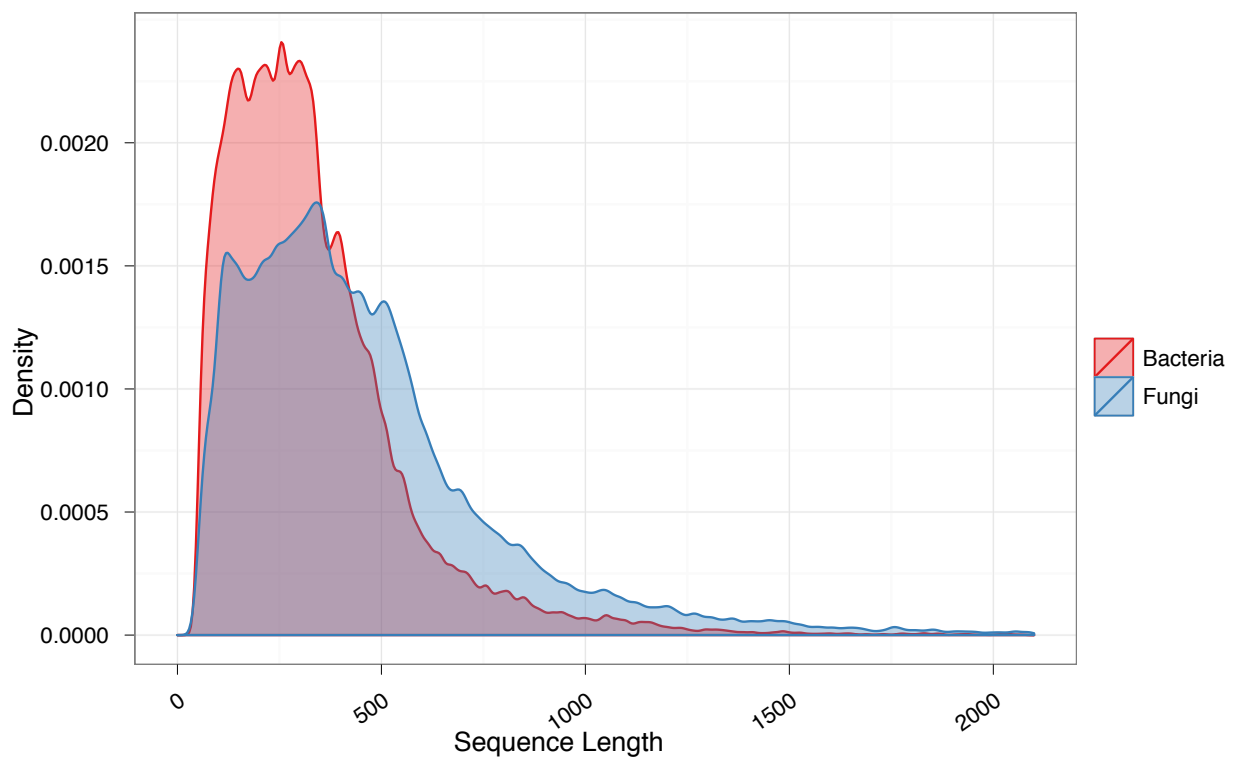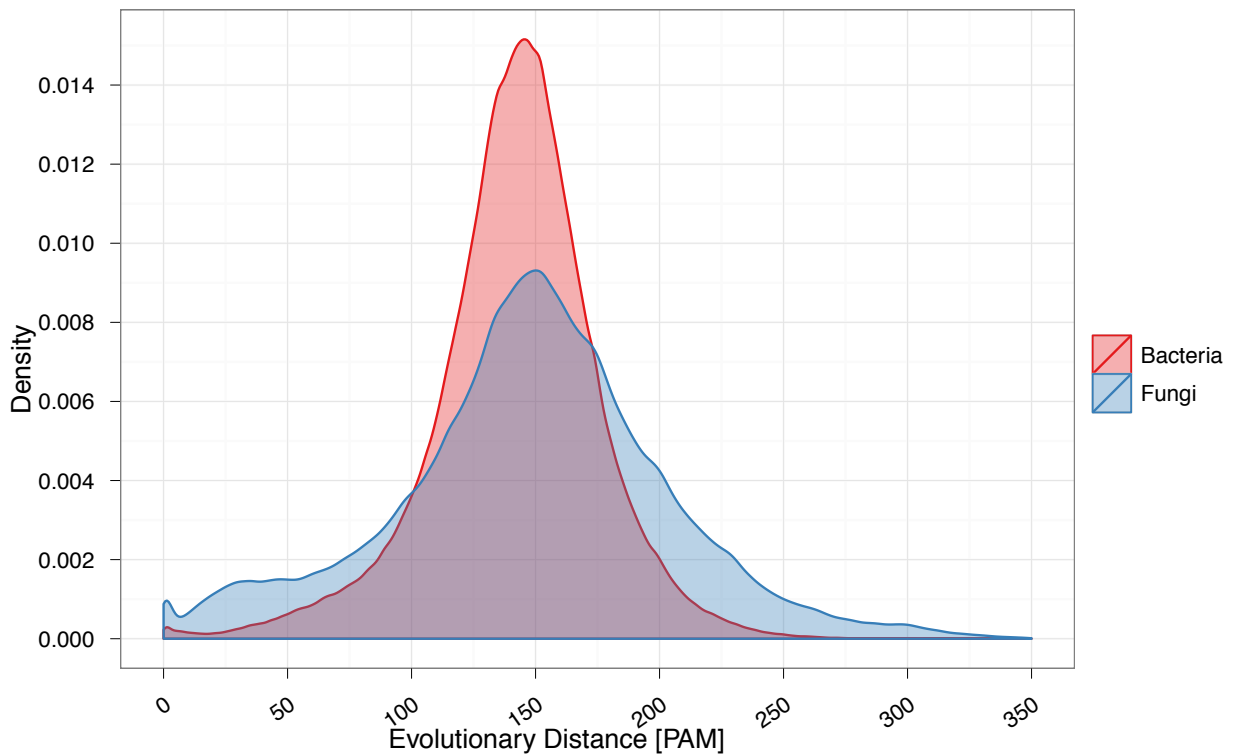# Supplementary materials for

