

IN, LASEC: Bachelor Project #1

Due on Spring 2016

Pr. Serge Vaudenay

Max Premi

Abstract

This Hill cipher is a polygraphic substitution cipher based on linear algebra, invented by Lester S. in 1929. Each letter is represented by a number modulo 26, from $A = 0$ to $Z = 25$. The algorithm breaks the plaintext into blocks of size d and then applies a matrix $d \times d$ to these blocks to yield ciphertext blocks. As it's a linear encryption, it can be simply broken with Know PlainText Attacks. The author takes the previous paper about a new Ciphertext-only Attacks on Cipher Hill, and try to improve it's complexity to get a better result that $O(d13^d)$.

The goal of this project is to actually study the algorithm to get the key matrix modulo 2 and then to improve the algorithm to get the key matrix modulo 26.

The project report is organized as follows: Section1 presents the Hill cipher and the work done in the previous report. In section2, the author studies the complexity and try to improve the algorithm to get the key matrix modulo 26 . Section 3 presents the possible enhancement of the FFT of algorithm 1. Experimental results and algorithm are presented at the end.

Contents

Abstract	2
Introduction	4
Key recovery modulo 26	4
Study of Algorithm to get K_{26}	5
Study of Faster Fourier Transform for Algorithm 1	6
Simple and practical algorithm for sparse Fourier transform	7
Deterministic Sparse Fourier Approximation via Fooling Arithmetic Progressions	7
Nearly optimal Sparse Fourier Transform	7
Experiment	8
Experiment 1:Probability of independent English letters	8
Experiment 2:Probability considering blocks of size d	8
Experiment 3:Blocksize and diffence of block in n ciphers	9
Algorithm	10
References	11

Introduction

The motivation of this project is, first and foremost, to improve the Linear Cipher only attack on the Hill cipher, by changing the recovering of the key modulo 26 and then see the possible algorithm to improve the FFT in recovery of matrix key modulo 2.

Indeed, it is known that a brute force attack can be done on the Hill cipher, as it is a Linear cipher, but to have a better complexity and less restrictive resources, improvement have been made.

Considering a matrix's size of d , it is now possible to get a matrix key with minimum length required on ciphertext of $n = 8.96d^2 - O(\log d)$.

This method has been then improved [2] using the divide-and-conquer technique, and eliminating repeated calculation while doing matrix multiplication, and have led to a ciphertext required length of $n = 8.96d^2$. Eventually, using the Chinese Remainder Theorem [2], the length has been brought to $n = 12.5d^2$, and the complexity to $O(d13^d)$.

By this same Chinese Remainder theorem, it is believed that we can find the key matrix modulo 2 first and then recover the matrix modulo 26 with a lower complexity [1].

It is shown in the previous paper that this matrix modulo 2 can be found in $O(d2^d)$.

Let's briefly describe how this attack works:

Let's consider X a random vector constituted of d letters, we can pick a fixed vector $\lambda \in \{0, 1\}^*$ and consider the dot-product $\lambda.X$ in $\mathbb{Z}/2\mathbb{Z}$.

Then with the aid of the bias(X) = $\varphi_X(\frac{2\pi}{p})$ in $\mathbb{Z}/26\mathbb{Z}$, we found correspondence between λ and μ (the last is the same vector but for the cipher text). It is needed to search $bias(\lambda.X)$ and we found $\mu = (K^T)^{-1} \times \lambda$.

Then with this formula and the approximation of all the vector μ , we get the vectors column of the key matrix in $\mathbb{Z}/2\mathbb{Z}$.

An algorithm to reorder them with the correlation, to find the last one and first one easily, and then recursively find all the vectors in the correct order.

All this process is described by algorithm 1 in the Annexe, and is done in a time $O(d2^d)$

This project will present possible improvement of this algorithm to get a lower complexity than the one mentioned before, with the help of Sparse Fourier Transform. Then, a possible enhancement to get the key matrix in modulo 26 will be discussed, as the one presented in the previous paper runs in $O(8^{nd})$.

Key recovery modulo 26

So now that we have the key matrix in $\mathbb{Z}/2\mathbb{Z}$, we can have the plain text in $\mathbb{Z}/2\mathbb{Z}$ using the linearity of the cipher.

