

HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment

Michael Remmert, Andreas Biegert, Andreas Hauser & Johannes Söding

Sequence-based protein function and structure prediction depends crucially on sequence-search sensitivity and accuracy of the resulting sequence alignments. We present an open-source, general-purpose tool that represents both query and database sequences by profile hidden Markov models (HMMs): ‘HMM-HMM-based lightning-fast iterative sequence search’ (HHblits; <http://toolkit.genzentrum.lmu.de/hhblits/>). Compared to the sequence-search tool PSI-BLAST, HHblits is faster owing to its discretized-profile prefilter, has 50–100% higher sensitivity and generates more accurate alignments.

Building protein multiple-sequence alignments (MSAs) by iterative sequence searches is of fundamental importance in computational biology, as MSAs are a key intermediate step in the sequence-based prediction of evolutionarily conserved properties, such as tertiary structure, functional sites or interaction interfaces. Sequence profiles and profile hidden Markov models (HMMs) are condensed representations of MSAs that specify for each sequence position the probability of observing each of the 20 amino acids in evolutionarily related proteins. PSI-BLAST¹, the most widely used iterative search tool, progressively refines a query sequence profile by adding statistically significant sequence matches to the profile for the next search iteration. The tools SAM2K (ref. 2) and HMMER3 (ref. 3) use profile HMMs for better sensitivity.

Profile-profile and HMM-HMM alignment are the most sensitive classes of sequence-search methods. They are the methods of choice for identifying and aligning templates for three-dimensional homology modeling⁴. Our HMM-HMM alignment method HHsearch⁵ is used by many of the best protein structure prediction servers, among which is HHpred⁶, the top-ranked server for template-based protein structure prediction in last year’s Critical Assessment of Techniques for Protein Structure Prediction exercise (http://predictioncenter.org/casp9/groups_analysis.cgi?type=server&tbd=on/). However, these methods

are generally too slow for iteratively searching through large sequence databases such as UniProt or NCBI’s nonredundant (nr) database. Here we present HMM-HMM-based lightning-fast iterative sequence search (HHblits), which extends HHsearch to enable fast, iterative sequence searches. The profile-profile alignment prefilter of HHblits reduces the number of full HMM-HMM alignments from many millions to a few thousand, making it faster than PSI-BLAST but still as sensitive as HHsearch (Supplementary Fig. 1).

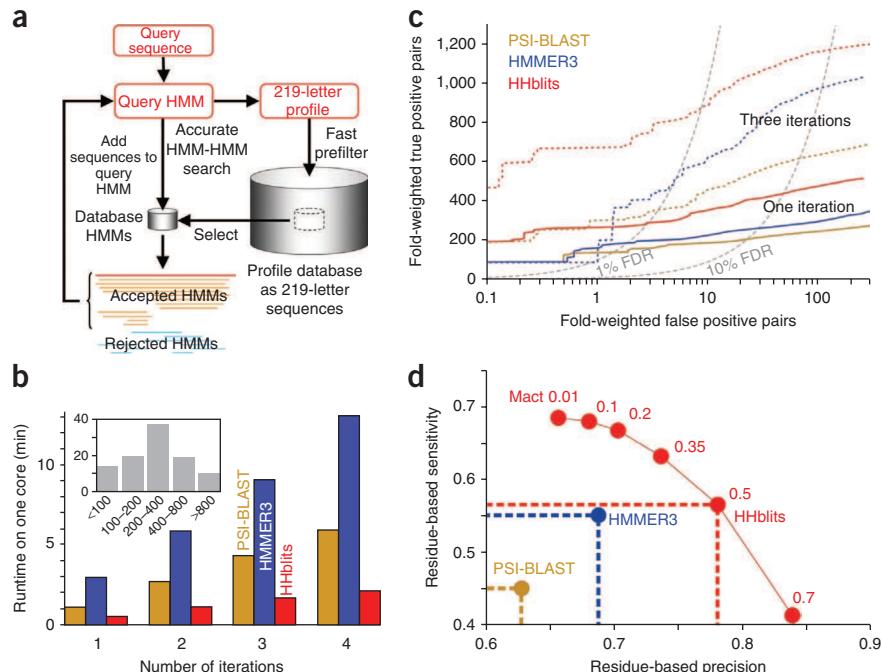
For iterative searches, HHblits needs a database of HMMs that covers the entire sequence space. We devised a very fast method, kClust (M. Hauser, C.E. Mayer and J.S., unpublished data), for clustering large sequence databases down to 20–30% maximum pairwise sequence identity while requiring almost full-length alignability (>80% coverage of longer sequences). This strict coverage criterion enriches for orthologous sequences with the same domain architecture⁷: of the UniProt20 clusters containing more than two Swiss-Prot sequences with enzyme commission numbers, 98.4% had all four enzyme commission digits conserved (Supplementary Fig. 2). kClust is sufficiently fast (~1,000 times faster than BLAST) to allow for regular reclustering of the updated UniProt and nr databases. UniProt20 (the version from July 2011) contained 15 million sequences in 2.6 million HMMs, with an average of 5.5 sequences per cluster.

HHblits first converts the query sequence (or MSA) to an HMM. This is conventionally done by adding pseudocounts of amino acids that are physicochemically similar to the amino acid in the query. In contrast, HHblits calculates pseudocounts that depend on the local sequence context (that is, the 13 positions around each residue). This method had improved the sensitivity and alignment quality of the resulting profile considerably⁸. HHblits then searches the HMM database and adds the sequences from HMMs below a defined expected value (*E* value) threshold to the query MSA, from which the HMM for the next search iteration is built (Fig. 1a and Supplementary Fig. 3). For speed and sensitivity, the prefilter is crucial. The key idea was to implement profile-profile comparison as a sequence-to-profile comparison by discretizing the vectors of 20 amino acid probabilities in each HMM column into an alphabet of 219 letters. Each letter represents a typical profile column (Supplementary Fig. 4). We approximate the database HMMs by sequences over this extended alphabet, ignoring the insertion and deletion probabilities of the HMMs (Supplementary Fig. 5). Before prefiltering, we calculate the score of each query HMM column with each of the 219 letters, which results in a 219-row extended sequence profile. The prefiltering consists of two steps (Supplementary Fig. 3): (i) a very fast gapless local alignment between the extended query profile and the extended database sequences and (ii) a gapped

Gene Center and Center for Integrated Protein Science Munich, Ludwig-Maximilians Universität München, Munich, Germany. Correspondence should be addressed to J.S. (soeding@genzentrum.lmu.de).

RECEIVED 29 JULY; ACCEPTED 1 DECEMBER; PUBLISHED ONLINE 25 DECEMBER 2011; DOI:10.1038/NMETH.1818

Figure 1 | Workflow and benchmark comparison. (a) HHblits can iteratively search for homologous sequences in large databases such as UniProt. The HHblits database is a clustered version in which each set of full-length alignable sequences is represented by an HMM. Sequences from matched HMMs with a statistically significant *E* value are added to the query MSA, from which a new HMM is calculated for the next search iteration. A prefilter reduces the number of full HMM-HMM alignments by ~2,500-fold. (b) Median run times for searches with 100 test sequences through the UniProt or UniProt20 database (the inset shows the test sequence length distribution). (c) True positive pairs (same SCOP fold) compared to false positive pairs (different SCOP fold) for one and three search iterations in an all-against-all comparison. FDR, false discovery rate. (d) Mean fraction of correctly aligned residue pairs out of all structurally alignable pairs (sensitivity) compared to the fraction of correctly aligned pairs out of all the aligned pairs (precision). The parameter mact controls the alignment greediness (**Supplementary Fig. 10**).



local alignment. For step ii, we modified the code from previous work⁹. Each of the two steps allows 1–5% of the sequences to pass. We implemented both filters with streaming SIMD extension 3 (SSE3) instructions, which are available on all modern Intel and Advanced Micro Devices (AMD) central processing units and process 16 single-byte operations per core and clock cycle⁹. The database HMMs whose extended sequences passed the prefilter are aligned to the query HMM, and *E* values are calculated. Statistically significant matches are realigned with a local maximum accuracy algorithm¹⁰.

A single search iteration with HHblits version 2.2.17 through UniProt20 (2.6 million clusters and 15 million sequences) for 100 randomly selected query sequences took a median 31 s and an average 1 min 13 s on a single Xeon 2.9 GHz core (Fig. 1b and **Supplementary Data 1**). For a single search iteration through UniProt (15 million sequences), PSI-BLAST needed 1 min 7 s (median) and 1 min 26 s (average) and HMMER3 needed 2 min 57 s (median) and 5 min 8 s (average). Additional iterations took roughly the same amount of time as the first iteration (**Supplementary Fig. 6**), and therefore overall, HHblits was

about twice (15%) as fast as PSI-BLAST and was 6× (median) and 4× (average) faster than HMMER3.

We compared the sensitivity of HHblits to that of PSI-BLAST and HMMER3 in detecting homologous proteins (to rank true positive, homologous pairs above false positive, unrelated pairs) (Fig. 1c). We performed an all-against-all comparison of 5,287 representative domain sequences from the Structural Classification of Proteins (SCOP) database¹¹. After one iteration, HHblits detected 107% more true positive pairs than PSI-BLAST and 53% more than HMMER3 at 1% false discovery rate, and after three iterations, the improvement was 147% over PSI-BLAST and 69% over HMMER3. We obtained similar values in a receiver operating curve 5 (ROC5) analysis (Online Methods and **Supplementary Fig. 7**). Furthermore, HHblits reported more reliable *E* values than PSI-BLAST (**Supplementary Fig. 8**).

To assess the quality of the pairwise alignments (Fig. 1d), we randomly selected from each SCOP superfamily up to ten pairs of domains with <30% sequence identity and a TM-align (Online Methods) structural similarity score of >0.6 (**Supplementary Data 2**). For each method, we built MSAs for the queries using

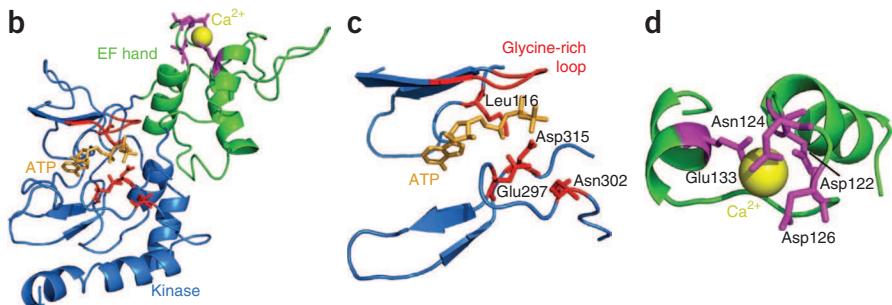
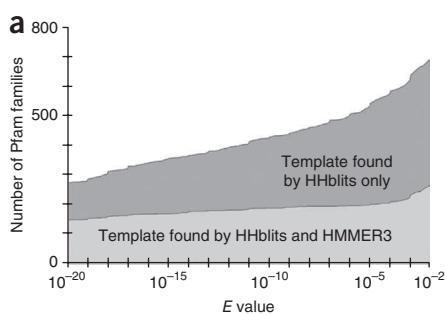


Figure 2 | Structure predictions for Pfam families and the modeling of human Pip49 (also known as FAM69B). (a) Families to which only HHblits and both HHblits and HMMER3 assigned a structural template below a given *E* value. (b) Homology model of human Pip49 kinase domain (blue) with the inserted EF hand (green). (c) Catalytic center showing the conserved residues (red) for protein kinase activity. (d) EF hand insertion with the conserved residues (magenta) for the predicted Ca²⁺-dependent activation.

two search iterations through UniProt and aligned the resulting query MSAs with their corresponding templates. We determined correctly aligned residues through comparison with the structural alignments. Compared to PSI-BLAST and HMMER3, HHblits sensitivity per residue using default parameters (mact 0.5) was 12 and 2 percentage points higher and the precision per residue was 15 and 10 percentage points higher, respectively (**Fig. 1d**). The higher precision of HHblits alignments explains its robustness against homologous overextension (tested on a benchmark with multidomain proteins; **Supplementary Fig. 9**), which is the main cause of corrupted PSI-BLAST alignments¹².

