiv. Link to article.

#### See Also

leduc2022\_plexDIA

### **Examples**

leduc2022\_pSCoPE()

liang2020\_hela 25

ing)	liang2020_hela	Liang et al. 2020 (Anal. Chem.): HeLa cells (MaxQuant preprocessing)
------	----------------	--

# **Description**

Single-cell proteomics data from HeLa cells using the autoPOTS acquisition workflow. The samples contain either no cells (blanks), 1 cell, 10 cells, 150 cells or 500 cells. Samples containing between 0 and 10 cells are isolated using micro-pipetting while samples containing between 150 and 500 cells were prepared using dilution of a bulk sample.

# Usage

liang2020\_hela

#### **Format**

A QFeatures object with 17 assays, each assay being a SingleCellExperiment object:

- HeLa\_\*: 15 assays containing PSM data.
- peptides: quantitative data for 48705 peptides in 15 samples (all runs are combined).
- proteins: quantitative data for 3970 protein groups in 15 samples (all runs combined).

Sample annotation is stored in colData(liang2020\_hela()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: The HeLa cells come from a commercially available cell line. Samples containing between 0 and 10 cells were isolated using micro-manipulation and the counts were validated using a microscope. Samples containing between 150 and 500 cells were prepared by diluting a bulk sample and the exact counts were evaluated by obtaining phtotmicrographs.
- Sample preparation performed using the autoPOTS worflow that relied on the OT-2 pipeting robot. Cell are lysed using sonication. Samples are then processed by successive incubation with DTT (reduction), then IAA (alkylation), then Lys-C and trypsin (protein digestion).
- **Separation**: Samples were injected on the column using a modified Ultimate WPS-3000 TPL autosampler coupled to an UltiMate 3000 RSLCnano pump. The LC column is a home-packed nanoLC column (45cm x 30um; 40nL/min)
- **Ionization**: Nanospray Flex ion source (2,000V)
- Mass spectrometry: Orbitrap Exploris 480. MS1 settings: accumulation time = 250 ms (0-10 cells) or 100 ms (150-500 cells); resolution = 120,000; AGC = 100\ duration = 90 s (0-10 cells) or 60 s (150-500 cells); accumulation time = 500 ms (0-1 cell), 250 ms (10 cells), 100 ms (150 cells) or 50 ms (500 cells); resolution = 60,000 (0-10 cells) or 30,000 (150-500 cells); AGC = 5E3 (0-1 cells) or 1E4 (10-500 cells).
- Data analysis: MaxQuant (v1.6.7.0) and the search database is Swiss-Prot (July 2020).

# **Data collection**

All data were collected from the PRIDE repository (accession ID: PXD021882).

The sample annotations were collected from the methods section and from table S3 in the paper.

The PSM data were found in the evidence.txt file. The data were converted to a QFeatures object using the scp::readSCP() function.

The peptide data were found in the peptides.txt file. The column names holding the quantitative data were adapted to match the sample names in the QFeatures object. The data were then converted to a SingleCellExperiment object and then inserted in the QFeatures object. Links between the PSMs and the peptides were added

A similar procedure was applied to the protein data. The data were found in the proteinGroups.txt file. The column names were adapted, the data were converted to a SingleCellExperiment object and then inserted in the QFeatures object. Links between the peptides and the proteins were added

#### Source

The PSM data can be downloaded from the PRIDE repository PXD021882 The source link is: http://ftp.pride.ebi.ac.uk/pride/data/archive/2020/12/PXD021882/

#### References

Liang, Yiran, Hayden Acor, Michaela A. McCown, Andikan J. Nwosu, Hannah Boekweg, Nathaniel B. Axtell, Thy Truong, Yongzheng Cong, Samuel H. Payne, and Ryan T. Kelly. 2020. "Fully Automated Sample Processing and Analysis Workflow for Low-Input Proteome Profiling." Analytical Chemistry, December. (link to article).

# **Examples**

liang2020\_hela()

petrosius2023\_AstralAML

Petrosius et al. 2023 (bioRxiv): AML hierarchy on Astral.

# Description

Single cell proteomics data from FACS sorted cells from the OCI-AML8227 model. The dataset contains leukemic stem cells (LSC; CD34+, CD38-), progenitor cells (CD34+, CD38+), CD38+ blasts (CD34-, CD38+) and CD38- blasts (CD34-, CD38-). It contains quantitative information at PSM, peptide and protein levels. Data was acquired using an Orbitrap Astral mass spectrometer. Direct DIA analysis was performed with Spectronaut version 17.

# Usage

petrosius2023\_AstralAML

#### **Format**

A QFeatures object with 217 assays, each assay being a SingleCellExperiment object:

- Assays 1-215: PSM data from the Spectronaut PEPQuant file with LFQ quantities from the FG.MS1Quantity column.
- peptides: Peptide data resulting from the PSM to peptide aggregation the 215 PSM assays. Resulting peptide assays were joined into a single assay.
- proteins: Protein data from the Spectronaut PGQuant file with LFQ quantities from the PG.Quantity column.

The colData(petrosius2023\_AstralAML()) contains cell type annotation, batch annotation and FACS data. The description of the rowData fields can be found in the Spectronaut user manual.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see *References*).

- Cell isolation: Cell sorting was done on a FACS Aria III or Aria II instrument, controlled by the DIVA software package and operated with a 100 microm nozzle. Cells were sorted at single-cell resolution, into a 384-well Eppendorf LoBind PCR plate containing 1 microL of lysis buffer.
- Sample preparation Single-cell protein lysates were digested overnight at 37°C with 2 ng of Trypsin supplied in 1 microL of digestion buffer. Digestion was stopped by the addition of 1 microL 1% (v/v) trifluoroacetic acid (TFA). All liquid dispensing was done using an I-DOT One instrument.
- Liquid chromatography: Chromatographic separation of peptides was conducted on a vanquish Neo UHPLC system connected to a 50 cm uPAC Neo Low-load and an EASY-spray. Autosampler and injection valves were configured to perform direct injections from a 384 well plate using a 25 uL injection loop on 11.8 min gradients.
- Mass spectrometry: Acquisition was conducted with an Orbitrap Astral mass spectrometer operated in positive mode with the FAIMSPro interface compensation voltage set to -45 V. MS1 scans were acquired with the Orbitrap at a resolution of 120,000 and a scan range of 400 to 900 m/z with normalized automatic gain control (AGC) target of 300 % and maximum injection time of 246 ms. Data independent acquisition of MS2 spectra was performed in the Astral using loop control set to 0.7 seconds per cycle with varying isolation window widths and injection times. Fragmentation of precursor ions was performed using higher energy collisional dissociation (HCD) using a normalized collision energy (NCE) of 25 %. AGC target was set to 800 %.
- Raw data processing: Raw files were processed using Spectronaut version 17. Direct DIA analysis was performed in pipeline mode. Pulsar searches were performed without fixed modifications. N-terminal acetylation and methionine oxidation were set as variable modifications. Quantification level was set to MS1 and the quantity type set to area under the curve.

### **Data collection**

The data were provided by the authors and is accessible at the Dataverse The dataset ('Astral AML single-cell data from Petrosius et al. 2023 preprint') contains the following files of interest:

28 petrosius2023\_mES

- 20240201\_130747\_PEPQuant (Normal).tsv: the PSM level data
- 20240201\_130747\_PGQuant (Normal).tsv: the protein level data
- index\_map.csv: FACS data.
- msRuns\_overview.csv: Sample annotations.

We added the FACS data to the sample annotations in a single table. Both annotations and PSM features tables are then combined in a single QFeatures object using the scp::readSCP() function.

The peptide data were obtained by aggregation of the PSM data to the peptide level. All of the resulting peptides assays were joined into a single assays. Individual peptides assays were discarded.

