

1 The importance of study design for detecting differentially 2 abundant features in high-throughput experiments

3 Luo Huaien^{1*}, Li Juntao^{1*}, Chia Kuan Hui Burton¹, Paul Robson², Niranjan Nagarajan^{1#}

4 ¹ *Computational and Systems Biology, Genome Institute of Singapore, Singapore 138672, Singapore*

5 ² *Stem-cell and Developmental Biology, Genome Institute of Singapore, Singapore 138672, Singapore*

6 *Contributed equally

7 #Corresponding Author: nagarajann@gis.a-star.edu.sg

8 ABSTRACT

9 The use of high-throughput experiments, such as RNA-seq, to simultaneously identify
10 differentially abundant entities across conditions has become widespread, but the systematic
11 planning of such studies is currently hampered by the lack of general-purpose tools to do so.
12 Here we demonstrate that there is substantial variability in performance across statistical
13 tests, normalization techniques and study conditions, potentially leading to significant
14 wastage of resources and/or missing information in the absence of careful study design. We
15 present a broadly applicable experimental design tool called EDDA, and the first for single-
16 cell RNA-seq, Nanostring and Metagenomic studies, that can be used to i) rationally choose
17 from a panel of statistical tests, ii) measure expected performance for a study and iii) plan
18 experiments to minimize mis-utilization of valuable resources. Using case studies from recent
19 single-cell RNA-seq, Nanostring and Metagenomics studies, we highlight its general utility
20 and, in particular, show a) the ability to correctly model single-cell RNA-seq data and do
21 comparisons with 1/5th the amount of sequencing currently used and b) that the selection of
22 suitable statistical tests strongly impacts the ability to detect biomarkers in Metagenomic
23 studies. Furthermore, we demonstrate that a novel mode-based normalization employed in
24 EDDA uniformly improves in robustness over existing approaches (10-20%) and increases
25 precision to detect differential abundance by up to 140%.

26 INTRODUCTION

27 The availability of high-throughput approaches to do counting experiments (e.g. by using
28 DNA sequencing) has enabled scientists in diverse fields (especially in Biology) to
29 simultaneously study a large set of entities (e.g. genes or species) and quantify their relative

30 abundance. These estimates are then compared across replicates and experimental conditions
31 to identify entities whose abundance is significantly altered. One of the most common
32 scenarios for such experiments is in the study of gene expression levels, where sequencing
33 (with protocols such as SAGE¹, PET² and RNA-Seq³) and probe-based approaches⁴ can be
34 used to obtain a digital estimate of transcript abundance in order to identify genes whose
35 expression is altered across biological conditions (e.g. cancer versus normal⁵). Other popular
36 settings where such differential abundance analysis is performed include the study of DNA-
37 binding proteins and histone modifications (e.g. using ChIP-Seq^{6,7}), RNA-binding proteins
38 (e.g. using RIP-Seq⁸ and CLIP-Seq⁹) and the profiling of microbial communities (using 16S
39 rRNA amplicon¹⁰ and shotgun sequencing¹¹).

40 Due to its generality, a range of software tools have been developed to do differential
41 abundance tests (DATs), often with specific applications in mind, including popular
42 programs such as edgeR¹², DEseq⁶, Cuffdiff^{13,14}, Metastats¹¹, baySeq¹⁵ and NOISeq¹⁶. The
43 digital nature of associated data has allowed for several model-based approaches including
44 the use of exact tests (e.g. Fisher's Exact Test¹¹), Poisson¹⁷ and Negative-Binomial^{6,12} models
45 as well as Bayesian¹⁵ and Non-parametric¹⁶ methods. Recent comparative evaluations of
46 DATs in a few different application settings (e.g. for RNA-Seq^{6,16,18-21} and Metagenomics¹¹)
47 have further suggested that there is notable variability in their performance, though a
48 consensus on the right DATs to be used remains elusive. In addition, it is not clear, which (if
49 any) of the DATs are broadly applicable across experimental settings despite the generality of
50 the statistical models employed. The interaction between modelling assumptions of a DAT
51 and the application setting, as defined by both **Experimental Choices** (e.g. number of
52 sequencing reads to produce for RNA-seq) as well as intrinsic **Experimental**
53 **Characteristics** (e.g. number of genes in the organism of interest), could be complex and not
54 predictable *a priori*. Correspondingly, only in very recent work, have experimental design
55 issues been discussed in a limited setting, i.e. using a *t*-test for RNA-seq analysis²². Also, as
56 experimental conditions can vary significantly and along several dimensions (**Table 1**), a
57 systematic assessment of DATs under all conditions is likely to be infeasible. As a result, the
58 choice of DAT as well as decisions related to experimental design (e.g. number of replicates
59 and amount of sequencing) are still guided by rules of thumb and likely to be far from
60 optimal.

61 In this study, we establish the strong and pervasive impact of experimental design decisions
62 on differential abundance analysis, with implications for study design in diverse disciplines.

63 In particular, we identified data normalization as a source of performance variability and
64 designed a robust alternative (**mode-normalization**) that uniformly improves over existing
65 approaches. We then propose a new paradigm for rational study design based on the ability to
66 model counting experiments in a wide spectrum of applications (**Figure 1**). The resulting
67 general-purpose tool called **EDDA** (for “Experimental Design in Differential Abundance
68 analysis”), is the first program to enable researchers to design experiments for single-cell
69 RNA-seq, NanoString assays and Metagenomic sequencing and we highlight its use through
70 case studies. EDDA provides researchers access to an array of popular DATs through an
71 intuitive online interface (<http://edda.gis.a-star.edu.sg>) and answers questions such as “How
72 much sequencing should I be doing?”, “Does the study adequately capture biological
73 variability?” and “Which test should I use to sensitively detect differential abundance in my
74 application setting?”. To provide full access to its functionality, EDDA is also available as a
75 user-friendly R package (on SourceForge: <https://sourceforge.net/projects/eddanorm/> and
76 Bioconductor: <http://www.bioconductor.org/packages/devel/bioc/html/EDDA.html>), and is easily
77 extendable to new DATs and simulation models.

78 **RESULTS**

79 In the following, we present and emphasize the somewhat under-appreciated impact of
80 various experimental conditions (grouped into two categories - Experimental Choices and
81 Experimental Characteristics, see **Table 1**) and various popular DATs (see **Table 2**) on the
82 ability to detect differential abundance. Our results highlight the importance of careful
83 experimental design and the interplay between experimental conditions and DATs in
84 dictating the success of such experiments. We first establish that the impact of experimental
85 choices on performance can be significant and in the next section explore their interaction
86 with various experimental characteristics. The results presented here are based on synthetic
87 datasets to allow controlled experiments and exploration of a wide-range of parameters, with
88 no emphasis on a particular application. In the following section, we discuss the validity of
89 the modelling assumptions and parameter ranges that we investigated and use this to motivate
90 the design of EDDA. We conclude by showcasing EDDA’s application in various settings.
91 For ease of reference, a schematic overview of the simulation model in EDDA (**Figure 1a**)
92 and a flowchart of how it can be used (**Figure 1b**) is provided in **Figure 1** with detailed
93 descriptions in the **Methods** section.

94 **Impact of Experimental Choices on Performance**

95 While the availability of high-throughput technologies (such as massively parallel DNA
96 sequencing) to do counting experiments has significantly increased the resolution of such
97 experiments, the cost of the experiment is often still an important factor and the number of
98 replicates and data-points that can be afforded may be less than optimal. Furthermore, in
99 many settings the number of replicates or the number of data-points possible may be
100 constrained due to technological limitations or uniqueness of samples and conditions. In such
101 settings, it would be ideal to understand and exploit the trade-off between the number of
102 replicates and data-points needed (e.g. by doing deeper RNA-seq for a few biological
103 replicates) for such analysis. In the following, we investigate these dimensions individually
104 and in conjunction.

105 *Number of Replicates* (NR): The performance of DATs (measured here by the area-under-
106 curve of the receiver operating characteristic or ROC curve; AUC) as a function of the
107 number of replicates in the study is highly non-linear, with significant improvements
108 obtainable until a saturation point (indicated by a larger marker size; **Figure 2a**). The
109 saturation point is likely to be largely determined by the intrinsic variability of replicates in
110 an experiment. However, clear differences are also seen across DATs as seen in **Figure 2a**,
111 where edgeR and DESeq achieved AUC greater than 0.95 with five replicates, while baySeq
112 required eight replicates and has substantially lower AUC with one replicate. While all DATs
113 seem to converge toward optimal performance (AUC=1) in this setting, the rate of
114 convergence varies markedly (e.g. note the curve for Cuffdiff). Note that the relative
115 performance of DATs is influenced by the experimental setting, especially in conditions
116 where the number of replicates available is small (as is typically the case), and thus the
117 choice of DAT is strongly dependent on the desired precision/recall trade-off
118 (**Supplementary Figure 1**).

119 *Number of Data-points* (ND): The number of data-points generated in an application setting
120 is often set to be the maximum possible given the resources available. However, this may
121 lead to misallocation of resources as suggested by **Figure 2b, 2c** and **2d**. In this setting,
122 increasing the number of data-points continues to improve AUC over a wide range of values
123 (for most DATs). At very high values AUC saturates, but not necessarily at 1 (**Figure 2b**,
124 **1c**). Oddly, for one of the DATs (baySeq), performance decreases with increase in ND – a
125 feature that is not *a priori* evident from its specification (**Figure 2b**). However, with more

126 replicates, much fewer data-points are needed to obtain high AUC values, suggesting that this
127 is a better trade-off in this setting (**Figure 2c, 2d**; not necessarily the case in other settings
128 e.g. when the number of replicates is already high or intrinsic variability across replicates is
129 low). Note that if the number of data-points is limited by the application, then there is
130 considerable variability in performance across DATs (**Supplementary Figure 2**) and some
131 DATs may have consistently lower AUCs (Cuffdiff and Metastats, in this setting) with high
132 ND as well. Thus, to meet experimental objectives, especially when high precision is desired,
133 informed decisions on statistical test and number of replicates to employ are needed (as
134 facilitated by experimental design tools such as EDDA).

135 **Interaction of Experimental Characteristics and Choices**

136 It is important to note that experimental choices alone do not dictate the ability to detect
137 signals of differential abundance and as we show here, the intrinsic characteristics of the
138 experiment are also an important variable that need to be taken into account. This excludes
139 the possibility of pre-computing recommendations for experimental choices and DATs to use
140 in various applications and emphasizes the need for a general-purpose experimental design
141 tool such as EDDA.

142 *Entity Count* (EC): Intuitively, the impact of the number of entities being profiled is expected
143 to be minimal and the *a priori* assumption is that scaling the number of data-points as a
144 function of the number of entities should lead to comparable performance. Our results
145 suggest that this is not quite true. As seen in **Figure 3a**, when the number of entities is low,
146 there is not only greater variability in performance, but average AUC is also lower. Some
147 statistical tests seem to be less appropriate when the number of entities is low (e.g. Cuffdiff),
148 while others exhibited greater robustness (edgeR and DESeq, followed by baySeq).

149 *Sample Variability* (SV): The intrinsic variability seen across replicates in an application
150 setting dictates the trade-off between the number of data-points and replicates needed in
151 complex ways as shown in **Figure 3b**. While the specific patterns will depend on the
152 application, even for a setting with large effect sizes as shown here, the specific tradeoffs
153 chosen by the various DATs vary (more sequencing for DESeq vs more replicates for
154 Cuffdiff) and the cost-effectiveness of a method ($= NR \times ND$ needed) can switch with
155 sample variability (e.g. Cuffdiff goes from being the least to the most cost-effective when
156 sample variability increases; **Figure 3b**).

157 *Abundance profile* (AP): The relative abundance of entities is often seen to follow a power
158 law distribution (**Supplementary Figure 3**), but the precise shape can vary and together with
159 the number of data-points generated, impact overall performance for an application. In
160 particular, testing differential abundance for rare entities (with low relative abundance) can
161 be difficult and could explain the variability in performance seen in **Figure 3c**. While all
162 methods have lower AUC for a rare-entity-enriched profile from Wu et al²³ (**Supplementary**
163 **Figure 3**), some methods seem to be more robust (e.g. baySeq) or tuned to detect rare entities
164 (e.g. Metastats), while others experience a larger relative drop in performance (e.g. Cuffdiff
165 or NOISeq), suggesting that DAT choices need to take abundance profiles into account.

166 *Perturbation profile*: The affect of specific profiles of differential abundance on prediction
167 performance is likely to be the least predictable from first-principles and this was also seen in
168 our experiments (**Figure 4**). Altering the fraction of differentially abundant entities alone
169 could reorder the performance of various statistical tests, as seen in **Figures 4a** and **4b**, where
170 baySeq went from being the worst performer to the best performer. Furthermore, switching
171 the distribution of fold-changes was also seen to affect results as seen in **Figures 4b** and **4c**,
172 with NOISeq now becoming the best performing DAT. Other parameters such as the
173 abundance profile also combine with the perturbation profile to influence relative
174 performance as seen in **Figures 4b** and **4d**, where DEseq went from being one of the best to
175 being the worst performer. Overall, no single DAT was found to outperform others
176 (**Supplementary Figure 4**), highlighting that specific experimental characteristics and
177 choices need to be taken into account while choosing an appropriate DAT.