As another measure of MSA quality, we sought to improve the accuracy of PSIPRED¹³ secondary structure prediction by running PSIPRED on MSAs generated by HHblits. Although PSIPRED had been trained on PSI-BLAST MSAs, HHblits MSAs improved the Q3 score (fraction of correctly predicted secondary structure states) for proteins from the PDBselect 2007 dataset (Online Methods) from 80.4% to 81.3% and the secondary structure segment overlap (SOV) score from 77.5% to 78.6% (**Supplementary Table 1**). These results, obtained without training a large parameter set, are among the best achieved at present¹⁴.

A potential drawback of HHblits is the requirement that its databases consist of MSAs and their HMMs instead of single sequences. Although we will regularly update standard HHblits databases such as UniProt20, nr20, PDB, SCOP and Pfam, customized databases, for example databases representing an organism's proteome, will need to be built specifically for HHblits.

To show the utility of HHblits, we predicted structures for Pfam families¹⁵ for which no template is known and also for which no template is known for any family from its Pfam clan (**Fig. 2a**). We jumpstarted two HHblits iterations through UniProt20 with the Pfam seed alignment and then searched the PDB70 database ([ftp://toolkit.genzentrum.lmu.de/HHblits/databases/](http://toolkit.genzentrum.lmu.de/HHblits/databases/)). HHblits assigned templates to 620 families with $E < 10^{-3}$, only 226 of which HMMER3 detected (41 families were found only by HMMER3 and not HHblits) (**Supplementary Table 2**).

As an example of these results, we describe the predictions for Pip49-C, the C-terminal part of the pancreatitis induced protein 49, a Pfam domain of unknown structure and function with a predicted N-terminal transmembrane helix. The 100 best HHblits matches in PDB70 were with protein kinases (best E value of 2×10^{-20}), even though the Pfam MSA is missing 70 N-terminal residues from the kinase domain. An HHblits search started with full-length human Pip49 (also known as FAM69B) (with two iterations through UniProt20 and one iteration through PDB70) detected many protein kinase domains, and, notably, a tandem Ca^{2+} -binding EF hand (E value = 0.09) inserted in the kinase domain. Based on our homology models (**Fig. 2b–d** and **Supplementary Data 3**)

and the conservation of key residues, we predict that Pip49 and its paralog FAM69A are membrane-bound protein kinases in the lumen of the endoplasmic reticulum that are activated by Ca^{2+} through structural rearrangement of their EF hand.

In conclusion, HHblits is an open-source, robust, general-purpose, iterative protein sequence search tool that is faster, considerably more sensitive and produces alignments of much better quality than PSI-BLAST. HHblits has the potential to improve many downstream analysis and prediction methods, such as a *de novo* protein structure prediction method requiring large and accurate MSAs¹⁶.

METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturemethods/>.

Note: Supplementary information is available on the Nature Methods website.

ACKNOWLEDGMENTS

We acknowledge financial support by the Deutsche Forschungsgemeinschaft (grant SFB646) and by a Gastprofessur grant from Ludwig-Maximilians-Universität München financed through the Excellence Initiative of the Bundesministerium für Bildung und Forschung.

AUTHOR CONTRIBUTIONS

M.R. performed research, J.S. initiated and guided research, A.B. generated the profile-column alphabet, A.H. contributed code for fast file access, and M.R. and J.S. wrote the manuscript.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/naturemethods/>.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

1. Altschul, S.F. et al. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
2. Karplus, K., Barrett, C. & Hughey, R. *Bioinformatics* **14**, 846–856 (1998).
3. Eddy, S.R. *Genome Inform.* **23**, 205–211 (2009).
4. Söding, J. & Remmert, M. *Curr. Opin. Struct. Biol.* **21**, 404–411 (2011).
5. Söding, J. *Bioinformatics* **21**, 951–960 (2005).
6. Söding, J., Biegert, A. & Lupas, A.N. *Nucleic Acids Res.* **33**, W244–W248 (2005).
7. Hegyi, H. & Gerstein, M. *Genome Res.* **11**, 1632–1640 (2001).
8. Biegert, A. & Söding, J. *Proc. Natl. Acad. Sci. USA* **106**, 3770–3775 (2009).
9. Farrar, M. *Bioinformatics* **23**, 156–161 (2007).
10. Biegert, A. & Söding, J. *Bioinformatics* **24**, 807–814 (2008).
11. Andreeva, A. et al. *Nucleic Acids Res.* **36**, D419–D425 (2008).
12. Gonzalez, M.W. & Pearson, W.R. *Nucleic Acids Res.* **38**, 2177–2189 (2010).
13. Jones, D.T. *J. Mol. Biol.* **292**, 195–202 (1999).
14. Aydin, Z., Singh, A., Bilmes, J. & Noble, W. *BMC Bioinformatics* **12**, 154 (2011).
15. Finn, R.D. et al. *Nucleic Acids Res.* **38**, D211–D222 (2010).
16. Marks, D.S. et al. *PLoS ONE* **6**, e28766 (2011).

ONLINE METHODS

HHblits server usage in a nutshell. The HHblits server (<http://hhblits.genzentrum.lmu.de/> or <http://toolkit.genzentrum.lmu.de/hhblits/>) takes as input a single sequence or an MSA and iteratively searches through the selected HMM databases (UniProt20, nr20, PDB, SCOP and Pfam) for a specified number of iterations. Two or more iterations only make sense when a database covering the entire sequence space (such as UniProt20 or the nr20) is selected. A larger number of search iterations increases the sensitivity of the alignment but also increases the risk of alignment corruption, for example, through homologous overextension¹². Owing to this trade-off, we recommend between one and four iterations. For optimum reliability it is advisable to first identify domains in the query sequence by performing a single HHblits iteration through the PDB70 database and then to cut the query sequence along domain boundaries into shorter segments, which are less prone to alignment corruption.

When the “realign with MAC” box is checked in HHblits, it calculates the more accurate maximum accuracy (MAC) HMM-HMM alignments after the Viterbi HMM-HMM comparisons. However, the Viterbi alignments are better suited for the calculation of scores, *E* values and probabilities, and, therefore, the results are a combination of Viterbi scores and *E* values and MAC alignments. The MAC threshold parameter ‘mact’ controls the alignment greediness during MAC realignment. At a mact value of 0.35, segments that have an average probability of being correct of below 35% will be omitted. A mact value of 0.01 will generate quasi-global MAC alignments for use in, for example, homology modeling, while still using the local Viterbi alignment for the scoring. When searching with single domains, a mact value of 0.2 can be sufficient; otherwise, higher mact values (for example, 0.5) are recommended. The MSA built by HHblits can be inspected and extracted under the “show alignment” tab on the results page. An MSA of consensus sequences of matched HMMs can be viewed with a Jalview applet on the results page, for example to check for alignment corruption. For more information, see the HHblits server help pages and the HHblits user guide.

HHblits command line usage in a nutshell. HHblits is available as source code and as executable RPM and DPKG packages for most Linux 64 bit platforms, MAC OS X and Berkeley Software Distribution (BSD) Unix at <http://hhblits.genzentrum.lmu.de/>. The command “\$ hhblits -i query.fasta -d /databases/UniProt20 -n 2 -mact 0.01 -oa3m query_msa.a3m” will run two search iterations through the UniProt20 database, starting from the input sequence (or MSA) in query.fasta. The mact value 0.01 generates quasi-global alignments. The human-readable output is written to query.hhr by default (this option can be changed using -o <file>), whereas the resulting MSA is written to query_msa.a3m in a3m format. This format can be transformed to other formats using reformat.pl. Custom databases (such as for a single genome) can be built by generating MSAs for each protein sequence using, for example, two HHblits iterations, adding secondary structure with the Perl script addss.pl and building the HHblits database files using create_db.pl and create cs_db.pl. For more information see the user guide in the HHblits package at <http://hhblits.genzentrum.lmu.de/> or <http://toolkit.genzentrum.lmu.de/hhblits/>.

Fast sequence clustering with kClust. We developed kClust to cluster large sequence databases for use with HHblits in a fraction of the time that would be necessary using BLAST and down to much lower sequence identities (20–30%) than is possible with CD-HIT¹⁷. kClust achieves its high speed and sensitivity with two new algorithms. First, a fast prefilter sums the similarity scores of all similar 6-mers between sequences *Q* and *T*. The score threshold is set stringently such that only $\sim 4 \times 10^{-6} \times L_Q L_T$ chance matches occur between sequences of lengths L_Q and L_T . Thus, the time to compare two sequences is reduced in comparison to classic dynamic programming approaches by a factor of $\sim 2.5 \times 10^5$. A more sensitive comparison is performed in the second step using dynamic programming on the set of similar 4-mers between *Q* and *T*. Here the threshold is set such that chance matches occur with a probability of $\sim 2 \times 10^{-3}$ per 4-mer pair. This allows us to achieve a speedup relative to the SSEARCH implementation of classic Smith-Waterman dynamic programming by a factor ~ 30 . kClust binaries and scripts for automatically generating MSAs from clusters are available at <ftp://toolkit.genzentrum.lmu.de/kClust/>.

Discretized profile-column alphabet. We discretized profile columns into an alphabet of 219 states (the number of printable ASCII characters), where each letter represents a typical profile column. This allows us to approximate any sequence profile by a sequence over this 219-letter extended alphabet. To compare two profiles, we first calculate the score S_{ik} of each query profile column *i* with each of the 219 letters *k*

$$S_{ik} = \log_2 \sum_{a=1}^{20} q_i(a) p_k(a) / f(a)$$

where $q_i(a)$ denotes the query profile at position *i*, $p_k(a)$ is the profile column represented by the letter $k \in \{1, \dots, 219\}$ and $f(a)$ is the background frequency of residue *a*. We thus obtain a 219-row extended sequence profile, which can be aligned to extended sequences representing the other profile using fast, standard dynamic programming techniques. We generated the 219-letter alphabet using the same method that was previously used for learning an optimal set of sequence context profiles⁸, but here we set the window size from 13 to 1 residue. We also set the window weights w_j to 100 to obtain a hard clustering. We initialized the 219 states randomly and maximized the likelihood that the 10 million training sequence profile columns were generated by the 219 profile columns. The best of several trials was used. The 10 million profile columns were randomly sampled from the MSAs in our clustered nonredundant database.

Pre-filtering. In the two prefilter steps, the extended query sequence profile is aligned to the extended database sequences. The first step calculates the score of the largest ungapped alignment. To pass this filter, the score has to be larger than $2.5 + \log_2(L_Q L_T)$ bits, where L_Q and L_T are the lengths of the query profile and database sequence, respectively. The log term is a standard length correction. The second step calculates a Smith-Waterman alignment with affine gap penalties (gap open: 5 bits, gap extend: 1 bit). From the bit score *S*, an approximate *E* value is calculated: $E = N_{db} L_Q L_T \times 2^{-S}$, where N_{db} is the number of sequences per HMMs in the database, and sequences pass if their *E* value is

below the prefilter threshold (E_{pre}) = 1,000. Each filter step leads to a 10- to 100-fold reduction of database sequences.

