

Visualisation d'arbres de grandes tailles

Rapport de PSTL

Érika Baëna
erika.baena@etu.upmc.fr

Diana Malabard
diana.malabard@etu.upmc.fr

Antoine Genitrini (encadrant)
antoine.genitrini@lip6.fr

7 mai 2014

Table des matières

1	État des lieux	2
2	Étude préliminaire	3
2.1	TikZ	3
2.1.1	500 nœuds	3
2.1.2	10000 nœuds	4
2.2	Asymptote	4
2.2.1	500 nœuds	4
2.2.2	10000 nœuds	4
2.3	NetworkX accompagné de Matplotlib	5
2.3.1	500 nœuds	5
2.3.2	10000 nœuds	5
2.4	Conclusion	5
3	Choix d'implémentation	6
3.1	Fonctionnement général	6
3.2	Plus en détails	6
3.2.1	Structure d'arbre	6
3.2.2	Fonctionnement des algorithmes de parsing	6
3.2.3	Fonctionnement de l'algorithme de calcul de coordonnées	7
3.2.4	Fonctionnement de la génération du code	7
3.3	Complexité	8
4	Conclusion	9
4.1	Bilan	9
4.2	Pour la suite	9
A	Images obtenues lors de l'étude préliminaire	10
B	Extraits de l'implémentation	16

Résumé

Des outils existent actuellement pour représenter des arbres de grande taille de façon efficace. Citons par exemple GraphViz. L'inconvénient d'un tel outil est qu'il ne prend pas en compte l'ordre des fils. Ceci pose problème lorsque l'on souhaite représenter des arbres dont l'ordre des fils est primordial : les arbres de recherche.

Ce projet consiste à fournir une alternative à Graphviz, afin de pouvoir visualiser n'importe quel type d'arbre, toujours de manière efficace, mais en conservant l'ordre des fils. Ce problème possède plusieurs problématiques. Tout d'abord, nous voulons que l'affichage d'un arbre soit faite de manière élégante. Ensuite, il faut que le calcul de la mise en page de l'arbre soit rapide.

Pour ce faire, nous avons donc dû étudier les algorithmes déjà existants pour la mise en page élégante des arbres de grande taille. Sachant ces algorithmes, nous avons conçu un algorithme permettant cette mise en page. Enfin, nous avons implémenté cet algorithme, de telle façon qu'il puisse être utilisé avec différentes sorties. Nous avons ici choisi de considérer trois sorties possibles : Tikz, pour pouvoir générer automatiquement un pdf ou intégrer le code à un document LaTeX ; Asymptote, une alternative à Tikz ; et NetworkX, pour pouvoir générer une image de l'arbre que l'on pourra ultérieurement insérer dans n'importe quel document.

Chapitre 1

État des lieux

Chapitre 2

Étude préliminaire

Le but de cette étude préliminaire est de trouver un outil adapté à la représentation d'arbres de grande taille. Étudions les performances de TikZ, Asymptote et NetworkX pour la génération d'un cas particulier d'arbres : les chaînes.

2.1 TikZ

TikZ est un package L^AT_EX permettant la création de graphiques.

On va utiliser le code Python suivant pour générer du code TikZ décrivant un arbre linéaire d'une taille passée en paramètre.

```
1 import sys
3 nbIte = int(sys.argv[1])
4 if (sys.argv[2] == "true"):
5     labels = True
6 else:
7     labels = False
8 i = 0
9
10 fileName = "testTikz%d" % (nbIte,)
11
12 if (labels):
13     fileName += ".tex"
14 else:
15     fileName += ".tex"
16
17 fichierTest = open(fileName, "w")
18
19 if (labels):
20     fichierTest.write("\\node (a%d) at (0,%d) {$%d$};\n" % (i, i, i))
21     i += 1
22     while i < nbIte:
23         fichierTest.write("\\node (a%d) at (0,%d) {$%d$};\n" % (i, i, i))
24         fichierTest.write("\\draw (a%d) — (a%d);\n" % (i-1, i))
25         i += 1
26 else:
27     i += 1
28     while i < nbIte:
29         fichierTest.write("\\draw (0,%d) — (0,%d);\n" % (i-1, i))
30         i += 1
31
32 fichierTest.close()
```

2.1.1 500 nœuds

On utilise le code précédent pour générer un arbre linéaire de taille 500. On insère le code obtenu dans un fichier L^AT_EX pour voir le résultat. Le fichier L^AT_EX compile et le résultat est aux figures A.1 et A.2.

