

TM-1236

Stochastic Motif Extraction using a Genetic
Algorithm with the MDL Principle

by

A. Konagaya & H. Kondou (NEC)

November, 1992

© 1992, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03)3456-3191 ~ 5
Telex ICOT J32964

Institute for New Generation Computer Technology

Stochastic Motif Extraction using a Genetic Algorithm with the MDL Principle

Akihiko Konagaya

Hiroyasu Kondou

C&C Systems Research Laboratories., NEC Corporation
1-1, Miyazaki 4-chome, Miyamae-ku, Kawasaki, Kanagawa 216, Japan

Abstract

This paper proposes a new methodology to extract "stochastic motifs" from protein sequences. Extracting motifs is not trivial because (1) almost all motifs have exceptions, (2) no quantitative criterion has been available so far for good motifs, and (3) combinatorial explosion may occur when searching for all motif candidates.

Instead of pursuing precise motifs, we are trying to extract stochastic motifs that inherently include exceptions, are more stable and suitable for representing important regions. As for the quantitative criterion, we adopt Rissanen's Minimum Description Length (MDL) principle to avoid overfitting to sample sequences. To avoid combinatorial explosion in motif extraction, we adopt a "genetic algorithm", a kind of probabilistic search algorithm based on the biological evolution process. Our experimental results demonstrate that the MDL principle greatly increases the convergence speed of a genetic algorithm when extracting stochastic motifs.

1 Introduction

Recently, some biologists have focused on searching for common patterns in protein sequences which have been preserved in the evolution process. Such patterns are called "motifs" and are considered to represent special biological functions (e.g. Serine proteases and Cysteine proteases) and/or special structures (e.g. Zinc fingers and Leucine zipper consensus)[1]. However, extracting motifs is not trivial because (1) almost all motifs have exceptions, (2) no quantitative criterion has been available so far for good motifs, and (3) combinatorial explosion may occur when searching for all motif candidates.

Common patterns in protein sequences are good approximations of protein functions. A good example is the well-known motif for the heme c binding site in a cytochrome c which plays an important role in the respiratory chain. Figure 1 shows some portions of known cytochrome c sequences for various species. Each character in the sequence corresponds to an amino acid. In most cytochrome c sequences, we can find the common pattern "CXXCH" which represents a cysteine,

Species	Sequence of Cytochrome c
Human	..FIMKCSQCCHTVEK..
Mouse	..FVQKCAQCCHTVEK..
Chicken	..FVQKCSQCCHTVEK..
Snake	..FSMKCGTCHTVEE..
Prawn	..FVQRCACQCHSAQA..
Yeast	..FKTRCLQCCHTVEK..
Homp	..FKTKCAECHTVGR..
Tetrahymena	..FDSQCSACHAIEG..
Rhodospila	..FHITICILCHTDIK..
Microblum	..VFKCKKICHQVGP..
Pseudomonas	..VEKQCMTCCHRAUK..

Figure 1: Some portions of cytochrome c sequences

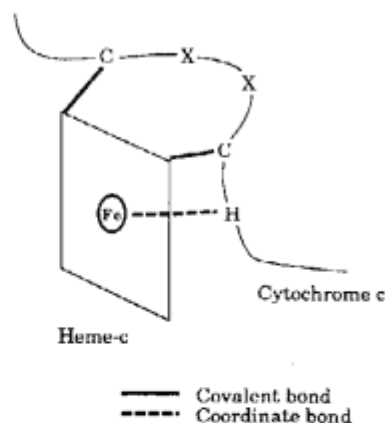


Figure 2: A heme c binding in a cytochrome c

followed by two arbitrary amino acids, followed by another cysteine, followed by a histidine. In this pattern the second "X" does not necessarily coincide with the first "X". The pattern "CXXCH" can be considered as a motif for a cytochrome c, because it corresponds to a protein function; two cysteines and one histidine bind to a heme which cytochrome c holds in the center (Figure 2).

As with other motifs, the pattern "CXXCH" also has exceptions. It does not exist in the cytochrome c of *Euglena*, and the pattern "CXXCH" exists in an adrenodoxin of a pig which is a different category from the cytochrome c. In this case, it would be possible to eliminate such exceptions by introducing more complex patterns. However, one should not expect that more complex patterns always represent protein functions more precisely. This is because more complex rules may cause overfitting to sample data and would not necessarily work better for the discrimination of unknown data, especially in the case of learning stochastic rules from noisy data[3]. This suggests that, instead of pursuing precise motifs, we should try to extract more stable motifs which may contain exceptions but work better for the prediction of unknown data. We call such motifs "stochastic motifs" in this paper. The following example gives the flavor of a stochastic motif. "If the pattern ... "CXXCH" ... is included in the sequence, then the sequence is cytochrome c with probability 130/227 and otherwise it belongs to other protein categories with probability 8072/8076." For the representation of a stochastic motif, we also propose a stochastic decision predicate, which consists of Horn clauses and their probability parameters.

To establish a quantitative criterion for stochastic motifs, Rissanen's MDL principle[2] is adopted. This is because overfitting may occur if we try to extract the stochastic motif that best fits the sample protein sequences. We can easily show that the best fitting stochastic motif is unstable in the sense that it varies according to the sampling of sequences. The MDL principle solves this problem by balancing between the complexity of a motif and its classification errors. It gives a strategy of selecting a "good" stochastic motif on the basis of the sum of the bit lengths required to encode a stochastic motif and its logarithmic likelihood to the sample protein sequences. That is, the principle enables us to compare a simple stochastic motif with classification errors and a complex stochastic motif without classification errors, quantitatively.

To avoid the combinatorial explosion in the motif extraction, we use "genetic algorithms", which are a kind of probabilistic search algorithm based on the biological evolution process. The virtue of genetic algorithms is that they offer an efficient generate-and-test search by means of simple genetic operators that simulate "crossover", "mutation" and "selection". Our experimental results demonstrate that the MDL principle plays an essential role for extracting stable stochastic motifs in terms of convergence speed of genetic algorithms. In fact, a genetic algorithm cannot find stable stochastic motifs without the bias to the complexity of stochastic motifs, that is, with a maximum likelihood method, as far as we have seen in our tests.

The organization of the rest of this paper is as follows. Section 2 gives a representation for stochastic motifs, which we call *Stochastic Decision Predicates*. Section 3 gives a strategy for selecting a good stochastic motif using the MDL principle. Section 4 gives an algorithm for finding optimal stochastic motifs. Section 5 presents the experimental results on extracting stochastic motifs based on our proposed methodology. Finally, in section 6 we discuss current difficulties and future work. This work has been done as a part of the fifth generation computer systems project for the evaluation of the parallel inference machines.

2 Stochastic Decision Predicates

There are many ways to represent stochastic motifs. As a first step for a stochastic representation of motifs, we devised the stochastic decision predicate, a natural extension of a decision list with probabilities. The stochastic decision predicate consists of Horn clauses with probability parameters as follows.

```

motif(S,cytochrome_c) with 137/244.
    :- contain(S,'CXXCH').
motif(S,others) with 9386/9389.

