

# OVERESTIMATING THE ADMISSIBLE HEURISTIC OF A\* FOR MULTIPLE SEQUENCE ALIGNMENT

Tristan CAZENAVE<sup>1</sup>

<sup>1</sup>*LIASD, Dept Informatique,  
Université Paris 8, France*  
E-mail: [cazenave@ai.univ-paris8.fr](mailto:cazenave@ai.univ-paris8.fr)

## Abstract

Multiple sequence alignment is an important problem in computational biology. A-star is an algorithm that can be used to find exact alignments. We present a simple modification of the A-star algorithm that improves much multiple sequence alignment, both in time and memory, at the cost of a small accuracy loss. It consists in overestimating the admissible heuristic. A typical speedup for random sequences of length two hundred fifty is 47 associated to a memory gain of 13 with an error rate of 0.09%. Concerning real sequences, the speedup can be greater than 13,000 and the memory gain greater than 150, the error rate being in the range from 0.08% to 0.71% for the sequences we have tested. Overestimation can align sequences that are not possible to align with the exact algorithm.

**Keywords:** multiple sequence alignment, A\* search, overestimation

## 1 Introduction

Multiple sequence alignment is one of the most important problem in computational biology. It is used to align DNA and protein sequences. The problem of aligning more than eight sequences takes too much memory for current exact algorithms such as A-star or dynamic programming. Biologists use programs that give an approximate answer to overcome the difficulty of finding exact alignment.

From a search point of view, the problem has properties that are different from other problems that are commonly solved with A-star such as the sliding-tile puzzle, or pathfinding on game maps. It has a branching factor in  $O(2^s)$ , when  $s$  is the number of sequences to align. The state space forms a lattice, and there are many paths that go through the same node.

We propose to improve the basic A-star algorithm with overestimation. We have found that it enables large speedups and memory gains at the cost of a low error rate.

The second section presents the multiple sequence alignment problem, the third section details the modifications to the basic A-star algorithm, the fourth section presents experimental results, the last section concludes.

## 2 Multiple Sequence Alignment

In this section we present the multiple sequence alignment problem, then we show how dynamic programming can be applied to it. We present the approximate algorithms currently used to solve the problem, and we give an overview of the exact algorithms that have been tested on it.

### 2.1 The problem

The multiple sequence alignment problem can be considered as a shortest path problem in a  $s$ -dimensional lattice [2]. Let's first consider the case of dimension two, it consists in aligning two sequences. We can write the letter of the first sequence on the horizontal axis, and the letters of the second sequence on the vertical axis. The path starts at the origin point of the matrix (point (0,0) at the upper left). For each point there are three possible moves: the diagonal, the horizontal, and the vertical moves. A diagonal move is equivalent to aligning two characters of the sequences, an horizontal move is equivalent to aligning a character of the horizontal sequence with a gap in the second sequence, a vertical move aligns a character of the vertical sequence with a gap in the horizontal sequence. All paths stop at the bottom right of the matrix after the last two characters of both sequences have been aligned.

The simple model to evaluate the cost of a move is: 0 for a match (aligning the two same characters), 1 for a mismatch, and 2 for a gap (a gap is represented with a -). The cost of a path is the sum of the costs of its moves.

Much more elaborate models of the cost of a mismatch such as the PAM matrices [3] have been developed and have been widely used, elaborate models of gap penalties have also been studied. In order to keep things simple and

as our interest is the search behavior of the alignment algorithms, we have restricted ourselves to the simple model.

For example, if the first sequence is ACGTGCGCT and the second sequence is ACAGTGCCT the best alignment is:

```
AC-GTGCGCT
ACAGTGC-CT
```

and it has a cost of four.

When aligning  $s$  sequences, the path goes through a  $s$ -dimensional lattice, the branching factor is  $2^s - 1$ , and the cost of a move is the sum of the costs of the moves for each pair of sequences.

## 2.2 Dynamic programming

Dynamic programming can be used to efficiently find solution to the problem of sequence alignment. However if the average length of the sequences to align is  $length$ , and the number of sequences is  $s$ , dynamic programming needs  $O(length^s)$  memory and time. A possible improvement trades off time for space [5] but it still requires  $O(length^{s-1})$  memory which is still too much for aligning many sequences.