The protein data were formatted from the 20240201\_130747\_PGQuant (Normal).tsv to a Single-CellExperiment object and the sample metadata were matched to the column names and stored in the colData. The object is then added to the QFeatures object and the rows of the peptide data are linked to the rows of the protein data based on the protein sequence information through an AssayLink object.

Note that the QFeatures object has not been further processed and has therefore not been normalized, log-transformed or batch-corrected.

#### Source

The PSM data, protein data and sample annotations can be downloaded from the dataset 'Astral AML single-cell data from Petrosius et al. 2023 preprint' in the Dataverse.

#### References

Valdemaras Petrosius, Pedro Aragon-Fernandez, Tabiwang N. Arrey, Nil Üresin, Benjamin Furtwängler, Hamish Stewart, Eduard Denisov, Johannes Petzoldt, Amelia C. Peterson, Christian Hock, Eugen Damoc, Alexander Makarov, Vlad Zabrouskov, Bo T. Porse and Erwin M. Schoof. 2023. "Evaluating the capabilities of the Astral mass analyzer for single-cell proteomics." biorxiv. https://doi.org/10.1101/2023.06.06.543943

# **Examples**

petrosius2023\_AstralAML()

petrosius2023\_mES Petrosius et al, 2023 (Nat. Comm.): Mouse embryonic stem cell (mESC) in different culture conditions

# Description

Profiling mouse embryonic stem cells across ground-state (m2i) and differentiation-permissive (m15) culture conditions. The data were acquired using orbitrap-based data-independent acquisition (DIA). The objective was to demonstrate the capability of their approach by profiling mouse embryonic stem cell culture conditions, showcasing heterogeneity in global proteomes, and highlighting differences in the expression of key metabolic enzymes in distinct cell subclusters.

petrosius2023\_mES 29

#### Usage

petrosius2023\_mES

#### **Format**

A QFeatures object with 605 assays, each assay being a SingleCellExperiment object:

- Assay 1-603: PSM data acquired with an orbitrap-based data-independent acquisition (DIA) protocol, hence those assays contain single column that contains the quantitative information.
- peptides: peptide data containing quantitative data for 9884 peptides and 603 single-cells.
- proteins: protein data containing quantitative data for 4270 proteins and 603 single-cells.

Sample annotation is stored in colData(petrosius2023\_mES()).

### **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Sample isolation: Cell sorting was done on a Sony MA900 cell sorter using a 130 microm sorting chip. Cells were sorted at single-cell resolution, into a 384-well Eppendorf LoBind PCR plate (Eppendorf AG) containing 1 microL of lysis buffer.
- Sample preparation: Single-cell protein lysates were digested with 2 ng of Trypsin supplied in 1 microL of digestion buffer which was carried out overnight at 37 °C, and subsequently acidified by the addition of 1 microL 1% (v/v) trifluoroacetic acid (TFA). All liquid dispensing was done using an I-DOT One instrument.
- Liquid chromatography: For the HRMS1-DIA experiments and the DIA isolation window survey, Evosep One liquid chromatography was used. The standard 31-minute or 58-minute pre-defined Whisper gradients were used with a flow rate of 100 nl/min for peptide elution.
- Mass spectrometry: The mass spectrometer was operated in positive mode with the FAIM-SPro interface compensation voltage set to -45 V. MS1 scans were carried out at 120,000 resolution with an automatic gain control (AGC) of 300% and maximum injection time set to auto. For the DIA isolation window survey, a scan range of 500–900 was used, and 400–1000 rest of the experiments. Higher energy collisional dissociation (HCD) was used for precursor fragmentation with a normalized collision energy (NCE) of 33% and the MS2 scan AGC target was set to 1000%.
- Raw data processing: The mESC raw data files were processed with Spectronaut 17.

#### **Data collection**

The data were provided by the Author and is accessible at the Dataverse The folder ('20240205\_111248\_mESC\_SNEcombine m2i/') contains the following files of interest:

- 20240205\_111251\_PEPQuant (Normal).tsv: the PSM level data
- 20240205\_111251\_Peptide Quant (Normal).tsv: the peptide level data
- 20240205\_111251\_PGQuant (Normal).tsv: the protein level data

The metadata were downloaded from the Zenodo repository.

30 schoof2021

• sample\_facs.csv: the metadata

We formatted the quantification table so that columns match with the metadata. Then, both tables are then combined in a single QFeatures object using the scp::readSCP() function.

The peptide data were formated to a SingleCellExperiment object and the sample metadata were matched to the column names and stored in the colData. The object is then added to the QFeatures object and the rows of the PSM data are linked to the rows of the peptide data based on the peptide sequence information through an AssayLink object.

The protein data were formated to a SingleCellExperiment object and the sample metadata were matched to the column names and stored in the colData. The object is then added to the QFeatures object and the rows of the peptide data are linked to the rows of the protein data based on the protein sequence information through an AssayLink object.

#### Source

The peptide and protein data can be downloaded from the Dataverse The raw data and the quantification data can also be found in the MassIVE repository MSV000092429: ftp://MSV000092429@massive.ucsd.edu/.

#### References

**Source article**: Petrosius, V., Aragon-Fernandez, P., Üresin, N. et al. "Exploration of cell state heterogeneity using single-cell proteomics through sensitivity-tailored data-independent acquisition." Nat Commun 14, 5910 (2023). (link to article).

# **Examples**

petrosius2023\_mES()

schoof2021	Schoof et al. 2021 (Nat. Comm.): acute myeloid leukemia differentiation
------------	---

# Description

Single-cell proteomics data from OCI-AML8227 cell culture to reconstruct the cellular hierarchy. The data were acquired using TMTpro multiplexing. The samples contain either no cells, single cells, 10 cells (reference channel) 200 cells (booster channel) or are simply empty wells. Single cells are expected to be one of progenitor cells (PROG), leukaemia stem cells (LSC), CD38- blast cells (BLAST CD38-) or CD38+ blast cells (BLAST CD38+). Booster are either a known 1:1:1 mix of cells (PROG, LSC and BLAST) or are isolated directly from the bulk sample. Samples were isolated and annotated using flow cytometry.

# Usage

schoof2021

schoof2021 31

#### **Format**

A QFeatures object with 194 assays, each assay being a SingleCellExperiment object:

• F\*: 192 assays containing PSM quantification data for 16 TMT channels. The quantification data contain signal to noise ratios as computed by Proteome Discoverer.

- proteins: quantitative data for 2898 protein groups in 3072 samples (all runs combined). The quantification data contain signal to noise ratios as computed by Proteome Discoverer.
- logNormProteins: quantitative data for 2723 protein groups in 2025 single-cell samples. This assay is the protein datasets that was processed by the authors. Dimension reduction and clustering data are also available in the reducedDims and colData slots, respectively

Sample annotation is stored in colData(schoof2021()). The cell type annotation is stored in the Population column. The flow cytometry data is also available: FSC-A, FSC-H, FSC-W, SSC-A, SSC-H, SSC-W, APC-Cy7-A (= CD34) and PE-A (= CD38).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- **Sample isolation**: cultured AML 8227 cells were stained with anti-CD34 and anti-CD38. The sorting was performed by FACSAria instrument and deposited in 384 well plates.
- Sample preparation: cells are lysed using freeze-boil and sonication in a lysis buffer (TFE) that also includes reduction and alkylation reagents (TCEP and CAA), followed by trypsin (protein) and benzonase (DNA) digestion, TMT-16 labeling and quenching, desalting using SOLAµ C18 plate, peptide concentration, pooling and peptide concentration again. The booster channel contains 200 cell equivalents.
- Liquid chromatography: peptides are separated using a C18 reverse-phase column (50cm x 75 µm i.d., Thermo EasySpray) combined to a Thermo EasyLC 1200 for 160 minute gradient with a flowrate of 100nl/min.
- Mass spectrometry: FAIMSPro interface is used. MS1 setup: resolution 60.000, AGC target of 300%, accumulation of 50ms. MS2 setup: resolution 45.000, AGC target of 150, 300 or 500%, accumulation of 150, 300, 500, or 1000ms.
- Raw data processing: Proteome Discoverer 2.4 + Sequest spectral search engine and validation with Percolator

# **Data collection**

All data were collected from the PRIDE repository (accession ID: PXD020586). The data and metadata were extracted from the SCeptre\_FINAL.zip file.

