# Algorithmique 3 : structures de données arborescentes Projet

## Distance de Jaccard

#### **Généralités**

La distance de Jaccard  $^1$  qui sépare deux ensembles A et B est notée  $J_{\delta}(A,B)$ . Elle est définie par :

$$J_{\delta}(A,B) = egin{cases} 0, & ext{si } A = \emptyset \land B = \emptyset \, ; \ 1 - \dfrac{|A \cap B|}{|A \cup B|}, & ext{sinon}. \end{cases}$$

Elle est à valeurs dans l'intervalle réel [0; 1]. Elle mesure la dissimilarité entre les ensembles : plus sa valeur est proche de 0, plus les ensembles sont semblables; plus sa valeur s'éloigne de 0, plus les ensembles sont dissemblables.

Considérons les deux ensembles  $A = \{a, b, c, d, e\}$  et  $B = \{b, c, f\}$ . Nous avons alors  $A \cap B = \{b, c\}$  et  $A \cup B = \{a, b, c, d, e, f\}$ , et donc  $|A \cap B| = 2$  et  $|A \cup B| = 6$ . Nous obtenons donc  $J_{\delta}(A, B) = 1 - 2/6 = 1 - 1/3 = 2/3$ .

Le projet consiste à écrire un programme en C dont le but est d'afficher la distance de Jaccard qui sépare les ensembles des mots qui figurent dans toute paire de fichiers textes signifiés sur la ligne de commande.

Un mot est, par défaut, une séquence de longueur maximale de caractères n'appartenant pas à la classe **isspace**.

Pour toute paire de fichiers textes, l'affichage est de la forme leur distance de Jaccard exprimée avec quatre décimales, le caractère de tabulation horizontale '\t', l'argument associé au premier fichier texte (ou deux caractères '"' consécutifs s'il s'agit de l'entrée standard), le caractère de tabulation horizontale, l'argument associé au deuxième fichier texte (même convention), le caractère de fichiers ou fichiers ou

L'exécutable doit être nommé jdis (pour Jaccard dissimilarity).

#### **Exemples basiques**

Trois fichiers textes sont tout d'abord créés, à savoir x0.txt, x1.txt et x2.txt. Quatre exé-commande cutions sont ensuite lancées, les trois premières sur deux fichiers, la dernière sur les trois.

```
$ cat > x0.txt
a ee iii oooo uuuuu yyyyyy
$ cat > x1.txt
a zz eee rrrr ttttt yyyyyy
$ cat > x2.txt
a bb ccc dddd eeeee ffffff ggggg hhhh iii jj k
a bb ccc dddd eeeee hhhh iii jj k
a bb ccc hhhh iii jj k
a bb ccc jj k
```

MÀJ 22-03 Simplification: au moins une paire de noms de fichiers ou assimilés doivent figurer sur la ligne de commande.

<sup>1.</sup> Paul Jaccard est un professeur suisse de botanique, spécialiste de physiologie végétale, à l'origine de deux mesures utilisées en statistique : l'« indice (de Jaccard) » (ou le « coefficient de similarité », et, originellement, le « coefficient de communauté », défini comme le quotient des cardinaux de l'intersection et de l'union) et la distance dont il est ici question.

```
a jj k
$ ./jdis x0.txt x1.txt
0.8000 \rightarrow x0.txt \rightarrow x1.txt
$ ./jdis x1.txt x2.txt
0.9375 \rightarrow x1.txt \rightarrow x2.txt
$ ./jdis x2.txt x0.txt
0.8667 \rightarrow x2.txt \rightarrow x0.txt
$ ./jdis x0.txt x1.txt x2.txt
0.8000 \rightarrow x0.txt \rightarrow x1.txt
0.8667 \rightarrow x0.txt \rightarrow x2.txt
0.9375 \rightarrow x1.txt \rightarrow x2.txt
```

Les ensembles associés à x0.txt et x1.txt ont en commun les 2 mots a et yyyyyy; ils possèdent à eux deux les 6 mots de x0.txt plus les 4 mots que possède seul x1.txt, à savoir zz, eee, rrrr et ttttt. Leur distance de Jaccard vaut donc 1-2/(6+4) = 1-2/10 = 1-1/5 = 4/5 = 0.8. C'est ce qu'affichent les première et quatrième exécutions de jdis.

Tous les mots différents de x2.txt figurent sur sa première ligne; 11 au total. Les 2 mots qu'ont en commun x0.txt et x2.txt sont a et iii. Les 4 autres mots de x0.txt n'apparaissent pas dans x2.txt. Leur distance de Jaccard vaut donc  $1-2/(4+11)=13/15\simeq0,866$  7. C'est ce qu'affichent les deux dernières exécutions.

Le seul mot qu'ont en commun x1.txt et x2.txt est a. Les 5 autres mots de x0.txt n'apparaissent pas dans x2.txt. Leur distance de Jaccard vaut donc 1 - 1/(5 + 11) = 15/16 = 0.937 5. C'est ce qu'affichent les deuxième et quatrième exécutions.

# **Exemples avec options**

Des options peuvent être ajoutées au programme, fournies sur la ligne de commande selon les canons usuels, qui modifient les productions des exécutions.

