

Fast gene set enrichment analysis

Gennady Korotkevich^{1,*}, Vladimir Sukhov^{1,2,*}, Nikolay Budin¹, Boris Shpak¹,
Maxim N. Artyomov³, and Alexey Sergushichev^{1,**}

¹Computer Technologies Laboratory, ITMO University, Saint Petersburg,
197101, Russia

²JetBrains Research, Saint Petersburg, Russia

³Washington University in St. Louis, St. Louis, MO, USA

* contributed equally

** corresponding author, e-mail: alserg@itmo.ru

Abstract

Gene set enrichment analysis (GSEA) is an ubiquitously used tool for evaluating pathway enrichment in transcriptional data. Typical experimental design consists in comparing two conditions with several replicates using a differential gene expression test followed by preranked GSEA performed against a collection of hundreds and thousands of pathways. However, the reference implementation of this method cannot accurately estimate small P-values, which significantly limits its sensitivity due to multiple hypotheses correction procedure.

Here we present FGSEA (Fast Gene Set Enrichment Analysis) method that is able to estimate arbitrarily low GSEA P-values with a high accuracy in a matter of minutes

20 or even seconds. To confirm the accuracy of the method, we also developed an exact
21 algorithm for GSEA P-values calculation for integer gene-level statistics. Using the
22 exact algorithm as a reference we show that FGSEA is able to routinely estimate P-
23 values up to 10^{-100} with a small and predictable estimation error. We systematically
24 evaluate FGSEA on a collection of 605 datasets and show that FGSEA recovers much
25 more statistically significant pathways compared to other implementations.

26 FGSEA is open source and available as an R package in Bioconductor ([http://](http://bioconductor.org/packages/fgsea/)
27 bioconductor.org/packages/fgsea/) and on GitHub (<https://github.com/ctlab/fgsea/>).

28 **1 Main**

29 Preranked gene set enrichment analysis (GSEA) [1] is a widely used method for analyzing
30 gene expression data, particularly for datasets with small number of replicates. It allows to
31 select from an *a priori* defined collection of pathways those which have non-random behavior
32 in a considered experiment (Fig 1a). The method uses an enrichment score (ES) statistic
33 which is calculated based on a vector of gene-level signed statistics, such as *t*-statistic from
34 a differential expression test. As the analytical form of the null distribution for the ES
35 statistic is not known, empirical null distribution has to be calculated. That can be done
36 in a straightforward manner by sampling random gene sets as was done in the reference
37 implementation [1] and reimplementations [2, 3]. In this case for each of the input pathways
38 a number of random gene sets of the same size are generated, and for each of them an ES
39 value is calculated. Then a P-value is estimated as the number of random gene sets with
40 the same or more extreme ES value divided by the total number of generated gene sets (a
41 formal definition is available in the section 2.1). Finally, a multiple hypothesis correction
42 procedure is applied to get adjusted P-values.

43 However, a large number of generated random gene sets can be required to reach a
44 given false discovery rate (FDR) level on some datasets. As an example, we calculated

45 GSEA P-values for Gene Ontology Biological Pathways collection (C5_BP subset of MSigDB
46 collection [1]) on four datasets from Gene Expression Omnibus (GEO) varying the sampling
47 depth, and calculated the number of pathways reaching FDR level of 0.01 after Benjamini-
48 Hochberg (BH) correction. Due to the properties of BH procedure, the dependence of the
49 number of significant pathways in these experiments has a phase transition behavior (Fig 1b):
50 for each of the datasets there exist a certain critical sampling depth after which the number
51 of significant pathways becomes non-zero and stays on the same level. This critical sampling
52 depth is different for different datasets, but ultimately can reach the order of M/α , where α
53 is the selected FDR threshold and M is the number of considered pathways.

54 To systematically assess the distribution of GSEA critical sampling depth on real datasets
55 we prepared a collection of 605 microarray datasets from Gene Expression Omnibus (GEO)
56 containing only two biological conditions. For each of these datasets we ran the differential
57 expression analysis and used the results as gene-level statistics. We discovered that more than
58 half of the datasets has the critical sampling depth of at least 10^4 and a noticeable portion
59 (10—20% depending on the collection) has the critical sampling depth of at least 10^5 (Fig 1c).
60 When a large pathway collection is considered (entire MSigDB collection) individual datasets
61 has values of critical sampling depths reaching $5 \cdot 10^6$. However, even running the reference
62 implementation with the sampling depth of $n = 10^4$ routinely is inconvenient and running
63 it with $n = 10^5$ can be impossible due to the time and memory consumption (Fig 1d): time
64 and memory requirements grow linearly with the number of samples and the collection size.

65 To improve applicability of preranked GSEA analysis we present a fast gene set enrich-
66 ment analysis (FGSEA) method for accurate and efficient estimation of GSEA P-values for
67 a collection of pathways. The method consist of two main procedures: *FGSEA-simple* and
68 *FGSEA-multilevel*. FGSEA-simple procedure allows to efficiently estimate P-values with a
69 limited accuracy but simultaneously for the whole *collection* of gene sets, while FGSEA-
70 multilevel procedure allows to accurately estimate arbitrarily low P-values but for *individual*

71 gene sets.

72 FGSEA-simple procedure is based on an idea that generated random gene set samples
73 can be shared between different input pathways. Indeed, consider M gene sets of the sizes
74 $K_1 \leq K_2 \leq \dots \leq K_M = K$ and a collection of n independent samples g_i of size K (Fig 2a).
75 As in the naive approach, due to g_i being independent samples of the size K the P-value
76 for the pathway M can be estimated as a proportion of samples g_i having the same or more
77 extreme ES value as the pathway M . However, for any other pathway j we can construct
78 a set of n independent samples of size K_j by considering the prefixes $g_{i,1..K_j}$. Again, given
79 a set of independent samples, the P-value can be estimated as a proportion of the samples
80 having the same or more extreme ES value.

81 The next important idea is that given a gene set sample g_i of the size K the ES values for
82 *all* the prefixes $g_{i,1..j}$ can be calculated in an efficient manner using a square root heuristic
83 (Fig 2b). Briefly, a variant of an enrichment curve is considered: the genes are enumerated
84 starting from the most up-regulated to the most down-regulated, with the curve going to
85 the right if the gene is not present in the pathway, and the curve goes upward if the gene is
86 present in the pathway. It can be shown that the enrichment score can be easily calculated if
87 curve point most distant from the diagonal is known. Let us split K genes from the gene set
88 into $b \approx \sqrt{K}$ consecutive blocks of size \sqrt{K} and consider what happens with the curve when
89 we change the prefix from $g_{i,1..j-1}$ to $g_{i,1..j}$ by adding gene $g_{i,j}$. The curve in the blocks to the
90 left of $g_{i,j}$ are not changed at all, while the blocks to the right of $g_{i,j}$ are uniformly shifted.
91 This observation allows us to consider the prefixes in an increasing order and update the
92 position of the most distant point in $O(\sqrt{K})$ time. Briefly, for the each block which is either
93 not changed or shifted the update procedure takes $O(1)$ time, while for the changed block the
94 update procedure is proportional to its size and takes $O(\sqrt{K})$ time. Finally, aggregating the
95 blocks takes additional $O(\sqrt{K})$ time. Overall this results in time complexity of $O(K\sqrt{K})$
96 to calculate ES values for all the prefixes. In total, the time complexity of the calculating

97 P-values for the set of M pathways is $n(K\sqrt{K} + M)$, which gives around $O(K\frac{\log K}{\sqrt{K}})$ speed
98 up compared to a naive approach. The full description of the algorithm is given in the
99 section 2.3.

100 As an example we ran FGSEA-simple and the reference implementations on the same
101 example dataset of genes differentially regulated on Th1 activation [4] against a set of 700 Re-
102 actome [5] pathways (see section 2.2) and compared the resulting nominal P-values (Fig 2c).
103 Both methods were ran with $n = 10000$ and the results are indistinguishable from each
104 other up to the random noise inherent to both methods. However, on this example the
105 reference implementation took about 420 seconds, while FGSEA-simple finished in about 4
106 seconds. The two order of magnitude speed-up is consistent with the theoretical one due to
107 the algorithm time complexity. Given a highly parallel implementation of FGSEA-simple,
108 its performance allows to routinely achieve sampling depth of 10^5 and accurately estimate
109 P-values as low as 10^{-5} .

110 However, accurately estimating P-values lower than 10^{-5} with FGSEA-simple can be
111 impractical or even infeasible. To estimate such low P-values we developed *FGSEA-multilevel*
112 method, which is based on an adaptive multi-level split Monte Carlo scheme [6]. The method
113 takes as an input an ES value $\gamma > 0$ and a gene set size K , and calculates the probability
114 $P_K(\text{ES} \geq \gamma)$ of a random gene set of size K to have an enrichment score no less than γ . The
115 method sequentially finds ES levels l_i for which the probability $P_K(\text{ES} \geq l_i)$ is approximately
116 equal to 2^{-i} (see Fig 3a for a toy example). The method stops when l_i becomes greater than
117 γ and the P-value can be crudely approximated as 2^{-i} .

118 The intermediate l_i thresholds are calculated as follows. First, a set of Z (an odd number,
119 parameter of the method) random gene sets of size K are generated uniformly and ES values
120 for them are calculated. The median value of the ES values is calculated and assigned to l_1 .
121 By construction, the probability $P_K(\gamma \geq l_1)$ of a random gene set to have an ES value no less
122 than l_1 can be approximated as $\frac{1}{2}$. Next $\frac{Z-1}{2}$ generated gene sets with the ES values less than

123 l_1 are discarded, while $\frac{Z-1}{2}$ gene sets with the ES values greater than l_1 are duplicated. This
124 results in a sample of Z gene sets with the ES values no less than l_1 , but the distribution is
125 non-uniform. However, it can be made into a uniform sample with a Metropolis algorithm.
126 On each Metropolis algorithm step each gene set sample is tried to be modified by swapping
127 a random gene from the set with a gene outside of the set. The change is accepted if an
128 enrichment score of the new set is no less than current threshold l_1 , otherwise the change is
129 rejected. Metropolis algorithm guarantees, that after enough steps the sample becomes close
130 to uniformly distributed. Thus, a median of the enrichment scores (l_2) would correspond to
131 probability of $\frac{1}{2}$ for a gene set to have an enrichment score no less than l_2 given it has an
132 enrichment score no less than l_1 :

$$P_K(\text{ES} \geq l_2 \mid \text{ES} \geq l_1) \approx \frac{1}{2}.$$

133 Which means

$$P_K(\text{ES} \geq l_2) \approx 2^{-2}.$$

134 The same procedure is applied to calculate the next l_i values.

135 The iterations stop when l_i becomes greater than γ . On this iteration the probability of
136 a random gene set to have a ES value no less than γ can be approximated as:

$$\frac{1}{2^{i-1}} \cdot \frac{\#\{\text{samples with ES} \geq \gamma\}}{Z}.$$

137 When estimating small P-values it becomes practical to carry out the estimation in log-
138 scale. In particular, the values become practically unbiased both in median and mean sense
139 and it becomes simple to estimate the approximation error and confidence intervals (see
140 section 2.5.4).

