

[Click here to view linked References](#)

1 **Comparing the utility of *in vivo* transposon mutagenesis approaches in yeast species to infer**
2 **gene essentiality**

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

Anton Levitan¹, Andrew N. Gale², Emma K. Dallon², Darby W. Kozan², Kyle W. Cunningham², Roded

Sharan³ and Judith Berman¹

¹School of Molecular Microbiology and Biotechnology, Faculty of Life Sciences, Tel Aviv University,

Tel Aviv, Israel

²Department of Biology, Johns Hopkins University, Baltimore, MD, USA

³Blavatnik School of Computer Science, Tel Aviv University, Tel Aviv, Israel

25

26 **ABSTRACT**

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

Background:

In vivo transposon mutagenesis coupled with deep sequencing enables large-scale genome-wide mutant screens for genes essential in different growth conditions. Six large scale studies have now been performed with three yeast species (*S. cerevisiae*, *S. pombe* and *C. albicans*), each mutagenized with two of three different heterologous transposons (*AcDs*, *Hermes*, and *PiggyBac*).

Results: We analyzed predictions of gene essentiality for each of the six studies and evaluated the ability of the data to predict gene essentiality using a machine-learning approach. Important data features included a sufficient number of independent insertions and the degree of random insertion distribution. All transposons showed some bias in insertion site preference, both because of jackpot events, specific insertion sequence preferences and preferences for short-range vs long range insertions. For *PiggyBac*, a stringent target sequence limited the ability to predict essentiality in genes with few or no target sequences. Furthermore, the machine learning approach is robust for predicting gene function in less well-studied species by leveraging cross-species orthologs. Finally, comparisons of isogenic diploid vs haploid *S. cerevisiae* isolates identified several genes that are haplo-insufficient, while most essential genes, as expected, were recessive.

Conclusions: We provide recommendations for the choice of transposons and the inference of gene essentiality in genome-wide studies of eukaryotic microbes such as yeasts, including species that have been less amenable to classical genetic studies. These include maximizing the

49 number of unique insertions, avoiding transposons with stringent target sequences and a method
50 for cross-species transfer learning.

52 INTRODUCTION

54 Work with model yeasts such as *Saccharomyces cerevisiae* and *S. pombe* has pioneered the
55 combination of genotype/phenotype comparisons at a genomic scale. For these yeasts, with
56 genome sequences available for over 20 years[1, 2] and facile gene replacement protocols, have
57 relied heavily on comprehensive collections of deletion mutants[3, 4] for high throughput
58 dissection of specific genotypes as well as for genetic interactions with drugs (reviewed in Lehár et
59 al.[5]) and for gene-gene interactions through systematic analysis of double and triple mutant
60 analysis (e.g., Reguly et al.[6], Kuzmin et al.[7]). In animals and plants that are less amenable to
61 such directed molecular manipulations, the use of heterologous transposons *in vivo* has facilitated
62 genetic analysis, within the limitations imposed by the transposon excision/insertion process[8, 9].
63 With the advent of deep sequencing, such studies have also become more facile and have been
64 performed in the two model yeasts as well[10, 11].

66 *In vivo* transposon mutagenesis generally involves the introduction of a heterologous DNA
67 transposon, along with the genes (e.g., the relevant transposase) required to induce its active
68 transposition into a clonal isolate of a species of interest. Upon induction, the transposase excises
69 the transposon from its original location (excision site) and inserts it into a single new position in
70 the genome. Each cell harbors, at most, a single transposition mutation because the frequency of
71 transposon excision and reinsertion is quite low. The transposon is usually engineered for facile

72 selection of excision and/or reinsertion events, allowing detection and enrichment of these rare
73 events.

74
75 *In vivo* transposon insertion provides several advantages, as it rapidly yields large numbers of
76 mutants in a single step and easily can be performed in parallel strains with different mutations or
77 genetic backgrounds. Because it does not require much prior knowledge, it can also be performed
78 in non-model species, where each experiment is likely to be highly informative. The only
79 transformation steps required are those used to engineer the starting strain. This bypassing the
80 problem of low transformation efficiency in many species. It also avoids the unintended genome
81 alterations (e.g., aneuploidies) that often accompany DNA transformation[12]. The sites of
82 transposon insertion throughout the genome can be identified *en masse* using very large
83 collections of independent insertion event clones, coupled with deep sequencing of the DNA
84 immediately adjacent to the new transposon locus. For example, a high throughput
85 genotype/phenotype analysis of 30 bacterial species grown under >170 different nutrient and
86 stress conditions recently assigned functions to thousands of genes including ~300-600 genes per
87 bacterium that are essential for viability[13].

88
89 Three different transposon systems have been used for *in vivo* mutagenesis in yeasts: *AcDs* from
90 *Zea mays*, *Hermes* from *Musca domestica* and *PiggyBac*, from *Trichoplusia ni*. *AcDs* has been used
91 primarily in plant species but was also engineered for increased efficiency in the model yeast *S.*
92 *cerevisiae*[14] and, later, in *C. albicans*[15]. *AcDs* does not display any insertion sequence
93 preference, although it has a higher frequency of insertions into intergenic regions than coding
94 regions and has a bias for reinsertion near the initial site of excision. *Hermes* has been used for
95 mutagenesis in *S. pombe* and *S. cerevisiae*[17, 18]; it prefers to insert at genomic positions with

96 the target sequence TnnnnA. *PiggyBac* (PB) has been used in mammalian systems such as rat,
97 mouse[19] and also in *S. pombe*[17]. *PB* has a strong preference for insertion at TTAA sequences,
98 which are generally more frequent in A-T-rich intergenic regions than within coding sequences.
99
100 Transposon insertion within an ORF is generally assumed to cause a loss-of-function mutation.
101 Identifying the phenotypes associated with loss-of-function mutations in specific genes allows the
102 prediction of genetic functions. Cells in which the transposon inserted into a gene essential for
103 viability will fail to grow and thus be lost from the population. By contrast, cells with mutations in
104 non-essential genes are expected to be well-represented in the cell population. Insertion of a
105 transposon carrying a strong promoter into an ORF could activate expression inappropriately; can
106 be useful for the study of gain-of-function mutations.
107
108 *In vivo* transposon mutagenesis studies of yeasts include analysis of *S. cerevisiae* with *Hermes* (this
109 study) or with the mini-Ds derivative of the *AcDs* system[10], in *S. pombe* with *Hermes*[11] and
110 *PB*[20] and in *C. albicans* with *AcDs*[21] and with *PB*[22] (Fig. 1A). In earlier work, we applied a
111 machine learning (ML) approach to infer gene essentiality from the *C. albicans AcDs* data. Here,
112 we modified the ML approach to predict the likelihood of essentiality for the complete set of
113 predicted open reading frames these sixth *in vivo* transposon datasets. We compared the
114 strengths and challenges of the different transposons in each species, with the goal of reaching
115 insights concerning the number of insertion events required for accurate predictions, the
116 distribution of mutations, and the degree to which different transposons, with different sequence
117 dependencies, provided similar or different conclusions. The goal was to provide metrics that
118 assist in determining the advantages and disadvantages of different transposon systems so as to

119 optimize the data produced in a given *in vivo* transposon system and to suggest approaches for
120 generating whole genome data in understudied yeast species.

121
122

122 **RESULTS AND DISCUSSION**

123
124

124 **A comparative analysis of the transposon mutagenesis studies**

125
126

126 We compared six *in vivo* transposon insertion mutagenesis experiments, produced using three
127 different heterologous transposons (*AcDs*, *Hermes* and *PiggyBac*) in three different species (*S.*
128 *cerevisiae*, *S. pombe* and *C. albicans*). Details of the datasets are provided in the methods section
129 and relevant parameters are highlighted in Table 1. The number of transposition events detected
130 in the different studies varied considerably, from over 500,000 unique insertion sites (hits) and 84
131 M total reads for the *C. albicans* *AcDs* (*CaAcDs*), to as few as 37,500 unique hits and 6.1 M reads in
132 the *S. pombe* *PiggyBac* (*SpPB*) data set (Table 1). The number of reads per hit also varied
133 considerably, from 41 to 170.

134
135

135 **Overview of ML approach for gene essentiality prediction**

136
137

137 We first mapped the hit and read frequency of the transposons in each of the three reference
138 genomes (Fig. 1C). Sites of transposon insertion were identified based upon targeted sequencing
139 of regions adjacent to the inserted transposon. Slightly different sequencing protocols were used
140 in the different studies, but all six essentially amplified Tn-adjacent sequences and mapped them.
141 Theoretically, essential genes have no hits and non-essential genes have many hits; however,
142 distinguishing essential and non-essential genes from the data is not entirely straight-forward

143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165

143 (e.g., Fig. 1C, Gene X_n). To address this ambiguity, we extended a previous approach for gene
144 essentiality prediction[21]. We first chose input features from transposon data that were likely to
145 be informative in the essential/nonessential decision: the number of unique insertion sites (hits)
146 per ORF, the degree to which those insertion sites were enriched in the population (reads), as well
147 as normalization factors that consider the insertion frequency as a function of chromosome
148 position.

149
150 Training sets were built using information from the two model species with gene essentiality data
151 available from classical genetic approaches (e.g., comprehensive ORF deletion analysis) (reviewed
152 in Giaever and Nislow[23] and Spirek et al.[24]) (Table S1 (training sets)). The specificity and
153 sensitivity of the approach was analyzed using the AUC (area under the receiver operating
154 characteristic curve (Fig. 1D). For *C. albicans*, which did not have extensive prior knowledge of
155 gene essentiality, we constructed a training set from a core set of genes whose orthologs were
156 known to be essential in both model yeasts. This approach is likely to be useful for other species
157 that lack sufficient prior knowledge of gene essentiality to construct a within-species training set.

158
159 We assessed the performance of each classifier by producing training sets using genes known or
160 inferred to be essential and non-essential, as described in the Methods. The classifiers were
161 trained and their performance, assessed using the AUC measure, was high across most examined
162 studies (>0.94). The one exception was the *SpPB* study, which had far fewer unique insertion sites
163 (Table 1) and had an AUC of 0.785. The highest AUC levels were seen with the *AcDs* in both *C.*
164 *albicans* and *S. cerevisiae*. Of note, these two studies also had the largest number of total hits and
165 reads. All the considered ML features for each ORF in every study and the predicted verdicts of

166 essentiality are provided in Tables S2-S7. Below, we describe the main insights gained from this
167 comparison.

167
2
3
168
4
5

169 **Insight #1: Optimize the number of independent insertion sites (hits) for highest quality**
170 **predictions of gene essentiality**

10
11
12
13

14 The total number of unique insertion sites (hits) and the performance (AUC values) were high
15
16 highly correlated, and this correlation was statistically significant (Fig. 2a; Pearson's $r = 0.892$; p-
17 value = 0.0169). By contrast, the total number of sequencing reads showed a weaker correlation
18
19 with the AUC that was not statistically significant (Fig. 2b; Pearson's $r = 0.636$; p-value = 0.1741). If
20
21 we disregard the worst performing *SpPB*, the correlation of the AUCs with the total number of hits
22 rises dramatically to Pearson's $r = 0.995$; p-value = 0.0003, and the correlation with the total
23
24 number of sequencing reads remains similarly weak: Pearson's $r = 0.652$; p-value = 0.2327. Thus, it
25 a library with many independent hits will improve performance and simply increasing the number
26 of sequencing reads is not likely to be sufficient to obtain optimal results. Increasing the number
27 of independent hits requires collection of sufficient numbers of independent colonies soon after
28 transposase induction and the resulting transposon excision and reinsertion. An advantage of
29 *Hermes* is that most insertions occur during stationary phase, so transposase-inducing conditions
30 can be tolerated throughout the growth period. Experimental designs that optimize isolation of
31 independent events are critical.

