1    **Extensive loss of cell cycle and DNA repair genes in an ancient lineage of bipolar budding**

2    **yeasts**

3

4    Jacob L. Steenwyk[1], Dana A. Opulente[2], Jacek Kominek[2], Xing-Xing Shen[1], Xiaofan Zhou[3],

5    Abigail L. Labella[1], Noah P. Bradley[1], Brandt F. Eichman[1], Neža Čadež[5], Diego Libkind[6],

6    Jeremy DeVirgilio[7], Amanda Beth Hulfachor[4], Cletus P. Kurtzman[7], Chris Todd Hittinger[2]*, and

7    Antonis Rokas[1]*

8

9    1. Department of Biological Sciences, Vanderbilt University, Nashville, TN 37235, USA

10   2. Laboratory of Genetics, Genome Center of Wisconsin, DOE Great Lakes Bioenergy Research

11   Center, Wisconsin Energy Institute, J. F. Crow Institute for the Study of Evolution, University of

12   Wisconsin–Madison, Wisconsin 53706, USA

13   3. Guangdong Province Key Laboratory of Microbial Signals and Disease Control, Integrative

14   Microbiology Research Centre, South China Agricultural University, 510642 Guangzhou, China

15   4. Laboratory of Genetics, Genome Center of Wisconsin, Wisconsin Energy Institute, J.F. Crow

16   Institute for the Study of Evolution, University of Wisconsin-Madison, Madison, WI 53706,

17   USA

18   5. University of Ljubljana, Biotechnical Faculty, Department of Food Science and Technology,

19   Jamnikarjeva 101, 1000 Ljubljana, Slovenia

20   6. Laboratorio de Microbiología Aplicada, Biotecnología y Bioinformática, Instituto Andino

21   Patagónico de Tecnologías Biológicas y Geoambientales (IPATEC), Universidad Nacional del

22   Comahue - CONICET, San Carlos de Bariloche, 8400, Río Negro, Argentina

23   7. Mycotoxin Prevention and Applied Microbiology Research Unit, National Center for

24   Agricultural Utilization Research, Agricultural Research Service, US Department of Agriculture,

25   Peoria, Illinois 61604, USA

26

27   *Correspondence: cthittinger@wisc.edu and antonis.rokas@vanderbilt.edu

28

29   Running title: Gene loss among *Hanseniaspora*

30

31   Keywords: phylogenomics, Saccharomycotina, biodiversity, carbohydrate metabolism,

32   genomics, DNA repair, cell cycle, genome stability, DNA damage response

33 **Abstract**

34 Cell cycle checkpoints and DNA repair processes protect organisms from potentially lethal

35 mutational damage. Compared to other budding yeasts in the subphylum Saccharomycotina, we

36 noticed that a lineage in the genus *Hanseniaspora* exhibited very high evolutionary rates, low

37 GC content, small genome sizes, and lower gene numbers. To better understand *Hanseniaspora*

38 evolution, we analyzed 25 genomes, including 11 newly sequenced, representing 18 / 21 known

39 species in the genus. Our phylogenomic analyses identify two *Hanseniaspora* lineages, the fast-

40 evolving lineage (FEL), which began diversifying ~87 million years ago (mya), and the slow-

41 evolving lineage (SEL), which began diversifying ~54 mya. Remarkably, both lineages lost

42 genes associated with the cell cycle and genome integrity, but these losses were greater in the

43 FEL. For example, all species lost the cell cycle regulator *WHI5*, and the FEL lost components of

44 the spindle checkpoint pathway (e.g., *MAD1*, *MAD2*) and DNA damage checkpoint pathway

45 (e.g., *MEC3*, *RAD9*). Similarly, both lineages lost genes involved in DNA repair pathways,

46 including the DNA glycosylase gene *MAG1*, which is part of the base excision repair pathway,

47 and the DNA photolyase gene *PHR1*, which is involved in pyrimidine dimer repair. Strikingly,

48 the FEL lost 33 additional genes, including polymerases (i.e., *POL4* and *POL32*) and telomere-

49 associated genes (e.g., *RIF1, RFA3, CDC13, PBP2*). Echoing these losses, molecular

50 evolutionary analyses reveal that, compared to the SEL, the FEL stem lineage underwent a burst

51 of accelerated evolution, which resulted in greater mutational loads, homopolymer instabilities,

52 and higher fractions of mutations associated with the common endogenously damaged base, 8-

53 oxoguanine. We conclude that *Hanseniaspora* is an ancient lineage that has diversified and

54 thrived, despite lacking many otherwise highly conserved cell cycle and genome integrity genes

55 and pathways, and may represent a novel system for studying cellular life without them.

56 **Introduction**

57 Genome maintenance is largely attributed to the fidelity of cell cycle checkpoints, DNA repair

58 pathways, and their interaction [1]. Dysregulation of these processes often leads to the loss of

59 genomic integrity [2] and hypermutation, or the acceleration of mutation rates [3]. For example,

60 improper control of cell cycle and DNA repair processes can lead to 10- to 100-fold increases in

61 mutation rate [4]. Furthermore, deletions of single genes can have profound effects on genome

62 stability. For example, the deletion of *MEC3*, which is involved in sensing DNA damage in the

63 G1 and G2/M cell cycle phases, can lead to a 54-fold increase in the gross chromosomal

64 rearrangement rate [5]. Similarly, nonsense mutations in mismatch repair proteins account for the

65 emergence of hypermutator strains in the yeast pathogens *Cryptococcus deuterogattii* [6] and

66 *Cryptococcus neoformans* [7,8]. Due to their importance in ensuring genomic integrity, most

67 genome maintenance-associated processes are thought to be evolutionarily ancient and broadly

68 conserved [9].

69

70 One such ancient and highly conserved process in eukaryotes is the cell cycle [10,11]. Landmark

71 features of cell cycle control include cell size control, the mitotic spindle checkpoint, the DNA

72 damage response checkpoint, and DNA replication [9]. Cell size is controlled, in part, through

73 the activity of *WHI5*, which represses the G1/S transition by inhibiting G1/S transcription [12].

74 Similarly, when kinetochores are improperly attached or are not attached to microtubules, the

75 mitotic spindle checkpoint helps to prevent activation of the anaphase-promoting complex

76 (APC), which controls the G1/S and G2/M transitions [9,13]. Additional key regulators in this

77 process are Mad1 and Mad2, which dimerize at unattached kinetochores and delay anaphase.

78 Failure of Mad1:Mad2 recruitment to unattached kinetochores results in failed checkpoint

4

79      activity [14]. Importantly, many regulators, including but not limited to those mentioned here,

80      are highly similar in structure and function between fungi and animals and are thought to have a

81      shared ancestry [10]. Interestingly, cell cycle initiation in certain fungi (including

82      *Hanseniaspora*) is achieved through SBF, a transcription factor that is functionally equivalent

83      but evolutionarily unrelated to E2F, the transcription factor that that initiates the cycle in

84      animals, plants, and certain early-diverging fungal lineages [11]. SBF is postulated to have been

85      acquired via a viral infection, suggesting that evolutionary changes in this otherwise highly

86      conserved process can and do rarely occur [11,15].

87

88      DNA damage checkpoints can arrest the cell cycle and influence the activation of DNA repair

89      pathways, the recruitment of DNA repair proteins to damaged sites, and the composition and

90      length of telomeres [16]. For example, *MEC3* and *RAD9*, function as checkpoint genes required

91      for arrest in the G2 phase after DNA damage has occurred [17]. Additionally, the deletions of

92      DNA damage and checkpoint genes have been known to cause hypermutator phenotypes in the

93      baker's yeast *Saccharomyces cerevisiae* [18]. Similarly, hypermutator phenotypes are associated

94      with loss-of-function mutations in DNA polymerase genes [19]. For example, deletion of the

95      DNA polymerase $\delta$ subunit gene, *POL32*, which participates in multiple DNA repair processes,

96      causes an increased mutational load and hypermutation in *S. cerevisiae*, in part, through the

97      increase of genomic deletions and small indels [18,20]. Likewise, the deletion of *MAG1*, a gene

98      encoding a DNA glycosylase that removes damaged bases via the multi-step base excision repair

99      pathway, can cause a 2,500-fold increased sensitivity to the DNA alkylating agent methyl

100     methanesulfonate [21].

101

102    In contrast to genes in multi-step DNA repair pathways, other DNA repair genes function

103    individually or are parts of simpler regulatory processes. For example, *PHR1*, a gene that

104    encodes a photolyase, is activated in response to and repairs pyrimidine dimers, one of the most

105    frequent types of lesions caused by damaging UV light [22,23]. Other DNA repair genes do not

106    interact with DNA but function to prevent the misincorporation of damaged bases. For example,

107    *PCD1* encodes a 8-oxo-dGTP diphosphatase [24], which suppresses $G \rightarrow T$ or $C \rightarrow A$

108    transversions by removing 8-oxo-dGTP, thereby preventing the incorporation of the base 8-oxo-

109    dG, one of the most abundant endogenous forms of an oxidatively damaged base [24–26].

110    Collectively, these studies demonstrate that the loss of DNA repair genes can lead to

111    hypermutation and increased sensitivity to DNA damaging agents.

112

113    Hypermutation phenotypes are generally short-lived because most mutations are deleterious and

114    are generally adaptive only in highly stressful or rapidly fluctuating environments [27]. For

115    example, in *Pseudomonas aeruginosa* infections of cystic fibrosis patients [28] and mouse gut-

116    colonizing *Escherichia coli* [29], hypermutation is thought to facilitate adaptation to the host

117    environment and the evolution of drug resistance. Similarly, in the fungal pathogens *C.*

118    *deuterogattii* [6] and *C. neoformans* [7,8], hypermutation is thought to contribute to within-host

119    adaptation, which may involve modulating traits such as drug resistance [6]. However, as

120    adaptation to a new environment nears completion, hypermutator alleles are expected to decrease

121    in frequency due to the accumulation of deleterious mutations that result as a consequence of the

122    high mutation rate [30,31]. In agreement with this prediction, half of experimentally evolved

123    hypermutating lines of *S. cerevisiae* had reduced mutation rates after a few thousand generations

124    [32], suggesting hypermutation is a short-lived phenotype and that compensatory mutations can

125    restore or lower the mutation rate. Additionally, this experiment also provided insights to how

126    strains may cope with hypermutation; for example, all *S. cerevisiae* hypermutating lines

127    increased their ploidy to presumably reduce the impact of higher mutation rates [32]. Altogether,

128    hypermutation can produce short-term advantages but causes long-term disadvantages, which

129    may explain its repeated but short-term occurrence in clinical environments [29] and its

130    sparseness in natural ones. While these theoretical and experimental studies have provided

131    seminal insights to the evolution of mutation rate and hypermutation, we still lack understanding

132    of the long-term, macroevolutionary effects of increased mutation rates.

133

134    Recently, multiple genome-scale phylogenies of species in the budding yeast subphylum

135    Saccharomycotina showed that certain species in the bipolar budding yeast genus *Hanseniaspora*

136    are characterized by very long branches [33–35], which are reminiscent of the very long

137    branches of fungal hypermutator strains [6–8]. Most of what is known about these cosmopolitan

138    apiculate yeasts relates to their high abundance on mature fruits and in fermented beverages [36],

139    especially on grapes and in wine must [37,38]. As a result, *Hanseniaspora* plays a significant

140    role in the early stages of fermentation and can modify wine color and flavor through the

141    production of enzymes and aroma compounds [39]. Surprisingly, even with the use of *S.*

142    *cerevisiae* starter cultures, *Hanseniaspora* species, particularly *Hanseniaspora uvarum,* can

143    achieve very high cell densities , in certain cases comprising greater than 80% of the total yeast

144    population, during early stages of fermentation [40], suggesting exceptional growth capabilities

145    in this environment.

146

147     To gain insight into the long branches and the observed fast growth of *Hanseniaspora*, we

148     sequenced and extensively characterized gene content and patterns of evolution in 25 genomes,

149     including 11 newly sequenced for this study, from 18 / 21 known species in the genus. Our

150     analyses delineated two lineages, the fast-evolving lineage (FEL), which has a strong signature

151     of acceleration in evolutionary rate at its stem branch, and the slow-evolving lineage (SEL),

152     which has a weaker signature of evolutionary rate acceleration at its stem branch. Relaxed

153     molecular clock analyses estimate that the FEL and SEL split ~95 million years ago (mya). The

154     degree of evolutionary rate acceleration is commensurate with the preponderance of loss of

155     genes associated with metabolic, cell cycle, and DNA repair processes. Specifically, compared to

156     *S. cerevisiae*, there are 748 genes that were lost from two-thirds of *Hanseniaspora* genomes with

157     FEL yeasts having lost an additional 661 genes and SEL yeasts having lost only an additional 23.

158     Both lineages have lost major cell cycle regulators, including *WHI5* and components of the APC,

159     while FEL species additionally lost numerous genes associated with the spindle checkpoint (e.g.,

160     *MAD1* and *MAD2*) and DNA damage checkpoint (e.g., *MEC3* and *RAD9*). Similar patterns are

161     observed among DNA repair-related genes; *Hanseniaspora* species have lost 14 genes, while the

162     FEL yeasts have lost an additional 33 genes. For example, both lineages have lost *MAG1* and

163     *PHR1*, while the FEL has lost additional genes including polymerases (i.e., *POL32* and *POL4*)

164     and multiple telomere-associated genes (e.g., *RIF1*, *RFA3, CDC13, PBP2*). Compared to the

165     SEL, analyses of substitution patterns in the FEL show higher levels of sequence substitutions,

166     greater instability of homopolymers, and a greater mutational signature associated with the

167     commonly damaged base, 8-oxo-dG [26]. Furthermore, we find that the transition to transversion

168     (or transition / transversion) ratios of the FEL and the SEL are both very close to the ratio

169     expected if transitions and transversions occur neutrally. These results are consistent with the

8

170     hypothesis that species in the FEL represent a novel example of diversification and long-term

171     evolutionary survival of a hypermutator lineage, which highlights the potential of *Hanseniaspora*

172     for understanding the long-term effects of hypermutation on genome function and evolution.

173

174     **Results**

175     **An exceptionally high evolutionary rate in the FEL stem branch**

176     Concatenation and coalescence analyses of a data matrix of 1,034 single-copy OGs (522,832

177     sites; 100% taxon-occupancy) yielded a robust phylogeny of the genus *Hanseniaspora* (Fig 1A,

178     Fig S2, Fig S3). Consistent with previous analyses [34,35,41], our phylogeny revealed the

179     presence of two major lineages, each of which was characterized by long stem branches; we

180     hereafter refer to the lineage with a longer stem branch as the fast-evolving lineage (FEL) and to

181     the other as the slow-evolving lineage (SEL). Relaxed molecular clock analysis suggests that the

182     FEL and SEL split 95.34 (95% credible interval (CI): 117.38 – 75.36) mya, with the origin of

183     their crown groups estimated at 87.16 (95% CI: 112.75 – 61.38) and 53.59 (95% CI: 80.21 –

184     33.17) mya, respectively (Fig 1A, Fig S4, File S2).

185

186     The FEL stem branch is much longer than the SEL stem branch in the *Hanseniaspora* phylogeny

187     (Fig 1) (see also phylogenies in: Shen et al., 2016, 2018). To determine whether this difference in

188     branch length was a property of some or all single-gene phylogenies, we compared the difference

189     in length of the FEL and SEL stem branches among all single-gene trees where each lineage was

190     recovered monophyletic ($n = 946$). We found that the FEL stem branch was nearly four times

191     longer ($0.62 \pm 0.38$ substitutions / site) than the SEL stem branch ($0.17 \pm 0.11$ substitutions /

192     site) (Fig 1B; $p < 0.001$; Paired Wilcoxon Rank Sum test). Furthermore, of the 946 gene trees

193     examined, 932 had a much longer FEL stem branch (0.46 ± 0.33 Δ substitutions / site), whereas

194     only 14 had a slightly longer SEL stem branch (0.06 ± 0.05 Δ substitutions / site).

195

196     **The genomes of FEL species have lost substantial numbers of genes**

197     Examination of GC content, genome size, and gene number revealed that the some of the lowest

198     GC content values, as well as the smallest genomes and lowest gene numbers, across the

199     subphylum Saccharomycotina are primarily observed in FEL yeasts (Fig S1). Specifically, the

200     average GC contents for FEL yeasts (33.10 ± 3.53%), SEL yeasts (37.28 ± 2.05%), and all other

201     Saccharomycotina yeasts (40.77 ± 5.58%) are significantly different ($\chi^2(2) = 30.00$, $p < 0.001$;

202     Kruskal-Wallis rank sum test). Further examination revealed only the FEL was significantly

203     different from other Saccharomycotina yeasts ($p < 0.001$; Dunn's test for multiple comparisons

204     with Benjamini-Hochberg multi-test correction). For genome size and gene number, FEL yeast

205     genomes have average sizes of 9.71 ± 1.32 Mb and contain 4,707.89 ± 633.56 genes,

206     respectively, while SEL yeast genomes have average sizes of 10.99 ± 1.66 Mb and contain

207     4,932.43 ± 289.71 genes. In contrast, all other Saccharomycotina have average genome sizes and

208     gene numbers of 13.01 ± 3.20 Mb and 5,726.10 ± 1,042.60, respectively. Statistically significant

209     differences were observed between the FEL, SEL, and all other Saccharomycotina (genome size:

210     $\chi^2(2) = 33.47$, $p < 0.001$ and gene number: $\chi^2(2) = 31.52$, $p < 0.001$; Kruskal-Wallis rank sum

211     test for both). Further examination revealed the only significant difference for genome size was

212     between FEL and other Saccharomycotina yeasts ($p < 0.001$; Dunn's test for multiple

213     comparisons with Benjamini-Hochberg multi-test correction), while both the FEL and SEL had

214     smaller gene sets compared to other Saccharomycotina yeasts ($p < 0.001$ and $p = 0.008$,

215     respectively; Dunn's test for multiple comparisons with Benjamini-Hochberg multi-test

10

216    correction). The lower numbers of genes in the FEL (especially) and SEL lineages were also

217    supported by gene content completeness analyses using orthologous sets of genes constructed

218    from sets of genomes representing multiple taxonomic levels across eukaryotes (Fig S5) from the

219    ORTHODB database [43].

220

221    To further examine which genes have been lost in the genomes of FEL and SEL species relative

222    to other representative Saccharomycotina genomes, we conducted HMM-based sequence

223    similarity searches using annotated *S. cerevisiae* genes as queries in HMM construction (see

224    *Methods*) (Fig S6). Because we were most interested in identifying genes absent from the FEL

225    and SEL, we focused our analyses on genes lost in at least two-thirds of each lineage (i.e., $\geq 11$

226    FEL taxa or $\geq 5$ SEL taxa). Using this criterion, we found that 1,409 and 771 genes have been

227    lost in the FEL and SEL, respectively (Fig 2A). Among the genes lost in each lineage, 748 genes

228    were lost across both lineages, 661 genes have been uniquely lost in the FEL, and 23 genes have

229    been uniquely lost in the SEL (File S3).