178 **Modelling Assumptions and Normalization**

179 In the absence of variability across replicates and experimental biases, counting experiments
180 of the sort studied here can naturally be modelled as samples from a Multinomial distribution.
181 To simulate technical and intrinsic variability, a common approach has been to model the
182 relative abundance of each entity across replicates using the Negative Binomial
183 distribution^{15,24}. In fact, in many studies this is the model from which counts are simulated for
184 each entity^{15,25}, independent of those for other entities (we refer to this as the **Negative**
185 **Binomial** model), and bypassing the joint simulation of counts from a Multinomial
186 distribution (referred to here as the **Full** model). In practice, both the Full and the Negative
187 Binomial model elicit similar performance for most experimental settings and most DATs.
188 However, for a few DATs (baySeq, NOISeq and Cuffdiff) we observed deterioration in

189 performance on the Full model when compared to a similar experiment using the Negative
190 Binomial model (**Supplementary Figure 5**) suggesting that the Full model is a better
191 measure of performance of a DAT.

192 By analysing several published and in-house datasets we established that, in general, for bulk
193 transcriptome sequencing (confirming earlier reports^{15,24}), Nanostring assays and Shotgun
194 Metagenomic sequencing (not shown in prior work), variability in replicates can be
195 adequately modelled using the negative binomial distribution (**Supplementary Figure 6**). An
196 exception to this rule was, however, seen in single-cell RNA-seq experiments in accordance
197 with observations of unusually high cell-to-cell variability in recent reports^{26,27}
198 (**Supplementary Figure 6c, d**). For cases where an appropriate model for variability across
199 replicates is not available (as in the single-cell case), we developed a **Model-Free** approach,
200 that uses sub-sampling (and appropriate scaling where needed) of existing datasets to provide
201 simulated datasets that match sample-to-sample variability in real datasets better (with the
202 drawback that it relies on the availability of a dataset with many replicates; see **Methods**).

203 To evaluate the data generation models used in this study (either model-based or model-free),
204 as well as establish their suitability for the design of EDDA, we first investigated
205 distributional properties of real and simulated datasets (**Supplementary Figure 7**). The
206 results here indicate that while overall both simulation approaches (where applicable) provide
207 good approximations and capture the general trend, the Model-Free approach more closely
208 mimics true sample variability (**Supplementary Figure 7**). We next tested the suitability of
209 an approach where simulated datasets are generated to mimic an existing pilot dataset and
210 employed to measure trends in performance. Our results confirmed that simulated data
211 generated by our simulation models enable reliable measurement of true performance for
212 DATs (relative-error in AUC < 6%) and monitoring of trends as a function of experimental
213 choices and characteristics (**Supplementary Figure 8**). In addition, experimental
214 recommendations from EDDA simulations were also found to match DAT recommendations
215 based on benchmarking on real datasets²⁸ (**Supplementary Figure 9**), suggesting that EDDA
216 can help avoid this step and still reliably guide experimental design.

217 In some experimental settings, variability in replicates can be extremely low and directly
218 simulating from the Multinomial distribution (a special case of the Full model that we refer to
219 as the **Multinomial** model; see **Methods**) is sufficient. In principal, with enough data-points,
220 statistical testing under the Multinomial model should be straightforward and we expect

221 various DATs to perform well. The few exceptions that we noted, suggest that aspects other
222 than statistical testing, such as data normalization, may play a role in their reduced
223 performance (**Supplementary Figure 5**).

224 An investigation of different normalization approaches (**Table 2**) under the various
225 experimental conditions explored in this study suggests that their robustness can vary
226 significantly as a function of the experimental setting. In particular, we observed a few
227 settings under which many of the existing approaches performed sub-optimally (**Figure 5a**)
228 and to address this we designed a new method (**mode-normalization**) that analyzes the
229 distribution of un-normalized fold-changes of entities using mode statistics to select a suitable
230 normalization factor (see **Methods** and **Supplementary Figure 10**). We compared mode-
231 normalization to the default normalization and a popular alternate (Upper-quartile
232 Normalization), for each DAT and across all the conditions tested here, to find that the use of
233 mode-normalization uniformly improved performance (on average, AUC by 9% and
234 precision by 14% at 5% FDR). Also, in cases where the performance of a few DATs dipped
235 under the Multinomial model, mode-normalization was able to rescue the AUC values
236 (**Supplementary Figure 5d**). In addition, we identified several examples where mode-
237 normalization significantly improved AUC values for all the DATs tested (improving
238 precision to detect differential abundance by up to 140% at 5% FDR), highlighting that
239 proper data normalization is a key step in attaining experimental goals (**Figure 5a**). As
240 depicted in **Figure 5a**, there are often cases where the default normalization of a DAT or a
241 popular alternate (Upper-quartile Normalization) lead to reduced performance while calling
242 differentially abundant entities, while mode normalization consistently achieves optimal
243 performance across DATs. Note that, no normalization can be expected to work under all
244 conditions and simulated datasets generated by EDDA can also be valuable to compare and
245 choose among alternative normalization techniques.

246 To further evaluate normalization methods on real datasets, we studied the consistency of
247 differential abundance predictions (against predictions on the full dataset) upon down-
248 sampling of data-points, using 3 deeply-sequenced RNA-seq datasets²⁸ (**Figure 5b**). These
249 results highlight the robustness of mode-normalization versus other popular approaches
250 (UQN and TMM; **Table 2**). Mode-normalization was found to improve robustness by 10-
251 20% across datasets and was the least affected by imbalances in sequencing depth across
252 conditions (**Figure 5b**).

253 Applications of EDDA and mode-normalization

254 The observed variability in the performance of DATs across experimental characteristics and
255 choices, and the demonstration that data from many kinds of high-throughput experiments
256 can be adequately modelled *in silico*, motivated the use of a *simulate-and-test* paradigm in
257 EDDA to guide experimental design (see **Figure 1a** and **Methods**). EDDA allows users fine-
258 scale control of all the variables discussed here (summarized in **Table 1**), but also provides
259 the option to directly learn experimental parameters (and models for the model-free
260 approach) from pilot or publicly-available datasets. Some of the commonly expected modes
261 of usage for EDDA are discussed in the methods section **EDDA Modules** and illustrated in
262 **Figure 1b**. Furthermore, to showcase the use of EDDA and mode-normalization, we present
263 results from EDDA analysis of several recently generated datasets in three different
264 experimental settings, each highlighting a different aspect of the utility of the package in a
265 practical scenario.

266 For the first case study, we analyzed data from a recent single-cell RNA-seq study of
267 circulating tumor cells (CTCs) from melanoma patients²⁹. The authors generated on average
268 1000 data-points per entity (> 20 million reads) and used a one-way ANOVA test (equivalent
269 to a *t*-test) to identify differentially abundant genes between CTCs and primary melanocytes.
270 We reanalyzed the data using EDDA to simulate synthetic datasets that mimicked real data
271 (with the **Model-Free** approach and a 96-cell dataset generated as a resource for this study)
272 and used them to test a panel of DATs (see **Methods**). The availability of new micro-fluidics
273 based systems to automate single-cell omics has highlighted the cost of sequencing as the
274 major bottleneck in studying a large number of cells. Strikingly, EDDA analysis revealed that
275 this study could have been conducted with 1/5th of the sequencing that was done (by reducing
276 sequencing depth to 1/10th and doubling the number of replicates) without affecting
277 performance in terms of identifying differentially abundant genes (**Figure 6a**). This was,
278 however, only possible if the appropriate DAT was used (edgeR and BaySeq, in this case),
279 with the choice of DAT playing a more significant role than the amount of sequencing done.
280 Using BaySeq with on average 100 reads per gene (i.e. 2 million reads per cell as opposed to
281 the 20 million reads used in this study) and increasing the number of replicates from 5 to 50
282 (and thus maintaining sequencing cost) would be expected to boost AUC from 0.86 (and 0.75
283 using the *t*-test) to 0.96 and sensitivity from 57% to 72% at 5% FDR, in this study (**Figure**
284 **6a**). Note that while practical considerations could limit the number of CTCs that can be
285 captured and studied, this should not be an issue for other cell types in this study.

286 In the second case study, we analyzed data from an in-house project (manuscript in
287 preparation) for the development of prognostic and predictive gene signatures in breast
288 cancer on the NanoString nCounter platform (NanoString, WA, USA). The NanoString
289 platform allows for the digital measurement of gene expression, similar to RNA-seq, but is
290 typically used to profile a small, selected set of genes in a large number of samples (107
291 genes and 306 samples in this study), making data normalization a critical step for robust
292 analysis. Using EDDA, we explored the impact of a range of normalization approaches,
293 including the one recommended in the NanoString data analysis guide (NanoString, WA,
294 USA). As shown in **Figure 6b**, the coefficient of variation (COV) of a panel of 6
295 housekeeping genes (*ACTG1*, *ACTB*, *EIF4G2*, *GAPDH*, *RPLP0* and *UBE2D4*) is
296 significantly lower, as expected, when the data is properly normalized and, in particular, this
297 is the case when using mode-normalization which produces the lowest average COV across
298 all the methods tested. We then investigated the effect of normalization on the power to
299 discriminate between ER positive and negative breast cancers using a panel of 8 known ER
300 signature genes³⁰. Not surprisingly, the ability to distinguish ER positive and negative breast
301 cancers improves significantly with proper normalization with mode-normalization providing
302 the largest F-score (see **Methods**) among all the approaches tested (**Figure 6c**).

303 In our third case study, we critically assessed the analysis done in a recent Metagenomic
304 study that looked into the association of markers in gut microflora with type 2 diabetes in
305 Chinese patients³¹. Due to the complexity of the microbial community in the gut, the authors
306 reported that they were able to assemble more than 4 million microbial genes overall.
307 Correspondingly, since on average ~20 million paired-end reads were generated per sample,
308 the sequencing done here is expected to provide shallow coverage of the gene set, on average.
309 The study involved a large number of cases (71) and controls (74) and the Wilcoxon test was
310 used to identify differentially abundant genes in a case-control comparison. Another batch of
311 100 cases and controls was then used to validate biomarkers identified from the first batch.
312 We used EDDA to generate virtual microbiome profiles and assessed the performance of the
313 Wilcoxon test in this setting, in addition to the default panel of DATs (**Table 2**). EDDA
314 analysis revealed that the Wilcoxon test was likely to have been too conservative in this
315 setting and could have been improved upon using DATs like Metastats, which was designed
316 for Metagenomic data and edgeR, which is commonly used for RNA-seq analysis (**Figure**
317 **6d**). In addition, while increased coverage is likely to improve the ability to detect true
318 differences in the microbiome, the gains are expected to be relatively modest (for edgeR,

319 ~1% increase in AUC with 10-fold increase in sequencing; **Figure 6d**). Correspondingly,
320 despite the shallow coverage employed in this study, it is likely to have captured a significant
321 fraction of the biomarkers that could have been determined with more sequencing. In
322 contrast, increasing the number of replicates is likely to have markedly improved the ability
323 to detect true differences in the microbiome, (with edgeR, ~7% increase in AUC by doubling
324 the number of replicates **Figure 6d**). Keeping sequencing cost fixed and using 300 replicates
325 and 5 reads per gene is thus expected to boost AUC from 0.73 (using the Wilcoxon test) in
326 the study to 0.97 (using edgeR; **Figure 6d**) and sensitivity from 32% to 86% at 5% FDR.

327 Based on this, we reanalyzed cases and controls from the first batch in this study to identify
328 an additional 37,664 differentially abundant genes (17% increase) using edgeR, of which a
329 greater fraction (27% increase over the original study) were also validated in the second batch
330 of samples (**Supplementary Table 2**). The newly identified genes highlighted previously
331 missing aspects of the role of the microbiome in Type 2 diabetes including the identification
332 of 24 additional gene families enriched in differentially abundant genes (**Supplementary**
333 **Table 3**). In particular, this analysis detected two bacterial genes identified as multiple sugar
334 transport system substrate-binding proteins as being abundant in cases vs controls
335 (**Supplementary Figure 11**), as well as the enrichment of two new families for multiple
336 sugar transport system permease proteins (K02025 and K02026). Strikingly, the newly
337 detected genes also enabled construction of an improved microbiome-based predictive model
338 for Type 2 diabetes (AUC of 0.96 vs 0.81 in the original study; see **Methods** and
339 **Supplementary Figure 11**), based on the selection of 50 marker genes (see **Methods** and
340 **Supplementary Table 4**), highlighting that improved differential abundance analysis based
341 on informed choices using EDDA can significantly impact major conclusions from a study.

342 **DISCUSSION**

343 The case studies highlighted in the previous section are not unique in any way and point to a
344 general trend in current design of high-throughput experiments, where commonly used rules
345 of thumb lead to suboptimal designs and poor utilization of precious experimental resources.
346 Considering that the market for sequencing based experiments alone is currently in the
347 billions of dollars, savings in research budgets worldwide would be substantial with even a
348 modest 10% improvement in study design. On the other end, with a fixed budget, optimizing
349 study design can ensure that key insights are not missed. In particular, in many scenarios
350 where either (a) effect sizes are small and fold-changes are marginal or (b) large effects on a

351 few entities mask subtle effects on other entities or (c) the goal is to understand coordinated
352 processes such as cellular pathways through enrichment analysis³², loss in sensitivity or
353 precision due to unguided experimental choices can be detrimental to the study. The use of a
354 personalized-benchmarking tool such as EDDA provides a measure of insurance against this.