32
33
34
35
36
37
38
39

40 **Insight #2: Avoid libraries with high levels of jackpot events**

41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

189 Jackpot events are the appearance of extraordinarily high numbers of reads in a very small
190 number of insertion sites. When the number of reads greatly exceeds the theoretical number of
191 cell divisions in the experiment, this is likely due to a transposition event happened prior to the
192 induction of transposon excision in the experiment. Jackpot events are a major pitfall in that much
193 sequencing capacity is wasted on detection of a single insertion site. Jackpot events with >1M
194 sequencing reads were present in 4 of the 6 data sets; *SpPB* and *ScHermes* had no major jackpot
195 events (no hits with ≥ 1000 -fold more reads than the average read/hit) (Table 1). Both of these
196 libraries also had far fewer total sequencing reads than the other studies (6M and 18M vs 24-84M
197 for the other libraries).

199 Of note, within a data set, some individual experiments had jackpot events that were far more
200 than others (Table 1), which would be expected if jackpot events occur stochastically. Importantly,
201 jackpot events were not clearly associated with one of the three species or with the transposon
202 type. This suggests that jackpots arise from technical, rather than biological issues.

204 Avoiding jackpot events is important because they reduce data quality considerably: the higher
205 the number of reads at a few jackpot sites, the lower the number of informative hits and reads.

206 Avoiding the selection of cells in which a transposon was already mobilized is key to ensuring that
207 the number of hits and reads provide good genome coverage. Dividing the cultures into dozens of
208 small cultures and then re-pooling these sub-cultures after transposase induction can effectively
209 dilute out most jackpot events. Preparing several independent libraries and sampling a few
210 sequences in each may also be worthwhile. For example, if a tested library shows one sequence
211 twice in one hundred colonies, it is likely to be a >1M jackpot event.

213 **Insight #3. Consider which features are most important in the analysis of a given transposon**

214
2

215 In decision tree based algorithms, such as Random Forest[25], every node is a condition of a split

5
216 of the data by a single feature. The splitting process continues until it reaches a stop condition

8
217 such as: all the features have been used, the obtained subset is very small or the training labels

10
218 are the same for the obtained subset. The goal is to reduce entropy (uncertainty) in the data.

13
219 Entropy is zero when all the labels in the obtained subset are the same; and is maximum when

16
220 half of the labels are the same in the obtained subset (in a binary classification). Each split of the

18
221 data by a given feature (node in the tree) reduces the entropy. The importance of a given feature

21
222 in the Random Forrest classifier is the calculated decrease in entropy contributed by that feature.

23
223 Here we describe the features of the classifiers, and discuss their relative importance.

24
224

28
225

31
226 For each ORF, we calculated predictive features including the number of hits, number of reads and

33
227 the length of each ORF, a neighborhood index, which normalizes for insertion bias due to genomic

36
228 position (e.g., proximity to the initial excision site in the genome) and a freedom index, which

38
229 reports the proportion of an ORF that is hit-free. The freedom index is especially useful for

41
230 identifying genes with essential domains, that are able to tolerate insertions outside of the

43
231 essential domain. The number of the transposon hits per transposon target sequences in an ORF

46
232 was an additional feature used in the analysis (Figure 3), where applicable (in *PB* and *Hermes*

48
233 studies). Furthermore, we calculated the number of hits and the number of reads normalized by

51
234 the length of each ORF. We compared the ‘feature importance’ for each library to ask whether

54
235 specific features were more important for the different classifiers and whether feature

56
236 importance was characteristic for a given transposon or yeast species.

59
60

61
62

63
64

65

237

238 The number of hits per ORF played an important role in determining essentiality, with essential
239 genes having far fewer hits than non-essential ORFs (~7 times less, on average, across the 6
240 datasets), consistent with the strong correlation between number of hits and the AUC (Fig. 2A).

241 The number of reads per ORF played a lesser role in these classifications, also consistent with the
242 correlation above (Fig. 2B). Gene length also affected the probability of transposon insertion in a
243 gene, and thus was a crucial normalization parameter for the numbers of hits and reads for every
244 ORF.

245
246 The Neighborhood Index (NI) feature made important contributions in all of the classifications
247 (except *SpPB*, which had far less data). Importantly, the NI did not differ considerably between the
248 different transposons, consistent with the idea that chromatin accessibility, 3D chromosome
249 organization and other factors that may bias the insertion site frequency in a given organism affect
250 the frequency of insertion of different transposons in a similar manner.

251
252 The Freedom Index (FI) was a major contributor to both *ScAcDs* and *CaAcDs* predictions while
253 results with the *PB* and *Hermes* datasets were mixed (Fig. 3). This is consistent with the idea that
254 *AcDs* does not have a specific target sequence and thus inserts throughout ORFs, while *PB* and
255 *Hermes* have fewer target sequences within ORFs. Thus, the FI more important in *AcDs*
256 experiments because hits occur more randomly throughout an ORF.

257
258 The importance of the number of hits in the proximal regulatory sequences (100 bp upstream to
259 the start codon) to the essential/non-essential predictions was only minor, but is highly variable.
260 For example, the impact of *ScHermes* was nearly twice that of *ScAcDs* for this feature. As

261
262
263
264
265

261 described further below, we posit that this difference is due to cryptic enhancer/promoter activity
262 in the miniDs transposon in *S. cerevisiae* that not seen with the *Hermes* transposon.

263

264

264 **Insight #4: Consider the effect of transposon-specific target sequence specificity**

265

266

266 Some transposons have preferred sites of insertion: *Hermes* prefers TnnnnA and *PiggyBac* inserts

267

267 primarily at TTAA sequences; *AcDs* does not have an insertion site preference. Theoretically, the

268

268 length of the insertion site sequence necessarily scales inversely with the number of potential

269

269 unique hit sites. However, it was not clear at what insertion sequence length the resolution of

270

270 studies of gene essentiality becomes limiting.

271

272

273

272 The feature importance of the number of hits per transposon target sequence in an ORF, which

274

274 should be a measure of library saturation, showed a varying degree of importance in the *PB* and

275

275 *Hermes* studies. Curiously, its importance wasn't dependent on the type of the transposon or the

276

276 target sequence prevalence in the genome. This likely because target sequences are preferred

277

277 sites of insertion, yet are not exclusive or absolute. For example, *PiggyBac* in *C. albicans* had 1.6-

278

278 fold more unique insertion sites than the theoretical number of target sequences in the *C. albicans*

279

279 genome. By contrast, for both *ScHermes* and *SpHermes*, the number of target sequences available

280

280 far outnumbered the number of unique hits. The proportion of target sequences not hit ranged

281

281 from 8.9% for *CaPB* to 85% for *SpHermes* (Table 1) and the proportion of hits not in target

282

282 sequences ranged from 14% in *SpPB* to ~50% in *CaPB* as well as both *Hermes* data sets. Thus, we

283

283 surmise that the preference for target sequences is only a minor limitation for both of these

284

284 transposons, except when the total number of hits is very low as in *SpPB*.

285

286

287

288

289

290

291

292

293

285 Another critical issue is the number of genes that lack any preferred target sequences within the
286 ORF; there are 228 and 185 ORFs without a single TTAA sequence in *C. albicans* and *S. pombe*,
287 respectively. These ORFs have a lower probability of acquiring insertions and, if the genes are
288 non-essential, they are much more likely to give false positive information (be predicted essential
289 for lack of insertions). Indeed, 155 ORFs without TTAA sequences were predicted essential in the
290 *CaPB* data and yet are predicted non-essential in the *CaAcDs* study. Similarly, 118 of the 185 ORFs
291 lacking TTAA sequences were predicted essential from the *SpPB* study, but were non-essential in
292 the *SpHermes* study. We assume that many of these genes are false positives, especially given that
293 127 of the 185 ORFs lacking TTAA, including 95 of the 118 aforementioned ORFs, were non-
294 essential in classical genetics studies of *S. pombe*.

295
296
297 Next, we asked if the number of target sequences within an ORF affected the *CaPB* classification
298 performance for that ORF. To address this, we compared the performance (AUC) to sets of ORFs
299 filtered to exclude ORFs with different numbers of target sequences (from 0 to 10) from the
300 training set used to train the classifier (Fig. 4). The AUC increased from ~0.94 for the entire
301 training set to >0.98 for the training set containing only genes with 10 or more target sites (~50%
302 of the genes in the training set). This suggests that studies using the *PiggyBac* transposon may
303 struggle to correctly infer gene essentiality for ORFs with low numbers of target sites.

304
305 **Insight #5. Consider whether the transposon can activate as well as disrupt gene expression.**

306
307 The prediction of essentiality was based upon the assumption that transposon insertion into an
308 ORF disrupted gene expression and produced loss-of-function allele. However, this is not

309 necessarily the case for all genes. For example, if an insertion allele removes a regulatory domain
310 from a protein, the protein may become hyperactive, and thus the result would be a gain-of-
311 function allele. Additionally, some transposons may introduce enhancer and promoter activities
312 that could increase gene expression in some species. The *miniDs* transposon used in the *ScAcDs*
313 data is likely to contain such activities.[10] Consistent with this idea, the *ScAcDs* dataset contains
314 an average of 1.89 insertions within the first 100 bp of essential genes whereas the other datasets
315 including *CaAcDs*, which has a transposon modified from the *miniDs*, contain significantly fewer
316 (0.82 insertions in the first 100 bp of Ess genes). Additionally, many essential genes of *S. cerevisiae*
317 appeared to tolerate *miniDs* insertions, but not *Hermes* insertions, at sites in the 5' UTR that are
318 very close to the start codon. Thus, the *miniDs* transposon in *S. cerevisiae* may facilitate
319 inappropriate activation of gene expression when inserted upstream or within certain genes.

320 321 **Cross-study analysis**

322
323 Knowing the full set of essential and non-essential genes in eukaryotic microbes, including
324 pathogens of humans, animals and plants, will improve our understanding of common and
325 species-specific properties of these understudied organisms. Furthermore, once a transposon
326 library has been collected, it can be screened under many other growth conditions to reveal
327 genotype/phenotype relationships. *In vivo* transposon analysis of gene essentiality is a practical
328 and feasible approach, because the cost in time and resources for obtaining libraries is far lower
329 than that for producing engineered deletion mutants, especially given that the amount of baseline
330 information (other than the genome sequence) about the organisms may be minimal. The only
331 technical hurdle is to introduce the heterologous transposon of interest, either on a plasmid
332 (where feasible) or into a useful locus within the genome of the relevant organism.

333

334 An additional challenge is that ML approaches require a high-quality training dataset (of gold-
2
3
335 standard essential and non-essential genes). For many non-model organisms, such training data is
4
5
336 too sparse to build a robust training set. For *C. albicans*, we circumvented the low numbers of
6
7
8
337 genes already known to be essential by relying upon genes that had been determined to be
9
10
11
338 essential from comprehensive classical genetic deletion studies in both model yeasts (*S. cerevisiae*
12
13
14
339 and *S. pombe*) and that had orthologs in *C. albicans*. Training on *S. cerevisiae* or *S. pombe*
15
16
340 orthologs with consistently essential orthologs yielded good performance predictions for *C.*
17
18
19
341 *albicans* (AUC: 0.940 to 0.993, Table 1). *CaAcDs* performance was lower when training only on the
20
21
22
342 66 genes known to be essential plus the set of presumed non-essential genes (those that had been
23
24
343 successfully deleted in *C. albicans* studies, AUC of ~0.92)[21].
25
26

27
28
29
344
30
345 Next, we considered the quality of the learning performance for each dataset, if we trained on
31
32
33
346 orthologs from one species and predicted essentiality of genes in a different organism (Figure 5).
34
35
347 The transfer learning performance of the classifications was of a comparable quality to the single
36
37
348 study classifiers for most *AcDs* and *Hermes* cases (Figure 5a and Figure 2). Furthermore, it
38
39
40
349 displayed a somewhat symmetrical property: in most cases, there were minor differences in
41
42
350 performance between train/test and test/train pairs (reducing the quality by ~ 0.5% to 5.7%) when
43
44
45
351 the tests were between or among *AcDs* and *Hermes* experiments. Conversely, when testing for
46
47
352 predictions from *PB* data that were trained on either *AcDs* or *Hermes*, the AUCs dropped more
48
49
50
353 dramatically (up to ~21.9%). Thus, *PB* data was less transferable than the *Hermes* and *AcDs* data.
51
52

53
54
55
354
56
355 The low *PB* transferability between *SpPB* and *CaPB* is likely due to the sparser target sequence
57
58
356 distribution relative to either *Hermes* or *AcDs*, which causes *PB* studies to produce false positives
59
60
61
62
63
64
65

357 as noted above, and thus might contribute to reduced performance in cross-study analyses. The
358 lower performance of the classifiers in the original *PB* single studies (Table 1), also may have
359 contributed to the reduced ability to predict essentiality in pools of *PB* mutants using cross-species
360 models.

