

METHODOLOGY ARTICLE

Open Access



# Estimating statistical significance of local protein profile-profile alignments

Mindaugas Margelevičius

## Abstract

**Background:** Alignment of sequence families described by profiles provides a sensitive means for establishing homology between proteins and is important in protein evolutionary, structural, and functional studies. In the context of a steadily growing amount of sequence data, estimating the statistical significance of alignments, including profile-profile alignments, plays a key role in alignment-based homology search algorithms. Still, it is an open question as to what and whether one type of distribution governs profile-profile alignment score, especially when profile-profile substitution scores involve such terms as secondary structure predictions.

**Results:** This study presents a methodology for estimating the statistical significance of this type of alignments. The methodology rests on a new algorithm developed for generating random profiles such that their alignment scores are distributed similarly to those obtained for real unrelated profiles. We show that improvements in statistical accuracy and sensitivity and high-quality alignment rate result from statistically characterizing alignments by establishing the dependence of statistical parameters on various measures associated with both individual and pairwise profile characteristics. Implemented in the COMER software, the proposed methodology yielded an increase of up to 34.2% in the number of true positives and up to 61.8% in the number of high-quality alignments with respect to the previous version of the COMER method.

**Conclusions:** The more accurate estimation of statistical significance is implemented in the COMER method, which is now more sensitive and provides an increased rate of high-quality profile-profile alignments. The results of the present study also suggest directions for future research.

**Keywords:** Homology search, Profile-profile alignment, Random profile model, Statistical significance, Protein structure prediction

## Introduction

Establishing homology between proteins is essential in evolutionary and highly important in structural and functional studies of proteins [1] and their complexes. Alignment, in general, represents the most fundamental way to deduce homology directly or indirectly. Sequence alignment has proved indispensable in annotating uncharacterized proteins [2] and paved the way for alignment of sequence families described by profiles, constituting the basis for inferring protein structure and function by homology.

In the presence of a large and steadily growing amount of sequence data, estimating the statistical significance of alignments plays a prominent role in alignment-based homology search algorithms [3]. For an alignment with a particular similarity score, it provides a probability or related quantity indicating how likely a chance alignment with the same or greater score is to be observed.

The limiting distribution of the ungapped local alignment score  $S$  (nonlattice) for large sequence lengths  $m$  and  $n$  has been proved [4–7] to be the type 1 (Gumbel-type) extreme value distribution (EVD) [8]

$$P(S \leq x) \approx \exp(-Kmn e^{-\lambda x}), \quad (1)$$

where  $\lambda$  and  $K$  ( $K = e^{\lambda\mu} / [mn]$ ;  $\mu$ , the location parameter under the standard parametrization) are parameters calculable from the score matrix used to align sequences and

Correspondence: [mindaugas.margelevicius@bti.vu.lt](mailto:mindaugas.margelevicius@bti.vu.lt)  
Institute of Biotechnology, Life Sciences Center, Vilnius University, Saulėtekio al. 7, 10257 Vilnius, Lithuania



satisfying the conditions of the negative expected score and existence of at least one positive score.

Ample empirical evidence, e.g., [9–13], has suggested that the statistical theory for ungapped alignments applies to gapped alignments, provided gap penalties, or costs, lead to alignment scores in the (local) region of logarithmic growth [14].

Importantly, factors, such as the length of sequences being compared [10, 15, 16] and their compositional similarity [9, 17], affect the distribution of alignment scores. While solutions to take them into account exist for sequence and profile-to-sequence alignments [9, 18] (Supplementary Section S1, Additional file 1), there is no established procedure for addressing them in profile-profile alignment [19–22].

This study aims at characterizing the distribution of profile-profile alignment scores [23, 24] to improve statistical accuracy and remote homology detection. Our focus is local alignment scores. Hence, we assume that the expected profile-profile alignment score per aligned pair is negative and the probability for a positive score is positive. These assumptions are satisfied in part when the score for a pair of positions of two profiles is (implicitly) log-odds [3, 25]. Importantly, the score preserves the form of log-odds when it linearly combines different components (e.g., the similarity of two amino acid frequency vectors with that of secondary structures) as long as each component itself represents a log-odds score. Such a construction of composite scores is typical, and, given appropriate gap penalties, we can expect that the assumptions will hold for most profile-profile alignment methods.

Here, we use the COMER method [26] as a means for producing profile-profile alignment scores the algorithms developed in this work employ. The COMER profile at each position encapsulates the transition probabilities of moving to and from the states of insertion and deletion, which for a pair of profiles transform to position-specific gap penalties. Comparing two COMER profiles also includes scoring predicted secondary structure (SS) and sequence contexts [27]. These properties make the characterization of the distribution of alignment scores a challenging task.

Moreover, the distribution of profile-profile alignment scores strongly depends on the extent of (dis-)similarity between unrelated profiles, or the null model of random sequence families. Real sequences do not conform to the canonical (and simplest) model, in which the amino acids in the sequence are iid, and exhibit a more complex structure with short- and long-range amino acid correlations [17, 25, 28–30]. It has been proved that in the limit of infinitely long sequences, ungapped alignment scores for random Markov-dependent sequences are still distributed according to the EVD [31]. Gapped alignments

of correlated random sequences accounted for by a null model have been shown numerically to correspond to an EVD too [32]. Notably, the values of the statistical parameters differed substantially from those obtained for iid sequences.

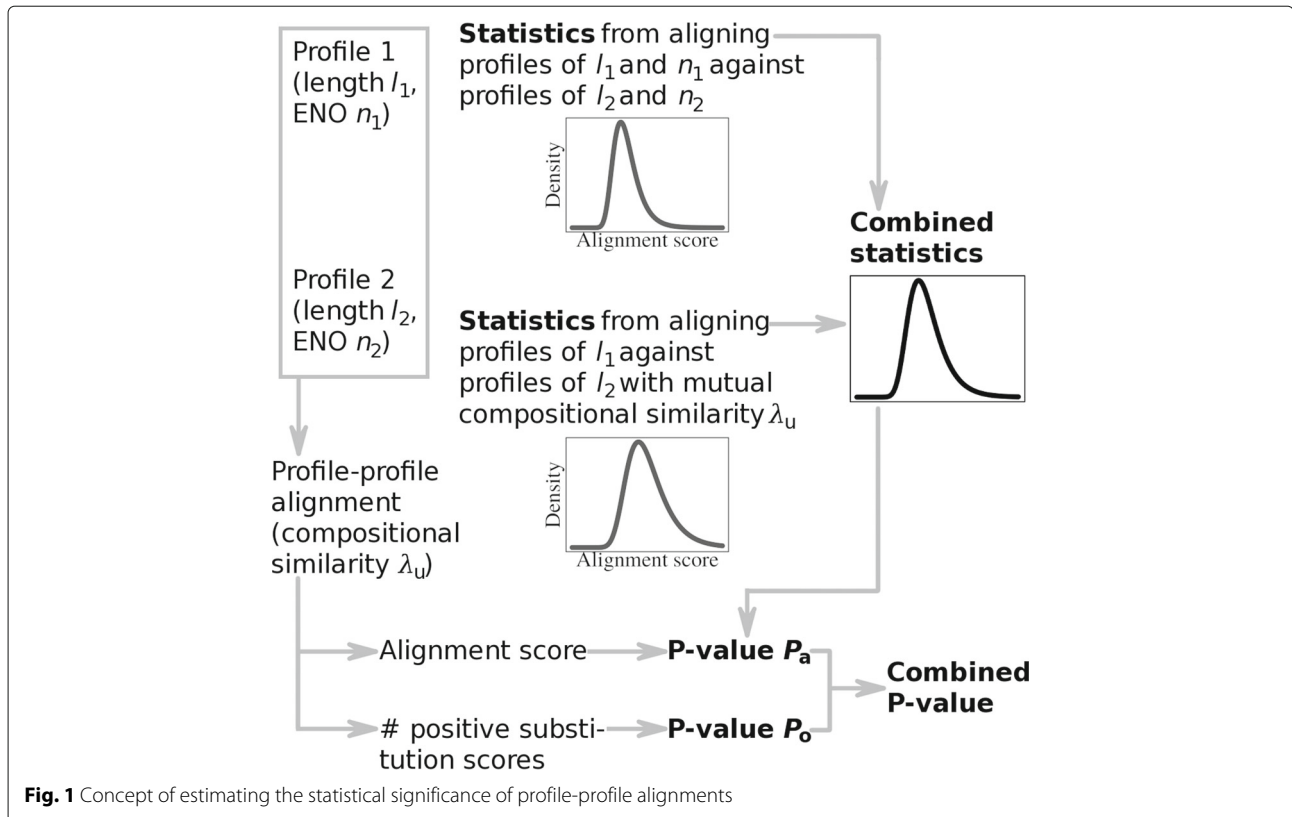
In this work, the null model of protein sequence family does not explicitly incorporate correlations between families and leaves the profile-profile scoring function unchanged. Instead, we randomly generate profiles with properties of real unrelated profiles and take into account the correlations between them by establishing the dependence of statistical parameters on the compositional similarity between the profiles. The null model affects the probability of a score of a pair of profile positions, yet the change in composition does not alter the type of alignment score distribution but values of statistical parameters. The approach proposed here predicts differences in these values.

