

C C A T T 0 1 0 0 0  
G A G G A 0 1 1 0 1  
G A A T T 0 0 1 1 0  
A C A A G 0 0 1 0 0  
T A C C A 0 0 1 1 0  
T T A C A 0 1 0 0 0  
A C C T C 0 0 0 1 0  
A A G G A 0 0 0 0 0  
G A T G A 0 1 1 0 0  
T A G A T 0 0 1 0 0  
G A T G A 1 0 1 0 0  
T G T A G 1 0 0 0 0  
T A G T A 0 0 0 0 0  
G A T A T 1 0 0 0 0  
G A G T G 1 0 0 0 0  
A G A T T 0 0 0 0 0  
G A G T A 0 0 0 0 0  
T G A T G 1 0 0 0 0  
A T T A G 0 0 0 0 0  
T A G A T 0 0 0 0 0  
T A G T A 1 0 0 0 0  
G A G A A 0 0 0 0 0  
G T A T A 0 0 0 0 0  
T A G A T 0 0 0 0 0  
A G A A 0 0 0 0 0  
G A G A 0 0 0 0 0  
A A T 0 0 0 0 0

# Bioinformatics Explained

## Bioinformatics explained: BLAST versus Smith-Waterman

July 4, 2007

## Bioinformatics explained: BLAST versus Smith-Waterman

Database similarity searches are fundamental in bioinformatics as they are some of the best ways of gaining information about unknown genes and proteins and thus being able to characterize and predict functionality based on homology of new sequencing results from the laboratory.

So, database searching and identification of homology is fundamental in bioinformatics today and different methods are in different ways relevant and important to research. Both time and accuracy is of great significance as handling of large-scale projects is required in a suitable time as well as precision should of course not be compromised. The time aspect is of considerable importance as ever growing amounts of data are entering the research laboratories. The growing amounts of data provide us with more and more information and it also becomes more and more important to get the best quality results - we cannot risk overlooking any information hidden in our data material.

The requirements for database similarity searches should of course be reflected by the algorithms used for performing these identifications of homology. This article will therefore look into some advantages and disadvantages of two of the most widely used algorithms within this field, the Smith-Waterman and the BLAST algorithms.

### Commonly used algorithms

Two of the most commonly used algorithms for database similarity searching are BLAST [Altschul et al., 1990] and Smith-Waterman [Smith and Waterman, 1981].

The Smith-Waterman algorithm is a method of database similarity searching which considers the best local alignment between a query sequence and sequences in the database being searched. The Smith-Waterman algorithm allows consideration of indels (insertions/deletions) and compares fragments of arbitrary lengths between two sequences and this way the optimal local alignments are identified [Smith and Waterman, 1981].

See *Bioinformatics explained: Smith-Waterman* on <http://www.clcbio.com/be/> to learn more about the Smith-Waterman algorithm.

BLAST (Basic Local Alignment Search Tool) also identifies homologous sequences by database searching [Altschul et al., 1990]. BLAST identifies the local alignments between sequences by finding short matches and from these initial matches (local) alignments are created. The BLAST algorithm is a development of the Smith-Waterman algorithm suggesting a time-optimized model contrary to the more accurate but timeconsuming calculations of the Smith-Waterman algorithm [Altschul et al., 1990].

See *Bioinformatics explained: BLAST* on <http://www.clcbio.com/be/> to learn more about BLAST.

### Comparison of Smith-Waterman to BLAST

Comparisons between Smith-Waterman and BLAST show notable differences in the two methods according to accuracy and speed of the searches.

An example from Shpaer et al. [Shpaer et al., 1996] compares accuracy between the two search methods by use of a large objective test data set. Accuracy is compared as a means of selectivity and sensitivity of the search and relates to both the number of results obtained and to the quality of these. For the comparison, BLASTP version 1.4 distributed by NCBI [Altschul et al., 1990] and the publicly available software implementation of Smith-Waterman SSEARCH from the FASTA distribution were running the algorithms using default parameter settings. Test data were generated from the PIR release 40 database

where homologous sequences are indicated by the classification of protein super families (SF). Query sequences were selected to represent all nontrivial super families in the PIR release 40 database and 502 protein sequences were chosen to act as queries against the test database consisting of all entries with a SF number assigned in the database [Shpaer et al., 1996]. The method of "missed@equivalence" point [Pearson, 1995, Pearson, 1991] was used to determine accuracy of the search algorithms. Results of this comparison of Smith-Waterman with BLAST on protein sequence similarity shows that the number of false positives as well as false negatives is significantly lower for Smith-Waterman than it is for BLAST and the risk that sequence similarity readily detected by Smith-Waterman will be missed using BLAST is considerable. In short, the conclusion of the comparison between the two sequence similarity search methods is that the Smith-Waterman algorithm performs significantly better than BLAST on accuracy [Shpaer et al., 1996]. Similar results have previously been obtained by Pearson [Pearson, 1995].