Both filters were implemented with SSE3 instructions that process 16 single bytes in parallel on 128-bit SIMD units present on each central processing unit (CPU) core. Each byte holds the score in units of 1/4 bits plus an offset of 50, which allows us to represent a score range between -12.5 and +51.5 bits. The algorithms were programmed such that the scores will saturate at 255 on overflow. Because any score larger than 51 bits will always pass the filter, this range is sufficient for prefiltering. The first step processes four or five cells of the dynamic programming matrix per CPU clock cycle, and the second step processes ~1.3 cells per clock cycle. The clustered UniProt database (version from 07/2011) contains 2.6 million sequences of average length 320 cells, and therefore the first prefilter search with a query profile of length 300 through UniProt takes about $300 \times 320 \times 2.6 \times 10^6 / (4.5 \times 2.9 \text{ GHz}) = 18$ s, which is about 25% of the average time needed for the entire HHblits search.

For sequences that pass the two prefilters, we calculate local alignments using SSE3 instructions to restrict the resulting HMM-HMM alignment to the region likely to contain the true alignments. For back-tracing, we need to prevent the score from saturating. Therefore, each score is held in 2 bytes in this step (again, in units of 1/4 bits), which yields a score range of -12.5 bits to +16,371.5 bits. Up to ten suboptimal alignments are extracted by masking all cells at a distance of <150 residues from the previously extracted alignments until the prefilter E value is above the E_{pre} value.

Viterbi alignment and E value calculation. To speed up the time-consuming HMM-HMM alignment steps, all cells with a distance of >200 residues to all alignments identified in the previous step are masked out. An HMM-HMM alignment is performed on the active cells using the Viterbi algorithm from HHsearch. The Viterbi algorithm determines the alignment with the maximum score. Even though it does not yield the most accurate alignments (see the maximum accuracy alignment section below), it yields reliable scores for ranking and P value calculation. From the Viterbi score S , a P value is calculated using an extreme value distribution: $P = 1 - \exp(-\exp(-\lambda(S - \mu)))$. The extreme value distribution parameters μ and λ are estimated from the four features L_Q , L_T , N_Q^{eff} and N_T^{eff} using two standard two-layer neural networks with four hidden nodes each. Here N_Q^{eff} and N_T^{eff} are the numbers of effective sequences in the query and template HMMs, respectively (defined in ref. 5). The Viterbi E value is calculated from the P value using $E = N_{\text{db}} P (E_{\text{pre}}/N_{\text{db}})^{\alpha}$, where $\alpha = 0.4 + 0.02 \times (N_T^{\text{eff}} - 1) \times (1 - 0.1 \times (N_Q^{\text{eff}} - 1))$. The term $(E_{\text{pre}}/N_{\text{db}})^{\alpha}$ is an empirical correction for the correlation between the prefiltering and Viterbi scores ($\alpha = 0$: perfect correlation, $\alpha = 1$: no correlation). The three coefficients for α were optimized to yield accurate E values (Supplementary Fig. 8).

Further speedups. Viterbi alignments are performed in the order of decreasing prefilter E values. We stop the time-consuming HMM-HMM comparisons when very few homologs are likely to have been observed among the last 200 HMM-HMM alignments. A coarse estimate for the probability of a match to be a true homolog is $1/(1+E)$ for a Viterbi E value of E . We average $1/(1+E)$ over the last 200 processed Viterbi alignments and skip all further database HMMs when this average drops below 0.01.

Maximum accuracy alignment. Whereas the Viterbi algorithm calculates the alignment with the best score, the maximum accuracy alignment (proposed in ref. 18) yields the global alignment with the maximum possible accuracy as defined by the sum of probabilities for each residue pair to be correctly aligned

$$\sum_{(i,j) \in \text{alignment}} P(i \text{ aligned to } j) \rightarrow \max$$

We extended this algorithm to local HMM-HMM comparison¹⁰, which produces the local alignment that maximizes the sum of probabilities for each residue pair to be correctly aligned minus the mact penalty

$$\sum_{(i,j) \in \text{alignment}} (P(i \text{ aligned to } j) - \text{mact}) \rightarrow \max$$

With the mact parameter, the alignment greediness can be controlled from nearly global, long, greedy alignments (mact near 0) to very precise and short alignments (mact near 1). To speed up the MAC alignment, cells at a city block distance of >200 from the optimal and all suboptimal Viterbi HMM-HMM alignments are masked.

Adding sequences from significant matches to the query HMM. Sequences from all HMMs below the Viterbi E value inclusion threshold (with a default value of 10^{-3}) are read from the alignment files of the clustered database and aligned to the query MSA according to the HMM-HMM maximum accuracy alignment. The query HMM is calculated from the query MSA.

Parameter optimization. We optimized the parameters (filter thresholds, gap costs, amino acid and transition pseudocount strengths and E value inclusion threshold) on an optimization set that had no member from the same fold as the sequences in the test set (see the sensitivity benchmarks section below). We varied the parameters in discrete steps one after another, performed an all-against-all search on the optimization set and tried to maximize the mean ROC5 value (see below). For the prefilter settings, we chose the best trade off between efficiency and sensitivity.

Sensitivity benchmarks. We filtered the sequences from SCOP 1.73 (ref. 11) to a maximum pairwise sequence identity of 20%. We assigned every fifth fold to the optimization set (1,329 sequences in 215 folds) and the other folds to the test set (5,287 sequences in 862 folds; Supplementary Data 4). SCOP is a hierarchically ordered database of protein domain sequences with known structure. We considered domains from the same fold as true positive, homologous pairs and domains from different folds as false positive, nonhomologous pairs. Exceptions to this were members of Rossman-like folds (c.2-c.5, c.27, c.28, c.30 and c.31) and the four- to eight-bladed β -propellers (b.66-b.70), which are probably related and which we treated as ‘unknown’. To prevent a few large folds from dominating the benchmark⁴, we weighed each hit with the value of one over the number of members in the query SCOP fold (‘fold-weighted true positives and false positives’). All but the last search iteration were performed against the UniProt database. The final iteration of PSI-BLAST and HMMER3 searches were performed against all UniProt and SCOP sequences. For HHblits, the final iteration was performed against the UniProt20/SCOP database, a UniProt20 database to which SCOP test sequences

had been added as singleton clusters: each SCOP sequence from the test set was either mapped to its UniProt20 cluster containing the test sequence or was added as a singleton cluster to UniProt20/SCOP if no matching cluster was found. All pairs of domains were ranked by *E* value for each of the tools, and the number of true positives versus false positives below a given *E* value were plotted. The ROC5 plots in **Supplementary Figures 7d** and **9b** assess how well a method ranks the matched proteins within each search. These plots show the fraction of queries with ROC5 scores above the threshold on the *x* axis. The ROC5 score is the area under the true positive versus false positive ROC curve up to the fifth false positive divided by the area under the optimal ROC curve.

Sensitivity benchmark for multidomain proteins. Because multi-domain protein sequences present particular challenges, such as homologous overextension¹², to iterative sequence search methods, we tested our tools on a benchmark set of multi-domain proteins. For each of the 5,287 sequences in our test set, we searched for a sequence in the nonredundant database that had a BLAST match to the SCOP sequence with an *E* value <10⁻⁴⁰, sequence coverage >95% sequence identity >60% and whose full-length sequence contained at least 100 additional residues. These criteria resulted in 2,343 multidomain proteins. For all extracted multi-domain proteins, we proceeded as described in the previous paragraph (two iterations through UniProt and one iteration through UniProt or UniProt/SCOP, respectively). We counted true positive and false positive pairs only if the alignment covered at least 50 residues of the SCOP domain in the nonredundant query sequence.

Alignment quality. To assess the accuracy of the pairwise alignments (**Fig. 1d**), we chose 4,128 query-template pairs by randomly selecting from each SCOP superfamily up to ten pairs with <30% sequence identity and a structural similarity TM-align¹⁹ score >0.6 (**Supplementary Data 2**). For each method, we built MSAs for the queries using two search iterations through UniProt and aligned the resulting query MSAs with their corresponding templates. For HHblits, we selected the template HMMs from the clustered UniProt20 that contained the SCOP template sequence (using the same procedure as described in section Sensitivity benchmarks). We determined correctly aligned residues by comparison with structural alignments from TM align¹⁹.

Improving PSIPRED secondary structure prediction. We compared the accuracy of the secondary structure prediction of PSIPRED¹³ using the PSIPRED procedure to generate sequence profiles (three iterations of PSI-BLAST on a filtered database), with the accuracy of PSIPRED run on profiles built from MSAs generated by HHblits. As a test set, we used PDBselect 2007 (ref. 20), which contains 3,649 sequences ranging from 30 to 1,040 amino acids in length. We built MSAs for each sequence using two and

three iterations of PSI-BLAST on the nr database filtered with pfilt from the PSIPRED package and using one, two and three iterations of HHblits through UniProt20. HHblits alignments with a diversity around 7.0 were generated by applying hhfilter from the HHblits package with the option ‘-neff 7’. For all MSAs, we performed the PSIPRED procedure with the default parameters and calculated the Q3 and SOV scores based on the known DSSP sequences (mapping E and B to strand, H, G and I to helix, and S, T and C to coil states).

Fold prediction for Pfam. For nearly half of all Pfam families in version 24.0 (5,716 out of 11,913), no structure is known, and the structures of any of the remotely related families in their Pfam clan are also unknown. We generated MSAs for these 5,716 Pfam families by using their seed alignments as input and performing two iterations with HHblits through the UniProt20 database. The PDB70 database of HHpred was searched with the resulting MSAs. For HMMER3, we scanned the PDB70 sequence database with the full HMMER3 models provided by Pfam.

Pip49/FAM69B modeling. We built an MSA for human Pip49/FAM69B (UniProt identifier Q5VUD6) by running two iterations of HHblits through the clustered UniProt database and adding the secondary structure prediction from PSIPRED to this MSA using the script addss.pl from the HHblits package (**Supplementary Data 5**). To identify structural homologs, the PDB database was scanned by HHblits with this MSA with a mact value of 0.2. From the list of PDB matches, we chose as templates a protein kinase with bound ATP (PDB identifier 1RDQ) and a Ca²⁺-bound EF hand (PDB identifier 3C1V). We used the corresponding HHblits alignments to create a homology model with MODELLER²¹ (**Supplementary Data 3**). Although many protein kinases contain EF hands downstream of their kinase domains¹⁵, Pip49 is the first one known in which an EF hand is inserted in the kinase domain, directly after the small N-terminal β sheet. We validated the presence of the EF hand insertion by building an MSA with two iterations of HHblits starting from the presumed inserted sequence and then searching the PDB70 database. This yielded highly significant matches with EF hands (best *E* value = 4 × 10⁻⁵). The previously reported transmembrane helix from residue position 31 to 51 could be confirmed by HMMTOP, MEMSAT-SVM and Phobius. The kinase domain is framed by two short domains with highly conserved cysteines that are likely to form disulfide bonds, which suggests that it resides in the lumen of the endoplasmic reticulum.

