

## GESTION DE VERSIONS DE GRANDS TEXTES

## I Différentiels par positions fixes

Dans toute cette partie, on suppose que les textes comparés ont la même taille.

**Q1** La fonction `textes_égaux` teste caractère par caractère que les listes de caractères `texte1` et `texte2` sont égales.

```
2 def textes_égaux(texte1, texte2):
3     n = len(texte1) # hypothèse de l'énoncé : len(texte1) == len(texte2)
4     for i in range(n):
5         if texte1[i] != texte2[i]:
6             return False
7     return True
```

La complexité de la fonction `textes_égaux` est en  $\mathcal{O}(n)$  où  $n$  est la longueur commune des deux textes.

**Q2** La fonction `distance` calcule ce que l'on appelle la distance de Hamming entre les listes de caractères `texte1` et `texte2`.

```
18 def distance(texte1, texte2):
19     n = len(texte1) # hypothèse de l'énoncé : len(texte1) == len(texte2)
20     compteur = 0
21     for i in range(n):
22         if texte1[i] != texte2[i]:
23             compteur += 1
24     return compteur
```

La complexité de la fonction `distance` est en  $\mathcal{O}(n)$  où  $n$  est la longueur commune des deux textes.

**Q3** On suppose dans cette question que les listes de caractères `texte1` et `texte2` peuvent être de longueurs différentes. Dans le cas où `texte1` et `texte2` sont toutes deux des listes vides, la fonction renvoie `True`.

```
36 def aucun_caractère_commun(texte1, texte2):
37     n1, n2 = len(texte1), len(texte2)
38     dico1 = {texte1[i]: True for i in range(n1)} # dictionnaire des
39     caractères de texte1
40     for i in range(n2):
41         if texte2[i] in dico1:
42             return False
43     return True
```

Complexité de la fonction `aucun_caractère_commun` :

La complexité de la fonction dépend de `len(t1)` et de `len(texte2)`.

- À l'extérieur de la boucle `for`, les instructions sont en  $\mathcal{O}(1)$ , à l'exception de la création du dictionnaire `dico1` en  $\mathcal{O}(\text{len}(\text{texte1}))$ .
- À chacun des `len(texte2)` tours de boucle, les instructions exécutées, en particulier le test `texte2[i] in dico1`, sont en  $\mathcal{O}(1)$ .  
La boucle `for` a donc une complexité en  $\mathcal{O}(\text{len}(\text{texte2}))$ .

Par somme, la complexité de la fonction `aucun_caractère_commun` est bien en  $\mathcal{O}(\text{len}(\text{texte1}) + \text{len}(\text{texte2}))$ , comme demandé par l'énoncé.

## Q4

```

93 def différentiel(texte1, texte2):
94     n = len(texte1) # hypothèse de l'énoncé : len(texte1) == len(texte2)
95     diff = [] # différentiel de texte2 vis-à-vis de texte1
96
97     i = 0 # indice de parcours de texte1
98     while i < n:
99         while i < n and texte1[i] == texte2[i]: # on avance jusqu'au
premier caractère différent, s'il existe
100             i += 1
101
102         # ici, i == n ou texte1[i] != texte2[i]
103         arg_début, arg_avant, arg_après = i, [], []
104         while i < n and texte1[i] != texte2[i]: # on avance jusqu'au
dernier caractère différent et on complète arg_avant et arg_après
105             arg_avant.append(texte1[i])
106             arg_après.append(texte2[i])
107             i += 1
108
109         # ici, i == n ou texte1[i] == texte2[i]
110         if arg_avant != [] : # si arg_avant (et donc arg_après) a été
alimenté, on crée une tranche
111             tr = tranche(arg_début, arg_avant, arg_après)
112             diff.append(tr)
113
114     return diff

```

Complexité de la fonction `differentiel` :

La complexité de la fonction dépend de  $n = \text{len}(\text{texte1}) = \text{len}(\text{texte2})$ .

- À l'extérieur de la boucle `while` principale, les instructions sont en  $\mathcal{O}(1)$ .
  - À chacun des  $n$  tours de boucle, les instructions exécutées sont en  $\mathcal{O}(1)$ .  
En effet, on réalise une première boucle `while` avec un indice  $i$  qui est incrémenté et vaut au plus  $n - 1$  et une seconde boucle `while` qui démarre avec la dernière valeur de l'indice  $i$ . Par ailleurs, toutes les instructions sont en  $\mathcal{O}(1)$ , y compris l'exécution de la fonction `tranche`.
- La boucle `while` principale a donc une complexité en  $\mathcal{O}(n)$ .

Par somme, la complexité de la fonction `differentiel` est en  $\boxed{\mathcal{O}(n) = \mathcal{O}(\text{len}(\text{texte1}))}$ , où  $n$  est la longueur commune des deux textes.

