



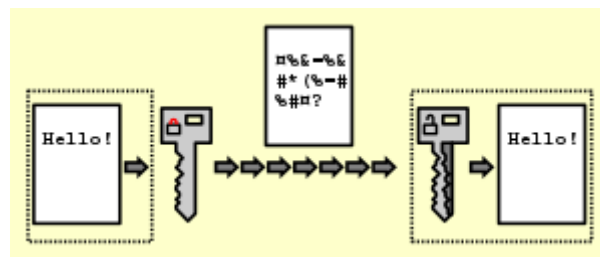
**DEEC**  
DEPARTAMENTO DE ENGENHARIA  
ELETROTÉCNICA E DE COMPUTADORES  
TÉCNICO LISBOA

# PROGRAMAÇÃO

## (LEEC 23/24)

### Enunciado do Projeto

### Cifras



## 1. Introdução

O trabalho que se descreve corresponde ao projeto da UC de Programação para a LEEC em 2023/2024. O projeto a realizar corresponde ao desenvolvimento de um programa para implementar mecanismos simples para cifrar e decifrar informação. O objetivo do projeto é promover e avaliar a prática dos conceitos e técnicas de programação estudados no âmbito desta UC. O projeto deverá ser realizado ao longo do período acompanhando a matéria teórica lecionada nesta UC. Conforme descrito na secção de avaliação, haverá uma entrega intermédia e uma entrega final, que corresponde à entrega do projeto completo. Não se pretende de maneira nenhuma que os métodos de cifra e decifra sejam seguros, mas apenas ilustrar alguns mecanismos básicos e como eles podem ser usados. Lembra-se ainda que pode ser considerado crime tentar decifrar dados de outros sem autorização. Para mais informações consulte o site do Centro Nacional de Cibersegurança [1].

## 2. Descrição do Problema a Resolver

O programa deve poder funcionar com dados lidos do teclado e o resultado do programa ser escrito no écran, ou funcionar com ficheiros, podendo ser parametrizado através de argumentos na linha de comando. De seguida descrevem-se as várias funcionalidades a implementar com exemplos.

### 2.1. Conjunto de Caracteres a Considerar

Na Tabela 1 são apresentados, a negrito, os 67 caracteres a considerar, por ordem: algarismos de 0..9; maiúsculas de A..Z; minúsculas de a..z; espaço em branco; ponto final; vírgula; ponto e vírgula; hífen.

A correspondente numeração também é indicada na Tabela 1, a vermelho, de 0..66.

**Tabela 1: Conjunto de caracteres a utilizar a negrito, com numeração a vermelho**

0	1	2	3	4	5	6	7	8	9	10	11	12	13	14
0	1	2	3	4	5	6	7	8	9	A	B	C	D	E
15	16	17	18	19	20	21	22	23	24	25	26	27	28	29
F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
30	31	32	33	34	35	36	37	38	39	40	41	42	43	44
U	V	W	X	Y	Z	a	b	c	d	e	f	g	h	i
45	46	47	48	49	50	51	52	53	54	55	56	57	58	59
j	k	l	m	n	o	p	q	r	s	t	u	v	w	x
60	61	62	63	64	65	66								
y	z		.	,	;	-								

### 2.2. Método 1: Cifra de César/Substituição

Na cifra de substituição [2], ou cifra de César, cada carácter do conjunto de caracteres a utilizar, definido na Tabela 1, é substituído por outro com um determinado offset fixo dado pela 1ª letra da senha definida. Por exemplo, se a 1ª letra da senha for o carácter '3', codifica-se cada carácter avançando 3 caracteres na tabela, como se a tabela comesse no carácter '3', conforme ilustrado na Fig. 1. Naturalmente que os caracteres que excedam 66 dão a volta módulo 67. Os caracteres que não constem na Tabela 1, não são modificados, como por exemplo, o carácter '!'. Os restantes caracteres da senha, depois do 1ª, não são usados neste método.

Dados	0	1	2	3	4	5	6	7	8	9	A	...	-
Cifrado	3	4	5	6	7	8	9	A	B	C	D	...	2

**Fig. 1: Exemplo da troca de caracteres da cifra de César/Substituição com senha começada por '3'.**

Matematicamente, o caracter cifrado é igual ao caracter dos dados mais o offset, módulo 67:

$$\text{cifrado} = (\text{dados} + \text{offset}) \% 67$$

Para decifrar a mensagem, basta subtrair o offset fixo a cada caracter cifrado para obter o caracter decifrado. Para não se obter números negativos, pode ser feito:

$$\text{dados} = (\text{cifrado} + 67 - \text{offset}) \% 67$$

Exemplo de cifra com offset '3', que numericamente corresponde a 3 pela Tabela 1:

dados: To Sherlock Holmes she is always

cifrado: Wr;Vkhurfn;Krophv;vkh;lv;dozd.v

Exemplo de cifra com offset 'P', que numericamente corresponde a 25 pela Tabela 1:

dados: To Sherlock Holmes she is always

cifrado: s8Kr1;B58.4Kg856;CKC1;K2CKz5GzIC

Note-se que se se utilizar offset '0', os dados não são alterados. Caso se cifre com '3' e depois se cifre novamente com '3', obtém-se os dados originais, porque  $3+64=67$  que, módulo 67, dá 0.

Naturalmente que não sabendo qual a senha/offset usado, o código pode ser atacado com um método de força bruta que consiste em experimentar os 67 offsets possíveis.

### 2.3. Método 2: Cifra de Vigenère/Polialfabética

No método 2, cifra de Vigenère [3], cada caracter do conjunto de caracteres a utilizar, definido na Tabela 1, é substituído por outro com um determinado offset variável, tal como na cifra de César, consoante a letra da senha onde se vai. O primeiro caracter cifrado usa o offset correspondente à primeira letra da senha. O segundo caracter cifrado usa o offset correspondente à segunda letra da senha, e assim sucessivamente. Depois de se usar o último caracter da senha, volta-se ao primeiro caracter da senha, como se a senha se repetisse indefinidamente.

Na Fig. 2 ilustra-se o offset a aplicar a cada caracter dos dados a cifrar quando a senha é "123". A preto está a informação a cifrar. A verde está a senha "123" repetida até ao fim dos dados. A azul, estão os caracteres cifrados com o offset indicado a verde para esse caracter.

Dados	T	o		S	h	e	r	l	o	c	k		H
Offset	1	2	3	1	2	3	1	2	3	1	2	3	1
Cifrado	U	q	;	T	j	h	s	n	r	d	n	;	l

Fig. 2: Exemplo da troca de caracteres da cifra de Vigenère com senha "123".