La longueur maximale prise en compte pour les mots peut être fixée. Au delà de cette longueur, tout mot est coupé. La longueur est fixée à 2 dans l'exemple ci-dessous, les avertissements signifiant la coupure de mots étant redirigés vers le trou noir :

```
$ ./jdis -i2 x0.txt x1.txt x2.txt 2>/dev/null
0.6667 \rightarrow x0.txt \rightarrow x1.txt
0.7857 \rightarrow x0.txt \rightarrow x2.txt
0.8667 \rightarrow x1.txt \rightarrow x2.txt
```

Les fichiers x0.txt et x1.txt partagent un mot de plus, ee; d'où la distance 1 - 3/(6+3) = $2/3 \simeq 0.666$  7. Même chose pour x0.txt et x2.txt; d'où la distance  $1 - 3/(3 + 11) = 11/14 \simeq$ 0,785 7. Même chose encore pour x1.txt et x2.txt; d'où la distance  $1-2/(4+11)=13/15\simeq$ 0.866 7.

Une option permet de modifier les productions : au lieu d'afficher les distances, c'est le graphe d'appartenance de chaque mot appartenant à la réunion des ensembles des textes qui est affiché.

<pre>\$ ./jdis -g -i2 x0.txt x1.txt x2.txt 2&gt;/dev/null</pre>				
	$\rightarrow$ x0.txt-	$\rightarrow$ x1.txt-	yx2.txt	
a	<b>&gt;</b> x	>x	×	
bb	<del></del>	<del></del>	×	
cc	<del></del>	<del></del>	×	
dd	<del></del>	<del></del>	×	
ee	<b>&gt;</b> x	<b>&gt;</b> x	×	

ff	<del>-</del>	<del>-</del>	—)x
gg	<u> </u>	<u> </u>	$\longrightarrow$ x
hh	<u> </u>	<u> </u>	$\longrightarrow$ $\mathbf{x}$
ii	x	<u> </u>	—)x
j j	<del>-</del>	<del>-</del>	—)x
k	<u> </u>	<u> </u>	— }x
00	x	<u> </u>	<del></del>
rr	<u> </u>	—>x——	
	/	/ **	/
tt	<u> </u>	,	
tt		,	<u></u> -
	/	x	<u></u> -

Selon que le mot figurant dans le colonne de gauche apparait ou non dans le fichier dont le nom est cité sur le première ligne dans les trois autres colonnes, un « x » ou un « - » est affiché. La tabulation sert de séparateur de colonnes. Les mots sont affichés dans l'ordre lexicographique engendré par la fonction standard  $\frac{\text{strcoll}}{\text{strcoll}}^2$ . On retrouve très facilement sur ce graphe les résultats annoncés plus haut. Pour les fichiers x0.txt et x1.txt par exemple : 3 mots en commun, 9 mots au total, soit 1-3/9=2/3.

Sur des textes plus « réels » maintenant, en ce sens où y figurent des mots de langue française primée. et des symboles de ponctuation :

```
$ cat > toto0.txt
```

La maitresse de Toto lui demande :

- Toto, conjugue-moi le verbe savoir a tous les temps.
- Je sais qu'il pleut, je sais qu'il fait beau, je sais qu'il neige.
- \$ cat > toto1.txt

La maitresse demande a Toto de conjuguer le verbe manger a la premiere personne du present, du futur et du passe compose.

#### Toto dit :

- Euh... je mange, je mangerai... euh... j'ai plus faim !
- \$ cat > toto2.txt

La maitresse demande a Toto de conjuguer le verbe marcher au present et a toutes les personnes.

#### Toto lui repond :

- Je marche... tu marches... il marche... nous...

La maitresse l'interrompt et lui demande d'aller plus vite.

#### Alors Toto accelere :

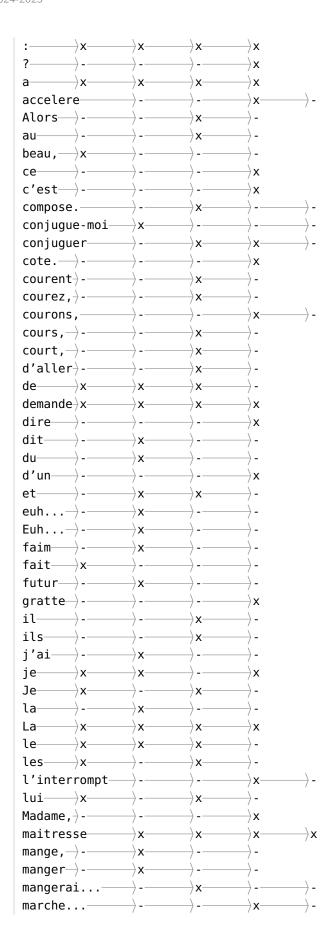
- Je cours, tu cours, il court, nous courons, vous courez, ils courent !
- \$ cat > toto3.txt

La maitresse demande a Toto :

- Toto, pourrais-tu me dire ce qu'est un oiseau migrateur ?
- Oui Madame, je peux : c'est un oiseau qui ne se gratte que d'un seul cote.
- \$ ./jdis -g toto0.txt toto1.txt toto2.txt toto3.txt

strcoll est imposée. L'option -l est supprimée. Les deux exemples qui suivent s'en trouvent modifiés en conséquence.

