# beadplexr: Reproducible and automated analysis of multiplex bead assays

Ulrik Stervbo [Corresp., 1, 2] , Timm H Westhoff [1] , Nina Babel [1, 2]

[1] Center for Translational Medicine, Medical Clinic I, Marien Hospital Herne, University Hospital of the Ruhr-University Bochum, Herne, Germany

[2] Berlin-Brandenburg Center for Regenerative Therapies, Charité – Universitätsmedizin, Berlin, Germany

Corresponding Author: Ulrik Stervbo
Email address: ulrik.stervbo-kristensen@charite.de

Multiplex bead assays are an extension of the commonly used sandwich ELISA. The advantage over ELISA is that they make simultaneous evaluation of several analytes possible. Several commercial assay systems, where the beads are acquired on a standard flow cytometer, exist. These assay systems come with their own software tool for analysis and evaluation of the concentration of the analyzed analytes. However, these tools are either tied to particular commercial software or impose other limitations to their licenses, such as the number of events which can be analyzed. In addition, all these solutions are 'point and click' which potentially obscures the steps taken in the analysis. Here we present beadplexer, an open-source R-package for the reproducible analysis of multiplex bead assay data. The package makes it possible to automatically identify bead clusters, and provides functionality to easily fit a standard curve and calculate the concentrations of the analyzed analytes. beadplexer is available from CRAN and from https://gitlab.com/ustervbo/beadplexr

PeerJ

# 1 beadplexr: Reproducible and automated
# 2 analysis of multiplex bead assays

3 Ulrik Stervbo[1,2], Timm H. Westhoff[1], Nina Babel[1,2]

4 [1]Center for Translational Medicine, Medical Clinic I, Marien Hospital Herne, University
5 Hospital of the Ruhr-University Bochum, Herne, Germany

6 [2]Berlin-Brandenburg Center for Regenerative Therapies, Charité – Universitätsmedizin,
7 Berlin, Germany

8 **Correspondence**
9 Ulrik Stervbo: ulrik.stervbo-kristensen@charite.de

10 **Abstract**
11 Multiplex bead assays are an extension of the commonly used sandwich ELISA. The
12 advantage over ELISA is that they make simultaneous evaluation of several analytes
13 possible. Several commercial assay systems, where the beads are acquired on a standard
14 flow cytometer, exist. These assay systems come with their own software tool for analysis
15 and evaluation of the concentration of the analyzed analytes. However, these tools are
16 either tied to particular commercial software or impose other limitations to their licenses,
17 such as the number of events which can be analyzed. In addition, all these solutions are
18 'point and click' which potentially obscures the steps taken in the analysis. Here we present
19 `beadplexer`, an open-source R-package for the reproducible analysis of multiplex bead assay
20 data. The package makes it possible to automatically identify bead clusters, and provides
21 functionality to easily fit a standard curve and calculate the concentrations of the analyzed
22 analytes. `beadplexer` is available from CRAN and from
23 https://gitlab.com/ustervbo/beadplexr.

## Introduction

24

25 The enzyme-linked immunosorbent assay (ELISA) is a commonly used method to

26 determine the concentration of soluble analytes such as cytokines (Elshal & McCoy, 2006).

27 The concentration of the analyte is determined from a standard curve, which is created

28 from standard samples with known concentrations. The ELISA is a single point assay and

29 query into several analytes can be time consuming or impossible when the sample is

30 limited. Development in polystyrene bead preparations made it possible to construct

31 assays that allow for query of several analytes at the same time. Similar to the ELISA, the

32 analytes of interest are captured by a primary antibody (Figure 1A). The captured analytes

33 are subsequently labelled with a secondary antibody which in turn is detected with a

34 fluorochrome conjugated tertiary antibody. The level of fluorochrome intensity is directly

35 related to the amount of bound tertiary antibody, and therefore also to the amount of

36 analyte present in the sample. In a multiplex bead assay, the primary antibody is fixed on a

37 polystyrene bead, and physical properties such as size and granularity as well as

38 fluorescent colors of the beads are used to distinguish the different analytes studied. The

39 data is usually collected using a standard flow cytometer.

40 The LEGENDplex system from BioLegend, the CBA system from BD Biosciences, and the

41 MACSPlex system from Miltenyi Biotec are all bead based multiplex systems (Morgan et al.,

42 2004; Miltenyi Biotec, 2014; Yu et al., 2015). The systems differ slightly in terms of physical

43 properties and colors used, and in the number of analytes that can be simultaneously

44 identified. The Bio-Plex system from Bio-Rad works in a similar manner as those described

45 here, but requires a dedicated instrument and does not produce files suitable for analysis

46 with `beadplexr`. The individual assays that can be analyzed with `beadplexr` are described in

47 the following.

48 **LEGENDplex:** Beads fall into two large groups based on size and granularity – as related to

49 the forward light scattering, FSC, and the perpendicular light scatter, SSC. Within each

50 group, individual analytes are discriminated by the intensity of Allophycocyanin (APC) of

51 the beads. The concentration of the analyte is related to the intensity of Phycoerythrin (PE).

52  **CBA:** All beads have similar size and granularity. The individual analytes are discriminated
53  by the intensity of APC and APC-Cy7 of the bead. The concentration of the analyte is related
54  to the intensity of PE.

55  **MACSPlex:** All beads have similar size and granularity. The individual analytes are
56  discriminated by the intensity of PE and Fluorescein isothiocyanate (FITC) of the bead. The
57  concentration of the analyte is related to the intensity of APC.

