

Using set theory to reduce redundancy in pathway sets

Ruth Stoney^{1*}, Jean-Marc Schwartz², David L Robertson³, Goran Nenadic¹

¹School of Computer Science, University of Manchester, M13 9PT, UK

²Faculty of Biology, Medicine and Health, University of Manchester, M13 9PT, UK

³MRC-University of Glasgow Centre for Virus Research, Garscube Campus, Glasgow, G61 1QH, UK.

Ruth Stoney: ruth.stoney@manchester.ac.uk

David L Robertson: David.L.Robertson@glasgow.ac.uk

Goran Nenadic: g.nenadic@manchester.ac.uk

Jean-Marc-Schwartz: jean-marc.schwartz@manchester.ac.uk

This work has been supported by the Biotechnology and Biological Sciences Research Council DTP [BB/J014478/1].

* corresponding author

1. Abstract

1.01 Background

The consolidation of pathway databases, such as KEGG[1], Reactome[2] and ConsensusPathDB[3], has generated widespread biological interest, however the issue of pathway redundancy impedes the use of these consolidated datasets. Attempts to reduce this redundancy have focused on visualizing pathway overlap or merging pathways, but the resulting pathways may be of heterogeneous sizes and cover multiple biological functions. Efforts have also been made to deal with redundancy in pathway data by consolidating enriched pathways into a number of clusters or concepts. We present an alternative approach, which generates pathway subsets capable of covering all of genes presented within either pathway databases or enrichment results, generating substantial reductions in redundancy.

1.02 Results

We propose a method that uses set cover to reduce pathway redundancy, without merging pathways. The proposed approach considers three objectives: removal of pathway redundancy, controlling pathway size and coverage of the gene set. By applying set cover to the ConsensusPathDB dataset we were able to produce a reduced set of pathways, representing 100% of the genes in the original data set with 74% less redundancy, or 95% of the genes with 88% less redundancy. We also developed an algorithm to simplify enrichment data and applied it to a set of enriched osteoarthritis pathways, revealing that within the top ten pathways, five were redundant subsets of more enriched pathways. Applying set cover to the enrichment results removed these redundant pathways allowing more informative pathways to take their place.

1.03 Conclusion

Our method provides an alternative approach for handling pathway redundancy, while ensuring that the pathways are of homogeneous size and gene coverage is maximised. Pathways are not altered from their original form, allowing biological knowledge regarding the data set to be directly applicable. We demonstrate the ability of the algorithms to prioritise redundancy reduction, pathway size control or gene set coverage. The application of set cover to pathway enrichment results produces an optimised summary of the pathways that best represent the differentially regulated gene set.

Keywords

Set cover, data redundancy, pathways, gene enrichment analysis

40

41 2. Background

42 Pathways are sets of genes corresponding to functionally related interacting
43 proteins. Pathway data is available from many databases dependent on biological
44 focus. The fragmented nature of pathways across multiple databases makes it
45 difficult to perform inclusive analysis of all known data. To address this issue,
46 many attempts have been made to consolidate pathway databases such as
47 ConsensusPathDB (CPDB) [4], PathwayCommons [5], The Human Pathway
48 Database (HPD) [6], Pathway Interaction Database (PID) [7], and NCBI Biosystems
49 [8]. Amalgamating multiple databases into a consistent searchable format
50 facilitates the use of these resources, however the arbitrary nature of pathway
51 boundaries results in overlap and redundancy. This redundancy greatly increases
52 the quantity and complexity of pathway data, which has led to the development of
53 a range of tools to assist in data simplification and interpretation [6, 7, 9–11].
54 Previous solutions presented to deal with redundancy include visualizing
55 redundancy between pathways to the user [6], merging pathways based on
56 similarity [10, 11] and even integrating full pathway sets into a non-redundant,
57 single unified pathway [12]. Reducing redundancy simplifies the pathway-related
58 descriptive space, allowing multiple resources to be combined while limiting the
59 number of pathway attributes assigned to each gene. The advantages are apparent,
60 with resources such as PathCards being integrated into the widely used
61 GeneCards[11].

62
63 Redundancy Control in Pathway Databases (ReCiPa) [10] uses a pathway merging
64 algorithm to combine pathways with high levels of overlap. Users select a
65 maximum overlap threshold and pathway pairs displaying greater levels of
66 overlap are merged. Within that study redundancy was observed within five large
67 databases (KEGG, Biocarta, CGP, NCI-PID, and Reactome). They proceeded to
68 merge pathways from the Molecular Signatures Database (MSigDB), whose overlap
69 exceeded 75%, reducing pathway redundancy.

70
71 Pathcards described a multistep procedure to reduce pathway redundancy, also
72 through pathway merging [11]. Two thresholds were calculated and sequential
73 merging steps were used to minimize overlap, while preventing the generated
74 super-pathways from becoming too large to be informative. By merging pathways
75 into super-pathways, Pathcards suggested many new molecular interactions. They
76 demonstrated that many of these newly generated interactions are supported by
77 high numbers of literature co-mentions and high experimental interactions scores
78 according to STRING. However, while the generation of potential interactions can
79 be highly beneficial, if the aim is to utilize previously validated data, merging
80 pathways introduces a source of uncertainty into the dataset.

81
82 A major application of pathway data sets is pathway enrichment analysis. Both
83 Pathcards and ReCiPa explored the capability of their reduced pathway dataset to

84 improve enrichment results. Enrichment analysis of 830 differential expression
85 sets was performed using the super-pathways generated within Pathcards. The
86 enrichment results from super-pathways tended to be more significant than the
87 enrichment scores of their constituent pathways. Similarly within the ReCiPa study
88 enrichment analysis was performed using genes differentially expressed in
89 obesity. After merging, the top 20 most significantly enriched pathways showed
90 less overlap and greater significance towards the disease, compared to the original
91 dataset.

92
93 Pathway Distiller implemented an alternative approach by removing redundancy
94 from enriched pathway sets following enrichment analysis [9]. Pathways may be
95 consolidated into pathway concepts based on gene expression profiles, gene
96 membership, protein-protein interaction data or shared Gene Ontology (GO)
97 terms. Each method provides varying, complementary views of the data, with
98 different pathway concepts generated. Consolidating enrichment output into a
99 reduced number of pathway concepts increases data manageability and
100 readability, by organizing redundant pathways into their major groups.

101
102 All of the approaches discussed to this point have used merging and consolidation
103 to address redundancy. Alexa et al. (2006) demonstrated that redundancy in GO
104 enrichment results could be reduced by selecting a subset of representative terms
105 [13]. Pathway enrichment analysis and GO enrichment analysis are similar
106 techniques in which sets of differentially expressed genes are compared to gene
107 sets associated with pathways or GO terms. Alexa et al. (2006) introduced two
108 algorithms, *elim* and *weight*, which use the Gene Ontology topology to select a
109 representative subset of highly enriched GO terms [13]. The enrichment set cover
110 algorithm presented in this paper shares some conceptual similarity with this
111 approach however, the implementation is different since there is no organized
112 topological hierarchy for combined pathway datasets and the rules governing the
113 Gene Ontology, such as the true path rule [14], do not apply.

114
115 Within this paper we show that set cover can be used to reducing redundancy by
116 selecting subsets of representative pathways. We describe a set of algorithms for
117 reducing redundancy in pathway datasets, as well as a separate algorithm for
118 reducing redundancy from pathway enrichment results. The proportional set
119 cover algorithm and hitting set cover algorithm aim to identify a minimum subset
120 of pathways required to cover the genes in highly redundant, consolidated
121 pathway databases. The generated set covers are not designed to depict the full
122 range of possible pathway boundaries and their accompanying cellular functions,
123 but rather they provide a simplified set of pathways to represent the actions of
124 genes within the dataset. Since the pathways are not merged database and
125 biological information remains directly applicable and functional specificity is not
126 lost through pathway size expansion. The proposed method also removes the risk
127 of biologically distinct pathways being merged. The algorithm's ability to remove
128 overlap is not limited by thresholds, conferring an advantage compared to
129 approaches such as Pathcards and ReCiPa in which redundancy between pathway

130 pairs can only be removed if the overlap exceeds the threshold. Set cover
131 algorithms also consider redundancy between multiple pathways, rather than just
132 comparing pathway pairs.