230

231    To identify the likely functions of genes lost from each lineage, we conducted GO enrichment

232    analyses. Examination of significantly over-represented GO terms for the sets of genes that have

233    been lost in *Hanseniaspora* genomes revealed numerous categories related to metabolism (e.g.,

234    MALTOSE METABOLIC PROCESS, GO:0000023, $p = 0.006$; SUCROSE ALPHA-GLUCOSIDASE

235    ACTIVITY, GO:0004575, $p = 0.003$) and genome-maintenance processes (e.g., MEIOTIC CELL

236    CYCLE, GO:0051321, $p < 0.001$) (File S4). Additional terms, such as CELL CYCLE, GO:0007049

237    ($p < 0.001$), CHROMOSOME SEGREGATION, GO:0007059 ($p < 0.001$), CHROMOSOME

238    ORGANIZATION, GO:0051276 ($p = 0.009$), and DNA-DIRECTED DNA POLYMERASE ACTIVITY,

11

239    GO:0003887 ($p < 0.001$), were significantly over-represented among genes absent only in the

240    FEL. Next, we examined in more detail the identities and likely functional consequences of

241    extensive gene losses across *Hanseniaspora* associated with metabolism, cell cycle, and DNA

242    repair.

243

244    *Metabolism-associated gene losses.*       Examination of the genes causing over-

245    representation of metabolism-associated GO terms revealed gene losses in the *IMA* gene family

246    and the *MAL* loci, both of which are associated with growth primarily on maltose but can also

247    facilitate growth on sucrose, raffinose, and melezitose [44,45]. All *IMA* genes have been lost in

248    *Hanseniaspora*, whereas *MALx3*, which encodes the *MAL*-activator protein [46] has been lost in

249    all but one species (*Hanseniaspora jakobsenii*; Fig 2B). Consistent with these losses,

250    *Hanseniaspora* species cannot grow on the carbon substrates associated with these genes (i.e.,

251    maltose, raffinose, and melezitose) with the exception of *H. jakobsenii*, which has weak/delayed

252    growth on maltose (Fig 2B; File S5). The growth of *H. jakobsenii* on maltose may be due to a

253    cryptic α-glucosidase gene or represent a false positive, as *MALx2* encodes the required enzyme

254    for growth on maltose and is absent in *H. jakobsenii*. Because these genes are also associated

255    with growth on sucrose in some species [44], we also examined their ability to grow on this

256    substrate. In addition to the *MAL* loci conferring growth on sucrose, the invertase Suc2 can also

257    break down sucrose into glucose and fructose [47]. We found that FEL yeasts have lost *SUC2*

258    and are unable to grow on sucrose, while SEL yeasts have *SUC2* and are able to grow on this

259    substrate (Fig 2B; File S5). Altogether, patterns of gene loss are consistent with known metabolic

260    traits.

261

262    Examination of gene sets associated with growth on other carbon substrates revealed that

263    *Hanseniaspora* species also cannot grow on galactose, consistent with the loss of one or more of

264    the three genes involved in galactose assimilation (*GAL1, GAL7,* and *GAL10*) from their

265    genomes (Fig 2C; File S5). Additionally, all *Hanseniaspora* genomes appear to have lost two

266    key genes, *PCK1* and *FBP1*, encoding enzymes in the gluconeogenesis pathway (Fig S7A and

267    S7C); in contrast, all *Hanseniaspora* have an intact glycolysis pathway (Fig S7B and S7D).

268

269    Manual examination of other metabolic pathways revealed that *Hanseniaspora* genomes are also

270    missing some of their key genes. For example, we found that THIAMINE BIOSYNTHETIC PROCESS,

271    GO:0009228 ($p$ = 0.003), was an over-represented GO term among genes missing in both the

272    FEL and SEL due to the absence of *THI* and *SNO* family genes. Further examination of genes

273    present in the thiamine biosynthesis pathway revealed extensive gene loss (Fig 2D), which is

274    consistent with their inability to grow on vitamin-free media [45] (File S5). Notably,

275    *Hanseniaspora* are still predicted to be able to import extracellular thiamine via Thi73 and

276    convert it to its active cofactor via Thi80, which may explain why they can rapidly consume

277    thiamine [39]. Similarly, examination of amino acid biosynthesis pathways revealed the

278    methionine salvage pathway was also largely disrupted by gene losses across all *Hanseniaspora*

279    (Fig 2E). Lastly, we found that *GDH1* and *GDH3* from the glutamate biosynthesis pathway from

280    ammonium are missing in FEL yeasts (File S3). However, *Hanseniaspora* have *GLT1,* which

281    enables glutamate biosynthesis from glutamine.

282

283    *Cell cycle and genome integrity-associated gene losses.*              Many genes involved in cell

284    cycle and genome integrity, including cell cycle checkpoint genes, have been lost across

13

285    *Hanseniaspora* (Fig 3). For example, *WHI5* and *DSE2*, which are responsible for repressing the

286    Start (i.e., an event that determines cells have reached a critical size before beginning division)

287    [48] and help facilitate daughter-mother cell separation through cell wall degradation [49], have

288    been lost in both lineages. Additionally, the FEL has lost the entirety of the DASH complex (i.e.,

289    *ASK1*, *DAD1*, *DAD2*, *DAD3*, *DAD4*, *DUO1, DAM1*, *HSK3, SPC19*, and *SPC34*), which forms

290    part of the kinetochore and functions in spindle attachment and stability, as well as chromosome

291    segregation, and the MIND complex (i.e., *MTW1*, *NNF1*, *NSL1*, and *DSN1*), which is required

292    for kinetochore bi-orientation and accurate chromosome segregation (File S3 and S4). Similarly,

293    FEL species have lost *MAD1* and *MAD2*, which are associated with spindle checkpoint processes

294    and have abolished checkpoint activity when their encoded proteins are unable to dimerize [14].

295    Lastly, components of the anaphase-promoting complex, a major multi-subunit regulator of the

296    cell cycle, are lost in both lineages (i.e., *CDC26* and *MND2*) or just the FEL (i.e., *APC2, APC4,*

297    *APC5,* and *SWM1*).

298

299    Another group of genes that have been lost in *Hanseniaspora* are genes associated with the DNA

300    damage checkpoint and DNA damage sensing. For example, both lineages have lost *RFX1*,

301    which controls a late point in the DNA damage checkpoint pathway [50], whereas the FEL has

302    lost *MEC3* and *RAD9*, which encode checkpoint proteins required for arrest in the G2 phase after

303    DNA damage has occurred [17]. Since losses in DNA damage checkpoints and dysregulation of

304    spindle checkpoint processes are associated with genomic instability, we next evaluated the

305    ploidy of  *Hanseniaspora* genomes [51]. Using base frequency plots, we found that the ploidy of

306    genomes of FEL species ranges between 1 and 3, with evidence suggesting that certain species,

307    such as *H. singularis, H. pseudoguilliermondii*, and *H. jakobsenii,* are potentially aneuploid (Fig

14

308    S8). In contrast, the genomes of SEL species have ploidies of 1-2 with evidence of potential

309    aneuploidy observed only in *H. occidentalis* var. *citrica.* Greater variance in ploidy and

310    aneuploidy in the FEL compared to the SEL may be due to the FEL's loss of a greater number of

311    components of the anaphase-promoting complex (APC), whose dysregulation is thought to

312    increase instances of aneuploidy [52].

313

314    *Pronounced losses of DNA repair genes in the FEL.*        Examination of other GO-enriched

315    terms revealed numerous genes associated with diverse DNA repair processes that have been lost

316    among *Hanseniaspora* species, and especially the FEL (Fig 4). We noted 14 lost DNA repair

317    genes across all *Hanseniaspora*, including the DNA glycosylase gene *MAG1* [53], the photolyase

318    gene *PHR1* that exclusively repairs pyrimidine dimers [23], and the diphosphatase gene *PCD1*, a

319    key contributor to the purging of mutagenic nucleotides, such as 8-oxo-dGTP, from the cell [24].

320    An additional 33 genes were lost specifically in the FEL such as *TDP1*, which repairs damage

321    caused by topoisomerase activity [54]; the DNA polymerase gene *POL32* that participates in

322    base-excision and nucleotide-excision repair and whose null mutants have increased genomic

323    deletions [20]; and the *CDC13* gene that encodes a telomere-capping protein [55].

324

325    **FEL gene losses are associated with accelerated sequence evolution**

326    *Loss of DNA repair genes is associated with a burst of sequence evolution.*        To examine

327    the mutational signatures of losing numerous DNA repair genes on *Hanseniaspora* substitution

328    rates, we tested several different hypotheses that postulated changes in the ratio of the rate of

329    nonsynonymous (dN) to the rate of synonymous substitutions (dS) (dN/dS or ω) along the

330    phylogeny (Table 1; Fig 5). For each hypothesis tested, the null was that the ω value remained

15

331      constant across all branches of the phylogeny. Examination of the hypothesis that the $\omega$ values of

332      both the FEL and SEL stem branches were distinct from the background $\omega$ value ($H_{FE-SE\ branch}$;

333      Fig 5B), revealed that 678 genes (68.55% of examined genes) significantly rejected the null

334      hypothesis (Table 1; $\alpha = 0.01$; LRT; median FEL stem branch $\omega = 0.57$, median SEL stem

335      branch $\omega = 0.29$, and median background $\omega = 0.060$). Examination of the hypothesis that the $\omega$

336      value of the FEL stem branch and the $\omega$ value of the FEL crown branches were distinct from the

337      background $\omega$ value ($H_{FE}$; Fig 5C) revealed 743 individual genes (75.13% of examined genes)

338      that significantly rejected the null hypothesis (Table 1; $\alpha = 0.01$; LRT; median FEL stem branch

339      $\omega = 0.71$, median FEL crown branches $\omega = 0.06$, median background $\omega = 0.063$). Testing the

340      same hypothesis for the SEL ($H_{SE}$; Fig 5D) revealed 528 individual genes (53.7% of examined

341      genes) that significantly rejected the null hypothesis (Table 1; $\alpha = 0.01$; LRT; median SEL stem

342      branch $\omega = 0.267$, median SEL crown branches $\omega = 0.074$, median background $\omega = 0.059$).

343      Finally, testing of the hypothesis that the FEL and SEL crown branches have $\omega$ values distinct

344      from each other and the background ($H_{FE-SE\ crown}$; Fig 5E) revealed 717 genes (72.5% of

345      examined genes) that significantly rejected the null hypothesis (Table 1; $\alpha = 0.01$; LRT; median

346      FEL crown branches $\omega = 0.062$, median SEL crown branches $\omega = 0.074$, median background $\omega =$

347      0.010). These results suggest a dramatic, genome-wide increase in evolutionary rate in the FEL

348      stem branch (Fig 5B and 5C), which coincided with the loss of a large number of genes involved

349      in DNA repair.

350

351      *The FEL has a greater number of base substitutions and indels.*      To better understand

352      the mutational landscape in the FEL and SEL, we characterized patterns of base substitutions

353      across the 1,034 OGs. Focusing on first ($n = 240,565$), second ($n = 318,987$), and third ($n =$

354     58,151) codon positions that had the same character state in all outgroup taxa, we first examined

355     how many of these sites had experienced base substitutions in FEL and SEL species (Fig 6A).

356     We found significant differences between the proportions of base substitutions in the FEL and

357     SEL ($F(1) = 196.88$, $p < 0.001$; Multi-factor ANOVA) at each codon position (first: $p < 0.001$;

358     second: $p < 0.001$; and third: $p = 0.02$; Tukey Honest Significance Differences post-hoc test).

359

360     Examination of whether the observed base substitutions were AT- (i.e., G|C → A|T) or GC- (i.e.,

361     A|T → G|C) biased revealed differences between the FEL and SEL ($F(1) = 447.1$, $p < 0.001$;

362     Multi-factor ANOVA), as well as between AT- and GC-bias ($F(1) = 914.5$, $p < 0.001$; Multi-

363     factor ANOVA) among sites with G|C ($n = 232,546$) and A|T ($n = 385,157$) pairs (Fig 6B).

364     Specifically, we observed significantly more base substitutions in the FEL compared to the SEL

365     and a significant bias toward A|T across both lineages ($p < 0.001$ for both tests; Tukey Honest

366     Significance Differences post-hoc test). Examination of transition / transversion ratios revealed a

367     lower transition / transversion ratio in the FEL ($0.67 \pm 0.02$) compared to the SEL ($0.76 \pm 0.01$)

368     (Fig 6C; $p < 0.001$; Wilcoxon Rank Sum test); this finding is in contrast to the transition /

369     transversion ratios found in most known organisms, whose values are substantially above 1.00

370     [56–59]. Altogether, these analyses reveal more base substitutions in the FEL and SEL across all

371     codon positions and a significant AT-bias in base substitutions across all *Hanseniaspora*.

372

373     Examination of indels revealed that the total number of insertions or deletions was significantly

374     greater in the FEL ($mean_{insertions} = 7521.11 \pm 405.34$; $mean_{deletions} = 3894.11 \pm 208.16$) compared

375     to the SEL ($mean_{insertions} = 6049.571 \pm 155.85$; $mean_{deletions} = 2346.71 \pm 326.22$) (Fig 6D; $p <$

376     0.001 for both tests; Wilcoxon Rank Sum test). The difference in number of indels between the

17

377    FEL and SEL remained significant after taking into account indel size (F(1) = 2102.87, $p <$

378    0.001; Multi-factor ANOVA). Further analyses revealed there are significantly more insertions

379    in the FEL compared to the SEL for insertion sizes 3-18 bp ($p < 0.001$ for all comparisons

380    between each lineage for each insertion size; Tukey Honest Significance Differences post-hoc

381    test), while there were significantly more deletions in the FEL compared to the SEL for deletion

382    sizes 3-21 bp ($p < 0.001$ for all comparisons between each lineage for each deletion size; Tukey

383    Honest Significance Differences post-hoc test). These analyses suggest that there are

384    significantly more indels in the FEL compared to the SEL and that this pattern is primarily

385    driven by short indels.

386

387    **Greater sequence instability in the FEL and signatures of endogenous and exogenous DNA**

388    **damage**

389    *The FEL has greater instability of homopolymers.*            Examination of the total proportion

390    of mutated bases among homopolymers (i.e., (substituted bases + deleted bases + inserted bases)

391    / total homopolymer bases) revealed significant differences between the FEL and SEL (Fig 6G;

392    F(1) = 27.68, $p < 0.001$; Multi-factor ANOVA). Although the FEL had a higher proportion of

393    mutations among homopolymers across all sizes of two ($n = 17,391$), three ($n = 1,062$), four ($n =$

394    104), and five ($n = 5$), significant differences were observed for homopolymers of length two and

395    three ($p = 0.02$ and $p = 0.003$, respectively; Tukey Honest Significance Differences post-hoc). To

396    gain more insight into the drivers differentiating mutational load in homopolymers, we

397    considered the additional factors of homopolymer sequence type (i.e., A|T or C|G) and mutation

398    type (i.e., base substitution, insertion, or deletion) (Fig S9). In addition to recapitulating

399    differences between the types of mutations that occur at homopolymers (F(2) = 1686.70, $p <$

18

400    0.001; Multi-factor ANOVA), we observed that base substitutions occurred more frequently than

401    insertions and deletions ($p < 0.001$ for both tests; Tukey Honest Significance Differences post-

402    hoc test). For example, among A|T and C|G homopolymers of length two and C|G

403    homopolymers of length three, base substitutions were higher in the FEL compared to the SEL ($p$

404    $= 0.009$, $p < 0.001$, and $p < 0.001$, respectively; Tukey Honest Significance Differences post-hoc

405    test). Additionally, there were significantly more base substitutions in A|T homopolymers of

406    length five in the FEL compared to the SEL ($p < 0.001$; Tukey Honest Significance Differences

407    post-hoc test). Altogether, these analyses reveal greater instability of homopolymers in the FEL

408    compared to the SEL due to more base substitutions.