Para decifrar a informação, basta decifrar cada caracter cifrado com a correspondente letra da senha, tal como no método 1, cifra de César.

Exemplo de cifra com senha "Programacao2024":

dados: To Sherlock Holmes she is always

cifrado: saj3T9YGL7T,Hqp6QbbeCLVFNjclYeIe

Este código pode ser atacado com análise estatística de frequência das letras. Mas se a senha for mais comprida do que mensagem a cifrar, pode ser impossível de quebrar. O conhecimento de algumas palavras cifradas pode ajudar.

#### 2.4. Método 3: Cifra de Hill/Matricial

No método 3, cifra de Hill [4], codificam-se pares de caracteres de dados, multiplicando por uma matriz de 2 linhas por 2 colunas, formada pelos 4 primeiros caracteres da senha a utilizar, módulo 67, produzindo pares de caracteres cifrados. Nesta cifra, é obrigatório filtrar os dados de entrada retirando os que não pertencem à Tabela 1 e formatar os dados de saída, conforme explicado na secção 2.5. No caso de se pretender cifrar um número ímpar de caracteres (porque por exemplo o ficheiro acaba), acrescenta-se um espaço em branco no fim dos caracteres a cifrar de forma a ter um número par de caracteres. Não se descarta este espaço em branco na descodificação.

Sendo a senha começada pelos 4 caracteres “abcd”, forma-se a matriz  $K = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$ .

Sendo o par de caracteres a cifrar  $m = [m_1 \ m_2]$ , cifra-se calculando o produto de matrizes  $m.K$ , módulo 67, que dá:

$$m.K = [m_1 \ m_2] \cdot \begin{bmatrix} a & b \\ c & d \end{bmatrix} = [m_1 \cdot a + m_2 \cdot c \quad m_1 \cdot b + m_2 \cdot d] \text{ módulo } 67$$

Se a senha for “1234”, teremos  $K = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ . Para cifrar as letras “To”, que se codificam como  $m = [29 \ 50]$ , conforme a Tabela 1, calcula-se o produto:

$$m.K = [29 \ 50] \cdot \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = [29 \cdot 1 + 50 \cdot 3 \quad 29 \cdot 2 + 50 \cdot 4] = [179 \ 258] \text{ módulo } 67 = [45 \ 57]$$

que corresponde às letras “jv”, conforme a Tabela 1.

Para decifrar, primeiro é necessário calcular a inversa da matriz  $K$ , calculada a partir do determinante da matriz  $K$  e da matriz adjunta de  $K$ . A matriz inversa é dada por:  $K^{-1} = \det(K)^{-1} \cdot \text{adj}(K)$ , sendo os cálculos feitos módulo 67. O determinante é calculado por:

$$\det(K) = \det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = (a \cdot d - b \cdot c) \text{ módulo } 67$$

Para que não dê um número negativo, sugere-se somar 67x67 antes de realizar a subtração.

No caso de o determinante dar zero, a matriz  $K$  não é invertível, pelo que, neste caso, na matriz  $K$  se incrementa sucessivamente as entradas  $a$  e  $d$  de uma unidade,  $b$  e  $c$  de duas unidades, módulo 67, recalculando-se o determinante até o determinante dar não nulo. Isto tem de ser feito tanto antes da cifra como antes da decifra. Por exemplo, se o utilizador der a senha “1212”, o determinante daria zero, e não seria possível decifrar os dados cifrados. Com este método, a senha seria convertida em “2433”, para a qual o determinante já seria -6 módulo 67 = 61, não nulo, permitindo cifrar e decifrar os dados. No caso de a senha ter menos de 4 caracteres, dá-se um erro ao utilizador indicando que a senha é demasiado curta. No caso de a senha ter mais de 4 caracteres, só se utilizam os primeiros 4 caracteres.

O inverso do determinante módulo 67 pode ser tirado da seguinte tabela, que foi calculada como explicado no anexo 1:

{ 0, 1, 45, 30, 67, 18, 15, 51, 78, 10, 9, 81, 52, 48, 70, 6, 39, 21, 5, 75, 49, 17, 85, 31, 26, 57, 24, 33, 35, 43, 3, 23, 64, 27, 55, 28, 47, 77, 82, 16, 69, 76, 53, 29, 87, 2, 60, 36, 13, 20, 73, 7, 12, 42, 61, 34, 62, 25,

66, 86, 46, 54, 56, 65, 32, 63, 58, 4, 72, 40, 14, 84, 68, 50, 83, 19, 41, 37, 8, 80, 79, 11, 38, 74, 71, 22, 59, 44, 88 }

A matriz adjunta é dada por:

$$\text{adj}(K) = \text{adj}\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \text{ módulo } 67$$

Para não ficar com números negativos, sugere-se somar 67 nos termos com '- '.

Se a senha for "1234", teremos  $K = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$ .  $\det(K) = 1 \cdot 4 - 2 \cdot 3$  módulo 67 =  $4 - 6$  módulo 67 =  $-2$  módulo 67 =  $67 - 2 = 65$ . Pela tabela do anexo 1,  $\det(K)^{-1}$  módulo 67 =  $(65)^{-1}$  módulo 67 = 33 porque  $65 \times 33$  módulo 67 =  $2145$  módulo 67 = 1.

Logo a matriz inversa será:

$$K^{-1} = \det(K)^{-1} \cdot \text{adj}(K) = 33 \cdot \begin{bmatrix} 4 & 67 - 2 \\ 67 - 3 & 1 \end{bmatrix} \text{ módulo } 67 = 33 \cdot \begin{bmatrix} 4 & 65 \\ 64 & 1 \end{bmatrix} \text{ mod } 67 = \begin{bmatrix} 65 & 1 \\ 35 & 33 \end{bmatrix}$$

Para decifrar as letras "jv", que correspondem aos códigos [45 57], multiplica-se pela matriz inversa para recuperar os dados originais:

$$\begin{aligned} m \cdot K \cdot K^{-1} &= \begin{bmatrix} 45 & 57 \end{bmatrix} \cdot \begin{bmatrix} 65 & 1 \\ 35 & 33 \end{bmatrix} = \begin{bmatrix} 45 \cdot 65 + 57 \cdot 35 & 45 \cdot 1 + 57 \cdot 33 \end{bmatrix} = \\ &= \begin{bmatrix} 4920 & 1926 \end{bmatrix} \text{ módulo } 67 = \begin{bmatrix} 29 & 50 \end{bmatrix} \end{aligned}$$

que resulta nas letras originais "To".

Note-se que a senha "1001" resulta na matriz identidade, pelo que os dados não são alterados quando se cifra com esta senha. Note-se que a senha "0110" resulta em trocar os caracteres de ordem dentro de cada par de caracteres.

