# Parallel Expansion and Divergence of an Adhesin Family in Pathogenic Yeasts

Rachel A. Smoak[*1], Lindsey F. Snyder[*2], Jan S. Fassler[^3], Bin Z. He[^3]


[1]Civil and Environmental Engineering, [2]Interdisciplinary Graduate Program in Genetics, [3]Biology
Department, University of Iowa, Iowa City, IA 52242
*These authors contributed equally
^Correspondence should be addressed to:
Bin Z. He, bin-he@uiowa.edu
Jan S. Fassler: jan-fassler@uiowa.edu

## Abstract

Opportunistic yeast pathogens arose multiple times in the Saccharomycetes class, including the
recently emerged, multidrug-resistant *Candida auris*. We show that homologs of a known yeast
adhesin family in *Candida albicans*, the Hyr/Iff-like (Hil) family, are enriched in distinct clades of
*Candida* species as a result of multiple, independent expansions. Following gene duplication,
the tandem repeat-rich region in these proteins diverged extremely rapidly and generated large
variations in length and β-aggregation potential, both of which were known to directly affect
adhesion. The conserved N-terminal effector domain was predicted to adopt a β-helical fold
followed by an α-crystallin domain, making it structurally similar to a group of unrelated bacterial
adhesins. Evolutionary analyses of the effector domain in *C. auris* revealed relaxed selective
constraint combined with signatures of positive selection, suggesting functional diversification
after gene duplication. Lastly, we found the Hil family genes to be enriched at chromosomal
ends, which likely contributed to their expansion via ectopic recombination and break-induced
replication. Combined, these results suggest that the expansion and diversification of adhesin
families generate variation in adhesion and virulence within and between species and are a key
step toward the emergence of fungal pathogens.

Running title: Parallel evolution of adhesins in yeasts
Keywords: *Candida auris*, adhesin, gene duplication, convergent evolution, opportunistic yeast
pathogen, natural selection

## Introduction

*Candida auris*, a newly emerged multidrug-resistant yeast pathogen, is associated with a high mortality rate (up to 60% in a multi-continent meta-analysis (Lockhart et al. 2017)) and has caused multiple outbreaks across the world (CDC global *C. auris* cases count, February 15th, 2021). As a result, it became the first fungal pathogen to be designated by CDC as an urgent threat (CDC 2019). The emergence of *C. auris* as a pathogen is part of a bigger evolutionary puzzle: *Candida* is a polyphyletic genus that contains most of the human yeast pathogens. Phylogenetically, species like *C. albicans*, *C. auris* and *C. glabrata* belong to distinct clades with close relatives that either don't or only rarely infect humans (Fig. 1A). This strongly suggests that the ability to infect humans evolved multiple times in yeasts (Gabaldón et al. 2016). Because many of the newly emerged *Candida* pathogens are either resistant or can quickly evolve resistance to antifungal drugs (Lamoth et al. 2018; Srivastava et al. 2018), it is urgent to understand how yeast pathogenesis arose and what increases their survival in the host. We reason that any shared genetic changes among independently derived *Candida* pathogens will reveal key factors for host adaptation.

Gene duplication and the subsequent functional divergence is a major source for the evolution of novel phenotypes (Zhang 2003; Qian and Zhang 2014; Kuang et al. 2016; Eberlein et al. 2017). In a genome comparison of seven pathogenic *Candida* species and nine low-pathogenic potential relatives, three of the top six pathogen-enriched gene families encode GlycosylPhosphatidylInositol (GPI)-anchored cell wall proteins, namely Hyr/Iff-like, Als-like and Pga30-like (Butler et al. 2009). The first two encode known fungal adhesins (Bailey et al. 1996; Hoyer 2001; Luo et al. 2010). These glycosylated cell wall proteins play key roles in fungal attachment to host endo- and epithelial cells, mediate biofilm formation and iron acquisition, and are well-established virulence factors (HOYER et al. 2008; de Groot et al. 2013; Lipke 2018). It has been suggested that expansion of the cell wall protein repertoire, particularly adhesins, is a key step towards the evolution of yeast pathogens (Gabaldón et al. 2016). This is supported by a study showing that several adhesin families independently expanded in *C. glabrata* and close relatives (Gabaldón et al. 2013). Interestingly, studies of pathogenic *Escherichia coli* found that multiple strains independently acquired genes mediating intestinal adhesion, giving credence to the hypothesis from a different kingdom (Reid et al. 2000).

Despite the importance of adhesins in both the evolution and virulence of *Candida* pathogens, there is a lack of detailed phylogenetic study for their evolutionary history (Hoyer 2001; Linder and Gustafsson 2008; Gabaldón et al. 2013). Even less is known about their sequence divergence and the role of natural selection in their evolution (Xie et al. 2011). In the

66  newly emerged *Candida auris*, individual adhesins have been characterized but there is little
67  information about their evolutionary relationship with homologs in other *Candida* species and
68  how their sequences diverged (Kean et al. 2018; Singh et al. 2019; Muñoz et al. 2021). In this
69  study we characterized the detailed evolutionary history of a yeast adhesin family and used *C.*
70  *auris* as a focal group to determine how adhesin sequences diverged under various natural
71  selection forces. To choose a candidate adhesin family in *C. auris*, we compared it with the well-
72  studied *C. albicans*, which belongs to the same CUG-Ser1 clade. Of the known adhesins in *C.*
73  *albicans*, *C. auris* lacks the Hwp family and has only three Als or Als-like proteins compared with
74  eight Als proteins in *C. albicans* (Muñoz et al. 2018). By contrast, *C. auris* has eight genes with
75  a Hyphal_reg_CWP (PF11765) domain found in the Hyr/Iff family in *C. albicans* (Muñoz et al.
76  2021). This family was one of the most highly enriched in pathogenic *Candida* species relative
77  to non-pathogenic ones (Butler et al. 2009). Transcriptomic studies identified two *C. auris*
78  Hyr/Iff-like (Hil) genes as being upregulated during biofilm formation and under antifungal
79  treatment (Kean et al. 2018). Interestingly, isolates from the less virulent *C. auris* Clade II lack
80  five of the eight Hil genes (Muñoz et al. 2021). It is currently not known whether the *C. auris HIL*
81  genes encode adhesins, how they relate to the *C. albicans* Hyr/Iff family genes and how their
82  sequences diverged after duplication.
83      We show that the Hil family independently expanded multiple times, including in *C. auris*
84  and *C. albicans*. Using *C. auris* as a focal species, we show in detail how sequence features
85  and predicted structures of the effector domain offer support for the hypothesis that its Hil family
86  members encode adhesins, while rates of nonsynonymous-to-synonymous substitutions reveal
87  varying strengths of selective constraint and positive selection acting on the effector domain
88  during the expansion of the family. The observed pattern of rapid divergence in the repeat-rich
89  central domain was found to be general across the entire family and led to large variations in
90  length and β-aggregation potential both between and within species, likely contributing to
91  phenotypic diversity in adhesion and virulence.

92  **Results**

93  **Phylogenetic distribution of the Hyr/Iff-like (Hil) family and its potential to encode**
94  **adhesin**

95  The Hyr/Iff family was first identified and characterized in *Candida albicans* (Bailey et al. 1996;
96  Richard and Plaine 2007). The family is defined by its N-terminal Hyphally regulated Cell Wall
97  Protein domain (Hyphal_reg_CWP, PF11765), followed by a highly variable central domain rich
98  in tandem repeats (Boisramé et al. 2011). Because the effector domain is more conserved than

99    the repeat region and plays a prominent role in mediating adhesion in known yeast adhesins

100   (Willaert 2018), here we define the Hyr/Iff-like (Hil) family as the group of evolutionarily related

101   proteins sharing the Hyphal_reg_CWP domain, different from a previous definition based on

102   sequence similarity in either the Hyphal_reg_CWP domain or the repeat region (Butler et al. 2009).

103        To determine the phylogenetic distribution of the Hil family and its association with the

104   pathogenic potential of species, we performed BLASTP searches using the Hyphal_reg_CWP

105   domain from three distantly related Hil homologs as queries (from *C. auris*, *C. albicans* and *C.*

106   *glabrata*). We scrutinized the database hits and searched additional assemblies to ensure that

107   their sequences are complete and accurate given the available genome assemblies (Text S1).

108   Using the criteria of *E*-value$<10^{-5}$ and query coverage>50%, we identified a total of 215 proteins

109   containing the Hyphal_reg_CWP domain from 32 species (Fig. 1A, Table S1). No credible hits

110   were identified outside the budding yeast subphylum even after a lower *E*-value cutoff of $10^{-3}$

111   was tested, suggesting that this family is specific to this group (Materials and Methods). Species

112   with eight or more Hil family genes fell largely within the Multi-Drug Resistant (MDR) and the

113   Candida/Lodderomyces (CaLo) clades, which include *C. auris* and *C. albicans*, respectively.

114   Only three such species were found outside of the two clades: *C. glabrata*, *M. bicuspidata* and

115   *K. africana*. *C. glabrata* is a major opportunistic pathogen that is more closely related to *S.*

116   *cerevisiae* than to most other *Candida* species (Dujon et al. 2004; Butler et al. 2009; Gabaldón

117   et al. 2013). *M. bicuspidata* is part of the CUG-Ser1 clade. While not a pathogen in humans, it is

118   a parasite of freshwater animals (Hall et al. 2010; Jiang et al. 2022). *K. africana* is not closely

119   related to any known yeast pathogen and its ecology is poorly understood (Gordon et al. 2011).

120        We then asked how many of the Hil family genes in each species are likely to encode

121   yeast adhesins. To get an initial estimate, we combined a Machine Learning tool for predicting

122   fungal adhesins (Chaudhuri et al. 2011) with predictions for the N-terminal signal peptide and C-

123   terminal GPI-anchor sequence, two features shared by the majority of known fungal adhesins

124   (Lipke 2018). Half of all Hil homologs passed all three tests (Fig. 1A). Notably, *M. bicuspidata*

125   has the largest Hil family among all species, but none of its 29 Hil genes passed all tests. We

126   found most of the identified hits in this species were short relative to the rest of the family (Fig.

127   S1), and 10 of the 29 hits were annotated as being incomplete in the RefSeq database. Further

128   analyses with a better assembled genome and functional studies are needed to determine if the

129   Hil family in this species has unique properties and functions.

130   **Independent expansion of the Hil family in multiple pathogenic *Candida* lineages**

131   Pathogenic yeast species have on average a larger Hil family and also more of its members

132   were predicted to encode adhesins than in low pathogenic-potential species (Fig. 1B, t-test with

133    unequal variance and Mann-Whitney U test both yielded $P < 0.005$, one-sided test). This naive

134    comparison doesn't account for phylogenetic relatedness between species and could result in a

135    false positive association (Levy et al. 2017; Bradley et al. 2018). To address this, we performed

136    phylogenetic logistic regression, which uses the known phylogeny to specify the residual

137    correlation structure among species with shared ancestry (Ives and Garland 2010). We tested

138    for associations between the pathogen status with either the total number of Hil homologs or the

139    number of putative adhesins in each species. Both tests were significant ($P = 0.005$ and 0.007,

140    respectively). Together, these results strongly support an enrichment of the Hil family and the

141    putative adhesins therein among the pathogenic yeast species.

142        Some adhesin families have undergone independent expansions even among closely

143    related species (Gabaldón et al. 2013). This would result in overestimation of the phylogenetic

144    signal in the above analysis. To further characterize the evolutionary history of the Hil family,

145    including among closely related *Candida* lineages, we reconstructed a species tree-aware

146    maximum likelihood phylogeny for the Hil family based on the Hyphal_reg_CWP domain

147    alignment (Fig. 1C, Fig. S2). We found that homologs from the MDR clade and the Candida /

148    Lodderomyces (CaLo) clade separated into two groups, suggesting that the duplications of the

149    Hil family genes in the two clades occurred independently. To better illustrate the history of gene

150    duplications in the Hil family, we reconciled the gene tree with the species tree and mapped the

151    number of duplications onto the species phylogeny (Materials and Methods). The result showed

152    that the Hil family has independently expanded multiple times, not only between clades but also

153    among closely related species within a clade, such as in *C. albicans* and *C. tropicalis* (Fig. 1D).

154    **Sequence features of the *C. auris* Hil family support their adhesin status**

155    Experiments have demonstrated that Hil family members function as adhesin in *C. albicans* and

156    more recently for one member in *C. glabrata* (Bailey et al. 1996; Boisramé et al. 2011; Reithofer

157    et al. 2021; Rosiana et al. 2021). To further evaluate the adhesin function of Hil family proteins,

158    we focused on *C. auris*, in which Hil family members were implicated in biofilm formation and

159    response to antifungal treatments, but still remain poorly characterized (Kean et al. 2018). We

160    named the eight *C. auris* Hil family proteins Hil1-Hil8 ordered by their length (Table S2). This

161    differs from the literature, which referred to them by their most closely related Hyr/Iff genes in *C.*

162    *albicans* (Kean et al. 2018; Jenull et al. 2021; Muñoz et al. 2021). The renaming was to avoid

163    the incorrect implication of one-to-one orthology between the two species (Fig. 1C).

164        To further assess the adhesin potential for the *C. auris* Hil family, we compared their

165    domain architecture and sequence features to those typical of known yeast adhesins, including

166    a signal peptide, an effector domain, a Ser/Thr-rich and highly glycosylated central domain with

167    tandem repeats and β-aggregation prone sequences and a GPI-anchor signal (Fig. 2A) (de

168    Groot et al. 2013; Lipke 2018). All eight *C. auris* Hil proteins followed this domain architecture

169    (Fig. 2B). Hil1-4 were additionally characterized by an array of regularly spaced β-aggregation

170    prone sequences (red ticks below the protein, Fig. 2B). All eight proteins also had elevated

171    Ser/Thr frequencies in their central domain and were predicted to be heavily O-glycosylated

172    (Fig. 2C). Predicted N-glycosylation was rare except in Hil5 and Hil6 (Fig. 2C). The overall

173    Ser/Thr frequencies in the Hil family proteins were significantly elevated compared with the rest

174    of the proteome (Fig. S3). All eight members were predicted to be fungal adhesins by

175    FungalRV, a support vector machine-based classifier that showed high sensitivity and specificity

176    in eight pathogenic fungi based on sequence features (Chaudhuri et al. 2011).

177    **Hyphal_reg_CWP domain in the Hil family is predicted to adopt a β-helical fold similar to**

178    **unrelated bacterial adhesin binding domains**

179    Crystal structures of the effector domain in several yeast adhesin families, including Als, Epa

180    and Flo, revealed carbohydrate or peptide binding activities supporting the proteins' adhesin

181    functions (Willaert 2018). The structure of the Hyphal_reg_CWP domain in the Hil family in this

182    study has not yet been experimentally determined. However, crystal structures for the effector

183    domains of two Adhesin-like Wall Proteins (Awp1 and Awp3b) in *C. glabrata*, which are distantly

184    related to those in the Hil family, were recently reported, and the predicted structure of one of *C.*

185    *glabrata*'s Hil family members (Awp2) was found to be highly similar to the two solved structures

186    (Reithofer et al. 2021). We used AlphaFold2 (Jumper et al. 2021) to predict the structures of the

187    effector domain for two *C. auris* Hil proteins, Hil1 and Hil7 (Fig. 3A, B). Both resemble the *C.*

188    *glabrata* Awp1 effector domain (Fig. 3C), consisting of a right-handed β-helix at the N-terminus

189    followed by an α-crystallin fold. There are three β-strands in each of the 9 rungs in the β-helix,

190    stacked into three parallel β-sheets (Fig. 3D). The α-crystallin domain consists of seven β-

191    strands forming two antiparallel β-sheets, adopting an immunoglobulin-like β-sandwich fold (Fig.

192    3E) (Koteiche and Mchaourab 1999; Stamler et al. 2005).

193          The β-strand-rich structure is typical of effector domains in known yeast adhesins, but

194    the β-helix fold at the N-terminus is uncommon (Willaert 2018). Proteins with a β-helix domain

195    often have carbohydrate-binding capabilities and act as enzymes, e.g., hydrolase and pectate

196    lyase (SCOP ID: 3001746). To gain further insight into Hyphal_reg_CWP domain's function, we

197    searched the PDB50 database for structures similar to what was predicted for *C. auris* Hil1

198    using DALI (Holm 2022). We identified a number of bacterial adhesins with a highly similar β-

199    helix fold but no α-crystallin domain (Table S3), e.g., Hmw1 from *H. influenzae* (PDB: 2ODL),

200    Tāpirins from *C. hydrothermalis* (PDB: 6N2C), TibA from enterotoxigenic *E. coli* (PDB: 4Q1Q)

201    and SRRP from *L. reuteri* (PDB: 5NY0). For comparison, the binding region of the Serine Rich

202    Repeat Protein 100-23 (SRRP$_{100\text{-}23}$) from *L. reuteri* was shown in Fig. 3F (Sequeira et al. 2018).

203    Together, these results strongly suggest that the Hyphal_reg_CWP domain in the *C. auris* Hil

204    family genes mediates adhesion. Additionally, the low sequence identity (12-15%) between the

205    yeast Hyphal_reg_CWP domain and the bacterial adhesins' binding regions further suggests

206    the two groups have convergently evolved a similar structure to achieve adhesion functions.

207    **Rapid divergence of the repeat-rich central domain in Hil family proteins in *C. auris***

208    While the overall domain architecture is well conserved, the eight Hil family proteins in *C. auris*

209    differ significantly in length and sequence of their central domains (Fig. 2B). While not involved

210    in ligand binding, central domains in yeast adhesins are known to play a critical role in mediating

211    adhesion: the length and stiffness of the central domain are essential for elevating and exposing

212    the effector domain (Frieman et al. 2002; Boisramé et al. 2011); and the tandem repeats and β-

213    aggregation sequences within them directly contribute to adhesion by mediating homophilic

214    binding and amyloid formation (Rauceo et al. 2006; Otoo et al. 2008; Frank et al. 2010; Wilkins

215    et al. 2018). Thus, divergence in the central domain of the Hil family has the potential to lead to

216    phenotypic diversity, as shown in *S. cerevisiae* (Verstrepen et al. 2004; Verstrepen et al. 2005).

