

63-11.2 : Fondements de la programmation

Fichiers (textes)

eddy.meylan@he-arc.ch

Les fichiers

- Un fichier est une collection ordonnée d'informations homogènes
- Leur nombre n'est à priori pas connu
- Ces données sont stockées ailleurs qu'en mémoire volatile (RAM)
 - Généralement sur un «disque», et dans un répertoire
- L'unité d'information est **l'enregistrement** (quelque soit son type)

- De nos jours, les fichiers ne sont plus utilisés comme support de stockage.
 - Les bases de données existent depuis fort longtemps comme support éprouvé de stockage et manipulation de données persistantes
- Les fichiers textes sont encore utilisés, plutôt comme support d'échange de données (xml, csv *encore...* ☺) ou de configuration/paramétrage (ini, cfg)
- La manipulation de fichiers depuis un langage de programmation est besogneux.
 - Chaque langage fourni via ses librairies des services facilitant la manipulation des fichiers et de leur contenu
- Dans le cadre de cette introduction, nous nous contenterons de manipuler des fichiers textes.
 - *Il en existe d'autres...*

- Un fichier séquentiel se caractérise par le fait que deux enregistrements contigus au niveau logique occupent deux emplacements consécutifs au niveau physique.
- L'accès séquentiel impose, pour lire ou écrire un enregistrement, d'avoir lu ou écrit les enregistrements précédents
- **Un fichier texte est un fichier séquentiel dont l'unité d'enregistrement est le caractère**

- Pour vous, un fichier texte, c'est ça :

Petit Papa Noël
Quand tu descendras du ciel
Avec des jouets par milliers
N'oublie pas mon petit soulier

- Mais sur votre machine, c'est plutôt ça :

```
1011100010001010100011101101000101010010100010111110100101000101001001000100100101010010
```

- Un fichier texte est constitué de caractères, organisés en lignes
 - Cette organisation est utile pour nous et/ou pour le programme manipulant un fichier texte
 - La machine elle ne stocke que des données binaires
- Un fichier texte ne contient pas le caractère en tant que tel, mais sa représentation numérique (Ascii, Unicode)
- La marque de fin de ligne est un caractère spécial, qui diffère selon les systèmes d'exploitation.
 - CR-LF pour DOS/Windows
 - LF pour Unix
- Ce sont les routines de lecture/écriture qui les traitent

- Un texte comme nous nous le représentons...

Petit Papa Noël
Quand tu descendras du ciel
Avec des jouets par milliers
N'oublie pas mon petit soulier

- Est mis en forme... Les sauts de lignes sont aussi des caractères

Petit Papa Noël<CR><LF>Quand tu descendras du ciel<CR><LF>Avec des jouets par milliers<CR><LF>N'oublie pas mon petit soulier

- Chaque caractère, y.c les espaces et les saut de lignes, sont «encodés»

P e t i t P a p a N o e l C R L F Q u a n t ...
80 101 116 105 116 32 80 97 112 97 32 78 111 101 108 13 10 81 117 97 110 100 32 ...

- Puis convertis en binaires...

- Pour transformer en

Petit Papa Noël
Quand tu descendras du ciel
Avec des jouets par milliers
N'oublie pas mon petit soulier

il faut déterminer

101110001000101010001110110100010101001010001011110100101000101001000100100101010010

- Combien de bit représente un caractère ?
 - » Pour le transformer en un entier non signé
 - Généralement 8 (ou 16) bits
- Quel caractère représente cet entier
 - » Utilisation d'une table d'encodage
 - A un nombre (entier) correspond un caractère
 - Convention commune

- ASCII - American Standard Code for Information Interchange
- Première norme définissant en informatique un jeu de caractère standard
 - années 60
- Créé pour la langue anglaise
- Codé sur 8 bits, mais n'en utilise que 7 (!)
 - $2^7 \rightarrow 128$ caractères

Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	@	96	60	`
1	1	[START OF HEADING]	33	21	!	65	41	A	97	61	a
2	2	[START OF TEXT]	34	22	"	66	42	B	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	C	99	63	c
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	e
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	'	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	H	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49	I	105	69	i
10	A	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	B	[VERTICAL TAB]	43	2B	+	75	4B	K	107	6B	k
12	C	[FORM FEED]	44	2C	,	76	4C	L	108	6C	l
13	D	[CARRIAGE RETURN]	45	2D	-	77	4D	M	109	6D	m
14	E	[SHIFT OUT]	46	2E	.	78	4E	N	110	6E	n
15	F	[SHIFT IN]	47	2F	/	79	4F	O	111	6F	o
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	P	112	70	p
17	11	[DEVICE CONTROL 1]	49	31	1	81	51	Q	113	71	q
18	12	[DEVICE CONTROL 2]	50	32	2	82	52	R	114	72	r
19	13	[DEVICE CONTROL 3]	51	33	3	83	53	S	115	73	s
20	14	[DEVICE CONTROL 4]	52	34	4	84	54	T	116	74	t
21	15	[NEGATIVE ACKNOWLEDGE]	53	35	5	85	55	U	117	75	u
22	16	[SYNCHRONOUS IDLE]	54	36	6	86	56	V	118	76	v
23	17	[END OF TRANS. BLOCK]	55	37	7	87	57	W	119	77	w
24	18	[CANCEL]	56	38	8	88	58	X	120	78	x
25	19	[END OF MEDIUM]	57	39	9	89	59	Y	121	79	y
26	1A	[SUBSTITUTE]	58	3A	:	90	5A	Z	122	7A	z
27	1B	[ESCAPE]	59	3B	;	91	5B	[123	7B	[
28	1C	[FILE SEPARATOR]	60	3C	<	92	5C	\	124	7C	\
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D]	125	7D]
30	1E	[RECORD SEPARATOR]	62	3E	>	94	5E	^	126	7E	^
31	1F	[UNIT SEPARATOR]	63	3F	?	95	5F	_	127	7F	_
											[DEL]

- Utilisation des 8 bits
 - $2^8 \rightarrow 256$ caractères
- Décliné en plusieurs variantes qui utilise des jeux de caractères différents pour la plage 128-255 ☹
 - ISO-8859
 - ISO-8859-1
 - ISO-8859-2
 - ISO-8859-15
 - ...
 - Window-1252