58  All multiplex systems come with their own analysis software. However, these solutions
59  might come with an added price tag because of binding to a particular piece of software, or
60  the license is valid only for a number of bead events. In this case, large data files with many
61  bead events or repeated re-evaluation of the acquired data might result an expiration of the
62  license. In addition, the usability and flexibility of the analysis solutions are restricted and
63  often impractical for experiments with a large number of samples. Currently no open
64  source alternative exists.

65  Here the general usage of the `beadplexr` package for R (R Core Team, 2018) is introduced. It
66  will be demonstrated how to load the files generated by the flow cytometer, identify bead
67  populations, draw standard curves and calculate concentration of the experimental
68  samples.

## Materials & Methods

70  The `beadplexr` package includes data from an unpublished "Human Growth Factor Panel
71  (13-plex)" LEGENDplex (BioLegend) experiment performed in our laboratory. The dataset
72  consists of eight controls samples and a serum sample from a single healthy volunteer. All
73  samples were processed in duplicates and per manufacturer's instructions. The data was
74  acquired on a CytoFLEX cytometer (Beckman Coulter). An example of a flow cytometry
75  data file is also included in the package. We utilize these data to illustrate the functionality
76  of the package.

77  The data here were analyzed with R, version 3.5.1, (R Core Team, 2018) and plots created
78  with `ggplot2` (Wickham, 2009) and `cowplot` (Wilke, 2017). The workflow and examples

79    presented here make use of or suggests the following R-packages: `devtools` (Wickham,

80    Hester & Chang, 2018), `dplyr` (Wickham et al., 2018), `hexbin` (Carr et al., 2018), `magrittr`

81    (Bache & Wickham, 2014), `purr` (Henry & Wickham, 2018), `stringr` (Wickham, 2018), and

82    `tidyr` (Wickham & Henry, 2018).


## 83    Results

### 84    Package overview

85    The released package can be installed from CRAN and the development version from

86    GitLab:

```
87    # Installing the package -------------------------------------------------
88    # From CRAN
89    install.packages("beadplexr")
90    # From GitLab using devtools
91    # install.packages("devtools")
92    # devtools::install_git("https://gitlab.com/ustervbo/beadplexr")
93    #
94    # Or with vignettes built
95    # devtools::install_git("https://gitlab.com/ustervbo/beadplexr",
96    #                       build_vignettes = TRUE)
97
```

98    The package provides several steps to extract the analyte concentration from the raw data

99    (Figure 1B). The functions for interacting with the data are flexible, but sensible defaults

100   make them accessible to the novice R-user. The workflow and examples presented here are

101   collected in Script S1, and a more detailed workflow is presented in the package vignette.

102   The latter can be viewed using the command `vignette("legendplex-analysis")`.

### 103   Reading FCS-files

104   `beadplexr` works with Flow Cytometry Standard (FCS) files (Seamer et al., 1997), which is

105   the usual output of a flow cytometer. The function `read_fcs()` loads the given FCS-file using

106   the functionality provided by the Bioconductor package `flowcore` (Ellis et al., 2017) and

107   performs the following steps:

108     1.  Apply an *arcsinh* transformation of the bead channels – this natural logarithm based

109         transformation generally performs well on all flow cytometry data (Finak et al.,

110         2010). Opposed to the traditionally used *log10* scaling of flow cytometry data, the

111         *arcsinh* can deal with the negative values produced by some newer digital flow

112         cytometers

113     2.  Remove boundary events of the size (FCS) and granularity (SSC) channels – events

114         outside the range of the detectors are registered with the maximum value possible.

115         These events can interfere with the clustering

116     3.  Optionally subset the channels to contain just bead events – similar to removal of

117         boundary events, this might improve identification of the bead clusters

118     4.  Convert the FCS-data to a `data.frame`

```
119   # Reading fcs-files -----------------------------------------------------
120   library(beadplexr)
121
122   # Get the path to the example fcs-file
123   .file_name <- system.file("extdata",
124                             "K2-C07-A7.fcs",
125                             package = "beadplexr")
126
127   # `read_fcs()` requires at least a path and file name of the file to load,
128   # by identifying the required forward and side scatter and the bead
129   # property channels, only the required data is returned.
130   #
131   # The argument `.filter` takes a named list, where each element is a size
132   # two vector, giving the lower and upper cut-off for the channel given in
133   # the element name
134   .data <- read_fcs(
135     .file_name = .file_name,
136     .fsc_ssc = c("FSC-A", "SSC-A"),
137     .bead_channels = c("FL6-H", "FL2-H"),
138     .filter = list(
139       "FSC-A" = c(3.75e5L, 5.5e5L),
140       "SSC-A" = c(4e5L, 1e6L),
141       "FL6-H" = c(7L, Inf)
```

142  )

143 )

144

145 Because of the variation in detector settings between flow cytometers, it is left to the user

146 to get the event filtering settings correct for an experiment. However, the event filtering

147 should remain stable once established. This, of course, requires that there is no change of

148 cytometer, and that there is no particular drift in the used cytometer. Visualizing the

149 populations greatly helps in setting the appropriate cut-offs (Figure 2). It is for this reason

150 that the `ggplot2` based convenience function `facs_plot()` is included.

### Naming the FCS-files

152 Each sample in a multiplex bead assay must have a unique and meaningful name. A later

153 step in the workflow separates standard samples from experimental samples. The standard

154 samples are in addition ordered in a way that calculation of dilution of standard

155 concentrations is possible. For the dataset included in the package, 'C' followed by an

156 integer denotes the standard (control) samples – as suggested in the LEGENDplex manual –

157 and 'S' followed by an integer denotes the experimental samples. The different parts of the

158 file name should be separated by a character not used in the IDs; this will make for easy

159 parsing of the file names.