To get the key matrix in $\mathbb{Z}/26\mathbb{Z}$, we can use the Chinese Remainder Theorem, but we would get a complexity proportional to $O(13^d)$. In the previous paper, it was believed that it's possible to get the key matrix in $\mathbb{Z}/26\mathbb{Z}$ without considering $\mathbb{Z}/13\mathbb{Z}$.

First of all, we create a hash table using long text, and search mapping between segments of reference text and plain text modulo 2.

$\#(\text{seg in reference}) = \text{len}(\text{reference text}) - n + 1$, with n the segment size.

Indeed, if the following text is taken as an example: *thisisatest*, with $n = 5$, we get the following segment: *thisi, hisis, isisa, sisat, isate, sates, atest* which is 7 segments $11 - 5 + 1 = 7$

It's the same idea for $\#(\text{seg in plain}) = \text{len}(\text{plaintext}) - n + 1$, with n the segment size.

Let a , b be random segment of length n and X the plaintext. We use Rényi entropy, with the following formula:

$$H_\alpha(X) = \frac{1}{1-\alpha} \log_2 \left(\sum_{i=1}^n \Pr(X=i)^\alpha \right)$$

When alpha has the value 2, the following result is obtained:

$$-\log_2 \left(\sum_{i=1}^n \Pr(X=i)^2 \right)$$

that gives us the probability that a segment equals another one as $\sum_{i=1}^n \Pr(X=i)^2 = \sum \Pr(a=b)^2$.

Rényi entropy represent more generally the quantity of information in the probability of a random variable's collision.

Then we define good matching : segments are equals before and after modulo 2, and bad matching segment which are not equal but equal modulo 2.

For good matching, we have $E(\#goodmatching) = (\#segmentsinreference) \times (\#segmentinplaintext) \times 2^{-H_2(X)}$, as the number of good matching is actually the collision between segment in plaintext and segment in reference text multiplied by the rényi entropy of this segment (which represents the rate of collision for a given block X).

Then the same is done for $E(\#allmatching)$, the difference is that it must be taken into account that we are in $\mathbb{Z}/2\mathbb{Z}$: $E(\#goodmatching) = (\#segmentsinreference) \times (\#segmentinplaintext) \times 2^{-H_2(X \bmod 2)}$. And indeed it is understandable that if 2 words modulo 2 are equals, these words are not always equals modulo 26.

Now let's consider $E(\#allmatching)$. As we only 2 values are possible it's easier than the previous, indeed $E(\#allmatching) = (\#segmentsinreference) \times (\#segmentinplaintext) \times 2^{-H_2(X \bmod 2)}$.

$H_2(X \bmod 2) = -\log_2(\sum_{i=0}^1 \Pr(X=i)^2)$, where $\Pr(X \bmod 2 = i)$ declined in $\Pr(X \bmod 2 = 0)$ and $\Pr(X \bmod 2 = 1)$

From the experiment, we always get 0.5^n for $X \bmod 2$ so $E(\#allmatching)$ is never supposed to be different than $(\#segmentsinreference) \times (\#segmentinplaintext) \times 0.5^n$.

Then to have an idea of the complexity, the ratio $\frac{E(\#goodmatchings)}{E(\#allmatchings)}$ is computed, to find a general expression to express the number of good matching in all matching: $\frac{1}{8^n}$

In the Experiment part, the computation of $E(\#allmatchings)$ are done again with the help of a Java program. To decrease the actual complexity, we need to increase the ratio of good matching as $E(\#allmatchings)$ can't be changed. So the only solution left is to try different assumptions and calculations for $E(\#goodmatchings)$.

The one that is interesting is to consider blocks of letters as being independent from each others, and look at the evolution of the ratio through the growing block size. This is done in Experiment 2.

We finally conclude that it's possible to have a correct ratio for large size block, but the actual algorithm depends too much on the blocksize, and it's therefore impossible to get a correct complexity with the found curve that looks more like. For example, $blocksize = 27$ gives ratio $\frac{1}{5}$ but still $\frac{1}{ratio^{blocksize}}$ is too high as $blocksize = 27$. Eventually, the following parts focus on other way to implement this recovery of key matrix modulo 26.

Study of Algorithm to get K_{26}

We are going to try to improve the complexity of algorithm 2 to find another complexity than $O(8^{nd})$ with n the segment size and $d \times d$ the matrix size.