```

The general form is the following.

```

motif(S,C1) (with p1) :- Q1(1) ∧ ... ∧ Qk1(1).
motif(S,C2) (with p2) :- Q1(2) ∧ ... ∧ Qk2(2).
.....
motif(S,Cm-1) (with pm-1) :- Q1(m-1) ∧ ... ∧ Qkm-1(m-1).
motif(S,Cm) (with pm) :- Q1(m) ∧ ... ∧ Qkm(m).

```

Here we call each "motif(S,C_i) (with p_i) :- Q₁⁽ⁱ⁾ ∧ ... ∧ Q_{k_i}⁽ⁱ⁾" a *stochastic clause*. The stochastic clause can be read as S is categorized into C_i with probability p_i if Q₁⁽ⁱ⁾, ..., Q_{k_i}⁽ⁱ⁾ are all true. We assume sequential interpretation of the stochastic clauses in this paper. That is, motif(S,C_i) is selected after motif(S,C₁), ..., motif(S,C_{i-1}) are examined. The body goals Q₁⁽ⁱ⁾ ∧ ... ∧ Q_{k_i}⁽ⁱ⁾ (i = 1, ..., m) represent a condition to discriminate a category C_i when S is given. Each goal Q_j⁽ⁱ⁾ consists of the disjunction of goals R_{1_j}⁽ⁱ⁾, ..., R_{h_j}⁽ⁱ⁾ where R_{h_j}⁽ⁱ⁾ represents some predicate that discriminates a category C_i, such as contain(S,σ) which is true when S contains a pattern σ.

2.1 Semantics of Stochastic Decision Predicate

The semantics of stochastic decision predicates are given from the viewpoint of computational learning theory of stochastic rules[3]. A stochastic decision predicate represents a probabilistic mapping from protein sequences to categories. The probabilistic mapping can be regarded as a conditional probability distribution over the categories when a sequence is given, by introducing a probability structure on the

sequence-category pairs. See the paper [4] for the formal approach to learning stochastic motifs.

3 The MDL Principle in Motif Extraction

In our methodology, the MDL principle gives a new quantitative criterion for “good” stochastic motifs. The most important point is that it enables us to avoid overfitting when extracting stochastic motifs.

For example, as we have shown in the previous section, the pattern “CXXCH” has exceptions in the cytochrome c. It is possible to avoid these exceptions by adding more conjunctions and disjunctions of patterns such as “AAQCH” and “PGTKM”. However, care must be taken so that the obtained result does not become too complex and overfit to the sample sequences. Therefore, we adopt the MDL principle to extract simple but stable stochastic motifs which may contain exceptions rather than precise motifs without exceptions.

The MDL principle originally comes from coding theory in communication. The basic idea is to optimize the number of bits when sending a piece of information, by means of encoding a rule and its exceptions in the piece of information. The MDL principle selects a rule such that minimizes the total bit length of the rule and the exceptions.

The flavor of the MDL principle is the following. Suppose there is a binary string “101101100”. Sending the string requires 9 bits if we do not use any rule. Less bits are sufficient if we compress the string using a rule and its exception. In this case, we can represent the string as three repeats of “10*” and exceptions “110” for the third bit of each repeat instead of * in the rule. The rule requires $\log 3^3 = 4.75^1$ bits since we have to choose on of 3^3 varieties that represent 3-character rules using three kinds of characters. The exception requires $\log 2^3 = 3.0$ bits. The total bits becomes 7.75 bits. We may find a more complex rule to reduce the number of exceptions, but such a rule might require a longer bit length. Therefore, it is important to balance the complexity of the rule and the number of exceptions to reduce the total bit length.

In our methodology, we apply the MDL principle for extracting stochastic motifs in the way proposed by Yamanishi for learning stochastic rules: Yamanishi’s MDL learning algorithm[3]. In his algorithm, the MDL principle selects a stochastic rule that balances the complexity of the stochastic rule and its likelihood of matching the sample data. The rest of this section follows his algorithm with slight modification which mainly comes from the difference of stochastic rule representation, that is, stochastic decision lists and stochastic decision predicates, and some practical reasons for applying the MDL learning algorithm to the motif extraction.

Our methodology selects a stochastic motif that balances the complexity of representation and likelihood of matching the sample sequences. The complexity of a stochastic motif representation is measured by

the description lengths to encode the probability parameters and the Horn clauses of a stochastic decision predicate. The likelihood of a stochastic motif is measured by the description length of likelihood, that is, by the logarithmic likelihood of categories when the sequences are given to the stochastic motif. The description lengths are calculated as follows.

3.1 Description Length of Likelihood

Let $\ell(L)$ be the description length of likelihood given by logarithmic likelihood of categories when sequences are given to a stochastic motif. The likelihood of the categories can be calculated using probabilities associated for categories on each Horn clause in the stochastic motif.

Let $(S_1, C_1), \dots, (S_N, C_N)$ be given N sample sequence and category pairs. Let E_j be the set of sequences which are false for the $1, \dots, j-1$ th clauses and are true for the j th clause. Let N_j be the number of sequences in E_j and let N_j^+ be the number of sequences which are in E_j and belong to the category of the j -th clause. Then the likelihood of the categories (C_1, \dots, C_N) when given sample sequences (S_1, \dots, S_N) with respect to a stochastic predicate with probabilities (p_1, \dots, p_m) , which we denote L , is calculated as follows.

$$L = \prod_{j=1}^m p_j^{N_j^+} (1 - p_j)^{N_j - N_j^+}.$$

The description length $\ell(L)$ is given by $-\log L$ which can be calculated, as follows:

$$\ell(L) = \sum_{i=1}^m N_i \{H(\tilde{p}_i) + D_{KL}(\tilde{p}_i \parallel \hat{p}_i)\} \quad (1)$$

where $\tilde{p}_i = N_i^+ / N_i$ and \hat{p}_i is an estimate of the true parameter p_i^* , which is set to be $\frac{N_i^+ + 1}{N_i + 2}$ (the Bayes estimator) to avoid the difficulties of calculating the description length when $N_i^+ = 0$ or N_i . In addition, $H(\tilde{p}_i)$ and $D_{KL}(\tilde{p}_i \parallel \hat{p}_i)$ are the entropy function and Kullback-Leibler divergence defined as follows.

$$H(\tilde{p}_i) = -\tilde{p}_i \log \tilde{p}_i - (1 - \tilde{p}_i) \log (1 - \tilde{p}_i)$$

$$D_{KL}(\tilde{p}_i \parallel \hat{p}_i) = \tilde{p}_i \log \frac{\tilde{p}_i}{\hat{p}_i} + (1 - \tilde{p}_i) \log \frac{1 - \tilde{p}_i}{1 - \hat{p}_i}$$

The description length $\ell(L)$ indicates the number of bits required to encode the distribution of positive examples and negative examples relative to the stochastic decision predicate. The length varies from near 0 bit², when $p_i = 0$ or 1.0 ($i = 1, \dots, m$), to N bits, when $p_i = 0.5$ ($i = 1, \dots, m$). The former occurs when the stochastic decision predicate completely discriminates the target categories in the given sequences. The latter occurs when the stochastic decision predicate does not contribute to any discrimination of the given sequences.