17. Li, W. & Godzik, A. *Bioinformatics* **22**, 1658–1659 (2006).
18. Holmes, I. & Durbin, R. *J. Comput. Biol.* **5**, 493–504 (1998).
19. Zhang, Y. & Skolnick, J. *Nucleic Acids Res.* **33**, 2302–2309 (2005).
20. Griep, S. & Hobohm, U. *Nucleic Acids Res.* **38**, D318–D319 (2009).
21. Sali, A. & Blundell, T.L. *J. Mol. Biol.* **234**, 779–815 (1993).

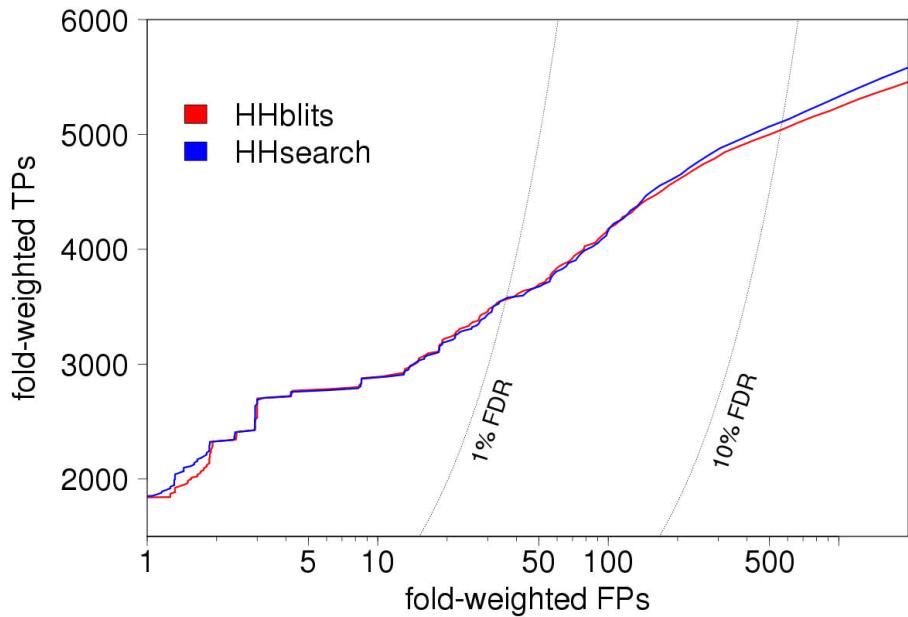
HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment

Michael Remmert, Andreas Biegert, Andreas Hauser & Johannes Söding

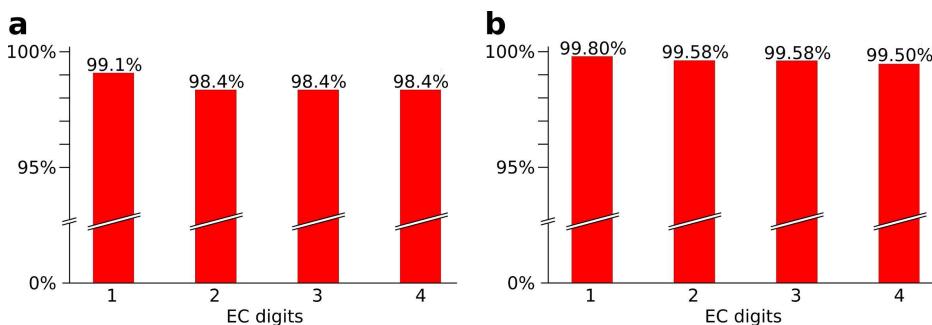
Supplementary Figure 1	Sensitivity / selectivity comparison between HHblits and HHsearch on the SCOP test set for a single search.
Supplementary Figure 2	Sequence clusters in the UniProt20 predominantly contain functionally closely related proteins.
Supplementary Figure 3	Detailed workflow of HHblits.
Supplementary Figure 4	Histogram representation of the amino acid distributions in the 219 profile columns of the extended profile column alphabet.
Supplementary Figure 5	Translation of an alignment into a sequence of column states.
Supplementary Figure 6	Average run times on 100 randomly selected sequences from the nr database, measured on an Intel Xeon X5570 at 2.93 GHz.
Supplementary Figure 7	Sensitivity and selectivity of homology detection.
Supplementary Figure 8	Accuracy of E-value estimation by HHblits and PSI-BLAST.
Supplementary Figure 9	Sensitivity and selectivity of homology detection for multi-domain proteins.
Supplementary Figure 10	Relationship between HHblits/HHsearch confidence estimates from the maximum accuracy algorithm and the probability for a residue pair to be correctly aligned.
Supplementary Table 1	Improvement of PSIPRED secondary structure prediction accuracy

	through HHblits multiple sequence alignments (MSAs).
Supplementary Table 2	List of 394 PFAM families for which no homologous template is known, HMMER3 has no match in the PDB below an E-value of $< 10^{-3}$ and for which HHblits has a match in the PDB database with E-value $< 10^{-3}$.

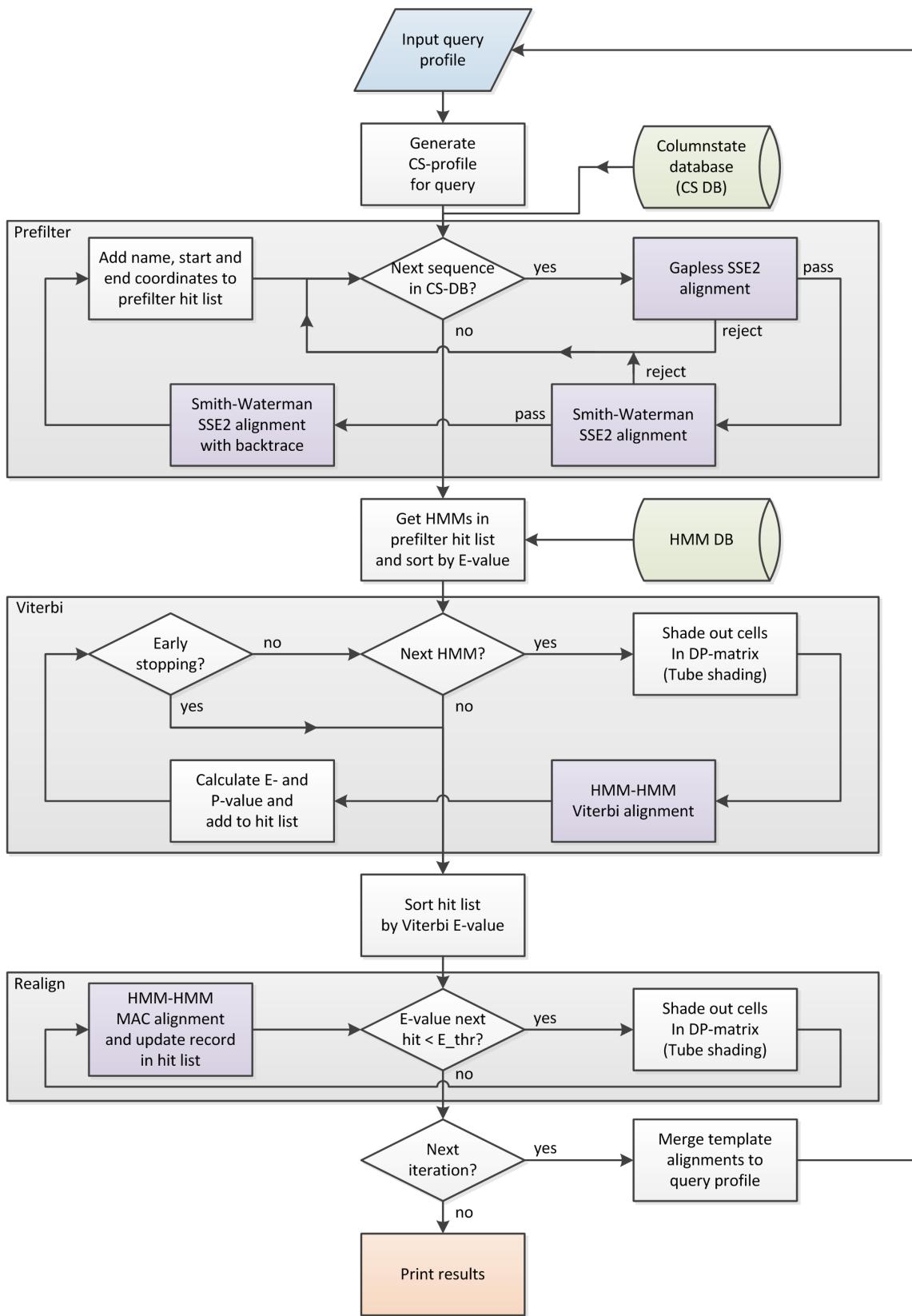
Note: Supplementary Data 1–5 are available on the Nature Methods website.



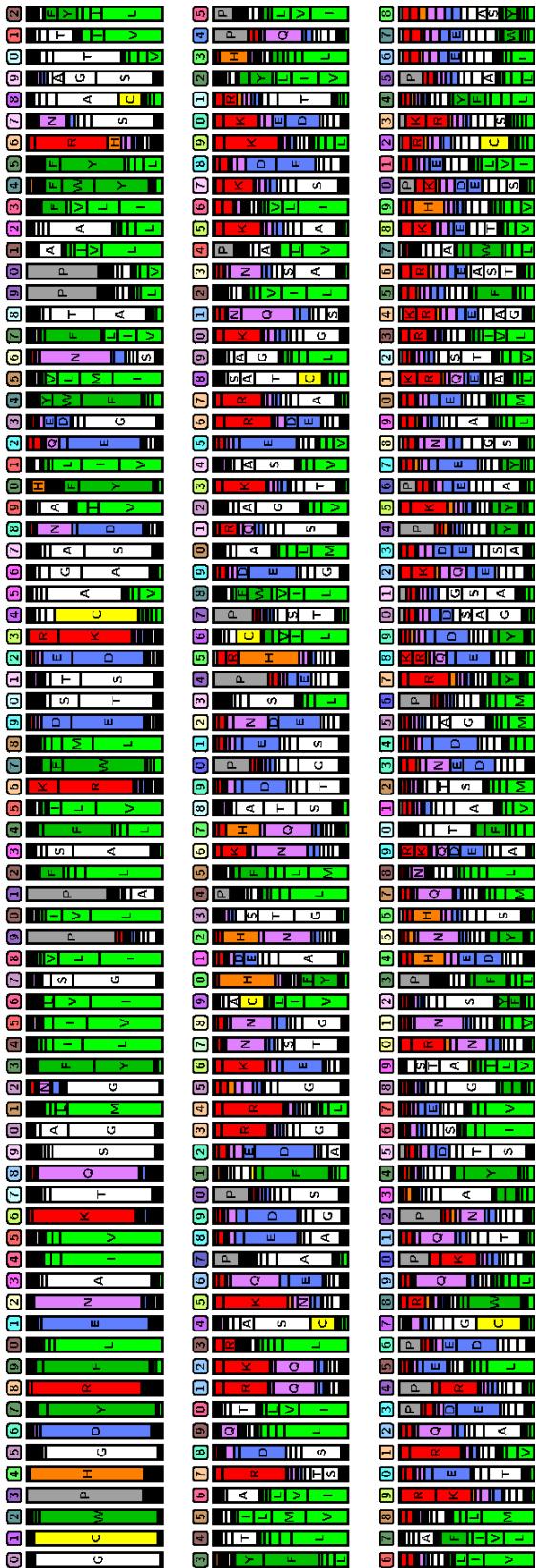
Supplementary Figure 1: Sensitivity / selectivity comparison between HHblits and HHsearch on the SCOP test set for a single search. The prefiltering in HHblits leads only to a very slight performance decrease with respect to HHsearch, even though the run time of HHblits is decreased dramatically.