Exemplo de cifra com senha "0110":

**dados:** To Sherlock Holmes she is always

**cifrado:** oTS ehlrco koHmlses ehi slaawsy

Exemplo de cifra com senha "Programacao2024":

**dados:** To Sherlock Holmes she is always

**cifrado:** 9J2ey6vQ1PTHirOIFXTyy6;gSdY.YUkL

Este código pode ser atacado com análise estatística de pares de letras. Se for usada uma matriz maior, fica mais difícil a análise.

## 2.5. Filtragem do ficheiro de entrada e formatação do ficheiro de saída

Se não for usada filtragem, os caracteres dos dados a cifrar que não constem na Tabela 1 são passados para o output sem qualquer modificação. Isto significa que, caso não se use filtragem, cifrando e decifrando com os métodos 1 e 2 recuperam-se os dados originais sem qualquer modificação.

Caso se use filtragem, os caracteres dos dados a cifrar que não constem na Tabela 1 (incluindo mudanças de linha) são eliminados. Neste caso, também é formatado o output cifrado que deve ser escrito em linhas com um máximo de 8 blocos de 6 caracteres cifrados, separados por um sublinhado

('\_'). No caso do output decifrado escreve-se em linhas de 48 caracteres sem espaços a separar os blocos.

No fim da escrita de um bloco de 6 caracteres cifrados escreve-se sempre o sublinhado separador ou uma mudança de linha, independentemente de haver mais dados a escrever ou de se ter chegado ao fim dos dados.

No fim dos dados deve-se mudar de linha, quer ela esteja incompleta ou completa, tanto nas operações de cifra como nas operações de decifra. Os sublinhados, não sendo caracteres da Tabela 1, são descartados na decifra da informação.

Naturalmente que, no caso de se usar filtragem e formatação dos dados, os dados cifrados não são recuperados no seu formato original exacto após a decifra.

A filtragem é obrigatória no método de cifra 3, de matriz, por uma questão de simplificação, porque os caracteres têm de ser tratados aos pares.

Segue-se um exemplo de cifra e decifra com formatação, com o método 1, e senha começada por 'P', (note-se que no fim de cada linha dos dados há um espaço em branco):

To Sherlock Holmes she is always the woman. I have seldom heard him  
mention her under any other name. In his eyes she eclipses and  
predominates the whole of her sex. It was not that he felt any emotion  
akin to love for Irene Adler. All emotions, and that one particularly,

```
s8Kr1;_B58.4K_g856;C_KC1;K2_CKz5Gz_ICKD1;_KG86z7_LKhK1z
F;KC;5_,86K1;_zB,K12_6K6;7D_287K1;_BKE7,;_BKz7IK_8D1;BK
7z6;LK_h7K12C_K;I;CK_C1;K;. _529C;C_Kz7,K9_B;,862_7zD;CK
D1;KG1_85;K8-_K1;BKC_;HLKhD_KGzCK7_8DKD1z_DK1;K-;5DKz7
IK;68D_287Kz4_27KD8K_58F;K-_8BKHB;_7;KZ,5;_BLKZ5_5K;68D
287CMK_z7,KD1_zDK87;_K9zBD2_.E5zB5_IMK
```

To Sherlock Holmes she is always the woman. I ha  
ve seldom heard him mention her under any other  
name. In his eyes she eclipses and predominates  
the whole of her sex. It was not that he felt an  
y emotion akin to love for Irene Adler. All emot  
ions, and that one particularly,

## 2.6. Cálculo de estatísticas de um ficheiro

Prende-se que o programa apresente estatísticas da frequência de caracteres da Tabela 1 encontrados nos dados de entrada. A frequência de um dado carácter é definida como o quociente entre o número de caracteres desse dado tipo encontrados na totalidade dos dados de entrada e o número total de caracteres da Tabela 1 encontrados na totalidade dos dados de entrada.

O programa deve escrever uma linha por cada carácter da Tabela 1, contendo a contagem de caracteres desse tipo, um tab, e a frequência desse carácter calculada como um `double` em percentagem. Depois, deve escrever uma linha com o total de caracteres da Tabela 1. Depois, deve escrever uma linha com a contagem de caracteres que não constam na Tabela 1 e a sua frequência no ficheiro completo. Finalmente, a última linha é a contagem de todos os caracteres do ficheiro.

Considere-se o exemplo do ficheiro "1line.txt", cujo conteúdo pode ser visto com o comando "hexdump":

```
$ hexdump -C lline.txt
00000000 54 6f 20 53 68 65 72 6c 6f 63 6b 20 48 6f 6c 6d |To Sherlock Holm|
00000010 65 73 20 73 68 65 20 69 73 20 61 6c 77 61 79 73 |es she is always|
00000020 20 74 68 65 20 77 6f 6d 61 6e 2e 0a                | the woman..|
0000002c
```

De seguida apresenta-se o output que se pretende que o programa tenha:

```
conta('0')=0 0.000000%
conta('1')=0 0.000000%
conta('2')=0 0.000000%
conta('3')=0 0.000000%
conta('4')=0 0.000000%
conta('5')=0 0.000000%
conta('6')=0 0.000000%
conta('7')=0 0.000000%
conta('8')=0 0.000000%
conta('9')=0 0.000000%
conta('A')=0 0.000000%
conta('B')=0 0.000000%
conta('C')=0 0.000000%
conta('D')=0 0.000000%
conta('E')=0 0.000000%
conta('F')=0 0.000000%
conta('G')=0 0.000000%
conta('H')=1 2.325581%
conta('I')=0 0.000000%
conta('J')=0 0.000000%
conta('K')=0 0.000000%
conta('L')=0 0.000000%
conta('M')=0 0.000000%
conta('N')=0 0.000000%
conta('O')=0 0.000000%
conta('P')=0 0.000000%
conta('Q')=0 0.000000%
conta('R')=0 0.000000%
conta('S')=1 2.325581%
conta('T')=1 2.325581%
conta('U')=0 0.000000%
conta('V')=0 0.000000%
conta('W')=0 0.000000%
conta('X')=0 0.000000%
conta('Y')=0 0.000000%
conta('Z')=0 0.000000%
conta('a')=3 6.976744%
conta('b')=0 0.000000%
conta('c')=1 2.325581%
conta('d')=0 0.000000%
conta('e')=4 9.302326%
conta('f')=0 0.000000%
conta('g')=0 0.000000%
conta('h')=3 6.976744%
conta('i')=1 2.325581%
conta('j')=0 0.000000%
conta('k')=1 2.325581%
conta('l')=3 6.976744%
conta('m')=2 4.651163%
conta('n')=1 2.325581%
conta('o')=4 9.302326%
conta('p')=0 0.000000%
conta('q')=0 0.000000%
conta('r')=1 2.325581%
conta('s')=4 9.302326%
conta('t')=1 2.325581%
conta('u')=0 0.000000%
conta('v')=0 0.000000%
conta('w')=2 4.651163%
conta('x')=0 0.000000%
conta('y')=1 2.325581%
```

```

conta('z')=0 0.000000%
conta(' ') =7 16.279070%
conta('.')=1 2.325581%
conta(',')=0 0.000000%
conta(';')=0 0.000000%
conta('-')=0 0.000000%
Total: 43 caracteres
conta(outros)=1 2.272727%
Total do ficheiro: 44 caracteres