409

410    *The FEL has a stronger signature of endogenous DNA damage from 8-oxo-dG.*   Examination

411    of mutational signatures associated with common endogenous and exogenous mutagens revealed

412    greater signatures of mutational load in the FEL compared to the SEL, as well as in both FEL

413    and SEL compared to the outgroup taxa. The oxidatively damaged guanine base, 8-oxo-dG, is a

414    commonly observed endogenous form of DNA damage that causes the transversion mutation of

415    G $\rightarrow$ T or C $\rightarrow$ A [26]. Examination of the direction of base substitutions among all sites with a

416    G base in all outgroup taxa revealed differences in the direction of base substitutions ($F(2) =$

417    $5,682$, $p < 0.001$; Multi-factor ANOVA). Moreover, there are significantly more base

418    substitutions at G sites associated with 8-oxo-dG damage in the FEL compared to the SEL (Fig

419    6H; $p < 0.001$; Tukey Honest Significance Differences post-hoc test). These analyses reveal that

420    FEL genomes have higher proportions of G site substitutions associated with the mutational

421    signature of a common endogenous mutagen.

422

423   Hanseniaspora *have a greater genomic signature of UV-damage.* Both the FEL and SEL have

424   lost *PHR1*, a gene encoding a DNA photolyase that repairs pyrimidine dimers, so we next

425   examined the genomes for evidence of a CC → TT dinucleotide substitution bias, an indirect

426   molecular signature of UV radiation damage (Fig 6I). To do so, we used a CC|GG and TT|AA

427   score, which quantifies the abundance of CC|GG and TT|AA dinucleotides in a genome and

428   corrects for the total number of dinucleotides and GC content in the same genome. When

429   comparing CC|GG scores between the FEL, SEL, and outgroup taxa, there were no significant

430   differences ($\chi^2(2) = 5.96$, $p = 0.051$; Kruskal-Wallis rank sum test). When comparing all

431   *Hanseniaspora* to the outgroup, we found that the CC|GG score was significantly lower in

432   *Hanseniaspora* ($p = 0.03$; Wilcoxon Rank Sum test). Examination of TT|AA scores revealed

433   significant differences between the three groups ($\chi2(2) = 8.84$, $p = 0.012$; Kruskal-Wallis rank

434   sum test), which was driven by differences between the FEL and SEL compared to the outgroup

435   ($p = 0.011$ and $0.016$, respectively; Dunn's test for multiple comparisons with Benjamini-

436   Hochberg multi-test correction). The same result was observed when comparing all

437   *Hanseniaspora* to the outgroup ($p < 0.001$; Wilcoxon Rank Sum test). Altogether, these analyses

438   suggest *Hanseniaspora* have a greater signal of UV damage compared to other budding yeasts.

439

440   Lastly, we examined if all of these mutations were associated with more radical amino acid

441   changes in the FEL compared to the SEL using two measures of amino acid change: Sneath's

442   index [60] and Epstein's coefficient of difference [61]. For both measures, we observed

443   significantly more radical amino acid substitutions in the FEL compared to the SEL (Fig S10; *p*

444   *< 0.001*; Wilcoxon Rank Sum test for both metrics). Altogether, these analyses reveal greater

20

445     DNA sequence instability in the FEL compared to the SEL, which is also associated with more

446     radical amino acid substitutions.

447

448     **<u>Discussion</u>**

449     The genus *Hanseniaspora* has been recently observed to exhibit the longest branches among

450     budding yeasts (Fig 1) [33–35], and their genomes have some of the lowest numbers of genes,

451     lowest GC contents, and smallest assembly sizes in the subphylum (Fig S1). Through the

452     analysis of the genomes of nearly every known *Hanseniaspora* species this study presents

453     multiple lines of evidence suggesting that one lineage of *Hanseniaspora*, which we have named

454     FEL, is a lineage of long-term, hypermutator species that have undergone extensive gene loss

455     (Figs. 1-4).

456

457     Evolution by gene loss is gaining increasing attention as a major mode of genome evolution

458     [34,62] and is mainly possible due to the dispensability of the majority of genes. For example,

459     90% of *E. coli* [63]*,* 80% of *S. cerevisiae* [64], and 73% of *Candida albicans* [65] genes are

460     dispensable in laboratory conditions*.* The loss of dispensable genes can be selected for [66] and

461     is common in lineages of obligate parasites or symbionts, such as in the microsporidia,

462     intracellular fungi which have lost key metabolic pathways such as amino acid biosynthesis

463     pathways [67,68], and myxozoa, a group of cnidarian obligate parasites that infect vertebrates

464     and invertebrates [69]. Similar losses are also increasingly appreciated in free-living organisms,

465     such as the budding yeasts (this study; Hittinger et al., 2004; Riley et al., 2016; Shen et al., 2018;

466     Slot and Rokas, 2010; Wolfe et al. 2015) and animals [62]. For example, a gene known to enable

21

467     sucrose utilization, *SUC2* [47], is lost in the FEL and reflects an inability to grow on sucrose,

468     while the *SUC2* is present in the SEL and reflects an ability to grow on sucrose (Fig 2).

469

470     However, *Hanseniaspora* species have experienced not just the typically observed losses of

471     metabolic genes (Figs. 2A and 2B), but more strikingly, the atypical loss of dozens of cell cycle

472     and DNA damage, response, and repair genes (Figs. 3 and 4). Losses of cell cycle genes are

473     extremely rare [11], and most such losses are known in the context of cancers [73]. Losses of

474     individual or a few DNA repair genes have also been observed in individual hypermutator fungal

475     isolates [6–8]. In contrast, the *Hanseniaspora* losses of cell cycle and DNA repair genes are not

476     only unprecedented in terms of the numbers of genes lost and their striking impact on genome

477     sequence evolution, but also in terms of the evolutionary longevity of the lineage.

478

479     *Missing checkpoint processes are associated with fast growth and bipolar budding.*

480     *Hanseniaspora* species lost numerous components of the cell cycle (Fig 3), such as *WHI5*, which

481     causes accelerated G1/S transitions in knock-out *S. cerevisiae* strains [12,48], as well as

482     components of APC (i.e., *CDC26* and *MND2*), which may accelerate the transition to anaphase

483     [13]. These and other cell cycle gene losses are suggestive of rapid cell division and growth and

484     consistent with the known ability of *Hanseniaspora* yeast of rapid growth in the wine

485     fermentation environment [40].

486

487     One of the distinguishing characteristics of the *Hanseniaspora* cell cycle is bipolar budding,

488     which is known only in the genera *Wickerhamia* (Debaryomycetaceae) and *Nadsonia*

489     (Dipodascaceae), as well as in *Hanseniaspora* and its sister genus *Saccharomycodes* (both in the

22

490    family Saccharomycodaceae) [45][74]. These three lineages are distantly related to one another

491    on the budding yeast phylogeny [34], so bipolar budding likely evolved three times

492    independently in Saccharomycotina, including in the last common ancestor of *Hanseniaspora*

493    and *Saccharomycodes*. Currently, there is only one genome available for *Saccharomycodes* [74],

494    making robust inferences of ancestral states challenging. Interestingly, examination of cell cycle

495    gene presence and absence in the only representative genome from the genus, *Saccharomycodes*

496    *ludwigii* [74], reveals that *CDC26*, *PCL1*, *PDS1*, *RFX1*, *SIC1*, *SPO12*, and *WHI5* are absent (File

497    S6), most of which are either absent from all *Hanseniaspora* (i.e., *CDC26*, *RFX1*, *SPO12*, and

498    *WHI5*) or just from the FEL (i.e., *PDS1* and *SIC1*). This evidence raises the hypothesis that

499    bipolar budding is linked to the dysregulation of cell cycle processes due to the absence of cell

500    cycle genes and in particular cell cycle checkpoints (Fig 3).

501

502    *Some gene losses may be compensatory.*    Deletion of many of the genes associated with DNA

503    maintenance that have been lost in *Hanseniaspora* lead to dramatic increases of mutation rates

504    and gross genome instability [12,13,20], raising the question of how these gene losses were

505    tolerated in the first place. Examination of the functions of the genes lost in *Hanseniaspora*

506    suggests that at least some of these gene losses may have been compensatory. For example,

507    *POL4* knock-out strains of *S. cerevisiae* can be rescued by the deletion of *YKU70* [75]*,* both of

508    which were lost in the FEL. Similarly, the loss of genes responsible for key cell cycle functions

509    (e.g., kinetochore functionality and chromosome segregation) appears to have co-occurred with

510    the loss of checkpoint genes responsible for delaying the cell cycle if its functions fail to

511    complete, which may have allowed *Hanseniaspora* cells to bypass otherwise detrimental cell

512    cycle arrest. Specifically, *MAD1* and *MAD2*, which help delay anaphase when kinetochores are

513     unattached [14]; the 10-gene DASH complex, which participates in spindle attachment, stability,

514     and chromosome segregation [76]; and the 4-gene MIND complex, which is required for

515     kinetochore bi-orientation and accurate chromosome segregation [77], were all lost in the FEL.

516

517     *Long-term hypermutation and the subsequent slowing of sequence evolution.*        Estimates of ω

518     suggest the FEL and SEL, albeit to a much lower degree in the latter, underwent a burst of

519     accelerated sequence evolution in their stem lineages, followed by a reduction in the pace of

520     sequence evolution (Fig 5). This pattern is consistent with theoretical predictions that selection

521     against mutator phenotypes will reduce the overall rate of sequence evolution [27], as well as

522     with evidence from experimental evolution of hypermutator lines of *S. cerevisiae* that showed

523     that their mutation rates were quickly reduced [32]. Although we do not know the catalyst for

524     this burst of sequence evolution, hypermutators may be favored in maladapted populations or in

525     conditions where environmental parameters frequently change [27,32]. Although the

526     environment occupied by the *Hanseniaspora* last common ancestor is unknown, it is plausible

527     that environmental instability or other stressors favored hypermutators in *Hanseniaspora*. Extant

528     *Hanseniaspora* species are well known to be associated with the grape environment [39,78,79].

529     Interestingly, grapes appear to have originated  [80] around the same time window that

530     *Hanseniaspora* did (Fig 1B), leading us to speculate that the evolutionary trigger of

531     *Hanseniaspora* hypermutation could have been adaptation to the grape environment.

532

533     *Losses of DNA repair genes are reflected in patterns of sequence evolution.*        Although the

534     relationship between genotype and phenotype is complex, the loss of genes involved in DNA

535     repair can have predictable outcomes on patterns of sequence evolution in genomes. In the case

24

536    of the observed losses of DNA repair genes in *Hanseniaspora*, the mutational signatures of this

537    loss and the consequent hypermutation can be both general (i.e., the sum total of many gene

538    losses), as well as specific (i.e., can be putatively linked to the losses of specific genes or

539    pathways). Arguably the most notable general mutational signature is that *Hanseniaspora*

540    genome sequence evolution is largely driven by random (i.e., neutral) mutagenic processes with

541    a strong AT-bias. For example, whereas the transition / transversion ratios of eukaryotic

542    genomes are typically within the 1.7 and 4 range [56–59], *Hanseniaspora* ratios are ~0.66-0.75

543    (Fig 6C), which are values on par with estimates of transition / transversion caused by neutral

544    mutations alone (e.g., 0.6-0.95 in *S. cerevisiae* [56,81], 0.92 in *E. coli* [82], 0.98 in *Drosophila*

545    *melanogaster* [83], and 1.70 in humans [84]). Similarly, base substitutions across *Hanseniaspora*

546    genomes are strongly AT-biased, especially in the FEL (Fig 6), an observation consistent with

547    the general AT-bias of mutations observed in diverse organisms, including numerous bacteria

548    [85], the fruit fly [83], *S. cerevisiae* [56], and humans [84].

549

550    In addition to these general mutational signatures, examination of *Hanseniaspora* sequence

551    evolution also reveals mutational signatures that can be linked to the loss of specific DNA repair

552    genes. For example, we found a higher proportion of base substitutions associated with the most

553    abundant oxidatively damaged base, 8-oxo-dG, which causes G → T or C → A transversions

554    [26], in the FEL compared to the SEL, which reflects specific gene losses. Specifically,

555    *Hanseniaspora* yeasts have lost *PCD1*, which encodes a diphosphatase that contributes to the

556    removal of 8-oxo-dGTP [24] and thereby reduces the chance of misincorporating this damaged

557    base. Once 8-oxo-dG damage has occurred, it is primarily repaired by the base excision repair

558    pathway [26]. Notably, the FEL is missing a key component of the base excision repair pathway,

25

559    a DNA polymerase $\delta$ subunit, encoded by *POL32*, which aids in filling the gap after excision

560    [86]. Accordingly, the proportion of G|C sites with substitutions indicative of 8-oxo-dG damage

561    (i.e., G → T or C → A transversions) is significantly greater in the FEL compared to the SEL

562    (Fig 5H). Similarly, the numbers of dinucleotide substitutions of CC → TT associated with UV-

563    induced pyrimidine dimers [87] are higher across *Hanseniaspora* compared to other yeasts due

564    to the loss of *PHR1*, which encodes a DNA photolyase that repairs pyrimidine dimers (Fig 5I)

565    [23].

566

567    Our analyses provide the first major effort to characterize the genome function and evolution of

568    the enigmatic genus *Hanseniaspora* and identify major and extensive losses of genes associated

569    with metabolism, cell cycle, and DNA repair processes. These extensive losses and the

570    concomitant acceleration of evolutionary rate mean that levels of amino acid sequence

571    divergence within each of the two *Hanseniaspora* lineages alone, but especially within the FEL,

572    are similar to those observed within plant classes and animal subphyla (Fig S11). These

573    discoveries set the stage for further fundamental molecular and evolutionary investigations

574    among *Hanseniaspora,* such as potential novel rewiring of cell cycle and DNA repair processes.

575

576    **Methods**

577    **DNA sequencing**       For each species, genomic DNA (gDNA) was isolated using a two-step

578    phenol:chloroform extraction previously described to remove additional proteins from the gDNA

579    [34]. The gDNA was sonicated and ligated to Illumina sequencing adaptors as previously

580    described [88], and the libraries were submitted for paired-end sequencing (2 x 250) on an

581    Illumina HiSeq 2500 instrument.

582

583     **Phenotyping** We qualitatively measured growth of species on five carbon sources (maltose,

584     raffinose, sucrose, melezitose, and galactose) as previously described in [34]. We used a minimal

585     media base with ammonium sulfate and all carbon sources were at a 2% concentration. Yeast

586     were initially grown in YPD and transferred to carbon treatments. Species were visually scored

587     for growth for about a week on each carbon source in three independent replicates over multiple

588     days. A species was considered to utilize a carbon source if it showed growth across $\geq$ 50% of

589     biological replicates. Growth data for *Hanseniaspora gamundiae* were obtained from Čadež et

590     al., 2019.

591

592     **Genome assembly and annotation** To generate *de novo* genome assemblies, we used paired-

593     end DNA sequence reads as input to iWGS, version 1.1 [89], a pipeline which uses multiple

594     assemblers and identifies the "best" assembly according to largest genome size and N50 (i.e., the

595     shortest contig length among the set of the longest contigs that account for 50% of the genome

596     assembly's length) [90] as described in [34]. More specifically, sequenced reads were first

597     quality-trimmed, and adapter sequences were removed used TRIMMOMATIC, version 0.33 [91],

598     and LIGHTER, version 1.1.1 [92]. Subsequently, KMERGENIE, version 1.6982 [93], was used to

599     determine the optimal *k*-mer length for each genome individually. Thereafter, six *de novo*

600     assembly tools (i.e., ABYSS, version 1.5.2 [94]; DISCOVAR, release 51885 [95]; MASURCA,

601     version 2.3.2 [96]; SGA, version 0.10.13 [97]; SOAPDENOVO2, version 2.04 [98]; and SPADES,

602     version 3.7.0 [99]) were used to generate genome assemblies from the processed reads. Using

603     QUAST, version 4.4 [100], the best assembly was chosen according to the assembly that

604     provided the largest genome size and best N50.

27

605

606    Annotations for eight of the *Hanseniaspora* genomes (i.e., *H. clermontiae*, *H. osmophila* CBS

607    313, *H. pseudoguilliermondii*, *H. singularis*, *H. uvarum* DSM2768, *H. valbyensis*, *H. vineae* T02

608    19AF, and *K. hatyaiensis*) and the four outgroup species (i.e., *Cy. jadinii*, *K. marxianus*, *S.*

609    *cerevisiae*, and *W. anomalus*) were generated in a recent comparative genomic study of the

610    budding yeast subphylum [34]. The other 11 *Hanseniaspora* genomes examined here were

611    annotated by following the same protocol as in [34].

612

613    In brief, the genomes were annotated using the MAKER pipeline, version 2.31.8 [101]. The

614    homology evidence used for MAKER consists of fungal protein sequences in the SwissProt

615    database (release 2016_11) and annotated protein sequences of select yeast species from

616    MYCOCOSM [102], a web portal developed by the US Department of Energy Joint Genome

617    Institute for fungal genomic analyses. Three *ab initio* gene predictors were used with the

618    MAKER pipeline, including GENEMARK-ES, version 4.32 [103]; SNAP, version 2013-11-29

619    [104]; and AUGUSTUS, version 3.2.2 [105], each of which was trained for each individual

620    genome. GENEMARK-ES was self-trained on the repeat-masked genome sequence with the

621    fungal-specific option ("–fugus"), while SNAP and AUGUSTUS were trained through three

622    iterative MAKER runs. Once all three *ab initio* predictors were trained, they were used together

623    with homology evidence to conduct a final MAKER analysis in which all gene models were

624    reported ("keep_preds" set to 1), and these comprise the final set of annotations for the genome.