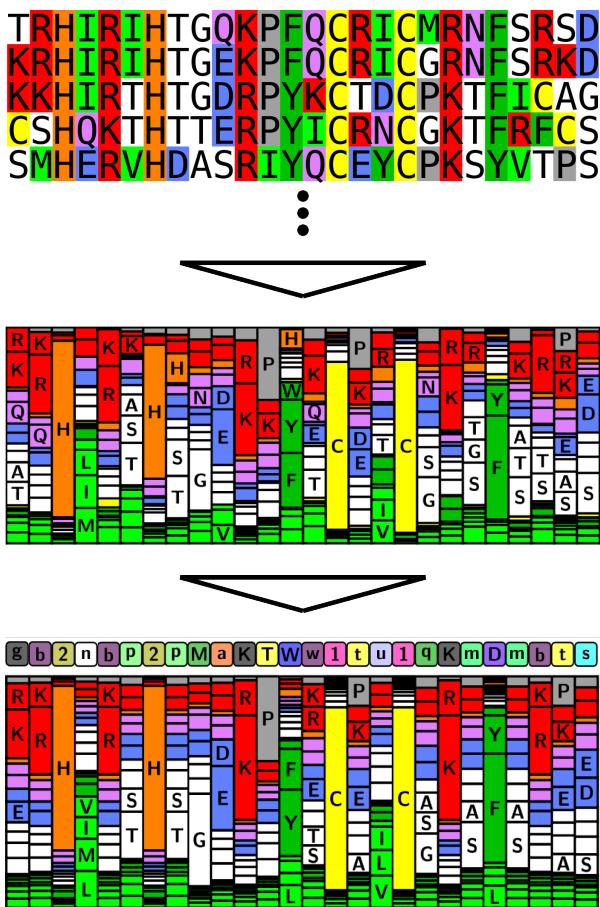
Supplementary Figure 2: Sequence clusters in the UniProt20 predominantly contain functionally closely related proteins. **(a)** Out of 10061 clusters containing at least two SwissProt sequences with annotated enzyme commission (EC) numbers, the graph gives the fraction of those clusters that contain only a single EC annotation up to digit 1 to 4. **(b)** Average fraction of sequences possessing the majority annotation in each cluster, averaged over all 10061 clusters with more than one EC-annotated sequence. Note that the actual functional homogeneity of UnitProt20 clusters may be lower than suggested by these numbers, since a part of SwissProt's EC annotations will have been based on evidence from annotations of homologs, among other information sources.



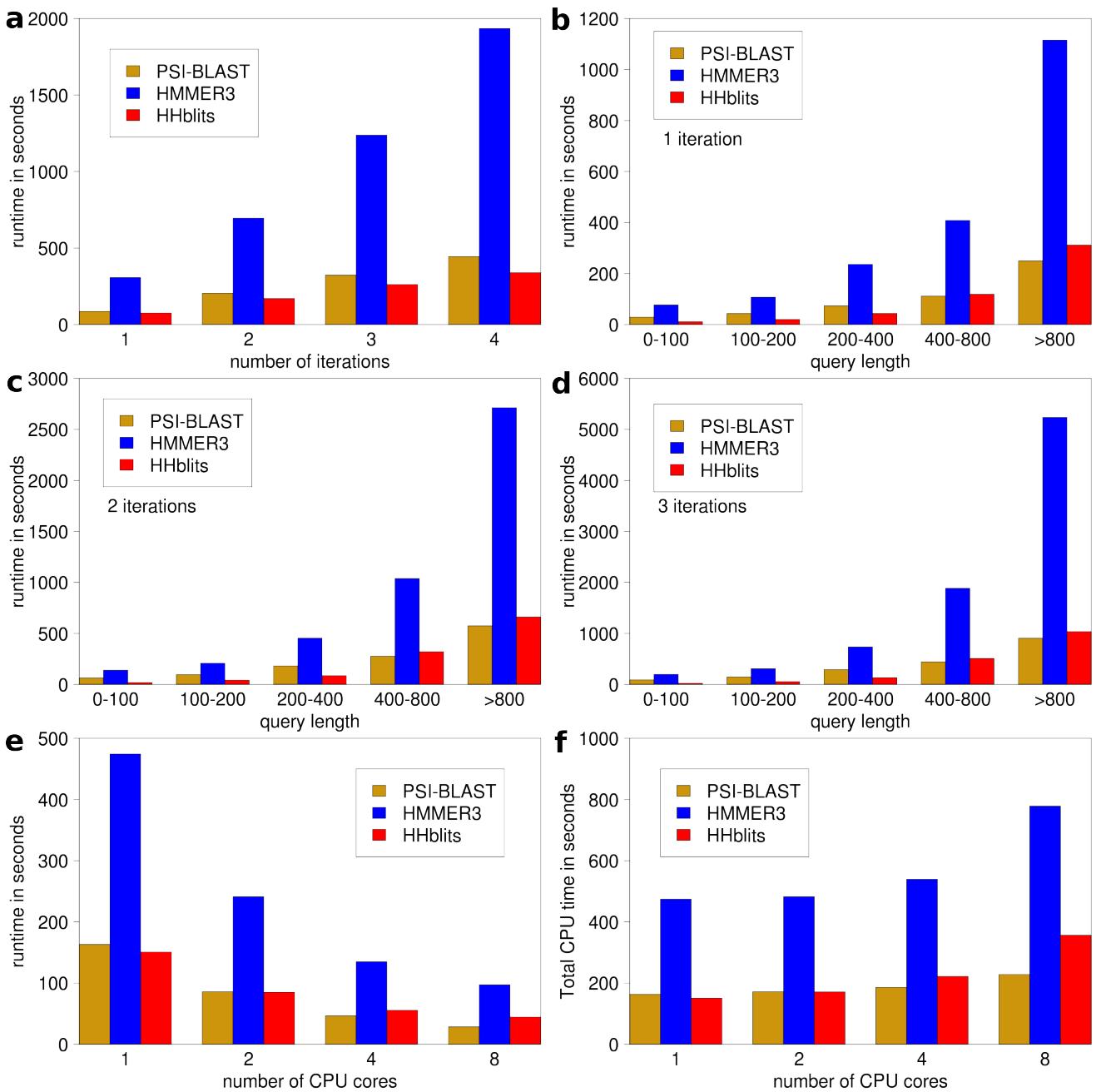
Supplementary Figure 3: Detailed workflow of HHblits. The prefilter starts with a column state (CS) database and a CS-profile for the query. A list of HMMs that pass the prefilter, sorted by their prefilter E-values, is given to the HMM-HMM comparison part. In the Viterbi algorithm, several additional filters are applied and the E-values and P-values of all hits are calculated. Afterwards, the alignments of the best hits are realigned by the Maximum Accuracy (MAC) algorithm. If a further iteration will be performed, the template alignments are merged to the query profile and the scheme is repeated. DP: dynamic programming.



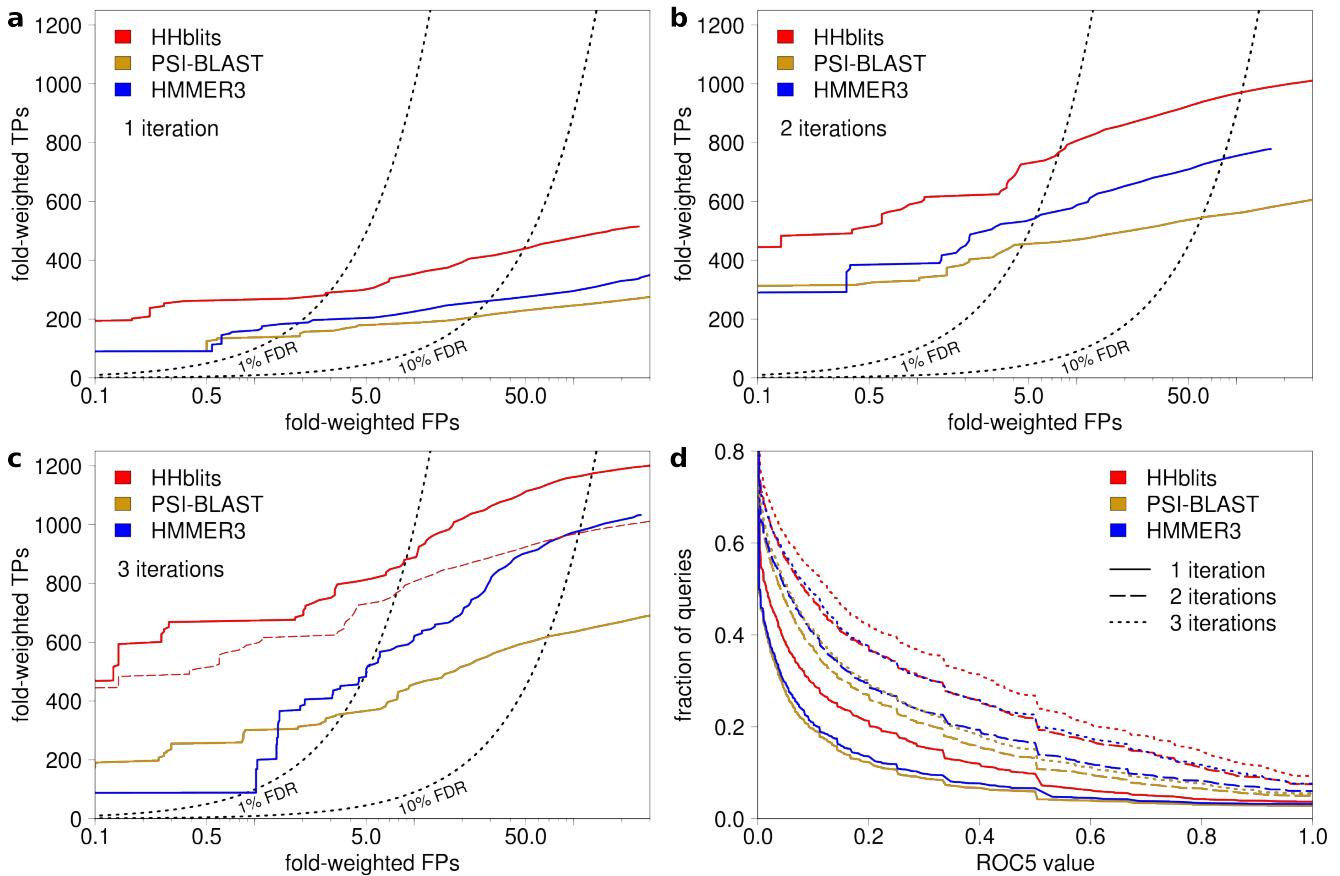
Supplementary Figure 4: Histogram representation of the amino acid distributions in the 219 profile columns of the extended profile column alphabet. The column vectors are ordered by entropy, starting with the almost pure states and ending with near-random background distributions.



