



HAL
open science

CeGAL: Redefining a Widespread Fungal-Specific Transcription Factor Family Using an In Silico Error-Tracking Approach

Claudine Mayer, Arthur Vogt, Tuba Uslu, Nicolas Scalzitti, Kirsley Chennen, Olivier Poch, Julie D. Thompson

► **To cite this version:**

Claudine Mayer, Arthur Vogt, Tuba Uslu, Nicolas Scalzitti, Kirsley Chennen, et al.. CeGAL: Redefining a Widespread Fungal-Specific Transcription Factor Family Using an In Silico Error-Tracking Approach. *Journal of Fungi*, 2023, 9 (4), pp.424. 10.3390/jof9040424 . hal-04253824

HAL Id: hal-04253824

<https://cnrs.hal.science/hal-04253824>

Submitted on 23 Oct 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Article

2 **CeGAL: redefining a widespread fungal-specific transcription**
 3 **factor family using an *in silico* error-tracking approach**

4 Claudine Mayer ^{1,2,*}, Arthur Vogt ¹, Tuba Uslu ¹, Nicolas Scalzitti ¹, Kirsley Chennen ¹, Olivier Poch ¹ and Julie D.
 5 Thompson ^{1,*}

6 ¹ Complex Systems and Translational Bioinformatics (CSTB), ICube laboratory, UMR7357, University of
 7 Strasbourg, 1 rue Eugène Boeckel, 67000, Strasbourg, France; claudine.mayer@unistra.fr, arthur.vogt@etu.unistra.fr, tuba.uslu@etu.unistra.fr, n.scalzitti@yahoo.com, kirsley.chennen@unistra.fr, olivier.poch@unistra.fr, thompson@unistra.fr
 8
 9
 10 ² Université Paris Cité, F-75013, Paris, France; claudine.mayer@unistra.fr
 11 * Correspondance: CM: claudine.mayer@unistra.fr; JD: thompson@unistra.fr

12 **Abstract:** In fungi, the most abundant transcription factor (TF) class contains a fungal-specific
 13 ‘GAL4-like’ Zn2C6 DNA binding domain (DBD), while the second class contains another fungal-
 14 specific domain, known as ‘fungal_trans’ or Middle Homology Domain (MHD), whose function
 15 remains largely uncharacterized. Remarkably, almost a third of MHD-containing TF in public se-
 16 quence databases apparently lack DNA binding activity, since they are not predicted to contain a
 17 DBD. Here, we reassess the domain organization of these ‘MHD-only’ proteins using an *in silico*
 18 error-tracking approach. In a large-scale analysis of ~17000 MHD-only TF sequences **present**
 19 **in all fungal phyla except Microsporidia and Cryptomycota**, we show that the vast
 20 majority (>90%) result from genome annotation errors and we were able to predict a new
 21 DBD sequence for 14261 of them. **Most of these sequences correspond to a Zn2C6 domain**
 22 **(82%), with a small proportion of C2H2 domains (4%) found only in Dikarya.** Our results
 23 contradict previous findings that the MHD-only TF are widespread in fungi. In contrast, we show
 24 that they are exceptional cases, and that the fungal-specific Zn2C6-MHD domain pair represents the
 25 canonical domain signature defining the most predominant fungal TF family. We call this family
 26 CeGAL, after the highly characterized members: Cep3, whose 3D structure is determined, and
 27 GAL4, a eukaryotic TF archetype. We believe that this will not only improve the annotation and
 28 classification of the Zn2C6 TF, but will also provide critical guidance for future fungal gene regula-
 29 tory network analyses.

30 **Keywords:** Fungal-specific Transcription Factors; genome annotation errors; large scale biocomputing
 31 analysis.

Citation: To be added by editorial staff during production.

Academic Editor: Firstname Last-name

Received: date
 Revised: date
 Accepted: date
 Published: date



Copyright: © 2023 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

33 **1. Introduction**

34 Transcription factors (TF) are essential for the regulation of expression pathways in
 35 eukaryotes by binding genomic DNA *via* a DNA binding domain (DBD) for example com-
 36 posed of a zinc finger structural motif [1]. The ‘classical’ zinc finger domain coordinates a
 37 single zinc atom with a combination of 4 amino acids, usually cysteine or histidine. How-
 38 ever, the Zn2C6 domain (also called Zn(II)2Cys6, Zn2/Cys6 or Zn(2)-Cys(6) binuclear
 39 cluster domain) is an atypical zinc finger, where the well-conserved
 40 CX₂CX₆CX_{5,16}CX₂CX_{6,8}C motif contains six conserved cysteines that coordinate
 41 two zinc atoms to establish correct folding of the zinc cluster domain [2,3]. The Zn2C6
 42 domain defines the GAL4-like Zn2C6-TF family, which is quasi-specific to fungi and ob-
 43 served ubiquitously in all fungal species, where it represents the most abundant TF family
 44 in each species [4–8]. Zn2C6-TF are involved in a wide range of functions from primary
 45 and secondary metabolisms to multidrug resistance and virulence [4–9]. In addition to the

46 DBD, which is generally localized in the N-terminal part, Zn2C6-TF contain a region for
47 activation of the transcriptional machinery. This region, sometimes called TAD (TransAc-
48 tivation Domain), is present in many eukaryotic TF from yeast to humans [10] and is gen-
49 erally found in the C-terminal part of the proteins.

50 Comparative sequence analyses of the Zn2C6-TF family, initiated in the 1990s, re-
51 vealed the existence of conserved regions between the Zn2C6 DBD and the TAD [11]. One
52 of these regions, named MHR (Middle Homology Region) [2], is composed of three con-
53 served motifs involving about 80 amino acids. The MHR (also known as Fungal_trans)
54 was extended to eight consecutive conserved motifs embedded in a large functional do-
55 main ranging from 225 to 405 residues [12], which is, like the Zn2C6 DBD, specific to fun-
56 gal species and represents the second largest fungal-specific TF class [7]. A mean second-
57 ary structure prediction performed on the eight motifs suggested that they are mainly
58 composed of α -helices. Ten years later, the crystal structure of Cep3 [13,14], a yeast ki-
59 netochore subunit present in the analysis [12], confirmed that the eight motifs are included
60 in an all-alpha domain, hereafter called MHD (Middle Homology Domain). To date, Cep3
61 remains the only experimental 3D structure known for a MHD-containing protein.

62 The functional role of the MHD remains largely elusive, although it has been postu-
63 lated that the fungal-specific Zn2C6-MHD TF might correspond to the metazoan nuclear
64 receptors, with the MHD echoing the metazoan Ligand Binding Domain involved in the
65 regulation of the TF activity (notably in an inhibitory function) and/or in the regulation
66 and recognition of effectors [15]. Furthermore, it has been postulated that the MHD might
67 also participate in DNA target discrimination [16,17].

68 In terms of protein domain organization, the domain pair or bigram [18] composed
69 of the Zn2C6 DBD combined with the MHD is the most frequent in the Zn2C6-TF family.
70 For example, of the 54 Zn2C6-TF from the *Saccharomyces cerevisiae* S288C strain, 44 (81.5%)
71 contain a Zn2C6-MHD domain pair [12]. Strikingly, in the InterPro protein family data-
72 base [19], approximately one third of the proteins exhibiting an MHD (InterPro ID:
73 IPR007219) are not predicted to contain a zinc finger motif of the Zn2C6 or C2H2 types
74 (InterPro ID: IPR001138 or IPR013087). These TF, which apparently lack a DNA binding
75 activity, represent the second largest fungal TF class after the Zn2C6 TF and will be called
76 'MHD-only' hereafter. Except for some rare exceptions [20], MHD-only TF have not been
77 confirmed experimentally and there is some debate about whether the MHD can act in-
78 dependently. For example, in all experimentally proven TF listed in the TRANSFAC da-
79 tabase (<https://genexplain.com/transfac/>), the MHD is always located downstream of a
80 DBD [5].

81 A recent genome-wide study of the complement of TF in the fungus *Aspergillus nid-*
82 *ulans* revealed numerous discrepancies between the predicted protein sequences and the
83 deduced sequences from experimental transcriptomic data, with approximately 30% of
84 the TF needing some type of correction [21]. Among the badly predicted TF, a large ma-
85 jority (78%) concern the Zn2C6- and/or MHD-containing sequences which frequently ex-
86 hibit non-predicted or non-processed introns leading to premature stop codons and erro-
87 neous sequences. It is interesting that most of the *A. nidulans* MHD-only proteins have a
88 domain with predicted DNA binding (mainly of the Zn2C6 type) after RNA sequence
89 analysis. These high-throughput experimental results prompted us to reassess the fungal-
90 specific MHD-containing TF family, by developing a domain-centric error-tracking ap-
91 proach that takes into account potential mispredictions of protein sequences.