**Q5** On crée une fonction auxiliaire pour concaténer des listes en se limitant à l'utilisation de `append`.

```

196 def ajout(texte1, texte2):
197     '''renvoie texte1 + texte2.'''
198     for i in range(len(texte2)):
199         texte1.append(texte2[i])

```

Dans la fonction `applique` :

- soit `diff` est vide, et on peut renvoyer directement `texte1`,
- soit `diff` n'est pas vide, et on construit `texte2` à partir d'une liste vide en lui ajoutant ses caractères un à un (ou par slicing) en utilisant les tranches données (dont le nombre est majoré par  $\lceil \frac{n}{2} \rceil$ , car chaque tranche est séparée d'au moins 1 caractère). Ceci garantit une complexité en  $\boxed{\mathcal{O}(n)}$ , où  $n$  est la longueur commune des deux textes.

```

202 def applique(texte1, diff):
203     if diff == []:
204         return texte1[:] # copie de texte1
205     else: # diff contient des tranches à modifier
206         texte2 = [] # on construit texte2 en partant de la liste vide
207
208         arg_début = 0
209         for tr in diff: # pour chaque tranche: d'une part, on récupère
                les caractères identiques précédents (s'ils existent), d'autre part,
                on ajoute les caractères donnés par la tranche
210             ajout(texte2, texte1[arg_début: début(tr)]) # caractères non
                modifiés
211             ajout(texte2, après(tr)) # caractères donnés par la tranche
212             arg_début = fin(tr)
213
214             ajout(texte2, texte1[arg_début:]) # caractères non modifiés é
                ventuels après la dernière tranche
215         return texte2

```

**Q6** Pour répondre aux contraintes données par l'énoncé, il suffit d'intervertir les textes avant et après dans chaque tranche de diff dans une nouvelle version de diff.

```

295 def inverse(diff):
296     diff_inv = []
297     for tr in diff:
298         diff_inv.append(tranche(début(tr), après(tr), avant(tr)))
299     return diff_inv

```

La complexité de la fonction `inverse` est en  $\mathcal{O}(\text{len}(\text{diff}))$ .

**Q7** La procédure `modifie(texte_versionné, texte)` :

- détermine le différentiel de `texte` vis-à-vis du texte courant contenu dans le dictionnaire `texte_versionné` (l'énoncé dit de faire l'hypothèse que `texte` a la même taille que le texte courant),
- empile ce nouveau différentiel dans l'historique des différentiels,
- remplace le texte courant par `texte`.

```

401 def modifie(texte_versionné, texte):
402     diff = différentiel(courant(texte_versionné), texte)
403     hist = historique(texte_versionné)
404     hist.append(diff)
405     remplace_courant(texte_versionné, texte)

```

La complexité de la fonction `modifie` est en  $\mathcal{O}(\text{len}(\text{diff}))$ .

Comme  $\text{len}(\text{diff})$  a une taille au plus égale à  $\lceil \frac{n}{2} \rceil$  (comme vu en **Q5**), la complexité de la fonction `modifie` est aussi en  $\mathcal{O}(n)$ , où  $n$  est la longueur de `texte`.

La fonction `annule(texte_versionné)` :

- dépile le dernier différentiel de l'historique de `texte_versionné`,
- récupère le texte courant de `texte_versionné`,
- détermine le texte précédent en appliquant l'inverse du dernier différentiel au texte courant,
- remplace le texte courant par le texte précédent,
- renvoie le texte précédent.

L'énoncé indique que la fonction **annule** n'est jamais appelée quand il n'y a pas eu de modification. On ne demande pas de programmation défensive.

```

408 def annule(texte_versionné):
409     dernier_diff = historique(texte_versionné).pop()
410     texte2 = courant(texte_versionné)
411     texte1 = applique(texte2, inverse(dernier_diff))
412     replace_courant(texte_versionné, texte1)
413     return texte1

```

La complexité de la fonction **annule** est en  $\mathcal{O}(n)$  où  $n = \text{len}(\text{texte})$ .

## II Différentiels sur des positions variables

Q8

```

536 def poids(diff):
537     weight = 0
538     for tr in diff:
539         weight += len(avant(tr)) + len(après(tr))
540     return weight

```

La complexité de la fonction **poids** est en  $\mathcal{O}(\text{len}(\text{diff}))$ .

Q9 Soit  $0 \leq i \leq \text{len}(\text{texte1})$  et  $0 \leq j \leq \text{len}(\text{texte2})$ .

Soit  $M[i][j]$ , le nombre minimal de suppressions et d'insertions de caractères pour passer de  $\text{texte1}[:i]$  à  $\text{texte1}[:j]$ .

Une relation de récurrence est donnée par,

si  $\text{texte1}[i] = \text{texte1}[j]$ ,

- alors  $M[i+1][j+1] = M[i][j]$
- sinon  $M[i+1][j+1] = \min(M[i+1][j], M[i][j+1]) + 1$

Explication de la relation de récurrence :

Soit  $M[i][j]$  donné pour passer de  $\text{texte1}[:i]$  à  $\text{texte2}[:j]$ . Si  $\text{texte1}[i]$  et  $\text{texte2}[j]$  sont **égaux**, alors on n'ajoute ni supprime de caractère pour passer de  $\text{texte1}[:i+1]$  à  $\text{texte2}[:j+1]$ . Dans ce cas, le coût d'édition de  $M[i+1][j+1]$  est égal à  $M[i][j]$ .

Si  $\text{texte1}[i]$  et  $\text{texte2}[j]$  sont **différents**, alors on a deux possibilités pour obtenir le nombre minimal de suppressions et insertions pour passer de  $\text{texte1}[:i+1]$  à  $\text{texte2}[:j+1]$  :

- $M[i+1][j]$  est le nombre minimal de suppressions et insertions pour passer de  $\text{texte1}[:i+1]$  à  $\text{texte2}[:j]$ . En ajoutant  $\text{texte2}[j]$  à  $\text{texte2}$ , on ajoute 1 à  $M[i+1][j]$  pour obtenir le nombre minimal de suppressions et insertions pour passer de  $\text{texte1}[:i+1]$  à  $\text{texte2}[:j+1]$ . Donc,  $M[i+1][j+1]$  vaut  $M[i+1][j] + 1$ .
- $M[i][j+1]$  est le nombre minimal de suppressions et insertions pour passer de  $\text{texte1}[:i]$  à  $\text{texte2}[:j+1]$ . En supprimant  $\text{texte1}[i]$  de  $\text{texte1}$ , on ajoute 1 à  $M[i][j+1]$  pour obtenir le nombre minimal de suppressions et insertions pour passer de  $\text{texte1}[:i+1]$  à  $\text{texte2}[:j+1]$ . Donc,  $M[i+1][j+1]$  vaut  $M[i][j+1] + 1$ .