## Speeding up all-against-all protein comparisons while maintaining sensitivity by considering subsequence-level homology

Lucas D. Wittwer, Ivana Piližota et al.

Supplementary figure 1: **Distribution of sequence length in bacteria and fungi datasets (in number of amino acids).**

Supplementary figure 2: **Distribution of evolutionary distances among inferred homologous pairs (based on full all-against-all in bacteria and fungi datasets).**

**Case studies of two missing homologs:**

To illustrate the nature of missing homologs, we provide two detailed descriptions of high-scoring homologous pairs that are missed by the new approach (1 representative, subsequence homology). In both cases, one sequence is added to a cluster with an alignment score slightly above the threshold while the other is not added due to a score just below the threshold:

1) Bacteria: Sequences CHIPD1706 and CHIPD2153 (we use OMA IDs unless stated otherwise) have an alignment score of 2238.183 (estimated PAM distance: 40.3). CHIPD1706 is member of a cluster with CHIPD533 as representative because the score 141.248 is above the threshold (137.75). CHIPD1706 however is not part of the cluster because the alignment score with the representative is 117.317 only and thus below the threshold. Supplementary figure 3 depicts a multiple sequence alignment of the these three sequences, and additionally three other cluster members. Supplementary figure 4 depicts a phylogenetic tree of the sequences and confirms that the terminal branch of CHIPD2153 (outside) is slightly longer than that of CHIPD1706 (inside).

CHIPD2153_(outside_cluster)   1  ____MKRPLLLLLLVVLSQGLH__AQQQPHYTQY ILNPFIINPAVAG   48
CHIPD1706_(inside_cluster)    1  ____MQLKGMITAIILVLALPVQVLAQQQPHYTQYVLNTFIINPAVAG    48
CHIPD2136                     1  __MYTTKKWFTVVLLLCLAAGSR__AQQSVQFSQYIFNGLAINPAYAG    48
CHIPD533_(representative)     1  MK__NKNILFVVGTVLIALMPGWVKAQVDPHFSQYYAYPMWLNPALTG    48
CHIPD1704                     1  MRTFTKALLML___CLVGLTGKKLQASDPHFSGYYVYPAWLNPAMTG     48
CHIPD1170                     1  ____MKKVLLFFTIALYLMNPARAQ___DPHFSQFFASPLTLNPAMTG    48

CHIPD2153_(outside_cluster)  49  _IENYWDVKASHRHQWTGLNGAPVTT_____YLTVHGPLRKSDYPVAS    96
CHIPD1706_(inside_cluster)   49  _IENYWDVKASHRHQWTGVNGSPVTT_____YLTIHGPLRKTDYPQAS    96
CHIPD2136                    49  _YKDVLHLNASYRQQWTGLEGAPRTG_____SISLDGPLNRGN_____   96
CHIPD533_(representative)    49  IVDGDYRVSANYRNQWVNI_GKPFST_____VGVSFDAAAAN_____  96
CHIPD1704                    49  VFDGDYRVSAIYRSQWGSV_SSPFKT_____YGIAGEVKTNN_____  96
CHIPD1170                    49  LFSGDFRVSGNYRSQWSSI_STPFTTGTAAVDFGILKNVLNYT_____  96

CHIPD2153_(outside_cluster)  97  VTGLTPPGDNPRGRAYWQEYTTPPAHAGVGMTILNDKTGPLNRFSISA   144
CHIPD1706_(inside_cluster)   97  ATGFNPEGSNPRGKAYWETYTAPPSHAGAGVTILNDKTGPLSRFSFSG   144
CHIPD2136                    97  _____KDANVGLGIQAMMDNLGPQSAISLYA   144
CHIPD533_(representative)    97  _____NIGVGLNITNMSAGDAGYNYLNAMA   144
CHIPD1704                    97  _____NINFGASVLNQKAGDGGYNYTTAYG   144
CHIPD1170                    97  _____DI_WGVGVMAMYDRTGGGALTSTYL   144

CHIPD2153_(outside_cluster) 145  TYA___HH_IPLSTRLS__VSGGISVGMQSVSVDAAKLQFQQPGDPV_   192
CHIPD1706_(inside_cluster)  145  TYA___HH_IGIAPGTS__LSAGISVGAQRISLDATALEFQNPSDPA_   192
CHIPD2136                   145  SYA___YR_IRLDEEDTRRLCFGLGVGATQYGMDGKDLIYETNGDRI_   192
CHIPD533_(representative)   145  SFS___YRGVRFGETGTSQLVFGIQAGMINRKIDPAKWQLGSQYDPVM   192
CHIPD1704                   145  SAS___YTGVRLGAFEQHRLVFGLQMGLIQRRFDPAKLHFGDQWNPIT   192
CHIPD1170                   145  SFSTAYHKG__LDPEGNHTLAVGLQATLVQKRLDQTKLIFENQIDN_N   192

