1    **Single Object Profiles Regression Analysis (SOPRA): A novel method for**
2    **analyzing high content cell-based screens**

3

4    Rajendra Kumar Gurumurthy[2#], Klaus-Peter Pleissner[1#], Cindrilla Chumduri [2], Thomas F.
5    Meyer[2], André P. Mäurer[1*]

6

7    [#]These authors contributed equally to the manuscript

8    [1]Steinbeis Center for Systems Biomedicine, Steinbeis Innovation gGmbH, 14612 Falkensee,
9    Germany

10   [2]Max Planck Institute for Infection Biology, Department of Molecular Biology, 10117 Berlin,
11   Germany

12

13   *To whom correspondence should be addressed:

14   Klaus-Peter Pleissner

15   Max Planck Institute for Infection Biology

16   10117 Berlin, Germany

17   E-Mail: pleissner@mpiib-berlin.mpg.de

18   Tel: +49 030 28 460 0

19

20   **Running title:** Single Object Profiles Regression Analysis

21

22

23

## 24 **Abstract**

25 <u>Motivation</u>

26      High content screening (HCS) experiments generate complex data from multiple object
27 features for each cell within a treated population. Usually these data are analyzed by using
28 population-averaged values of the features of interest, increasing the amount of false positives and
29 the need for intensive follow-up validation. Therefore, there is a strong need for novel approaches
30 with reproducible hit prediction by identifying significantly altered cell populations.

31 <u>Results</u>

32 Here we describe SOPRA, a workflow for analyzing image-based HCS data based on regression
33 analysis of non-averaged object features from cell populations, which can be run on hundreds of
34 samples using different cell features. Following plate-wise normalization the values are counted
35 within predetermined binning intervals, generating unique frequency distribution profiles
36 (histograms) for each population, which are then normalized to control populations. Statistically
37 significant differences are identified using a regression model approach. Significantly changed
38 profiles can be used to generate a heatmap from which altered cell populations with similar
39 phenotypes are identified, enabling detection of siRNAs and compounds with the same 'on-target'
40 profile, reducing the number of false positive hits. A screen for cell cycle progression was used to
41 validate the workflow, which identified statistically significant changes induced by siRNA-mediated
42 gene perturbations and chemical inhibitors of different cell cycle stages.

43

44

## Background

46  The availability of robotic liquid handling combined with automated fluorescence microscopy and
47  high-performance image computing has enabled rapid advances in the development of high-
48  throughput screening. Numerous studies have demonstrated the power of high-throughput image-
49  based assays for characterizing drug effects (Perlman, et al., 2004), identifying active small molecules
50  (Tanaka, et al., 2005) and classifying sub-cellular protein localization (Boland and Murphy, 2001;
51  Conrad, et al., 2004), including genome-wide siRNA-mediated loss-of-function screens (Neumann, et
52  al., 2006) or gene deletion (Ohya, et al., 2005) libraries. For each single cell within a cellular sample
53  population, it is possible to achieve quantitative measurements of phenotypes such as expression
54  level and localization of proteins, post-translational modifications and even cellular or sub-cellular
55  morphologies.

56       Analyzing cellular populations in the early drug discovery process allows the complexity of
57  living systems to be addressed and produces vast amounts of data that are more meaningful than
58  those obtained from isolated proteins (Taylor, 2007). In combination with advanced bioinformatics
59  tools treatments can be identified which lead to altered cell populations, and therefore might be
60  relevant drugs or drug targets.

61       Nonetheless, several limitations in data analysis have restricted the full potential of high-
62  throughput image-based assays so far (Lang, et al., 2006; Zhou and Wong, 2006). The usual course of
63  events for a HCS analysis workflow starts with the extraction of image feature data, followed by
64  normalization and statistical analysis, including final hit selection (Buchser, et al., 2004). A wide
65  variety of microscopes, image-analysis and data-analysis software packages are available to address
66  these issues (Gough and Johnston, 2007). However, distributions of multidimensional, multivariate
67  phenotypic measurements from cellular populations are mostly transformed into single population-
68  averaged values such as mean or median values. These population-averaged values are used for
69  plate-wise or batch-wise normalizations, as well as for statistical analysis for hit selection
70  (Birmingham, et al., 2009; Singh, et al., 2014), which leads to a substantial loss of information.
71  Population-averaged values can indicate whether the value of the measured phenotype increases or
72  decreases upon treatment, but do not reflect the detailed response of a cellular population to a
73  certain treatment or gene depletion. Therefore, these population-averaged values are limiting the
74  power of the statistical approaches that are widely used, such as Z-Score or percent-of-control (POC)
75  analysis, making it impossible to identify more distinct reactions of a cell population. This loss of

76    information also hampers the differentiation of treatments or gene depletions with the same 'on-
77    target' effect from those with 'off-target' effects, which is extremely important for RNAi gene
78    perturbation experiments, where multiple siRNAs are used per gene.

79    Some publications have described methods for non-averaged cell population data analysis from high-
80    content image-based screens. Knapp et al. (Knapp, et al., 2011) showed considerable effects of
81    population context on observed phenotypes when using non-averaged population data for the
82    normalization steps, but still used population-averaged values for hit detection. Another method
83    uses multivariate cell classification based on phenotypic changes for hit identification (Loo, et al.,
84    2007), which results in a drug effect score, and a vector, indicating the simultaneous phenotypic
85    changes induced by the drug. Another publication used multi-parametric phenotypic profiles to
86    cluster genes based on morphological changes of individual cells (Fuchs, et al., 2010). Yet another
87    group has proposed the use of Ripley's K-function to identify knockdowns resulting in perturbation of
88    this cell clustering (Suratanee, et al., 2010). Also the Kolomogorow-Smirnov (KS) test has been used
89    to score the difference between control and samples populations (Gorenstein, et al., 2010).
90    However, all these methods have limitations that prevent them from being widely used for large-
91    scale high content cell population analysis. Multivariate classification methods are mostly based on
92    the analysis of predominantly redundant image features, spatial clustering requires a subjective and
93    work intensive classification step for the cellular populations and KS only uses one unique value to
94    identify cell population with altered distributions.

95    Here we present a new approach called Single Object Profiles Regression Analysis (SOPRA) that
96    overcomes many of these limitations by analyzing non-averaged cell population data. It uses a
97    classification free regression analysis of normalized frequency distribution profiles of cell
98    populations. SOPRA can be used to analyze data derived from various high-throughput techniques,
99    such as images from automated microscopy or single cell data from FACS analysis. The regression
100   workflow consists of i) a pre-processing step, ii) data gathering and normalization steps, iii)
101   identification of significant profiles, iv) post-processing. The normalization is performed in a plate-
102   wise and bin-wise fashion, resulting in a unique normalized frequency distribution profile for each
103   feature of a cell population. Finally, normalized distribution profiles that exhibit statistically
104   significant changes are identified by using a p-value and R-squared (RSQ)-value derived from the
105   regression analysis with the R-package maSigPro (Conesa, et al., 2006). Additionally, normalized
106   feature profiles that have been identified as significantly altered can be further clustered in a
107   heatmap according to their similarity. This can be used to identify treatments with the same 'on-
108   target' effects. Most loss-of-function screens use multiple siRNAs for the same gene, which should

109   end up in the same cluster if they have a similar cell population phenotype. The more siRNAs for the
110   same gene are identified as having a similar cell population profile, the more reliably this gene can be
111   regarded as a hit. Beyond this, the derived values of a regression analysis of distribution profiles of
112   cellular features are not affected by experimental bias  to the same degree as population-averaged
113   approaches (Sacher, et al., 2008), leading to more reproducible results. We used a cell-based
114   chemical compound and RNAi screen of cell cycle progression to validate the SOPRA workflow. The
115   cellular features 'Area', 'Total Intensity DAPI' and 'Mean Intensity DAPI' were extracted for each
116   nucleus using image analysis software and subjected to the SOPRA workflow. We found that SOPRA
117   can be used to identify statistically significant changes of frequency distribution profiles within
118   cellular populations, whether induced by gene perturbation through siRNAs or by chemical inhibitor
119   treatment. Taken together, SOPRA is a novel object-based data analysis workflow based on
120   regression analysis of cellular feature distribution profiles to identify significantly changed cell
121   populations from high-throughput data sets.