92 As a starting point, we collected proteins containing an MHD from three sequence
93 databases with different levels of human expertise involved in the genome annotation
94 process. First, the Saccharomyces Genome Database (SGD) is dedicated to the budding
95 yeast *S. cerevisiae* [22] and provides comprehensive information including protein se-
96 quences from a collection of 48 *S. cerevisiae* strain genomes. Second, the UniProtKB/Swiss-
97 Prot database is the expertly curated component of UniProtKB [23]. Third, the Uni-
98 ProtKB/TrEMBL database contains computer-generated annotations for all translations of
99 the EMBL nucleotide sequence entries. We then focused our analysis on the MHD-only

TF by applying a specific error-tracking protocol that uses different DBD-MHD combinations to identify potentially mispredicted genes in available fungal genomic sequences, and especially mispredictions that affected the protein domain organization.

Our large-scale analysis of almost 17000 MHD-only TF allowed us to verify that at least 90% of them possess upstream genomic sequence regions coding for a DBD, mostly of the Zn2C6 type. These results suggest that the vast majority of the MHD-only TF sequences present in public databases result from errors, and that the Zn2C6-MHD domain pair represents a canonical domain signature defining the most predominant family of TF composed of two fungal-specific domains.

2. Materials and Methods

2.1. Collection of SGD sequences

The 44 proteins from the *S. cerevisiae* S288C strain with a Zn2C6-MHD domain organization [12] were identified in the SGD database (Table S1), and their annotated orthologs in the 47 available strains (Table S2) were downloaded from the SGD web site (http://sgd-archive.yeastgenome.org/sequence/strains/strain_alignments.tar). Genome assemblies for the 47 strains were also downloaded from the SGD web site (<http://sgd-archive.yeastgenome.org/sequence/strains>). Ortholog sequences that did not contain the conserved CX₂CX₆CX_{5,16}CX₂CX_{6,8}C motif were considered to be potentially erroneous.

For each potentially erroneous sequence, we performed a TBLASTN alignment of the S288C reference protein sequence with the corresponding genome assembly. We then tried to identify the causes of the erroneous sequences. First, if TBLASTN hits were found on multiple scaffolds (with percent identity >95% and length>20 amino acids), we assumed that the misprediction was due to a genome assembly issue. If multiple TBLASTN hits (with percent identity >95% and length>20 amino acids) were found on a single scaffold, but in different reading frames, we assumed that the misprediction was due to a sequence insertion leading to a frameshift error. If a TBLASTN hit was found with percent identity >95% and coverage =100%, we assumed that the misprediction was due to a wrongly predicted start codon.

In order to propose a corrected sequence, a protein sequence was then reconstructed from the TBLASTN hits found on the same scaffold. This corrected protein sequence was searched for the conserved CX₂CX₆CX_{5,16}CX₂CX_{6,8}C motif.

2.2 Collection of UniprotKB sequences

Fungal proteins were identified in the UniProt 2022_01 database [23], by querying for proteins annotated with the InterPro entry IPR007219: *Transcription_factor_dom_fun* or *Fungal_trans*, which covers the Middle Homology Domain (MHD) specific to these transcription factors. Domain architectures of all proteins containing an MHD were then extracted from the InterPro v86.0 database [19]. The 37646 UniProt sequences annotated with an MHD, but no DBD, were considered to be potentially erroneous (MHD-only). Potentially erroneous sequences from reviewed UniProt/Swiss-Prot entries and unreviewed UniProt/TrEmbl entries were processed separately: the 12 UniProt/Swiss-Prot proteins with no DBD were analyzed manually, while the 37634 UniProt/TrEMBL sequences were input to the error identification protocol (Figure 1) and described in detail in the following sections.

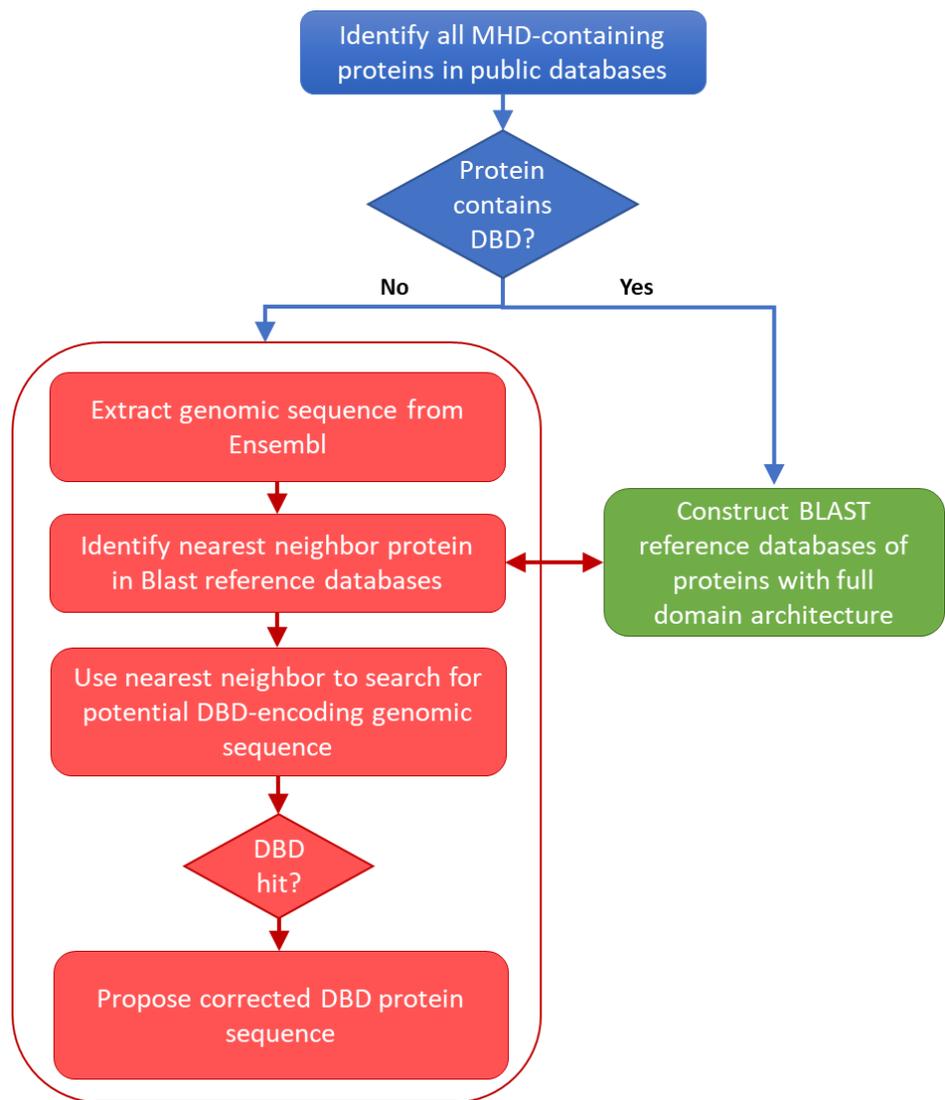


Figure 1. Schema of the protocol used to locate potential errors in sequences retrieved from public databases.

2.3 Construction of BLAST databases of proteins with full domain architecture

UniProt sequences containing an MHD (IPR007219) in combination with a DBD were used to construct BLAST reference databases. Two BLAST databases were constructed: one for each of the two main DBD types, namely Zn2C6 (IPR001138) and C2H2 (IPR013087), found in this TF family to which the well-studied GAL4 protein belongs. The Zn2C6 BLAST reference database contains 80456 sequences, while the C2H2 BLAST reference database contains 6314 sequences.

2.4 Extraction of genomic sequences

For all potentially erroneous sequences in UniProt, the corresponding genomic DNA sequences were extracted from the Ensembl database [24], when an Ensembl cross-reference was available in the UniProt database. To improve detection of missing DBD, the full length gene sequences were retrieved together with an additional 1000 nucleotides upstream of the 5' end of the gene. For the 37634 potential error sequences, 16760 genomic DNA sequences were found in the Ensembl database.

2.5 Identification of nearest neighbor reference sequences

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

For each MHD-only sequence, BLASTP searches were performed in the two BLAST reference databases containing Zn2C6 and C2H2 sequences in combination with an MHD. The nearest neighbor sequences with the required domain combination (*i.e.* DBD and MHD) were selected, if a BLASTP hit was identified with E-value < 0.005. Figure 2 shows the E-value distribution of BLASTP hits obtained for all neighbor sequence searches.