217    To determine how the central domain sequences evolved in the *C. auris* Hil family, we

218    used dot plots both to reveal the tandem repeat structure within each protein and to examine the

219    similarity among the paralogs. A "dot" on the x-y plot indicates that the corresponding segments

220    (window size = 50 a.a.) from the two proteins on the x- and y-axes share similarity, with the gray

221    scale being proportional to the degree of similarity (Brodie et al. 2004). We found that *C. auris*

222    Hil1, 2, 3 and 4 share a ~44 aa repeat unit, whose copy number varies between 15 and 46,

223    driving differences in their protein lengths (Fig. 4A). These repeats have conserved periodicity

224    as well as sequence (Fig. 4B, Fig. S4). There are two interesting features of this 44 aa repeat

225    unit: a) it contains a heptapeptide "GVVIVTT" that is predicted to be strongly β-aggregation

226    prone, which explains the large number of regularly spaced β-aggregation motifs in Hil1-Hil4

227    (Fig. 2B); b) it is predicted to form three β-strands in the same orientation (Fig. 4B), raising an

228    interesting question of whether the tandem repeats may adopt a β-structure similar to that of the

229    effector domain. Hil7 and Hil8 encode the same repeat unit but have only one copy (Fig. 4A, red

230    boxes). By contrast, Hil5 and Hil6 encode very different low complexity repeats with a unit

231    length of ~5 aa. Their copy numbers range between 15 to 49 (Fig. 4C, D) have relatively low

232    Ser/Thr frequencies (Fig. 2C). Another consequence of encoding only one or zero copies of the

233    44 aa repeat unit found in Hil1-Hil4 is that Hil5-Hil8 are predicted to have 2-4 β-aggregation

234    prone sequences in contrast to 21-50 in Hil1-Hil4. For comparison, characterized yeast

235    adhesins contain 1-3 such sequences at a cutoff of >30% β-aggregation potential predicted by

236    TANGO (Fernandez-Escamilla et al. 2004; Ramsook et al. 2010; Lipke 2018). The variable

237    lengths, Ser/Thr frequencies and distribution of β-aggregation sequences, all resulting from the

238    evolution of the tandem repeats, suggest the intriguing possibility that the 8 different Hil proteins

239    in *C. auris* are non-redundant, playing distinct roles in cell adhesion and other cell-wall related

240    phenotypes.

241        Because tandem repeats are prone to recombination-mediated expansions and

242    contractions, we asked if there are variable numbers of tandem repeats (VNTR) among strains

243    in *C. auris*, which could generate diversity in cell adhesive properties as shown in *S. cerevisiae*

244    (Verstrepen et al. 2005). To answer this question, we identified homologs of Hil1-Hil4 in nine *C.*

245    *auris* strains from three geographically-stratified clades (Muñoz et al. 2018; Muñoz et al. 2021).

246    The genomes of these strains were *de novo* assembled using long-read technologies (Table

247    S4), which allowed us to confidently assess copy number variations within tandem repeats. We

248    identified a total of eight indel polymorphisms in Hil1-Hil4 (Table S5, example alignments in Fig.

249    S5). Except for one 16 aa deletion that is in a single Clade III strain, all seven other indels span

250    one or multiples of the repeat unit and affect all strains within a clade. This is consistent with

251    them being driven by recombination between repeats. The agreement within clades additionally

252    show the indels are not due to sequencing / assembly artifacts, which are not expected to follow

253    the clade labels. As previously reported, Clade II strains lack five of the eight Hil family proteins,

254    including Hil1-4 (Muñoz et al. 2021). Our phylogenetic analysis further showed that this was due

255    to gene losses within Clade II (Fig. S6). The potential relationship between the Hil family size

256    and the virulence profiles of Clade II strains is discussed later.

257    **Natural selection on the effector domain during the Hil family expansion in *C. auris***

258    Gene duplication provides raw materials for natural selection and is often followed by a period of

259    relaxed functional constraints on one or both copies, allowing for sub- or neo-functionalization

260    (Zhang 2003; Innan and Kondrashov 2010). Positive selection can be involved in this process,

261    which can lead to a ratio of nonsynonymous to synonymous substitution rates dN/dS > 1 (Yang

262    1998). Here we ask if the Hyphal_reg_CWP domain in *C. auris* Hil1-Hil8 experienced relaxed

263    selective constraints and/or positive selection following gene duplications, the latter of which

264    would suggest functional diversification. We chose to focus on the Hyphal_reg_CWP domain

265    because of its functional importance and because the high-quality alignment in this domain

266    allowed us to make confident evolutionary inferences (Fig. S7).

267        Because gene conversion between paralogs can cause distinct genealogical histories for

268    different parts of the alignment and mislead evolutionary inferences (Casola and Hahn 2009),

269    we first identified putatively non-recombining partitions using GARD (Kosakovsky Pond et al.

270    2006) (Fig. S8), and chose two partitions, P1-414 and P697-981, for maximum-likelihood based

271    analyses using PAML (Yang 2007) (Fig. 5A).

272         We first tested if a subset of the sites evolved under positive selection consistently on *all*

273    branches. We found moderate evidence supporting the hypothesis for the P697-981 partition,

274    where the M8 vs. M7 and M8 vs. M8a tests were significant at a 0.01 level, but the conservative

275    test M2a vs M1a was not (Table S6). All three tests were insignificant for the P1-414 partition.

276    Next, we tested for elevated dN/dS on selected branches of the tree, sign of relaxed selective

277    constraints or positive selection. We first estimated the dN/dS for each branch using a free-ratio

278    model and designated those with dN/dS greater than 10 as the "foreground" (Fig. 5B, C, "FG").

279    We found strong evidence for the FG branches to have a higher dN/dS than the remainder of

280    the tree (log-likelihood ratio test $P < 0.01$, Fig. 5D). There is no evidence, however, for the

281    dN/dS across the entire domain on the FG branches to be greater than 1 (Fig. 5D, *a*, row 2). We

282    then tested the more realistic scenario, where a subset of the sites on the FG branches were

283    subject to positive selection. Using the branch-site test 2 as defined in (Zhang et al. 2005), we

284    found evidence for positive selection on a subset of the sites on the FG branches for both

285    partitions (log-likelihood ratio test $P < 0.01$), and identified residues in both as candidate targets

286    of positive selection with a posterior probability greater than 0.99 (Fig. 5D). We conclude that

287    there is strong evidence for relaxed selective constraint on the Hyphal_reg_CWP domain on

288    some branches following gene duplications; there is also evidence for positive selection acting

289    on a subset of the sites on those branches. However, as the free-ratio model estimates were

290    noisy and the Empirical Bayes method used to identify the residues under selection lacks power

291    (Zhang et al. 2005) and can produce false positives (Nozawa et al. 2009), the specific branches

292    and residues implicated must be interpreted with caution.

293    **The yeast Hil family has adhesin-like domain architecture with rapidly diverging central**

294    **domain sequences**

295    We next examined the entire yeast Hil family to reveal the broader patterns of their evolution.

296    We found that the Hil family in general has elevated Ser/Thr content compared with the rest of

297    the proteome (Fig. S9). Moreover, the majority encode tandem repeats in the central domain

298    (Fig. 6A) and contain predicted β-aggregation prone sequences (Fig. 6B). Together, these

299    features further suggest that most yeast Hil family members encode fungal adhesins. While

300    these key features typical of yeast adhesins are conserved, the yeast Hil family exhibits extreme

301    variation in protein length, tandem-repeat content as well as in β-aggregation potential (Fig. 6A,

302    B, S10), extending the pattern seen in *C. auris* (Fig. 2). The length of the protein outside of the

303  Hyphal_reg_CWP domain has a mean ± standard deviation of 822.4±785.8 aa and a median of

304  608.5 aa. This large variation in protein length is almost entirely driven by the tandem repeats

305  (Fig. 6C, linear regression slope = 1.0, $r^2$ = 0.83). A subset of the Hil proteins (vertical bar in Fig.

306  6A, B) stand out in that they are both longer than the rest of the family (1745 vs 770 aa, median

307  protein length) and have an unusually large number of β-aggregation prone motifs (25 vs 6,

308  median number of TANGO hits per protein). The motifs in this group of proteins are regularly

309  spaced as a result of being part of the tandem repeat unit (median absolute deviation, or MAD,

310  of distances between adjacent TANGO hits less than 5 aa, Fig. 6D). The motif "GVVIVTT" and

311  its variants account for 61% of all hits in this subset and are not found in significant number in

312  the rest of the family. Together, these observations combined with previous experimental

313  studies showing a direct impact of adhesin length and β-aggregation potential on their function

314  (Verstrepen et al. 2005; Lipke et al. 2012) lead us to propose that the rapid divergence of the Hil

315  family following the parallel expansion led to functional diversification in adhesion in pathogenic

316  yeasts and may have contributed to their enhanced virulence.

317  **The yeast Hil family genes are preferentially located near chromosome ends**

318  Several well-characterized yeast adhesin families, including the Flo family in *S. cerevisiae* and

319  the Epa family in *C. glabrata*, are enriched in the subtelomeres (Teunissen and Steensma 1995;

320  De Las Peñas et al. 2003; Xu et al. 2020; Xu et al. 2021). This region is associated with high

321  rates of SNPs, indels and copy number variations, and can undergo ectopic recombination that

322  enables the spread of genes between chromosome ends or their losses (Mefford and Trask

323  2002; Anderson et al. 2015). To determine if the Hil family is also enriched in the subtelomeric

324  region, we compared their chromosomal locations with the background gene density distribution

325  (Fig. 7A) in species with a chromosomal level assembly (Table S7). To account for the shared

326  evolutionary history, we selected one species per closely related group such that the Hil family

327  homologs in these species were mostly derived through independent duplications based on our

328  gene tree (Fig. S2). The result showed that the Hil family genes are indeed enriched at

329  chromosomal ends (Fig. 7B). A goodness-of-fit test confirmed that the difference between the

330  chromosomal locations of the Hil family and the genome background is highly significant (*P* =

331  $1.3 \times 10^{-12}$). As ectopic recombination between subtelomeres has been suggested to underlie the

332  spread of gene families (Anderson et al. 2015), we hypothesize that the enrichment of the Hil

333  family towards the chromosome ends is both a cause and consequence of its parallel expansion

334  in different *Candida* lineages.

## Discussion

The repeated emergence of human pathogens in the Saccharomycetes class poses serious health threats, as many emerging pathogenic species are multi-drug resistant or quickly gain resistance (Lamoth et al. 2018; Srivastava et al. 2018). This raises an evolutionary question: are there shared genomic changes in independently derived *Candida* pathogens, which could be key factors in host adaptation? Yeast adhesin families were among the most enriched gene families in pathogenic lineages relative to the low pathogenic potential relatives (Butler et al. 2009). It has been proposed that expansion of adhesin families could be a key step in the emergence of novel yeast pathogens (Gabaldón et al. 2016). However, detailed phylogenetic studies supporting this hypothesis are rare (Gabaldón et al. 2013), and far less is known about how their sequences diverge and what selective forces are involved during the expansions (Xie et al. 2011; Muñoz et al. 2021). In this study, we found that the Hyr/Iff-like (Hil) family, defined by the conserved Hyphal_reg_CWP domain, is significantly enriched among distantly related pathogenic clades (Fig. 1A, B). This resulted from independent expansion of the family in these clades, including among closely related species (Fig. 1C, D). We also showed that the protein sequences diverged extremely rapidly after duplications, driven mostly by the evolution of the tandem repeats and resulting in large variations in protein length, Ser/Thr content and β-aggregation potential (Fig. 2B, C, Fig. 6). Our evolutionary analyses revealed evidence of relaxed selective constraint and a potential role of positive selection acting on the Hyphal_reg_CWP domain during the family's expansion in *C. auris* (Fig. 5). We also found the Hil family to be strongly enriched near chromosomal ends (Fig. 7). Overall, our results support the hypothesis that expansion and diversification of adhesin families is a key step towards the emergence of yeast pathogens.

### Genome assembly quality limits gene family evolution studies

Like any study of multi-gene family evolution, our work relies on and is limited by the quality of the genome assemblies. Two additional challenges in our study are due to the fact that Hil family genes are rich in tandem repeats (Fig. 2B, 6A), and many are located near chromosome ends (Fig. 7B), both of which pose problems for genome assemblies. For example, we found significant disagreement in length for 8 of the 16 Hil proteins in *C. tropicalis* between a long-read assembly and the RefSeq assembly, consistent with a recent study (Oh et al. 2020) (Table S8); in *C. glabrata*, we identified 13 Hil family genes in a long-read assembly (GCA_010111755.1) vs 3 in the RefSeq assembly (GCF_000002545.3); 12 of the 13 genes were in the subtelomeres (Xu et al. 2020). However, similar analyses in additional species didn't reveal these problems,

368 suggesting that the issues were at least in part due to difficulties in some genomes (Text S1).

369 Nonetheless, we acknowledge the possibility of missing homologs and inaccurate sequences,

370 especially in the tandem-repeat region. We thus believe the expected improvements in genome

371 assemblies due to advances in long-read sequencing technologies will be crucial for future

372 studies of the adhesin gene family in yeasts. It is worth noting that our main conclusions about

373 the parallel expansion of the Hil family and its rapid divergence patterns are robust with respect

374 to isolated problems as described above. Also, the long-read technology-based and *de novo*

375 assembled genomes for *C. auris* strains allowed us to confidently assess variation in the Hil

376 family size and tandem repeat copy number between paralogs and among individual strains

377 (Table S4). The accuracy of the tandem-repeat sequences in multiple strains in this species is

378 supported by the conservation of repeat copy numbers within clades (Table S5).

379 **Evidence for adhesin functions in the Hil family**

380 A few members of the Hil family, e.g., Iff4 in *C. albicans* and Awp2 in *C. glabrata* were shown to

381 mediate adhesiveness to polystyrene (Fu et al. 2008; Kempf et al. 2009; Reithofer et al. 2021).

382 While further experimental studies are needed to establish the adhesin functions of other Hil

383 family members, our work provides bioinformatic support for this hypothesis (Fig. 2, 6). The

384 predicted β-helix fold of the Hyphal_reg_CWP domain (Fig. 3), while unusual among

385 characterized yeast adhesins (Willaert 2018), is found in many virulence factors residing on the

386 surface of bacteria or viruses as well as enzymes that degrade or modify polysaccharides

387 (Table S3) (Kajava and Steven 2006). The elongated shape and rigid structure of the β-helix are

388 consistent with the functional requirements of adhesins, including the need to protrude from the

389 cell surface and the capacity for multiple binding sites along its length that facilitate adhesion. In

390 a bacterial adhesin – the serine rich repeat protein (SRRP) from the Gram-positive bacterium, *L.*

391 *reuteri* – a protruding, flexible loop in the β-helix was proposed to serve as a binding pocket for

392 its ligand (Sequeira et al. 2018). Such a feature is not apparent in the predicted structure of the

393 Hyphal_reg_CWP domain. Further studies are needed to elucidate the mechanism of action of

394 this domain and its potential substrates.

395 The cross-kingdom similarity in adhesin effector domain structure is intriguing in several

396 ways. First, it suggests convergent evolution in bacteria and yeasts. Second, it suggests that

397 what is known about the structure-function relationship in bacteria can provide insight into the

398 Hyphal_reg_CWP domain in yeast. Notably, the *Lr*SRRP shows a pH-dependent substrate

399 specificity that is potentially adapted to distinct host niches (Sequeira et al. 2018). Finally, the

400 similar structure and function of the bacterial and yeast adhesins could mediate cross-kingdom

401 interactions in natural and host environments (Uppuluri et al. 2018).

402    However, not all Hil family homologs are likely to encode adhesins. Sequence features
403    suggest some Hil family proteins may have non-adhesin functions. For example, 39 of 193 Hil
404    proteins (homologs labeled as incomplete were excluded) have the requisite signal sequence
405    (SP+), but lack a GPI anchor attachment site (gpi-, Fig. S1). One, Iff11 in *C. albicans*, was
406    shown to be secreted, and a null mutant of it was found to be hypersensitive to cell wall-
407    damaging agents and less virulent in a murine systemic infection model (Bates et al. 2007).
408    Moreover, 75% of these "SP+, gpi-" proteins are shorter than 600 amino acids, in contrast to
409    only 4% of the 117 proteins having both a signal peptide and a GPI anchor attachment site.
410    Such short, secreted proteins with tandem repeat sequences identical or similar to those
411    present in the cell-wall associated Hil protein counterparts may serve an important regulatory
412    function by bundling with wall associated adhesins as previously suggested for similar subclass
413    of proteins within the Als family (Oh et al. 2019). It is possible that the Hil family has evolved
414    diverse functions broadly related to cell adhesion.

415    **Ongoing diversification of the Hil family within species**

416    In addition to the parallel expansion and the subsequent rapid sequence divergence in the Hil
417    family between species, we and others also revealed population level variation in both the family
418    size and sequences within *C. auris* (Fig. S5, S6, Table S5) (Muñoz et al. 2021). Notably, among
419    the four geographically stratified clades, Clade II strains lost five of the eight Hil family members
420    (Fig. S6). Besides missing members of the Hil family, Clade II strains also lack seven of the
421    eight members of another GPI-anchor family that is specific to *C. auris* (Muñoz et al. 2021).
422    These coincide with the finding that Clade II strains were mostly associated with ear infections
423    (57/61 isolates according to (Kwon et al. 2019)) rather than hospital outbreaks, as reported for
424    strains from the other clades, and that they were generally less resistant to antifungal drugs
425    (Kwon et al. 2019; Welsh et al. 2019). This raises the question of whether the smaller adhesin
426    repertoire in Clade II strains limits their adhesive capability and results in a different pathology.
427    Similar expansion and contraction of adhesin families have been shown for the *C. glabrata* Hil
428    family (AWP Cluster V) and Epa family (Marcet-Houben et al. 2022), suggesting that dynamic
429    evolution of adhesin families in pathogenic yeasts could be a common pattern. Variation in the
430    tandem repeat copy number in Hil1-Hil4 among *C. auris* strains is also intriguing (Fig. S5). Prior
431    studies of the *S. cerevisiae* Flo proteins have shown that protein length directly impacts cellular
432    adhesion phenotypes (Verstrepen et al. 2005) and thus population level variation in adhesin
433    length could further contribute to phenotypic diversity. Lastly, scans for selective sweeps in *C.*
434    *auris* identified Hil and Als family members as being among the top 5% of all genes, suggesting

435    that adhesins are targets of natural selection in the recent evolutionary history of this newly

436    emerged pathogen (Muñoz et al. 2021).