The algorithm used is not the only parameter in the measure of accuracy. Also, scoring matrices and penalties can affect this, and for instance optimization of scoring matrices can improve accuracy of the search significantly. This is the case for both the Smith-Waterman and the BLAST methods.

### Data example of results obtained by Smith-Waterman and BLAST searches

Smith-Waterman has more than one time been compared to BLAST and has been proved more accurate, returning more qualified hits. An example is illustrated below.

#### Software

BLAST and Smith-Waterman output is compared by BLASTP implementation [Altschul et al., 1990] in CLC Combined Workbench 3.0 and CLC Bioinformatics Cell 1.04 accelerating Smith-Waterman searches and integrating the algorithm with CLC Combined Workbench. Both algorithms run as local database searches and are run with the BLOSUM62 matrix and an expect value of 10.

#### Test data

The family of Glutathione S-transferases (GST) is a diverse enzyme family having members within the species of plants, animals, and prokaryotes. Elongation factors have shown homology with the GST family and Elongation factor 1 $\gamma$  (EF1 $\gamma$ ) has been shown to contain GST-related domains [Koonin et al., 1994] as well as GST activity has been proved for EF1 $\gamma$  expressed in *Escherichia coli* [Kobayashi et al., 2001]. For comparison of output from BLAST and from Smith-Waterman searches respectively, the human EF1 $\gamma$  (NCBI accession P26641) acts as query sequence in a match against a database containing 100 randomly chosen GST sequences identified by a simple search for this protein family at NCBI.

#### Results

CLC Combined Workbench returns the results from both search methods in a similar way making it easy to compare the quality of the output. Results are shown as graphic or tabular views (see 1).

The Smith-Waterman search returns a result list of 85 hits from the match between human EF1 $\gamma$  against the 100 members of the GST family. Comparatively, the BLAST search only returns a list of 49 hits. There is obviously a difference in identified homology between sequences for the two search methods as illustrated in figure 1 showing partly the tabular output of the two searches.