²It is not appropriate to neglect the value of Kullback-Leibler divergence when the value of entropy function is small.

¹“log” denotes logarithm with base 2.

3.2 Description Length of Probabilities

Let $\ell(P)$ be the description length of the probabilities $\hat{P} = (\hat{p}_1, \dots, \hat{p}_m)$ for a stochastic decision predicate. Since the accuracy (variance) of the maximum likelihood estimator is $O(1/\sqrt{N})$, the description length $\ell(P)$ is given by:

$$\ell(P) = \sum_{i=1}^m \frac{\log N_i}{2} \quad (2)$$

3.3 Description Length of Horn Clauses

Let $\ell(M)$ be the description length of the Horn clauses M . $\ell(M)$ significantly depends on the encoding scheme from Horn clauses to binary strings. The scheme ought to be designed so that the description length can reflect the complexity of the Horn clauses.

In the motif extraction system, $\ell(M)$ is given by:

$$\begin{aligned} \ell(M) = & \sum_{i=1}^m \left[\log^* \left(\sum_{j=1}^{k_i} h_j \right) + \left(\sum_{j=1}^{k_i} h_j - 1 \right) \right. \\ & + \sum_{j=1}^{k_i} \sum_{l=1}^{h_j} \left\{ \log \left(\frac{L_i^j(i)}{X_i^j(i)} \right) \right. \\ & \left. \left. + (L_i^j(i) - X_i^j(i)) * \log(|\mathcal{A}| - 1) \right\} + \log r \right] \end{aligned} \quad (3)$$

where $L_i^j(i)$ and $X_i^j(i)$ are the number of amino acids and of variables, respectively, in the pattern in the l -th predicate in the j -th disjunction region of the i -th clause. On the righthand of (3), the first term denotes the description length of the number of *contain* predicates in the i -th clause. For any $d > 0$, $\log^* d$ denotes $\log d + \log \log d + \dots$ where the sum is taken over all positive terms (Rissanen's integer coding scheme [5]). The second term of (3) denotes the description length to encode the disjunctions and conjunctions occurring in the i -th clause. The third term denotes the description length of the positions of variables in the pattern σ appearing in the predicate '*contain*(S, σ).'. The fourth term denotes the description length required to describe amino acids (not variables) included in the pattern σ appearing in the predicate '*contain*(S, σ).'. $|\mathcal{A}|$ is 20 for amino acids. The last term $\log r$ denotes the description length of the category C appearing in the predicate '*motif*(S, C)'.

3.4 Description Length of Stochastic Motif

By summing (1), (2), and (3), we have the following description length $\ell(T)$ of a stochastic motif represented by a decision predicate:

$$\begin{aligned} \ell(T) \\ \stackrel{\text{def}}{=} \ell(L) + \lambda \{ \ell(P) + \ell(M) \} \end{aligned} \quad (4)$$

where λ is the adjustment parameter. The MDL principle asserts that one should select the stochastic motif which minimizes the description length $\ell(T)$. Notice

here that it is still computationally intractable to find the stochastic motif that minimizes the description length $\ell(T)$ when all possible combinations of Horn clauses are large. Next, we will discuss algorithms to avoid this combinatorial explosion of the search space.

4 Genetic Algorithms

Genetic algorithms are stochastic search algorithms based on the biological evolution process[6]. As in figure 3, genetic algorithms simulate the survival of the fittest in a population of individuals which represent points in a search space. The individuals are represented by binary strings. A function, often called a fitness function, gives values to the binary strings. The aim of a genetic algorithm is to find a global optimum of the fitness function when given an initial population of individuals by applying genetic operators in each generation. The genetic operators consist of the following: crossover, mutation and selection.

Crossover

The crossover operator produces two descendants by exchanging part of two individuals. This operator aims to make a better individual by replacing a part of an individual with a better part of another individual. For example, crossover of the strings "000110" and "110111" at the third position produces the strings "000111" and "110110". The candidates of the crossover operation and the crossover position are randomly chosen.

Mutation

The mutation operator changes certain bit(s) in an individual. For example, the string "000110" becomes "001110" if mutation occurs at the third bit. This operation aims to escape from search spaces from which individuals cannot escape by means of only the crossover operator.

Selection

The selection operator chooses good individuals in a population according to their fitness values and the given selection strategy. This operator aims to increase better individuals in the population while maintaining certain diversity. It simulates the survival of the fittest principle. The operator first calculates the relative fitness of all individuals. Then, several lesser individuals are discarded and the same number of better individuals are duplicated according to their relative fitness values. In case of roulette wheel selection strategy, it selects the next individuals with the probabilities in proportion to their relative fitness values. So, better individuals have a higher chance of remaining or being duplicated but this is not guaranteed.

One of interesting characteristics of our genetic algorithm is in its use of the MDL principle to calculate the fitness value of an individual motif. The MDL length gives the natural relative fitness values in the population, although the smaller the better in this case.

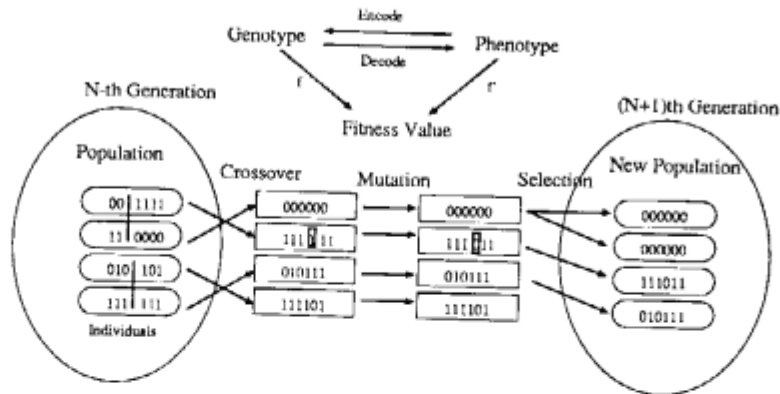


Figure 3: Mechanism of Simple Genetic Algorithms

5 Evaluation

5.1 The Experimental Motif Extraction System

The overview of our experimental motif extraction system is the following. The target hypothesis space is the domain of stochastic decision predicates. The search strategy is based on the MDL principle. The search algorithm is an asynchronous parallel genetic algorithm which consists of the set of subpopulations in which individuals migrate asynchronously. In each subpopulation, individuals represent stochastic decision predicates in the target hypothesis space, and fitness function calculates the corresponding description lengths of the stochastic decision predicates.

