

# Hidden Markov Models for Gene Sequence Classification: Classifying the VSG genes in the *Trypanosoma brucei* Genome

**Andrea Mesa**

ANDREAGMESA@CCEE.EDU.UY

*Facultad de Ciencias Económicas y Administración,  
Universidad de la República,  
Montevideo 11200, Uruguay*

**Sebastián Basterrech**

SEBASTIAN.BASTERRECH.TISCORDIO@VSB.CZ

*National Supercomputer Center  
VŠB-Technical University of Ostrava  
Ostrava, 708 00, Czech Republic*

**Gustavo Guerberoff**

GGUERBER@FING.EDU.UY

*Instituto de Matemática y Estadística, Facultad de Ingeniería  
Universidad de la República,  
Montevideo 11200, Uruguay*

**Fernando Alvarez-Valin**

FALVAREZ@FCIEN.EDU.UY

*Sección Biomatemática - Facultad de Ciencias  
Universidad de la República,  
Montevideo 11400, Uruguay*

**Editor:** –

## Abstract

The article presents an application of *Hidden Markov Models (HMMs)* for pattern recognition on genome sequences. We apply HMM for identifying genes encoding the *Variant Surface Glycoprotein (VSG)* in the genomes of *Trypanosoma brucei (T. brucei)* and other African trypanosomes. These are parasitic protozoa causative agents of sleeping sickness and several diseases in domestic and wild animals. These parasites have a peculiar strategy to evade the host's immune system that consists in periodically changing their predominant cellular surface protein (VSG). The motivation for using patterns recognition methods to identify these genes, instead of traditional homology based ones, is that the levels of sequence identity (amino acid and DNA sequence) amongst these genes is often below of what is considered reliable in these methods. Among pattern recognition approaches, HMM are particularly suitable to tackle this problem because they can handle more naturally the determination of gene edges. We evaluate the performance of the model using different number of states in the Markov model, as well as several performance metrics. The model is applied using public genomic data. Our empirical results show that the VSG genes on *T. brucei* can be safely identified (high sensitivity and low rate of false positives) using HMM.

**Keywords:** Hidden Markov Model, Classification, Gene Sequence Classification, *Trypanosoma brucei*, Variant Surface Glycoprotein

## 1. Introduction

Nowadays, the community disposes of a huge amount of genomic sequential data. Machine learning techniques are indispensable in the analysis of sequences due to the enormous size of the available genome databases. In this article we apply *Hidden Markov Models (HMMs)* for identifying a specific pattern in a hyper-variable gene sequences. Markov models have been widely applied for modelling several complex systems that evolve in time. In particular the HMM have proven to be a powerful tool for classification tasks on time series data and spatial sequential data. During the last 20 years, the model has been used for analysing biological sequences, and for solving problems like gene finding and protein analysis. We apply a HMM for identifying genes encoding the *Variant Surface Glycoprotein (VSG)* in the genome of *Trypanosoma brucei (T. brucei)*. This approach can be naturally be extended to other African trypanosomes that also have this type of of hyper-variable genes. These are parasitic protozoa causative agent of *sleeping sickness* in humans and several diseases in domestic and wild animals. These parasites have a peculiar strategy to evade the host's immune system that consists in periodically changing their predominant cellular surface protein (VSG).

There are several methods that have been applied for detecting coding regions in a genome, among which the most used are those based on sequence comparison and *ab initio* methods (Wang et al., 2004). Homology gene finding methods are based in given a DNA sequence identifying its homologous sequences in other genomes stored in databases. In *ab initio* gene prediction the goal is finding either the nearby presence or statistical properties of the coding region i.e. use gene structure as a template to detect genes. Although, for the particular case of the VSG genes these techniques sometimes do not yield reliable results, due to the high variability of the VSG genes. It means that, the levels of sequence identity (amino acid and DNA sequence) amongst the VSG genes is often below of what is considered reliable in the traditional homology based methods. Therefore in this article, we study the possibility of locating these genes looking for a statistical signature using HMM. In the following, we begin by specifying the motivations for studying this particular genoma. Next, we clarify the contributions and goals of this article.

### 1.1 Motivation

The *African Trypanosomiasis* also called *sleeping sickness* is caused by *T. brucei* parasite. There are two subspecies that causes disease in humans: *T. brucei gambiense* and *T. brucei rhodesiense*. A third sub-species named *T. brucei brucei* exists although it rarely infects humans. Around the 90% of the sleeping sickness is caused by *T. brucei gambiense* (Organization, 2006). The disease infections occur in regions of the sub-Saharan Africa. The disease is propagated by the bite of an infected tsetse fly that acts as disease vector. Even though over the last years the number of cases have decreased, in 2009 were reported 9878 cases (Organization, 2006). This decreasing tendency has continued, and it has been reported 7216 cases in 2012.

Cell surface of the trypanosoma consists of a uniform layer of *Variant Surface Glycoprotein (VSG)* that blocks the trypanosomes recognition by the immune system of the mammal host. Once the individual is infected, the trypanosoma expresses a particular VSG and the host immune system reacts by generating a response to this protein layer, which causes a

decrease in the number of trypanosomes. Some trypanosomes can evade the recognition by the immune system. By expressing an alternative VSG for which have not yet developed antibodies. This allows increasing the population of infectious agents sustaining a long term infection, thus increasing the chances of transmission. Therefore, various infections occur each of which is caused by cells expressing a different VSG variant.

## 1.2 Objectives and Contributions

HMM is a probabilistic model used for the analysis of sequential data. There is a general consensus about the power of HMM for classifying time-series data and sequential data. The goal is developing a machine learning tool based on HMM that identifies homogeneous zones with the VSG genes in the genome of *T. brucei*. Among pattern recognition approaches, HMM are particularly suitable to tackle this problem because they can handle more naturally the determination of gene edges. As far as our knowledge is concerned, this article presents the first mathematical approach based on HMM for modelling such sequence.

The main contribution is proposing an algorithm for classifying zones of a genome on two classes, which is based on the following three features: gene location, correlation between symbols in nucleotide sequence and its hidden information included in that sequence. The proposed approach is an hybrid procedure that combines the gene location and the power of the HMM technique for modelling the underlying structure of the sequence. We analyse the performance of our procedure on a public data from *GenBank* database (Benson et al., 2009). The performance of the algorithm is evaluated using different measures from the confusion matrix. Despite its simplicity, according to our experimental results the algorithm is robust and performant in accuracy (high sensitivity and low rate of false positives) and time for detecting VSG genes on *T. brucei*.

## 1.3 Organization

This article is organized as follows. Next section is a general introduction to HMMs and its applications. Section 3 introduces the algorithms used for estimating the parameters of the model. In 3.1 we present the Baum-Welch Algorithm, and in 3.2 we introduce the Viterbi algorithm. We present the methodology used for solving our problem in 4. The experimental results are presented in 5. Finally, we conclude the article and discuss future research lines.

## 2. Hidden Markov Model Description

In this Section, we describe in brief the *Hidden Markov Model (HMM)*. We discuss the general assumptions and properties of a HMM, and we present the approach to use it on the gene identification problem. The end of this section presents some related works that have used HMM for classification on sequential data.

### 2.1 General Introduction

Markov models are one of the most powerful tools for stochastic modelling of dynamic systems. The theory of Markov is rich and quite complete, and the associated algorithmic tools are very efficient. The Markov models are characterised for their high ability to capture all kinds of behaviors. The model provides time-scale invariability in recognition and learning

capabilities (Yamato et al., 1992). A HMM is a particular time-discrete Markov model where the states are not directly visible. The theory of HMM was developed by Baum in the late 60s (Baum and Eagon, 1967; Baum et al., 1970), and in 1989, Rabiner (Rabiner, 1989) publishes a tutorial describing Baum theory applied to speech recognition that is a reference in the study of such models. The family of HMMs has also shown to be effective for analysing biological sequence data (Yoon, 2009; Durbin et al., 1998). The applications include: sequence alignment (Pachter et al., 2002), gene and protein structure predictions (Munch and Krogh, 2006; Won et al., 2007), modelling DNA sequencing errors (Lottaz et al., 2003), and for analysing RNA structure (Yoon and Vaidyanathan, 2008; Harmanci et al., 2007). A HMM is composed by the following elements (Rabiner and H., 1986): a finite set of states, a finite alphabet, and probabilities of state transition and symbol emission. At each time  $t$  a new state is entered in the system following a transition probability distribution that verifies the Markovian property, after each transition an output symbol is produced according to some probability distribution that depends on the current state.

