1    sppIDer: a species identification tool to investigate hybrid genomes with high-throughput

2    sequencing

3    Quinn K. Langdon[1,2]; David Peris[1,2,3,4]; Brian Kyle[2]; Chris Todd Hittinger[1,2,3]

4    1 Laboratory of Genetics, J. F. Crow Institute for the Study of Evolution, Genome Center of

5    Wisconsin, University of Wisconsin–Madison, Madison, WI USA

6    2 Wisconsin Energy Institute, University of Wisconsin–Madison, Madison, WI USA

7    3 DOE Great Lakes Bioenergy Research Center, University of Wisconsin–Madison, Madison,

8    WI USA

9    4 Department of Food Biotechnology, Institute of Agrochemistry and Food Technology (IATA),

10   CSIC, Valencia Spain

11   **Abstract:**

12          The genomics era has expanded our knowledge about the diversity of the living world,

13   yet harnessing high-throughput sequencing data to investigate alternative evolutionary

14   trajectories, such as hybridization, is still challenging. Here we present sppIDer, a pipeline for

15   the characterization of interspecies hybrids and pure species, that illuminates the complete

16   composition of genomes. sppIDer maps short-read sequencing data to a combination genome

17   built from reference genomes of several species of interest and assesses the genomic contribution

18   and relative ploidy of each parental species, producing a series of colorful graphical outputs

19   ready for publication. As a proof-of-concept, we use the genus *Saccharomyces* to detect and

20   visualize both interspecies hybrids and pure strains, even with missing parental reference

21   genomes. Through simulation, we show that sppIDer is robust to variable reference genome

22   qualities and performs well with low-coverage data. We further demonstrate the power of this

23   approach in plants, animals, and other fungi. sppIDer is robust to many different inputs and

24   provides visually intuitive insight into genome composition that enables the rapid identification

25   of species and their interspecies hybrids. sppIDer exists as a Docker image, which is a reusable,

26   reproducible, transparent, and simple-to-run package that automates the pipeline and installation

27   of the required dependencies (https://github.com/GLBRC/sppIDer).

28

29   **Introduction:**

30          Interspecies hybrids play a large role in both natural and in industrial settings (Dunn and

31   Sherlock 2008; Soltis et al. 2015; Payseur and Rieseberg 2016; Peris et al. 2017c). However,

1

32  identification and characterization of the genomic contributions of hybrids can be difficult. High-

33  throughput sequencing can be used to address many of the barriers to identifying and

34  characterizing hybrids. With the influx of sequencing data, the quality and number of reference

35  genomes available is increasing at a rapid pace. Population genomic, ecological diversity, and

36  gene expression projects are underway in many fields. These studies are yielding a high volume

37  of short-read data, but determining the best way to leverage these data can be challenging. A key

38  goal of the modern genomic era is to be able to integrate and synthesize these data to further our

39  understanding of natural diversity (Richards 2017), including addressing key questions about the

40  frequency and genomic identities of hybrid and admixed lineages in the wild.

41      The number of reference genomes available has rapidly increased, but it is not complete

42  in most clades. To avoid the drawbacks of limited reference genomes, several new phylogenetic

43  approaches have been developed that do not require sequence alignments or whole-genome

44  assemblies, such as phylogeny-building approaches using kmers (Fan et al. 2015), de novo

45  identification of phylogenetically informative regions (Schwartz et al. 2015), and local

46  assemblies of target genes (Allen et al. 2015; Johnson et al. 2016). These methods can accurately

47  reconstruct known and simulated phylogenies of pure lineages. However, these methods have not

48  been tested on hybrid or admixed lineages. As hybrids are the result of an outcrossing event

49  between two independently evolving lineages, their origin is inherently not tree-like. Therefore,

50  placing hybrids on a bifurcating tree will not reflect the topology observed with pure lineages.

51  Placing hybrids on a phylogenetic network is more apt, but it is still untested with alignment-free

52  phylogenetic approaches. Other species identification methods based on local assembly of target

53  genes could lead to erroneous identification, depending on which parent the gene of interest is

54  retained from in the hybrid, or could lead to the assembly of a chimeric gene if the hybrid has

55  retained copies from multiple parents. Therefore, in organisms with alternative evolutionary

56  trajectories, such as hybrids with complex genomes, applying alignment-free phylogenetic

57  methods is difficult and could potentially result in imprecise conclusions.

58      Other methods to detect interspecies hybrids have been adapted from methods developed

59  for intraspecies diversity, such as $F_{ST}$, STRUCTURE analysis, phylogenetic discordance, linkage

60  disequilibrium, and PCA approaches (Payseur and Rieseberg 2016). There are numerous

61  drawbacks to using these methods to detect interspecies hybrids. For example, most definitions

62  of speciation require the cessation of gene flow and the accumulation of sequence divergence

2

63 well beyond the levels observed between populations, which are therefore beyond the

64 expectations of most of these approaches. Many of these methods were also developed for

65 diploid obligately outcrossing species, which makes problematic their application to

66 allopolyploids or species that primarily undergo selfing or other forms of inbreeding. Indeed, the

67 basic assumptions of these methods, including gene flow, demographic history, and natural

68 selection, are violated by most interspecies hybrids.

69 Here we present sppIDer as a novel, assumption-free method that rapidly provides visual

70 and intuitive insight into ancestry genome-wide, which will aid in the discovery and

71 characterization of interspecies hybrids. This method maps short-read data to combination

72 genome, built from available reference genomes chosen by the user. sppIDer allows for the

73 analysis and visualization of the genomic makeup of a single organism of interest, facilitating the

74 rapid discovery of hybrids and individuals with other unique genomic features, such as

75 aneuploidies and introgressions. Therefore, sppIDer is an unbiased method that provides unique

76 and intuitive insights into complex genomic ancestry and regions of differing evolutionary

77 history, which can complement existing methods in the characterization of hybrids.

78

79 **New Approaches**

80 Here we describe and make available a user-friendly short-read data analysis pipeline that

81 utilizes existing bioinformatic tools and custom scripts to determine species identity, hybrid

82 status, and chromosomal copy-number variants (CCNVs). Short-reads are mapped to a

83 combination reference genome of multiple species of interest, and the output is parsed for where,

84 how well, and how deeply the reads map across this combination genome. A colorful automated

85 output allows end-users to rapidly and intuitively assess the genomic contribution, either from a

86 single species or multiple species, and relative ploidy of an organism. Figure 1 illustrates the

87 basic workflow in a flow chart of each step. An upstream step creates a combination reference

88 genome, which is a concatenation of reference genomes of interest, before the main pipeline is

89 run. The main pipeline starts with mapping short-read data to this combination reference

90 genome. Then, this output is parsed for percentage and quality of reads that map to each

91 individual reference genome within the combination reference and percentage of unmapped

92 reads; this summary is then plotted so these metrics can be visualized. In parallel, the mapping

93 output is analyzed for depth of coverage. Reads with a mapping quality (MQ) greater than three

94    are retained and sorted into the combination reference genome order; then, coverage across the

95    combination reference genome is computed. A custom script then calculates the mean coverage

96    for each species, and the combination reference genome broken into windows. The output of

97    these analyses is then plotted so that coverage across the combination reference genome can be

98    visualized.

99         We have given this computational pipeline and wrapper a portmanteau of the pluralized

100    abbreviation of species (spp.) and identifier (IDer), to reflect its ability to identify hybrids of

101    multiple species. sppIDer also detects CCNVs, such as those caused by aneuploidy and other

102    genomic changes that do not meet the textbook definition of aneuploidy, including interspecies

103    loss-of-heterozygosity events, interspecies unbalanced translocations, and other differences in

104    relative ploidy. sppIDer is provided as an open source Docker (http://www.docker.com), which

105    organizes the pipeline and all the dependencies into a reusable, reproducible, transparent, and

106    simple-to-run package (https://github.com/GLBRC/sppIDer).

107         Here we present several applications of sppIDer in yeast, plant, and animal genomes.

108    Through simulations, we show that sppIDer can detect hybrids of closely or distantly related

109    species, and of recent or ancient origin. We use the genus *Saccharomyces* to 1) detect both

110    interspecies hybrids and pure strains; 2) detect hybrids, even with missing reference genomes;

111    and 3) determine how divergent lineages and poor-quality data and reference genomes affect

112    sppIDer's performance. Next, we test sppIDer's utility in non-*Saccharomyces* systems: another

113    yeast genus, *Lachancea*; an animal genus, *Drosophila*; and a plant genus, *Arabidopsis*. Finally,

114    we test an extension for non-nuclear DNA using mitochondrial genome data. Overall, sppIDer is

115    robust to many different inputs and can be used across organisms to provide rapid insight into the

116    species identity, hybrid status, and CCNVs of an organism.

117

118    **Results and Discussion:**

119    <u>Species and interspecies hybrid identifications:</u>

120         To test sppIDer, we first used the well-studied genus *Saccharomyces* (Hittinger 2013).

121    Seven of the eight species have reference genomes scaffolded at a near-chromosomal level, and

122    there are many interspecies hybrids (Goffeau et al. 1996; Fischer et al. 2000; Dunn and Sherlock

123    2008; Liti and Carter et al. 2009; Scannell and Zill et al. 2011; Liti et al. 2013; Baker et al. 2015;

124    Naseeb et al. 2017; Peris et al. 2017c). To test sppIDer's species-level classification ability for a

125    natural isolate, we used the short-read data available for a *Saccharomyces eubayanus* strain

126    isolated in New Zealand (P1C1) (Gayevskiy and Goddard 2016). The reads from this wild *S.*

127    *eubayanus* strain mapped preferentially to the *S. eubayanus* reference genome (Figure 2a), as

128    seen by normalized coverage only being above zero for the *S. eubayanus* genome. This strain

129    belongs to the same diverse lineage as the reference strain (Peris and Langdon et al. 2016), but as

130    the first isolate from Oceania, these results show that sppIDer can easily classify, to the species

131    level, a divergent wild strain isolated from a novel environment. To test sppIDer's utility for

132    industrial strains, we used short reads from an ale strain, Fosters O (Gonçalves et al. 2016). This

133    test shows that this brewing strain is a pure species; the *S. cerevisiae* genome is the only genome

134    that had normalized coverage above zero. However, normalized coverage differed within the *S.*

135    *cerevisiae* genome (Figure 2b & Figure S1a), implying aneuploidies. Coverage was lower for

136    chromosomes VII and XIV and increased for chromosome XIII, in comparison with the genome-

137    wide average coverage, indicating that there are more copies of chromosome XIII and fewer

138    copies of chromosome VIII and XIV. Additionally, we could detect regions of CCNV within a

139    chromosome, such as the small region within chromosome VII where the normalized coverage

140    returned to the genome average.