355 With the recent, dramatic expansion in the number of high-throughput applications (largely
356 based on DNA sequencing) as well as end-users (often non-statisticians), differential
357 abundance testing is now frequently done by non-experts in settings different from the
358 original benchmarks for a method. This can make it difficult to determine if a particular
359 analysis was appropriate or lead to incorrect results. One possible approach that could
360 account for this is to use multiple DATs to get a consensus set (also available as an option on
361 the EDDA web-server) but this can result in overly conservative predictions. For example, in
362 a recent analysis of RNA-seq data from two temporally-separated mouse brain samples using
363 edgeR and DESeq (with default parameters), we found that the intersection of differentially
364 expressed genes (at 10% FDR) contained less than 10% of the union. Breakdown of the
365 results showed that while edgeR was primarily reporting up-regulated genes (998 out of
366 1189), DESeq was largely reporting down-regulated genes (875 out of 878), with no
367 indication as to which analysis was more appropriate. EDDA simulations and analysis were
368 then used to clarify that results from edgeR were more reliable here (FPR of 3.8% vs 9.2% at
369 10% FDR) and could be improved further using mode-normalization (FPR of 1.5%).
370 Furthermore, the bias towards detecting up or down-regulated genes was intrinsic to the tests
371 here (not affected by normalization as we originally suspected) and hence reporting the union
372 of results was more appropriate. Examples such as this are not uncommon in the analysis of
373 high-throughput datasets and experimental design tools such as EDDA can help provide
374 informed answers to researchers.

375 We hope the results in this study serve to further highlight the still under-appreciated
376 importance of proper normalization for differential abundance analysis with high-throughput
377 datasets^{20,33,34}. Normalization based on mode-statistics provides an intuitive alternative to
378 existing approaches, exhibiting greater robustness to experimental conditions in general,
379 >20% improved AUC performance in some conditions, as well as the ability to detect cases
380 where proper normalization may not be feasible.

381 EDDA was designed to provide an easy-to-use and general-purpose platform for
382 experimental design in the context of differential abundance analysis. To our knowledge, it is

383 the first method that allows users to plan single-cell RNA-seq, Nanostring assays and
384 Metagenomic sequencing experiments, where the larger number of samples involved could
385 lead to important experimental tradeoffs. The combination of model-based and model-free
386 simulations in EDDA allows for greater flexibility and, in particular, we provide evidence
387 that the commonly used Negative Binomial model may not be appropriate for single-cell
388 RNA-seq, but a model-free approach (leveraging on a 96 cell dataset generated in this study)
389 is better suited. Model-free simulations using EDDA can thus serve as a basis for refining
390 new statistical tests and clustering techniques for single-cell RNA-seq. Note that a common
391 assumption in EDDA and most statistical testing packages is that deviations from the
392 multinomial model due to experimental biases can be corrected for and hence these issues
393 were ignored in this study^{14,35}.

394 The basis for EDDA is a simple *simulate-and-test* paradigm as the diversity of statistical tests
395 precludes more sophisticated approaches (e.g. deriving closed-form or numerical bounds on
396 expected performance). Given the simplicity of this approach, it is even more surprising that
397 the field has until now relied on rules of thumb. In light of this, the main contribution of this
398 work should be seen as the demonstration that significant variability can be observed across
399 all experimental dimensions and, therefore, lack of experimental design tailored to a
400 particular application setting can lead to substantial wastage of resources and/or loss of
401 detection power. We hope that the availability of EDDA through an intuitive, easy-to-use,
402 point-and-click web-based interface will thus encourage a wide-spectrum of researchers to
403 employ experimental design in their studies.

404 **METHODS**

405 **Single-cell Library Preparation and Sequencing**

406 ATCC® CCL-243™ cells (*i.e.* K562 cells) were thawed and maintained following vendor's
407 instructions using IMDM medium (ATCC® 30-2005™) supplemented with 10% FBS
408 (GIBCO® 16000-077™). The cells were fed every 2 to 3 days by dilution and maintained
409 between 2×10^5 and 1×10^6 cells/ml in 10 to 15 ml cultures kept in T25 flasks placed
410 horizontally in an incubator set at 37 °C and 5% CO₂. Cells were slowly frozen two days
411 after feeding at a concentration of 4 million cells per ml in 100µl aliquots of complete
412 medium supplemented with 5% DMSO (ATCC® 4X). The cryo-vials containing the frozen
413 aliquots were kept in the vapor phase of liquid nitrogen until ready to use. On the day of the
414 C1™ experiment, a 900 µl aliquot of frozen complete medium was thawed and brought to
415 room temperature. The cryo-vial was retrieved from the cryo-storage unit and placed in direct
416 contact with dry ice until the last minute. As soon as the cryo-vial was taken out of dry ice,
417 the cells were thawed as quickly as possible at a temperature close to 37°C (in about 30
418 seconds).

419 The room temperature complete medium was slowly added to the thawed cells directly in the
420 cryo-tube and mixed by pipetting four to five times with a 1000 µl pipette tip. This cell
421 suspension was mixed with C₁™ cell suspension reagent at the recommended ratio of 3:2
422 immediately before loading 5 µl of this final mix on the C₁™ IFC. The C1 Single-Cell Auto
423 Prep System (Fluidigm) was used to capture individual cells and to perform cell lysis and
424 cDNA amplification following the chip manufacturer's protocol for single-cell mRNA-seq
425 (PN 100-5950). Briefly, chemistry provided by the SMARTer Ultra Low RNA Kit (Clontech)
426 was used for reverse transcription and subsequent amplification of polyadenylated transcripts
427 using the C1™ script 1772x/1773x. After harvest of the amplified cDNA from the chip, 96-
428 way bar-coded Illumina sequencing libraries were prepared by tagmentation with the Nextera
429 XT kit (Illumina) following the manufacturer's protocol with modifications stated in
430 Fluidigm's PN 100-5950 protocol. The 96 pooled libraries were 51-base single-end
431 sequenced over 3 lanes of a Hi-Seq 2000.

432 Raw reads for all libraries are available for download from NCBI using the following link:
433 <http://www.ncbi.nlm.nih.gov/bioproject/238846>.

434 **Simulation of Count Data in EDDA**

435 *Abundance Profile (AP)*: When provided with sample data, EDDA uses the entity count and
436 the sample abundance profile from the data (when multiple samples are provided, the counts
437 are aggregated to get an average frequency profile) to do simulations. Users can also
438 explicitly provide a profile or choose from among pre-defined profiles including **BP** (for
439 BaySeq Profile; the profile used in simulations by Hardcastle et al.¹⁵), **HBR** (a profile derived
440 from a “human brain reference” dataset³⁶) and the profile from **Wu et al**²³. In order to
441 simulate with entity counts that differ from the original profile, EDDA allows users to sample
442 entities (without replacement for subsampling and with replacement for over-sampling) from
443 the middle 80% (entities are ordered by relative abundance), top 10% and bottom 10%
444 independently. This procedure is designed to maintain the dynamic range of the original
445 profile. In addition, to avoid working with entities with very low counts, EDDA allows users
446 to filter out those with counts below a minimum threshold for all replicates (default of 10).

447 *Perturbation Profile*: If EDDA is provided with sample data under two conditions then the
448 profile of differential abundance seen there is assigned to genes by keeping the relationship of
449 mean expression and fold change. Specifically, EDDA applies a DAT (DESeq and FDR-
450 cutoff of 5% by default, after mode normalization) to the sample data to identify
451 differentially abundant entities and their corresponding fold-changes ($f_i = x_i^2/x_i^1$ where x_i^1
452 and x_i^2 are the mean relative abundance of entity i under the first and second condition
453 respectively). Given d differentially abundant entities, with fold-changes f_1 to f_d , a set of d
454 entities from the first condition are perturbed by these fold-changes to obtain the abundance
455 profile for the second condition (i.e. $x_i^2 = \sqrt{f_j} \times x_i^1$ and $x_i^1 = x_i^1/\sqrt{f_j}$) while retaining the
456 correspondence between mean expression level and fold change. In addition, to account for
457 undetected entities an additional fraction of entities is randomly selected (from those that fail
458 the FDR cutoff and with fold-change > 1.5) and their observed fold-changes used to perturb
459 abundance profiles as before. The fraction of entities was determined to ensure that the
460 overall count matched the expected number of differentially abundant entities in the dataset
461 (estimated from the expected number of true positives at each gene's FDR level). In the
462 absence of sample data, EDDA also allows users to specifically set the percentage of entities
463 with increased and decreased entity counts (**PDA**), ranges for the fold-changes (**FC**) and the
464 distribution to sample fold-changes from (e.g. Log-normal, Normal or Uniform).

465 *Simulation Model (SM)*: The default model for simulations in EDDA is the **Full** model where
466 the mean abundance for each entity (under each condition) and the dispersion value provided
467 (**SV**) is used to compute means for the replicates using a Negative Binomial distribution (this
468 is done emulating the procedure in baySeq¹⁵ where each entity has a dispersion sampled from
469 a gamma distribution). When sample data is available, EDDA estimates dispersion values by
470 using the procedure in DESeq (alternately, edgeR) to fit the empirical distribution. Entity
471 abundances for each replicate are then normalized to get a frequency vector (that sums to 1)
472 to simulate count data from a Multinomial distribution (where the total count is sampled from
473 *Uniform*[$0.9 \times ND$, $1.1 \times ND$] and *ND* is the number of data-points specified by the user).
474 EDDA also allows users a simpler **Negative Binomial** model (**NB**) where the counts are
475 directly obtained from the Negative Binomial sampling described above and a **Multinomial**
476 model where the abundance profile (normalized to 1) is used to directly simulate from the
477 Multinomial distribution.

478 *Model-Free Approach*: To create **Model-Free** simulations, EDDA requires sample data as
479 input (ideally with many replicates and data-points). For RNA-seq, single-cell RNA-seq and
480 Metagenomic simulations, EDDA is packaged with sample datasets discussed in this study to
481 be used as input. If enough replicates are available in the sample dataset (i.e. greater than NR
482 \times desired number of simulations), EDDA subsamples counts from the entity in the sample
483 dataset with the closest average count and scales counts as needed. To simulate more
484 replicates than the number available in input data, EDDA groups entities according to their
485 average count to sub-sample entity counts. This approach was validated using RNA-seq
486 data²³ where more than 90% of genes had similar expression variability as the 10 closest
487 genes (in terms of average count; Kolmogorov–Smirnov test p -value < 0.05) as opposed to
488 2% of genes in the case of random groupings. After simulating variability in counts across
489 replicates using the model-free approach, EDDA also provides users the option to convert
490 counts back to a relative abundance profile for multinomial sampling of counts with a desired
491 number of data-points.

492 **Mode Normalization**

493 In principle, the ideal normalization factor for detecting differential abundance would be
494 based on counts for an entity that is not differentially abundant (or the sum of counts for all
495 such entities, see **Figure 5**). The idea behind **mode-normalization** is to identify such entities
496 under the assumption that non-differentially-abundant entities will tend to have similar un-

497 normalized fold changes (**UFCs**, computed as ratio of average counts across conditions). In
498 methods such as DESeq, a related idea is implemented using a quantity called the size-factor
499 (= ratio of observed count to a pseudo-reference sample computed by taking the geometric
500 mean across samples) and by taking the median (or upper-quartile) size factor under the
501 assumption that it would typically come from a non-differentially-abundant entity.

502 Mode-normalization in EDDA is based on calculating UFCs for all entities and determining
503 the approximate modes for their empirical distribution (**Supplementary Figure 10**).
504 Specifically, we used a kernel density estimation approach³⁷ to smooth the empirical
505 distribution and to compute local maxima for it. In cases where the number of maxima is not
506 as expected (i.e. 3, corresponding to entities with decreased, unchanged and increased relative
507 abundance), the bandwidth for smoothing was decreased as needed (starting from 0.5, in
508 steps of 0.02, till the number of maxima is as close to 3 as possible). If the final smoothed
509 distribution was uni- or tri-modal then the mode in the middle (presumably composed of non-
510 differentially-abundant entities) was chosen and the normalization factor was calculated from
511 the geometric mean of 10 entity counts around the mode. For bimodal distributions, selecting
512 the correct mode is potentially error-prone and we flag this to the user, picking the mode with
513 the narrowest peak (as given by the width of the peak at half the maximum value) and
514 calculating the normalization factor as before.

515 **Parameter/DAT settings and EDDA extensions**

516 EDDA is designed to be a general-purpose experimental design tool (that is easily extendable
517 due to its implementation in R) and correspondingly it provides significant flexibility in user
518 settings. In addition, we investigated the question of which parameter values are typically
519 seen in common applications (e.g. RNA-seq, Nanostring analysis and Metagenomics) and
520 used these to guide the evaluations presented in this study as detailed below.