Let  $\mathcal{A}$  be a finite alphabet and  $\mathcal{E}$  be a finite set of states or labels. Given a sequence of observations  $x = (x_1, x_2, \dots, x_n)$ ,  $x_j \in \mathcal{A}$  and a collection of states  $y = (y_1, y_2, \dots, y_n)$ ,  $y_j \in \mathcal{E}$  for all  $j$ , where each  $x_j$  has been emitted by its state  $y_j$ . A HMM describes how the sequence  $x$  and its associated state  $y$  progresses in time. Note that, the index  $j$  ( $j = 1, \dots, n$ ) can be thought of as a time index. We say that the states are *hidden* due to the sequence of observations is known but its related states are unknown. The model works as follows: it starts in an initial state chosen according to some probability, it emits an observation, moves to a new state emits a new observation and so on until to reach some final conditions.

The model parameters are the probability of emitting symbols and the probabilities of the state transitions. Let  $\mathbb{P}(x_i|y_i = k)$  be the emission probability that the symbol  $x_i$  is generated when the sequence is at the state  $k$  at time  $i$ . The probability of transition from the state  $k$  to state  $l$  is given by  $\mathbb{P}(y_{i+1} = l|y_i = k)$ , i.e. the conditional probability that the sequence is in the state  $l$  at time  $i + 1$  when at time  $i$  was in the state  $k$ . For simplicity the above probabilities are respectively written by  $\mathbb{P}(y_{i+1}|y_i)$  and  $\mathbb{P}(x_i|y_i)$ . At  $n$  steps the model generates a sequence  $x = (x_1, x_2, \dots, x_n)$  for which passes through a sequence of states  $y = (y_1, y_2, \dots, y_n)$ . For a fixed length  $n$  the model defines a probability distribution over all sequences of observations  $x$  and all possible states  $y$ . Then, the probability that the model transits the path  $y$  and generates the sequence  $x$  is given by

$$\mathbb{P}(x, y) = \prod_{i=1}^n \mathbb{P}(y_i|y_{i-1}) \mathbb{P}(x_i|y_i), \quad (1)$$

where the first transition  $\mathbb{P}(y_1|y_0)$  is given by  $\mathbb{P}(y_1) = \mathbb{P}(y_1|y_0)$ . Figure 1 illustrates the relationship among states and observations in a HMM. We can see the conditional independence among the events in the system. The probability of the state  $y_i$  only depends on the previous state  $y_{i-1}$ , and the observation generated at time  $i$  only depends on the state at time  $y_i$ .

In this article, we denote the unknown model parameters using greek letters, the matrix of transition probabilities by  $\boldsymbol{q}$ , the emission matrix by  $\boldsymbol{\varepsilon}$ , and the collection of the whole parameter of the model by  $\boldsymbol{\theta}$ . We use bold fonts for matrices and vectors and normal fonts for scalars. Those parameters are estimated in a training step by a particular case of

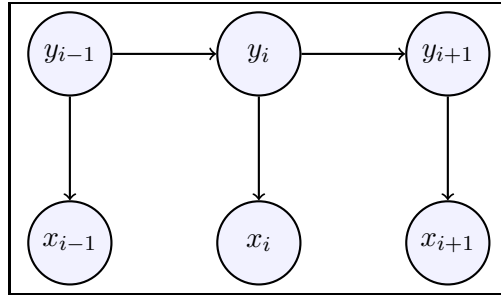


Figure 1: Usual representation of a HMM. Figure illustrates the conditional independence relationship among the states and the observations in the system.

the Expectation Maximization algorithm named the Baum-Welch algorithm (Welch, 2003; Hastie et al., 2001). The model is based in the following assumptions:

- Markov assumption: the transition probability among the states is defined by  $\mathbb{P}(y_{i+1}|y_i)$ . This means that the next state depends only on the current state that is called Markov assumption of first order.
- Homogeneity: it is assumed that the state transition probabilities are homogeneous, i.e. they are independent of its location in the sequence. As a consequence, the probability of moving from one state to the next state is independent on the position of the sequence in which this transition occurs (it is independent on time). Besides, it is assumed that the emission probabilities are also homogeneous.
- Conditional independence among the observations: each observation  $x_i$  only depends on the state  $y_i$ , that is  $x_i$  is generated by the state  $y_i$  with  $\mathbb{P}(x_i|y_i)$ .

## 2.2 Applying HMM in Biological Contexts

A DNA is a macromolecule composed of simple units called nucleotides. Each nucleotide is composed of deoxyribose, phosphate and one of the four nitrogenous bases: *adenine* (*A*), *cytosine* (*C*), *guanine* (*G*), *thymine* (*T*). The DNA is represented by a double helical structure wherein each complementary base pair (*A-T*, *C-G*) are connected by hydrogen bonds. The two DNA strands are in opposite directions to each other, positive or 5' – 3' (left to right) is complemented by the reverse or complementary sequence. The selection of a HMM is given by the alphabet, the number of states and the pattern of transitions of the Markov chain. For this problem, the alphabet is defined as  $\mathcal{A} = \{A, C, G, T\}$ . The number of states depends on the problem, for example can be  $\mathcal{E} = \{\text{gene, intergenic region}\}$  if the goal is determining genes and intergenic regions, and can be  $\mathcal{E} = \{\text{exon, intron, intergenic region}\}$  if also there is interest in identifying non-coding (introns) regions. Although, the hidden space can be easily generalised for more states in the case that we are interested in identifying a higher number of regions in the sequence. The emission and transitions probabilities are numerically computed.

More formal, given an alphabet  $\mathcal{A}$ , a set of states  $\mathcal{E}$  and a DNA sequence  $x$  composed by  $x_1, x_2, \dots, x_n$  where  $x_j \in \mathcal{A}$ , the task is to assign to each symbol  $x_j$  for  $j = 1, \dots, n$  one

label  $y_i$  in  $\mathcal{E}$ . In the specific problem of VSG-gene identification, we define the alphabet as  $\mathcal{A} = \{A, C, G, T\}$ , and we analyse the HMM model for two situations: considering two states ( $\mathcal{E} = \{\text{VSG-gene, no VSG-gene}\}$ ) and considering three states ( $\mathcal{E} = \{\text{VSG-gene, no VSG-gene, other structures}\}$ ). The sequence  $x$  was taken from the public *GenBank* database (Benson et al., 2009).

### 2.3 Related Works

HMMs have been applied in several areas of pattern recognition due to their powerful for modelling sequential data. In particular, they have been especially used for developing speech recognition devices (Rabiner, 1989). Basically, this problem consists in predicting a spoken word from a speech signal and translate into text. For example, hybrid techniques that use ideas from the Neural Network domain and HMMs has been applied for speech recognition systems (Trentin and Gori, 2001). Another tool used for solving this problem has been Gaussian Mixture HMM (Rabiner, 1989). In the case of phone speech recognition a hybrid model that employs deep learning with HMMs was studied in (Dahl et al., 2011). HMM has been also applied for acoustic modeling using Mixture of Gaussian Distributions (Juang et al., 1986; Rabiner and Schafer, 2007). In addition, the model was successfully applied in the Human-Computer Interface area. In (Yamato et al., 1992) time sequential images expressing human actions were encoded in a sequence of symbols, which were modeled by a HMM. In (Liu and Wang, 2011) the authors use the sequence of facial temperature for emotion recognition, the transitions between emotional states are modelled by HMM.