Overall, this work purposes a procedure for estimating the statistical significance of profile-profile alignments, regarding the effects and factors that influence alignment score distribution, such as edge effects, compositional similarity, and profile attributes (length and the effective number of observations). The procedure builds on an algorithm proposed for generating profiles that resemble real sequence families. The results of the implemented procedure give rise to an analysis of its impact on sensitivity and alignment quality, which we also provide in this paper.

## Methods

The concept of estimating the statistical significance of profile-profile alignments, as proposed here, is illustrated in Fig. 1.

The significance of the alignment between a profile of length  $l_1$  and effective number of observations (ENO)  $n_1$  and another profile of length  $l_2$  and ENO  $n_2$  follows from combining the significance of the alignment score and that of the number of positive substitution scores in the alignment. The distribution of alignment scores depends on both the lengths of profiles being compared and their ENOs [22] (Supplementary Section S2.1 in Additional file 1 defines ENO, which represents the median number of residues per profile position). The composition of profiles is another factor that affects the distribution of alignment scores, and it is a measure independent of profile ENO. Therefore, we incorporate a measure of compositional similarity calculated between profiles into a statistical model for the distribution of alignment scores. The significance of the alignment score of profile 1 and 2 is then calculated using statistics obtained from aligning unrelated profiles of the same length and ENO and having the same mutual compositional similarity.



The parameters of alignment score distributions are estimated by simulation. We develop algorithms to generate random profiles that express properties of, and their alignment scores distribute similarly to alignment scores of, real unrelated profiles. Dividing random profiles according to (1) their length and ENO and (2) length and mutual compositional similarity provides two categories of profiles. Two categories of distributions of alignment scores, illustrated in Fig. 1, result from aligning the profiles in each category. We develop conditional mean estimators to combine the statistics of these distributions, which captures the dependence of the distribution of alignment scores on profile length, ENO, and the compositional similarity between profiles. Finally, we derive a statistic based on the number of positive profile-profile substitution scores and combine its statistical significance and the significance of the alignment score to obtain the final estimate. These steps are described in detail below. More details can be found in Supplementary Section S4, Additional file 1.

### Compositional similarity

The statistical parameter  $\lambda_u \equiv \lambda$  of the limiting distribution of sequence alignment scores (1) has several related meanings [5, 33, 34].  $\lambda_u$  can also be regarded as a measure of compositional similarity [18]. The value of  $\lambda_u$  found as the positive solution to

$$\sum_k p(s_k) \exp(\lambda_u s_k) = 1, \quad (2)$$

where  $\{s_k\}_k$  represent different values of scores in the substitution matrix and  $p(s_k)$  is the probability of  $s_k$ , will decrease as the number of positive substitution scores increases ( $p(s_k)$  increases for all  $k : s_k > 0$ ). Therefore, low values of  $\lambda_u$  may indicate an increased probability for a high-scoring alignment to occur by chance due to compositionally biased regions in two sequences.

The parameter  $\lambda_u$  calculated for a pair of profiles [19, 21] has the same dependence on composition. We take into account compositional similarity between profiles by specifying the dependence of the distribution of profile-profile alignment scores on it.

### Alignment scores of real unrelated profiles

We analyze the distribution of alignment scores of real unrelated profiles (constructed from multiple sequence alignments, MSAs, of real sequences) for two reasons. One is to determine the type of distribution that describes them. The other reason is that these data provide a reference point for generating random profiles whose alignment scores would follow the same type of distribution and be similarly distributed.

We found that alignment scores of real unrelated profiles follow an EVD and 85% of goodness-of-fit tests do not

reject the null hypothesis. Supplementary Sections S2.3 and S3.1 in Additional file 1, respectively, describe how real unrelated profiles were obtained and detail the comparison of profiles of different length, ENO, and compositional similarity (see also Additional files 2, 3, 8 and 9).

### Profile simulation

It is impractical to use real unrelated profiles to produce sufficient data for any values of profile length, ENO, and compositional similarity. Instead, random profiles are generated.

It has been shown that using a null model to generate realistic random profiles improves remote homology detection [22]. However, using the earlier procedure for generating random profiles may lead to highly correlated profiles when the comparison of generated profiles includes the scoring of predicted SSs (Supplementary Section S3.2, Additional file 1). Alignment scores of random profiles may not then represent a distribution obtained by aligning unrelated profiles. Therefore, the aim is to develop an algorithm for generating random profiles that exhibit properties of real profiles, and the degree of correlation between the random profiles can be controlled.

We achieve that by using a modest number of “seed” profiles constructed for a diverse set of real sequences. Random profiles then result from adding noise to profiles/MSAs generated using these seed profiles as a model. Seeds provide properties of real profiles, whereas noise and randomization represent a means of controlling correlation between random profiles.

An important feature of the algorithm that we propose (Algorithm 1) is that random profiles result from concatenating fixed-length fragments sampled randomly from a set of generated MSAs regardless of SS predictions the fragments entail. Supplementary Section S4.1 in Additional file 1 provides additional details, and Supplementary Algorithm S2 defines how MSAs with added noise are produced using a profile model in step 2 of Algorithm 1.

Algorithm 1 shows that the fragment length  $s$  and the level  $r$  of noise added to  $S \times M$  generated MSAs are the parameters that determine the degree of similarity among resulting random profiles. Too large  $s$  and/or too small  $r$  lead to highly correlated random profiles, while too small  $s$  or too large  $r$  result in divergent profiles and useless alignment statistics.

### Optimizing $s$ and $r$

Based on the results obtained for real unrelated profiles, we find the optimal  $s$  and  $r$  using two criteria: (1) the goodness of fit of the EVD to the empirical distribution of alignment scores and (2) the distance between the distribution function obtained for real unrelated profiles and

---

**Algorithm 1** Generating random profiles of length  $l$  and ENO  $n$  given noise level  $r$  and fragment length  $s$

---

Input:  $S$  sequences.

Output:  $R$  random profiles.

1. Make profiles (seeds) for  $S$  diverse sequences chosen at random for which profiles are obtained to be of a sufficiently large ENO (e.g., 12).
  2. Using each of the  $S$  seed profiles as a model, generate  $M$  MSAs of ENO  $n$  with noise  $r$ .
  3. Using each set of  $S \times M$  generated MSAs of ENO  $n$  as source MSAs, generate a set of  $R$  random MSAs of length  $l$  in the following way:
    - (a) choose a number  $j$  at random uniformly between 1 and  $S \times M$ ;
    - (b) randomly select a fragment of length  $s$  of source MSA  $j$ ;
    - (c) copy and add the selected fragment to a random MSA being generated;
    - (d) repeat steps (3a)–(3c) until the length of the random MSA becomes  $l$ .
  4. Construct profiles from the  $R$  generated random MSAs.
- 

the distribution function obtained from simulations. In this way, profiles generated using the optimal  $s$  and  $r$  possess features characteristic to real profiles and correlations between them do not dominate.