521 For RNA-seq experiments, ECs in the range 1000 (microbial genomes) to 30000 (mammalian
522 genomes) are common with NR and ND in the range [1,10] and [10,10000] per entity,
523 respectively. For Nanostring and Metagenomic experiments (species profile), EC can be
524 significantly lower (in the range [10, 1000]) as well as significantly higher (>1 million for
525 Metagenomic gene profile). For the RNA-seq datasets analyzed in this study, SV was found
526 to be in the range [0.1, 0.9] (**Supplementary Table 1**) though much higher variability was
527 seen in Metagenomic data (> 5). Abundance profiles are best learnt from pilot data but in
528 their absence, the sample profiles provided with EDDA should serve as a useful range of

529 proxies. For RNA-seq experiments, PDA values in the range [5%, 30%] can be expected
530 while Nanostring and Metagenomic experiments can have even higher percentages. Fold
531 change distributions were typically observed to be well-approximated by a Log-normal
532 model (**Log-normal** $[\mu, \sigma]$) but other models are also feasible in EDDA (e.g. **Normal** $[\mu, \sigma]$
533 and **Uniform** $[a, b]$).

534 For DATs such as DESeq, edgeR, baySeq, NOISeq and Metastats that are implemented in R,
535 EDDA is set to call corresponding R functions, running them with default parameters and
536 normalization options unless otherwise specified. The results in this study were obtained with
537 the following versions of the various packages: DESeq (v1.7.6), edgeR (v2.4.6), baySeq
538 (v1.8.3), NOISeq (version as of 20th April 2011), Metastats (version as of 14 April 2009) and
539 R (v2.14.0). For Cuffdiff, the relevant C++ code was extracted from Cufflinks (v2.1.1) and
540 incorporated into EDDA as a pre-compiled dynamically-linked library using Rcpp³³
541 (<http://cran.r-project.org/web/packages/Rcpp/index.html>).

542 In its current form, EDDA installs the DATs listed in Table 2 by default. In addition, EDDA
543 is designed to support the easy integration of new DATs and a step-by-step guide to do so
544 (with the Wilcoxon test used here as an example) is provided as part of the package (see
545 **Supplementary Text**). EDDA is also designed to be extendable in terms of simulation
546 models and a guide for this is also provided in the installation package (see **Supplementary**
547 **Text**).

548 **EDDA modules**

549 For expert users the full functionality of EDDA and mode-normalization are available in a
550 package written in the statistical computing language R that can be freely downloaded from
551 public websites such as SourceForge and Bioconductor. In addition, to enable easy access for
552 those who are unfamiliar with the R environment, we designed web-based modules that
553 encapsulate typical use cases for EDDA (<http://edda.gis.a-star.edu.sg>; also see **Figure 1a**)
554 including modules for:

555 a) **Differential Abundance Testing**: This module is meant to enable users to easily run
556 a panel of DATs on any given dataset, to assess the variability of results across DATs,
557 compute the intersection and union of these results and correspondingly select a more
558 robust or comprehensive set of calls for downstream analysis. The assumption here is

559 that a user has already generated all their data and would like a limited comparison of
560 results from various DATs.

561 b) **Performance Evaluation:** The purpose of this module is to allow users to evaluate
562 the relative performance of various DATs based on the characteristics of their
563 experimental setting. A salient feature of this module is that users can adjust the
564 stringency thresholds for the DATs and immediately assess the impact on
565 performance, without re-running the DATs. The expected use case for this module is
566 when users have pilot data and would like to do a systematic evaluation of the DATs.

567 c) **Experimental Design:** This module allows users to specify desired performance
568 targets and the range of experimental choices that are feasible, to identify
569 combinations that can meet the targets as well as the appropriate DATs that can be
570 used to achieve them. Ideally, users in the planning stages of an experiment would use
571 this module to optimize their experimental design.

572 **Pre-processing of RNA-seq datasets**

573 Count data for RNA-seq datasets in ENCODE (for GM12892 NR=3, MCF7 NR=3 and h1-
574 hESC NR=4) were obtained directly from <http://bitbucket.org/soccin/seqc>. Count data for
575 Pickrell et al³⁸ (NR=69) was obtained from (<http://bowtie-bio.sourceforge.net/recount/>).
576 Reads from each library of the K562 single cell RNA-seq dataset were mapped uniquely and
577 independently using TopHat³⁹ (version 2.0.7) against the human reference genome (hg19).
578 Raw counts for each gene were then extracted using Human Gencode 19 annotations and
579 htseq-count (<http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>).

580 **Single-cell RNA-seq and Metagenomic Analysis**

581 RNA-seq count data from single-cell experiments in Ramskold et al²⁹ was obtained from
582 Supplementary Table 4 in the manuscript (15 thousand genes and 10 samples - 6 putative
583 circulating tumor cells and 4 from the melanoma cell line SKMEL5; RPKM values were
584 converted back to raw counts). Metagenomic count data from the study by Qin et al³¹ was
585 obtained from www.gigadb.org (1.14 million genes and 145 samples). The fit of the
586 metagenomic count data to the Negative Binomial distribution was assessed using a
587 Kolmogorov–Smirnov test, where < 0.1% of the genes failed the test at a *p*-value threshold of
588 0.01. Characteristics of both datasets were learned by EDDA and used to generate simulated
589 datasets (Model-Free for single-cell and with the Full Model for metagenomic data).

590 To build a predictive model for type 2 diabetes from the data in Qin et al³¹, we followed the
591 procedure described there to identify 50 marker genes from the top 1000 differentially
592 abundant genes (based on edgeR *p*-values) by employing the maximum relevance minimum
593 redundancy (mRMR) feature selection framework⁴⁰. The identified marker genes were
594 combined into a “T2D index” (= mean abundance of positive markers - mean abundance of
595 negative markers) for each sample, which was then used to rank samples and compute ROC
596 curves as in the original study.

597 **Nanostring Analysis**

598 Nanostring count data was obtained from an in-house preliminary study of prognostic and
599 predictive gene signatures for breast cancer (manuscript in preparation). Briefly, expression
600 levels of 107 genes of interest from 369 patients in different stages of breast cancer and with
601 known estrogen receptor (ER) alpha status and clinical outcomes were quantified using the
602 NanoString nCounter System (NanoString, WA, USA). The raw data was normalized by
603 different generally applicable methods (e.g. median-normalization as implemented in DESeq,
604 mode-normalization from the EDDA package and UQN as implemented in edgeR; see **Table**
605 **2**) as well as the recommended standard from the NanoString data analysis guide (normalized
606 by positive and negative controls, followed by global normalization). Note that as this dataset
607 has multiple categories we extended the standard two-condition version of mode-
608 normalization by randomly labelling samples as controls or cases to identify the top 10 genes
609 that are consistently chosen. In order to measure the impact of normalization on the ability to
610 separate patients based on their ER status a standard F-score was calculated, as the ratio of
611 between-group variance to within-group variance of mean counts for the 8 ER signature
612 genes (formally $F\text{-score} = F_{\text{Between}}/F_{\text{Within}}$, where $F_{\text{Between}} = |X_1|(E(X_1) - E(X))^2 +$
613 $|X_2|(E(X_2) - E(X))^2$ and $F_{\text{within}} = (|X_1|\text{Var}(X_1) + |X_2|\text{Var}(X_2))/(|X_1|+|X_2| - 2)$ and
614 $X = X_1 \cup X_2$, for mean counts X_1 and X_2 in the two groups).

615 **AVAILABILITY**

616 As an open-source R package at <https://sourceforge.net/projects/eddanorm/> or
617 <http://www.bioconductor.org/packages/devel/bioc/html/EDDA.html> and as web-modules at
618 <http://edda.gis.a-star.edu.sg>.

619 **ACKNOWLEDGEMENTS**

620 We would like to thank Shyam Prabhakar, Swaine Chen, Denis Bertrand, Sun Miao, Li Yi
621 and Chng Kern Rei for providing valuable feedback on drafts of the manuscript, Christopher
622 Wong for sharing the Nanostring dataset with us and Lili Sun, Naveen Ramalingam and Jay
623 West for technical assistance in generating the single-cell RNA-seq dataset.

624 **FUNDING**

625 This work was done as part of the IMaGIN platform (project No. 102 101 0025), supported
626 by a grant from the Science and Engineering Research Council as well as IAF grants
627 IAF311009 and IAF111091 from the Agency for Science, Technology and Research
628 (A*STAR), Singapore.

629 **AUTHORS CONTRIBUTIONS**

630 NN and CKHB initiated the project. LH and CKHB implemented EDDA with additional
631 inputs from LJ. LH designed and implemented mode-normalization with inputs from NN. PR
632 coordinated the single-cell RNA-seq data generation. LJ conducted the single-cell and
633 metagenomic analysis. LH conducted the nanostring analysis. NN, LH, LJ and CHKB wrote
634 the manuscript with inputs from all authors.

635 **COMPETING FINANCIAL INTERESTS**

636 None

637 **TABLES**

638 **Table 1. Experimental conditions affecting differential abundance tests (DATs).**

| | | Abbr. | Description | Notes |
|-------------------------------------|------------------------------------|--------------|--|---|
| Experimental Choices | Number of Replicates | NR | Number of technical or non-technical replicates for the two groups compared in the test | For simplicity, in many cases, we assume NR to be the same in both groups |
| | Number of Data-points | ND | Number of data-points generated in the counting experiment | e.g. reads generated in an RNA-seq experiment |
| Experimental Characteristics | Entity Count | EC | Number of entities in the counting experiment | e.g. number of genes in an RNA-seq experiment |
| | Sample Variability | SV | Variability across replicates (see Methods) | e.g. biological variability in RNA-seq datasets |
| | Abundance Profile | AP | Relative abundance of the entities in the first group | Typically follows a power-law distribution |
| | Perturbation Profile | PDA, FC | Perturbations to the abundance profile of the first group to obtain the profile for the second group (see Methods) | Used to generate the differentially abundant entities (PDA = Percentage of entities, FC = fold-change distribution) |
| Test Settings | Biases in Data Generation | | Deviations from multinomial sampling due to biases inherent in the experimental protocol | These are often corrected for in a pre-processing step e.g. composition-bias in RNA-seq data ^{35,41} |
| | Differential Abundance Test | DAT | See Table 2 | |

640 **Table 2. Description of various software packages for conducting differential**
 641 **abundance tests.** The various normalization approaches used are: TMM – Trimmed Mean of
 642 M values⁴², RPM – Normalization by read count (RNA-seq), RPKM – Normalization by read
 643 count and gene length (RNA-seq)³, UQN – Upper Quartile Normalization²⁰.

| Name | Statistical Testing | Normalization Approach | Target Application Areas |
|--------------------------------|--|-------------------------|---|
| edgeR ¹² | Negative Binomial Model, Conditional Maximum Likelihood to estimate parameters, Exact Test | TMM, UQN | SAGE ¹ , MPSS ⁴³ , PMAGE ⁴⁴ , miRAGE ⁴⁵ , SACO ⁴⁶ etc. |
| DESeq ⁶ | Negative Binomial Model, local regression to estimate parameters, Exact Test | Normalization by median | RNA-seq ³ , HITS-CLIP ⁹ , ChIP-seq ⁴⁷ etc. |
| baySeq ¹⁵ | Negative Binomial Model, empirical Bayes to estimate parameters | RPM, TMM | DNA-seq, RNA-seq ³ , SAGE ¹ etc. |
| NOISeq ¹⁶ | Non-parametric approach | RPM, RPKM, UQN | ChIP-seq ⁴⁷ , RNA-seq ³ etc. |
| Cuffdiff ¹⁴ | <i>t</i> -test | RPM, RPKM, UQN | RNA-seq ³ |
| Metastats ¹¹ | Non-parametric <i>t</i> -test, Fisher's Exact Test for small counts | RPM | Metagenomics |

644

645 **FIGURE LEGENDS**

646 **Figure 1. A schematic overview of EDDA and its usage.** (a) Overview of the simulate-and-
647 test approach used in EDDA involving 4 key steps (see **Methods** for details). The figure also
648 indicates the various parameters studied here, that impact the ability to predict differentially
649 abundant entities (also see **Table 1**), in red and green text. (b) A flowchart describing how
650 EDDA can be used to optimize design and performance in a counting experiment. EDDA
651 functionality is depicted in green boxes and rhombi indicate decision nodes.

652 **Figure 2. Performance of DATs as a function of Experimental Choices.** (a) AUC as a
653 function of the number of replicates. AUC as a function of the number of data-points with (b)
654 1 replicate (c) 2 replicates and (d) 5 replicates. Note that the number of data-points reported is
655 the average per entity. Points with large markers indicate the saturation point where the AUC
656 is within 5% of the maximum observed AUC. Results reported are based on an average over
657 10 simulations with the following parameters, EC=1000, PDA=26%, FC=Uniform[3, 5],
658 ND=1000 per entity, AP=BP, SM=NB and SV=0.85 (see **Table 1** and **Methods** for key to
659 abbreviations used here).

660 **Figure 3. Performance of DATs as a function of Experimental Characteristics.** (a) AUC
661 as a function of entity count (small=100, medium=1,000, large=30,000). Simulations were
662 done with the same parameters as in **Figure 2a**, except with ND=100 per entity, AP=HBR,
663 SM=Full and FC=Log-normal[2, 1] (b) AUC as a function of SV (and NR, ND per entity),
664 where low, medium and high were set to SV=0.05, 0.5 and 0.85 respectively. Note that the
665 plotted point shows the smallest NR and ND values (small NR was given priority in the case
666 of ties) that achieved an AUC target of 0.95. Synthetic data was generated with EC=1000,
667 PDA=10%, FC=Uniform[8, 8], AP=BP and SM=NB. (c) AUC as a function of abundance
668 profile. Simulations were done with the same parameters as in **Figure 2a** and with NR=5 and
669 ND=500 per entity. The inset shows the fraction of rare entities (mean count < 10) in each
670 abundance profile. Results reported are based on an average over 10 simulations and error
671 bars in subfigures a and c represent one standard deviation.