CHIPD2153_(outside_cluster) 193  _____VASSALLNKWRP___EVNAGLLLY_____GPDFYLGAAAQNI   240
CHIPD1706_(inside_cluster)  193  _____IASSTVLGRWRP___DVSAGLYLY_____SSQYFAGISAQNV    240
CHIPD2136                   193  _____IPDGSA_KATTP___DARVGIYYY_____TPSVYIGVSVLDL   240
CHIPD533_(representative)   193  GFDPSKPSGENISTTSSNSFDAAAGVMFFDGNPNHQFNPFAGFSAGHL   240
CHIPD1704                   193  GYNPGQATNDMFNKTTSAAFDAGAGVLYYDAQPGKKYNLYGGFSVMHL   240
CHIPD1170                   193  GYNPAIPNGETFVNPTISYLDPNIGIL_YNGLVGESSNIYLGASYYHI   240

CHIPD2153_(outside_cluster) 241  VPQ__EVAYDNGKVVGDSL____YRGKLVPHLFFSGGYRLWLSEDFTM   288
CHIPD1706_(inside_cluster)  241  VPS__GIGFDDGKVKGDSI____YRGKLVPHLFATAGYRLWVNEEVSL   288
CHIPD2136                   241  LSKYTSSGYKWRGYTYESI___RRKQ___HLYVTAGYMFNVNDEISL   288
CHIPD533_(representative)   241  TQP_____QDPFVSAGSNKRLPVRYIGHGGTRIKLNEIFSL    288
CHIPD1704                   241  NKP_____SDQF_SATGDARIPMRTTAHAGVRVTISELFSL    288
CHIPD1170                   241  TQP_____TETF_MAQNNNRLTSRYTVHGGGSVPVNGANRI    288

CHIPD2153_(outside_cluster) 289  LPSVMVRLVTA_____APVSYDVNAKFMYRDRMWVGTSYRVK__   336
CHIPD1706_(inside_cluster)  289  LPSVMIKYVTA_____LPVSFDVNAKLQYRDRIWVGGSYRYN__   336
CHIPD2136                   289  KPSVLFKSDFS_____GPAGLDATLMMHIDELLWVGGSYRTNLS   336
CHIPD533_(representative)   289  TPHGLYMRQGNAHETVVGLYGQAYLNEEFD_____FLLGANYRIN__   336
CHIPD1704                   289  TPNVLYMKQGAASEKMLGAYGQYAVSAETD_____VMLGANYRLK__   336
CHIPD1170                   289  HFSAIFMKQSTASEISFG__GAYGFNLNGDDENPTTFYLGSWYRVK__   336

CHIPD2153_(outside_cluster) 337  _____DGFAAMVGV___NISSTINIGYAYDYTTSSLNAVS   384
CHIPD1706_(inside_cluster)  337  _____DGIAAMVGI___NVNATFNIGYSYDYTTSGLNIAS   384
CHIPD2136                   337  VLNKKSIVNNTALDKANAISGILEYYISPKYRIGYSYDYSMNKLAGIQ   384
CHIPD533_(representative)   337  _____DSAIPFAGF__HFKN_FVLGLSYDVNASNLRRLV   384
CHIPD1704                   337  _____DAFSVYTGV___SYKS_FMLGLSYDVNSSDLGKIT   384
CHIPD1170                   337  _____DAINPYLGL___EIGG_FTFGASYDTNVSTLRPAS   384