We performed extensive data wrangling to combine all the metadata available from different files into a single table available using colData(schoof2021).

The PSM data were found in the bulk\_PSMs.txt file. Contaminants were defined based on the protein accessions listed in contaminant.txt. The data were converted to a QFeatures object using the scp::readSCP() function.

The protein data were found in the bulk\_Proteins.txt file. Contaminants were defined based on the protein accessions listed in contaminant.txt.The column names holding the quantitative data

32 scpdata

were adapted to match the sample names in the QFeatures object. Unnecessary feature annotations (such as in which assay a protein is found) were removed. Feature names were created following the procedure in SCeptre: features names are the protein symbol (or accession if missing) and if duplicated symbols are present (protein isoforms), they are made unique by appending the protein accession. Contaminants were defined based on the protein accessions listed in contaminant.txt. The data were then converted to a SingleCellExperiment object and inserted in the QFeatures object.

The log-normalized protein data were found in the bulk.h5ad file. This dataset was generated by the authors by running the notebook called bulk.ipynb. The bulk.h5ad was loaded as an AnnData object using the scanpy Python module. The object was then converted to a SingleCellExperiment object using the zellkonverter package. The column names holding the quantitative data were adapted to match the sample names in the QFeatures object. The data were then inserted in the QFeatures object.

The script to reproduce the QFeatures object is available at system.file("scripts", "make-data\_schoof2021.R", package = "scpdata")

#### Source

The PSM and protein data can be downloaded from the PRIDE repository PXD020586 The source link is: https://www.ebi.ac.uk/pride/archive/projects/PXD020586

#### References

Schoof, Erwin M., Benjamin Furtwängler, Nil Üresin, Nicolas Rapin, Simonas Savickas, Coline Gentil, Eric Lechman, Ulrich auf Dem Keller, John E. Dick, and Bo T. Porse. 2021. "Quantitative Single-Cell Proteomics as a Tool to Characterize Cellular Hierarchies." Nature Communications 12 (1): 745679. (link to article).

#### **Examples**

schoof2021()

scpdata

Single-Cell Proteomics Data Package

# Description

The scpdata package distributes mass spectrometry-based single-cell proteomics datasets. The datasets were collected from published work and formatted to a standardized data framework. The scp frameworks stores the expression data for different MS levels (identified spectrum, peptide, or protein) in separate assays. Each assay is an object of class SingleCellExperiment that allows easy integration with state-of-the-art single-cell analysis tools. All assays are contained in a single object of class QFeatures. An overview of the data structure is shown provided in the scp package.

The scpdata() function returns a summary table with all currently available datasets in the package. More information about the data content and the data collection can be found in the corresponding manual pages.

specht2019v2 33

# Usage

```
scpdata()
```

# Value

A DataFrame table containing a summary of the available datasets.

# Author(s)

Christophe Vanderaa

#### See Also

More information about the data manipulation can be found in the scp package.

# **Examples**

```
## List available datasets and their metadata
scpdata()

## Load data using the ExperimentHub interface
hub <- ExperimentHub()

## Not run:

## Download the data set of interest using ExperimentHub indexing
hub[["EH7711"]]

## Download the same data set using the build-in function
leduc2022()

## End(Not run)</pre>
```

specht2019v2

Specht et al. 2019 - SCoPE2 (biorRxiv): macrophages vs monocytes (version 2)

# Description

Single cell proteomics data acquired by the Slavov Lab. This is the version 2 of the data released in December 2019. It contains quantitative information of macrophages and monocytes at PSM, peptide and protein level.

# Usage

```
specht2019v2
```

34 specht2019v2

#### **Format**

A QFeatures object with 179 assays, each assay being a SingleCellExperiment object:

• Assay 1-63: PSM data for SCoPE2 sets acquired with a TMT-11plex protocol, hence those assays contain 11 columns. Columns hold quantitative information from single-cell channels, carrier channels, reference channels, empty (blank) channels and unused channels.

- Assay 64-177: PSM data for SCoPE2 sets acquired with a TMT-16plex protocol, hence those assays contain 16 columns. Columns hold quantitative information from single-cell channels, carrier channels, reference channels, empty (blank) channels and unused channels.
- peptides: peptide data containing quantitative data for 9208 peptides and 1018 single-cells.
- proteins: protein data containing quantitative data for 2772 proteins and 1018 single-cells.

The colData(specht2019v2()) contains cell type annotation and batch annotation that are common to all assays. The description of the rowData fields for the PSM data can be found in the MaxQuant documentation.

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: flow cytometry (BD FACSAria I).
- **Sample preparation** performed using the SCoPE2 protocol. mPOP cell lysis + trypsin digestion + TMT-11plex or 16plex labelling and pooling.
- **Separation**: online nLC (DionexUltiMate 3000 UHPLC with a 25cm x 75um IonOptick-sAurora Series UHPLC column; 200nL/min).
- Ionization: ESI (2,200V).
- Mass spectrometry: Thermo Scientific Q-Exactive (MS1 resolution = 70,000; MS1 accumulation time = 300ms; MS2 resolution = 70,000).
- Data analysis: DART-ID + MaxQuant (1.6.2.3).

#### Data collection

The PSM data were collected from a shared Google Drive folder that is accessible from the SlavovLab website (see Source section). The folder contains the following files of interest:

- ev\_updated.txt: the MaxQuant/DART-ID output file
- annotation\_fp60-97.csv: sample annotation
- batch\_fp60-97.csv: batch annotation

We combined the sample annotation and the batch annotation in a single table. We also formatted the quantification table so that columns match with those of the annotation and filter only for single-cell runs. Both table are then combined in a single QFeatures object using the scp::readSCP() function.

The peptide data were taken from the Slavov lab directly (Peptides-raw.csv). It is provided as a spreadsheet. The data were formatted to a SingleCellExperiment object and the sample metadata were matched to the column names (mapping is retrieved after running the SCoPE2 R script) and

specht2019v3 35

stored in the colData. The object is then added to the QFeatures object (containing the PSM assays) and the rows of the peptide data are linked to the rows of the PSM data based on the peptide sequence information through an AssayLink object.

The protein data (Proteins-processed.csv) is formatted similarly to the peptide data, and the rows of the proteins were mapped onto the rows of the peptide data based on the protein sequence information.

#### Source

The data were downloaded from the Slavov Lab website via a shared Google Drive folder. The raw data and the quantification data can also be found in the massIVE repository MSV000083945: ftp://massive.ucsd.edu/MSV000083945.

#### References

Specht, Harrison, Edward Emmott, Aleksandra A. Petelski, R. Gray Huffman, David H. Perlman, Marco Serra, Peter Kharchenko, Antonius Koller, and Nikolai Slavov. 2019. "Single-Cell Mass-Spectrometry Quantifies the Emergence of Macrophage Heterogeneity." bioRxiv. (link to article).

# **Examples**

specht2019v2()

specht2019v3	Specht et al. 2019 - SCoPE2 (biorRxiv): macrophages vs monocytes
	(version 3)

# **Description**

Single cell proteomics data acquired by the Slavov Lab. This is the version 3 of the data released in October 2020. It contains quantitative information of macrophages and monocytes at PSM, peptide and protein level.