Dynamic programming is very useful to compute pairwise alignments. The memory needed to compute the pairwise alignment is manageable and it gives precise information to the admissible heuristic of A-star.

Table 1 gives the dynamic programming table corresponding to the previous two sequences. For each pair of index in the two sequences, the table gives the cost of the optimal alignment starting at the corresponding position. Each entry is computed as the minimum among the entry to the left plus 2, the entry below plus 2 and the entry to the lower left diagonal plus the cost of aligning the two corresponding letters. The table is computed starting at the lower right and going from the right to the left and then bottom-up. The upper left cell contains four which is the cost of the optimal alignment.

## 2.3 Approximate algorithms

The programs currently used by biologist such as CLUSTAL W [14] and DCA [13] find sub-optimal alignments. They consist in series of progressive pairwise alignments.

Other approximate algorithms use metaheuristics or evolutionary methods [10].

	A	C	G	T	G	C	G	C	T	
A	4	6	7	8	9	10	12	14	16	18
C	4	4	5	6	7	8	10	12	14	16
A	5	3	4	4	5	7	8	10	12	14
G	6	4	2	3	3	5	6	8	10	12
T	8	6	4	2	3	3	4	6	8	10
G	10	8	6	4	2	2	2	4	6	8
C	13	11	9	7	5	2	1	2	4	6
C	14	12	10	8	6	4	2	0	2	4
T	16	14	12	10	8	6	4	2	0	2
	18	16	14	12	10	8	6	4	2	0

**Table 1.** Example of a dynamic programming table.

## 2.4 Exact algorithms

A-star was applied to the optimal alignment of multiple sequences by Ikeda and Imai [7]. The admissible heuristic is computed using the dynamic programming tables for the pairwise alignments. Because of the large branching factor of the problem, and the large number of open nodes, A-star cannot align more than seven sequences due to memory limits. To overcome this difficulty and reduce the memory requirements, A-star with partial expansion was proposed [16]. It consists in not memorizing in the open list child nodes that have a  $f$  value greater than the  $f$  value of their parent plus a threshold. Experimental results show that Partial Expansion A-star can align seven sequences using fewer stored nodes than A-star, and can align some eight sequences problems. However the gain in memory is acquired at the cost of a greater search time.

Another refinement was proposed to reduce both the memory and the time requirements, using an octree to represent a three-way heuristic [11]. A close approach is the use of external memory pattern databases using structured duplicate detection [18]. It reduces the memory requirements of the pattern databases by an average factor of 83 times, and makes Sweep-A-star [17] run 16% faster than using traditional pattern databases.

Other researchers have compared A-star and dynamic programming [6].

### 3 A-star

This section presents modifications to the A-star algorithm. In all the paper, the admissible heuristic we have used for A-star is the sum of the pairwise alignments given by the 2-dimensional dynamic programming tables.

In A-star, the value associated to a node (the  $f$  value) is the sum of the cost of the path from the origin to the node (the  $g$  value), and of the admissible heuristic that gives a lower bound on the cost of the remaining path to the goal (the  $h$  value). We have for the node  $n$ :  $f(n) = g(n) + h(n)$ .

The main loop of the A-star algorithm works as follow:

```
void AStar() {
    bool outOfMemory = false;
    Node * node = nodeWithSmallestf();
    while(!node->final() && !outOfMemory){
        outOfMemory = develop(node);
        node = nodeWithSmallestf();
    }
}
```

The function *develop(node)* puts the node in the closed list and puts all the children of the node in the open list, and *nodeWithSmallestf()* sends back a node of the open list with a minimal  $f$ .

The first subsection deals with the efficient choice of the best open node. The second subsection is about efficient duplicate node detection. The third subsection explains how we have used overestimation.

#### 3.1 Choosing the best open node

Naive implementations of A-star use a list to store the open and the closed nodes. In this case, when the program has to find the best open node, it has to go through all the list to find the node with the minimum  $f$ . The cost of using a list is linear in the size of the list.

A more elaborate, and commonly used, implementation of A-star uses a priority queue to represent the open list. A priority queue uses a heap to maintain the nodes sorted. The insertion of a new node, as well as the finding of the best node require a logarithmic time in the size of the list.