361
362 Reduced differences cross-study performance could also be due to differences between the
363 importance of different features in the classifiers for the different datasets. To test this possibility,
364 we correlated the vector of the relative feature importance for each study with the feature
365 importance in all the other studies (Figure 5b). The analysis distinguished 3 groups within the 6
366 studies, based on the correlation coefficient values for feature importance between members of
367 the group: *CaAcDs* and *ScAcDs* (Pearson's $r = 0.902$); *Sp Hermes* and *CaPB* (Pearson's $r = 0.898$);
368 and *SpPB* and *ScHermes* (Pearson's $r = 0.929$). Notably, the quality of the transfer learning
369 predictions appears to be independent of both the transposon type and the organism studied,
370 with the exception of the *AcDs* studies. We presume that this is due to the lack of a specific target
371 sequence for the *AcDs* transposon system.

372
373 **Insight #6: As necessary, construct training sets using genes with orthologs in models where**
374 **essentiality is known and then validate the training set manually.**

375
376 We suggest that an initial training set of orthologous genes known to be essential and non-
377 essential in related model organisms can be used to facilitate analysis of a transposon insertion
378 study in a non-model organism with sparse essentiality information. An important caveat is that
379 differences between gene function in different species can alter gene essentiality of a small
380 number of these orthologs; thus, it is important to visually inspect this orthologous training set

381 before applying it. The goal is to remove any genes with insertion patterns that are highly
382 contradictory to the ‘essentiality label’ that the orthologs provided. For example, for *C. albicans*,
383 the entire orthologous training set was reviewed in an unprejudiced fashion by three independent
384 inspectors, who visually reviewed the insertion patterns in the *CaAcDs* data and manually labeled
385 each gene as essential, non-essential or ambiguous. When all three inspectors classified a gene as
386 non-essential (e.g., many insertions throughout an ORF within a genome region that had many
387 insertions outside of that ORF) and the orthologs were labeled ‘essential’ in the two model yeasts,
388 we removed that gene from the training set.

389
390 Once a training set has been established, and the features for the ORFs have been calculated, the
391 Random Forest classifier can be run in a cross-validation scheme and the AUC can be calculated
392 using the essentiality labels. This provides an efficient approach to obtain information about all of
393 the genes in a species that has been sequenced but not subjected to much molecular
394 manipulation. Clearly, the same approach can be used to compare the essentiality of the same
395 sets of genes grown in different conditions as well, potentially providing large amounts of
396 phenotypic data across an entire set of ORFs. If applied to a species that had not been the subject
397 of many genetic studies, such data would represent a treasure-trove of information about genes
398 that had not been previously studied and the phenotypes associated with loss-of-function of those
399 genes.

400 401 **A Comparative Analysis of Gene Essentiality Predictions**

402
403 An important issue is whether different transposon insertion studies in the same organism had
404 similar or different predictions from one another and from the known essentiality status of

405 deletion mutants, which are by definition 'loss-of-function' null alleles. For *S. cerevisiae*, the
406 classifiers displayed a high degree of agreement on the final verdicts of gene essentiality (Figure
407 6a), while more discrepancies were evident for the *C. albicans* and *S. pombe* studies. Both
408 *PiggyBac* studies predicted a much higher number of essential genes than the *AcDs* or *Hermes*
409 studies (Figure 6b and 6c) as expected from the paucity of target sequences that are likely to give
410 false positive predictions discussed above. For example, the *CaPB* study had an average 5.84
411 target sites per kb in genes likely to be false positives vs. 10.62 target sites per kb in all the genes
412 (Mann Whitney U: p-value < 2.38*e⁻⁷⁸). Importantly, when compared to the set of essential genes
413 for each species determined by deletion analysis, the transposon studies also did quite well, with
414 only 20 to 35% of the genes in disagreement. In some cases, such discrepancies were found to be
415 due to issues with the deletion collection isolates. For example, ~8% of the original *S. cerevisiae*
416 deletion collection carried aneuploidies or gene amplifications,[26] and ~10% of *S. pombe* deletion
417 strains retained a wild-type copy of the ORF that had been targeted for deletion. Extra copies of
418 the 'deleted' gene reduces the apparent number of essential genes.

420 **Gene essentiality in haploid versus diploid strains of *S. cerevisiae***

421
422 *S. cerevisiae* is readily grown in both the diploid and haploid states, which allows identification of
423 the haplo-insufficient subset of genes among the set of essential genes. Based on gene knockout
424 studies, only 2 genes (*NDC1*, *MLC1*) were classified as haplo-insufficient,[27, 28] while all other
425 essential genes were haplo-proficient (i.e. heterozygous knockouts in diploids were viable). To
426 determine whether additional haplo-insufficient genes exist in *S. cerevisiae*, we collected
427 *ScHermes* insertions in diploid strain BY4743 and compared them to haploid strains BY4741 and
428 BY4742. *ScHermes* transposon mutagenesis libraries were used with the classifier that had been

429 trained on the *SpHermes* haploid training set data, applying the same threshold for classification
430 (Figure 8). The classifier identified 155 genes as “essential in both haploid and diploid”, a number
431 far higher than expected. Upon closer analysis, 98 contained regions of poor mapping due to
432 duplications elsewhere in the genome, 50 were categorized as dubious in the Saccharomyces
433 Genome Database (yeastgenome.org), one (*LEU2*) had been deleted in the strains studied, and the
434 two known haploinsufficient genes (*NDC1*, *MLC1*) had been identified, providing support for this
435 approach. Upon visual inspection of data of the remaining four genes, one essential gene (*BCY1*)
436 appeared haploinsufficient, whereas another (*RPC10*) contained numerous insertions in its 5’ UTR
437 in diploids but not haploids, suggesting that it might not be essential (Fig. 9). The other two ORFs
438 predicted to be haploinsufficient are very small (165-225 bp) and also are within regions of sparse
439 insertion density. Thus, we have lower confidence in the data for these two genes. Thus,
440 transposon mutagenesis of a diploid strain successfully revealed the two known haploinsufficient
441 genes and one new one (*BCY1*), which is known to be essential in the conditions employed in the
442 screen, but not essential in other culture conditions[29].

443
444 The classifier also identified 29 genes as “haploinsufficient” in diploids and not essential in
445 haploids. Of these, 20 could be dismissed based on their annotation as dubious or the presence of
446 duplicated (unmappable) segments. All of the remaining 9 genes were small (87-528 bp) and were
447 found in regions of sparse insertion density. These genes are annotated in SGD as non-essential
448 and are probably false positives.

449
450 These observations raise an important issue about data quality control. It is important to filter
451 dubious and uninformative ORFs from the data set (as was done in the analysis of *CaAcDs*[21]. This
452 includes genes with repeated domains or duplicate copies in the genome that prevent

453 unambiguous mapping of short Illumina reads. Furthermore, predicting the essential/non-
454 essential status for short ORFs and especially those located in regions with sparse intergenic
455 insertions is more likely to be problematic.

456
457 **Insight #7: Prediction quality increases considerably when uninformative data such as**
458 **mitochondrial genome sequences and duplicated genes that are difficult to map.**

459
460 **In summary,** we suggest a number of metrics for the inference of gene essentiality using *in vivo*
461 transposon mutagenesis studies in yeasts, including those with little available genetic data.
462 Maximizing the total number of unique transposon insertions is the most critical factor in
463 achieving optimal performance of the classification. It can be attained by collecting many
464 independent insertion clones, striving to reduce the possible jackpot events in the study,
465 maximizing the depth of the coverage and by utilizing a transposon with a fairly permissive target
466 sequence or no preferred target sequence. Furthermore, while transposons with relatively
467 stringent target sequences have some advantages for screens that identify individual mutants, for
468 determining gene essentiality they are less robust, as the low number of potential target
469 sequences, and especially the lack of any target sequences, in certain ORFs will increase the
470 likelihood of falsely classifying non-essential genes as essential. Additionally, we think that
471 transposon mutagenesis is an ideal approach to gain large amounts of useful genotype/phenotype
472 data understudied organisms: the cross-species learning methodology allows inference of gene
473 essentiality based on conserved orthologs, especially when coupled with visual screening of the
474 data. Finally, *in vivo* transposon mutagenesis is an incredibly useful tool for high throughput
475 genomic studies, not only of gene essentiality *per se*, but also of genes required under specific

476 selective conditions. We hope that the recommendations provided here will facilitate future work
477 to understand gene in a wide range of yeast species.

2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

478 MATERIALS AND METHODS

479 Data Acquisition

480 Experimental

481 *Sc Hermes* data was obtained as follows: The haploid and diploid strains of *S. cerevisiae* were
482 transformed with plasmid pSG36.[30] A single colony was suspended in 100 mL synthetic
483 complete (SC) medium lacking uracil and containing 2% galactose, divided into twenty 16 x 150
484 mm glass culture tubes, and shaken for 3 days at 30°C. This protocol yielded $\sim 5 \times 10^6$ cells
485 bearing transposon insertions per mL ($\sim 3\%$ of all cells). To enrich for cells bearing transposon
486 insertions, the twenty cultures were pooled, centrifuged, and the cell pellet was resuspended in
487 600 mL SC medium containing 2% glucose, 0.1 mg/mL nourseothricin, and 1 mg/mL 5-fluoroorotic
488 acid, and then shaken overnight at 30°C. The cells were pelleted, resuspended in 600 mL of the
489 same medium, and cultured as before. Finally, 60 mL of these enriched cells were pelleted,
490 resuspended in 600 mL of the same medium, and cultured as before. These highly enriched cells
491 were pelleted, resuspended in 15% glycerol, and frozen in aliquots at -80°C. To extract genomic
492 DNA, 100 mg of thawed cell pellets were washed three times in 1 mL deionized water and
493 extracted using Quick-DNA Fungal/Bacterial Miniprep kit (Zymo Research). A total of 2.4 μ g of
494 purified gDNA was fragmented by sonication in four separate tubes using a Diagenode Picoruptor.
495 The fragmented DNA was then end repaired, ligated to splinkerette adapters, size selected with
496 AMPure xp beads, and PCR amplified in separate reactions using transposon-specific and adapter-
497 specific primers as detailed previously.[31] Samples were then PCR amplified to attach Illumina P5
498 and P7 (indexed) adapters, purified with AMPure xp beads, mixed with phiX-174, loaded into
499 MiSeq instrument (Illumina) and 75 bp of each end was sequenced using primers specific for
500 Hermes right inverted repeat and P7. Detailed protocols and primer sequences are available upon
501 request. De-multiplexed reads were mapped to the *S. cerevisiae* S288C reference genome using

502 Bowtie2, and any mapped reads with a quality score < 20 or a mismatch at nucleotide +1 were
503 removed. This process was repeated a total of 3 times in diploid strain BY4743, 2 times in haploid
504 strain BY4741, and 1 time in haploid strain BY4742. The diploid and haploid datasets were
505 combined prior to analyses. The *S.cerevisiae* *Hermes* data (mapped reads and counts) are
506 available at <http://genome-euro.ucsc.edu/s/CunninghamLab/Hermes%20Vs%20AcDs> . FastQ files
507 are available from Sequence Read Archive (SRA) and the ArrayExpress Experiment Archive
508 (ArrayExpress), which are core repositories of the European Nucleotide Archive (ENA) at accession
509 number [XXXX] <https://www.ebi.ac.uk/ena/about/data-repositories> (to be updated when
510 accession number is issued).