141          To test sppIDer's ability to delineate hybrids, we used short-read data from two *S.*

142    *cerevisiae* X *S. eubayanus* lager yeast lineages, Saaz (strain CBS1503) and Frohberg (strain

143    W34/70). These results recapitulated the known relative ploidy and rearrangements, where

144    ploidy differs both within and between genomes. Specifically, the Frohberg lineage contains

145    approximately two copies of each chromosome from both *S. cerevisiae* and *S. eubayanus*. Thus,

146    what was observed matched this expectation, where the average normalized coverage across both

147    the *S. cerevisiae* and *S. eubayanus* genomes were approximately at the same level, but there were

148    clear fluctuations, indicating ploidy changes (Figure 2c & Figure S1b). In our test with a

149    representative of the Saaz lineage, we observed that the *S. cerevisiae* genome had an average

150    normalized coverage of ~0.5, that fluctuated from none to two, and the *S. eubayanus* genome had

151    an average normalized coverage of 1.5, that fluctuated from zero to three (Figure 2d & Figure

152    S1c). These results match with previous observations that the Saaz lineage is approximately

153    haploid for the *S. cerevisiae* genome and diploid for the *S. eubayanus* genome. Additionally,

154    from the sppIDer plots, we also easily inferred the previously described aneuploidies and

155    translocations (Figure 2c-d) (Dunn and Sherlock 2008; Okuno et al. 2016).

156      As an additional hybrid test, we used short-read data from the wine strain Vin7, a *S.*
157   *cerevisiae* X *Saccharomyces kudriavzevii* hybrid. From the normalized coverage plot (Figure
158   2e), we could determine that Vin7 has retained complete copies of both parental genomes, but at
159   different ploidy levels. Specifically, the normalized coverage for *S. cerevisiae* was around two
160   across the genome, while the normalized coverage for *S. kudriavzevii* was consistently around
161   one across the genome. Here we could infer that this strain has double the number of copies of *S.*
162   *cerevisiae* chromosomes as it does of *S. kudriavzevii* chromosomes. Although exact ploidy
163   cannot be measured without direct measures of DNA content, the inferred ploidy is consistent
164   with previous studies (Borneman et al. 2012; Peris et al. 2012; Borneman et al. 2016).
165      As a final test of interspecies hybrids, we used data from the cider strain CBS2834
166   (Almeida et al. 2014). Here sppIDer detected large genetic contributions from *S. cerevisiae*, *S.*
167   *kudriavzevii*, and *Saccharomyces uvarum,* as well as introgressed contributions from *S.*
168   *eubayanus* (Figure 2f & Figure S1d). Although the *S. eubayanus* genetic contribution is quite
169   small, seen on chromosomes XII and XIV, it was still easily detected by sppIDer. These
170   examples show that sppIDer can easily detect higher-order interspecies hybrids, even those with
171   minor contributions from several species.
172
173   <u>Testing the limits of sppIDer with a simulated phylogeny:</u>
174      To test sppIDer's performance with hybrids of varying levels of parental divergence, we
175   used a simulated phylogeny. To build this phylogeny we started with the *S. cerevisiae* reference
176   genome and produced a phylogeny of 10 species through several rounds of simulating short-read
177   sequencing data, applying a set mutation rate, and assembling those reads. For these simulated
178   genomes, sister species were ~4% diverged, and the most distantly related species were ~20%
179   diverged (Figure 3a). This simulated phylogeny allowed us to test pseudo-hybrids from closely
180   and distantly related lineages. Further, the iterative process of phylogeny building allowed us to
181   create ancient pseudo-hybrids that simulated the result from hybridization of a common ancestor
182   predating a lineage split. sppIDer accurately mapped pure lineages to their corresponding
183   reference genome (Figure 3b). For all 10 species, >90% of the reads mapped to their
184   corresponding reference genome. The read simulation and assembly process resulted in varying
185   quality final references, but despite differences in genome quality, all reads still mapped
186   accurately and were not biased to the best reference genome.

187          To determine sppIDer's applicability to hybrids of both closely and distantly related

188    parents and of recent and ancient origin, we tested sppIDer with pseudo-hybrids of different

189    combinations of simulated species. sppIDer accurately detected all true hybrid parents. When

190    pseudo-hybrids were between sister species, <0.01% of the reads mapped promiscuously to other

191    species (Figure 3c). When we used more divergent pseudo-hybrids, sppIDer still detected the

192    true parents, with <5% of the reads mapped promiscuously to the sister species (Figure 3d).

193    Additionally, we simulated ancient pseudo-hybrids, between common ancestors before lineage

194    splits, and found that sppIDer mapped the reads of these hybrids to the references of the lineages

195    that descended from the ancestors that hybridized (Figure 3e). With complete knowledge of this

196    simulated phylogeny, we were able to test many different potential hybrid arrangements and

197    found that sppIDer detected the true parents of all hybrids.

198          Finally, we tested a scenario, which is common in biology, of incomplete knowledge of

199    the clade of interest. This dearth could due to many variables, such as a described species lacking

200    a reference genome or a species being unknown to science altogether. To test the effect of

201    missing a species, we removed one species' reference genome from the combination reference

202    genome, then mapped pure lineage and pseudo-hybrid reads to this permuted genome. With

203    reads of a simulated pseudo-hybrid of sister species, G and H, we observed that, when one parent

204    genome was missing, the reads mapped primarily to the reference genome of the remaining

205    parent, reference H, with slightly increased promiscuous mapping of reads to the next-closest

206    clade, references I and J (Figure 3f). Therefore, with incomplete reference genome knowledge,

207    detecting hybrids of closely related species is limited. However, we could still detect hybrids of

208    more distantly related species, such as a pseudo-hybrid of E and G and a pseudo-hybrid of the

209    common ancestor of A and the common ancestor of G and H (Figure S2), though our inference

210    of parentage was biased by the availability of reference genomes. Therefore, with incomplete

211    knowledge of reference genomes, hybrid detection is limited, and the inference of true parentage

212    can suffer in specific cases, but generally, distant and ancient hybrids can be detected.

213

214    <u>Hybrid detection with missing reference genomes:</u>

215          To empirically address how sppIDer would be affected by missing reference genomes,

216    such as for hybrids whose parents are themselves unknown (Hoot et al. 2004; Pryszcz et al.

217    2014), we focused again on the genus *Saccharomyces*. Specifically, we used the *S. cerevisiae* X

7

218    *S. kudriavzevii* hybrid (Vin7) and the *S. cerevisiae* X *S. eubayanus* Frohberg lager yeast

219    (W34/70) as examples. We tested the performance of sppIDer on short-read data from both

220    hybrids by removing the *S. cerevisiae* reference genome and, in a separate test, removing the

221    reference genome of the other parent. Our expectation was that reads would map to the genome

222    of the sister species, if it were available, or that they would fail to map or be distributed across

223    other genomes, if there were no close relatives.

224      When we removed the *S. eubayanus* reference genome for the lager example, the

225    proportion of reads that failed to map increased, as did those reads that mapped to *S. uvarum*, its

226    sister species (~93% identical in DNA sequence, Libkind and Hittinger et al. 2011), albeit with a

227    decreased mapping quality (MQ) (Figure 4c). We then tested sppIDer on Vin7 and W34/70 when

228    the *S. cerevisiae* reference genome was removed (Figure 4a & d). In both examples, the

229    proportion of reads that mapped to *Saccharomyces paradoxus*, *S. cerevisiae's* sister species

230    (~87% identical in DNA sequence), increased (Figure 4a & d). Thus, the absence of a reference

231    genome for one of the parents of a hybrid led to increased mapping to its sister species, instead.

232    We also tested removing the *S. kudriavzevii* reference genome for Vin7. Since there is not a

233    sister species closely related to *S. kudriavzevii*, the number of unmapped reads increased, and the

234    remaining reads mapped to the reference genomes other species of the genus in approximately

235    equal proportions (Figure 4f).

236      From these tests, we would have easily inferred that W34/70 was a hybrid, regardless of

237    whether either parent genome was withheld (the actual state of affairs for *S. eubayanus* before

238    Libkind and Hittinger et al. 2011). Using the coverage plots, we were still even able to infer the

239    same CCNVs for W34/70 that we observed with the full suite of reference genomes. With Vin7,

240    we still easily inferred its hybrid status without including the *S. cerevisiae* genome. Without the

241    *S. kudriavzevii* reference genome, Vin7 produced an unusually high number of unmapped reads

242    without a decrease in mapping quality to *S. cerevisiae*, a result that should spur the investigator

243    to perform more detailed analyses to search for evidence of contributions by an unknown

244    species, such as de novo genome assembly and phylogenetics. Therefore, even without a full

245    complement of reference genomes, sppIDer can still be useful for rapid inference of interspecies

246    hybrids.

247

248    Hybrid detection with simulated low-quality reference genomes:

249        To test a scenario where not all of the reference genomes are ideal, we used iWGS (Zhou

250    et al. 2016) to independently simulate reads and then assemble de novo genomes for *S.*

251    *cerevisiae*, *S. kudriavzevii*, *S. uvarum*, and *S. eubayanus*. These simulations resulted in reference

252    genomes with many more scaffolds and with a lower N50 than the published genomes (Table

253    S1). These low-quality genomes were independently swapped for the high-quality references in

254    the combination reference genome and tested with short-read data. We started by testing

255    simulated pseudo-lager short reads where we expected reads to map both to the *S. cerevisiae* and

256    *S. eubayanus* reference genomes. Whether we swapped in the low-quality *S. cerevisiae* reference

257    (Figure S3a) or the low-quality *S. eubayanus* reference (Figure S3b), the reads still mapped

258    equally to the references that were used to simulate the reads with minimal promiscuously

259    mapped reads to their sister species reference genomes.

260        We next tested the limits of sppIDer with the empirical data for CBS2834 because it has

261    the most complex arrangement of contributions from four species. Tests with each simulated

262    low-quality reference genome independently showed that we could indeed recapitulate the same

263    inference of ancestry and that roughly the same proportion of reads mapped to each reference

264    genome (Figure S4) as with high-quality reference genomes (Figure S1d). Here, the inference of

265    approximate ploidy became more difficult, and visually interpreting translocations between

266    species was impossible. When both high-quality *S. cerevisiae* and *S. kudriavzevii* reference

267    genomes were used, we could infer translocations between these two genomes on chromosomes

268    IV, X, and XV due to mid-chromosome ploidy changes that are compensated for in the other

269    genome. There were more promiscuously mapped reads to the high-quality reference genomes of

270    the sister species, but not at the same level as mapped to the true parent reference genomes.