The search time depends considerably on the size of the hypothesis space. A large hypothesis space makes it difficult for us to find the optimal stochastic decision predicate in a reasonable time. Therefore, as the first step of motif extraction, we restricted the stochastic predicates to the following forms.

```

motif(S,proteinClass) with p1
:- contain(S,pattern1) and
   contain(S,pattern2) ...
motif(S,others) with p2.

```

That is, we use a predicate *motif* which discriminates the target protein category *proteinClass* from other proteins (*others*) in the database. The discrimination conditions are represented by the conjunction of a predicate *contain*. As the pattern candidates in the *contain* predicate, we adopt 128 patterns that occur frequently in the target proteins.

The mapping from a stochastic decision predicate to a binary string is the following. Each bit corresponds to one of the 128 patterns. A bit 1 represents the occurrence of the pattern in a discrimination condition, and a bit 0 represents the pattern does not occur in the discrimination condition. For example, suppose we use 3-bit length binary strings whose first, sec-

ond, third bits correspond to the pattern "CXXCH", "PXLXG", "GXKM", respectively. Then, the binary string "100" represents the following stochastic decision predicate.

```

motif(S,proteinClass) with p1
:- contain(S,"CXXCH").
motif(S,others) with p2.

```

The binary string "011" represents the following stochastic decision predicate.

```

motif(S,proteinClass) with p1
:- contain(S,"PXLXG") & contain(S,"GXKM").
motif(S,others) with p2.

```

According to this mapping, 128 bits binary strings can express 2^{128} kinds of stochastic decision predicates. As for the genetic operators, we adopt one-point crossover, one-point mutation and roulette wheel selection as described in section 4. The values of other runtime parameters are: the adjustment parameter is 1.0, the number of subpopulations is 63, the subpopulation size is 16, the crossover rate is 1.0, the mutation rate is 0.01 and the migration rate is 0.5, that is, one individual per two generations in average.

5.2 Experimental Results

Table 1 contains some of the stochastic motifs extracted by our experimental system when applied to the protein categories that have more than 10 entries in the Protein Identification Resources (PIR32.0) which currently has 9633 entries³. The rest of results are presented in the appendix.

In table 1, the column *PC* is the super family number of the protein category in PIR32.0. The column *StochasticMotif* is the conjunctions of patterns extracted by our system. The columns $\ell(T)$, $\ell(M)$, $\ell(P)$ and $\ell(L)$ are description lengths of a stochastic motif,

³ Annotated and classified entries by homology in pirl.dat.

Table 1: Results of Stochastic Motif Extraction

PC	Stochastic Motif	$\ell(T)$	$\ell(M)$	$\ell(P)$	$\ell(L)$	E	N_1^+	N_1	N_2^+	N_2
1	CXXCH	309.3441	18.248	10.564	980.693	140	137	244	9356	9389
13	CXXCH&GXGXXC	94.6111	36.383	9.114	47.114	17	16	32	9600	9601
14	IXXXWY&WGXT	47.3127	56.255	8.346	2.715	11	11	11	9622	9622
23	PXXGXXC&CXGXXA	98.7051	40.868	9.179	45.688	33	31	35	9596	9598
23	GXXXC&CXXGXC&TXSC	447.698	55.201	9.435	383.059	91	50	50	9542	9583
29	HXXV&PXXXXXXG	180.516	38.724	9.068	135.124	35	35	30	9593	9603
29	CXRD&RXXXXL&LXXXXY	102.494	56.008	8.877	37.609	23	21	23	9608	9610
33	CXXCXC&GHE	74.458	35.061	8.877	30.520	20	19	23	9609	9610
38	GXXHD&HGD&RPR	72.659	47.879	8.399	16.381	13	12	12	9494	9495
63	AXCXXN&WXXNE	63.554	37.575	9.275	16.503	41	40	40	9592	9593
80	LXXRXN&PXPXXN&PXXXXN	116.070	59.137	8.616	48.316	17	16	16	9614	9617
105	AYXS&MXXYG&YSS	74.436	49.879	8.201	16.355	10	9	9	9628	9624
117	HXXMXP&IPF	46.198	35.061	8.408	2.725	12	12	12	9621	9621

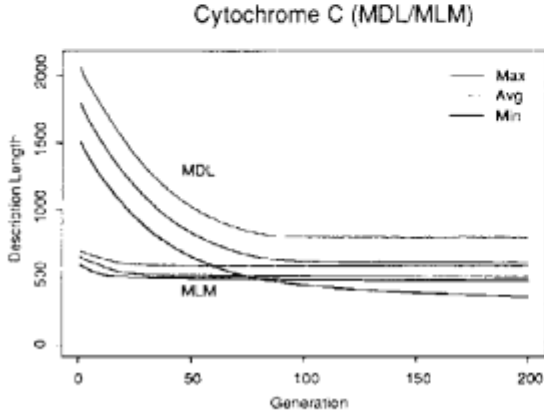


Figure 4: Average description lengths of the best stochastic motif encountered in each generation

Horn clauses, probabilities, and a logarithmic likelihood to the sample sequences.

The column E is the number of target protein sequences in the protein sequence database (PIR). The column N_1, N_2 is the number of protein sequences that become *true* in the first, second clause of a stochastic decision predicate. The column N_1^+, N_2^+ is the number of protein sequences which belong to the target protein category in N_1, N_2 , respectively.

The correspondence between the obtained stochastic motifs and biologically meaningful regions remains as future research issues.

5.3 Comparison of the MDL principle and the Maximum likelihood method

To demonstrate the effectiveness of the MDL principle, various indexes including prediction errors, convergence speed are compared to the maximum likelihood method (MLM). In MLM, good individuals are selected using only the description length of likelihood ($\ell(L)$) without consideration for the complexity of a stochastic decision predicate ($\ell(M) + \ell(P)$).

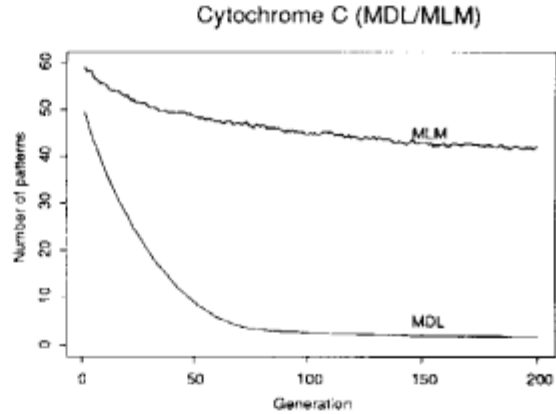


Figure 5: Average number of patterns of the best stochastic motif encountered in each generation

Table 2: Prediction errors for Cytochrome C by Cross Validation Method

	MDL-GA	MLM-GA
$\sum_{i=1}^{10} E_i^+$	3	57
$\sum_{i=1}^{10} E_i^-$	96	0
Total	99	57

Using cross validation technique ([7] p.75-76), the prediction errors can be counted as follows. Let S_i be a disjoint subgroup of protein sequences S for certain N where $S = \cup_{i=1}^N S_i$. Let S'_i be a sample set which removes the i th subgroup from the original protein sequences ($S'_i = S - S_i$). Then, let M_i be a stochastic motif extracted from the sample set S'_i , and count the number of prediction errors E_i^+ and E_i^- using the subgroup S_i as a test set, where E_i^+ shows the number of protein sequences that belong to the target protein category but are not true for the first clause of the stochastic motif M_i . E_i^- shows the number of protein sequences that do not belong to the target protein category but are true for the first clause of the stochastic motif M_i .