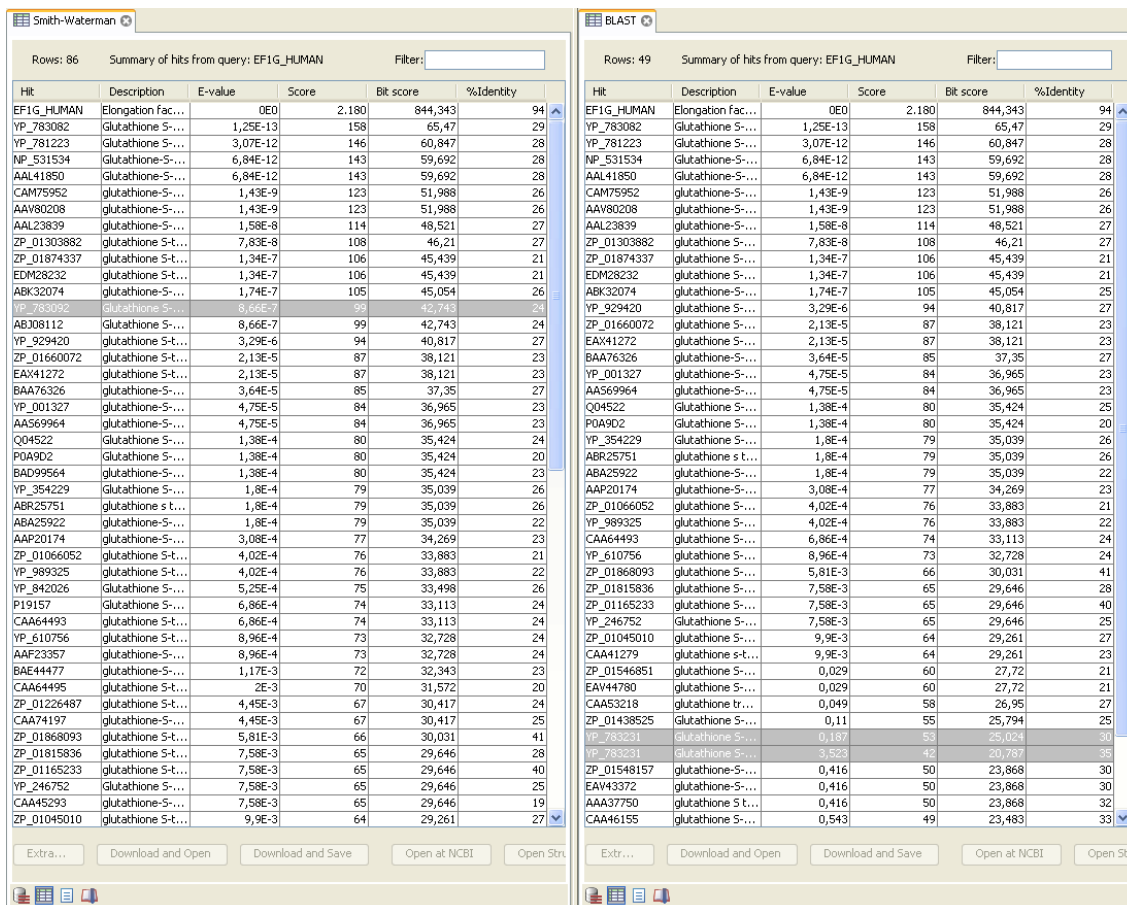


Figure 1: Tabular view of BLAST and Smith-Waterman searches of human Elongation factor  $1\gamma$  matched against a database containing 100 Glutathione S-transferase family members. The results show a significant difference in the number of hits returned by the algorithms. In the left side of the picture is shown the result list returned by the Smith-Waterman search. The row greyed out shows an example of a sequence match identified by Smith-Waterman but not identified by BLAST. The greyed out rows of the BLAST results in the right side of the picture show an example of BLAST returning more hits for a specific comparison of two sequences.

Searching using the Smith-Waterman algorithm clearly identifies a more detailed pattern concerning similarity between query and database sequences. For instance, the sequence returned in row 13 of the result list from the Smith-Waterman search is not identified by BLAST even though the same database has been searched. From this search for matches between the human EF1 $\gamma$  and the 100 GST family members, we expect to identify some degree of similarity as the elongation factor is distantly related to the protein family represented in the database, and as the database consists of members of the same family, a significant number of homologous are expected.

Furthermore, the results returned by the BLAST search contains in some cases more hits for each sequence, e.g. you get more possible alignments for the same sequence match of query and database sequence. This is explained by the fact that BLAST identifies local alignments by initially finding short matches between two sequences not taking the entire sequence into account whereas Smith-Waterman identifies the optimal local alignment for each two sequences compared. This exemplifies how a match identified by BLAST is not necessarily the optimal alignment between query and database sequence whereas the Smith-Waterman algorithm returns one - the optimal - local alignment for each pair of compared sequences.

Elongation factors seem to share similarities with the Glutathione S-transferase family at sequence level. The query sequence is though only distantly related to the Glutathione S-transferase proteins in the test database and we see how a BLAST search might miss important information and identification of sequence relatedness which is readily detected from the Smith-Waterman search.

### **Smith-Waterman: Advantages and disadvantages in review**

Sequence similarity searches performed using the Smith-Waterman algorithm guarantees you *the* optimal local alignments between query and database sequences. Thus, you are ensured the best performance on accuracy and the most precise results - aspects of significant importance when you cannot afford to miss any information gained from the similarity search as e.g. when searching for remote homology. The Smith-Waterman algorithm being the most sensitive algorithm for detection of sequence similarity has however some costs. Time is a considerable disadvantage and performing a Smith-Waterman search is both time consuming and computer power intensive.

### **BLAST: Advantages and disadvantages in review**

The algorithm behind BLAST increases speed of the database searches compared to the Smith-Waterman algorithm. Similarity between two sequences using BLAST is determined by identifying initial short matches and starting local alignments from these matches. Some matches between query sequences and database sequences may be missed by BLAST, and the method does not *guarantee* identification of the optimal alignment between query and database sequence. This might be a disadvantage but often the results obtained by use of the BLAST algorithm may be sufficient and BLAST has, however, some qualities appealing to research. The method is fast, results are returned in a short time and the tool has become the *de facto* standard within database similarity searching.