271    These tests with simulated low-quality de novo genomes showed that, both with simulated and

272    empirical data, proper hybrid genome contributions can still be identified, and ploidy shifts still

273    detected, despite the poor-quality reference genomes, but the inference of translocations and

274    ploidies of specific chromosomes becomes difficult.

275

276    <u>Hybrid detection with low-coverage and long-read data:</u>

277        To further explore the power of sppIDer, we wanted to test how little coverage was

278    needed to still detect the proper ancestry (Figure S5). Using data simulated at varying coverages,

279    we found that only 0.5X coverage was needed to recover the true ancestry for a single species

280    (Figure S5a-b), single species with aneuploidies (Figure S5c-d), and interspecies hybrids (Figure

281    S5e-f). We also tested empirical data by down-sampling the FASTQ files of CBS2834 and found

282    that we could still detect contributions from the four species at as low as ~0.05X coverage

283    (Figure S5g), but we lost the ability to infer ploidy at around ~0.5X coverage (Figure S5h).

284    These low coverage tests show how powerful sppIDer is, even with scant data, which could be a

285    boon in many systems with large genomes or when sequencing resources are limited.

286          We also tested sppIDer with simulated PacBio long-read data from the *S. cerevisiae*

287    genome and a hybrid pseudo-lager genome with equal contributions from the *S. cerevisiae* and *S.*

288    *eubayanus* reference genomes. We found that we could still easily determine the species

289    contribution for each (Figure S6), suggesting sppIDer's utility will continue if long-read

290    technologies eventually supplant short-read sequencing technologies.

291

292    <u>Divergent lineages and poor-quality data:</u>

293          Since sppIDer relies on reference genomes, we recognized that it might be biased in its

294    ability to work with lineages that were highly divergent from the reference genome, as might be

295    the case in many systems. We tested this scenario with an example from *S. paradoxus*, one of the

296    most diverse *Saccharomyces* species (Liti and Carter et al. 2009; Leducq et al. 2016). Compared

297    to a representative of the reference genome's lineage, fewer reads from the divergent lineage

298    (~96% identical) mapped and with poorer quality (Figure S7a-b). We also tested this effect in *S.*

299    *kudriavzevii* using poor-quality data (36-bp reads from a first-generation Illumina Genome

300    Analyzer run by Hittinger et al. 2010) and found qualitatively similar results, but many more

301    unmapped reads. Thus, while divergence from the reference genome affected map-ability,

302    sppIDer still worked generally as expected. However, when mapping percentage and quality

303    decline substantially, such as seen in these test cases, sppIDer can provide an early indication

304    that the organism may be highly divergent from the reference genome, which may merit further

305    investigation.

306

307    <u>Comparison to alignment-free phylogenetic methods:</u>

308          Alignment and assembly (AA)-free phylogeny-building methods are gaining popularity,

309    but they have not previously been applied to hybrid data. Therefore, we also tested how AA-free

310    phylogenetic methods, such as AAF (Fan et al. 2015) or SISRS (Schwartz et al. 2015),

311    performed in detecting and visualizing hybrids compared to sppIDer. We found that these

312    methods performed well when given only pure lineages, but when hybrids were included, they

313    either failed completely or produced incorrect phylogenies. We tested both our simulated

314    phylogeny and empirical *Saccharomyces* data. For the simulated data, both AAF and SISRS

315    produced the correct phylogeny when given the 10 simulated species. However, when given any

316    hybrid data, AAF failed to produce the correct phylogeny and instead clustered the hybrid with

317    its parents, while SISRS failed to complete at all. With the empirical data, we saw similar results;

318    with AAF, we could recapitulate the phylogeny of the genus *Saccharomyces* when using only

319    pure samples, but when we included any hybrid, an incorrect phylogeny was produced (Figure

320    S8a and c). SISRS had similar issues with producing the correct phylogeny with hybrids, but its

321    output allowed for more nuanced network visualizations. For CBS2834, the SISRS output

322    allowed us to infer the shared background with *S. cerevisiae*, *S. kudriavzevii*, *S. uvarum*, and *S.*

323    *eubayanus* (Figure S8d), but the proportion of contribution from each species was difficult to

324    estimate compared to the sppIDer output. Overall, we found that these methods have serious

325    limitations when used with hybrids, but they could be used as a complement to sppIDer to make

326    inferences about pure parental lineages.

327        Methods that assemble targeted genes from short read-data, such as aTRAM and

328    HybPiper, can be used with poor-quality references and/or references that may be missing genes

329    of interest. We tested these tools with a panel of loci that can be used to delineate the *S.*

330    *eubayanus* populations (Peris and Langdon et al. 2016). HybPiper and aTRAM were able to

331    match short-reads to a locus of interest 59% or 34% of the time, respectively, but they could only

332    assemble these reads 23% or 28% of the time, respectively. Neither method could assemble one

333    locus for all 15 strains tested, including both hybrids and non-hybrids (Table S1). While these

334    methods can be powerful when applied in a targeted manner to pure strains, they fail when

335    applied to hybrid data.

336

337    <u>Non-*Saccharomyces* examples</u>

338    *Lachancea*: refining the interpretation of voucher specimens

339        With the publication of 10 high-quality *Lachancea* genome sequences (Vakirlis et al.

340    2016) and another two recently-described and fully-sequenced species (González et al. 2013;

341    Freel et al. 2015; Sarilar et al. 2015; Freel et al. 2016), this genus is becoming a powerful yeast

342    model. As molecular techniques improve, initial identifications in culture and museum

343    collections can yield new interpretations. For example, the strain CBS6924 was initially

344    identified as *Lachancea thermotolerans*, but recent evidence suggested it as a candidate for a

345    novel species (*Lachancea fantastica* nom. nud. Vakirlis et al. 2016). Its closest relative,

346    *Lachancea lanzarotensis,* was also recently described (González et al. 2013). To test sppIDer's

347    utility for determining whether a strain or voucher specimen is or is not properly classified, we

348    tested mapping reads from CBS6924 to a combination genome with all *Lachancea* reference

349    genomes (Figure S9a), then removing the '*L. fantastica*' reference genome (Figure S9b), and then

350    removing both the '*L. fantastica*' and *L. lanzarotensis* reference genomes (Figure S9c). When

351    both reference genomes were removed, the reads were spread across many genomes, and the

352    initial classification of the strain as *L. thermotolerans* would have been easily falsified. By

353    including the *L. lanzarotensis* reference genome, most reads mapped to that reference, but still

354    poorly enough to warrant additional investigation. When the '*L. fantastica*' reference genome was

355    included, CBS6924 reads mapped unambiguously to this reference. These results demonstrate

356    sppIDer's utility outside of the genus *Saccharomyces* to aid in reclassifying provisional species

357    identifications of voucher specimens from culture and museum collections.

358

359    *Drosophila*

360         To test sppIDer with larger genomes, we examined the animal genus *Drosophila*, which

361    is a large genus with many available reference genomes (Adams et al. 2000; *Drosophila* 12

362    Genomes Consortium 2007; Alekseyenko et al. 2013; Sanchez-Flores et al. 2016), as well as

363    several species still lacking reference genomes. The difference in reference genome qualities led

364    us to remove contigs <10Kb from our combination reference genome; this paring down sped up

365    computation time, reduced memory usage, and improved the visualization, but otherwise did not

366    affect the results (Figure S10a). To test the ability of sppIDer to distinguish closely related

367    species, we started with the *Drosophila yakuba* species complex (Turissini et al. 2015), where *D.*

368    *yakuba* has a sequenced reference genome available, but its close relative *Drosophila santomae*

369    and more distant relative *Drosophila teissieri* do not. Here we observed that short reads from a

370    *D. yakuba* (Comeault et al. 2016) representative mapped well to the *D. yakuba* reference

371    genome. As we moved from the close relative *D. santomae* (Figure 5b) to a more distant one, *D.*

372    *teissieri* (Figure 5c), the mapping percentage and quality decreased with increased promiscuous

373    mapping to other relatives (Figure 5a-c). Thus, as in yeasts, sppIDer can classify pure species

374    and their close relatives well and provide insight to guide downstream analyses.

375         We also used *Drosophila* short-read data to test sppIDer's ability to detect hybrids in

376    non-fungal systems. In this case, we used genomic data from a pure parent and RNA-seq data

377    from a $F_1$ interspecies hybrid (Coolon et al. 2014). We found that sppIDer could easily detect

378    hybrids in an animal model, but as expected, detection of CCNVs using RNA-seq was not

379    possible (Figure S11).

380

381    *Arabidopsis*

382         The study of hybrid speciation and allopolyploidy in plants has a long history (Rieseberg

383    1997; Soltis et al. 2015), and we choose *Arabidopsis* as our plant test case because it has

384    reference genomes available for *Arabidopsis halleri*, *Arabidopsis thaliana*, and *Arabidopsis*

385    *lyrata* (Swarbreck et al. 2008). There are drastic differences in the quality of reference genomes

386    available: the *A. thaliana* reference has seven scaffolds with an N50 of 23,459,830, whereas the

387    *A. halleri* reference has 282,453 scaffolds with an N50 of 17,686. To control for this limitation,

388    we again removed contigs <10KB from our combination reference genome, which helped with

389    run time and memory usage but did not affect the conclusions (Figure S10b). These tests in

390    *Arabidopsis* provide an empirical illustration of sppIDer's performance with differing quality

391    reference genomes. *Arabidopsis* also provides a useful test of detecting hybrids in a plant system,

392    as there are two well-supported allotetraploid species in the genus, *Arabidopsis suecica* and

393    *Arabidopsis kamchatica* (Shimizu-Inatsugi et al. 2009; Schmickl et al. 2010). First, we tested

394    short-read data from a divergent lineage of *A. thaliana* (Durvasula et al. 2017) and found that the

395    reads mapped well to the *A. thaliana* reference genome (Figure 5d). As expected, reads from the

396    interspecies hybrid *A. kamchatica* (Novikova et al. 2016) mapped both to *A. lyrata* and to *A.*

397    *halleri* (Figure 5e), approximately equally, confirming that *A. kamchatica* indeed has genomic

398    contributions from these two species and that sppIDer can detect hybrids, even when the

399    combination reference genome contains reference genomes of substantially varying quality.