Since the late 80s, HMM has been successfully applied in the area of computational biology (Durbin et al., 1998). In (Churchill, 1989) the author evaluated HMM for modeling five DNA sequences: the yeast mitochondrial, human and mouse mitochondrial, a human X chromosomal fragment, and the bacteriophage lambda genome. Other applications in this first stage have been presented in (Krogh et al., 1994) and (Baldi et al., 1994). In (Krogh et al., 1994) the authors applied HMM for searching common patterns between a protein sequence and multiply aligned sequences. They shown the effectiveness of the model for searching a sequence in a given protein family. In (Baldi et al., 1994) the HMM was used to capture statistical characteristics in three protein families: globins, immunoglobulins, and kinases, wherein the model was applied for classification, multiple alignments, and motif detection. Another protein structure modeling using HMM was presented in (Stultz et al., 1993). A classification problem was studied in (Henderson et al., 1996) where the authors applied HMM for segmenting uncharacterised genomic DNA sequences into exons, introns and intergenic regions. The DNA sequencing errors were modeling in (Lottaz et al., 2003) where the authors improved the detection of translation start and stop sites, and thus the location of the coding regions.

Other applications correspond to the task of searching structures on genomes. The model was used for gene searching in bacterial genomes (Lukashin and Borodovsky, 1998), and in (Delcher et al., 1999) was shown the effectiveness of HMM for finding gene on eukaryotic genomes. The model was also used for identifying homologous protein and nucleotide sequences (Finn et al., 2011). A variation of HMM named *Profile-HMM* have been applied for analysing genomic sequences (Eddy et al., 1995; Eddy, 1996, 1998; Gough et al., 2001).

The huge interest in the community concerning HMMs for solving biological computational problems has caused that several researchers develop specific software for gene sequence and prediction signals, these include (Lukashin and Borodovsky, 1998; Delcher et al., 1999; Finn et al., 2011; Flicke, 2007).

Other learning techniques have also been used for solving classification tasks and clustering problems in the field of computational biology. In (Decaprio et al., 2007) the authors described a gene predictor based on *Conditional Random Fields (CRF)* that incorporates information of *Expressed Sequence Tags (ESTs)* for predicting genes on fungal genome. In (Gross et al., 2007) the authors combined CRF with *Support Vector Machines (SVM)* and tested for the human genome, the model incorporates data from other genomes called informants, making a multiple alignment based on some assumptions about the respective evolutionary processes. Another application of the HMM combined with SVM analyses the genome of a nematode (Schweikert et al., 2009).

In biological sequence analysis *Artificial Neural Network (ANN)* have been also used. In (Rebello et al., 2011) ANN was applied for predicting genes on a *Streptococcus* genome, and they were also used on the prediction of splice sites of a gene in (Johansen et al., 2009). *Decision trees* were performed for finding genes in vertebrate DNA sequences (Salzberg et al., 1996), and a variation (called *Randomized oblique decision trees*) combined with other learning techniques were applied for gene modeling in (Allen et al., 2004).

In the particular case of *T. brucei*, an evolutionary analysis considering the family members of this gene was introduced in (Alvarez et al., 1996). Although, in this article the authors analyse the sequence doing pairwise comparisons between *T. brucei* and other protein coding genes. An example of using a system based on HMM to search genes on a particular *T. brucei* chromosoma is presented in (El-sayed et al., 2003), although in this article the sequence has only one VSG cluster.

During the last decades, HMM have been used in a large number of applications for analysing sequential data. Here, we presented a summary that is not exhaustive, a more complete overview about HMM applications in biological sequences can be seen in (Durbin et al., 1998; Choo et al., 2004; Yoon, 2009).

### 3. Training and Inference

Let  $\theta$  be the vector with the parameters of the model (the probabilities of transitions and the probabilities of emissions). The training process is based in the *Maximum-likelihood Estimation (MLE)* that consists in adjusting parameters to maximize the likelihood for a given sequence  $x$  generated by the model. That is finding  $\theta$  such that  $\mathbb{P}_\theta(x)$  is maximised. For any two states  $k$  and  $l$ , and a symbol  $b$  we find the parameters  $p_{kl} = \mathbb{P}(y_{i+1} = l | y_i = k)$  and  $e_k(b) = \mathbb{P}(x_i = b | y_i = k)$  as follows:

$$\hat{p}_{kl} = \frac{P_{kl}}{\sum_{l'} P_{kl'}}, \quad \hat{e}_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')},$$

where  $P_{kl}$  is the number of transitions from  $k$  to  $l$ , and  $E_k(b)$  is the number of times that the state  $k$  emits the symbol  $b$  in the current training sequence. The adjustable vector parameter  $\theta$  is a collection composed by  $\hat{p}_{kl}$  and  $\hat{e}_k$  for all  $k$  and  $l$ .

A computational efficient algorithm for setting the parameter is *Baum-Welch (BM)* Algorithm (Welch, 2003), that is a particular case of the well-known *Expectation-Maximization (EM)* (Baum et al., 1970; Dempster et al., 1977; Hastie et al., 2001). The BW algorithm is based in two auxiliary recursive procedures called *forward* and *backward* algorithms. The forward algorithm is a dynamic programming algorithm that computes in an efficient way the *forward probability*, which is given by the following expression (Yoon, 2009)

$$\begin{aligned}\alpha_k(n) &= \mathbb{P}(x_1, \dots, x_n, y_n = k) \\ &= \mathbb{P}(x_n | y_n = k) \sum_j \alpha_j(n-1) \mathbb{P}(y_n = k | y_{n-1} = j).\end{aligned}$$

The time-complexity of the forward algorithm is  $nK^2$  where  $K$  is the number of states and  $n$  is the size of the sequence  $x$ . Algorithm 1 presents the procedure used for compute the forward probabilities  $\alpha_k(n)$  for all states  $k$ .

---

**Algorithm 1:** Forward Algorithm.

---

```
// Initialization
1  $\alpha_k(1) = \mathbb{P}(x_1, y_1 = k) = \mathbb{P}(y_1 = k) \mathbb{P}(x_1 | y_1 = k), \forall k = 1, \dots, K;$ 
2 for ( $i = 1$  to  $n$ ) do
3    $\left[ \alpha_k(i) = \mathbb{P}(x_i | y_i = k) \sum_j \alpha_j(i-1) \mathbb{P}(y_i = k | y_{i-1} = j); \right.$ 
4  $\mathbb{P}(x) = \sum_k \alpha_k(n);$ 
5 Return  $\mathbb{P}(x)$  and  $\alpha_k(n)$  for all  $k;$ 
```

---

The backward probability is also recursively defined

$$\begin{aligned}\beta_k(n) &= \mathbb{P}(x_{n+1}, \dots, x_n | y_n = k) \\ &= \sum_l \mathbb{P}(x_{n+1} | y_{n+1} = l) \beta_l(n+1) \mathbb{P}(y_{n+1} = l | y_n = k).\end{aligned}$$

Algorithm 2 illustrates the procedure used for compute the backward probabilities  $\beta_k(n)$  for all states  $k$ . In a similar way that the forward case, this backward probability can be recursively computed in a  $nK^2$  time-complexity (Yoon, 2009).