625

626    *Data acquisition*        All publicly available *Hanseniaspora* genomes, including multiple strains

627    from a single species, were downloaded from NCBI (https://www.ncbi.nlm.nih.gov/ File S1).

28

628    These species and strains include *H. guilliermondii* UTAD222 [78], *H. opuntiae* AWRI3578, *H.*

629    *osmophila* AWRI3579, *H. uvarum* AWRI3580 [106], *H. uvarum* 34-9, *H. vineae* T02-19AF

630    [107], *H. valbyensis* NRRL Y-1626 [33], and *H. gamundiae* [41]. We also included

631    *Saccharomyces cerevisiae* S288C, *Kluyveromyces marxianus* DMKU3-1042, *Wickerhamomyces*

632    *anomalus* NRRL Y-366-8, and *Cyberlindnera jadinii* NRRL Y-1542, four representative

633    budding yeast species that are all outside the genus *Hanseniaspora* [34], which we used as

634    outgroups. Together with publicly available genomes, our sampling of *Hanseniaspora*

635    encompasses all known species in the genus (or its anamorphic counterpart, *Kloeckera*), except

636    *Hanseniaspora lindneri*, which likely belongs to the FEL based on a four-locus phylogenetic

637    study [108], and *Hanseniaspora taiwanica*, which likely belongs to the SEL based on neighbor-

638    joining analyses of the LSU rRNA gene sequence [109].

639

640    **Assembly assessment and identification of orthologs**        To determine genome assembly

641    completeness, we calculated contig N50 [90] and assessed gene content completeness using

642    multiple databases of curated orthologs from BUSCO, version 3 [110]. More specifically, we

643    determined gene content completeness using orthologous sets of genes constructed from sets of

644    genomes representing multiple taxonomic levels, including Eukaryota (superkingdom; 100

645    species; 303 BUSCOs), Fungi (kingdom; 85 species; 290 BUSCOs), Dikarya (subkingdom; 75

646    species; 1,312 BUSCOs), Ascomycota (phylum; 75 species; 1,315 BUSCOs), Saccharomyceta

647    (no rank; 70 species; 1,759 BUSCOs), and Saccharomycetales (order; 30 species; 1,711

648    BUSCOs).

649

29

650    Genomes sequenced in the present project were sequenced at an average depth of 63.49 ± 52.57

651    (File S1). Among all *Hanseniaspora*, the average scaffold N50 was 269.03 ± 385.28 kb, the

652    average total number of scaffolds was 980.36 ± 835.20 (398.32 ± 397.97 when imposing a 1kb

653    scaffold filter), and the average genome assembly size was 10.13 ± 1.38 Mb (9.93 ± 1.35 Mb

654    when imposing a 1kb scaffold filter). Notably, the genome assemblies and gene annotations

655    created in the present project were comparable to publicly available ones. For example, the

656    genome size of publicly available *Hanseniaspora vineae* T02 19AF is 11.38 Mb with 4,661

657    genes, while our assembly of *Hanseniaspora vineae* NRRL Y-1626 was 11.15 Mb with 5,193

658    genes.

659

660    We found that our assemblies were of comparable quality to those from publicly available

661    genomes. For example, *Hanseniaspora uvarum* NRRL Y-1614 (N50 = 267.64 kb; genome size =

662    8.82 Mb; number of scaffolds = 258; gene number = 4,227), which was sequenced in the present

663    study, and *H. uvarum* AWRI3580 (N50 = 1,289.09 kb; genome size = 8.81 Mb; number of

664    scaffolds = 18; gene number = 4,061), which is publicly available [106] had similar single-copy

665    BUSCO genes present in the highest and lowest ORTHODB [43] taxonomic ranks (Eukaryota and

666    Saccharomycetales, respectively). Specifically, *H. uvarum* NRRL Y-1614 and *H. uvarum*

667    AWRI3580 had 80.20% (243 / 303) and 79.87% (242 / 303) of universally single-copy

668    orthologs in Eukaryota present in each genome respectively, and 52.31% (895 / 1,711) and

669    51.49% (881 / 1,711) of universally single-copy orthologs in Saccharomycetales present in each

670    genome, respectively.

671

672    To identify single-copy orthologous genes (OGs) among all protein coding sequences for all 29

673    taxa, we used ORTHOMCL, version 1.4 [111]. ORTHOMCL clusters genes into OGs using a

674    Markov clustering algorithm (van Dongen, 2000; https://micans.org/mcl/) from gene similarity

675    information acquired from a blastp 'all-vs-all' using NCBI's BLAST+, version 2.3.0 (Fig S2;

676    Madden, 2013) and the proteomes of species of interest as input. The key parameters used in

677    blastp 'all-vs-all' were: e-value = 1e$^{-10}$, percent identity cut-off = 30%, percent match cutoff =

678    70%, and a maximum weight value = 180. To conservatively identify OGs, we used a strict

679    ORTHOMCL inflation parameter of 4.

680

681    To identify additional OGs suitable for use in phylogenomic and molecular sequence analyses,

682    we identified the single best putatively orthologous gene from OGs with full species

683    representation and a maximum of two species with multiple copies using PHYLOTREEPRUNER,

684    version 1.0 [114]. To do so, we first aligned and trimmed sequences in 1,143 OGs out a total of

685    11,877 that fit the criterion of full representation and a maximum of two species with duplicate

686    sequences. More specifically, we used MAFFT, version 7.294b [115], with the BLOSUM62 matrix

687    of substitutions [116], a gap penalty of 1.0, 1,000 maximum iterations, the 'genafpair' parameter,

688    and TRIMAL, version 1.4 [117], with the 'automated1' parameter to align and trim individual

689    sequences, respectively. The resulting OG multiple sequence alignments were then used to infer

690    gene phylogenies using FASTTREE, version 2.1.9 [118], with 4 and 2 rounds of subtree-prune-

691    regraft and optimization of all 5 branches at nearest-neighbor interchanges, respectively, as well

692    as the 'slownni' parameter to refine the inferred topology. Internal branches with support lower

693    than 0.9 Shimodaira-Hasegawa-like support implemented in FASTTREE [118] were collapsed

694    using PHYLOTREEPRUNER, version 1.0 [114], and the longest sequence for species with multiple

31

695     sequences per OG were retained, resulting a robust set of OGs with every taxon being

696     represented by a single sequence. OGs were realigned (MAFFT) and trimmed (TRIMAL) using the

697     same parameters as above.

698

699     **Phylogenomic analyses**     To infer the *Hanseniaspora* phylogeny, we performed

700     phylogenetic inference using maximum likelihood [119] with concatenation [120,121] and

701     coalescence [122] approaches. To determine the best-fit phylogenetic model for concatenation

702     and generate single-gene trees for coalescence, we constructed trees per single-copy OG using

703     RAxML, version 8.2.8. [123], where each topology was determined using 5 starting trees.

704     Single-gene trees that did not recover all outgroup species as the earliest diverging taxa when

705     serially rooted on outgroup taxa were discarded. Individual OG alignments or trees were used for

706     species tree estimation with RAxML (i.e., concatenation) using the LG [124] model of

707     substitution, which is the most commonly supported model of substitution (874 / 1,034;  84.53%

708     genes), or ASTRAL-II, version 4.10.12 (i.e., coalescence; Mirarab and Warnow, 2015). Branch

709     support for the concatenation and coalescence phylogenies was determined using 100 rapid

710     bootstrap replicates [126] and local posterior support [122], respectively.

711

712     Several previous phylogenomic studies have shown that the internal branches preceding the

713     *Hanseniaspora* FEL and SEL are long [33,35]. To examine whether the relationship between the

714     length of the internal branch preceding the FEL and the length of the internal branch preceding

715     the SEL was consistent across genes in our phylogeny, we used NEWICK UTILITIES, version 1.6

716     [127] to remove the 88 single-gene trees where either lineage was not recovered as monophyletic

717     and calculated their difference for the remaining 946 genes.

718

719     **Estimating divergence times**      To estimate divergence times among the 25 *Hanseniaspora*

720      genomes, we used the Bayesian method MCMCTree in the PAML, version 4.9 [128], and the

721      concatenated 1,034-gene matrix. The input tree was derived from the concatenation-based ML

722      analysis under a single LG+G4 [124] model (Figure 1A). The in-group root (i.e., the split

723      between the FEL and SEL) age was set between 0.756 and 1.177 time units (1 time unit = 100

724      million years ago [mya]), which was adopted from a recent study [34].

725

726      To infer the *Hanseniaspora* timetree, we first estimated branch lengths under a single LG+G4

727      [124] model with codeml in the PAML, version 4.9 [128], package and obtained a rough mean of

728      the overall mutation rate. Next, we applied the approximate likelihood method [129,130] to

729      estimate the gradient vector and Hessian matrix with Taylor expansion (option usedata = 3).

730      Last, we assigned (a) the gamma-Dirichlet prior for the overall substitution rate (option

731      rgene_gamma) as G(1, 1.55), with a mean of 0.64, (b) the gamma-Dirichlet prior for the rate-

732      drift parameter (option sigma2 gamma) as G(1, 10), and (c) the parameters for the birth-death

733      sampling process with birth and death rates $\lambda=\mu=1$ and sampling fraction $\rho=0$. We employed the

734      independent-rate model (option clock=2) to account for the rate variation across different

735      lineages and used soft bounds (left and right tail probabilities equal 0.025) to set minimum and

736      maximum values for the in-group root mentioned above. The MCMC run was first run for

737      1,000,000 iterations as burn-in and then sampled every 1,000 iterations until a total of 30,000

738      samples was collected. Two separate MCMC runs were compared for convergence, and similar

739      results were observed.

740

741     **Gene presence and absence analysis**        To determine the presence and absence of genes in

742     *Hanseniaspora* genomes, we built hidden Markov models (HMMs) for each gene present in

743     *Saccharomyces cerevisiae* and used the resulting HMM profile to search for the corresponding

744     homolog in each *Hanseniaspora* genome, as well as outgroup taxa. More specifically, for each of

745     the 5,917 verified open reading frames from *S. cerevisiae* [131] (downloaded Oct 2018 from the

746     *Saccharomyces* genome database), we searched for putative homologs in NCBI's Reference

747     Sequence Database for Fungi (downloaded June 2018) using NCBI's BLAST+, version 2.3.0

748     [113], blastp function, and an e-value cut-off of $1e^{-3}$ as recommended for homology searches

749     [132]. We used the top 100 hits for the gene of interest and aligned them using MAFFT, version

750     7.294b [115], with the same parameters described above. The resulting gene alignment was then

751     used to create an HMM profile for the gene using the hmmbuild function in HMMER, version

752     3.1b2 [133]. The resulting HMM profile was then used to search for each individual gene in each

753     *Hanseniaspora* genome and outgroup taxa using the hmmsearch function with an expectation

754     value cutoff of 0.01 and a score cutoff of 50. This analysis was done for the 5,735 genes with

755     multiple blast hits allowing for the creation of a HMM profile. To evaluate the validity of

756     constructed HMMs, we examined their ability to recall genes in *S. cerevisiae* and found that we

757     recovered all nuclear genes. Altogether, our ability to recall 99.63% of genes demonstrates the

758     validity of our pipeline for the vast majority of genes and for nuclear genes in particular.

759

760     To determine if any functional categories were over- or under-represented among genes present

761     or absent among *Hanseniaspora* species, we conducted gene ontology (GO) [134] enrichment

762     analyses using GOATOOLS, version 0.7.9 [135]. We used a background of all *S. cerevisiae*

763     genes and a *p*-value cut-off of 0.05 after multiple-test correction using the Holm method [136].

764    Plotting gene presence and absence among pathways was done by examining depicted pathways

765    available through the KEGG project [137] and the *Saccharomyces* Genome Database [131].

766

767    We examined the validity of the gene presence and absence pipeline by examining under-

768    represented terms and the presence or absence of essential genes in *S. cerevisiae* [138]. We

769    hypothesized that under-represented GO terms will be associated with basic molecular processes

770    and that essential genes will be under-represented among the set of absent genes. In agreement

771    with these expectations, GO terms associated with basic biological processes and essential *S.*

772    *cerevisiae* genes are under-represented among genes that are absent across *Hanseniaspora*

773    genomes. For example, among all genes absent in the FEL and SEL, the molecular functions

774    BASE PAIRING, GO:0000496 ($p < 0.001$); GTP BINDING, GO:0005525 ($p < 0.001$); and

775    ATPASE ACTIVITY, COUPLED TO MOVEMENT OF SUBSTANCES, GO:0043492 ($p <$

776    0.001), are significantly under-represented (File S4). Similarly, *S. cerevisiae* essential genes are

777    significantly under-represented ($p < 0.001$; Fischer's exact test for both lineages) among lost

778    genes with only 3 and 2 *S. cerevisiae* essential genes having been lost from the FEL and SEL

779    genomes, respectively.

780

781    **Ploidy estimation**    To determine ploidy, we leveraged base frequency distributions at variable

782    sites, which we generated by mapping each genome's reads to its assembly. To ensure high-

783    quality read mapping, we first quality-trimmed reads suing TRIMMOMATIC, version 0.36 [91],

784    using the parameters leading:10, trailing:10, slidingwindow:4:20, and minlen:50. Reads were

785    subsequently mapped to their respective genome using BOWTIE2, version 1.1.2 [139], with the

786    "sensitive" parameter and converted the resulting file to a sorted bam format using SAMTOOLS,

35

787    version 1.3.1 [140]. We next used NQUIRE [141], which extracts base frequency information at

788    segregating sites with a minimum frequency of 0.2. Prior to visualization, we removed

789    background noise by utilizing the Gaussian Mixture Model with Uniform noise component

790    [141].

791

792    **Molecular evolution and mutation analysis**      *Molecular sequence rate analysis along the*

793    *phylogeny.*      To determine the rate of sequence evolution over the course of

794    *Hanseniaspora* evolution, we examined variation in the rate of nonsynonymous (dN) to the rate

795    of synonymous (dS) substitutions (dN/dS or $\omega$) across the species phylogeny. We first obtained

796    codon-based alignments of the protein sequences used during phylogenomic inference by

797    threading nucleotides on top of the amino acid sequence using PAL2NAL, version 14 [142], and

798    calculated $\omega$ values under the different hypotheses using the CODEML module in PAML, version

799    4.9 [128]. For each gene tested, we set the null hypothesis ($H_o$) where all internal branches

800    exhibit the same $\omega$ (model = 0) and compared it to four different alternative hypotheses. Under

801    the $H_{FE\text{-}SE\ branch}$ hypothesis, the branches immediately preceding the FEL and SEL were assumed

802    to exhibit distinct $\omega$ values from the background (model = 2) (Fig 5Bi). Under the $H_{FE}$

803    hypothesis, the branch immediately preceding the FEL was assumed to have a distinct $\omega$ value,

804    all FEL crown branches were assumed to have their own collective $\omega$ value, and all background

805    branches were assumed to have their own collective $\omega$ value (model = 2) (Fig 5Ci). The $H_{SE}$

806    hypothesis assumed the branch preceding the lineage had its own $\omega$ value, all SEL crown

807    branches had their own collective $\omega$ value, and all background branches were assumed to have

808    their own collective $\omega$ value (model = 2) (Fig 5Di). Lastly, the $H_{FE\text{-}SE\ crown}$ hypothesis assumed

809    that all FEL crown branches had their own collective $\omega$ value, all SEL crown branches had their

36

810    own collective ω value, and the rest of the branches were assumed to have their own collective ω

811    value (model = 2) (Fig 5Ei). To determine if each of the alternative hypotheses was significantly

812    different from the null hypothesis, we used the likelihood ratio test (LRT) ($\alpha = 0.01$). A few

813    genes could not be analyzed due to fatal interruptions or errors during use in PAML, version 4.9

814    [128], which have been reported by other users [143]; these genes were removed from the

815    analysis. Thus, this analysis was conducted for 989 genes for three tests ($H_{FE-SE\ branch}$, $H_{FE}$, and

816    $H_{SE}$ hypotheses) and 983 genes for one test ($H_{FE-SE\ crown}$ hypothesis).

817

818    *Examination of mutational signatures*          To conservatively identify base substitutions,

819    insertions, and deletions found in taxa in the FEL or SEL, we examined the status of each

820    nucleotide at each position in codon-based and amino acid-based OG alignments. We examined

821    base substitutions, insertions, and deletions at sites that are conserved in the outgroup (i.e., all

822    outgroup taxa have the same character state for a given position in an alignment). For base

823    substitutions, we determined if the nucleotide or amino acid residue in a given *Hanseniaspora*

824    species differed from the conserved outgroup nucleotide or amino acid residue at the same

825    position. To measure if amino acid substitutions in each lineage were conservative or radical

826    (i.e., a substitution to a similar amino acid residue versus a substitution to an amino acid residue

827    with different properties), we used Sneath's index of dissimilarity, which considers 134

828    categories of biological activity and chemical change to quantify dissimilarity of amino acid

829    substitutions, and Epstein's coefficient of difference, which considers differences in polarity and

830    size of amino acids to quantify dissimilarity. Notably, Sneath's index is symmetric (i.e.,

831    isoleucine to leucine is equivalent to leucine to isoleucine), whereas Epstein's coefficient is not

832    (i.e., isoleucine to leucine is not equivalent to leucine to isoleucine). For indels, we used a sliding

37

833     window approach with a step size of one nucleotide. We considered positions where a nucleotide

834     was present in all outgroup taxa but a gap was present in *Hanseniaspora* as deletions, and

835     positions where a gap was present in all outgroup taxa and a nucleotide was present in

836     *Hanseniaspora* species as insertions. Analyses were conducted using custom PYTHON, version

837     3.5.2 (https://www.python.org/), scripts, which use the BIOPYTHON, version 1.70 [144], and

838     NUMPY, version 1.13.1 [145], modules.

839

840     We discovered that all *Hanseniaspora* species lack the *PHR1* gene, which is associated with the

841     repair of UV radiation damage. UV exposure induces high levels of CC → TT dinucleotide

842     substitutions [87]. If *Hanseniaspora* have a reduced capacity to repair UV radiation damage,

843     they would be expected to contain fewer CC|GG dinucleotides and more TT|AA ones. To test

844     whether this was the case, we created a CC or GG (hereby denoted as CC|GG) score, which was

845     calculated using the following formula:

846     $CC|GG\ score = \frac{CC|GG}{D} \times \frac{1}{G|C}$     where $D = \frac{GS}{2}$

847     where CC|GG is the number of observed CC or GG dinucleotides in a genome, D is the number

848     of dinucleotides in the genome, GS is the genome size, and G|C is GC-content. Similarly, we

849     created a TT|AA score calculated the following formula:

850     $TT|AA\ score = \frac{TT|AA}{D} \times \frac{1}{A|T}$     where $D = \frac{GS}{2}$

851     where TT|AA is the number of TT or AA dinucleotides in a genome, D is the number of

852     dinucleotides in the genome, GS is the genome size, and A|T is AT-content.