### Identification of analyte MFI

161 The mean fluorescence intensity (MFI) of each analyte relates directly to the concentration

162 of the analyte in the sample (Figure 1A). The first step to calculate the analyte

163 concentration is to identify the bead populations representing the analytes and calculate

164 the MFIs of these.

165 `beadplexr` makes use of structured Panel Information to provide analyte metadata such as

166 name and start concentration for each standard sample, as well as the name of the panel,

167 the fold dilution of the standards, and the units of the analytes. The desired Panel

168 Information is loaded using the `load_panel()` function by passing the name or a name

169 pattern to the function. The package itself comes with a set of LEGENDplex Panel

170 Information, which are documented in the help files to `load_panel()`. The Panel Information

171 file itself is in YAML format, and the `load_panel()` function can also load a Panel

172  Information file located outside the package. The latter is useful in the cases of custom

173  panels. The Panel Information is not required, but makes sense if the assay is repeated

174  across several projects.

175

176  ```
# Libraries ---------------------------------------------------------
```
177

178  ```
library(beadplexr)
```
179  ```
library(ggplot2)
```
180  ```
library(cowplot)
```
181  ```
library(dplyr)
```
182  ```
library(purrr)
```
183  ```
library(tidyr)
```
184  ```
library(readr)
```
185  ```
library(stringr)
```
186

187

188  ```
# Load data ---------------------------------------------------------
```
189

190  ```
data(lplex)
```
191  ```
# Load one of the panels distributed with the package, see ?load_panel() for
```
192  ```
# the included panels
```
193  ```
panel_info <- load_panel(.panel_name = "Human Growth Factor Panel (13-plex)")
```
194

195  Analytes of any assay system are identified using the function `identify_analyte()`, which

196  identifies analyte clusters and assign an analyte ID to each cluster. The function takes a

197  `data.frame` with events and a character vector giving the name of column(s) where the

198  analytes can be discriminated. An identifier for each analyte is passed in the argument

199  `.analyte_id`, which is simply a character vector giving the ID of the analyte.

200  `identify_analyte()` sorts the clusters based on their centers and use this ranking to assign

201  the analyte IDs. The order of analyte IDs given in `.analyte_id` is therefore important and must

202  match the expected order of analytes. An optional argument is `.trim` which allows the

203  removal events in the periphery of a cluster. The value of the argument gives the fraction of

204    the most distant points to be removed. Distance based trimming is non-trivial since the

205    possible numerical range depends on the detection range of the flow cytometer.

206    The function `identify_analyte()` interfaces several methods for unsupervised clustering,

207    which are passed in the `.method` argument.  The default clustering method is clustering

208    large applications (`clara`) from the package `cluster` (Maechler et al., 2017). The method

209    selects a number of subsets of fixed size and applies the partitioning around medoids

210    (`pam`)-algorithm to each subset. The objective of the pam-algorithm is to minimize the

211    dissimilarity between the representative of $k$ clusters and the members of each cluster

212    (Kaufman & Rousseeuw, 2009). The best resulting set of medoids (cluster centers) is that

213    with the lowest average dissimilarity of all points in the original dataset to the medoids.

214    Though similar to pam in algorithm type, the Base-R included `kmeans` works on minimizing

215    the distance to the cluster representative (Zaki & Wagner Meira, 2014).

216    The `dbscan` method in the `fpc` package differs from `clara` and `kmeans` in that `dbscan` identifies

217    clusters based local density (Hennig, 2015). The function requires a neighborhood size and

218    minimum number of events in each neighborhood to evaluate whether points can be

219    considered as belonging to a cluster (Zaki & Wagner Meira, 2014). If the bead populations

220    have different local densities, there is no guarantee that the correct number of clusters will

221    returned. This problem does not exist for `Mclust` from the `mclust` package, which fits a

222    Gaussian mixture model using the EM-algorithm (Scrucca et al., 2016).  This algorithm

223    iteratively optimizes the individual parameters of k normal distributions (Zaki & Wagner

224    Meira, 2014). This way the relationship between a cluster and a set of data points is given

225    by a set of probability scores.

226    We have found that `dbscan()` is the best clustering method for the forward-side scatter

227    population identification. However, it can be difficult to get the parameters *event count* and

228    *neighborhood size* correct. The reason for this difficulty lies in the sensitivity of the method

229    to the choice of *neighborhood size*; if it is too large clusters might be merged, and if it is too

230    small everything might be classified as noise. In our experience, the clustering function

231    `clara()` is a great all-rounder although the subsampling performed by the function can lead

232    to slight differences between each run. Using the same value for `set.seed()` at the

233    beginning of each session will alleviate this and make each run reproducible.

234    Different flow cytometers perform differently in terms of separation of the individual bead

235    populations. This is due to factors such as detector settings and age of the cytometer and its

236    light sources. The consequence is that the populations of interest might be closer together

237    or further apart. Another consequence might be an increased in the noise of the detectors

238    of the flow cytometer. Collectively these differences in the data constitution means that one

239    clustering function might perform better on one dataset while be inferior on another. As

240    with analysis of all flow cytometric data the optimal solution is a matter of taste, but the

241    better clustering function is the one that separates the populations well, without including

242    too much noise.