514 **Publicly available databases**

515 The rest of the datasets analyzed here were obtained from previously published studies. *ScAcDs*
516 was published by Michel et al, 2017[10], from which both WT1 & WT2 were combined for the
517 analysis. The data was downloaded: [https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-](https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-4885/samples/)
518 [4885/samples/](https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-4885/samples/). *SpPB* was published by Li et al, 2011[20] and the data was obtained from the SRA
519 database: SRR089408. *Sp Hermes* was published by Guo et al, 2013[22] and the data was obtained
520 from the SRA database: SRR327340. *CaAcDs* was published by Segal et al, 2018[21] and the data
521 was obtained from the SRA database, where SRR7824843, SRR7824841 and SRR7824838 files,
522 were combined for the analysis. *CaPB* was published by Gao et al, 2018[22] and the data was
523 obtained from the SRA database, where all the following files were for the analysis: DMSO
524 (untreatment): SRR7704188, SRR7704193, SRR7704196; 5-FOA (untreatment: SRR7704189,
525 SRR7704194, SRR7704200; No drug screen: SRR7704195. All the SRR files were obtained using

526 fastq-dump, with the following bash command: fastq-dump --gzip --skip-technical --readids --

527 dumpbase --split-files --clip <SRR*****>

528

529

529 Data Processing

530 The .fastq files downloaded with fastq-dump, were processed using cutadapt to filter out reads

531 not containing partial transposon sequences. Reads with transposon sequences were trimmed to

532 remove the transposon sequences for alignment purposes, as follows: cutadapt --cores=8 -m 2 -g

533 <primer sequence> <input fastq filename> -o <output fastq filename> --discard-untrimmed --

534 overlap <overlap length>. In the analysis of the *Sc Hermes* study all the sequencing reads

535 contained the transposon and the reads start at the first genomic base, thus required no filtering.

536 In the analysis of the *Sp PiggyBac* we filtered the reads containing the transposons from the rest,

537 by identifying the ACGCAGACTATCTTTCTAGGG sequence, cutting it out and aligning only the

538 remaining part of the relevant reads. In the analysis of the *Ca AcDs* we filtered the reads

539 containing the transposons from the rest, by identifying the

540 GTATTTTACCGACCGTTACCGACCGTTTTTCATCCCTA sequence, cutting it out and aligning only the

541 remaining part of the relevant reads, starting 37bp downstream (Segal et al, 2018[21]). In the

542 analysis of the *Ca PiggyBac*, we filtered the reads containing the transposons from the rest, by

543 identifying the TGCATGCGTCAATTTTACGCAGACTATCTTTCTA sequence, cutting it out and aligning

544 only the remaining part of the relevant reads, starting 3bp downstream. In the analysis of the *Sc*

545 *AcDs* study we used the published transposon hitmaps of WT1 and WT2. In the analysis of the *Sp*

546 *Hermes* study we used the published transposon hitmaps (Segal et al, 2018[21]).

547

548

548 Alignment of reads and mapping the transposon hits

549

550

551

552

553

554

555

556

557

558

549 bowtie2 indices were created for each organism and gffutils databases were created for each
550 organism's genetic features, using the latest versions of the reference genomes (fasta) and the
551 genomic feature files (gff), which were downloaded from the respective official sources for the
552 three organisms: *S. cerevisiae* from <https://downloads.yeastgenome.org>, *S. pombe* from
553 <ftp://ftp.pombase.org/pombe/> and *C. albicans* <http://www.candidagenome.org/download/>.
554 Sequencing reads were aligned using bowtie2 with the default settings. The resulting sam files
555 were converted to bam using samtools. bam files were sorted using samtools and indexed via
556 pysam. Transposon hits and their corresponding reads were mapped to the respective genomes
557 and counted in each genomic feature. Transposon target sites were found in every genome using
558 Biopython and counted in each genetic feature.

560 **Gene essentiality classification**

561 Table 2 summarizes the features for machine learning classification that were engineered from the
562 mapped transposon hits, reads and the transposon target sequences in the genomes. Random
563 Forrest classification was performed, using Python's scikit-learn library with the default
564 parameters, except the n_estimators parameter that was increased to 200, and the random_state
565 parameter was fixed at 0, for reproducibility purposes. The results were validated using a 5-fold
566 cross-validation. Essentiality labels for the training set of each organism were obtained previously
567 (Shiftman et al, 2018) and are provided in Table S1.

568 Thresholds for the essentiality predictions in each classification were chosen as follows: Two
569 metrics were evaluated (Figure 8): 1) Minimum of the Euclidean distance between (0, 1) and the
570 receiver operating characteristic (ROC) curve. 2) Maximum of the vertical distance between the
571 line describing a random choice (a straight line from (0, 0) to (1, 1)) and the ROC curve. The first
572 method was chosen, and we verified that the second metric is reasonably close, to eliminate any

573 possible artifacts. We predicted the essentiality of all the available genes for each organism based
574 on their respective features, and using the aforementioned method to choose the threshold for
575 each binary classification.

576
577 Figures were generated using Python's matplotlib and seaborn libraries. The schematics were
578 drawn using Inkscape. Mann Whitney U p-values and Pearson's correlation coefficients and p-
579 values were calculated using Python's Scipy.

580

581 **FUNDING ACKNOWLEDGEMENTS**

582 European Research Council Advanced Award 340087 (RAPLODAPT) to J.B., the Israel Science

583 Foundation (grants no. 715/18, 757/12) to R.S. by NIH R21-AI130722 and T32-GM007231 to KWC.

584 A.L. is supported by a fellowship from the Edmond J. Safra Center for Bioinformatics.

585

586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665

586 **References**

- 587 1. Goffeau A, Barrell BG, Bussey H, Davis RW, Dujon B, Feldmann H, Galibert F, Hoheisel JD,
1
2
388 Jacq C, Johnston M, et al: **Life with 6000 genes.** *Science* 1996, **274**:546, 563-547.
3
4
589 2. Wood V, Gwilliam R, Rajandream MA, Lyne M, Lyne R, Stewart A, Sgouros J, Peat N, Hayles
6
7
590 J, Baker S, et al: **The genome sequence of *Schizosaccharomyces pombe*.** *Nature* 2002,
9
10
591 **415**:871-880.
11
12
592 3. Giaever G, Chu AM, Ni L, Connelly C, Riles L, Veronneau S, Dow S, Lucau-Danila A, Anderson
14
15
593 K, Andre B, et al: **Functional profiling of the *Saccharomyces cerevisiae* genome.** *Nature*
16
17
594 2002, **418**:387-391.
19
20
595 4. Kim DU, Hayles J, Kim D, Wood V, Park HO, Won M, Yoo HS, Duhig T, Nam M, Palmer G, et
21
22
596 al: **Analysis of a genome-wide set of gene deletions in the fission yeast**
24
25
597 ***Schizosaccharomyces pombe*.** *Nat Biotechnol* 2010, **28**:617-623.
26
27
598 5. Lehar J, Stockwell BR, Giaever G, Nislow C: **Combination chemical genetics.** *Nat Chem Biol*
28
29
30
599 2008, **4**:674-681.
32
33
600 6. Reguly T, Breitkreutz A, Boucher L, Breitkreutz BJ, Hon GC, Myers CL, Parsons A, Friesen H,
34
35
601 Oughtred R, Tong A, et al: **Comprehensive curation and analysis of global interaction**
37
38
602 **networks in *Saccharomyces cerevisiae*.** *J Biol* 2006, **5**:11.
39
40
603 7. Kuzmin E, VanderSluis B, Wang W, Tan G, Deshpande R, Chen Y, Usaj M, Balint A, Mattiazzi
41
42
43
604 Usaj M, van Leeuwen J, et al: **Systematic analysis of complex genetic interactions.** *Science*
45
46
605 2018, **360**.
47
48
606 8. Munoz-Lopez M, Garcia-Perez JL: **DNA transposons: nature and applications in genomics.**
50
51
607 *Curr Genomics* 2010, **11**:115-128.
52
53
608 9. Kawakami K, Largaespada DA, Ivics Z: **Transposons as tools for functional genomics in**
54
55
56
609 **vertebrate models.** *Trends Genet* 2017, **33**:784-801.
58
59
60
61
62
63
64
65

- 610 10. Michel AH, Hatakeyama R, Kimmig P, Arter M, Peter M, Matos J, De Virgilio C, Kornmann B:
611 **Functional mapping of yeast genomes by saturated transposition.** *Elife* 2017, **6**.
612 11. Guo Y, Park JM, Cui B, Humes E, Gangadharan S, Hung S, FitzGerald PC, Hoe KL, Grewal SI,
613 Craig NL, Levin HL: **Integration profiling of gene function with dense maps of transposon**
614 **integration.** *Genetics* 2013, **195**:599-609.
615 12. Bouchonville K, Forche A, Tang KE, Selmecki A, Berman J: **Aneuploid chromosomes are**
616 **highly unstable during DNA transformation of *Candida albicans*.** *Eukaryot Cell* 2009,
617 **8**:1554-1566.
618 13. Price MN, Wetmore KM, Waters RJ, Callaghan M, Ray J, Liu H, Kuehl JV, Melnyk RA, Lamson
619 JS, Suh Y, et al: **Mutant phenotypes for thousands of bacterial genes of unknown**
620 **function.** *Nature* 2018, **557**:503-509.
621 14. Lazarow K, Du ML, Weimer R, Kunze R: **A hyperactive transposase of the maize**
622 **transposable element activator (Ac).** *Genetics* 2012, **191**:747-756.
623 15. Mielich K, Shtifman-Segal E, Golz JC, Zeng G, Wang Y, Berman J, Kunze R: **Maize**
624 **transposable elements Ac/Ds as Insertion mutagenesis tools in *Candida albicans*.** *G3*
625 *(Bethesda)* 2018, **8**:1139-1145.
626 16. Betrán E, Long M: **Dntf-2r, a young *Drosophila* retroposed gene with specific male**
627 **expression under positive Darwinian selection.** *Genetics* 2003, **164**:977-988.
628 17. Park JM, Evertts AG, Levin HL: **The Hermes transposon of *Musca domestica* and its use as**
629 **a mutagen of *Schizosaccharomyces pombe*.** *Methods* 2009, **49**:243-247.
630 18. Edskes HK, Mukhamedova M, Edskes BK, Wickner RB: **Hermes transposon mutagenesis**
631 **shows [URE3] prion pathology prevented by a Ubiquitin-targeting protein: Evidence for**
632 **carbon/nitrogen assimilation cross talk and a second function for Ure2p in**
633 ***Saccharomyces cerevisiae*.** *Genetics* 2018, **209**:789-800.