Par définition,  $M[i+1][j+1]$  est la plus petite de ces deux valeurs.

**Q10** On utilise une programmation dynamique **ascendante** afin de générer la matrice de distance d'édition,

```

578 def mini(a, b):
579     '''renvoie le plus petit des deux nombres a et b.'''
580     if a < b: return a
581     return b

584 def levenshtein(texte1, texte2):
585     n, m = len(texte1), len(texte2)
586     M = [[0 for j in range(m + 1)] for i in range(n + 1)] # structure de
                    # mémorisation
587     for i in range(n + 1):
588         for j in range(m + 1):
589             if i == 0 : M[i][j] = j
590             elif j == 0 : M[i][j] = i
591             else: # i > 0 et j > 0
592                 if texte1[i-1] == texte2[j-1] :
593                     M[i][j] = M[i-1][j-1]
594                 else:
595                     M[i][j] = mini(M[i-1][j], M[i][j-1]) + 1
596     return M

```

La complexité de la fonction levenshtein est en  $\mathcal{O}(\text{len}(\text{texte1}) \times \text{len}(\text{texte2}))$ . Elle est donc bien polynomiale en  $n = \max(\text{len}(\text{texte1}), \text{len}(\text{texte2}))$ .

**Q11** On utilise la fonction auxiliaire inv qui inverse les termes d'une liste et la fonction concaténation(texte1, texte2) qui renvoie texte1 + texte 2.

```

617 def inv(liste):
618     ''' inverse les termes d'une liste.'''
619     n = len(liste)
620     inv_liste = []
621     for i in range(n):
622         inv_liste.append(liste[n-1-i])
623     return inv_liste
624
625 def concaténation(texte1, texte2):
626     '''on concatène texte1 avec texte2.'''
627     arg = []
628     ajout(arg, texte1)
629     ajout(arg, texte2)
630     return arg

```

La fonction differentiel parcourt la matrice des distances de Levenshtein en partant de « la fin » (i.e.  $M[\text{len}(\text{texte1})][\text{len}(\text{texte2})]$ ) et en remontant jusqu'au début (i.e.  $M[0][0]$ ). Elle crée progressivement les tranches du différentiel en démarrant par la dernière tranche (d'où l'utilisation de la fonction inv pour inverser in fine les tranches du différentiel).

Pour tout  $0 < i \leq \text{len}(\text{texte1})$  et  $0 < j \leq \text{len}(\text{texte2})$ ,

- si  $\text{texte1}[i-1]$  est différent de  $\text{texte2}[j-1]$ ,
  - si  $M[i][j-1] \leq M[i-1][j]$ , alors on alimente inv\_arg\_après (i.e. arg\_après en sens inverse) avec  $\text{texte2}[j-1]$ . En cas d'égalité, on fait également ce choix, cf. FIGURE 3 de l'énoncé,
  - sinon  $M[i][j-1] > M[i-1][j]$ , alors on alimente inv\_arg\_avant (i.e. arg\_avant en sens inverse) avec  $\text{texte1}[i-1]$ .
- si  $\text{texte1}[i-1]$  est égal à  $\text{texte2}[j-1]$ , on insère une nouvelle tranche dans le différentiel à condition que inv\_arg\_avant ou inv\_arg\_après soit non vide : ceci

permet entre-autres de n'ajouter la tranche qu'une seule fois si l'égalité de caractères se poursuit. Lors de l'insertion de la tranche, on doit veiller à ce que les éléments de `arg_avant` et `arg_après` soient dans le sens de lecture (d'où l'utilisation de la fonction `inv` pour l'inversion).

Une fois arrivé au bord, on complète `arg_avant` ou `arg_après` de la tranche respectivement par `texte2[:j]` (quand `i` est nul) ou `texte1[:i]` (quand `j` est nul).

```

633 def différentiel(texte1, texte2, M):
634     n1, n2 = len(texte1), len(texte2)
635     i, j = n1, n2
636
637     diff = []
638     inv_arg_avant, inv_arg_après = [], [] # la collecte des lettres est
639     inversée car on part de la fin
640     while i != 0 or j != 0:
641         if i == 0 : # bord haut, on complète inv_arg_après par texte2[:j]
642             arg_après = concaténation(texte2[:j], inv(inv_arg_après))
643             diff.append(tranche(i, inv(inv_arg_avant), 0, arg_après))
644             j = 0
645         elif j == 0 : # bord gauche, on complète inv_arg_avant par
646             texte1[:i]
647             arg_avant = concaténation(texte1[:i], inv(inv_arg_avant))
648             diff.append(tranche(0, arg_avant, j, inv(inv_arg_après)))
649             i = 0
650         else: # i > 0 et j > 0
651             if texte1[i-1] == texte2[j-1] : # on remonte en diagonale
652                 if inv_arg_avant != [] or inv_arg_après != []: # on crée
653                 la tranche une seule fois dès qu'elle est non vide
654                     diff.append(tranche(i, inv(inv_arg_avant), j, inv(
655                     inv_arg_après)))
656                     inv_arg_avant, inv_arg_après = [], []
657                     i, j = i-1, j-1
658             else: # texte1[i-1] != texte2[j-1]
659                 if M[i-1][j] < M[i][j-1]: # on monte d'une ligne et on
660                 complète inv_arg_avant
661                     inv_arg_avant.append(texte1[i-1])
662                     i -= 1
663             else: # M[i][j-1] <= M[i-1][j], on se décale à gauche et
664             on complète inv_arg_après
665                 inv_arg_après.append(texte2[j-1])
666                 j -= 1
667
668     return inv(diff) # la collecte des tranches est inversée car on part
669     de la fin

```

Explication de la complexité :

La complexité de la fonction `différentiel` est en  $\mathcal{O}(\text{len}(\text{texte1}) + \text{len}(\text{texte2}))$ .