243    The function `identify_legendplex_analyte()` can be applied to each sample individually in a

244    loop. However, it is more prudent to apply the function to all samples at the same time

245    because the clustering decision will be identical for each sample. In addition, clustering on

246    all the samples is 1.4 times faster than clustering on each sample individually.

```
247   # Identify analytes ---------------------------------------------------
248
249   # The function `identify_legendplex_analyte()` used here is convenience
250   # around the clustering work horse `identify_analyte`. The
251   # `identify_legendplex_analyte()` identifies the bead populations according
252   # to size and granularity, and for each of the two populations the individual
253   # bead populations are identified
254   #
255   # The function requires a named list with analytes from the Panel
256   # Information, and a list with a list of key-value pairs giving the arguments
257   # for the bead identification on the forward and side scatter, and a list of
258   # key-value pairs giving arguments for the bead identification in each
259   # subpopulation in the APC channel.
260   #
261   # The argument .trim gives the fraction of events furthest from the centers of
262   # the groups that should be removed. The population center is found by a
263   # Gaussian kernel estimate. In this case we remove 1% and 3% of the of the
```

```
264  # events based on their distance to the group center.
265  #
266  # The inner lists can be named, but this is not required.
267  args_ident_analyte <- list(fs = list(.parameter = c("FSC-A", "SSC-A"),
268                                        .column_name = "Bead group",
269                                        .trim = 0.01),
270                             analyte = list(.parameter = "FL6-H",
271                                            .column_name = "Analyte ID",
272                                            .trim = 0.03))
273
274  # The FCS-data is a list of samples, which we combine before cluster
275  # identification.
276  analytes_identified <- lplex %>%
277    bind_rows(.id = "Sample") %>%
278    identify_legendplex_analyte(.analytes = panel_info$analytes,
279                                .method_args = args_ident_analyte)
280
```

281  The analyte IDs for the "Human Growth Factor Panel (13-plex)" bead group A are A4, A5,
282  A6, A7, A8, A10 and for group B the analyte IDs are B2, B3, B4, B5, B6, B7, B9. In this case,
283  the beads are arranged from low to high, that is the lowest analyte ID has lowest intensity
284  in the APC channel (Figure 3).

285  This initial and crucial step of the analysis has been successfully performed with data from
286  a CBA experiment (C. McGuckin, CTIBIOTECH, Lyon, France, unpublished) and from a
287  MACSPlex experiment (Miltenyi Biotec, Bergisch Gladbach, Germany, unpublished) using
288  the function `identify_analyte()`.

289  With the analytes identified and the bead populations documented, the MFI of each analyte
290  can finally be calculated. The function `calc_analyte_mfi()` gives the possibility to calculate
291  geometric, harmonic, and arithmetic mean of the in intensity of each respective analyte
292  reporter, such as PE in a LEGENDplex assay. Since the reporter intensities are usually log-
293  transformed only the geometric mean is relevant, but harmonic and arithmetic mean are
294  included to accommodate for special cases.

```
295  # Calculate analyte MFI ---------------------------------------------
```

```
296
297   # The mean fluorescence intensity is calculated for each sample and analyte.
298   # The function `calc_analyte_mfi()` provides three ways of calculating the
299   # MFI: geometric, harmonic, and arithmetic mean.
300   analyte_mfi <- analytes_identified %>%
301     filter(!is.na(`Analyte ID`)) %>%
302     # Call `calc_analyte_mfi()` for each sample
303     group_by(Sample) %>%
304     do(calc_analyte_mfi(., .parameter = "FL2-H",
305                      .column_name = "Analyte ID",
306                      .mean_fun = "geometric")) %>%
307     # Later we will fit the standard curve on a log-log scale, so we transform
308     # here
309     mutate(`FL2-H` = log10(`FL2-H`))
```

## Calculation of standard and experimental samples

The calculation of the concentration of the analytes of the experimental samples requires two steps:

1.  Create a standard curve by fitting a model to the MFI of the standard analytes and their known concentrations
2.  Estimate the concentration of each sample analyte from the fitted model

The samples in the dataset included in the package can be distinguished by the presence of 'C' or 'S', respectively. The sample type indicating letter is then followed by one or more integers. Using this naming scheme, it is easy to separate standard samples from the experimental samples. It is also easy to order the standard samples for concentration assignment. In this case the naming scheme suggested in the LEGENDplex assay protocol is followed: 7 indicates the highest concentration of the standard analyte, 1 indicates the lowest concentration, and 0 indicates blank.

The order of the standard samples is crucial for the function `calc_std_conc()` to correctly calculate the concentration of an analyte in each standard sample. The function requires a

326   vector which gives the order of the standard samples, a start concentration for the analyte,

327   and a dilution factor. The standard samples are ordered numerically from high to low and

328   assigned a standard concentration, such that the first sample is given the start

329   concentration and the second to last sample the lowest concentration, and the very last

330   sample the concentration 0, as this is assumed to be for background measurement.

331   The start concentration is stored in the Panel Information for each analyte separately, as

332   the start concentration might differ from analyte to analyte. The dilution factor is also given

333   in the Panel Information. It will always be the same for all standard analytes and is usually

334   4, meaning that the concentration of each standard analyte is 4 times lower than the