- 634 19. Zhao S, Jiang E, Chen S, Gu Y, Shangguan AJ, Lv T, Luo L, Yu Z: **PiggyBac transposon vectors:**
635 **the tools of the human gene encoding.** *Transl Lung Cancer Res* 2016, **5**:120-125.
- 636 20. Li J, Zhang JM, Li X, Suo F, Zhang MJ, Hou W, Han J, Du LL: **A piggyBac transposon-based**
637 **mutagenesis system for the fission yeast *Schizosaccharomyces pombe*.** *Nucleic Acids Res*
638 2011, **39**:e40.
- 639 21. Segal ES, Gritsenko V, Levitan A, Yadav B, Dror N, Steenwyk JL, Silberberg Y, Mielich K,
640 Rokas A, Gow NAR, et al: **Gene essentiality analyzed by in vivo transposon mutagenesis**
641 **and machine learning in a stable haploid isolate of *Candida albicans*.** *MBio* 2018, **9**.
- 642 22. Gao J, Wang H, Li Z, Wong AH, Wang YZ, Guo Y, Lin X, Zeng G, Liu H, Wang Y, Wang J:
643 ***Candida albicans* gains azole resistance by altering sphingolipid composition.** *Nat*
644 *Commun* 2018, **9**:4495.
- 645 23. Giaever G, Nislow C: **The yeast deletion collection: a decade of functional genomics.**
646 *Genetics* 2014, **197**:451-465.
- 647 24. Spirek M, Benko Z, Carnecka M, Rumpf C, Cipak L, Batova M, Marova I, Nam M, Kim DU,
648 Park HO, et al: ***S. pombe* genome deletion project: an update.** *Cell Cycle* 2010, **9**:2399-
649 2402.
- 650 25. Breiman L: **Random Forests.** *Machine Learning* 2001, **45**:5-32.
- 651 26. Hughes TR, Roberts CJ, Dai H, Jones AR, Meyer MR, Slade D, Burchard J, Dow S, Ward TR,
652 Kidd MJ, et al: **Widespread aneuploidy revealed by DNA microarray expression profiling.**
653 *Nat Genet* 2000, **25**:333-337.
- 654 27. Chial HJ, Giddings TH, Jr., Siewert EA, Hoyt MA, Winey M: **Altered dosage of the**
655 ***Saccharomyces cerevisiae* spindle pole body duplication gene, NDC1, leads to aneuploidy**
656 **and polyploidy.** *Proc Natl Acad Sci U S A* 1999, **96**:10200-10205.

- 657 28. Stevens RC, Davis TN: **Mlc1p is a light chain for the unconventional myosin Myo2p in**
658 ***Saccharomyces cerevisiae*. *J Cell Biol* 1998, **142**:711-722.**
- 659 29. Matsumoto K, Uno I, Oshima Y, Ishikawa T: **Isolation and characterization of yeast**
660 **mutants deficient in adenylate cyclase and cAMP-dependent protein kinase. *Proc Natl***
661 ***Acad Sci U S A* 1982, **79**:2355-2359.**
- 662 30. Gangadharan S, Mularoni L, Fain-Thornton J, Wheelan SJ, Craig NL: **DNA transposon**
663 **Hermes inserts into DNA in nucleosome-free regions in vivo. *Proc Natl Acad Sci U S A***
664 **2010, **107**:21966-21972.**
- 665 31. Bronner IF, Otto TD, Zhang M, Udenze K, Wang C, Quail MA, Jiang RH, Adams JH, Rayner JC:
666 **Quantitative insertion-site sequencing (QIseq) for high throughput phenotyping of**
667 **transposon mutants. *Genome Res* 2016, **26**:980-989.**
- 668 32. **UCSC Genome Browser on *S. cerevisiae* Apr. 2011 (SacCer_Apr2011/sacCer3) Assembly**
669 **[\[http://genome-euro.ucsc.edu/cgi-](http://genome-euro.ucsc.edu/cgi-bin/hgTracks?db=sacCer3&lastVirtModeType=default&lastVirtModeExtraState=&virtModeType=default&virtMode=0&nonVirtPosition=&position=chrIX%3A183118%2D203117&hgslid=232744249_LxCF1Ofg7R9Jcsy2DHsyDa984DWn)**
670 **[bin/hgTracks?db=sacCer3&lastVirtModeType=default&lastVirtModeExtraState=&virtMode](http://genome-euro.ucsc.edu/cgi-bin/hgTracks?db=sacCer3&lastVirtModeType=default&lastVirtModeExtraState=&virtModeType=default&virtMode=0&nonVirtPosition=&position=chrIX%3A183118%2D203117&hgslid=232744249_LxCF1Ofg7R9Jcsy2DHsyDa984DWn)**
671 **[Type=default&virtMode=0&nonVirtPosition=&position=chrIX%3A183118%2D203117&hgslid=232744249](http://genome-euro.ucsc.edu/cgi-bin/hgTracks?db=sacCer3&lastVirtModeType=default&lastVirtModeExtraState=&virtModeType=default&virtMode=0&nonVirtPosition=&position=chrIX%3A183118%2D203117&hgslid=232744249_LxCF1Ofg7R9Jcsy2DHsyDa984DWn)**
672 **[_LxCF1Ofg7R9Jcsy2DHsyDa984DWn\]](http://genome-euro.ucsc.edu/cgi-bin/hgTracks?db=sacCer3&lastVirtModeType=default&lastVirtModeExtraState=&virtModeType=default&virtMode=0&nonVirtPosition=&position=chrIX%3A183118%2D203117&hgslid=232744249_LxCF1Ofg7R9Jcsy2DHsyDa984DWn)**

675 **Figure Legends**

676 **Figure 1. Overview of data acquisition and analysis: transposition events to gene essentiality.**

- 677 a. Three yeast species analyzed (Sp, *S. pombe*; Sc, *S. cerevisiae*; and Ca, *C. albicans*) by in vivo
678 transposition in this study and which transposons (PB, PiggyBac; AcDs and Hermes) were
679 used to mutagenize which species. Note that each species was analyzed with two different
680 transposon systems.

- 682 b. Comparison of genome insertion sites for transposition events that initiate from an
683 extrachromosomal plasmid (left, red region of plasmid circle) or a specific chromosomal
684 locus (right, red bar on a given chromosome). Horizontal lines represent multiple copies of
685 the same genome, each of which underwent a single insertion event (green arrow) per
686 genome. While transposition is generally random, a bias for loci in close proximity to the
687 initial transposon insertion site demands normalization of the final data.
- 688 c. Mapped Tnseq analysis of the pool of transposition events yields the chromosomal
689 insertion sites (brown vertical lines) in the reference chromosomes relative to the ORFs
690 (purple regions). A close up of a small region of a single chromosome (olive horizontal line)
691 including 5 ORFs is illustrated.
- 692 d. A training set is constructed using known or inferred labels (non-essential, blue; essential,
693 red) together with extracted features calculated from the data and its position relative to
694 ORFs. Features are defined in Table 2.
- 695 e. The training set features, as well as features for all ORFs are used as input for Random
696 Forest classification (black rectangle); output is a prediction of essentiality (red or blue as
697 in d), for which an optimal threshold is determined and applied to designate all genes in
698 one of the two categories.

700 *Figure 2.* Contribution of unique hits and total number of reads to the quality of ML predictions for
701 gene essentiality/non-essentiality.

702 Performance of the classifier vs (a) the total number of unique insertion sites (hits) and (b) the
703 total number of sequencing reads, in each study (organism abbreviations as in Fig. 1a; Ac, AcDs; H,
704 Hermes; PB, PiggyBac).

706 *Figure 3.* Feature importances in the different classifiers.

707 Importance of each feature used in the Random Forest classifier of essentiality for each dataset.

708 Features are described in Table 2; Neighborhood index generally normalizes for non-random

709 insertion frequencies across the genome; Freedom index reports on the largest proportion of an

710 ORF that has no hits, which is a measure of domains that may be essential.[21]

712 *Figure 4.* Analysis of the ability of the CaPB classifier to infer gene essentiality in genes with

713 increasing number of target sequences.

714 When only ORFs with a specific number of target sites are considered (\geq x-axis), AUC rises

715 accordingly (red), but the number of ORFs that can be analyzed necessarily decreases (numbers

716 above red dots). This demonstrates the importance of the prevalence of the transposon target

717 sequences in ORFs, for the quality of gene essentiality inference, using in-vivo transposon

718 mutagenesis studies. X axis: Minimum number of target sequences per ORF needed for inclusion

719 in the classification. Y-axis (red): CaPB classifier AUC.

721 *Figure 5.* Analysis of ROC AUC values for Random Forest classification trained on data from a

722 different organism and/or transposon in all possible combinations.

- 723 a. For each ROC AUC value in the table, training was performed on 80% of the original
724 training set used in the training species/transposon described in the rows. This training
725 data was then used to predict the essentiality of the remaining 20% of the training set in
726 the species/transposon described in the columns. The train/test split ratio was similar to
727 the 5-fold cross-validation performed in the single study analyses.

728 b. For each study, the vector of the relative feature importance was correlated with the
729 feature importance in all the other studies. Pearson r correlation coefficient values are
730 presented.

731

732 *Figure 6: Comparison with known essentials genes.*

733 a. Comparison of the essentiality verdicts in *S. cerevisiae*, based on the known essential genes
734 from the literature, ScAcDs and ScHermes classifiers.

735 b. Comparison of the essentiality verdicts in *S. pombe*, based on the known essential genes
736 from the literature, SpPB and SpHermes classifiers.

737 c. Comparison of the essentiality verdicts in *C. albicans*, based on the known essential Sp and
738 Sc orthologs from the literature, CaAcDs and CaPB classifiers.

739 Classification thresholds differ slightly from the previously published analyses [21] based on
740 threshold selection applied systematically to all 6 studies (described in detail in Methods).

741

742 *Figure 7. Gene essentiality in haploid and diploid S. cerevisiae.*

743 Comparison of essential genes in haploid and diploid *S. cerevisiae* analyzed with Sc Hermes. RF
744 classifier was trained on the haploid ScHermes study and predicted gene essentiality in a diploid
745 strain, using the same threshold for the final verdict. Mitochondrial genes were not considered.

746

747 *Figure 8. Threshold optimization.*

748 Two metrics were evaluated: 1) Minimum of the Euclidean distance between (0, 1) and the
749 receiver operating characteristic (ROC) curve. 2) Maximum of the vertical distance between the
750 line describing a random choice (a straight line from (0, 0) to (1, 1)) and the ROC curve.

751

752
753
754
755
756
757
758
759
760
761
762
763
764
765

752 *Figure 9. Suspected haploinsufficient genes in S. cerevisiae.*

753 Four genes suspected to be haploinsufficient in *S. cerevisiae*: *NDC1*, *MLC1*, *RPC10* and *BCY1*, as

2
3
754 they appear in the UCSD genome browser [32]. *NDC1* and *MLC1* were known to be

5
6
755 haploinsufficient. *BCY1* appears to be a previously unknown haploinsufficient gene. *RPC10* might

7
8
756 not be essential as it sustained hits in the 5' UTR in diploids but not haploids.