We want to turn the problem in another way, meaning instead of looking at all possible matching and do all the decryption possible with d matching, try to found the number of good matching we need so that an

algorithm can find the key matrix by solving equations.

So the problem can be turned like this: find the number y of matching needed to have a set of linear equation of the first order, to find the matrix coefficient in $\mathbb{Z}/26\mathbb{Z}$.

To be clear, let's recall what we are given. We got cipher Y_1, Y_2, \dots, Y_n in $\mathbb{Z}/26\mathbb{Z}$ but also in $\mathbb{Z}/2\mathbb{Z}$, the matrix K_2 thanks to algorithm 1, and from these two we get the plaintext in $\mathbb{Z}/2\mathbb{Z}$, X_1, X_2, \dots, X_n . We now create a matrix K_{26}^{-1} of the same size and with element $x_1, x_2, \dots, x_{d^2} \in \mathbb{Z}^{13}$.

Let's say we want to find $2d$ good matching so good equation in this. Thanks to previous study, it is known that we need $2d8^d$ matching to have the number wanted. Thus the probability that there are $2d$ good equation in the matrix, if coefficient are chosen randomly among the 13 possible values, is :

$$\binom{2d8^d}{2d} \frac{1}{13^d} \left(1 - \frac{1}{13^d}\right)^{2d8^d - 2d}$$

d is considered big enough so that $(1 - \frac{1}{13^d})^{2d8^d - 2d}$ is close to 1. From $\binom{2d8^d}{2d}$ it can be said: $\frac{8^d 2d!}{2d!(8^d 2d - 2d)!} = \frac{(8^d 2d)^{2d}}{2d!} \approx \frac{(8^d 2d)^{2d}}{2d^d}$

The previous result was obtain using Striling's approximation, $n! \simeq \sqrt{2\pi n} \left(\frac{n}{e}\right)^n$.

From here we multiply with other terms and find this: $\left(\frac{8}{13}\right)^{2d^2} 2^{2d-1} d^d$.

Let's go back to K_{26} and let's take a plaintext block X . From the linearity of Hill, $Y = K_{26} \times X$. If we pick a Y , we actually know from K_2 the value of plaintext in $\mathbb{Z}/2\mathbb{Z}$ with $Y = K_2 \times Y_i$. We can pick one of the given cipher in the list Y_1, Y_2, \dots, Y_n called from now on $Y_i = (a_0, a_1, \dots, a_d)$ and multiply it by K_{26}^{-1} and get a vector X_i constituted of d equations with d^2 unknown, x_1, \dots, x_d^2 of this form:

$$\begin{pmatrix} x_1 a_1 + x_2 a_2 + x_3 a_3 + \dots + x_d a_d \equiv y \pmod{2} \\ \dots \\ x_{d^2-d} a_{d^2-d} + x_{d^2-d+1} a_{d^2-d+1} + \dots + x_{d^2} a_{d^2} \equiv y \pmod{2} \end{pmatrix}$$

Moreover one theorem of linear algebra says:

the congruence $x_1 a_1 + x_2 a_2 + x_3 a_3 + \dots + x_n a_n \equiv b \pmod{m}$, with $(a_1, \dots, a_n, m) = c|b$ has $c|m|^{n-1}$ distinct solutions represented as matrix of size n .

Brought back to our equation, the number of distinct solutions is $c|2|^{d-1}$ with $c \in \mathbb{Z}/26\mathbb{Z}$ it will generally be 1.

We actually know the solution modulo 2 for x_1, x_2, \dots, x_d^2 thanks to the recovery of the matrix modulo 2. So there is only one solution modulo 2 for each equations. To get solution in $\mathbb{Z}/26\mathbb{Z}$, we need to use the matching given above to find possible value for the matrix, and verify by looking at the plaintext, meaning if it actually makes sense or not.

Study of Faster Fourier Transform for Algorithm 1

With a fast Fourier Transform (FFT) the complexity is $O(N \log N)$ for N the input size.

A general algorithm for computing the DFT must take time proportional to its output size N . However, in some cases, most of the Fourier coefficients of a signal are small or equal to zero, meaning, the output of the DFT is sparse.

For sparse signals, the n lower bound for the complexity of DFT no longer applies. If a signal has a small number k of non-zero Fourier coefficients the output of the Fourier transform can be represented succinctly using only k coefficients.