We calculate the supremum class upper tail Anderson-Darling statistic  $AD_{\text{up}}$  [35] (Supplementary Section S5.1, Additional file 1) for testing the goodness of fit of the EVD to the data (the first criterion). We normalize it,  $AD_{\text{up}}^* = AD_{\text{up}}/\sqrt{N}$ , by the square root of the number of alignment scores,  $\sqrt{N}$ , to make it independent of sample size. The distance between two empirical distribution functions (the second criterion) is measured by the two-sample Kolmogorov-Smirnov statistic  $D$ .

The whole procedure for finding the optimal  $s$  and  $r$  can be described as follows.

1. Choose values for  $s$  and  $r$ .
2. Using Algorithm 1, generate random profiles of all ENOs and lengths considered.
3. Align simulated profiles.
4. Fit the EVD in the right tail of alignment score distributions.
5. Calculate the  $AD_{\text{up}}^*$  statistic for each alignment score distribution.
6. Train models for predicting the statistical parameters for profiles with given attributes (ENO and length) and compositional similarity.

7. Predict the statistical parameters and calculate the  $p$ -value of each alignment of real unrelated profiles.
8. For each combination of pair values of profile ENO and length, calculate the distance  $D$  between the empirical distribution function obtained for real unrelated profiles and the distribution function obtained in step 7.
9. Repeat steps 1–8 until the optimal  $s$  and  $r$  with respect to  $AD_{up}^*$  and  $D$  have been found.

We discuss prediction of statistical parameters (step 7) below. Here we note that predictions are used to incorporate compositional similarity into the statistical model of profile-profile alignments.

The results (Fig. 2a) obtained using  $S = 1012$  seed profiles ( $M = 1$ ) suggest that considering a small number of different values suffices to find the optimal  $s$  and  $r$ : Large  $s$  increases correlation between profiles, whereas large  $r$  makes them unalignable. Figure 2a shows that the best balance between the two criteria is achieved for  $s = 9$  and  $r = 0.03$ .

Sections S6.1 and S6.2 (Additional file 1 and Additional files 4, 5, 6, 10 and 11) present additional results from simulation experiments. They also show that seeding from three different profiles representing different SCOPe [36] classes leads to similar results (Fig. 2b). Using one seed facilitates profile simulation.

### Conditional mean estimators

Let us consider two profiles, one of length  $l_1$  and ENO  $n_1$  and another one of length  $l_2$  and ENO  $n_2$ . The profiles share compositional similarity  $\lambda_u$ . Then, the significance

of their alignment score is determined from combining statistics obtained (1) from aligning simulated profiles of length  $l_1$  and ENO  $n_1$  against profiles of length  $l_2$  and ENO  $n_2$  and (2) from aligning profiles of length  $l_1$  against profiles of length  $l_2$  with mutual compositional similarity  $\lambda_u$ , as shown in Fig. 1. Here we define this statistical combination.

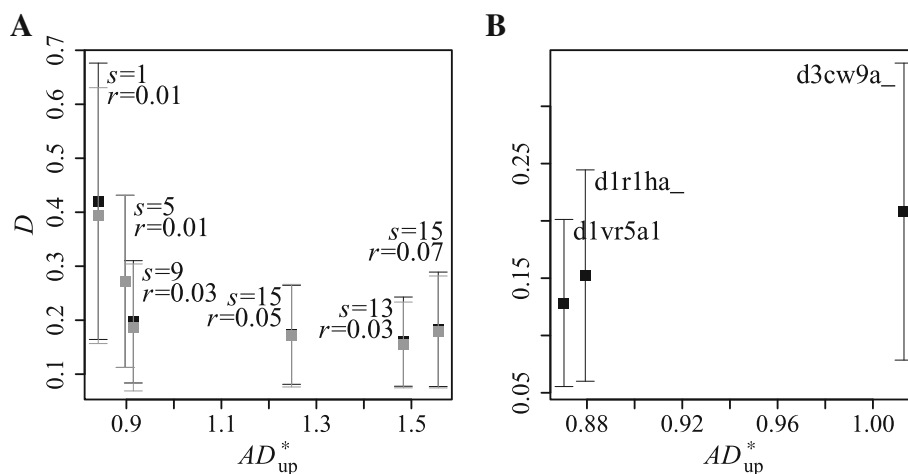
Let A and B index two distributions. The first is obtained for profiles described by the set of attributes  $\{n_1, l_1; n_2, l_2\}$  (i.e., profiles of length  $l_1$  and ENO  $n_1$  have been aligned against profiles of length  $l_2$  and ENO  $n_2$ ) and the other for profiles characterized by the set of attributes  $\{\lambda_u; l_1; l_2\}$ . Assume that distributions A and B belong to the family of EVDs. Let  $\hat{\mu}^A$  and  $\hat{\sigma}^A$  respectively denote the estimates of the location and scale parameters of the EVD corresponding to distribution A. Similarly, let  $\hat{\mu}^B$  and  $\hat{\sigma}^B$  be the estimates of the statistical parameters of distribution B. Then, the conditional mean estimator of the location parameter given two sets of parameters specifying its distribution in settings A and B is

$$\hat{\mu} = a\hat{\mu}^A + (1 - a)\hat{\mu}^B \quad (0 < a < 1), \tag{3}$$

and the corresponding conditional mean estimator of the scale parameter is

$$\hat{\sigma} = b\hat{\sigma}^A + (1 - b)\hat{\sigma}^B \quad (0 < b < 1), \tag{4}$$

where  $a$  and  $b$  are parameters that depend on the parameters specifying the distributions of the location and scale parameters, respectively, in settings A and B (see Supplementary Appendix A in Additional file 1 for details and a proof).



**Fig. 2** Distance between distributions against goodness of fit for different values of  $s$  and  $r$ . Each point represents the average over all distributions obtained for different pair values of profile attributes. Vertical bars represent one standard deviation. Gray and black colors show the results obtained with and without considering the compositional similarity between profiles, respectively. (a): Results obtained using  $S = 1012$  seed profiles. (b): Results for three different seed profiles ( $S = 1$ ) used to generate  $M = 1000$  source MSAs with  $s = 9$  and  $r = 0.03$ . In this case, the results obtained with and without considering compositional similarity coincide

### Prediction of statistical parameters

In simulations (Optimizing  $s$  and  $r$ ), random profiles are generated of predefined lengths  $l \in L$  and ENOs  $n \in N$  and discretized values of mutual compositional similarity  $\lambda_u$  rounded to the nearest multiple of 0.1. The sets  $L = \{50, 100, 200, 400, 600, 800\}$  and  $N = \{2, 4, 6, \dots, 14\}$  represent these predefined values.

Statistical significance for any length  $l$ , ENO  $n$ , and  $\lambda_u$  has to be estimated based on the statistics obtained from the observed distributions. As the statistics depend non-linearly on  $l$ ,  $n$ , and  $\lambda_u$  (Supplementary Section S6, Additional file 1), we use low-complexity artificial neural network (NN) models to predict statistical parameters. The models are trained on the estimates obtained (1) from aligning simulated profiles of length  $l_1 \in L$  and ENO  $n_1 \in N$  against profiles of length  $l_2 \in L$  and ENO  $n_2 \in N$  ( $n_2 < n_1$ ) and (2) from aligning profiles of length  $l_1 \in L$  against profiles of length  $l_2 \in L$  with  $\lambda_u$  discretized. We denote the predictions of the location parameter in these two settings by  $\hat{\mu}^A(n_1, l_1; n_2, l_2)$  and  $\hat{\mu}^B(\lambda_u; l_1; l_2)$ , respectively. The notation for the scale parameter is similar. (See also Supplementary Section S4.4, Additional file 1.)