Table 2 shows the prediction errors for cytochrome c by the cross validation method when divided into 10 subgroups. The best scored stochastic motif is selected from 50 trials for each subgroups. Each trial requires 200 genetic algorithm generations.

The results show that the stochastic motifs obtained using a genetic algorithm with the MDL principle (MDL-GA) are more stable than the ones obtained using a genetic algorithm with the ML method (MLM-GA). As seen in table 2, the stochastic motifs obtained by the genetic algorithm with the ML method is apparently overfitted to the sample protein sequences. It shows strong discrimination performance for the sample protein sequences ($\sum_{i=1}^{10} E_i^-$), but shows weak predictive performance for the test sequences ($\sum_{i=1}^{10} E_i^+$).

Contrary to our expectations, the result does not come from the intrinsic difference between MDL and MLM, but comes from the difference of convergence speed between MDL-GA and MLM-GA. As in figure 4, MDL-GA shows good convergence speed compared to MLM-GA. That is, MLM-GA is too slow to give us the global optimum in the search space within reasonable time. The difference of the convergence speed comes from the bias caused by the MDL principle. As shown in figure 5, MDL-GA rapidly decreases the number of patterns in the best stochastic motif encountered, while MLM-GA gradually decreases. This is natural since the description length of Horn clauses basically corresponds to the number of patterns. In other words, the MDL principle gives a bias for GA to select individuals with fewer patterns.

Figure 6 shows the effectiveness of the bias for the convergence speed of a genetic algorithm with the MDL principle by changing the adjustment parameter (λ) from 0.5 to 2.0. The histogram in the fig-

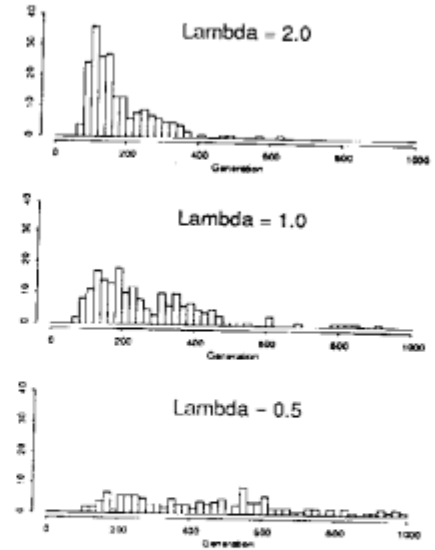


Figure 6: Comparison of convergence speed by the distributions of generations in which the optimal solution is found

ure 6 shows the distribution of generations in which the optimal solution (*CXXCH*) is found. In case of $\lambda = 0.0$, that is, the genetic algorithm with the maximum likelihood method, no optimum stochastic motif is found so far as 10000 generations. In addition, the same stochastic motif has not been extracted even in 10000 trials.

6 Discussion

The following work remains to deal with actual protein sequences on the basis of our methodology.

- The extension of stochastic decision predicate form: In our experience, the number of categories for discrimination is limited to two, that is, the target category and the others. A stochastic decision predicate over two categories can be constructed by concatenating the obtained stochastic clauses for each protein category and recalculating the probabilistic parameter, although it causes another combinatorial problem; in the order of protein categories. Another interesting extension is providing other predicates, such as a distance between patterns.
- Disjunction of patterns: In the current implementation, no form is provided for the disjunction of patterns on the mapping from stochastic decision predicates to binary strings on the genetic algorithm. For example, the pattern "*CXXCH* \vee *AXXCH*" may be more appropriate since it eliminates three exceptions caused by *Euglinae*.

Finding the pattern "AXXCH" is possible if we use (multiple) alignment information of homologous protein sequences.

- More complex patterns: The patterns we used in our experiments are too simple to reflect protein functions. For example, it is a well known fact that in the heme-c binding motif "CXXCH", neither histidine, cysteine, proline nor tryptophan occur in "XX" and small amino acids tend to occur there. To represent such information, more complex patterns are required. Our early experience shows that hidden markov models (HMM) seems to be appropriate for this purpose.
- The handling of category hierarchy: The current MDL principle might select too simple stochastic motifs which have nothing to do with the protein categories. For example, the MDL principle might select only "PGTKM" instead of "CXXCH \wedge PGTKM" for a mitochondria cytochrome c, a subcategory of a cytochrome c. Such selection is tolerable for the purpose of database search, but less effective in the sense that it might lose biological meaning. Such over-simplification can be avoided by adding constraint that reflects category hierarchy.
- Reducing hypothesis space: Since the MDL principle has a bias against selecting complex patterns, it is possible to eliminate complex patterns, for example, more than five patterns from the hypothesis space. One may think it would be faster to search all candidates less than four patterns than to use a genetic algorithm. However, genetic algorithms are also faster if we change their mapping so that it only represents combinations of less than four patterns. In addition, we might bias to the description length of Horn clauses. If this is true, we have to change the adjustment parameter, and also have to search a larger hypothesis space which may include complex patterns more than five patterns. In that case, genetic algorithms would be more powerful tools than conventional search algorithms.
- The handling of point mutations and experimental ambiguity: For example, actual amino acid sequences contain mutation information and special characters that represent ambiguous elements, such as B for asparagine or aspartic acid, and Z for glutamine and glutamic acid. The disjunctive form of stochastic decision predicates may help to some extent. However, such information should be counted for the calculation of description lengths of the stochastic motifs.

7 Conclusion

We have proposed a new methodology for extracting stochastic motifs from protein sequences. Our proposed methodology is characterized by the stochastic representation of motifs using stochastic decision predicates, quantitative criterion using the MDL principle and fast search algorithms using genetic algorithms.

Our experimental results show that the methodology actually produces stable motifs from real protein sequences. The effectiveness of the MDL principle has been statistically proven and compared to the maximum likelihood method, although data are limited to cytochrome c in this paper. We believe the methodology can also be applied to the various kind of discrimination problems on genetic information such as protein sequences. This work has been done as a part of fifth generation computer systems project for the evaluation of parallel inference machines.