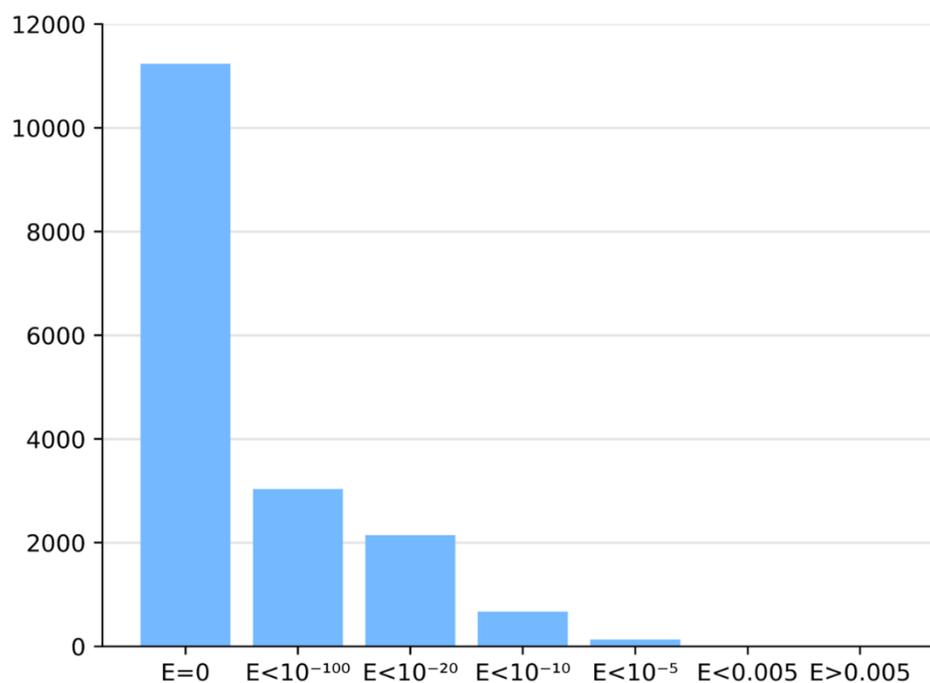


Figure 2. Histogram of BLASTP E-values in neighbor identification step.

2.6 Identification of missing DBD sequences

For each MHD-only sequence with a BLASTP hit to a nearest neighbor reference protein, two complementary approaches were implemented to search for the missing DBD sequence. First, a local TBLASTN search was performed in the genomic sequence of the potential error sequence, using the protein DBD sequence segment of the nearest neighbor as a query. TBLASTN alignments with E-value < 0.0001 were taken into account. Second, a global pairwise alignment was performed between the genomic sequence of the MHD-only sequence and the full-length protein sequence of the nearest neighbor, using the ProSplign software developed by the NCBI (<https://www.ncbi.nlm.nih.gov/sutils/static/prosplign/prosplign.html>). ProSplign is a tool for protein to genomic sequence alignment, and is an integral component of the NCBI Eukaryotic Genome Annotation Pipeline. Genes are first localized on the genomic sequence in a compartmentalization step that starts with computing protein-to-genomic blast hits. These give initial insight into the structure of compartments. Hits are separated into two same-strand sets and then compartments are identified within each strand. To do so, the optimization problem is formally defined in terms of genomic sequence coverage and then solved with a dynamic programming algorithm. ProSplign has been shown to produce accurate spliced alignments and is able to compute alignments of distantly related proteins with low similarity.

Pairwise alignments obtained from TBLASTN and ProSplign were analyzed to identify potential DBD-encoding sequence segments in the erroneous sequences. Finally, the potential DBD-encoding sequence segments were compared to an HMM representing the DBD downloaded from the Pfam protein family database [25]: PF00172 for the Zn2C6 DBD and PF00096 for the C2H2 DBD. To do this, the hhmsearch program from the

HMMER suite [26] was used and sequences with an E-value < 0.1 were considered as hits. In addition to the hmmsearch E-value and to eliminate a number of false positive hits, we also checked for conserved amino acids: for Zn2C6 DBD, two occurrences of the pattern C-x(2)-C (where x is any amino acid) were required, while for C2H2 DBD, one occurrence for each of the C-x(2,4)-C and H-x(3,5)-H patterns were required.

In order to determine whether our sequence curation protocol over-predicts DBDs in the potentially erroneous sequences, we used the same error identification protocol to search for Zn2C6 DBD in the C-terminal region of the proteins (*i.e.* downstream of the MHD). The results of this analysis are described in the Supplementary methods.

3. Results

3.1. MHD-containing proteins in the SGD database

We first analyzed the MHD-containing proteins in the SGD, a database that provides comprehensive integrated information for *S. cerevisiae*. The reference (S288C) strain of *S. cerevisiae* contains 44 proteins with an MHD (Table S1), all of them exhibiting a Zn2C6-MHD domain pair [12]. In addition to S288C, the SGD contains genome assemblies and annotations for a further 47 strains of *S. cerevisiae* (Table S2). For each of these 47 strains, we extracted the annotated orthologs of the 44 S288C Zn2C6-MHD proteins. If all orthologs are conserved in all strains, we would expect a total of 2068 orthologs (44 orthologs from each of the 47 strains), but only 1793 orthologous sequences were found in the SGD (Table 1). In other words, 275 (13%) orthologous sequences were not predicted. Furthermore, for the 1793 predicted sequences, 253 (14%) of them did not contain the conserved CX₂CX₆CX_{5,16}CX₂CX_{6,8}C motif and were considered to be potentially mispredicted genes.

Table 1. Domain annotations of MHD containing proteins in the SGD database.

Domain annotation	S288C reference strain	47 other strains
Zn2C6-MHD	44 (100%)	1540 (86%)
MHD-only	0 (0%)	253 (14%)
Total	44	1793

Proportion of sequences with each domain combination, with respect to the total number of sequences (in parentheses).

To investigate the causes of the 253 genes with potential errors, we used the 44 S288C protein sequences to search the corresponding genome assemblies in the SGD using TBLASTN (Table 2).

Table 2. Probable sources of protein sequence prediction errors in the SGD database.

Error type	Probable cause of error	MHD-only
Genome sequence error	Frameshift	190 (75%)
	2 or more scaffolds	9 (4%)
Gene prediction error	Wrong start codon	46 (18%)
Undetermined	Undetermined	8 (3%)
Total		253 (100%)

Proportion of each cause of error with respect to the total number of errors detected (in parentheses).

In the majority of cases, significant hits to genomic regions were found and the protein sequence errors could be linked to genome sequencing or assembly issues. Indeed, for 190 (75%) of the 253 proteins, the S288C protein sequence matched to multiple segments of a single genome scaffold with sequence identity of at least 95%, although the matching segments were found in different reading frames. These frameshifts were

238 mainly due to the insertion of one or two bases in the genome sequence of the *S. cerevisiae*
239 strain, as compared to the S288C sequence. A further nine sequences were found split over
240 multiple scaffolds. For 46 (18%) of the 253 proteins, the S288C protein sequence matched
241 the genome scaffold with a coverage of 100% and sequence identity of at least 95%, and
242 we concluded that the absence of the Zn2C6 domain was due to a wrongly predicted start
243 codon.

244 We then tried to correct the 253 erroneous sequences by reconstructing the protein
245 sequence from the TBLASTN genome hits. For 243 (95%) sequences, a complete Zn2C6
246 domain could be found upstream of the MHD domain (Table S3). The full-length se-
247 quences are provided as a Fasta file. After taking into account the detected gene prediction
248 errors, only ten of the 253 MHD-only proteins remained for which a Zn2C6 domain was
249 not found. These included the nine gene sequences split over multiple scaffolds, which
250 could not be resolved due to the genome assembly issues, and one sequence with a
251 frameshift error (although manual analysis of this genome sequence indicates a frameshift
252 error affecting one of the conserved cysteines in the Zn2C6 domain).

253 In order to validate our predictions for the 253 erroneous sequences, we searched for
254 RNA-seq datasets in the NCBI Gene Expression Omnibus (GEO) project corresponding to
255 the different *S. cerevisiae* strains, focusing on the strains with the highest number of poten-
256 tially mispredicted genes. This transcriptome analysis is described in Supplementary
257 methods. For each mispredicted gene with a coverage of at least 30 reads, the aligned
258 reads in the region of the gene were manually reviewed, confirming that all the predicted
259 DBD sequences were expressed at similar levels as the MHD portion (Figure S1).

260 In summary, no convincing evidence of MHD-only proteins was found in any of the
261 47 *S. cerevisiae* strains analyzed here and all the identified MHD located in reliable genome
262 sequence scaffolds were associated with upstream Zn2C6 domains.