```

Note-se que havendo 43 caracteres da Tabela 1 no ficheiro (linha “Total:”) e sendo a contagem do carácter ‘H’ de 1, a sua frequência é  $1/43$ , ou em percentagem  $1/43 \cdot 100\% = 2.325581\%$ .

Há 1 caracteres no ficheiro que não constam na Tabela 1: um carácter ‘\n’ no fim do ficheiro. A frequência de outros caracteres é assim  $1/44 \cdot 100\% = 2.272727\%$ , conforme apresentado em cima.

## 2.7. Ataque método 1: baseado em dicionário

Como dito acima, o método de cifra 1 pode ser atacado com um mecanismo de força bruta que consiste em testar todas os 67 offsets possíveis.

Deve ser implementado um método de ataque baseado em dicionário que testa todos os offsets possíveis, verificando para cada um quantas palavras de um dicionário dado consegue encontrar. Aquele que corresponder a um maior número de palavras válidas é considerado o offset correto.

Para começar, o programa tem de ler todas as palavras de um ficheiro de dicionário fornecido para estruturas de dados em memória, alocadas dinamicamente, para depois a pesquisa de palavras ser rápida. Cada linha do ficheiro de dicionário tem uma única palavra que começa no 1º carácter da linha e termina com mudança de linha (‘\n’), espaço em branco (‘ ’), tab (‘\t’) ou sinal de dividir (‘/’). Se for encontrado o carácter plica numa palavra, este carácter é descartado, pois não faz parte dos caracteres da Tabela 1, sendo a palavra armazenada sem a plica, a menos que já constasse no dicionário. Se, depois disto, for encontrado algum carácter que não consta da Tabela 1, a palavra é descartada e não armazenada no dicionário em memória. Por exemplo, “Asunción”, “Asuncións”, “Atatürk”, etc., seriam descartadas.

Assim, por exemplo, com a linha:

```
aardvark/SM
```

é guardado no dicionário “aardvark” sem o hífen e parando no ‘/’.

Com as linhas:

```

aardvark's
aardvarks
abacus's

```

são guardadas no dicionário “aardvarks”, apenas uma única vez, e “abacuss”.

Sugere-se que o dicionário seja armazenado numa árvore, com um nó por letra, para tornar as pesquisas mais rápidas. O anexo 2 ilustra como isso poderia ser feito.

De seguida, lê-se os dados a analisar para memória. Novamente, deve ser usada memória dinâmica para armazenamento dos dados. Cria-se uma cópia dos dados descartando todos os caracteres que não façam parte da Tabela 1, tal como no método de cifra 1.

Por fim, percorrem-se todos valores possíveis para o offset da cifra de César.



Para cada valor de offset possível, percorrem-se todos os caracteres do ficheiro, considerando esse carácter como possível carácter inicial para procurar palavras.

Procuram-se palavras no dicionário começando nesse carácter, contando quantas palavras se encontram no dicionário a começar nesse carácter, e considerando como *score* o somatório dos quadrados dos tamanhos das palavras encontradas no dicionário.

Avança-se para o offset seguinte, percorrendo todos os caracteres do ficheiro, substituindo cada um pelo próximo módulo 67.

O offset que tiver resultado no maior *score* referente às palavras encontradas no dicionário é a solução. O resultado é apresentado ao utilizador, bem como a correspondente descodificação dos dados. Para descodificar os dados, é necessário percorrer todos os caracteres dos dados, antes do descarte dos caracteres que não fazem parte da Tabela 1, convertendo apenas os que façam parte da Tabela 1, e transcrevendo os outros sem modificação.

Exemplo. Considere que tem o seguinte texto, de um exemplo da secção 2.2:

s8Kr1;B58.4Kg856;CKC1;K2CKz5GzIC

e o dicionário:

and  
any  
as  
he  
is  
not  
of  
she  
the

Com offset 35, encontra-se apenas uma única palavra, estando listada debaixo do carácter do texto onde foi encontrada:

MhtLaXkehVdtAhefXltlaXtbltTepTrltmaXtphfTgu  
he

Note-se que não é necessário as palavras estarem terminadas por um espaço para serem contabilizadas.

Com offset 42, encontram-se as seguintes 4 palavras, estando listadas debaixo do carácter do texto onde foram encontradas:

To Sherlock Holmes she is always  
he  
she  
he  
is

No fim, este é o offset para o qual se encontram mais palavras. O correspondente *score* é  $2 \times 2 + 3 \times 3 + 2 \times 2 + 2 \times 2 = 21$ . O offset 42, pela Tabela 1 corresponde à letra 'g', mas a correspondente cifra terá sido feita com  $(67-42) \% 67 = 25$ , que pela Tabela 1 corresponde à letra 'P'. Neste caso, o ataque é bem sucedido. O programa escreve o resultado, bem como o texto descodificado, conforme se apresenta:

offset máximo 42, letra 'g', cifrado com letra 'P', total 4 palavras, score 21.  
To Sherlock Holmes she is always

## 2.8. Ataque método 2: baseado em probabilidades para cifra de César

As probabilidades de ocorrência das letras em diferentes línguas são conhecidas [5]. A probabilidade de ocorrência dos algarismos, em muitos casos, segue a lei de Benford [6].