Instead of a simple queue, we have used an array of stacks to maintain the open list. The index in the array is the value of  $f$  for the nodes stored in the corresponding stack. The insertion of an element is performed in constant

time, just pushing it on the top of the stack that corresponds to its  $f$  value. Finding the best node is also performed in almost constant time. The smallest  $f$  value (*currentf*) over all the nodes is maintained, and updated each time a node is inserted in the open list. When retrieving the best open node, *currentf* is used to check if the stack at this index has an element. If so the first element is popped and returned as the best node. If the stack is empty, *currentf* is incremented and the test is performed again for the next index in the array of stacks.

### 3.2 Duplicate node detection

In the multiple sequence alignment problem, it is very important not to expand again nodes that are already present in the open or closed lists, with a smaller or equal  $g$ .

A possible implementation is to go through all the closed and open nodes, to verify if the node to insert in the open list is not already present with a smaller or equal  $g$ . When implemented this way, the duplicate node detection takes most of the time of the algorithm.

We propose an implementation of duplicate node detection which uses a transposition table. Transposition tables are often used in game programs so as to memorize the results of search at nodes of the tree [1]. In order to hash position, we have used Zobrist hashing [19]. A position in the state space is defined by its coordinates. There are as many coordinates as there are sequences. For each possible sequence, and each possible coordinate in this sequence, a 64-bit random number has been computed once for all. These static random numbers are computed using the rand function to set each one of the 64 bits of the random number. The Zobrist hashing of a position is the XOR of all the random numbers that correspond to the coordinates of the position. The XOR is used because:

- it is a very fast operation on bits,
- it is incremental: in order to undo the XOR with a number, the only operation needed is to XOR again with this number. When a node is expanded, it is a move to a neighboring position. The program only has to XOR the random numbers of the old coordinates that change, and to XOR also the random numbers of the new coordinates.
- The XOR of random values that have a uniform repartition gives a random value that has a uniform repartition. It is important to lower the collision probability.

Each position in the search space is associated to a 64-bit hashcode. The lowest bits of this hashcode are used to index the position in the transposition table. An entry of the transposition table contains a list of structures. Each structure contains a hashcode and a  $g$  value. When the algorithm detects duplicate, it goes through the list and verifies if an entry with the same hashcode as the current position, and a less or equal  $g$  is present. In that case, the node is cut.

### 3.3 Consequences on the implementation of A-star

In this subsection the implementation of A-star with array of stacks and hashtables is detailed.

The function *nodeWithSmallestf()* sends back a node of the open list with a minimal  $f$ . As the open list is structured as an array of stacks, it consists in sending back the head of the non empty array that corresponds to the minimal  $f$ :

```
Node * nodeWithSmallestf () {
    while (Open [currentf].next == NULL &&
           currentf < MaxLength)
        currentf++;
    return Open [currentf].next;
}
```

The function *develop(node)* puts the node in the closed list and puts all the children of the node in the open list. The function starts with removing the node from the head of the open structure, then it finds the successors, and checks that memory is not exhausted. For all the successors, if the successors is not already visited (i.e. the *elementTable* function sends back true when the node is already in the hash table), it is inserted in the open structure using the *insert(succ)* function. In the end the node is added to the closed list. The function is:

```

bool develop (Node * node) {
    Open [node->f ()].next = node->next;
    if (not getSuccessors (node, listSuccessors))
        return false;
    for all successors succ in listSuccessors
        if (not elementTable (succ, succ->g ()))
            insert (succ)
    node->next = Closed.next;
    Closed.next = node;
    return true;
}

```

The *insert(node)* function inserts a node in the open list. It simply consists in pushing the node in the array that is indexed by its f value:

```

bool insert (Node *node) {
    int f = node->f ();
    Node * tmp = & Open [f];
    node->next = tmp->next;
    tmp->next = node;
    if (currentf > f)
        currentf = f;
    addToHashTable (node);
}

```

The hashtable essentially consists in two functions. The first one detects if an element is in the table. It first finds the index of the node taking the lower bits of its hashcode, then for all nodes stored at this index it verifies if the node has a smaller or equal g and if it has the same position as the parameter node. When it is the case it returns true as the node is present in the hashtable. The pseudo-code is:

```

bool elementTable (Node *node, int g) {
    Transpo * t = &table [node->hash () & SizeTable];
    for all nodes n in t
        if (n->g <= g)
            if (n->position == node->position)
                return true;
    return false;
}

```

The other function of the hashtable consists in adding a node to the table. It also starts finding the index of the node in the hashtable, and then simply adds the node to the list of nodes at this index:

```
void addToHashTable (Node * node) {
    Transpo * t = &table [node->hash () & SizeTable];
    t->addNodeToList (node);
}
```

### 3.4 Overestimation

Overestimation has been used for Sokoban in the program Rolling Stones, adding all the patterns that match instead of only selecting the ones that ensure admissibility [8]. This use of overestimation has helped Rolling Stones solve 52 problems instead of 47 without overestimation.