### Adjustment of predictions

Statistical accuracy does not necessarily correlate with increased sensitivity [18, 22] and high-quality alignment rate. Therefore, to simultaneously improve all these qualities, we introduce simple adjustments to the predicted location  $\hat{\mu}^\nu(\cdot)$  and scale parameters  $\hat{\sigma}^\nu(\cdot)$  ( $\nu = A, B$ ), as shown in Fig. 3.

Observing a high correlation between the scale and the location parameters (Supplementary Section S6.3, Additional file 1), we adjust the scale parameters as follows:

$$\tilde{\sigma}^A = g_s \hat{\mu}^A + g_i, \quad \tilde{\sigma}^B = g_c \left( \exp(\hat{\sigma}^B) - 1 \right). \quad (5)$$

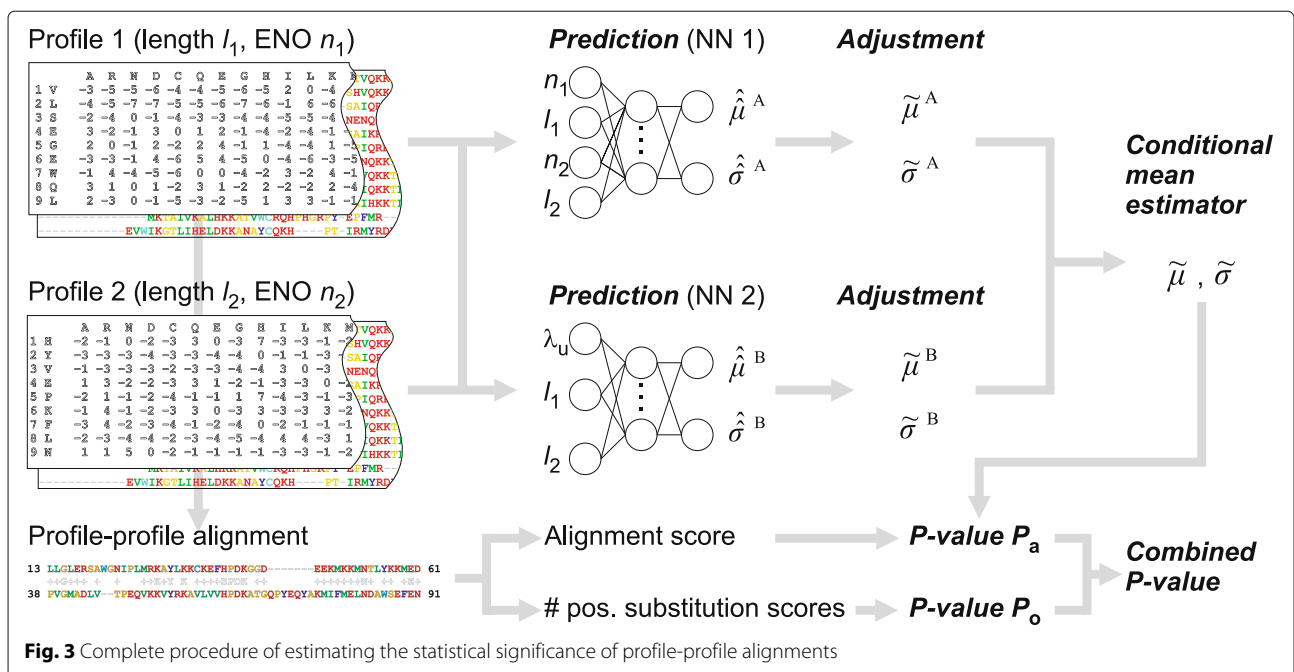
The first equation eliminates the need to predict  $\hat{\sigma}^A$ . The second equation, although expressed non-linearly in  $\hat{\sigma}^B$ , employs one coefficient whose optimal value implies almost the same result as in the case of expressing the scale parameter linearly in  $\hat{\mu}^B$ .

We adjust the location parameters similarly:

$$\tilde{\mu}^A = h_s \hat{\mu}^A + h_i, \quad \tilde{\mu}^B = \hat{\mu}^B + h_c. \quad (6)$$

Note that the adjustments  $\tilde{\mu}^\nu$  and  $\tilde{\sigma}^\nu$  ( $\nu = A, B$ ) retain the dependence of the statistical parameters on the profile length and ENO and compositional similarity. They only scale and shift predictions made by the trained models.

Conditional mean estimators (3) and (4) are used to combine  $\tilde{\mu}^\nu$  and  $\tilde{\sigma}^\nu$  ( $\nu = A, B$ ). The parameters  $a$  and  $b$  and the coefficients  $W = \{g_s, g_i, g_c, h_s, h_i, h_c\}$  in (5) and (6), which we refer to as the adjustment parameters, are optimized with respect to statistical accuracy and alignment quality and sensitivity (Supplementary Section S6.5, Additional file 1). Here we note that the optimal balance is achieved with the compositional similarity between profiles having a large impact ( $a = b = 0.35$ ) on conditional mean estimates.



**Fig. 3** Complete procedure of estimating the statistical significance of profile-profile alignments

### Number of positive substitution scores

The number of positive profile-profile substitution scores in the alignment provides additional information to the alignment score. For example, the same alignment score can be the result of many weakly positive substitution scores or a few high scores. Therefore, it can be a useful indicator for estimating the significance of profile-profile alignments.

The distribution of the number of positive substitution scores can be accurately approximated by the negative binomial distribution (NBD) (Supplementary Section S6.4 in Additional file 1 and Additional files 7 and 12). However, aiming to reduce the overall complexity of the calculation of statistical significance, we propose a derived statistic  $\omega_n$ . It depends on the lengths of profiles being compared and its value (consequently, its significance) decreases as the alignment search space (the product of the profile lengths) increases:

$$\omega_n = \left\lfloor \frac{c_0 \omega}{(l_1 l_2)^{\frac{3}{2}}} \right\rfloor. \quad (7)$$

Here  $\lfloor \cdot \rfloor$  denotes rounding to the nearest integer,  $l_1$  and  $l_2$  are the lengths of two profiles compared,  $\omega$  is the number of positive substitution scores in the alignment, and  $c_0 = 10^5$  is the constant that is equal to the value of the denominator when  $l_1$  and  $l_2$  are slightly less than 50. Hence, the  $\omega_n$  statistic corresponds to the number of positive substitution scores when the search space approximately equals that of two profiles of length 50. The number of positive substitution scores becomes less informative as the search space increases, and, consequently, the value of the  $\omega_n$  statistic decreases. The distribution of  $\omega_n$  is approximately negative binomial (Supplementary Section S6.4, Additional file 1).

The complete procedure for statistical significance estimation is shown in Fig. 3.  $p$ -values of alignment score and  $\omega_n$ ,  $P_a$  and  $P_o$ , are combined using the empirical Brown's method [37] (Supplementary Section S4.5, Additional file 1).

### $E$ -value and its correction

The expected number of local alignments with a score greater than or equal to  $x$ ,  $E = Kl_1 l_N e^{-\lambda x}$ , increases as the size of the search space increases [5, 7, 12, 38]. ( $E$ -value and  $p$ -value  $P$  are related by the equation  $P = 1 - \exp(-E)$ .) Let  $E_N$  be the  $E$ -value of an alignment with score  $x$  obtained from searching a database of size  $l_N$ . Let also  $E_{N'}$  be the  $E$ -value of an alignment with the same score  $x$  (and profile composition) obtained by searching a database of size  $l_{N'}$  with the same query. Then, the relationship between  $E_N$  and  $E_{N'}$  is [18]

$$E_N = \frac{l_N}{l_{N'}} E_{N'}. \quad (8)$$

We use this relationship when compensating for a limited number of randomly generated profiles used in simulations. We refer to the results obtained by calculating the corrected  $E$ -value as the alternative model.

### Results

This section presents results from the application of the proposed procedure (Fig. 3) for estimating the statistical significance of profile-profile alignments. The NN models predicting statistical parameters were trained on the estimates obtained for profiles generated by Algorithm 1 with a noise level  $r = 0.03$  and a fragment length  $s = 9$ , as determined in [Optimizing  \$s\$  and  \$r\$](#) .