## 122   Results

123   SOPRA utilizes a data gathering step combined with plate-wise and a so-called bin-wise normalization
124   methods, as well as a two-step regression approach that first adjusts a global regression model with
125   defined variables in order to identify profiles exhibiting statistically significant changes (Conesa, et
126   al., 2006). The SOPRA workflow consists of several steps as outlined in Figure 1. *A: High Content*
127   *Screen.* This first step includes screening, image analysis and data extraction. B: Preparation of screen
128   description files and the single cell data files. *C: Preprocessing (optional).* The derived data files for
129   various image features at single cell level are subjected to a preprocessing step to exclude all data
130   from flagged wells that should be excluded from the analysis. *D1: Data Gathering and Plate-Wise*
131   *Normalization.* In this step each single cell object is annotated with additional information such as
132   RNA.ID, plate number, well number, replicate number, well content and gene symbol. If the imaging
133   software supports a gating procedure for objects that do not meet certain criteria, such as cell size,
134   these can also be flagged and excluded from subsequent analysis steps. The measured value of each
135   cell for the feature of interest is then normalized to the median of the objects in the neutral control
136   wells. *D2: Data Gathering and Frequency Distribution Profiles (Histogram) Generation.* Next, the
137   common binning axis of the distribution profiles is generated by determining the minimum and
138   maximum limits of the measured feature across all the data of the screen to avoid strong relative
139   differences at the tails of the distribution. The data are divided into equally spaced binning intervals,
140   which is sufficient for population data that follows a given order of regression model (such as
141   quadratic). A pseudo count of one is added to each bin to avoid bins with zero objects, and the

142    relative frequency for each treatment is calculated by dividing the number of objects in each bin by

143    the population size (sum of all objects in all bins). Next, a bin-wise normalization step is performed by

144    dividing the relative frequency of each bin for each treatment by the median of the corresponding

145    bin of the control wells, such as the 'AllStars' control.

146    Binning creates equal-length bins to which data are assigned. The default number of bins (the

147    binning level) is 7.

148    For variable x, assume that the data set is $\{x_i\}$, let $x_1$, $x_2$,....$x_m$ represent the ordered values of the

149    variable. Let the x$^{th}$ percentile be $min(x)$ and $max(x)$. The range of the variable is $range(x) = max(x) -$

150    $min(x)$. For binning, the width of binning interval is $L = \frac{max(x) - min(x)}{n}$. The split points are $s_k =$

151    $min(x) + L * k$ , where $k = 1, 2,...$ , $numbin-1$ and $numbin$ is $n$. For each bin a pseudocount of 1 is

152    added $Countp(X_{ik}) = Count(X_{ik}) + 1$.

153    The output data, consisting of a normalized frequency distribution profile for each cell population

154    and the annotation data are stored in the file 'AllDataTable'. *E: Determination of Statistically*

155    *Significant Altered Profiles.* A regression analysis is performed using the Bioconductor R-package

156    maSigPro (Conesa, et al., 2006) to identify significantly changed normalized distribution profiles. *F:*

157    *Postprocessing (optional).* The gene, cluster and frequency of the gene within the cluster are listed

158    for all significantly changed normalized distribution profiles identified by the maSigPro analysis.

159    To generate the data for the cell cycle progression screen we seeded HeLa cells in 384-well plates

160    either transfected or treated with the inhibitors in three independent biological replicates (Figure

161    2A). We used 166 different siRNAs to target 107 genes, from which 54 had been reported to interfere

162    with cell cycle progression (Kittler, et al., 2007) (Supplementary Figure 1, Supplementary Table 1).

163    Additionally, we used cells treated with the chemical inhibitors aphidicolin or nocodazole, which lead

164    to G1/S and G2/M cell cycle arrest, respectively. Cells left untreated (Mock) or treated with siRNAs

165    against 'AllStars' or 'Luciferase' were used as negative controls. While images from AllStars- and

166    Luciferase-treated cell populations showed an unaltered, normal phenotype, treatment of cells with

167    aphidicolin (A1-4) or nocodazole (N1-4) resulted in an altered phenotype as a consequence of G1/S

168    or G2/M phase arrest, respectively (Figure 2B). On day 4, cells were fixed, nuclei stained with

169    Hoechst (Figure 2A), images acquired using automated microscopy and automated image analysis

170    (Olympus Scan^R) was performed for extracting the image features 'Area', 'Total Intensity DAPI' and

171    'Mean Intensity DAPI' for each nucleus (Supplementary Figure 2) in tab-delimited files using a Scan^R

172    export script. The cell population distribution profiles for the control as well as the chemically or

173    siRNA-treated samples behave differently for the extracted object features (Figure 2C). They show a

174    strong shift towards smaller nuclei for nocodazole, and towards larger nuclei for aphidicolin-treated

175    samples for the feature 'Area', while for the feature 'Mean Intensity DAPI' the influence of these two

176    chemical treatments on the mean intensity is the opposite. Interestingly, for the feature 'Total

177    Intensity DAPI' a strong shift towards higher values was observed for nocodazole-treated samples,

178    while aphidicolin treatment did not alter the profile compared to that of the 'Allstars', 'Luciferase' or

179    'Mock'-treated wells. Distribution profiles of the cell populations treated with different siRNAs

180    (samples) showed no clear tendency (Figure 2C).

181    We then calculated p-values and RSQ-values using maSigPro regression analysis, as described, to

182    identify significantly altered distribution profiles compared to the neutral controls. The maSigPro

183    package computes a regression fit for each frequency distribution profile , and uses a linear step-up

184    (BH) false discovery rate (FDR) procedure (Benjamini and Hochberg, 1995). Here, we used a level of

185    0.05 for FDR control. Once statistically significant distribution profiles have been found, a variable

186    selection procedure is applied to find significant variables for each profile. The final step is to

187    generate lists of statistically significant profiles. As expected, cell populations treated with the

188    'AllStars' or 'Luciferase' controls usually had high p-values and low RSQ- values. Only two (10%) and

189    four (20%) out of 20 cellular populations treated with the neutral controls 'Allstars' or 'Luciferase',

190    respectively, were identified to be significantly changed for at least one of the three cellular features

191    used (Supplementary Table 2; Figure 3A - Plate2, Well 207). When hits were only considered positive

192    if at least two of the image features were identified as significantly changed, none of the neutral

193    controls were identified as a hit. In contrast, cells treated with aphidicoline (A1-4) or nocodazole (N1-

194    4) showed significant changes, indicated by low p-values and high RSQ-values for all of the three

195    extracted cellular features (Figure 4A). All 28 profiles for each of the aphidicolin conditions A2 (4

196    µg/ml/24 h), A3 (2 µg/ml/12 h) and A4 (4 µg/ml/12 h), for each of the nocodazole conditions N1 (50

197    ng/ml/24 h), N2 (75 ng/ml/24 h), N3 (50 ng/ml/12 h) and N4 (75 ng/ml/12 h) and 27 out of 28

198    profiles for the aphidicolin condition A1 (2 µg/ml/24 h) were identified as significantly changed hits

199    (Supplementary Figure 3). Interestingly, aphidicolin-treated samples showed marked differences for

200    the cellular features 'Area' and 'Mean Intensity' and only slight changes for the cell feature 'Total

201    Intensity DAPI' (Figure 3A – A1: Plate1, Well 208, A2: Plate2, Well 353, A3: Plate1, Well 44, A4:

202    Plate2, Well 213) , while nocodazole-treated samples showed strong changes in all three cellular

203    features used (Figure 3A -N4: Plate1, Well 47, N1: Plate 2, Well 354). In total, using these thresholds

204    for the p-value and the RSQ-value, 359 normalized distribution profiles were identified as

205    significantly altered for each of the cellular features 'Area' and 'Mean Intensity DAPI' and 335

206    normalized distribution profiles for the cellular feature 'Total Intensity DAPI'. This resulted in a total

207    of 448 significantly changed cell populations; with 247 profiles significantly changed for all three, 111

208    profiles for two and 90 profiles for only one of the analyzed cellular features. Next, for the 448

209    profiles identified as significantly changed, a k-means clustering approach was performed (Figure 4B

210    and 4C). The normalized distribution profiles for the features 'Area', 'Mean Intensity DAPI' and 'Total

211    Intensity DAPI' were arranged in four, three and two profile clusters, respectively. Cluster numbers

212    were selected to give high cluster reproducibility. Finally, for all clustered profiles a heatmap was

213    defined, based on the k-means clustering result arranged as a vector (consisting of zeros and ones

214    such as 0001-010-10 for a profile resulting in cluster 4 for 'Area', cluster 2 for 'Mean Intensity DAPI'

215    and cluster 1 for 'Total Intensity DAPI'). The heatmap was sorted using a hierarchical clustering

216    (hclust) algorithm to identify cell populations with similar distribution profiles (Figure 3A). Finally, a

217    dendrogram cut-off value of 1.8 was used to generate three main groups in the matrix.