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Table 1. Statistics of the transposon data sets

	Ca AcDs	Sc AcDs	Sc Hermes	Sp Hermes	Ca PiggyBac	Sp PiggyBac
Transposon target sequence	-	-	TnnnnA	TnnnnA	TTAA	TTAA
Initial transposon insertion	Genome	Plasmid	Plasmid	Plasmid	Genome	Genome
Number of target sequences (10^3)	-	-	1154.84	1302.41	120.27	111.37
Total number of unique hits (10^3)	588.97	514.89	444.41	382.82	191.49	37.5
Target Sequences without a hit (10^3)	-	-	924.56	1110.2	10.69	79.09
Percent of target sequences without a hit	-	-	80.06%	85.24%	8.89%	71.02%
Percent of hits in target sequences	-	-	51.82%	50.21%	56.74%	86.06%
Total number of reads (10^6)	84.16	47.1	18.22	23.92	32.58	6.14
Average number of reads per hit	143	91	41	62	170	164
Standard deviation reads per hit	9250	4357	137	5069	4584	481
Highest reads per hit (10^3)	3254.83	2210.55	11.72	2355.61	1301.82	48.31
Highest reads per hit / average reads per hit	22761	24292	286	37994	7658	295
Number of over 106 reads per hit	10	2	0	2	1	0
Average number of hits per gene	29.11	44.05	30.36	24.64	11.81	1.32
Number of genes with 0 hits	300	261	332	99	397	2580
Average number of target sequences per gene	-	-	123.64	131.44	10.62	8.39
Number of genes with 0 target sequences	-	-	0	0	228	186
ROC AUC	0.993	0.985	0.972	0.962	0.94	0.785

Table 2. Classification features

Feature	Description
Hits	Number of transposon hits within the ORF
Reads	Number of reads associated with the transposon hits within the ORF
Neighborhood Index (NI)	Number of transposon hits within the ORF, normalized by length of the ORF and the surrounding 10 kbp
Freedom Index (FI)	Length of the largest hit-free region in the ORF, divided by the ORF's length
Hits 100 upstream	Number of transposon hits within the upstream region of the ORF
Hits per Target Seqs	Number of transposon hits divided by the number of transposon target sequences within an ORF
Reads per Length	Number of transposon hits divided by the length of the ORF
Hits per Length	Number of reads associated with the transposon hits divided by the length of the ORF

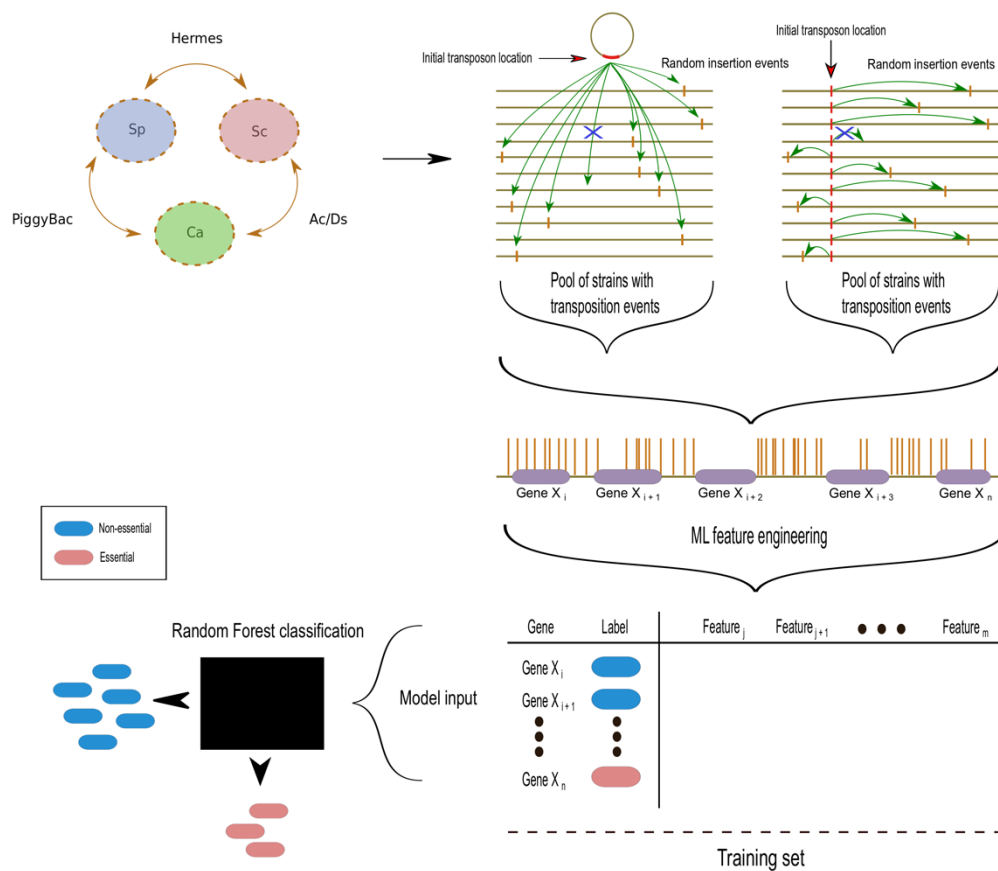


Figure 1. Overview of data acquisition and analysis: transposition events to gene essentiality.

- Three yeast species analyzed (Sp, *S. pombe*; Sc, *S. cerevisiae*; and Ca, *C. albicans*) by in vivo transposition in this study and which transposons (PB, PiggyBac; AcDs and Hermes) were used to mutagenize which species. Note that each species was analyzed with two different transposon systems.
- Comparison of genome insertion sites for transposition events that initiate from an extrachromosomal plasmid (left, red region of plasmid circle) or a specific chromosomal locus (right, red bar on a given chromosome). Horizontal lines represent multiple copies of the same genome, each of which underwent a single insertion event (green arrow) per genome. While transposition is generally random, a bias for loci in close proximity to the initial transposon insertion site demands normalization of the final data.
- Mapped Tnseq analysis of the pool of transposition events yields the chromosomal insertion sites (brown vertical lines) in the reference chromosomes

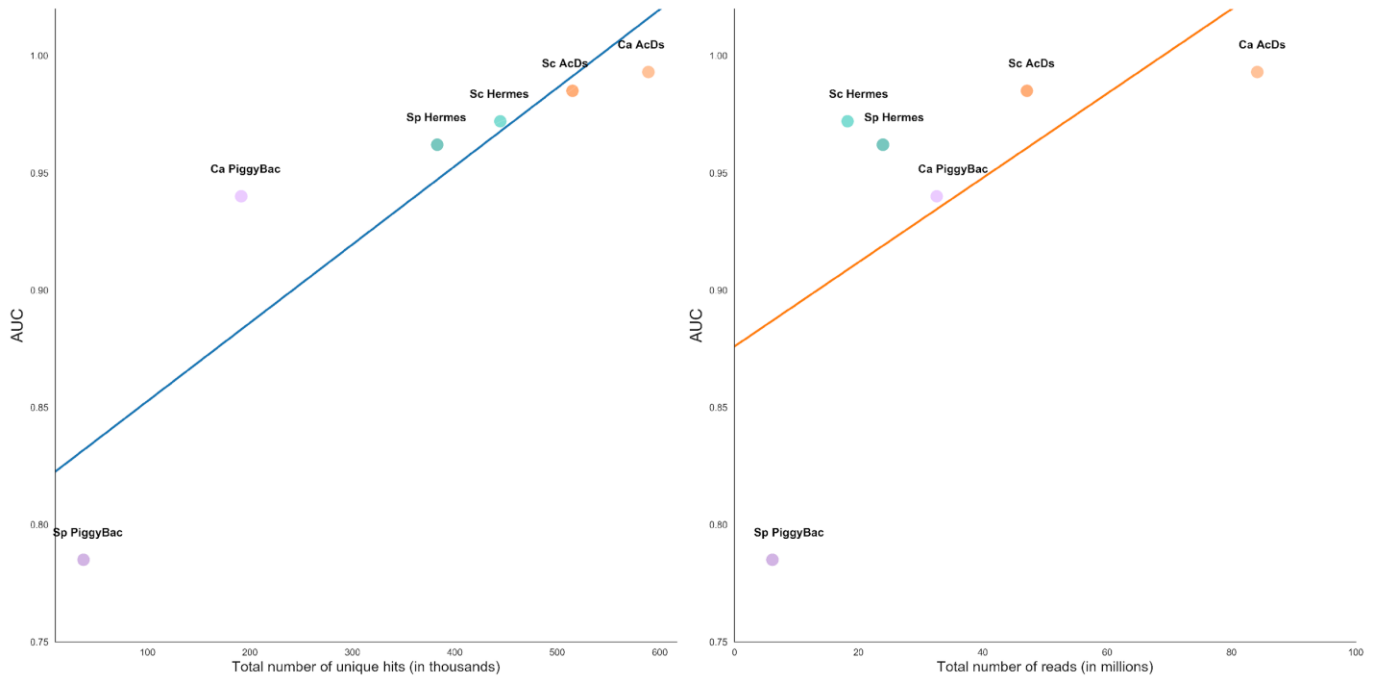
relative to the ORFs (purple regions). A close up of a small region of a single chromosome (olive horizontal line) including 5 ORFs is illustrated.

- D. A training set is constructed using known or inferred labels (non-essential, blue; essential, red) together with extracted features calculated from the data and its position relative to ORFs. Features are defined in Table 2.
- E. The training set features, as well as features for all ORFs are used as input for Random Forest classification (black rectangle); output is a prediction of essentiality (red or blue as in d), for which an optimal threshold is determined and applied to designate all genes in one of the two categories.

F.

Figure 2. Contribution of unique hits and total number of reads to the quality of ML predictions for gene essentiality/non-essentiality.

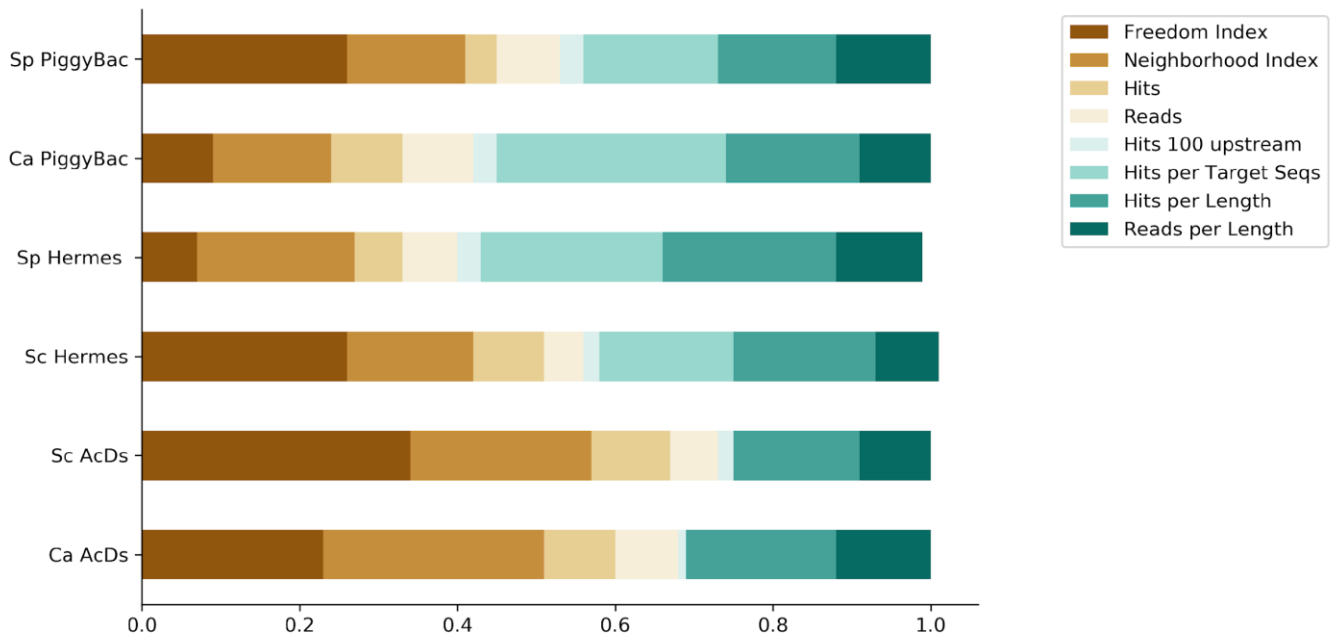
- Performance of the classifier vs (a) the total number of unique insertion sites (hits) and (b) the total number of sequencing reads, in each study (organism abbreviations as in Fig. 1a; Ac, AcDs; H, Hermes; PB, PiggyBac).



●

Figure 3. Feature importances in the different classifiers.

- Importance of each feature used in the Random Forest classifier of essentiality for each dataset. Features are described in [Table XXX](#); Neighborhood index generally normalizes for non-random insertion frequencies across the genome; Freedom index reports on the largest proportion of an ORF that has no hits, which is a measure of domains that may be essential.[21]



●

Figure 4. Analysis of the ability of the CaPB classifier to infer gene essentiality in genes with increasing number of target sequences.