Pearl has introduced  $\epsilon$ -admissible search [12] which finds solutions with bounded costs.

A related but different approach is to use likely-admissible heuristics [4]. It consists in relaxing the admissibility requirement in a probabilistic sense. Instead of providing an upper bound to the cost, it guarantees to end-up with optimal solutions with a given probability.

In order to overestimate the length of the remaining path, we have used the following  $f$ :  $f(n) = g(n) + w \times h(n)$  where  $w$  is a real number greater than one. Overestimation speeds up A-star at the cost of making it inexact.

The main property of this overestimation is that it more easily develops nodes that have a low  $h$  value first, i. e. nodes that are closer to the goal than in usual A-star. Nodes that have been discarded early in the search stay discarded longer than in usual A-star. So overestimation prefers paths where some search has already been invested, and paths that have a low admissible heuristic at the beginning of the path.

## 4 Experimental results

Experiments use machines with 1 GB of RAM. Given the available memory, we have chosen a limit of 10,000,000 nodes for A-star.

### 4.1 Generation of random test data

In order to test the different algorithms, we have generated random sequences of bases (i.e. strings composed of letters in the {A,C,G,T} alpha-

**Table 2.** Comparison of priority queues and array of stacks.

$s$	<i>algorithm</i>	$\Sigma$ <i>time</i>	$\Sigma$ <i>nodes</i>
5	array	38.54s	3,154,269
5	queue	45.38s	3,310,618
6	array	744.42s	23,045,612
6	queue	888.92s	23,914,925

bet). The tests use sets of strings of length fifty, one hundred, two hundred or two hundred and fifty. This methodology is similar to Korf and Zhang testing methodology [9]. Generating random problems allows to generate a large number of problems, and to easily replicate experiments. For each length, we have generated one hundred problems. Each problem is composed of ten strings.

## 4.2 Array of stacks

Table 2 gives the time and number of nodes used by A-star with a STL priority queue, and A-star with an array of stacks.  $s$  is the number of sequences to align. Each line describes the result of solving one hundred problems with sequences of length fifty.

For five sequences the speed of the array of stacks is 81,844 nodes per second while the speed of the priority queue is 72,953 nodes per second. The array of stacks is 12% faster.

On more complex problems with more nodes, the comparison is even better for the array of stacks: for six sequences it develops 30,957 nodes per second versus 26,903 nodes per second for the priority queue. The array of stacks is 15% faster.

## 4.3 Duplicate node detection

Table 3 gives the time used to solve the one hundred problems with sequences of length fifty, using transposition tables, and using lists. The transposition table uses 65,535 entries. The index of a position is the last 16 bits of its hashcode.

Even for problems with a small number of nodes such as the alignment of four sequences, the transposition table algorithm clearly out performs the list

**Table 3.** Comparison of list and transposition table.

$s$	<i>algorithm</i>	$\Sigma$ <i>time</i>
4	list	124.54s
4	transposition	2.03s

**Table 4.** Upper bounds for one hundred sets of sequences of length fifty.

$s$	$\Sigma$ <i>score</i>
5	37,341
6	56,041
7	78,561
8	104,711
9	135,545
10	168,401

implementation.

On the more complex problem of aligning five sequences, the list implementation takes 2203 seconds for solving the first six problems, when the transposition table implementation takes 2.05 seconds. Lists become even worse for more than five sequences since the number of nodes grows and the list implementation takes time proportional to the square of the number of nodes.