437        Diversification of the adhesin repertoire within a strain can arise from a variety of

438    molecular mechanisms. For example, chimeric proteins generated through recombination

439    between Als family members or between an Als protein's N terminal effector domain and an

440    Hyr/Iff protein's repeat region have been shown (Butler et al. 2009; Zhao et al. 2011; Oh et al.

441    2019). Some of the adhesins with highly diverged central domains may have arisen in this

442    manner (Fig. S10). Gene conversion between members of the same family can also drive the

443    evolution of adhesin families within a species, as shown in *S. cerevisiae* and *C. glabrata*

444    (Verstrepen et al. 2004; Marcet-Houben et al. 2022). Evidence of this in the Hil family was

445    revealed in our analysis of recombination within the effector domain in *C. auris* (Fig. S8).

446    **Special properties of the central domain in *C. auris* Hil1-Hil4 and related Hil proteins**

447    A subset of Hil proteins represented by *C. auris* Hil1-Hil4 (Fig. 6A, B, vertical bar) stand out in

448    that they are much longer on average and encode a large number of β-aggregation prone

449    sequences compared with the rest of the family (Fig. 6B, D). Behind these properties is a

450    conserved ~44 aa repeat unit containing a highly β-aggregation prone sequence ("GVVIVTT"

451    and its variants) (Fig. 4B). β-aggregation prone sequences and the amyloid-like interaction they

452    mediate have been extensively studied, especially in the Als protein family in *C. albicans*: they

453    were experimentally shown to mediate aggregation (Otoo et al. 2008; Ramsook et al. 2010) and

454    were crucial for forming protein clusters on cell surfaces known as nanodomains in response to

455    physical tension or sheer forces (Alsteens et al. 2010; Lipke et al. 2012). Recently, they were

456    also shown to mediate cell-cell *trans* interactions via homotypic protein binding (Dehullu et al.

457    2019; Ho et al. 2019). This may underlie biofilm formation and kin discrimination (Smukalla et al.

458    2008; Brückner et al. 2020; Lipke et al. 2021). Most known yeast adhesins, including the Als

459    family proteins, encode between one and three β-aggregation prone sequences (Ramsook et al.

460    2010). *C. auris* Hil1-Hil4 and their close relatives are unusual in that they have as many as 50

461    such sequences, with each predicted by TANGO to have ~90% probability of aggregation,

462    whereas the positive threshold for the algorithm is only >5% over 5-6 residues (Fernandez-

463    Escamilla et al. 2004). The structural implications of the vast number of β-aggregation prone

464    motifs may be that such tandem repeat domains are constitutively amyloid in nature, rather than

465    requiring force or other stimuli as required by the Als proteins. The functional implications are

466    unclear without the requisite experimental tests. However, we speculate that variations in

467    protein length and β-aggregation potential resulting from the central domain divergence could

468    directly impact the adhesion functions as previously suggested (Verstrepen et al. 2005;

469    Boisramé et al. 2011; Lipke et al. 2012).

470    **Structural predictions of the tandem repeat region in *C. auris* Hil1 and Hil2**

471    Given the large number of ~44 aa repeats in the central domain of *C. auris* Hil1-Hil4 and the

472    prediction that each repeat encodes 3-4 short consecutive β-strands (Fig 4B), we wondered

473    what structural properties this region may have and how these features might contribute to the

474    adhesion function. We explored this question using threading based structural prediction tools

475    such as I-TASSER (Yang et al. 2015) and pDOMThreader (Lobley et al. 2009). For the tandem

476    repeat region in the central domain of Hil1, I-TASSER identified (S)-layer protein (SLP)

477    structures (e.g., RsaA from *C. crescentus*, SbsA and SbsC from *G. sterotherophilus*) as among

478    the top structural analogs. These β-strand-rich structures are known to self-assemble to form a

479    2-dimensional array on the surface of bacteria, mediating a range of functions including

480    adhesion to host cells in pathogens (Fagan and Fairweather 2014). pDOMThreader analyses of

481    the central domains in Hil1 and Hil2 identified a different set of templates, namely bacterial self-

482    associating proteins including Ag43a from uropathogenic *E. coli*, pertactin from *B. pertussis* and

483    the *H. influenzae* hap adhesin. Interestingly, these proteins have β-helical structures like the

484    Hyphal_reg_CWP domain, with the β-helices being involved in cell-cell interaction via an

485    interface along the long solenoidal axis for homotypic interactions, and mediating bacterial

486    clumping (Heras et al. 2014) and lead to biofilm formation in *H. influenzae* (Meng et al. 2011).

487    We speculate that the long repeat regions in Hil1 and Hil2 may similarly mediate cell-cell

488    interactions in *C. auris*.

489        The possibility that the central domains in Hil1 and Hil2 form a β-helix is interesting in

490    that β-helix is one of the commonly described structural motifs in functional amyloids, e.g., HET-

491    s from the fungus *Podospora anserina* (Wasmer et al. 2008). Such a solenoid-type amyloid is

492    distinguished from other amyloid types in that the β-helices formed by repeats within the same

493    protein, rather than among distinct monomeric proteins, are suggested to be stabilized not only

494    by polar zippers and hydrophobic contacts, but also by electrostatic interactions between the

495    alternating β-strands (Willbold et al. 2021). Other examples of amyloid forming proteins with a

496    predicted β-helix structure include the imperfect repeat domain in the human premelanosome

497    protein Pmel17 (Louros et al. 2016) and the extracellular curli proteins of Enterobacteriaceae

498    that are involved in biofilm formation and adhesion to host cells (Shewmaker et al. 2009). The

499    proposed solenoidal structure of the central domain of Hil1-Hil4 like proteins, if true, would have

500    two significant implications. First, it confers the necessary rigidity and extended conformation

501    required for cell wall anchored adhesins to extend into the surrounding extracellular milieu.

502 Second, the numerous β-strand rich repeats each containing a highly amyloid prone heptameric
503 sequence, and capable of wrapping into a solenoidal shaped stack, is likely to substantially
504 reduce the rate-limiting nucleation step, which limits the formation of, e.g., an Aβ amyloid fiber.
505 This would allow the formation of extracellular extensions at low protein concentrations without
506 the need for an extensive fiber lengthening process via the incorporation of additional
507 monomeric units. Finally, the observation of solenoid-mediated intercellular interactions in the
508 Hap adhesins suggests that Hil proteins may likewise have a biofilm related function.

509 **Genomic context**

510 As reported by (Muñoz et al. 2021), we found that the Hil family genes are preferentially located
511 near chromosomal ends in *C. auris* and also in other species examined (Fig. 7). This is similar
512 to previous findings for the Flo and Epa families (Teunissen and Steensma 1995; De Las Peñas
513 et al. 2003; Xu et al. 2020; Xu et al. 2021), as well as the Als genes in some species (Oh et al.
514 2021). This location bias of the Hil and other adhesin families is likely a key mechanism for their
515 dynamic expansion and sequence evolution via ectopic recombination (Anderson et al. 2015)
516 and by Break-Induced Replication (Bosco and Haber 1998; Sakofsky and Malkova 2017; Xu et
517 al. 2021). Another potential consequence of the Hil family genes being located in subtelomeres
518 is that they may be subject to epigenetic silencing as an additional regulatory mechanism, which
519 can be derepressed in response to stress (Ai et al. 2002). Such epigenetic regulation of the
520 adhesin genes was found to generate cell surface heterogeneity in *S. cerevisiae* (Halme et al.
521 2004) and lead to hyperadherent phenotypes in *C. glabrata* (Castaño et al. 2005).

522 **Concluding remarks**

523 To address the lack of candidate adhesins in *C. auris*, we identified and characterized the
524 Hyr/Iff-like (Hil) family in this species and all yeasts. Based on our results, we hypothesize that
525 expansion and diversification of adhesin gene families is a key step towards the evolution of
526 fungal pathogenesis and that variation in the adhesin repertoire contributes to within and
527 between species differences in the adhesive and virulence properties. Future experimental tests
528 of these hypotheses will be important biologically for improving our understanding of the fungal
529 adhesin repertoire, biotechnologically for inspiring additional nanomaterials, and biomedically for
530 advancing the development of *C. auris*-directed therapeutics.

531

532 **Materials and Methods**

533 **RESOURCE AVAILABILITY**

534  **Lead contact**

535  Further information and requests for resources and reagents should be directed to and will be

536  fulfilled by the Lead Contact, Bin Z. He (bin-he@uiowa.edu).

537  **Data and code availability**

538  All raw data and code for generating the intermediate and final results are available at the

539  GitHub repository at https://github.com/binhe-lab/C037-Cand-auris-adhesin. Upon publication,

540  this repository will be digitally archived with Zenodo and a DOI will be minted and provided to

541  ensure reproducibility.

542  **Software and algorithms list**

| NAME | REFERENCE | WEB OR DOWNLOAD URL |
|---|---|---|
| AlphaFold2 | (Jumper et al. 2021) | https://github.com/sokrypton/ColabFold (links to DeepMind Google Colab Notebook) |
| BLAST+ v2.12.0 | (Camacho et al. 2009) | https://blast.ncbi.nlm.nih.gov/ |
| ClipKit | (Steenwyk et al. 2020) | https://github.com/JLSteenwyk/ClipKIT |
| Clustal Omega v1.2.4 | (Sievers et al. 2011) | http://www.clustal.org/omega/ |
| Custom R, Python and shell scripts | This study | https://github.com/binhe-lab/C037-Cand-auris-adhesin |
| DALI | (Holm 2022) | http://ekhidna2.biocenter.helsinki.fi/dali/ |
| EMBOSS v6.6.0.0 | (Rice et al. 2000) | http://emboss.open-bio.org/ |
| FungalRV | (Chaudhuri et al. 2011) | http://fungalrv.igib.res.in/ |
| GeneRax v2.0.1 | (Morel et al. 2020) | https://github.com/BenoitMorel/GeneRax |
| HmmerWeb (hmmscan) | (Potter et al. 2018) | https://www.ebi.ac.uk/Tools/hmmer/search/hmmscan |
| Jalview v2.11 | (Waterhouse et al. 2009) | https://www.jalview.org/ |
| JDotter | (Brodie et al. 2004) | https://4virology.net/virology-ca-tools/jdotter/ |
| NetNGlyc v1.0 | (Gupta and Brunak 2002) | https://services.healthtech.dtu.dk/service.php?NetNGlyc-1.0 |

| NAME | REFERENCE | WEB OR DOWNLOAD URL |
|---|---|---|
| NetOGlyc v4.0 | (Steentoft et al. 2013) | https://services.healthtech.dtu.dk/service.php?NetOGlyc-4.0 |
| PAL2NAL.pl | (Suyama et al. 2006) | http://www.bork.embl.de/pal2nal/ |
| PAML v4.9e | (Yang 2007) | http://abacus.gene.ucl.ac.uk/software/paml.html |
| PredGPI | (Pierleoni et al. 2008) | http://gpcr.biocomp.unibo.it/predgpi/ |
| PSIPred | (Buchan and Jones 2019) | http://bioinf.cs.ucl.ac.uk/psipred/ |
| PyMol v2.5.2 | (Schrödinger, LLC 2021) | https://pymol.org/ |
| R package - ggtree v3.2.1 | (Yu 2020) | https://github.com/YuLab-SMU/ggtree |
| R package - phylolm | (Ho and Ané 2014) | https://cran.r-project.org/web/packages/phylolm/index.html |
| R package - rentrez v1.2.3 | (Winter 2017) | https://github.com/ropensci/rentrez |
| R package - treeio v1.18.1 | (Wang et al. 2020) | https://github.com/YuLab-SMU/treeio |
| R v4.1.0 | (R Core Team) | https://cran.r-project.org/ |
| RAxML v8.0.0 | (Stamatakis 2014) | https://cme.h-its.org/exelixis/web/software/raxml/ |
| RAxML-NG v1.1.0 | (Kozlov et al. 2019) | https://github.com/amkozlov/raxml-ng |
| RStudio v1.4 | (RStudio Team 2021) | https://www.rstudio.com/ |
| SignalP 6.0 | (Teufel et al. 2022) | http://www.cbs.dtu.dk/services/SignalP/ |
| TANGO v2.3.1 | (Fernandez-Escamilla et al. 2004) | http://tango.crg.es/ |
| XSTREAM | (Newman and Cooper 2007) | https://amnewmanlab.stanford.edu/xstream/download.jsp |

543

**METHOD DETAILS**

**Identify Hyr/Iff-like (Hil) family homologs in yeasts and beyond**

To identify the Hil family proteins in yeasts and beyond, we used the Hyphal_reg_CWP domain sequence from three distantly related Hil homologs as queries, namely, *C. albicans* Hyr1 (XP_722183.2), *C. auris* Hil1 (XP_028889033) and *C. glabrata* CAGL0E06600g (XP_722183.2). We performed BLASTP searches in the RefSeq protein database with an *E*-value cutoff of $1\times10^{-5}$, a minimum query coverage of 50% and with the low complexity filter on. All hits were from Ascomycota (yeasts) and all but one were from the Saccharomycetes class (budding yeast). A single hit was found in the fission yeast *Schizosacchromyces cryophilus*. Using that hit as the query, we searched all fission yeasts in the nr protein database, with a relaxed E-value cutoff of $10^{-3}$ and identified no additional hits. We thus excluded that one hit from downstream analyses. To supplement the RefSeq database, which lacks some yeast species such as those in the Nakaseomyces genus, we searched the Genome Resources for Yeast Chromosomes (GRYC, http://gryc.inra.fr/). Using the same criteria, we recovered 16 additional sequences. To allow for gene tree and species tree reconciliation, we excluded three species that are not part of the 322 species yeast phylogeny (Shen et al. 2018) and not a member of the Multidrug-Resistant clade (Muñoz et al. 2018). Further details, including additional quality control steps taken to ensure that the homolog sequences are accurate and complete, can be found in Text S1. In total, we curated a list of 215 Hil family homologs from 32 species.

**Gene family enrichment analysis**

To determine if the Hil family is enriched in the pathogenic yeasts, we performed two analyses. In the first analysis, we divided the species into pathogens vs low-pathogenic potential groups and performed a t-test with unequal variance (also known as Welch's test) as well as a non-parametric Mann-Whitney U test to compare the Hil family size in the two groups. For both tests, we used either the total size of the family, or the number of putative adhesins as the random variable, and the results were consistent. We excluded homologs from *M. bicuspidata* because 10 of its 29 Hil family proteins were annotated as incomplete in the RefSeq protein database, and also because as a parasite of freshwater crustaceans, it does not fit into either the human pathogen or the low-pathogenic potential group. *S. cerevisiae* was included in the comparison as an example of species with zero members of the Hil family. We chose *S. cerevisiae* because we could be confident about its lack of a Hil family homolog thanks to its well assembled and well annotated genome.

576        In the second test, we used phylogenetic logistic regression (Ives and Garland 2010) to

577    account for the phylogenetic relatedness between species. We used the `phyloglm` function in

578    the `phylolm` package in R, with {method = "logistic_IG10", btol = 50, boot = 100}. The species

579    tree, including the topology and branch lengths, were based on the 322 species phylogeny from

580    (Shen et al. 2018), supplemented by the phylogenetic relationship for the MDR clade based on

581    (Muñoz et al. 2018). The *P*-values based on phylogenetically specified residual correlations

582    were reported.

583    **Phylogenetic analysis of the Hil family and inference of gene duplications and losses**

584    To infer the evolutionary history of the Hil family, we reconstructed a maximum-likelihood tree

585    based on the alignment of the Hyphal_reg_CWP domain. First, we used hmmscan (HmmerWeb

586    version 2.41.2) to identify the location of the Hyphal_reg_CWP domain in each Hil homolog. We

587    used the "envelope boundaries" to define the domain in each sequence, and then aligned their

588    amino acid sequences using Clustal Omega with the parameter {--iter=5}. We then trimmed the

589    alignment using ClipKit with its default smart-gap trimming mode (Steenwyk et al. 2020).

590    RAxML-NG v1.1.0 was run in the MPI mode with the following parameters on the alignment:

591    "raxml-ng-mpi --all --msa INPUT --model LG+G --seed 123 --bs-trees autoMRE". The resulting

592    tree was corrected using GeneRax, which seeks to maximize the joint likelihood of observing

593    the alignment given the gene family tree (GFT) and observing the GFT given the species

594    phylogeny, using the parameter {--rec-model UndatedDL}. The species tree used is the same as

595    the one used for the phylogenetic logistic regression above. In addition to correcting the gene

596    family tree, GeneRax also reconciled it with the species tree and inferred duplication and loss

597    event counts on each branch. Tree annotation and visualization were done in R using the treeio

598    and ggtree packages (Wang et al. 2020; Yu 2020).

599        To infer the phylogenetic tree for the Hil family homologs in various *C. auris* strains and

600    infer gains and losses within species, we identified orthologs of the *HIL* genes in representative

601    strains from the four major clades of *C. auris* (B8441, B11220, B11221, B11243) (Muñoz et al.

602    2018). Orthologs from two MDR species, *C. haemuloni* and *C. pseudohaemulonis*, and from *D.*

603    *hansenii* were included to help root the tree. The gene tree was constructed as described

604    above. To root the tree, we first inferred a gene tree without the outgroup (*D. hansenii*)

605    sequences in the alignment. Then, the full alignment with the outgroup sequences along with

606    the gene tree from the first step were provided to RAxML to run the Evolutionary Placement

607    Algorithm (EPA) algorithm (Berger et al. 2011), which identified a unique root location. To

608    reconcile the gene tree with the species tree, we performed maximum likelihood-based gene

609    tree correction using GeneRax (v2.0.1) with the parameters: {--rec-model UndatedDL}. The

610    species tree was based on (Muñoz et al. 2018).

611    **Prediction for fungal adhesins and adhesin-related sequence features**

612    **1)** The potential of Hil homologs encoding fungal adhesins was assessed using FungalRV, a

613    Support Vector Machine-based fungal adhesin predictor (Chaudhuri et al. 2011). Proteins

614    passing the recommended cutoff of 0.511 were considered positive. **2)** Signal Peptide was

615    predicted using the SignalP 6.0 server, with the "organism group" set to Eukarya. The server

616    reported the proteins that had predicted signal peptides. No further filtering was done. **3)** GPI-

617    anchor was predicted using PredGPI using the General Model. Proteins with a false positive

618    rate of 0.01 or less were considered as containing a GPI-anchor. **4)** Tandem repeats were

619    identified using XSTREAM with the following parameters: {-i.7 -I.7 -g3 -e2 -L15 -z -Asub.txt -B -

620    O}, where the "sub.txt" was provided by the software package. **5)** β-aggregation prone

621    sequences were predicted using TANGO v2.3.1 with the following parameters: {ct="N" nt="N"

622    ph="7.5" te="298" io="0.1" tf="0" stab="-10" conc="1" seq="SEQ"}. **6)** Serine and Threonine

623    content in proteins were quantified using `freak` from the EMBOSS suite, with a sliding window

624    of 100 or 50 aa and a step size of 10 aa. To compare with proteome-wide distribution of Ser/Thr

625    frequency, the protein sequences for *C. albicans* (SC5314), *C. glabrata* (CBS138) and *C. auris*

626    (B11221) were downloaded from NCBI Assembly database (IDs in Table S7) and the frequency

627    of serine and threonine residues were counted for each protein. **7)** O-linked and N-linked

628    glycosylations were predicted using NetOGlyc (v4.0) and NetNGlyc (v1.0) servers.