## Usage

specht2019v3

#### **Format**

A QFeatures object with 179 assays, each assay being a SingleCellExperiment object:

• Assay 1-63: PSM data for SCoPE2 sets acquired with a TMT-11plex protocol, hence those assays contain 11 columns. Columns hold quantitative information from single-cell channels, carrier channels, reference channels, empty (blank) channels and unused channels.

36 specht2019v3

• Assay 64-177: PSM data for SCoPE2 sets acquired with a TMT-16plex protocol, hence those assays contain 16 columns. Columns hold quantitative information from single-cell channels, carrier channels, reference channels, empty (blank) channels and unused channels.

- peptides: peptide data containing quantitative data for 9208 peptides and 1018 single-cells.
- proteins: protein data containing quantitative data for 2772 proteins and 1018 single-cells.

The colData(specht2019v2()) contains cell type annotation and batch annotation that are common to all assays. The description of the rowData fields for the PSM data can be found in the MaxQuant documentation.

#### **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- **Cell isolation**: flow cytometry (BD FACSAria I).
- **Sample preparation** performed using the SCoPE2 protocol. mPOP cell lysis + trypsin digestion + TMT-11plex or 16plex labeling and pooling.
- **Separation**: online nLC (DionexUltiMate 3000 UHPLC with a 25cm x 75um IonOptick-sAurora Series UHPLC column; 200nL/min).
- **Ionization**: ESI (2,200V).
- Mass spectrometry: Thermo Scientific Q-Exactive (MS1 resolution = 70,000; MS2 accumulation time = 300ms; MS2 resolution = 70,000).
- Data analysis: DART-ID + MaxQuant (1.6.2.3).

# **Data collection**

The PSM data were collected from a shared Google Drive folder that is accessible from the SlavovLab website (see Source section). The folder contains the following files of interest:

- ev\_updated\_v2.txt: the MaxQuant/DART-ID output file
- annotation\_fp60-97.csv: sample annotation
- batch\_fp60-97.csv: batch annotation

We combined the sample annotation and the batch annotation in a single table. We also formatted the quantification table so that columns match with those of the annotation and filter only for single-cell runs. Both table are then combined in a single QFeatures object using the scp::readSCP() function.

The peptide data were taken from the Slavov lab directly (Peptides-raw.csv). It is provided as a spreadsheet. The data were formatted to a SingleCellExperiment object and the sample metadata were matched to the column names (mapping is retrieved after running the SCoPE2 R script) and stored in the colData. The object is then added to the QFeatures object (containing the PSM assays) and the rows of the peptide data are linked to the rows of the PSM data based on the peptide sequence information through an AssayLink object.

The protein data (Proteins-processed.csv) is formatted similarly to the peptide data, and the rows of the proteins were mapped onto the rows of the peptide data based on the protein sequence information.

williams2020\_lfq 37

## Note

Since version 2, a serious bug in the data were corrected for TMT channels 12 to 16. Many more cells are therefore contained in the data. Version 2 is maintained for backward compatibility. Although the final version of the article was published in 2021, we have kept specht2019v3 as the data set name for consistency with the previous data version specht2019v2.

#### Source

The data were downloaded from the Slavov Lab website via a shared Google Drive folder. The raw data and the quantification data can also be found in the massIVE repository MSV000083945: ftp://massive.ucsd.edu/MSV000083945.

#### References

Specht, Harrison, Edward Emmott, Aleksandra A. Petelski, R. Gray Huffman, David H. Perlman, Marco Serra, Peter Kharchenko, Antonius Koller, and Nikolai Slavov. 2021. "Single-Cell Proteomic and Transcriptomic Analysis of Macrophage Heterogeneity Using SCoPE2." Genome Biology 22 (1): 50. (link to article).

# **Examples**

specht2019v3()

williams2020\_lfq

Williams et al. 2020 (Anal. Chem.): MCF10A cell line

# **Description**

Single-cell label free proteomics data from a MCF10A cell line culture. The data were acquired using a label-free quantification protocole based on the nanoPOTS technology. The objective was to test 2 elution gradients for single-cell applications and to demonstrate successful use of the new nanoPOTS autosampler presented in the article. The samples contain either no cells, single cells, 3 cells, 10 cells 50 cells.

## **Usage**

williams2020\_lfq

# **Format**

A QFeatures object with 9 assays, each assay being a SingleCellExperiment object:

• peptides\_[30 or 60]min\_[intensity or LFQ]: 3 assays containing peptide intensities or LFQ normalized quantifications (see References) for either a 30min or a 60 min gradient.

38 williams2020\_lfq

• proteins\_[30 or 60]min\_[intensity or iBAQ or LFQ]: 6 assays containing protein intensities, iBAQ normalized or LFQ normalized quantifications (see References) for either a 30min or a 60 min gradient.

Sample annotation is stored in colData(williams2020\_lfq()).

## **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Sample isolation: cultured MCF10A cells were isolated using flow-cytometry based cell sorting and deposit on nanoPOTS microwells
- Sample preparation: cells are lysed using using a DDM+DTT lysis buffer. Alkylation was then performed using an IAA solution. Proteins are digested with Lys-C and trypsin followed by acidification with FA. Sample droplets are then dried until LC-MS/MS analysis.
- Liquid chromatography: peptides are loaded using the new autosampler described in the paper. Samples are loaded using a a homemade miniature syringe pump. The samples are then desalted and concentrated through a SPE column (4cm x 100µm i.d. packed with 5µm C18) with microflow LC pump. The peptides are then eluted from a long LC column (60cm x 50 µm i.d. packed with 3µm C18) coupled to a nanoflox LC pump at 150nL/mL with either a 30 min or a 60 min gradient.
- Mass spectrometry: MS/MS was performed on an Orbitrap Fusion Lumos Tribrid MS coupled to a 2kV ESI. MS1 setup: Orbitrap analyzer at resolution 120.000, AGC target of 1E6, accumulation of 246ms. MS2 setup: ion trap with CID at resolution 60.000, AGC target of 2E4, accumulation of 120ms (50 cells) or 250ms (0-10 cells).
- **Raw data processing**: preprocessing using Maxquant v1.6.2.10 that use Andromeda search engine (with UniProtKB 2016-21-29), MBR and LFQ normalization were enabled.

# **Data collection**

All data were collected from the MASSIVE repository (accession ID: MSV000085230).

The peptide and protein data were extracted from the Peptides\_[...].txt or ProteinGroups[...].txt files, respectively, in the MCF10A\_LC\_[30 or 60]minutes folders.

The tables were duplicated so that peptide intensisities, peptide LFQ, protein intensities, protein LFQ and protein intensities are contained in separate tables. Tables are then converted to Single-CellExperiment objects. Sample annotations were infered from the sample names and from the paper. All data is combined in a QFeatures object. AssayLinks were stored between peptide assays and their corresponding proteins assays based on the leading razor protein (hence only unique peptides are linked to proteins).

The script to reproduce the QFeatures object is available at system.file("scripts", "make-data\_williams2020\_lfq.R' package = "scpdata")

# Suggestion

See QFeatures:: joinAssays if you want to join the 30min and 60min assays in a single assay for an integrated analysis.

williams2020\_tmt 39

#### Source

The PSM and protein data can be downloaded from the MASSIVE repository MSV000085230.

#### References

**Source article**: Williams, Sarah M., Andrey V. Liyu, Chia-Feng Tsai, Ronald J. Moore, Daniel J. Orton, William B. Chrisler, Matthew J. Gaffrey, et al. 2020. "Automated Coupling of Nanodroplet Sample Preparation with Liquid Chromatography-Mass Spectrometry for High-Throughput Single-Cell Proteomics." Analytical Chemistry 92 (15): 10588–96. (link to article).