672 **Figure 4. ROC plots under various Perturbation Profiles.** Statistical tests in the legend are
673 listed from best to worst (in terms of AUC values) for each setting. Note the striking re-
674 ordering of test performance across subfigures with slight changes in experimental
675 conditions. Simulations in EDDA were done with (a) PDA=10% (b) PDA=[50% UP, 25%
676 DOWN] (c) PDA=[50% UP, 25% DOWN], FC= Log-normal[2, 1] (d) PDA=[50% UP, 25%

677 DOWN], AP = HBR. Unless stated otherwise, common parameters include NR=3, EC=1000,
678 FC=Uniform[3, 7], ND=500 per entity, AP=BP, SM=Full and SV=0.5.

679 **Figure 5. Comparison of different data normalization approaches.** (a) On simulated
680 benchmark datasets using EDDA. Note that we used normalization by the sum of counts for
681 all non-differentially-abundant entities (non-DA) as a measure of ideal performance here. In
682 general, Upper Quartile Normalization (UQN) improved over the Default normalization
683 (improvement shown by the checked box) but for cases where it did not we mark its
684 performance with a solid line. Mode normalization always improved over the performance
685 from UQN and the Default normalization for the DAT (improvement shown by the solid
686 box). The parameters for the various experiments include for A: PDA=[26% UP, 10%
687 DOWN], AP=Wu et al B: PDA=[35% UP, 15% DOWN] C: PDA=[40%
688 DOWN],SM=Mutinomial and D: the same as in **Figure 4b**. Unless stated otherwise, common
689 parameters include NR=3, EC=1000, FC=Log-normal[1.5, 1], ND=500 per entity, AP=HBR,
690 SM=Full and SV=0.5. (b) Comparison on real datasets, highlighting the robustness of mode
691 normalization. Shown here is the overlap fraction (used to measure robustness) of the top 500
692 differentially abundant genes (sorted by p -value using edgeR) in comparisons involving both
693 the original full-size libraries versus those where one of the libraries is down-sampled (5% or
694 10% of original size). Barplots show the average of 5 runs and error bars represent one
695 standard deviation.

696 **Figure 6. Case studies of EDDA usage.** (a) Predicted performance of DATs as a function of
697 replicates (NR) and sequencing depth (ND) for a single cell RNA-seq experiment (Ramskold
698 et al²⁹). Note that the vertical black dashed line indicates design choices in the original study.
699 (b) Coefficient of variation for a panel of housekeeping genes from a NanoString experiment
700 under different normalization approaches. (c) Separability of ER positive and negative
701 samples under different normalization approaches. (d) Predicted performance of DATs as a
702 function of replicates (NR) and sequencing depth (ND) for a Metagenomics experiment (Qin
703 et al³¹). Note that the vertical black dashed line indicates design choices in the original study.

704 **REFERENCES**

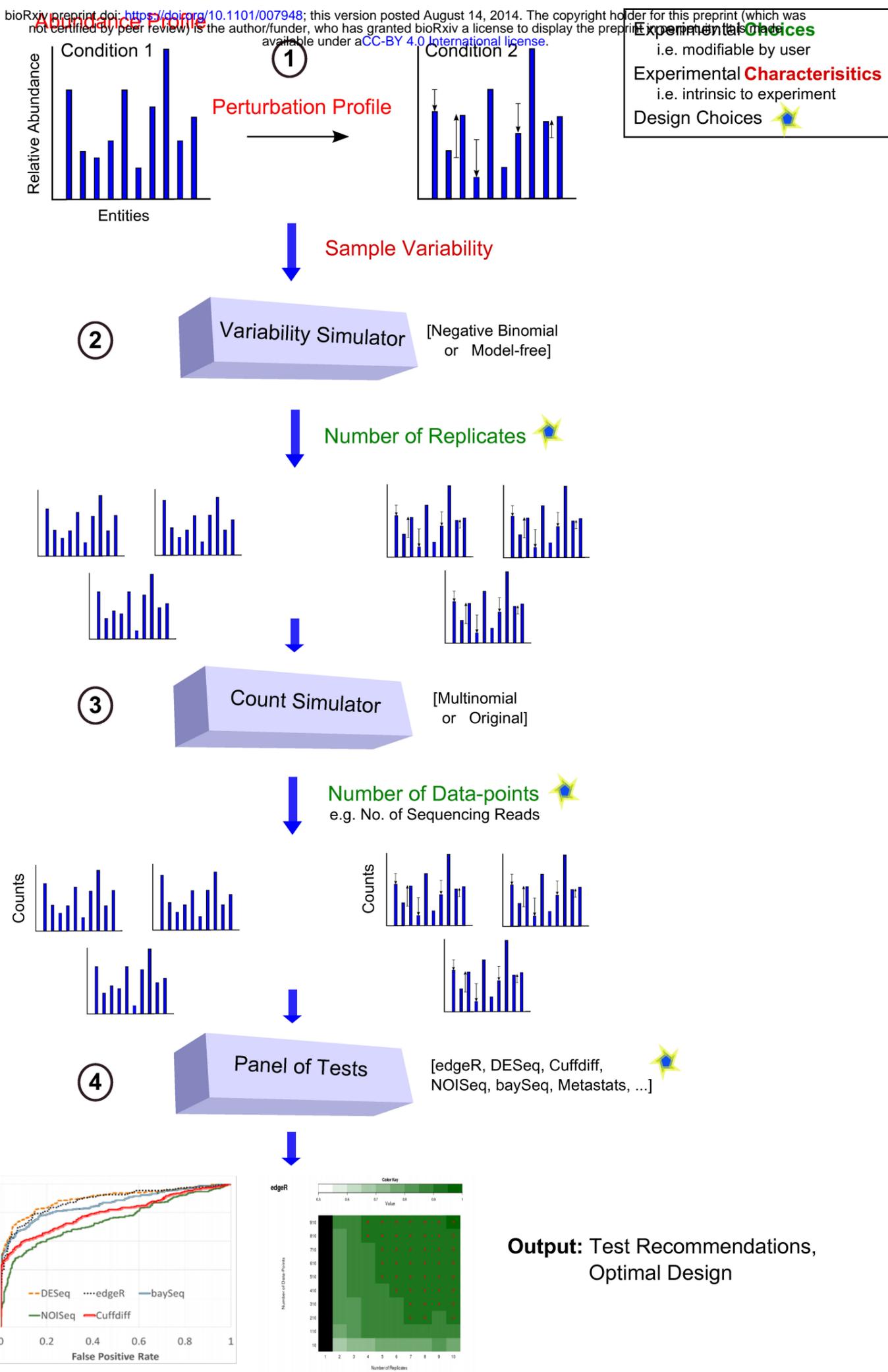
- 705 1. Velculescu, V.E., Zhang, L., Vogelstein, B. & Kinzler, K.W. Serial analysis of gene expression.
706 *Science* **270**, 484-487 (1995).
- 707 2. Ng, P. et al. Gene identification signature (GIS) analysis for transcriptome characterization
708 and genome annotation. *Nat Methods* **2**, 105-111 (2005).
- 709 3. Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying
710 mammalian transcriptomes by RNA-Seq. *Nat Methods* **5**, 621-628 (2008).
- 711 4. Geiss, G.K. et al. Direct multiplexed measurement of gene expression with color-coded
712 probe pairs. *Nat Biotechnol* **26**, 317-325 (2008).
- 713 5. Meyerson, M., Gabriel, S. & Getz, G. Advances in understanding cancer genomes through
714 second-generation sequencing. *Nat Rev Genet* **11**, 685-696.
- 715 6. Anders, S. & Huber, W. Differential expression analysis for sequence count data. *Genome*
716 *biology* **11**, R106-R106.
- 717 7. Zhang, Y. et al. Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137 (2008).
- 718 8. Zhao, J. et al. Genome-wide identification of polycomb-associated RNAs by RIP-seq. *Mol Cell*
719 **40**, 939-953.
- 720 9. Licatalosi, D.D. et al. HITS-CLIP yields genome-wide insights into brain alternative RNA
721 processing. *Nature* **456**, 464-469 (2008).
- 722 10. Ong, S.H. et al. Species identification and profiling of complex microbial communities using
723 shotgun Illumina sequencing of 16S rRNA amplicon sequences. *PLoS One* **8**, e60811.
- 724 11. White, J.R., Nagarajan, N. & Pop, M. Statistical methods for detecting differentially abundant
725 features in clinical metagenomic samples. *PLoS Comput Biol* **5**, e1000352 (2009).
- 726 12. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor package for differential
727 expression analysis of digital gene expression data. *Bioinformatics (Oxford, England)* **26**, 139-
728 140.
- 729 13. Trapnell, C. et al. Transcript assembly and quantification by RNA-Seq reveals unannotated
730 transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**, 511-515.
- 731 14. Trapnell, C. et al. Differential analysis of gene regulation at transcript resolution with RNA-
732 seq. *Nat Biotechnol* **31**, 46-53.
- 733 15. Hardcastle, T.J. & Kelly, K.a. baySeq: empirical Bayesian methods for identifying differential
734 expression in sequence count data. *BMC bioinformatics* **11**, 422-422.
- 735 16. Tarazona, S., Garc a-Alcalde, F., Dopazo, J.n., Ferrer, A. & Conesa, A. Differential expression
736 in RNA-seq: a matter of depth. *Genome research* **21**, 2213-2223.
- 737 17. Wang, L., Feng, Z., Wang, X. & Zhang, X. DEGseq: an R package for identifying differentially
738 expressed genes from RNA-seq data. *Bioinformatics* **26**, 136-138.
- 739 18. Nookaew, I. et al. A comprehensive comparison of RNA-Seq-based transcriptome analysis
740 from reads to differential gene expression and cross-comparison with microarrays: a case
741 study in *Saccharomyces cerevisiae*. *Nucleic Acids Res* **40**, 10084-10097.
- 742 19. Sonesson, C. & Delorenzi, M. A comparison of methods for differential expression analysis of
743 RNA-seq data. *BMC Bioinformatics* **14**, 91.
- 744 20. Bullard, J.H., Purdom, E., Hansen, K.D. & Dudoit, S. Evaluation of statistical methods for
745 normalization and differential expression in mRNA-Seq experiments. *BMC bioinformatics* **11**,
746 94-94.
- 747 21. Seyednasrollah, F., Laiho, A. & Elo, L.L. Comparison of software packages for detecting
748 differential expression in RNA-seq studies. *Brief Bioinform.*
- 749 22. Busby, M.A., Stewart, C., Miller, C.A., Grzeda, K.R. & Marth, G.T. Scotty: a web tool for
750 designing RNA-Seq experiments to measure differential gene expression. *Bioinformatics* **29**,
751 656-657.
- 752 23. Wu, J.Q. et al. Transcriptome sequencing revealed significant alteration of cortical promoter
753 usage and splicing in schizophrenia. *PloS one* **7**, e36351-e36351.