Considerando as frequências das letras em inglês indicadas na Wikipédia [5], considerando que a probabilidade do espaço é o dobro da letra mais frequente (o 'e'), assumindo que as maiúsculas são 30 vezes menos prováveis que as minúsculas, considerando a probabilidade de cada algarismo dada pela lei de Benford considerando apenas o 1º e 2º algarismo dos números e que a probabilidade do '1' é igual à do 'z', considerando que a probabilidade da vírgula é igual à do 'p' e a do ponto final é igual à do 'g', a probabilidade do ponto e vírgula é igual à do 'z' e a probabilidade do hífen é igual à do 'j', obtém-se a seguinte tabela para as probabilidades dos caracteres da Tabela 1, após normalização (para que a soma dê 1):

```
{ 0.000160079, 0.000554985, 0.000381091, 0.000306662, 0.000263793, 0.000235223,
  0.000214441, 0.00019842, 0.000185551, 0.000174894, 0.002049944, 0.00037499,
  0.000699981, 0.00107497, 0.003174913, 0.000549985, 0.000499986, 0.001524958,
  0.001749952, 3.7499E-05, 0.000192495, 0.000999973, 0.000599984, 0.001674954,
  0.001874949, 0.000474987, 2.37493E-05, 0.001499959, 0.001574957, 0.002274938,
  0.000699981, 0.000244993, 0.000599984, 3.7499E-05, 0.000499986, 1.84995E-05,
  0.061498311, 0.011249691, 0.020999423, 0.032249114, 0.095247384, 0.016499547,
  0.014999588, 0.045748744, 0.052498558, 0.001124969, 0.005774841, 0.029999176,
  0.017999506, 0.05024862, 0.056248455, 0.014249609, 0.00071248, 0.044998764,
  0.047248702, 0.068248126, 0.020999423, 0.007349798, 0.017999506, 0.001124969,
  0.014999588, 0.000554985, 0.190494768, 0.014249609, 0.014999588, 0.000554985,
  0.001124969 }
```

As probabilidades a usar devem ser sempre as indicadas acima.

Deve ser implementado um método de ataque baseado em probabilidades que testa todos os offsets possíveis da cifra de César, verificando qual corresponde a um menor erro quadrático médio entre as frequências de ocorrência das letras nos dados e as probabilidades tabeladas. Aquele offset que corresponder ao menor erro quadrático médio é considerado o offset correto.

Para realizar este ataque, o programa começa por ler o ficheiro a analisar para memória, tal como no método de ataque 1. Novamente, deve ser usada memória dinâmica para armazenamento dos dados.

Depois, calculam-se as frequências de cada carácter da Tabela 1, tal como indicado no cálculo de estatísticas explicado na secção 2.6.

De seguida, percorrem-se todos valores possíveis para o offset, calculando o erro quadrático médio entre as frequências de ocorrência das letras calculadas a partir dos dados (*freq*) e as probabilidades tabeladas (*prob*) com esse offset, conforme a seguinte fórmula, e fazendo as contas com variáveis do tipo `double`:

$$erro(offset) = \sum_{i=0}^{66} \frac{(freq[i] - prob[(i + offset) \% 67])^2}{prob[(i + offset) \% 67]}$$

O offset que corresponder ao menor erro quadrático médio é o seleccionado.

Considere-se como exemplo, o ficheiro `1line.txt` apresentado anteriormente e cifrado com o método de César com o resultado indicado na secção 2.2. Percorrendo todos os offsets possíveis de 0 a 66, obtinha-se, sucessivamente (para simplificar, não se apresentam todos os resultados):

```
offset 0, letra '0', erro quadrático médio: 272.997165.
```

offset 1, letra 'l', erro quadrático médio: 161.850810.  
offset 2, letra '2', erro quadrático médio: 171.115201.  
offset 3, letra '3', erro quadrático médio: 192.723549.  
offset 4, letra '4', erro quadrático médio: 148.931592.  
offset 5, letra '5', erro quadrático médio: 166.647241.  
offset 6, letra '6', erro quadrático médio: 1251.019047.  
offset 7, letra '7', erro quadrático médio: 357.302215.  
offset 8, letra '8', erro quadrático médio: 208.437593.  
offset 9, letra '9', erro quadrático médio: 107.105496.  
offset 10, letra 'A', erro quadrático médio: 246.629764.  
offset 11, letra 'B', erro quadrático médio: 440.690860.  
offset 12, letra 'C', erro quadrático médio: 166.700277.

offset 38, letra 'c', erro quadrático médio: 14.811947.  
offset 39, letra 'd', erro quadrático médio: 210.456887.  
offset 40, letra 'e', erro quadrático médio: 53.883851.  
offset 41, letra 'f', erro quadrático médio: 315.642658.  
offset 42, letra 'g', erro quadrático médio: 1.325340.  
offset 43, letra 'h', erro quadrático médio: 7.772516.  
offset 44, letra 'i', erro quadrático médio: 37.368196.

offset 66, letra '-', erro quadrático médio: 845.552643.

Neste caso, o menor erro quadrático médio corresponde ao offset 42, que pela Tabela 1 corresponde ao carácter 'g'. O carácter usado na cifra é  $(67-42)\%67 = 25$ , que corresponde ao carácter 'P'.

Para decifrar os dados soma-se a cada carácter o offset 42, módulo 67. Neste caso, o ataque é bem sucedido e recupera-se a mensagem original. O programa não precisa de escrever o erro quadrático médio para cada valor de offset, bastando escrever o resultado do menor erro como se segue:

offset com menor erro 42, letra 'g', cifrado com letra 'P', erro quadrático médio:  
1.325340.  
To Sherlock Holmes she is always the woman.

## 2.9. Ataque método 3: baseado em probabilidades para cifra de Vigenère

O ataque anterior pode ser sucessivamente aplicado a cada uma das letras da senha usada na cifra de Vigenère. Naturalmente que, para que o ataque seja bem sucedido, é necessário ter texto suficiente para que as frequências das letras calculadas para cada carácter da senha tenham uma precisão razoável. Com o texto de apenas 43 caracteres do exemplo anterior, com a senha usada como exemplo na secção 2.3, que tem 15 caracteres, cada carácter da senha foi usado para cifrar apenas cerca de 3 caracteres de dados, o que é demasiado pouco para que as frequências de ocorrência das letras se assemelhem às probabilidades de ocorrência das letras no texto em inglês.

Não se sabendo o tamanho da senha usada na cifra, tenta-se sucessivamente com tamanhos de senha crescentes, começando com 1 carácter apenas. O tamanho máximo da senha é um dado passado ao programa. Por omissão, considera-se o tamanho máximo da senha 20 caracteres. Note-se que o caso de senha de 1 carácter corresponde à cifra de César, cujo ataque já foi explicado na secção 2.8.

Para uma senha de tamanho  $n$  caracteres, vai-se sucessivamente tentar descobrir cada um dos  $n$  caracteres da senha. Para cada carácter da senha, primeiro determina-se a frequência de ocorrência de cada carácter nos dados, considerando os dados que seriam cifrados especificamente com esse carácter da senha. Escolhe-se para esse carácter da senha o carácter que corresponde ao offset que dá o menor erro quadrático médio. Quando se tiver determinado os  $n$  caracteres da senha, calcula-se o erro quadrático médio para as frequências dos caracteres descodificados no ficheiro completo com a senha escolhida. No fim, escolhe-se a senha que resulta num menor erro quadrático médio para a totalidade dos dados descodificados.