Acknowledgements

The authors wish to express their sincere gratitude to Dr. K. Nitta of ICOT, and to Dr. N. Koike of NEC Corporation for their encouragement and support. The authors thank Dr. K. Yamanishi for his technical advices especially for the MDL principle. The authors also thank Mr. K. Yamagishi, Mr. S. Oyanagi, Ms A. Ikeda of NSIS, and Ms Y. Kobayashi and Ms K. Hikita of C&C Systems Research Laboratories, NEC Corporation for their great contribution to our analysis on real protein sequences.

References

- [1] Aitken, Alastair, (1990). *Identification of Protein Consensus Sequences*, Ellis Horwood Series in Biochemistry and Biotechnology.
- [2] Rissanen, J.(1978). Modeling by shortest data description. *Automatica*, 14, 465-471.
- [3] Yamanishi, K.(1990). A learning criterion for stochastic rules. *Proceedings of the 3-rd Annual Workshop on Computational Learning Theory*, (pp. 67-81), Rochester, NY: Morgan Kaufmann. Its full version is to appear in *Jr. on Machine Learning*.
- [4] Yamanishi, K. & Konagaya, A.(1991). Learning Stochastic Motifs from Genetic Sequences. in *Proc. of the Eighth International Workshop of Machine Learning*.
- [5] Rissanen, J.(1983). A universal prior for integers and estimation by minimum description length. *Annals of Statistics*, 11, 416-431.
- [6] Goldberg, D.E., (1989). *Genetic Algorithms in Search, Optimization, and Machine Learning*, Addison-Wesley Publishing Company, Inc.
- [7] Breiman, L., Friedman, J.H., Olshen, R.A., & Stone, C.J.(1984). Classification and regression trees. *Wadsworth Statistics/Probability Series*.

APPENDIX: Continued from Table 1

PC	Stochastic Motif	$\ell(T)$	$\ell(M)$	$\ell(P)$	$\ell(L)$	E	N_1^+	N_1	N_2^+	N_2
118	CSE& QWY	41.255(29.932,	8.569,	2.754),	13,	13,	15,	9618,	9618
119	HXXG& WLF	42.906(31.932,	8.277,	2.698),	10,	10,	10,	9623,	9623
124	QXHXXP& HXGD	59.337(36.253,	8.616,	14.468),	12,	12,	16,	9617,	9617
162	CDP& CXXFXP	46.036(35.061,	8.277,	2.698),	10,	10,	10,	9623,	9623
166	HXXXXH& KXAD	51.720(36.253,	8.520,	6.947),	13,	13,	14,	9619,	9619
177	ENV& PCXXXS& QXRX	67.494(52.201,	8.466,	6.827),	12,	12,	13,	9620,	9620
201	FCT& RXMM	43.451(31.932,	8.735,	2.780),	12,	12,	13,	9614,	9614
251	EXXPF& GPM	51.247(35.233,	8.466,	9.626),	11,	11,	13,	9620,	9620
252	EXXXYW& PXXYW	49.214(37.375,	8.845,	2.783),	22,	22,	22,	9611,	9611
255	DFG& HXUAXXXN	246.366(33.735,	9.985,	202.641),	107,	94,	108,	9512,	9525
261	DANV& WXP	54.242(33.932,	8.677,	11.433),	21,	21,	23,	9610,	9610
304	FXQF& PYH									
	APY& FXQF	52.539(31.932,	8.937,	11.691),	33,	33,	25,	9608,	9608
306	YNGXXV& FXSXY& LYXXI	66.646(55.525,	8.569,	2.754),	15,	15,	15,	9618,	9618
348	YQCE& YQXXC	78.021(33.253,	9.516,	35.251),	57,	55,	56,	9575,	9577
360	NYC& RAH	40.930(29.932,	8.346,	2.713),	11,	11,	11,	9622,	9622
404	CKXXXT& TFXH	54.163(36.253,	9.293,	8.617),	40,	40,	41,	9522,	9592
408	LXXWXX& NDD									
	KNW& LXXWXX									
	EPK& LXXWXX	45.387(34.253,	8.408,	2.725),	12,	12,	12,	9621,	9621
407	DXRXD& FXNHE& NND	43.259(52.201,	8.346,	2.713),	11,	11,	11,	9622,	9622
415	QXXXWXX& LXXXXXXC	78.105(40.190,	8.937,	26.579),	27,	25,	25,	9606,	9606
420	KSC& YCXNI	44.994(33.253,	8.457,	2.804),	25,	25,	25,	9608,	9608
421	MXPN& NQR	61.233(31.932,	8.359,	20.732),	15,	14,	15,	9617,	9617
458	GWDA& CXXDXG	124.490(34.253,	9.435,	80.802),	40,	37,	50,	9580,	9580
473	NSW& YWXXN	45.815(34.253,	8.776,	2.785),	20,	20,	20,	9613,	9613
476	FXXXFD& WXXXC	80.402(38.575,	8.700,	33.326),	19,	17,	18,	9613,	9613
483	DXGA& GAD& GDXDXXQ	98.757(51.686,	8.478,	38.662),	10,	8,	12,	9619,	9621
513	ELXXAD& DXIA& YXPT	56.543(52.679,	9.158,	34.506),	35,	33,	34,	9597,	9599
518	GFA& HAP& VLXXA	116.004(48.879,	8.408,	58.776),	14,	10,	12,	9617,	9621
520	ARXP& IGXGA& AXGCG	123.084(52.879,	9.157,	60.069),	16,	14,	35,	9598,	9600
521	RLXXN& FXLXXL& LXXI	122.429(55.008,	8.677,	38.337),	20,	17,	23,	9607,	9610