335   previous concentration. This generally gives a good range of standard concentrations.

```
336   # Helper function to extract the sample number ----------------------------
337
338   #' Cast sample ID to numeric
339   #'
340   #' @param .s A string with the sample ID pattern to be cast
341   #' @param .pattern A string giving the pattern
342   #'
343   #' @return
344   #' A numeric
345   #'
346   as_numeric_sample_id <- function(.s, .pattern = c("C[0-9]+", "S[0-9]+")){
347     .pattern <- match.arg(.pattern)
348
349     # Extract the pattern defined just above, remove the first element, and
350     # cast to a numeric
351     .s %>%
352       str_extract(.pattern) %>%
353       str_sub(start = -1L) %>%
354       as.numeric()
355   }
356
357   # Split in standard and sample --------------------------------------------
358
359   # We need to fit a standard curve on the standard samples, and use this curve
```

```
360   # to calculate the concentration of the experimental samples. Here we split
361   # the data set in two: one with the standard samples and one with the
362   # experimental samples.
363   #
364   # We need to order the standard samples from high to low in order to
365   # calculate the concentration of the analytes in the standard sample.
366   # Incorporating the information into the sample name in terms of an easily
367   # parsable pattern is a good practice.
368
369   # All standard samples have the pattern C[number]
370   standard_data <- analyte_mfi %>%
371     ungroup() %>%
372     filter(str_detect(Sample, "C[0-9]+")) %>%
373     mutate(`Sample number` = as_numeric_sample_id(Sample, . pattern = "C")) %>%
374     select(-Sample)
375
376   # All non-standards are experimental samples... we could also filter on
377   # S[number]
378   experiment_data <- analyte_mfi %>%
379     ungroup() %>%
380     filter(!str_detect(Sample, "C[0-9]+")) %>%
381     mutate(`Sample number` = as_numeric_sample_id(Sample, . pattern = "S")) %>%
382     select(-Sample)
383
384   # To the standard data we have to add additional information such the start
385   # concentration of each standard analyte and the dilution factor, as well as
386   # as the analyte names (analyte IDs by themselves do not make much sense).
387   #
388   # The concentration of the standard samples is calculated using
389   # `calc_std_conc()`, which take a vector of sample numbers for ordering, a
390   # start concentration and a dilution factor.
391   standard_data <- standard_data %>%
392     left_join(as_data_frame_analyte(panel_info$analytes), by = "Analyte ID") %>%
393     rename(`Analyte name` = name) %>%
394     group_by(`Analyte ID`, `Analyte name`) %>%
395     mutate(
396       Concentration = calc_std_conc(
```

```
397        `Sample number`,
398        concentration,
399        .dilution_factor = panel_info$std_dilution
400      )
401    ) %>%
402    # Later we will fit the standard curve on a log-log scale, so we transform
403    # here
404    mutate(Concentration = log10(Concentration)) %>%
405    select(-concentration, -`Bead group`)
406
```

407 The next step is to fit a standard curve for each analyte. With the standard curve we can

408 calculate the concentration of the experimental samples (the purpose of the initial work),

409 we can check the quality of the measurements and the standard curve, and plot the

410 experimental samples on the standard curve (beadplexr provides easy access to all of this).

411 The latter is to allow for visual verification that the experimental samples are within the

412 linear part of the standard curve.

413 However, in each case we need to ensure that the correct standard curve is used with the

414 correct experimental data, which means we have to juggle at least three structures: A

415 data.frame with the standard data, a data.frame with the experimental sample data, and the

416 models for each analyte (probably a list). It quickly becomes tedious to ensure that

417 everything is in the correct order - and it is most certainly error prone. To circumvent this,

418 we can use the nest() and its inverse unnest() functions of the tidyr package. nest() relies

419 the fact that a data.frame in R is in fact a list, and uses this to pack a data.frame into a

420 single cell of a data.frame.

```
421  combine # Nest standard and experimental data ------------------------
422
423  # Nested data.frames is a great way of combining and working with complex
424  # data structures.
425  #
426  # First we pack all the standard data in to a data.frame with a set of
427  # data.frames
428  standard_data <- standard_data %>%
429    nest(-`Analyte ID`, .key = "Standard data")
```

```
430
431   # The the same for all the experimental data
432   experimental_data <- experiment_data %>%
433     nest(-`Analyte ID`, .key = "Experimental data")
434
435   # Since both structures are data.frames we can easily combine them
436   plex_data <- inner_join(standard_data, experiment_data, by = "Analyte ID")
437
```

With everything in a neatly arranged `data.frame` we can now focus on the actual task at hand, namely calculation of the standard curve for each analyte. For this we use the function `fit_standard_curve()`, which interfaces the `drm()` function from the `drc` package (Ritz et al., 2015). The `drm()` function specializes in fitting various biological response-models, and the `drc` package provides several response-models, such as the four- and five-parameter log-logistic model. `fit_standard_curve()` is designed to be used in the piped workflow, and takes a `data.frame` with MFIs and concentrations and returns the model as a `drc` object. The four-parameter log-logistic model is widely used in analysis of ELISA data. Since the five-parameter model yields better fits, because of the increased flexibility, this is the default function (Gottschalk & Dunn, 2005).

```
448   # Calculate standard curves ----------------------------------------------
449
450   # For each of the analytes we calculate the standard curve. Working with
451   # nested data.frames means that we have to loop over each row to calculate
452   # the standard curve using the data.frame in "Standard data"
453   #
454   # When clustering is performed with mclust, the package mclust is
455   # loaded in the background (an unfortunate necessity). The mclust
456   # package also has a function called `map`, so an unlucky side effect
457   # of clustering with mclust, is that we need to be specify which map
458   # function we use.
459
460   plex_data <- plex_data %>%
461     group_by(`Analyte ID`) %>%
462     mutate(`Model fit` = purrr::map(`Standard data`,
463                                      fit_standard_curve,
```

464                          `.parameter = "FL2-H"))`

465

466     We can plot the standard curve using the built in `plot_std_curve()` function (Figure 4A).

467     With the standard curve created we can calculate the concentrations of the experimental

468     samples using the function `calculate_concentration()`, which requires a `data.frame` with

469     the MFIs in a column, and the fitted model. It can be helpful to apply

470     `calculate_concentration()` to the standard samples, as this can be used to verify that the

471     standard measurements were all fine, and that the estimation of the sample concentrations

472     therefore is trustworthy.