133

134 We also developed the enrichment set cover algorithm for handling pathway
135 enrichment data and applied it to a set of enriched osteoarthritis pathways [15]. In
136 contrast to the approaches used by ReCiPa and Pathcards, the enrichment set
137 cover algorithm is designed to be used following enrichment analysis, which
138 should be performed using the full pathway dataset. Redundancy is then removed
139 from the enriched pathway set by selecting the pathway with the lowest p-value to
140 cover each differentially regulated gene. Enriched pathways are not merged or
141 altered and the number of enriched pathways required to cover the dataset is
142 reduced. The resulting pathways set can therefore be used as an optimized
143 summary output, conveniently showing the most important pathways for
144 describing the differentially regulated gene set. By increasing the number of
145 differentially regulated genes covered by the most highly enriched pathways,
146 researchers examining the top 10 or 20 pathways are provided with a more
147 inclusive portrayal of the gene set.

148

149 3. Approach

150 We downloaded pathway data from ConsensusPathDB (CPDB), an opensource
151 online collection of pathways, that incorporates 32 sources including KEGG,
152 Wikipathways, PDB, Reactome. CPDB makes these resources available as a single
153 download, which we acquired on 24/09/2015 containing 4,011 pathways. We
154 applied the set cover algorithm to the CPDB data set, analyzing it's effectiveness at:
155 reducing pathway overlap; reducing pathway size variability; and preserving the
156 maximum number of genes in the data set. We found that standard set cover
157 caused unacceptable increases in pathway size, therefore we modified the
158 algorithm and assessed the modified algorithms capability to meet the previous
159 three objectives.

160

161 Set cover is a well-defined algorithm in computer science for handling overlapping
162 sets of sets. For example, set cover is used by CLASS, a bioinformatics program that
163 maps RNA sequence data to transcripts [16]. Set cover has also been used to
164 predict protein-protein interactions based on binding domains [17], to reduce the
165 complexity of SNP sets [18] and to minimize the number of probes needed to
166 analyze DNA [19].

167

168 Set cover algorithms deal with elements and sets, which relate to genes and
169 pathways respectively. All the unique genes in the data set are collectively referred
170 to as the universe. The aim is to produce a reduced selection of sets (pathways),
171 which collectively cover all the elements (genes) in the universe (dataset). This

172 subset of the original data is called the cover set [20]. Each time a pathway is
173 added to the cover set the genes in the pathway become covered (Figure 1). Direct
174 application of set cover lead to extremely large, functionally non-specific pathways
175 dominating the cover set, therefore we implemented the proportional set cover
176 and hitting set cover algorithms to better control pathway size, while reducing
177 redundancy and covering the dataset.

178

179 When dealing with enrichment analysis data the aim is to reduce redundancy
180 between pathways, while preserving the order of enrichment significance denoted
181 by the p-values. We designed an algorithm that would select the set of pathways
182 with the lowest p-values capable of covering all the genes in the dataset. This
183 ensures that the filtered results return the most enriched pathways available for
184 each gene.

185

186 4. Methods

187 4.01 **Overlap score**

188 To measure overlap across different algorithms we measured the mean number of
189 pathways that each gene appears in. Within the raw data genes appeared in a
190 mean of 12.4 pathways. We refer to this metric as the overlap score.

191 4.02 **Set cover**

192 We applied the set cover algorithm to the data set, which generates a subset of
193 pathways called a cover set, in which all the genes in the data set are represented
194 or "covered". Set cover begins by first assigning values to each pathway (v_i). The
195 set cover values correspond to the number of uncovered genes each pathway
196 contains (Equation 1).

197

$$v_i = |s_i \cap \mathbf{R}|$$

198

199 where (s_i) is the pathway's gene set and \mathbf{R} is the set of all uncovered genes.

200

201 At the beginning of the algorithm all the genes in the dataset are uncovered so the
202 algorithm selects the largest pathway. The genes from the selected pathway are
203 then covered, so it is unnecessary to cover them again using additional pathways.
204 The algorithm then recalculates how many uncovered genes each pathway
205 contains and continues to add the pathway with the maximum value to the set
206 cover until all genes in the data set are covered.

207

208

209 *Algorithm 1 Set cover (in separate file)*

210 where R is the set of uncovered genes, U is all the genes in the dataset, C is the
211 covered genes, SC is the set cover result, GC is the gene coverage (see Section 4.03)
212 and s_i is a pathway.

213

214 Application of the set cover algorithm was effective in reducing overlap between
215 the pathways; however, it selected very large pathways with reduced
216 informativeness (maximum size 2320, standard deviation 160, almost double the
217 standard deviation on the original dataset 86.9). We therefore explored methods
218 that avoid preferential selection of large pathways.

219

220 4.03 Gene Set Coverage

221 As the set cover algorithm approaches completion and the final sets are added to
222 the cover set, increases in data coverage are gained at the expense of redundancy
223 reduction. This is because the final sets required to cover the few remaining genes
224 tend to have the most overlap with other pathways already in the set cover. In
225 addition, fewer pathways are available to cover the final few genes, restricting
226 options to control pathway size. To allow a user-defined compromise between the
227 gene coverage, pathway redundancy and pathway size we introduce the Gene
228 Coverage (GC) parameter. Setting GC below 100% allows the algorithm to finish
229 before the final elements have been covered. We experimented setting GC to 90,
230 95, 99 and 100% of the number of genes in the data set.

231

232 4.04 Proportional set cover

233 When reducing pathway redundancy there are three competing aims: reducing
234 redundancy; controlling pathway size; and covering the entire gene set. The
235 proportional set cover algorithm was generated to focus on controlling pathway
236 size.

237

238 To control the size of the pathways we altered the scoring mechanism to rank
239 pathways based on the proportion of uncovered genes they contained, rather than
240 the absolute number (Equation 2). This works because larger pathways are more
241 likely to have a proportion of their genes covered when other pathways are
242 selected. Additionally this mechanism directly penalizes overlap, which the
243 standard algorithm does not. At the beginning of the proportional set cover
244 algorithm none of the genes are covered so the proportion of uncovered genes in
245 every pathway is 1. This would result in the starting pathway being selected at
246 random. To ensure that pathway size variability is controlled as strictly as
247 possible, we implemented the second part of Equation 2, which ensures that
248 pathways of mean pathway size are preferentially selected when multiple
249 pathways with the same proportion of uncovered genes are available.

250

$$v_i = \frac{|s_i \cap R|}{|s_i|} + \frac{1}{\text{abs}(|s_i| - \overline{|s_i|}) * k}$$

251

252 where s_i is the pathway's gene set, $\overline{|s_i|}$ is the mean pathway length, R is the
253 uncovered genes set and k is a large constant to limit the influence of the second
254 term (taken equal to 10,000).

255

256 4.05 Hitting set cover

257 The set-covering problem can be reformulated into the equivalent set-hitting
258 problem. In this formulation genes and pathways are visualized as bi-partite graph
259 in which the pathways are connected to the genes that they contain. In this
260 depiction it is clear that some genes are only linked to a single pathway, which
261 must be selected if the gene is to be covered. The importance of pathways can
262 therefore be considered as a factor of how infrequent their genes are. The hitting
263 set cover is therefore designed to reduce redundancy as much as possible without
264 directly selecting for pathway size.