#### 4.4 Straight alignment

In order to find an upper bound to the cost of an alignment and better evaluate overestimation, we have tested the algorithm which consists of aligning all the sequences without introducing gaps. The tests were run for one hundred sets of sequences of length fifty and one hundred sets of sequences of length one hundred. The same sets of sequences are used for testing overestimation. The results are given in table 4 for sequences of length fifty, and in table 5 for sequences of length one hundred. The  $\Sigma$  *score* field gives the sum of the scores of the straight alignment for all the random sequences.

**Table 5.** Upper bounds for one hundred sets of sequences of length one hundred.

$s$	$\Sigma$ score
4	44,894
5	75,061
6	112,500
7	157,548
8	209,945

#### 4.5 Overestimation for random sequences

We tested overestimation for different numbers of sequences, different weights, and different lengths of sequences. Results for sequences of length fifty are given in table 6. The first column gives the number of sequences to align, the second column gives the weight used for overestimation (1.00 corresponds to the exact algorithm), the third column gives the cumulated time used to solve one hundred problems, the fourth column gives the sum of the scores (length of the shortest path) found for each problem, the fifth column gives the sum of the nodes used for solving each problem.

The results for the exact algorithm are not given for eight or more sequences, since the node limit is reached for these problems before the problem is solved.

We can observe that the 1.05 weight is a safe weight. It significantly reduces the time and the number of nodes, while finding alignments that are better than straight alignments and quite close to optimal alignments. The 1.10 and the 1.20 weights sometimes give worse results than the straight alignment, and should be avoided.

Table 7 has been created using table 6. For each number of sequences, and each weight, the speedup, the error and the memory gain are given. The error is calculated dividing the sum of the lengths of the paths found with overestimation by the sum of the lengths of the shortest paths found by the exact algorithm. The memory gain is computed dividing the number of nodes of the exact algorithm by the number of nodes of the approximate algorithm.

We can observe in this table that the memory gains increase fast with the number of sequences, and that the speedups increase even faster with the number of sequences.

We tested the program on sequences of length one hundred. Results are

**Table 6.** Results for one hundred sets of sequences of length fifty.

$s$	$w$	$\Sigma$ time	$\Sigma$ score	$\Sigma$ nodes
5	1.00	38.54s	36,654	3,154,269
5	1.05	2.59s	36,747	556,463
5	1.10	0.57s	37,036	212,282
5	1.20	0.30s	37,320	161,983
6	1.00	744.42s	55,362	23,045,612
6	1.05	24.97s	55,477	2,556,103
6	1.10	1.76s	55,929	496,354
6	1.20	0.85s	56,161	328,479
7	1.00	30,844.41s	77,982	168,829,955
7	1.05	268.40s	78,147	12,052,417
7	1.10	6.47s	78,592	1,247,218
7	1.20	2.17s	78,767	653,935
8	1.05	5,639.71s	104,396	58,990,176
8	1.10	26.17s	104,895	3,151,376
8	1.20	4.62s	104,914	1,305,738
9	1.10	94.02s	134,856	8,153,718
9	1.20	16.68s	134,765	2,586,087
10	1.10	572.74s	168,920	21,489,700
10	1.20	44.12s	168,592	5,193,800

**Table 7.** Gains over the exact algorithm for sequences of length fifty.

$s$	$w$	speedup	error	memory gain
5	1.05	14.88	0.25%	5.67
5	1.10	67.61	1.04%	14.86
5	1.20	128.47	1.82%	19.47
6	1.05	29.81	0.21%	9.02
6	1.10	422.97	1.02%	46.42
6	1.20	875.78	1.44%	70.15
7	1.05	114.92	0.21%	14.01
7	1.10	4,767.30	0.78%	135.36
7	1.20	14,214.00	1.01%	258.17

**Table 8.** Results for one hundred sets of sequences of length one hundred.

<i>s</i>	<i>w</i>	$\Sigma$ <i>time</i>	$\Sigma$ <i>score</i>	$\Sigma$ <i>nodes</i>
4	1.00	65.31s	42,605	5,403,857
4	1.05	3.18s	42,696	732,666
4	1.10	0.58s	43,069	243,311
4	1.20	0.36s	43,587	162,110
5	1.00	4547.52s	72,152	88,548,072
5	1.05	109.81s	72,246	8,045,147
5	1.10	2.90s	72,986	788,188
5	1.20	0.84s	74,167	340,647
6	1.05	6250.38s	109,398	93,687,353
6	1.10	17.59s	110,401	2,865,184
6	1.20	1.91s	112,185	700,769
7	1.10	216.68s	155,251	14,090,458
7	1.20	4.47s	157,674	1,415,208
8	1.10	7,288.41s	207,858	89,290,030
8	1.20	10.43s	210,554	2,839,425