853

854     **Data Availability**

38

855    Data matrices, species-level and single-gene phylogenies, dN/dS results, and HMMs will be

856    made available through the figshare repository upon publication.

857

858

**Acknowledgements**

859

860    We thank members of the Rokas and Hittinger laboratories for helpful suggestions and

861    discussion.

862

**Financial Statement**

863

40

888    **References**

889    1.    Lindahl T. Quality Control by DNA Repair. Science (80- ). 1999;286: 1897–1905.

890          doi:10.1126/science.286.5446.1897

891    2.    Hakem R. DNA-damage repair; the good, the bad, and the ugly. EMBO J. 2008;27: 589–

892          605. doi:10.1038/emboj.2008.15

893    3.    Broustas CG, Lieberman HB. DNA Damage Response Genes and the Development of

894          Cancer Metastasis. Radiat Res. 2014;181: 111–130. doi:10.1667/RR13515.1

895    4.    Pal C, Maciá MD, Oliver A, Schachar I, Buckling A. Coevolution with viruses drives the

896          evolution of bacterial mutation rates. Nature. 2007;450: 1079–1081.

897          doi:10.1038/nature06350

898    5.    Myung K, Datta A, Kolodner RD. Suppression of Spontaneous Chromosomal

899          Rearrangements by S Phase Checkpoint Functions in Saccharomyces cerevisiae. Cell.

900          2001;104: 397–408. doi:10.1016/S0092-8674(01)00227-6

901    6.    Billmyre RB, Clancey SA, Heitman J. Natural mismatch repair mutations mediate

902          phenotypic diversity and drug resistance in Cryptococcus deuterogattii. Elife. 2017;6.

903          doi:10.7554/eLife.28802

904    7.    Boyce KJ, Wang Y, Verma S, Shakya VPS, Xue C, Idnurm A. Mismatch Repair of DNA

905          Replication Errors Contributes to Microevolution in the Pathogenic Fungus Cryptococcus

906          neoformans. Alspaugh JA, editor. MBio. 2017;8. doi:10.1128/mBio.00595-17

907    8.    Rhodes J, Beale MA, Vanhove M, Jarvis JN, Kannambath S, Simpson JA, et al. A

908          Population Genomics Approach to Assessing the Genetic Basis of Within-Host

909          Microevolution Underlying Recurrent Cryptococcal Meningitis Infection. G3

910          Genes|Genomes|Genetics. 2017;7: 1165–1176. doi:10.1534/g3.116.037499

911    9.      Barnum KJ, O'Connell MJ. Cell Cycle Regulation by Checkpoints. 2014. pp. 29–40.

912           doi:10.1007/978-1-4939-0888-2_2

913    10.    Cross FR, Buchler NE, Skotheim JM. Evolution of networks and sequences in eukaryotic

914           cell cycle control. Philos Trans R Soc Lond B Biol Sci. 2011;366: 3532–44.

915           doi:10.1098/rstb.2011.0078

916    11.    Medina EM, Turner JJ, Gordân R, Skotheim JM, Buchler NE. Punctuated evolution and

917           transitional hybrid network in an ancestral cell cycle of fungi. Elife. 2016;5.

918           doi:10.7554/eLife.09492

919    12.    Costanzo M, Nishikawa JL, Tang X, Millman JS, Schub O, Breitkreuz K, et al. CDK

920           Activity Antagonizes Whi5, an Inhibitor of G1/S Transcription in Yeast. Cell. 2004;117:

921           899–913. doi:10.1016/j.cell.2004.05.024

922    13.    Castro A, Bernis C, Vigneron S, Labbé J-C, Lorca T. The anaphase-promoting complex: a

923           key factor in the regulation of cell cycle. Oncogene. 2005;24: 314–325.

924           doi:10.1038/sj.onc.1207973

925    14.    Heinrich S, Sewart K, Windecker H, Langegger M, Schmidt N, Hustedt N, et al. Mad1

926           contribution to spindle assembly checkpoint signalling goes beyond presenting Mad2 at

927           kinetochores. EMBO Rep. 2014;15: 291–298. doi:10.1002/embr.201338114

928    15.    Hendler A, Medina EM, Kishkevich A, Abu-Qarn M, Klier S, Buchler NE, et al. Gene

929           duplication and co-evolution of G1/S transcription factor specificity in fungi are essential

930           for optimizing cell fitness. Snyder M, editor. PLOS Genet. 2017;13: e1006778.

931           doi:10.1371/journal.pgen.1006778

932    16.    Zhou B-BS, Elledge SJ. The DNA damage response: putting checkpoints in perspective.

933           Nature. 2000;408: 433–439. doi:10.1038/35044005

934   17.   Weinert TA, Kiser GL, Hartwell LH. Mitotic checkpoint genes in budding yeast and the

935         dependence of mitosis on DNA replication and repair. Genes Dev. 1994;8: 652–665.

936         doi:10.1101/gad.8.6.652

937   18.   Serero A, Jubin C, Loeillet S, Legoix-Né P, Nicolas AG. Mutational landscape of yeast

938         mutator strains. Proc Natl Acad Sci. 2014;111: 1897–1902. doi:10.1073/pnas.1314423111

939   19.   Campbell BB, Light N, Fabrizio D, Zatzman M, Fuligni F, de Borja R, et al.

940         Comprehensive Analysis of Hypermutation in Human Cancer. Cell. 2017;171: 1042–

941         1056.e10. doi:10.1016/j.cell.2017.09.048

942   20.   Huang M-E, de Calignon A, Nicolas A, Galibert F. POL32 , a subunit of the

943         Saccharomyces cerevisiae DNA polymerase δ, defines a link between DNA replication

944         and the mutagenic bypass repair pathway. Curr Genet. 2000;38: 178–187.

945         doi:10.1007/s002940000149

946   21.   Xiao W, Chow BL, Hanna M, Doetsch PW. Deletion of the MAG1 DNA glycosylase

947         gene suppresses alkylation-induced killing and mutagenesis in yeast cells lacking AP

948         endonucleases. Mutat Res - DNA Repair. 2001; doi:10.1016/S0921-8777(01)00113-6

949   22.   Sebastian J, Sancar GB. A damage-responsive DNA binding protein regulates

950         transcription of the yeast DNA repair gene PHR1. Proc Natl Acad Sci. 1991;88: 11251–

951         11255. doi:10.1073/pnas.88.24.11251

952   23.   Sebastian J, Kraus B, Sancar GB. Expression of the yeast PHR1 gene is induced by DNA-

953         damaging agents. Mol Cell Biol. 1990;10: 4630–7.

954   24.   Nunoshiba T. A novel Nudix hydrolase for oxidized purine nucleoside triphosphates

955         encoded by ORFYLR151c (PCD1 gene) in Saccharomyces cerevisiae. Nucleic Acids Res.

956         2004;32: 5339–5348. doi:10.1093/nar/gkh868

957    25.    Cartwright JL, Gasmi L, Spiller DG, McLennan AG. The Saccharomyces cerevisiae

958        PCD1 gene encodes a peroxisomal nudix hydrolase active toward coenzyme A and its

959        derivatives. J Biol Chem. 2000; doi:10.1074/jbc.M005015200

960    26.    De Bont R. Endogenous DNA damage in humans: a review of quantitative data.

961        Mutagenesis. 2004;19: 169–185. doi:10.1093/mutage/geh025

962    27.    Ram Y, Hadany L. The evolution of stress-induced hypermutation in asexual populations.

963        Evolution (N Y). 2012;66: 2315–2328. doi:10.1111/j.1558-5646.2012.01576.x

964    28.    Oliver A. High Frequency of Hypermutable Pseudomonas aeruginosa in Cystic Fibrosis

965        Lung Infection. Science (80- ). 2000;288: 1251–1253. doi:10.1126/science.288.5469.1251

966    29.    Giraud A, Matic I, Tenaillon O, Clara A, Radman M, Fons M, et al. Costs and Benefits of

967        High Mutation Rates: Adaptive Evolution of Bacteria in the Mouse Gut. Science (80- ).

968        2001;291: 2606–2608. doi:10.1126/science.1056421

969    30.    Taddei F, Radman M, Maynard-Smith J, Toupance B, Gouyon PH, Godelle B. Role of

970        mutator alleles in adaptive evolution. Nature. 1997;387: 700–702. doi:10.1038/42696

971    31.    Sniegowski PD, Gerrish PJ, Lenski RE. Evolution of high mutation rates in experimental

972        populations of E. coli. Nature. 1997;387: 703–705. doi:10.1038/42701

973    32.    McDonald MJ, Hsieh Y-Y, Yu Y-H, Chang S-L, Leu J-Y. The Evolution of Low

974        Mutation Rates in Experimental Mutator Populations of Saccharomyces cerevisiae. Curr

975        Biol. 2012;22: 1235–1240. doi:10.1016/j.cub.2012.04.056

976    33.    Riley R, Haridas S, Wolfe KH, Lopes MR, Hittinger CT, Göker M, et al. Comparative

977        genomics of biotechnologically important yeasts. Proc Natl Acad Sci. 2016;113: 9882–

978        9887. doi:10.1073/pnas.1603941113

979    34.    Shen X-X, Opulente DA, Kominek J, Zhou X, Steenwyk JL, Buh K V., et al. Tempo and

980  Mode of Genome Evolution in the Budding Yeast Subphylum. Cell. 2018;

981  doi:10.1016/j.cell.2018.10.023

982  35.  Shen X-X, Zhou X, Kominek J, Kurtzman CP, Hittinger CT, Rokas A. Reconstructing the

983  Backbone of the Saccharomycotina Yeast Phylogeny Using Genome-Scale Data. G3

984  Genes|Genomes|Genetics. Genetics Society of America; 2016;6: 3927–3939.

985  doi:10.1534/g3.116.034744

986  36.  Albertin W, Setati ME, Miot-Sertier C, Mostert TT, Colonna-Ceccaldi B, Coulon J, et al.

987  Hanseniaspora uvarum from Winemaking Environments Show Spatial and Temporal

988  Genetic Clustering. Front Microbiol. 2016;6. doi:10.3389/fmicb.2015.01569

989  37.  Jordão A, Vilela A, Cosme F. From Sugar of Grape to Alcohol of Wine: Sensorial Impact

990  of Alcohol in Wine. Beverages. 2015;1: 292–310. doi:10.3390/beverages1040292

991  38.  Montero CM, Dodero MCR, Sanchez DAG, Barroso CG. Analysis of low molecular

992  weight carbohydrates in food and beverages: A review. Chromatographia. 2004;

993  doi:10.1365/s10337-003-0134-3

994  39.  Martin V, Valera M, Medina K, Boido E, Carrau F. Oenological Impact of the

995  Hanseniaspora/Kloeckera Yeast Genus on Wines—A Review. Fermentation. 2018;4: 76.

996  doi:10.3390/fermentation4030076

997  40.  Langenberg A-K, Bink FJ, Wolff L, Walter S, von Wallbrunn C, Grossmann M, et al.

998  Glycolytic Functions Are Conserved in the Genome of the Wine Yeast Hanseniaspora

999  uvarum, and Pyruvate Kinase Limits Its Capacity for Alcoholic Fermentation. Dudley EG,

1000  editor. Appl Environ Microbiol. 2017;83. doi:10.1128/AEM.01580-17

1001  41.  Čadež N, Bellora N, Ulloa R, Hittinger CT, Libkind D. Genomic content of a novel yeast

1002  species Hanseniaspora gamundiae sp. nov. from fungal stromata (Cyttaria) associated with

1003    a unique fermented beverage in Andean Patagonia, Argentina. Yurkov AM, editor. PLoS

1004    One. 2019;14: e0210792. doi:10.1371/journal.pone.0210792

1005  42.  Zhou X, Shen X-X, Hittinger CT, Rokas A. Evaluating Fast Maximum Likelihood-Based

1006    Phylogenetic Programs Using Empirical Phylogenomic Data Sets. Mol Biol Evol.

1007    2018;35: 486–503. doi:10.1093/molbev/msx302

1008  43.  Waterhouse RM, Tegenfeldt F, Li J, Zdobnov EM, Kriventseva E V. OrthoDB: a

1009    hierarchical catalog of animal, fungal and bacterial orthologs. Nucleic Acids Res.

1010    2013;41: D358–D365. doi:10.1093/nar/gks1116

1011  44.  Opulente DA, Rollinson EJ, Bernick-Roehr C, Hulfachor AB, Rokas A, Kurtzman CP, et

1012    al. Factors driving metabolic diversity in the budding yeast subphylum. BMC Biol.

1013    2018;16: 26. doi:10.1186/s12915-018-0498-3

1014  45.  Kurtzman CP, Fell JW. The Yeasts - A Taxonomic Study. 4th ed. Kurtzman CP, Fell JW,

1015    editors. Elsevier Science; 1998.

1016  46.  Charron MJ, Read E, Haut SR, Michels CA. Molecular evolution of the telomere-

1017    associated MAL loci of Saccharomyces. Genetics. 1989;122: 307–16.

1018  47.  Koschwanez JH, Foster KR, Murray AW. Sucrose Utilization in Budding Yeast as a

1019    Model for the Origin of Undifferentiated Multicellularity. Keller L, editor. PLoS Biol.

1020    2011;9: e1001122. doi:10.1371/journal.pbio.1001122

1021  48.  Jorgensen P. Systematic Identification of Pathways That Couple Cell Growth and Division

1022    in Yeast. Science (80- ). 2002;297: 395–400. doi:10.1126/science.1070850

1023  49.  Colman-Lerner A, Chin TE, Brent R. Yeast Cbk1 and Mob2 Activate Daughter-Specific

1024    Genetic Programs to Induce Asymmetric Cell Fates. Cell. 2001;107: 739–750.

1025    doi:10.1016/S0092-8674(01)00596-7

1026    50.    Lubelsky Y, Reuven N, Shaul Y. Autorepression of Rfx1 Gene Expression: Functional

1027           Conservation from Yeast to Humans in Response to DNA Replication Arrest. Mol Cell

1028           Biol. 2005;25: 10665–10673. doi:10.1128/MCB.25.23.10665-10673.2005

1029    51.    Galgoczy DJ, Toczyski DP. Checkpoint Adaptation Precedes Spontaneous and Damage-

1030           Induced Genomic Instability in Yeast. Mol Cell Biol. 2001;21: 1710–1718.

1031           doi:10.1128/MCB.21.5.1710-1718.2001

1032    52.    Kim IY, Kwon HY, Park KH, Kim DS. Anaphase-Promoting Complex 7 is a Prognostic

1033           Factor in Human Colorectal Cancer. Ann Coloproctol. 2017;33: 139–145.

1034           doi:10.3393/ac.2017.33.4.139

1035    53.    Xiao W, Chow BL. Synergism between yeast nucleotide and base excision repair

1036           pathways in the protection against DNA methylation damage. Curr Genet. 1998;

1037           doi:10.1007/s002940050313

1038    54.    Nitiss KC, Malik M, He X, White SW, Nitiss JL. Tyrosyl-DNA phosphodiesterase (Tdp1)

1039           participates in the repair of Top2-mediated DNA damage. Proc Natl Acad Sci. 2006;103:

1040           8953–8958. doi:10.1073/pnas.0603455103

1041    55.    Lustig AJ. Cdc13 subcomplexes regulate multiple telomere functions. Nat Struct Biol.

1042           2001;8: 297–9. doi:10.1038/86157

1043    56.    Zhu YO, Siegal ML, Hall DW, Petrov DA. Precise estimates of mutation rate and

1044           spectrum in yeast. Proc Natl Acad Sci. 2014;111: E2310–E2318.

1045           doi:10.1073/pnas.1323011111

1046    57.    Marth GT, Yu F, Indap AR, Garimella K, Gravel S, Leong W, et al. The functional

1047           spectrum of low-frequency coding variation. Genome Biol. 2011;12: R84.

1048           doi:10.1186/gb-2011-12-9-r84

1049    58.    Wang J, Raskin L, Samuels DC, Shyr Y, Guo Y. Genome measures used for quality

1050            control are dependent on gene function and ancestry. Bioinformatics. 2015;31: 318–23.

1051            doi:10.1093/bioinformatics/btu668

1052    59.    Vignal A, Milan D, SanCristobal M, Eggen A. A review on SNP and other types of

1053            molecular markers and their use in animal genetics. Genet Sel Evol. 2002;34: 275–305.