264 3.2. MHD-containing proteins in the Uniprot database

265 We then queried the UniProt database for all proteins annotated with the MHD (In-
266 terPro ID: IPR007219), resulting in a total of 126861 proteins, with 126691 in the unre-
267 viewed TrEMBL section and 170 in the reviewed Swiss-Prot section. The MHD containing
268 proteins have a wide range of domain architectures, with 1905 different architectures
269 listed in the InterPro database, although the most frequent domain pairs are: (i) MHD with
270 a Zn2C6 DBD (IPR001138), (ii) MHD-only, and (iii) MHD with one or two C2H2 DBD
271 (IPR013087), as shown in Figure 3 and Table S4.
272

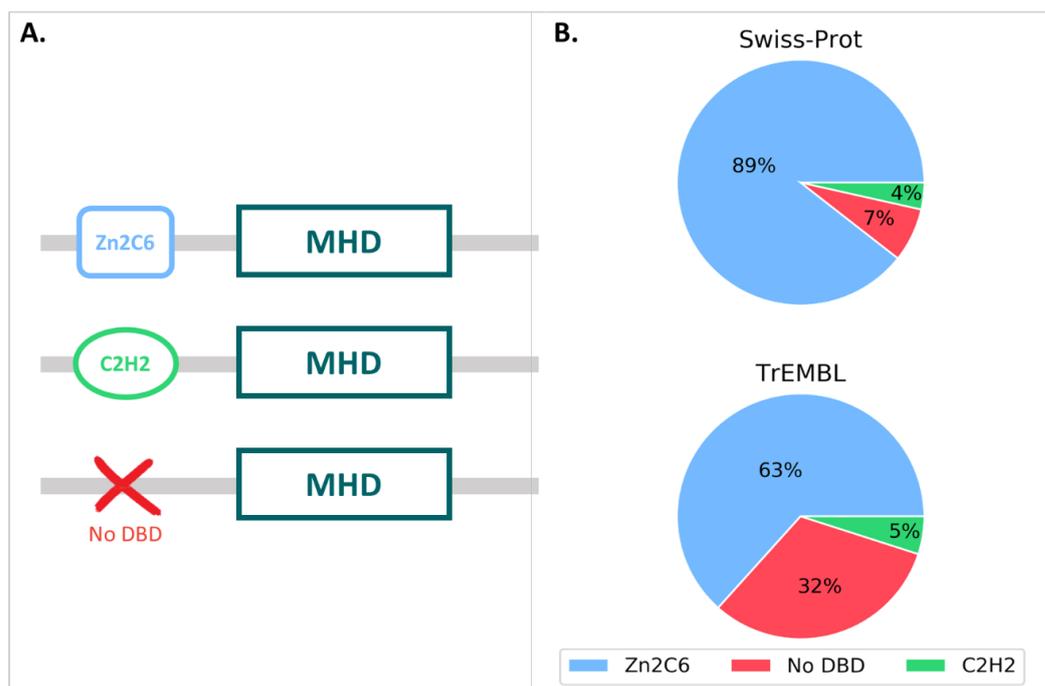


Figure 3. (A) MHD domain pairs found in the UniProtKB. (B) Proportion of each domain pair found in the Swiss-Prot and TrEMBL sections, with respect to the total number of MHD-containing sequences.

For the 170 proteins from the Swiss-Prot section, nearly 90% contain the Zn2C6-MHD domain pair, although this combination is found in a smaller proportion of the TrEmbl proteins with only 63.4%. Conversely, the proportion of MHD-only proteins lacking an annotated DBD is much higher in TrEmbl (31.6%) than in Swiss-Prot (7.1%).

3.3. Manual analysis of the 12 Swiss-Prot MHD-only sequences

Since Swiss-Prot entries are curated by experts, we manually investigated the twelve MHD-only sequences in this database (Table S5). Where possible, we extracted the corresponding genome sequence from either ENSEMBL [24] or GENBANK [27] databases, and translated the genome sequence in the three frames to search for potential DBD encoding regions. This was not possible for four of the twelve sequences. For Q5AR44 and A0A5C1RF03, the gene is located at the start of a contig and the genome region upstream of the annotated gene is not available. For B8NJC5, according to the ENSEMBL database, the upstream gene codes for a small protein coding a Zn2C6 DBD. Finally, A6SSW6 is a shorter protein of length 179 (compared to > 500 for the other Swiss-Prot proteins), and the MHD hit in the InterPro database is a partial domain. The genome sequence (FR718884) is annotated as "possibly a relic of a transcription factor".

For all remaining eight proteins, a potential DBD sequence is found either within the existing annotated gene *via* alternative splicing, or the proximal 5' region (<1000 nt) *via* an alternative start codon or in a new exon (Figure 4).

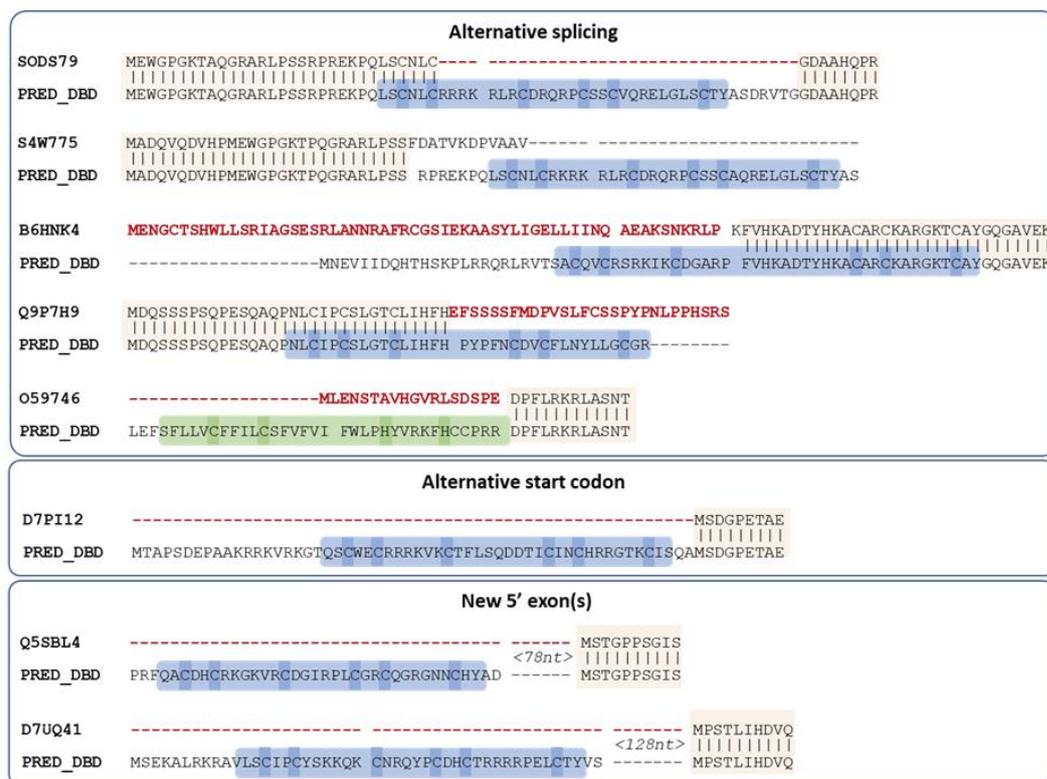


Figure 4. Proposed new sequences for missing DBD of Swiss-Prot MHD proteins. Conserved regions between the existing Swiss-Prot sequence and the proposed sequence are indicated by vertical lines. Regions in the existing Swiss-Prot sequence shown in red are replaced in the predicted sequence, while the predicted DBD is outlined in blue (Zn2C6) or green (C2H2), with cysteines/histidines corresponding to potential zinc binding amino acids highlighted in blue or green. Spaces in the sequences indicate annotated or predicted splice sites.

3.4. Automatic analysis of 16760 TrEMBL MHD-only sequences

Based on the manual analysis of Swiss-Prot described in the previous section, an automatic protocol was developed to analyze potentially erroneous sequences retrieved from the UniProt/TrEMBL database. The first step in the protocol involved identifying the corresponding genomic sequences in the ENSEMBL database. This resulted in a set of 16760 sequences that were used as input for the main error detection step (see Methods). Two different methods were implemented to locate genomic regions within or upstream of the gene that could encode the missing DBD, using either a local or global alignment approach. Figure 5 shows the number of DBD identified by the two methods.

299
 300
 301
 302
 303
 304
 305
 306
 307
 308
 309
 310
 311
 312
 313
 314
 315
 316

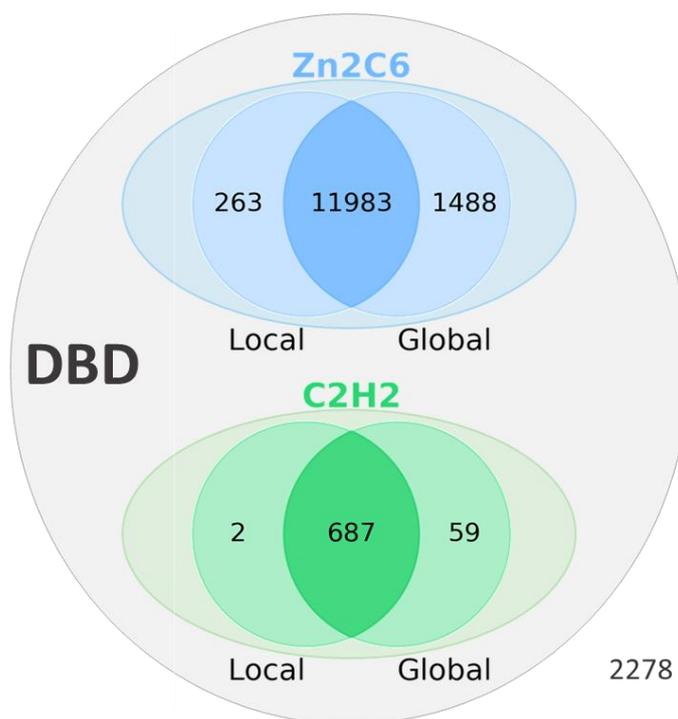


Figure 5. Results of the error identification step in 16760 MHD-only TF sequences from the TrEMBL database. Number of DBD (Zn2C6 or C2H2) identified by local and global alignment methods. **A total of 14482 MHD-only TF sequences could be attributed to gene prediction errors, while no DBD could be identified for the remaining 2278 sequences.**