En effet, au sein de la boucle `while`, à chaque test, on décrémente l'indice `i` ou `j` (ou les deux). Ainsi la boucle fait exactement  $\text{len}(\text{texte1}) + \text{len}(\text{texte2})$  tours.

Concernant la fonction `inv`,

- elle permute au plus une seule fois chaque caractère de `texte1` et `texte2`, ce qui se déroule en  $\mathcal{O}(\text{len}(\text{texte1}) + \text{len}(\text{texte2}))$ ,
- elle inverse une seule fois les  $\text{len}(\text{diff})$  tranches de `diff`, avec  $\text{len}(\text{diff})$  qui est un  $\mathcal{O}(\text{len}(\text{texte1}) + \text{len}(\text{texte2}))$ .

Concernant la fonction `concaténation`, elle n'est utilisée qu'une seule fois lors de l'arrivée au bord.

- Pour  $i = 0$ , la complexité du slicing `texte2[:j]` est en  $\mathcal{O}(j)$ ; à ce texte slicé, de taille au plus égale à `len(texte2)`, est concaténé l'inverse de `inv_arg_après`, de taille au plus égale à `len(texte2)`. La complexité de l'instruction `arg_après = concaténation(texte2[:j], inv(inv_arg_après))` est donc en  $\mathcal{O}(\text{len}(\text{texte2}))$ .
- Pour  $j = 0$ , de façon symétrique, la complexité de l'instruction `arg_avant = concaténation(texte1[:i], inv(inv_arg_avant))` est en  $\mathcal{O}(\text{len}(\text{texte1}))$ .

D'où le résultat sur la complexité.

Explication des propriétés du différentiel fournit par la fonction `différentiel` :

Par construction,

- Comme `i` et `j` correspondent respectivement à `arg_début_avant` et `arg_début_après` d'une tranche, et sont décrémentés (l'un ou les deux) lors de la constitution d'une tranche, les suites des `arg_début_avant` et `arg_début_après` sont croissantes.
- On ajoute un caractère à `arg_avant` ou `arg_après` uniquement s'ils sont différents. L'indice `arg_début_avant`, respectivement, `arg_début_après` ne change pas à l'ajout d'un caractère respectivement dans `arg_après` et `arg_avant` et, dès qu'un caractère pourrait être en commun, on ajoute une tranche : ceci fait que l'on n'a aucun caractère en commun à `arg_avant` et `arg_après`.
- On ajoute systématiquement un caractère à `arg_avant` ou `arg_après` à chaque tour de boucle, sauf si le caractère est le même. Si tous les caractères sont les mêmes, `diff` reste vide. Ainsi, l'une des deux listes `arg_avant` ou `arg_après` est toujours non vide.

On a donc bien les propriétés d'un différentiel.

**Q12** L'algorithme suivant exploite le fait que les dates de début et fin des tranches sont globalement croissantes (voire strictement entre `fin_avant` d'une tranche et `début_avant` de la tranche suivante).

```

870 def conflit(diff1, diff2):
871     m1, m2 = len(diff1), len(diff2)
872     i, j = 0, 0
873     while i < m1 and j < m2:
874         tr1, tr2 = diff1[i], diff2[j]
875         if (fin_avant(tr2) >= début_avant(tr1)) and (début_avant(tr2) <=
fin_avant(tr1)): # existence d'un conflit
876             return True
877         else:
878             if début_avant(tr2) >= début_avant(tr1):
879                 i += 1
880             else:
881                 j += 1
882
883     return False # pas de conflit détecté

```



Complexité de la fonction `conflit` :

La complexité de la fonction dépend de  $m1 = \text{len}(\text{diff1})$  et de  $m2 = \text{len}(\text{diff2})$ .

- À l'extérieur de la boucle `while`, les instructions sont en  $\mathcal{O}(1)$ .
- Dans le pire des cas, on exécute  $(m1 - 1) + m2$  (ou  $(m2 - 1) + m1$  par symétrie) tours de boucle. Les instructions exécutées sont en  $\mathcal{O}(1)$ . La boucle `while` a donc une complexité en  $\mathcal{O}(m1 + m2)$ .

Par somme, la complexité de la fonction `conflit` est donc bien en

$\mathcal{O}(\text{len}(\text{diff1}) + \text{len}(\text{diff2}))$ , comme demandé par l'énoncé.

**Q13** L'application d'une tranche de `diff2` à gauche d'une tranche de `diff1` sur le texte `texte` (tel qu'exposé dans l'énoncé) ne pose aucune difficulté : cela revient à les appliquer l'une à la suite de l'autre, à partir de `texte` (i.e. `fusionne(diff1, diff2)` renvoie `diff2`).