1054            doi:10.1051/gse:2002009

1055    60.    Sneath PH. Relations between chemical structure and biological activity in peptides. J

1056            Theor Biol. 1966;12: 157–95.

1057    61.    Epstein CJ. Non-randomness of Ammo-acid Changes in the Evolution of Homologous

1058            Proteins. Nature. 1967;215: 355–359. doi:10.1038/215355a0

1059    62.    Albalat R, Cañestro C. Evolution by gene loss. Nat Rev Genet. 2016;17: 379–391.

1060            doi:10.1038/nrg.2016.39

1061    63.    Baba T, Ara T, Hasegawa M, Takai Y, Okumura Y, Baba M, et al. Construction of

1062            Escherichia coli K-12 in-frame, single-gene knockout mutants: the Keio collection. Mol

1063            Syst Biol. 2006;2. doi:10.1038/msb4100050

1064    64.    Giaever G, Chu AM, Ni L, Connelly C, Riles L, Véronneau S, et al. Functional profiling

1065            of the Saccharomyces cerevisiae genome. Nature. 2002;418: 387–391.

1066            doi:10.1038/nature00935

1067    65.    Segal ES, Gritsenko V, Levitan A, Yadav B, Dror N, Steenwyk JL, et al. Gene

1068            Essentiality Analyzed by In Vivo Transposon Mutagenesis and Machine Learning in a

1069            Stable Haploid Isolate of Candida albicans. Di Pietro A, editor. MBio. 2018;9.

1070            doi:10.1128/mBio.02048-18

1071    66.    Koskiniemi S, Sun S, Berg OG, Andersson DI. Selection-Driven Gene Loss in Bacteria.

1072   Casadesús J, editor. PLoS Genet. 2012;8: e1002787. doi:10.1371/journal.pgen.1002787

1073 67. Keeling PJ, Slamovits CH. Simplicity and complexity of microsporidian genomes.

1074   Eukaryot Cell. 2004;3: 1363–9. doi:10.1128/EC.3.6.1363-1369.2004

1075 68. Katinka MD, Duprat S, Cornillot E, Méténier G, Thomarat F, Prensier G, et al. Genome

1076   sequence and gene compaction of the eukaryote parasite Encephalitozoon cuniculi.

1077   Nature. 2001;414: 450–453. doi:10.1038/35106579

1078 69. Chang ES, Neuhof M, Rubinstein ND, Diamant A, Philippe H, Huchon D, et al. Genomic

1079   insights into the evolutionary origin of Myxozoa within Cnidaria. Proc Natl Acad Sci U S

1080   A. 2015;112: 14912–7. doi:10.1073/pnas.1511468112

1081 70. Hittinger CT, Rokas A, Carroll SB. Parallel inactivation of multiple GAL pathway genes

1082   and ecological diversification in yeasts. Proc Natl Acad Sci. 2004;101: 14144–14149.

1083   doi:10.1073/pnas.0404319101

1084 71. Slot JC, Rokas A. Multiple GAL pathway gene clusters evolved independently and by

1085   different mechanisms in fungi. Proc Natl Acad Sci. 2010;107: 10136–10141.

1086   doi:10.1073/pnas.0914418107

1087 72. Wolfe KH, Armisén D, Proux-Wera E, ÓhÉigeartaigh SS, Azam H, Gordon JL, et al.

1088   Clade- and species-specific features of genome evolution in the Saccharomycetaceae.

1089   Nielsen J, editor. FEMS Yeast Res. 2015;15: fov035. doi:10.1093/femsyr/fov035

1090 73. Hartwell L. Defects in a cell cycle checkpoint may be responsible for the genomic

1091   instability of cancer cells. Cell. 1992. doi:10.1016/0092-8674(92)90586-2

1092 74. Tavares MJ, Güldener U, Esteves M, Mendes-Faia A, Mendes-Ferreira A, Mira NP.

1093   Genome Sequence of the Wine Yeast Saccharomycodes ludwigii UTAD17. Cuomo CA,

1094   editor. Microbiol Resour Announc. 2018;7. doi:10.1128/MRA.01195-18

1095    75.    Sterling CH. DNA Polymerase 4 of Saccharomyces cerevisiae Is Important for Accurate

1096            Repair of Methyl-Methanesulfonate-Induced DNA Damage. Genetics. 2005;172: 89–98.

1097            doi:10.1534/genetics.105.049254

1098    76.    Jenni S, Harrison SC. Structure of the DASH/Dam1 complex shows its role at the yeast

1099            kinetochore-microtubule interface. Science (80- ). 2018;360: 552–558.

1100            doi:10.1126/science.aar6436

1101    77.    Dimitrova YN, Jenni S, Valverde R, Khin Y, Harrison SC. Structure of the MIND

1102            Complex Defines a Regulatory Focus for Yeast Kinetochore Assembly. Cell. 2016;167:

1103            1014–1027.e12. doi:10.1016/j.cell.2016.10.011

1104    78.    Seixas I, Barbosa C, Salazar SB, Mendes-Faia A, Wang Y, Güldener U, et al. Genome

1105            Sequence of the Nonconventional Wine Yeast Hanseniaspora guilliermondii UTAD222.

1106            Genome Announc. American Society for Microbiology; 2017;5: e01515-16.

1107            doi:10.1128/genomeA.01515-16

1108    79.    Chavan P, Mane S, Kulkarni G, Shaikh S, Ghormade V, Nerkar DP, et al. Natural yeast

1109            flora of different varieties of grapes used for wine making in India. Food Microbiol.

1110            2009;26: 801–808. doi:10.1016/j.fm.2009.05.005

1111    80.    Wikstrom N, Savolainen V, Chase MW. Evolution of the angiosperms: calibrating the

1112            family tree. Proc R Soc B Biol Sci. 2001;268: 2211–2220. doi:10.1098/rspb.2001.1782

1113    81.    Lynch M, Sung W, Morris K, Coffey N, Landry CR, Dopman EB, et al. A genome-wide

1114            view of the spectrum of spontaneous mutations in yeast. Proc Natl Acad Sci. 2008;

1115            doi:10.1073/pnas.0803466105

1116    82.    Lynch M. The Origins of Genome Architecture. Journal of Heredity. 2007.

1117            doi:10.1093/jhered/esm073

1118   83.   Keightley PD, Trivedi U, Thomson M, Oliver F, Kumar S, Blaxter ML. Analysis of the

1119         genome sequences of three Drosophila melanogaster spontaneous mutation accumulation

1120         lines. Genome Res. 2009;19: 1195–1201. doi:10.1101/gr.091231.109

1121   84.   Lynch M. Rate, molecular spectrum, and consequences of human mutation. Proc Natl

1122         Acad Sci. 2010;107: 961–968. doi:10.1073/pnas.0912629107

1123   85.   Hershberg R, Petrov DA. Evidence That Mutation Is Universally Biased towards AT in

1124         Bacteria. Nachman MW, editor. PLoS Genet. 2010;6: e1001115.

1125         doi:10.1371/journal.pgen.1001115

1126   86.   Seeberg E, Eide L, Bjørås M. The base excision repair pathway. Trends Biochem Sci.

1127         1995;20: 391–397. doi:10.1016/S0968-0004(00)89086-6

1128   87.   Huang MN, Yu W, Teoh WW, Ardin M, Jusakul A, Ng AWT, et al. Genome-scale

1129         mutational signatures of aflatoxin in cells, mice, and human tumors. Genome Res.

1130         2017;27: 1475–1486. doi:10.1101/gr.220038.116

1131   88.   Hittinger CT, Gonçalves P, Sampaio JP, Dover J, Johnston M, Rokas A. Remarkably

1132         ancient balanced polymorphisms in a multi-locus gene network. Nature. 2010;464: 54–58.

1133         doi:10.1038/nature08791

1134   89.   Zhou X, Peris D, Kominek J, Kurtzman CP, Hittinger CT, Rokas A. in silico Whole

1135         Genome Sequencer &amp; Analyzer (iWGS): A Computational Pipeline to Guide the

1136         Design and Analysis of de novo Genome Sequencing Studies. G3

1137         Genes|Genomes|Genetics. 2016; doi:10.1534/g3.116.034249

1138   90.   Yandell M, Ence D. A beginner's guide to eukaryotic genome annotation. Nat Rev Genet.

1139         2012;13: 329–42. doi:10.1038/nrg3174

1140   91.   Bolger AM, Lohse M, Usadel B. Trimmomatic: A flexible trimmer for Illumina sequence

1141    data. Bioinformatics. 2014;30: 2114–2120. doi:10.1093/bioinformatics/btu170

1142    92.    Song L, Florea L, Langmead B. Lighter: fast and memory-efficient sequencing error

1143    correction without counting. Genome Biol. 2014;15: 509. doi:10.1186/s13059-014-0509-9

1144    93.    Chikhi R, Medvedev P. Informed and automated k-mer size selection for genome

1145    assembly. Bioinformatics. 2014;30: 31–37. doi:10.1093/bioinformatics/btt310

1146    94.    Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM, Birol I. ABySS: A parallel

1147    assembler for short read sequence data. Genome Res. 2009;19: 1117–1123.

1148    doi:10.1101/gr.089532.108

1149    95.    Weisenfeld NI, Yin S, Sharpe T, Lau B, Hegarty R, Holmes L, et al. Comprehensive

1150    variation discovery in single human genomes. Nat Genet. 2014;46: 1350–1355.

1151    doi:10.1038/ng.3121

1152    96.    Zimin A V., Marçais G, Puiu D, Roberts M, Salzberg SL, Yorke JA. The MaSuRCA

1153    genome assembler. Bioinformatics. 2013;29: 2669–2677.

1154    doi:10.1093/bioinformatics/btt476

1155    97.    Simpson JT, Durbin R. Efficient de novo assembly of large genomes using compressed

1156    data structures. Genome Res. 2012;22: 549–556. doi:10.1101/gr.126953.111

1157    98.    Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically

1158    improved memory-efficient short-read de novo assembler. Gigascience. 2012;1: 18.

1159    doi:10.1186/2047-217X-1-18

1160    99.    Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes:

1161    A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing. J

1162    Comput Biol. 2012;19: 455–477. doi:10.1089/cmb.2012.0021

1163    100.    Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for

1164        genome assemblies. Bioinformatics. 2013;29: 1072–1075.

1165        doi:10.1093/bioinformatics/btt086

1166   101.  Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database management

1167        tool for second-generation genome projects. BMC Bioinformatics. 2011;12: 491.

1168        doi:10.1186/1471-2105-12-491

1169   102.  Grigoriev I V., Nikitin R, Haridas S, Kuo A, Ohm R, Otillar R, et al. MycoCosm portal:

1170        gearing up for 1000 fungal genomes. Nucleic Acids Res. 2014;42: D699–D704.

1171        doi:10.1093/nar/gkt1183

1172   103.  Ter-Hovhannisyan V, Lomsadze A, Chernoff YO, Borodovsky M. Gene prediction in

1173        novel fungal genomes using an ab initio algorithm with unsupervised training. Genome

1174        Res. 2008;18: 1979–1990. doi:10.1101/gr.081612.108

1175   104.  Korf I. Gene finding in novel genomes. BMC Bioinformatics. 2004;5: 59.

1176        doi:10.1186/1471-2105-5-59

1177   105.  Stanke M, Waack S. Gene prediction with a hidden Markov model and a new intron

1178        submodel. Bioinformatics. 2003;19: ii215-ii225. doi:10.1093/bioinformatics/btg1080

1179   106.  Sternes PR, Lee D, Kutyna DR, Borneman AR. Genome Sequences of Three Species of

1180        Hanseniaspora Isolated from Spontaneous Wine Fermentations. Genome Announc.

1181        2016;4: e01287-16. doi:10.1128/genomeA.01287-16

1182   107.  Giorello FM, Berná L, Greif G, Camesasca L, Salzman V, Medina K, et al. Genome

1183        Sequence of the Native Apiculate Wine Yeast Hanseniaspora vineae T02/19AF. Genome

1184        Announc. 2014;2: e00530-14. doi:10.1128/genomeA.00530-14

1185   108.  Cadez N. Phylogenetic placement of Hanseniaspora-Kloeckera species using multigene

1186        sequence analysis with taxonomic implications: descriptions of Hanseniaspora

1187    pseudoguilliermondii sp. nov. and Hanseniaspora occidentalis var. citrica var. nov. Int J

1188    Syst Evol Microbiol. 2006;56: 1157–1165. doi:10.1099/ijs.0.64052-0

1189    109.   Chang C-F, Huang L-Y, Chen S-F, Lee C-F. Kloeckera taiwanica sp. nov., an

1190    ascomycetous apiculate yeast species isolated from mushroom fruiting bodies. Int J Syst

1191    Evol Microbiol. 2012;62: 1434–1437. doi:10.1099/ijs.0.034231-0

1192    110.   Waterhouse RM, Seppey M, Simão FA, Manni M, Ioannidis P, Klioutchnikov G, et al.

1193    BUSCO Applications from Quality Assessments to Gene Prediction and Phylogenomics.

1194    Mol Biol Evol. 2018;35: 543–548. doi:10.1093/molbev/msx319

1195    111.   Li L, Stoeckert CJ, Roos DS. OrthoMCL: Identification of ortholog groups for eukaryotic

1196    genomes. Genome Res. 2003;13: 2178–2189. doi:10.1101/gr.1224503

1197    112.   van Dongen S. Graph clustering by flow simulation. Graph Stimul by flow Clust.

1198    2000;PhD thesis: University of Utrecht. doi:10.1016/j.cosrev.2007.05.001

1199    113.   Madden T. The BLAST sequence analysis tool. BLAST Seq Anal Tool. 2013; 1–17.

1200    114.   Kocot KM, Citarella MR, Moroz LL, Halanych KM. PhyloTreePruner: A phylogenetic

1201    tree-based approach for selection of orthologous sequences for phylogenomics. Evol

1202    Bioinforma. 2013;2013: 429–435. doi:10.4137/EBO.S12813

1203    115.   Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:

1204    Improvements in Performance and Usability. Mol Biol Evol. 2013;30: 772–780.

1205    doi:10.1093/molbev/mst010

1206    116.   Mount DW. Using BLOSUM in Sequence Alignments. Cold Spring Harb Protoc.

1207    2008;2008: pdb.top39-pdb.top39. doi:10.1101/pdb.top39

1208    117.   Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. trimAl: a tool for automated

1209    alignment trimming in large-scale phylogenetic analyses. Bioinformatics. 2009;25: 1972–

1210     1973. doi:10.1093/bioinformatics/btp348

1211  118.  Price MN, Dehal PS, Arkin AP. FastTree 2 - Approximately maximum-likelihood trees

1212     for large alignments. PLoS One. 2010;5. doi:10.1371/journal.pone.0009490

1213  119.  Felsenstein J. Evolutionary trees from DNA sequences: A maximum likelihood approach.

1214     J Mol Evol. 1981;17: 368–376. doi:10.1007/BF01734359

1215  120.  Philippe H, Delsuc F, Brinkmann H, Lartillot N. Phylogenomics. Annu Rev Ecol Evol

1216     Syst. 2005;36: 541–562. doi:10.1146/annurev.ecolsys.35.112202.130205

1217  121.  Rokas A, Williams BL, King N, Carroll SB. Genome-scale approaches to resolving

1218     incongruence in molecular phylogenies. Nature. 2003;425: 798–804.

1219     doi:10.1038/nature02053

1220  122.  Edwards SV. Is a new and general theory of molecular systematics emerging? Evolution

1221     (N Y). 2009;63: 1–19. doi:10.1111/j.1558-5646.2008.00549.x

1222  123.  Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of

1223     large phylogenies. Bioinformatics. 2014;30: 1312–1313.

1224     doi:10.1093/bioinformatics/btu033

1225  124.  Le SQ, Gascuel O. An Improved General Amino Acid Replacement Matrix. Mol Biol

1226     Evol. 2008;25: 1307–1320. doi:10.1093/molbev/msn067

1227  125.  Mirarab S, Warnow T. ASTRAL-II: coalescent-based species tree estimation with many

1228     hundreds of taxa and thousands of genes. Bioinformatics. 2015;31: i44–i52.

1229     doi:10.1093/bioinformatics/btv234

1230  126.  Stamatakis A, Hoover P, Rougemont J. A Rapid Bootstrap Algorithm for the RAxML

1231     Web Servers. Renner S, editor. Syst Biol. 2008;57: 758–771.

1232     doi:10.1080/10635150802429642

1233    127.    Junier T, Zdobnov EM. The Newick utilities: high-throughput phylogenetic tree

1234            processing in the UNIX shell. Bioinformatics. 2010;26: 1669–1670.

1235            doi:10.1093/bioinformatics/btq243

1236    128.    Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol Biol Evol.

1237            2007;24: 1586–1591. doi:10.1093/molbev/msm088

1238    129.    dos Reis M, Donoghue PCJ, Yang Z. Bayesian molecular clock dating of species

1239            divergences in the genomics era. Nat Rev Genet. 2016;17: 71–80. doi:10.1038/nrg.2015.8

1240    130.    Reis M d., Yang Z. Approximate Likelihood Calculation on a Phylogeny for Bayesian

1241            Estimation of Divergence Times. Mol Biol Evol. 2011;28: 2161–2172.

1242            doi:10.1093/molbev/msr045

1243    131.    Cherry JM, Hong EL, Amundsen C, Balakrishnan R, Binkley G, Chan ET, et al.