<sup>2.</sup> Après avoir pris soin de faire appel, en début d'exécution, à la fonction setlocale de l'en-tête standard <locale.h> avec comme premier paramètre la macroconstante LC\_COLLATE définie dans le même en-tête



$marcher  angle  ext{}$	\x
, , ,	/ /
marches	→ - → X - → -
me	→
migrateur	→ - → ×
ne	→ -
$\text{neige.} \rightarrow \text{x} \longrightarrow \text{-}$	<u> </u>
nous	
nous>	—→ <b>x</b> ———— -
oiseau $\longrightarrow$ - $\longrightarrow$ -	>>x
0ui	
passe	<del></del>
personne $\longrightarrow$ - $\longrightarrow$	
personnes. $\longrightarrow$ - $\longrightarrow$	
peux	x
pleut, $\rightarrow x \longrightarrow -$	, - , , , , , , , , , , , , , , , , , ,
plus—>-—>x—	
pourrais-tu	→ - — → - — → <b>x</b>
premiere	×
present	→ <b>x</b>
present,	×
que	→
qu'est→	→ - ×
. , , ,	/ /
· / /	→×x
qu'il—>x——>-	
repond—>>-	
sais—>x——>-	<del></del>
savoir—x————	<del></del>
se	<u></u> →>x
seul	
temps. $\rightarrow$ x $\longrightarrow$ -	<del></del>
Toto x x	$\rightarrow$ <b>x</b> $\longrightarrow$ <b>x</b>
Toto,	
tous————————————————————————————————————	<del></del>
$toutes \!\!  o \!\!\! - \!\!\! \longrightarrow \!\!\! - \!\!\!\! - \!\!\!\! \longrightarrow \!\!\!\! - \!\!\!\! - \!\!\!\!\! - \!\!\!\!\! \longrightarrow \!\!\!\! - \!\!\!\!\! - \!\!\!\!\!\! - \!\!\!\!\! \longrightarrow \!\!\!\!\! - \!\!\!\!\!\!\!\!$	
tu	→ x
un	
$verbe \longrightarrow x \longrightarrow x \longrightarrow$	
vite. $\longrightarrow$ - $\longrightarrow$ -	
vous	
<pre>\$ ./jdis toto0.txt tot</pre>	to1.txt toto2.txt toto3.txt
0,7556→toto0.txt	$\longrightarrow$ toto1.txt
0,7547→toto0.txt	
0,8043→toto0.txt	•
0,7586→toto1.txt	— ⊤toto2.txt
0,8491→toto1.txt——	— toto3.txt
0,8906→toto2.txt	
7,000 / COLOZI CAC	, 131331177

**MÀJ** 08-04 En début d'extrait des deux exemples, ajouts de l'obtention des textes.

La vérification « à l'œil et à la main » afin de retrouver les résultats à partir du graphe est plus fastidieuse. Mais le graphe, au format CSV, peut par exemple être édité dans un tableur; des colonnes « B » et « C », il vient 11 mots en commun, « x et x », pour 45 au total, « x ou x », soit  $1 - 11/45 = 34/45 \simeq 0.755 6.$ 

Des résultats sensiblement plus cohérents peuvent toutefois être obtenus en tentant de se MÀJ 10-04 rapprocher de la notion de « mot » en supprimant les caractères de ponctuation de la classe ispunct : (Voir le

tum de la question 7.) (Sébastien)

	/toto0.	txt	$\longrightarrow$ /toto1.	txt	$\longrightarrow$ /toto2.txt $-\!\!\!-$	/toto3.tx
			—)x——	$ ightarrow {\sf x}$	, -	, -
accelere		· -		√x		
,	\ <b>-</b>	×		→ <b>-</b>	/	
(	· -	<i>-</i>	x	→ <b>-</b>		
Alors—	· -	<i>-</i>	x	→ <b>-</b>		
au	· -	<i>-</i>	—) <b>x</b> ——	→ <b>-</b>		
. ′	×	<i>-</i>		→ <b>-</b>		
c	> <b>-</b>	\( \) -		$\stackrel{'}{\rightarrow}$ x		
ce	·	·		$\rightarrow$ x		
compose	\ \	√ ×		→ <u>-</u>		
compose, conjugue		}x		<i>_</i>		
conjugue		\	x	√ ×		
cote		\		->x	/	
courent)		\	x	_\		
,		/ <b>-</b>	→x			
courez—)		7 -	/			
courons		<b>} -</b>	→) <b>x</b>	<b>→</b> -		
,	<b>-</b>	<i>-</i>	→ x	<b>→</b> -		
/	<b>-</b>	<b>}</b> -	—	<del>-</del>		
- /	-	<b>-</b>	x	$\rightarrow$ x		
/	⟩x	} <b>x</b>	→ x	<b>→</b> -		
demande		} <b>x</b>	—) <b>x</b> ——	$\rightarrow$ x		
dire /	<b>-</b>	<del>}</del> -	<u> </u>	$\rightarrow$ x		
dit——	<b>-</b>	}x	<u> </u>	<del>-</del>		
du	<b>-</b>	}x	<u> </u>	<del>-</del>		
est——	<b>-</b>	<del>-</del>	<u> </u>	$\rightarrow$ x		
et———	<b>-</b>	}x	—}x——	<b>→</b> -		
euh——)	<b>-</b>	}x	<u> </u>	<del>-</del>		
Euh———	<b>-</b>	}x	<u> </u>	<del>-</del>		
faim	<b>-</b>	<b>x</b> ——		<del>-</del>		
fait——)	\x	<del>-</del>	<u> </u>	<del>-</del>		
futur——		/	<u> </u>	<del>-</del>		
gratte )	<b>-</b>	<del>}</del>	<u> </u>	$\rightarrow$ x		
il	\rangle x	<b>-</b>	x	<del>-</del>		
ils	<b>-</b>	<del>-</del>	x	<del>-</del>		
interron	npt	<b>-</b>	<u> </u>	→x	<del></del>	
j—	<b>-</b>	}x	<u> </u>	<del>-</del>		
je ,	x	x	<u> </u>	$\rightarrow$ x		
Je	\x	<del>-</del>	x	<del>-</del>		
1	<b>-</b>	ý <b>-</b> —	x	<u> </u>		
la	· -	, x	<u> </u>	<u>-</u>		
La—	x	×	x	$\stackrel{'}{ ightarrow}$ <b>x</b>		