**LFQ normalization**: Cox, Jürgen, Marco Y. Hein, Christian A. Luber, Igor Paron, Nagarjuna Nagaraj, and Matthias Mann. 2014. "Accurate Proteome-Wide Label-Free Quantification by Delayed Normalization and Maximal Peptide Ratio Extraction, Termed MaxLFQ." Molecular & Cellular Proteomics: MCP 13 (9): 2513–26. (link to article).

**iBAQ normalization**: Schwanhäusser, Björn, Dorothea Busse, Na Li, Gunnar Dittmar, Johannes Schuchhardt, Jana Wolf, Wei Chen, and Matthias Selbach. 2011. "Global Quantification of Mammalian Gene Expression Control." Nature 473 (7347): 337–42. (link to article).

# **Examples**

williams2020\_lfq()

williams2020\_tmt

Williams et al. 2020 (Anal. Chem.): 3 AML cell line

# **Description**

Single-cell label data from three acute myeloid leukemia cell line culture (MOLM-14, K562, CMK). The data were acquired using a TMT-based quantification protocole and the nanoPOTS technology. The objective was to demonstrate successful use of the new nanoPOTS autosampler presented in the source article. The samples contain either carrier (10 ng), reference (0.2ng), empty or single-cell samples..

# Usage

williams2020\_tmt

# Format

A QFeatures object with 4 assays, each assay being a SingleCellExperiment object:

- peptides\_[intensity or corrected]: 2 assays containing peptide reporter ion intensities or corrected reporter ion intensities as computed by MaxQuant.
- proteins\_[intensity or corrected]: 2 assays containing protein reporter ion intensities or corrected reporter ion intensities as computed by MaxQuant.

Sample annotation is stored in colData(williams2020\_tmt()).

40 williams2020 tmt

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Sample isolation: cultured MOLM-14, K562 or CMK cells were isolated using flow-cytometry based cell sorting and deposit on nanoPOTS microwells
- Sample preparation: cells are lysed using using a DDM lysis buffer. Proteins are digested with trypsin followed by TMT labelling and quanching with HA. The samples are then acidified with FA, pooled in a single samples (adding carrier and reference peptide mixtures), and dried until LC-MS/MS analysis.
- Liquid chromatography: peptides are loaded using the new autosampler described in the paper. Samples are loaded using a a homemade miniature syringe pump. The samples are then desalted and concentrated through a SPE column (4cm x 100μm i.d. packed with 5μm C18) with microflow LC pump. The peptides are then eluted from a long LC column (60cm x 50 μm i.d. packed with 3μm C18) coupled to a nanoflox LC pump at 150nL/mL (elution time is not explicated).
- Mass spectrometry: MS/MS was performed on an Orbitrap Fusion Lumos Tribrid MS coupled to a 2kV ESI. MS1 setup: Orbitrap analyzer at resolution 120.000, AGC target of 1E6, accumulation of 246ms. MS2 setup: Orbitrap with HCD at resolution 120.000, AGC target of 1E6, accumulation of 246ms.
- **Raw data processing**: preprocessing using Maxquant v1.6.2.10 that use Andromeda search engine (with UniProtKB 2016-21-29).

# **Data collection**

All data were collected from the MASSIVE repository (accession ID: MSV000085230).

The peptide and protein data were extracted from the Peptides\_AML\_SingleCell.txt or ProteinGroups\_AML\_SingleCell files, respectively, in the AML\_SingleCell folders.

The tables were duplicated so that intensisities and corrected intensities are contained in separate tables. Tables are then converted to SingleCellExperiment objects. Sample annotations were inferred from the sample names, from table S2 and from the Experimental Section of the paper. All data is combined in a QFeatures object. AssayLinks were stored between peptide assays and their corresponding proteins assays based on the leading razor protein (hence only unique peptides are linked to proteins).

The script to reproduce the QFeatures object is available at system.file("scripts", "make-data\_williams2020\_tmt.R package = "scpdata")

# Source

The PSM and protein data can be downloaded from the MASSIVE repository MSV000085230.

#### References

**Source article**: Williams, Sarah M., Andrey V. Liyu, Chia-Feng Tsai, Ronald J. Moore, Daniel J. Orton, William B. Chrisler, Matthew J. Gaffrey, et al. 2020. "Automated Coupling of Nanodroplet Sample Preparation with Liquid Chromatography-Mass Spectrometry for High-Throughput Single-Cell Proteomics." Analytical Chemistry 92 (15): 10588–96. (link to article).

woo2022\_lung 41

# **Examples**

williams2020\_tmt()

woo2022\_lung

Woo et al. 2022 (Cell Syst.): 26 primary human lung cells

# **Description**

Single-cell proteomics data from dissociated primary human lung cells. The data were acquired using the TIFF (transfer identification based on FAIMS filtering) acquisition method. The data contain 26 single cells.

# Usage

woo2022\_lung

## **Format**

A QFeatures object with 5 assays, each assay being a SingleCellExperiment object:

- peptides\_[intensity or LFQ]: 2 assays containing peptide quantities or normalized quantities using the maxLFQ method as computed by MaxQuant.
- proteins\_[intensity or iBAQ or LFQ]: 3 assays containing protein quantities or normalized proteins using the iBAQ or maxLFQ methods as computed by MaxQuant.

Sample annotation is stored in colData(woo\_lung()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Sample isolation: primary human lung cells were dissociated following the protocol in Bandyopadhyay et al., 2018. The cells were sorted using the Influx II cell sorter and deposited on a nanoPOTS chip.
- Sample preparation: cells are lysed using using a DDM+DTT lysis and reduction buffer. The proteins are alkylated with IAA and digested with LysC and trypsin. Samples are then acidified with FA, vacuum dried and stored in freezer until data acquisition.
- Liquid chromatography: peptides are loaded using an in-house autosampler (Williams et al. 2020). The samples are concentrated through a SPE column (4cm x 100μm i.d. packed with 5μm C18) with microflow LC pump. The peptides are then eluted from an LC column (25cm x 50 μm i.d. packed with 1.7μm C18) from a 60 min gradient (100nL/min).
- Mass spectrometry: MS/MS was performed on an Orbitrap Fusion Lumos Tribrid MS with FAIMSpro coupled to a 2.4 kV ESI. FAIMS setup: 4-CV method (-45, -55, -65, -75 V). MS1 setup: resolution = 120.000, range = 350-1500 m/z,AGC target of 1E6, accumulation of 254ms. MS2 setup: 30% HCD, resolution AGC 2E4, accumulation of 254ms.

• **Raw data processing**: preprocessing using Maxquant v1.6.2.10 that use Andromeda search engine (with UniProtKB 2016-21-29). MBR was enabled.

# **Data collection**

All data were collected from the MASSIVE repository (accession ID: MSV000085937).

The peptide and protein data were extracted from the peptides\_nondepleted\_Lung\_scProteomics.txt or proteinGroups\_nondepleted\_Lung\_scProteomics.txt files, respectively, in the NonDepleted\_Lung\_SingleCellPro folders.

The tables were split so that intensities, maxLFQ, and iBAQ data are contained in separate tables. Tables are then converted to SingleCellExperiment objects. Sample annotations were inferred from the sample names. All data is combined in a QFeatures object. AssayLinks were stored between peptide assays and their corresponding proteins assays based on the leading razor protein (hence only unique peptides are linked to proteins).

The script to reproduce the QFeatures object is available at system.file("scripts", "make-data\_woo2022\_lung.R", package = "scpdata")

#### Source

The peptide and protein data can be downloaded from the MASSIVE repository MSV000085937

## References

**Source article**: Woo, Jongmin, Geremy C. Clair, Sarah M. Williams, Song Feng, Chia-Feng Tsai, Ronald J. Moore, William B. Chrisler, et al. 2022. "Three-Dimensional Feature Matching Improves Coverage for Single-Cell Proteomics Based on Ion Mobility Filtering." Cell Systems 13 (5): 426–34.e4. (link to article).