400    Thus, sppIDer can accurately detect interspecies hybrid in a plant model and will likely become

401    more generally useful in other plant systems, where allopolyploidy is frequent (Soltis et al.

402    2015), as more reference genomes become available.

403

404    <u>mitoSppIDer</u>

405         Applications of sppIDer with non-nuclear sequencing data are also of considerable

406    interest. Organelle genomes (e.g. mitochondria, chloroplast) have a different mode of

407    inheritance, and increasing data suggest widespread reticulation and cases where their ancestries

408    differ from the nuclear genomes (Peris et al. 2014; Wu et al. 2015; Leducq et al. 2017; Peris et

409    al. 2017a; Peris et al. 2017c; Sulo et al. 2017). We developed mitoSppIDer as an extension to

410    explore these non-nuclear inherited elements. Since mitochondrial genomes are generally small,

411    the coding regions can be easily visualized, which allows precise mapping of introgressions in

412    both coding and non-coding regions. However, more cautious interpretation is warranted,

413    because mitochondrial reads are often at low and variable abundance, and quality can differ

414    between DNA isolations and sequencing runs. Again, we tested using the genus *Saccharomyces*

415    because of the availability of mitochondrial reference genomes (Foury et al. 1998; Procházka et

416    al. 2012; Baker et al. 2015). We first tested mitoSppIDer with a strain of *S. uvarum* (ZP1021)

417    (Almeida et al. 2014) and found that, of the reads that mapped to any mitochondrial genome,

418    >99% mapped to the *S. uvarum* mitochondrial genome (Figure S12a). Next, we examined Vin7,

419    a hybrid strain of *S. cerevisiae* X *S. kudriavzevii*, and mitoSppIDer revealed that this strain

420    inherited the mitochondrial genome of *S. kudriavzevii* with intergenic introgressions from

421    multiple non-*S. kudriavzevii* mitochondrial genomes (Figure S12b) (Peris et al. 2017c). As with

422    conventional sppIDer, mitoSppIDer rapidly highlights interesting regions for further analysis,

423    such as detailed phylogenetic analyses of introgression candidates.

424

425    <u>Summary</u>

426         Altogether, these tests show the versatility of sppIDer across clades: in fungi, plants, and

427    animals. sppIDer allows for the rapid exploration and visualization of short-read sequencing data

428    to answer a variety of questions, including species identification; determination of the genome

429    composition of natural, synthetic, and experimentally evolved interspecies hybrids; and inference

430    of CCNVs (Brickwedde et al. 2017; Gorter de Vries et al. 2017; Peris et al. 2017b). With

431    examples from the genus *Saccharomyces*, sppIDer could detect contributions from up to four

432    species and recapitulated the known relative ploidy and aneuploidies of brewing strains. From a

433    simulated phylogeny, we found that sppIDer accurately detected hybrids from a range of

434    divergences in the parents and even detected ancient hybrids. In systems with low-quality or

435    varying quality references genomes, sppIDer performs well without much promiscuous mapping

436    between varying reference qualities, but its ability to infer translocations and CCNVs is limited.

437    Even in systems missing reference genomes, sppIDer still enables rapid inferences by using the

438    reference genomes of closely related species, with the caveat that mapping quality declines with

439    sequence divergence. Additionally, sppIDer works on long-read data and with coverage as low

440    as 0.5X. Finally, sppIDer can be extended to non-nuclear data, allowing for the exploration of

441    alternative evolutionary trajectories of mitochondria or chloroplasts. As more high-quality

442    reference genomes become available across the tree of life, we expect sppIDer will become an

443    increasingly useful and versatile tool to quickly provide a first-pass summary and intuitive

444    visualization of the genomic makeup in diverse organisms and interspecies hybrids.

445

**Methods:**

446

447        The sppIDer workflow to identify pure species, interspecies hybrids, and CCNVs consists

448    of one main pipeline that utilizes common bioinformatics programs, as well as several custom

449    summary and visualization scripts (Figure 1). An upstream step is required to prepare the

450    combination reference genome to test the desired comparison species. The inputs for the main

451    sppIDer pipeline are this combination reference genome and short-read FASTQ file(s) from the

452    organism to test. The output consists of several plots showing to which reference genomes the

453    short-reads mapped, how this mapping varies across the combination reference genome, and

454    several text files of summary information. Additionally, the pipeline retains all the intermediate

455    files used to make the plots and summary files; these contain much more detailed information

456    and may be useful as inputs to various other potential downstream analyses. We are releasing

457    sppIDer as a Docker, which runs as an isolated, self-contained package, without the need to

458    download dependencies and change environmental settings. Packaging complex bioinformatics

459    pipelines as Docker containers increases their reusability and reproducibility, while simplifying

460    their ease of use (Boettiger 2015; Di Tommaso et al. 2015). sppIDer can be found here

461    (https://github.com/GLBRC/sppIDer), where a transparent Dockerfile lays out the technical

462    prerequisites, platform, how they work in combination, and is a repository for all the custom

463    scripts. A manual for sppIDer can be found both at the GitHub page and at

464    http://sppider.readthedocs.io.

465

15

466     The pipeline:

467         Before running the main sppIDer script, a combination reference genome must first be

468     created and properly formatted (top of Figure 1). This is a separate script,

469     combineRefGenomes.py, that takes multiple FASTA-formatted reference genomes and a key

470     listing the reference genomes to use and a unique ID for each. The script concatenates the

471     reference genomes together in the order given in the input key, outputting a combination

472     reference FASTA where the chromosomes/scaffolds are renamed to reflect their reference

473     unique ID and their numerical position within the reference-specific portion of the combination

474     output. For reference genomes that contain many short and uninformative scaffolds, there is an

475     option to remove scaffolds below a desired base-pair length. This option improves speed,

476     memory usage, and visual analysis for large genomes with many scaffolds and low N50 values.

477     Setting a threshold usually does not affect the conclusions (Figure S10), but we recommend

478     trying different thresholds to determine how much information is lost. The choice of reference

479     genomes to concatenate is completely at the discretion of the user and their knowledge of the

480     system to which they are applying sppIDer. We recommend choosing multiple phylogenetically

481     distinct lineages or species, where gene flow and incomplete lineage sorting are limited, from a

482     single genus. We caution that, for ease of analysis and interpretation, less than 30 reference

483     genomes should be used at once. To illustrate the power of sppIDer, for our examples, we used

484     all available species-level reference genomes for the genera tested, but we excluded lineages and

485     strains within species. However, sppIDer could be applied iteratively with different combinations

486     of reference genomes that are more targeted for a particular lineage or question. For example,

487     with an experimentally evolved hybrid, just the parental genomes could be included to detect

488     CCNVs that occurred during the evolution, but with a suspected hybrid isolated from the wild or

489     industry, all potential parent species reference genomes should be included.

490         The main body of sppIDer (Figure 1b) uses a custom `python 2` (Python Software

491     Foundation) script to run published tools and custom scripts to map short-reads to a combination

492     reference genome and parse the output. The first step uses the `mem` algorithm in `BWA` (Li and

493     Durbin 2009) to map the reads to the combined concatenated reference genome. Two custom

494     scripts use this output to count and collect the distribution of mapping qualities (MQ) for the

495     reads that map to each reference genome and produce plots of percentage and MQ of reads that

496     map to each reference genome. The BWA output is also used by `samtools view` and `sort`

16

497 (Li et al. 2009) to keep only reads that map with a MQ > 3, a filter that removes reads that map

498 ambiguously. From here, the number of reads that map to each base pair can be analyzed using

499 `bedtools genomeCoverageBed` (Quinlan and Hall 2010), for smaller genomes using the

500 per-basepair option (`-d`) and, for large genomes, the `-bga` option. The depth of coverage output

501 is used by an R (Wickham 2009; R Core Team 2013) script that determines the mean coverage of

502 the combined reference genome that is subdivided into 10,000 windows of equal size. Finally, a

503 plot for the average coverage for each component reference genome and a second plot of average

504 coverage for the windows are produced.

505

506 <u>The metrics:</u>

507 Several different metrics are used to summarize the data. Depth of coverage is a count of

508 how many reads cover each base pair or region of the genome. Coverage can vary greatly from

509 sequencing run to sequencing run; hence, a $\log_2$ conversion is used to normalize to the mean

510 coverage. As discussed in the Results, depth of coverage plots can be used to infer the species,

511 the parents of hybrids, and ploidy changes either between or within a genome. sppIDer also

512 reports the percentage of reads that map to each reference genome. Finally, sppIDer uses the

513 established MAPPing Quality (MQ) score introduced in Li et al. (2008) to bin reads by their

514 map-ability on a 0-60 scale. A score of zero is used for reads where it is unlikely that their

515 placement is correct, so sppIDer reports these as "unmapped", along with reads that cannot be

516 mapped and therefore do not receive a MQ score. The mapping quality scale can therefore

517 provide a rough assessment of data quality, as well as divergence to the provided reference

518 genomes.

519

520 <u>Tested reference genomes and data:</u>

521 For the *Saccharomyces* tests, we used reference genomes that are scaffolded to a

522 chromosomal level. In some cases, there is only one reference genome available per species, and

523 for the others, we used the first available near-complete reference; see Table S1 for those used.

524 For systems with multiple reference genomes available, the choice could be more targeted, such

525 as utilizing lineage specific references or references that contains unplaced scaffolds with genes

526 of interest. Alternatively, for systems where few genomes are available, we have shown here that

527 a close relative works as a proxy. For the *Saccharomyces* references, each ordered "ultra-

528    scaffolds" genome was downloaded from http://www.saccharomycessensustricto.org/ or for *S.*

529    *arboricola* and *S. eubayanus* from NCBI. The published *S. uvarum* genome (Scannell and Zill et

530    al 2011) had chromosome X swapped with chromosome XII, which was fixed manually. These

531    genomes were concatenated together using the python script combineRefGenomes.py, creating a

532    combination reference FASTA with all *Saccharomyces* species. This combination reference

533    genome can then be used repeatedly to test any dataset of interest. For the *Saccharomyces* tests,

534    we used publicly available FASTQ data from a number of publications, all available on NCBI

535    (Table S1 contains all accession numbers). Using the data for each strain separately and the

536    combination reference genome created above, we then called sppIDer.py with, --out uniqueID, --

537    ref SaccharomcyesCombo.fasta, --r1 read1.fastq, and optionally --r2 read2.fastq. sppIDer is

538    written to test one sample's FASTQ file(s) against one combination reference genome at a time,

539    but this could be easily parallelized.