546	WXXWXX& PXXLXXP	69.211(39.383,	9.676,	40.852),	73,	70,	70,	9560,	9563
547	KDD& PKLE& YGR									
	NVPE& NYGE& YGR									
	QPP& KDD& NVF									
	FEK& KNYE& NVF	56.076(44.557,	8.739,	2.780),	19,	19,	19,	9614,	9614
552	KXNM& KWR	43.134(31.932,	8.466,	2.730),	13,	13,	13,	9620,	9620
561	GDXX& QKRF& YPG	89.635(44.557,	8.346,	32.732),	12,	10,	11,	9620,	9622
577	HXXW& HWN	44.224(33.253,	8.277,	2.698),	10,	10,	10,	9623,	9623
616	FHR& PXXXXNF	71.520(35.061,	8.411,	27.648),	19,	18,	21,	9611,	9612
617	HSE& YEXXW	44.519(33.253,	8.520,	2.746),	14,	14,	14,	9619,	9619
661	QXRXF& GXXXXXS& NXXGXH	97.367(57.330,	9.293,	30.744),	39,	38,	41,	9591,	9592
694	CXXFX& FYXXC	88.032(38.575,	8.677,	40.581),	24,	20,	20,	9608,	9610
696	PXGG& TYXXXC	101.289(36.253,	8.569,	56.466),	18,	14,	15,	9614,	9618
697	WAXFXXXXXT& VXXMM	126.723(38.575,	8.739,	89.145),	25,	18,	19,	9607,	9614
704	TXCXL& YXXPW	131.305(36.575,	8.659,	46.071),	19,	13,	17,	9610,	9616
708	CXKXN& CXKXGC	46.464(30.384,	8.659,	16.442),	18,	17,	17,	9613,	9616
719	CXKXLE& CXKXQV	109.800(38.575,	8.116,	63.108),	14,	8,	8,	9620,	9625
731	PXXDXW& WKR	67.962(36.253,	8.908,	2.800),	24,	24,	24,	9609,	9609
733	FXXXYX& PFXXXV& YXPT	120.311(50.810,	9.238,	63.258),	43,	38,	38,	9590,	9595
790	QSE& RCF	41.134(29.932,	8.466,	2.730),	13,	13,	13,	9620,	9620
791	PRW& HFXW	89.537(31.932,	8.659,	28.946),	19,	17,	17,	9614,	9616
796	CXGY& CXXYC	60.305(25.253,	8.619,	16.436),	17,	16,	16,	9616,	9617
797	SWXXP& YXXLC	86.021(37.575,	8.569,	40.776),	18,	15,	15,	9615,	9618
802	QXXXXT& GIST& STC	98.444(49.646,	8.408,	40.749),	15,	12,	12,	9618,	9621
806	HXXXXP& MXXXXXXAY	174.417(41.180,	8.992,	194.244),	40,	21,	27,	9587,	9606
809	PXXPG& NXXTR& YNXXXXR	81.207(56.330,	8.466,	16.411),	11,	13,	13,	9619,	9620
812	CXKXGC& LCG	100.669(34.253,	9.688,	58.960),	64,	62,	66,	9563,	9565
813	GWXD& WXP	71.215(29.932,	8.408,	32.878),	13,	11,	12,	9619,	9621
830	CXXXXG& CXXXXXXN	188.469(40.866,	9.981,	137.670),	109,	98,	100,	9522,	9533
833	CXXXXX& CXKXLL	152.972(39.383,	8.659,	104.930),	26,	17,	17,	9607,	9616
851	KSC& YXXCR	57.903(35.253,	8.277,	16.373),	11,	10,	10,	9623,	9623
872	AWXXV& CAW	83.175(33.253,	9.293,	46.629),	39,	37,	41,	9590,	9592
886	DXXXVXO	692.147(30.095,	10.871,	661.181),	266,	237,	279,	9223,	9254
892	CXXXXP& PXXXXXXW& CXXXXXXP	249.788(64.079,	9.477,	217.203),	74,	63,	53,	9559,	9580
893	CXXXX& VXXXXXXP	167.048(30.646,	9.934,	117.478),	69,	67,	69,	9537,	9534
902	PXTXXX& HXXXXV& PXTXXXXXXP	703.051(55.915,	10.878,	632.268),	456,	382,	343,	9178,	9250
905	MQXP& MQI									
	MXHE& MQI	48.698(31.932,	8.592,	27.774),	16,	16,	27,	9606,	9606
908	AKKE& KKA& KXKXK& KXGXO	104.297(70.314,	8.577,	25.106),	22,	21,	23,	9609,	9610
909	AGXXF& FXV& LFP	69.235(51.879,	9.114,	8.243),	31,	31,	33,	9601,	9601
910	FAND& NFXXXD& VTK	62.681(50.879,	8.992,	2.809),	27,	27,	27,	9606,	9606
911	AIH& AQD& PFQ									
	DFK& FQR& PFQ									
	DFK& DTN& FQR									
	AQD& DFK& IAQ									
	AIH& AQD& FQR									
	AIH& FQR& QDF									
912	IRK& KRN& VKR	56.118(44.557,	8.776,	2.785),	20,	20,	20,	9613,	9613
917	PXRR& RPXXXX& SRXXX	59.032(44.557,	8.700,	2.774),	18,	18,	18,	9615,	9615
927	GFXX& HNP& RXPXT	116.476(54.201,	9.238,	53.037),	21,	20,	38,	9594,	9595
943	GK& KXKX& DAXXXG	86.517(51.201,	8.408,	24.908),	10,	9,	12,	9620,	9621
952	DXPXP& PXT& VTP	115.949(53.008,	8.992,	53.949),	13,	11,	27,	9604,	9606
954	QXKX& PXXAXH& PNE	85.984(50.879,	8.408,	6.466),	11,	11,	12,	9621,	9621
957	CXXXX& CXXXX	63.466(52.201,	8.520,	2.746),	14,	14,	14,	9619,	9619
963	EPXXT& MXXH	91.684(36.575,	9.463,	43.446),	13,	13,	32,	9581,	9581
976	ROXXM& MNP	72.881(33.253,	8.700,	28.927),	15,	14,	16,	9616,	9616
979	EXXXP& FDXG& IGXXV	74.898(33.253,	8.520,	33.125),	18,	13,	14,	9617,	9619
984	FPD& GUR& VIK	92.478(55.523,	8.114,	28.839),	10,	8,	8,	9623,	9625
		66.637(44.557,	8.520,	13.840),	10,	10,	14,	9619,	9619