### 3.1 Description of the Baum-Welch Algorithm

The BW algorithm is iterative, at the first iteration is initialised the transition probability matrix  $\boldsymbol{\rho}$  and the emission matrix  $\boldsymbol{\varepsilon}$ . The initialisation can be done using a uniform distribution or in an arbitrary way based on *a priori* known information. In the following we present in brief the main mathematical expressions of the Baum-Welch method. Let  $x$  be the current training sequence, the transition probability from the state  $k$  to the state  $l$  at the position  $i$  of  $x$  is computed as (Yoon, 2009)

$$\mathbb{P}(y_i = k, y_{i+1} = l | x) = \frac{\alpha_k(i) \mathbb{P}(y_{i+1} = l | y_i = k) \mathbb{P}(x_{i+1} | y_{i+1} = l) \beta_l(i+1)}{\mathbb{P}(x)},$$



---

**Algorithm 2:** Backward Algorithm.

---

```
// Initialization
1  $\beta_k(n) = 1 \forall k = 1, \dots, K;$ 
2 for ( $i = n - 1$  to 1) do
3    $\left[ \beta_k(i) = \sum_l \mathbb{P}(x_{i+1}|y_{i+1} = l)\beta_l(i + 1)\mathbb{P}(y_{i+1} = l|y_i = k); \right.$ 
4  $\mathbb{P}(x) = \sum_k \mathbb{P}(y_1|y_0 = k)\mathbb{P}(x_1|y_1 = k) \beta_k(1);$ 
5 Return  $\mathbb{P}(x)$  and  $\beta_k(n)$  for all  $k;$ 
```

---

where  $\mathbb{P}(x) = \sum_k \alpha_k(n)$  is computed using the forward algorithm (Rabiner, 1989). The probability of generating the symbol  $b$  when the state is  $k$  is computed as

$$\frac{\sum_{\{i:x_i=b\}} \alpha_k(i)\beta_k(i)}{\mathbb{P}(x)}. \quad (2)$$

Let  $T$  be the number of sequences in the training set, and we denote the training sequences by  $x^{(t)}$  for  $t = 1, \dots, T$ . The probabilities of transition and emission are computed as

$$\hat{p}_{kl} = \frac{P_{kl}}{\sum_{l'} P_{kl'}} \quad \text{and} \quad \hat{e}_k(b) = \frac{E_k(b)}{\sum_{b'} E_k(b')},$$

where

$$P_{kl} = \sum_t \frac{1}{\mathbb{P}(x^{(t)})} \sum_i \alpha_k^t(i) \mathbb{P}(y_{i+1} = l | y_i = k) \mathbb{P}(x_{i+1}^{(t)} | y_{i+1} = l) \beta_l^t(i + 1), \quad (3)$$

$$E_k(b) = \sum_t \frac{1}{\mathbb{P}(x^{(t)})} \sum_{\{i:x_i^{(s)}=b\}} \alpha_k^t(i) \beta_k^t(i). \quad (4)$$

It has been proven in (Hastie et al., 2001) that the likelihood function increases at each step of the BW algorithm. That is  $P_{\hat{\theta}^{(h+1)}}(x) \geq P_{\hat{\theta}^{(h)}}(x)$  for all step  $h$ . The algorithm ends when the difference between the likelihood function at the step  $h + 1$  and the value at the step  $h$  is lower than some arbitrary threshold, that is  $dist(P_{\hat{\theta}^{(h+1)}}(x), P_{\hat{\theta}^{(h)}}(x)) < \epsilon$  where  $dist(\cdot)$  is an arbitrary distance. In our numerical experiences, we used Euclidean distance and  $\epsilon = 0.00001$ , in addition we control the number of iterations (the maximum number of iterations was 500).

### 3.2 Description of the Viterbi Algorithm

Once the model parameters are estimated, the Viterbi algorithm is useful for finding the *best* sequence of states  $y$  for generating a specific sequence  $x$ . The method assumes that the *best* sequence is that that generates the sequence  $x$  with higher probability, that is

$$y^* = \arg \max_y \frac{\mathbb{P}(x, y)}{\mathbb{P}(x)}. \quad (5)$$

Note that the denominator in (5) does not depend of  $y$ . So, the goal is to find  $y^*$  such that  $\mathbb{P}(x, y)$  is maximised. Viterbi algorithm is iterative. Given the samples until the position  $i$ , let  $V_i(k)$  be the probability of the best sequence of states (*path*) that generates the state  $k$ . At the initial iteration, the method sets  $V_1(k)$  with the probability that observation  $x_1$  is generated. Next,  $V_i(l)$  is computed for each position  $i = 2, \dots, n$  and each state  $k = 1, \dots, K$ . At the step  $i$ , the algorithm finds the sequence of states  $y_1, y_2, \dots, y_i$  that generates the sequence  $x_1, x_2, \dots, x_i$  with highest probability, for each location  $i$  and state  $k$  ( $y_i = k$ ). The algorithm computes the  $V_{i+1}(l)$  for the location  $i + 1$  and all possible states  $l$  following

$$V_{i+1}(l) = \max_k \left( V_i(k) \mathbb{P}(y_{i+1} = l | y_i = k) \right) \mathbb{P}(x_{i+1} | y_{i+1} = l). \quad (6)$$

We define an auxiliary variable  $I_i(l)$  that contains the state that maximises the expression (6) at each iteration. The probability  $\mathbb{P}(x, y^*)$  is given by the maximum value of  $V_n(k)$  for all the states  $k$ . The most probable path  $y^*$  is computed using a backtracking procedure from ( $n$  till 1) and  $y_{i-1}^* = I_i(y_i^*)$ .

### 3.3 Technical Issues

We present here some technical issues related to the performance and numeric stability of the algorithms during the training process

- Algorithmic stability: The product among probabilities should be controlled, because it can occur underflow. This is a common problem when we are working with power of Markov matrices and products of probabilities. The stability can be improved using scaling coefficients in the forward and backward probabilities and logarithmic transformations in the Viterbi algorithm (Durbin et al., 1998; Rabiner, 1989).
- Performance of learning on imbalanced dataset: we say that a dataset is imbalanced if the classes are not *approximately equally represented* (Chawla et al., 2002). This imbalance can provoke suboptimal classification performance, due to the learning tool overweight the large classes ignoring the small ones (Chawla et al., 2002). Some techniques have been introduced for improving the predictors in the imbalanced context, which can be performed for solving the gene classification problem (Chawla et al., 2004).

## 4. Methodology

In this section we describe the methodology used for applying HMM for solving the problem of the labelling areas of *T. brucei* genome as VSG gene or non-VSG gene. The alphabet is composed by the four nucleotide types  $\mathcal{A} = \{A, C, G, T\}$ . Concerning the number of states in the problem, we begin by adjusting the model with two states, which represents a binary situation: yes or not *VSG* region. Let  $S$  be the non-redundant collection of sequences of *T. brucei* genome. The set  $S$  has 9 sequences of different lengths and different numbers of clusters of VSG genes. As usual, we generate a random partition of the set  $S$  in two subsets: a training set  $S^T$  and a validation set  $S^V$ . The sequences in  $S^T$  are used for estimating the transition probability matrices and the parameters of the emission distribution. We apply

BW Algorithm for performing that estimation. Note that, only the nucleotide sequence is taken account in the BW Algorithm, the information about the VSG genes is not used. Next, we apply the Viterbi Algorithm on the sequences in  $S^V$  in order to estimate the most likely sequence of hidden states. We determine those regions in the testing sequences that are predicted as VSG genes. Finally we evaluate the model using the predicted values with the VSG annotated genes of the sequences in  $S^V$ .

In our experiments we find that the predictor can be improved defining an interval of control, which we call *location control filter*. We can define a range where the VSG are presented, which is based on the minimum and maximum lengths of the labeled VSG. Once this range is defined, we can use it as filter for the predictions given by the BW and Viterbi algorithms. The model assigns a label VSG to a sub-sequence of nucleotide if verifies the following two conditions: it is predicted by the HMM model such as VSG and it passes the filtering operation. We call *filter operation* if the subsequence has location inside the range defined by the training set. As a consequence, it is possible to adjust the predictions given by the model discarding those that do not pass the filter. Obviously, the range depends of the training set, then it can change according different training data. In spite of that, this is a common practice in machine learning. For instance, when the input patterns are normalised. This normalisation more often takes as reference the maximum and minimum value of the input variables. Here, the approach is similar, we consider the minimum VSG location and the maximum VSG location using the training data, out of this range the genes are considered non-VSG. Figure 2 presents a hierarchical schema among the techniques that compose our approach.

We evaluate the learning capability of the model using the following quantitative metrics. Table 1 shows the error types (confusion matrix) that can happen during the parameter estimations, four situations are identified:

- *True positive (TP)* refers when the VSG was correctly classified.
- *False positive (FP)* refers to an incorrectly identification, in our problem that occurs when the HMM predicts that some nucleotides are part of the VSG gene when this is not correct.
- *False negative (FN)* refers when the model incorrectly rejected a VSG gene.
- *True negative (TN)* refers when the model correctly rejected a VSG gene.