ASCII de base												Utilisation du 8ième bit											
0	<NUL>	32	<SPC>	64	@	96	`	128	À	160	†	192	¿	224	#								
1	<SOH>	33	!	65	A	97	a	129	Á	161	°	193	¡	225	·								
2	<STX>	34	"	66	B	98	b	130	Â	162	¢	194	ª	226	,								
3	<ETX>	35	#	67	C	99	c	131	Ã	163	£	195	£	227	„								
4	<EOT>	36	\$	68	D	100	d	132	Ä	164	§	196	§	228	‰								
5	<ENQ>	37	%	69	E	101	e	133	Å	165	•	197	•	229	À								
6	<ACK>	38	&	70	F	102	f	134	Ä	166	¶	198	Δ	230	É								
7	<BEL>	39	'	71	G	103	g	135	Å	167	ß	199	«	231	Ê								
8	<BS>	40	(72	H	104	h	136	Ä	168	®	200	»	232	Ë								
9	<TAB>	41)	73	I	105	i	137	Ä	169	©	201	...	233	Ê								
10	<LF>	42	*	74	J	106	j	138	Ä	170	™	202	...	234	Ë								
11	<VT>	43	+	75	K	107	k	139	Ä	171	‘	203	À	235	Ì								
12	<FF>	44	,	76	L	108	l	140	Ä	172	’	204	Á	236	Í								
13	<CR>	45	-	77	M	109	m	141	Ä	173	•	205	Â	237	Î								
14	<SO>	46	.	78	N	110	n	142	Ä	174	Æ	206	Æ	238	Ò								
15	<SI>	47	/	79	O	111	o	143	Ä	175	Ø	207	Ø	239	Ó								
16	<DLE>	48	0	80	P	112	p	144	Ä	176	∞	208	-	240	•								
17	<DC1>	49	1	81	Q	113	q	145	Ä	177	±	209	-	241	•								
18	<DC2>	50	2	82	R	114	r	146	Ä	178	≤	210	-	242	•								
19	<DC3>	51	3	83	S	115	s	147	Ä	179	≥	211	-	243	•								
20	<DC4>	52	4	84	T	116	t	148	Ä	180	¥	212	-	244	•								
21	<NAK>	53	5	85	U	117	u	149	Ä	181	µ	213	-	245	•								
22	<SYN>	54	6	86	V	118	v	150	Ä	182	¶	214	÷	246	•								
23	<ETB>	55	7	87	W	119	w	151	Ä	183	Σ	215	•	247	•								
24	<CAN>	56	8	88	X	120	x	152	Ä	184	Π	216	ϕ	248	•								
25		57	9	89	Y	121	y	153	Ä	185	η	217	Υ	249	•								
26	<SUB>	58	:	90	Z	122	z	154	Ä	186	ƒ	218	/	250	•								
27	<ESC>	59	;	91	[123	{	155	Ä	187	•	219	€	251	•								
28	<FS>	60	<	92	\	124		156	Ä	188	•	220	-	252	•								
29	<GS>	61	=	93]	125	}	157	Ä	189	Ω	221	>	253	•								
30	<RS>	62	>	94	^	126	~	158	Ä	190	æ	222	fi	254	•								
31	<US>	63	?	95	_	127		159	Ä	191	•	223	fi	255	•								

- Standard informatique développé par le Consortium Unicode
 - Le standard actuel
- Objectif : Définir un jeu de caractère universel
 - Prendre en compte chaque symbole/caractère de n'importe quel système d'écriture
 - Supprimer la profusion de pages de codes spécifiques
 - Augmenter l'interopérabilité et l'échange de données dans le monde
 - Simplifier le développement des logiciels

- Attribue à chaque symbole/caractère
 - Un NOM
 - Un NOMBRE (points de code)
- Exemple avec 'a'
 - LETTRE MINUSCULE LATINE A
 - U+0061
- Exemple avec 'é'
 - LETTRE MINUSCULE LATINE E
ACCENT AIGU
 - U+00E9
- *Unicode va bien plus loin, ce n'est pas qu'une nouvelle table de caractère (mais nous n'irons pas plus loin ici)*

	000	001	002	003	004	005	006	007	ation
0	NUL 0000	DLE 0010	SP 0020	0 0030	@ 0040	P 0050	` 0060	p 0070	
1	SOH 0001	DC1 0011	! 0021	1 0031	A 0041	Q 0051	a 0061	q 0071	
2	STX 0002	DC2 0012	" 0022	2 0032	B 0042	R 0052	b 0062	r 0072	
3	ETX 0003	DC3 0013	# 0023	3 0033	C 0043	S 0053	c 0063	s 0073	
4	EOT 0004	DC4 0014	\$ 0024	4 0034	D 0044	T 0054	d 0064	t 0074	
5	ENQ 0005	NAK 0015	% 0025	5 0035	E 0045	U 0055	e 0065	u 0075	
8	TH 0008	DC 0008	.. 0008	, 0008	È 00C8	Ø 00D8	è 00E8	ø 00F8	
9	TH 0009	XXX 0009	© 00A9	1 00B9	É 00C9	Ù 00D9	é 00E9	ù 00F9	
A	TV 000A	ICU 000A	à 00AA	ó 00BA	Ê 00CA	Ú 00DA	ê 00EA	ú 00FA	