**MÀJ** 08-04 Dans le deuxième exemple, suppression de l'incongru - g. (Jugurta)

1	\	\.,	\	
les————————————————————————————————————	<b>7</b> -	} <b>x</b>	7 <b>-</b> \	
lui————————————————————————————————————	<del>-</del>	→ X	<del>)</del> -	
Madame $\rightarrow$ -	<del>} -</del>	<del>}</del> -	×	
maitresse	→x	}x	<b>x</b> ——>	X
$mange \longrightarrow -$	→x	<del>} -</del>	<del>)</del> -	
manger $\longrightarrow$ - $\longrightarrow$	→x	<b>-</b>	<del>)</del> -	
mangerai———	<del>-</del>	}x	<del>}</del>	-
marche $\longrightarrow$ - $\longrightarrow$	<del>-</del>	}x	<del>}</del> -	
marcher $ angle$ - $$	<del>-</del>	}x	<del>}</del> -	
marches $ angle$	<del>-</del>	}x	<del>\</del> -	
me	<del>-</del>	<del>}</del> -	×	
migrateur-	<b>-</b>	<u> -                                   </u>	·	x
moi————————————————————————————————————	· -	· -	<del>)</del> -	
ne	ý <b>-</b>	· -	×	
neige— $\stackrel{/}{x}$	\( \rightarrow -	\( \rightarrow -	) <b>-</b>	
nous—>-	<i>-</i>	×	<i>→</i> _	
oiseau→-	\	\	×	
013cdd /- 0ui	\	\	×	
7	\ <b>\</b>	/ <del>-</del>	/ <b>^</b>	
passe	→ <b>X</b>	7 <b>-</b>	7 <b>-</b> \	
personne	<b>} -</b>	→ X	<i>-</i>	-
personnes	<del>\</del> -	<i>-</i>	$\rangle$ x $\longrightarrow$	-
peux————————	<del>} -</del>	<del>}</del> -	X	
pleut—x——	<del>}</del> -	<del>}</del> -	<del>)</del> -	
plus $\longrightarrow$ - $\longrightarrow$	→x	}x	<del>)</del> -	
pourrais———	<del>-</del>	-	<del></del>	X
premiere	<del>-</del>	}x	<b>}</b>	-
present  angle	→x	<b></b> x	<del>}</del> -	
qu	· -	<del>} -</del>	×	
que	<del>-</del>	<del></del>	×	
qui	<del>-</del>	· -	×	
repond $\longrightarrow$ - $\longrightarrow$	<b>-</b>	x	-	
sais———————————————————————————————————	· -	· -	<del>-</del>	
savoir—x——	ý <b>-</b> —	ý <b>-</b>	<del>)</del> -	
,	· -	\	×	
seul	<u> </u>	\	×	
temps————————————————————————————————————	\	\	) <u>-</u>	
Toto x	√x	√ x	×	
`	→	\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	\ <b>^</b>	
/	→ <b>-</b>	/ <b>-</b>	7 <b>-</b> \	
100100	/	}x	7 <b>-</b> \ <b>.</b>	
/	→ <b>-</b>	}x	<del>)</del> X	
/	/	/	<b>) X</b>	
,	→x	}x	<del>)</del> -	
vite——>-	<del>}</del> -	}x	<del>)</del> -	
vous—>-	<del>-</del>	}x	<del>)</del> -	
\$ ./jdis -p tot				t toto3.txt
$0,7955 \rightarrow /toto0.$				
0,7647 $\rightarrow$ /toto0.	txt	/toto2.	txt	
0,8444 $\rightarrow$ /toto0.		,		
0,7857 $\longrightarrow$ /toto1.	txt	/toto2.	txt	

```
0,8824→/toto1.txt------/toto3.txt
0,8833—/toto2.txt———/toto3.txt
```