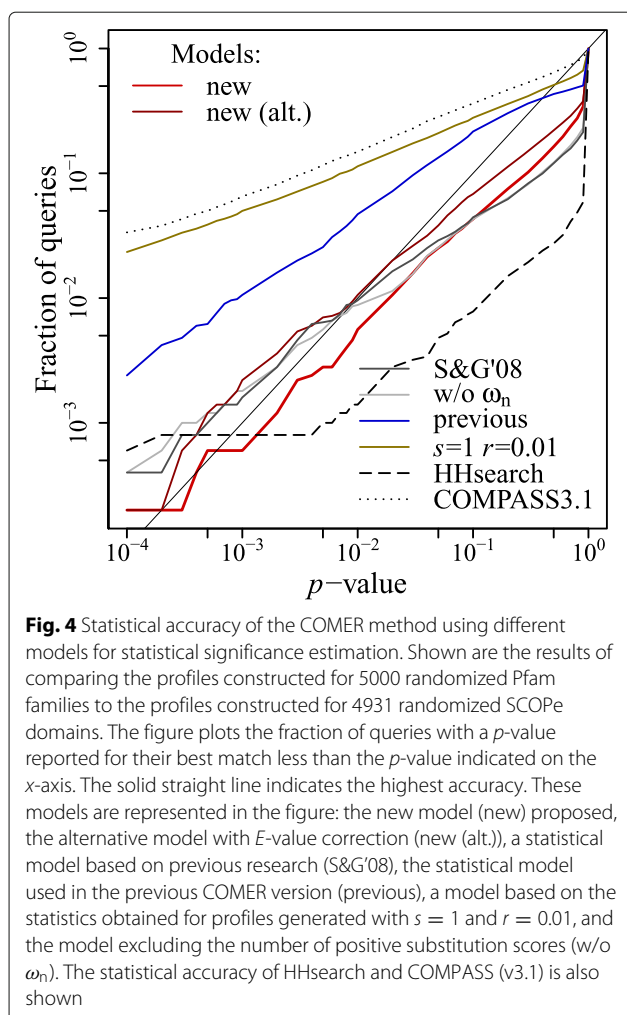
For comparison, we have implemented and present results from the application of the method for estimating statistical significance proposed by [22] (Supplementary Section S5.3.4, Additional file 1). We refer to the COMER version implementing this method as S&G'08 in the text. We also show the results of the application of the COMPASS [19] (v3.1) profile-profile alignment method, where the approach [22] applies to the scoring function for which it was originally developed.

### Statistical accuracy

We used COMER (or another tool) to search and align the profiles constructed for 5000 simulated Pfam [39] (v30.0) MSAs against the database of simulated profiles representing 4931 SCOPe (v2.03) domains. The profiles were randomized using Algorithm 1 to preserve length, ENO, and composition inherent in each of the Pfam MSAs and profiles constructed for the SCOPe domains (see Supplementary Section S5.2, Additional file 1). Then, we calculated the fraction of queries with a  $p$ -value reported by COMER (or another tool) for their best match less than the specified  $p$ -value. The expected number of such queries is the product of the given  $p$ -value and the total number of queries. Therefore, the correspondence between the fraction of queries and the given  $p$ -value represents the theoretical result.

The results are shown in Fig. 4. The new model for statistical significance estimation is more accurate than the model [21] implemented in the previous version [27] of the COMER method. The results also show that the statistics obtained from aligning profiles generated by randomly permuting MSA columns ( $s = 1$ ,  $r = 0.01$ ) lead to overestimation of statistical significance. This result can be accounted for by unrealistic representation of protein sequence families (profiles).

Finally, although the models with and without the  $\omega_n$  statistic taken into consideration achieve comparable statistical accuracy, the sensitivity and high-quality alignment rate obtained using the latter model (w/o  $\omega_n$ ) are lower. The same applies to the model implemented based on the previous research (S&G'08) (see below).



### Sensitivity and alignment quality

We compared the performance of the new COMER version with that of its previous version [27], where the new and previous versions differed only in how they estimated the statistical significance of the same alignments. For reference, we also provide results for three other profile-profile alignment methods, HHsearch [40] (v3.0.0), FFAS [41], and COMPASS (v3.1) [22].

4900 protein domains of the test dataset from the SCOPe database (v2.03) filtered to 20% sequence identity was used to evaluate performance. The test and training datasets shared no common folds.

Profiles were constructed using two categories of MSAs. The MSAs for each domain sequence were obtained by running PSI-BLAST [33] (v2.2.28+) for six iterations and HHblits [42] for three iterations, respectively, against a filtered UniProt database [43].

A pair of aligned domains (profiles) that belonged to the same SCOPe superfamily or shared statistically significant structural similarity (DALI [44] Z-score  $\geq 2$ ) was

considered a true positive (TP). Aligned pairs that did not meet the above criteria but belonged to the same SCOPe fold were considered to have an unknown relationship and were ignored. Other aligned pairs were considered false positives (FPs). The sensitivity was also summarized using the  $ROC_n$  score, which is the normalized area under the ROC curve up to  $n$  FPs.

Alignment quality was evaluated by generating, using MODELLER [45] (v9.4), protein structural models for each produced alignment. Statistically significant similarity between a model and the real structure, TM-score  $\geq 0.4$  [46], was considered to correspond to a high-quality alignment (HQA). An alignment with a TM-score  $< 0.2$ , a characteristic value for a random pair, was assumed to be of low-quality (LQA).

Two evaluation modes were used to evaluate alignment quality. The local mode penalizes alignment overextension, whereas the global mode penalizes too short alignments (e.g., a few amino acids in length). These evaluation modes were also used to evaluate the quality of maximally extended alignments produced by COMER and HHsearch (option -mact set to 0). The evaluation of maximally extended alignments reveals how the quality of local alignments changes with their extension. It provides an indication of the quality of the match between the query and a database protein, which is also important in protein homology modeling.

(Details about the evaluation setting can be found in Supplementary Section S5.3, Additional file 1).

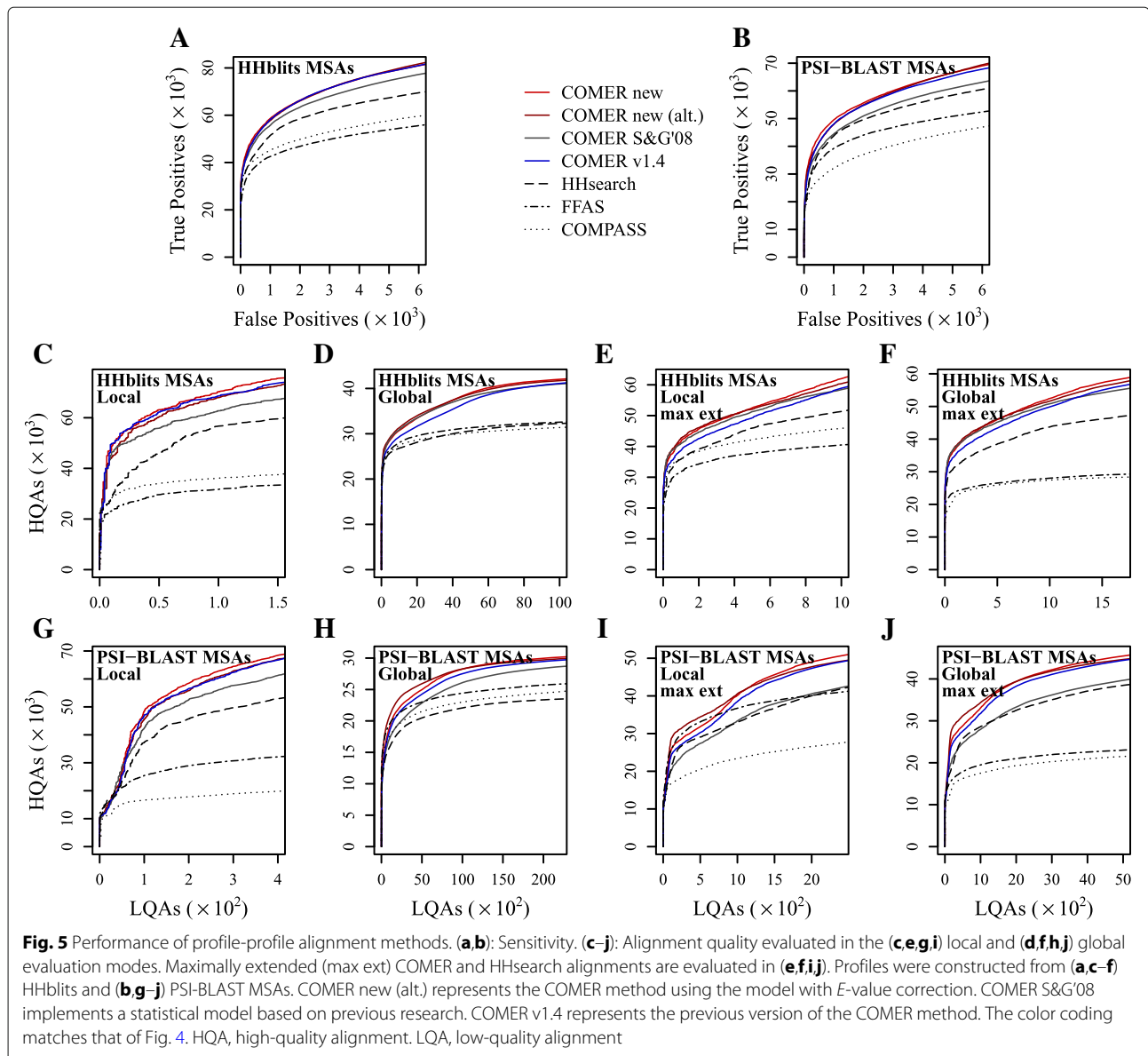
Figure 5 and Table 1 reveal that the new statistical model leads to consistent improvement in both sensitivity and HQA rate and that most of the improvements are statistically significant.