141 The full formal description of the algorithm is available in the section 2.5.

142 For the example dataset we show that P-values are as low as 10^{-26} for some of the path-
143 ways and the results are consistent with FGSEA-simple P-values ran on 10^8 permutations
144 (Fig 3b). Note, that FGSEA-multilevel calculation with sample size of $Z=101$ took only
145 10 seconds working on a single thread while 10^8 permutations on FGSEA-simple took 40
146 minutes working in 32 threads.

147 To further prove the approximation quality of FGSEA-multilevel algorithm we developed
148 an exact method for calculating GSEA P-values, but limited to integer weights. The method
149 is based on dynamic programming, the full description is given in section 2.4. The complexity
150 of the algorithm is $O(NKT^2)$, where N is the number of genes, K is the size of gene sets
151 and T is the sum of the top K absolute values of gene-level statistics. With a number of
152 optimizations this method allows to calculate P-values for rounded weights in the example
153 dataset in a couple of hours.

154 When run on the same integer weights FGSEA-multilevel and the exact method give
155 highly concordant results (Fig 3c). Additionally, using the exact P-values, real approximation
156 errors can be compared with the estimated ones. We show, that the FGSEA-multilevel error
157 estimation are highly concordant with the real errors (Fig 3d) for a wide range of P-values
158 (from 10^{-4} to 10^{-100}), gene set sizes (from 15 to 250) and sample sizes (from 101 to 1001).

159 In practice FGSEA-multilevel method is combined with FGSEA-simple. First, for all the
160 input pathways FGSEA-simple method can be run with a limited sample size. Next, for the
161 pathways that have high relative error after FGSEA-simple (i.e. pathways with low p-values)
162 FGSEA-multilevel method is executed. As many of the pathways in an input collection
163 usually are not enriched, they have a relatively high P-value and will be batch-processed
164 with a highly efficient FGSEA-simple algorithm with deterministic time boundaries. The
165 more interesting pathways with lower P-values will then be processed with FGSEA-multilevel
166 algorithm individually and the amount of processing time will depend on their P-values.

167 As FGSEA allows to practically estimate the P-values for a large collections of gene sets, it

168 can lead to a large number of statistically significant hits with high overlaps. To deal with this
169 issue and make the representation of FGSEA results more concise we developed a procedure
170 to filter the redundant gene sets. The procedure is similar to GO Trimming method [7]
171 but is based on the Bayesian network construction approaches. It considers the significant
172 pathways one by one and tries to remove gene sets that do not provide new information given
173 some other pathway already present in the output. In this case, we consider a pathway P_1
174 to give a new information given a pathway P_2 if the P-value of pathway P_1 in the universe
175 of genes just from P_2 , or just from genes outside of P_2 , is less than some threshold. This
176 procedure allows to filter redundant pathways without requirement of having any explicit
177 hierarchy of pathways. The full description of the procedure is given in section 2.6. The
178 table resulting from running FGSEA on the example dataset with filtering of redundant hits
179 is shown on Fig 3e.

180 Finally, we have explored FGSEA performance on the collection of 605 curated GEO
181 datasets described earlier and the C5_BP pathway collection. Notably, it took less than a
182 minute per dataset of running time to finish FGSEA analysis with the multilevel algorithm
183 (Fig 4a) on a laptop with 4-core Intel Core i5 processor, with a median time of 8 seconds.
184 Besides the multilevel algorithm, the same analysis was carried out with FGSEA-simple
185 algorithm with the sampling depth values of $n = 10^3$, 10^4 and 10^5 , and the reference imple-
186 mentation (Broad GSEA) with the sampling depth of 10^4 . For all these methods we compared
187 the number of pathways reaching FDR level of 0.01: BH-adjusted P-values were used for
188 FGSEA and reported Q-values (*Broad Q-values*) were used for the reference implementation
189 (Fig 4b). The results reiterate that even sampling depth of 10^5 is not enough to detect statis-
190 tically significant enriched pathways for some of the datasets when BH-adjustment procedure
191 is used.

192 The *ad hoc* procedure implemented in Broad GSEA aggregates ES values generated
193 across different pathways increasing the sensitivity on some of the datasets compared to BH-

194 adjusted P-values for the same sampling depth of 10^4 , however this increase in sensitivity
195 comes with overall more conservative behavior. The total number of pathways reaching
196 FDR level of 0.01 for Broad Q-values is 39467, which is only 60% of 65690 pathways for
197 BH-adjustment P-values with the sampling depth of 10^4 and 48% of 81628 pathways for the
198 multilevel algorithm. We further characterized this behavior by directly comparing Broad
199 Q-values with BH-adjusted P-values and have shown that Broad Q-values are individually
200 more conservative (Fig 4c) in a pathway size dependent manner (Fig 4d).

201 To conclude, here we have presented FGSEA method for fast preranked gene set enrich-
202 ment analysis. The method allows to routinely estimate even very low P-values and can be
203 used with conjunction with standard multiple hypothesis testing correction methods, such as
204 Benjamini-Hochberg procedure. This, in turn, leads to better sensitivity and the ability to
205 detect significant pathways in hard cases, where other implementations fail. FGSEA method
206 is freely available as an R package at Bioconductor (<http://bioconductor.org/packages/fgsea>)
207 and on GitHub (<https://github.com/ctlab/fgsea>).

208 2 Methods

209 2.1 Formal definitions

210 The preranked gene set enrichment analysis takes as input two objects: an array of gene-level
211 statistic values S for the genes $U = \{1, 2, \dots, N\}$ and a list of query gene sets (pathways)
212 P . The goal of the analysis is to determine which of the gene sets from P has a non-random
213 behavior.

214 The statistic array S of the size $|S| = N$ for each gene $i \in U$ contains a value $S_i \in \mathbb{R}$
215 that characterizes the gene behavior in a considered biological process. Commonly, if $S_i > 0$
216 the expression of gene i goes up on treatment compared to control and $S_i < 0$ means that
217 the expression goes down. Absolute values $|S_i|$ represent magnitude of the change. Array
218 S is sorted in a decreasing order: $S_i > S_j$ for $i < j$. The value of N in practice is about
219 10000–20000.

220 The list of gene sets $P = \{P_1, P_2, \dots, P_M\}$ of length M usually contains groups of genes
221 that are commonly regulated in some biological process. We assume that the gene sets P_i are
222 ordered by their size (denoted as K_i): $K_1 \leq K_2 \leq \dots \leq K_M = K$. Usually only relatively
223 small gene sets are considered with $K \approx 500$ genes.

224 To quantify a co-regulation of genes in a gene set p Subramanian *et al.*[1] introduced a
225 gene set enrichment score function $s_r(p)$ that uses gene rankings (values of S). The more
226 positive is the value of $s_r(p)$ the more enriched the gene set is in the positively-regulated
227 genes (with $S_i > 0$). Accordingly, negative $s_r(p)$ corresponds to enrichment in the negatively
228 regulated genes.

229 Value of $s_r(p)$ can be calculated as follows. Let $k = |p|$, $\text{NS} = \sum_{i \in p} |S_i|$. Let also ES be

230 an array specified by the following formula:

$$ES_i = \begin{cases} 0 & \text{if } i = 0, \\ ES_{i-1} + \frac{1}{NS} |S_i| & \text{if } 1 \leq i \leq N \text{ and } i \in p, \\ ES_{i-1} - \frac{1}{N-k} & \text{if } 1 \leq i \leq N \text{ and } i \notin p. \end{cases}$$

231 The value of $s_r(p)$ corresponds to the largest by the absolute value entry of ES:

$$s_r(p) = ES_{i^*}, \text{ where } i^* = \arg \max_i |ES_i|.$$

For convenience, we also introduce the following notation:

$$s_r^+(p) = ES_{i^+}, \quad i^+ = \arg \max_i ES_i,$$

$$s_r^-(p) = ES_{i^-}, \quad i^- = \arg \min_i ES_i.$$

232 From these two values it easy to find value of $s_r(p)$, which is equal to $s_r^+(p)$ if $|s_r^+(p)| > |s_r^-(p)|$
 233 or $s_r^-(p)$ otherwise.

234 Often we will consider only the positive values of the gene set enrichment score function
 235 since:

$$\forall p \in P, \quad s_r(p) = -s'_r(p'),$$

236 where $p' = (p'_1, p'_2, \dots, p'_k) = (N - p_1 + 1, N - p_2 + 1, \dots, N - p_k + 1)$ and s'_r corresponds to
 237 the gene set enrichment score function for array S' such that $S'_i = S_{N-i+1}$.

238 Next, following Subramanian *et al* for a pathway p we define GSEA P-value as:

$$P_{value}(p) = \begin{cases} \frac{P_k(s_r(q) \geq s_r(p))}{P_k(s_r(q) \geq 0)} & \text{if } s_r(p) > 0, \\ \frac{P_k(s_r(q) \leq s_r(p))}{P_k(s_r(q) \leq 0)} & \text{if } s_r(p) < 0, \end{cases}$$

239 where q is a random gene set of size k .

240 **2.2 The datasets**

241 A collection of 605 curated datasets was generated from a set of all curated datasets (GDS)
242 in Gene Expression Omnibus. Only the datasets with two biological conditions were kept.
243 Differential expression was done using limma. The moderated t-statistic was used for gene
244 ranking. Only top 10000 genes by average expression were used in ranking. The final
245 rankings are available at https://ctlab.itmo.ru/files/software/fgsea/geo_ranks/.

246 Pathway collections C2_REACTOME, C5_BP, C2, C5, MSigDB were obtained via
247 msigdb package. LINCS perturbation collection was downloaded from Enrichr web-site. The
248 corresponding gmt files are available at <https://ctlab.itmo.ru/files/software/fgsea/gmts/>.

249 As the example ranking Th0 vs Th1 comparison was used from dataset GSE14308 [4].
250 The differential expression was calculated using limma [8]. Only top 12000 genes by mean
251 expression were used. Limma t-statistic was used as gene-level statistic. The script to gen-
252 erate rankings is available on GitHub: [https://github.com/ctlab/fgsea/blob/master/inst/](https://github.com/ctlab/fgsea/blob/master/inst/gen_gene_ranks.R)
253 [gen_gene_ranks.R](https://github.com/ctlab/fgsea/blob/master/inst/gen_gene_ranks.R).