By integrating the results of the local and global alignment searches, DBD sequences could be proposed for 14482 (86%) of the 16760 MHD-only sequences tested (Table S6). The proposed DBD sequences are distributed in 476 fungal species or strains and are provided as a Fasta file and in the Supplementary Data file with additional information concerning the nearest neighbor used for blast searches, the description, the pathogenicity (against animals or plants) and a complete taxonomic description. Most of these sequences correspond to a Zn2C6 domain (82%), with a smaller proportion of C2H2 domains (4%), which correlates well with the proportions found in the manually curated Swiss-Prot section. To verify the quality of the computer-predicted DBD sequences, we took advantage of a previous study performed on the *Aspergillus flavus* TF proteome by Chang and Ehrlich [28]. By manual analysis of the genomic region, the authors identified an upstream Zn2C6 domain for 67% of the studied MHD-only TF. Of the 85 DBD sequences we predicted here for *A. flavus* MHD-only TF (Supplementary File 2), 59 sequences (69.4%) were strictly identical to the DBD sequences detected by Chang and Ehrlich, thus highlighting the accuracy of our automatic error-checking protocol.

For the 2278 (13.6%) proteins with no DBD identified by our automatic protocol, we then investigated potential causes for the erroneous sequences. Partial hits, with hmm-search scores below the defined threshold and part of the conserved Zn2C6 or C2H2 motifs (see Methods), were found in 905 sequences. These might indicate complex exon/intron structures that were partly mispredicted by our protocol (an example is shown in the following section) or might be caused by genome sequencing or assembly issues. For example, undefined regions in the genomic sequences, represented by 'N' characters, were found in 1363 of the 2278 proteins. Other reasons for not identifying a DBD include (i) the DBD is located more than 1000 nucleotides upstream of the gene, (ii) the related sequence is not conserved enough to allow protein-DNA alignment of the DBD.

317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350

3.5. Reassessment of domain pairs in MHD-containing sequences

If the missing DBD sequences proposed here were integrated into the public databases, the number of MHD-only proteins would be reduced from 12 to 4 for Swiss-Prot, and from 16760 to 2278 for TrEMBL (Figure 6 and Table S6). More importantly perhaps, this would lead to a significant difference in the distribution of domain pairs present in MHD-containing sequences. In the public databases, this distribution is 65%, 5% and 30% for Zn2C6-MHD, C2H2-MHD and MHD-only respectively. However, our error-tracking protocol indicates that the true distribution is closer to 90%, 6% and 4% for Zn2C6-MHD, C2H2-MHD and MHD-only respectively.

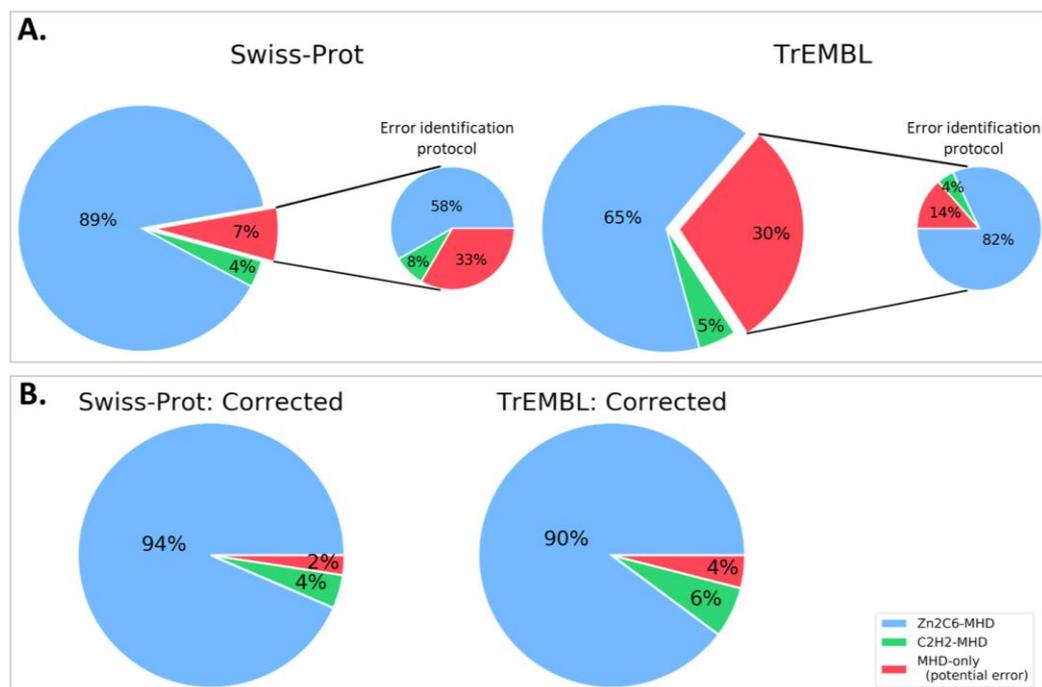


Figure 6. Proportion of sequences with Zn2C6-MHD (blue) or C2H2-MHD (green) domain combinations in A. public databases: Swiss-Prot and TrEMBL (sequences mapped to ENSEMBL only) and B. after applying our error identification protocol. Proportion of potentially erroneous sequences lacking a DBD is shown in red.

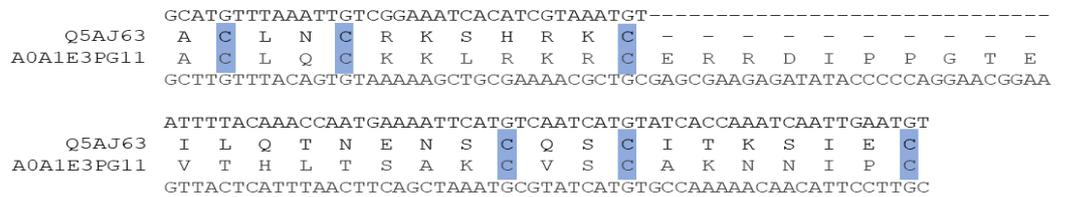
Concerning the phylogenetic distribution of the DBD-MHD domain pair, at the phylum level, MHD domains are present in all phyla except Microsporidia and Cryptomycota (Figure S2). The Zn2C6-MHD domain pair is also present in all phyla except Microsporidia and Cryptomycota, and it is therefore difficult to determine the origin or emergence of this domain pair. In contrast, the C2H2-MHD domain pair is found only in Dikarya (Basidiomycota and Ascomycota).

Interestingly, it has been shown previously that there is a significant difference in the TF repertoire of ascomycete and basidiomycete fungi [8], and in particular that the Zn2C6 family (33%) is much more prevalent than the C2H2 (10%) in ascomycete TF, compared to basidiomycete TF (20% and 15% for Zn2C6 and C2H2 respectively). Despite this overall enrichment of C2H2 in the basidiomycete TF, within the MHD sequences, the proportions of Zn2C6 and C2H2 are similar in both clades (90% and 7% for Ascomycota compared to 92% and 3% for Basidiomycota) (Table S7).

The level of protein sequence errors is of course dependent on the quality of the genome sequencing, assembly and annotation. A small number of well-characterized organisms had no MHD-only sequences in the Uniprot database, including model organisms like *S. cerevisiae*, *Yarrowia lipolytica*, or *Ustilago maydis*. Nevertheless, some model organisms had a small number of MHD-only sequences, for example *Schizosaccharomyces pombe* has 27 proteins annotated with an MHD, of which two proteins had no DBD, or *Candida*

385
386
387

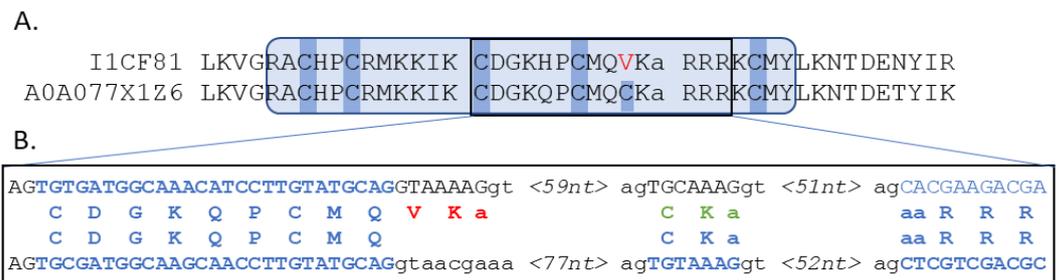
albicans with 28 MHD-containing proteins, of which only one has no DBD: Q5AJ63_CANAL (Figure 7).