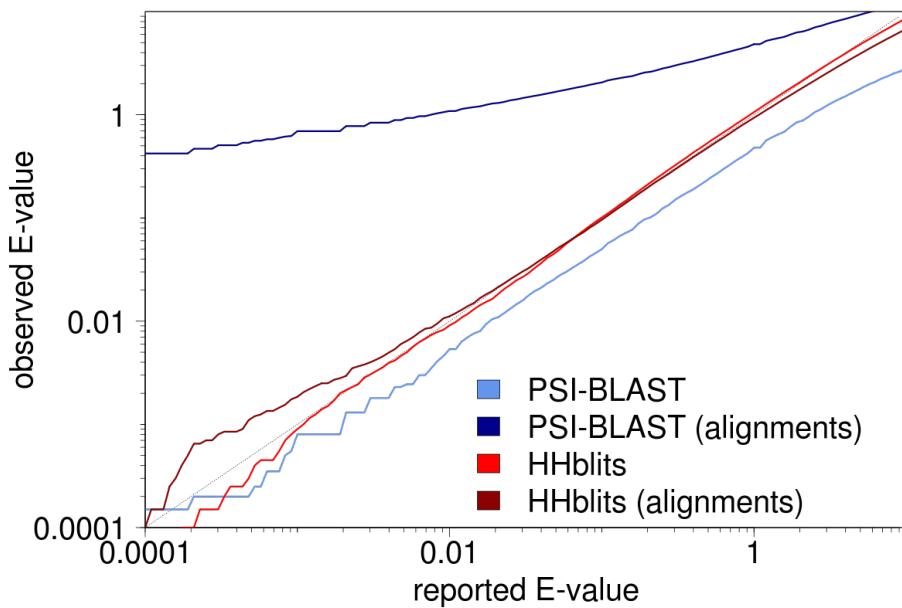
Supplementary Figure 5: Translation of an alignment into a sequence of column states. In a first step, a sequence profile is generated from a given sequence alignment by adding a small amount of context-specific pseudocounts (histogram below alignment). Afterwards, each profile column is translated into the one column state that best describes the observed counts (colored boxes with rounded corners). Below each column state letter its characteristic profile column is shown. Clearly, the column state sequence represents the given sequence profile very well.



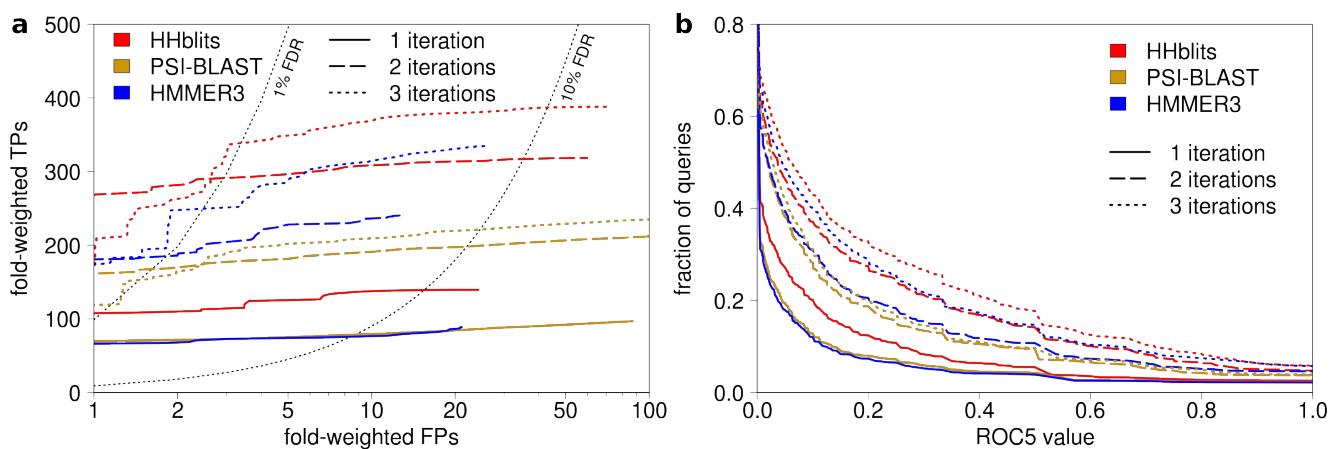
Supplementary Figure 6: Average run times on 100 randomly selected sequences from the nr database, measured on an Intel Xeon X5570 at 2.93 GHz. (a) Average run times for 1 to 4 iterations. (b)-(d) Run times for various bins of query length for (b) one, (c) two, and (d) three search iterations. HHblits has a very good run time for queries with a sequence length below 400 residues and clearly outperforms PSI-BLAST in this range of query lengths. In the range of 400 to 800 residues both methods have a similar run time and only for proteins with a length > 800 residues the run time of HHblits is slightly worse to that of PSI-BLAST. HMMER3 scales in a similar way as HHblits with the query length, always by a factor 3 to 5 slower. (e), (f) Run time for two search iterations on 1 to 8 CPU cores. (e) shows the wall time, whereas (f) gives the total CPU time, equal to the wall time times the number of CPUs. All three tools scale well with an increasing number of CPUs.



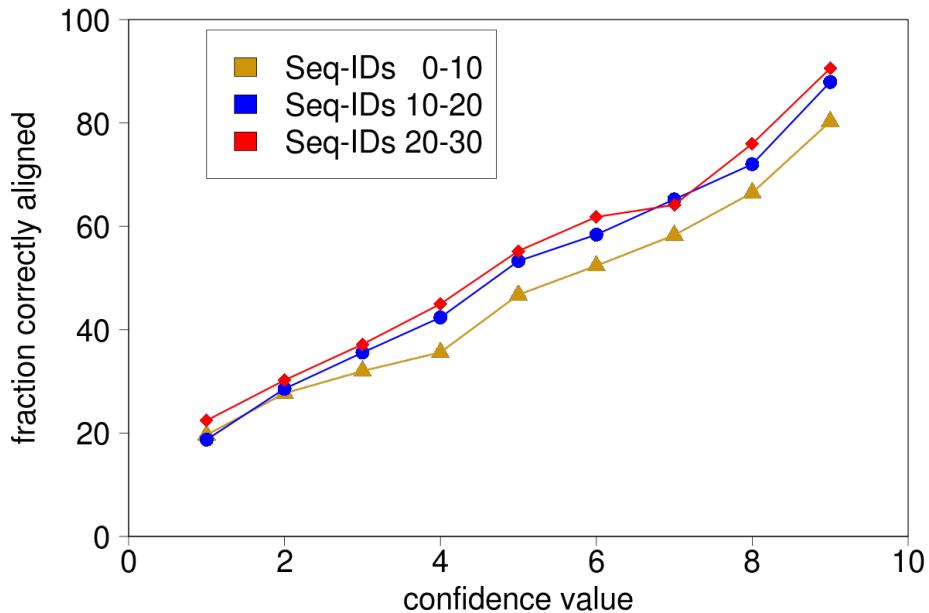
Supplementary Figure 7: Sensitivity and selectivity of homology detection. **(a)-(c)** ROC (receiver operating characteristic) plots for (a) one, (b) two, and (C) three iterations on the test set (5287 sequences from the SCOP 1.73 database). All but the last search iteration are performed against the UniProt. The last search iteration is done through a combined database containing the UniProt and the SCOP sequences (See Online Methods). TPs are defined as pairs from the same SCOP folds, FPs as pairs from different folds, with the exception of Rossman folds and β propellers. At a false discovery rate (FDR) of 10% HHblits detects in the first iteration twice as many TPs as PSI-BLAST and 68% more than HMMER3. In (c), the light red curve (two iterations of HHblits) shows a clear improvement over three iterations of PSI-BLAST and HMMER3. **(d)** Fraction of queries with ROC5 value above the threshold on the x-axis. The ROC5 value is the area under the ROC curve up to the 5'th FP, normalized to yield a theoretical maximum of 1. The ROC5 plot is more robust to overfitting than the ROC curves in (a-c). (see Söding & Remmert, Curr. Opin. Struct. Biol. 2011).



Supplementary Figure 8: Accuracy of E-value estimation by HHblits and PSI-BLAST. We generated a version of UniProt (and the clustered UniProt) with randomly shuffled residues (MSA columns) and randomly selected 20 000 proteins from the nr database to search through the scrambled database. Any match is therefore a false positive. We counted the number of matches below a given reported E-value. Dividing this number through the number of total searches (20 000) yields the empirical, observed E-value. Reported and observed E-values should be similar. Both PSI-BLAST (light blue) and HHblits (red) report reliable E-values when started with a single sequence. When searches are jump-started with a MSAs (obtained from one iteration of HHblits), PSI-BLAST produces a great excess of false positive matches at E-values below 1 (dark blue).



Supplementary Figure 9: Sensitivity and selectivity of homology detection for multi-domain proteins. True positive pairs (TPs) and false positive pairs (FPs) are counted only if the alignment covers at least 50 residues of the SCOP domain in the query NR protein. (a) ROC (receiver operating characteristic) plot showing TPs versus FPs detected at the same E-value thresholds, for 1, 2 and 3 search iterations. After three iterations, HHblits has significantly fewer FPs at high confidence (false discovery rate FDR < 1%) than PSI-BLAST and HMMER3. (b) Fraction of queries with ROC5 value above the threshold on the x-axis. The ROC5 value is the area under the ROC curve up to the 5'th FP, normalized to yield a theoretical maximum of 1.



Supplementary Figure 10: Relationship between HHblits/HHsearch confidence estimates from the maximum accuracy algorithm and the probability for a residue pair to be correctly aligned. The confidence values have an excellent correlation with the fraction of correctly aligned columns, and are nearly independent of the sequence identity between query and template sequences.

Supplementary Table 1: Improvement of PSIPRED secondary structure prediction accuracy through HHblits multiple sequence alignments (MSAs). Performance is measured on the sequences from the PDBselect 2007 data set by Segment Overlap score (SOV) and 3-state accuracy (Q3). For each sequence in this set, MSAs are generated by 2 and 3 iterations of PSI-BLAST and 1, 2 and 3 iterations HHblits. Even 1 iteration of HHblits yield better performance than the standard PSIPRED, which uses 3 iterations of PSI-BLAST. The best results with an improvement of more than 1% are achieved by performing up to 3 iterations of HHblits and filtering the generated MSAs to a diversity of $N_{eff} = 7$.

input alignments	SOV	Q3
2 iterations PSI-BLAST	74.64%	77.31%
3 iterations PSI-BLAST	77.52%	80.38%
1 iteration HHblits	77.87%	80.71%
2 iterations HHblits	78.31%	80.99%
3 iterations HHblits	78.12%	80.83%
HHblits diversity 7	78.62%	81.31%

Supplementary Table 2: List of 394 PFAM families for which no homologous template is known, HMMER3 has no match in the PDB below an E-value of $< 10^{-3}$ and for which HHblits has a match in the PDB database with E-value $< 10^{-3}$. In each row, the best HHblits match is given with its E-value and the coverage of the PFAM query. The last column specifies the HMMER3 E-value for the best match, an '-' indicates that HMMER3 has no matches up to the default reporting E-value threshold of 10.