254 Reactome [5] database was used as an example collection via reactome.db R package.
255 For the analysis only the pathways of the size from 15 to 500 were used. The script to
256 generate pathway collection is available on GitHub: [https://github.com/ctlab/fgsea/blob/](https://github.com/ctlab/fgsea/blob/master/inst/gene_reactome_pathways.R)
257 [master/inst/gene_reactome_pathways.R](https://github.com/ctlab/fgsea/blob/master/inst/gene_reactome_pathways.R)

258 **2.3 FGSEA-simple: an algorithm for fast calculation of GSEA P-** 259 **values simultaneously for many pathways**

260 In this section we describe an algorithm for fast estimation of GSEA P-values simultaneously
261 for a collection of pathways P . There, for each pathway p a set of n uniformly random gene

262 sets q_i are considered. Then P-value is estimated as:

$$\frac{\#\{q_i \mid s_r(q_i) \geq s_r(p)\} + 1}{\#\{q_i \mid s_r(q_i) \geq 0\} + 1}$$

263 for positively enriched pathway p and as:

$$\frac{\#\{q_i \mid s_r(q_i) \leq s_r(p)\} + 1}{\#\{q_i \mid s_r(q_i) \leq 0\} + 1}$$

264 for negatively enriched pathway. These two formulas follow Subramanian *et al.* implementa-
265 tion, except of +1 terms, which are recommended by Phipson and Smyth [9]. Otherwise, the
266 nominal P-values from FGSEA-simple and reference implementation are indistinguishable,
267 however FGSEA-simple works orders of magnitude faster.

268 2.3.1 Cumulative statistic calculation for the mean statistic

269 Let first describe the idea of the proposed algorithm on a simple mean statistic s_m :

$$s_m(p) = \frac{1}{|p|} \sum_{i \in p} S_i.$$

270 The main idea of the algorithm is to reuse sampling for different query gene sets. This
271 can be done due to the fact that for an estimation of null distributions samples have to be
272 independent only for a specific gene set size, while they can be dependent between different
273 sizes.

274 Instead of generating nM independent random gene sets: n for each of M input gene
275 sets, we will generate only n random gene sets of size K . Let π_i be an i -th random gene set
276 of size K . From that gene set we can generate gene sets for a all the query pathways P_j by
277 using its prefix: $\pi_{i,j} = \pi_i[1..K_j]$.

278 The next step is to calculate the enrichment scores for all gene sets $\pi_{i,j}$. Instead of

279 calculating enrichment scores separately for each gene set we will calculate simultaneously
280 scores for all $\pi_{i,j}$ for a fixed i . Using a simple procedure it can be done in $\Theta(K)$ time.

281 Let us find enrichment scores for all prefixes of π_i . This can be done by element-wise
282 dividing of cumulative sums array by the length of the corresponding prefix:

$$s_m(\pi_i[1..k]) = \frac{1}{k} \sum_{i \in \pi_i[1..k]} S_i.$$

283 Selecting only the required prefixes takes an additional $\Theta(m)$ time.

284 The described procedure allows to find P-values for all query gene sets in $\Theta(n(K + m))$
285 time. This is about $\min(K, m)$ times faster than the straightforward procedure.

286 **2.3.2 Cumulative statistic calculation for enrichment score**

287 For the enrichment score S_r we use the similar idea as above: we will also be sampling only
288 gene sets of size K and from that sample will calculate statistic values for all the other sizes.
289 However, calculation of the cumulative statistic values for the subsamples is more complex
290 in this case. In this section we only be considering the positive mode of enrichment statistic
291 s_r^+ .

It is helpful to look at enrichment score from a geometric point of view. Let us consider
for a pathway p of size $|p| = k$ a graph of $N + 1$ points (Fig. S1) with the coordinates (x_i, y_i)
for $0 \leq i \leq N$ such that:

$$(x_0, y_0) = (0, 0), \tag{1}$$

$$x_i = x_{i-1} + [i \notin p], \quad \forall i \in 1..N, \tag{2}$$

$$y_i = y_{i-1} + [i \in p] \cdot |S_i| \quad \forall i \in 1..N. \tag{3}$$

292 The calculation of s_r^+ corresponds to finding the point farthest up from a diagonal

293 $((x_0, y_0), (x_N, y_N))$. Indeed, it is easy to see that $x_N = N - |p| = N - k$ and $y_N = \sum_{j \in p} |S_j| =$
 294 NS, while the individual enrichment scores ES_i can be calculated as $ES_i = \frac{1}{NS}y_i - \frac{1}{N-k}x_i$.
 295 Value of ES_i is proportional to the directed distance from the line going through (x_0, y_0) and
 296 (x_N, y_N) to the point (x_i, y_i) .

297 Let us fix a sample π of size K . To efficiently calculate cumulative values $s_r^+(\pi[1..k])$ for
 298 all $k \leq K$ we need a fast method of updating the farthest point when a new gene is added.
 299 In that case we can add genes from π one by one and calculate values $s_r^+(\pi[1..k])$ from the
 300 corresponding maximal distances.

Because we are calculating values for $\pi[1..k]$ for $k \leq K$ we know in advance which K
 genes will be added. This allows us to consider $K + 1$ points instead of $N + 1$ for each
 iteration k . Let array o of size K contain the sorted order of genes in π : that is, π_{o_1} is the
 minimal among π , π_{o_2} is the second minimal and so on. The coordinates can be calculated
 as follows:

$$(x_0^k, y_0^k) = (0, 0), \tag{4}$$

$$x_i^k = x_{i-1}^k + \pi_{o_i} - \pi_{o_{i-1}} - [o_i \leq k], \quad \forall i \in 1..K, \tag{5}$$

$$y_i^k = y_{i-1}^k + [o_i \leq k] \cdot |S_{o_i}|, \quad \forall i \in 1..K, \tag{6}$$

301 where we set π_{o_0} to be zero.

It can be shown that finding the farthest up point among (4)–(6) is equivalent to finding
 the farthest up point among (1)–(3) with (x_i^k, y_i^k) being equal to $(x_{\pi_{o_i}}, y_{\pi_{o_i}})$ calculated for

$p = \pi[1..k]$. Consider $x_{\pi_{o_i}} - x_{\pi_{o_{i-1}}}$. By the definition of x it is equal to:

$$x_{\pi_{o_i}} - x_{\pi_{o_{i-1}}} = \sum_{j=1}^{\pi_{o_i}} [j \notin \pi[1..k]] - \sum_{j=1}^{\pi_{o_{i-1}}} [j \notin \pi[1..k]] = \sum_{j=\pi_{o_{i-1}}+1}^{\pi_{o_i}} [j \notin \pi[1..k]] = \pi_{o_i} - \pi_{o_{i-1}} - \sum_{j=\pi_{o_{i-1}}+1}^{\pi_{o_i}} [j \in \pi[1..k]].$$

302 By the definition of o , in the interval $[\pi_{o_{i-1}} + 1, \pi_{o_i} - 1]$ there are no genes from π and,
 303 thus, from $\pi[1..k]$. Thus we can replace the sum with its last member:

$$x_{\pi_{o_i}} - x_{\pi_{o_{i-1}}} = \pi_{o_i} - \pi_{o_{i-1}} - [\pi_{o_i} \in \pi[1..k]] = \pi_{o_i} - \pi_{o_{i-1}} - [o_i \leq k].$$

304 We got the same difference as in (5).

305 Now consider $y_{\pi_{o_i}} - y_{\pi_{o_{i-1}}}$. By the definition of y it is equal to:

$$y_{\pi_{o_i}} - y_{\pi_{o_{i-1}}} = \sum_{j=1}^{\pi_{o_i}} [j \in \pi[1..k]] \cdot |S_j| - \sum_{j=1}^{\pi_{o_{i-1}}} [j \in \pi[1..k]] \cdot |S_j| = \sum_{j=\pi_{o_{i-1}}+1}^{\pi_{o_i}} [j \in \pi[1..k]] \cdot |S_j|.$$

306 Again, in the interval $[\pi_{o_{i-1}} + 1, \pi_{o_i} - 1]$ there are no genes from $\pi[1..k]$. Thus we can
 307 replace the sum with only the last member:

$$y_{\pi_{o_i}} - y_{\pi_{o_{i-1}}} = [\pi_{o_i} \in \pi[1..k]] \cdot |S_{o_i}| = [o_i \leq k] \cdot |S_{o_i}|.$$

308 We got the same difference as in (6).

309 We do not need to consider other points, because points from o_{i-1} to $o_i - 1$ have the same
 310 y coordinate and o_{i-1} is the leftmost of them. Thus, when at least one gene is added the diag-
 311 onal $((x_0, y_0), (x_N, y_N))$ is not horizontal and o_{i-1} is the farthest point among $o_{i-1}, \dots, o_i - 1$.

312 Now let consider what happens with the enrichment score graph when gene π_k is added

313 to the query set $\pi[1..k-1]$ (Fig. S2). Let r_k be a rank of gene π_k among genes π , then
314 coordinate of points (x_i, y_i) for $i < r_k$ do not change, while all (x_i, y_i) for $i \geq r_k$ are changed
315 on $(\Delta_x, \Delta_y) = (-1, |S_{\pi_k}|)$.

316 To make fast incremental updates we will decompose the problem into multiple smaller
317 ones. For simplicity we assume that $K+1$ is an exact square of an integer b . Let split $K+1$
318 points into b consecutive blocks of the size b : $\{(x_0^k, y_0^k), \dots, (x_{b-1}^k, y_{b-1}^k)\}, \{(x_b^k, y_b^k), \dots, (x_{2b-1}^k, y_{2b-1}^k)\}$
319 and so on.

320 For each of b blocks we will store and update the farthest up point from the diagonal.
321 When we know for each block its farthest point we can find the globally farthest point by a
322 simple pass in $O(b)$ time.

323 Next, we show how to update the farthest points in blocks in amortized time $O(b)$. This
324 taken together with one $O(b)$ pass will get us an algorithm to update the globally farthest
325 point in amortized $O(b)$ time.

326 Below we use $c = \lfloor r_k/b \rfloor$ as an index of a block where gene π_k belongs, where r_k is the
327 ranking of the genes from π , i.e. $r_{o_i} = i$.

328 First, we describe the procedure to update point coordinates. We will store x_i coordinates
329 using two vectors: B of size b and D of size $K+1$, such that $x_i = B_{i/b} + D_i$. When gene
330 π_k is added all x_i for $i \geq r_k$ are decremented by one. To reflect this we will decrement all
331 B_j for $j > c$ and decrement all D_i for $r_k \leq i < cb$. The update takes $O(b)$ time. After this
332 update procedure we can get value x_i in $O(1)$ time. The same procedure is applied for y
333 coordinates.