540         For the tests to determine if hybrids could be detected with missing reference genomes,

541    new combination reference genomes without one species' genomes were created by removing

542    the desired species' reference name from the reference genome key before running

543    combineRefGenomes.py. Since both Vin7 and W34/70 contain contributions from *S. cerevisiae,*

544    the combination reference genome lacking the *S. cerevisiae* reference was tested for each set for

545    FASTQ files for Vin7 and W34/70. The same process was followed to remove the *S.*

546    *kudriavzevii* reference genome from the combination reference to test Vin7, as well as to remove

547    the *S. eubayanus* reference genome from the combination reference to test W34/70.

548         For the *Lachancea* test, all of the genomes were available and downloaded from

549    http://gryc.inra.fr/. The FASTQ data for CBS6924 was downloaded from NCBI. A combination

550    reference genome with all available genomes was created and used. Then, sequentially, the

551    "*Lachancea fantastica*" and *Lachancea lanzarotensis* genomes were removed by modifying the

552    input key and rerunning combineRefGenomes.py. The FASTQ data for CBS6924 was tested

553    against all three of these combination reference genomes. See Table S1 for the full accessions.

554         For the non-*Saccharomyces* tests, we used the most complete reference genome available

555    for each species in the genus (accessions Table S1). Therefore, there is quite a bit of variation

556    between different references. For the *Drosophila* and *Arabidopsis* genomes, we tested removing

557    contigs, using the --trim option, with combineRefGenomes.py, as well as not removing contigs,

558    and found the cleanest results when we removed contigs less than 10 Kb. The combined

18

559    reference genomes of both *Drosophila* and *Arabidopsis* were both larger than four gigabases;

560    therefore, the --byGroup option was used with sppIDer.py to speed up processing and reduce

561    memory usage. The data we tested came from a variety of publications, but we targeted data of

562    divergent or hybrid lineages. See Table S1 for complete information.

563         For the mitoSppIDer test, we used the complete species-level *Saccharomyces*

564    mitochondrial reference genomes available on NCBI, which do not necessarily correspond to the

565    same strain that was used to build nuclear genomic reference (Table S1). Again,

566    combineRefGenomes.py was used to concatenate these references. An additional script,

567    combineGFF.py, was used to create a combination GFF file that was used to denote the coding

568    regions on the output plots. mitoSppIDer.py has an additional flag for the GFF file, but it

569    otherwise runs in a similar manner to sppIDer.py; the same input FASTQ file(s) can even be

570    used. Whole genome sequencing data contains varying amounts of mitochondrial sequences;

571    therefore, using the raw FASTQ data works sufficiently, even when many of the genomic reads

572    will be classified as "unmapped".

573

574    <u>Simulations:</u>

575         To create the simulated low-quality de novo genomes, we used the software `iWGS` (Zhou

576    et al. 2016) to simulate 100bp paired-end reads with an average inter-read insert size of 350bp

577    (sd 10) at 2X coverage from the reference genomes of *S. cerevisiae, S. kudriavzevii, S. uvarum,*

578    and *S. eubayanus*. For the simulated de novo *Saccharomyces* genomes the N50 scores ranged

579    from 1254-1274 and the number of scaffolds ranged from 10023-10426 (Table S1).

580         To simulate short-read data, we used `DWGSIM` (https://github.com/nh13/DWGSIM),

581    which allowed us to vary the coverage, error rate, and mutation rate as needed. The *S. cerevisiae*

582    reference genome was used to simulate single species reads and a concatenation of the *S.*

583    *cerevisiae* and *S. eubayanus* reference genomes was used for hybrid pseudo-lager reads. As a test

584    of an aneuploid genome, we also manually manipulated the *S. cerevisiae* reference genome so

585    that it contained zero copies of chromosomes I and III and duplicate copies of chromosome XII,.

586    All simulated reads were 100bp paired-end reads with an average insert size of 500bp. For the

587    coverage tests, we varied the coverage from 0.01-10X. For the short reads used against the low-

588    quality de novo genomes, we used 10X coverage and a 3% mutation rate. To simulate PacBio-

589   style long reads, we used iWGS on the hybrid pseudo-lager concatenated genome with the

590   default settings of 30X coverage, average read accuracy of 0.9, and SD of read accuracy 0.1.

591        To make our simulated phylogeny, we used the *S. cerevisiae* reference genome as a base

592   and simulated reads with DWGSIM at a 2% mutation rate as 100bp paired-end reads with an

593   average insert size of 500bp at 10X coverage. iWGS was used to assemble these reads. The

594   resulting assembly was again simulated with a 2% mutation rate, and those reads were

595   assembled. This procedure was followed for 6 rounds with one lineage being independently

596   simulated twice each round to produce a speciation event. This simulation resulted in 10 species

597   in the phylogenetic arrangement shown in Figure 3a. Summaries of the final assemblies can be

598   found in Table S1, but the median of the final assemblies was 5100 scaffolds, N50 of 1335, and

599   total length of 6.4MB. Each simulated species was ~12% diverged from *S. cerevisiae,* the most

600   closely related species were ~4% diverged, and the most distantly related species were ~20%

601   diverged. The reads used to produce the final assemblies were used to test whether sppIDer

602   mapped each set of reads to their corresponding reference genomes. The reads of different

603   references were concatenated to simulate pseudo-hybrids of different divergences. To simulate

604   ancient hybrids, the reads from earlier rounds of simulation, before speciation events, were

605   concatenated and tested against the final assemblies with sppIDer. As with the empirical data, to

606   simulate a missing reference genome, that reference was removed from the input key prior to

607   running combineRefGenomes.py.

608

609   <u>Alignment-free phylogenetic methods:</u>

610        We tested four alignment-free phylogenetic methods: two that build phylogenies using

611   short-read data, SISRS (Schwartz et al. 2015) and AAF (Fan et al. 2015), and two that assemble

612   targeted loci from short-read data, aTRAM (Allen et al. 2015) and HybPiper (Johnson et al.

613   2016). We simulated 10X coverage paired-end, 100bp data for each *Saccharomyces* reference

614   genome at a mutation rate of 0 with DWGSIM to use as input for these methods. For SISRS, we

615   used the default settings with a genome size of 12Mb, first using only the reference

616   *Saccharomyces* data, then including empirical data for hybrids. SISRS failed at the missing data

617   filtering step when data from the lager strain W34/70 was used, even when we allowed for all but

618   one sample to have missing data. SISRS nexus outputs were visualized with SplitsTree

619   (Huson and Bryant 2006). For AAF, we found that a *k* of 17 accurately recapitulated the

20

620    *Saccharomyces* phylogeny, even with the inclusion of empirical data from other pure lineages.

621    Once we determined the optimal *k*, we tested including empirical hybrid data. We also used AAF

622    with our simulated phylogeny, which constructed the tree that matched the simulations with the

623    default *k* of 25. The output of AAF was visualized with `iTol` (Letunic and Bork 2016).

624         For the targeted loci methods, we used 13 loci that can delineate *S. eubayanus*

625    populations (Peris and Langdon et al. 2016), as well as the ITS sequences for *S. cerevisiae*

626    (AY046146.1) (Kurtzman and Robnett 2003) and *S. eubayanus* (JF786673.1) (Libkind and

627    Hittinger et al. 2011) as bait, all obtained from NCBI. We tested the simulated *Saccharomyces*

628    reads, as well as the empirical data for P1C1, Fosters O, CBS1503, CBS2834, Vin7, and

629    W34/70. For `aTRAM,` we used the default settings and the option for the Velvet assembler. For

630    `HybPiper,` we used the default settings and the SPADES assembler.

631

632    **Acknowledgments**

645

646    **Reference:**

647    Adams MD, Celniker SE, Holt RA, Evans CA, Gocayne JD, Amanatides PG, Scherer SE, Li

648         PW, Hoskins RA, Galle RF, et al. 2000. The Genome Sequence of *Drosophila*

649         *melanogaster*. Science. 287:2185–2196.

650    Alekseyenko AA, Ellison CE, Gorchakov AA, Zhou Q, Kaiser VB, Toda N, Walton Z, Peng S,

651  Park PJ, Bachtrog D, et al. 2013. Conservation and de novo acquisition of dosage

652    compensation on newly evolved sex chromosomes in *Drosophila*. Genes Dev. 27:853–858.

653 Allen JM, Huang DI, Cronk QC, Johnson KP. 2015. aTRAM - automated target restricted

654    assembly method: a fast method for assembling loci across divergent taxa from next-

655    generation sequencing data. BMC Bioinformatics 16:1–7.

656 Almeida P, Gonçalves C, Teixeira S, Libkind D, Bontrager M, Masneuf-Pomarède I, Albertin W,

657    Durrens P, Sherman DJ, Marullo P, et al. 2014. A Gondwanan imprint on global diversity

658    and domestication of wine and cider yeast *Saccharomyces uvarum*. Nat. Commun. 5:4044.

659 Baker E, Wang B, Bellora N, Peris D, Hulfachor AB, Koshalek JA, Adams M, Libkind D,

660    Hittinger CT. 2015. The genome sequence of *Saccharomyces eubayanus* and the

661    domestication of lager-brewing yeasts. Mol. Biol. Evol. 32:2818–2831.

662 Boettiger C. 2015. An introduction to Docker for reproducible research, with examples from the

663    R environment. Unpubl. Data [Internet]. Available from: https://arxiv.org/abs/1410.0846v1

664 Borneman AR, Desany BA, Riches D, Affourtit JP, Forgan AH, Pretorius IS, Egholm M,

665    Chambers PJ. 2012. The genome sequence of the wine yeast VIN7 reveals an allotriploid

666    hybrid genome with *Saccharomyces cerevisiae* and *Saccharomyces kudriavzevii* origins.

667    FEMS Yeast Res. 12:88–96.

668 Borneman AR, Forgan AH, Kolouchova R, Fraser JA, Schmidt SA. 2016. Whole Genome

669    Comparison Reveals High Levels of Inbreeding and Strain Redundancy Across the

670    Spectrum of Commercial Wine Strains of *Saccharomyces cerevisiae*. G3 6:957–971.

671 Brickwedde A, van den Broek M, Geertman J-MA, Magalhães F, Kuijpers NGA, Gibson B,

672    Pronk JT, Daran J-MG. 2017. Evolutionary Engineering in Chemostat Cultures for

673    Improved Maltotriose Fermentation Kinetics in *Saccharomyces pastorianus* Lager Brewing

674    Yeast. Front. Microbiol. 8:1–15.

675 Comeault AA, Venkat A, Matute DR. 2016. Correlated evolution of male and female

676    reproductive traits drive a cascading effect of reinforcement in *Drosophila yakuba*. Proc. R.