Hence, we can find Fourier Transform algorithm whose run time is sub-linear in the signal size N .

what we want is to enhance the possible FFT on a table called n_y which contains the number of times k where each cipher y appears. So it is a table containing numbers $\in \mathbb{N}$ of size $N = 2^d$.

Simple and practical algorithm for sparse Fourier transform

This algorithm considers a complex vector x of length l .

It computes the k -sparse Fourier transform in $O(\sqrt{kl} \log^{3/2} l)$, if x is sparse then find it takes exactly $O(k \log^2 l)$, but in general estimate x is approximately $O(\sqrt{lk})$

So this algorithm is better if the ratio $\frac{l}{k} \in [2 \times 10^3, 10^6]$, but it's clearly not the best one as recently found are supposed to find it in a lower complexity ($k \log(l)$).

Now let's consider the input of this DFT to be n_y with size 2^d . This table contains integer whose the sum is equal to the number of cipher given. As a result, we can't say that the resulting DFT of this vector will be sparse or not. It'll almost never be as we always take a large number of cipher such that it is bigger than 2^d .

Deterministic Sparse Fourier Approximation via Fooling Arithmetic Progressions

Here if we gave a threshold $\tau \in (0, 1]$ and an oracle access to a function f , it outputs the τ -significant Fourier Coefficient. This is called SFT and runs in $\log(N), \frac{1}{\tau}$.

An oracle access to a function take as input x and return the $f(x)$ of the function f .

This algorithm is robust to random noise and local (meaning it runs in polynomial time)

It's based on partition of set by binary search, we have at the beginning 4 intervals, then testing for the first two if the norm of f Fourier Transform squared is equals to the set_i oracle output squared.

Meaning more explicitly : $f(J_i)^2 = \sum_{\alpha \in J_i} |f(\alpha)|^2$ If this pass, it will output yes, and we'll be able to continue the algorithm by replacing the J and insert the J_i

The heart of the code is actually to decide which intervals potentially contain a significant Fourier coefficient. Yes if weight on J , exceeds significant threshold τ , NO if J larger.

The threshold τ can be chosen, with the fact that a α is a τ -significant Fourier coefficient iff $|\hat{f}|^2 \geq \tau \|f\|_2^2$ where $\hat{f} = \langle f, X_\alpha \rangle$ and $X_\alpha = e^{2\pi i \alpha x / N}$.

Considering the table n_y of size n with all entries=1, we get $\tau \|f\|_2^2 = \tau \times n$

And as \hat{f} got lot of small coefficient and large one on the extremities, they will be some coefficient that will satisfy this equation for small τ . So the complexity will depends on $\frac{1}{\tau}$ and it'll clearly be too big.

Let's take the same example as the previous subsection, the FFT with 200 one, we get some coefficient equals 1 (the smallest ones) and so τ needs to be inferior to $\frac{1}{800}$, so if an input more important is taken $\frac{1}{\tau}$ will be too big to match the previous complexity.

Nearly optimal Sparse Fourier Transform

We want here to compute the k -sparse approximation to the discrete Fourier transform of an 2^d -dimensional signal.

In the case where the input has at most k non-zero Fourier coefficient, we got $O(k \cdot \log(2^d))$ time, else we have $O(k \cdot \log(2^d) \cdot \log(\frac{2^d}{k}))$

The basis is still the same, if a signal has a small number k of non-zero Fourier coefficient, the output of this DFT can be represented succinctly using only k coefficient.

What is required, is that the input size n is a power of 2 which is complete in this case.

This algorithm has a better performance if and only if $k < O(\frac{2^d}{\log(2^d)})$. This won't be the case here and the formula for the k superior to this limit perform too poorly for worst case: $O(\sqrt{(2^d k) \log^{\frac{3}{2}}(2^d)})$

Experiment

Experiment 1: Probability of independent English letters

From the frequency letter given by Wikipédia, in english we got the following result :

Proba sum = 0.9999999999999999

Sum of probability squared = 0.06549717159999999, which corresponds to $(\sum_{i=0}^{25} \Pr(i = y)^2)^n, y \in \{alphabet\}$

Sum of probability that gives 0 modulo 2 squared = 0.32298762240000006 which corresponds to $(\sum_{i=0}^{25} \Pr(i = 0)^2)^n, i \in \{alphabet \bmod 2\}$

Sum of probability that gives 1 modulo 2 squared = 0.18634762239999997 which corresponds to $(\sum_{i=0}^{26} \Pr(i = 1)^2)^n, i \in \{alphabet \bmod 2\}$

Ration of good matching and all matching = 0.1285934407027314ⁿ

So $\frac{1}{7,77644^n}$.