629    **Structural prediction and visualization for the Hyphal_reg_CWP domain**

630    To perform structural predictions using AlphaFold2, we used the Google Colab notebook

631    (https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.i

632    pynb) authored by the DeepMind team. This is a reduced version of the full AlphaFold version 2

633    in that it searches a selected portion of the environmental BFD database, and doesn't use

634    templates. The Amber relaxation step is included, and no other parameters other than the input

635    sequences are required. DALI was used to search for similar structures in the PDB50 database.

636    Model visualization and annotation were done in PyMol v2.5.2. Secondary structure prediction

637    for *C. auris* Hil1's central domain was performed using PSIPred.

638    **Dotplot**

639    To determine the self-similarity and similarity between the eight *C. auris* Hil proteins, we made

640    dotplots using JDotter (Brodie et al. 2004). The window size and contrast settings were labeled

641    in the legends for the respective plots. The self-alignment for *C. auris* Hil1 tandem repeats was

642    visualized using Jalview v2.11.

**Identification of intraspecific tandem repeat copy number variations among *C. auris* strains**

644    To identify polymorphisms in Hil1-Hil4 in diverse *C. auris* strains, we downloaded the genome

645    sequences for the following strains from NCBI: Clade I - B11205, B13916; Clade II - B11220,

646    B12043, B13463; Clade III - B11221, B12037, B12631, B17721; Clade IV - B11245, B12342

647    (Table S4). The amino acid sequences for Hil1-Hil4 from the strain B8441 were used as the

648    query to search the nucleotide sequences of the above assemblies using TBLASTN, with the

649    following parameters {-db_gencode 12 -evalue 1e-150 -max_hsps 2}. Orthologs in each strain

650    were curated based on the BLAST hits to either the Hyphal_reg_CWP domain alone or the

651    entire protein query. All Clade II strains had no hits for Hil1-Hil4. Several strains in Clade I, III

652    and IV were found to lack one or more Hil proteins (Table S5). But upon further inspection, it

653    was found that they had significant TBLASTN hits for part of the query, e.g., the central domain,

654    and the hits were located at the end of a chromosome, suggesting the possibility of incomplete

655    or misassembled sequences. Further experiments will be needed to determine if those *HIL*

656    genes are present in those strains.

**Estimation of dN/dS ratios and model comparisons**

658    We used `codeml` in PAML (v4.9e) to perform evolutionary inferences on the Hyphal_reg_CWP

659    domain in *C. auris*. We first used Clustal Omega to align the amino acid sequences for the

660    Hyphal_reg_CWP domain from Hil1-Hil8 from *C. auris* similar to how we generated the multiple

661    sequence alignment for all Hil proteins. A closely related outgroup (XP_018709340.1 from *M.*

662    *bicuspidata*) was included to root the tree. We then generated a coding sequence alignment

663    from the protein alignment using PAL2NAL (Suyama et al. 2006). We used GARD (Kosakovsky

664    Pond et al. 2006) to analyze the coding sequence alignment to detect gene conversion events.

665    The web service of GARD on datamonkey.org was run with the following parameters: {data

666    type: nucleotide, run mode: normal, genetic code: yeast alternative nuclear, site-to-site rate

667    variation: general discrete, rate classes: 3}. Based on the results, we identified two putatively

668    non-recombining partitions, P1 = 1-414 and P2 = 697-981 (the numbers refer to the alignment

669    columns). We then separately analyzed the two partitions in PAML. To test hypotheses about

670    positive selection on a subset of the sites on all branches, we compared models M2a vs M1a,

671    M8 vs M7 and M8a vs M8. The first 4 models were specified by: {seqtype = 1, CodonFreq = 1,

672    model = 0, NSsites = 0,1,2,7,8, icode = 8, fix_kappa = 0, kappa = 2, fix_omega = 0, omega =

673    0.4, cleandata = 1}. The model M8a is additionally specified by { seqtype = 1, CodonFreq = 1,

674   model = 0, NSsites = 8, fix_omega = 1 and omega = 1, cleandata = 1}. To test hypotheses for

675   variable dN/dS on different branches (no variation across sites), we used {model = 0 or 1 or 2,

676   NSsites = 0}, with the rest being the same as the site tests. Model = 0 specified the single ratio

677   model, model = 1 the free ratio model and model = 2 the user-defined model. For the user-

678   defined model, we first used estimates from the free ratio model to designate a set of branches

679   with dN/dS > 10 as the foreground and then tested if their dN/dS was significantly different from

680   the rest of the tree by comparing a two-ratio model with the single-ratio model. Since the results

681   were significant, we further tested if the foreground dN/dS was significantly greater than 1, by

682   comparing the two-ratio model to a constrained version of the model where omega was fixed at

683   1. For branch-site test, we used {model = 2, NSsites = 2, fix_omega = 0, omega = .4} as the

684   alternative model and {model = 2, NSsites = 2, fix_omega = 1, omega = 1} as the null to test for

685   positive selection on a subset of the sites on the foreground branches. Sites under positive

686   selection were identified using the Bayes Empirical Bayes (BEB) procedure, with a posterior

687   probability threshold of 0.99.

688   **Chromosomal locations of Hil family genes**

689   To compare the chromosomal locations of the Hil family genes to the background distribution,

690   we selected eight species whose genomes were assembled to a chromosomal level and are not

691   within a closely related group, including *C. albicans*, *D. hansenii*, *C. orthopsilosis*, *K. africana*, *K.*

692   *lactis*, *N. dairenensis*, *C. auris* and *C. glabrata* (Table S7). We did not include some species,

693   e.g., *C. dubliniensis*, to minimize statistical dependence due to shared ancestry. The RefSeq

694   assembly for *C. auris* was included even though it was at a scaffold level because a recent

695   study showed that seven of its longest scaffolds were chromosome-length, allowing the

696   mapping of the scaffolds to chromosomes (Muñoz et al. 2021, Supplementary Table 1). To

697   determine the chromosomal locations of the Hil homologs in these eight species, we used

698   Rentrez v1.2.3 (Winter 2017) in R to retrieve their chromosome ID and coordinates. To calculate

699   the background gene density on each chromosome, we downloaded the feature tables for the

700   eight assemblies from the NCBI assembly database and calculated the location of each gene as

701   its start coordinate divided by the chromosome length. To compare the chromosomal location of

702   the Hil family genes to the genome background, we divided each chromosome into five equal-

703   sized bins based on the physical distance to the nearest chromosomal end. We calculated the

704   proportion of genes residing in each bin for the Hil family or for all protein coding genes. To

705   determine if the two distributions differ significantly from one another, we performed a

706   goodness-of-fit test using either a Log Likelihood Ratio (LLR) test or a Chi-Squared test, as

707   implemented in the XNomial package in R (Engels 2015). The LLR test *P*-value was reported.

**Reference**

Ai W, Bertram PG, Tsang CK, Chan TF, Zheng XFS. 2002. Regulation of subtelomeric silencing during stress response. *Mol. Cell* 10:1295–1305.

Alsteens D, Garcia MC, Lipke PN, Dufrêne YF. 2010. Force-induced formation and propagation of adhesion nanodomains in living fungal cells. *Proc. Natl. Acad. Sci. U. S. A.* 107:20744–20749.

Anderson MZ, Wigen LJ, Burrack LS, Berman J. 2015. Real-Time Evolution of a Subtelomeric Gene Family in Candida albicans. *Genetics* 200:907–919.

Bailey DA, Feldmann PJ, Bovey M, Gow NA, Brown AJ. 1996. The Candida albicans HYR1 gene, which is activated in response to hyphal development, belongs to a gene family encoding yeast cell wall proteins. *J. Bacteriol.* 178:5353–5360.

Bates S, de la Rosa JM, MacCallum DM, Brown AJP, Gow NAR, Odds FC. 2007. Candida albicans Iff11, a secreted protein required for cell wall structure and virulence. *Infect. Immun.* 75:2922–2928.

Berger SA, Krompass D, Stamatakis A. 2011. Performance, accuracy, and Web server for evolutionary placement of short sequence reads under maximum likelihood. *Syst. Biol.* 60:291–302.

Boisramé A, Cornu A, Da Costa G, Richard ML. 2011. Unexpected role for a serine/threonine-rich domain in the Candida albicans Iff protein family. *Eukaryot. Cell* 10:1317–1330.

742   Bosco G, Haber JE. 1998. Chromosome break-induced DNA replication leads to nonreciprocal
743       translocations and telomere capture. *Genetics* 150:1037–1047.

744   Bradley PH, Nayfach S, Pollard KS. 2018. Phylogeny-corrected identification of microbial gene
745       families relevant to human gut colonization. *PLOS Comput. Biol.* 14:e1006242.

746   Brodie R, Roper RL, Upton C. 2004. JDotter: a Java interface to multiple dotplots generated by
747       dotter. *Bioinforma. Oxf. Engl.* 20:279–281.

748   Brückner S, Schubert R, Kraushaar T, Hartmann R, Hoffmann D, Jelli E, Drescher K, Müller DJ,
749       Oliver Essen L, Mösch H-U. 2020. Kin discrimination in social yeast is mediated by cell
750       surface receptors of the Flo11 adhesin family.Brakhage AA, Barkai N, editors. *eLife*
751       9:e55587.

752   Buchan DWA, Jones DT. 2019. The PSIPRED Protein Analysis Workbench: 20 years on.
753       *Nucleic Acids Res.* 47:W402–W407.

754   Butler G, Rasmussen MD, Lin MF, Santos MAS, Sakthikumar S, Munro CA, Rheinbay E,
755       Grabherr M, Forche A, Reedy JL, et al. 2009. Evolution of pathogenicity and sexual
756       reproduction in eight Candida genomes. *Nature* 459:657–662.

757   Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
758       BLAST+: architecture and applications. *BMC Bioinformatics* 10:421.

759   Casola C, Hahn MW. 2009. Gene conversion among paralogs results in moderate false
760       detection of positive selection using likelihood methods. *J. Mol. Evol.* 68:679–687.

761   Castaño I, Pan S-J, Zupancic M, Hennequin C, Dujon B, Cormack BP. 2005. Telomere length
762       control and transcriptional regulation of subtelomeric adhesins in Candida glabrata. *Mol.*
763       *Microbiol.* 55:1246–1258.

764   CDC. 2019. Antibiotic resistance threats in the United States, 2019. *US Dep. Health Hum. Serv.*
765       *CDC* [Internet]. Available from: https://stacks.cdc.gov/view/cdc/82532

766   Chaudhuri R, Ansari FA, Raghunandanan MV, Ramachandran S. 2011. FungalRV: adhesin
767       prediction and immunoinformatics portal for human fungal pathogens. *BMC Genomics*
768       12:192.

769   De Las Peñas A, Pan S-J, Castaño I, Alder J, Cregg R, Cormack BP. 2003. Virulence-related
770       surface glycoproteins in the yeast pathogen Candida glabrata are encoded in
771       subtelomeric clusters and subject to RAP1- and SIR-dependent transcriptional silencing.
772       *Genes Dev.* 17:2245–2258.

773   Dehullu J, Valotteau C, Herman-Bausier P, Garcia-Sherman M, Mittelviefhaus M, Vorholt JA,
774       Lipke PN, Dufrêne YF. 2019. Fluidic Force Microscopy Demonstrates That Homophilic
775       Adhesion by Candida albicans Als Proteins Is Mediated by Amyloid Bonds between
776       Cells. *Nano Lett.* 19:3846–3853.

777   Dujon B, Sherman D, Fischer G, Durrens P, Casaregola S, Lafontaine I, de Montigny J, Marck
778       C, Neuvéglise C, Talla E, et al. 2004. Genome evolution in yeasts. *Nature* 430:35–44.

779    Eberlein C, Nielly-Thibault L, Maaroufi H, Dubé AK, Leducq J-B, Charron G, Landry CR. 2017.
780          The Rapid Evolution of an Ohnolog Contributes to the Ecological Specialization of
781          Incipient Yeast Species. *Mol. Biol. Evol.* 34:2173–2186.

782    Engels B. 2015. XNomial: Exact Goodness-of-Fit Test for Multinomial Data with Fixed
783          Probabilities. Available from: https://CRAN.R-project.org/package=XNomial

784    Fagan RP, Fairweather NF. 2014. Biogenesis and functions of bacterial S-layers. *Nat. Rev.*
785          *Microbiol.* 12:211–222.

786    Fernandez-Escamilla A-M, Rousseau F, Schymkowitz J, Serrano L. 2004. Prediction of
787          sequence-dependent and mutational effects on the aggregation of peptides and
788          proteins. *Nat. Biotechnol.* 22:1302–1306.

789    Frank AT, Ramsook CB, Otoo HN, Tan C, Soybelman G, Rauceo JM, Gaur NK, Klotz SA, Lipke
790          PN. 2010. Structure and Function of Glycosylated Tandem Repeats from Candida
791          albicans Als Adhesins. *Eukaryot. Cell* 9:405–414.

792    Frieman MB, McCaffery JM, Cormack BP. 2002. Modular domain structure in the Candida
793          glabrata adhesin Epa1p, a beta1,6 glucan-cross-linked cell wall protein. *Mol. Microbiol.*
794          46:479–492.

795    Fu Y, Luo G, Spellberg BJ, Edwards JE, Ibrahim AS. 2008. Gene overexpression/suppression
796          analysis of candidate virulence factors of Candida albicans. *Eukaryot. Cell* 7:483–492.

797    Gabaldón T, Martin T, Marcet-Houben M, Durrens P, Bolotin-Fukuhara M, Lespinet O, Arnaise
798          S, Boisnard S, Aguileta G, Atanasova R, et al. 2013. Comparative genomics of emerging
799          pathogens in the Candida glabrata clade. *BMC Genomics* 14:623.

800    Gabaldón T, Naranjo-Ortíz MA, Marcet-Houben M. 2016. Evolutionary genomics of yeast
801          pathogens in the Saccharomycotina. *FEMS Yeast Res.* 16.

802    Gordon JL, Armisén D, Proux-Wéra E, ÓhÉigeartaigh SS, Byrne KP, Wolfe KH. 2011.
803          Evolutionary erosion of yeast sex chromosomes by mating-type switching accidents.
804          *Proc. Natl. Acad. Sci.* 108:20024–20029.

805    de Groot PWJ, Bader O, de Boer AD, Weig M, Chauhan N. 2013. Adhesins in human fungal
806          pathogens: glue with plenty of stick. *Eukaryot. Cell* 12:470–481.

807    Gupta R, Brunak S. 2002. Prediction of glycosylation across the human proteome and the
808          correlation to protein function. *Pac. Symp. Biocomput. Pac. Symp. Biocomput.*:310–322.

809    Hall SR, Becker CR, Duffy MA, Cáceres CE. 2010. Variation in resource acquisition and use
810          among host clones creates key epidemiological trade-offs. *Am. Nat.* 176:557–565.

811    Halme A, Bumgarner S, Styles C, Fink GR. 2004. Genetic and Epigenetic Regulation of the FLO
812          Gene Family Generates Cell-Surface Variation in Yeast. *Cell* 116:405–415.

813    Heras B, Totsika M, Peters KM, Paxman JJ, Gee CL, Jarrott RJ, Perugini MA, Whitten AE,
814          Schembri MA. 2014. The antigen 43 structure reveals a molecular Velcro-like

815  mechanism of autotransporter-mediated bacterial clumping. *Proc. Natl. Acad. Sci. U. S.*
816  *A.* 111:457–462.

817  Ho L si T, Ané C. 2014. A linear-time algorithm for Gaussian and non-Gaussian trait evolution
818  models. *Syst. Biol.* 63:397–408.

819  Ho V, Herman-Bausier P, Shaw C, Conrad KA, Garcia-Sherman MC, Draghi J, Dufrene YF,
820  Lipke PN, Rauceo JM. 2019. An Amyloid Core Sequence in the Major Candida albicans
821  Adhesin Als1p Mediates Cell-Cell Adhesion. *mBio* 10:e01766-19.

822  Holm L. 2022. Dali server: structural unification of protein families. *Nucleic Acids Res.*:gkac387.

823  Hoyer LL. 2001. The ALS gene family of Candida albicans. *Trends Microbiol.* 9:176–180.

824  HOYER LL, GREEN CB, OH S-H, ZHAO X. 2008. Discovering the Secrets of the Candida
825  albicans Agglutinin-Like Sequence (ALS) Gene Family—a Sticky Pursuit. *Med. Mycol.*
826  *Off. Publ. Int. Soc. Hum. Anim. Mycol.* 46:1–15.

827  Innan H, Kondrashov F. 2010. The evolution of gene duplications: classifying and distinguishing
828  between models. *Nat. Rev. Genet.* 11:97–108.

829  Ives AR, Garland T. 2010. Phylogenetic logistic regression for binary dependent variables. *Syst.*
830  *Biol.* 59:9–26.

831  Jenull S, Tscherner M, Kashko N, Shivarathri R, Stoiber A, Chauhan M, Petryshyn A, Chauhan
832  N, Kuchler K. 2021. Transcriptome Signatures Predict Phenotypic Variations of Candida
833  auris. *Front. Cell. Infect. Microbiol.* [Internet] 11. Available from:
834  https://www.ncbi.nlm.nih.gov/labs/pmc/articles/PMC8079977/

835  Jiang H, Bao J, Xing Y, Li X, Chen Q. 2022. Comparative Genomic Analyses Provide Insight
836  Into the Pathogenicity of Metschnikowia bicuspidata LNES0119. *Front. Microbiol.*
837  13:939141.

838  Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates
839  R, Žídek A, Potapenko A, et al. 2021. Highly accurate protein structure prediction with
840  AlphaFold. *Nature*:1–11.

841  Kajava AV, Steven AC. 2006. The turn of the screw: variations of the abundant beta-solenoid
842  motif in passenger domains of Type V secretory proteins. *J. Struct. Biol.* 155:306–315.