# **Examples**

woo2022\_lung()

woo2022\_macrophage

Woo et al. 2022 (Cell Syst.): LPS-treated macrophages

# Description

Single-cell data from macrophages subjected to 3 LPS treatments. The data were acquired using the TIFF (transfer identification based on FAIMS filtering) acquisition method. The data contain 155 single cells: 54 control cells (no treatment), 52 cells treated with LPS during 24h and 49 cells treated with LPS during 49h.

# Usage

woo2022\_macrophage

#### **Format**

A QFeatures object with 5 assays, each assay being a SingleCellExperiment object:

- peptides\_[intensity or LFQ]: 2 assays containing peptide quantities or normalized quantities using the maxLFQ method as computed by MaxQuant.
- proteins\_[intensity or iBAQ or LFQ]: 3 assays containing protein quantities or normalized proteins using the iBAQ or maxLFQ methods as computed by MaxQuant.

Sample annotation is stored in colData(woo\_macrophage()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- **Sample isolation**: cultured RAW 264.7 cells treated or not with 100 ng/ul LPS. The cells were sorted using the Influx II cell sorter and deposited on a nanoPOTS chip.
- Sample preparation: cells are lysed using using a DDM+DTT lysis and reduction buffer. The proteins are alkylated with IAA and digested with LysC and trypsin. Samples are then acidified with FA, vacuum dried and stored in freezer until data acquisition.
- Liquid chromatography: peptides are loaded using an in-house autosampler (Williams et al. 2020). The samples are concentrated through a SPE column (4cm x 100μm i.d. packed with 5μm C18) with microflow LC pump. The peptides are then eluted from an LC column (25cm x 50 μm i.d. packed with 1.7μm C18) from a 60 min gradient (100nL/min).
- Mass spectrometry: MS/MS was performed on an Orbitrap Fusion Lumos Tribrid MS with FAIMSpro coupled to a 2.4 kV ESI. FAIMS setup: 4-CV method (-45, -55, -65, -75 V). MS1 setup: resolution = 120.000, range = 350-1500 m/z, AGC target of 1E6, accumulation of 254ms. MS2 setup: 30% HCD, resolution AGC 2E4, accumulation of 254ms.
- **Raw data processing**: preprocessing using Maxquant v1.6.2.10 that use Andromeda search engine (with UniProtKB 2016-21-29). MBR was enabled.

# **Data collection**

All data were collected from the MASSIVE repository (accession ID: MSV000085937).

The peptide and protein data were extracted from the peptides\_RAW\_LPS\_scProteomics.txt or proteinGroups\_RAW\_LPS\_scProteomics.txt files, respectively, in the RAW\_LPS\_SingleCellProteomics folders.

The tables were split so that intensities, maxLFQ, and iBAQ data are contained in separate tables. Tables are then converted to SingleCellExperiment objects. Sample annotations were inferred from the sample names. All data is combined in a QFeatures object. AssayLinks were stored between peptide assays and their corresponding proteins assays based on the leading razor protein (hence only unique peptides are linked to proteins).

The script to reproduce the QFeatures object is available at system.file("scripts", "make-data\_woo2022\_macrophage package = "scpdata")

# Source

The peptide and protein data can be downloaded from the MASSIVE repository MSV000085937

44 zhu2018MCP

## References

**Source article**: Woo, Jongmin, Geremy C. Clair, Sarah M. Williams, Song Feng, Chia-Feng Tsai, Ronald J. Moore, William B. Chrisler, et al. 2022. "Three-Dimensional Feature Matching Improves Coverage for Single-Cell Proteomics Based on Ion Mobility Filtering." Cell Systems 13 (5): 426–34.e4. (link to article).

# **Examples**

woo2022\_macrophage()

zhu2018MCP

Zhu et al. 2018 (Mol. Cel. Prot.): rat brain laser dissections

# Description

Near single-cell proteomics data of laser captured micro-dissection samples. The samples are 24 brain sections from rat pups (day 17). The slices are 12 um thick squares of either 50, 100, or 200 um width. 5 samples were dissected from the corpus callum (CC), 4 samples were dissected from the corpus collosum (CP), 13 samples were extracted from the cerebral cortex (CTX), and 2 samples are labeled as (Mix).

# Usage

zhu2018MCP

## **Format**

A QFeatures object with 4 assays, each assay being a SingleCellExperiment object:

- peptides: quantitative information for 13,055 peptides from 24 samples
- proteins\_intensity: protein intensities for 2,257 proteins from 24 samples
- proteins\_LFQ: LFQ intensities for 2,257 proteins from 24 samples
- proteins\_iBAQ: iBAQ values for 2,257 proteins from 24 samples

Sample annotation is stored in colData(zhu2018MCP()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the original article (see References).

• Cell isolation: brain patches were collected using laser-capture microdissection (PALM MicroBeam) on flash frozen rat (*Rattus norvergicus*) brain tissues. Note that the samples were stained with H&E before dissection for histological analysis. DMSO is used as sample collection solution

zhu2018MCP 45

• **Sample preparation** performed using the nanoPOTs device: DMSO evaporation + protein extraction (DMM + DTT) + alkylation (IAA)

- Lys-C digestion + trypsin digestion.
- Separation: nanoLC (Dionex UltiMate with an in-house packed 60cm x 30um LC columns; 50nL/min)
- **Ionization**: ESI (2,000V)
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid (MS1 accumulation time = 246ms; MS1 resolution = 120,000; MS1 AGC = 3E6). The MS/MS settings depend on the sample size, excepted for the AGC = 1E5. 50um (time = 502ms; resolution = 240,000), 100um (time = 246ms; resolution = 120,000), 200um (time = 118ms; resolution = 60,000).
- Data analysis: MaxQuant (v1.5.3.30) + Perseus (v1.5.6.0) + Origin Pro 2017

## Data collection

The data were collected from the PRIDE repository (accession ID: PXD008844). We downloaded the MaxQuant\_Peptides.txt and the MaxQuant\_ProteinGroups.txt files containing the combined identification and quantification results. The sample annotations were inferred from the names of columns holding the quantification data and the information in the article. The peptides data were converted to a SingleCellExperiment object. We split the protein table to separate the three types of quantification: protein intensity, label-free quantitification (LFQ) and intensity based absolute quantification (iBAQ). Each table is converted to a SingleCellExperiment object along with the remaining protein annotations. The 4 objects are combined in a single QFeatures object and feature links are created based on the peptide leading razor protein ID and the protein ID.

#### **Source**

The PSM data can be downloaded from the PRIDE repository PXD008844. FTP link ftp://ftp.pride.ebi.ac.uk/pride/data/archi

# References

Zhu, Ying, Maowei Dou, Paul D. Piehowski, Yiran Liang, Fangjun Wang, Rosalie K. Chu, William B. Chrisler, et al. 2018. "Spatially Resolved Proteome Mapping of Laser Capture Microdissected Tissue with Automated Sample Transfer to Nanodroplets." Molecular & Cellular Proteomics: MCP 17 (9): 1864–74 (link to article).

# **Examples**

zhu2018MCP()

46 zhu2018NC\_hela

zhu2018NC\_hela

Zhu et al. 2018 (Nat. Comm.): HeLa titration

## **Description**

Near single-cell proteomics data of HeLa samples containing different number of cells. There are three groups of cell concentrations: low (10-14 cells), medium (35-45 cells) and high (137-141 cells). The data also contain measures for blanks, HeLa lysates (50 cell equivalent) and 2 cancer cell line lysates (MCF7 and THP1, 50 cell equivalent).