473     After calculating the concentrations we can plot the known standard concentrations versus

474     the estimated standard concentrations using the function `plot_target_est_conc()` (Figure

475     4B) and visualize where the samples fall on the standard curve with `plot_estimate()`

476     (Figure 4C).

```
477   # Calculate experimental sample concentrations ----------------------------
478
479   # Using the standard curve just calculated, we can back-calculate the
480   # concentration of the standard concentrations, and more importantly the
481   # concentration of the experimental samples
482   plex_data <- plex_data %>%
483   mutate(`Standard data` =
484           purrr::map2(`Standard data`, `Model fit`,
485                       calculate_concentration, .parameter = "FL2-H")) %>%
486     mutate(`Experimental data` =
487             purrr::map2(`Experimental data`, `Model fit`,
488                       calculate_concentration, .parameter = "FL2-H"))
489
490
491   # Add concentration plots -------------------------------------------------
492
493   # We can also loop over each row and add plots to the data.frame
494   plex_data <-  plex_data %>%
495     mutate(`Std curve` =
496             purrr::pmap(list(.data = `Standard data`,
```

```
497                          .model = `Model fit`,
498                          .title = `Analyte name`),
499                     plot_std_curve, .parameter = "FL2-H")) %>%
500    mutate(`Std conc` =
501           purrr::map(`Standard data`,
502                      plot_target_est_conc)) %>%
503    mutate(`Est curve` =
504           purrr::pmap(list(`Experimental data`,
505                            `Standard data`,
506                            `Model fit`,
507                            `Analyte name`),
508                     plot_estimate, .parameter = "FL2-H"))
509
```

510  Lastly we fulfill the purpose of all the previous actions and extract the concentration of

511  each analyte for each sample.

```
512  # Extract analyte concentration -------------------------------------------
513
514  plex_data %>%
515    unnest(`Experimental data`) %>%
516    # Make the names a little more telling and transform them back to useful
517    # concentrations
518    rename(`Concentration (pg/ml)` = Calc.conc,
519           `Concentration error` = `Calc.conc error`) %>%
520    mutate(`Concentration (pg/ml)` = 10^ `Concentration (pg/ml)`,
521           `Concentration error` = 10^`Concentration error`)
```

## Discussion

523  Multiplex bead assays make simultaneous evaluation of several analytes possible. Because

524  of this, they are an attractive alternative to the commonly used sandwich ELISA.

525  Commercial systems are available for acquisition on a standard flow cytometer, but these

526  commercial systems make use of their own proprietary software for the data analysis. This

527  can impose different limitations to the analysis. The R-package `beadplexr`, released under

528  the MIT license, is meant as an open-source alternative to these commercial systems. The

529  package is available from CRAN and from https://gitlab.com/ustervbo/beadplexr.

530 A critical step in the analysis multiplex bead assays is the identification of bead populations
531 corresponding to each analyte. A single function in `beadplexr` acts as an interface to several
532 common, and tested, clustering functions, making it easy to find the best suited clustering
533 function. Future versions of the package will see improvements in this part, with inclusion
534 of other clustering methods and perhaps a heuristic for automatic method selection.

535 Flow cytometry data are inherently noisy. `beadplexr` only provides a rudimentary function
536 for removing points with no neighbors and lets the clustering functions determine which
537 events are considered noisy though the `.trim` argument. However, a very noisy data set
538 might make it difficult for an optimal identification of the bead clusters in the first place.
539 De-noising multidimensional data is not trivial, but work is planned in this direction for a
540 future release.

## Conclusion

542 The R-package `beadplexr` provides a frame work for easy and reproducible analysis of
543 multiplex bead assays for the experienced and the novice user alike.

## Acknowledgments

## References

550 Bache SM., Wickham H. 2014. *magrittr: A Forward-Pipe Operator for R*.

551 Carr D., Lewin-Koh  ported by N., Maechler M., Sarkar  contains copies of lattice functions
552 written by D. 2018. *hexbin: Hexagonal Binning Routines*.

553 Ellis B., Haaland P., Hahne F., Meur NL., Gopalakrishnan N., Spidlen J., Jiang M. 2017.
554 *flowCore: flowCore: Basic structures for flow cytometry data*.

555  Elshal MF., McCoy JP. 2006. Multiplex Bead Array Assays: Performance Evaluation and

556      Comparison of Sensitivity to ELISA. *Methods (San Diego, Calif.)* 38:317–323. DOI:

557      10.1016/j.ymeth.2005.11.010.

558  Finak G., Perez J-M., Weng A., Gottardo R. 2010. Optimizing transformations for automated,

559      high throughput analysis of flow cytometry data. *BMC Bioinformatics* 11:546. DOI:

560      10.1186/1471-2105-11-546.

561  Gottschalk PG., Dunn JR. 2005. The five-parameter logistic: a characterization and

562      comparison with the four-parameter logistic. *Analytical Biochemistry* 343:54–65.

563      DOI: 10.1016/j.ab.2005.04.035.

564  Hennig C. 2015. *fpc: Flexible Procedures for Clustering*.

565  Henry L., Wickham H. 2018. *purrr: Functional Programming Tools*.

566  Kaufman L., Rousseeuw PJ. 2009. *Finding Groups in Data: An Introduction to Cluster

567      Analysis*. John Wiley & Sons.

568  Maechler M., Rousseeuw P., Struyf A., Hubert M., Hornik K. 2017. *cluster: Cluster Analysis

569      Basics and Extensions*.

570  Miltenyi Biotec 2014. *Data acquisition and analysis without the MACSQuant® Analyzer --

571      General instructions for MACSPlex Cytokine Kits*. Miltenyi Biotec.

572  Morgan E., Varro R., Sepulveda H., Ember JA., Apgar J., Wilson J., Lowe L., Chen R., Shivraj L.,

573      Agadir A., Campos R., Ernst D., Gaur A. 2004. Cytometric bead array: a multiplexed

574      assay platform with applications in various areas of biology. *Clinical Immunology*

575      110:252–266. DOI: 10.1016/j.clim.2003.11.017.