334 Second, for each block we will maintain an upper part of its convex hull. Having convex
335 hull is useful because the farthest point in block always lays on its convex hull. All blocks
336 except c have the points either not changed or shifted simultaneously on the same value. That
337 means that the lists of points on the convex hulls for these blocks remain unchanged. For
338 the block c we can reconstruct convex hull from scratch using Graham scan algorithm [10].

339 Because the points are already sorted by x coordinate, this reconstruction takes $O(b)$ time.
340 In total, it takes $O(b)$ time to update the convex hulls.

341 Third, the farthest points in blocks can be updated using the stored convex hulls. Con-
342 sider a block where the convex hull was not changed (every block except, possibly, block c).
343 Because diagonal always rotates in the same counterclockwise direction, the farthest point
344 in block on iteration k either stays the same or moves on the convex hull to the left of the
345 farthest point on the $(k - 1)$ -th iteration. Thus, for each such block we can compare current
346 farthest point with its left neighbor on the convex hull and update the point if necessary. It
347 is repeated until the next neighbor is closer to the diagonal than the current farthest point.
348 In the block c we just find the farthest point in a single pass by the points on the convex
349 hull.

350 To show that the updating the farthest points takes $O(b)$ amortized time we will use
351 potential method. Let a potential after adding k -th gene Φ_k be a sum of relative indexes
352 of the farthest points for all the blocks. As there are b blocks of size b the sum of relative
353 indexes lies between 0 and b^2 . Thus, $\Phi_k = O(b^2)$. For an update of all $b - 1$ blocks except c
354 we need to make $t_k = b - 1 + z$ operations of comparing two points, where z is the number
355 of times the farthest points were updated. This can take up to $\Theta(b^2)$ time in the worst case.
356 However, it can be noticed, that potential change $\Phi_k - \Phi_{k-1}$ is equal to $-z + O(b)$: the sum
357 of indexes is decreased by a number of times the farthest points were updated plus $O(b)$ for
358 the block c where the index can go from 0 to $b - 1$. This gives an amortized cost of k -th
359 iteration to be $a_k = t_k + \Phi_k - \Phi_{k-1} = b - 1 + z - z + O(b) = O(b)$. The total real cost of K
360 iterations is $\sum_{k=1}^K a_k + \Phi_0 - \Phi_K = O(Kb) + O(b^2) = O(Kb)$, which means amortized cost of
361 one iteration to be $O(b)$.

362 Taken together the algorithm allows to find all cumulative enrichment scores $s_r(\pi[1..k])$
363 in $O(Kb) = O(K\sqrt{K})$ time. The straightforward implementation of calculating cumulative
364 values from scratch would take $O(K^2 \log K)$ time. Thus, we have improved the performance

365 $O(K \frac{\log K}{\sqrt{K}})$ times.

366 **2.3.3 Implementation details**

367 We also implemented an optimization so that the algorithm does not build convex hull from
368 scratch for a changed block c , but only updates the changed points. This does not influence
369 the asymptotic performance, but decreases the constant factor.

370 First, we start updating the convex hull from position r_k and not from the start. To be
371 able to do this, we have an array `prev` that for each gene $g \in \pi$ stores the previous point on
372 the convex hull if g were the last gene in the block. This actually is the same as the top of
373 the stack in Graham algorithm and represent the algorithms state for any given point. As
374 all points h to the left of g are not changed `prevh` also remains unchanged and need not to
375 be recalculated.

376 Second, we stop updating the hull, when we reach the point on the previous iteration
377 convex hull. We can do this because every point to the left of g is rotated counterclockwise
378 of any point to the right of g , which means that the first point on the convex hull right of g
379 on $(k - 1)$ -th iteration remains being a convex hull point at k -th iteration.

380 **2.4 An algorithm for exact calculation of GSEA P-values for integer** 381 **gene-level statistics**

382 In this section we describe a polynomial algorithm to calculate GSEA P-value exactly, but
383 only for the case when gene-level statistics are integer numbers: $S_i \in \mathbb{Z}$. For simplicity we
384 will consider a problem of calculating the following probability:

$$P(s_r^+(q) \geq \gamma),$$

385 where q is a random gene set of size k . We also assume $\gamma > 0$.

386 Let denote the sum of k largest absolute values of gene ranks by T . The algorithm will
 387 be polynomial in terms of N , k and T .

388 2.4.1 The basic algorithm

389 Let us consider a gene set $q = \{q_1, q_2, \dots, q_k\}$. Recall the formula for $s^+(q)$:

$$s^+(q) = \text{ES}_{i^+}, \text{ where } i^+ = \arg \max_i \text{ES}_i,$$

$$\text{ES}_i = \begin{cases} 0, & \text{if } i = 0, \\ \text{ES}_{i-1} + \frac{1}{\text{NS}}|S_i|, & \text{if } 1 \leq i \leq N \text{ and } i \in q, \\ \text{ES}_{i-1} - \frac{1}{N-k}, & \text{if } 1 \leq i \leq N \text{ and } i \notin q. \end{cases}$$

390 First, let rewrite the formula for ES_i in an equivalent fashion, grouping positive and
 391 negative summands:

$$392 \text{ES}_i = \frac{1}{\text{NS}} \left(\sum_{j=1}^i [j \in q] |S_j| \right) - \frac{1}{N-k} \left(i - \sum_{j=1}^i [j \in q] \right).$$

393 Then for calculating ES_i the following values are sufficient:

- 394 • i : the index of the current gene;
- 395 • $c = \sum_{j=1}^i [j \in q]$: the number of genes included into the set q among genes 1.. i ;
- 396 • $s = \sum_{j=1}^i [j \in q] |S_j|$: the sum of the absolute values of gene-level statistics for genes
 397 included in the set among genes 1.. i ;
- 398 • $\text{NS} = \sum_{j=1}^N [j \in q] |S_j|$: the sum of the absolute values of gene-level statistics for *all* genes
 399 in the set.

400 Knowing the values above, ES_i can be calculated as $\text{ES}_i = \frac{s}{\text{NS}} - \frac{i-c}{N-k}$.

401 Notice that NS can take only integer values from 0 to T (for a set of genes with the
 402 largest absolute values of gene-level statistics). Let us split the desired probability to a sum
 403 of independent probabilities based on the value of NS:

$$404 \quad P(s^+(q) \geq \gamma) = \sum_{\text{NS}=0}^T P \left(\left(\sum_{j \in q} |S_j| = \text{NS} \right) \wedge (s^+(q) \geq \gamma) \right).$$

405 Our algorithm will be based on *dynamic programming*. For each possible value of NS
 406 we will process the genes one by one in increasing order of index and calculate an array
 407 $f_{\text{NS}}(i, c, s)$. The value $f_{\text{NS}}(i, c, s)$ will contain the probability for a uniformly random gene
 408 set q' of c genes selected from genes 1.. i to simultaneously have the following two properties:

- 409 1. the sum of the absolute values of gene-level statistics of genes from q' is equal to s ;
- 410 2. $\text{ES}_j < \gamma$ holds for all $j \leq i$, where the values of ES are calculated for the gene set q'
 411 but using the selected values of NS and k , not the ones calculated for the set q' .

412 Suppose that we have calculated all values of $f_{\text{NS}}(i, c, s)$, then

$$413 \quad P \left(\left(\sum_{j \in q} |S_j| = \text{NS} \right) \wedge (s^+(q) < \gamma) \right) = f_{\text{NS}}(N, k, \text{NS})$$

414 and

$$415 \quad P(s^+(q) < \gamma) = \sum_{\text{NS}=0}^T f_{\text{NS}}(N, k, \text{NS}).$$

416 Finally, the sought probability is equal to:

$$417 \quad P(s^+(q) \geq \gamma) = 1 - P(s^+(q) < \gamma) = 1 - \sum_{\text{NS}=0}^T f_{\text{NS}}(N, k, \text{NS}).$$

418 Let us find a formula for $f_{\text{NS}}(i, c, s)$. The base case of dynamic programming is $i = 0$ for
 419 all NS:

$$f_{\text{NS}}(0, c, s) = \begin{cases} 1, & \text{if } c = s = 0, \\ 0, & \text{otherwise.} \end{cases}$$

420 Suppose we want to calculate $f_{\text{NS}}(i, c, s)$ for some $i > 0$. First, calculate

$$421 \quad \text{ES}_i = \frac{s}{\text{NS}} - \frac{i-c}{N-k}$$

422 and compare it to γ . If $\text{ES}_i \geq \gamma$, then $f_{\text{NS}}(i, c, s) = 0$ by definition.

423 Otherwise, condition “ $\text{ES}_j < \gamma$ holds for all $j \leq i$ ” can be simplified to “ $\text{ES}_j < \gamma$ holds
424 for all $j \leq i - 1$ ”. This observation allows us to use values of f that have already been
425 calculated. Consider two cases:

- 426 1. Gene i does not belong to the set q' . As q' is a set of c genes chosen uniformly at random
427 from i genes, this case happens with the probability $\frac{i-c}{i}$. The conditional probability
428 that such set satisfies the two necessary properties is $f_{\text{NS}}(i-1, c, s)$. Indeed, any set of
429 size c with the sum of absolute values of gene-level statistics values equal to s , chosen
430 among genes 1.. $i-1$ and satisfying the conditions on ES, is a valid set chosen among
431 genes 1.. i . Similarly, if a set does not satisfy the condition on ES_j for some $j \leq i-1$,
432 this set should not be counted towards $f_{\text{NS}}(i, c, s)$ since obviously $j \leq i$.
- 433 2. Gene i belongs to the set. This case happens with the probability $\frac{c}{i}$. The probability
434 that this set satisfies the necessary conditions is $f_{\text{NS}}(i-1, c-1, s-S_i)$. Indeed, any
435 set of size $c-1$ with the sum of absolute values of gene-level statistics equal to $s-S_i$,
436 chosen among genes 1.. $i-1$ and satisfying the conditions on ES, can be extended with
437 gene i , thus forming a set of size c satisfying both necessary properties. Similarly, if a
438 set does not satisfy the condition on ES_j for some $j \leq i-1$, adding gene i will not fix
439 the situation.

440 Then we can calculate $f_{\text{NS}}(i, c, s)$ using the law of total probability:

$$f_{\text{NS}}(i, c, s) = \frac{i-c}{i} f_{\text{NS}}(i-1, c, s) + \frac{c}{i} f_{\text{NS}}(i-1, c-1, s-S_i),$$

441 in the case when $i > 0$ and $\text{ES}_i < \gamma$.