We measure the quality of our estimation using two well-known measures of binary prediction quality: the *True Positive Rate (TPR)* (also named *sensitivity*) and the *Positive Predictive Value (PPV)* (also named *precision*). The TPR is the proportion of nucleotides correctly predicted over the total number of nucleotides that are part of VSG gene. It is defined as:

$$TPR = \frac{TP}{TP + FN}. \quad (7)$$

The PPV measures the proportion of correctly classify nucleotides over the total of VSG identifies. It is given by:

$$PPV = \frac{TP}{TP + FP}. \quad (8)$$

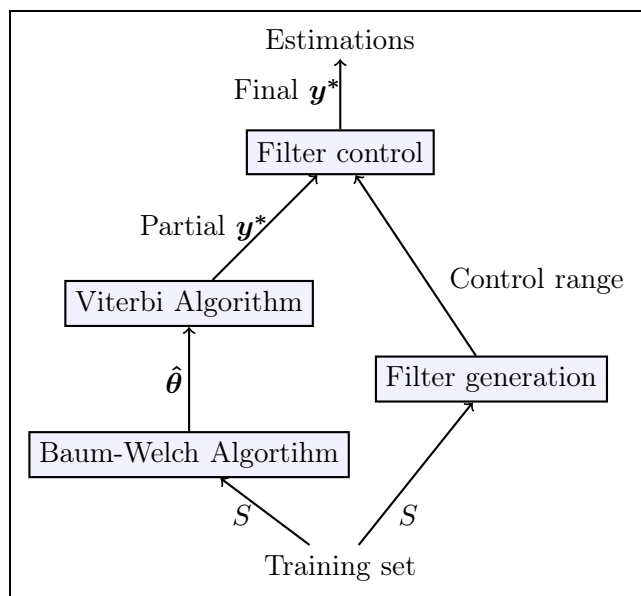


Figure 2: Schema of the procedure for gene sequence classification using HMM. A training set composed by sequences is the input of the Baum-Welch algorithm, then the Viterbi algorithm is applied producing a *partial* estimation. The training sequence is also used for generating an interval control corresponding the locations of the VSG genes. Next, the location interval of control is applied as filter on the *partial* estimation to produce the final estimation.

Table 1: Confusion matrix: type of estimation errores.

Prediction / Target	Positive	Negative
Positive	True Positive (TP)	False Positive (FP)
Negative	False Negative (FN)	True Negative (TN)

The following hypothesis are considered. We assume that exists independence among sequences, as well as among genes into the sequences. That is, we dispose of a non-redundant collection of sequences of *T. brucei* genomic data. A common problem when we are working on pattern recognition for genome sequences is that we can not affirm the independence between genes and segment of genes in the sequence. Some parts in the sequence can be strongly correlated each of other. The correlation among several sequences arises from the process of phylogenetic inertia, basically this means that the sequences have the same origin. In spite of that, the process of divergence evolution can vanish the phylogenetic information. Even two sequences that diverged from a common ancestral sequence can absolutely lost the information from their ancestors if the sequences diverged enough. Often, it is used the percentage of identical base pair (sequence identity) between two sequences as correlation assessment between that sequences. As a consequence, in order to assume independence among sequences and to produce a non-redundant collection, the sequence percent identity is used as independence assessment reference. The VSG genes diverge from a common

ancestor, despite that their amino acid identity is low. An usual assumption is considering the VSG genes as non-redundant set, and assuming the hypothesis of independence among the sequences.

## 5. Results

### 5.1 Data description

We use a set of 9 genomic sequence segments from chromosome 9 of *T. brucei* genome where some VSG gene clusters are located. The data set was taken from the *GenBank* database (Benson et al., 2009). Table 2 shows size of sequences and number of annotated VSG genes for each one. Since in these genomic sequences the location of most VSG had been previously identified using traditional homology-based annotation methods, our predictions obtained with the HMM approach can be compared with these annotations obtained by traditional based methods.

Table 2: Size of the sequence segments and number of annotated VSG genes in each sequence of chromosome 9 of *T. brucei*.

Sequence Id	Size	Number of VSG genes
$S_1$	75462	16
$S_2$	67530	21
$S_3$	19745	7
$S_4$	89317	10
$S_5$	102962	3
$S_6$	17110	5
$S_7$	88923	8
$S_8$	24879	0
$S_9$	65933	7

### 5.2 Applying the HMM model with 2 states

We start by analysing the performance of the HMM model with 2 states. Without loss of generality we index the sequences in  $S$  from 1 to 9 as  $S_1, S_2, \dots, S_9$ . We create three partitions of  $S$  using a uniform random selection. Each partition have two groups, one is used for training and another one for validations. Table 3 shows the training and validation sets for each partition of  $S$ . In case of HMM with 2 states the BW algorithm converges after 160 iterations for the three training sets. We compute the Viterbi path and the model accuracy.

We reached the following estimated matrices when we used the sequences of the partition  $A$

$$\hat{\boldsymbol{\theta}} = \begin{pmatrix} 0.9969 & 0.0031 \\ 0.0029 & 0.9971 \end{pmatrix},$$

where the  $(i, j)$  element represents the transition probability from state  $i$  to  $j$  for  $i, j = 1, 2$  and

Table 3: Generation of the training and validation sets.

Partition Id	Training set	Validation Set
<i>A</i>	$\{S_1, S_2, S_4, S_5, S_7, S_8\}$	$\{S_3, S_6, S_9\}$
<i>B</i>	$\{S_1, S_2, S_3, S_4, S_6, S_9\}$	$\{S_5, S_7, S_8\}$
<i>C</i>	$\{S_3, S_5, S_6, S_7, S_8, S_9\}$	$\{S_1, S_2, S_4\}$

$$\hat{\epsilon} = \begin{pmatrix} 0.3815 & 0.1911 & 0.2018 & 0.2255 \\ 0.1932 & 0.2213 & 0.2031 & 0.3824 \end{pmatrix}$$

where  $(i, 1)$ ,  $i = 1, 2$  represents the probability of state  $i$  emits the symbol  $A$ ,  $(i, 2)$  the probability that state  $i$  emits the second symbol  $C$  and so on. The matrix of transition probabilities is diagonal dominant, therefore the chance that the model jumps from one states to another one is very low.

Table 4 shows some metrics of the accuracy of our model. The rows of the table show the accuracy metric reached for the model with 2 states. The first column specify the metric, the next three columns shows the TPR, PPV, and number of VSG non-detected (ND-VSG) using the partition  $A$ . The other following columns present the results using  $B$  and  $C$  sets. The model predicts all VSG genes in almost all cases, therefore TPR is high. In the worst case the TPR is higher than 0.81. The PPV for the sequences  $S_9$  of the partition  $A$  and  $S_5$  of the partition  $B$  are small, as a consequence the model prediction has a relatively high false positive error.

Table 4: Sensitivity (TPR), precision (PPV) and number of non-detected VSG (ND-VSG) reached by a HMM with 2 states

		<i>A</i>			<i>B</i>			<i>C</i>	
	$S_3$	$S_6$	$S_9$	$S_5$	$S_7$	$S_8$	$S_1$	$S_2$	$S_4$
TPR	0.9709	0.9634	0.9838	0.9581	0.9485	-	0.9759	0.8104	0.9959
PPV	0.5214	0.4652	0.2781	0.1769	0.4775	-	0.3171	0.4455	0.3085
ND-VSG	0	0	0	0	0	0	8	4	0

### 5.3 Application of the HMM with 3 states

In order to improve the accuracy in predictions, we analyse a HMM with 3 states. Our purpose is to investigate the existence of another group that does not appear in the HMM with 2 states. We use training and validation set presented in Table 3. The BW algorithm converges for the three partitions in less than 400 iterations.