218    As a result, the aphidicolin-treated samples A1 and A2 grouped differently from the aphidicolin-

219    treated samples A3 and A4, (Figure 3A, sidebar), while the nocodazole-treated samples N1 and N2, as

220    well as N3 and N4, grouped together. Further, the significant distribution profiles of samples treated

221    with siRNA were more dispersed in the heatmap, depending on the individual feature distribution.

222    Thus, with this two-step method - first identifying statistically significant normalized profiles for each

223    analyzed image feature, then using a heatmap to generate profile groups – we were able to

224    differentiate cell populations showing a similar distribution among the cluster profiles. Taken

225    together, the SOPRA workflow was responsive enough to distinguish not only nocodazole-treated

226    from aphidicolin-treated samples, but also to differentiate between samples that were treated with

227    the same concentration but for different durations (A1/A2 vs A3/A4).

228    As laid out above, analyzed features for siRNAs that target the same gene and have the same 'on-

229    target' phenotype should end up in the same cluster and also in the same heatmap group. Therefore,

230    we further analyzed if individual siRNAs for the same gene were represented in the same or different

231    heatmap groups. The individual siRNAs of 38 and 42 genes appeared exclusively in profile groups 1 or

232    2, respectively, strengthening the 'on-target' specificity of these siRNAs. In contrast, individual siRNAs

233    of 21 genes were represented in both of these groups, indicating less stringent 'on-target' specificity

234    or other influences, such as experimental variation. For the SOPRA workflow, two different siRNAs

235    were used for each gene in duplicate, therefore hits were classified as medium or weak hits if the

236    two siRNAs did not show the same cluster profile and were not grouped in the same heatmap group.

237  To assess the reproducibility of plate replicates and SOPRA workflow the RSQ-values for different
238  groups of replicates were determined and the correlation matrix between these groups was
239  calculated. Firstly, we defined replicate group R1 containing replicate r2, r3 and r4 (i.e. without
240  replicate 1), replicate group R2 containing r1, r3, and r4 (without replicate 2) and so on. Running
241  SOPRA for cell feature 'Total Intensity DAPI' with pre-defined maSigPro parameters alfa=1, Q=1 and
242  RSQ=0 one gets the p-values and RSQ-values for each replicate group. The correlation matrix
243  between the RSQ-values for sample data from the different groups R1- R4 is calculated and
244  visualized (Figure 5).

245  Furthermore, we used Receiver Operating Characteristic (ROC) to assess the statistical performance
246  of SOPRA workflow in comparison to other approach, such as "Kolmogorov-Smirnov (KS) test" which
247  uses probability density and "t.test" for assessment of population differences. The RO curve for cell
248  feature 'Total Intensity DAPI' is depicted in Figure 6 and show that the SOPRA method lies between
249  the other two RO curves.

250  To benchmark the efficiency of this method in gene perturbation hit prediction, we tested whether
251  the results of the SOPRA workflow could be validated by either the original cell cycle data from Kittler
252  *et al.* (Kittler, et al., 2007) or FACS data generated by our group (Figure 3B). We selected 46 of the
253  genes (hits and non-hits) analyzed with SOPRA and performed FACS analysis for cell cycle profiles
254  with one siRNA per gene. A hit was scored as positive for a particular method if at least one other
255  method also leads to the same (positive or negative) result (Supplementary Table 3). Out of the 46
256  genes analyzed, 30 genes from the Kittler *et al.* study were validated with at least one of the other
257  methods (SOPRA or FACS), while for the SOPRA and FACS analyses 36 and 38 genes, respectively,
258  were validated by one of the other two methods. Taken together, SOPRA and FACS analysis scored
259  best in their ability to predict hits, compared to the data published by Kittler *et al.* (Kittler, et al.,
260  2007).

261  Thus, the SOPRA workflow offers a unique and fast analysis approach, based on measured single
262  features of cell populations, comparable to or better than published methods. In contrast to FACS
263  data analysis it does not need manual intervention or thresholding, such as cell gating. SOPRA is
264  therefore well suited for high-throughput and high-content data, as it can be easily run on multiple
265  features from an identical cell population.

266  ## Conclusions

267  Most methods published for analyzing high-content microscopic screens use population-averaged
268  values or manually performed cell classification steps for normalization and hit classification. The
269  SOPRA workflow represents a novel approach for analyzing large microscopy-based high-content
270  screens using non-averaged data of cell populations for normalization and hit determination. The
271  workflow generates frequency distribution profiles of cellular features normalized to a neutral
272  control for each treatment. These normalized distribution profiles are used for hit identification by
273  regression analysis to identify significantly altered profiles using the R-package maSigPro, as originally
274  described for the analysis of single series time course gene-expression data.

275  RNAi screens are frequently performed with multiple siRNAs per target gene; however the use of
276  population-averaged values often leads to the identification of 'off-target' effects as hits, since
277  population averaged values can only monitor major variations of the phenotype such as up- or down-
278  regulation compared to a control. In contrast, non- averaged data can indicate more diverse changes
279  of a cell population upon treatment; thus different siRNAs targeting the same gene should have a
280  similar 'on-target' effect on the distribution profile of the measured cellular features and
281  consequently these are more likely to be 'true' hits. The SOPRA workflow we describe here has the
282  power to cluster all significantly altered normalized distribution profiles, identifying siRNAs with
283  similar 'on-target' profiles for the same gene via a heatmap approach. Therefore, the SOPRA
284  workflow can be used to avoid false-positive hits or 'off-target' effects, leading to more reliable HCS
285  hit results, reducing time and work intensive validation steps.

286  In principle, the SOPRA workflow can be used to analyze single cell population data from various
287  sources such as microscopy or FACS. In this study, we performed a microscopy-based high content
288  screen of the effect of siRNA-mediated gene knockdown of selected genes taken from a published
289  cell cycle data set from Kittler et al (Kittler, et al., 2007), as an example to demonstrate the utility of
290  the SOPRA workflow.

291  We were able to show that the false positive detection rate (detection of neutral controls as
292  significantly changed) can be reduced considerably when taking into account more than one cellular
293  feature. As described using the generated cell cycle data, we were able to demonstrate that the
294  SOPRA workflow led to no false-positive hits among the neutral controls, when at least two of the
295  image features were taken into account. For the cell populations treated with the cell cycle
296  inhibitors, a very high hit detection rate of 99.55% was achieved (223 of 224 cell population profiles).
297  We also used siRNA knockdowns in this screen, which produce less significant phenotypic effects
298  compared to small chemical compounds. Nevertheless, analysis of changed cell populations based on

10

299  gene perturbation with siRNA using SOPRA still achieved a hit detection rate comparable to a manual

300  FACS analysis with commercial software, which requires predetermined gating or thresholding.

301  Taken together, SOPRA is a novel analysis workflow that uses a unique analysis approach for non-

302  averaged high-throughput data from cellular features, based on regression analysis of normalized

303  frequency distribution profiles of cell populations. It offers an easy to handle workflow and can be

304  run on hundreds of cell populations using multiple features. In particular, treated cell populations are

305  defined as significantly changed on two measurements - the p-value and the RSQ-value - followed by

306  a clustering step to identify treatments with the same normalized density profiles. A following

307  heatmap analysis enabled us to filter out most hits that are likely to be false positive. Thus, SOPRA is

308  a unique tool ideal for high content analysis of cell population data.