Toute la difficulté consiste à faire l'inverse : appliquer des tranches de `diff2` à droite de tranches de `diff1` ; il faut alors décaler en conséquence les positions `début_avant(tr2)` et `début_après(tr2)` : on utilise alors la fonction auxiliaire `applique_ecart_de_position`.

```

955 def applique_ecart_de_position(tr, écart_position):
956     ''' décale les positions début_avant(tr) et début_après(tr) de tr de
        écart_position. '''
957     return tranche(début_avant(tr) + écart_position, avant(tr), dé
        but_après(tr) + écart_position, après(tr))

```

On fusionne des différentiels sur des positions variables. On a vérifié auparavant qu'ils sont sans conflit.

L'algorithme suivant commence par traiter le cas où `diff1` ou `diff2` sont vides.

Dans le cas où ce n'est pas le cas, `diff1` et `diff2` contiennent tous les deux au moins un « groupe » de tranches (éventuellement réduit à 1 tranche) que nous considérons dans l'ordre suivant :

Groupe de `diff2` (éventuellement vide) - Groupe de `diff1` - Groupe de `diff2` (éventuellement vide si le premier Groupe de `diff2` ne l'est pas) -...- Groupe de `diff1` (éventuellement vide).

Nommons `diff`, initialement vide, le résultat de la fusion de `diff1` et `diff2`.

- On commence par vérifier s'il existe un groupe de tranches de `diff2` qui précède un groupe de tranches de `diff1`. En effet, si tel est le cas, aucun écart de position n'est à effectuer, et le groupe de tranches de `diff2` est récupéré tel quel dans `diff`.
- À ce stade, on est assuré qu'il existe un groupe de tranches de `diff1`, éventuellement suivi d'un groupe de tranches de `diff2`.

On itère sur les groupes de tranches de `diff2` ce qui suit :

- Soit il n'y a plus de groupe de tranches de `diff2` et le traitement s'arrête en renvoyant `diff`.
- Sinon, on est assuré qu'il existe un groupe de tranches de `diff1`, suivi d'un groupe de tranches de `diff2`.

On parcourt le groupe de tranches de `diff1` afin de calculer l'écart de position à apporter au groupe de tranches de `diff2` situé « à sa droite » (écart « cumulé » lors d'itérations).

Puis, on parcourt le groupe de tranches de `diff2` pour lui appliquer l'écart de position avec la fonction `applique_ecart_de_position`. Les nouvelles tranches portant les modifications sont récupérées dans `diff`.



- On a alors plusieurs cas :
  - Soit il s'agissait du dernier groupe de tranches de `diff2` et le traitement s'arrête en renvoyant `diff`.
  - Soit il reste au moins un groupe de tranches de `diff2`. Cela signifie qu'il y a forcément un groupe de tranches de `diff1` qui le précède car les différentiels sont sans conflit. Ceci permet d'itérer.

```

960 def fusionne(diff1, diff2):
961     m1, m2 = len(diff1), len(diff2)
962
963     if diff1 == [] : return diff2 # on applique diff2 sur le texte
964     elif diff2 == []: return [] # on conserve les modifs de diff1
965
966     # ici, diff1 et diff2 sont non vides
967     diff = [] # création d'un nouveau différentiel
968     i, j = 0, 0
969
970     # il existe un groupe de tranche de diff2 qui précède un groupe de
    tranches de diff1
971     while j < m2 and début_avant(diff1[i]) > début_après(diff2[j]):
972         diff.append(diff2[j])
973         j += 1
974
975     # il existe un groupe de tranches de diff1, éventuellement suivi d'
    un groupe de tranches de diff2
976     écart_position = 0
977     # s'il n'y a plus de groupe de tranches de diff2, le traitement s'
    arrête en renvoyant diff
978     while j < m2:
979         # on calcule l'écart de position lié aux tranches du groupe de
    diff1
980         while i < m1 and début_après(diff1[i]) < début_avant(diff2[j]):
981             écart_position += len(après(diff1[i])) - len(avant(diff1[i])
    )
982             i += 1
983
984     # on applique l'écart de position à toutes les tranches du
    groupe de diff2 mises dans diff
985     while j < m2 :
986         tr = applique_ecart_de_position(diff2[j], écart_position)
987         diff.append(tr)
988         j += 1
989
990     return diff

```

On parcourt itérativement les tranches de `diff1` et `diff2`, en n'utilisant que des instructions en  $\mathcal{O}(1)$ , donc la complexité de la fonction `fusionne` est bien en  $\mathcal{O}(\text{len}(\text{diff1}) + \text{len}(\text{diff2}))$ .

### III Calcul de différentiels par calcul de plus courts chemins

#### Q14

```

1138 def successeurs(texte1, texte2, sommet):
1139     n1, n2 = len(texte1), len(texte2)
1140     i, j = sommet
1141     succ = []
1142     if i < n1:
1143         succ.append(((i+1, j), 1))
1144     if j < n2:
1145         succ.append(((i, j+1), 1))
1146     if i < n1 and j < n2 and texte1[i] == texte2[j]:
1147         succ.append(((i+1, j+1), 0))
1148     return succ

```

Montrons par récurrence que, pour tout  $0 \leq i \leq \text{len}(\text{texte1})$  et  $0 \leq j \leq \text{len}(\text{texte2})$ , la longueur d'un plus court chemin de  $(0,0)$  à  $(i,j)$  est égal à  $M[i][j]$  avec la pondération utilisée pour les arcs du graphe :

On fait une récurrence double sur  $0 \leq i \leq \text{len}(\text{texte1})$  et  $0 \leq j \leq \text{len}(\text{texte2})$ .