- 754 24. Robinson, M.D. & Smyth, G.K. Small-sample estimation of negative binomial dispersion, with
755 applications to SAGE data. *Biostatistics (Oxford, England)* **9**, 321-332 (2008).
- 756 25. Lu, J., Tomfohr, J.K. & Kepler, T.B. Identifying differential expression in multiple SAGE
757 libraries: an overdispersed log-linear model approach. *BMC Bioinformatics* **6**, 165 (2005).
- 758 26. Brennecke, P. et al. Accounting for technical noise in single-cell RNA-seq experiments. *Nat*
759 *Methods* **10**, 1093-1095.
- 760 27. Deng, Q., Ramskold, D., Reinius, B. & Sandberg, R. Single-cell RNA-seq reveals dynamic,
761 random monoallelic gene expression in mammalian cells. *Science* **343**, 193-196.
- 762 28. Rapaport, F. et al. Comprehensive evaluation of differential gene expression analysis
763 methods for RNA-seq data. *Genome Biol* **14**, R95.
- 764 29. Ramskold, D. et al. Full-length mRNA-Seq from single-cell levels of RNA and individual
765 circulating tumor cells. *Nat Biotechnol* **30**, 777-782.
- 766 30. Yu, K., Ganesan, K., Miller, L.D. & Tan, P. A modular analysis of breast cancer reveals a novel
767 low-grade molecular signature in estrogen receptor-positive tumors. *Clin Cancer Res* **12**,
768 3288-3296 (2006).
- 769 31. Qin, J. et al. A metagenome-wide association study of gut microbiota in type 2 diabetes.
770 *Nature* **490**, 55-60.
- 771 32. Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach for
772 interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-15550
773 (2005).
- 774 33. Garmire, L.X. & Subramaniam, S. Evaluation of normalization methods in mammalian
775 microRNA-Seq data. *RNA* **18**, 1279-1288.
- 776 34. Dillies, M.A. et al. A comprehensive evaluation of normalization methods for Illumina high-
777 throughput RNA sequencing data analysis. *Brief Bioinform* **14**, 671-683.
- 778 35. Jones, D.C., Ruzzo, W.L., Peng, X. & Katze, M.G. A new approach to bias correction in RNA-
779 Seq. *Bioinformatics* **28**, 921-928.
- 780 36. Au, K.F., Jiang, H., Lin, L., Xing, Y. & Wong, W.H. Detection of splice junctions from paired-
781 end RNA-seq data by SpliceMap. *Nucleic Acids Res* **38**, 4570-4578.
- 782 37. Rosenblatt, M. Remarks on Some Nonparametric Estimates of a Density Function. *The*
783 *Annals of Mathematical Statistics* **27** (1956).
- 784 38. Pickrell, J.K. et al. Understanding mechanisms underlying human gene expression variation
785 with RNA sequencing. *Nature* **464**, 768-772.
- 786 39. Kim, D. et al. TopHat2: accurate alignment of transcriptomes in the presence of insertions,
787 deletions and gene fusions. *Genome Biol* **14**, R36.
- 788 40. Peng, H., Long, F. & Ding, C. Feature selection based on mutual information: criteria of max-
789 dependency, max-relevance, and min-redundancy. *IEEE Trans Pattern Anal Mach Intell* **27**,
790 1226-1238 (2005).
- 791 41. Roberts, A., Trapnell, C., Donaghey, J., Rinn, J.L. & Pachter, L. Improving RNA-Seq expression
792 estimates by correcting for fragment bias. *Genome Biol* **12**, R22.
- 793 42. Robinson, M.D. & Oshlack, A. A scaling normalization method for differential expression
794 analysis of RNA-seq data. *Genome biology* **11**, R25-R25.
- 795 43. Brenner, S. et al. Gene expression analysis by massively parallel signature sequencing (MPSS)
796 on microbead arrays. *Nat Biotechnol* **18**, 630-634 (2000).
- 797 44. Kim, J.B. et al. Polony multiplex analysis of gene expression (PMAGE) in mouse hypertrophic
798 cardiomyopathy. *Science* **316**, 1481-1484 (2007).
- 799 45. Cummins, J.M. et al. The colorectal microRNAome. *Proc Natl Acad Sci U S A* **103**, 3687-3692
800 (2006).
- 801 46. Impey, S. et al. Defining the CREB regulon: a genome-wide analysis of transcription factor
802 regulatory regions. *Cell* **119**, 1041-1054 (2004).
- 803 47. Robertson, G. et al. Genome-wide profiles of STAT1 DNA association using chromatin
804 immunoprecipitation and massively parallel sequencing. *Nat Methods* **4**, 651-657 (2007).