## Methods

309

310  <u>Cell Cycle perturbation screen</u>

311  We generated a set of screening plates consisting of siRNAs (Qiagen, Germany) targeting proteins

312  responsible and not responsible for cell cycle progression, as well as the neutral siRNAs 'AllStars' and

313  'Luciferase', and wells without treatment (Mock) (Supplementary Figure 1). On day one cells were

314  seeded in 96-well plates and transfected using Hiperfect (Qiagen, Germany). The chemical cell cycle

315  inhibitors nocodazole and aphidicolin were added as positive controls at the described time points

316  and concentrations. On day four cells were fixed using 4% PFA and stained with Hoechst 33342 (5

317  μg/ml, Sigma). The plates were imaged using an automated microscope (IX-81, Olympus, Germany)

318  and analyzed using the Scan^R software with an image analysis assay designed in-house

319  (Supplementary Figure 2).

320  Using a Scan^R single cell export script, single cell data was exported and are downloadable from

321  https://transfer.mpiib-berlin.mpg.de/s/AibR4AHLCR9xzDB?path=%2F.  The  SOPRA  project

322  description  (Supplementary  File  1)  is  also  available  from  GitHub

323  https://github.com/kppleissner/SOPRA/ .

324  <u>Cell Cycle FACS validation</u>

325  For FACS analysis of cell cycle profiles, $1 \times 10^5$ cells were seeded into each well of a 12-well plate 24 h

326  before transfection. Cells were then transfected with Hiperfect transfection reagent (Qiagen)

327  according to the manufacturer's guidelines. In brief, 150 ng of specific siRNA was added to RPMI

328  without serum and incubated with 6 μl Hiperfect in a total volume of 100 μl. After 10 to 15 min, the

329  liposome-siRNA mixture was added to the cells with 1 ml of cell culture medium (RPMI (Gibco)

330     supplemented with 10% fetal calf serum (FCS) (Biochrome), 2 mM glutamine, and 1 mM sodium

331     pyruvate), to give a final siRNA concentration of 10 nM. After 1 day, cells were trypsinized and

332     seeded into new 6-well plates. Three days after transfection, cells were detached from the plate with

333     the addition of trypsin-EDTA for 5 minutes, spun down for 10 minutes at 500 x $g$ and resuspended in

334     0.5 ml PBS. The resuspended cells were then added to 70% ethanol for fixation and left at −20°C

335     overnight. Cells were collected by centrifugation, resuspended, rinsed in PBS and re-collected by

336     centrifugation. Pelleted cells were resuspended in 500 µl PBS containing a final concentration of 20

337     µg/ml propidium iodide and 200 µg/ml RNAse A and left in the dark for 30 minutes at room

338     temperature. Cell Cycle analysis was then performed using a Becton Dickinson FACsort flow

339     cytometer and BD CellQuest Pro Software (BD Biosciences).

340     SOPRA

341     The SOPRA workflow (Figure 1) consists of several steps and requires a variety of input files. The

342     *'Single Cell Feature Files'* contain the features for every single cell measured, while the files

343     *'PlateConf_LookUp', 'PlateList'* and *'ScreenLog'* contain information about well content, plate content

344     and flagged wells. In the first step, the data is gathered, including flagging of wells and single objects

345     within wells. In the next step, a plate-wise median normalization is performed and the limits for the

346     binning intervals are defined. Subsequently, the single objects within each binning interval (bin) are

347     counted, and a bin-wise normalization is performed. Derived frequency distribution profiles of

348     measured features are then subjected to the regression analysis using R-Package maSigPro.

349     Significantly different profiles can be identified using the calculated p-value, RSQ-value and alpha-

350     value for each sample profile. The significant profiles can be clustered using different clustering

351     algorithms. Finally, a post processing step (optional) can be performed in order to convert siRNA into

352     gene names, cluster membership and frequency. The SOPRA workflow is written as a Shiny

353     application in R. A detailed project description with specific instructions for how to run the workflow

354     is available from GitHub.

355

## Figures and Files

357     **Figure 1 –SOPRA workflow of high-throughput data sets**

358     **(A)** High content screening data is generated and used to prepare single object data files and input

359     data files. **(B)** Screen description and the single cell data files are generated manually. **(C)** Wells that

360     should be omitted are flagged and **(D1)** the single object data is filtered, normalized to the median of

361   the controls and a common binning axis for all plates is determined. **(D2)** For each measured feature

362   the frequency distribution profile (histogram) is generated for each sample well, which is then

363   normalized for each bin to the median distribution profile of the controls. **(E)** Significantly changed

364   normalized distribution profiles are determined using regression analysis and **(F)** a post processing

365   step is performed to determine the number of screening hits.

366   **Figure 2 – Schematic representation of the microscopic cell cycle screening assay**

367   **(A)** Cells were seeded in 384-well plates and treated with siRNAs or chemical cell cycle inhibitors at

368   different concentrations and time points to inhibit cell cycle progression. Cells were fixed, stained

369   with Hoechst and subjected to automated microscopy and image analysis. **(B)** Treatment with the

370   control-siRNAs AllStars and luciferase did not lead to any changes of the cell population. Treatment

371   with aphidocoline A1 (2 µg/ml, 24 h), A2 (4 µg/ml, 24 h), A3 (2 µg/ml, 12 h), A4 (4 µg/ml, 12 h) and

372   nocodazole N1 (50 ng/ml, 24 h), N2 (75 ng/ml, 24 h), N3 (50 ng/ml, 12 h), N4 (75 ng/ml, 12 h)

373   resulted in cell populations arrested at various stages of the cell cycle. **(C)** Distribution profiles were

374   generated for each well from the data exported for the features 'Area', 'Mean Intensity DAPI' and

375   'Total Intensity DAPI' for all nuclei.

376   **Figure 3 – Heatmap analysis and examples of significantly altered distribution profiles**

377   **(A)** The normalized regression profiles for different treatment conditions for aphidicolin (A1-A4) and

378   nocodazole (N1 and N4), as well as Luciferase are displayed. A heatmap was generated showing the

379   distribution of all cell populations with at least one significantly changed profile for the features

380   'Area', 'Mean Intensity DAPI' and 'Total Intensity DAPI' among the SOPRA cluster profiles. Wells

381   treated with aphidicolin or nocodazole are displayed in different shades of green or blue in the row

382   sidebar. Wells Mock-treated or treated with siRNA against Luciferase or AllStars are indicated in red,

383   orange and yellow, respectively. Wells treated with siRNA against specific genes are displayed in grey

384   in the row sidebar. The heatmap is clustered using hierarchical clustering, and a dendogram, cut-off

385   of 1.8 performed resulting in the heatmap groups (1), (2) and (3). The Venn diagram displays the

386   distribution of the significantly changed profiles for each treatment among the heatmap groups (1)-

387   (3). **(B)** Examples of profiles for the features 'Area', 'Total Intensity DAPI' and 'Mean Intensity DAPI'

388   of cell populations significantly changed upon siRNA treatment, as well as the corresponding

389   microscopic and FACS images.

390   **Figure 4 – Results of SOPRA regression analysis and cluster profiles**

391   **(A)** Calculated RSQ and p-values of each well for the features 'Area', 'Mean Intensity DAPI' and

392   'Total Intensity DAPI' using the maSigPro package. **(B)** Data visualization by cluster analysis.

13

393        Normalized distribution profiles of all significantly altered normalized profiles for the three

394        image features were clustered using k-means with 4, 3 and 2 clusters, respectively. The

395        average feature profile is shown (black line) together with the individual profiles of the cell

396        populations in the cluster (grey lines) or **(C)** as the mean of 3 replicates.

397

398    **Figure 5 – Reproducibility assessment between replicates**

399    Correlation matrix between the RSQ-values for sample data from the different replicate groups R1-

400    R4 for cell feature 'Total Intensity DAPI'

401

402    **Figure 6 – Assessment of diagnostic quality by Reciever Operating Curve (ROC)**

403    Receiver Operating Characteristic (ROC) serves to assess the SOPRA workflow in comparison to other

404    statistical approaches, such as "t.test" and "Kolmogorov-Smirnov (KS) test". The RO curves for cell

405    feature 'Total Intensity DAPI' shows that the SOPRA method lies between the t.test and KS-test.