265

266 We calculated the frequency of each gene in the data set (F), then assigned the
267 gene's value $gv(j)$ as $1/F$. We then assigned a value v_i to each pathway defined as
268 the sum of each uncovered gene's scores divided by the number of genes in the
269 pathway (Equation 3).

270

$$gv(j) = 1 / F(j)$$
$$v_i = \frac{\sum_{j \in s_i \cap R} gv(j)}{|s_i|}$$

271

272 where $gv(j)$ is the value of a gene, $F(j)$ is the number of pathways a gene is in,
273 $j \in s_i \cap R$ means for each uncovered gene in the pathway and $|s_i|$ is the length of
274 the pathway.

275

276 4.06 Set cover for pathway enrichment analysis

277 Pathway analysis is a frequently used method; therefore a modified set cover
278 algorithm to address this situation could be highly useful. The universe represents
279 differentially expressed genes and the sets are enriched pathways generated
280 through enrichment analysis. Enrichment analysis results represent entirely
281 different input data compared to the pathway datasets used in the previous
282 algorithms, as the enriched pathways already have scores (p-values). We wish to
283 reduce redundancy (gene overlap) between enriched pathways and it is essential
284 that the pathways with the lowest possible p-values are selected. Equation 4
285 allows the pathways with the lowest p-values to be selected, unless all of their
286 genes are covered by other enriched pathways with even lower p-values.

287

$$s_i \cap R = \theta \rightarrow b = 0, \quad s_i \cap R \neq \theta \rightarrow b = 1 \\ v_i = (1 - pvalue_i) * b$$

288

289 where s_i is the enriched pathway's gene set, R is the uncovered gene set, b is a
290 binomial operator, $pvalue_i$ is the pathway's p-value and v_i is the pathway's set
291 cover value.

292 We generated the enriched data set by applying GOseq [21] to expression data
293 from the damaged cartilage in osteoarthritis patients and controls [15].

294

295 5. Results

296 We started with the large, extensively redundant CPDB data set and used set cover
297 to reduce pathway overlap, while controlling pathway size and seeking to cover as
298 much of the data set as possible. We describe the ability of the standard set cover
299 algorithm and two modified algorithms, in conjunction with the GC parameter, to
300 meet these objectives.

301

302 5.01 Pathway redundancy varies between different algorithms

303 The original pathway data set contained 11,196 genes and 3,305 pathways; the
304 starting overlap score (see methods) was 12.4. The standard set cover algorithm
305 reduced overall redundancy from 12.4 to 4.1, a 73% reduction (since a completely
306 discrete pathway set would have a score of 1). The overlap score for proportional
307 set cover was 4.36, slightly higher than the standard set cover algorithm, but still
308 representing a 70% reduction in overlap from the original data. The hitting set
309 cover algorithm was designed to select pathways that contained rare genes within
310 the data set, resulting in the greatest reduction in overlap (overlap score of 3.95
311 equivalent to a 74% reduction).

312

313 After application of the set cover algorithms the distribution of the remaining
314 overlap between pathways varied greatly. Figure 2 shows the Jaccard similarity
315 between pairs of pathways, in the outputs produced by each of the three
316 algorithms. The standard set cover algorithm produced the lowest maximum
317 overlap (Jaccard similarity = 0.68) between the pathway pairs. However, compared
318 to the original data, a higher proportion of pathway pairs in the set cover output
319 showed Jaccard similarities between 10-30%. Proportional set cover had the
320 greatest maximum Jaccard similarity at 0.93, out of the set cover algorithms. The
321 hitting set cover algorithm produced a maximum Jaccard similarity between two
322 pathways of 0.82, despite having the lowest overlap score.

323

324 **Gene Coverage can be lowered to reduce redundancy**

325 For each of the algorithms it is possible to use the *GC* parameter to prioritize
326 reductions in redundancy over gene coverage by stopping any algorithm before all
327 of the genes in the dataset have been covered. Figure 3 shows improved ability of
328 the set cover algorithms to reduce pathway overlap for different values of *GC*. If
329 99% of the genes are required then the hitting set algorithm achieves the lowest
330 overlap score of 3.24, equivalent to an 80% reduction in overlap. Redundancy can
331 be further reduced if only 95% of the genes are covered, with the proportional and
332 hitting set algorithms producing an overlap score of 2.41, equivalent to a 88%
333 reduction in redundancy. Both the proportional set cover and the hitting set cover
334 are more effective at reducing redundancy than the standard set cover if *GC* is set
335 to less than 100%.
336

337 **Pathway size is affected by the set cover algorithm and Gene Coverage**
338 **setting**

339 When *GC* was set to 100% the standard set cover algorithm represented all of the
340 genes in the dataset using only 524 pathways (16% of the original pathway set).
341 However, many of these were very large increasing the mean size to 87.2
342 (standard deviation 160.1). These pathways have reduced informativeness since
343 functional specificity is lost. Figure 4A illustrates the tendency of this algorithm to
344 select extremely large pathways.
345

346 The proportional set cover algorithm was designed to preferentially select
347 moderately sized pathways. This returned a cover set of 1,336 pathways with
348 controlled size variation (mean of 36.5, standard deviation 55.1) shown in Figure
349 4A. The hitting set cover algorithm was less able to control pathway size than the
350 proportional set cover algorithm, returning 957 pathways with a mean size of 46.2
351 (standard deviation 61.7).
352

353 Figures 4B – D show that as *GC* is reduced the tendency of the standard set cover to
354 select very large pathways becomes more exaggerated. Decreasing *GC* also
355 improves the ability of the proportional set cover algorithm to select moderately
356 sized pathways. The hitting set algorithm also tends to select smaller pathways
357 when *GC* is reduced, since larger pathways often contain more frequent genes.
358 Reducing *GC* affects pathway size since in the later stages of the algorithm, fewer
359 pathways are available to cover the remaining genes, reducing the available
360 options. Therefore, lowering *GC* has the ability to help control pathway size when
361 the proportional set cover and hitting set cover algorithms are used.
362

363 *Since the databases that contribute to CPDB contain pathways of different sizes, the set cover generated may*
364 *preferentially select pathways from some databases more than others.*

365 Table 1 shows the proportion of pathways that come from each database in the
 366 cover set generated by each algorithm. All algorithms generate set covers with
 367 reduced INOH and SMPDB pathways, showing that SMPBD's focus on small
 368 molecules and INOH's ontology-based approach tend to be ill-suited to the
 369 generation of discrete pathway protein sets. The standard set cover algorithm
 370 generates sets containing large pathways, preferentially selecting pathways from
 371 KEGG (median size 65, see Table 1) and Netpath (median size 51); while
 372 proportional set cover tends to select smaller pathways from Reactome (median
 373 size 17), HumanCyc (median size 5) and Signalink (median size 32), whilst
 374 avoiding NetPath.
 375

376 *Table 1. Proportion of pathways from CPDB databases. Median size represents the*
 377 *median sizes of the pathways in the CPDB dataset. CPDB % represents the proportion*
 378 *of the pathways in the unaltered dataset that came from each database. The*
 379 *following columns represent the proportion of pathways in the set cover generated*
 380 *by the standard set cover algorithm, the hitting set cover algorithm and the*
 381 *proportional set cover algorithm. Different results are obtained by altering the*
 382 *proportion of the gene set covered, shown in subcolumns below the algorithm header.*