La distance entre deux histoires de Toto est bien sûr égale à la tête à Toto s'il s'agit de la même histoire:

```
$ ./jdis toto0.txt ./././toto0.txt
```

Et elle est très grande, très proche de 1 donc, entre une histoire de Toto et une production 3 du Grand Homme (où « . . . » signifie le chemin donnant accès au fichier en question) :

```
$ ./jdis toto0.txt .../lesmiserables.txt
$ ./jdis -p toto0.txt .../lesmiserables.txt
```

#### Mise en œuvre

Votre programme doit être écrit en C, pouvoir être compilé avec gcc branche 13.1 (minimum) et MÀJ 22-03 supporter les options de compilation usuelles : -std=c2x, -Wall, -Wconversion, -Werror, -Wextra, Précision de -Wpedantic, -Wwrite-strings et -02.

MÀJ 08-04

Sans puis

avec -p,

les bonnes distances

et avec

Vos sources C doivent être mis en forme par uncrustify selon les fichiers de configuration fournis. Aucune ligne ne doit avoir une longueur strictement supérieure à 80.

Les types, sous-programmes et paramètres doivent être correctement nommés. Toute duplication de code est à éviter. En particulier, les éventuels caractères, chaines et nombres magiques seront nommés.

Vous éviterez l'utilisation de variables globales.

L'exécutable doit être nommé jdis.

À la base de votre travail, vous utiliserez nécessairement soit l'en-tête "hashtable.h", soit l'en-tête "bst.h", sans les modifier. La partie implantation associée à l'en-tête "bst.h" se doit de traiter les arbres binaires de recherche comme des AVL.

Vous pouvez développer et utiliser des modules additionnels, par exemple pour gérer les options ou la partie récupération de mots dans un flot texte.

Toutes les zones mémoires explicitement allouées doivent être désallouées avant la fin de l'exécution, même en cas d'erreur. Jamais deux chaines de caractères de même valeur ne doivent être allouées.

Tout message d'avertissement ou d'erreur doit être envoyé vers la sortie erreur.

Une option courte ou une option longue doit obligatoirement être supportée : « -? » ou « --help », qui affiche l'aide selon le format usuel ainsi que les éventuelles limitations de l'exécutable.

Il est possible de ne développer que des options courtes ou que des options longues.

Les options suivantes, courtes ou longues, peuvent être développées :

-q ou --qraph: pour supprimer l'affichage normal et le remplacer par l'affichage des mots appartenant à la réunion des ensembles de mots de tous les fichiers texte en signifiant pour chaque mot son appartenance aux ensembles. L'affichage se fait sur une colonne de plus que le nombre de fichiers. La première ligne affiche les noms des fichiers, chaque nom étant précédé du caractère de tabulation horizontale. Pour les autres lignes, la première colonne est réservée aux mots, les

<sup>3.</sup> Entre autres, les fichiers textes lesmiserables.txt , toto0.txt, toto1.txt, toto2.txt et toto3.txt sont MAJ 22-03 disponibles sur l'espace de cours sous UniversiTICE. Ajout de la note.

autres au codage de l'appartenance. Les colonnes sont séparées par le caractère de tabulation. Selon que le mot appartient ou non à l'ensemble des mots du fichier spécifié en haut de la colonne, un « x » ou un « - » est affiché. Les mots sont affichés dans l'ordre lexicographique;

- -i VALUE ou --initial=VALUE : pour fixer la longueur maximale des mots à VALUE. Lorsque VALUE vaut 0, la longueur n'a aucune limite. La valeur par défaut de VALUE est 0;
- -p ou --punctuation-like-space : pour donner aux symboles de ponctuation le même rôle que les caractères de la classe isspace.

Toutes les paramètres figurant sur la ligne de commande doivent être pris en compte. En par- MÀJ 22-03 ticulier, l'exécutable doit pourvoir différencier parmi les paramètres ceux qui relèvent de la catégorie Ajout du option et ceux qui n'en relèvent pas — des noms de fichiers donc, ou assimilés. Face à une option paragraphe. qu'il ne supporte pas, l'exécutable ne doit pas planter mais prévoir une fin d'exécution normale avec un message d'erreur approprié.

À l'instar de certaines commandes Linux, la possibilité de lire les mots sur l'entrée standard au lieu de les lire dans un fichier de nom donné sur la ligne de commande peut être développée. Pour le signifier explicitement à la commande, il suffit d'utiliser « - » en lieu et place d'un nom de fichier. Dans l'affichage de la ligne en-tête, « "" » doit figurer à l'emplacement du nom de fichier manguant.

Enfin, « - - » donné seul doit pouvoir spécifier que l'argument qui suit sur la ligne de commande doit être considéré comme un nom de fichier.

#### Modalités

Vous pouvez réaliser ce projet à deux, pas plus.

Vous soutiendrez le projet individuellement.

Votre projet sera compilable et exécutable sur les machines des salles de travaux pratiques du département d'informatique et sous Xubuntu.

Votre programme devra répondre aux diverses exigences signifiées dans la section « Mise en œuvre ». Il sera en partie testé avec des outils automatiques.