406

407    **Code availability and implementation**

408    Source code of SOPRA shiny application (ui.R , server.R) , single cell data (96-wells plate) data for

409    testing and SOPRA project description (folder: Manual) are freely available from GitHub

410    https://github.com/kppleissner/SOPRA/ .

411

412    **Supporting Data ZIP File for 384-wells plates (Single Cell Features)** :

413    Due to large size of files the 384-wells plate data couldn`t be uploaded to GitHub and therefore are

414    available as ***384_Plates_for_SOPRA.zip*** from MPI-IB Cloud tranfer server via this URL

415    https://transfer.mpiib-berlin.mpg.de/s/AibR4AHLCR9xzDB?path=%2F .

416    The ***384_wells_Plates_for_SOPRA.zip*** file contains data based on a cell cycle screen analyzed with

417    Scan^R (Olympus). Following cell features were measured: 'Area' ,' Mean Instensity DAPI' and 'Total

418    Intensity DAPI'. In general, any file of the correct format can be used for SOPRA. For each plate – or

419    part of a plate – one file is needed. The folders also contain the descriptive files '*PlateConf_LookUp*',

420    '*PlateList*' and '*ScreenLog*'.

421

## Author contributions

423    RKG- Conceived, designed, performed and analyzed the screen and wrote the manuscript

424    KPP- Wrote R-Scripts for SOPRA workflow, replication and ROC analysis, shiny interface, project
425    description and realized storage of SOPRA into GitHub

426    CC – Performed the FACS validation of hits

427    TFM – Supervised the project

428    APM- Conceived the project, conceived, designed and analyzed the screen, wrote the R-Scripts , user
429    interface UI.R and server.R in the Shiny environment, the SOPRA-Project Description and the
430    manuscript.