Assim, por exemplo, para uma senha de 2 caracteres, primeiro calcula-se a frequência dos caracteres de dados cifrados com o 1º caracter da senha, isto é, o 1º, 3º, 5º, 7º, etc., caracteres dos dados e calcula-se o erro quadrático médio tal como na secção 2.8, escolhendo o primeiro caracter da senha. Depois calcula-se a frequência dos caracteres de dados cifrados com o 2º caracter da senha, isto é, o 2º, 4º, 6º, 8º, etc., caracteres dos dados, e calcula-se o erro quadrático médio tal como na secção 2.8, escolhendo o segundo caracter da senha. Por fim, calcula-se as frequências dos caracteres para o ficheiro completo descodificado com a senha formada com os 2 caracteres escolhidos, e o correspondente erro quadrático médio entre estas frequências e as probabilidades tabeladas para o ficheiro completo.

Depois, passa-se a senhas de 3 caracteres, em que se avança nos caracteres dos dados de 3 em 3, e assim sucessivamente até chegar ao tamanho máximo de senha.

Considerando o texto de 4 linhas da secção 2.5, obter-se-ia:

```
min    tamanho chave 1: "m" erro 22.905399
min    tamanho chave 2: "mk" erro 19.326523
min    tamanho chave 3: "knS" erro 17.810620
min    tamanho chave 4: "mYmy" erro 15.001547
min    tamanho chave 5: "msZeq" erro 10.371095
        tamanho chave 6: "cnnksd" erro 12.902389
        tamanho chave 7: "zxUdmkn" erro 14.611716
min    tamanho chave 8: "nrs--pmz" erro 9.730406
min    tamanho chave 9: "ktSXnsdnn" erro 7.917868
        tamanho chave 10: "mzmeqmsjoq" erro 8.813991
        tamanho chave 11: "kmmmlm,mYyno" erro 11.243709
        tamanho chave 12: "itoyxdllncss" erro 8.603000
        tamanho chave 13: "tstjlkxsb mpm" erro 8.302794
        tamanho chave 14: "9yndmnqlmwkhSz" erro 8.377635
min    tamanho chave 15: "Programacao2024" erro 0.237347
        tamanho chave 16: "8hmnQwzznrs-oolz" erro 8.986137
        tamanho chave 17: " mfnaZrxpq9hV uq" erro 7.578374
        tamanho chave 18: " ,nPwspsnysfVnsxqz" erro 6.425133
        tamanho chave 19: "dsx2rp nqculitdmssv" erro 7.102089
        tamanho chave 20: "WyZermlj;qmz-gsZs9n-" erro 5.783108
To Sherlock Holmes she is always the woman. I have seldom heard him
mention her under any other name. In his eyes she eclipses and
predominates the whole of her sex. It was not that he felt any emotion
akin to love for Irene Adler. All emotions, and that one particularly,
```

O programa deve escrever no standard output, uma linha com a melhor chave de cada tamanho de chave, começando em 1, até ao tamanho máximo de chave indicado (por omissão 20 caracteres). Para cada tamanho de chave, deve escrever “min” se foi achado um novo mínimo, um tab, “tamanho chave” seguido do tamanho de chave, dois pontos, a chave formada pelos melhores caracteres para esse tamanho de chave dentro de aspas, “erro” e o erro quadrático médio para o ficheiro inteiro descodificado com essa chave. Quando chegar ao tamanho máximo de chave indicado, descodifica o ficheiro com a chave encontrada com o menor erro. Neste caso, seria a chave de tamanho 15 caracteres e o ataque era bem sucedido.

### 3. Interface de entrada/saída e parametrização do programa

É possível o programa escrever informação para depuração de erros para “stderr”. Tal informação será descartada durante as verificações automáticas. Naturalmente que se escrever muita informação em “stderr”, o seu programa ficará mais lento.

Por omissão, o programa deve ler os dados do “stdin” e escrever o resultado da execução em “stdout”. O programa pode ser parametrizado de forma a alterar este comportamento por omissão.

De seguida descrevem-se os parâmetros de linha de comando que deverão ser interpretados pelo programa a desenvolver.

### 3.1. Parametrização do programa

O programa deverá ser invocado na linha de comando da seguinte forma:

**\$ ./cifras [OPTIONS]**

'\$' é a prompt do Linux e "./" representa a directoria corrente.

**cifra** designa o nome do ficheiro executável contendo o programa desenvolvido.

**[OPTIONS]** designa a possibilidade de o programa ser invocado com diferentes opções de funcionamento

As **opções** de funcionamento são identificadas sempre como strings começadas com o caractere '-', e podem aparecer por qualquer ordem. De seguida, descrevem-se as várias opções disponíveis:

- h            mostra a ajuda para o utilizador
  - i filename   nome do ficheiro de entrada, em alternativa a stdin
  - o filename   nome do ficheiro de saída, em alternativa a stdout
  - s senha      senha a usar para cifrar/decifrar, em alternativa a "Programacao2024"
  - f            filtra o ficheiro de entrada e formata o ficheiro de saída
  - c *nn*        a operação a realizar deve ser cifrar, com o método *nn*
  - d *nn*        a operação a realizar deve ser decifrar, com o método *nn*
  - e            a operação a realizar deve ser calcular estatísticas
  - a *nn*        a operação a realizar deve ser atacar, com o método *nn*
  - n *nn*        para o ataque método 3, a dimensão máxima da chave deve ser *nn*, em vez de 20
  - w filename   nome do ficheiro de dicionário a usar
- A opção -h, quando invocada, deverá imprimir para stdout uma mensagem de ajuda de execução da linha de comandos. Para ver um exemplo deste tipo de mensagens, executar no terminal: "gedit -h". Também pode apresentar esta mensagem de ajuda no caso não haver nenhum parâmetro válido.
  - Para a opção -n, o valor por omissão deve ser 20.
  - Para a opção -w, o ficheiro de dicionário por omissão a usar deve ser /usr/share/dict/words que contém um dicionário de Inglês.

### 3.2. Exemplos de Invocação do Programa

De seguida apresentam-se alguns exemplos válidos e inválidos de invocação do programa com especificação de parâmetros em linha de comando.