Another site [8], with a total number of 100000 letters composed with texts from Edgar Allan Poe, Arthur Conan Doyle, and 4 articles from encyclopedia Encarta 95:

proba sum = 0.99990000000000001

sum of probability squared = 0.06609151 which corresponds to $(\sum_{i=0}^{25} \Pr(i = y)^2)^n, y \in \{alphabet\}$

sum of probability that gives 0 modulo 2 squared = 0.32001649 which corresponds to $(\sum_{i=0}^{25} \Pr(i = 0)^2)^n, i \in \{alphabet \bmod 2\}$

sum of probability that gives 1 modulo 2 squared = 0.18852964 which corresponds to $(\sum_{i=0}^{25} \Pr(i = 1)^2)^n, i \in \{alphabet \bmod 2\}$

Ration of good matching and all matching = 0.12996168115565054ⁿ

So $\frac{1}{7,69457^n}$.

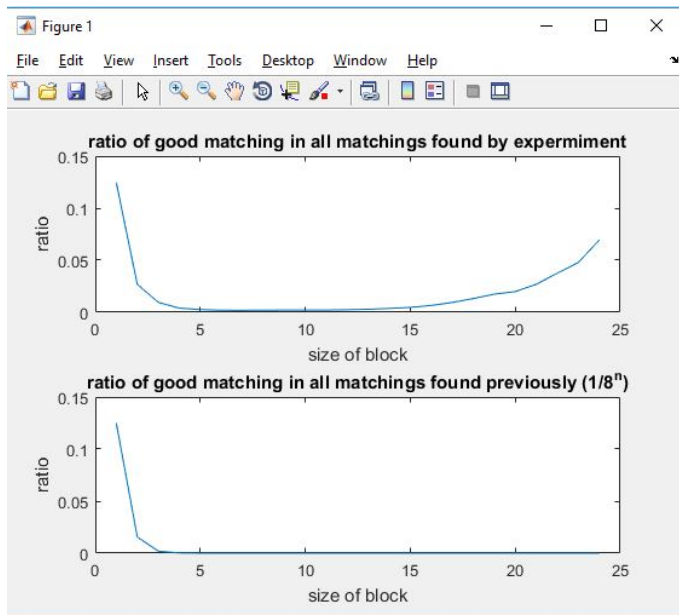
Experiment 2: Probability considering blocks of size d

So calculation are done on a text of approximately 860000 characters to see the evolution of the ratio good matching/bad matching.

A program is ran to see the evolution for a block size between 1 and 25, and give the ratio, thanks to the probability that a block appears. It is completely heuristic as it's just counting the number of block that appears and do some manipulation with it. So the basic is to choose a block size, then it'll count every different blocks that appears modulo 26 and modulo 2. Then it'll compute the probability that a good matching happen with the following : $\sum_{X \in block} (\frac{\#X-1}{\#block-1})^2$

The exact same thing is done with X in modulo 2, to get the probability of all matching, and then we compute the ratio $\frac{E(\#goodmatchings)}{E(\#allmatchings)}$

With this, the evolution of the ratio in function of the block size looks like this:



So we can see that the ratio follow the $\frac{1}{8^n}$ until blocksize 15, then it goes up again to almost match $1/8$ for blocksize 24. With the current algorithm and this size of block the number of iteration would be 8^d and as $d = 24$ is really large it's still not effective. Even if the ratio do not behave as previously thought, the complexity stay too high for reasonable blocksize between 8 and 14, and if the blocksize is increased to a certain point, the ratio is good but the complexity depends on a ratio power the blocksize, so it will not be good enough to be taken in account.