431

## Acknowledgements

438

# References

440　Benjamini, Y. and Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach
441　to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)* 1995:289-300.
442　Birmingham, A., *et al.* Statistical methods for analysis of high-throughput RNA interference screens.
443　*Nat Methods* 2009;6(8):569-575.
444　Boland, M.V. and Murphy, R.F. A neural network classifier capable of recognizing the patterns of all
445　major subcellular structures in fluorescence microscope images of HeLa cells. *Bioinformatics*
446　2001;17(12):1213-1223.
447　Buchser, W., *et al.* Assay Development Guidelines for Image-Based High Content Screening, High
448　Content Analysis and High Content Imaging. In: Sittampalam, G.S., *et al.*, editors, *Assay Guidance*
449　*Manual.* Bethesda (MD); 2004.
450　Conesa, A., *et al.* maSigPro: a method to identify significantly differential expression profiles in time-
451　course microarray experiments. *Bioinformatics* 2006;22(9):1096-1102.
452　Conrad, C., *et al.* Automatic identification of subcellular phenotypes on human cell arrays. *Genome*
453　*Res.* 2004;14(6):1130-1136.
454　Fuchs, F., *et al.* Clustering phenotype populations by genome-wide RNAi and multiparametric
455　imaging. *Mol. Syst. Biol.* 2010;6:370.
456　Gorenstein, J., *et al.* Reducing the multidimensionality of high-content screening into versatile
457　powerful descriptors. *Biotechniques* 2010;49(3):663-665.
458　Gough, A.H. and Johnston, P.A. Requirements, features, and performance of high content screening
459　platforms. *Methods Mol. Biol.* 2007;356:41-61.
460　Kittler, R., *et al.* Genome-scale RNAi profiling of cell division in human tissue culture cells. *Nat. Cell*
461　*Biol.* 2007;9(12):1401-1412.
462　Knapp, B., *et al.* Normalizing for individual cell population context in the analysis of high-content
463　cellular screens. *BMC Bioinformatics* 2011;12:485.
464　Lang, P., *et al.* Cellular imaging in drug discovery. *Nat Rev Drug Discov* 2006;5(4):343-356.
465　Loo, L.H., Wu, L.F. and Altschuler, S.J. Image-based multivariate profiling of drug responses from
466　single cells. *Nat Methods* 2007;4(5):445-453.
467　Neumann, B., *et al.* High-throughput RNAi screening by time-lapse imaging of live human cells. *Nat*
468　*Methods* 2006;3(5):385-390.
469　Ohya, Y., *et al.* High-dimensional and large-scale phenotyping of yeast mutants. *Proc. Natl. Acad. Sci.*
470　*U. S. A.* 2005;102(52):19015-19020.
471　Perlman, Z.E., *et al.* Multidimensional drug profiling by automated microscopy. *Science*
472　2004;306(5699):1194-1198.
473　Sacher, R., Stergiou, L. and Pelkmans, L. Lessons from genetics: interpreting complex phenotypes in
474　RNAi screens. *Curr. Opin. Cell Biol.* 2008;20(4):483-489.
475　Singh, S., Carpenter, A.E. and Genovesio, A. Increasing the Content of High-Content Screening: An
476　Overview. *J Biomol Screen* 2014;19(5):640-650.
477　Suratanee, A., *et al.* Detecting host factors involved in virus infection by observing the clustering of
478　infected cells in siRNA screening images. *Bioinformatics* 2010;26(18):i653-i658.
479　Tanaka, M., *et al.* An unbiased cell morphology-based screen for new, biologically active small
480　molecules. *PLoS Biol.* 2005;3(5):e128.
481　Taylor, D.L. Past, present, and future of high content screening and the field of cellomics. *Methods*
482　*Mol. Biol.* 2007;356:3-18.
483　Zhou, X. and Wong, S.T. Informatics challenges of high-throughput microscopy. *Signal Processing*
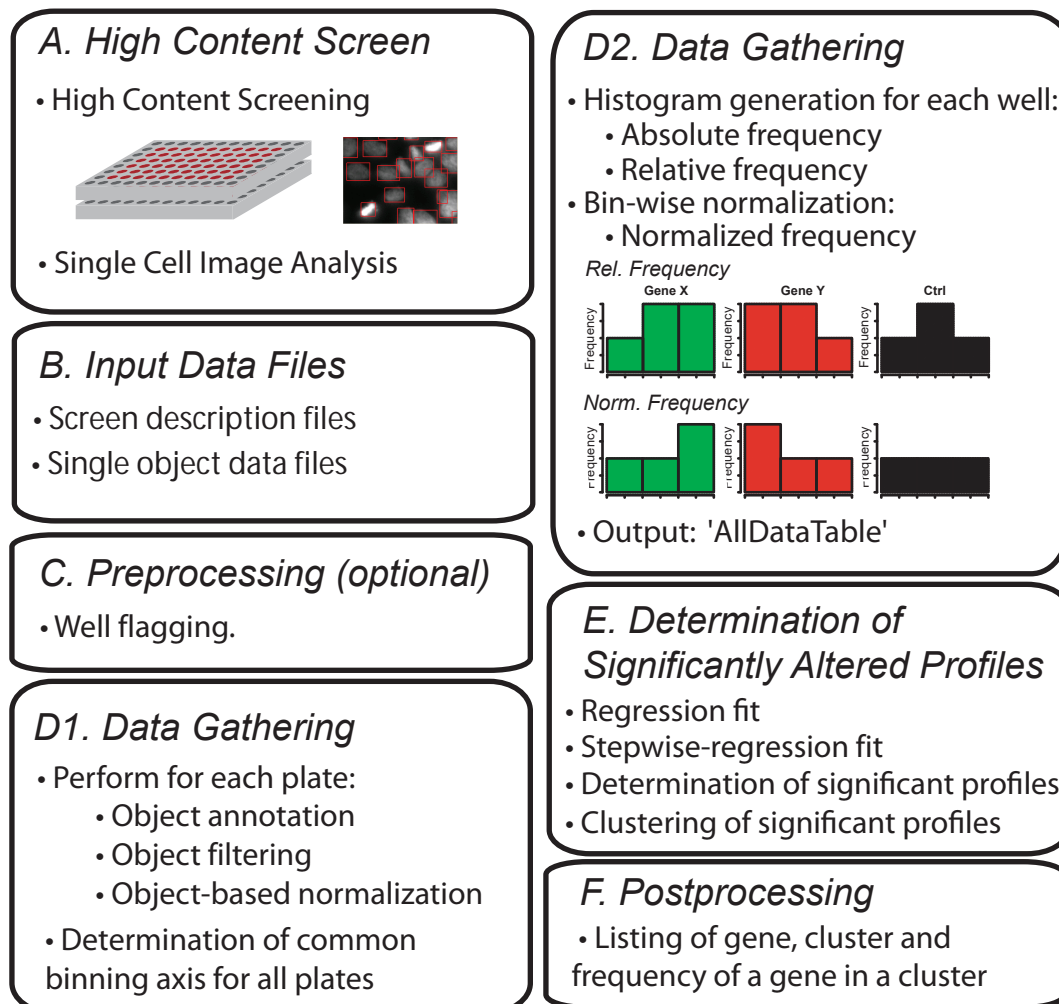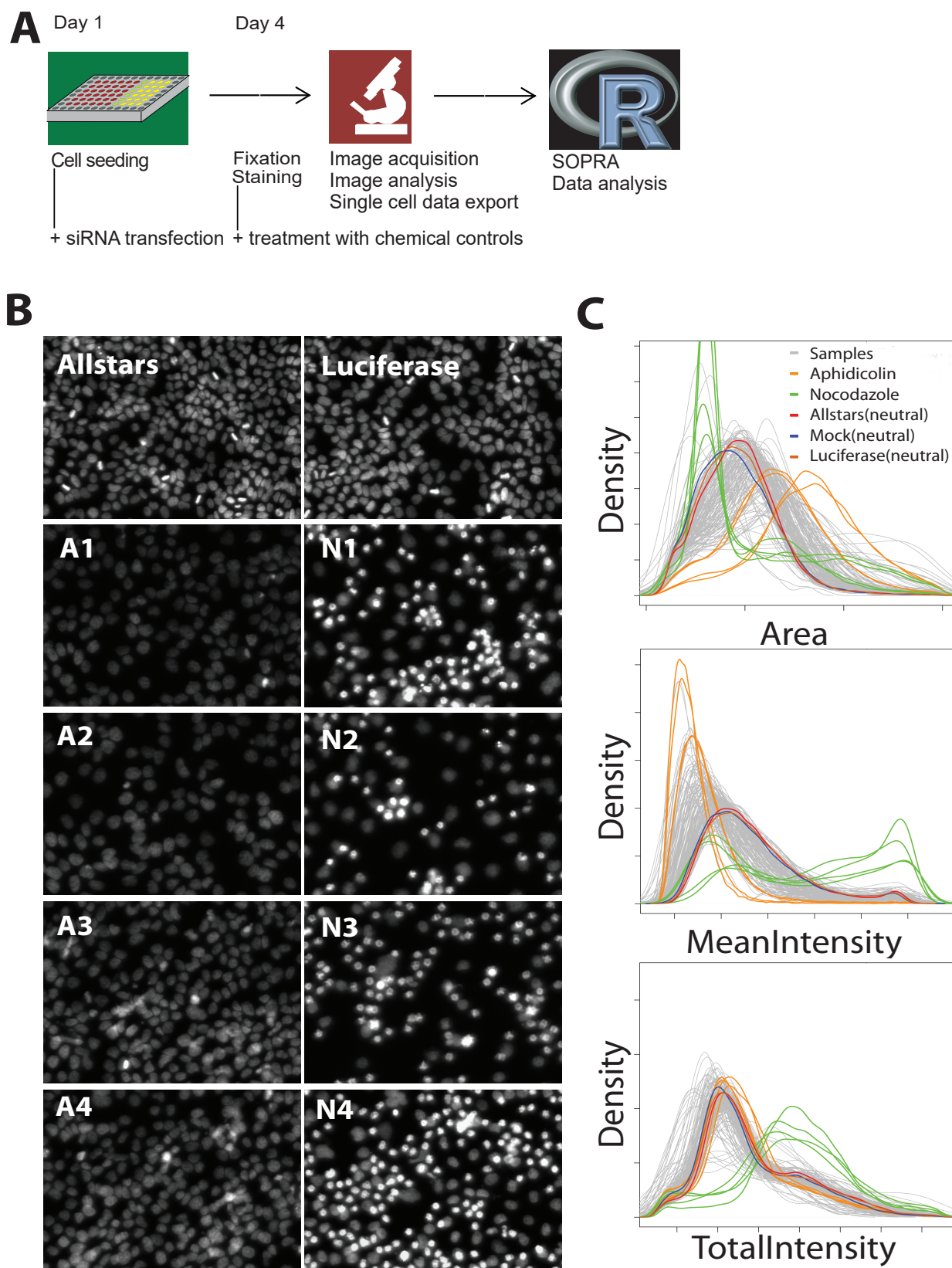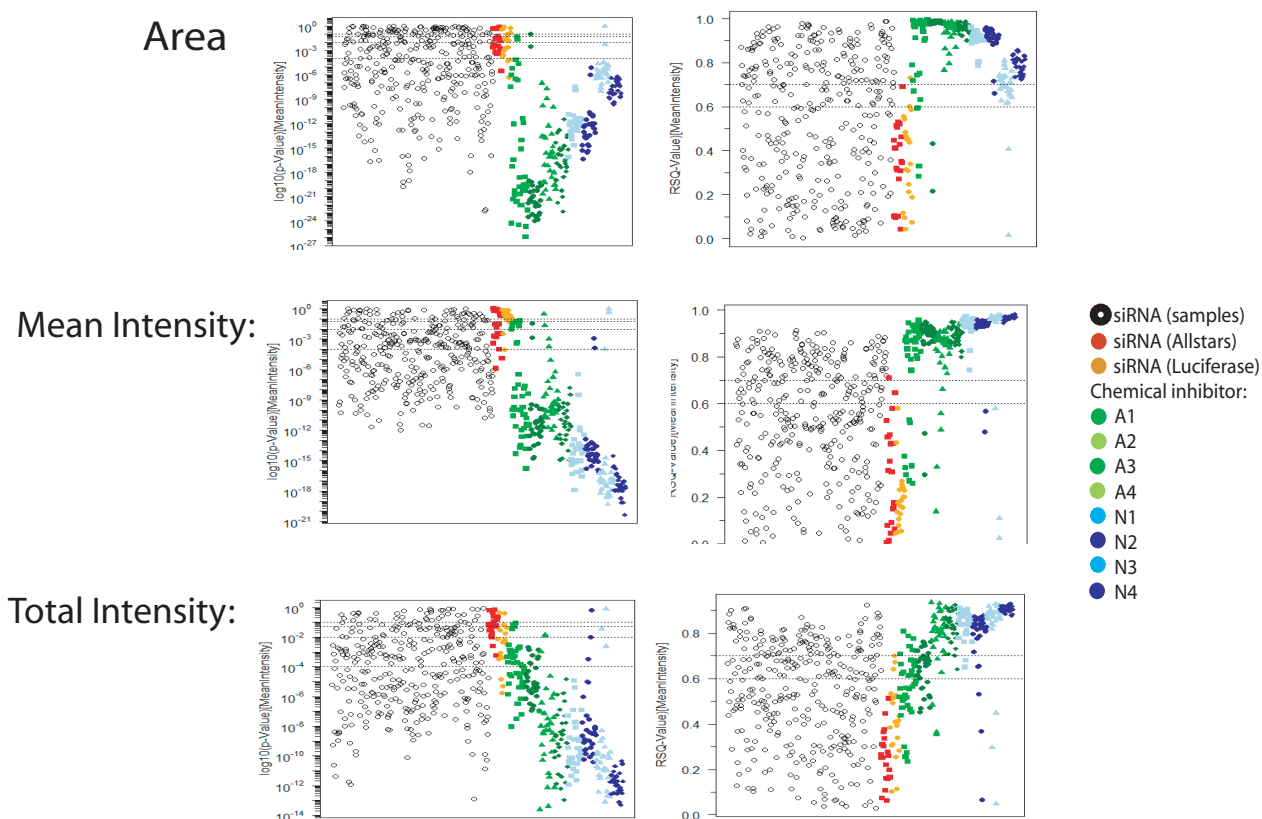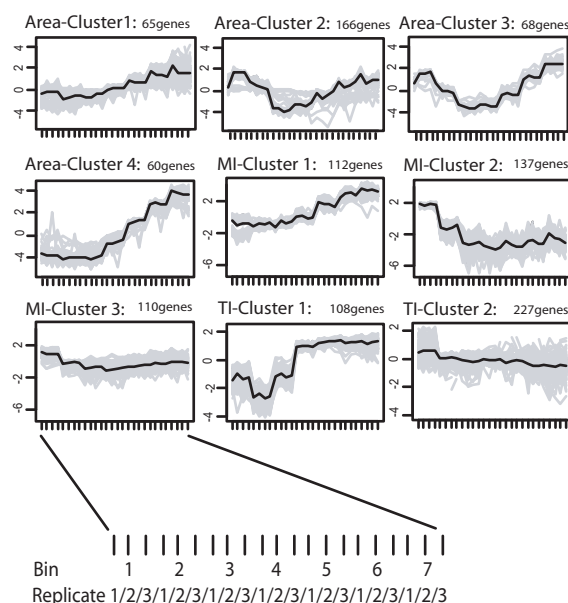484　*Magazine, IEEE* 2006;23(3):63-72.
485