Votre projet devra être rendu avec tous les fichiers sources nécessaires, un fichier makefile pour produire l'exécutable (make sans paramètre ou avec les paramètres all ou jdis) et réaliser MÀJ 22-03 le nettoyage d'après production (make avec le paramètre clean), un rapport de développement. Le Précisions rapport sera fourni sous forme électronique (fichier au format pdf) qui précisera l'objet de chacun pour des modules, décrira l'implantation (dessins souhaités), commentera des exemples, explicitera les éventuelles limitations. Aucun autre fichier que ceux-là ne doit être joint; en particulier, aucun des textes ayant servi aux exemples.

Vous remettrez votre projet au plus tard le jeudi 15 mai 2025 à 18 h 47 dans une archive au MÀJ 25-04 format tar.gz de nom identifiant\_projet.tar.gz, où identifiant est l'identifiant de 8 caractères Précision de de l'un des membres du groupe.

la date.

La notation est sensible à la propreté du code, à son homogénéité, aux performances de MÀJ 22-03 l'exécutable, à la clarté des explications fournies dans le rapport et lors de la soutenance, ainsi qu'à Ajout du un passage sans encombre au révélateur valgrind.

paragraphe.

#### Hors jeux

La note attribuée au projet est 0/20 en cas :

- de demande expresse de la part de l'étudiant·e,
- d'absence de dépôt,
- d'erreur d'identifiant dans le nom de l'archive,

MÀJ 28-04 Ajout du corps de la section.

- d'archive non conforme : format, contenu (fichier manquant, tel un source, le fichier makefile ou le rapport; fichier en trop, tel un exécutable, un fichier objet, un fichier temporaire, un fichier texte qui n'est ni un source ni un fichier makefile),
  - d'erreur apparaissant à compilation,
  - de non respect de la mise en forme standard,
- pour toute recopie dans un source quelconque d'une structure normalement cachée dans la partie implantation d'un module,
- hors interface "holdall.h", d'utilisation d'une interface dont les types élémentaires « à la base de votre travail » ne sont pas const void \*;
- d'utilisation d'un forceur de type pour contrecarrer l'aspect non mutable d'une donnée (sauf pour se conformer au type de retour d'une fonction et, dans ce cas, avec une utilisation de ce genre de forceur de type immédiatement après un return),
  - pour tout plagiat (recopie totale ou partielle d'un source tiers, dont IA<sup>4</sup>),
  - pour méconnaissance du projet présenté lors de la soutenance.

#### Grille d'évaluation

En dehors des cas rédhibitoires précédemment cités, la réalisation du développement et de sa Ajout du présentation est notée sur un maximum de 20 points répartis comme suit :

MÀJ 28-04 section.

critère	points
pour 2 fichiers uniquement	6
pour un nombre de fichiers ou équivalents supérieur ou égal à 2 et éventuellement	12
inférieur à une constante d'au moins 16	
option -i	1
option -p	1
utilisation de « - » et « »	1
rapport et soutenance	5

Des pourcentages de points peuvent être retirés en fonction des critères suivants :

critère	pénalités
absence d'aide ou aide en inadéquation avec ce qui est développé	-20 %
absence de l'affichage normal ou de l'affichage associé à l'option « -g »	−50 %
plantage dans la phase de reconnaissance des options (une option non supportée doit donner lieu à un message d'erreur et à l'arrêt de l'exécutable)	-20 %
mots découpés en tranches de longueur maximale constante	−50 %
mots découpés en tranches de longueur maximale fixée par l'utilisateurice (ce qui correspondrait à une mauvaise interprétation du rôle de « -i »)	−40 %
longueur des mots limitée par une constante $C$ (ce qui correspondrait à « -i $C$ » avec $C$ non nul par défaut au lieu du « -i $0$ » par défaut)	−30 <b>%</b>
longueur des mots limitée par une constante fixée par l'utilisateurice sans que l'option « -i 0 » ne soit fonctionnelle	−20 <b>%</b>
gestion non correcte d'opérations sur les fichiers ou plantage consécutif à cette mauvaise gestion (toutes les opérations sur le fichiers doivent être testées et correctement traitées)	−30 %
erreurs signalées par valgrind, hors limitation propres de valgrind	-50 %

<sup>4.</sup> Toute expression de la forme x ou !x là où, de manière standard, est écrit en C, à l'inversion des objets de la comparaison près, x != 0, x != '\0', x != nullptr ou x == 0, x == '\0', x == nullptr sera considérée au moins comme suspecte.

désallocations non faites signalées par valgrind, même en cas d'échec messages d'erreur inadaptés (les textes des messages d'erreur doivent décrire les erreurs de la manière la plus appropriée possible et doivent être envoyés vers la sortie erreur)	-25 % -10 %
complexités non satisfaisantes à l'exécution	−25 <b>%</b>
utilisation de variables globales (hors errno)	-10 %
présence de nombres ou chaines magiques	-10 %
répétitions de parties de code (le code doit être factorisé)	−20 <b>%</b>
fonctions auxiliaires non spécifiées ou de spécification incohérente ou incomplète	−20 %
code mal présenté, malgré l'utilisation de uncrustify	-25 %

## Foire aux questions

▶ 1 Bien que vous n'avez pas encore publié le barème pour le projet, je voudrais savoir si l'on MÀJ 22-03 peut s'attendre à des points (bonus) si les deux implémentations (ABR et table de hachage) Ajout 1-1. sont réalisés. Si mes souvenirs ne me trompent pas, c'était le cas l'année dernière.?