576  R Core Team 2018. *R: A Language and Environment for Statistical Computing*. Vienna,

577      Austria.

578    Ritz C., Baty F., Streibig JC., Gerhard D. 2015. Dose-Response Analysis Using R. *PLOS ONE* 10.

579    Scrucca L., Fop M., Murphy TB., Raftery AE. 2016. mclust 5: clustering, classification and

580        density estimation using Gaussian finite mixture models. *The R Journal* 8:205–233.

581    Seamer LC., Bagwell CB., Barden L., Redelman D., Salzman GC., Wood JCS., Murphy RF. 1997.

582        Proposed new data file standard for flow cytometry, version FCS 3.0. *Cytometry*

583        28:118–122. DOI: 10.1002/(SICI)1097-0320(19970601)28:2<118::AID-

584        CYTO3>3.0.CO;2-B.

585    Wickham H. 2009. *ggplot2: elegant graphics for data analysis*. Springer.

586    Wickham H. 2018. *stringr: Simple, Consistent Wrappers for Common String Operations*.

587    Wickham H., François R., Henry L., Müller K. 2018. *dplyr: A Grammar of Data Manipulation*.

588    Wickham H., Henry L. 2018. *tidyr: Easily Tidy Data with "spread()" and "gather()" Functions*.

589    Wickham H., Hester J., Chang W. 2018. *devtools: Tools to Make Developing R Packages Easier*.

590    Wilke CO. 2017. *cowplot: Streamlined Plot Theme and Plot Annotations for "ggplot2."*

591    Yu Y., Wang C., Clare S., Wang J., Lee S-C., Brandt C., Burke S., Lu L., He D., Jenkins NA.,

592        Copeland NG., Dougan G., Liu P. 2015. The transcription factor Bcl11b is specifically

593        expressed in group 2 innate lymphoid cells and is essential for their development.

594        *The Journal of Experimental Medicine* 212:865–874. DOI: 10.1084/jem.20142318.

595    Zaki MJ., Wagner Meira J. 2014. *Data Mining and Analysis: Fundamental Concepts and*

596        *Algorithms*. Cambridge University Press.

597

598    **Figure legends**

599    **Figure 1: Overview of assay principle and the package workflow.**

600    A) Schematic overview of the principle of a LEGENDplex assay. B) Steps in analysis of a

601    multiplex bead assay with accompanying visualizations.

602    **Figure 2: Visualization of flow cytometry data.**

603    Size (FSC) and granularity (SSC) can be used distinguish the two LEGENDplex bead

604    populations. A) Common monochrome scatter-plot created with `facs_plot(.x = "FSC-H",`

605    `.y = "SSC-H", .beads = "Bead group")` on the sample 'K3-C0-1.fcs'. High density regions

606    are obscured in this type of plots. B) Pseudo-colored scatter -plot created with

607    `facs_hexbin(.x = "FSC-H", .y = "SSC-H", .beads = "Bead group", .bins = 75)` on the

608    same sample as in A). The number of events in discrete bins is indicated by color. The

609    coloring is according to the standard blue-green-yellow-red scheme, where blue indicates a

610    low number of events, and red indicates a high number. The Pseudo-colored scatter -plot

611    requires the R-package `hexbin` to be installed.

612    **Figure 3: Bead identification and visualization of LEGENDplex data.**

613    Populations identified in the sample 'K3-C0-1.fcs'. A) Identification of the two bead

614    populations 'A' and 'B' according to size and granularity: The two clusters were identified

615    using `.method = clara` and noisy data points were excluded by `.trim = 0.01`. B-C)

616    Identification of analytes of the bead population 'A' and 'B': The 1 dimensional clusters

617    along the APC channel were identified using `.method = clara` and noisy data points were

618    excluded by `.trim = 0.03`. Noisy data points are assigned the group 'NA'.

619    **Figure 4: Visualization of standard and experimental samples for Angiopoietin-2.**

620    The dataset included in `beadplexr` is from a 13-plex assay. Here we use Angiopoietin-2 to

621    illustrate the visualizations. A) A log-log plot of the standard curve of Angiopoietin-2. Each

622    point is a single measurement (each in duplicate). The standard concentration is diluted in

623    steps of four fold dilution from 50,000.0 to 12.21 pg/ml. The intensity of the analyte is

624    measured in the PE channel. The full line indicates the best fit, and gray the confidence

625    interval. B) Correlation between the standard concentration (x-axis) and the calculated

626    concentration of the standard samples (y-axis).  The back calculation is done using the fit in

627    A) and the MFI of the samples. C) Using the fit in A) the concentration of an experimental

628    sample is calculated. Visual inspection of the position of the experimental samples on the

629    standard curve can reveal samples that are close to the upper or lower bound of the
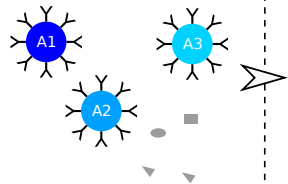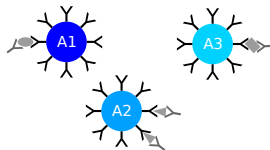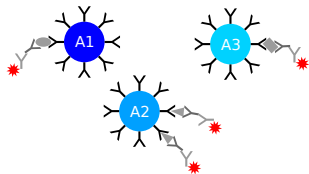
630    standard curve.