677    Soc. B 283.

678 Coolon JD, Mcmanus CJ, Stevenson KR, Graveley BR, Wittkopp PJ. 2014. Tempo and mode of

679    regulatory evolution in *Drosophila*. Genome Res. 24:797–808.

680 Di Tommaso P, Palumbo E, Chatzou M, Prieto P, Heuer ML, Notredame C. 2015. The impact of

681    Docker containers on the performance of genomic pipelines. PeerJ 3:e1273.

682    Dunn B, Sherlock G. 2008. Reconstruction of the genome origins and evolution of the hybrid
683        lager yeast *Saccharomyces pastorianus*. Genome Res. 18:1610–1623.

684    *Drosophila* 12 Genomes Consortium. 2007. Evolution of genes and genomes on the *Drosophila*
685        phylogeny. Nature 450:203–218.

686    Durvasula A, Fulgione A, Gutaker RM, Irez S, Flood PJ, Neto C, Tsuchimatsu T, Burbano HA,
687        Picó FX, Alonso-Blanco C, et al. 2017. African genomes illuminate the early history and
688        transition to selfing in *Arabidopsis thaliana*. PNAS 114:5213–5218.

689    Fan H, Ives AR, Surget-Groba Y, Cannon CH. 2015. An assembly and alignment-free method of
690        phylogeny reconstruction from next-generation sequencing data. BMC Genomics 16:1–18.

691    Fischer G, James SA, Roberts IN, Oliver SG, Louis EJ. 2000. Chromosomal evolution in
692        *Saccharomyces*. Nature 405:451–454.

693    Foury F, Roganti T, Lecrenier N, Purnelle B. 1998. The complete sequence of the mitochondrial
694        genome of *Saccharomyces cerevisiae*. FEBS Lett. 440:325–331.

695    Freel KC, Charron G, Leducq J-B, Landry CR, Schacherer J. 2015. *Lachancea quebecensis* sp.
696        nov., a yeast species consistently isolated from tree bark in the Canadian province of
697        Quebec. Int. J. Syst. Evol. Microbiol. 65:3392–3399.

698    Freel KC, Friedrich A, Sarilar V, Devillers H, Neuvéglise C, Schacherer J. 2016. Whole-Genome
699        Sequencing and Intraspecific Analysis of the Yeast Species *Lachancea quebecensis*. GBE
700        8:733–741.

701    Gayevskiy V, Goddard MR. 2016. *Saccharomyces eubayanus* and *Saccharomyces arboricola*
702        reside in North Island native New Zealand forests. Environ. Microbiol. 18:1137–1147.

703    Goffeau A, Barrell BG, Bussey H, Davis RW, Dujon B, Feldmann H, Galibert F, Hoheisel JD,
704        Jacq C, Johnston M, et al. 1996. Life with 6000 Genes. Science. 274:546+563-567.

705    Gonçalves M, Pontes A, Almeida P, Barbosa R, Serra M, Libkind D, Hutzler M, Gonçalves P,
706        Sampaio JP. 2016. Distinct Domestication Trajectories in Top- Fermenting Beer Yeasts and
707        Wine Yeasts Distinct Domestication Trajectories in Top-Fermenting Beer Yeasts and Wine
708        Yeasts. Curr. Biol. 26:1–12.

709    González SS, Alcoba-Flórez J, Laich F. 2013. *Lachancea lanzarotensis* sp. nov., an
710        ascomycetous yeast isolated from grapes and wine fermentation in Lanzarote , Canary
711        Islands. Int. J. Syst. Evol. Microbiol. 63:358–363.

712    Gorter De Vries AR, Pronk JT, Daran J-MG. 2017. Industrial Relevance of Chromosomal Copy

23

713   Number Variation in *Saccharomyces* Yeasts. Appl. Environ. Microbiol. 83:1–15.

714 Hittinger CT. 2013. *Saccharomyces* diversity and evolution: a budding model genus. Trends

715   Genet. 29:309–317.

716 Hittinger CT, Gonçalves P, Sampaio JP, Dover J, Johnston M, Rokas A. 2010. Remarkably

717   ancient balanced polymorphisms in a multi-locus gene network. Nature 464:54–58.

718 Hoot SB, Napier NS, Taylor WC. 2004. Reveling Unknown or Extinct Lineages within *Isoetes*

719   (Isoetaceae) Using DNA Sequences from Hybrids. Am. J. Bot. 91:899–904.

720 Huson DH, Bryant D. 2006. Application of Phylogenetic Networks in Evolutionary Studies.

721   Mol. Biol. Evol. 23:254–267.

722 Johnson MG, Gardner EM, Liu Y, Medina R, Goffinet B, Shaw J, Zerega NJC, Wickett NJ.

723   2016. HybPiper: Extracting Coding Sequence and Introns for Phylogenetics from High-

724   Throughput Sequencing Reads Using Target Enrichment. Appl. Plant Sci. 4.

725 Kurtzman CP, Robnett CJ. 2003. Phylogenetic relationships among yeasts of the

726   '*Saccharomyces* complex' determined from multigene sequence analyses. FEMS Yeast Res.

727   3:417–432.

728 Leducq J-B, Henault M, Charron G, Nielly-Thibault L, Terrat Y, Fiumera HL, Shapiro BJ,

729   Landry CR. 2017. Mitochondrial Recombination and Introgression during Speciation by

730   Hybridization. Mol. Biol. Evol. 34:1947–1959.

731 Leducq J-B, Nielly-Thibault L, Charron G, Eberlein C, Verta J-P, Samani P, Sylvester K,

732   Hittinger CT, Bell G, Landry CR. 2016. Speciation driven by hybridization and

733   chromosomal plasticity in a wild yeast. Nat. Microbiol. 1:1–10.

734 Letunic I, Bork P. 2016. Interactive tree of life ( iTOL ) v3: an online tool for the display and

735   annotation of phylogenetic and other trees. Nucleic Acids Res. 44:W242–W245.

736 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows – Wheeler

737   transform. Bioinformatics 25:1754–1760.

738 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,

739   Subgroup 1000 Genome Project Data Processing. 2009. The Sequence Alignment/Map

740   format and SAMtools. Bioinformatics 25:2078–2079.

741 Li H, Ruan J, Durbin R. 2008. Mapping short DNA sequencing reads and calling variants using

742   mapping quality scores. Genome Res. 18:1851–1858.

743 Libkind D, Hittinger CT, Valério E, Gonçalves C, Dover J, Johnston M, Gonçalves P, Sampaio

744    JP. 2011. Microbe domestication and the identification of the wild genetic stock of lager-

745    brewing yeast. Proc. Natl. Acad. Sci. U.S.A. 108:14539–14544.

746  Liti G, Carter DM, Moses AM, Warringer J, Parts L, James S a, Davey RP, Roberts IN, Burt A,

747    Koufopanou V, et al. 2009. Population genomics of domestic and wild yeasts. Nature

748    458:337–341.

749  Liti G, Nguyen Ba AN, Blythe M, Müller CA, Bergström A, Cubillos FA, Dafnhis-Calas F,

750    Khoshraftar S, Malla S, Mehta N, et al. 2013. High quality de novo sequencing and

751    assembly of the *Saccharomyces arboricolus* genome. BMC Genomics 14.

752  Naseeb S, James SA, Alsammar H, Michaels CJ, Gini B, Nueno-palop C, Bond CJ, Mcghie H,

753    Roberts IN, Delneri D. 2017. *Saccharomyces jurei* sp. nov., isolation and genetic

754    identification of a novel yeast species from *Quercus robur*. Microbiology 67:2046–2052.

755  Novikova PY, Hohmann N, Nizhynska V, Tsuchimatsu T, Ali J, Muir G, Guggisberg A, Paape

756    T, Schmid K, Fedorenko OM, et al. 2016. Sequencing of the genus *Arabidopsis* identifies a

757    complex history of nonbifurcating speciation and abundant trans-specific polymorphism.

758    Nat. Genet. 48:1077–1082.

759  Okuno M, Kajitani R, Ryusui R, Morimoto H, Kodama Y, Itoh T. 2016. Next-generation

760    sequencing analysis of lager brewing yeast strains reveals the evolutionary history of

761    interspecies hybridization. DNA Res. 1:1–14.

762  Payseur BA, Rieseberg LH. 2016. A genomic perspective on hybridization and speciation. Mol.

763    Ecol. 25:2337-2360.

764  Peris D, Arias A, Orlic S, Belloch C, Pérez-Través L, Querol A, Barrio E. 2017a. Molecular

765    Phylogenetics and Evolution Mitochondrial introgression suggests extensive ancestral

766    hybridization events among *Saccharomyces* species. Mol. Phylogenet. Evol. 108:49–60.

767  Peris D, Langdon QK, Moriarty R V, Sylvester K, Bontrager M, Charron G, Leducq J, Landry

768    CR, Libkind D, Hittinger CT. 2016. Complex Ancestries of Lager-Brewing Hybrids Were

769    Shaped by Standing Variation in the Wild Yeast *Saccharomyces eubayanus*. PLoS Genet.

770    12.

771  Peris D, Lopes CA, Arias A, Barrio E. 2012. Reconstruction of the Evolutionary History of

772    *Saccharomyces cerevisiae* x *S. kudriavzevii* Hybrids Based on Multilocus Sequence

773    Analysis. PLoS One 7.

774  Peris D, Moriarty R V, Alexander WG, Baker E, Sylvester K, Sardi M, Langdon QK, Libkind D,

775   Wang QM, Bai FY, et al. 2017b. Biotechnology for Biofuels Hybridization and adaptive
776       evolution of diverse *Saccharomyces* species for cellulosic biofuel production. Biotechnol.
777       Biofuels 10:1–19.

778   Peris D, Pérez-Torrado R, Hittinger CT, Barrio E, Querol A. 2017c. On the origins and industrial
779       applications of *Saccharomyces cerevisiae* × *Saccharomyces kudriavzevii* hybrids. Yeast.