2.1.2 10000 nœuds

On utilise maintenant le même code mais pour avoir un arbre de 10000 nœuds. Le fichier \LaTeX ne compile plus. On obtient l'erreur **dimension too large** à la ligne :

```
\node (a576) at (0,576) {$576$};
```

Si l'arbre est trop grand, essayons de réduire sa taille : on diminue l'échelle, on diminue la distance entre deux points et on supprime les labels. L'erreur persiste au même endroit. TikZ limite notre arbre à 575 nœuds à la verticale. Peut-être la limite serait-elle différente si les nœuds étaient répartis autrement.

2.2 Asymptote

Asymptote est un langage de description de dessins vectoriels. Un package permet de le compiler dans un fichier \LaTeX mais le code asymptote peut aussi être autonome.

On va utiliser le code Python suivant pour générer du code Asymptote décrivant un arbre linéaire d'une taille passée en paramètre.

```
1 import sys
3 nbIte = int(sys.argv[1])
4 if (sys.argv[2] == "true"):
5     labels = True
6 else:
7     labels = False
9 i = 0
11 fileName = "testAsymptote%d" % (nbIte,)
12 if (labels):
13     fileName += ".tex"
14 else:
15     fileName += ".tex"
17 fichierTest = open(fileName, "w")
18 fichierTest.write("\\begin{asy}\n")
19 fichierTest.write("size(20cm,20cm);\n")
21 if (labels):
22     fichierTest.write("label(\"a%d\", (0, %d), E);\n" % (i, i))
23 i += 1
25 while i < nbIte:
26     if (labels):
27         fichierTest.write("label(\"a%d\", (0, %d), E);\n" % (i, i))
28         fichierTest.write("draw((0, %d) -- (0, %d));\n" % ((i-1), i))
29         i += 1
31 fichierTest.write("\\end{asy}\n")
33 fichierTest.close()
```

2.2.1 500 nœuds

On commence doucement en générant un arbre de 500 nœuds. On insère le code obtenu dans un fichier \LaTeX comme précédemment. Le fichier compile et le résultat obtenu est celui des figures A.3 et A.4.

2.2.2 10000 nœuds

On recommence en mettant la barre à 10000 nœuds. Le fichier compile sans problème et le résultat est celui des figures

2.3 NetworkX accompagné de Matplotlib

NetworkX est une bibliothèque Python pour l'étude des graphes, conçue pour fonctionner sur des grands graphes.

Matplotlib est aussi une bibliothèque Python mais qui permet quant à elle de générer une image 2D dans différents formats de sortie possible comme par exemple un png, un pdf ou un svg.

On va utiliser le code Python suivant pour générer du code NetworkX décrivant un arbre linéaire d'une taille passée en paramètre.

```
import matplotlib.pyplot as plt
2 import networkx as nx
import sys
4
nbIte = int(sys.argv[1])
longueur_arete = float(sys.argv[2])
6 G = nx.path_graph(nbIte)
8 pos={x: (5, x*longueur_arete) for x in G.nodes()}
nx.draw(G, pos, node_size=5, with_labels=False)
10 plt.savefig("networkx_%d_nodes.png" % nbIte)
```

2.3.1 500 nœuds

De même que précédemment, on génère d'abord un arbre de 500 nœuds. On choisit le format de sortie png. Le résultat est visible aux figures

2.3.2 10000 nœuds

Passons maintenant à 10000 nœuds. Le résultat est aux figures

2.4 Conclusion

TikZ permet une représentation claire d'un arbre avec ses labels. En effet, même avec 500 nœuds, les labels sont lisibles si on zoome suffisamment. Cependant, une limite a rapidement été atteinte. TikZ serait préférable pour la représentation de petits arbres avec (ou sans!) labels.

Asymptote permet de représenter de grands arbres. Son point faible est la représentation des labels. Cependant, les labels sont compilés avec L^AT_EX_ε qui peut permettre d'avoir des labels un peu plus évolués qu'une chaîne de caractères, une fois la taille des labels maîtrisée.

NetworkX et Matplotlib permettent plusieurs formats de sortie différents. Cela pourrait être utile aux non-utilisateurs de L^AT_EX. Cependant, l'affichage des labels n'est pas non plus très optimal.

Notons tout de même que cette étude préliminaire ne prend pas en compte le temps de calcul des coordonnées, ce qui est le cœur de notre projet et qui est expliqué dans le chapitre suivant.