# Figure 1



**A. High Content Screen**
- High Content Screening
- Single Cell Image Analysis

**B. Input Data Files**
- Screen description files
- Single object data files

**C. Preprocessing (optional)**
- Well flagging.

**D1. Data Gathering**
- Perform for each plate:
  - Object annotation
  - Object filtering
  - Object-based normalization
- Determination of common binning axis for all plates

**D2. Data Gathering**
- Histogram generation for each well:
  - Absolute frequency
  - Relative frequency
- Bin-wise normalization:
  - Normalized frequency

*Rel. Frequency*

*Norm. Frequency*

- Output: 'AllDataTable'

**E. Determination of Significantly Altered Profiles**
- Regression fit
- Stepwise-regression fit
- Determination of significant profiles
- Clustering of significant profiles

**F. Postprocessing**
- Listing of gene, cluster and frequency of a gene in a cluster
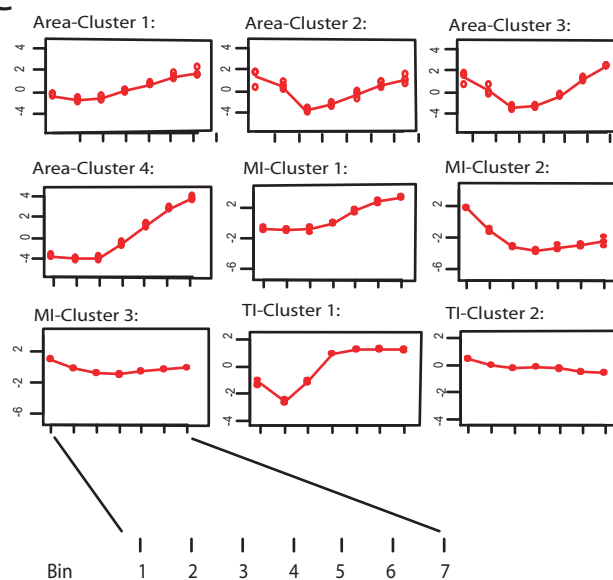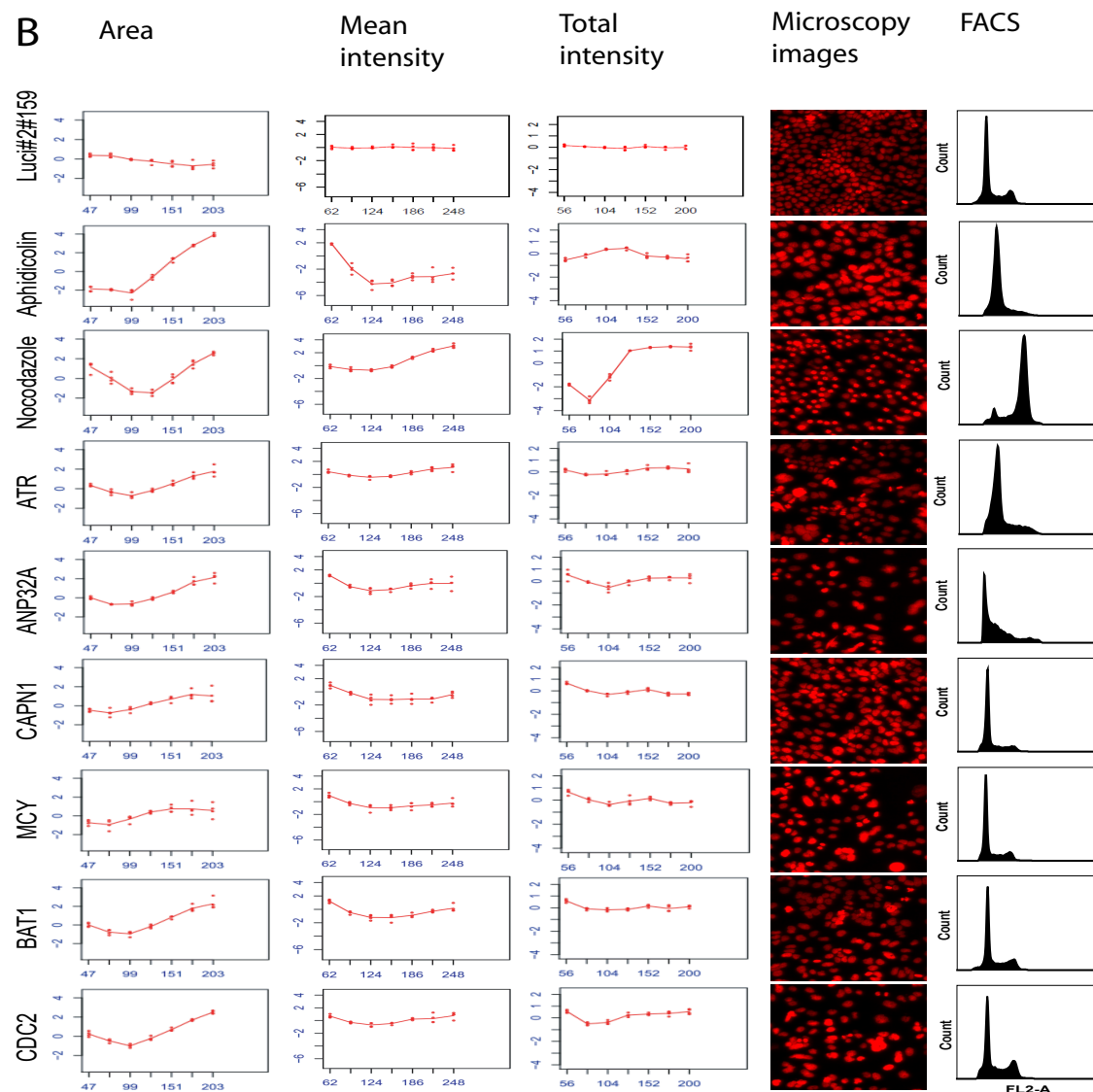
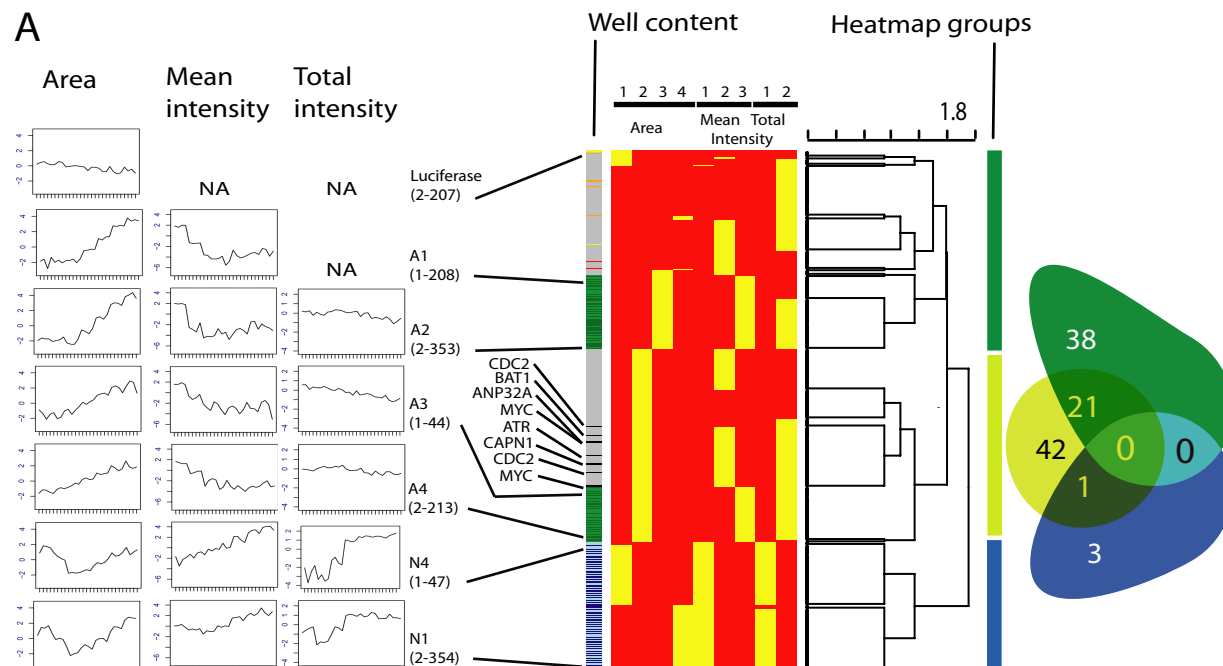# Figure 2