442 Putting all the cases together, we arrive to the final formula for $f_{NS}(i, c, s)$:

$$f_{NS}(i, c, s) = \begin{cases} 1, & \text{if } i = c = s = 0, \\ 0, & \text{if } i = 0 \text{ and either } c \neq 0 \text{ or } s \neq 0, \\ 0, & \text{if } i > 0 \text{ and } \frac{s}{NS} - \frac{i-c}{N-k} \geq s^+(p), \\ \frac{i-c}{i} f_{NS}(i-1, c, s) + \\ \frac{c}{i} f_{NS}(i-1, c-1, s-S_i), & \text{otherwise.} \end{cases}$$

443 The overall complexity of the algorithm is $O(NkT^2)$. The values of f can be evaluated
 444 sequentially in increasing order of i . It is enough to evaluate $f_{NS}(i, c, s)$ for $0 \leq i \leq N$,
 445 $0 \leq c \leq k$, and $0 \leq s \leq NS \leq T$. Each value of f can be evaluated in constant time.

446 2.4.2 Optimizations and implementation details

447 While the algorithm described above is polynomial, a number of further optimizations are
 448 required to make execution on real size inputs feasible.

449 First, let note that the following property holds: $f_{NS_2}(i, c, s) \geq f_{NS_1}(i, c, s)$ as long as
 450 $NS_2 \geq NS_1$. Indeed, ES values calculated using different values of NS are decreasing when
 451 NS is increased. That means all gene sets counted towards $f_{NS_1}(i, c, s)$ should also be counted
 452 towards $f_{NS_2}(i, c, s)$ if $NS_2 \geq NS_1$.

453 Following the observation above, instead of calculating values of $f_{NS}(i, c, s)$ we will con-
 454 sider the values $g(i, c, s, b) = f_{b+1}(i, c, s) - f_b(i, c, s)$. These values will contain the probability
 455 of a random gene set q of size k selected uniformly from genes $1..N$ to satisfy simultaneously
 456 the following three properties:

- 457 1. set q contains exactly c genes from the genes $1..i$.

- 458 2. the sum of the absolute values of gene-level statistics of the first c genes from q is equal
459 to s ;
- 460 3. $ES_j < \gamma$ holds for all $j \leq i$, where the values of ES are calculated for the gene set q
461 using $NS = b + 1$ (and for all higher values of NS);
- 462 4. $ES_j \geq \gamma$ holds for at least one $j \leq i$, where the values of ES are calculated for the gene
463 set q using $NS = b$ (and for all lower values of NS).

464 The sought probability can be calculated from values of g as follows:

$$465 \quad P(s^+(q) \geq \gamma) = 1 - P(s^+(q) < \gamma) = 1 - \sum_{NS=0}^T \sum_{b=0}^{NS} g(N, k, NS, b).$$

466 To calculate the values of g we will use the forward dynamic programming algorithm.
467 In this algorithm we expand a tree of reachable dynamic programming states, starting from
468 $g(0, 0, 0, 0)$ which is equal to 1.

469 The states will be considered by “levels” in an increasing order of i . The values $g(i +$
470 $1, c, s, b)$ from $(i + 1)$ -th level are calculated based on level i . Note, that the sum of values
471 on i -th level is always equal to 1.

472 To calculate all values from the $(i + 1)$ -th level all non-zero values from the i -th level are
473 considered sequentially. Let consider state (i, c, s, b) and let define $p = (k - c)/(N - i)$ – the
474 probability that gene $i + 1$ will be added to the set. The corresponding set $G(i, c, s, b)$ can
475 be divided into two groups.

- 476 1. The gene sets from $G(i, c, s, b)$ that do not include gene $i + 1$. These gene sets are
477 included into gene sets $G(i + 1, c, s, b)$ on the level $i + 1$. Thus the corresponding
478 probability $g(i, c, s, b) \cdot (1 - p)$ is added to the value of $g(i + 1, c, s, b)$.
- 479 2. The gene sets from $G(i, c, s, b)$ that do include gene $i + 1$. These gene sets are included
480 into $G(i + 1, c + 1, s' = s + |S_{i+1}|, b')$ where b' is an updated bound. To calculate

481 b' let note that ES_j will be greater or equal to γ iff $\frac{s'}{NS} - \frac{(i+1)-(c+1)}{N-k} \geq \gamma$ which is
482 equivalent to $NS \leq \frac{s'}{\frac{i-j}{N-k} + \gamma}$. Thus $b' = \max\left(b, \left\lfloor \frac{s'}{\frac{i-j}{N-k} + \gamma} \right\rfloor\right)$ The probability that is
483 added to $g(i+1, c+1, s', b')$ is equal to $g(i, c, s, b) \cdot p$.

484 While the asymptotic number of states remains to be $O(NkT^2)$ the forward dynamic
485 programming allows to consider only “reachable” gene stats with $g(i, c, s, b) > 0$. In practice
486 the number of reachable stats can be several orders of magnitude smaller than the total
487 states.

488 Furthermore, for the algorithm we can consider only states with $g(i, c, s, b) > \varepsilon$ to be
489 reachable for some small value of ε . If we do not consider the unreachable states we would
490 not be able to calculate the desired probability exactly. However, if we calculate the value
491 of δ as a sum of all the skipped states values, the desired probability will be calculated with
492 the absolute error no more than δ .

493 The algorithm implementation with few other optimizations is available at: [https://](https://github.com/ctlab/fgsea/blob/master/inst/exact/exact.cpp)
494 github.com/ctlab/fgsea/blob/master/inst/exact/exact.cpp.

495 **2.5 FGSEA-multilevel: an algorithm for calculation of arbitrar-** 496 **ily low P-values using adaptive multilevel split Monte Carlo** 497 **scheme**

498 In this section we describe FGSEA-multilevel algorithm that can accurately estimate GSEA
499 P-value for a pathway p of size k even when the true P-value is very small.

500 Let $\gamma = s_r(p) > 0$ be the enrichment score of the query pathway p for which we want to
501 calculate the following value:

$$\frac{P(s_r(q) \geq \gamma)}{P(s_r(q) \geq 0)},$$

502 where q is a random gene set of size k . This probability can be rewritten as follows:

$$\frac{P(s_r(q) \geq \gamma)}{P(s_r(q) \geq 0)} = \frac{P(s_r^+(q) \geq \gamma) \cdot P(s_r(q) \geq 0 \mid s_r^+(q) \geq \gamma)}{P(s_r(q) \geq 0)}.$$

503 First, we focus on determining the probability $P(s_r^+(q) \geq \gamma)$. This probability can be
504 extremely small, so using a naive sampling gives a bad estimation. We use the adaptive
505 multilevel split Monte Carlo method [6] to solve this problem.

To estimate the probability $P(s_r^+(q) \geq \gamma)$ we split the enrichment scores into levels
 $0 = l_0 < l_1 < \dots < l_t = \gamma$. Then we can define the following probabilities:

$$\begin{aligned} P(s_r^+(q) \geq l_1 \mid s_r^+(q) \geq l_0) &= \alpha_1, \\ P(s_r^+(q) \geq l_2 \mid s_r^+(q) \geq l_1) &= \alpha_2, \\ &\dots \\ P(s_r^+(q) \geq l_t \mid s_r^+(q) \geq l_{t-1}) &= \alpha_t. \end{aligned}$$

506 Now the probability $P(s_r^+(q) \geq \gamma)$ can be rewritten as $\prod_{i=1}^t \alpha_i$.

507 To estimate α_i we can draw a sample $\{q_1^i, q_2^i, \dots, q_Z^i\}$ of size Z from a conditional distri-
508 bution $P(\cdot \mid s_r^+(q) \geq l_{i-1})$. Then

$$\alpha_i \approx \frac{Z_i}{Z},$$

509 where Z_i is the number of elements in the set $\{q_j^i \mid s_r^+(q_j^i) \geq l_i\}$.

510 Below we show how levels l_i can be chosen and how to sample from the corresponding
511 conditional distributions.

512 2.5.1 Choosing the enrichment score levels

513 We propose to chose value for a level l_i as a median of the enrichment scores for the q_j^i
514 sample. For simplicity Z is required to be an odd number.

515 Then the procedure for estimating probability $P(s_r^+(q) \geq \gamma)$ consists of repetition of the

516 following steps:

- 517 1. On iteration $i \geq 1$ sample Z gene sets q_j^i of size k from the distribution $P(\cdot \mid s_r^+(q) \geq l_{i-1})$.
- 518 2. Set the level \tilde{l}_i to be equal to the median of values $s_r^+(q_j^i)$.
- 519 3. If $\tilde{l}_i > \gamma$ then stop the iterations and set $l_i = \gamma$ and $t = i$, otherwise set $l_i = \tilde{l}_i$.

520 As a result, by construction, $\alpha_i \approx 1/2$ for $1 \leq i \leq t - 1$. The value of α_t can be
 521 approximated as Z_t/Z (which is always $\geq 1/2$). Together we get the following expression for
 522 estimating the desired probability:

$$P(s_r^+(q) \geq \gamma) \approx 2^{-(t-1)} \cdot \frac{Z_t}{Z}.$$

523 2.5.2 The conditional sampling implementation

524 To generate a uniform sample q_j^i from the conditional distribution $P(\cdot \mid s_r^+(q) \geq l_{i-1})$ we
 525 use the Metropolis algorithm.

526 First, we generate a sample $q_1^1, q_2^1, \dots, q_Z^1$ of size Z from the distribution $P(\cdot \mid s_r^+(q) \geq l_0)$.
 527 Since $l_0 = 0$ and values of s_r^+ are always non-negative it can be done by generating a uniformly
 528 random subset of size k from the genes $\{1, 2, \dots, N\}$.

529 Now let consider a sample $q_1^{i-1}, q_2^{i-1}, \dots, q_Z^{i-1} \sim P(\cdot \mid s_r^+(q) \geq l_{i-2})$ at a step $i > 1$. The
 530 sample can be sorted in an increasing order of enrichment score values: $s_r^+(q_{(1)}^{i-1}) \leq s_r^+(q_{(2)}^{i-1}) \leq$
 531 $\dots \leq s_r^+(q_{(Z)}^{i-1})$. Let $d = \lceil Z/2 \rceil$. The level l_{i-1} is the median of the values $s_r^+(q_j^{i-1})$ and, thus,
 532 is equal to $l_{i-1} = s_r^+(q_{(d)}^{i-1})$.

533 Let first populate q_j^i in the following way:

$$q_j^i = \begin{cases} q_{(Z+1-j)}^{i-1}, & \text{if } j < d, \\ q_{(j)}^{i-1}, & \text{otherwise.} \end{cases}$$

534 This gives us a sample from the conditional distribution $P(\cdot \mid s_r^+(q) \geq l_{i-1})$, however it is

535 not uniform.

536 To make the sample uniform we apply a number of the Metropolis algorithm iterations.

537 On each iteration for each gene set q_j^i we apply the following steps:

538 1. Choose a random gene $g \in q_j^i$.

539 2. Choose a random gene $g' \notin q_j^i$.

540 3. Consider $\tilde{q}_j^i = q_j^i \setminus \{g\} \cup \{g'\}$. If $s_r^+(\tilde{q}_j^i) \geq l_{i-1}$ then we replace q_j^i with \tilde{q}_j^i .