**Table 9.** Gains over the exact algorithm for sequences of length one hundred.

<i>s</i>	<i>w</i>	<i>speedup</i>	<i>error</i>	<i>memory gain</i>
4	1.05	20.54	0.21%	7.38
4	1.10	112.60	1.09%	22.21
4	1.20	181.42	2.30%	33.34
5	1.05	41.41	0.13%	11.01
5	1.10	1568.11	1.15%	112.34
5	1.20	5413.71	2.79%	259.94

given in table 8. We can observe that the 1.10 weight is always better than the straight alignment, while the 1.20 weight becomes worse for seven and eight sequences. The 1.05 weight gives interesting speedups and memory gains for alignments that are close to optimal.

Table 9 gives the gains and the error calculated with table 8. If we compare table 7 with table 9, we can observe that the gains for four sequences in table 9 are similar to the gain for five sequences in table 7. The branching factor is fifteen for four sequences, and thirty one for five sequences, the average length of the shortest path is four hundred twenty six for four sequences of length one hundred, and three hundred sixty seven for five sequences of length fifty. The gains are slightly greater for four sequences of length one hundred than for five sequences of length fifty, with length of the shortest paths which are also slightly greater.

Concerning five sequences of length one hundred, and seven sequences of length fifty, the average length of the shortest paths are respectively seven hundred twenty two and seven hundred eighty, when the gains in memory are equivalent, and the speedups are three times greater for the seven sequences.

The error rates are more important for five sequences of length one hundred than for seven sequences of length fifty, even if the speedup are lower. It is interesting as it shows that speedup and error rates are not always correlated, and that there are portions of the space of problems (denoted by the length of the sequences and the number of sequences) that are more favorable to overestimation than others.

As the overestimation has a better behavior for sequences of length one hundred than for sequences of length fifty, we have tested the algorithm on sequences of length two hundred (tables 10 and 11), and on sequences of length two hundred and fifty (tables 12 and 13).

**Table 10.** Results for one hundred sets of sequences of length two hundred.

$s$	$w$	$\Sigma$ time	$\Sigma$ score	$\Sigma$ nodes
4	1.00	3,664.79s	84,102	91,691,701
4	1.05	97.87s	84,197	7,638,509
4	1.10	2.28s	85,190	611,965
4	1.20	0.99s	86,631	336,527
4	straight	0.06s	90,104	0

**Table 11.** Gains over the exact algorithm for sequences of length two hundred.

$s$	$w$	speedup	error	memory gain
4	1.05	37.44	0.11%	12.00
4	1.10	1607.36	1.29%	149.83
4	1.20	3701.81	3.00%	272.46

For sequences of length two hundred, all weights give much better results than the straight alignment. Moreover, the speedups and the memory gains are also better than for four sequences of length one hundred.

Concerning sequences of length two hundred fifty, the speedups and memory gains are even better, and all the weights give alignments much better than the straight one. For five sequences, the overestimation finds alignments that are much better than the straight one when the exact algorithm exhausts memory. When we have tested the first alignment of five sequences with a weight of 1.05, the node limit was reached and A-star stopped with no solution after 2,818 seconds. Comparatively, a weight of 1.10 found a path of length 1807 in 0.39 seconds and 57,044 nodes (a straight alignment gave 1865). For this problem the speedup was therefore of much more than 7,225 and the memory gain of much more than 175.

In conclusion, in our experiments on random sequences, overestimation gives better results for sequences that are difficult to align.