- Généralement, la table d'encodage utilisée actuellement est UTF-8
- Encode chaque point de code sur un nombre variable d'octets (de 1 à 4 octets)
 - 0 à 127 : encodage sur 1 seul octet
 - 128 à 2047 : encodage sur 2 octets
 - 2048 à 65536 : encodage sur 3 octets
 - ...
 - *UTF-8 encapsule parfaitement la table ASCII !*
- 'é' : Unicode : point de code U+00E9 (233 en décimal)
 - Donc est codé sur 2 octets, qui sont C3 A9 en UTF-8
 - En ASCII étendu, le 'é' n'est codé que sur un seul octet, puisque ASCII (même étendu) code tout sur un octet seulement

- Lorsque vous manipulez du texte, il faut TOUJOURS connaître
 - Le jeu de caractère
 - L'encodage
 - Sinon, comment interpréter les octets ????
- Un classique
 Ã© au lieu de é
 Les caractères Ã et © sont selon Unicode aux points de code C3 et A9...

- Un fichier séquentiel texte ne supporte que 3 modes :
 - Lecture (consultation d'un fichier existant)
 - Ecriture (création d'un nouveau fichier)
 - Ajout (écriture à la fin du fichier existant)
- L'accès à un fichier se fait par un flux
 - une variable représentant le flux (sur le fichier)
 - Un flux est défini en entrée ou en sortie
 - C'est à l'ouverture du fichier que le lien physique se fait
- La syntaxe est passablement différente selon les langages, mais le principe reste (généralement) le même

- Avant toute action, un fichier doit être ouvert
 - On lie un fichier physique à un flux et on définit le «sens» du flux

```
OPEN varFlux, "nom fichier", {lecture|écriture|Ajout}
```

- `nom fichier` est le nom physique du fichier
 - avec le chemin complet ou relatif
- L'ouverture en lecture d'un fichier inexistant génère une erreur
- L'ouverture en écriture d'un fichier écrase le fichier s'il existe déjà

- Après traitement, un fichier doit être fermé

```
CLOSE varFlux
```

- Entraîne la fermeture du fichier et la libération de la mémoire de la variable associée
 - Et éventuellement les verrous sur le fichier

- La lecture et l'écriture dans un fichier texte se fait par des instructions d'entrées/sorties, en y spécifiant le flux

- Lecture

```
Input (varFlux, maVar (E/S))
```

- » Lit dans le flux varFlux et affect la variable maVar

- Ecriture

```
Output (varFlux, maVar (E) )
```

- » Ecrit le flux varFlux le contenu de la variable maVar

- Les lectures/écritures peuvent se faire
 - Par caractère
 - Par ligne

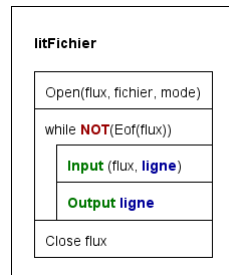
- Quelque soit le mode d'ouverture du fichier, il y a toujours un "curseur" qui pointe sur un endroit du fichier physique.
 - C'est la tête de lecture/écriture.
- A l'ouverture du fichier, elle se trouve naturellement au début du fichier, avant le premier enregistrement
 - Peut varier selon le langage
- Lorsque l'on lit un enregistrement, on lit l'enregistrement là où se trouve la tête de lecture.
 - La tête lit autant d'octets que nécessaire
 - Lorsque l'on vient de lire, la tête avance à l'enregistrement suivant

- La détection de la fin de fichier permet de savoir si la tête de lecture se trouve à la fin du fichier
 - Le curseur pointe «dans le vide»
- Cela se fait par la fonction `EOF`, qui retourne un booléen indiquant si la tête se trouve à la fin du fichier ou non.

```
Eof(varFlux) : Booléen
```

- `EOF` : *End Of File*

- Principe de lecture séquentielle de tous les enregistrements d'un fichier



- Ici la lecture se fait ligne par ligne (avec une String) et les affiche à l'écran