Using the developed statistical model yielded an increase of up to 34.2% and 27.4% in the number of TPs and up to 43.9% and 61.8% in the number of HQAs. Table 2 shows that these percentage increases were achieved at a low false discovery rate (FDR). Hence, the largest relative improvements are expected for relationships detected at a high confidence level.

Table 2 also shows that the relative increases were greater when the profiles were constructed from the PSI-BLAST MSAs. A similar trend was also observed for the number of HQAs examined as a function of the number of alignments of inferior quality (IQAs) (Supplementary Section S7.1, Additional file 1). In this evaluation setting, the developed statistical model showed an increase of up to 102.8% and 193.3% in the number of HQAs.

MSAs obtained from a PSI-BLAST search contain more alignment errors than those built using HHblits. Larger relative improvements in performance achieved when using PSI-BLAST MSAs for profile construction, therefore, show a certain degree of robustness that the new model exhibit. We attribute this characteristic to the





combination of statistical parameters dependent upon both profile length and ENO and compositional similarity. If a high alignment score of two unrelated profiles arises due to false positives in the corresponding MSAs, a high compositional similarity between the profiles counterbalances the statistical significance estimated solely based on profile length and ENO.

By contrast, the version based on the previous approach for estimating statistical significance (S&G'08) displays a large decrease in sensitivity and HQA rate with respect to the previous COMER version (v1.4) for PSI-BLAST MSAs. Among other differences from the new model, that model does not depend explicitly on the compositional similarity between profiles, which in part accounts for this decrease. Although version S&G'08 achieved similar

statistical accuracy (Fig. 4), the sensitivity and the rate of HQAs were consistently lower than those obtained using the new statistical model. (We also provide statistical analysis with respect to FPs found among the top-ranked alignments of the queries from the test dataset, but these results should be interpreted with caution [see Supplementary Section S7.2, Additional file 1]).

#### Application to pairwise profile HMM alignments

We applied the methodology for estimating statistical significance (Fig. 3) to pairwise profile HMM alignments produced by HHsearch. An improvement in both statistical accuracy and sensitivity and HQA rate (Supplementary Section S7.3, Additional file 1) confirms the effectiveness of the developed methodology.

**Table 1** Area under the ROC curve and improvement of the COMER method

Input	Evaluation	x	COMER new		COMER new (alt.)		COMER v1.4
			ROC <sub>x</sub>	Z (p-value)	ROC <sub>x</sub>	Z (p-value)	ROC <sub>x</sub>
HHblits MSAs	Sensitivity	6000	0.837	0.9 (0.350)	0.836	0.6 (0.571)	0.835
PSI-BLAST MSAs	Sensitivity	6000	0.705	5.6 (1.7 × 10 <sup>-8</sup> )	0.698	2.6 (0.009)	0.690
HHblits MSAs	Local	150	0.782	0.9 (0.350)	0.753	-0.9 (0.352)	0.768
	Global	10000	0.456	8.6 (<3 × 10 <sup>-16</sup> )	0.456	7.3 (3.8 × 10 <sup>-13</sup> )	0.438
	Local (max ext)	1000	0.628	9.0 (<3 × 10 <sup>-16</sup> )	0.624	7.1 (1.7 × 10 <sup>-12</sup> )	0.591
	Global (max ext)	1700	0.603	7.6 (3.9 × 10 <sup>-14</sup> )	0.598	7.2 (7.7 × 10 <sup>-13</sup> )	0.574
PSI-BLAST MSAs	Local	400	0.642	1.6 (0.103)	0.623	0.2 (0.874)	0.621
	Global	22000	0.330	3.7 (2.2 × 10 <sup>-4</sup> )	0.334	4.9 (1.2 × 10 <sup>-6</sup> )	0.321
	Local (max ext)	2500	0.496	5.3 (1.2 × 10 <sup>-7</sup> )	0.499	6.3 (2.8 × 10 <sup>-10</sup> )	0.477
	Global (max ext)	5000	0.471	5.1 (3.0 × 10 <sup>-7</sup> )	0.474	5.6 (2.1 × 10 <sup>-8</sup> )	0.457

The sensitivity and alignment quality (Local and Global evaluation modes) of versions of the COMER method are evaluated. Maximally extended (max ext) COMER alignments are included in the evaluation. Profiles were constructed from HHblits and PSI-BLAST MSAs. ROC<sub>x</sub> is the ROC score calculated up to x false positives (FPs; Sensitivity) or low-quality alignments (LQAs; alignment quality in the Local and Global modes). The number of FPs or LQAs, x, depends on the evaluation mode (see also Fig. 5). Z is the difference between the areas (ROC<sub>x</sub> scores) obtained for a new and the previous (v1.4) versions of the COMER method, divided by the estimated standard error. The statistical significance of Z is indicated in parentheses

#### Example of reduced significance for a false positive

Reranking alignments using the proposed estimation of statistical significance has been shown to increase sensitivity and the rate of HQAs. This is demonstrated by an example.

Two domains d2hi7b1 (a.29.15.1) and d2cfqa\_ (f.38.1.2) are alpha-helical proteins but represent different SCOPe classes. They have different topology and do not share statistically significant structural similarity.

Alpha-helical structure implies a high compositional similarity  $\lambda_u = 0.296$  between the profiles constructed for the two domains (low values of  $\lambda_u$  correspond to high

compositional similarity; Supplementary Sections S6.1 and S6.2, Additional file 1). Significance estimation dependent upon compositional similarity allowed the new COMER version to correctly remove the alignment from the list of statistically significant alignments. In contrast, the alignment was considered significant by the previous COMER version. Note that both versions estimated the significance of the same alignment with the same score.