Non pour cette année. Mais si vous tenez à faire voir votre maitrise des deux implantations. utilisez la même syntaxe que l'année dernière :

- -a: même chose que --words-processing=avl-binary-tree;
- -h : même chose que --words-processing=hash-table;
- -w TYPE ou --words-processing=TYPE : pour signifier le traitement des mots selon la structure de données spécifiée par TYPE. Les valeurs possibles pour TYPE sont explicites : avl-binary-tree et hash-table. La valeur par défaut de TYPE est hash-table.
- ▶ 2 A-t-on le droit d'utiliser des modules déjà existant comme « getopt »? Oui vous pouvez utiliser getopt.

MÀJ 28-03 Ajout 2-3.

Attention cependant à ne pas introduire de caractères ou chaines magiques. Par exemple les « g » et « graph » pour le remplacement de l'affichage normal ne doivent apparaitre qu'une fois et une seule dans tout votre code, aide comprise.

▶ 3 Je ne comprends pas vraiment ce que fait 'setlocale(...)', pourquoi on doit la mettre pour strcoll,... J'ai cherché dans l'internet et les manuels et après quelques tests les résultats sont les mêmes : la ponctuation reste (l'option -p ne fonctionne pas, et même en mettant par défaut la locale, elle ne change rien). J'ai essayé de mettre en second paramètre nullptr, "", "fr-FR" mais rien...

Réponse par deux exemples reproductibles. Dans le sous-dossier courstd/7/set\_mp/ d'Algo2 :

```
$ ./main2 < .../toto3.txt 2>/dev/null
:
?
La
Madame,
0ui
Toto
Toto,
c'est
ce
cote.
d'un
demande
dire
gratte
jе
maitresse
migrateur
ne
oiseau
peux
```

Algorithmique 3 2024-2025

```
pourrais-tu
qu'est
que
qui
se
seul
un
```

Le même résultat est obtenu après l'inclusion de l'en-tête <locale.h>, après tout appel à setlocale ou après le remplacement de strcmp par strcoll.

Si, en revanche, la variable LANG de l'environnement vaut fr\_FR.UTF-8, et que <locale.h> est inclus, et que setlocale est appelée avec les paramètres LC\_COLLATE et "", et que cet appel est un succès, et que strcoll est utilisée en lieu et place de strcmp, on obtient :

```
$ ./main2 < .../toto3.txt 2>/dev/null
:
?
а
ce
c'est
cote.
demande
dire
d'un
gratte
jе
La
Madame,
maitresse
migrateur
ne
oiseau
Oui
peux
pourrais-tu
que
qu'est
qui
se
seul
Toto
Toto,
```

# ▶ 4 Est-ce que l'aide peut être dans un fichier texte?

**MÀJ** 07-04

Non. L'aide est partie intégrante d'un source C. Il ne doit y avoir qu'un seul endroit à modifier Ajout 4-6. s'il venait l'envie de modifier le caractère ou la chaine de caractères qui donne accès à une option.

▶ 5 Lorsque l'on alloue de la mémoire pour stocker un mot, il faut une taille maximum (pour le buffer), donc les mots qui font plus qu'un certaine taille doivent être coupés même si il n'y a pas de coupage spécifié? Cela n'est pas précisé dans le projet mais il est évident qu'il y a des chaînes très longues que l'on ne pourra pas stocker complètement.

Si la longueur maximale est spécifiée, les mots dont la longueur excède cette limite doivent être coupés (et les caractères de ces mots situés au delà de la limite ignorés). Sinon, la longueur maximale n'a pas de limite; tous les mots qui peuvent être stockés en mémoire doivent l'être; si un mot ne peut être stocké entièrement, l'exécutable doit prendre fin immédiatement avec un message d'erreur adapté.

▶ 6 J'aurais une question à vous posez au sujet de l'option "--" pour le projet d'Algorithmique. Est - ce que c'est seulement l'argument qui suit qui est considéré comme un fichier ou est - ce que c'est tous les arguments qui suivent les deux tirets?

L'argument qui suit. Avec « ./jdis -- --help -- - -- toto0.txt », quatre fichiers doivent être traités : « --help », « - », « -- » et « toto0.txt ».

▶ 7 j'ai quelques interrogations concernant le projet d'algo 3. Pour l'exécutable déposé, MÀJ 10-04 est-il possible de savoir quelle est l'implémentation choisie entre les AVL et les tables de ha- Ajout 7-8. chage? De plus j'aurais voulu savoir s'il était possible d'avoir des seuils concernant le temps d'execution sur des exemples donnés de notre exécutable afin de savoir si notre implémentation est performante ou non?

ps : je pense qu'il manque la fin de la phrase en haut de la page 6 du sujet : "Des résultats sensiblement plus cohérents peuvent toutefois être obtenus en tentant de se rapprocher de la notion de « mot » en supprimant [la ponctuation.]"?