541 The iterations are repeated until the total number of successful replacements becomes
542 greater or equal to $k \cdot Z$. In practice, this number of steps is enough to get a sufficiently
543 uniform sample to obtain a good estimation of probability, without a significant increase in
544 the running time of the algorithm.

545 2.5.3 Estimating the P-value

546 In order to estimate the desired P-value we also need to calculate the probabilities $P(s_r(q) \geq 0)$

547 and $P(s_r(q) \geq 0 \mid s_r^+(q) \geq \gamma)$.

548 To calculate the probability $P(s_r(q) \geq 0)$ we generate gene sets $q_1, q_2, \dots, q_{Z'}$, where
549 each sample q_i is selected uniformly at random from all the subsets of size k from the
550 set $\{1, 2, \dots, N\}$. The samples are generated until the number of samples q_i with $s_r(q_i) \geq 0$
551 becomes equal to Z . Then the probability $P(s_r(q) \geq 0)$ is estimated as follows

$$P(s_r(q) \geq 0) \approx \frac{Z}{Z'}$$

552 To determine the remaining probability $P(s_r(q) \geq 0 \mid s_r^+(q) \geq \gamma)$ we calculate the number
553 of gene sets in $\{q_j^t \mid s_r^+(q_j^t) \geq \gamma\}$ with value of the enrichment score function s_r is greater
554 than zero. After that, the probability can be estimated as follows:

$$P(s_r(q) \geq 0 \mid s_r^+(q) \geq \gamma) \approx \frac{\#\{q_j^t \mid s_r^+(q_j^t) \geq \gamma \wedge s_r(q_j^t) \geq 0\}}{\#\{q_j^t \mid s_r^+(q_j^t) \geq \gamma\}}$$

555 2.5.4 Estimating log-probability

556 To properly estimate a logarithm of the desired probability let note that the j -th order statis-
557 tic of a standard uniform sample of size Z is a random variable from the beta distribution
558 Beta($j, Z + 1 - j$). Therefore, we can use the properties of the beta distribution and make
559 correct transition to the logarithm of probability. So for the median value of sample of odd
560 size Z we have:

$$E[\log \alpha_i] = \psi\left(\frac{Z+1}{2}\right) - \psi(Z+1) \text{ for } i \in \{1, 2, \dots, t-1\},$$

561 where ψ is digamma function. In the same way, we can calculate the expectation of the
562 logarithm α_t :

$$E[\log \alpha_t] = \psi(Z_t + 1) - \psi(Z + 1).$$

563 Then the logarithm of probability $P(s_r^+(q) \geq \gamma)$ is estimated as

$$\log P(s_r^+(q) \geq \gamma) \approx (t-1) \cdot \left(\psi\left(\frac{Z+1}{2}\right) - \psi(Z+1) \right) + \psi(Z_t + 1) - \psi(Z + 1).$$

564 Similarly, we can estimate the variance of the estimates $\text{var}[\log \alpha_i] = \psi_1\left(\frac{Z+1}{2}\right) - \psi_1(Z+1)$,
565 where ψ_1 is trigamma function. From this we can approximate a standard error of our esti-
566 mator as:

$$\text{sd} = \sqrt{t \cdot \left(\psi_1\left(\frac{Z+1}{2}\right) - \psi_1(Z+1) \right)}.$$

567 The same approach with digamma functions is used to calculate the logarithm of the
568 probabilities $P(s_r(q) \geq 0 | s_r^+(q) \geq \gamma)$ and $P(s_r(q) \geq 0)$.

569 2.5.5 Optimizations and implementation details

570 The most demanding step of FGSEA-multilevel algorithm is to check whether the newly
571 obtained gene set has the enrichments score value of at least l_i . Importantly, this does not
572 require full computation of enrichment score value. It is enough to show that there is at least
573 one gene that has the running enrichment score value of at least l_i or, in other words, that
574 at least one point on the enrichment score graph is farther away from the diagonal (Fig. S1)
575 than some value calculated from l_i .

576 Similarly to FGSEA-simple algorithm, a square root decomposition is used to split all the
577 genes into approximately \sqrt{K} blocks. The block boundaries are determined automatically on
578 each level based on the existing sample in such a way, that each block contains approximately
579 \sqrt{K} genes. In particular this decomposition enables adding and removing genes in $O(\sqrt{K})$
580 time while keeping the gene set sorted.

581 Unlike FGSEA-simple algorithm we will not maintain the convex hull but will apply a
582 few heuristics to do the required check faster.

583 First we check a "candidate" point which has a high chance to be at the required distance
584 from the diagonal. If it is so, we do not have to continue the check. The "candidate" gene
585 is carried over from the previous iterations, as the point where the successful check has been
586 interrupted.

587 Next we go block by block. At the beginning we construct a "rectangle" upper bound
588 on the enrichment score value at the block, which can be obtained by moving all the genes
589 of the block to its start. If this upper bound does not satisfy our criterion we can skip the
590 block. Otherwise, we go gene by gene and calculate the enrichment score values until it
591 reaches the required value or the end of the block is reached. In the former case the check is
592 interrupted with a successful result.

593 2.5.6 Comparison with the exact method

594 To compare FGSEA-multilevel and the exact method on the same dataset we used rounded
595 values of the gene-level statistics from the example data (section 2.2) as input data for both
596 algorithms. Both algorithms calculated the probability $P(s_r^+(q) \geq \gamma)$.

597 The results of the algorithms for the pathways from the example data are shown on
598 Fig 3c. The exact algorithm was run with $\varepsilon = 10^{-40}$, all the probabilities were obtained with
599 accuracy of at least six significant digits. For FGSEA-multilevel $Z = 101$ was used.

600 We also calculated empirical estimation errors and compared it to the theoretical ones
601 (Fig 3d). For this we generated 100 independent estimates for a range of ES values (corre-
602 sponding to P-values of 10^{-4} to 10^{-100} , gene set sizes (from 15 to 250) and sample size (from
603 101 to 1001). The raw values are available in the Supplementary Table.

604 2.6 Filtering redundant pathways

605 In this section we describe an algorithm to filter redundant pathways from the results of
606 FGSEA.

607 Let consider two pathways p_1 and p_2 that both have a significant GSEA P-value. There
608 are two situations in which we will consider p_2 to be non-redundant given p_1 :

- 609 1. If pathway p_2 is enriched even if we do not consider the genes from p_1 at all. Formally,
610 we calculate GSEA P-value for gene set $p_2 \setminus p_1$ and gene-level statistics vector $S[U \setminus p_1]$
611 for all the genes except p_1 . If the P-value is less than a pre-defined threshold, then
612 pathway p_2 is considered as non-redundant given p_1 .
- 613 2. If pathway p_2 is enriched even if we consider only genes from p_1 . Formally, we calculate
614 GSEA P-value for gene set $p_2 \cap p_1$ and gene-level statistics vector $S[p_1]$ for the genes
615 from p_1 . Again, if the P-value is less than a pre-defined threshold, then pathway p_2 is
616 considered as non-redundant given p_1 .

617 Otherwise pathway p_2 is considered to be redundant.

618 The filtering procedure starts with a set of significantly enriched pathways P_{sig} selected
619 by the user: for example the pathways with GSEA P-values less than 0.01 after Benjamini-
620 Hochberg correction, sorted by P-value. The output of the procedure is a list $P_{main} \subset P_{sig}$
621 of pathways that are pairwise non-redundant. At the same time, all the other pathways
622 $P_{red} = P_{sig} \setminus P_{main}$ are redundant given some pathway from P_{sig} .

623 The procedure itself is similar to Sieve of Eratosthenes algorithm. The pathways are
624 considered one by one and some of them are marked as redundant. For a pathway p we first
625 check if it is already marked as redundant, if yes, we go to the next pathway. Otherwise,
626 we first run FGSEA-simple algorithm on a vector of statistics $S[U \setminus p]$ and all the pathway
627 currently not marked as redundant (including the ones that already have been considered,
628 but excluding pathway p). Then, similarly, we run FGSEA-simple algorithm on a vector of
629 statistics $S[p]$. Pathways that do not achieve non-redundant P-value threshold in both tests
630 are marked as redundant.

631 References

- 632 [1] Aravind Subramanian, Pablo Tamayo, Vamsi K Mootha, Sayan Mukherjee, Benjamin L
633 Ebert, Michael A Gillette, Amanda Paulovich, Scott L Pomeroy, Todd R Golub, Eric S
634 Lander, and Jill P Mesirov. Gene set enrichment analysis: a knowledge-based approach
635 for interpreting genome-wide expression profiles. *Proceedings of the National Academy
636 of Sciences of the United States of America*, 102(43):15545–50, 2005.
- 637 [2] G. Yu, L. G. Wang, G. R. Yan, and Q. Y. He. DOSE: an R/Bioconductor package for
638 disease ontology semantic and enrichment analysis. *Bioinformatics*, 31(4):608–609, Feb
639 2015.

- 640 [3] L. Varemo, J. Nielsen, and I. Nookaew. Enriching the gene set analysis of genome-
641 wide data by incorporating directionality of gene expression and combining statistical
642 hypotheses and methods. *Nucleic Acids Res.*, 41(8):4378–4391, Apr 2013.
- 643 [4] Gang Wei, Lai Wei, Jinfang Zhu, Chongzhi Zang, Jane Hu-Li, Zhengju Yao, Kairong
644 Cui, Yuka Kanno, Tae-Young Roh, Wendy T Watford, Dustin E Schones, Weiqun Peng,
645 Hong-Wei Sun, William E Paul, John J O’Shea, and Keji Zhao. Global mapping of
646 H3K4me3 and H3K27me3 reveals specificity and plasticity in lineage fate determination
647 of differentiating CD4+ T cells. *Immunity*, 30(1):155–67, 2009.
- 648 [5] G Joshi-Tope, M Gillespie, I Vastrik, P D’Eustachio, E Schmidt, B de Bono, B Jassal,
649 G R Gopinath, G R Wu, L Matthews, S Lewis, E Birney, and L Stein. Reactome: a
650 knowledgebase of biological pathways. *Nucleic acids research*, 33(Database issue):D428–
651 32, 2005.
- 652 [6] Zdravko I Botev and Dirk P Kroese. An efficient algorithm for rare-event probability
653 estimation, combinatorial optimization, and counting. *Methodology and Computing in
654 Applied Probability*, 10(4):471–505, 2008.
- 655 [7] S. G. Jantzen, B. J. Sutherland, D. R. Minkley, and B. F. Koop. GO Trimming:
656 Systematically reducing redundancy in large Gene Ontology datasets. *BMC Res Notes*,
657 4:267, Jul 2011.
- 658 [8] M. E. Ritchie, B. Phipson, D. Wu, Y. Hu, C. W. Law, W. Shi, and G. K. Smyth. limma
659 powers differential expression analyses for RNA-sequencing and microarray studies. *Nu-
660 cleic Acids Research*, pages gkv007–, 2015.
- 661 [9] Belinda Phipson and Gordon K Smyth. Permutation p-values should never be zero: cal-
662 culating exact p-values when permutations are randomly drawn. *Statistical applications
663 in genetics and molecular biology*, 9(1), 2010.