PFAM-ID	HHblits			HMMER	PFAM-ID	HHblits			HMMER
	hit	E-value	cov(%)	E-value		hit	E-value	cov(%)	E-value
PF03115	3hag	2e-101	53.24	-	PF12243	3d9j	1.e-26	99.27	-
PF11838	2xdt	1.1e-55	99.77	1	PF09960	2w3z	1.6e-26	43.59	1.2
PF09562	2oa9	1.2e-54	98.85	-	PF10962	1v9m	2.1e-26	87.63	-
PF10991	1je5	4e-54	96.02	-	PF07014	2nw8	2.3e-26	95.20	-
PF04412	1c96	2.6e-53	83.81	0.006	PF10100	3c7a	2.7e-26	92.56	0.25
PF09863	3iv3	1.2e-50	95.77	0.0048	PF04841	1got	7.2e-26	70.80	-
PF12043	3c5n	1.9e-49	98.81	0.26	PF07608	2zyz	1e-25	96.71	-
PF11047	3cxb	2e-49	73.10	0.049	PF07379	3ci0	1.3e-25	79.43	-
PF07718	1r4x	3.8e-49	86.38	0.051	PF03687	3bry	3.2e-25	87.35	0.11
PF07632	2mas	2.1e-47	95.94	-	PF07592	3hot	3.6e-25	93.57	-
PF08010	2b3w	4.2e-45	96.58	0.14	PF05651	2a2l	7.2e-25	82.22	0.0057
PF11329	3eu8	9.5e-44	98.64	0.31	PF12260	2acx	1.1e-24	89.76	-
PF03813	1q79	1e-43	47.19	0.017	PF06230	1zc2	1.5e-24	69.67	-
PF05316	3bbn	3.4e-41	90.71	-	PF05213	1vgj	3.9e-24	66.67	0.0084
PF09520	1wte	5.9e-40	97.33	-	PF11768	1vyh	4.8e-24	61.72	0.097
PF12264	1lqqp	7.5e-38	90.26	0.043	PF04551	1tx2	4.9e-24	71.68	0.07
PF11161	2ra9	7.2e-37	77.46	-	PF05550	2wur	5.5e-24	74.40	-
PF09739	3f8t	2.8e-36	70.51	-	PF03290	1th0	7.4e-24	53.85	1.5
PF10963	3fgx	4.5e-36	100.00	0.0027	PF09927	2jxp	8.9e-24	99.15	6.2
PF05864	2waq	5.8e-36	87.30	0.022	PF09807	3bs4	9.5e-24	97.60	0.087
PF04486	1tuw	8.2e-36	77.39	-	PF06199	2k4q	1.1e-23	100.00	0.013
PF05714	1w33	1e-35	86.26	0.078	PF08472	1tp6	1.2e-23	84.21	-
PF05428	3kq4	1e-34	92.90	0.06	PF08553	1got	1.7e-23	42.52	0.15
PF07588	1koe	2.1e-34	95.83	0.059	PF11340	3cz8	1.8e-23	86.14	0.17
PF09536	2kii	1.1e-33	95.60	-	PF09892	3na6	2.3e-23	84.54	0.091
PF03662	1qw9	1.5e-33	99.38	-	PF09674	1kea	5.4e-23	81.06	0.012
PF11443	2nvo	1.6e-33	92.88	-	PF05551	1a73	5.6e-23	52.24	0.91
PF05538	2odj	1.7e-33	95.05	-	PF08695	2ciu	7.4e-23	76.56	0.053
PF07520	1yuw	2.9e-33	58.91	-	PF04405	2k5e	7.9e-23	98.21	0.0013
PF06124	2r41	4.1e-33	100.00	-	PF10179	1fnh	8.7e-23	95.62	0.66
PF05291	2ilr	4.7e-33	60.86	0.073	PF02677	1wy5	9.2e-23	94.12	9.9
PF06045	1nkg	2.3e-31	92.65	-	PF04708	3guv	1e-22	62.02	-
PF06787	2a9s	4.6e-31	99.38	-	PF07006	2k3d	1.2e-22	66.40	-
PF10023	1z5h	8.3e-31	93.51	0.061	PF10250	2hhc	1.8e-22	94.17	0.099
PF05986	3ghm	1.4e-30	100.00	-	PF11813	3h2d	1.9e-22	71.67	-
PF05482	1tr2	4e-30	85.30	-	PF08928	2fef	2.3e-22	100.00	0.029
PF11686	1se7	1.6e-29	100.00	-	PF09859	3itq	2.5e-22	83.73	-
PF07307	3nf2	1.8e-29	80.95	0.027	PF10107	3fov	8.5e-22	48.75	0.89
PF07528	1px5	2e-29	94.19	0.46	PF09843	2zws	8.6e-22	73.60	-
PF03254	2de0	2.4e-29	74.48	0.5	PF08470	2vxr	1.2e-21	59.88	-
PF07395	1lrz	5.9e-29	98.86	-	PF11017	2hcy	2.4e-21	98.75	0.061
PF10287	3iln	6.3e-29	91.81	-	PF04937	2bw3	2.6e-21	99.35	0.21
PF01531	2hhc	7.1e-29	92.88	0.12	PF10222	1h54	2.7e-21	56.50	-
PF06074	3kdr	1.8e-28	56.66	0.14	PF03336	2jig	2.9e-21	46.15	0.23
PF06128	1yyh	2e-28	99.65	0.17	PF10677	2flc	3.4e-21	97.85	0.48
PF02088	1dec	1.2e-27	100.00	0.011	PF06420	1h2i	3.9e-21	56.50	-
PF05136	3kdr	1.2e-27	90.17	-	PF06674	3tdt	4.7e-21	46.13	3
PF07756	2qc0	1.5e-27	97.69	-	PF09796	2fyu	7.8e-21	87.10	0.33
PF04681	1z3q	1.8e-27	86.45	-	PF06951	1lw	1e-20	50.56	-
PF09865	2w7q	1.8e-27	96.79	-	PF06437	2fue	1.2e-20	64.90	-
PF08189	2b5b	3e-27	94.87	0.031	PF10748	2ivw	2.2e-20	72.39	-
PF10770	2plg	5.3e-27	82.88	0.014	PF07894	1byr	2.7e-20	58.80	0.0045
PF11841	3dad	6e-27	99.36	-	PF00609	2bon	4e-20	100.00	-
PF06245	1vk1	6.1e-27	52.51	0.78	PF09517	1yd6	5.6e-20	69.01	0.77
PF05060	1fo8	7.7e-27	75.42	0.062	PF11039	2vzy	5.9e-20	98.01	0.11

Supplementary Table 2 continued

PFAM-ID	HHblits			HMMER		PFAM-ID	HHblits			HMMER	
	hit	E-value	cov(%)	E-value			hit	E-value	cov(%)	E-value	
PF03351	1d7b	9.8e-20	96.80	0.0047		PF09337	3nnq	5.1e-15	100.00	0.022	
PF12541	logo	1e-19	79.50	-		PF08424	3dss	5.1e-15	90.61	0.031	
PF11680	3k44	1.2e-19	63.64	0.9		PF08170	3gir	8.5e-15	88.78	0.046	
PF12439	1v7w	1.5e-19	99.55	-		PF08379	3isr	1.1e-14	100.00	-	
PF05176	3gkn	2e-19	62.70	-		PF06805	3dwg	1.3e-14	42.70	-	
PF09810	3l0a	2.1e-19	54.39	0.011		PF09363	1umd	2.4e-14	72.41	-	
PF10738	3lyd	2.3e-19	55.51	0.0022		PF10826	2fe3	2.7e-14	85.19	0.46	
PF05046	2ogh	2.7e-19	88.89	0.043		PF07611	3bma	3.4e-14	80.00	0.021	
PF05342	3n6z	3.4e-19	48.98	-		PF08757	3dnlu	3.5e-14	58.54	-	
PF10288	1ni5	3.4e-19	98.95	0.14		PF07461	1tvg	4.2e-14	31.39	0.029	
PF09778	3erv	4.4e-19	99.12	0.0024		PF09941	1vet	7.7e-14	91.67	-	
PF04788	3bk5	5.1e-19	84.50	-		PF10141	2zxr	8.4e-14	75.13	0.11	
PF03214	1qg8	6.2e-19	34.00	-		PF09366	2pcs	9.6e-14	93.71	-	
PF06544	2iyg	6.3e-19	93.42	0.39		PF05272	2dhr	1e-13	64.00	1.2	
PF11854	2guf	7.5e-19	82.01	-		PF07618	1y6u	1.2e-13	98.25	0.46	
PF01973	2p2v	1e-18	85.88	0.0029		PF11824	2okx	1.3e-13	77.90	0.038	
PF07845	1m2d	1e-18	78.95	0.012		PF10483	3bs4	3.4e-13	79.18	-	
PF08885	1ivn	1.1e-18	79.34	0.051		PF10029	3c12	3.5e-13	81.51	0.095	
PF08642	1jbi	1.4e-18	85.27	3.5		PF11863	3hx1	4e-13	95.55	0.016	
PF07607	1z5h	1.4e-18	96.00	0.39		PF06147	3g27	4.2e-13	42.93	0.1	
PF06241	1lnq	1.4e-18	78.16	-		PF10743	1y6u	4.2e-13	70.93	0.047	
PF07076	2qcp	1.4e-18	83.54	0.26		PF10246	1k0r	4.3e-13	82.86	-	
PF10365	3km5	1.5e-18	89.51	0.49		PF11288	3fak	5.7e-13	66.99	0.0043	
PF07959	1yp2	1.5e-18	65.83	0.38		PF05227	1vls	8.2e-13	97.10	0.17	
PF10934	2ia7	1.7e-18	89.32	-		PF11845	3b42	8.3e-13	59.88	2.2	
PF02411	2h3o	3.5e-18	54.78	-		PF10302	2bps	8.8e-13	34.29	0.019	
PF10012	3h96	3.7e-18	80.00	0.028		PF11312	3mgg	9.4e-13	55.44	-	
PF07115	2ia7	4.1e-18	90.09	0.0018		PF11071	1s2d	1.8e-12	99.29	0.3	
PF07293	3i9v	6.3e-18	97.44	-		PF10037	1xi4	2.3e-12	75.06	-	
PF06477	2ag4	7.4e-18	97.32	-		PF04781	1elw	2.9e-12	96.72	0.27	
PF06021	1sqh	1.1e-17	99.51	-		PF09530	2i71	4.1e-12	74.59	0.023	
PF03018	2brj	1.2e-17	77.78	0.0043		PF07800	3knv	4.9e-12	80.62	0.1	
PF04114	3k9t	1.3e-17	44.98	-		PF08156	3id6	6.1e-12	100.00	2.8	
PF05565	2p2u	1.3e-17	74.84	0.016		PF06381	3kdr	8.8e-12	95.51	0.05	
PF07905	2ioj	1.3e-17	91.06	0.012		PF04244	2wq7	9.2e-12	69.78	-	
PF11751	3bry	1.4e-17	89.78	-		PF10127	3c18	1.2e-11	79.20	-	
PF09565	2c11	1.5e-17	61.20	-		PF08130	1w9n	1.2e-11	51.79	-	
PF08734	2zbc	1.8e-17	82.42	0.0045		PF11959	3hft	1.4e-11	84.09	0.27	
PF11019	3mc1	3.2e-17	75.37	0.0029		PF02066	1m0j	1.5e-11	51.85	3.2	
PF06044	2jne	3.4e-17	18.43	-		PF04986	1omh	1.6e-11	33.68	-	
PF04865	3h2t	5e-17	74.60	-		PF04082	2veq	1.6e-11	29.30	-	
PF12362	2aya	5.2e-17	84.62	-		PF10138	3ibz	1.6e-11	78.49	0.14	
PF05991	1exn	6.8e-17	98.14	-		PF06075	2b29	2.1e-11	20.55	-	
PF06622	1o9y	7.3e-17	24.59	-		PF03490	2plc	2.5e-11	72.55	-	
PF08521	2kse	1.3e-16	97.95	0.05		PF01927	3ga8	2.9e-11	37.58	0.035	
PF08480	1ru4	1.4e-16	99.46	-		PF09345	1h4x	2.9e-11	84.00	0.13	
PF03302	1yy9	1.8e-16	77.92	-		PF11761	3eeq	2.9e-11	100.00	0.1	
PF07506	1zx4	2.1e-16	98.83	0.0047		PF06881	2e31	3.4e-11	70.09	0.011	
PF01185	2fmc	5.1e-16	79.63	3.1		PF04492	3e6c	3.5e-11	82.00	0.02	
PF10087	2iw1	6.2e-16	95.65	0.057		PF06890	2p5z	4.2e-11	46.03	-	
PF11814	3erv	6.2e-16	90.00	0.15		PF08303	1yj5	4.7e-11	98.82	0.011	
PF03452	1xhb	7.3e-16	93.31	-		PF03281	1px5	6.3e-11	97.49	0.12	
PF10703	2p8g	8.9e-16	30.65	0.12		PF08497	2yx5	1.5e-10	47.75	-	
PF02413	2kz6	1.2e-15	57.14	0.71		PF10030	2jyx	1.5e-10	63.83	-	
PF07327	1wqj	1.3e-15	55.14	0.35		PF11997	3c48	1.7e-10	99.64	-	
PF11356	2ivw	1.5e-15	54.17	0.22		PF02697	3fmt	1.7e-10	86.67	0.0022	
PF12055	1k1x	2.3e-15	58.25	-		PF02474	1m4i	2.1e-10	73.10	-	