Chapitre 3

Choix d'implémentation

3.1 Fonctionnement général

L'utilisateur fournit un fichier et précise son type. Il peut aussi choisir d'afficher des labels sur les nœuds ainsi que le format de sortie souhaité. Par défaut, l'application génère un png sans labels. L'application fonctionne ensuite en trois étapes :

1. Parsing du fichier d'entrée selon le type indiqué pour obtenir une représentation d'arbre selon notre structure interne.
2. Calcul des coordonnées de chaque nœud.
3. Génération d'une image selon le type de sortie choisie.

```
erika@erika-K53SD:~/Documents/Cours/M1S2/PSTL/Implementation$ python3 treeDisplay.py -h
usage: treeDisplay.py [-h] [-L] [-O {tikz,asy,png,pdf,eps}] [-N NAME]
                    {str,arb,xml,dot} src

positional arguments:
  {str,arb,xml,dot}    The type of the given file.
  src                  The file which contains a tree description.

optional arguments:
  -h, --help            show this help message and exit
  -L, --labels          Print labels on the output tree, deprecated with
                        NetworkX and Asymptote
  -O {tikz,asy,png,pdf,eps}, --output {tikz,asy,png,pdf,eps}
                        The type of the output file
  -N NAME, --name NAME  The name of the output file
erika@erika-K53SD:~/Documents/Cours/M1S2/PSTL/Implementation$
```

3.2 Plus en détails

3.2.1 Structure d'arbre

```

1  Éë
2  This file is part of TreeDisplay .
3
4      TreeDisplay is free software: you can redistribute it and/or modify
5      it under the terms of the GNU General Public License as published by
6      the Free Software Foundation, either version 3 of the License, or
7      (at your option) any later version.
8
9      Foobar is distributed in the hope that it will be useful,
10     but WITHOUT ANY WARRANTY; without even the implied warranty of
11     MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
12     GNU General Public License for more details.
13
14     You should have received a copy of the GNU General Public License
15     along with Foobar. If not, see <http://www.gnu.org/licenses/>.
```

3.2.2 Fonctionnement des algorithmes de parsing

Mots bien parenthésés Le parser de mots bien parenthésés (cf. B) respecte la grammaire suivante :

```
ARBRE : '(' ID NOEUDS ')'  
NOEUDS : ARBRE NOEUDS | e  
ID : [a-zA-Z1-9]* | e
```

Dot Le parser dot parse une sous-partie du langage dot, à savoir :

```
DOT : STRICT GRAPH ID '{' SEQINST '}'  
STRICT : strict | e  
GRAPH = diagraph | graph  
SEQINST : INST SEQINST | e  
INST : REF [label = "ID"] ';' | REF LINK REF  
LINK : -- | ->  
REF : [0-9]*  
ID : [a-zA-Z1-9]* | e
```

XML Le parser XML utilise la librairie XML de Python. Il respecte la grammaire suivante :

```
XML : <?xml version="1.0"?><tree> NOEUDS </tree>  
NOEUDS : NOEUD NOEUDS | e  
NOEUD : <node type=TAG id=ID> NOEUDS </node> | <leaf type=TAG id=ID />  
TAG : "Leaf" | "BinNode"
```

3.2.3 Fonctionnement de l'algorithme de calcul de coordonnées

On décide que la distance minimale entre deux nœuds est de 1 unité.

L'ordonnée d'un nœud est triviale : c'est sa profondeur.

L'abscisse d'un nœud est un peu plus complexe et demande donc plus de réflexion. Pour un nœud donné, on commence par placer ses fils. On centre ensuite ce nœud au milieu de ses fils en faisant la moyenne de l'abscisse de ses deux fils extrêmes. Si un nœud n'a pas de fils, on le place à 1 de son frère gauche. D'un point de vue de l'architecture, on a donc besoin d'une structure qui, à profondeur p mémorise la prochaine place disponible (ou au choix la dernière place utilisée). Si un père a des fils, on le centre au milieu de ses fils. Par ce calcul, on peut se retrouver avec un nœud qui est trop proche de son frère gauche [Mettre un exemple]. Pour cela, on compare la position calculée avec la première position disponible et on prend le max. Si la position calculée n'est pas celle retenue, le père n'est plus centré au milieu de ses fils. On mémorise donc le décalage effectué pour ce père pour ensuite l'appliquer à ses sous-arbre dans un second temps.

1. Pour un nœud donné, on commence par placer ses fils.
2. On centre ensuite ce nœud au milieu de ses fils.
3. Si un nœud collisionne avec son frère gauche, on le décale vers la droite et on mémorise ce décalage car il faudra ensuite décaler ses sous-arbres. On applique ce décalage dans une seconde passe pour des raisons de complexité. En effet, si on faisait chaque décalage lorsqu'il se présentait, le décalage serait quadratique alors que dans le cas choisi, on est linéaire.

3.2.4 Fonctionnement de la génération du code

TikZ

Asymptote

Autre

3.3 Complexité

On suppose que la taille des labels est bornée. On note n le nombre de nœuds dans l'arbre. Dans ce cas :

1. Le parsing d'un fichier est en $O(n)$.
2. Le calcul des coordonnées est en $O(n)$ (on effectue 2 passes sur l'arbre).
3. La génération du fichier de sortie est en $O(n)$.

On a donc une complexité générale en $O(n)$ où n est le nombre de nœud de l'arbre.