843  Kean R, Delaney C, Sherry L, Borman A, Johnson EM, Richardson MD, Rautemaa-Richardson
844  R, Williams C, Ramage G. 2018. Transcriptome Assembly and Profiling of Candida auris
845  Reveals Novel Insights into Biofilm-Mediated Resistance. *mSphere* [Internet] 3.
846  Available from: https://msphere.asm.org/content/3/4/e00334-18

847  Kempf M, Cottin J, Licznar P, Lefrançois C, Robert R, Apaire-Marchais V. 2009. Disruption of
848  the GPI protein-encoding gene IFF4 of Candida albicans results in decreased adherence
849  and virulence. *Mycopathologia* 168:73–77.

850 Kosakovsky Pond SL, Posada D, Gravenor MB, Woelk CH, Frost SDW. 2006. Automated
851      Phylogenetic Detection of Recombination Using a Genetic Algorithm. *Mol. Biol. Evol.*
852      23:1891–1901.

853 Koteiche HA, Mchaourab HS. 1999. Folding pattern of the alpha-crystallin domain in alphaA-
854      crystallin determined by site-directed spin labeling. *J. Mol. Biol.* 294:561–577.

855 Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. 2019. RAxML-NG: a fast, scalable and
856      user-friendly tool for maximum likelihood phylogenetic inference. *Bioinforma. Oxf. Engl.*
857      35:4453–4455.

858 Kuang MC, Hutchins PD, Russell JD, Coon JJ, Hittinger CT. 2016. Ongoing resolution of
859      duplicate gene functions shapes the diversification of a metabolic network. *eLife*
860      5:e19027.

861 Kwon YJ, Shin JH, Byun SA, Choi MJ, Won EJ, Lee D, Lee SY, Chun S, Lee JH, Choi HJ, et al.
862      2019. Candida auris Clinical Isolates from South Korea: Identification, Antifungal
863      Susceptibility, and Genotyping. *J. Clin. Microbiol.* 57:e01624-18.

864 Lamoth F, Lockhart SR, Berkow EL, Calandra T. 2018. Changes in the epidemiological
865      landscape of invasive candidiasis. *J. Antimicrob. Chemother.* 73:i4–i13.

866 Levy A, Salas Gonzalez I, Mittelviefhaus M, Clingenpeel S, Herrera Paredes S, Miao J, Wang K,
867      Devescovi G, Stillman K, Monteiro F, et al. 2017. Genomic features of bacterial
868      adaptation to plants. *Nat. Genet.* 50:138–150.

869 Linder T, Gustafsson CM. 2008. Molecular phylogenetics of ascomycotal adhesins—A novel
870      family of putative cell-surface adhesive proteins in fission yeasts. *Fungal Genet. Biol.*
871      45:485–497.

872 Lipke PN. 2018. What We Do Not Know about Fungal Cell Adhesion Molecules. *J. Fungi Basel*
873      *Switz.* 4.

874 Lipke PN, Garcia MC, Alsteens D, Ramsook CB, Klotz SA, Dufrêne YF. 2012. Strengthening
875      relationships: amyloids create adhesion nanodomains in yeasts. *Trends Microbiol.*
876      20:59–65.

877 Lipke PN, Mathelié-Guinlet M, Viljoen A, Dufrêne YF. 2021. A New Function for Amyloid-Like
878      Interactions: Cross-Beta Aggregates of Adhesins form Cell-to-Cell Bonds. *Pathogens*
879      10:1013.

880 Lobley A, Sadowski MI, Jones DT. 2009. pGenTHREADER and pDomTHREADER: new
881      methods for improved protein fold recognition and superfamily discrimination.
882      *Bioinforma. Oxf. Engl.* 25:1761–1767.

883 Lockhart SR, Etienne KA, Vallabhaneni S, Farooqi J, Chowdhary A, Govender NP, Colombo
884      AL, Calvo B, Cuomo CA, Desjardins CA, et al. 2017. Simultaneous Emergence of
885      Multidrug-Resistant Candida auris on 3 Continents Confirmed by Whole-Genome
886      Sequencing and Epidemiological Analyses. *Clin. Infect. Dis. Off. Publ. Infect. Dis. Soc.*
887      *Am.* 64:134–140.

888    Louros NN, Baltoumas FA, Hamodrakas SJ, Iconomidou VA. 2016. A β-solenoid model of the
889        Pmel17 repeat domain: insights to the formation of functional amyloid fibrils. *J. Comput.*
890        *Aided Mol. Des.* 30:153–164.

891    Luo G, Ibrahim AS, Spellberg B, Nobile CJ, Mitchell AP, Fu Y. 2010. Candida albicans Hyr1p
892        Confers Resistance to Neutrophil Killing and Is a Potential Vaccine Target. *J. Infect. Dis.*
893        201:1718–1728.

894    Marcet-Houben M, Alvarado M, Ksiezopolska E, Saus E, de Groot PWJ, Gabaldón T. 2022.
895        Chromosome-level assemblies from diverse clades reveal limited structural and gene
896        content variation in the genome of Candida glabrata. *BMC Biol.* 20:226.

897    Mefford HC, Trask BJ. 2002. The complex structure and dynamic evolution of human
898        subtelomeres. *Nat. Rev. Genet.* 3:91–102.

899    Meng G, Spahich N, Kenjale R, Waksman G, St Geme JW. 2011. Crystal structure of the
900        Haemophilus influenzae Hap adhesin reveals an intercellular oligomerization mechanism
901        for bacterial aggregation. *EMBO J.* 30:3864–3874.

902    Morel B, Kozlov AM, Stamatakis A, Szöllősi GJ. 2020. GeneRax: A Tool for Species-Tree-
903        Aware Maximum Likelihood-Based Gene  Family Tree Inference under Gene
904        Duplication, Transfer, and Loss. *Mol. Biol. Evol.* 37:2763–2774.

905    Muñoz JF, Gade L, Chow NA, Loparev VN, Juieng P, Berkow EL, Farrer RA, Litvintseva AP,
906        Cuomo CA. 2018. Genomic insights into multidrug-resistance, mating and virulence in
907        Candida auris and related emerging species. *Nat. Commun.* 9:5346.

908    Muñoz JF, Welsh RM, Shea T, Batra D, Gade L, Howard D, Rowe LA, Meis JF, Litvintseva AP,
909        Cuomo CA. 2021. Clade-specific chromosomal rearrangements and loss of subtelomeric
910        adhesins in Candida auris. *Genetics* [Internet]. Available from:
911        https://doi.org/10.1093/genetics/iyab029

912    Newman AM, Cooper JB. 2007. XSTREAM: A practical algorithm for identification and
913        architecture modeling of tandem repeats in protein sequences. *BMC Bioinformatics*
914        8:382.

915    Nozawa M, Suzuki Y, Nei M. 2009. Reliabilies of identifying positive selection by the branch-
916        site and the site-prediction methods. *Proc. Natl. Acad. Sci.* 106:6700–6705.

917    Oh S-H, Isenhower A, Rodriguez-Bobadilla R, Smith B, Jones J, Hubka V, Fields C, Hernandez
918        A, Hoyer LL. 2020. Pursuing Advances in DNA Sequencing Technology to Solve a
919        Complex Genomic Jigsaw Puzzle: The Agglutinin-Like Sequence (ALS) Genes of
920        Candida tropicalis. *Front. Microbiol.* 11:594531.

921    Oh S-H, Schliep K, Isenhower A, Rodriguez-Bobadilla R, Vuong VM, Fields CJ, Hernandez AG,
922        Hoyer LL. 2021. Using Genomics to Shape the Definition of the Agglutinin-Like
923        Sequence (ALS) Family in the Saccharomycetales. *Front. Cell. Infect. Microbiol.*
924        11:794529.

925    Oh S-H, Smith B, Miller AN, Staker B, Fields C, Hernandez A, Hoyer LL. 2019. Agglutinin-Like
926        Sequence (ALS) Genes in the Candida parapsilosis Species Complex: Blurring the

927        Boundaries Between Gene Families That Encode Cell-Wall Proteins. *Front. Microbiol.*
928        10:781.

929    Otoo HN, Lee KG, Qiu W, Lipke PN. 2008. Candida albicans Als Adhesins Have Conserved
930        Amyloid-Forming Sequences. *Eukaryot. Cell* 7:776–782.

931    Pierleoni A, Martelli PL, Casadio R. 2008. PredGPI: a GPI-anchor predictor. *BMC Bioinformatics*
932        9:392.

933    Potter SC, Luciani A, Eddy SR, Park Y, Lopez R, Finn RD. 2018. HMMER web server: 2018
934        update. *Nucleic Acids Res.* 46:W200–W204.

935    Qian W, Zhang JG. 2014. Genomic evidence for adaptation by gene duplication. *Genome*
936        *Res.*:gr.172098.114.

937    R Core Team. R: A Language and Environment for Statistical Computing. Vienna, Austria: R
938        Foundation for Statistical Computing Available from: https://www.R-project.org

939    Ramsook CB, Tan C, Garcia MC, Fung R, Soybelman G, Henry R, Litewka A, O'Meally S, Otoo
940        HN, Khalaf RA, et al. 2010. Yeast cell adhesion molecules have functional amyloid-
941        forming sequences. *Eukaryot. Cell* 9:393–404.

942    Rauceo JM, De Armond R, Otoo H, Kahn PC, Klotz SA, Gaur NK, Lipke PN. 2006. Threonine-
943        rich repeats increase fibronectin binding in the Candida albicans adhesin Als5p.
944        *Eukaryot. Cell* 5:1664–1673.

945    Reid SD, Herbelin CJ, Bumbaugh AC, Selander RK, Whittam TS. 2000. Parallel evolution of
946        virulence in pathogenic Escherichia coli. *Nature* 406:64–67.

947    Reithofer V, Fernández-Pereira J, Alvarado M, de Groot P, Essen L-O. 2021. A novel class of
948        Candida glabrata cell wall proteins with β-helix fold mediates adhesion in clinical
949        isolates. *PLoS Pathog.* 17:e1009980.

950    Rice P, Longden I, Bleasby A. 2000. EMBOSS: the European Molecular Biology Open Software
951        Suite. *Trends Genet. TIG* 16:276–277.

952    Richard ML, Plaine A. 2007. Comprehensive Analysis of Glycosylphosphatidylinositol-Anchored
953        Proteins in Candida albicans. *Eukaryot. Cell* 6:119–133.

954    Rosiana S, Zhang L, Kim GH, Revtovich AV, Uthayakumar D, Sukumaran A, Geddes-McAlister
955        J, Kirienko NV, Shapiro RS. 2021. Comprehensive genetic analysis of adhesin proteins
956        and their role in virulence of Candida albicans. *Genetics* [Internet]. Available from:
957        https://doi.org/10.1093/genetics/iyab003

958    RStudio Team. 2021. RStudio: Integrated Development Environment for R. Boston, MA:
959        RStudio, PBC Available from: http://www.rstudio.com/

960    Sakofsky CJ, Malkova A. 2017. Break induced replication in eukaryotes: mechanisms,
961        functions, and consequences. *Crit. Rev. Biochem. Mol. Biol.* 52:395–413.

962    Schrödinger, LLC. 2021. The PyMOL Molecular Graphics System, Version 2.5.2.

963 Sequeira S, Kavanaugh D, MacKenzie DA, Šuligoj T, Walpole S, Leclaire C, Gunning AP,
964       Latousakis D, Willats WGT, Angulo J, et al. 2018. Structural basis for the role of serine-
965       rich repeat proteins from Lactobacillus reuteri in gut microbe–host interactions. *Proc.*
966       *Natl. Acad. Sci.* 115:E2706–E2715.

967 Shen X-X, Opulente DA, Kominek J, Zhou X, Steenwyk JL, Buh KV, Haase MAB, Wisecaver
968       JH, Wang M, Doering DT, et al. 2018. Tempo and Mode of Genome Evolution in the
969       Budding Yeast Subphylum. *Cell* [Internet]. Available from:
970       http://www.sciencedirect.com/science/article/pii/S0092867418313321

971 Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M,
972       Söding J, et al. 2011. Fast, scalable generation of high-quality protein multiple sequence
973       alignments using Clustal Omega. *Mol. Syst. Biol.* 7:539.

974 Singh S, Uppuluri P, Mamouei Z, Alqarihi A, Elhassan H, French S, Lockhart SR, Chiller T, Jr
975       JEE, Ibrahim AS. 2019. The NDV-3A vaccine protects mice from multidrug resistant
976       Candida auris infection. *PLOS Pathog.* 15:e1007460.

977 Smukalla S, Caldara M, Pochet N, Beauvais A, Guadagnini S, Yan C, Vinces MD, Jansen A,
978       Prevost MC, Latgé J-P, et al. 2008. FLO1 is a variable green beard gene that drives
979       biofilm-like cooperation in budding yeast. *Cell* 135:726–737.

980 Srivastava V, Singla RK, Dubey AK. 2018. Emerging virulence, drug resistance and future anti-
981       fungal drugs for Candida pathogens. *Curr. Top. Med. Chem.*

982 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
983       large phylogenies. *Bioinformatics* 30:1312–1313.

984 Stamler R, Kappé G, Boelens W, Slingsby C. 2005. Wrapping the α-Crystallin Domain Fold in a
985       Chaperone Assembly. *J. Mol. Biol.* 353:68–79.

986 Steentoft C, Vakhrushev SY, Joshi HJ, Kong Y, Vester-Christensen MB, Schjoldager KT-BG,
987       Lavrsen K, Dabelsteen S, Pedersen NB, Marcos-Silva L, et al. 2013. Precision mapping
988       of the human O-GalNAc glycoproteome through SimpleCell technology. *EMBO J.*
989       32:1478–1488.

990 Steenwyk JL, Buida TJ, Li Y, Shen X-X, Rokas A. 2020. ClipKIT: A multiple sequence alignment
991       trimming software for accurate phylogenomic inference. *PLoS Biol.* 18:e3001007.

992 Suyama M, Torrents D, Bork P. 2006. PAL2NAL: robust conversion of protein sequence
993       alignments into the corresponding codon alignments. *Nucleic Acids Res.* 34:W609-612.

994 Teufel F, Almagro Armenteros JJ, Johansen AR, Gíslason MH, Pihl SI, Tsirigos KD, Winther O,
995       Brunak S, von Heijne G, Nielsen H. 2022. SignalP 6.0 predicts all five types of signal
996       peptides using protein language models. *Nat. Biotechnol.* 40:1023–1025.

997 Teunissen AW, Steensma HY. 1995. Review: the dominant flocculation genes of
998       Saccharomyces cerevisiae constitute a new subtelomeric gene family. *Yeast Chichester*
999       *Engl.* 11:1001–1013.

1000 Uppuluri P, Lin L, Alqarihi A, Luo G, Youssef EG, Alkhazraji S, Yount NY, Ibrahim BA, Bolaris
1001     MA, Edwards JE, et al. 2018. The Hyr1 protein from the fungus Candida albicans is a
1002     cross kingdom immunotherapeutic target for Acinetobacter bacterial infection. *PLoS*
1003     *Pathog.* [Internet] 14. Available from:
1004     https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5963808/

1005 Verstrepen KJ, Jansen A, Lewitter F, Fink GR. 2005. Intragenic tandem repeats generate
1006     functional variability. *Nat. Genet.* 37:986–990.

1007 Verstrepen KJ, Reynolds TB, Fink GR. 2004. Origins of variation in the fungal cell surface. *Nat.*
1008     *Rev. Microbiol.* 2:533–540.

1009 Wang L-G, Lam TT-Y, Xu S, Dai Z, Zhou L, Feng T, Guo P, Dunn CW, Jones BR, Bradley T, et
1010     al. 2020. Treeio: An R Package for Phylogenetic Tree Input and Output with Richly
1011     Annotated and Associated Data. *Mol. Biol. Evol.* 37:599–603.

1012 Wasmer C, Lange A, Melckebeke HV, Siemer AB, Riek R, Meier BH. 2008. Amyloid Fibrils of
1013     the HET-s(218–289) Prion Form a β Solenoid with a Triangular Hydrophobic Core.
1014     *Science* [Internet]. Available from:
1015     https://www.science.org/doi/abs/10.1126/science.1151839

1016 Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview Version 2--a
1017     multiple sequence alignment editor and analysis workbench. *Bioinforma. Oxf. Engl.*
1018     25:1189–1191.

1019 Welsh RM, Sexton DJ, Forsberg K, Vallabhaneni S, Litvintseva A. 2019. Insights into the
1020     Unique Nature of the East Asian Clade of the Emerging Pathogenic Yeast Candida
1021     auris. *J. Clin. Microbiol.* 57:e00007-19.

1022 Wilkins M, Zhang N, Schmid J. 2018. Biological Roles of Protein-Coding Tandem Repeats in the
1023     Yeast Candida Albicans. *J. Fungi* 4:78.

1024 Willaert R. 2018. Adhesins of Yeasts: Protein Structure and Interactions. *J. Fungi* 4:119.

1025 Willbold D, Strodel B, Schröder GF, Hoyer W, Heise H. 2021. Amyloid-type Protein Aggregation
1026     and Prion-like Properties of Amyloids. *Chem. Rev.* 121:8285–8307.

1027 Winter DJ. 2017. rentrez: an R package for the NCBI eUtils API. *R J.* 9:520–526.

1028 Xie X, Qiu W-G, Lipke PN. 2011. Accelerated and adaptive evolution of yeast sexual adhesins.
1029     *Mol. Biol. Evol.* 28:3127–3137.

1030 Xu Z, Green B, Benoit N, Schatz M, Wheelan S, Cormack B. 2020. De novo genome assembly
1031     of Candida glabrata reveals cell wall protein complement and structure of dispersed
1032     tandem repeat arrays. *Mol. Microbiol.*

1033 Xu Z, Green B, Benoit N, Sobel JD, Schatz MC, Wheelan S, Cormack BP. 2021. Cell wall
1034     protein variation, break-induced replication, and subtelomere dynamics in Candida
1035     glabrata. *Mol. Microbiol.*

1036    Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y. 2015. The I-TASSER Suite: protein structure
1037        and function prediction. *Nat. Methods* 12:7–8.

1038    Yang Z. 1998. Likelihood ratio tests for detecting positive selection and application to primate
1039        lysozyme evolution. *Mol. Biol. Evol.* 15:568–573.

1040    Yang Z. 2007. PAML 4: Phylogenetic Analysis by Maximum Likelihood. *Mol. Biol. Evol.*
1041        24:1586–1591.

1042    Yu G. 2020. Using ggtree to Visualize Data on Tree-Like Structures. *Curr. Protoc. Bioinforma.*
1043        69:e96.