388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403

Figure 7. Predicted DBD for *Candida albicans* sequence Q5AJ63_CANAL, aligned with neighbor sequence A0A1E3PG11_9ASCO (hmmsearch E-value = 6.5e-10). Conserved cysteines characterizing the Zn2C6 DBD are highlighted in blue.

At the other extreme, *Rhizopus delemar* has 34 MHD-containing proteins of which only 6 are also annotated with a DBD, i.e. 82% are MHD-only proteins. According to our protocol, DBD could be detected for a further 24 MHD proteins and only 4 (12%) lack a DBD. Further manual analysis of these 4 MHD-only proteins showed that two genes (I1BR31, I1C782) have regions coding complete Zn2C6 domains within 1500 nt upstream of the 5' end, while one gene (I1CF81) has a partial hit within the default threshold of 1000 nt upstream of the 5' end (-490 to -192). According to the public databases, I1CF81 is coded by a gene with one exon and contains one known domain, the MHD. The partial DBD hit for I1CF81 in fact contains a misprediction of a short exon coding for 2 amino acids, as shown in Figure 8.



404
405
406
407
408
409
410
411

Figure 8. (A) Protein alignment of query I1CF81 with the neighbor A0A077X1Z6, showing the partial hit identified by the protocol, where the predicted sequence presents five of the six conserved cysteines that characterize the Zn2C6 DBD (hmmsearch E-value = 1.5e-11). Exon/intron boundaries are indicated by gaps in the sequences. (B) Genome-level comparison of query I1CF81 with the neighbor A0A077X1Z6, showing correctly predicted amino acids (blue), protocol mispredicted amino acids (red) and alternative manual prediction (green).

412
413
414
415
416
417
418
419
420
421
422
423

4. Discussion

In this work, we used a domain-centric *in silico* approach to show that the second most abundant fungal-specific TF family in the public databases, namely the MHD-only TF, results largely from genome annotation errors leading to unpredicted DBD. Taking advantage of the high quality sequences of the *S. cerevisiae* reference strain S288C and the availability of numerous other fungal genomes, we defined an error-tracking strategy involving increasing levels of difficulty: starting with the analysis of the MHD-only sequences present in 48 closely related *S. cerevisiae* strains, followed by the MHD-only sequences in the expert curated Swiss-Prot database and finally, in the automatically generated TrEMBL database.

The reference *S. cerevisiae* S288C strain has no MHD-only TF coding genes and we showed that, for 95% of the MHD-only TF genes observed in the other *S. cerevisiae* strains,

424 a complete Zn2C6 domain is located upstream of the MHD-coding genomic region. This
425 highlights an unexpectedly high rate of gene prediction errors in such closely related ge-
426 nomes. This high error rate was confirmed for the MHD-only TF present in the expert
427 curated Swiss-Prot database, since a DBD could be identified for all the proteins whose
428 corresponding genomic sequence was available. Finally, concerning the TrEMBL proteins,
429 our error-tracking protocol showed that 89% of the MHD-only TF exhibit upstream ge-
430 nomic sequence regions coding for a DBD. These analyses, showing that MHD-only TF
431 sequences result predominantly from prediction inaccuracies, are in line with the error
432 rate of 66% observed in the manual analysis of MHD-only TF in *Aspergillus flavus* [28].
433 Similarly, a recent high throughput RNA-sequencing experiment in fungal species iden-
434 tifies a large proportion of prediction errors in TF sequences [21].

435 The high rate of wrongly predicted TF sequences (at least 82%) is particularly sur-
436 prising given that (i) fungal genome sequences are generally of better quality with fewer
437 genome assembly errors, thanks to their relatively small, compact genomes and the low
438 level of repetitive sequences in most fungi [29], (ii) fungi serve as model eukaryotic organ-
439 isms and a wide range of diverse genomes have been sequenced and annotated (SGD,
440 Génolevures: genolevures.org, 1000 Fungal Genomes Project: mycocosm.jgi.doe.gov), (iii)
441 genome annotation in the fungi is facilitated by the relatively streamlined gene structures
442 and transcriptional processes in these organisms with few and typically short introns
443 rarely implicated in alternative splicing. Our results clearly indicate that all these fungal
444 features, which should promote gene prediction quality, do not limit the error rates at
445 least in the studied TF family. Most importantly, the true number of MHD-only genes
446 remains to be determined, if MHD can indeed act independently [5,7].

447 The notion of errors in public protein databases is a recurrent problem [30–32] and
448 substantial efforts have been invested to identify and correct genome annotation errors
449 [33–35]. Some important causes of erroneous protein sequences have been identified, in-
450 cluding the genome sequence quality and gene structure complexity [36], as well as re-
451 dundant or conflicting information in different resources or in the literature [32,37]. Con-
452 sequently, it has been estimated that 40 to 60% of the protein sequences in public databases
453 are erroneous [38–40]. Typical errors include missing exons, non-coding sequence reten-
454 tion in exons, wrong exon and gene boundaries, fragmenting genes and merging neigh-
455 boring genes. Our results confirm that genome sequence quality and gene structure com-
456 plexity are major drawbacks for correct annotation and provide further evidence of the
457 potential of domain-centric approaches to improve automated methods to identify and
458 correct mispredicted protein sequences [39,41–43].

459 It has been established that some domains always co-occur leading to the concept of
460 associated domains in proteins also been called ‘supra-domains’ [44], DASSEM units [45],
461 or domain co-occurrence [46,47].

462 Our results indicate that the fungal-specific MHD forms part of a synergistic domain
463 pair with a zinc finger DBD, mostly of the fungal-specific Zn2C6 type. As a consequence,
464 the protein sequences exhibiting the ‘supra-domain’ Zn2C6-MHD architecture may define
465 the most widely distributed and abundant fungal TF family that we propose to name
466 CeGAL after its most characterized members: Cep3, whose 3D structure has been deter-
467 mined, and GAL4, the archetypal fungal TF. This definition will clarify the classification
468 of fungal TF and will provide better discrimination of the so-called GAL4-like regulators
469 defined according to the presence of a Zn2C6 domain and which include TF with diverse
470 domain architectures.

471 Finally, the Zn2C6-MHD combination within the CeGAL family members may have
472 significant consequences for the fungal scientific community. As Zn2C6-TF function
473 mainly as homodimers or heterodimers, this implies that the number of sequence-specific
474 TF and the array of control DNA-sequences in target genes need to be reconsidered, as
475 well as the degree of combinatorial regulation involved in the wide range of fungal pro-
476 cesses controlled by **these** TF [4–8]. More importantly, this will also contribute to a better

understanding of fungal Gene Regulatory Networks (GRN), that aim to define the complete set of regulatory interactions between TFs and their target genes at a species level. Classically, GRN analyses combine an initial step of genome-wide characterization of TF families with experimental data related to transcriptional effects of TF deletion/overexpression, chromatin immunoprecipitation (ChIP)-based TF binding data, protein-protein interactions or pairs of genes involved in genetic interactions [48]. However, as GRN studies generally exclude proteins lacking a DBD, the complete repertoire of fungal TF is frequently underestimated. In light of the impressive improvement of the TF specificity prediction tools [49–51], we believe that the definition of the CeGAL family combined with the 14 000 DBD sequences provided in this study will permit more robust GRN analyses.

Supplementary Materials: The following supporting information can be downloaded at: www.mdpi.com/xxx/s1,

SupplementaryMethods.docx: Supplementary methods;

SupplementaryData.xlsx: Proposed DBD sequences for 14482 MHD-only sequences, with additional information concerning the nearest neighbor used for blast searches, the description, the pathogenicity (against animals or plants) and a complete taxonomic description;

Figure S1: Screenshots of IGV browser window displaying the genome region around the mispredicted genes in A. AWRI796 and B. VL3;

Figure S2: Phylogenetic distribution of the DBD-MHD domain pair and classification at the phylum level;

Table S1: List of the 44 Zn2C6-MHD containing genes coded by the *Saccharomyces cerevisiae* S288C genome;

Table S2: List of the 47 genomes of *Saccharomyces cerevisiae* strains extracted from the SGD;

Table S3: Results of the analysis of 253 MHD-containing sequences in the SGD database;

Table S4: Frequency of MHD domain pairs in the UniProtKB Swiss-Prot and TrEMBL sections.;

Table S5: Swiss-Prot entries for 12 MHD-only proteins, and cross-references to genome and literature databases.

Table S6: Number of MHD-containing sequences with different domain combinations in public databases: Swiss-Prot and TrEMBL (sequences mapped to ENSEMBL only) and number of DBD identified by our error identification protocol;

Table S7: Comparison of MHD-containing sequences from Ascomycota and Basidiomycota.