**Table 12.** Results for one hundred sets of sequences of length two hundred fifty.

<i>s</i>	<i>w</i>	$\Sigma$ <i>time</i>	$\Sigma$ <i>score</i>	$\Sigma$ <i>nodes</i>
4	1.00	12,998.15s	104,565	217,402,086
4	1.05	274.72s	104,660	16,514,033
4	1.10	4.17s	106,002	1,019,825
4	1.20	1.27s	107,831	426,185
4	straight	0.08s	112,060	0
5	1.10	64.88s	179,024	7,012,805
5	1.20	3.44s	183,058	933,319
5	straight	0.15s	187,103	0

**Table 13.** Gains over the exact algorithm for sequences of length two hundred fifty.

<i>s</i>	<i>w</i>	<i>speedup</i>	<i>error</i>	<i>memory gain</i>
4	1.05	47.31	0.09%	13.16
4	1.10	3,117.06	1.37%	213.18
4	1.20	10,234.76	3.12%	510.11

## 4.6 Overestimation for real sequences

We have tested overestimation on real sequences with another machine, a Pentium 2.8 GHz with 1 GB of RAM. The sequences are taken from BaliBase [15]. The results are given in table 14. The name of the sequences are the names in BaliBase. Some of the test sequences have been created by merging the five sequences of 1aho-ref1 and the five sequences of 1csp-ref1, which gives ten sequences. The test4-ref1 file contains all the ten sequences, test3-ref1 contains the first nine sequences (the five sequences of 1aho-ref1 followed by the first four sequences of 1csp-ref1), test2-ref1 contains the first eight sequences, and test-ref1 the first seven.

The *score* field gives the length of the shortest path found for the alignment of the sequences, the *time* field gives the time in seconds, and the *nodes* field the number of nodes used by A-star.

The memory gains and speedups are given in table 15. An impressive result is the alignment in test2-ref1 that cannot be found by A-star, when overestimation finds it in less than 13,000 times less time and less than 150 times less memory.

Another excellent result is the alignment of ten sequences for test4-ref1. We do not have error rates for these two results since the optimal alignment is too hard to find, however it has been calculated for test-ref1 and in this case it is very small (0.08%).

The experiments with real sequences show that overestimation works very well for difficult to align sequences, and is associated with a small error rate.

## 5 Conclusion and Future Work

Overestimation of the admissible heuristic of A-star applied to the multiple sequence alignment problem gives large speedups and memory gain for small error rates. In our tests, it works better with difficult sequences than with easy ones.

Future works include combining our improvement with other heuristics used for exact algorithms such as pattern databases [18], partial expansion [16] and dynamic programming [6].

Another interesting future work is to use neural networks or genetic programming to learn an overestimation function more complex than a simple multiplication by a constant.

Moreover overestimation of an admissible heuristic could certainly be useful in other domains. We will investigate its use in other applications of A\* search.

**Table 14.** Results for real sequences.

<i>name</i>	<i>s</i>	<i>w</i>	<i>score</i>	<i>time</i>	<i>nodes</i>
test4-ref1	10	1.00	unknown	>2,664.1382	>10,000,000
test4-ref1	10	1.05	2,634	15.1980	1,326,779
test3-ref1	9	1.00	unknown	>2,536.0442	>10,000,000
test3-ref1	9	1.05	2,081	0.8571	210,605
test2-ref1	8	1.00	unknown	>2,096.7981	>10,000,000
test2-ref1	8	1.05	1,624	0.1540	63,843
test-ref1	7	1.00	1,220	210.5230	3,373,102
test-ref1	7	1.05	1,221	0.0686	34,972
1bbt3-ref1	5	1.00	unknown	>903.0281	>10,000,000
1bbt3-ref1	5	1.05	1,847	5.4853	721,808
1aboA-ref1	5	1.00	701	3.9505	475,654
1aboA-ref1	5	1.05	706	0.0173	4,808
1ad2-ref1	4	1.00	985	0.1180	45,011
1ad2-ref1	4	1.05	985	0.0060	4,723
1aab-ref1	4	1.00	403	0.0550	25,017
1aab-ref1	4	1.05	405	0.0039	3,422
1aho-ref1	5	1.00	488	0.0403	14,410
1aho-ref1	5	1.05	489	0.0047	2,449
1ar5A-ref1	4	1.00	748	0.0185	5,107
1ar5A-ref1	4	1.05	749	0.0234	3,052
1csp-ref1	5	1.00	412	0.0144	3,127
1csp-ref1	5	1.05	413	0.0120	2,246