# Usage

zhu2018NC\_hela

## **Format**

A QFeatures object with 4 assays, each assay being a SingleCellExperiment object:

- peptides: quantitative information for 37,795 peptides from 21 samples
- proteins\_intensity: protein intensities for 3,984 proteins from 21 samples
- proteins\_LFQ: LFQ intensities for 3,984 proteins from 21 samples
- proteins\_iBAQ: iBAQ values for 3,984 proteins from 21 samples

Sample annotation is stored in colData(zhu2018NC\_hela()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the original article (see References).

- **Cell isolation**: HeLa cell concentration was adjusted by serial dilution and cell counting was performed manually using an inverted microscope.
- Sample preparation performed using the nanoPOTs device. Protein extraction using RapiGest (+ DTT) + alkylation (IAA) + Lys-C digestion + cleave RapiGest (formic acid).
- **Separation**: nanoACQUITY UPLC pump (60nL/min) with an Self-Pack PicoFrit 70cm x 30um LC columns.
- **Ionization**: ESI (1,900V).
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid. MS1 settings: accumulation time = 246ms; resolution = 120,000; AGC = 1E6. MS/MS settings, depend on the sample size, excepted for the AGC = 1E5. Blank and approx. 10 cells (time = 502ms; resolution = 240,000), approx. 40 cells (time = 246ms; resolution = 120,000), approx. 140 cells (time = 118ms; resolution = 60,000).
- Data analysis: MaxQuant (v1.5.3.30) + Perseus + OriginLab 2017

zhu2018NC\_islets 47

## **Data collection**

The data were collected from the PRIDE repository (accession ID: PXD006847). We downloaded the CulturedCells\_peptides.txt and the CulturedCells\_proteinGroups.txt files containing the combined identification and quantification results. The sample annotations were inferred from the names of columns holding the quantification data and the information in the article. The peptides data were converted to a SingleCellExperiment object. We split the protein table to separate the three types of quantification: protein intensity, label-free quantification (LFQ) and intensity based absolute quantification (iBAQ). Each table is converted to a SingleCellExperiment object along with the remaining protein annotations. The 4 objects are combined in a single QFeatures object and feature links are created based on the peptide leading razor protein ID and the protein ID.

## **Source**

The PSM data can be downloaded from the PRIDE repository PXD006847. FTP link: ftp://ftp.pride.ebi.ac.uk/pride/data/arch

#### References

Zhu, Ying, Paul D. Piehowski, Rui Zhao, Jing Chen, Yufeng Shen, Ronald J. Moore, Anil K. Shukla, et al. 2018. "Nanodroplet Processing Platform for Deep and Quantitative Proteome Profiling of 10-100 Mammalian Cells." Nature Communications 9 (1): 882 (link to article).

# See Also

The same experiment was conducted on HeLa lysates: zhu2018NC\_lysates.

# **Examples**

zhu2018NC\_hela()

zhu2018NC\_islets

Zhu et al. 2018 (Nat. Comm.): human pancreatic islets

# Description

Near single-cell proteomics data human pancreas samples. The samples were collected from pancreatic tissue slices using laser dissection. The pancreata were obtained from organ donors through the JDRFNetwork for Pancreatic Organ Donors with Diabetes (nPOD) program. The sample come either from control patients (n=9) or from type 1 diabetes (T1D) patients (n=9).

# Usage

zhu2018NC\_islets

48 zhu2018NC\_islets

## **Format**

A QFeatures object with 4 assays, each assay being a SingleCellExperiment object:

- peptides: quantitative information for 24,321 peptides from 18 islet samples
- proteins\_intensity: quantitative information for 3,278 proteins from 18 islet samples
- proteins\_LFQ: LFQ intensities for 3,278 proteins from 18 islet samples
- proteins\_iBAQ: iBAQ values for 3,278 proteins from 18 islet samples

Sample annotation is stored in colData(zhu2018NC\_islets()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: The islets were extracted from the pacreatic tissues using laser-capture microdissection.
- Sample preparation performed using the nanoPOTs device. Protein extraction using RapiGest (+ DTT) + alkylation (IAA) + Lys-C digestion + cleave RapiGest (formic acid)
- **Separation**: nanoACQUITY UPLC pump with an Self-Pack PicoFrit 70cm x 30um LC columns; 60nL/min)
- **Ionization**: ESI (1,900V)
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid. MS1 settings: accumulation time = 246ms; resolution = 120,000; AGC = 1E6. MS/MS settings: accumulation time = 118ms; resolution = 60,000; AGC = 1E5.
- Data analysis: MaxQuant (v1.5.3.30) + Perseus + OriginLab 2017

## Data collection

The data were collected from the PRIDE repository (accession ID: PXD006847). We downloaded the Islet\_t1d\_ct\_peptides.txt and the Islet\_t1d\_ct\_proteinGroups.txt files containing the combined identification and quantification results. The sample types were inferred from the names of columns holding the quantification data. The peptides data were converted to a SingleCell-Experiment object. We split the protein table to separate the three types of quantification: protein intensity, label-free quantitification (LFQ) and intensity based absolute quantification (iBAQ). Each table is converted to a SingleCellExperiment object along with the remaining protein annotations. The 4 objects are combined in a single QFeatures object and feature links are created based on the peptide leading razor protein ID and the protein ID.

# Source

The PSM data can be downloaded from the PRIDE repository PXD006847. The source link is: ftp://ftp.pride.ebi.ac.uk/pride/data/archive/2018/01/PXD006847

## References

Zhu, Ying, Paul D. Piehowski, Rui Zhao, Jing Chen, Yufeng Shen, Ronald J. Moore, Anil K. Shukla, et al. 2018. "Nanodroplet Processing Platform for Deep and Quantitative Proteome Profiling of 10-100 Mammalian Cells." Nature Communications 9 (1): 882 (link to article).

zhu2018NC\_lysates 49

# **Examples**

zhu2018NC\_islets()

zhu2018NC\_lysates

Zhu et al. 2018 (Nat. Comm.): HeLa lysates

# **Description**

Near single-cell proteomics data of HeLa lysates at different concentrations (10, 40 and 140 cell equivalent). Each concentration is acquired in triplicate.

# Usage

zhu2018NC\_lysates

#### **Format**

A QFeatures object with 4 assays, each assay being a SingleCellExperiment object:

- peptides: quantitative information for 14,921 peptides from 9 lysate samples
- proteins\_intensity: quantitative information for 2,199 proteins from 9 lysate samples
- proteins\_LFQ: LFQ intensities for 2,199 proteins from 9 lysate samples
- proteins\_iBAQ: iBAQ values for 2,199 proteins from 9 lysate samples

Sample annotation is stored in colData(zhu2018NC\_lysates()).

# **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the original article (see References).

- Cell isolation: HeLas were collected from cell cultures.
- Sample preparation performed in bulk (5E5 cells/mL). Protein extraction using RapiGest (+ DTT) + dilution to target concentration + alkylation (IAA) + Lys-C digestion + trypsin digestion + cleave RapiGest (formic acid).
- **Separation**: nanoACQUITY UPLC pump (60nL/min) with an Self-Pack PicoFrit 70cm x 30um LC columns.
- Ionization: ESI (1,900V).
- Mass spectrometry: Thermo Fisher Orbitrap Fusion Lumos Tribrid. MS1 settings: accumulation time = 246ms; resolution = 120,000; AGC = 1E6. MS/MS settings, depend on the sample size, excepted for the AGC = 1E5. Blank and approx. 10 cells (time = 502ms; resolution = 240,000), approx. 40 cells (time = 246ms; resolution = 120,000), approx. 140 cells (time = 118ms; resolution = 60,000).
- Data analysis: MaxQuant (v1.5.3.30) + Perseus + OriginLab 2017.