CHIPD2153_(outside_cluster) 385  KGTHEILIGFLIGNR_____YGDLCPRNNF                 416
CHIPD1706_(inside_cluster)  385  HGSHEMVLGIMIGNR_____FADLCPRNMW                 416
CHIPD2136                   385  TGSHELSIGILFNSK_____LFSTSNPRYF_                416
CHIPD533_(representative)   385  NGSNSFELSLSFISRKKKVYSEENFFCPRL__                 416
CHIPD1704                   385  RGNNSFELSLSFIGRKSVKTPAGDFVCPRL__                 416
CHIPD1170                   385  NYRGGIELSLIYIHRRN_EGSKYRTLCPRF__                 416
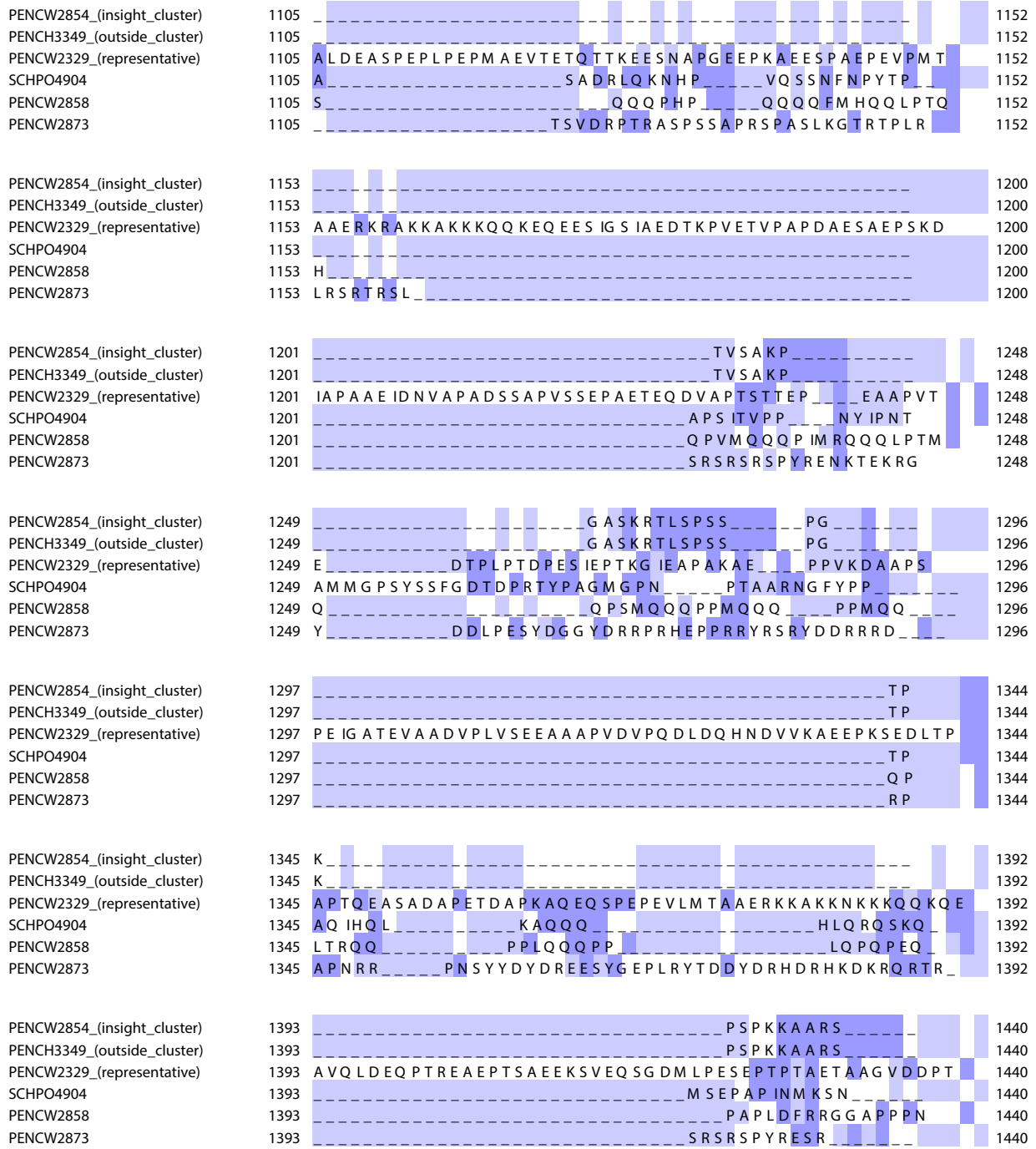
Supplementary figure 3: **Case study #1: Multiple sequence alignment of the cluster to which CHIPD2153 should be included to recover the missing pair CHIPD1706-CHIPD2153.** The corresponding tree is provided in Supplementary figure 4. Alignments drawn with JalView (Waterhouse et al. 2009 DOI:10.1093/bioinformatics/btp033)

Supplementary figure 4: **Case study #1: Distance tree of the cluster and the missing sequence (CHIPD2153).** The corresponding MSA is provided in Supplementary figure 3.

2) Fungi: The homology between sequences PENCW2854 and PENCH3349 is missed despite them being nearly identical (alignment score of 5606.4 and estimated distance of 0.38 PAM units). PENCH2854 is member of two clusters—the first cluster with representative PENCH2329 (score 138.9) and the second cluster with representative PENCW2605 (score 146.6). The alignment scores of PENCH3349 with the two representatives are below the threshold (124.6 and 128.9 respectively). Supplementary figure 5 and 6 provide representative subsets of the multiple sequence alignments for the two clusters.