1244            Saccharomyces Genome Database: the genomics resource of budding yeast. Nucleic

1245            Acids Res. 2012;40: D700–D705. doi:10.1093/nar/gkr1029

1246    132.    Pearson WR. An Introduction to Sequence Similarity ("Homology") Searching. Curr

1247            Protoc Bioinforma. 2013;42: 3.1.1-3.1.8. doi:10.1002/0471250953.bi0301s42

1248    133.    Eddy SR. Accelerated Profile HMM Searches. Pearson WR, editor. PLoS Comput Biol.

1249            2011;7: e1002195. doi:10.1371/journal.pcbi.1002195

1250    134.    GeneOntologyConsortium. The Gene Ontology (GO) database and informatics resource.

1251            Nucleic Acids Res. 2004;32: 258D–261. doi:10.1093/nar/gkh036

1252    135.    Klopfenstein D V., Zhang L, Pedersen BS, Ramírez F, Warwick Vesztrocy A, Naldi A, et

1253            al. GOATOOLS: A Python library for Gene Ontology analyses. Sci Rep. 2018;8: 10872.

1254            doi:10.1038/s41598-018-28948-z

1255    136.    Holm S. A simple sequentially rejective multiple test procedure. Scand J Stat. 1979;6: 65–

1256      70.

1257    137.   Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. KEGG as a reference

1258           resource for gene and protein annotation. Nucleic Acids Res. 2016;44: D457–D462.

1259           doi:10.1093/nar/gkv1070

1260    138.   Winzeler EA, Shoemaker DD, Astromoff A, Liang H, Anderson K, Andre B, et al.

1261           Functional characterization of the S-cerevisiae genome by gene deletion and parallel

1262           analysis. Science (80- ). 1999;285: 901–906. doi:10.1126/science.285.5429.901

1263    139.   Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.

1264           2012;9: 357–359. doi:10.1038/nmeth.1923

1265    140.   Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence

1266           Alignment/Map format and SAMtools. Bioinformatics. 2009;25: 2078–2079.

1267           doi:10.1093/bioinformatics/btp352

1268    141.   Weiß CL, Pais M, Cano LM, Kamoun S, Burbano HA. nQuire: a statistical framework for

1269           ploidy estimation using next generation sequencing. BMC Bioinformatics. 2018;19: 122.

1270           doi:10.1186/s12859-018-2128-z

1271    142.   Suyama M, Torrents D, Bork P. PAL2NAL: Robust conversion of protein sequence

1272           alignments into the corresponding codon alignments. Nucleic Acids Res. 2006;34.

1273           doi:10.1093/nar/gkl315

1274    143.   Liu L, Zhang J, Rheindt FE, Lei F, Qu Y, Wang Y, et al. Genomic evidence reveals a

1275           radiation of placental mammals uninterrupted by the KPg boundary. Proc Natl Acad Sci.

1276           2017;114: E7282–E7290. doi:10.1073/pnas.1616744114

1277    144.   Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely

1278           available Python tools for computational molecular biology and bioinformatics.

1279    Bioinformatics. 2009;25: 1422–1423. doi:10.1093/bioinformatics/btp163

1280    145.   Van Der Walt S, Colbert SC, Varoquaux G. The NumPy array: A structure for efficient

1281    numerical computation. Comput Sci Eng. 2011;13: 22–30. doi:10.1109/MCSE.2011.37

1282

1283

1284    **Tables**

1285    **Table 1.** Rate of sequence evolution hypotheses and results.

| Hypotheses for inter-lineage comparisons | Parameters | Fraction of genes significantly different than $H_O$ | Median ω values | | |
|---|---|---|---|---|---|
| | | | $\omega_{background}$ | $\omega_1$ | $\omega_2$ |
| **$H_O$: Uniform rate for all branches** **Figure S9A** | Single ω value | N/A | N/A | N/A | N/A |
| **$H_{FE-SE\ branch}$: Unique rates for FEL and SEL stem** **Figure S9B** | $\omega_{background} \neq \omega_1 \neq \omega_2$ $\omega_1$ = FEL stem branch $\omega_2$ = SEL stem branch | 678 genes (68.55% of examined genes) | 0.060 | 0.566 | 0.293 |
| **$H_{FE}$: Unique rates for FEL stem and FEL crown** **Figure S9C** | $\omega_{background} \neq \omega_1 \neq \omega_2$ $\omega_1$ = FEL stem branch $\omega_2$ = FEL crown branches | 743 genes (75.13% of examined genes) | 0.063 | 0.711 | 0.061 |
| **$H_{SE}$: Unique rates for SEL stem and SEL crown** **Figure S9D** | $\omega_{background} \neq \omega_1 \neq \omega_2$ $\omega_1$ = SEL stem branch $\omega_2$ = SEL crown branches | 528 genes (53.7% of examined genes) | 0.059 | 0.267 | 0.074 |
| **$H_{FE-SE\ crown}$: Unique rates for FEL crown and SEL crown** **Figure S9E** | $\omega_{background} \neq \omega_1 \neq \omega_2$ $\omega_1$ = FEL crown branches $\omega_2$ = SEL crown branches | 717 genes (72.5% of examined genes) | 0.010 | 0.062 | 0.074 |

1286

1287    **Main Figure Legends**

1288    **Fig 1. The evolutionary history and timeline of *Hanseniaspora* diversification and the**

1289    **stability of a long internode branch.**                 (A) Using 1,034 single-copy orthologous

1290    genes (SCOG), the evolutionary history of *Hanseniaspora* in geologic time revealed two well-

1291    support lineages termed the fast evolving and slow evolving lineages (FEL and SEL,

1292    respectively), which began diversifying around 87.2 and 53.6 million years ago (mya) after

1293    diverging 95.3 mya. (B) Among single-gene phylogenies where the FEL and SEL were

1294    monophyletic (*n* = 946), internode branch lengths leading up to each lineage revealed

1295    significantly longer internode branches leading up to the FEL ($0.62 \pm 0.38$ base substitutions /

1296    site) compared to the SEL ($0.17 \pm 0.11$ base substitutions / site) ($p < 0.001$; Paired Wilcoxon

1297    Rank Sum test). (C) Examination of the difference between internode branch lengths per single-

1298    gene tree revealed 932 single-gene phylogenies had a longer branch length in the FEL compared

1299    to the SEL (depicted in orange with values greater than 0), while the converse was only observed

1300    in 14 single-gene phylogenies (depicted in blue with values less than 0). Across all single-gene

1301    phylogenies, the average difference between the internode branch length leading up to the two

1302    lineages was 0.45.

1303

1304    **Fig 2. Gene presence and absence analyses reflect phenotype and reveal disrupted**

1305    **pathways.**                 (A) Examination of gene presence and absence (see *Methods*) revealed

1306    numerous genes that had been lost across *Hanseniaspora*. Specifically, 1,409 have been lost in

1307    the FEL, and 771 genes have been lost in the SEL. A Euler diagram represents the overlap of

1308    these gene sets. Both lineages have lost 748 genes, the FEL has lost an additional 661, and the

1309    SEL has lost an additional 23. (B) The *IMA* gene family (*IMA1-5*) encoding α-glucosidases,

60

1310    *MAL* (*MALx1-3*) loci, and *SUC2* are associated with growth on maltose, sucrose, raffinose, and

1311    melezitose. The *IMA* and *MAL* loci are largely missing among *Hanseniaspora* with the exception

1312    of homologs *MALx1*, which encode diverse transporters of the major facilitator superfamily

1313    whose functions are difficult to predict from sequence; as expected, *Hanseniaspora* spp. cannot

1314    grow on maltose, raffinose, and melezitose with the sole exception of *Hanseniaspora jakobsenii*,

1315    which has delayed/weak growth on maltose and is the only *Hanseniaspora* species with

1316    (*MALx3*), which encodes a homolog of the *MAL*-activator protein. (C) The genes involved with

1317    galactose degradation are largely missing among *Hanseniaspora* species, which correlates with

1318    their inability to grow on galactose. Genes that are present are depicted in white, and genes that

1319    are absent are depicted in black. The ability to grow, grow with delayed/weak growth on a given

1320    substrate, or the inability to grow is specified using white, grey, and black circles, respectively;

1321    dashes indicate no data. (D) Most genes involved in the thiamine biosynthesis pathway are

1322    absent among all *Hanseniaspora*. (E) Many genes involved in the methionine salvage pathway

1323    are absent among all *Hanseniaspora*. Absent genes are depicted in purple.

1324

1325    **Fig 3. Gene presence and absence in the budding yeast cell cycle.**            Examination

1326    of genes present and absent in the cell cycle of budding yeasts revealed numerous missing genes.

1327    Many genes are key regulators, such as *WHI5*; participate in spindle checkpoint processes and

1328    segregation, such as *MAD1* and *MAD2*; or DNA damage checkpoint processes, such as *MEC3*,

1329    *RAD9*, and *RFX1*. Genes missing in both lineages, the FEL, or the SEL are colored purple,

1330    orange, or blue, respectively. The "e" in the PHO cascade represents expression of Pho4:Pho2.

1331    Dotted lines with arrows indicate indirect links or unknown reactions. Lines with arrows indicate

1332    molecular interactions or relations. Circles indicate chemical compounds such as DNA.

1333

1334 **Fig 4. A panoply of genome maintenance and DNA repair genes are missing among**

1335 ***Hanseniaspora*, especially in the FEL.** Genes annotated as DNA repair genes

1336 according to gene ontology (GO:0006281) and child terms were examined for presence and

1337 absence in at least two-thirds of each lineage, respectively (268 total genes). 47 genes are

1338 missing among the FEL species, and 14 genes are missing among the SEL. Presence and absence

1339 of genes was clustered using hierarchical clustering (cladogram on the left) where each gene's

1340 ontology is provided as well. Genes with multiple gene annotations are denoted as such using the

1341 'multiple' term.

1342

1343 **Fig 5. dN/dS ($\omega$) analyses supports a historical burst of accelerated evolution in the FEL.**

1344 (A) The null hypothesis ($H_O$) that all branches in the phylogeny have the same $\omega$ value.

1345 Alternative hypotheses (B-E) evaluate $\omega$ along three sets of branches. (Bi) The alternative

1346 hypothesis ($H_{FE-SE}$ branch) examined $\omega$ values along the branch leading up the FEL and the SEL.

1347 (Bii) 311 supported $H_O$ and 678 genes supported $H_{FE-SE}$ branch. (Biii) Among the genes that

1348 supported $H_{FE-SE}$ branch, we examined the distribution of the difference between $\omega_1$ and $\omega_2$ as

1349 specified in part Bi. Here, a range of $\omega_1 - \omega_2$ of -3.5 to 3.5 is shown in the histogram.

1350 Additionally, we report the median $\omega_1$ and $\omega_2$ values, which are 0.57 and 0.29, respectively.

1351 (Biv) Among all genes examined, 0.39 genes significantly rejected $H_O$ and were faster in the

1352 FEL than the SEL, and 0.30 genes were faster in the SEL than the FEL. (Ci) The alternative

1353 hypothesis ($H_{FE}$) examined $\omega$ values along the branch leading up to the FEL and all branches

1354 thereafter ($FEL_{crown}$). (Cii) 246 genes supported $H_O$, and 743 genes supported $H_{FE}$. (Ciii) Among

1355 the genes that supported $H_{FE}$, we examined the distribution of the difference between $\omega_1$ and $\omega_2$

62

1356    as specified in part Ci. The median $\omega_1$ and $\omega_2$ values were 0.71 and 0.06, respectively. (Civ)

1357    Among all genes, 0.73 genes significantly rejected $H_O$ and were faster in the FEL than the

1358    FEL$_{crown}$, and 0.02 genes were faster in the FEL$_{crown}$ than the FEL. (Di) The alternative

1359    hypothesis (H$_{SE}$) examined $\omega$ values along the branch leading up to the SEL and all branches

1360    thereafter (SEL$_{crown}$). (Dii) 455 genes supported $H_O$, and 528 genes supported H$_{SE}$. (Diii) Among

1361    the genes that supported H$_{SE}$, we examined the distribution of the difference between $\omega_1$ and $\omega_2$

1362    as specified in part Di. The median $\omega_1$ and $\omega_2$ values were 0.27 and 0.07, respectively. (Div)

1363    Among all genes, 0.49 genes significantly rejected $H_O$ and were faster in the SEL than the

1364    SEL$_{crown}$, and 0.05 genes were faster in the SEL$_{crown}$ than the SEL. (Ei) The alternative

1365    hypothesis (H$_{FE-SE\ crown}$) examined $\omega$ values in the crown of the FEL$_{crown}$ and SEL$_{crown}$. (Eii) 272

1366    genes supported $H_O$, and 717 genes supported H$_{FE-SE\ crown}$. (Eiii) Among the genes that supported

1367    H$_{FE-SE\ crown}$, we examined the distribution of the difference between $\omega_1$ and $\omega_2$ as specified in part

1368    Di. The median $\omega_1$ and $\omega_2$ values were 0.06 and 0.07, respectively. (Eiv) Among all genes, 0.22

1369    genes significantly rejected $H_O$ and were faster in the FEL$_{crown}$ compared to the SEL$_{crown}$, and

1370    0.51 genes were faster in the SEL$_{crown}$ than the FEL$_{crown}$.

1371

1372    **Fig 6. Analyses of base substitutions and indels reveal a higher mutational load in the FEL**

1373    **compared to the SEL.**                (A) Analyses of substitutions at evolutionarily tractable

1374    sites among codon-based alignments revealed a higher number of base substitutions in the FEL

1375    compared to the SEL ($F(1) = 196.88$, $p < 0.001$; Multi-factor ANOVA) and an asymmetric

1376    distribution of base substitutions at codon sites ($F(2) = 1691.60$, $p < 0.001$; Multi-factor

1377    ANOVA). A Tukey Honest Significance Differences post-hoc test revealed a higher proportion

1378    of substitutions in the FEL compared to the SEL at evolutionarily tractable sites at the first ($n =$

63

1379    240,565; $p < 0.001$), second ($n = 318,987$; $p < 0.001$), and third ($n = 58,151$; $p = 0.02$) codon

1380    positions. (B) Analyses of the direction of base substitutions (i.e., G|C → A|T or A|T → G|C)

1381    reveals significant differences between the FEL and SEL ($F(1) = 447.1$, $p < 0.001$; Multi-factor

1382    ANOVA) and differences between the directionality of base substitutions ($F(1) = 914.5$, $p <$

1383    $0.001$; Multi-factor ANOVA). A Tukey Honest Significance Differences post-hoc test revealed a

1384    significantly higher proportion of substitutions were G|C → A|T compared to A|T → G|C among

1385    evolutionarily tractable sites that are G|C ($n = 232,546$) and A|T ($n = 385,157$) ($p < 0.001$),

1386    suggesting a general AT-bias of base substitutions. Additionally, there was a significantly higher

1387    proportion of evolutionary tractable sites with base substitutions in the FEL compared to the SEL

1388    ($p < 0.001$). More specifically, a higher number of base substitutions were observed in the FEL

1389    compared to the SEL for both G|C → A|T ($p < 0.001$) and A|T → G|C mutations ($p < 0.001$), but

1390    the bias toward AT was greater in the FEL. (C) Examinations of transition / transversion ratios

1391    revealed a lower transition / transversion ratio in the FEL compared to the SEL ($p < 0.001$;

1392    Wilcoxon Rank Sum test). (D) Comparisons of insertions and deletions revealed a significantly

1393    greater number of insertions ($p < 0.001$; Wilcoxon Rank Sum test) and deletions ($p < 0.001$;

1394    Wilcoxon Rank Sum test) in the FEL ($\bar{x}_{insertions} = 7521.11 \pm 405.34$; $\bar{x}_{deletions} = 3894.11 \pm 208.16$)

1395    compared to the SEL ($\bar{x}_{insertions} = 6049.571 \pm 155.85$; $\bar{x}_{deletions} = 2346.71 \pm 326.22$). (E and F)

1396    When adding the factor of size per insertion or deletion, significant differences were still

1397    observed between the lineages ($F(1) = 2102.87$, $p < 0.001$; Multi-factor ANOVA). A Tukey

1398    Honest Significance Differences post-hoc test revealed that most differences were caused by

1399    significantly more small insertions and deletions in the FEL compared to the SEL. More

1400    specifically, there were significantly more insertions in the FEL compared to the SEL for sizes 3-

1401    18 ($p < 0.001$ for all comparisons between each lineage for each insertion size), and there were

64

1402    significantly more deletions in the FEL compared to the SEL for sizes 3-21 ($p < 0.001$ for all

1403    comparisons between each lineage for each deletion size). Black lines at the top of each bar show

1404    the 95% confidence interval for the number of insertions or deletions for a given size. (G)

1405    Evolutionarily conserved homopolymers of sequence length two ($n = 17,391$), three ($n = 1,062$),

1406    four ($n = 104$), and five ($n = 5$) were examined for substitutions and indels. Statistically

1407    significant differences of the proportion mutated bases (i.e., (base substitutions + deleted bases +

1408    inserted bases) / total homopolymer bases) were observed between the FEL and SEL (F(1) =

1409    27.68, $p < 0.001$; Multi-factor ANOVA). Although the FEL had more mutations than the SEL

1410    for all homopolymers, a Tukey Honest Significance Differences post-hoc test revealed

1411    differences were statistically significant for homopolymers of two ($p = 0.02$) and three ($p =$

1412    0.003). Analyses of homopolymers using additional factors of mutation type (i.e., base

1413    substitution, insertion, deletion) and homopolymer sequence type (i.e., A|T and C|G

1414    homopolymers) can be seen in Fig S9. (H) G $\rightarrow$ T or C $\rightarrow$ A mutations are associated with the

1415    common and abundant oxidatively damaged base, 8-oxo-dG. When examining all substituted G

1416    positions for each species and their substitution direction, we found significant differences

1417    between different substitution directions (F(2) = 5682, $p < 0.001$; Multi-factor ANOVA). More

1418    importantly, a Tukey Honest Significance Differences post-hoc test revealed an over-

1419    representation of G $\rightarrow$ T or C $\rightarrow$ A in the FEL compared to the SEL ($p < 0.001$). (I) CC $\rightarrow$ TT

1420    dinucleotide substitutions are associated with UV damage. Using a CC|GG (left) and TT|AA

1421    (right) score, which is an indirect proxy for UV mutation damage where less UV damage would

1422    result in a higher CC|GG score and more UV damage would result in a higher TT|AA score, we

1423    found no significant differences when comparing CC|GG scores between the FEL, SEL, and

1424    outgroup taxa ($\chi2(2) = 5.964$, $p = 0.05$; Kruskal-Wallis rank sum test); however, when

65

1425    comparing the outgroup taxa to all *Hanseniaspora*, a significant difference was observed ($p =$

1426    0.03; Wilcoxon Rank Sum test). When examining TT|AA scores, we found significant

1427    differences between the FEL, SEL, and outgroup taxa ($\chi2(2) = 8.84$, $p = 0.01$; Kruskal-Wallis

1428    rank sum test). A post-hoc Dunn's test using the Benjamini-Hochberg method for multi-test

1429    correction revealed significant differences between the FEL and SEL compared to the outgroup

1430    taxa ($p = 0.01$ and $0.02$, respectively). A significant difference between all *Hanseniaspora* and

1431    the outgroup taxa were also observed ($p < 0.001$; Wilcoxon Rank Sum test). Results from the

1432    Kruskal-Wallis rank sum test and the Wilcoxon Rank Sum test are differentiated using lines and

1433    asterisks that are red and black, respectively.