- 664 [10] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Intro-*
665 *duction to Algorithms, Second Edition*, volume 7. 2001.

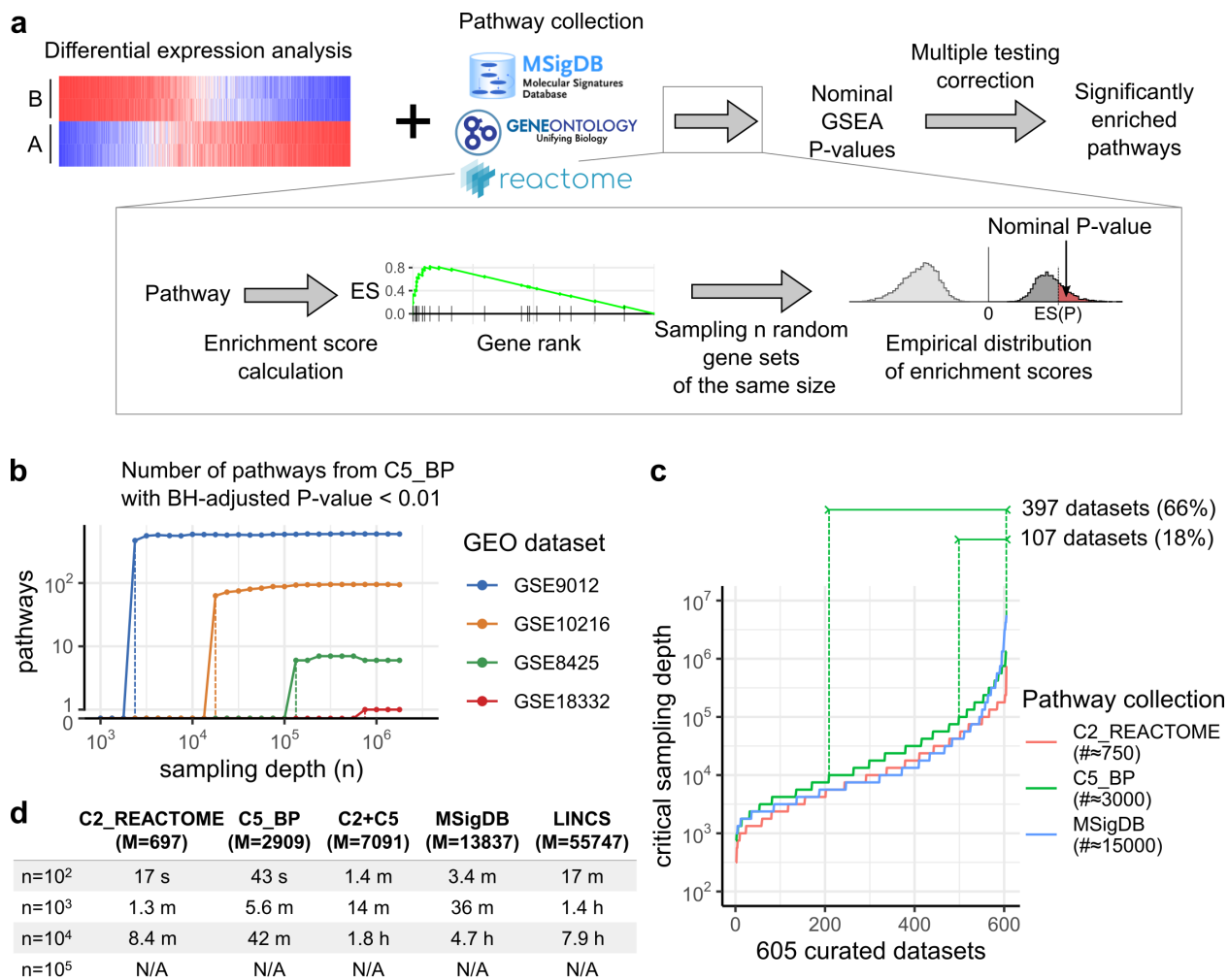


Figure 1: Gene Set Enrichment Analysis (GSEA) sensitivity depends on the ability to reach high sampling depth. **a**, Overview of preranked GSEA method. **b**, The number of pathways reaching Benjamini-Hochberg adjusted P-value threshold depends on the sampling depth in a phase transition manner. The critical sampling depth, where the transition happens, varies with the dataset. **c**, Distribution of critical sampling depth values across 605 curated datasets from Gene Expression Omnibus, as estimated for three pathway collections. **d**, Time required for the reference GSEA implementation to estimate P-values for different gene set collections and sampling depth values for dataset GSE22293. N/A values indicate out of memory errors.

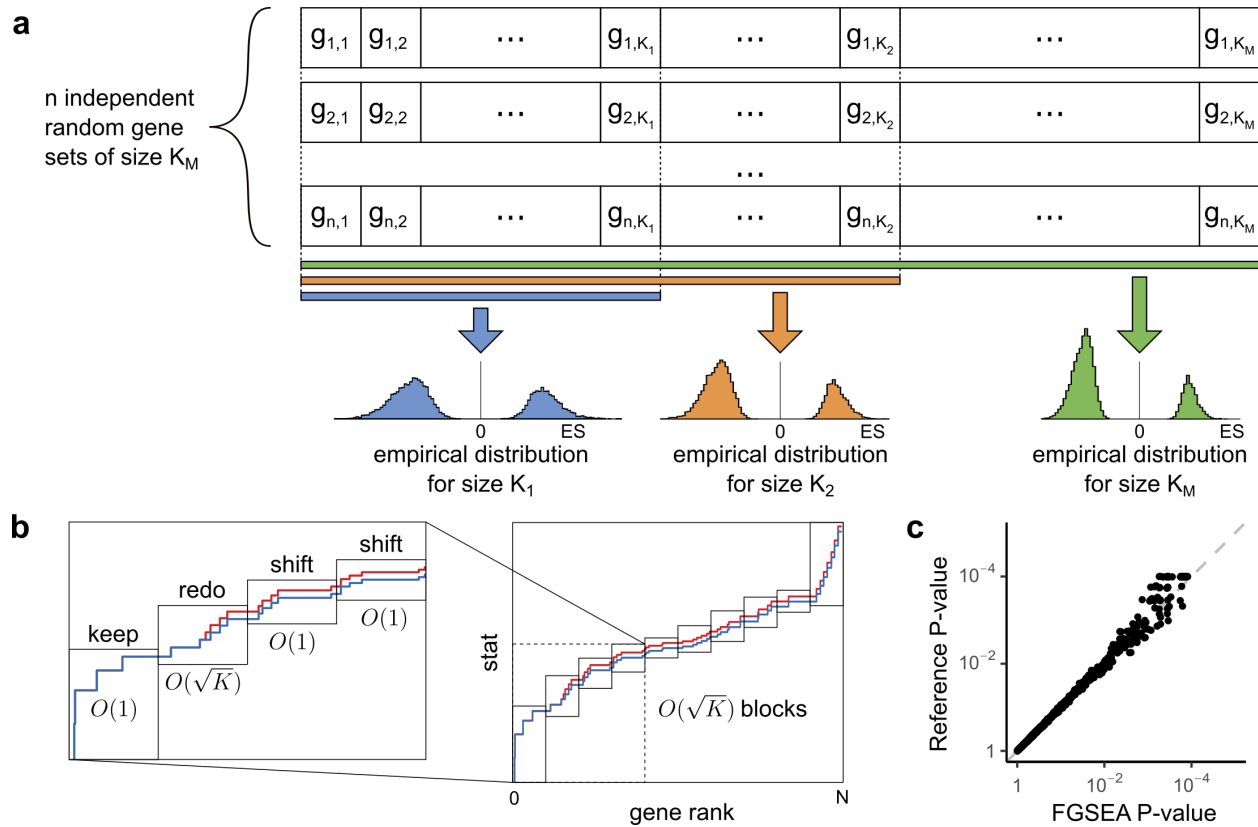


Figure 2: Preranked gene set enrichment analysis can be sped up by sharing sampling information between different gene set sizes. **a**, It is sufficient to generate n independent samples of size K_M to calculate empirical distribution for any sizes $K_j \leq K_M$ by considering only the prefix of the samples of the size K_j . **b**, For a given gene sample enrichment scores for all the prefixes can be efficiently calculated by employing a square root heuristic. **c**, The P-values calculated with the FGSEA-simple method are consistent with the reference implementation, but the results are obtained hundreds times faster.