Furthermore, we note that both clusters are very large: each contain >1000 sequences but only a small fraction of all member pairs are significant (7.8% and 26.96%). As mentioned in the Discussion section, splitting such clusters might improve the sensitivity and runtime of the method.

```
PENCW2854_(insight_cluster)  1105  _____  1152
PENCH3349_(outside_cluster)  1105  _____  1152
PENCW2329_(representative)   1105  A L D E A S P E P L P E P M A E V T E T Q T T K E E S N A P G E E P K A E E S P A E P E V P M T  1152
SCHPO4904                    1105  A _____ S A D R L Q K N H P _____ V Q S S N F N P Y T P __  1152
PENCW2858                    1105  S _____ Q Q Q P H P P ____ Q Q Q Q F M H Q Q L P T Q  1152
PENCW2873                    1105  _____ T S V D R P T R A S P S S A P R S P A S L K G T R T P L R  1152

PENCW2854_(insight_cluster)  1153  _____  1200
PENCH3349_(outside_cluster)  1153  _____  1200
PENCW2329_(representative)   1153  A A E R K R A K K A K K K Q Q K E Q E E S I G S I A E D T K P V E T V P A P D A E S A E P S K D  1200
SCHPO4904                    1153  _____  1200
PENCW2858                    1153  H _____  1200
PENCW2873                    1153  L R S R T R S L _____  1200

PENCW2854_(insight_cluster)  1201  _____ T V S A K P _____  1248
PENCH3349_(outside_cluster)  1201  _____ T V S A K P _____  1248
PENCW2329_(representative)   1201  I A P A A E I D N V A P A D S S A P V S S E P A E T E Q D V A P T S T T E P _____ E A A P V T  1248
SCHPO4904                    1201  _____ A P S I T V P P ____ N Y I P N T  1248
PENCW2858                    1201  _____ Q P V M Q Q Q P I M R Q Q Q L P T M  1248
PENCW2873                    1201  _____ S R S R S R S P Y R E N K T E K R G  1248

PENCW2854_(insight_cluster)  1249  _____ G A S K R T L S P S S _____ P G _____  1296
PENCH3349_(outside_cluster)  1249  _____ G A S K R T L S P S S _____ P G _____  1296
PENCW2329_(representative)   1249  E _____ D T P L P T D P E S I E P T K G I E A P A K A E ___ P P V K D A A P S  1296
SCHPO4904                    1249  A M M G P S Y S S F G D T D P R T Y P A G M G P N _____ P T A A R N G F Y P P _____  1296
PENCW2858                    1249  Q _____ Q P S M Q Q Q P P M Q Q Q ____ P P M Q Q _____  1296
PENCW2873                    1249  Y _____ D D L P E S Y D G G Y D R R P R H E P P R R Y R S R Y D D R R R D ____  1296

PENCW2854_(insight_cluster)  1297  _____ T P  1344
PENCH3349_(outside_cluster)  1297  _____ T P  1344
PENCW2329_(representative)   1297  P E I G A T E V A A D V P L V S E E A A A P V D V P Q D L D Q H N D V V K A E E P K S E D L T P  1344
SCHPO4904                    1297  _____ T P  1344
PENCW2858                    1297  _____ Q P  1344
PENCW2873                    1297  _____ R P  1344

PENCW2854_(insight_cluster)  1345  K _____ ___  1392
PENCH3349_(outside_cluster)  1345  K _____ ___  1392
PENCW2329_(representative)   1345  A P T Q E A S A D A P E T D A P K A Q E Q S P E P E V L M T A A E R K K A K K N K K K Q Q K Q E  1392
SCHPO4904                    1345  A Q I H Q L _____ K A Q Q Q _____ H L Q R Q S K Q _  1392
PENCW2858                    1345  L T R Q Q _____ P P L Q Q Q P P _____ L Q P Q P E Q _  1392
PENCW2873                    1345  A P N R R _____ P N S Y Y D Y D R E E S Y G E P L R Y T D D Y D R H D R H K D K R Q R T R _  1392

PENCW2854_(insight_cluster)  1393  _____ P S P K K A A R S _____  1440
PENCH3349_(outside_cluster)  1393  _____ P S P K K A A R S _____  1440
PENCW2329_(representative)   1393  A V Q L D E Q P T R E A E P T S A E E K S V E Q S G D M L P E S E P T P T A E T A A G V D D P T  1440
SCHPO4904                    1393  _____ M S E P A P I N M K S N ___  1440
PENCW2858                    1393  _____ P A P L D F R R G G A P P P N  1440
PENCW2873                    1393  _____ S R S R S P Y R E S R _____  1440
```