780   Peris D, Sylvester K, Libkind D, Gonçalves P, Sampaio JP, Alexander WG, Hittinger CT. 2014.
781       Population structure and reticulate evolution of *Saccharomyces eubayanus* and its lager-
782       brewing hybrids. Mol. Ecol. 23:2031–2045.

783   Procházka E, Franko F, Poláková S, Sulo P. 2012. A complete sequence of *Saccharomyces*
784       *paradoxus* mitochondrial genome that restores the respiration in S. cerevisiae. FEMS Yeast
785       Res. 12:819–830.

786   Pryszcz LP, Németh T, Gácser A, Gabaldón T. 2014. Genome Comparison of *Candida*
787       *orthopsilosis* Clinical Strains Reveals the Existence of Hybrids between Two Distinct
788       Subspecies. Genome Biol. Evol. 6:1069–1078.

789   Python Software Foundation. Python Language Reference, version 2.7.

790   Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
791       features. Bioinformatics 26:841–842.

792   R Core Team. 2013. R: A Language and Environment for Statistical Computing.

793   Richards S. 2017. It's more than stamp collecting: how genome sequencing can unify biological
794       research. Trends Genet. 31:411–421.

795   Rieseberg LH. 1997. Hybrid Origins of Plant Species. Annu. Rev. Ecol. Evol. Syst. 28:359–389.

796   Sanchez-Flores A, Peñaloza F, Carpinteyro-Ponce J, Nazario-Yepiz N, Abreu-Goodger C,
797       Machado CA, Markow TA. 2016. Genome Evolution in Three Species of Cactophilic
798       *Drosophila*. G3 6:3097–3105.

799   Sarilar V, Devillers H, Freel KC, Schacherer J, Neuvéglise C. 2015. Draft Genome Sequence of
800       *Lachancea lanzarotensi*s CBS 12615 T , an Ascomycetous Yeast Isolated from Grapes.
801       Genome Announc. 3:1–2.

802   Scannell DR, Zill OA, Rokas A, Payen C, Dunham MJ, Eisen MB, Rine J, Johnston M, Hittinger
803       CT. 2011. The Awesome Power of Yeast Evolutionary Genetics: New Genome Sequences
804       and Strain Resources for the *Saccharomyces sensu stricto* Genus. G3 1:11–25.

805   Schmickl R, Jørgensen MH, Brysting AK, Koch MA. 2010. The evolutionary history of the

806  *Arabidopsis lyrata* complex: a hybrid in the amphi-Beringian area closes a large distribution

807  gap and builds up a genetic barrier. BMC Evol. Biol. 10:1–18.

808  Schwartz RS, Harkins KM, Stone AC, Cartwright RA. 2015. A composite genome approach to

809  identify phylogenetically informative data from next-generation sequencing. BMC

810  Bioinformatics 16:1–10.

811  Shimizu-Inatsugi R, Lihová J, Iwanaga H, Kudoh H, Marhold K, Savolainen O, Watanabe K,

812  Yakubov V V., Shimizu KK. 2009. The allopolyploid *Arabidopsis kamchatica* originated

813  from multiple individuals of *Arabidopsis lyrata* and *Arabidopsis halleri*. Mol. Ecol.

814  18:4024–4048.

815  Soltis PS, Marchant DB, Van de Peer Y, Soltis DE. 2015. Polyploidy and genome evolution in

816  plants. Curr. Opin. Genet. Dev. 35:119–125.

817  Sulo P, Szabóová D, Bielik P, Poláková S, Šoltys K, Jatzová K, Szemes T. 2017. The

818  evolutionary history of *Saccharomyces* species inferred from completed mitochondrial

819  genomes and revision in the 'yeast mitochondrial genetic code.' DNA Res. 24:571–583.

820  Swarbreck D, Wilks C, Lamesch P, Berardini TZ, Garcia-hernandez M, Foerster H, Li D, Meyer

821  T, Muller R, Ploetz L, et al. 2008. The *Arabidopsis* Information Resource (TAIR): gene

822  structure and function annotation. Nucleic Acids Res. 36:1009–1014.

823  Turissini DA, Liu G, David JR, Matute DR. 2015. The evolution of reproductive isolation in the

824  *Drosophila yakuba* complex of species. J. Evol. Biol. 28:557–575.

825  Vakirlis N, Sarilar V, Drillon G, Fleiss A, Agier N, Meyniel J-P, Blanpain L, Carbone A,

826  Devillers H, Dubois K, et al. 2016. Reconstruction of ancestral chromosome architecture

827  and gene repertoire reveals principles of genome evolution in a model yeast genus. Genome

828  Res. 26:918–932.

829  Wickham H. 2009. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York

830  Wu B, Buljic A, Hao W. 2015. Extensive Horizontal Transfer and Homologous Recombination

831  Generate Highly Chimeric Mitochondrial Genomes in Yeast. Mol. Biol. Evol. 32:2559–

832  2570.

833  Zhou X, Peris D, Kominek J, Kurtzman CP, Hittinger CT, Rokas A. 2016. In Silico Whole

834  Genome Sequencer and Analyzer ( iWGS ): a Computational Pipeline to Guide the Design

835  and Analysis of de novo Genome Sequencing Studies. G3 6:3655–3662.

836

**Figure legends**

**Figure 1.** Workflow of sppIDer. (a) An upstream step concatenates all the desired reference genomes (represented by colored bars). Generally, references should be distinct species (see methods for advice about choosing references). This combination reference genome can be used for many analyses. (b) The main sppIDer pipeline. First, reads (short lines) are mapped. This output is used to parse for quality and percentage (left) or for coverage (right). On the left, quality (high MQ black lines versus low MQ light lines) is parsed, and the percentage of reads that map to each genome or do not map (grey bar) is calculated. To determine coverage, only MQ>3 reads (black lines) are kept and sorted into the combination reference genome order. These reads are then counted, either for each base pair or, for large genomes (combination length >4Gb), in groups. Then, the combination reference genome is broken into equally sized pieces, and the average coverage is calculated. (c) Several plots are produced. Shown here are examples of Percentage Mapped and Mapping Quality plots, a plot showing average coverage by species, and two ways to show coverage by windows with species side-by-side or stacked. *Scer = S. cerevisiae, Spar = S. paradoxus, Smik = Saccharomyces mikatae, Skud = S. kudriavzevii, Sarb = Saccharomyces arboricola, Suva = S. uvarum, Seub = S. eubayanus*.

**Figure 2.** Normalized coverage plots of *Saccharomyces* test cases. (a) Reads from a New Zealand isolate of *S. eubayanus*, P1C1, mapped to the *S. eubayanus* reference genome (magenta). (b) Reads from an ale strain, FostersO, mapped to the *S. cerevisiae* reference genome (red), with visually detectable aneuploidies. (c) Reads from a hybrid Frohberg lager strain, W34/70, mapped to both the *S. cerevisiae* and *S. eubayanus* reference genomes in an average approximately 1:1 ratio with visually detectable translocations and aneuploidies. (d) Reads from a hybrid Saaz lager strain, CBS1503, mapped to both *S. cerevisiae* and *S. eubayanus* reference genomes in an average approximately 1:2 (respectively) ratio with visually detectable translocations and aneuploidies. (e) Reads from a wine hybrid strain, Vin7, mapped to *S. cerevisiae* and *S. kudriavzevii* (green) reference genomes in an average approximately 2:1 (respectively) ratio. (f) Reads from a hybrid cider-producing strain, CBS2834, mapped to four reference genomes: *S. cerevisiae, S. kudriavzevii, S. uvarum* (purple), and *S. eubayanus*.

**Figure 3.** Simulated phylogeny of 10 species and sppIDer's detection of hybrids from this phylogeny. (a) Phylogeny built with AAF. (b) Reads from G mapped to the G reference genome.

868    (c) Reads from a pseudo-hybrid of the closely related species G and H mapped to the G and H

869    references. (d) Reads from more distant pseudo-hybrid of E and G mapped to references E and

870    G. (e) Reads of ancient pseudo-hybrid of A and a common ancestor of G and H mapped to the

871    references of A, G, and H, which are the lineages that descended from the hybrid's parents. (f)

872    Without the G reference genome, reads from a pseudo-hybrid of the closely related species G

873    and H mapped to the H reference genome, with some mapped promiscuously to references I and

874    J.

875    **Figure 4.** Comparison of the percentage of reads that mapped when different reference genomes

876    were excluded, compared to when all possible reference genomes for *Saccharomyces* were

877    available (middle panels). (a) When the *S. cerevisiae* reference genome was not provided and

878    reads from a Frohberg lager strain, W34/70, were mapped, more reads failed to map (grey) or

879    mapped to the *S. paradoxus* reference genome (yellow). (b) When the full array of

880    *Saccharomyces* genomes was provided, reads for the lager strain mapped to both *S. cerevisiae*

881    and *S. eubayanus*. (c) When the *S. eubayanus* reference genome was removed, more reads from

882    the lager strain failed to map or mapped to the *S. uvarum* reference genome (purple). (d) With

883    the removal of the *S. cerevisiae* reference genome, reads from the *S. cerevisiae* X *S. kudriavzevii*

884    hybrid strain Vin7, which would normally map to *S. cerevisiae*, instead failed to map or mapped

885    to *S. paradoxus*. (e) When all genomes were used, reads mapped to both *S. cerevisiae* and *S.*

886    *kudriavzevii*. (f) With the removal of the *S. kudriavzevii* reference genome, reads that would

887    normally map to *S. kudriavzevii* instead failed to map or were distributed across all other

888    genomes.