**Exemplos válidos de invocação do programa:**

**Exemplo 1:** ./cifras -h

**Exemplo 2:** ./cifras -c 1

**Exemplo 3:** ./cifras -d 1 -s Programacao

**Exemplo 4:** ./cifras -c 1 -i 1line.txt -o 1line-cesar.txt

**Exemplos inválidos de invocação do programa:**

**Exemplo 1:** `./cifras -c 20`

**Exemplo 2:** `./cifras -c 1 -s`

**Exemplo 3:** `./cifras -c 3 -s 00`

No exemplo 3, a senha é demasiado curta, sendo o mínimo 4 caracteres para o método 3.

**Nota:** Pode recorrer à função de biblioteca `getopt()` para efetuar a validação dos parâmetros de entrada ou implementar uma função para realizar essa validação. Consulte a página do manual com o comando “man 3 getopt” num terminal Linux ou numa pesquisa no Google.

No caso de invocação inválida do programa, o programa deve escrever uma mensagem de erro em `stdout` e terminar.

#### 4. Processo de Submissão

Os trabalhos submetidos irão ser **compilados com as seguintes opções:**

**-Wall -O3 -g**

e correr na máquina virtual fornecida.

O projeto deve ser submetido num **ficheiro ZIP** com:

- (1) código comentado com as funcionalidades indicadas no enunciado (ficheiros `.h` e `.c`);
- (2) Makefile para gerar o executável;

Para além disso, será necessário preencher uma **ficha de auto-avaliação no GoogleForms**.

Por fim, o código deve ser estruturado de forma lógica em vários ficheiros (`*.c` e `*.h`). As funções devem ter um cabeçalho curto, mas explicativo e o código deve estar corretamente indentado e com comentários que facilitem a sua legibilidade.

É importante reforçar que certos modos de operação do programa serão avaliados de forma automática, pelo que é imperativo que o programa respeite a impressão para `stdout`/ficheiro, a leitura de `stdin`/ficheiro. A falha na execução dos mesmos poderá levar a uma penalização na nota final.

#### 5. Processo de Avaliação

O projeto será avaliado em 2 fases. Na primeira fase, o programa só precisa de ler dados do `stdin` e escrever o resultado no `stdout`, implementando-se apenas a cifra e decifra com os métodos 1 e 2. Na segunda fase deverão estar operacionais todas as funcionalidades descritas no enunciado, à excepção do método 3, que funciona como um bónus caso seja implementado. A não utilização de memória dinâmica na fase 2 funciona como uma penalização, caso não seja usada.

Fase	Implementação	Cotação [%]	Cotação [valores]
1	Método 1: Cifra de César/Substituição, cifra e decifra <code>stdin/stdout</code>	15	3
1	Método 2: Cifra de Vigenère/Polialfabética, cifra e decifra <code>stdin/stdout</code>	10	2
1	Filtragem de dados e formatação do output	5	1
1	Parâmetros da linha de comando	5	1
2	Cálculo de estatísticas	5	1

2	Ataque método 1: baseado em dicionário	20	4
2	Ataque método 2: baseado em probabilidades para cifra de César	10	2
2	Ataque método 3: baseado em probabilidades para cifra de Vigenère	10	2
2	Verificação de erros com Valgrind	10	2
2	Qualidade do código: comentários, indentação, estruturação em funções, sem warnings, nomes de variáveis apropriados, etc.	10	2
2	Método 3: Cifra de Hill/Matricial (bónus)	+5	+1
2	Não utilização de memória dinâmica (penalização)	-20	-4

Avaliação oral (realizada após a entrega do projeto):

- Grupos com **nota superior a 17** realizarão obrigatoriamente oral para defender a nota (a nota atribuída antes da oral será o ponto de partida e limite máximo da nota após a discussão oral)
- Grupos com **nota entre 8 e 17** só realização oral se o docente de laboratório ou o avaliador de projeto considerar necessário.
- Grupos com **nota inferior a 8** podem fazer oral para obter nota mínima de projeto.

## 6. Código de Honestidade Académica

Espera-se que os alunos conheçam e respeitem o Código de Honestidade Académica que rege esta disciplina e que pode ser consultado na página da cadeira. O projeto é para ser planeado e executado por grupos de dois alunos e é nessa base que será avaliado. Quaisquer associações de grupos ou outras, que eventualmente venham a ocorrer, serão obviamente interpretadas como violação do Código de Honestidade Académica e terão como consequência a anulação do projeto aos elementos envolvidos.

Lembramos igualmente que a verificação de potenciais violações a este código é feita de forma automática com recurso a sofisticados métodos de comparação de código, que envolvem não apenas a comparação direta do código, mas também da estrutura do mesmo.

## 7. Referências bibliográficas

- [1] Centro Nacional de Cibersegurança. <https://www.cncs.gov.pt/>
- [2] Cifra de Substituição. [https://pt.wikipedia.org/wiki/Cifra\\_de\\_substitui%C3%A7%C3%A3o](https://pt.wikipedia.org/wiki/Cifra_de_substitui%C3%A7%C3%A3o)
- [3] Cifra de Vigenère. [https://pt.wikipedia.org/wiki/Cifra\\_de\\_Vigen%C3%A8re](https://pt.wikipedia.org/wiki/Cifra_de_Vigen%C3%A8re)
- [4] Cifra de Hill. [https://en.wikipedia.org/wiki/Hill\\_cipher](https://en.wikipedia.org/wiki/Hill_cipher)
- [5] Letter Frequency [https://en.wikipedia.org/wiki/Letter\\_frequency](https://en.wikipedia.org/wiki/Letter_frequency)
- [6] Lei de Benford [https://pt.wikipedia.org/wiki/Lei\\_de\\_Benford](https://pt.wikipedia.org/wiki/Lei_de_Benford)

## Anexo 1: Cálculo de inverso módulo 67

Na figura seguinte é ilustrada a tabela de multiplicação de dois números módulo 67. Um dos números está nas colunas e o outro número está nas linhas. O resultado da multiplicação módulo 67 é a entrada na tabela correspondente à linha e à coluna dos números a multiplicar, sendo este resultado calculado multiplicando os dois números e obtendo o resto da divisão inteira por 67. Quando o resultado deste produto dá 1, a entrada está com fundo vermelho. Por exemplo, multiplicando  $4 \times 17$ , obtém-se 68. Como  $68 \% 67 = 1$ , o resultado é 1.