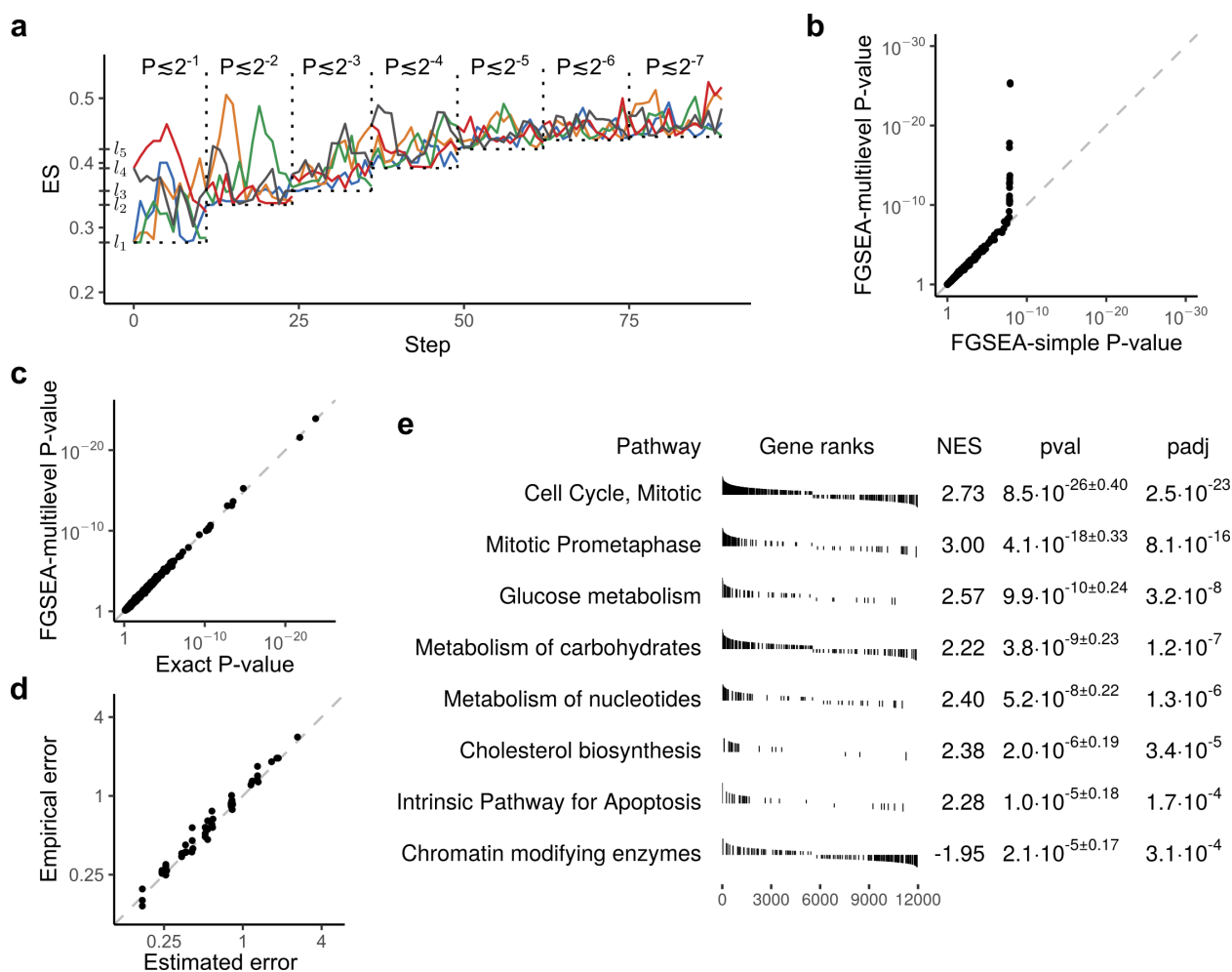


Figure 3: Adaptive multilevel Monte Carlo sampling scheme can be used to calculate arbitrarily low P-values. **a**, A toy illustration of the multilevel split Monte Carlo scheme for sample size of $Z = 5$. **b**, Comparison of GSEA P-values as calculated by FGSEA-simple method run with the sampling depth $n = 10^8$ and FGSEA-multilevel with the sample size of $Z=101$. **c**, Comparison of P-values as calculated with an exact method and FGSEA-multilevel method. Both methods were run on gene-level statistic values rounded to integers. **d**, Comparison between estimated and an observed error of \log_2 P-values for different P-values (from 10^{-4} to 10^{-100}), gene set sizes (from 15 to 250) and sample sizes (from 101 to 1001). **d**, An example of FGSEA results as run with FGSEA-multilevel method for Th0 vs Th1 comparison and Reactome pathways. The analysis was run with sample size of $Z = 101$. Redundant pathways were filtered.

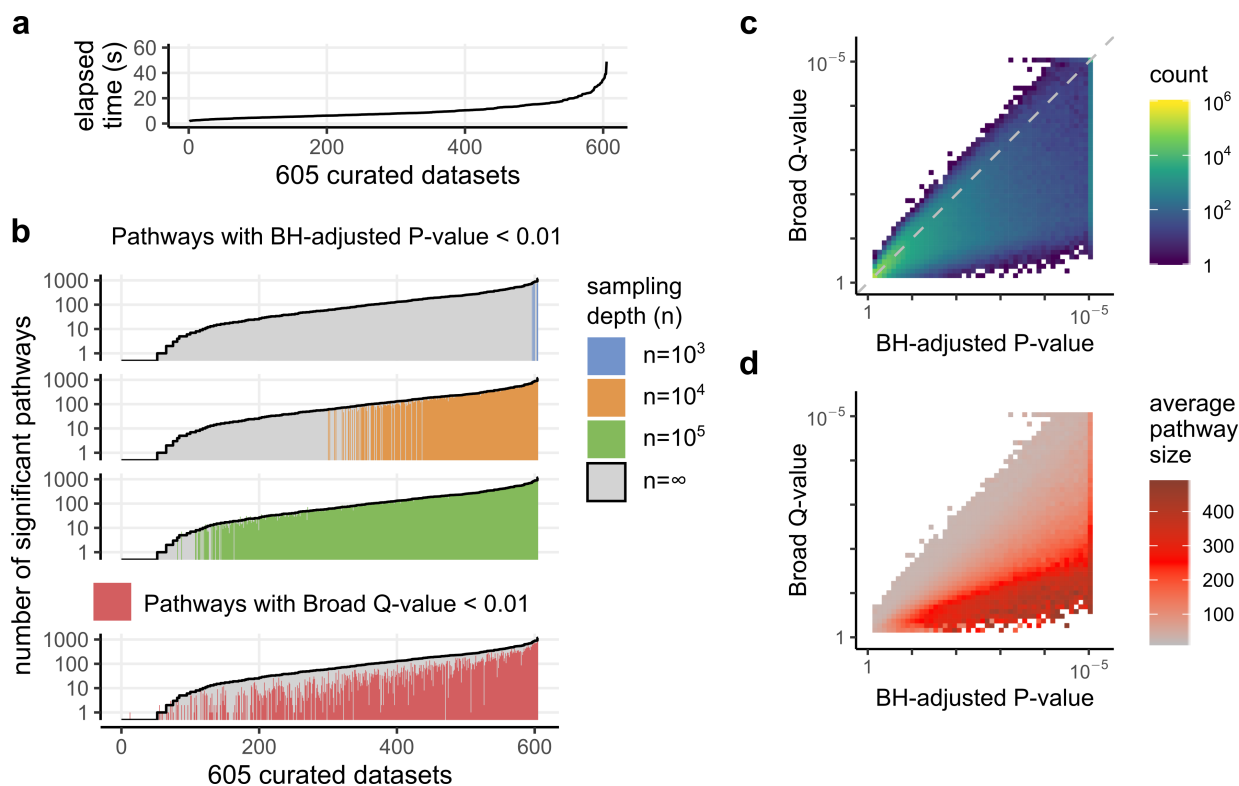


Figure 4: FGSEA method with the multilevel approach detects more statistically significant pathways compared to other GSEA implementations in a fraction of time. All plots represent data from analysis of the collection of 605 curated GEO datasets with C5_BP used as a gene set collection. **a**, Wall clock time of running FGSEA. **b**, The number of detected statistically significant pathways for different adjusted P-value calculation procedures. Pathways with small BH-adjusted P-value after FGSEA-multilevel are shown in grey. Pathways with small BH-adjusted P-value after FGSEA-simple are shown in blue, orange and green, depending on the sampling depth. Pathways with small Q-values as reported by Broad GSEA are shown in red. **c**, Comparison of BH-adjusted P-values for FGSEA-multilevel and Q-values reported by Broad GSEA. For illustration purposes all values are capped at 10⁻⁵. **d**, Same as **c** but average pathway sizes are shown.

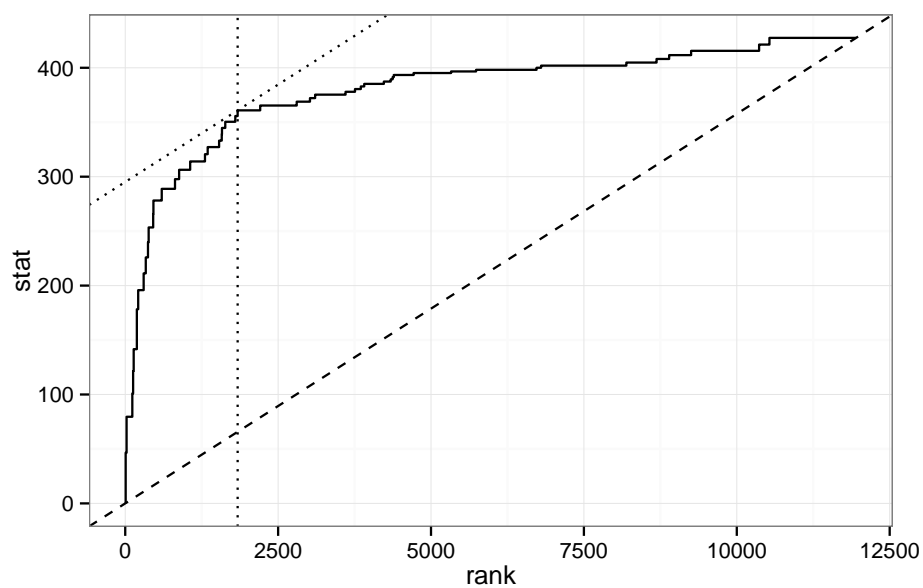


Figure S1: A graph that corresponds to a calculation of enrichment score. Each breakpoint on a graph corresponds to a gene present in the pathway. Dotted lines cross at a point which is the farthest up from a diagonal (dashed line). This point correspond to gene i^+ , where the maximal value of ES_i is reached.

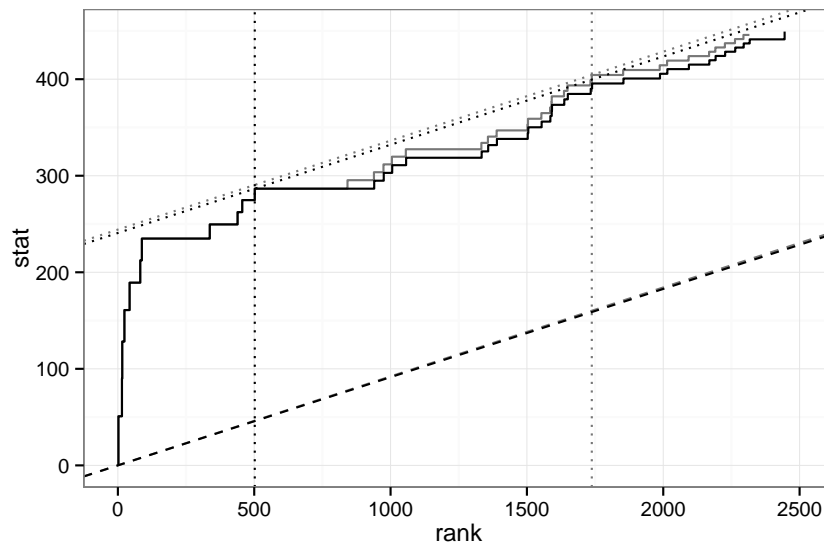


Figure S2: Update of an enrichment score graph when gene $\pi_k \approx 800$ is added. Only a fragment is shown. Black graph corresponds to a graph for gene set $\pi[1..k-1]$, gray graph corresponds to $\pi[1..k]$. A part of the graph to the left of $x = x_{r_k}$ does not change and the other part is shifted to the top-left corner. The diagonal $((x_0, y_0), (x_N, y_N))$ is rotated counterclockwise.