Dans la version 25.04 de l'exécutable déposé, c'est une table de hachage qui traite les mots. Pour les performances, quelques dixièmes de seconde sur « Les Trois Misérables », autrement dit les textes abeeeillmrsss.txt, lesmiserables.txt, sssrmllieeeba.txt. Mais les temps sont à mesurer sur les machines des salles de TP.

PS: fin de phrase restaurée.

▶ 8 Dans le cadre de l'optimisation de mon code pour le projet, je souhaitais savoir si je pouvais limiter le nombre de fichiers traités à 64. Cette limitation me permettrait de réduire les allocations dynamiques de mémoire et d'améliorer la complexité globale du programme.

Est-ce que cette limite est autorisée dans le cadre du projet?

64 fichiers, c'est précisément la limite de la version 25.04 de l'exécutable déposé. Toute limite est à clairement signifier dans l'aide et dans le rapport.

▶ 9 Pour lire un mot, ça passe nécessairement par un tampon. Donc le mot sera coupé si MÀJ 23-04 Ajout 9-10. sa longueur dépasse la capacité du tampon. Et c'est une erreur.

Deux fois « non ». Si une limite à la longueur des mots lus est fixée — et précisée dans l'aide comme dans le rapport — il ne s'agit pas d'une erreur : la coupure doit donner lieu à un avertissement envoyé sur la sortie erreur. Si aucune limite n'est fixée, la zone mémoire qui sert au tampon doit être réallouée pour y contenir tout le mot; il n'y a d'erreur que si une réallocation échoue par manque de place dans le tas. Pour tester la lecture et l'éventuelle mémorisation de mots longs, récupérez le source projet/glongwords.c qui figure dans l'archive algo3\_src.tar.gz sur UniversiTICE, produisez l'exécutable, dirigez sa sortie vers un fichier, exécutez votre jdis sur ce fichier.

▶ 10 Je ne comprends pas bien ce que vous voulez que l'on fasse avec LC\_ALL=C que vous donnez dans l'aide de votre exécutable.

Rien... Mais vous pouvez tester avec l'ordre traditionnel associé à strcmp ou avec celui associé à strcoll lorsque la variable LANG de l'environnement vaut fr\_FR.UTF-8 (c'est le défaut sur les machines des salles de TP). Pour de gros fichiers, ça doit être sensiblement plus rapide dans le premier cas (« LC\_ALL=C ./jdis ... », ou alors « export LC\_ALL=C » puis « ./jdis ... ») que dans le second, pour un résultat différent.

▶ 11 Bonjour, je voudrais savoir si je peux imposer le placement des options (au moins cer- MÀJ 25-04 taines) avant de mettre les noms de fichiers ou l'utilisation de l'entrée standard? La coupure Ajout 11-(et prochainement la ponctuation) nécessitent que je mette ces options avant tout ce qui est 14. manipulation de données. À part ces options cités précédemment qui doivent impérativement être en premiers, les options "-" (l'entrée standard), "--" (les noms de fichiers), "-?" (l'aide) et "-g" (le graphe) n'ont pas d'ordre imposés, donc ne contredisent pas vos exemples. Merci, cordialement.

Oui. S'agissant d'une limitation du programme, elle doit être indiquée dans l'aide et dans le rapport. De manière générale, tous les exécutables testés automatiquement le seront avec les options placées avant les noms des fichiers.

▶ 12 J'avais une petite question (ou plutôt une demande de confirmation) à propos du projet.

Il est indiqué que nous ne pouvons pas modifier les fichiers hashtable.h et bst.h. En revanche, est-il possible de modifier hashtable.c et bst.c? Et si nous reprenons le fichier hashtable.c du TP1 sans le modifier, cela aura-t-il un impact sur notre note finale? Doit-on nécessairement proposer une implémentation personnelle, étant donné que celle du TP1 a été réalisée par vous?

#### Merci d'avance pour votre réponse.

Vous pouvez modifier les parties implantation des modules hashtable et holdall si vous le souhaitez; mais ça ne peut être que pour les rendre plus performantes encore. Quant à l'impact sur la note (au projet?), cela dépendra de la pertinence des modifications.

▶ 13 Est-il possible d'avoir des indications sur la réelle possibilité de réaliser le projet avec un AVL? et si oui, de combien est-ce moins performant qu'avec une table de hachage?

Une nouvelle version (25.04.1) de l'exécutable donné en démonstration (dépôt UniversiTICE) inclus la possibilité d'un traitement des mots à l'aide d'un AVL. Sur des textes dont le cardinal de l'ensemble des mots est conséquent, le temps d'exécution semble être de 2 à 3 fois supérieur; ce qui est conforme à l'attendu.

▶ 14 Est-ce que l'ouput doit être le même que celui de votre exécutable ou pas? Genre ça passe avec des blancs au lieu de tabulations par exemple?

Les formats des sorties sont spécifiés. Vos sorties (hors l'aide) doivent être exactement les mêmes que celles de l'exécutable disponible sur UniversiTICE. Faites appel à diff.