631

**Figure 1**(on next page)

Overview of assay principle and the package workflow

A) Schematic overview of the principle of a LEGENDplex assay. B) Steps in analysis of a multiplex bead assay with accompanying visualizations.
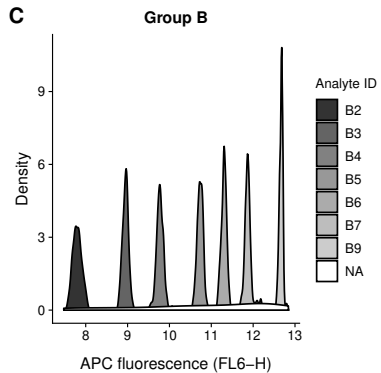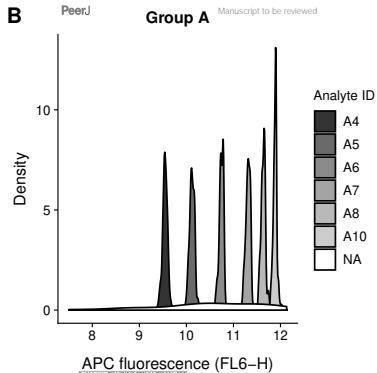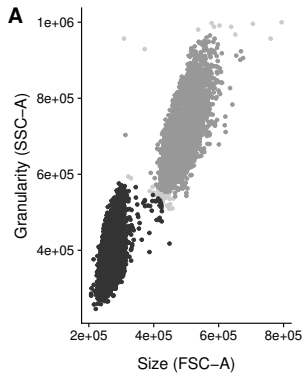
**Figure 2**(on next page)

Visualization of FACS data

Size (FSC) and granularity (SSC) can be used distinguish the two LEGENDplex bead

populations. A) Common monochrome scatter-plot created with facs_plot(.x = "FSC-H", .y =

"SSC-H", .beads = "Bead group") on the sample 'K3-C0-1.fcs'. High density regions are

obscured in this type of plots. B) Pseudo-colored scatter -plot created with facs_hexbin(.x =

"FSC-H", .y = "SSC-H", .beads = "Bead group", .bins = 75) on the same sample as in A). The

number of events in discrete bins is indicated by color. The coloring is according to the

standard blue-green-yellow-red scheme, where blue indicates a low number of events, and

red indicates a high number. The Pseudo-colored scatter -plot requires the R-package hexbin

to be installed.

**Figure 3**(on next page)

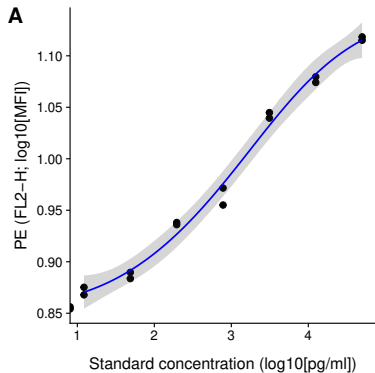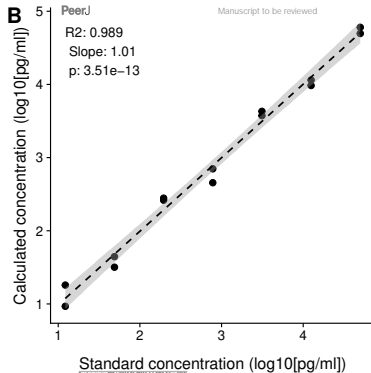Bead identification and visualization of LEGENDplex data

Populations identified in the sample 'K3-C0-1.fcs'. A) Identification of the two bead populations 'A' and 'B' according to size and granularity: The two clusters were identified using .method = clara and noisy data points were excluded by .trim = 0.01. B-C) Identification of analytes of the bead population 'A' and 'B': The 1 dimensional clusters along the APC channel were identified using .method = clara and noisy data points were excluded by .trim = 0.03. Noisy data points are assigned the group 'NA'.

**A**

Granularity (SSC-A) vs Size (FSC-A)

Bead group
- A
- B
- NA

**B** Group A

Density vs APC fluorescence (FL6-H)

Analyte ID
- A4
- A5
- A6
- A7
- A8
- A10
- NA

**C** Group B

Density vs APC fluorescence (FL6-H)

Analyte ID
- B2
- B3
- B4
- B5
- B6
- B7
- B9
- NA

**Figure 4**(on next page)

Visualization of standard and test samples for Angiopoietin-2

The dataset included in beadplexr is from a 13-plex assay. Here we use Angiopoietin-2 to illustrate the visualizations. A) A log-log plot of the standard curve of Angiopoietin-2. Each point is a single measurement (each in duplicate). The standard concentration is diluted in steps of four fold dilution from 50,000.0 to 12.21 pg/ml. The intensity of the analyte is measured in the PE channel. The full line indicates the best fit, and gray the confidence interval. B) Correlation between the standard concentration (x-axis) and the calculated concentration of the standard samples (y-axis). The back calculation is done using the fit in A) and the MFI of the samples. C) Using the fit in A) the concentration of an experimental sample is calculated. Visual inspection of the position of the experimental samples on the standard curve can reveal samples that are close to the upper or lower bound of the standard curve.

**A**

PE (FL2-H; log10[MFI])

Standard concentration (log10[pg/ml])

**B**

R2: 0.989
Slope: 1.01
p: 3.51e−13

Calculated concentration (log10[pg/ml])

Standard concentration (log10[pg/ml])

**C**

PE (FL2-H; log10[MFI])

Standard concentration (log10[pg/ml])