# Figure 3

A

# Figure 4

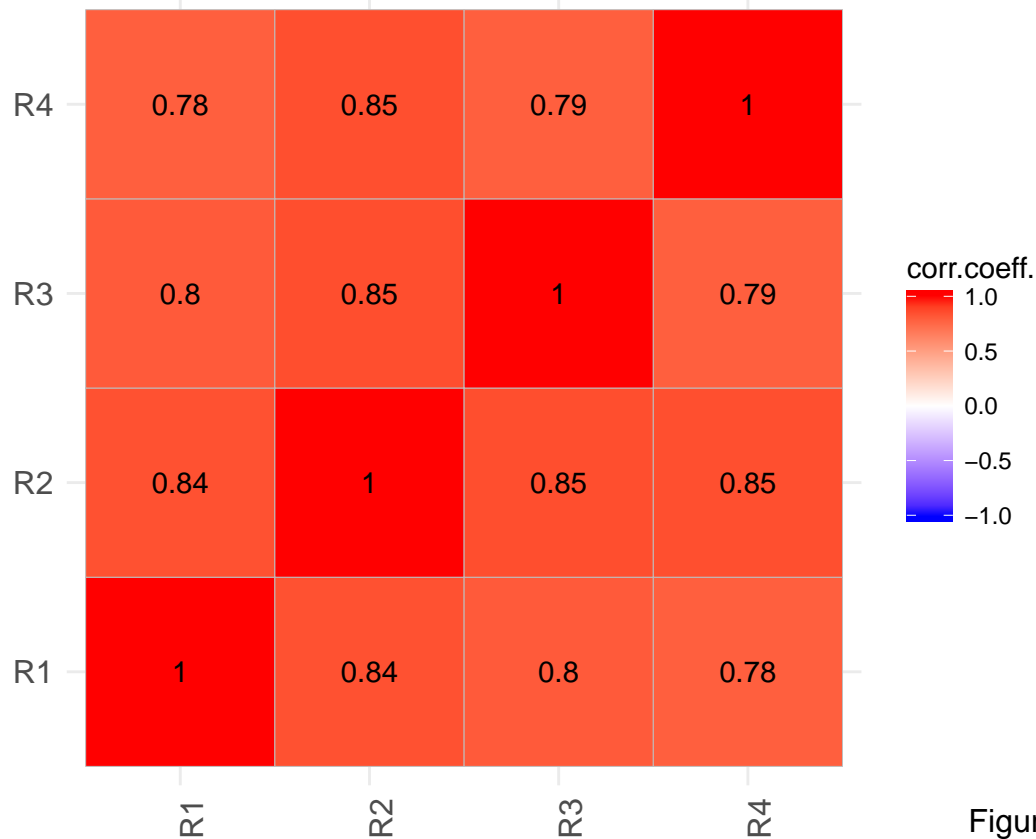Correlation matrix based on RSQ−values of samples for feature TotalIntensity
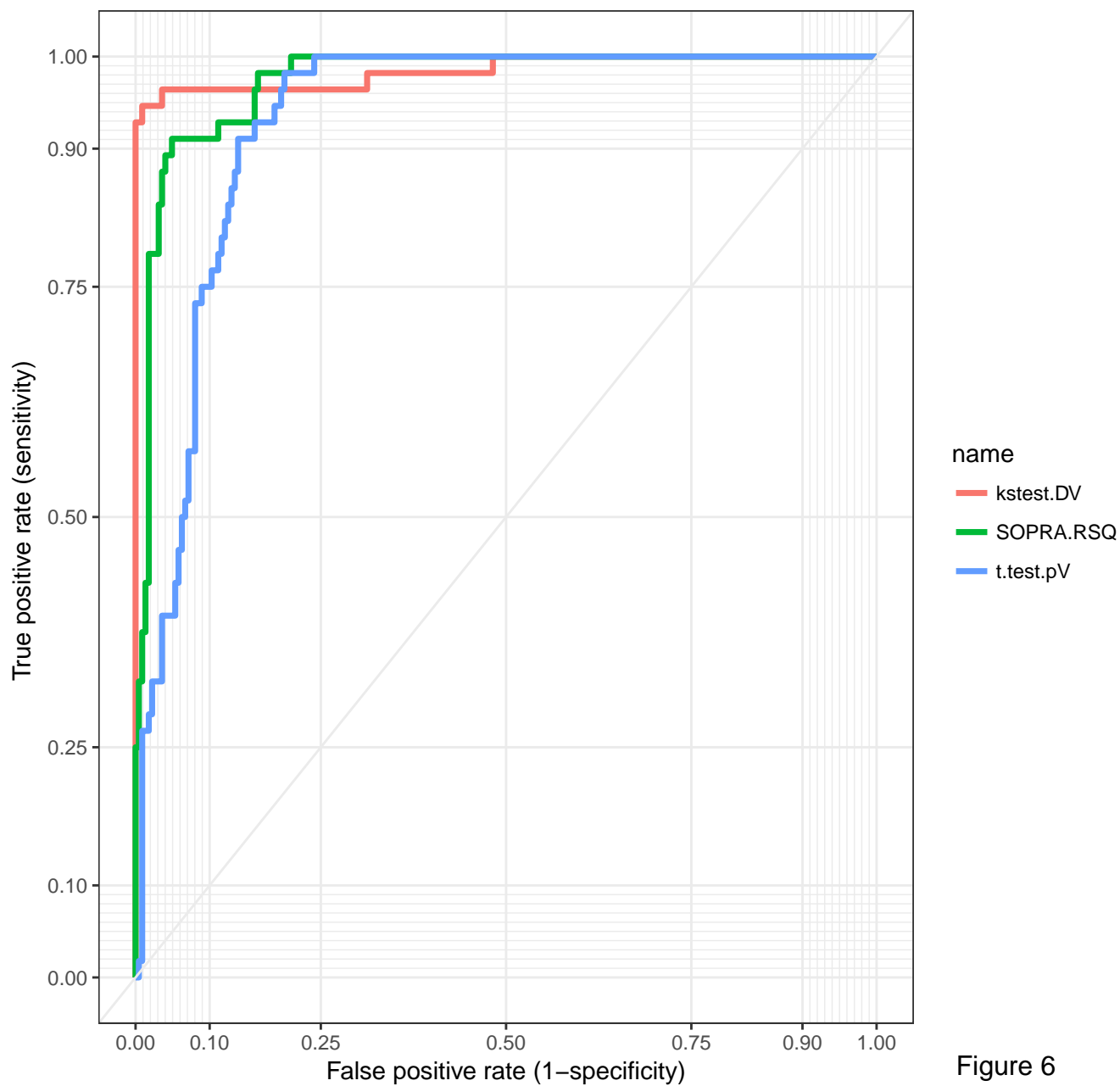
Figure 5

ROC for TotalIntensity: Positive−vs−Neutral

Figure 6