### **Should I use Smith-Waterman or BLAST for my sequence similarity searches?**

Both the BLAST and the Smith-Waterman algorithms have qualities for research. BLAST is a very popular algorithm giving you results shortly after starting the search. Smith-Waterman ensures you that you get the optimal local alignment and thereby that you will not miss any important information from your sequence similarity search. The Smith-Waterman algorithm is, however, time consuming and has strong requirements for computer power. So, BLAST can be a good method for initial screening of sequence data when getting an indication of results in a short time is more important than getting the most accurate results. The Smith-Waterman algorithm should, on the other hand, be the choice for your database similarity searches when getting precise results is more important than time.

As demands of improved performance in less time are growing, the similarity search algorithms are still being developed and optimized and there are ways of accelerating the Smith-Waterman algorithm, e.g. based on FPGA chips or by the SIMD technology. By such development the Smith-Waterman algorithm becomes a more reasonable choice for practical database similarity searches as time is reduced significantly.

### **Other useful resources**

Public available Smith-Waterman implementation from the Japanese Institute for Bioinformatics Research and Development <http://www-bt1s.jst.go.jp/cgi-bin/Tools/SSEARCH/index.cgi>

Public available Smith-Waterman implementation from the FASTA distribution [http://fasta.bioch.virginia.edu/fasta\\_www2/fasta\\_www.cgi?rm=select&pgm=sw](http://fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi?rm=select&pgm=sw)

The BLAST web page hosted at NCBI <http://www.ncbi.nlm.nih.gov/BLAST>

*Bioinformatics explained: Smith-Waterman* <http://www.clcbio.com/be/>

*Bioinformatics explained: BLAST* <http://www.clcbio.com/be/>

CLC Combined Workbench <http://www.clcbio.com/combined>

CLC Bioinformatics Cell <http://www.clccell.com>

### **Creative Commons License**

All CLC bio's scientific articles are licensed under a Creative Commons Attribution-NonCommercial-NoDerivs 2.5 License. You are free to copy, distribute, display, and use the work for educational purposes, under the following conditions: You must attribute the work in its original form and "CLC bio" has to be clearly labeled as author and provider of the work. You may not use this work for commercial purposes. You may not alter, transform, nor build upon this work.



See <http://creativecommons.org/licenses/by-nc-nd/2.5/> for more information on how to use the contents.

## References

- [Altschul et al., 1990] Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. (1990). Basic local alignment search tool. *J Mol Biol*, 215(3):403–410.
- [Kobayashi et al., 2001] Kobayashi, S., Kidou, S., and Ejiri, S. (2001). Detection and Characterization of Glutathione S-Transferase Activity in Rice EF-1betabetagamma and EF-1gamma Expressed in Escherichia coli. *Biochemical and Biophysical Research Communications*, 288(3):509–514.
- [Koonin et al., 1994] Koonin, E., MUSHEGIAN, A., TATUSOV, R., ALTSCHUL, S., BRYANT, S., BORK, P., and VALENCIA, A. (1994). Eukaryotic translation elongation factor 1 {gamma} contains a glutathione transferase domain—Study of a diverse, ancient protein superfamily using motif search and structural modeling. *Protein Science*, 3(11):2045.
- [Pearson, 1991] Pearson, W. (1991). Searching protein sequence libraries: comparison of the sensitivity and selectivity of the Smith-Waterman and FASTA algorithms. *Genomics*, 11(3):635–50.
- [Pearson, 1995] Pearson, W. (1995). Comparison of methods for searching protein sequence databases. *Protein Science*, 4(6):1145.
- [Shpaer et al., 1996] Shpaer, E. G., Robinson, M., Yee, D., Candlin, J. D., Mines, R., and Hunkapiller, T. (1996). Sensitivity and selectivity in protein similarity searches: a comparison of smith-waterman in hardware to blast and fasta. *Genomics*, 38(2):179–191.
- [Smith and Waterman, 1981] Smith, T. F. and Waterman, M. S. (1981). Identification of common molecular subsequences. *J Mol Biol*, 147(1):195–197.