Para achar o inverso de um número, é necessário encontrar qual o número que multiplicado por ele dá 1, módulo 67. Assim, para achar o inverso de um número, procura-se na linha correspondente a esse número qual a coluna onde o produto deu 1, sendo essa coluna correspondente ao inverso pretendido. Por exemplo, o inverso de 5 módulo 67 é 27. Na linha 5, o “1” aparece na coluna correspondente ao 27, pelo que 27 é o inverso de 5. Isso pode ser verificado, pois multiplicando  $5 \times 27 = 135$ , que módulo 67 dá 1. Assim, na 1ª coluna, a amarelo, estão indicados os inversos de cada um dos números de 0 a 66, que se transcrevem de seguida:

{ 0, 1, 34, 45, 17, 27, 56, 48, 42, 15, 47, 61, 28, 31, 24, 9, 21, 4, 41, 60, 57, 16, 64, 35, 14, 59, 49, 5, 12, 37, 38, 13, 44, 65, 2, 23, 54, 29, 30, 55, 62, 18, 8, 53, 32, 3, 51, 10, 7, 26, 63, 46, 58, 43, 36, 39, 6, 20, 52, 25, 19, 11, 40, 50, 22, 33, 66 }

Como 67 é um número primo, apenas 0 não tem inverso.





## Anexo 2: Representação de um dicionário em árvore

Representando o dicionário numa árvore é possível pesquisar palavras muito rapidamente, com um tempo de pesquisa proporcional ao número de caracteres da palavra e independente da dimensão do dicionário. Naturalmente que uma árvore gasta muita memória, mas o dicionário de qualquer língua não é assim tão grande que não permita um computador actual guardá-lo sem problemas.

Em cada nível da árvore representa-se a letra de uma determinada posição dentro da palavra. No primeiro nível da árvore representa-se a primeira letra de uma palavra. No segundo nível da árvore representa-se a segunda letra da palavra e assim sucessivamente. Em cada nível, é necessário para cada letra saber se é uma palavra válida até aí, ou não, para o que basta um booleano: True, se a palavra é válida; False, se a palavra não é válida. Também é necessário saber se as letras até aí formam um prefixo de uma palavra mais longa, caso em que se tem um ponteiro para o nível seguinte da árvore.

Na Fig. 3 ilustra-se parte de uma árvore com as seguintes palavras válidas no nível 1: "a", "e". No nível 1, sabe-se que "b", "c", etc., não são palavras válidas. Sabe-se ainda que há palavras começadas por "a", pois há um ponteiro do "a" para o nível 2.

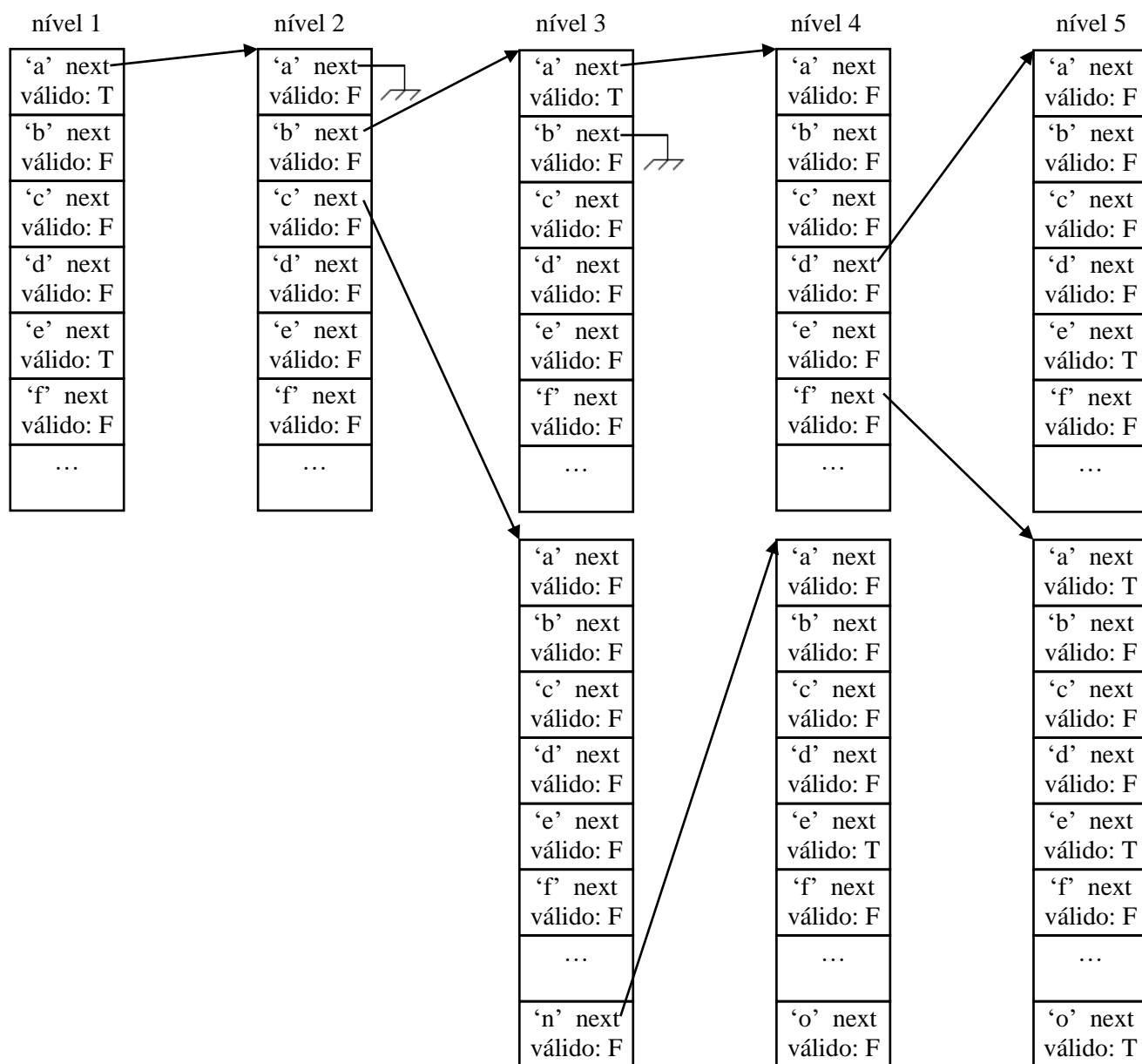


Fig. 3: Ilustração de parte de uma estrutura em árvore com palavras.

No nível 2, sabe-se que as palavras “aa”, “ab”, “ac”, etc., não são válidas e não há palavras começadas por “aa”, mas há palavras começadas por “ab” e “ac”.

No nível 3, sabe-se que “aba” é uma palavra válida, e há mais palavras começadas por “aba”; “abb” não é palavra válida e nenhuma palavra começa por “abb”; etc.

No nível 4, sabe-se que “acne” é uma palavra válida, e há palavras começadas por “abad” e “abaf”.

No nível 5, sabe-se que “abade”, “abafa”, “abafe” e “abafo” são palavras válidas.