- When only ORFs with a specific number of target sites are considered (\geq x-axis), AUC rises accordingly (red), but the number of ORFs that can be analyzed necessarily decreases (numbers above red dots). This demonstrates the importance of the prevalence of the transposon target sequences in ORFs, for the quality of gene essentiality inference, using in-vivo transposon mutagenesis studies.
- X axis: Minimum number of target sequences per ORF needed for inclusion in the classification.
- Y-axis (red): CaPB classifier AUC.

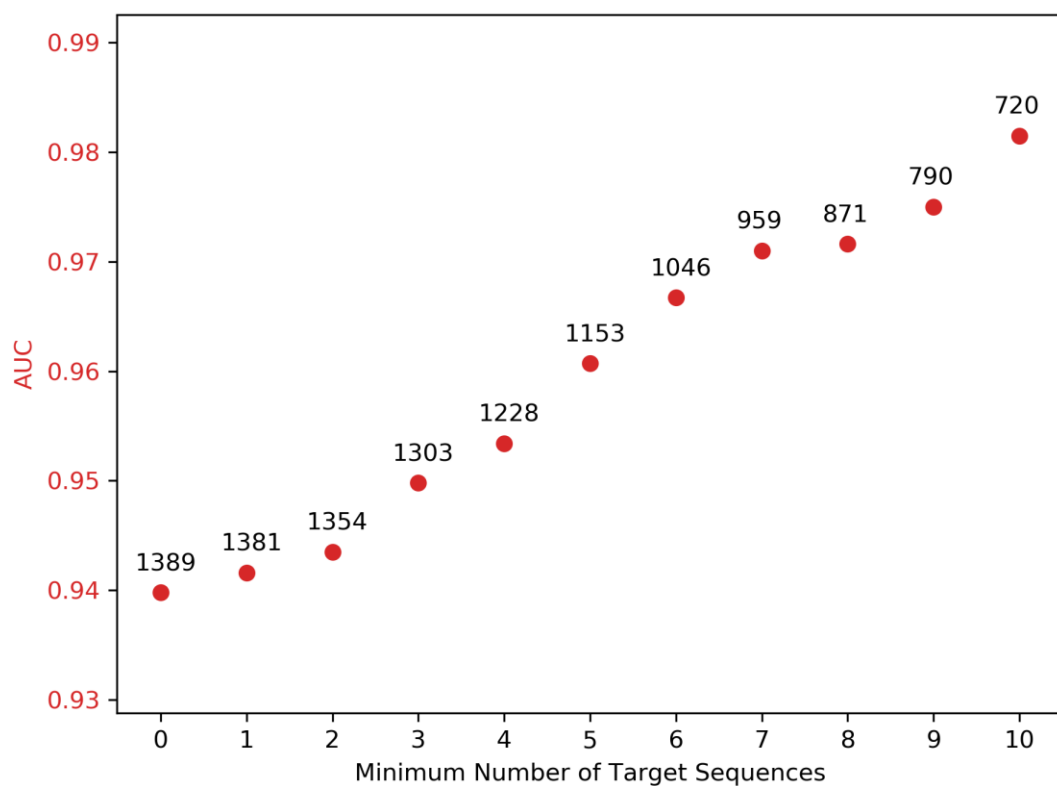
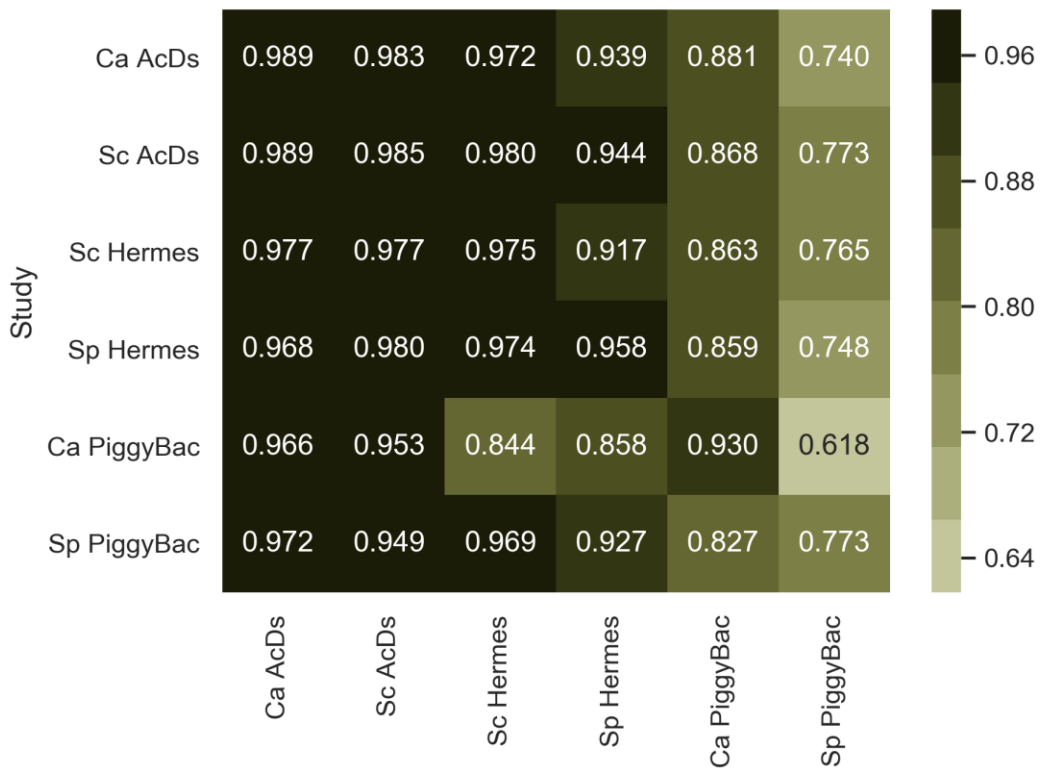


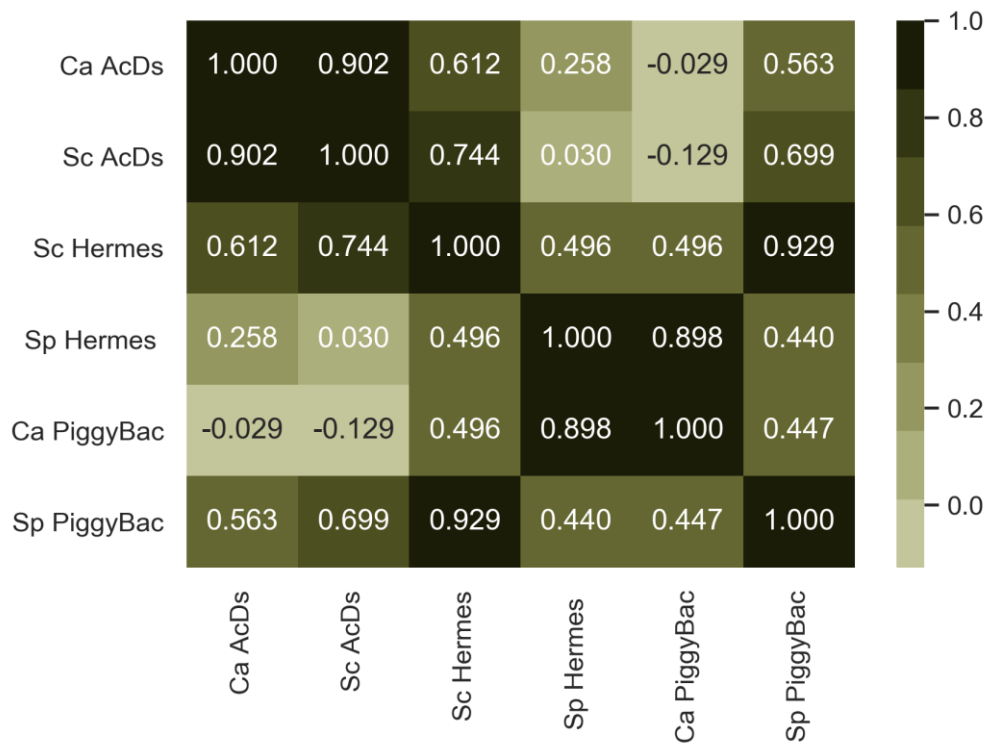
Figure 5. Analysis of ROC AUC values for Random Forest classification trained on data from a different organism and/or transposon in all possible combinations.

- A.** For each ROC AUC value in the table, training was performed on 80% of the original training set used in the training species/transposon described in the rows. This training data was then used to predict the essentiality of the remaining 20% of the training set in the species/transposon described in the columns. The train/test split ratio was similar to the 5-fold cross-validation performed in the single study analyses.
- B.** For each study, the vector of the relative feature importance was correlated with the feature importance in all the other studies. Pearson r correlation coefficient values are presented.

a. Transfer learning AUCs



b. Correlation of feature importance vectors

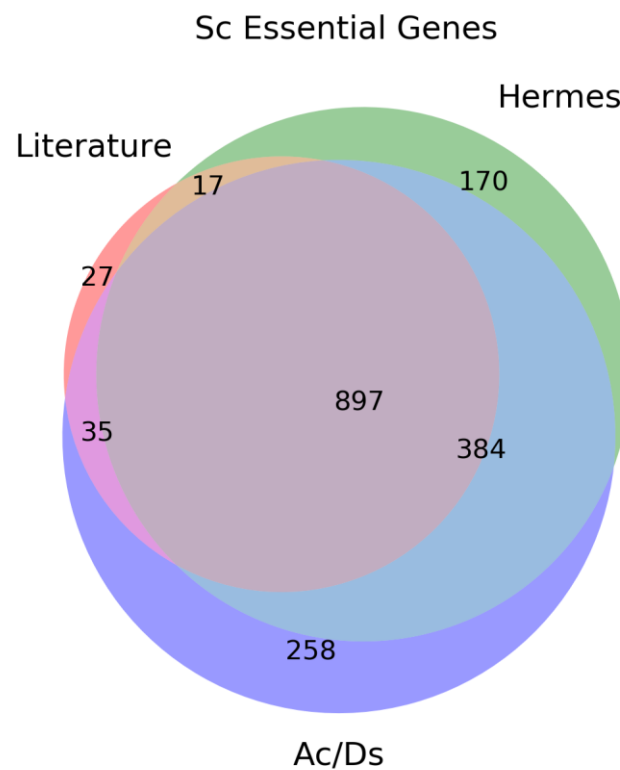


C.

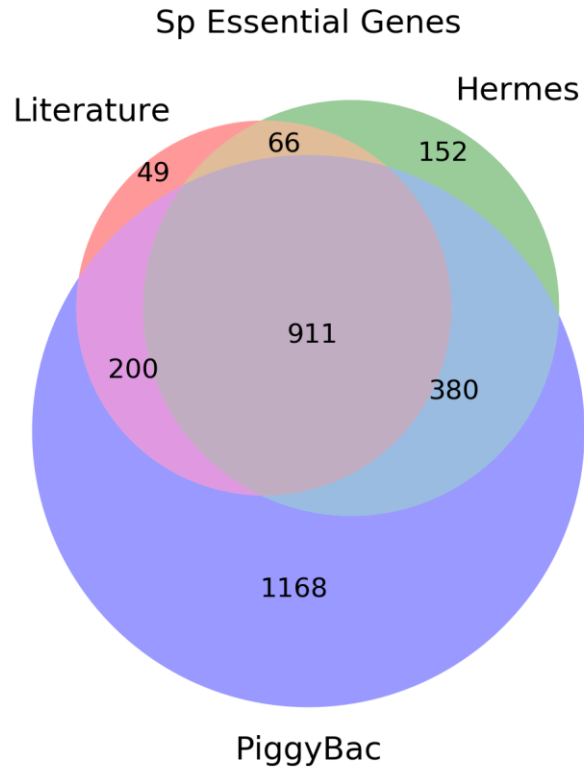
Figure 6: Comparison with known essentials genes.