383

	Median size	CPDB %	Standard set cover				Hitting set cover		
			100%	99%	95%	90%	100%	99%	95%
BioCarta	15.0	6.3	6.3	4.6	0.5	0.0	4.7	4.8	5.4
EHMN	32.5	1.6	3.2	3.4	2.6	1.0	2.1	2.3	1.8
HumanCyc	5.0	8.2	6.5	7.7	2.6	0.0	10.1	10.9	12.9
INOH	34.5	2.3	1.7	1.9	1.0	1.0	0.8	0.6	0.3
KEGG	65.0	7.2	29.0	30.5	37.6	40.4	15.8	15.0	13.5
NetPath	51.0	0.9	2.1	2.4	3.6	5.1	1.1	1.2	1.1
PharmGKB	13.0	2.8	3.1	2.9	0.5	0.0	2.0	2.1	2.4
PID	35.0	5.2	15.6	13.9	10.3	6.1	9.5	9.8	9.4
Reactome	17.0	39.6	4.2	5.3	10.8	21.2	36.1	35.1	34.7
Signalink	32.0	0.4	1.0	1.2	1.0	0.0	0.6	0.7	0.7
SMPDB	11.0	16.7	1.7	1.4	0.5	0.0	1.6	1.5	1.4
Wikipathways	26.0	8.8	25.6	24.9	28.9	25.3	15.6	16.0	16.2

384

385 **Reducing redundancy in pathway enrichment analysis**

386 To demonstrate the ability of the set cover algorithm to handle enrichment data,
 387 we applied the enrichment set cover algorithm to an osteoarthritis data set,
 388 retrieved from Dunn et al. (2016) [15]. From the osteoarthritis data set, 58.3% of
 389 the differentially expressed genes could be mapped to a CPDB pathway, which was

390 a 17% improvement on the G0seq [21] implemented data set. We retrieved 42
391 enriched pathways with a p-value lower than 0.05, following the Benjamini-
392 Hochberg correction for multiple testing. Set cover for enrichment analysis
393 reduced the number of pathways required to cover the differentially expressed
394 genes to 23 (supplementary table 1).

395

396 The heat map in Figure 5A shows the asymmetric overlap between the top ten
397 pathways before application of the algorithm. The p-values from pathway
398 enrichment determine the order in which pathways were considered for inclusion
399 in the cover set. Pathways were omitted if all of the differentially expressed genes
400 that they covered were also covered by more enriched pathways. Note that overlap
401 tends to be higher in the bottom left triangle as pathways added later were often
402 smaller subcomponents of larger pathways. We can see that 'extracellular matrix
403 organization', the most enriched pathway, was placed in the cover set first. Next
404 was 'collagen biosynthesis and modifying enzymes'; however, all of the
405 differentially expressed genes in this pathway are also covered by the larger
406 pathway 'extracellular matrix organization', as indicated by the red cell in the
407 'collagen biosynthesis and modifying enzymes' row, 'extracellular matrix
408 organization' column. The corresponding cell in the 'extracellular matrix
409 organization' row reveals that 24% of the differentially expressed genes in
410 'extracellular matrix organization' are also in 'collagen biosynthesis and modifying
411 enzymes'.

412

413 Figure 5B shows overlap between the top ten pathways after application of the
414 enrichment set cover algorithm. Because the differentially expressed genes
415 covered by the 'collagen biosynthesis and modifying enzymes' pathway are a
416 subset of those covered by the 'extracellular matrix organization' pathway, the
417 'collagen biosynthesis and modifying enzymes' pathway is removed from the cover
418 set (Figure 5B). The second pathway in this list therefore becomes 'GPCR signaling
419 g alpha q'. The 'collagen formation' and 'class b 2 secretin family receptors'
420 pathways are also removed because the differentially expressed genes they cover
421 are additionally covered by the more enriched pathways 'extracellular matrix
422 organization' and 'signal transduction' pathways (respectively). Additionally,
423 'GPCR signaling pertussis toxin' and 'GPCR signaling cholera toxin' are absent from
424 the returned list, as all of their differentially expressed genes are found in 'GPCR
425 signaling g alpha q' or 'signal transduction'.

426

427 Some pathways in the enrichment set cover do still show high levels of overlap, for
428 example 'wnt signalling network' is included despite 89% of its differentially
429 expressed genes being covered by 'signal transduction'. This is acceptable because
430 'signal transduction' is more highly enriched than 'wnt signalling network', yet the
431 'wnt signalling network' is worth including as it contains three differentially
432 expressed genes that are not in 'signal transduction'. If 'wnt signalling network'

433 had been excluded then these genes would not have been described by the most
434 significant pathway available to represent them. The unmodified top ten enriched
435 pathways only cover 78.0% of the enriched genes. Using the set cover enrichment
436 algorithm increases this figure to 85.2% without disrupting the pathway order
437 given by the enrichment p-values.
438

439 6. Discussion and conclusion

440 We described algorithms suitable for reducing overlap in large pathway data sets
441 allowing multiple databases to be amalgamated without excessive redundancy
442 impeding the usefulness of the resource. Standard set cover is the best algorithm
443 to reduce the number of pathways required to cover the data set, but significantly
444 increases pathway size, which can be controlled by proportional set cover or
445 hitting set cover. The proportional set cover is the best algorithm for controlling
446 pathway size and the hitting set cover is the preferred choice for covering all of the
447 genes in the dataset with minimal pathway redundancy. We showed that reducing
448 the *GC* parameter allows further reductions in pathway redundancy; for example, if
449 only 95% of the genes in the CPDB dataset were covered redundancy can be
450 reduced by up to 88%. In addition reducing *GC* increases pathway size control
451 when the proportional set cover and hitting set cover algorithms are used.
452

453 For pathway enrichment analysis we aimed to reduce redundancy while selecting
454 the most significantly enriched pathways based on p-values. As an application we
455 used the modified set cover algorithm to reduce the results of enrichment analysis
456 from a large osteoarthritis data set. We found that 5 out of the 10 top ranking
457 pathways could be omitted as they were subsets of more highly enriched
458 pathways. Overlap between pathways returned from enrichment data is not
459 always immediately obvious and requires further consideration. By reducing this
460 redundancy, data interpretation is made more intuitive. Reducing redundancy also
461 allows the user to explore substantially more of the data set using the same
462 number of pathways.
463

464 The enrichment set cover algorithm presented within this study differs from
465 existing methods implemented by ReCiPa and Pathcards, since enrichment
466 analysis is performed prior to reduction of redundancy. This is because the
467 different sets of pathway boundaries available in the full dataset may optimally fit
468 the differentially expressed genes. For example, comparison of the 'apoptosis'
469 taken from KEGG, Reactome and Wikipathways, reveals that many of the proteins
470 are specific to a single database [22]. This is due to the vague definition of pathway
471 boundaries, as well as differing experimental focus on cellular contexts, such as
472 tissues or disease states. Following enrichment analysis the pathways that are
473 most significantly enriched are selected to represent the differentially expressed

474 genes and superfluous pathways are removed. This prevents the top results from
475 being dominated by large numbers of highly similar pathways.

476

477 Set cover uses greedy heuristic methods, which provide good approximations of
478 the optimal solution in a time effective manner. These methods are extremely
479 efficient and can be run in a matter of minutes, however it should be noted that
480 they do not guarantee an optimal solution. This is particularly true for the
481 proportional set cover algorithm where the randomness of early selections
482 influences the result. However, all possible outcomes result in reduced
483 redundancy. The enrichment set cover algorithm is exempt from these
484 considerations unless multiple pathways have identical p-values.

485

486 We have provided a method to dramatically reduce redundancy in pathways
487 facilitating a more concise portrayal of cellular processes, while avoiding the issues
488 introduced by pathway merging. Our algorithms are publicly available and have
489 wide applicability to analysis of pathway datasets from any organism.

490

491 7. List of abbreviations

492 CPDB Consensus PathwayDB

493 GC Gene cover

494 SNP Single nucleotide polymorphism

495

496 8. Declarations

497 Ethics approval and consent to participate NA

498 Consent for publication NA

499 Availability of data and materials: [https://github.com/RuthStoney/set-cover-and-](https://github.com/RuthStoney/set-cover-and-set-packing-to-reduce-redundancy-in-pathway-data)
500 [set-packing-to-reduce-redundancy-in-pathway-data](https://github.com/RuthStoney/set-cover-and-set-packing-to-reduce-redundancy-in-pathway-data)

501 Competing interests: The authors declare that they have no competing interests

502 Funding: This work has been supported by the Biotechnology and Biological

503 Sciences Research Council DTP [BB/J014478/1].

504 Author contributions: All authors contributed to the design of the study. R.S.

505 performed the analysis and wrote the manuscript. All authors edited the

506 manuscript and approved its final version.