Experiment 3:Blocksize and diffence of block in n ciphers

blocksize	Probability that 2 blocks are different	# maximal of ciphers for $\text{Pr} > 0.95$
1	(0.9367250929)	less than 1
2	(0.99290120378)	7
3	(0.99868328054)	38
4	(0.99971599715)	180
7	(0.99998010694)	2578
10	(0.99999683297)	16195
12	(0.99999899833)	51207
14	(0.99999959261)	125907
16	(0.99999985551)	354995

Indeed, it's obvious that the bigger the block is, the lower is the chance that another one exist with same value. For blocksize 1, you cannot hop to have a good probability to get another block different as there are only 26 possibilities. You just have to resolve the following: $\text{Probability}^{\# \text{ciphers}} > 0.95 \Rightarrow n < \frac{\log(0.95)}{\log(\text{Probability})}$.

Algorithm

You hash a reference text.

You take the key matrix that you get from algorithm 1, find plain text in $\mathbb{Z}/2\mathbb{Z}$, and create an array.

find the list of all matching: find all pairs (seg, str) such that seg is a segment of plaintext modulo 2 and $str \in hash(seg)$ and save it in a list.

```

1: repeat
2:   select d matching form list (you'll get a  $d \times d$  key matrix)
3:   for each of these matchings  $(seg_i, str_i)$  do
4:     extract  $block_i$  from  $seg_i$  and  $str'_i$  from  $str_i$ ,
5:     then find  $ciphertext_i$  such that  $K^{-1} \times ciphertext_i \bmod 2 = block_i$ 
6:   end for
7:   solve  $ciphertext_i = K * str'_i$  for  $i=1$  to  $d$ 
8:   compute  $K^{-1} * ciphertext$ 
9: until decryption make sense
number of iteration is  $\frac{1}{ratio^{nd}} = 8^{nd}$ 

```

The following algorithm is to recover the key matrix in $\mathbb{Z}/2\mathbb{Z}$

```

1: Part1:
Require: Ciphertext  $Y_1, Y_2, \dots, Y_n$ 
Ensure:  $K(\bmod 2)$ 
2: for all  $\mu$  do
3:   compute  $S_n(\mu) = \sum_y (-1)^{\mu \cdot y} \times n_y$  where  $n_y = \#\{k; Y_k = y\}$ 
4: end for
5: set all  $\mu$  to the  $d$  values of  $\mu$  with largest  $S_n(\mu) = bias(\mu.Y)$ 
6: Part2:
7: for all  $(i, i')$  do
8:   compute  $n_{00}(i, i') = \#\{k < n : (\mu_i.Y_k, \mu_{i'}.Y_{k+1}) = (0, 0)\}$ 
9: end for
10: set  $(i_d, i_1)$  to the first pair with lowest  $n_{00}$ 
11: Part3:
12: for all  $t = 2$  to  $d - 1$  do
13:   for all  $i \notin \{i_1, i_2, \dots, i_{t-1}, i_d\}$  do
14:     compute  $n_{00}(i, i') = \#\{k : (\mu_{i_{t-1}}^T Y_k, \mu_i^T Y_k) = (0, 0)\}$ 
15:   end for
16:   take  $i$  such that  $n_{00}$  is minimum and set  $i_t = i$ 
17: end for
18: set  $\mu = (\mu_{i_1}, \mu_{i_2}, \dots, \mu_{i_d})$  and  $K = (\mu^{-1})^T$ 
19: output  $K$ 

```

Here to be faster we store n_y in a table and we do a FFT on this table to get S_n . With this operation the total complexity drop from $O(d^2 \times 2^d)$ to $O(d \times 2^d)$ But it seems with some other techniques we could do better.

References

- [1] Alina, Matyukhina. *Cryptanalysis of the Hill Cipher*.
- [2] S. Shazaei, S. Ahmadi. *Ciphertext- only attack on $d \times d$ Hill in $O(d13^d)$* .
- [3] Akavia, A. *Deterministic Sparse Fourier Approximation via Fooling Arithmetic Progressions*.
- [4] Akavia, A., Goldwasser, S., Safra, S. *Proving Hard-Core Predicates Using List Decoding*.
- [5] Hassanieh, H., Indyk, P., Katabi, D., Price, E. *Nearly optimal sparse Fourier transform*.
- [6] Hassanieh, H., Indyk, P., Katabi, D., Price, E. *Simple and practical algorithm for sparse Fourier transform*.
- [7] Iwen, M.A. *Combinatorial Sublinear-Time Fourier Algorithms*.
- [8] <http://www.nymphomath.ch/crypto/stat/anglais.html>