The estimated matrix of transitions when was used the sequences of the partition  $A$  was,

$$\hat{\rho} = \begin{pmatrix} 0.9956 & 0.0028 & 0.0016 \\ 0.0021 & 0.9976 & 0.0003 \\ 0.0017 & 0.0003 & 0.9980 \end{pmatrix},$$

where the  $(i, j)$  element represents the transition probability from state  $i$  to  $j$  for  $i, j = 1, 2, 3$ . As well as in the case of a model with two states, that matrix is diagonal dominant. The estimated matrix of emissions when was used the sequences of the partition  $A$  was,

$$\hat{\epsilon} = \begin{pmatrix} 0.3610 & 0.1455 & 0.1388 & 0.3547 \\ 0.1712 & 0.2346 & 0.2267 & 0.3675 \\ 0.3610 & 0.2330 & 0.2369 & 0.1691 \end{pmatrix},$$

where the first column of the matrix corresponds to the base  $A$ , the second to the base  $C$ , the third one to the base  $G$  and the last one to the base  $T$ , so the  $(i, 1)$  represents the probability of state  $i$  emits the symbol  $A$  and so on,  $i = 1, 2, 3$ .

Table 5 presents the accuracy of the model with 3 states. The first three columns shows TPR, PPV and number of non-detected VSG genes for partition  $A$ . The next group of three columns present accuracy for partition  $B$  and the last columns corresponds to partition  $C$ . We can affirm that the model with 3 states fits better than the 2-state model. The TPR metric is slightly lower in this case, but the precision increases for all sequences.

Table 5: Sensitivity (TPR), precision (PPV) and number of non-detected VSG (ND-VSG) reached by a HMM with 3 states.

	$A$			$B$			$C$		
	$S_3$	$S_6$	$S_9$	$S_5$	$S_7$	$S_8$	$S_1$	$S_2$	$S_4$
TPR	0.9472	0.9540	0.9559	0.9535	0.9405	-	0.9559	0.9280	0.9403
PPV	0.6712	0.7080	0.4022	0.3985	0.6629	-	0.5027	0.5435	0.4557
ND-VSG	0	0	0	0	0	0	0	0	0

Table 6 shows the locations of annotated VSG (left column) and predicted VSG (right column) by the HMM model with 3 states using the partition  $A$  and sequence 9. The model correctly predicts seven VSG genes and misclassifies other seven ones. The length of VSG genes correctly predicted in the sequence 9 varies from a minimum of 1115 up to 2217. The VSG genes wrong predicted by the model have lengths: 113, 510, 653, 250, 206, 5600 and 2432. Note that all the bad predictions are out of the range [1115, 2217]. This situation is repeated in a large number of wrong predicted genes on the validation sequences. As a consequence, if we consider the lengths of the annotated VSG genes (in the training set) we can define a bounded range by the minimum and maximum length of well-predicted VSG genes. Then, it is possible to adjust the predictions provided by the model discarding those that are outside of this range. Therefore, we can use the training set for both: to set the parameters of the model and to define a range for the VSG length. Figure 3 illustrates how the model is improved when a location filter is performed. Besides, it shows an additional challenge of the learning model due to the imbalance among the classes (Chawla et al., 2004). There are three histograms, the horizontal axe indicates the classes and the vertical axe indicates the percentage of genes. The data is the testing and corresponds to the sequence 9. There are two classes: the non-annotated VSG (denoted by 0) and the annotated VSG (denoted by 1). The left histogram shows the testing data, in the center is presented the estimations produced by the HMM without control of the gene locations, and the right histogram presents the estimations produced by the HMM with location control. The serie of histograms shows how is improved the percentage of well-classified genes for the sequence 9 using the location filter, the increment is larger than 10% of cases.

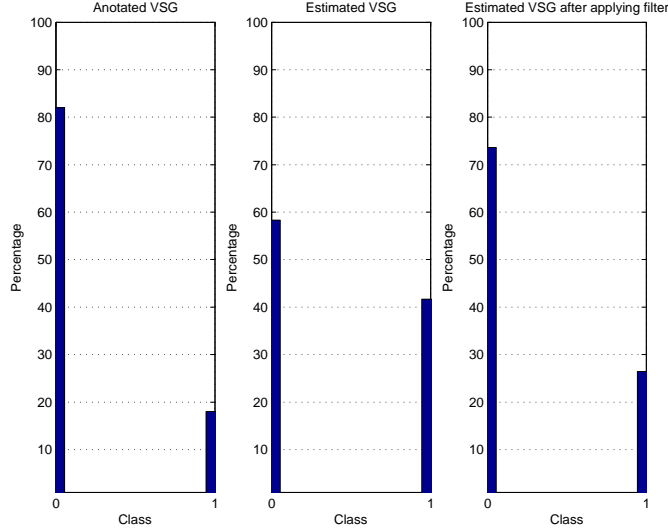


Figure 3: Example of the improvement for controlling the gene location. The data corresponds to the sequence 9. Left histogram was generated with the original data, in the center the histogram was done with the estimations produced by the HMM, and the right histogram corresponds to the estimations with a location control.

The VSG genes in the sequence 9 are located between the first 22000 nucleotides (Table 6). In this region, the most likely path of states determined by the Viterbi algorithm shows that only states 1 and 3 are presented (the state 3 corresponds to the VSG genes), with only the exception between the coordinates 6277 – 6429. From 27560 until the end of the sequence, the changes are between states 1 and 2 (with only the exception between 50143 and 52575 locations). This indicates that in addition to the VSG genes, the model detects other regions (corresponding to state 2) with particular characteristics that make them distinguishable from those predicted as states 1 and 3. We observe that those particular regions located at the second middle of the sequence and relating to the state 2 correspond to the VSG genes in the complementary sequence (Table 7). Table 7 presents an example about the location of annotated (left column) and predicted (right column) VSG genes by the 3-states HMM for the complementary sequence 9. Figure 4 was generated with the Artemis Software (Rutherford et al., 2000; Carver et al., 2012). It shows the location of annotated VSG genes and predicted by the 3-state model for sequence 9. The annotated VSG genes are indicated with the label “vsg” and the predictions are labeled as “pred”. The figure is useful for explain the model accuracy in finding the VSG on the positive strand (left to right direction, located at the top of the figure) and complementary strand (opposite direction, situated on the bottom) for the sequence 9. Table 8 shows the performance reached by the HMM with 3 states for the complementary sequence. The table presents the TPR and PPV values and the number of not detected VSG. The three groups of columns correspond to the *A*, *B* and *C* learning partitions. We can see that all VSG genes located on the complementary sequence are well-identified by the model.



Table 6: Location of the VSG genes annotates and predicted by the model with 3 states. The data corresponds to the sequence 9.

Location of the VSG gene	Location of the VSG prediction
1001 - 2502	658 - 2427
	3056 - 3169
4223 - 5692	3423 - 5640
6846 - 7375	6429 - 7544
8875 - 10387	8626 - 10315
	10829 - 11339
12360 - 13819	11555 - 13748
	14192 - 14845
16342 - 17943	16184 - 17849
	18148 - 18398
	18712 - 18918
20343 - 21634	19746 - 21585
	21676 - 27276
	50143 - 52575

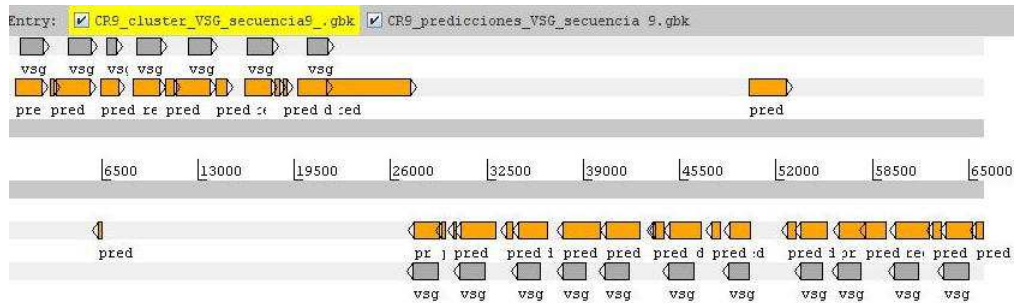


Figure 4: Annotated VSG genes and their predictions by an HMM with 3 states for the sequence 9. The Figure was generated using Artemis software (Carver et al., 2012).