Supplementary figure 5: **Case study #2: Representative extract of the multiple sequence alignment of the first cluster to which PENCH3349 should be included to recover the missing pair PENCW2854-PENCH3349.** Alignments drawn with JalView (Waterhouse et al. 2009 DOI:10.1093/bioinformatics/btp033)

```
PENCW2858                    289  N F A K D D D S S S D I D E A D G Y S V A Q F Q R T M N P A F R T S S P _ _ _ _ Q P S T F D S   336
PENCW2854_(inside_cluster)   289  A S T P P P P S P S T L _ _ _ _ _ _ _ _ _ _ _ R _ _ V P _ _ _ R A P _ _ _ _ R H G A K Y D D   336
PENCH3349_(outside_cluster)  289  A S T P P P P S P S T L _ _ _ _ _ _ _ _ _ _ _ R _ _ V P _ _ _ _ R A P _ _ R H G A K Y D D     336
PENCW2605_(representative)   289  Y I V K D P S A P P A L _ _ _ _ _ _ _ _ _ Q Y A P P _ _ _ Q R P M V H Y G H P N N F Q Q       336
PENCW2911                    289  G V T P D _ S T T A A L _ _ _ _ _ _ _ _ _ _ A A N V P _ _ _ _ K E P _ _ _ _ _ N G D L           336
SCHPO3779                    289  A A T N A P Q Q H P Q L _ _ _ _ _ _ _ _ _ _ Q R M M P I L S S N Q P I Q Q L P L P N Q A S P     336

PENCW2858                    337  H H D _ P H S D L A A Q M G H P S P P V V N N A P P P Q Q Q P L P S Q Q Q P H P Q Q Q Q F M H Q   384
PENCW2854_(inside_cluster)   337  Y E P Y P T R Y S A R L A G Q R G S R V A Q T T P P P R H A T V S A K _ _ P G A S K R T L S P     384
PENCH3349_(outside_cluster)  337  Y E P Y P T R Y S A R L A G Q R G S R V A Q T T P P P R H A T V S A K _ _ P G A S K R T L S P     384
PENCW2605_(representative)   337  Y R P Y P T P _ _ _ _ _ S N Q Q R P A Q Y S P A P _ P G R P G Y T G S _ _ _ H P S P Q Q A R G P   384
PENCW2911                    337  P G S F P E T _ _ _ _ _ P G Q E S E Q T F S V A P I P A S G G Y _ _ _ _ _ _ _ _ _ _ _ _ _         384
SCHPO3779                    337  Y I P V P L Q _ _ _ _ _ Q Q Q Q S Q P Q _ _ _ _ _ _ _ _ _ _ _ _ _ Q _ _ _ Q P Q Q Q Q H Q Q P   384

PENCW2858                    385  Q L P T Q H Q P V M Q Q Q P I M R Q Q Q L P T M Q Q P S M _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ Q Q Q   432
PENCW2854_(inside_cluster)   385  S S P G T P K P S P K K A A R S H R N A A S R I S P F D S _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _     432
PENCH3349_(outside_cluster)  385  S S P G T P K P S P K K A A R S H R N A A S R I S P F D S _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _     432
PENCW2605_(representative)   385  N Q G N A P S P Q G G Q K P P G Q K G Q P A K P S P D P V I Q M L A T R A A A D P E L K A L M     432
PENCW2911                    385  _ _ G N P V S L K P G E E V P H H E T L H _ _ _ _ _ N N T V D S N A T L D K E S Y E K G Q T L   432