- Soit  $i = 0$  et  $j \geq 0$  : il existe un unique chemin menant de  $(0,0)$  à  $(i,j)$  : il s'agit donc du plus court chemin à partir de  $(0,0)$ . Comme sa longueur est  $j$ , elle est bien égale à  $M[0][j]$ .
- Soit  $0 \leq i < \text{len}(\text{texte1})$ . On suppose que pour tout  $0 \leq j \leq \text{len}(\text{texte2})$ , la longueur d'un plus court chemin de  $(0,0)$  à  $(i,j)$  est égale à  $M[i][j]$  (HR1).
  - Soit  $j = 0$  : il existe un unique chemin menant de  $(0,0)$  à  $(i+1,0)$  : il s'agit donc du plus court chemin à partir de  $(0,0)$ . Comme sa longueur est  $i+1$ , elle est bien égale à  $M[i+1][0]$ .
  - Soit  $0 \leq j < \text{len}(\text{texte2})$  : on suppose que la longueur d'un plus court chemin de  $(0,0)$  à  $(i+1,j)$  est égale à  $M[i+1][j]$  (HR2).  
Un plus court chemin pour aller de  $(0,0)$  à  $(i+1,j+1)$  passe par  $(i,j)$  ou  $(i+1,j)$  ou  $(i,j+1)$ .
    - Si le chemin passe par  $(i,j)$ , il existe un arc  $(i,j) \rightarrow (i+1,j+1)$ . Dans ce cas, un plus court chemin a pour poids celui d'un chemin le plus court passant par  $(i,j)$ , de poids  $M[i][j]$  (hypothèse de récurrence (HR1)), auquel on ajoute le poids de l'arc  $(i,j) \rightarrow (i+1,j+1)$ , à savoir 0. Dans ce cas, un plus court chemin a pour poids  $M[i][j]$ .
    - Si le chemin passe par  $(i+1,j)$ , le chemin a pour poids celui d'un chemin le plus court passant par  $(i+1,j)$ , de poids  $M[i+1][j]$  (hypothèse de récurrence (HR2)), auquel on ajoute le poids de l'arc  $(i+1,j) \rightarrow (i+1,j+1)$ , à savoir 1. Dans ce cas, un plus court chemin a pour poids  $M[i+1][j] + 1$ .
    - Si le chemin passe par  $(i,j+1)$ , le chemin a pour poids celui d'un chemin le plus court passant par  $(i,j+1)$ , de poids  $M[i][j+1]$  (hypothèse de récurrence (HR1)), auquel on ajoute le poids de l'arc  $(i,j+1) \rightarrow (i+1,j+1)$ , à savoir 1. Dans ce cas, un plus court chemin a pour poids  $M[i][j+1] + 1$ .

Notons que dans le premier cas, on ne supprime, ni n'ajoute de caractère, c-à-d.  $\text{texte1}[i]$  est égal à  $\text{texte2}[j]$  et dans les 2 derniers cas, on ajoute un caractère à  $\text{texte1}$  ou (exclusif) à  $\text{texte2}$ , donc  $\text{texte1}[i]$  et  $\text{texte2}[j]$  sont différents. On retrouve bien que le poids d'un plus court chemin allant de  $(0,0)$  à  $(i+1,j+1)$  est  $M[i+1][j+1]$  (par définition des éléments de la matrice de Levenshtein, cf. question 9). Ceci achève la récurrence.

**Q15** L'algorithme de Dijkstra permet de calculer le plus court chemin du sommet d'entrée  $(0, 0)$  à tout sommet  $(i, j)$  du graphe, donc avec  $0 \leq i \leq \text{len}(\text{texte1})$  et  $0 \leq j \leq \text{len}(\text{texte2})$ . Or, nous venons de prouver qu'il s'agit aussi de la distance d'édition entre  $\text{texte1}[i]$  et  $\text{texte2}[j]$ .

En particulier, lorsque l'algorithme s'arrête, il est arrivé au sommet de sortie, de coordonnées  $(\text{len}(\text{texte1}), \text{len}(\text{texte2}))$ . Il donne alors la distance d'un chemin le plus court entre le sommet d'entrée et le sommet de sortie, c.à.d. la distance d'édition entre  $\text{texte1}[\text{len}(\text{texte1})] = \text{texte1}$  et  $\text{texte2}[\text{len}(\text{texte2})] = \text{texte2}$ .

D'où le résultat voulu.

Le dictionnaire renvoyé, `dist_final`, a pour clés les sommets extraits par la file de priorité utilisée par l'algorithme de Dijkstra (ces clés sont stockées dans le dictionnaire `vue`, alimenté à chaque extraction d'un sommet de la file de priorité). Pour chacun de ses sommets, il dispose de la plus courte distance du sommet considéré au sommet de départ, et donc par ricochet, la distance d'édition entre  $\text{texte1}[i]$  et  $\text{texte2}[j]$  (où  $(i, j)$  sont les coordonnées du sommet considéré).

**Q16** Soit  $|S|$  le nombre de sommets du graphe et  $|A|$  son nombre d'arêtes.

La complexité de l'algorithme de Dijkstra avec la file de priorité donnée est un  $\mathcal{O}((|S| + |A|) \log |S|)$ .

Comme  $|A|$  peut être majoré ici par  $3|S|$ , la complexité est en  $\mathcal{O}(|S| \log |S|)$ .

Comme  $|S| = \text{len}(\text{texte1}) \times \text{len}(\text{texte2})$ , la complexité de l'algorithme de Dijkstra est ici en  $\mathcal{O}((\text{len}(\text{texte1}) \times \text{len}(\text{texte2})) \log(\text{len}(\text{texte1}) \times \text{len}(\text{texte2})))$ .