#### Discussion

Profiles represent sequence families, and this fact alone suggests that a profile contains information whose content

**Table 2** Increase in the number of true positives or high-quality alignments.

Input	Evaluation	COMER new			COMER new (alt.)		
		TPs (+%)	TP <sub>Sv1.4</sub>	FDR	TPs (+%)	TP <sub>Sv1.4</sub>	FDR
HHblits MSAs	Sensitivity	34902 ( 7.4)	32483	0.001	33474 ( 3.1)	32483	0.001
PSI-BLAST MSAs	Sensitivity	27666 (34.2)	20612	0.002	26266 (27.4)	20612	0.002
HHblits MSAs	Local	75622 ( 2.9)	73526	0.002	0 ( 0)	0	0
	Global	16699 (15.9)	14414	0.001	18670 (29.5)	14414	0.001
	Local (max ext)	46608 ( 8.8)	42857	0.005	38090 (10.8)	34367	0.001
	Global (max ext)	43603 ( 7.6)	40531	0.008	42165 ( 8.9)	38704	0.006
PSI-BLAST MSAs	Local	61222 ( 4.4)	58627	0.004	13423 (13.4)	11842	0.001
	Global	10248 (43.9)	7123	0.001	11523 (61.8)	7123	0.001
	Local (max ext)	25208 ( 8.9)	23146	0.004	28819 (24.5)	23146	0.004
	Global (max ext)	25161 ( 7.3)	23453	0.007	15588 (24.6)	12515	0.003

Shown are the results of the evaluation of the sensitivity and alignment quality (Local and Global evaluation modes) of versions of the COMER method using profiles constructed from HHblits and PSI-BLAST MSAs. TPs stands for the number of true positives (Sensitivity) or high-quality alignments (in the Local and Global evaluation modes) at a specified false discovery rate (FDR) for a new version of the COMER method. TP<sub>Sv1.4</sub> represents the same number for the previous COMER version. The percentage improvement with respect to TP<sub>Sv1.4</sub> is given in parentheses. 0 indicates no improvement

largely depends on the family the profile describes. It means that the degree of (dis)similarity (distance) between unrelated sequence families strongly affects the distribution of scores of alignments between profiles that represent these families. The challenging aspects of characterizing the distribution of profile-profile alignment scores, however, are not limited to diversity across unrelated profiles. Similarity between SS, context and other predictions that accompany profiles complicate the characterization of alignment score distribution too.

The importance of a null model of profiles motivated us to develop an algorithm for generating random profiles. According to the algorithm, profiles are generated by concatenating fixed-length fragments sampled randomly from real unrelated profiles (seeds) with added noise. Controlling the noise level and fragment length allows to produce random profiles that are neither overly divergent nor share an excessive similarity. In this way, random profiles possess features characteristic to real profiles, but correlations that do not allow them to be considered unrelated do not dominate. Alignments between randomly generated profiles allowed us to determine the dependence of statistical parameters on profile length and ENO and also on compositional similarity between profiles.

The results suggest two important implications. First, profiles generated using long fragments (fragment length 9) represent real unrelated profiles much more accurately than do those obtained by randomly sampling profile columns. Statistics obtained from their alignments lead to higher statistical accuracy. On the other hand, randomly sampling of columns destroys higher-order dependencies inherent in real proteins and leads to the opposite result.

Second, the compositional similarity between profiles has a large impact on estimating the statistical significance of alignments. Improvements in sensitivity and HQA rate can be accounted for in part by that a high compositional similarity may indicate that the proteins share common structural elements or the profiles involve false positives. Significance estimation dependent upon compositional similarity, therefore, has a positive effect.

In fact, the issues of profile composition and random profile model are factors that hinder the application of techniques such as importance sampling [47] to accurately estimating statistical parameters. While composition can be measured in different ways, choosing a random profile model is more complicated because every profile represents a model. For example, the composition of two profiles, one of them obtained by randomly rearranging the positions of the other, will be the same. However, the distribution of alignment scores of profiles generated using the profile with rearranged positions as a seed (model) will differ from that obtained for profiles generated using the other profile as a seed. In this study, close match to the distribution of alignment scores of real unrelated reference

profiles constituted one of the criteria for random profile model selection. However, a supplementary approach might be beneficial.

Based on the above considerations and the results obtained by generating random profiles using one seed profile, we suspect that further improvements may result from introducing a new measure of profile “sequence affinity” (as opposed to the quantitative measure of the effective number of observations). The benefit may come from including into the model the statistical parameters dependent upon the new measure through, e.g., the application of the conditional mean estimator. The present study already demonstrated the conditional mean estimator to be both easily interpretable and effective when the calculation of compositional similarity or divergence between profiles was in use. And it provides possibilities for further improvements.

## Conclusions

We developed a methodology for estimating the statistical significance of profile-profile alignments such that improved statistical accuracy accompanies both increased sensitivity to homologous proteins and rate of high-quality alignments. The combination of statistics dependent upon different profile measures, an integral part of the methodology, may prove useful for future research, including developments of sensitive iterative search methods based on profile-profile comparison, where the importance of controlling false positives is particularly stressed.

## Additional files

**Additional file 1:** Supplementary Materials. Methodological details, derivations, evaluation description, additional simulation and application results. (PDF 5,053 kb)

**Additional file 2:** Figure S1. Distributions of alignment scores of real unrelated profiles for different pair values of profile ENO and length. (PDF 579 kb)

**Additional file 3:** Figure S3. Distributions of alignment scores obtained from aligning pairs of real profiles with mutual compositional similarity  $\lambda$  and different values of length  $l$ . (PDF 467 kb)

**Additional file 4:** Figure S8. Distributions of alignment scores of simulated profiles for different pair values of profile ENO and length. (PDF 6,262 kb)

**Additional file 5:** Figure S9. Distributions of alignment scores obtained from aligning pairs of simulated profiles with mutual compositional similarity  $\lambda$  and different values of length  $l$ . (PDF 1,872 kb)

**Additional file 6:** Figure S10. Observed  $p$ -values corresponding to the empirical distribution function obtained for real unrelated profiles against estimated  $p$ -values using predicted values of the statistical parameters. (PDF 322 kb)

**Additional file 7:** Figure S13. Distributions of the number of positive substitution scores observed in alignments of simulated profiles of different values of ENO and length. (PDF 5,102 kb)

**Additional file 8:** Table S1. Goodness of fit of the EVD to the distribution of alignment scores of real unrelated profiles. (PDF 35 kb)

**Additional file 9:** Table S2. Goodness of fit of the EVD to the distribution of alignment scores of real unrelated profiles. (PDF 35 kb)

**Additional file 10:** Table S3. Goodness of fit of the EVD to the distribution of alignment scores of profiles generated using  $S = 1012$  seed profiles with  $s = 9$  and  $r = 0.03$ . (PDF 65 kb)

**Additional file 11:** Table S4. Goodness of fit of the EVD to the distribution of alignment scores of profiles generated using  $S = 1012$  seed profiles with  $s = 9$  and  $r = 0.03$ . (PDF 44 kb)

**Additional file 12:** Table S5. Goodness of fit of the NBD to the distribution of the number of positive substitution scores observed in alignments of simulated profiles. (PDF 67 kb)

### Abbreviations

ENO: Effective number of observations; EVD: Extreme value distribution; FDR: False discovery rate; FP: False positive; HMM: Hidden Markov model; HQA: High-quality alignment; IQA: Alignment of inferior-quality; LQA: Low-quality alignment; MSA: Multiple sequence alignment; NBD: Negative binomial distribution; NN: (Artificial) neural network; ROC: Receiver operating characteristic; SS: Secondary structure; TP: True positive

### Acknowledgements

Not applicable.

### Authors' contributions

MM conceived, designed and carried out the research, analyzed data, and wrote the manuscript. The author read and approved the final manuscript.

### Funding

This research was funded by the European Regional Development Fund according to a supported activity under Measure No. 01.2.2-LMT-K-718 (Grant No. 01.2.2-LMT-K-718-01-0028). The funding bodies had no role in the design of the study or the collection/analysis/interpretation of data or in writing the manuscript.

### Availability of data and materials

A new version (v1.5.1) of the COMER software that employs the implementation of the new method for the estimation of statistical significance is available at <https://sourceforge.net/projects/comer>. The COMER software is also available on Github at <https://github.com/minmarg/comer> and as a Docker image (<https://hub.docker.com/r/minmar/comer>). Option SSEMODEL allows the user to choose between implemented statistical models. The software package contains programs to generate random MSAs and profiles. Other programs and scripts used in simulations and performance evaluation, dataset information and alignment data are available at <https://sourceforge.net/projects/comer/files/comer-pub-data-1.05>.