Author Contributions: Conceptualization, C.M., N.S., O.P. and J.D.T.; Methodology, C.M., A.V., T.U., N.S., O.P. and J.D.T.; Software, A.V., T.U. and J.D.T.; Validation, C.M., A.V., T.U., N.S., O.P. and J.D.T.; Formal analysis, C.M., O.P. and J.D.T.; Investigation, C.M., A.V., T.U., N.S., O.P. and J.D.T.; Resources, C.M., A.V., T.U., N.S., O.P. and J.D.T.; Data curation, C.M., A.V., T.U., N.S., O.P. and J.D.T.; Writing—original draft preparation, C.M., O.P. and J.D.T.; Writing—review and editing, C.M., N.S., O.P. and J.D.T.; Visualization, C.M., A.V., T.U., N.S., O.P. and J.D.T.; Supervision, C.M., O.P. and J.D.T.; Project administration, C.M., O.P. and J.D.T.; Funding acquisition, C.M., O.P. and J.D.T. All authors have read and agreed to the published version of the manuscript.

Funding: This research was funded by French Infrastructure Institut Français de Bioinformatique (IFB) ANR-11-INBS-0013, and Institute funds from the French Centre National de la Recherche Scientifique, the University of Strasbourg.

Institutional Review Board Statement: Not applicable

Data Availability Statement: The genomic and protein sequences supporting the conclusions of this article are available in public databases: SGD, UniprotKB, Ensembl and GENBANK. The full-length sequences for the corrected MHD-containing proteins from the SGD database, and the corrected sequences for the missing DBD in the MHD-containing proteins from the TrEMBL database are provided as supplementary datafiles.

Acknowledgments: We thank the members of the BiGEst bioinformatics platform for their assistance. This work was supported by French Infrastructure Institut Français de Bioinformatique (IFB) ANR-11-INBS-0013, and Institute funds from the French Centre National de la Recherche Scientifique, the University of Strasbourg.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Lehninger, A.; Nelson, D.; Cox, M. *Principles of Biochemistry*; 2nd ed.; Worth: New York, U.S, 1993;
2. Schjerling, P.; Holmberg, S. Comparative Amino Acid Sequence Analysis of the C6 Zinc Cluster Family of Transcriptional Regulators. *Nucleic Acids Res.* **1996**, *24*, 4599–4607.
3. Vallee, B.L.; Coleman, J.E.; Auld, D.S. Zinc Fingers, Zinc Clusters, and Zinc Twists in DNA-Binding Protein Domains. *Proc. Natl. Acad. Sci.* **1991**, *88*, 999–1003, doi:10.1073/pnas.88.3.999.
4. MacPherson, S.; Larochelle, M.; Turcotte, B. A Fungal Family of Transcriptional Regulators: The Zinc Cluster Proteins. *Microbiol. Mol. Biol. Rev.* **2006**, *70*, 583–604, doi:10.1128/MMBR.00015-06.
5. Shelest, E. Transcription Factors in Fungi. *FEMS Microbiol. Lett.* **2008**, *286*, 145–151, doi:10.1111/j.1574-6968.2008.01293.x.
6. Shelest, E. Transcription Factors in Fungi: TFome Dynamics, Three Major Families, and Dual-Specificity TFs. *Front. Genet.* **2017**, *8*.
7. Tianqiao, S.; Xiong, Z.; You, Z.; Dong, L.; Jiaoling, Y.; Junjie, Y.; Mina, Y.; Huijuan, C.; Mingli, Y.; Xiayan, P.; et al. Genome-Wide Identification of Zn2Cys6 Class Fungal-Specific Transcription Factors (ZnFTFs) and Functional Analysis of UvZnFTF1 in *Ustilaginoidea Virens*. *Rice Sci.* **2021**, *28*, 567–578, doi:10.1016/j.rsci.2021.03.001.
8. Todd, R.B.; Zhou, M.; Ohm, R.A.; Leeggangers, H.A.; Visser, L.; de Vries, R.P. Prevalence of Transcription Factors in Ascomycete and Basidiomycete Fungi. *BMC Genomics* **2014**, *15*, 214, doi:10.1186/1471-2164-15-214.
9. John, E.; Singh, K.B.; Oliver, R.P.; Tan, K. Transcription Factor Control of Virulence in Phytopathogenic Fungi. *Mol. Plant Pathol.* **2021**, *22*, 858–881, doi:10.1111/mpp.13056.
10. Piskacek, M.; Havelka, M.; Rezacova, M.; Knight, A. The 9aaTAD Transactivation Domains: From Gal4 to P53. *PLOS ONE* **2016**, *11*, e0162842, doi:10.1371/journal.pone.0162842.
11. Todd, R.B.; Andrianopoulos, A. Evolution of a Fungal Regulatory Gene Family: The Zn(II)2Cys6 Binuclear Cluster DNA Binding Motif. *Fungal Genet. Biol. FG B* **1997**, *21*, 388–405, doi:10.1006/fgbi.1997.0993.
12. Poch, O. Conservation of a Putative Inhibitory Domain in the GAL4 Family Members. *Gene* **1997**, *184*, 229–235, doi:10.1016/S0378-1119(96)00602-6.
13. Bellizzi, J.J.; Sorger, P.K.; Harrison, S.C. Crystal Structure of the Yeast Inner Kinetochores Subunit Cep3p. *Structure* **2007**, *15*, 1422–1430, doi:10.1016/j.str.2007.09.008.
14. Purvis, A.; Singleton, M.R. Insights into Kinetochores–DNA Interactions from the Structure of Cep3Δ. *EMBO Rep.* **2008**, *9*, 56–62, doi:10.1038/sj.embor.7401139.
15. Näär, A.M.; Thakur, J.K. Nuclear Receptor-like Transcription Factors in Fungi. *Genes Dev.* **2009**, *23*, 419–432, doi:10.1101/gad.1743009.
16. Turcotte, B.; Liang, X.B.; Robert, F.; Soontorngun, N. Transcriptional Regulation of Nonfermentable Carbon Utilization in Budding Yeast. *FEMS Yeast Res.* **2009**, *10*, 2–13, doi:10.1111/j.1567-1364.2009.00555.x.
17. Mollapour, M.; Piper, P.W. Activity of the Yeast Zinc-Finger Transcription Factor War1 Is Lost with Alanine Mutation of Two Putative Phosphorylation Sites in the Activation Domain. *Yeast Chichester Engl.* **2012**, *29*, 39–44, doi:10.1002/yea.1915.
18. Yu, L.; Tanwar, D.K.; Penha, E.D.S.; Wolf, Y.I.; Koonin, E.V.; Basu, M.K. Grammar of Protein Domain Architectures. *Proc. Natl. Acad. Sci.* **2019**, *116*, 3636–3645, doi:10.1073/pnas.1814684116.