## 6. Conclusions and Future Work

In this article we tackle a classification problem using a *Hidden Markov Model (HMM)*. Our approach identified with relatively *success* a particular type of genes on *Trypanosoma brucei* (*T. brucei*) genome, those encoding *Variant Surface Glycoproteins (VSG)* from African trypanosomes. *T. brucei* uses the expression of different VSG proteins to evade host immune mechanisms provoking the Trypanosomiasis disease. The results shows that the VSG genes have characteristics that make them detectable as statistically homogeneous zones in *T. brucei* genome. We perform a HMM with 2 and 3 states. The model with 3 states outperforms the 2-states model, identifying VSG in the two opposite strand directions. The HMM with 3 states is efficient in searching VSG genes, according our experiments in all cases the model reached a sensitivity over the 90% and the totality of VSG genes were well-identified.

Table 7: Location of the VSG gene and their predictions for the HMM with 3 states for the complementary sequence. The data corresponds to the sequence 9.

Location of the VSG gene	Location of the VSG prediction
27500 - 29086	27560 - 29217
	29404 - 29594
	30216 - 30348
30672 - 32290	30752 - 32972
	33791 - 34145
34509 - 35973	34619 - 36483
37592 - 39046	37648 - 40083
40462 - 41986	40528 - 42688
	43715 - 43783
	43853 - 44335
44756 - 46342	44816 - 46881
	47692 - 48183
48819 - 50113	48910 - 50143
	52782 - 53268
53634 - 55059	53685 - 55339
56172 - 57665	56223 - 57969
	58034 - 59358
59975 - 61544	60042 - 62274
	62542 - 63076
63341 - 64935	63466 - 65166
	65468 - 65933

Table 8: Sensitivity (TPR), precision (PPV) and number of non-detected VSG (ND-VSG) reached by a HMM with 3 states for the complementary sequences.

	<i>A</i>			<i>B</i>			<i>C</i>		
	$S_3$	$S_6$	$S_9$	$S_5$	$S_7$	$S_8$	$S_1$	$S_2$	$S_4$
TPR	-	-	0.9508	0.9485	0.9568	0.9540	0.9226	-	0.9385
PPV	-	-	0.6157	0.5568	0.4833	0.5422	0.5233	-	0.6005
ND-VSG	0	0	0	0	0	0	0	0	0

In addition, we propose a simple method for identifying the misclassified genes by the HMM model. We define a filter using the location of the VSG in the training data. Despite the simplicity of the approach presented in this article, according the empirical results for a particular genome the performance for gene detection was very high. This approach can be of great utility to apply to other African trypanosomes whose silent repertoires of VSG genes are not well-characterized experimentally as well as to others families of hyper-variable genes.

The present work was restricted to detect VSG genes within regions that are already known to be VSG gene clusters. Therefore the problem was delimited to differentiate only

three states, VSG genes, VSG genes located on the complementary strand of DNA and intergenic regions. An interesting issue for future work is the following: it would be of great interest to distinguish also between genomic regions containing VSG clusters from other genomic regions containing regular protein coding genes (encoding for normal functions in the cell).

## References

- Jonathan E. Allen, Mihaela Pertea, and Steven L. Salzberg. Computational Gene Prediction Using Multiple Sources of Evidence. *Genome Research*, 14(1):142–148, January 2004.
- Fernando Alvarez, Maria Noel Cortinas, and Hector Musto. The Analysis of Protein Coding Genes Suggests Monophyly of Trypanosoma. *Molecular Phylogenetics and Evolution*, 5(2):333 – 343, 1996. ISSN 1055-7903. doi: <http://dx.doi.org/10.1006/mpev.1996.0028>.
- Pierre Baldi, Yves Chauvin, Tim Hunkapiller, and Marcella A McClure. Hidden Markov models of biological primary sequence information. *Proceedings of the National Academy of Sciences*, 91(3):1059–1063, 1994.
- L. E. Baum, T. Petrie, G. Soules, and N. Weiss. A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *The Annals of Mathematical Statistics*, 41(1):164–171, 1970.
- L.E. Baum and J.A. Eagon. An inequality with applications to statistical estimation for probabilistic functions of Markov processes and to a model for ecology. *Bulletin of the American Mathematical Society*, 73(3):360–363, 1967.
- D.A. Benson, I. Karsch-Mizrachi, D.J. Lipman, J. Ostell, and E.W. Sayers. 37(database issue):d26-31. Technical report, GenBank, January 2009. URL <http://www.ncbi.nlm.nih.gov/genbank>.
- Tim Carver, Simon R Harris, Matthew Berriman, Julian Parkhill, and Jacqueline A McQuillan. Artemis: an integrated platform for visualization and analysis of high-throughput sequence-based experimental data. *Bioinformatics (Oxford, England)*, 28(4), February 2012. ISSN 1367-4803. doi: 10.1093/bioinformatics/btr703. URL <http://europepmc.org/articles/PMC3278759>.
- Nitesh V. Chawla, Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16:321–357, 2002.
- Nitesh V. Chawla, Nathalie Japkowicz, and Aleksander Kotcz. Editorial: Special Issue on Learning from Imbalanced Data Sets. *SIGKDD Explor. Newsl.*, 6(1):1–6, June 2004. ISSN 1931-0145. doi: 10.1145/1007730.1007733. URL <http://doi.acm.org/10.1145/1007730.1007733>.
- Khar Heng Choo, Joo Chuan Tong, and Louxin Zhang. Recent applications of hidden Markov models in computational biology. *Genomics Proteomics Bioinformatics*, 2(2): 84–96, 2004.

- Gary A Churchill. Stochastic models for heterogeneous DNA sequences. *Bulletin of mathematical biology*, 51(1):79–94, 1989.
- G.E. Dahl, Dong Yu, Li Deng, and A. Acero. Large vocabulary continuous speech recognition with context-dependent DBN-HMMS. In *Acoustics, Speech and Signal Processing (ICASSP), 2011 IEEE International Conference on*, pages 4688–4691, May 2011. doi: 10.1109/ICASSP.2011.5947401.
- D. Decaprio, J.P. Vinson, M.D. Pearson, P. Montgomery, M. Doherty, and J.E. Galagan. Conrad: Gene prediction using conditional random fields. *Neural Networks*, 17(9):1389–1398, 2007.
- Arthur L Delcher, Douglas Harmon, Simon Kasif, Owen White, and Steven L Salzberg. Improved microbial gene identification with GLIMMER. *Nucleic acids research*, 27(23):4636–4641, 1999.
- A.P. Dempster, N.M. Laird, and D.B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- R. Durbin, S. Eddy, A. Krogh, and G. Mitchinson. *Biological sequence analysis. Probabilistic models of proteins and nucleic acids*. Cambridge University Press, 1998.
- Sean R Eddy. Hidden Markov models. *Current opinion in structural biology*, 6(3):361–365, 1996.
- Sean R. Eddy. Profile hidden Markov models. *Bioinformatics*, 14(9):755–763, 1998.
- Sean R Eddy, Graeme Mitchison, and Richard Durbin. Maximum discrimination hidden Markov models of sequence consensus. *Journal of Computational Biology*, 2(1):9–23, 1995.
- Najib M. A. El-sayed, Elodie Ghedin, Jinming Song, Annette Macleod, Frederic Bringaud, Christopher Larkin, David Wanless, Jeremy Peterson, Lihua Hou, Sonya Taylor, Alison Tweedie, Nicolas Biteau, Hanif G. Khalak, Xiaoying Lin, Tanya Mason, Anjana J. Simpson, Samir Kaul, Hong Zhao, Grace Pai, Susan Van Aken, Teresa Utterback, Brian Haas, Hean L. Koo, Lowell Umayam, Bernard Suh, Caroline Gerrard, Vanessa Leech, Rong Qi, Shiguo Zhou, David Schwartz, Tamara Feldblyum, Steven Salzberg, Andrew Tait, C. Michael, R. Turner, Elisabetta Ullu, Owen White, Sara Melville, Mark D. Adams, Claire M. Fraser, and John E. Donelson. The sequence and analysis of *Trypanosoma brucei* chromosome II. *Nucleic Acids Res*, 16(31):4856–4863, 2003.
- Robert D. Finn, Jody Clements, and Sean R. Eddy. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Research*, 39:W29–W39, July 2011. doi: 10.1093/nar/gkr367.
- P. Flickek. Gene prediction: compare and contrast. *Genome Biology*, 8(12):233.1–233.3, 2007. doi: 10.1186/gb-2007-8-12-233.