1044    Zhang J. 2003. Evolution by gene duplication: an update. *Trends Ecol. Evol.* 18:292–298.

1045    Zhang J, Nielsen R, Yang Z. 2005. Evaluation of an Improved Branch-Site Likelihood Method
1046        for Detecting Positive Selection at the Molecular Level. *Mol. Biol. Evol.* 22:2472–2479.

1047    Zhao X, Oh S-H, Coleman DA, Hoyer LL. 2011. ALS51, a newly discovered gene in the
1048        Candida albicans ALS family, created by intergenic recombination: analysis of the gene
1049        and protein, and implications for evolution of microbial gene families. *FEMS Immunol.*
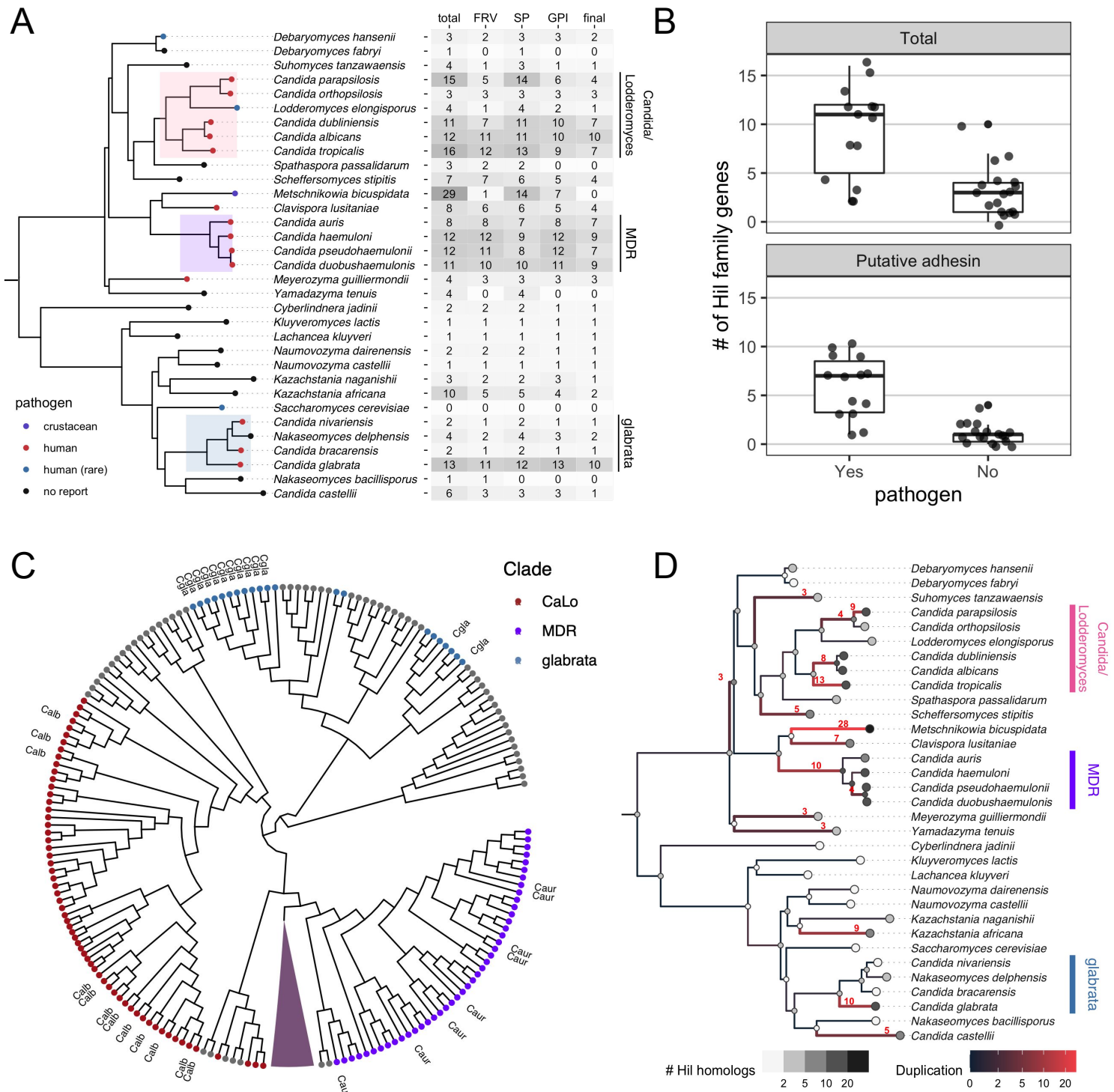1050        *Med. Microbiol.* 61:245–257.

1051

**Figure 1. Phylogenetic distribution of the yeast Hil family and its parallel expansion in independently derived pathogenic *Candida* species.** Legend on the next page.

**Figure 1. Phylogenetic distribution of the yeast Hil family and its parallel expansion in independently derived pathogenic *Candida* species.** A) Species tree is based on the phylogeny for 332 yeast species from (Shen *et al.* 2018), except for three species in the MDR clade other than *C. auris*, whose phylogenetic relationships are based on (Muñoz *et al.* 2018). The tip colors show the pathogenic status of the species. The highlighted clades are enriched in known human pathogens. In the table, the first column shows the total number of Hil family homologs per species. The number of homologs that pass each of the three tests for determining their adhesin status are shown in the next three columns. FRV = FungalRV, SP = Signal Peptide and GPI = GPI-anchor. See Materials and Methods for details. The number of homologs passing all three tests is shown in the "final" columns. (B) Boxplots comparing the number of Hil homologs (upper) or the number of putative adhesins passing all three tests (lower) per species between known pathogens and low pathogenic-potential species. Individual species numbers are shown as dots on top of the boxplot. Homologs from *M. bicuspidata* were excluded (see text). Both comparisons are significant at a 0.005 level by either a t-test with unequal variance or Mann-Whitney U test. (C) Maximum likelihood tree based on the Hyphal_reg_CWP domain of the Hil family was constructed using RAxML-NG and corrected with GeneRax based on the species tree. The tree is shown as a cladogram. All 29 homologs in *M. bicuspidata* formed a single group, which is shown as a triangle in dark plum. Homologs from the species in the three highlighted clades in (A) are colored accordingly. CaLo = Candida/Lodderomyces. Homologs from *C. albicans*, *C. auris* and *C. glabrata* are labeled as Calb, Caur and Cgla, respectively. (D) Species tree showing the number of inferred duplication events on each branch. The gray colors of the tip and internal nodes represent the identified and inferred number of Hil homologs, respectively. The branch color shows the inferred number of duplication events, with 3 or more duplications also shown as a number next to the branch.
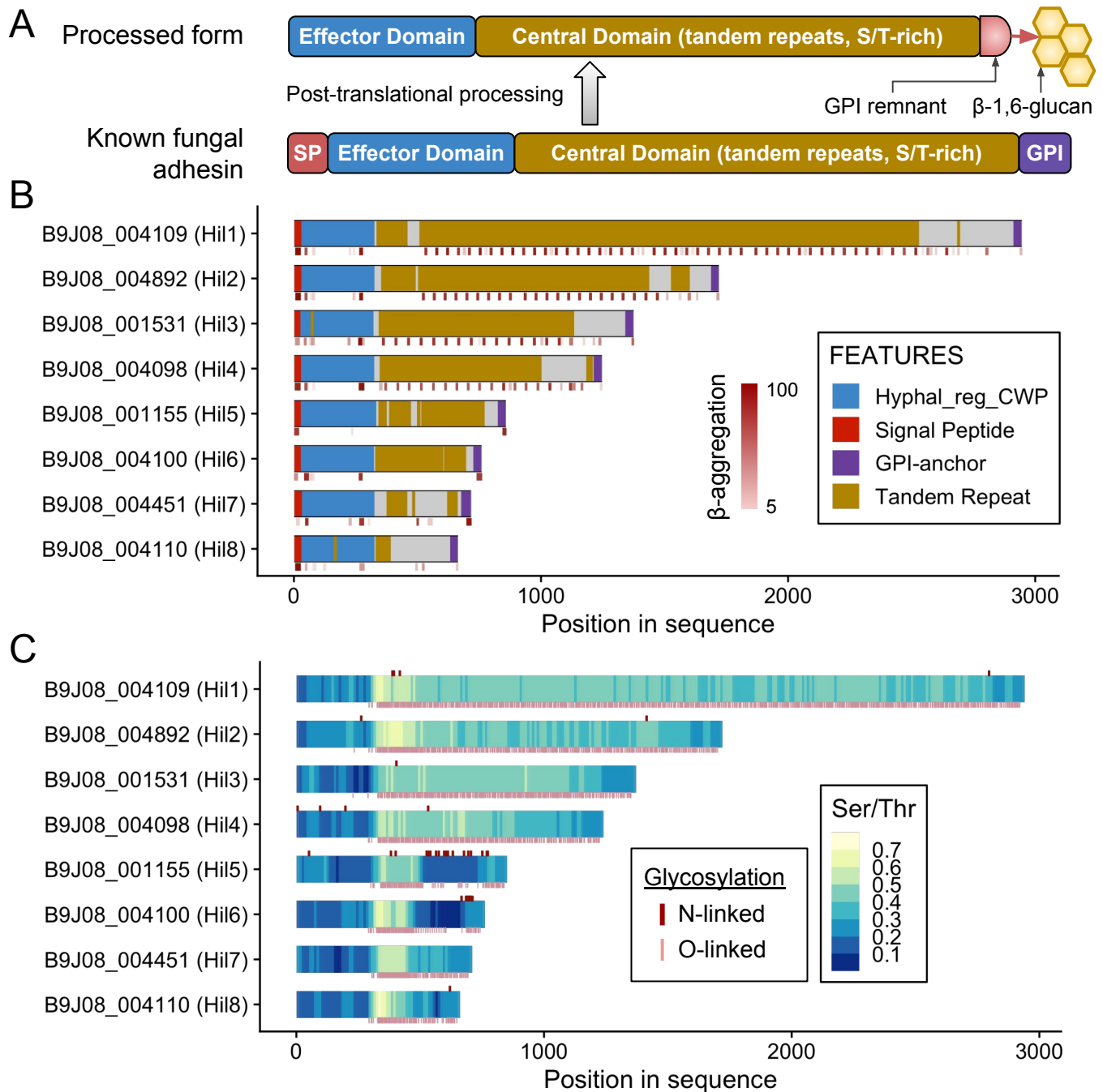
**Figure 2. Domain architecture and adhesin-associated features of the *C. auris* Hil family.** (A) Diagram depicting the domain organization of a typical yeast adhesin before and after the post-translational processing, adapted from (de Groot *et al.* 2013). (B) Domain features of the eight Hil proteins in *C. auris* (strain B8441). Gene IDs and names are labeled on the left. The short stripes below each diagram are the TANGO predicted β-aggregation prone sequences, with the intensity of the color corresponding to the score of the prediction. (C) Serine and Threonine (Ser/Thr) frequencies in each protein are plotted in 50 aa sliding windows with step size of 10 aa. N-linked and O-linked glycosylation sites were predicted by NetNGlyc 1.0 and NetOGlyc 4.0, respectively, and are shown as short ticks above and below each protein schematic.
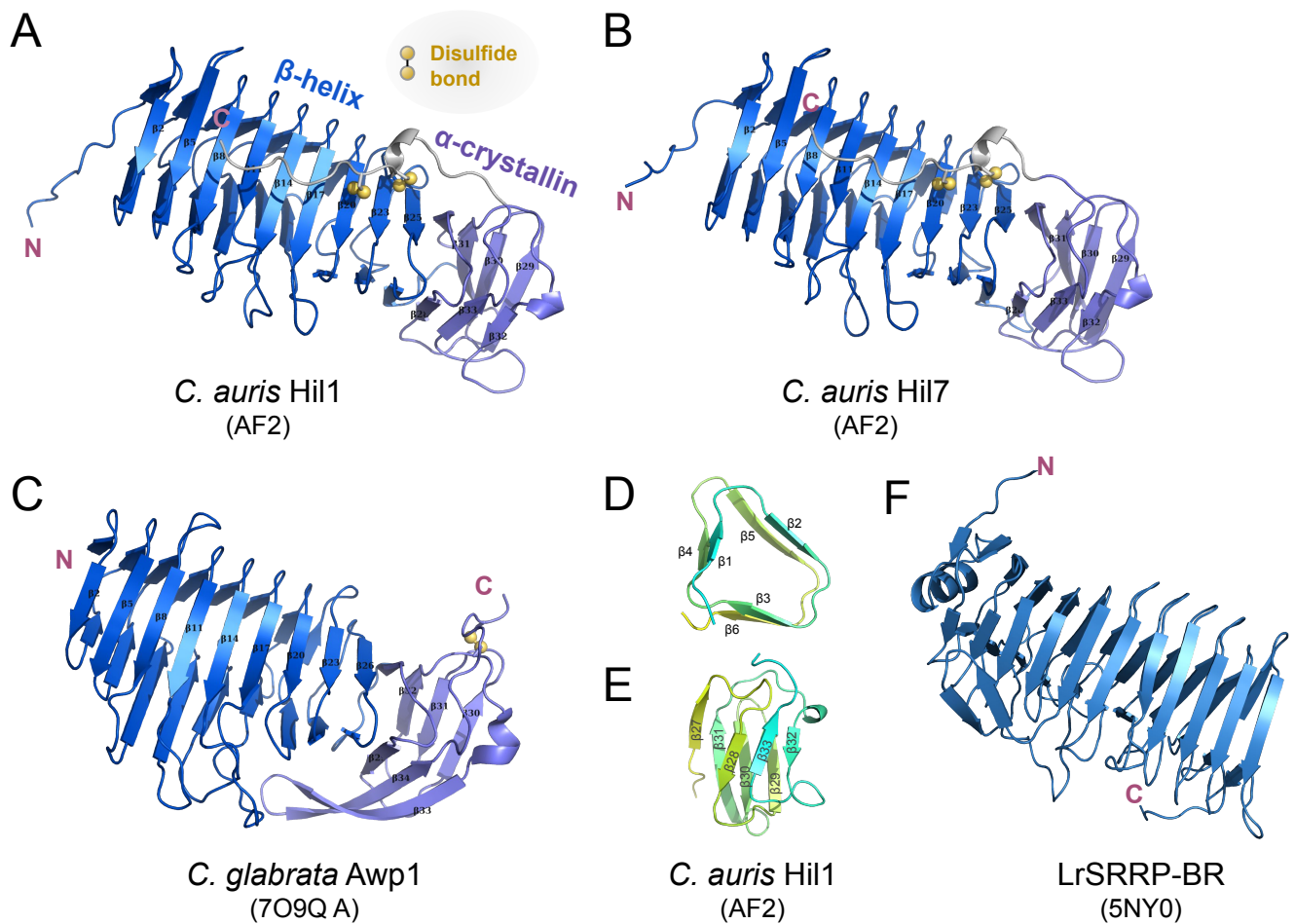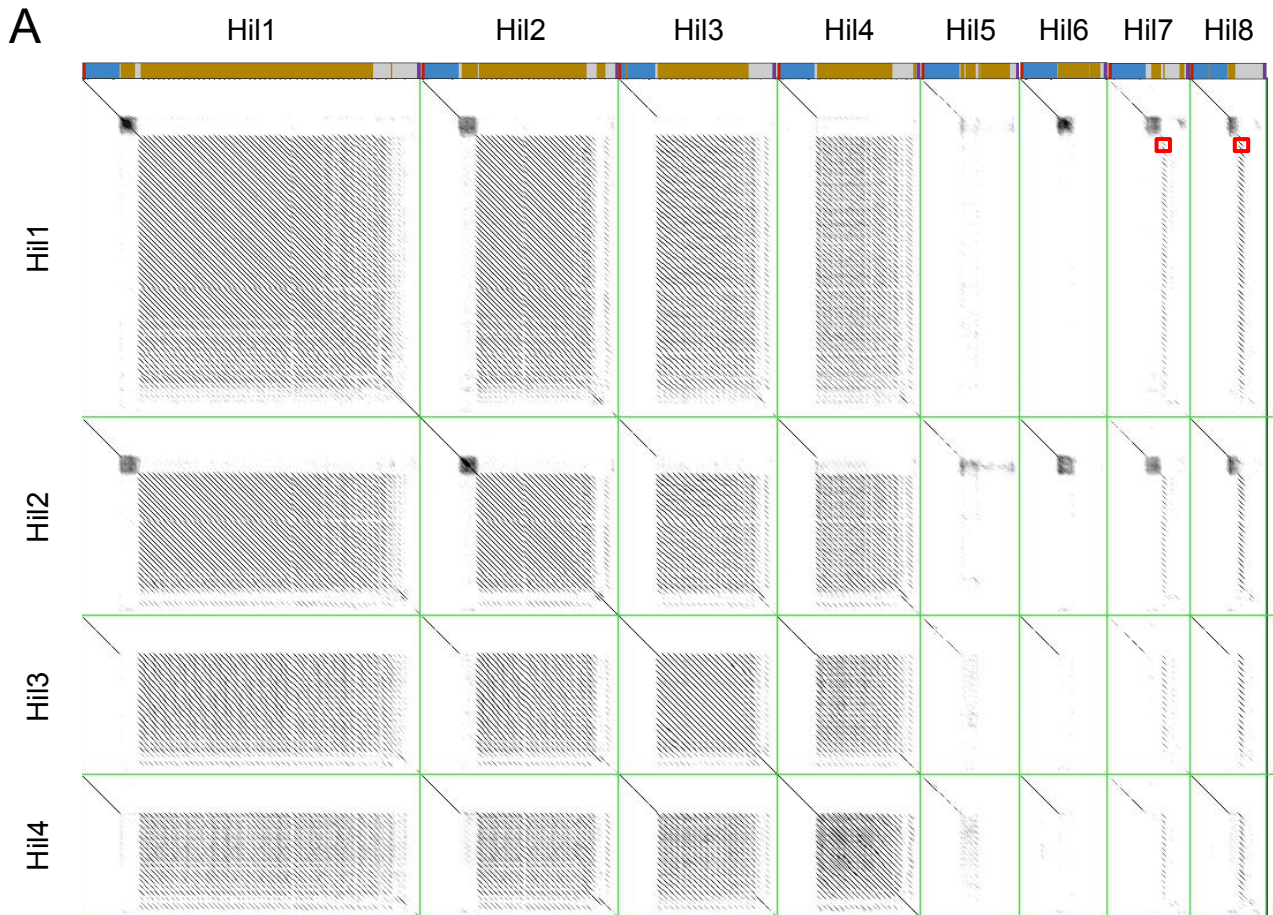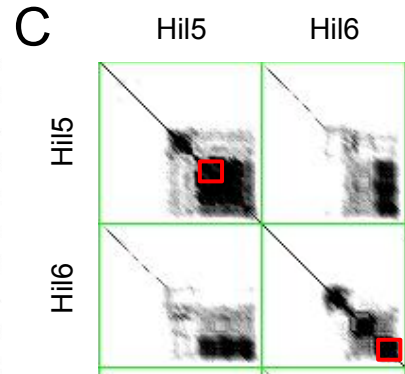
**Figure 3. Predicted structures of the Hyphal_reg_CWP domain in two *C. auris* Hil proteins are similar to yeast and bacterial adhesins.** (A) and (B) AlphaFold2 (AF2) predicted structures of the Hyphal_reg_CWP domains from *C. auris* Hil1 and Hil7, which consist of a β-helix followed by a α-crystallin domain, with the C-terminal loop linked to the β-helix via two disulfide bonds. (C) Crystal structure of the *C. glabrata* Awp1 effector domain, which is highly similar to *C. auris* Hil1 and Hil7, but with the disulfide bond in a different location. (D) cross-section of the first two rungs of the β-helix in (A), showing the three β-strands per rung. The cyan-to-yellow gradient follows the N- to C-terminus. (E) α-crystallin domain in (A), showing the seven β-strands forming two antiparallel β-sheets. Color is the same as in (D). (F) crystal structure of the Serine Rich Repeat Protein Binding Region (SRRP-BR) from the gram-positive bacterium *L. reuteri*, which adopts a β-helix fold.