Supplementary Table 2 continued

PFAM-ID	HHblits			HMMER		
	hit	E-value	cov(%)	hit	E-value	cov(%)
PF12226	3iyo	2.5e-10	43.04	-		
PF10567	1l3k	2.7e-10	69.36	-		
PF09826	1fwx	4.3e-10	81.37	0.014		
PF03158	2xeh	5.9e-10	75.65	-		
PF08371	3hs1	6.6e-10	86.42	-		
PF02666	2gpr	9.4e-10	90.23	-		
PF06676	2waq	1.4e-09	37.86	0.59		
PF06322	3e7l	1.5e-09	67.19	0.09		
PF08685	1z3u	1.9e-09	24.00	-		
PF07508	2r0q	1.9e-09	56.14	-		
PF09317	2z1q	2e-09	48.75	-		
PF07328	2ba3	2.2e-09	27.89	-		
PF10908	1zat	2.9e-09	75.70	-		
PF10474	2fji	3e-09	96.58	-		
PF02521	3jty	3.5e-09	57.42	-		
PF08499	3g4g	3.7e-09	88.52	0.025		
PF06669	3d9x	3.8e-09	97.14	-		
PF11954	1ei5	4e-09	63.87	0.028		
PF07202	1h3i	4.5e-09	71.98	0.22		
PF12010	1j1n	5.1e-09	94.20	0.12		
PF04936	3hot	7.7e-09	78.49	8.1		
PF08116	1c6w	7.8e-09	87.10	0.017		
PF12000	3fro	8.5e-09	69.59	-		
PF11766	1n67	8.6e-09	97.18	-		
PF11711	2qv7	1.2e-08	32.09	-		
PF06883	1twf	1.7e-08	100.00	2		
PF08074	1ofc	1.7e-08	41.07	1.3		
PF08465	1p6x	2.1e-08	96.97	-		
PF05263	2o8x	2.5e-08	48.89	0.095		
PF12128	1w1w	3.2e-08	6.19	0.0075		
PF09889	1lv3	3.3e-08	46.55	0.059		
PF04450	1z5h	4.3e-08	86.00	0.023		
PF11308	2zxq	4.5e-08	49.71	-		
PF09597	2e8n	4.7e-08	98.25	2.3		
PF09824	2p4w	5.3e-08	78.75	0.0017		
PF10780	1s3a	5.3e-08	100.00	-		
PF11658	3lxq	6.1e-08	58.70	-		
PF06011	1nep	7.2e-08	22.92	-		
PF04407	3dcn	8.4e-08	95.43	0.0061		
PF11853	2odj	8.5e-08	69.22	0.049		
PF09582	1p90	9.6e-08	50.46	-		
PF07610	2qsv	1e-07	100.00	0.015		
PF11325	2vw9	1.2e-07	98.85	-		
PF10686	2nx2	1.2e-07	88.73	0.017		
PF01439	2kak	1.3e-07	93.67	0.34		
PF08737	2fau	1.3e-07	77.88	-		
PF05380	1rw3	1.7e-07	86.67	-		
PF07699	2hey	1.9e-07	100.00	0.0045		
PF06378	1h2i	2.5e-07	80.50	6.1		
PF05610	2apn	3.6e-07	89.47	0.19		
PF10758	3hxl	4.1e-07	99.45	-		
PF04189	1yb2	4.4e-07	77.70	-		
PF09855	2k4x	5.5e-07	82.81	0.22		
PF06355	1gwy	6.1e-07	76.52	-		
PF05895	2fl8	6.6e-07	28.31	-		
PF04917	1oqw	6.8e-07	15.71	0.0036		
PF04572	2vk9	7e-07	80.00	0.087		
PF10144	3b42	7.5e-07	54.29	0.061		
PF08405	3i86	8.4e-07	12.01	-		
PF07087	1lw1	9.7e-07	77.17	-		
PF09970	2fc1	9.7e-07	77.30	-		
PF06977	1npe	1e-06	89.81	0.0024		
PF07429	2gek	1.1e-06	80.06	-		
PF10349	2hth	1.3e-06	32.41	-		
PF10116	3e20	1.3e-06	98.57	-		
PF08417	3gke	1.3e-06	59.43	-		
PF10711	2vxz	1.9e-06	82.65	0.014		
PF12303	2wg3	2.1e-06	61.70	0.0061		
PF07813	3epv	2.2e-06	88.46	0.0022		
PF10373	1ya0	3.5e-06	72.60	-		
PF02681	2ipb	4.2e-06	90.54	0.27		
PF11897	2gj4	4.2e-06	42.86	0.008		
PF10987	3h35	4.5e-06	72.65	0.1		
PF07617	3ia8	4.7e-06	98.18	-		
PF03345	2gk3	4.9e-06	54.17	-		
PF03850	3ibs	5.2e-06	80.92	0.011		
PF07919	2icn	5.3e-06	62.70	0.069		
PF10781	1whg	5.3e-06	98.39	-		
PF07107	3ec9	5.5e-06	69.39	0.0012		
PF05782	1kxp	5.5e-06	52.22	-		
PF12578	1lw3	6.2e-06	66.29	0.29		
PF04377	3gkr	6.4e-06	95.35	-		
PF11903	1baz	6.6e-06	47.95	-		
PF12073	1pj1	6.8e-06	94.23	0.26		
PF04312	1hjr	8.3e-06	76.98	0.025		
PF07480	3epv	8.6e-06	94.74	0.58		
PF12525	1v9n	8.7e-06	86.67	0.28		
PF04413	1vgv	8.8e-06	86.34	-		
PF05091	3fq1	9.4e-06	53.89	0.39		
PF10367	1chc	9.7e-06	29.36	6.2		
PF06353	2fph	1e-05	36.81	-		
PF06239	1xi4	1e-05	65.67	-		
PF08192	1hpg	1.1e-05	13.75	-		
PF08579	1xi4	1.2e-05	82.50	-		
PF05510	1u2c	1.3e-05	39.78	-		
PF11833	1faf	1.3e-05	22.06	3.9		
PF06823	2l1s	1.4e-05	80.33	-		
PF09984	3b42	1.5e-05	93.29	-		
PF09352	2qgp	1.6e-05	43.68	-		
PF06375	2g30	1.9e-05	34.78	-		
PF03406	1h6w	2.1e-05	90.70	-		
PF12340	3ly5	2.1e-05	77.73	0.036		
PF10952	2fbn	2.2e-05	88.03	-		
PF10240	2qp2	2.2e-05	37.80	-		
PF02754	3cf4	2.5e-05	80.95	0.91		
PF10941	1uf3	2.5e-05	55.08	0.054		
PF05689	1f00	2.5e-05	99.45	-		
PF01941	2p02	2.6e-05	91.86	0.62		
PF11839	1jcd	2.7e-05	39.76	-		
PF06956	1xmx	2.7e-05	75.94	-		
PF07585	2x55	2.8e-05	62.71	-		
PF06448	1lsh	3.2e-05	39.33	-		
PF10865	1ilo	3.2e-05	51.28	1		
PF10505	3fq1	3.5e-05	72.90	-		

Table Supplementary 2 continued

PFAM-ID	HHblits		HMMER	
	hit	E-value	cov(%)	E-value
PF11865	2qk1	3.9e-05	96.89	0.65
PF12222	1pgs	4e-05	69.09	0.21
PF00746	2ww8	4.1e-05	92.50	0.36
PF09759	1xqr	4.9e-05	72.63	1.2
PF11834	2dnf	5.1e-05	98.48	5.4
PF10126	3dfe	5.5e-05	92.86	0.068
PF07878	1nla	5.7e-05	92.00	-
PF07505	3c8f	7.2e-05	81.89	0.1
PF04155	1yo3	7.6e-05	72.97	-
PF10904	1j8b	7.6e-05	63.37	-
PF10706	2fc1	8.7e-05	88.51	-
PF01963	2g5g	8.9e-05	97.76	0.32
PF11336	2o4v	9.3e-05	77.48	0.014
PF03249	1p4t	9.4e-05	17.98	-
PF04599	1rxw	9.8e-05	65.00	-
PF01696	1pcl	0.0001	49.61	-
PF06702	1cja	0.00011	54.02	-
PF10122	2jr6	0.00013	80.39	10
PF11112	1z4h	0.00015	86.84	-
PF11849	3e0y	0.00016	94.77	-
PF06904	1lbu	0.00018	54.82	0.015
PF05918	1b3u	0.00019	60.90	0.0019
PF09854	2qgp	0.00019	23.82	0.021
PF10497	1wil	0.00025	61.76	0.71
PF07802	2k3j	0.00028	81.43	0.64
PF04305	3chh	0.00028	73.09	0.092
PF08736	2i1j	0.00028	51.06	-
PF06974	2jgp	0.0003	98.04	-
PF08288	3fro	0.0003	82.22	-
PF00242	1wz4	0.00031	7.33	-
PF09538	1vd4	0.00036	23.02	2.9
PF10673	3lub	0.00037	66.21	-
PF12215	2cqs	0.00038	63.13	-
PF11379	2cqy	0.00038	22.10	-
PF03258	2w7a	0.00045	46.67	0.094
PF10115	2kon	0.00046	76.34	-
PF08498	1wg8	0.00049	76.12	-
PF05444	3laq	0.00049	87.82	2.9
PF05869	3lkd	0.00055	81.40	0.06
PF11001	1wij	0.00058	65.61	0.0013
PF10309	3d45	0.00063	96.67	0.036
PF04049	3kae	0.00065	86.05	-
PF11006	2pxg	0.00067	87.21	1.6
PF07295	1lko	0.00068	25.85	0.096
PF10407	2ns5	0.00069	94.67	2.9
PF10165	1xm9	0.00073	77.28	-
PF08749	2plg	0.00074	92.41	-
PF04904	1rg6	0.00076	78.05	0.13
PF07409	2ia7	0.00078	63.25	0.049
PF12416	2dmh	0.00085	97.09	-
PF09576	1v54	0.00086	91.23	-
PF08004	2cob	0.00086	44.27	-
PF09415	1b67	0.00087	91.78	0.028
PF09894	1iru	0.0009	92.23	4.3
PF07855	2vfx	0.00091	95.73	0.16
PF10790	2al3	0.00095	92.11	0.019