507

508 Acknowledgements: We would like to thank Jamie Soul and Sara Dunn from the
509 University of Manchester, for providing the osteoarthritis data.

510

511 9. References

512 1. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic*
513 *Acids Res.* 2000;28:27–30.

- 514 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102409&tool=pmcentrez&rendertype=abstract>.
515
516 2. Fabregat A, Sidiropoulos K, Garapati P, Gillespie M, Hausmann K, Haw R, et al.
517 The reactome pathway knowledgebase. *Nucleic Acids Res.* 2016;44:D481–7.
518 3. Kamburov A, Pentchev K, Galicka H, Wierling C, Lehrach H, Herwig R.
519 ConsensusPathDB : toward a more complete picture of cell biology. 2011;39
520 November 2010:712–7.
521 4. Kamburov A, Wierling C, Lehrach H, Herwig R. ConsensusPathDB--a database for
522 integrating human functional interaction networks. *Nucleic Acids Res.* 2009;37
523 Database issue:D623-8. doi:10.1093/nar/gkn698.
524 5. Cerami EG, Gross BE, Demir E, Rodchenkov I. Pathway Commons , a web
525 resource for biological pathway data. 2011;39 November 2010:685–90.
526 6. Chowbina SR, Wu X, Zhang F, Li PM, Pandey R, Kasamsetty HN, et al. HPD: an
527 online integrated human pathway database enabling systems biology studies. *BMC*
528 *Bioinformatics.* 2009;14 Suppl 11:S5. doi:10.1186/1471-2105-10-S11-S5.
529 7. Schaefer CF, Anthony K, Krupa S, Buchoff J, Day M, Hannay T, et al. PID: The
530 pathway interaction database. *Nucleic Acids Res.* 2009;37 SUPPL. 1:674–9.
531 8. Geer LY, Marchler-Bauer A, Geer RC, Han L, He J, He S, et al. The NCBI BioSystems
532 database. *Nucleic Acids Res.* 2009;38 SUPPL.1:492–6.
533 9. Doderer MS, Anguiano Z, Suresh U, Dashnamoorthy R, Bishop AJR, Chen Y.
534 Pathway Distiller - multisource biological pathway consolidation. *BMC Genomics.*
535 2012;13 Suppl 6 Suppl 6:S18. doi:10.1186/1471-2164-13-S6-S18.
536 10. Vivar JC, Pemu P, McPherson R, Ghosh S. Redundancy control in pathway
537 databases (ReCiPa): an application for improving gene-set enrichment analysis in
538 Omics studies and “Big data” biology. *Omics.* 2013;17:414–22.
539 doi:10.1089/omi.2012.0083.
540 11. Belinky F, Nativ N, Stelzer G, Zimmerman S, Iny Stein T, Safran M, et al.
541 PathCards: multi-source consolidation of human biological pathways. *Database.*
542 2015;2015:bav006-bav006. doi:10.1093/database/bav006.
543 12. Yu N, Seo J, Rho K, Jang Y, Park J, Kim WK, et al. hiPathDB: A human-integrated
544 pathway database with facile visualization. *Nucleic Acids Res.* 2012;40:797–802.
545 13. Alexa A, Rahnenführer J, Lengauer T. Improved scoring of functional groups
546 from gene expression data by decorrelating GO graph structure. *Bioinformatics.*
547 2006;22:1600–7.
548 14. Camon E, Magrane M, Barrell D, Lee V, Dimmer E, Maslen J, et al. The Gene
549 Ontology Annotation (GOA) Database: sharing knowledge in Uniprot with Gene
550 Ontology. *Nucleic Acids Res.* 2004;32:D262–6.
551 15. Dunn SL, Soul J, Anand S, Schwartz JM, Boot-Handford RP, Hardingham TE.
552 Gene expression changes in damaged osteoarthritic cartilage identify a signature
553 of non-chondrogenic and mechanical responses. *Osteoarthr Cartil.* 2016;24:1431–
554 40. doi:10.1016/j.joca.2016.03.007.
555 16. Song L, Florea L. CLASS: constrained transcript assembly of RNA-seq reads.
556 *BMC Bioinformatics.* 2013;14 Suppl 5 Suppl 5:S14. doi:10.1186/1471-2105-14-S5-
557 S14.
558 17. Huang C, Morcos F, Kanaan SP, Wuchty S, Chen DZ. Predicting Protein-Protein
559 Interactions from Protein Domains Using a Set Cover Approach. *Quality.*
560 2007;4:78–87.
561 18. Ao SI, Yip K, Ng M, Cheung D, Fong PY, Melhado I, et al. CLUSTAG: Hierarchical
562 clustering and graph methods for selecting tag SNPs. *Bioinformatics.*

563 2005;21:1735–6.
564 19. Borneman J, Chrobak M, Della Vedova G, Figueroa a, Jiang T. Probe selection
565 algorithms with applications in the analysis of microbial communities.
566 Bioinformatics. 2001;17 Suppl 1:S39-48.
567 <http://www.ncbi.nlm.nih.gov/pubmed/11472991>.
568 20. Kordalewski D. New Greedy Heuristics For Set Cover and Set Packing. 2013;
569 April:49. <http://arxiv.org/abs/1305.3584>.
570 21. Young MD, Wakefield MJ, Smyth GK, Oshlack A. Gene ontology analysis for
571 RNA-seq: accounting for selection bias. Genome Biol. 2010;11:R14.
572 doi:10.1186/gb-2010-11-2-r14.
573 22. Herwig R, Hardt C, Lienhard M, Herwig R, Hardt C, Lienhard M, et al. Analyzing
574 and interpreting genome data at the network level with ConsensusPathDB
575 Analyzing and interpreting genome data at the network level with
576 ConsensusPathDB. 2016; October.
577
578

579 10. Figure legends

580
581 *Figure 1. Set cover A) A simple set of overlapping sets. B) The red set with 8*
582 *uncovered elements is selected first. C) The blue set with 3 elements is selected*
583 *second. D) The orange set then covers all the elements in the universe.*

584 *Figure 2. Jaccard coefficient between pathway pairs in the cover set results produced*
585 *by each algorithm.*

586
587 *Figure 3. Redundancy in set cover outputs given different GC values.*

588
589 *Figure 4. Pathway sizes in cover set when GC is set to A) 100%, B) 99%, C) 95% and*
590 *D) 90%. The boxes indicate the 25th and 75th percentiles and the whiskers indicate*
591 *the 5th and 95th percentiles.*

592
593 *Figure 5. Pathway redundancy heat maps. (A) Pathway overlap for top ten enriched*
594 *pathways. (B) Pathway overlap for top ten enriched pathways after application of set*
595 *cover. The values represent asymmetric overlap, i.e. for each pathway shown on the*
596 *left axis, values represent the proportion of genes that are also included in the*
597 *pathway shown on the bottom axis.*