**Figure 4. Dotplot shows the tandem repeat structure within and similarity between *C. auris* Hil proteins.** (A) Hil1-Hil4 are compared to all eight Hil proteins in *C. auris* including themselves in dotplots (JDotter, Brodie et al 2004) with a sliding window of 50 aa and Grey Map set to 60-245 (min-max). A schematic for each protein is shown above each column (colors same as in Fig. 2). The regions highlighted by the red boxes in row 1 reveal the presence of a single copy of the 44 aa repeat unit in Hil7 and Hil8. (B) Wrapped sequence of aa 543-982 from Hil1 showing the conserved period and sequence of the 44 aa tandem repeat. The magenta and plum fonts indicate motifs predicted by TANGO to have strong (probability > 90%) or moderate (30-90%) β-aggregation potential. The yellow highlighted regions are predicted to form β-strands by PSIPred, with cartoons shown above. (C) Dotplots between Hil5 and Hil6 with the same settings as in (A), showing the low complexity repeats unique to these two. Regions within the two red boxes are shown in (D), with residue numbers shown on both ends. The rectangles delineate individual repeats, with the number of copies for each repeat shown to the right.

**Figure 5. Maximum-likelihood-based analyses for selective pressure variation and role of positive selection on the Hyphal_reg_CWP domain in *C. auris*.** (A) Schematic showing the putative non-recombining partitions within the Hyphal_reg_CWP domain determined by GARD (see Fig. S8). The two partitions labeled in gray were studied separately. The numbers refer to the columns in the coding sequence alignment. (B, C) Phylogenetic trees were reconstructed for the two partitions and are shown as a cladogram. The vertical bar next to the Hil1/Hil2 pair indicates the difference in topology between the two trees. Branch colors are based on the dN/dS values estimated from a free-ratio model in PAML. "FG" designate foreground branches, whose dN/dS were greater than 10, except for branch 14..16 in (C), which was selected instead of 16..Hil4 because this would require fewer evolutionary changes in selective forces. We also analyzed the scenario with 16..Hil4 as the foreground and the conclusions remained the same with slightly different *P*-values. (D) Summary of the maximum-likelihood-based tests for selective force heterogeneity and for positive selection. "Insig." means *P*-value > 0.05. In the branch-site test, $p(\omega>1)$ is the total proportion of sites with dN/dS > 1 on the FG branches, and $\omega2$ their estimated dN/dS. The listed sites were identified as being under positive selection with a posterior probability greater than 0.99 by the Bayes Empirical Bayes (BEB). The one-letter code and number refer to the amino acid in the OG sequence and the alignment column (Fig S7).

**Figure 6. Evolution of protein length and β-aggregation potential in the yeast Hil family.** (A) Domain schematic shows that most homologs have a signal peptide at the N-terminus, then the Hyphal_reg_CWP domain and a highly repetitive region central domain, followed by the C-terminal GPI-anchor peptide. Homologs from *M. bicuspidata* were not included because many were annotated as incomplete. They were also excluded from other results in this figure. (B) Distribution of TANGO predicted β-aggregation sequences. The score for each sequence is shown as a color gradient and represents the median of the per-residue probability of aggregation. A vertical bar marks a group of MDR clade sequences that have a large number of β-aggregation prone sequences arranged in regular intervals. (C) X-Y plot showing the relationship between total protein length and tandem repeat sequence length for Hil family homologs. The linear regression line is shown in blue, with coefficients and r² values below. (D) The species tree on the left is the same as in Figure 1. The middle panel shows the number of Hil homologs per species. *M. bicuspidata* homologs were excluded; *S. cerevisiae* was included in the species tree but no Hil homolog was identified in it (see text). The right panel shows the number of predicted β-aggregation prone motifs per Hil homolog. Only motifs with a median probability >= 30% were counted. Proteins are colored in gold if they have five or more such motifs and if the Median Absolute Deviation (MAD) of the inter-motif distances is < 5 aa.

**Figure 7. Hil family genes are preferentially located near the chromosome ends.** (A) Schematic of the analysis: each chromosome (chr) is folded in half and divided into five equal-length "bins", ordered by their distance to the nearest telomere. The cumulative bar graph on the right summarizes the gene density distribution in the five bins. (B) Folded gene density distribution for six species with a chromosomal level assembly and more than two Hil family genes. The bin colors are as shown in (A). The Hil homologs in each species are plotted as a separate group. A goodness-of-fit test comparing the distribution of the Hil family genes to the genome background yielded a *P*-value of $1.3 \times 10^{-12}$.

A

B

| length | sp-gpi- | sp-GPI+ | SP+gpi- | SP+GPI+ | total |
|--------|---------|---------|---------|---------|-------|
| >600 | 3 | 14 | 9 | 112 | 138 |
| 251-600 | 12 | 5 | 27 | 4 | 48 |
| 0-250 | 3 | 0 | 3 | 1 | 7 |

**Supplementary Figure 1. Hil family proteins' length distribution and grouping by signal peptide (SP) and GPI-anchor signal presence.** (A) Histogram showing the distribution of protein lengths for Hil family proteins from 32 yeast species. Top: protein sequence records were labeled as complete or "NA"; bottom: proteins labeled as incomplete (no-right, no-left, no-ends). Most of the short sequences (<600 aa, dashed vertical line) came from the species *M. bicuspidata* (red) (B) Summary of the number of Hil family proteins predicted to have a signal peptide (SP+) and GPI-anchor signal (GPI+), grouped by protein length. Proteins labeled as incomplete were excluded from this table.

**Supplementary Figure 2. Maximum likelihood tree for the Hil family genes.** This tree is identical to the one shown in Fig 1C, except that it is shown in a rectangular format with the sequence names in the form of refseqID_species_name. We identified and named the *C. auris* homologs as Hil1-8 to be consistent with the latter figures. Note that all 29 *M. bicuspidata* homologs form a single clade, which is collapsed for the ease of viewing. Their sequence IDs can be found in Table S1.

**Supplementary Figure 3. Comparison of the Ser/Thr frequencies in *C. auris* Hil family members with all protein-coding genes in *C. auris*.** B8441 strain genome is used for this analysis. The frequency of Ser or Thr residues as a percent of the entire protein length is plotted as a histogram for all protein-coding genes. Red ticks indicate the eight Hil genes. A student's t-test was used to assess the significance of the difference in Ser/Thr frequencies between the Hil family proteins vs the rest of the proteome.

```
 1  DSGVVIVTTDSDGSLTTTTSVIPPPFTTYTSSWVTTNSAGETET
 2  DSGVVVVTTNSEGELTTSTSVIPPPYTTYTSTWTTTDGNGDVET
 3  DSGVVIVTTGSDGSLTTTTSVIPPPFTTFTSTWTTTNTDGETET
 4  DSGIVVVTTDSNGQLTTSTSIIPPPFTTYTSTWTSSQSDGSEVT
 5  DSGVVIVTTDSDGSLTTTTSVIPPPFTTYTSTWATTNSNGETET
 6  NSGVVVVTTGSDGELTTTTSTIPPPFTTYTSTWISTNSNGATET
 7  DSGVVVVTTDSDGALTTTTSIIPQPFTTYTSTWTSTNSDGDTET
 8  DSGVVVVTTNSDGALTTSTSVIPQPFTTFTSTWTSSNSNGAVQT
 9  DSGVVIVTTGSDGSLTTTTSVIPPPFTTYTSTWTSSNSDGETET
10  DSGVVVVTTDSNGELTTTTSIIPPPFTTFTSTWTSTKSDGAVET
11  DSGVIIVTTNSEGDLTTTTSIIPPPYTTYTSTWTTTDSNGVTET
12  DSGVVVVTTDSDGQLTTATSIIPPPFTTYTSTWTTTNSDGSEET
13  DSGVVIVTAGTDGSLTTTTSVIPPPFTTYTSTWITTNSNGAVET
14  DSGIIVVTTNSGGSLTTSTSVLPTPFTTYTSTWTTSDGDGNVQT
15  DSGVVIVTTGSDGALSTTTSVIPPPFTTYTSTWISTNSDGETET
16  DSGVVVVTTDSNGALTTTTSIIPPPFTTFTSTWTTTDENGATET
17  DSGVVVVTTGTDGSLTTTTSVIPPPYTTFTSTWTTSNSNGDIET
18  DSGVVIVTTNSDGSLTTTTSVIPPPYTTFTTTWATTNSDGTTET
19  DSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTSNKSDGAVET
20  DSGVVIVTTDSDGALTTTTSIIPQPFTTYTSTWTSTNSNGAIET
21  ESGVVVVTTDSNGALTTSTSVIPLPLTTFTTTWTTTNSAGETET
22  DSGVVVVETNSNGALTTTTSTFPEPFTTFTSTWTTTDDSGAIAT
23  DSGVVIVTTGSDGSLSTTTSVIPPPFTTYTTTWTSTNSNGGVET
24  DSGVVIVTTNSDGALETTTSVIDPPFNTYTSTWTTTDADGAIET
25  DSGVVVVTTGSDGSLTTTTSVIPHPFTTYTSTWVTTGSDGDTET
26  DSGVIVVTTDSDGALTTSTSLLPVPFTTYTSTWTITNSDGSQAT
27  DSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTGANGGEET
28  DSGVIIVTTDSDGQLATTTSVIPPPFTTFTSTWTTTNSDGNQAT
29  DSGVVIVTTDSDGQLTTTTSVIPPPFTTYTSTWTTTDGNGAEET
30  DSGVIIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTGADGSEET
31  DSGVIIVTTDSAGQLTTTTSVIPPPFTTFTSTWTTTDGNGNEGT
32  DSGVIIVTTDSDGALTTTTAVIPPAAGSGTDALSSSINDVPYTTYTSTWTTTDGNGNIET
33  DSGVVIVTTDSQGSLTTTTSIIDSPFTTYTSTWATTDNNGNVET
34  DSGVVIVTTDSNGQLTTTTSVIDSPYTTYTTSWPTTDANGGVET
35  DSGVVIVTTDSDGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVET
36  DSGVVIVTTDSNGQLTTTTSVIDSPYTTYTTSWPTTGADGAVET
37  NSGVVIVTTDSDGQLTTTTSVIDSPYTTYTSIWTTTDSVGNVET
38  DSGVVIVTTDSDGQVTTTTSRFENSPSDLTEYTTTWASTDSDGNIKT
39  DSGVVVVTTDSAGSTTTSTSTFDTPYTTFTSTWTTTNGNGDVKT
40  DSGVIIVTTDSVGQLTTTTSQFDSQQSGLTDYTTTWTTTDRNGNPST
41  ASGVVVVTTDSDGQITSTTSQFSDKSSGLTDYTTTWTTTDTDGSVVT
42  DSGVVIVTTDSAGSLTTTTSVFDTPITTFTSTWTTTNADGSIET
43  DSGVIIVTADSNGQLTTTTSQSDNRPSGLTDYTTTWTTTNTDGAIET
44  DSGVVVVTTDSQGQLTTITSVIESPVTASSGSSDKPSGITEFTTTWTTTDANGIAHT
45  DSGVVIVTTDADGSLTTTTSQIDNVSSGLTEFTSSWTTTLSDGSVET
46  DSGLVIVTTDSNGQLKTTTSQFEDIPSGLSEFTTSWTTTDADGDTRI
47  DSGVVIVTTDSDNRLTTTTSQFASVDPTDFTSYITSWTATNGDGSIET
48  DSGAVIVTTNSDGQLVTTTSVISSSHGAVSTSES
49  DS-NVIVTTDSEGSLTTSTVTLCPQCTHFTSTWTTSNSEGAIET
50  DSGVVVVTTDSVGSLTTYTKDCPEASGELSTFISTYTTTDTDGNIKTT
```

**Supplementary Figure 4. Tandem repeats in the *C. auris* Hil1 central domain.** The majority of the 50 tandem repeat copies have a conserved 44 aa period. Dark and light orange highlights show sequences predicted by TANGO to have strong (>90%) or moderate (30-90%) β-aggregation potentials.
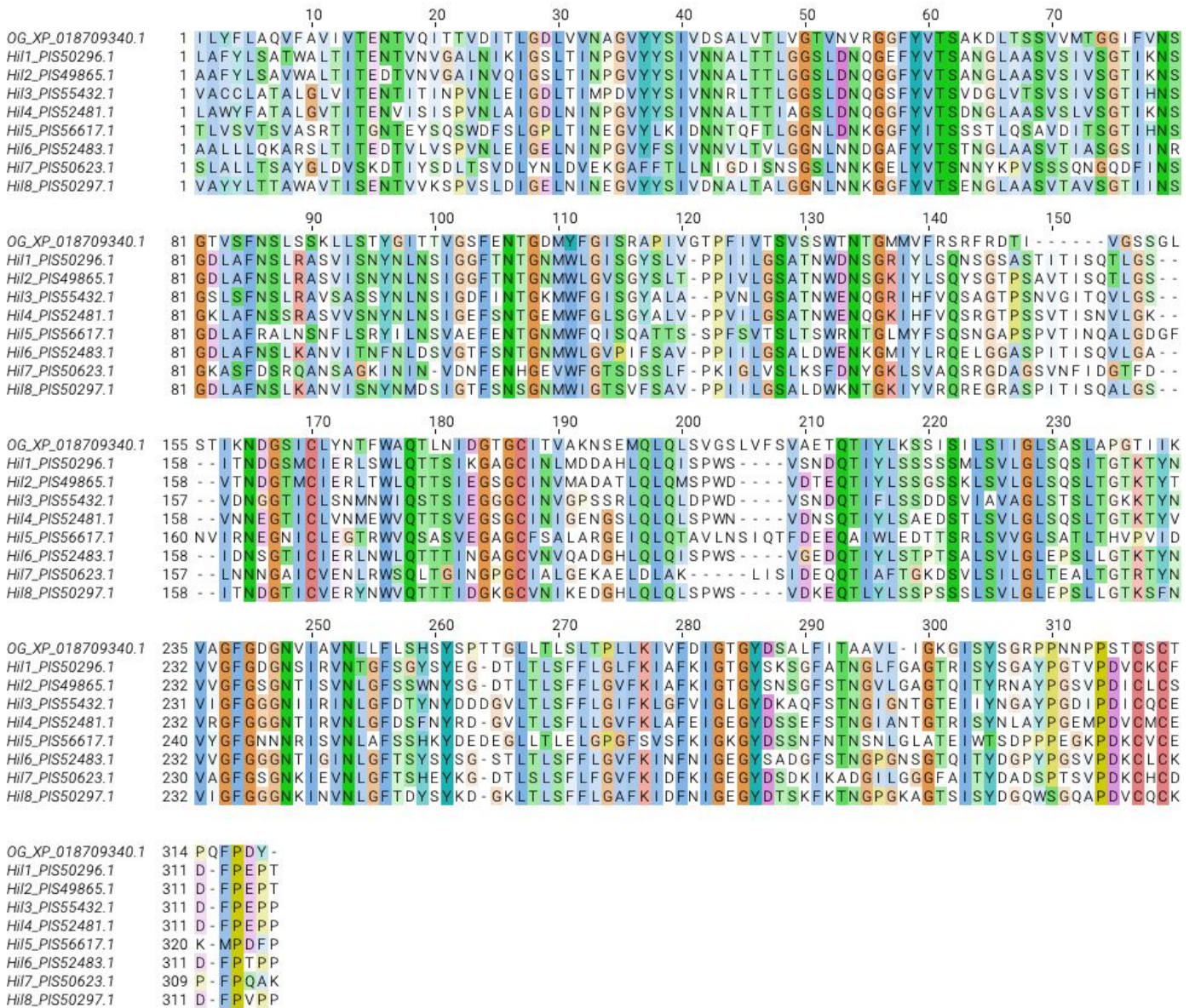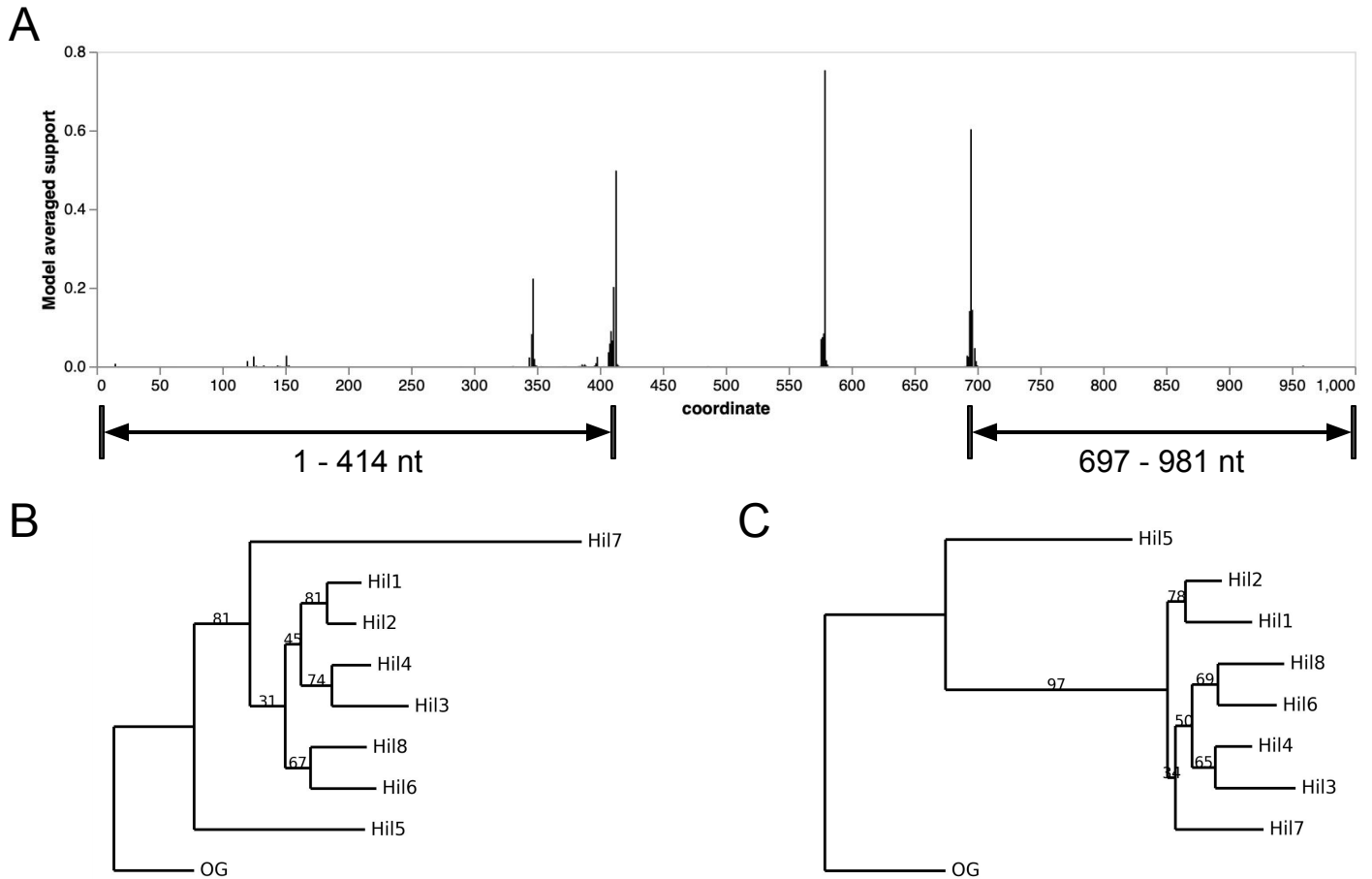
**A**