- Julian Gough, Kevin Karplus, Richard Hughey, and Cyrus Chothia. Assignment of Homology to Genome Sequences using a Library of Hidden Markov Models that Represent all Proteins of Known Structure. *J. Mol. Biol.*, 2001. doi: 10.1006/jmbi.2001.5080.
- S.S. Gross, Do C.B., Sirota M., and Batzoglou S. CONTRAST: a discriminative, phylogeny-free approach to multiple informant de novo gene prediction. *Genome Biology*, 8(12): R269.1–R269.16, 2007. doi: 10.1186/gb-2007-8-12-r269.
- A. O. Harmanci, G. Sharma, and D. H. Mathews. Efficient pairwise RNA structure prediction using probabilistic alignment constraints in Dynalign. *BMC Bioinform.*, 8(130), 2007. doi: 10.1186/1471-2105-8-130.
- T. Hastie, R. Tibshirani, and J. Friedman. *The elements of Statistical Learning*. Spring series in statistics. Springer-Verlag, New York, USA, 2001. ISBN 0387952845.
- J. Henderson, S. Salzberg, and K. Fasman. Finding Genes in Human DNA with a Hidden Markov Model. *Journal of Computational Biology*, 4(2):127–141, 1996.
- Ø. Johansen, T. Ryen, T. Eftesøl, T. Kjosmoen, and P. Ruoff. Splice site prediction using artificial neural networks. In F. Masulli, R. Tagliaferri, and G. M. Verkhivker, editors, *Computational Intelligence Methods for Bioinformatics and Biostatistics*, volume 5488 of *Lecture Notes in Computer Science*, pages 102–113. Springer Berlin Heidelberg, 2009. ISBN 978-3-642-02503-7. doi: 10.1007/978-3-642-02504-4\_9. URL [http://dx.doi.org/10.1007/978-3-642-02504-4\\_9](http://dx.doi.org/10.1007/978-3-642-02504-4_9).
- BH. Juang, S.E. Levinson, and M.M. Sondhi. Maximum Likelihood Estimation for Multivariate Mixture Observations of Markov Chains. *IEEE Trans. Information Theory*, It-32(2):307–309, March 1986.
- Anders Krogh, Michael Brown, I Saira Mian, Kimmen Sjölander, and David Haussler. Hidden Markov Models in computational biology: Applications to protein modeling. *Journal of molecular biology*, 235(5):1501–1531, 1994.
- Zhilei Liu and Shangfei Wang. Emotion Recognition Using Hidden Markov Models from Facial Temperature Sequence. In *Proceedings of the 4th International Conference on Affective Computing and Intelligent Interaction - Volume Part II, ACII'11*, pages 240–247, 2011. ISBN 978-3-642-24570-1.
- Claudio Lottaz, Christian Iseli, C. Victor Jongeneel, and Philipp Bucher. Modeling sequencing errors by combining Hidden Markov models. *Bioinformatics*, 19(suppl 2):ii103–ii112, 2003.
- Alexander V Lukashin and Mark Borodovsky. Genemark HMM: new solutions for gene finding. *Nucleic acids research*, 26(4):1107–1115, 1998.
- K. Munch and A. Krogh. Automatic Generation of Gene Finders for Eukaryotic Species. *BMC Bioinform.*, 7(263), 2006. doi: 10.1186/1471-2105-7-263.

- World Health Organization. Trypanosomiasis, human african (Sleeping sickness). Technical Report Fact sheet Number 259, World Health Organization, 2006. URL <http://www.who.int/mediacentre/factsheets/fs259/en/>. Accessed on 05/02/2015.
- L. Pachter, M. Alexandersson, and S. Cawley. Applications of Generalized Hidden Markov Models to Aligment and Gene Finding Problems. *J. Comput. Biol.*, 9:389–399, 2002.
- L. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, Feb 1989. ISSN 0018-9219. doi: 10.1109/5.18626.
- L. R. Rabiner and Juang B. H. An Introduction to Hidden Markov Models. *IEEE ASSP MAGAZINE*, 3(1):4–16, 1986 1986.
- Lawrence R. Rabiner and Ronald W. Schafer. Introduction to Digital Speech Processing. *Found. Trends Signal Process.*, 1(1):1–194, January 2007. ISSN 1932-8346. doi: 10.1561/2000000001.
- Santhosh Rebello, Uma Maheshwari, Safreena Venex DSouza, and Rashmi Venex DSouza. Back propagation neural network method for predicting Lac gene structures in Streptococcus pyogenes M Group A Streptococcus strains. *International Journal of Biotechnology and Molecular Biology Research*, 2(4):61–72, 2011.
- K Rutherford, J Parkhill, J Crook, T Horsnell, P Rice, MA Rajandream, and B Barrell. Artemis: sequence visualization and annotation. *Bioinformatics (Oxford, England)*, 16(10), October 2000. ISSN 1367-4803. doi: 10.1093/bioinformatics/16.10.944.
- Steven Salzberg, Xin Chen, John Henderson, and Kenneth Fasman. Finding genes in DNA using decision trees and dynamic programming. In *In ISMB-96: Proc. Fourth International Conference Intelligent Systems for Molecular Biology*, pages 201–210. AAAI Press, 1996.
- Gabriele Schweikert, Alexander Zien, Georg Zeller, Jonas Behr, Christoph Dieterich, Cheng Soon S. Ong, Petra Philips, Fabio De Bona, Lisa Hartmann, Anja Bohlen, Nina Krüger, Sören Sonnenburg, and Gunnar Rätsch. mGene: accurate SVM-based gene finding with an application to nematode genomes. *Genome research*, 19(11): 2133–2143, November 2009. ISSN 1549-5469. doi: 10.1101/gr.090597.108. URL <http://dx.doi.org/10.1101/gr.090597.108>.
- Collin M Stultz, James V White, and Temple F Smith. Structural analysis based on state-space modeling. *Protein Science*, 2(3):305–314, 1993.
- E. Trentin and M. Gori. A survey of hybrid ann/hmm models for automatic speech recognition. *Neurocomputing*, 37:91–126, 2001.
- Zhuo Wang, Yazhu Chen, and Yixue Li. A brief review of computational gene prediction methods. *Genomics Proteomics Bioinformatics*, 2(4):216–221, November 2004. ISSN 1672-0229. URL <http://view.ncbi.nlm.nih.gov/pubmed/15901250>.

- Lloyd Welch. Hidden Markov Models and the Baum-Welch Algorithm. *IEEE Information Theory Society Newsletter, The Shannon Lecture by Welch*, 53(4):1,10–13, Dec 2003.
- K. J. Won, T. Hamelryck, A. Pregel-Bennett, and A. Krogh. An evolutionary method for learning HMM structure: prediction of protein secondary structure. *BMC Bioinform.*, 8(357), 2007.
- J. Yamato, Jun Ohya, and K. Ishii. Recognizing human action in time-sequential images using hidden Markov model. In *Computer Vision and Pattern Recognition, 1992. Proceedings CVPR '92., 1992 IEEE Computer Society Conference on*, pages 379–385, Jun 1992. doi: 10.1109/CVPR.1992.223161.
- Byung-Jun Yoon. Hidden Markov models and their applications in biological sequence analysis. *Current Genomics*, 10(6):402–415, 2009.
- Byung-Jun Yoon and P. P. Vaidyanathan. Structural alignment of RNAs using profile-caHMMs and its application to RNA homology search: Overview and new results. *IEEE Trans. Automat. Contr. (Joint Special Issue on Systems Biology with IEEE Transactions on Circuits and System: Part-I)*, 53:10–25, 2008.