PC	StochasticMotif	$\ell(T)$	$\ell(M)$	$\ell(P)$	$\ell(L)$	E	N_1^+	N_1	N_2^+	N_2
987	GAA& FXLP& MGF	106.943	46.557	8.277	52.109	14	10	10	9619	9623
989	ADN& NAXV& NXGA	89.082	48.557	8.739	41.766	17	18	19	9612	9614
1013	FXXV& KXXRXG	117.105	37.573	9.375	70.154	11	8	46	9584	9587
1049	GPXXXR& GXM& MGP	75.877	50.879	8.569	16.429	16	13	13	9617	9618
1057	GXXXXFXR&LXVK&RXYXXP	161.615	57.271	9.316	94.828	64	56	56	9569	9577
1058	VPXXXXW&YRXXG	49.902	38.383	8.739	2.740	19	19	19	9614	9614
1060	EIXXV& LEXE& IXXVR	119.473	53.201	9.137	87.133	32	29	32	9507	9509
1072	FXXWXP& PXPXXH	81.457	38.575	8.008	33.974	25	23	24	9607	9609
1075	MPXT& RDXT	45.691	33.932	9.137	2.493	33	33	33	9600	9600
1079	PXXXXXXXG&DXXXXXG&DXXXXX	355.502	63.878	9.666	262.737	86	63	69	9541	9564
1087	EQ& QLXP	140.231	31.932	9.939	98.340	14	10	104	9525	9529
1110	NPXXY& WXPY	103.798	35.253	8.346	60.198	12	8	11	9618	9622
1116	WND& WDXXC	119.043	34.253	8.739	76.050	20	15	19	9609	9614
1120	FXXQX& WXXXX	76.313	36.375	8.776	28.961	22	20	20	9611	9613
1148	DXXXXXGXW&TDY	168.462	36.324	9.088	123.070	21	13	30	9593	9603
1155	CXXCR& CXY	50.720	35.253	8.520	6.947	13	13	14	9619	9619
1170	CXXXC& CXC	89.548	35.253	8.908	48.387	21	19	24	9607	9608
1191	RRM& WFQ									
	KEF& WFQ									
	RMK& WFQ									
1211	EXXQ& GDXXXXP	67.689	29.932	8.845	23.312	10	10	22	9611	9611
1214	FNQ& FPN& QLXQ	132.821	37.061	9.589	87.241	17	14	62	9568	9571
	AYXQ& FNQ& GGA									
	AAV& FNQ& PFXQ	57.691	46.557	8.408	2.725	12	12	12	9621	9621
1222	SWXF& WVXXXXS	74.666	37.061	8.659	28.946	19	17	17	9614	9616
1237	KQP& MXXC	43.065	31.932	8.408	2.725	12	12	12	9621	9621
1266	GXXQR& GXXXXQ& DEF	128.421	52.201	8.616	67.605	20	15	16	9612	9617
1341	FM& WFT									
	FM& IVP	81.197	29.932	8.520	2.745	14	14	14	9619	9619
1361	CXFP& FXXXXWP	48.263	37.061	8.466	3.736	13	13	13	9620	9620
1415	GXXDXG& HVD& PGM	60.307	48.879	8.659	2.769	17	17	17	9616	9616
1574	HXXXXXQ&HXXXX	49.441	38.343	8.346	2.713	12	11	11	9622	9622
1676	HXXDK& HFXXC	73.323	35.253	9.068	29.002	10	10	30	9603	9603
1680	MW& PLC	41.063	29.932	8.408	2.725	12	12	12	9621	9621
1681	DAXFP& QXXXXXW	56.476	38.383	8.466	9.628	11	11	13	9620	9620
1682	GEXW& YXRH	45.532	33.932	8.811	2.785	21	21	21	9612	9612
1683	GXXXXXG&YXXC	75.537	38.383	8.577	28.078	21	20	23	9609	9610
1687	QW& WXXXX	44.854	33.253	8.811	2.789	21	21	21	9612	9612
1688	PXXXXXK&WXXXXX	84.571	40.190	9.565	35.410	21	20	26	9606	9607
1691	YXXLXL& YPXXXXL& VXXXXXL	85.032	38.383	8.459	18.063	19	19	17	9616	9616
1705	GXXXXC&LXXEK& RXXXX	94.071	38.383	9.138	32.231	25	23	24	9599	9599
1707	CPXC& LXXC	122.789	35.253	9.016	78.421	20	16	26	9601	9601
1712	KAI& NXXFP& PXXTG	62.335	51.201	8.408	2.735	12	12	12	9621	9621
1717	PPXXXXP& TPW	57.787	34.253	8.659	14.874	15	15	17	9616	9616
1718	CPXXXX& PXXXXL	48.176	37.576	8.811	2.789	21	21	21	9612	9612
1719	FP& PXP& SPN	57.841	48.157	8.589	2.734	18	18	18	9618	9618
1721	CR& WEE	46.224	29.932	8.466	6.827	12	12	13	9620	9620
1779	GXXXXV& CXXXXXP	70.832	39.383	8.408	23.041	11	10	12	9620	9621
1826	SPF& WXXXX	44.728	33.253	8.700	2.774	18	18	18	9615	9615
1910	CXXXXH&GXXF	61.617	37.061	8.201	16.353	10	9	9	9623	9624
1917	HXXXXY& LWR	92.483	34.253	8.346	49.883	10	7	11	9619	9622
2052	CXY& VFL									
	GXXW& KFH	43.463	31.932	8.776	2.785	20	20	20	9613	9613
2054	GW& GYQ	40.909	29.932	8.277	2.698	10	10	10	9623	9623
2056	DXGW& GXXXX	72.991	35.253	8.776	28.961	22	20	20	9611	9613
2059	CXTXY& WGN	44.226	33.253	8.277	2.698	10	10	10	9623	9623
2055	MGQ& VXXX	49.302	33.253	8.700	7.348	17	17	16	9613	9613
2060	GXXGH& CXXXXX	88.040	38.383	8.845	40.812	25	22	22	9608	9611
2074	NXPXXXXK&RXXXXX&RXXXXX&YXDD	183.907	39.856	9.219	94.432	45	37	37	9588	9596
2078	KPC& WGC	41.493	29.932	8.776	2.785	20	20	20	9613	9613
2082	QCF& GXXC	117.542	31.932	9.137	76.514	35	30	33	9595	9600
2086	KXXXXW&YXWL	55.448	37.061	8.520	5.867	12	12	14	9619	9619
2092	UXGH& PXXXXM& VPL	61.307	49.879	8.659	2.769	17	17	17	9616	9616
2097	CFXXX& HXXXX	61.459	38.376	8.466	16.413	14	13	13	9619	9620
2098	RBW& RXXXXXW	55.478	35.061	8.616	15.802	11	11	16	9617	9617
2111	WXXXXP& WXXXXP	47.034	36.375	8.346	2.713	11	11	11	9622	9622
2112	CXXYXXY&PXXXXL	84.849	39.383	8.520	6.947	13	13	14	9619	9619
2122	SXMG& YAM	51.640	33.253	8.520	9.867	12	12	14	9619	9619
2135	NCXXXXC&QXTA	63.087	37.061	8.908	17.119	20	20	24	9609	9609
2140	CXPXM& FMV	104.826	33.253	8.408	63.173	17	17	12	9616	9621
2146	FXXXXCF&SXXXXL& GXXXXXW	183.648	59.137	9.589	114.921	72	62	62	9581	9571
2148	NPXXXX& QXXXXMXC	30.184	38.383	8.992	2.809	27	27	27	9606	9606
2149	CAXC& QMXXXXN	48.036	37.061	8.277	2.698	10	10	10	9623	9623
2151	HXXXXL&LHXXXXW	57.158	39.383	8.408	9.388	10	10	12	9621	9621
2153	CXXXXH& DRI									
	CXXXXH& SPQ	44.892	33.253	8.445	2.793	22	22	22	9611	9611
2155	MXKM& RMXX& TQXB	62.125	50.557	8.845	2.793	22	22	22	9611	9611
2201	WXXM& YIN									
	GXYQ& WKF									
	PYQ& WXXM									
	CXY& PYQ	43.790	31.932	9.044	2.814	29	26	29	9604	9604
2202	FXPXXXXW&NGN	46.036	35.061	8.277	2.698	10	10	10	9623	9623
2204	NWK& WWQ									
	TLC& WWQ									
	NYG& WWQ									
	CDY& RWW	41.672	29.932	8.527	2.804	25	25	25	9608	9608
2212	CKY& KYM	40.940	29.932	8.277	2.698	10	10	10	9623	9623
2253	UTR& NRX& RFP	57.532	46.557	8.277	2.698	10	10	10	9623	9623