SCHPO3779                    385  Q Q P Q P P Q _ Q P L Q Q Q Q Q Q R Q L H S G I Q Q P V S T I V S _ _ _ _ _ _ _ _ Q N G T Y Y   432

PENCW2858                    433  P P M Q Q Q P P M Q Q Q P L T R Q Q P P L Q Q Q P P L Q P Q P E Q P A P L D F R R G G A P P P N   480
PENCW2854_(inside_cluster)   433  _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ E L S L P R T S A H H I P Q T _ _ _ _ _ _ _ _ _ _           480
PENCH3349_(outside_cluster)  433  _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ E L S L P R T S A H H I P Q T _ _ _ _ _ _ _ _ _ _           480
PENCW2605_(representative)   433  R V V A S S Q A S Q E Q L R A F _ _ _ _ _ _ Q A H I D E L _ _ N A I I K S _ _ _ _ _ _ R _ _ _   480
PENCW2911                    433  P L G A G T Q N T D T N P G A F _ _ _ _ _ _ _ A F A I P P V T N N M I P E S _ _ _ _ _ _ S L P I   480
SCHPO3779                    433  S I P A V N H P M A G Q P I A I A P V P A P N Q A A L P P I P P Q A L P A N G T P N T L A S P V   480

PENCW2858                    481  Y D P E Q H G E I G A V P H N A Y P T _ _ _ _ _ _ _ _ D G M T M F C R A G P P S E R S S A T S A   528
PENCW2854_(inside_cluster)   481  _ _ _ _ _ _ _ _ S A _ _ _ _ E Q A L P T _ _ _ _ _ _ _ _ P A K T P S K K K I I N S D S S T S R S L   528
PENCH3349_(outside_cluster)  481  _ _ _ _ _ _ _ _ S A _ _ _ _ E Q A L P T _ _ _ _ _ _ _ _ P A K T P S K K K I I N S D S S T S R S L   528
PENCW2605_(representative)   481  E Q Q E Q R Q Q S _ _ _ _ S A G Q P L _ _ _ _ _ A Q Q P T P T S Q T Q T Q K Q Q D S G T K S V S   528
PENCW2911                    481  G G P A Q R A A T _ _ _ _ S A T E P T Y T G G P A Q Q A A M T D P G Y H I Q S A G P N S T T A A   528
SCHPO3779                    481  T L P A A N S A V _ _ _ _ Q N A Q P V _ _ _ _ _ _ _ _ P M T S S P A M A V V P Q N K T A A T S T   528

PENCW2858                    529  Y R P S S R D S Q S E V S N P T S I S S V E _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ E T   576
PENCW2854_(inside_cluster)   529  F _ _ _ _ _ _ _ _ _ _ _ _ _ P S T _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ S T     576
PENCH3349_(outside_cluster)  529  F _ _ _ _ _ _ _ _ _ _ _ _ _ P S T _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ S T     576
PENCW2605_(representative)   529  L E K _ _ _ _ _ _ _ _ _ S P P T S _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ Q T       576
PENCW2911                    529  L A A _ _ _ _ _ _ _ _ _ G V P L E _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ S K       576
SCHPO3779                    529  L A A Q Q G A N V L P P N A P E S V R H L I S L N E E T W I Q I G R L A E L F D D Q D K A L S A   576

PENCW2858                    577  P T K K S V V K P A S S A P V P S S _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ G D D       624