| B8441_I | 991 | PPPYTTY- - - - - - - - - - - - 44 aa = 1 repeat - - - - - - - - TSTWTTTNSDGSEETDSGVVIVTTGTDGSLTTTTS | 1032 |
| B11205_I | 980 | PPPYTTY- - - - - - - - - - - - TSTWTTTNSDGSEETDSGVVIVTTGTDGSLTTTTS | 1021 |
| B13916_I | 980 | PPPYTTY- - - - - - - - - - - - TSTWTTTNSDGSEETDSGVVIVTTGTDGSLTTTTS | 1021 |
| B11221_III | 1009 | PPPYTTYTSTWTTTDSNGVTETDSGVVVVTTDSDGQLTTATSIIPPPFTTYTSTWTTTNSDGSEETDSGVVIVTAGTDGSLTTTTS | 1094 |
| B17721_III | 1009 | PPPYTTYTSTWTTTDSNGVTETDSGVVVVTTDSDGQLTTATSIIPPPFTTYTSTWTTTNSDGSEETDSGVVIVTAGTDGSLTTTTS | 1094 |
| B12037_III | 1033 | PPPYTTYTSTWTTTDSNGVTETDSGVVVVTTDSDGQLTTATSIIPPPFTTYTSTWTTTNSDGSEETDSGVVIVTAGTDGSLTTTTS | 1118 |
| B12342_IV | 967 | PPPYTTYTSTWTTIDSNGVTETDSGVVVVTTDSDGQLTTTTSIIPPPFTTYTSTWTTTNSDGSEETDSGVVIVTTGTDGSLTTTTS | 1052 |

**B**

| B8441_I | 603 | SQATDSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTNSDGSSETDSGVIIVTTDSEGQLTTTTSMIPPPFTTYTSTWTTTDA | 688 |
| B11205_I | 603 | SQATDSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTNSDGSSETDSGVIIVTTDSEGQLTTTTSMIPPPFTTYTSTWTTTDA | 688 |
| B13916_I | 603 | SQATDSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTNSDGSSETDSGVIIVTTDSEGQLTTTTSMIPPPFTTYTSTWTTTDA | 688 |
| B11221_III | 603 | SQATDSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTNSDGSSETDSGVIIVTTDSEGQLTTTTSMIPPPFTTYTSTWTTTDA | 688 |
| B17721_III | 603 | SQATDSGVVIVTTDSEGQLTTTTSVIPPPFTTYTSTWTTTNSDGSSETDSGVIIVTTDSEGQLTTTTSMIPPPFTTYTSTWTTTDA | 688 |
| B12342_IV | 592 | SQATDSGVVIVTTDSE- - - - - - - - - - 44 aa = 1 repeat - - - - - - - - - - GQLTTTTSVIPPPFTTYTSTWTTTDA | 633 |
| B11245_IV | 592 | SQATDSGVVIVTTDSE- - - - - - - - - - GQLTTTTSVIPPPFTTYTSTWTTTDA | 633 |

**C**

| B8441_I | 775 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGADGSEETDSGVIIVTTDSAGQLTTTTSVISPPFTTFTSTWT | 860 |
| B11205_I | 775 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGADGSEETDSGVIIVTTDSAGQLTTTTSVISPPFTTFTSTWT | 860 |
| B13916_I | 775 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGADGSEETDSGVIIVTTDSAGQLTTTTSVISPPFTTFTSTWT | 860 |
| B11221_III | 775 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGAD- - - - - - - - - - - - - - - - AGQLTTTTSVISPPFTTFTSTWT | 844 |
| B17721_III | 775 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGADGSEETDSGVIIVTTDSAGQLTTTTSVISPPFTTFTSTWT | 860 |
| B12342_IV | 720 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGADGSEETDSGVIIVTTDSAGQLTTTTSVIPPPFTTFTSTWT | 805 |
| B11245_IV | 720 | DGNGAEETDSGVIIVTTDSEGQLTTTTSIIDPPFTTYTSTWTTTGADGSEETDSGVIIVTTDSAGQLTTTTSVIPPPFTTFTSTWT | 805 |

**D**

| B8441_I | 995 | VIDSPYTTYTTSWPTTDANGGVETDSGVVIVTTDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPYTTYTTSWPTTDANGNVETDSGVVIVT | 1114 |
| B11205_I | 995 | VIDSPYTTYTTSWPTTDANGGVETDSGVVIVTTDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPYTTYTTSWPTTDANGNVETDSGVVIVT | 1114 |
| B13916_I | 979 | VIDSPYTTYTTSWPTTDANGGVETDSGVVIVTTDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPYTTYTTSWPTTDANGNVETDSGVVIVT | 1098 |
| B11221_III | 995 | VIDSPYTTYTTSWPTTDANGGVETDSGVVIVTTDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPYTTYTTSWPTTDANGNVETDSGVVIVT | 1114 |
| B17721_III | 940 | VIDSPY- - - - - - - - - - - - - - - - - - - - - 220 aa = 5 repeats - - - - - - - - - - - - - - - - - - - - - - | 945 |
| B12342_IV | 1115 | TDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPFTTYTTSWPTTDANGNVETDSGVVIVTTDSNGQLSTTTSVIDSPYTTYTTSWPTTDANG | 1234 |
| B11245_IV | 1115 | TDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPFTTYTTSWPTTDANGNVETDSGVVIVTTDSNGQLSTTTSVIDSPYTTYTTSWPTTDANG | 1234 |
| | 1115 | TDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPFTTYTTSWPTTDANGNVETDSGVVIVTTDSNGQLSTTTSVIDSPYTTYTTSWPTTDANG | 1234 |
| | 1099 | TDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPFTTYTTSWPTTDANGNVETDSGVVIVTTDSNGQLSTTTSVIDSPYTTYTTSWPTTDANG | 1218 |
| | 1115 | TDSNGQLSTTTSVIDSPFTTYTTSWPTTDGNGAVETDSGVVIVTTDSNGQLTTTTSVIDSPFTTYTTSWPTTDANGNVETDSGVVIVTTDSNGQLSTTTSVIDSPYTTYTTSWPTTDANG | 1234 |
| | 946 | - - - - - - - - - - - - - - - - - - continue 220 aa = 5 repeats - - - - - - - - - TTYTTSWPTTDANG | 959 |
| | 946 | - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - TTYTTSWPTTDANG | 959 |

**Supplementary Figure 5. Examples of tandem repeat copy number variation in Hil1-Hil4 among the *C. auris* strains.** (A) A 44 aa indel in Hil1 removes exactly one repeat in all three Clade I strain orthologs. (B) A similar indel polymorphism of exactly one repeat length in Hil2 affecting the Clade IV strains. (C) An indel polymorphism in Hil2 that affects one Clade III strain and spans 16 aa, not a full repeat, but includes a predicted strong β-aggregation prone sequence "GVIIVTT". (D) An indel polymorphism in Hil2 that spans 220 aa or five full repeats affecting the Clade IV strains. Similar patterns were observed in Hil3 and Hil4.
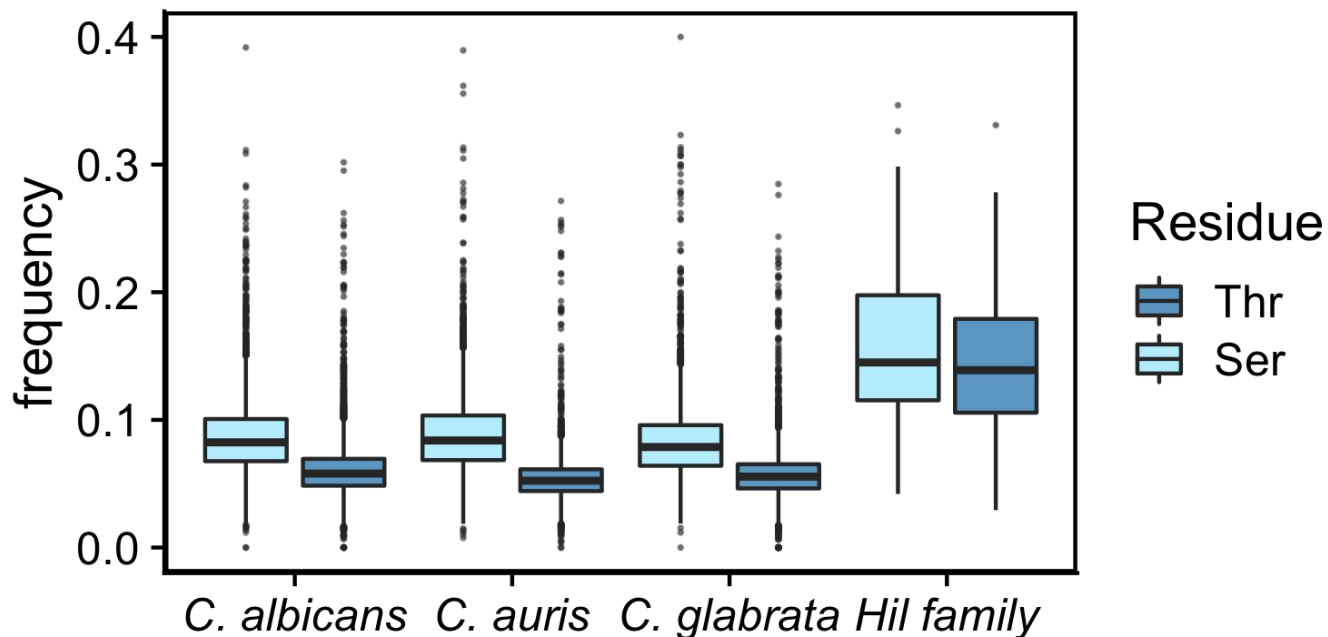
**Supplementary Figure 6. Reconciled Hil family gene tree based on the Hyphal_reg_CWP domain alignment in the four clades of *C. auris* strains and two closely related species.** The tree is rooted by the two homologs from the outgroup *D. hansenii*. The gene tree was corrected with the species/strain tree based on (Muñoz et al 2018) using GeneRax (v2.0.4). Hil genes lost in *C. auris* Clade II strains are labeled with an asterisk next to the Hil1-8 group labels.
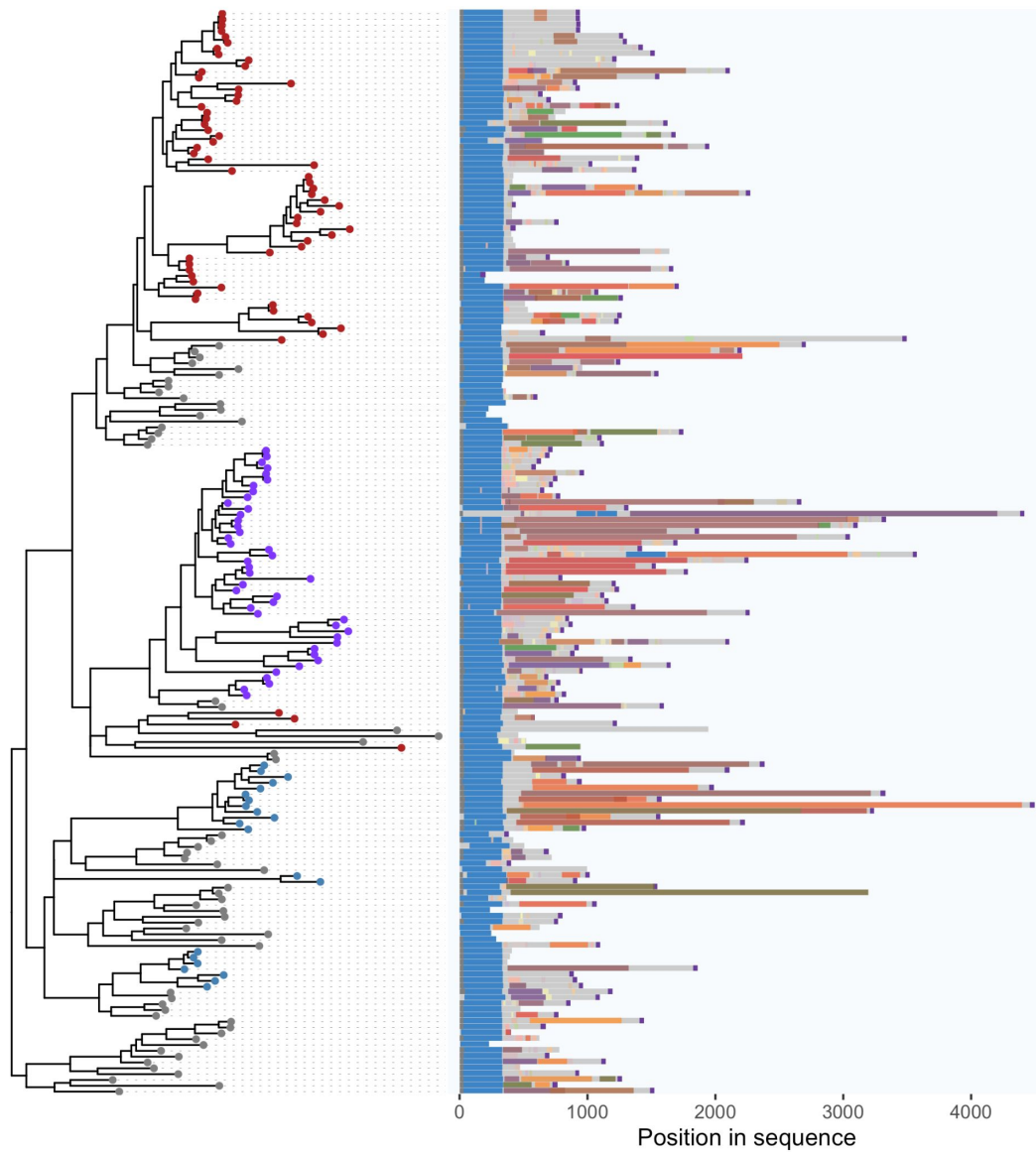
**Supplementary Figure 7. Multiple sequence alignment of *C. auris* Hil1-8 Hyphal_reg_CWP domain.** The domain regions were identified using HMMScan against the Pfam-A database. The sequences were aligned with clustalo v1.2.4 and the result visualized in Jalview v2.11.1.4 with the ClustalW color scheme. OG = outgroup from *M. bicuspidata*

**Supplementary Figure 8. Detecting intra-domain recombination and identifying non-recombining partitions in the Hyphal_reg_CWP domain using GARD.** (A) Model averaged support for breakpoint locations along the Hyphal_reg_CWP domain alignment for the eight Hil proteins in *C. auris* and an outgroup sequence from *M. bicuspidata* (protein ID: XP_018709340.1) to root the gene tree. Based on the GARD output, we chose the N- and C-terminal partitions for downstream analyses, i.e., coordinates 1-414 nt and 697-981 nt. (B) A maximum likelihood tree for partition 1-414 was constructed using RAxML-NG v1.1.0. Branch length is proportional to the amount of sequence divergence. OG stands for outgroup. Bootstrap support for internal splits are shown as a percentage and are based on 1000 replicates or until bootstrapping converges. (C) tree for 697-981nt, same format as in (B)

**Supplementary Figure 9. Yeast Hil family proteins have on average higher Ser/Thr frequencies than the rest of the proteome.** Proteome-wide distribution of Thr/Ser frequencies per protein from three species, compared with the yeast Hil family proteins (*M. bicuspidata* homologs were excluded because a large number of them were incomplete). The boxes represent the interquartile range (IQR), the middle thick line the median, the whiskers the 1.5 x IQR and the dots outliers outside that range.

**Supplementary Figure 10. Domain schematic for the Yeast Hil family showing rapidly evolving tandem repeat sequences in the central domain of the proteins.** Same as Fig. 6A except that tandem repeats belonging to different sequence clusters as determined by XSTREAM (Newman and Cooper 2007) are shown in different colors.