- A. Comparison of the essentiality verdicts in *S. cerevisiae*, based on the known essential genes from the literature, ScAcDs and ScHermes classifiers.
 - B. Comparison of the essentiality verdicts in *S. pombe*, based on the known essential genes from the literature, SpPB and SpHermes classifiers.
 - C. Comparison of the essentiality verdicts in *C. albicans*, based on the known essential Sp and Sc orthologs from the literature, CaAcDs and CaPB classifiers.
- Classification thresholds differ slightly from the previously published analyses[21]based on threshold selection applied systematically to all 6 studies (described in detail in Methods).

a.



b.



c.

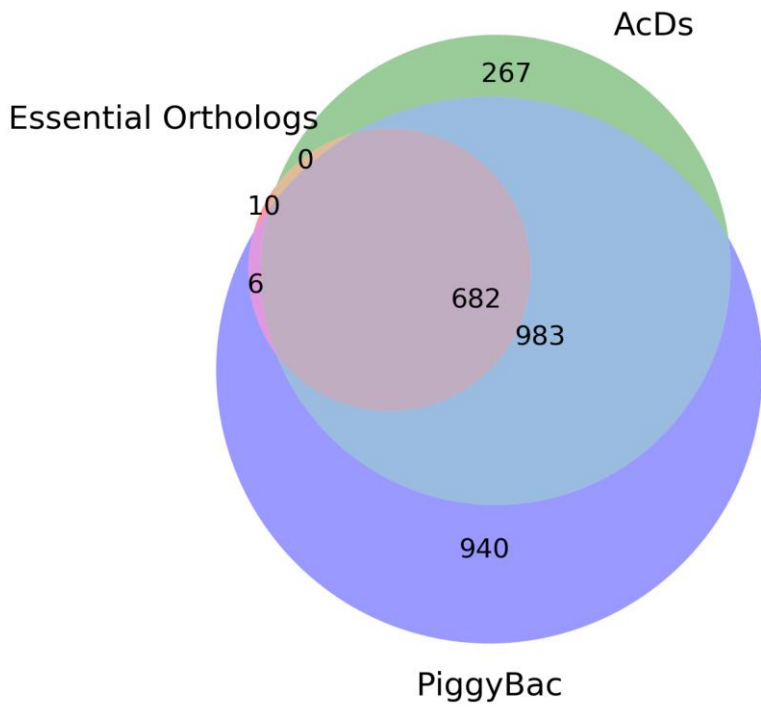


Figure 7. Gene essentiality in haploid and diploid *S. cerevisiae*.

- Comparison of essential genes in haploid and diploid *S. cerevisiae* analyzed with Sc Hermes. RF classifier was trained on the haploid ScHermes study and predicted gene essentiality in a diploid strain, using the same threshold for the final verdict. Mitochondrial genes were not considered.

Sc Essential Genes in Diploid and Haploid

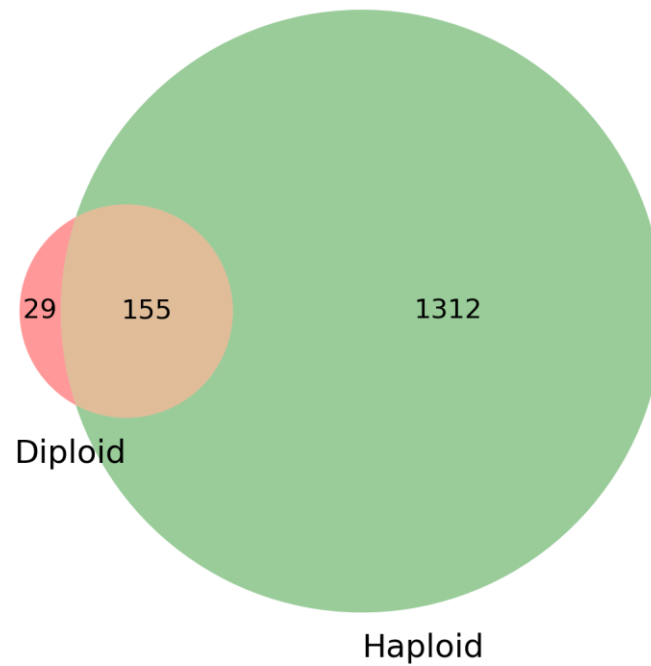
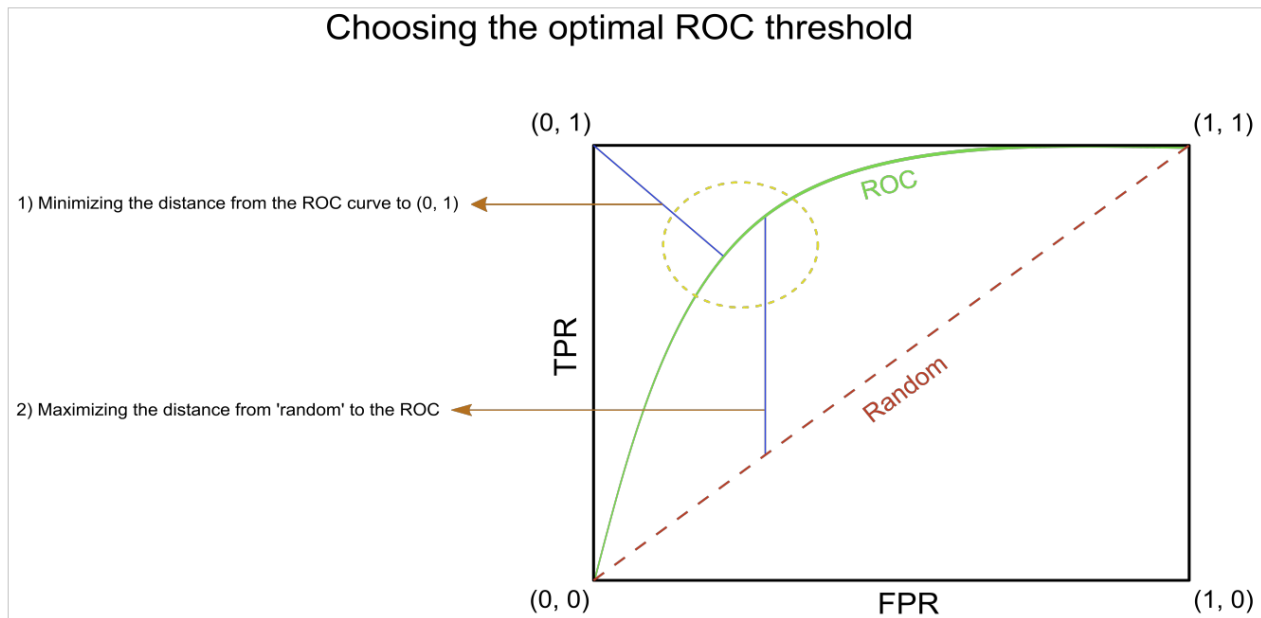
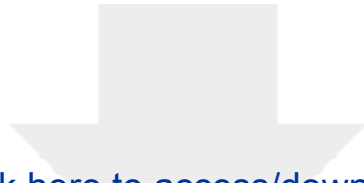


Figure 8. Threshold optimization.

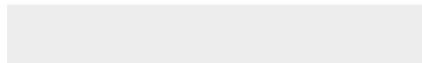
- Two metrics were evaluated: 1) Minimum of the Euclidean distance between (0, 1) and the receiver operating characteristic (ROC) curve. 2) Maximum of the vertical distance between the line describing a random choice (a straight line from (0, 0) to (1, 1)) and the ROC curve.



-

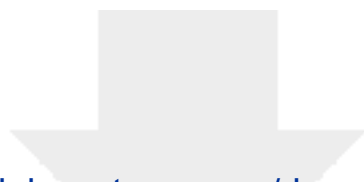


Click here to access/download
Supplementary Material
Supplementary Information v3.docx

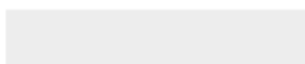


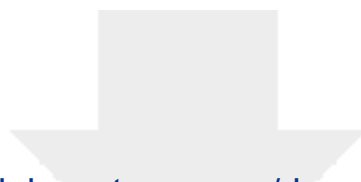


Click here to access/download
Supplementary Material
Table S1 Training Sets.csv



Click here to access/download
Supplementary Material
Table S2 CaAcDs ml predictions.csv

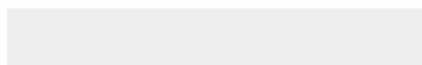
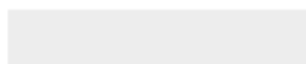




[Click here to access/download](#)

Supplementary Material

Table S3 ScAcDs ml predictions.csv

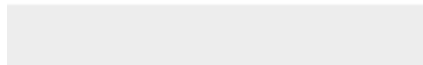




[Click here to access/download](#)

Supplementary Material

Table S4 ScHermes ml predictions.csv

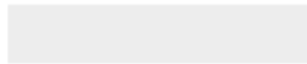


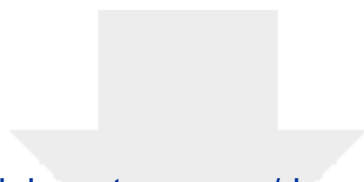


[Click here to access/download](#)

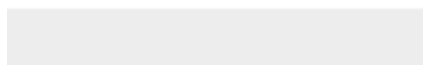
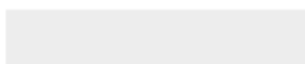
Supplementary Material

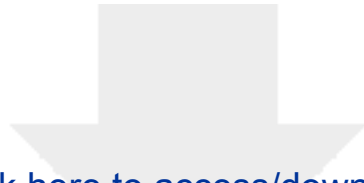
Table S5 SpHermes ml predictions.csv



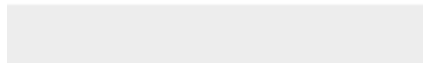



Click here to access/download
Supplementary Material
Table S6 CaPB ml predictions.csv



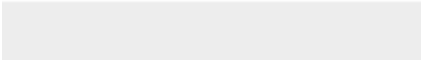



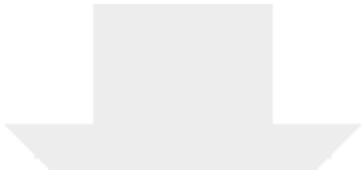
Click here to access/download
Supplementary Material
Table S7 SpPB ml predictions.csv






Click here to access/download
Supplementary Material
Figure S1 Ca AcDs.png





Click here to access/download
Supplementary Material
Figure S2 Sc AcDs.png





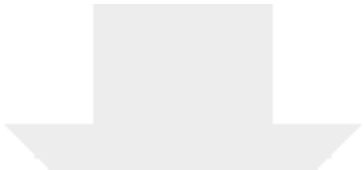
Click here to access/download
Supplementary Material
Figure S3 Sc Hermes.png






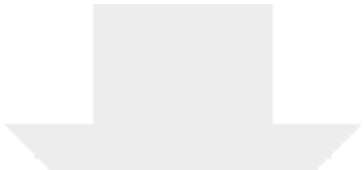
Click here to access/download
Supplementary Material
Figure S4 Sp Hermes.png






Click here to access/download
Supplementary Material
Figure S5 Ca PB.png





Click here to access/download
Supplementary Material
Figure S6 Sp PB.png





[Click here to access/download](#)

Supplementary Material

Figure S7 ScHermes target sequence frequency.png

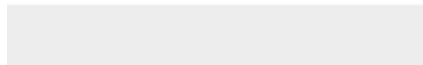




[Click here to access/download](#)

Supplementary Material

Figure S8 SpHermes target sequence frequency.png

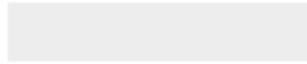




[Click here to access/download](#)

Supplementary Material

Figure S9 CaPB target sequence frequency.png





[Click here to access/download](#)

Supplementary Material

Figure S10 SpPB target sequence frequency.png