L'algorithme de Dijkstra, avec une structure de file de complexité quasiment optimale, comme celle proposée par l'énoncé, dispose d'une complexité plus importante que celle de la programmation dynamique de la partie II pour le calcul de la distance d'édition.

Par contre, dans le cas où l'algorithme de Dijkstra ne parcourt pas toute la matrice (par exemple, dans le cas d'un long préfixe commun, tel que `aaaaa` et `aaaab`), il peut être plus intéressant à utiliser que de la programmation dynamique dont l'algorithme parcourt toute la matrice de Levenshtein.

**Q17** On peut utiliser pour l'heuristique `h(texte1, texte2, sommet)` la fonction suivante :

Soient  $n_1 = \text{len}(\text{texte1})$  et  $n_2 = \text{len}(\text{texte2})$ .

Soit  $(i, j)$  les coordonnées du sommet `sommet`.

`h(texte1, texte2, sommet)` renvoie  $\max(n_1 - i, n_2 - j) - \min(n_1 - i, n_2 - j)$ .

La complexité de la fonction `h` est trivialement en  $\mathcal{O}(1)$ .

Montrons qu'il s'agit d'une heuristique admissible :

On considère un nouveau graphe dont toutes les diagonales  $(i, j) \rightarrow (i + 1, j + 1)$  existent (on a relâché une contrainte de l'énoncé). Un chemin du sommet  $(i, j)$  à la sortie sera donc de longueur inférieure à la longueur pondérée d'un plus court chemin du graphe initial.

En effet, le chemin pourra éventuellement emprunter une diagonale, là où il ne le peut pas dans le graphe initial.

Dans ce nouveau graphe, les longueurs de plus court chemins s'obtiennent en considérant « le bord » le plus proche de  $(i, j)$ ,

- s'il s'agit du bord de droite (i.e.  $n_2 - j < n_1 - i$ ), la longueur du plus court chemin est  $(n_1 - i) - (n_2 - j)$ ,
- sinon, il s'agit du bord du bas (i.e.  $n_1 - i \leq n_2 - j$ ), la longueur du plus court chemin est  $(n_2 - j) - (n_1 - i)$ ,

Justifions sur un exemple le gain de temps calcul en comparant les dictionnaires `dist_final` renvoyés par les deux algorithmes.

On considère `texte1 = ['A', 'B', 'C']` et `texte2 = ['B', 'X']`

avec  $n_1 = \text{len}(\text{texte1}) = 3$  et  $n_2 = \text{len}(\text{texte2}) = 2$ .

— Algorithme de Dijkstra

- Au départ, la file extrait le sommet  $(0,0)$ , de distance 0, aux sommets qui seront récupérés en fin par `dist_final`.

L'algorithme ajoute à la file les voisins de  $(0,0)$ , à savoir  $(0,1)$  et  $(1,0)$ , avec leur distance mise à jour.

- Puis, la file compare les distances de ses éléments  $(0,1)$  et  $(1,0)$  pour extraire celui de plus faible distance. Comme ces deux sommets ont la même distance 1, la file extrait le plus petit sommet selon l'ordre lexicographique, c.à.d. le sommet  $(0,1)$  et l'ajoute avec sa distance aux sommets qui seront récupérés en fin par `dist_final`.

Puis, l'algorithme ajoute à la file les voisins de  $(0,1)$ , à savoir  $(0,2)$  et  $(1,1)$ , avec leur distance mise à jour.

- Puis, la file compare les distances de ses éléments  $(1,0)$ ,  $(0,2)$  et  $(1,1)$  pour extraire celui de plus faible distance. Il s'agit de  $(1,0)$ , de distance 1. La file l'ajoute avec sa distance aux sommets qui seront récupérés en fin par `dist_final`. Puis, l'algorithme ajoute à la file les voisins de  $(1,0)$ , à savoir  $(2,0)$  et  $(2,1)$  (notons que  $(1,1)$  s'y trouve déjà), avec leur distance mise à jour.

La file contient à ce moment là :  $((2,1),1), ((2,0),2), ((1,1),2), ((0,2),2)$ .

On itère jusqu'au sommet d'arrivée. Alors, `dist_final` est égal à

$\{(0,0) : 0, (0,1) : 1, (1,0) : 1, (2,1) : 1, (0,2) : 2, (1,1) : 2, (2,0) : 2, (2,2) : 2, (3,1) : 2, (1,2) : 3, (3,0) : 3, (3,2) : 3\}$ .

— Algorithme  $A^*$

- Au départ, la file extrait le sommet  $(0,0)$ , de coût heuristique nul, et l'ajoute, avec sa distance de 0, aux sommets qui seront récupérés en fin par `dist_final`. L'algorithme ajoute à la file les voisins de  $(0,0)$ , à savoir  $(0,1)$  et  $(1,0)$ , avec leurs coûts heuristiques mise à jour.

- Pour  $(0,1)$  :

$d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 1 + \max(3-0, 2-1) - \min(3-0, 2-1)$ ,  
donc 3.

- Pour  $(1,0)$  :

$d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 1 + \max(3-1, 2-0) - \min(3-1, 2-0)$ ,  
donc 1.

- Puis, la file compare les coûts heuristiques de ses éléments  $(1,0)$ ,  $(0,1)$  pour extraire celui de plus faible coût. Elle extrait le sommet  $(1,0)$ , de coût 1, et l'ajoute, avec sa distance de 1, aux sommets qui seront récupérés en fin par `dist_final`.

Puis, l'algorithme ajoute à la file les voisins de  $(1,0)$ , à savoir  $(1,1)$ ,  $(2,0)$  et  $(2,1)$ , avec leurs coûts heuristiques.