50 zhu2019EL

# **Data collection**

The data were collected from the PRIDE repository (accession ID: PXD006847). We downloaded the Vail\_Prep\_Vail\_peptides.txt and the Vail\_Prep\_Vail\_proteinGroups.txt files containing the combined identification and quantification results. The sample annotations were inferred from the names of columns holding the quantification data and the information in the article. The peptides data were converted to a SingleCellExperiment object. We split the protein table to separate the three types of quantification: protein intensity, label-free quantitification (LFQ) and intensity based absolute quantification (iBAQ). Each table is converted to a SingleCellExperiment object along with the remaining protein annotations. The 4 objects are combined in a single QFeatures object and feature links are created based on the peptide leading razor protein ID and the protein ID.

## Source

The PSM data can be downloaded from the PRIDE repository PXD006847. The source link is: ftp://ftp.pride.ebi.ac.uk/pride/data/archive/2018/01/PXD006847

#### References

Zhu, Ying, Paul D. Piehowski, Rui Zhao, Jing Chen, Yufeng Shen, Ronald J. Moore, Anil K. Shukla, et al. 2018. "Nanodroplet Processing Platform for Deep and Quantitative Proteome Profiling of 10-100 Mammalian Cells." Nature Communications 9 (1): 882 (link to article).

## See Also

The same experiment was conducted directly on HeLa cells samples rather than lysates. The data is available in zhu2018NC\_hela.

# **Examples**

zhu2018NC\_lysates()

zhu2019EL

Zhu et al. 2019 (eLife): chicken utricle cells

# **Description**

Single-cell proteomics data from chicken utricle acquired to study the hair-cell development. The cells are isolated from peeled utrical epithelium and separated into hair cells (FM1-43 high) and supporting cells (FM1-43 low). The sample contain either 1 cell (n = 28), 3 cells (n = 7), 5 cells (n = 8) or 20 cells (n = 14).

# Usage

zhu2019EL

zhu2019EL 51

## **Format**

A QFeatures object with 62 assays, each assay being a SingleCellExperiment object:

• XYZw: 60 assays containing PSM data. The sample are annotated as follows. X indicates the experiment, either 1 or 2. Y indicated the FM1-43 signal, either high (H) or low (L). Z indicates the number of cells (0, 1, 3, 5 or 20). w indicates the replicate, starting from a, it can go up to i.

- peptides: quantitative data for 3444 peptides in 60 samples (all runs are combined).
- proteins\_intensity: protein intensities for 840 proteins from 24 samples
- proteins\_iBAQ: iBAQ values for 840 proteins from 24 samples

Sample annotation is stored in colData(zhu2019EL()).

## **Acquisition protocol**

The data were acquired using the following setup. More information can be found in the source article (see References).

- Cell isolation: The cells were taken from the utricles of E15 chick embryos. Samples were stained with FM1-43FX and the cells were dissociated using enzymatic digestion. Cells were FACS sorted (BD Influx) and split based on their FM1-43 signal, while ensuring no debris, doublets or dead cells are retained.
- Sample preparation performed using the nanoPOTs device. Cell lysis and protein extraction and reduction are performed using dodecyl beta-D-maltoside + DTT + ammonium bicarbonate. Protein were then alkylated using IAA. Protein digestion is performed using Lys-C and trypsin. Finally samples acidification is performed using formic acid.
- Separation: Dionex UltiMate pump with an C18-Packed column (50cm x 30um; 60nL/min)
- Ionization: ESI (2,000V)
- Mass spectrometry: Orbitrap Fusion Lumos Tribrid. MS1 settings: accumulation time = 246ms; resolution = 120,000; AGC = 3E6. MS/MS settings: accumulation time = 502ms; resolution = 120,000; AGC = 2E5.
- Data analysis: Andromeda & MaxQuant (v1.5.3.30) and the search database is NCBI GRCg6a.

## **Data collection**

All data were collected from the PRIDE repository (accession ID: PXD014256).

The sample annotation information is provided in the Zhu\_2019\_chick\_single\_cell\_samples\_CORRECTED.xlsx file. This file was given during a personal discussion and is a corrected version of the annotation table available on the PRIDE repository.

The PSM data were found in the evidence.txt (in the Experiment 1+ 2) folder. The PSM data were filtered so that it contains only samples that are annotated. The data were then converted to a OFeatures object using the scp::readSCP() function.

The peptide data were found in the peptides.txt file. The column names holding the quantitative data were adapted to match the sample names in the QFeatures object. The data were then converted to a SingleCellExperiment object and then inserted in the QFeatures object. Links between the PSMs and the peptides were added

52 zhu2019EL

A similar procedure was applied to the protein data. The data were found in the proteinGroups.txt file. We split the protein table to separate the two types of quantification: summed intensity and intensity based absolute quantification (iBAQ). Both tables are converted to SingleCellExperiment objects and are added to the QFeatures object as well as the AssayLink between peptides and proteins.

#### Source

The PSM data can be downloaded from the PRIDE repository PXD014256. The source link is: ftp://ftp.pride.ebi.ac.uk/pride/data/archive/2019/11/PXD014256

# References

Zhu, Ying, Mirko Scheibinger, Daniel Christian Ellwanger, Jocelyn F. Krey, Dongseok Choi, Ryan T. Kelly, Stefan Heller, and Peter G. Barr-Gillespie. 2019. "Single-Cell Proteomics Reveals Changes in Expression during Hair-Cell Development." eLife 8 (November). (link to article).

# **Examples**

zhu2019EL()

# **Index**

* datasets	khan2023, <u>18</u>
brunner2022, 2	
cong2020AC, 4	leduc2022, <u>20</u>
derks2022, 6	leduc2022_plexDIA, 20, 24
dou2019_boosting,8	leduc2022_pSCoPE, 22, 22
dou2019_lysates, 10	liang2020_hela, <u>25</u>
dou2019_mouse, 12	
gregoire2023_mixCTRL, 14	petrosius2023_AstralAML, 26
guise2024, 16	petrosius2023_mES, 28
khan2023, 18	2 22 24 22 24 52
leduc2022_plexDIA, 20	QFeatures, 3–22, 24–32, 34–52
leduc2022_pSCoPE, 22	QFeatures::addAssayLink, 9, 11, 13
liang2020_hela,25	QFeatures::joinAssays, $9$ , $11$ , $13$
petrosius2023_AstralAML, 26	schoof2021, 30
petrosius2023_mES, 28	school 2021, 30 scp::readSCP(), 5, 9, 11, 13, 15, 17, 19, 24,
schoof2021, 30	26, 28, 30, 31, 34, 36, 51
specht2019v2, 33	scpdata, 32
specht2019v3, 35	scpdata, 32 scpdata-package (scpdata), 32
williams2020_lfq, 37	SingleCellExperiment, 3, 4, 6–14, 16–22,
williams2020_tmt, 39	24–32, 34–52
woo2022_lung, 41	specht2019v2, 33
woo2022_macrophage, 42	specht2019v3, 35
zhu2018MCP, 44	Specific 2019 v 3, 33
zhu2018NC_hela, 46	williams2020_lfq,37
zhu2018NC_islets,47	williams2020_tmt, 39
zhu2018NC_lysates, 49	woo2022_lung, 41
zhu2019EL, 50	woo2022_macrophage, 42
,	WOOZOZZ_IIIdel Ophiage, 12
AssayLinks, 38, 40, 42, 43	zhu2018MCP, 44
brunner2022, 2	zhu2018NC_hela, 46, 50
	zhu2018NC_islets,47
2022010 1	zhu2018NC_lysates, 47, 49
cong2020AC, 4	zhu2019EL, 50
derks2022, 6	,
dou2019_boosting, 8, 12, 14	
dou2019_lysates, 9, 10, 14	
dou2019_mouse, 9, 12, 12	
4042013_110400, 7, 12, 12	
gregoire2023_mixCTRL, 14	
guise2024, 16	