- 571 19. Blum, M.; Chang, H.-Y.; Chuguransky, S.; Grego, T.; Kandasamy, S.; Mitchell, A.; Nuka, G.; Paysan-Lafosse, T.;
572 Qureshi, M.; Raj, S.; et al. The InterPro Protein Families and Domains Database: 20 Years On. *Nucleic Acids Res.* **2021**,
573 *49*, D344–D354, doi:10.1093/nar/gkaa977.
- 574 20. Zhang, W.-Q.; Gui, Y.-J.; Short, D.P.G.; Li, T.-G.; Zhang, D.-D.; Zhou, L.; Liu, C.; Bao, Y.-M.; Subbarao, K.V.; Chen,
575 J.-Y.; et al. Verticillium Dahliae Transcription Factor VdFTF1 Regulates the Expression of Multiple Secreted
576 Virulence Factors and Is Required for Full Virulence in Cotton. *Mol. Plant Pathol.* **2018**, *19*, 841–857,
577 doi:10.1111/mpp.12569.
- 578 21. Etxebeste, O. Transcription Factors in the Fungus Aspergillus Nidulans: Markers of Genetic Innovation, Network
579 Rewiring and Conflict between Genomics and Transcriptomics. *J. Fungi* **2021**, *7*, 600, doi:10.3390/jof7080600.
- 580 22. Engel, S.R.; Wong, E.D.; Nash, R.S.; Aleksander, S.; Alexander, M.; Douglass, E.; Karra, K.; Miyasato, S.R.; Simison,
581 M.; Skrzypek, M.S.; et al. New Data and Collaborations at the Saccharomyces Genome Database: Updated Reference
582 Genome, Alleles, and the Alliance of Genome Resources. *Genetics* **2021**, *220*, iyab224, doi:10.1093/genetics/iyab224.
- 583 23. The UniProt Consortium UniProt: The Universal Protein Knowledgebase in 2021. *Nucleic Acids Res.* **2021**, *49*, D480–
584 D489, doi:10.1093/nar/gkaa1100.
- 585 24. Howe, K.L.; Achuthan, P.; Allen, J.; Allen, J.; Alvarez-Jarreta, J.; Amode, M.R.; Armean, I.M.; Azov, A.G.; Bennett,
586 R.; Bhai, J.; et al. Ensembl 2021. *Nucleic Acids Res.* **2021**, *49*, D884–D891, doi:10.1093/nar/gkaa942.
- 587 25. Mistry, J.; Chuguransky, S.; Williams, L.; Qureshi, M.; Salazar, G.A.; Sonnhammer, E.L.L.; Tosatto, S.C.E.; Paladin,
588 L.; Raj, S.; Richardson, L.J.; et al. Pfam: The Protein Families Database in 2021. *Nucleic Acids Res.* **2021**, *49*, D412–
589 D419, doi:10.1093/nar/gkaa913.
- 590 26. Johnson, L.S.; Eddy, S.R.; Portugaly, E. Hidden Markov Model Speed Heuristic and Iterative HMM Search
591 Procedure. *BMC Bioinformatics* **2010**, *11*, 431, doi:10.1186/1471-2105-11-431.
- 592 27. Sayers, E.W.; Cavanaugh, M.; Clark, K.; Pruitt, K.D.; Schoch, C.L.; Sherry, S.T.; Karsch-Mizrachi, I. GenBank. *Nucleic*
593 *Acids Res.* **2021**, *50*, D161–D164, doi:10.1093/nar/gkab1135.
- 594 28. Chang, P.-K.; Ehrlich, K.C. Genome-Wide Analysis of the Zn(II)2Cys₆ Zinc Cluster-Encoding Gene Family in
595 Aspergillus Flavus. *Appl. Microbiol. Biotechnol.* **2013**, *97*, 4289–4300, doi:10.1007/s00253-013-4865-2.
- 596 29. Galagan, J.E.; Henn, M.R.; Ma, L.-J.; Cuomo, C.A.; Birren, B. Genomics of the Fungal Kingdom: Insights into
597 Eukaryotic Biology. *Genome Res.* **2005**, *15*, 1620–1631, doi:10.1101/gr.3767105.
- 598 30. International Society for Biocuration. Biocuration: Distilling Data into Knowledge. *PLOS Biol.* **2018**, *16*, e2002846,
599 doi:10.1371/journal.pbio.2002846.
- 600 31. Gabrielsen, A.M. Openness and Trust in Data-Intensive Science: The Case of Biocuration. *Med. Health Care Philos.*
601 **2020**, *23*, 497–504, doi:10.1007/s11019-020-09960-5.
- 602 32. Chen, Q.; Britto, R.; Erill, I.; Jeffery, C.J.; Liberzon, A.; Magrane, M.; Onami, J.; Robinson-Rechavi, M.; Sponarova, J.;
603 Zobel, J.; et al. Quality Matters: Biocuration Experts on the Impact of Duplication and Other Data Quality Issues in
604 Biological Databases. *Genomics Proteomics Bioinformatics* **2020**, *18*, 91–103, doi:10.1016/j.gpb.2018.11.006.
- 605 33. Salzberg, S.L. Next-Generation Genome Annotation: We Still Struggle to Get It Right. *Genome Biol.* **2019**, *20*, 92,
606 doi:10.1186/s13059-019-1715-2.
- 607 34. Ejigu, G.F.; Jung, J. Review on the Computational Genome Annotation of Sequences Obtained by Next-Generation
608 Sequencing. *Biology* **2020**, *9*, 295, doi:10.3390/biology9090295.
- 609 35. Zerbino, D.R.; Frankish, A.; Flicek, P. Progress, Challenges, and Surprises in Annotating the Human Genome. *Annu.*
610 *Rev. Genomics Hum. Genet.* **2020**, *21*, 55–79, doi:10.1146/annurev-genom-121119-083418.
- 611 36. Scalzitti, N.; Jeannin-Girardon, A.; Collet, P.; Poch, O.; Thompson, J.D. A Benchmark Study of Ab Initio Gene
612 Prediction Methods in Diverse Eukaryotic Organisms. *BMC Genomics* **2020**, *21*, 293, doi:10.1186/s12864-020-6707-9.

- 613 37. Poux, S.; Magrane, M.; Arighi, C.N.; Bridge, A.; O'Donovan, C.; Laiho, K. Expert Curation in UniProtKB: A Case
614 Study on Dealing with Conflicting and Erroneous Data. *Database J. Biol. Databases Curation* **2014**, *2014*, bau016,
615 doi:10.1093/database/bau016.
- 616 38. Prodocimi, F.; Linard, B.; Pontarotti, P.; Poch, O.; Thompson, J.D. Controversies in Modern Evolutionary Biology:
617 The Imperative for Error Detection and Quality Control. *BMC Genomics* **2012**, *13*, 5, doi:10.1186/1471-2164-13-5.
- 618 39. Meyer, C.; Scalzitti, N.; Jeannin-Girardon, A.; Collet, P.; Poch, O.; Thompson, J.D. Understanding the Causes of
619 Errors in Eukaryotic Protein-Coding Gene Prediction: A Case Study of Primate Proteomes. *BMC Bioinformatics* **2020**,
620 *21*, 513, doi:10.1186/s12859-020-03855-1.
- 621 40. Zhang, D.; Guelfi, S.; Garcia-Ruiz, S.; Costa, B.; Reynolds, R.H.; D'Sa, K.; Liu, W.; Courtin, T.; Peterson, A.; Jaffe,
622 A.E.; et al. Incomplete Annotation Has a Disproportionate Impact on Our Understanding of Mendelian and
623 Complex Neurogenetic Disorders. *Sci. Adv.* **2020**, *6*, eaay8299, doi:10.1126/sciadv.aay8299.
- 624 41. Nagy, A.; Patthy, L. MisPred: A Resource for Identification of Erroneous Protein Sequences in Public Databases.
625 *Database* **2013**, *2013*, bat053, doi:10.1093/database/bat053.
- 626 42. Evans, T.; Loose, M. AlignWise: A Tool for Identifying Protein-Coding Sequence and Correcting Frame-Shifts. *BMC*
627 *Bioinformatics* **2015**, *16*, 376, doi:10.1186/s12859-015-0813-8.
- 628 43. Drăgan, M.-A.; Moghul, I.; Priyam, A.; Bustos, C.; Wurm, Y. GeneValidator: Identify Problems with Protein-Coding
629 Gene Predictions. *Bioinformatics* **2016**, *32*, 1559–1561, doi:10.1093/bioinformatics/btw015.
- 630 44. Vogel, C.; Berzuini, C.; Bashton, M.; Gough, J.; Teichmann, S.A. Supra-Domains: Evolutionary Units Larger than
631 Single Protein Domains. *J. Mol. Biol.* **2004**, *336*, 809–823, doi:10.1016/j.jmb.2003.12.026.
- 632 45. McLaughlin, W.A.; Chen, K.; Hou, T.; Wang, W. On the Detection of Functionally Coherent Groups of Protein
633 Domains with an Extension to Protein Annotation. *BMC Bioinformatics* **2007**, *8*, 390, doi:10.1186/1471-2105-8-390.
- 634 46. Bernardes, J.; Zaverucha, G.; Vaquero, C.; Carbone, A. Improvement in Protein Domain Identification Is Reached
635 by Breaking Consensus, with the Agreement of Many Profiles and Domain Co-Occurrence. *PLOS Comput. Biol.* **2016**,
636 *12*, e1005038, doi:10.1371/journal.pcbi.1005038.
- 637 47. Menichelli, C.; Gascuel, O.; Bréhélin, L. Improving Pairwise Comparison of Protein Sequences with Domain Co-
638 Occurrence. *PLOS Comput. Biol.* **2018**, *14*, e1005889, doi:10.1371/journal.pcbi.1005889.
- 639 48. Monteiro, P.T.; Oliveira, J.; Pais, P.; Antunes, M.; Palma, M.; Cavalheiro, M.; Galocha, M.; Godinho, C.P.; Martins,
640 L.C.; Bourbon, N.; et al. YEASTRACT+: A Portal for Cross-Species Comparative Genomics of Transcription
641 Regulation in Yeasts. *Nucleic Acids Res.* **2020**, *48*, D642–D649, doi:10.1093/nar/gkz859.
- 642 49. Weirauch, M.T.; Yang, A.; Albu, M.; Cote, A.G.; Montenegro-Montero, A.; Drewe, P.; Najafabadi, H.S.; Lambert,
643 S.A.; Mann, I.; Cook, K.; et al. Determination and Inference of Eukaryotic Transcription Factor Sequence Specificity.
644 *Cell* **2014**, *158*, 1431–1443, doi:10.1016/j.cell.2014.08.009.
- 645 50. Si, J.; Zhao, R.; Wu, R. An Overview of the Prediction of Protein DNA-Binding Sites. *Int. J. Mol. Sci.* **2015**, *16*, 5194–
646 5215, doi:10.3390/ijms16035194.
- 647 51. Lambert, S.A.; Yang, A.W.H.; Sasse, A.; Cowley, G.; Albu, M.; Caddick, M.X.; Morris, Q.D.; Weirauch, M.T.; Hughes,
648 T.R. Similarity Regression Predicts Evolution of Transcription Factor Sequence Specificity. *Nat. Genet.* **2019**, *51*, 981–
649 989, doi:10.1038/s41588-019-0411-1.
- 650

651 **Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual au-
652 thor(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to
653 people or property resulting from any ideas, methods, instructions or products referred to in the content.