PENCW2854_(inside_cluster)   577  S K R R K M D N P K Y S A E T P A I _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _       624
PENCH3349_(outside_cluster)  577  S K R R K M D N P K Y S A E T P A I _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _       624
PENCW2605_(representative)   577  P S K G S K Q S P V K T E A Q P Q T P _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _         624
PENCW2911                    577  E K Q T N G D K P V _ _ E E V P Q R _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _ _         624
SCHPO3779                    577  Y E S A L R Q N P Y S I P A M L Q I A T I L R N R E Q F P L A I E Y Y Q T I L D C D P K Q G E I   624
```
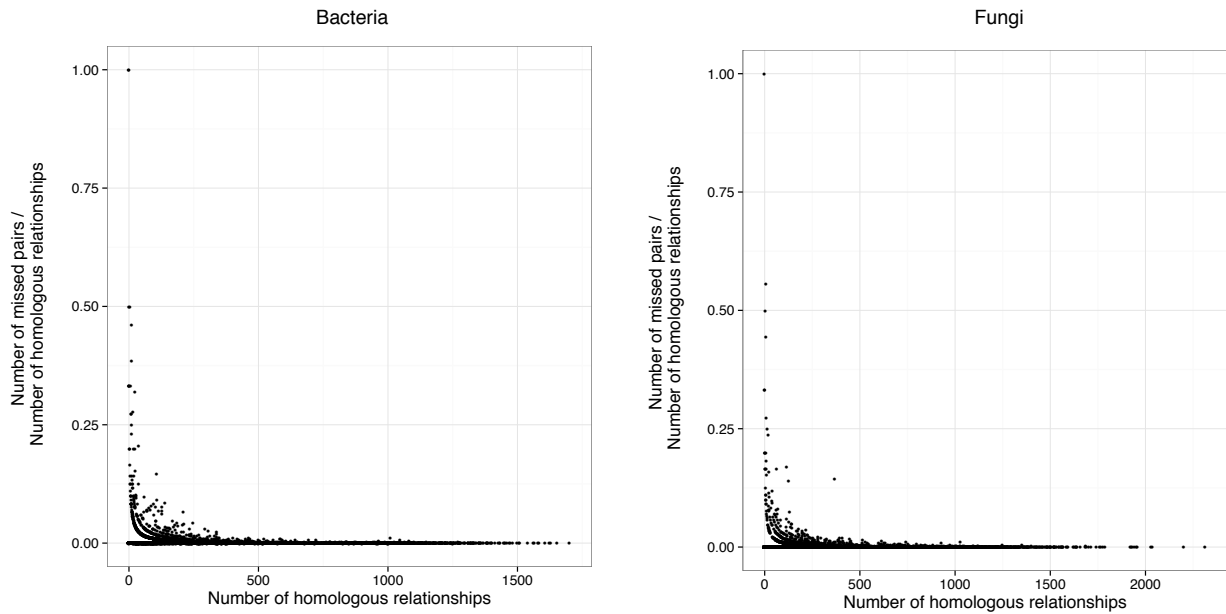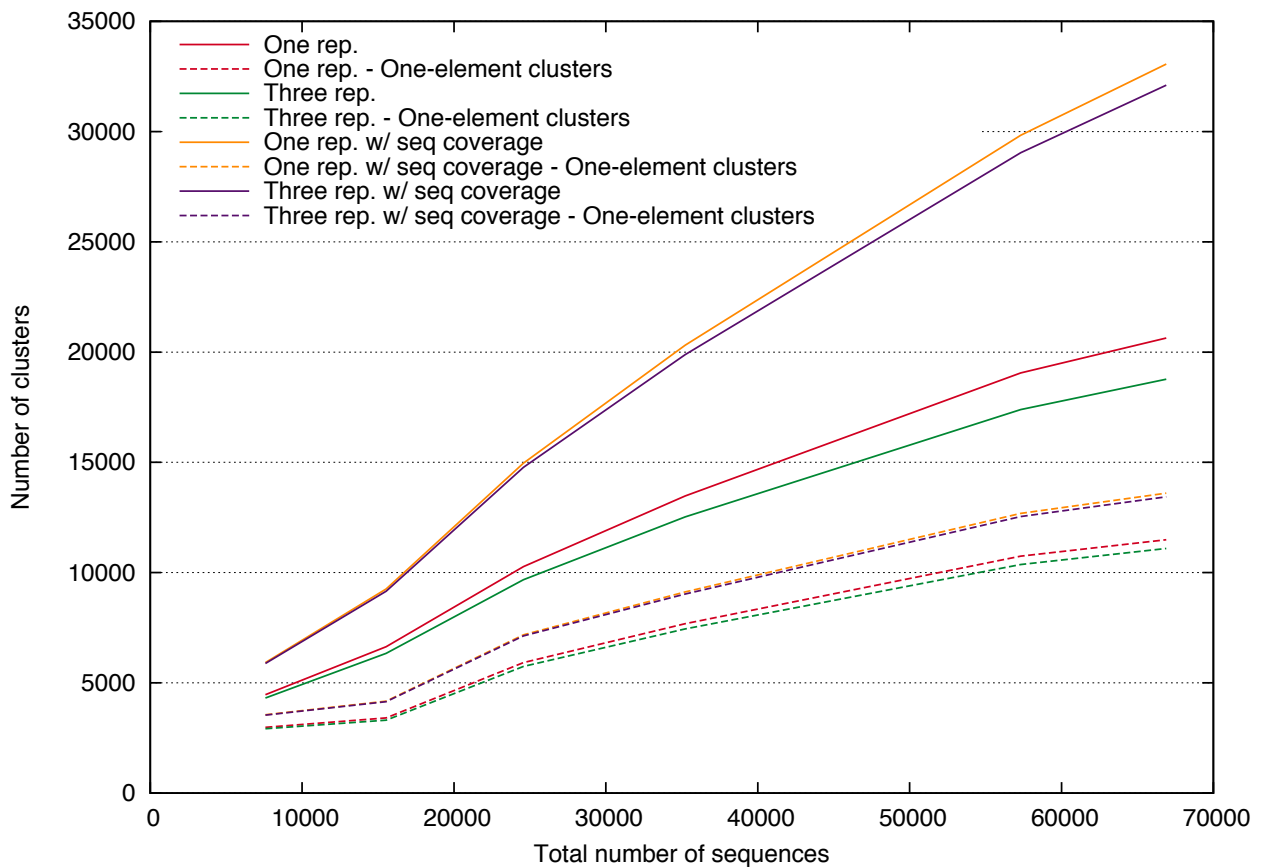
Supplementary figure 6: **Case study #2: Representative extract of the multiple sequence alignment of the second cluster to which PENCH3349 should be included to recover the missing pair PENCW2854-PENCH3349.** Alignments drawn with JalView (Waterhouse et al. 2009 DOI:10.1093/bioinformatics/btp033)
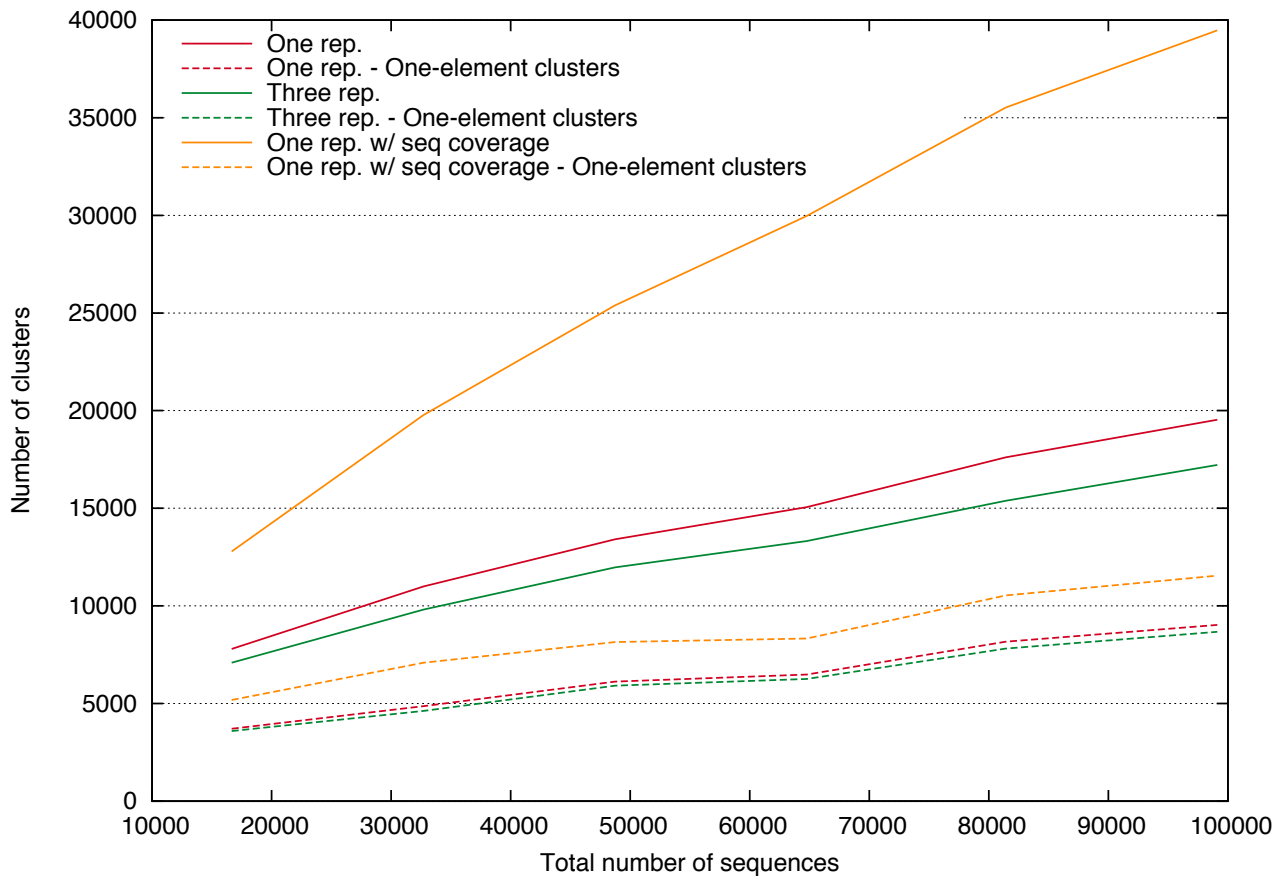
Supplementary figure 7: **Fraction of missing homologous relationships for each gene, as a function of the number of homologous relationships, for the full bacteria (14 proteomes) and fungi (12 proteomes) datasets with one representative and subsequence homology.** In large gene families (>100 members), virtually all homologous relationships are recovered compared to full all-against-all. This is also largely true for small families. The plot suggests that errors are distributed quite evenly across all types of genes.



Supplementary figure 8: **Growth of number of clusters on bacteria dataset.** No tapering is observed in the growth in the number of clusters generated by the new method

Supplementary figure 9: **Growth of number of clusters on fungi dataset.** No tapering is observed in the growth in the number of clusters generated by the new method.

Supplementary figure 10: **Distribution of cluster size for the full bacteria (14 proteomes) and fungi (12 proteomes) datasets with one representative and subsequence homology.** The distribution is heavily skewed toward small clusters.



Supplementary figure 11: **Histogram of the number of clusters overlapping with each cluster (top row) and of the number of clusters in which each sequence is involved (bottom row) for the full bacteria (14 proteomes) and fungi (12 proteomes) datasets with one representative and subsequence homology.** Some clusters overlap with thousands of other clusters, which suggest potential to merge some of them (see *Discussion*).