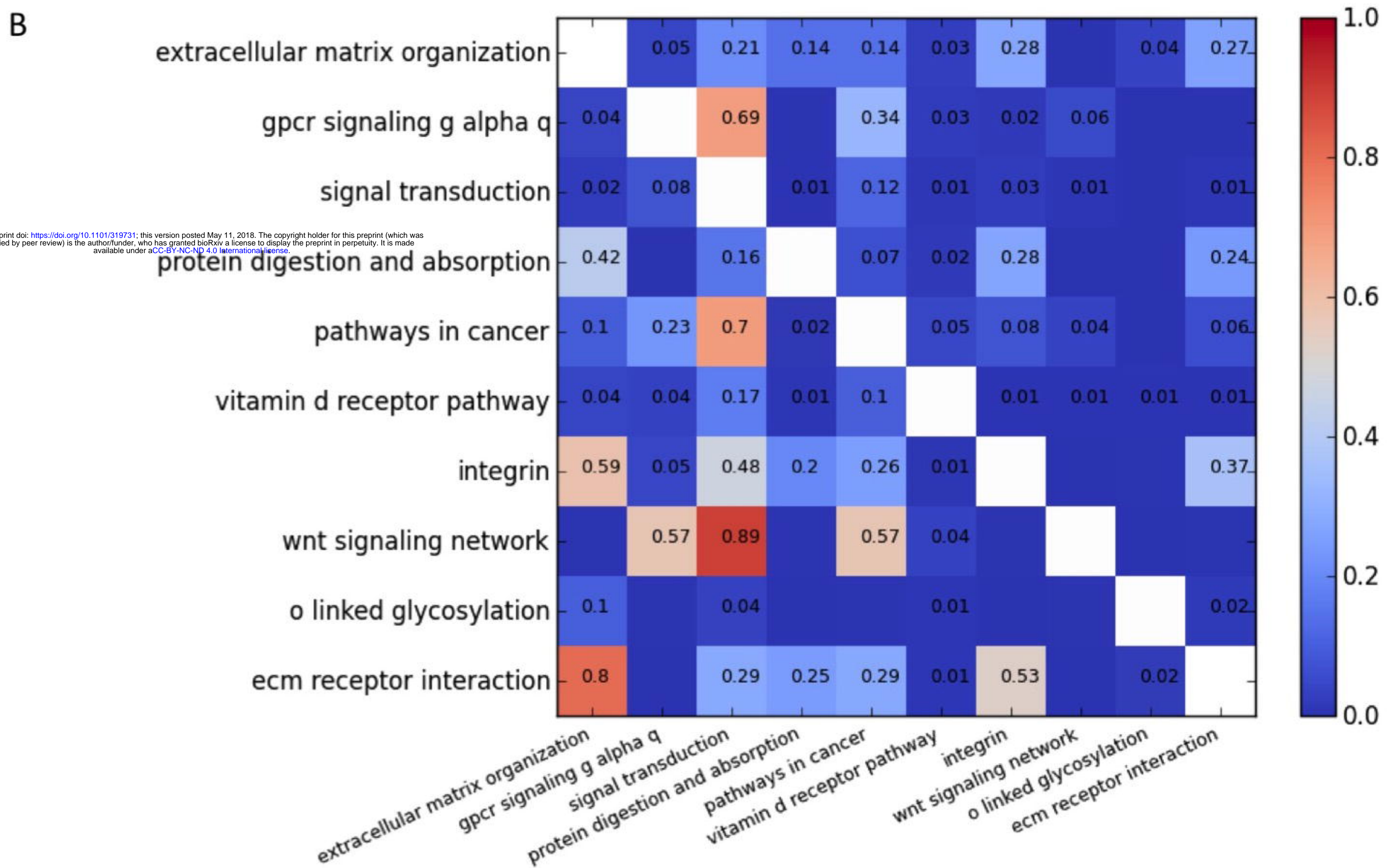
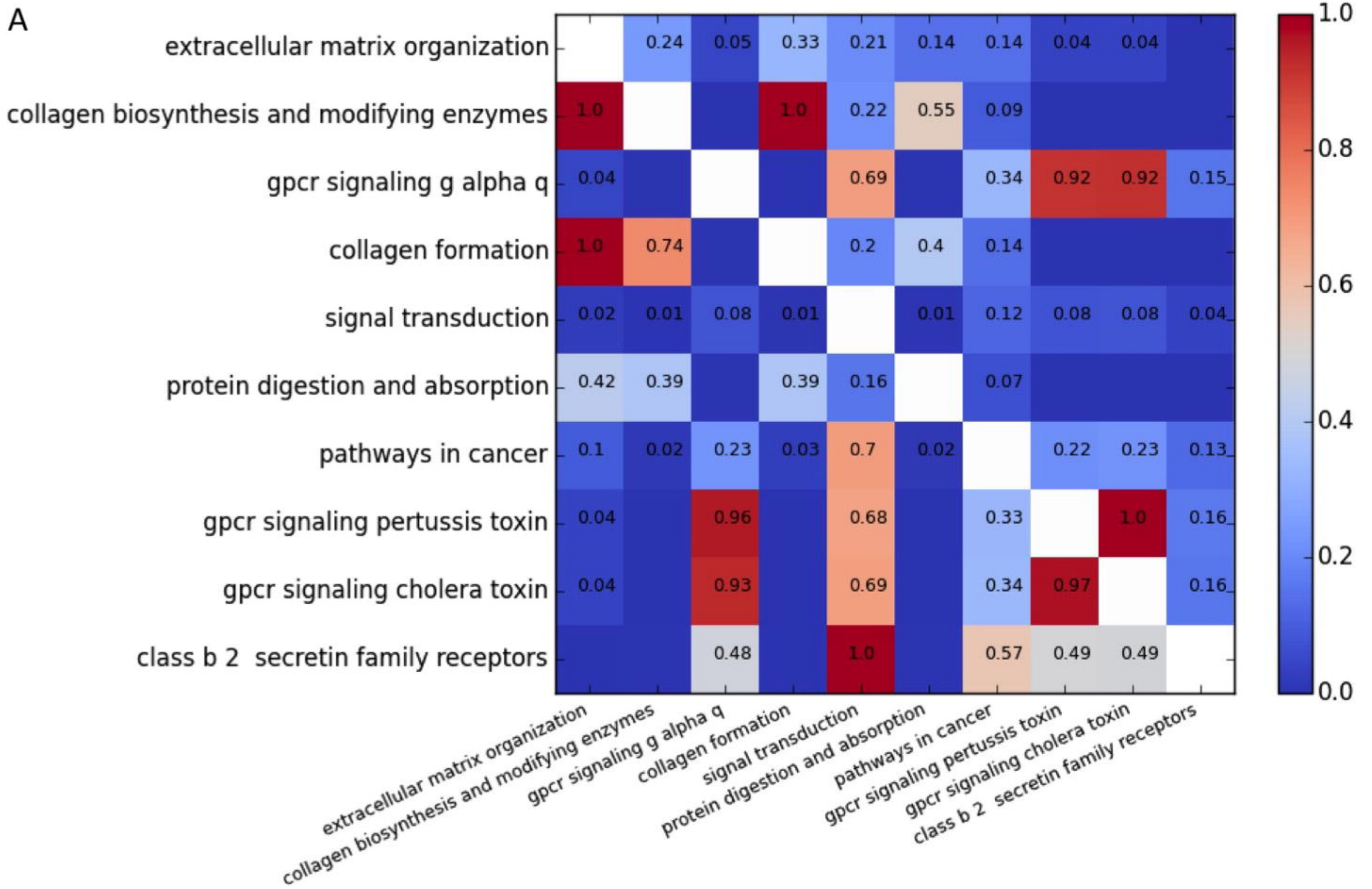
598