### Ethics approval and consent to participate

Not applicable.

### Consent for publication

Not applicable.

### Competing interests

The author declares no competing interests.

Received: 23 March 2019 Accepted: 23 May 2019

Published online: 13 August 2019

### References

- Wang S, Fei S, Wang Z, Li Y, Xu J, Zhao F, Gao X. PredMP: a web server for de novo prediction and visualization of membrane proteins. *Bioinformatics*. 2019;35(4):691–3.
- Li Y, Wang S, Umarov R, Xie B, Fan M, Li L, Gao X. DEEP: sequence-based enzyme EC number prediction by deep learning. *Bioinformatics*. 2018;34(5):760–9.
- Karlin S. Statistical signals in bioinformatics. *Proc Natl Acad Sci USA*. 2005;102(38):13355–62.
- Karlin S, Dembo A, Kawabata T. Statistical composition of high-scoring segments from molecular sequences. *Ann Stat*. 1990;18(2):571–81.
- Karlin S, Altschul SF. Methods for assessing the statistical significance of molecular sequence features by using general scoring schemes. *Proc Natl Acad Sci USA*. 1990;87(6):2264–8.
- Karlin S, Brendel V. Chance and statistical significance in protein and DNA sequence analysis. *Science*. 1992;257(5066):39–49.
- Dembo A, Karlin S, Zeitouni O. Limit distribution of maximal non-aligned two-sequence segmental score. *Ann Probab*. 1994;22(4):2022–39.
- Kotz S, Nadarajah S. *Extreme Value Distributions: Theory and Applications*. London: Imperial College Press; 2000.
- Mott R. Maximum-likelihood estimation of the statistical distribution of Smith-Waterman local sequence similarity scores. *Bull Math Biol*. 1992;54(1):59–75.
- Altschul SF, Gish W. Local alignment statistics. *Methods Enzymol*. 1996;266:460–80.
- Pearson WR. Empirical statistical estimates for sequence similarity searches. *J Mol Biol*. 1998;276(1):71–84.
- Waterman MS, Vingron M. Sequence comparison significance and poisson approximation. *Stat Sci*. 1994;9(3):367–81.
- Waterman MS, Vingron M. Rapid and accurate estimates of statistical significance for sequence data base searches. *Proc Natl Acad Sci USA*. 1994;91(11):4625–8.
- Arratia R, Waterman MS. A phase transition for the score in matching random sequences allowing deletions. *Ann Appl Probab*. 1994;4(1):200–25.
- Spang R, Vingron M. Statistics of large-scale sequence searching. *Bioinformatics*. 1998;14(3):279–84.
- Altschul SF, Bunschuh R, Olsen R, Hwa T. The estimation of statistical parameters for local alignment score distributions. *Nucleic Acids Res*. 2001;29(2):351–61.
- Mott R. Accurate formula for P-values of gapped local sequence and profile alignments. *J Mol Biol*. 2000;300(3):649–59.
- Yu YK, Gertz EM, Agarwala R, Schäffer AA, Altschul SF. Retrieval accuracy, statistical significance and compositional similarity in protein sequence database searches. *Nucleic Acids Res*. 2006;34(20):5966–73.
- Sadreyev R, Grishin N. COMPASS: a tool for comparison of multiple protein alignments with assessment of statistical significance. *J Mol Biol*. 2003;326(1):317–36.
- Poleksic A. Island method for estimating the statistical significance of profile-profile alignment scores. *BMC Bioinformatics*. 2009;10:112.
- Margelevičius M, Venclovas Č. Detection of distant evolutionary relationships between protein families using theory of sequence profile-profile comparison. *BMC Bioinformatics*. 2010;11:89.
- Sadreyev RI, Grishin NV. Accurate statistical model of comparison between multiple sequence alignments. *Nucleic Acids Res*. 2008;36(7):2240–8.
- Edgar RC, Sjölander K. A comparison of scoring functions for protein sequence profile alignment. *Bioinformatics*. 2004;20(8):1301–8.
- Wang G, Dunbrack RL. Scoring profile-to-profile sequence alignments. *Protein Sci*. 2004;13(6):1612–26.
- Meng L, Sun F, Zhang X, Waterman MS. Sequence alignment as hypothesis testing. *J Comput Biol*. 2011;18(5):677–91.
- Margelevičius M. Bayesian nonparametrics in protein remote homology search. *Bioinformatics*. 2016;32(18):2744–52.
- Margelevičius M. A low-complexity add-on score for protein remote homology search with COMER. *Bioinformatics*. 2018;34(12):2037–45.
- Yu YK, Hwa T. Statistical significance of probabilistic sequence alignment and related local hidden Markov models. *J Comput Biol*. 2001;8(3):249–82.
- Metzler D. Robust E-values for gapped local alignments. *J Comput Biol*. 2006;13(4):882–96.
- Eddy SR. A probabilistic model of local sequence alignment that simplifies statistical significance estimation. *PLoS Comput Biol*. 2008;4(5):1000069.
- Karlin S, Dembo A. Limit distributions of maximal segmental score among Markov-dependent partial sums. *Adv Appl Probab*. 1992;24(1):113–40.
- Messer PW, Bunschuh R, Vingron M, Arndt PF. Effects of long-range correlations in DNA on sequence alignment score statistics. *J Comput Biol*. 2007;14(5):655–68.
- Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*. 1997;25(17):3389–402.
- Schäffer AA, Aravind L, Madden TL, Shavirin S, Spouge JL, Wolf YI, et al. Improving the accuracy of PSI-BLAST protein database searches with

- composition-based statistics and other refinements. *Nucleic Acids Res.* 2001;29(14):2994–3005.
35. Chernobai A, Rachev ST, Fabozzi FF. Composite goodness-of-fit tests for left-truncated loss samples. In: Lee CF, Lee J, editors. *Handbook of Financial Econometrics and Statistics*. New York: Springer; 2015. p. 575–96.
  36. Fox NK, Brenner SE, Chandonia JM. SCOPe: Structural classification of proteins—extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res.* 2013;42(D1):304–9.
  37. Poole W, Gibbs DL, Shmulevich I, Bernard B, Knijnenburg TA. Combining dependent p-values with an empirical adaptation of Brown's method. *Bioinformatics.* 2016;32(17):430–6.
  38. Spang R, Vingron M. Limits of homology detection by pairwise sequence comparison. *Bioinformatics.* 2001;17(4):338–42.
  39. Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, et al. The Pfam protein families database: towards a more sustainable future. *Nucleic Acids Res.* 2016;44(D1):279–85.
  40. Söding J. Protein homology detection by HMM-HMM comparison. *Bioinformatics.* 2005;21(7):951–60.
  41. Jaroszewski L, Li Z, Cai XH, Weber C, Godzik A. FFAS server: novel features and applications. *Nucleic Acids Res.* 2011;39:38–44.
  42. Remmert M, Biegert A, Hauser A, Söding J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat Methods.* 2012;9(2):173–5.
  43. Suzek BE, Wang Y, Huang H, McGarvey PB, Wu CH. the UniProt Consortium: UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics.* 2015;31(6):926–32.
  44. Holm L, Kääriäinen S, Rosenström P, Schenkel A. Searching protein structure databases with DALI-Lite v.3. *Bioinformatics.* 2008;24(23):2780–1.
  45. Šali A, Blundell TL. Comparative protein modelling by satisfaction of spatial restraints. *J Mol Biol.* 1993;234(3):779–815.
  46. Zhang Y, Skolnick J. Scoring function for automated assessment of protein structure template quality. *Proteins.* 2004;57(4):702–10.
  47. Park Y, Sheetlin S, Spouge JL. Estimating the Gumbel scale parameter for local alignment of random sequences by importance sampling with stopping times. *Ann Stat.* 2009;37(6A):3697–714.

## Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