1434

1435  **Supplementary Figure Legends**

1436  **Fig S1. *Hanseniaspora* have among the smallest genome sizes, lowest number of genes, and**

1437  **lowest percent GC content in the budding yeast subphylum Saccharomycotina.**

1438        (A) The genus *Hanseniaspora* (family Saccharomycodaceae) includes the smallest

1439  budding yeast genome. The FEL, SEL, and all of Saccharomycotina have an average genome

1440  size of 9.71 ± 1.32 Mb (min: 8.10; max: 14.05), 10.99 ± 1.66 Mb (min: 7.34; max: 12.17), 12.80

1441  ± 3.20 Mb (min: 7.34; max: 25.83), respectively. (B) The genus *Hanseniaspora* includes the

1442  budding yeast genome with the fewest genes. The FEL, SEL, and all of Saccharomycotina have

1443  an average number of genes per genome of 4,707.89 ± 633.56 (min: 3,923; max: 6,380),

1444  4,932.43 ± 289.71 (min: 4,624; max: 5,349), and 5,657.66 ± 1,044.78 (min: 3,923; max: 12,786),

1445  respectively. (C) The genus *Hanseniaspora* has among the lowest GC-content values in budding

1446  yeast genomes. The FEL, SEL, and all of Saccharomycotina GC-content values were 33.10 ±

1447  3.53% (min: 26.32; max: 37.17), 37.28 ± 2.05% (min: 34.82; max: 39.93), and 40.30 ± 5.71%

1448  (min: 25.2; max: 53.98), respectively. Families of Saccharomycotina are depicted on the y-axis.

1449  Median values are depicted with a line, and dashed lines indicate plus or minus one standard

1450  deviation from the median. To the right of each figure, boxplots depict the median and standard

1451  deviations of each grouping. The grey represents all of Saccharomycotina. Blue represents the

1452  SEL, and orange represents the FEL.

1453

1454  **Fig S2. Phylogenomics method pipeline.**        Using 25 *Hanseniaspora* proteomes and the

1455  proteomes of 4 outgroup taxa, 11,877 orthologous groups (OGs) of genes were identified. 1,143

1456  OGs with few paralogs were identified has having few paralogs – that is, ≥ 90% of species do

1457  not have paralogs and have one gene in the OG. The sequences of the 1,143 OGs were

67

1458    individually aligned, trimmed, had their evolutionary history inferred, and paralogs were

1459    trimmed based on tree topology. Using the resulting 1,142 OGs with paralogs trimmed,

1460    sequences were realigned and trimmed and had their evolutionary history inferred. If the

1461    outgroup taxa were not the earliest diverging taxa after serially rooting on the outgroup taxa, the

1462    OG was removed resulting in 1,034 OGs. Among these 1,034 OGs of genes, a concatenated

1463    1,034-gene matrix was constructed and used for reconstructing evolutionary history. Similarly,

1464    evolutionary history was inferred using coalescence of the 1,034 OG single-gene phylogenies.

1465

1466    **Fig S3. Concatenation and coalescence produce nearly identical and well-supported**

1467    **phylogenies that support two distinct lineages.**            (Left) Concatenation supports one

1468    lineage with a long internode branch leading to the clade, which we term the fast-evolving

1469    lineage (FEL) and another lineage with a much shorter internode branch length leading to the

1470    clade (SEL). (Right) Coalescence supports monophyly of the FEL and SEL. Minor discrepancies

1471    are observed between the topologies. Bipartitions without full support have their support values

1472    depicted. Support for concatenation and coalescence was determined using 100 rapid bootstrap

1473    replicates and local posterior support, respectively.

1474

1475    **Fig S4. Internode key to accompany divergence time estimate file per internode.**

1476        Internode identifiers for timetree analysis in Fig 1B. Associated mean divergence time

1477    and credible intervals can be found in File S2.

1478

1479    **Fig S5. BUSCO analyses reveals extensive gene 'missingness' across various taxonomic**

1480    **ranks.**          BUSCO analyses of *Hanseniaspora* proteomes using the Eukaryota ($n_{BUSCOs}$ =

68

1481    303), Fungi ($n_{BUSCOs}$ = 290), Dikarya ($n_{BUSCOs}$ = 1,312), Ascomycota ($n_{BUSCOs}$ = 1,315),

1482    Saccharomyceta ($n_{BUSCOs}$ =1,759), and Saccharomycetales ($n_{BUSCOs}$ = 1,711) orthoDB databases

1483    revealed numerous BUSCO genes are missing among *Hanseniaspora* genomes, in particular the

1484    FEL.

1485

1486    **Fig S6. A liberal targeted gene searching pipeline and the number of missing genes in at**

1487    **least two-thirds of FEL and SEL taxa.**    (A) A FASTA file for gene *X*, where gene *X*

1488    is the FASTA entry of a verified ORF in the *Saccharomyces cerevisiae* proteome, is used as a

1489    query to search for putative homologs in the Fungal reference sequence (refseq) database. The

1490    top 100 putative homologs were subsequently aligned. From the alignment, a Hidden Markov

1491    Model (HMM) was made. Using the HMM, gene *X* was searched for in the genome of each

1492    species from the FEL, SEL, and outgroup individually using a liberal e-value cut-off of 0.01 and

1493    a score of > 50. This pipeline yields presence and absence information of gene *X* among FEL,

1494    SEL, and outgroup taxa. This method was subsequently applied to all verified ORF in the *S.*

1495    *cerevisiae* proteome.

1496

1497    **Fig S7. Gene presence and absence reveals a putatively diminished gluconeogenesis**

1498    **pathway.**    Gene presence and absence analysis of genes that participate in the

1499    gluconeogenesis (A) and glycolysis (B) pathway reveal key missing genes in the

1500    gluconeogenesis pathway, suggestive of a diminished capacity for gluconeogenesis. More

1501    specifically, *PCK1*, which encodes the enzyme that converts oxaloacetic acid to

1502    phosphoenolpyruvate, and *FBP1*, which encodes the enzyme that converts fructose-1,6-

1503    bisphosphate to fructose-6-phospbate, are missing among all *Hanseniaspora* species.

69

1504

1505     **Fig S8. Base frequency plots reveal diversity in ploidy of *Hanseniaspora* species.**

1506     (A) A lack of Gaussian distributions suggests *H. occidentalis* var. *occidentalis*, *H.*

1507     *uvarum* CBS 314, and *H. guilliermondii* CBS 465 are haploid. (B) A single Gaussian distribution

1508     suggests *H. occidentalis* var. *citrica*, *H. osmophila* CBS 313, *H. meyeri*, *H. clermontiae*, *H.*

1509     *nectarophila*, *H. thailandica*, *H. pseudoguilliermondii*, *H. singularis*, and *K. hatyaiensis* are

1510     diploids. (C) Two Gaussian distributions suggest *H. lachancei* and *H. jakobsenii* are triploid. (D)

1511     Analyses of *H. vineae* CBS 2171, *H. valbyensis*, *Hanseniaspora* sp. CRUB 1602, and *H.*

1512     *opuntiae* base frequency distributions were ambiguous. Certain FEL species, such as *H.*

1513     *singularis, H. pseudoguilliermondii*, and *H. jakobsenii,* are potentially aneuploid, while evidence

1514     of aneuploidy in the SEL is observed in only *H. occidentalis* var. *citrica.*

1515

1516     **Fig S9. Analyses of homopolymers by sequence length, type, and type of mutation.**

1517     Significant differences among the proportion of mutated bases among homopolymers of

1518     various lengths were observed (Figure 5). Addition of variables (i.e., sequence type (A|T or C|G)

1519     and mutation type (base substitution, insertion, and deletion)) allowed for further determination

1520     of what types of mutations caused differences between the FEL and SEL. As shown in Figure 5,

1521     we observed significant differences in the numbers of mutations between the FEL and SEL (F =

1522     27.06, $p < 0.001$; Multi-factor ANOVA) as well as in the type of mutations (F = 1686.70, $p <$

1523     0.001; Multi-factor ANOVA). A Tukey Honest Significance Differences post-hoc test revealed

1524     that the proportion of nucleotides that underwent base substitutions was significantly greater than

1525     insertions ($p < 0.001$) and deletions ($p < 0.001$). We next focused on significant differences

1526     observed between the FEL and SEL when considering all factors. We observed significant

1527    differences between the FEL and SEL at A|T and C|G homopolymers with a length of 2 ($p$ =

1528    0.009 and $p < 0.001$, respectively), C|G homopolymers of length 3 ($p < 0.001$), and A|T

1529    homopolymers of length 5 ($p < 0.001$).

1530

1531    **Fig S10. Metrics reveal more radical amino acid substitutions in the FEL compared to**

1532    **SEL.**            Using Sneath's index and Epstein's coefficient of difference, the average

1533    difference among amino acid substitutions were determined among sites where the outgroup taxa

1534    had all the same amino acid. Using either metric, amino acid substitutions were significantly

1535    more drastic in the FEL compared to the SEL ($p < 0.001$; Wilcoxon Rank Sum test for both

1536    metrics).

1537

1538    **Fig S11. Mean protein similarity reveals immense diversity in *Hanseniaspora*.**

1539            The FEL spans a large amount of mean protein similarity when comparing various

1540    species to *H. uvarum*. Similarly, but to a lesser degree, the same is true for the SEL when

1541    comparing various species to *H. vineae*. The diversity observed in these lineages is roughly on

1542    par with genus-level differences within the family Saccharomycetaceae, humans to zebrafish,

1543    and thale cress (*Arabidopsis thaliana*) to Japanese rice.

Figure 1

A



B



C

Figure 2

**A**

FEL
n = 1409

SEL
n = 771

661    748    23

**D**

pyridoxal 5'phosphate
L-histidine

*THI5, THI11,*
*THI12, THI13*

hydroxymethylpyrimidine

*THI21,*
*THI20*

hydroxymethylpyrimidine phosphate

*THI21,*
*THI20*

4-amino-5-
hydroxymethyl-2-
methylpyrimidine-
pyrophosphate

NAD L-glycine

*THI4*

adenylated thiazole

*THI4*

4-methyl-5-
(β-hydroxyethyl)
thiazole phosphate

*THI6*

thiamine-phosphate

$H_2O$

thiamine ← *THI73* Extracellular thiamine

*THI80*

thiamine-pyrophosphate

**Color key**
Genes that are present
Genes lost in *Hanseniaspora*

**B**

*H. occidentalis var. occidentalis*
*H. occidentalis var. citrica*
*H. osmophila AWR13579*
*H. osmophila CBS 313*
*H. vineae NRRL Y-1626*
*H. vineae T02 19AF*
*H. sp. CRUB 1602*
*H. jakobsenii*
*H. singularis*
*H. valbyensis*
*H. guilliermondii CBS 465*
*H. guilliermondii UTAD222*
*H. lachancei*
*H. pseudoguilliermondii*
*H. opuntiae*
*H. meyeri*
*H. thailandica*
*H. clermontiae*
*H. nectarophila*
*H. uvarum DSM2768*
*H. uvarum 34-9*
*H. uvarum CBS 314*
*H. uvarum AWRI3580*
*Wickerhamomyces anomalus*
*Cyberlindnera jadinii*
*Kluyveromyces marxianus*
*Saccharomyces cerevisiae*
*H. gamundiae*

Sucrose, Raffinose, Maltose, Melezitose metabolism

IMA1
IMA2
IMA3
IMA4
IMA5
MALx1
MALx2
MALx3
SUC2

Maltose
Sucrose
Raffinose
Melezitose

**C**

Galactose metabolism

GAL1
GAL7
GAL10
Galactose

**Gene presence or absence**  □ gene present  ■ gene absent

**Ability to grow on substrate**  ○ growth  ● no growth  ◗ weak/delayed growth

**E**

spermidine

formate

$O_2$

*ADI1*

1,2-dihydroxy-3-keto-5-methylthiopentene

phosphate

$H_2O$  *UTR4*

5-(methylthio)2,3-dioxopentyl phosphate

*MDE1*

$H_2O$

5-methylthioribulose-1-phosphate

*MRI1*

5-methylthioribose-1-phosphate

spermine

adenine  phosphate

*MEU1*

2-oxo-4-methylthiobutanoate

putrescine

standard α amino acid
keto acid

*BAT2*
*BAT1*
*ARO8*
*ARO9*

L-methionine

$H_2O$
ATP
pyrophosphate
phosphate

*SAM2*
*SAM1*

S-adenosyl-L-methionine

*SPE2* → $CO_2$

S-adenosylmethioninamine

*SPE4*    *SPE3*

5'-methylthioadenosine

spermidine

Figure 3

**Spindle checkpoint**

**DNA damage checkpoint**

Unattached kinetochores    Misaligned spindles    DNA damage sensing

Second messenger signaling pathway

Nutrients low    Pheromone (mating signal)

Ubiquitin mediated proteolysis

cAMP low

MAPK signaling pathway

SCB

SCF Grr1

SCF Cdc4

Mps1

**Mad1**

**Mad2**

Dam1    Bub3 Cdc20 Mad3

Bub3    Bub3 Bub1 → Cdc20 Mad3

Rad17    **Mec3** Ddc1    Rad24

Mec1    Ddc2

**Rad9**    **Mrc1**

Chk1

Mcm1    Yhp1 Yox1

DNA Fus3

Far1

APC/C Cdc20

SCF Cdc4

Mitotic exit

**APC/C** Cdc20

**Mad2**

**Pds1** Securin

Rad53

DNA    Swi4 Swi6

Cln1/2 Cdc28

Clb5/6 Cdc28

Clb3/4 Cdc28

Clb1/2 Cdc28

Esp1 Separin

Dun1

**Rfx1** Cyc8 Tup1

DNA

Start    Cln3 Cdc28

**Whi5**

**Sic1**

Mih1 Cak1

Slk19

Cdc5

**Bfa1** Bub2    Lte1

Cell cycle arrest

Postreplicative DNA repair

Mbp1 Swi6

**Tah11** - Cdc6

Cks1

Swi5

Cdc55    **Spo12** Fob1    Tem1 Cdc15

Dbf2,20    Mob1

Transcription of target genes

DNA

Cdc45-MCM-**ORC**

SCF Met30

Swe1

PP2A

Net1

Phosphate ○ →high Pho81

Pho80 Pho85    Pcl1/2 Pho85

Cdc7 Dbf4

DNA ○

Cdr1    Cdr2 Pom1

Kcc4 Hsl7 Hsl1 Cdc5

Gin4

**Sic1**

**APC/C** Cdh1

Cdc14

Release from the nucleolus

Pho4 Pho2 — e → **Pho5**

S-phase proteins    DNA biosynthesis

Cohesin Condensin

Scc2 **Scc4**

→ Chromosome segregation

**G1**    **S**    **G2**    **M**

**Complexes**

|  |  |  |  |  |
|---|---|---|---|---|
| ORC (Origin Recognition Complex) | MCM (Mini-Chromosome Maintenance) complex | Anaphase-promoting complex (APC) | Condensin Smc2 Smc4 Ycs4 Bm1 Ycg1 | Cohesin Smc1 Smc3 **Mcd1** Irr1 |
| Orc1 Orc2 **Orc3** Orc4 Orc5 **Orc6** | Mcm2 Mcm3 Mcm4 Mcm5 Mcm6 Mcm7 | Cdc27 Apc11 Cdc23 **Apc2** **Apc4** **Swm1** Ama1 Apc1 **Cdc26** Cdh1 **Mnd2** **Apc5** Cdc20 Doc1 Cdc16 |  |  |

**Color key**

Genes that are present

Genes lost in *Hanseniaspora*

Genes lost in FEL yeast

Genes lost in SEL yeast

Unknown presence or absence

Figure 4

Figure 3

Figure 6