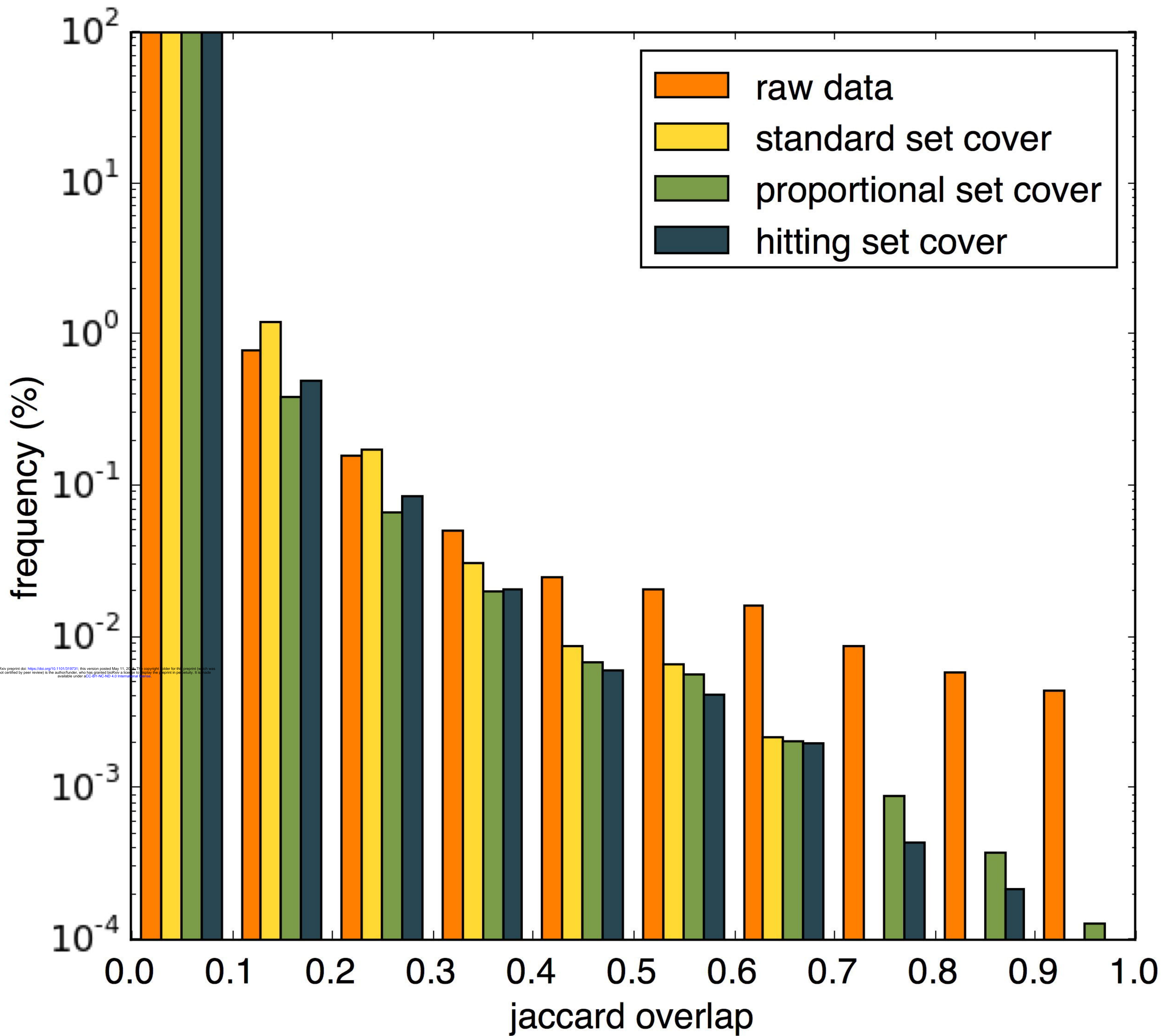
599 11. Additional material

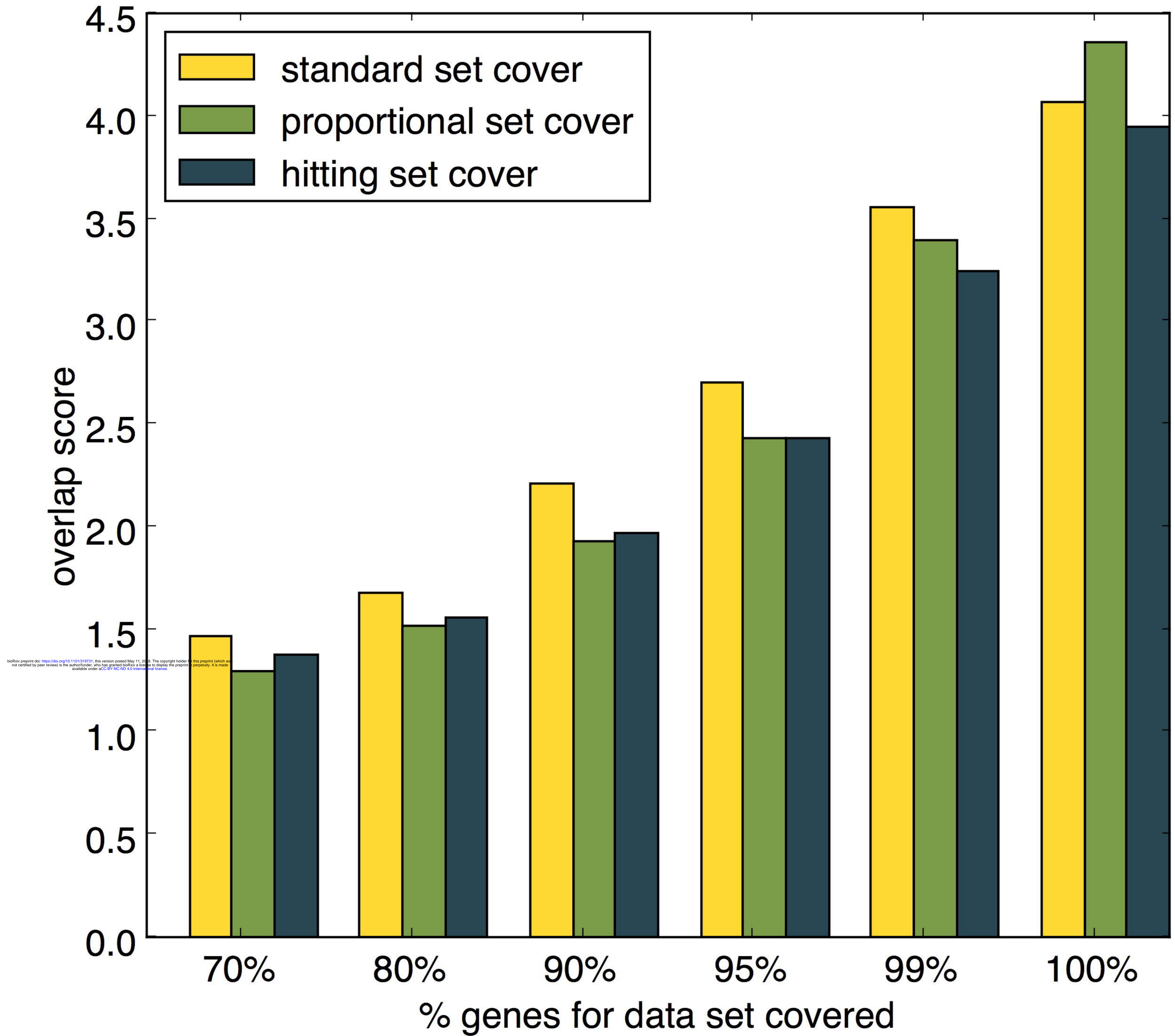
600 Supplementary table 1: Enriched pathways from the osteoarthritis dataset (p-
601 value<0.05). The set cover column indicated the 23 pathways that were included
602 in the set cover. Found in additional file 1.docx.