805

806

a



b

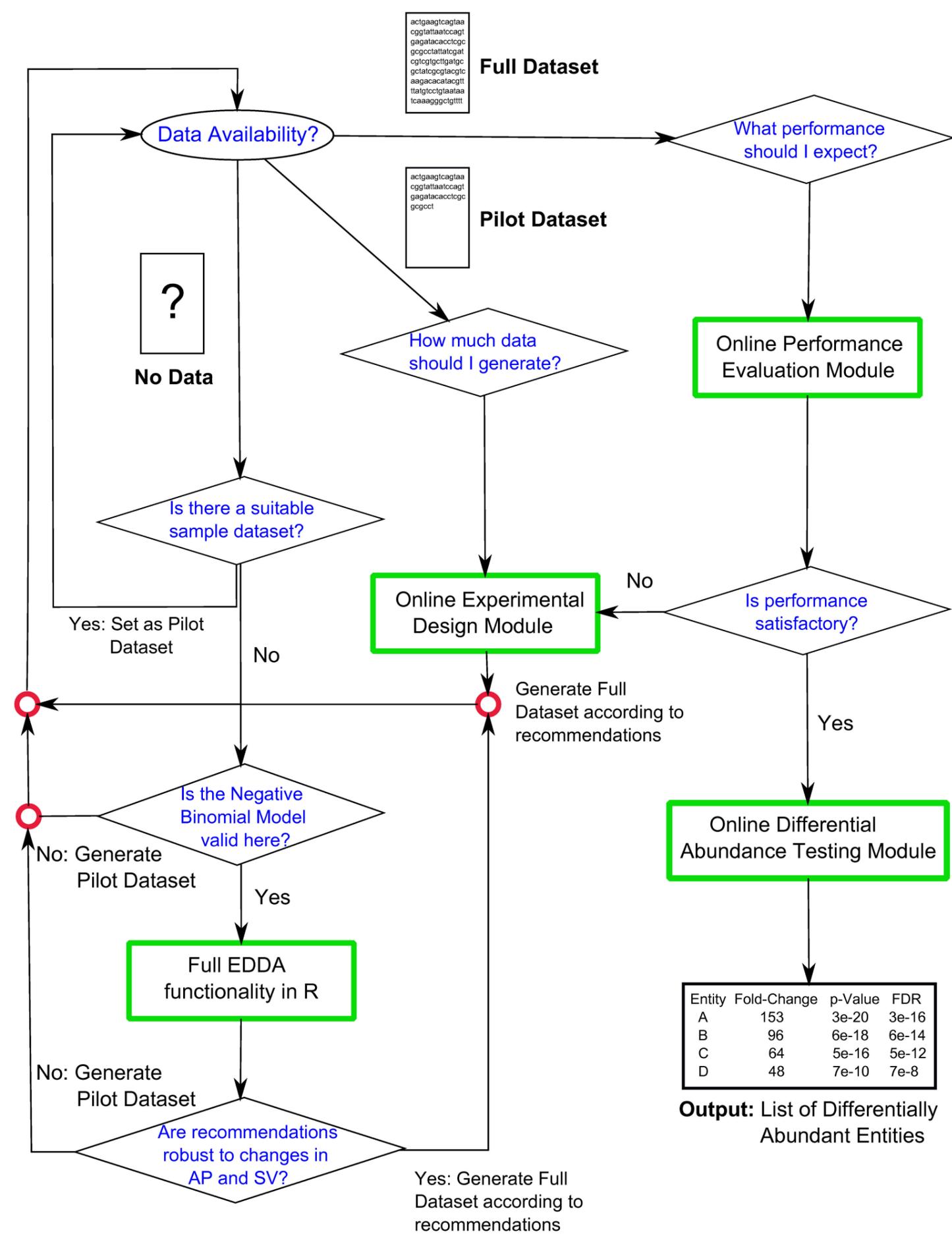


Figure 1

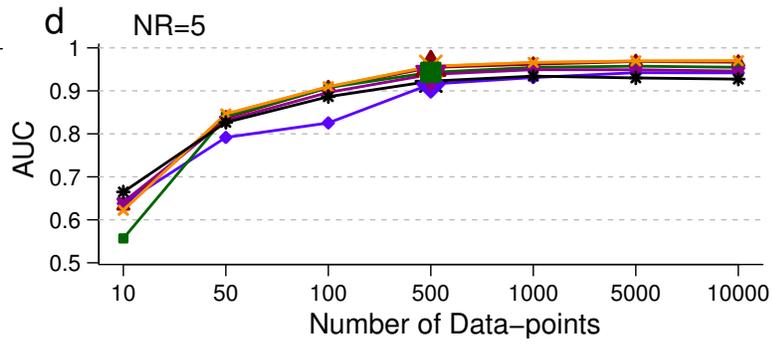
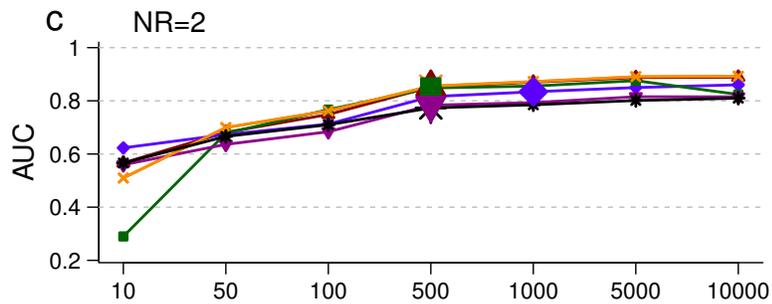
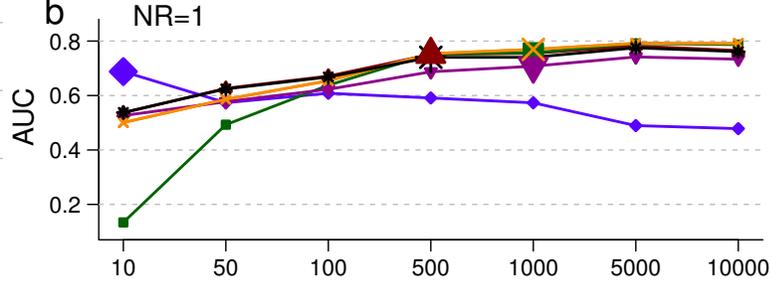
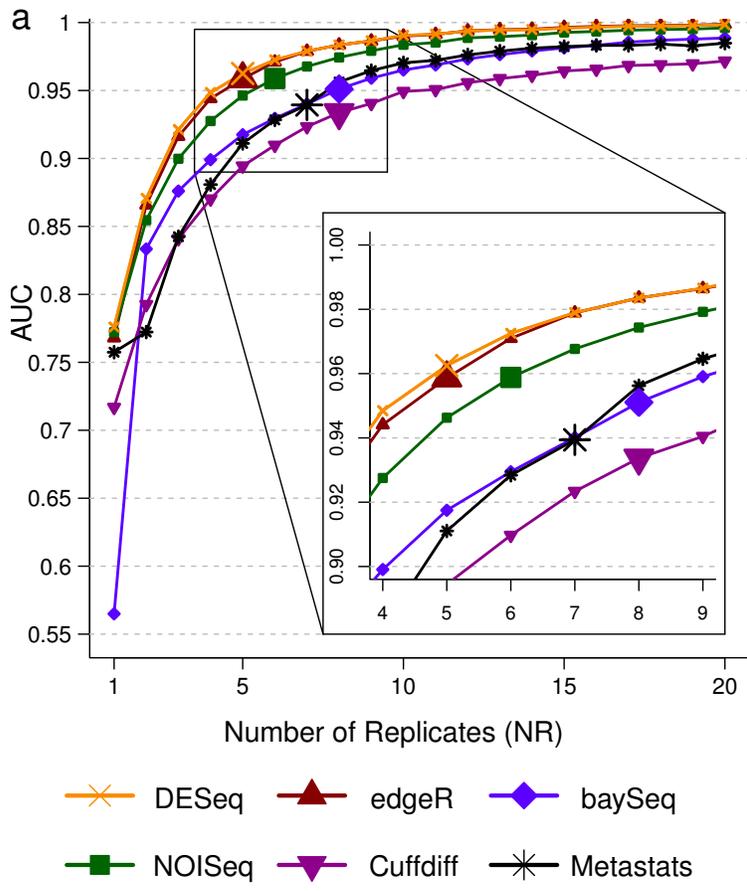
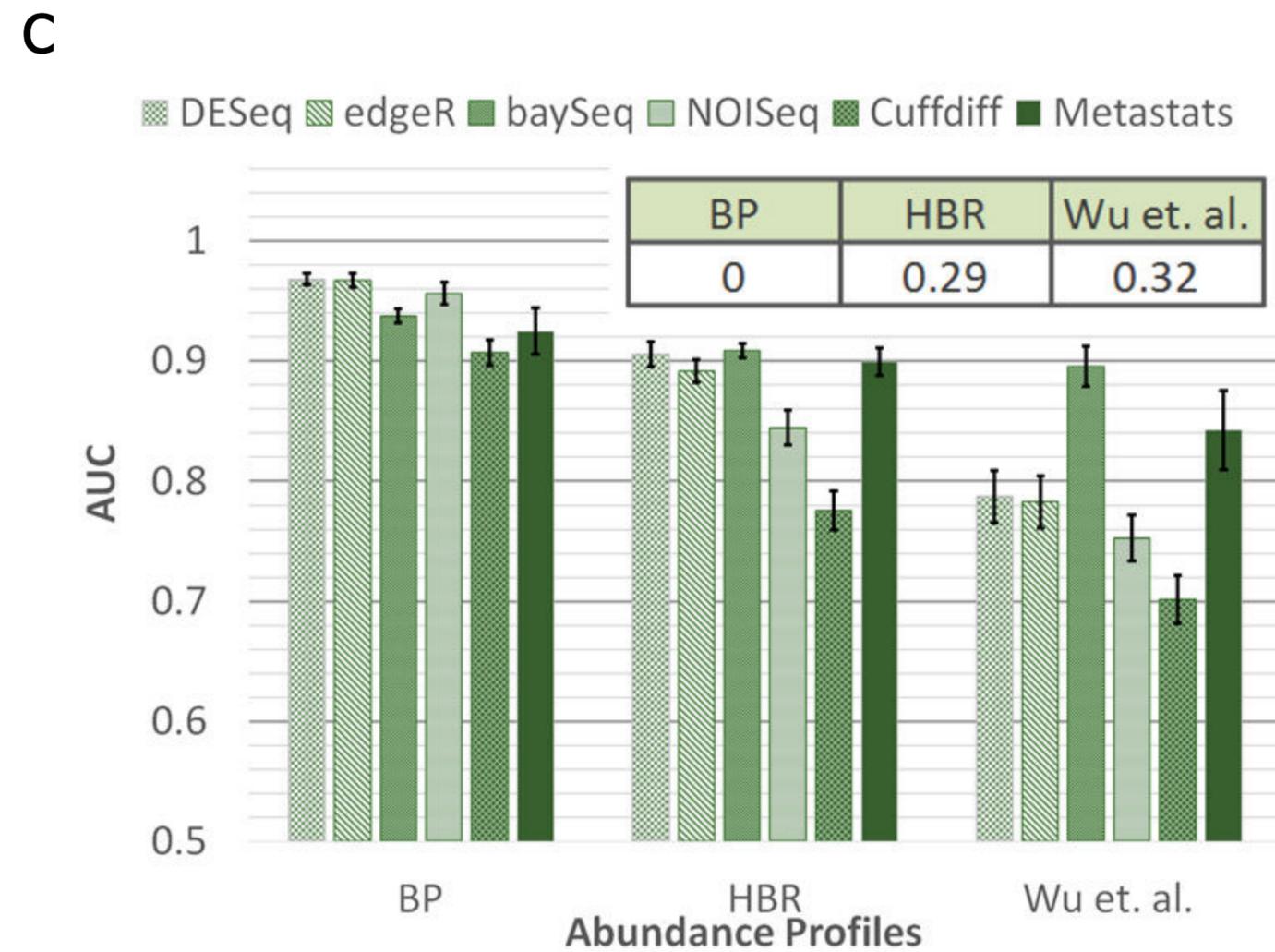
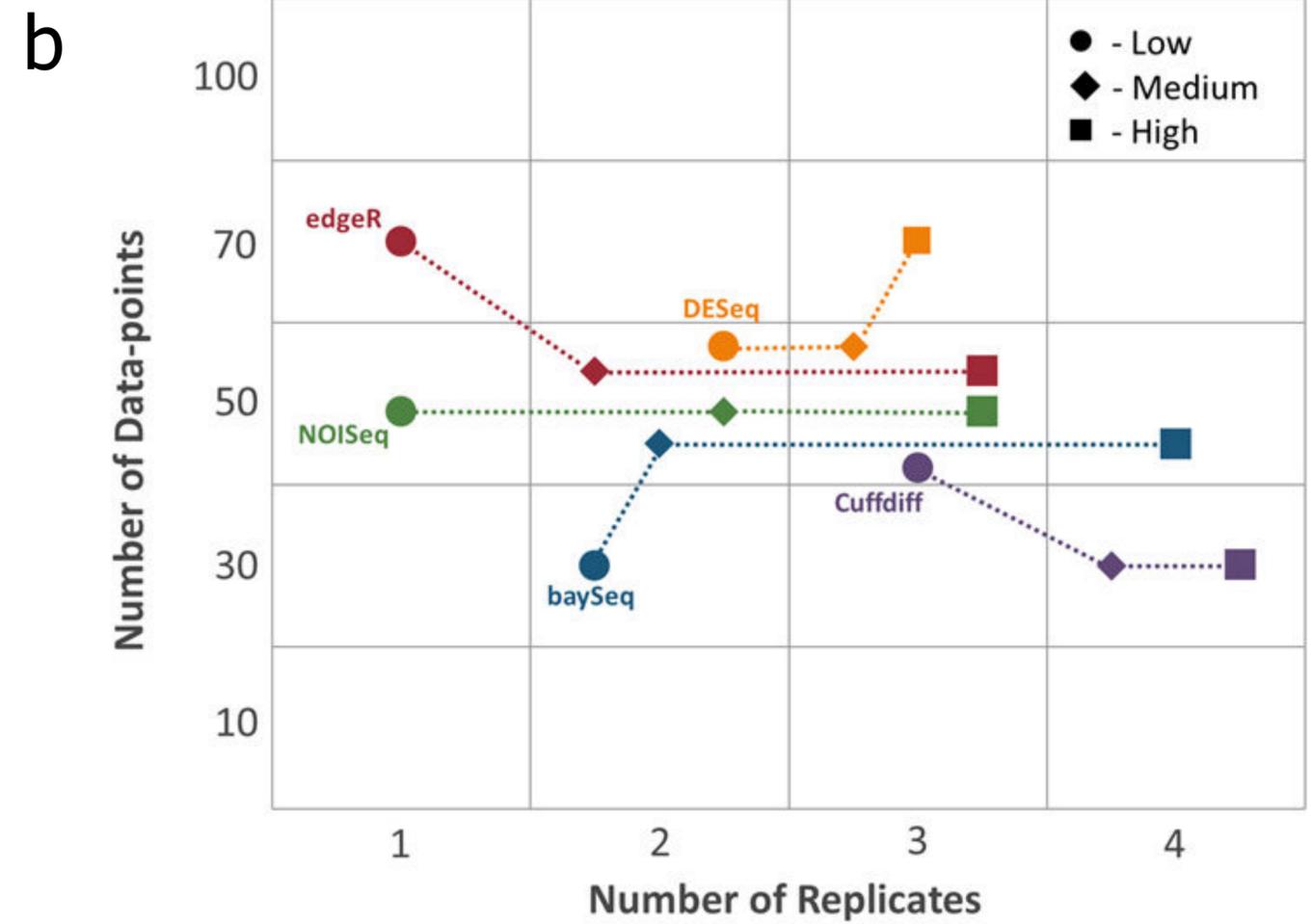
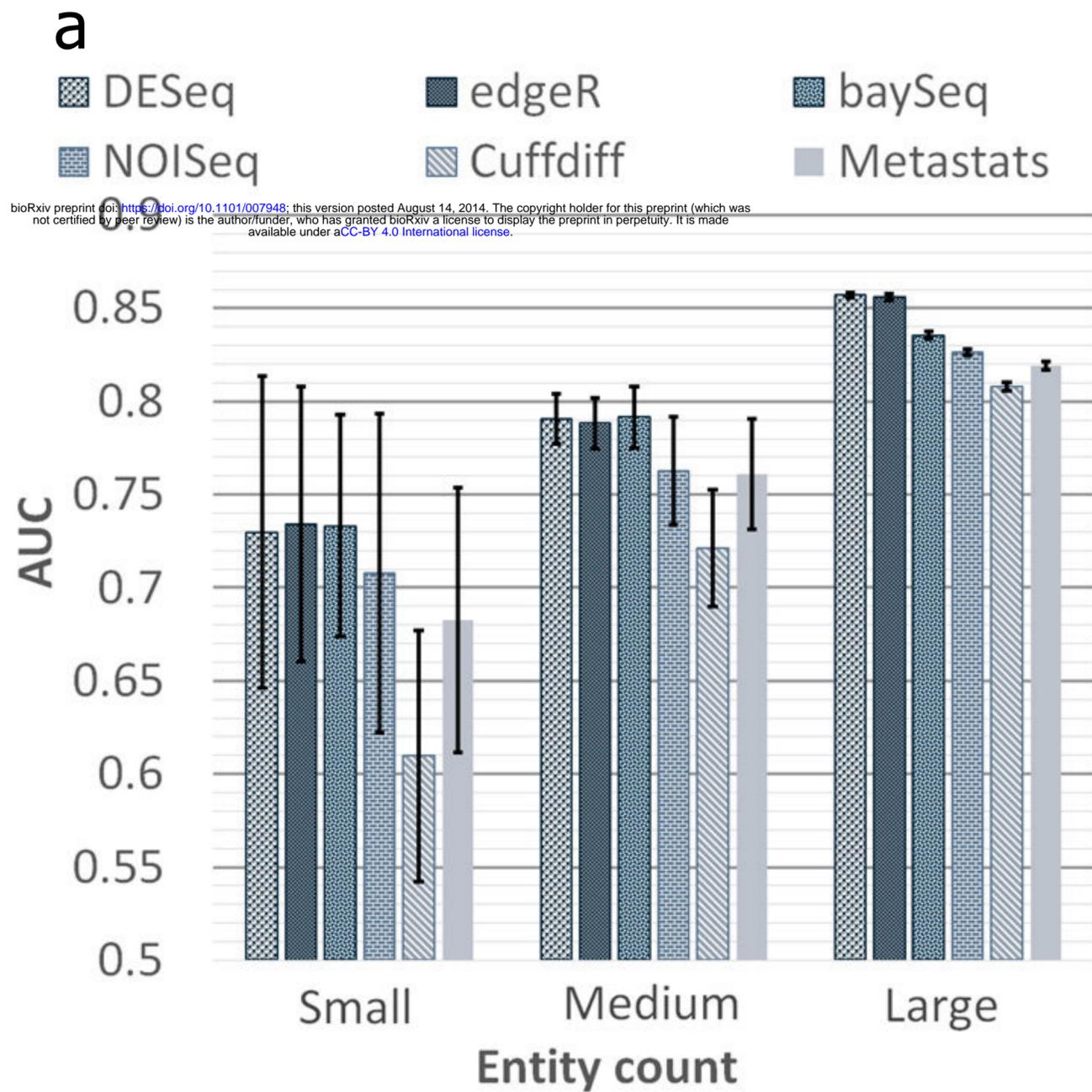
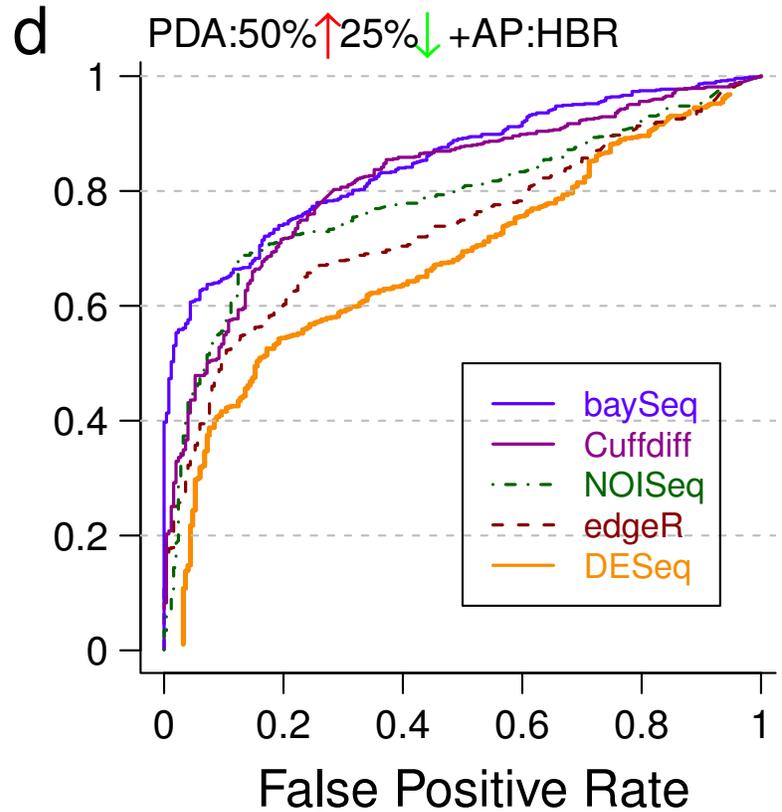
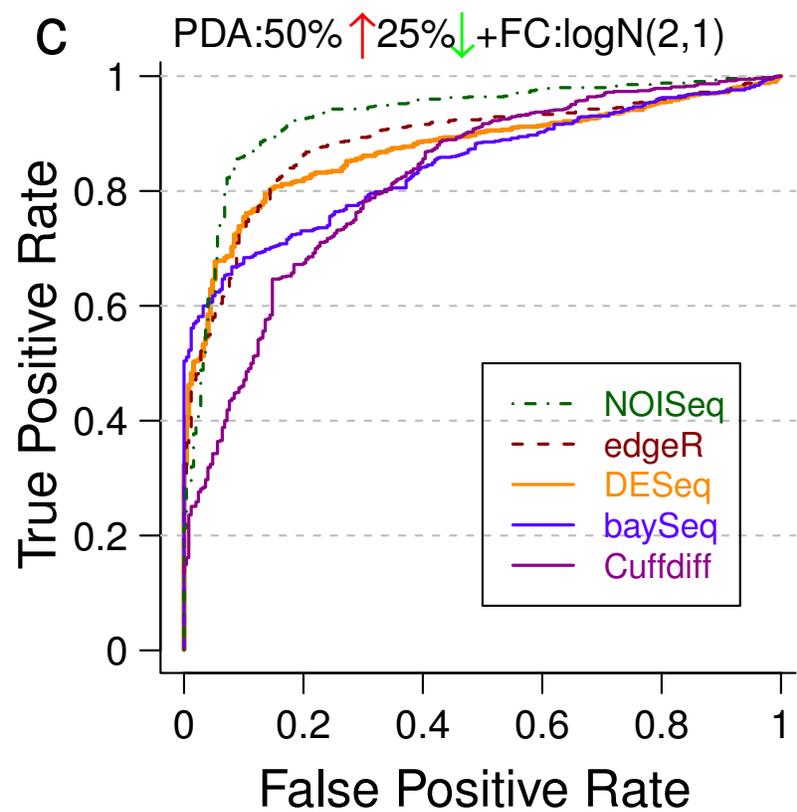
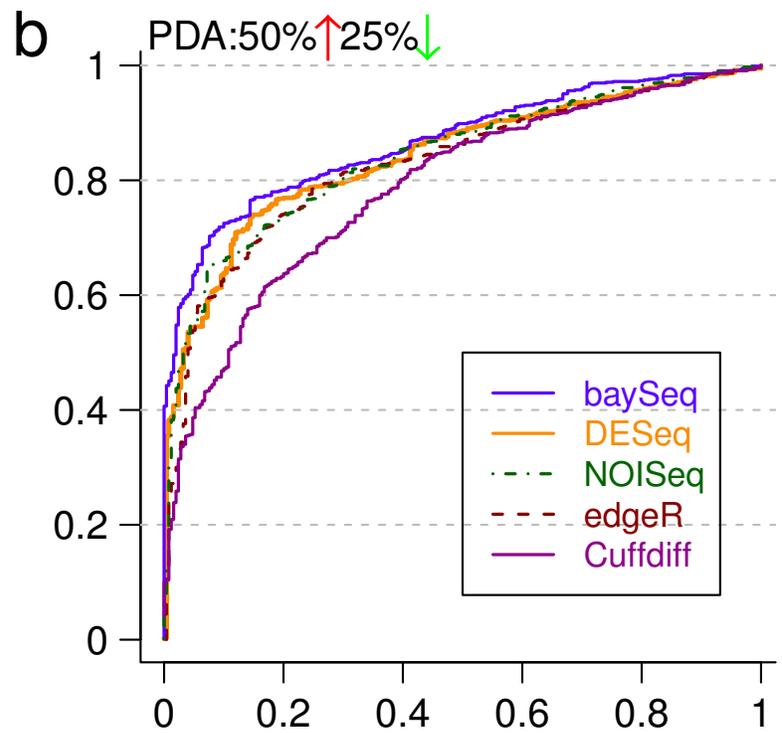
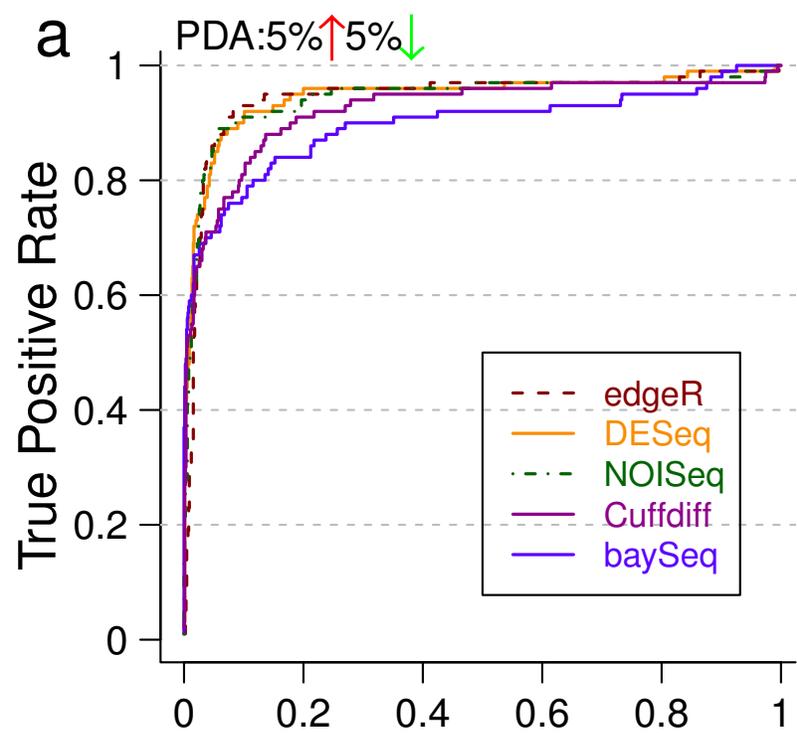
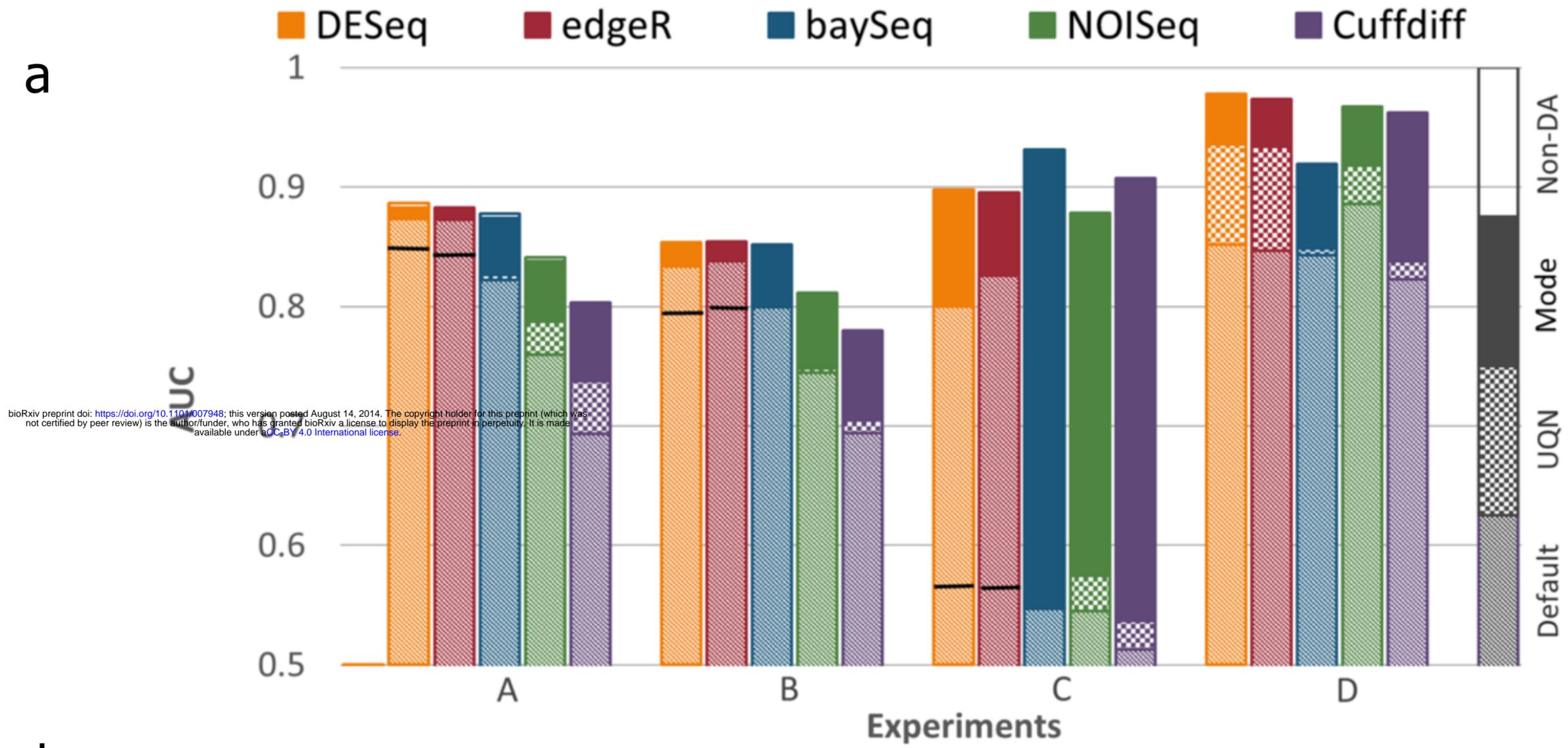