**Table 15.** Gains over the exact algorithm for real sequences.

<i>name</i>	<i>s</i>	<i>w</i>	<i>error</i>	<i>speedup</i>	<i>memory gain</i>
test4-ref1	10	1.05	unknown	>175.29	>7.54
test3-ref1	9	1.05	unknown	>2,958.81	>47.48
test2-ref1	8	1.05	unknown	>13,610.39	>156.63
test-ref1	7	1.05	0.08%	3,068.84	96.45
1bbt3-ref1	5	1.05	unknown	>164.62	>13.85
1aboA-ref1	5	1.05	0.71%	228.35	98.93
1ad2-ref1	4	1.05	0%	19.67	9.53
1aab-ref1	4	1.05	0.50%	14.10	7.31
1aho-ref1	5	1.05	0.20%	8.57	5.88
1ar5A-ref1	4	1.05	0.13%	1.26	1.67
1csp-ref1	5	1.05	0.24%	1.20	1.39

## References

- [1] D. Breuker, “Memory versus search in games,” University of Maastricht,” PhD thesis, October 1998.
- [2] H. Carrillo, D. Lipman, “The multiple sequence alignment problem in biology,” *SIAM Journal Applied Mathematics*, vol. 48, pp. 1073–1082, 1988.
- [3] M. O. Dayhoff, R. Schwartz, B. C. Orcutt, “A model of Evolutionary Change in Proteins,” *Atlas of protein sequence and structure*, vol. 5, pp. 345–358, 1978.
- [4] M. Ernandes, M. Gori, “Likely-admissible and sub-symbolic heuristics,” in *ECAI 2004*. Valencia, Spain: IOS Press, 2004, pp. 613–617.
- [5] D. S. Hirschberg, “A linear space algorithm for computing maximal common subsequences,” *Communications of the ACM*, vol. 18, no. 6, pp. 341–343, 1975.
- [6] H. Hohwald, I. Thayer, R. Korf, “Comparing best-first search and dynamic programming for optimal multiple sequence alignment,” in *IJCAI-03*, 2003, pp. 1239–1245.

- [7] T. Ikeda, T. Imai, "Fast A\* algorithms for multiple sequence alignment," in *Genome Informatics Workshop 94*, 1994, pp. 90–99.
- [8] A. Junghanns, J. Schaeffer, "Domain-dependent single-agent search enhancements," in *IJCAI-99*, 1999, pp. 570–575.
- [9] R. E. Korf, W. Zhang, "Divide-and-conquer frontier search applied to optimal sequence alignment," in *AAAI-00*, 2000, pp. 910–916.
- [10] Pawel Kupis, Jacek Mandziuk, "Multiple Sequence Alignment with Evolutionary-Progressive Method," in *ICANNGA (1)*, 2007, pp. 23–30.
- [11] M. McNaughton, P. Lu, J. Schaeffer, D. Szafron, "Memory-efficient A\* heuristics for multiple sequence alignment," in *AAAI-02*, 2002, pp. 737–743.
- [12] J. Pearl, *Heuristics: Intelligent Search Strategies for Computer Problem Solving*. Addison-Wesley, Reading, MA, 1984.
- [13] K. Reinert, J. Stoye, T. Will, "An iterative method for faster sum-of-pairs multiple sequence alignment," *Bioinformatics*, vol. 16, no. 9, pp. 808–814, 2000.
- [14] J. Thompson, D. Higgins, T. Gibson, "CLUSTAL W: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice," *Nucleic Acids Research*, vol. 22, pp. 4673–4680, 1994.
- [15] J. Thompson, F. Plewniak, and O. Poch, "BaliBase: a benchmark alignment database for the evaluation of multiple sequence alignment programs," *Bioinformatics*, vol. 15, pp. 87–88, 1999.
- [16] T. Yoshizumi, T. Miura, T. Ishida, "A\* with partial expansion for large branching factor problems," in *AAAI-00*, 2000, pp. 923–929.
- [17] R. Zhou, E. Hansen, "Sweep A\*: Space-efficient heuristic search in partially ordered graphs," in *Proceedings of 15th IEEE International Conference on Tools with Artificial Intelligence*, 2003, pp. 427–434.
- [18] R. Zhou, E. Hansen, "External-memory pattern databases using structured duplicate detection," in *AAAI-05*, Pittsburgh, PA, July 2005, pp. 1398–1405.
- [19] A. Zobrist, "A new hashing method with applications for game playing," *ICCA Journal*, vol. 13, no. 2, pp. 69–73, 1990.