- Le langage Pascal met à disposition du programmeur des procédures et fonctions permettant de manipuler des fichiers
 - Doit connaître le type de données d'un fichier pour pouvoir traiter ses enregistrements
 - Dans ce cours, nous nous intéressons uniquement aux fichiers texte,
- Pour les fichiers textes, les entrées sorties se font par les procédures standards
 - Write, Writeln, Read, Readln

- Le flux est représenté par une variable
- A sa déclaration, on définit le type de fichier

```
VAR
  fLec  : Text;
...
```

- Text est un type
 - » Il est équivalent à FILE OF Char

- Assign() établit un lien entre la variable de flux et un nom de fichier physique.
 - Avec le chemin complet ou relatif
- Reset() ouvre le flux en lecture
 - Conserve les données existantes
 - Positionne la tête sur le 1^{er} enregistrement.

```
Assign(fLec, CnomFichier);
Reset(fLec);
```

- Si ce fichier n'existe pas, une erreur est levée

- Ouverture d'un fichier en écriture

```
Assign(fEcr, CnomFichier);  
Rewrite(fEcr);
```

- Deux modes possibles (pour l'écriture)
 - **Rewrite()** ouvre le fichier :
 - Si le fichier n'existe pas, il est créé
 - Si le fichier existe déjà, il est écrasé (perte du contenu)
 - **Append()** ouvre le fichier :
 - Si le fichier n'existe pas une erreur est levée
 - Si le fichier existe déjà, le contenu est ajouté au fichier

- Lire un caractère

```
Read(varFlux, unChar);
```

- Lire une ligne (et le caractère de fin de ligne)

```
Readln(varFlux, uneString);
```

- La détection de la marque de fin de ligne et de fin de fichier sont respectivement **Eoln()** et **Eof()**

- Ecrire un caractère

```
Write(varFlux, unChar);
```

- Ecrire une ligne

```
Writeln(varFlux, uneString);
```

- Et sauter à la ligne...*(écrit le caractère de fin de ligne)*
- La lecture et l'écriture consistent simplement à lire ou écrire dans un flux, ce flux pouvant être le flux standard (la console) ou le flux fichier
- Caractère(s) de fin de ligne
 - Utile en cas de formatage de texte pour insérer manuellement une marque de fin de ligne
 - Unix LF (#10) Win CRLF(#13#10) MacOS (#13)

- Détection de la fin de fichier

```
Eof(varFlux): Boolean
```

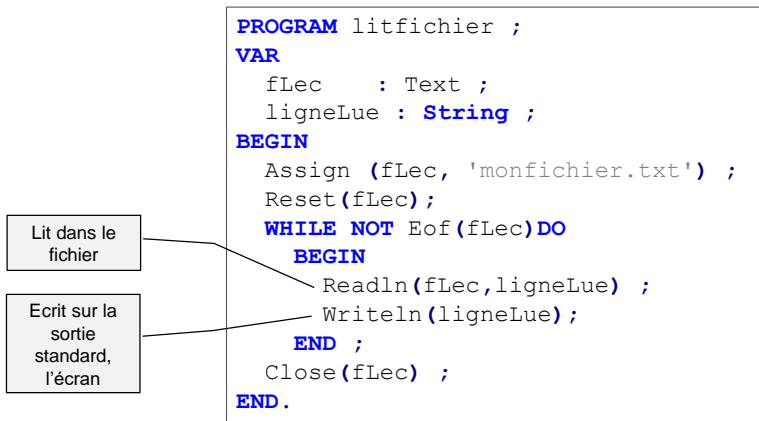
- Détection de la fin de ligne

```
Eoln(varFlux): Boolean
```