603

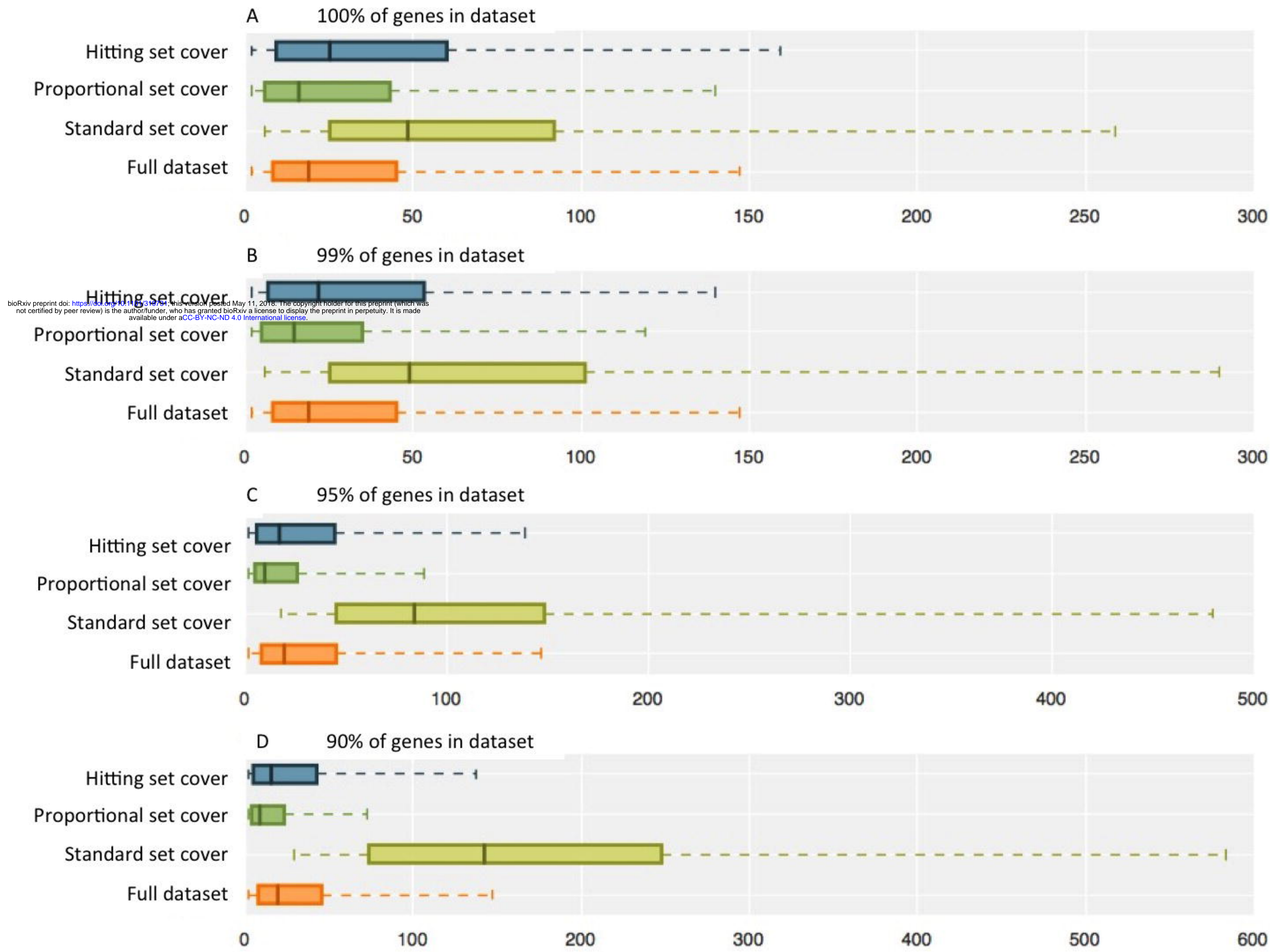
604



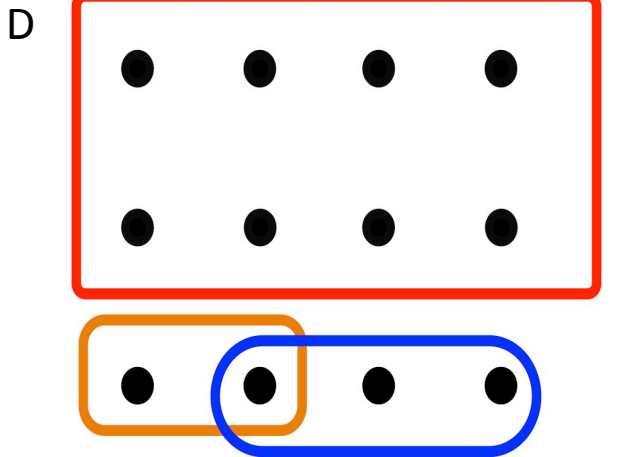
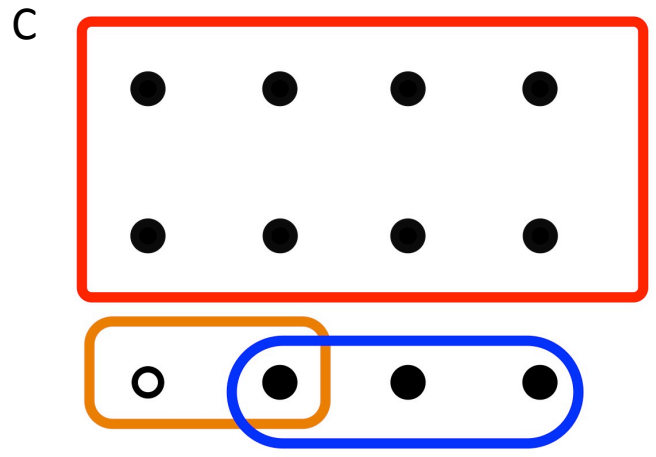
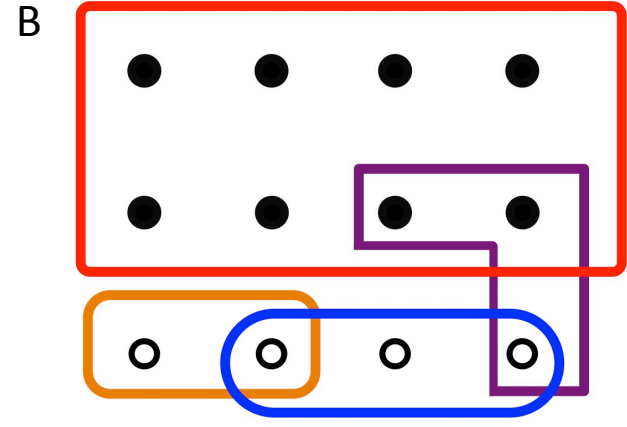
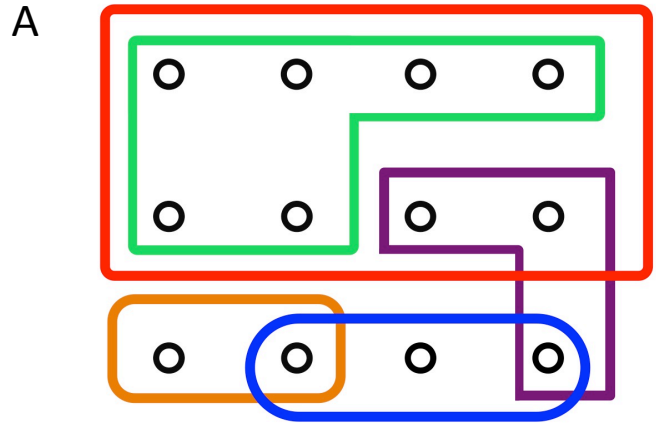




bioRxiv preprint doi: <https://doi.org/10.1101/319731>; this version posted May 11, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



Set cover



Set Cover

Start with $\mathbf{R} = \mathbf{U}$, $\mathbf{C} = \emptyset$ and $\mathbf{SC} = \emptyset$

while $|\mathbf{C}|/|\mathbf{U}| * 100 < GC$ **do**

Select set s_i that maximizes v_i

Add s_i to \mathbf{SC}

Add the elements in s_i to \mathbf{C}

Delete the elements in s_i from \mathbf{R}

end

Return \mathbf{SC}