- Pour  $(1,1)$  :

$d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 2 + \max(3-1, 2-1) - \min(3-1, 2-1)$ ,  
donc 3.

- Pour  $(2,0)$  :

$d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 2 + \max(3-2, 2-0) - \min(3-2, 2-0)$ ,  
donc 3.

- Pour (2, 1) :  
 $d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 1 + \max(3-2, 2-1) - \min(3-2, 2-1)$ ,  
 donc 1.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (0, 1), (1, 1), (2, 0) et (2, 1). Elle extrait le sommet (2, 1), de coût 1, et l'ajoute, avec sa distance de 1, aux sommets qui seront récupérés en fin par `dist_final`.

Puis, l'algorithme ajoute à la file les voisins de (2, 1), à savoir (2, 2) et (3, 1), avec leurs coûts heuristiques.

- Pour (2, 2) :  
 $d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 2 + \max(3-2, 2-2) - \min(3-2, 2-2)$ ,  
 donc 3.
- Pour (3, 1) :  
 $d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 2 + \max(3-3, 2-1) - \min(3-3, 2-1)$ ,  
 donc 3.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (0, 1), (1, 1), (2, 0), (2, 2) et (3, 1). Comme tous les sommets ont le même coût 3, la file extrait le plus petit sommet selon l'ordre lexicographique, càd. le sommet (0, 1) et l'ajoute, avec sa distance de 1, aux sommets qui seront récupérés en fin par `dist_final`.

Puis, l'algorithme ajoute à la file les voisins de (0, 1), à savoir (0, 2) (et pas (1, 1) car il y est déjà), avec son coût heuristique.

- Pour (0, 2) :  
 $d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 2 + \max(3-0, 2-2) - \min(3-0, 2-2)$ ,  
 donc 5.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (1, 1), (2, 0), (2, 2), (3, 1), (0, 2). Comme, en dehors de (0, 2), les sommets ont le même coût 3, la file extrait le plus petit sommet selon l'ordre lexicographique, càd. le sommet (1, 1) et l'ajoute, avec sa distance de 2, aux sommets qui seront récupérés en fin par `dist_final`.

Puis, l'algorithme ajoute à la file les voisins de (1, 1), à savoir (1, 2) (car (2, 1) y a déjà été inséré), avec sa distance mise à jour.

- Pour (1, 2) :  
 $d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 3 + \max(3-1, 2-2) - \min(3-1, 2-2)$ ,  
 donc 5.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (2, 0), (2, 2), (3, 1), (0, 2), (1, 2). Comme, en dehors de (0, 2) et (1, 2), les sommets ont le même coût 3, la file extrait le plus petit sommet selon l'ordre lexicographique<sup>1</sup>, càd. le sommet (2, 0) et l'ajoute, avec sa distance de 2, aux sommets qui seront récupérés en fin par `dist_final`.

Puis, l'algorithme ajoute à la file les voisins de (2, 0), à savoir (3, 0) (car (2, 1) y a déjà été inséré) avec sa distance mise à jour.

- Pour (3, 0) :  
 $d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 3 + \max(3-3, 2-0) - \min(3-3, 2-0)$ ,  
 donc 5.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (2, 2), (3, 1), (0, 2), (1, 2) et (3, 0). Les sommets (2, 2) et (3, 1) ont le même coût 3, la file extrait

---

1. avec l'ordre lexicographique donné...qui n'en est pas un !

le plus petit sommet selon l'ordre lexicographique, c.à.d. le sommet (2, 2) et l'ajoute, avec sa distance de 2, aux sommets qui seront récupérés en fin par `dist_final`. Puis, l'algorithme ajoute à la file les voisins de (2, 2), à savoir (3, 2) avec sa distance mise à jour.

- Pour (3, 2) :  

$$d + h(\text{texte1}, \text{texte2}, \text{voisin}) = 3 + \max(3-3, 2-2) - \min(3-3, 2-2),$$
 donc 3.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (3, 1), (0, 2), (1, 2), (3, 2) et (3, 0). Les sommets (3, 1) et (3, 2) ont le même coût 3, la file extrait le plus petit sommet selon l'ordre lexicographique, c.à.d. le sommet (3, 1) et l'ajoute, avec sa distance de 2, aux sommets qui seront récupérés en fin par `dist_final`. Aucun autre élément n'est ajouté à la file car le sommet (3, 2) y a déjà été ajouté.
- Puis, la file compare les coûts heuristiques de ses éléments, à savoir (0, 2), (1, 2), (3, 2) et (3, 0). Elle extrait le sommet (3, 2) et l'ajoute, avec sa distance de 3, aux sommets qui seront récupérés en fin par `dist_final`.  
 Aucun autre élément n'est ajouté à la fin car le sommet (3, 2) n'a pas de voisin (sommet d'arrivée).

Ainsi, `dist_final` est égal à

$\{(0, 0) : 0, (1, 0) : 1, (2, 1) : 1, (0, 1) : 1, (1, 1) : 2, (2, 0) : 2, (2, 2) : 2, (3, 1) : 2, (3, 2) : 3\}$ .

On a un gain de temps dans l'algorithme  $A^*$ , car seuls 9 sommets ont été extraits, au lieu de 12 sommets.

*NB : En cas d'égalité, un `heapdict` fournit le dernier élément inséré, ce qui fournit un meilleur résultat, avec seulement 5 sommets.*

Ainsi, *`dist_final` est égal à  $\{(0, 0) : 0, (1, 0) : 1, (2, 1) : 1, (2, 2) : 2, (3, 2) : 3\}$ .*