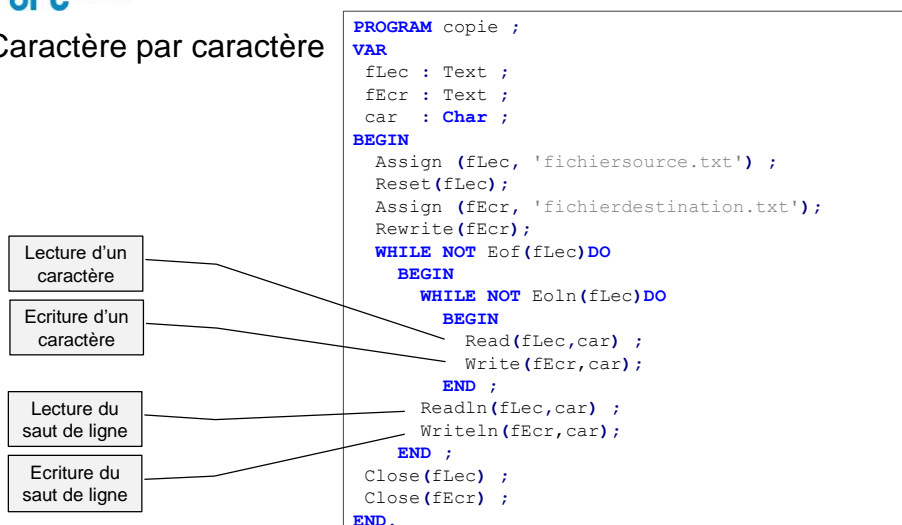
- Fermeture du flux sur le fichier

```
Close(varFlux);
```

- Lecture ligne par ligne



- Caractère par caractère



- Un fichier texte est une suite de caractères !
- On veut (parfois) stocker une structure dans un fichier texte, on va donc «formater» la structure en une suite caractères
 - La manipulation d'un fichier «formaté» nécessite de connaître sa structure (format)
- Le fichier formaté est généralement organisé en lignes :
 - Chaque ligne contient un enregistrement
 - Toutes les lignes ont le même format
- La transformation de ligne en enregistrement (et l'inverse) se fait avec des traitements de chaînes (String)
 - Les nombres (Entier et/ou Réel) sont stockés en String et nécessitent des fonctions de conversion

- Colonnes fixes
 - Chaque colonne à une longueur fixe
- Séparateur
 - Chaque colonne est séparé par un caractère connu (ici l'espace)
 - » Ne doit pas être contenu «dans» les colonnes
- CSV (peut se lire avec Excel)
 - Le séparateur est (généralement) le ;
 - La première ligne contient le nom des colonnes
- Lignes
 - Chaque ligne est un champs
 - » Tous les enregistrements doivent avoir les mêmes attributs

111Levert Albert
222Hyacinthe Facteur

111 Levert Albert
222 Hyacinthe Facteur

numero;nom;prn
111;Levert;Albert
222;Hyacinthe;Facteur

111
Levert
Albert
222
Hyacinthe
Facteur

- Pour manipuler des fichiers textes, il est (très) souvent nécessaire de «traiter» des chaines de caractères
- `nomChaine[position]`
 - Extraction d'un caractère de la chaine à une position
 - » En Pascal, commence à 1 (!)
- `Length()`
 - Retourne la longueur de la chaine
- `Pos()`
 - Retourne la position d'une chaine dans une autre
 - » Très utile pour détecter les séparateurs (et leurs positions)
- `Copy()`
 - Retourne une sous chaine extraite d'une autre
 - » Très utile pour «découper» une ligne

- `maChaine := 'Salut les lapins' ;`
 - `maChaine[2]` contient 'a'
 - `Length(maChaine)` retourne 16
 - `Pos('les',maChaine)` retourne 7
 - `Pos('zut',maChaine)` retourne 0
 - `Copy(maChaine,7,3)` retourne 'les'

- `Rename(f, 'nomfichier');`
 - Renomme le fichier (doit être «assigné» et fermé)
- `Erase(f);`
 - Efface le fichier (doit être «assigné» et fermé)
- `Getdir`
 - Retourne le répertoire de travail courant
- `Chdir()`
 - Change de répertoire de travail
- `Mkdir()`
 - Crée un répertoire
- `Rmdir()`
 - Supprime un répertoire

Supplément gratuit
pour les curieux