Figure 2

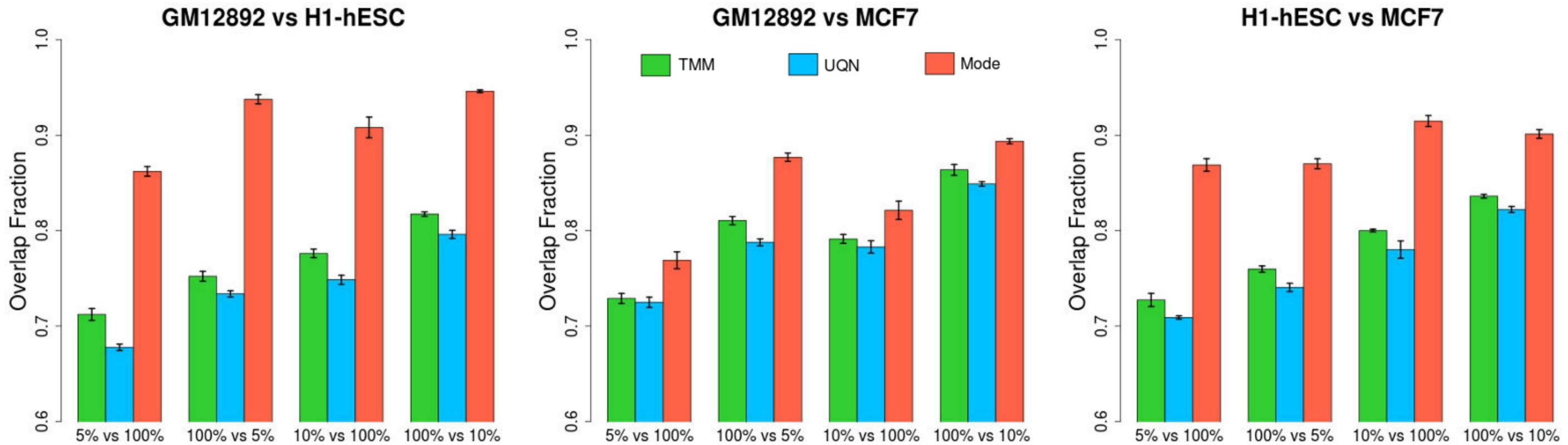




a



b



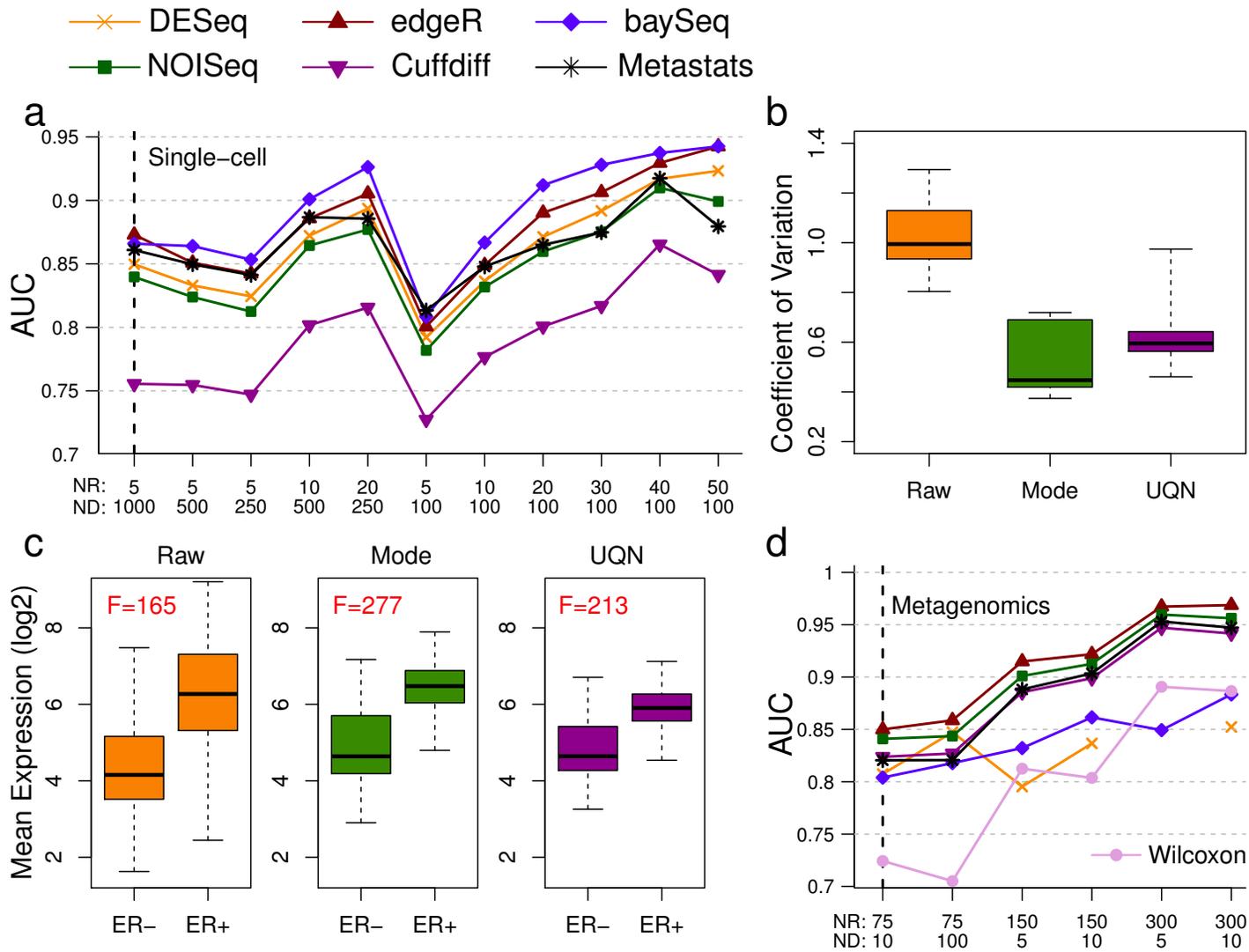


Figure 6