889    **Figure 5.** Examples using animal and plant genomes. (a) Reads from a *D. yakuba* individual

890    mapped primarily (>99%) to the *D. yakuba* reference genome. (b) Reads from the sister species

891    *D. santomae* mapped best to the *D. yakuba* reference genome with some mapped promiscuously

892    to other reference genomes. (c) Reads from the more distantly related species *D. teissieri* mapped

893    mostly to the *D. yakuba* reference genome, but with more reads not mapped and mapped

894    promiscuously to other related reference genomes. (d) Reads from an *Arabidopsis thaliana*

895    accession from Tanzania mapped back to the European reference genome for *A. thaliana*. The

896    repetitive nature of centromeres causes the coverage to fluctuate around those regions. (e) Reads

897    from the hybrid species *A. kamchatica* mapped to the two parental reference genomes: *A. halleri*
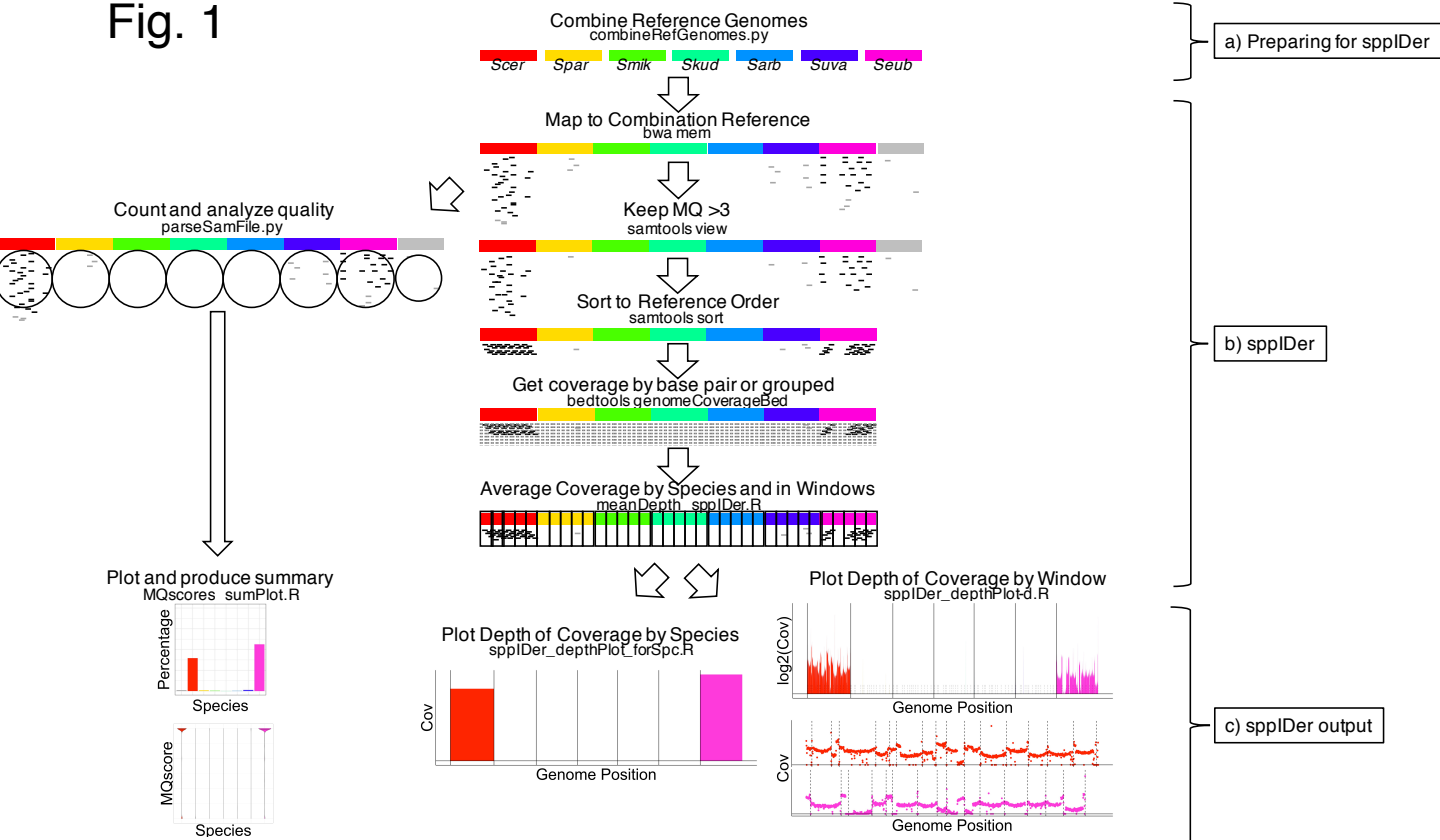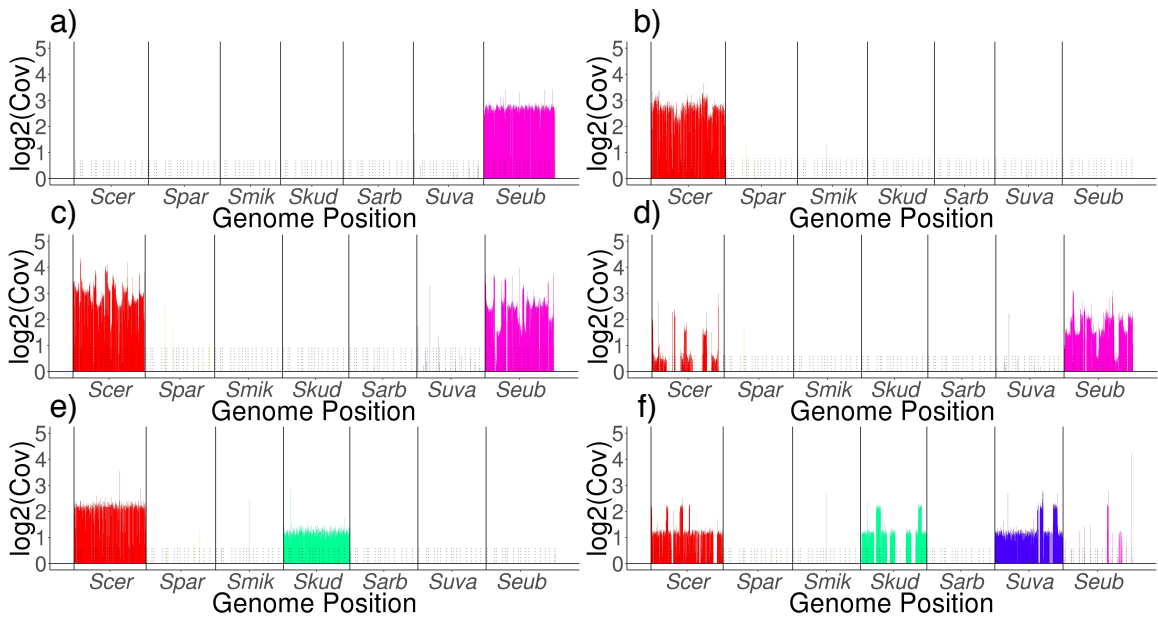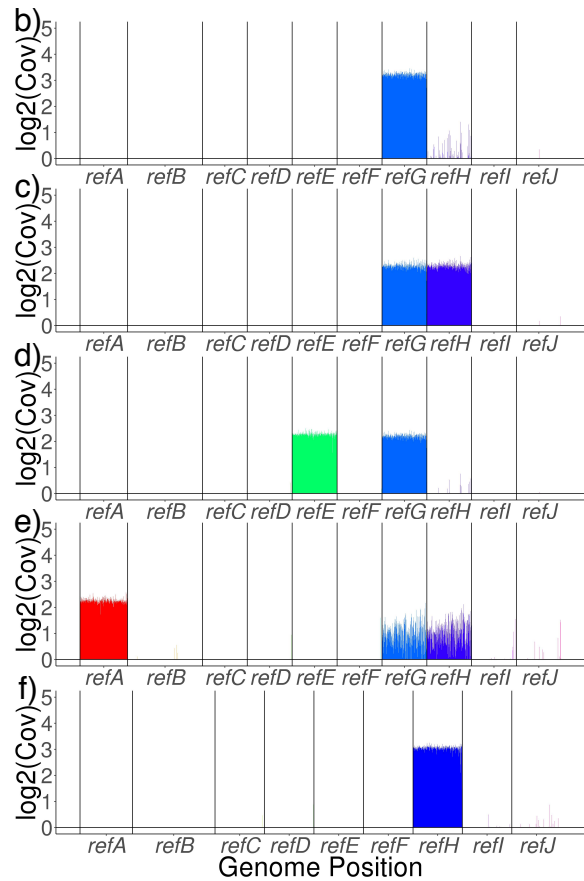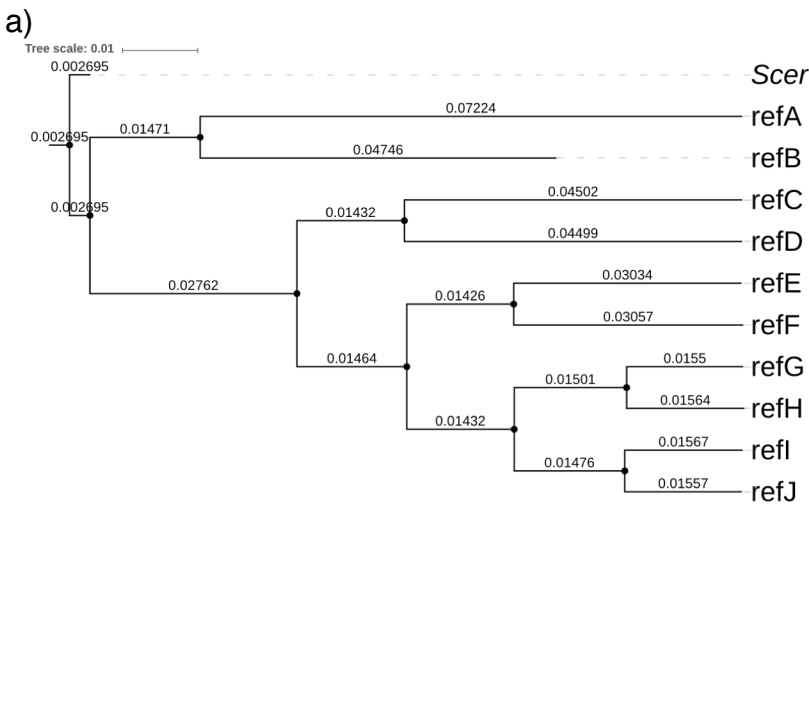
898    and *A. lyrata*.

# Fig. 1



**Combine Reference Genomes**
combineRefGenomes.py

*Scer*  *Spar*  *Smik*  *Skud*  *Sarb*  *Suva*  *Seub*

a) Preparing for sppIDer

**Map to Combination Reference**
bwa mem

**Keep MQ >3**
samtools view

**Sort to Reference Order**
samtools sort

**Get coverage by base pair or grouped**
bedtools genomeCoverageBed

**Average Coverage by Species and in Windows**
meanDepth_sppIDer.R

b) sppIDer

**Count and analyze quality**
parseSamFile.py

**Plot and produce summary**
MQscores_sumPlot.R

**Plot Depth of Coverage by Species**
sppIDer_depthPlot_forSpc.R

**Plot Depth of Coverage by Window**
sppIDer_depthPlot-d.R

c) sppIDer output

# Fig. 2

Fig. 3



Fig. 4

Fig. 5