Notons que nous avons supposé que la taille des labels était bornée. Or nous n'avons aucune prise sur la taille des labels du fichier qui nous est passé en entrée. Dans ce cas, même si on borne la taille des labels pour l'affichage, la complexité est dominée par le parsing du fichier d'entrée car on doit lire tous les caractères du fichier.

Chapitre 4

Conclusion

4.1 Bilan

L'application `treeDisplay` permet de calculer les coordonnées des nœuds d'un arbre en temps linéaire par rapport à ce nombre de nœuds.

4.2 Pour la suite

La structure choisie permet aussi de représenter des graphes. On peut donc envisager par la suite d'implémenter un algorithme de calcul de coordonnées pour les graphes. Les modules de parsing et de génération sont pleinement réutilisables.

La complexité en mémoire est actuellement égale à celle en temps. Cependant, on utilisant une table de hachage pour les sous-arbres, on peut devenir sous-linéaire. L'idée consiste à calculer des coordonnées relatives et lorsqu'un arbre comporte deux (ou plus) sous-arbres identiques. Le sous-arbre en question n'est sauvegarder en mémoire qu'une seule fois et la deuxième fois on ne sauvegarde qu'une référence. On calcule ensuite les coordonnées absolues lors de la génération.

Annexe A

Images obtenues lors de l'étude préliminaire

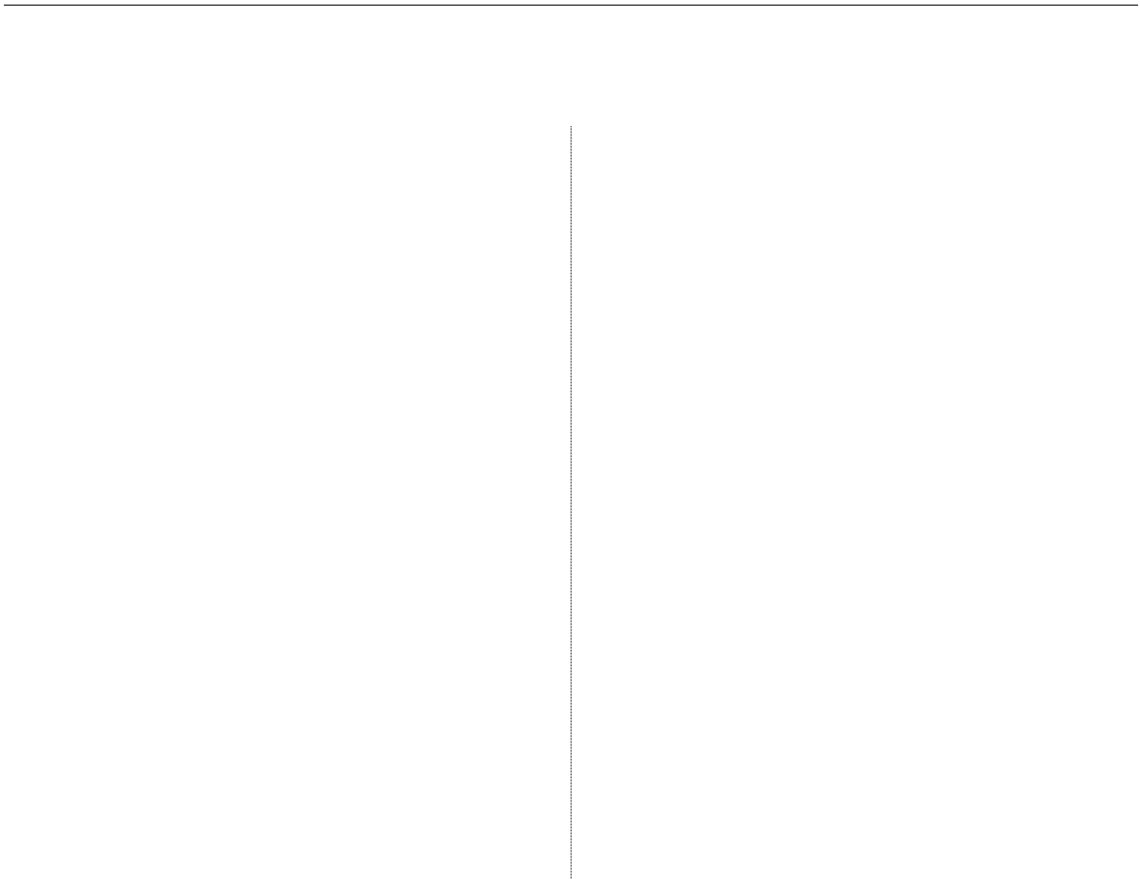


FIGURE A.1 – Arbre linéaire de taille 500 avec labels obtenu avec TikZ.

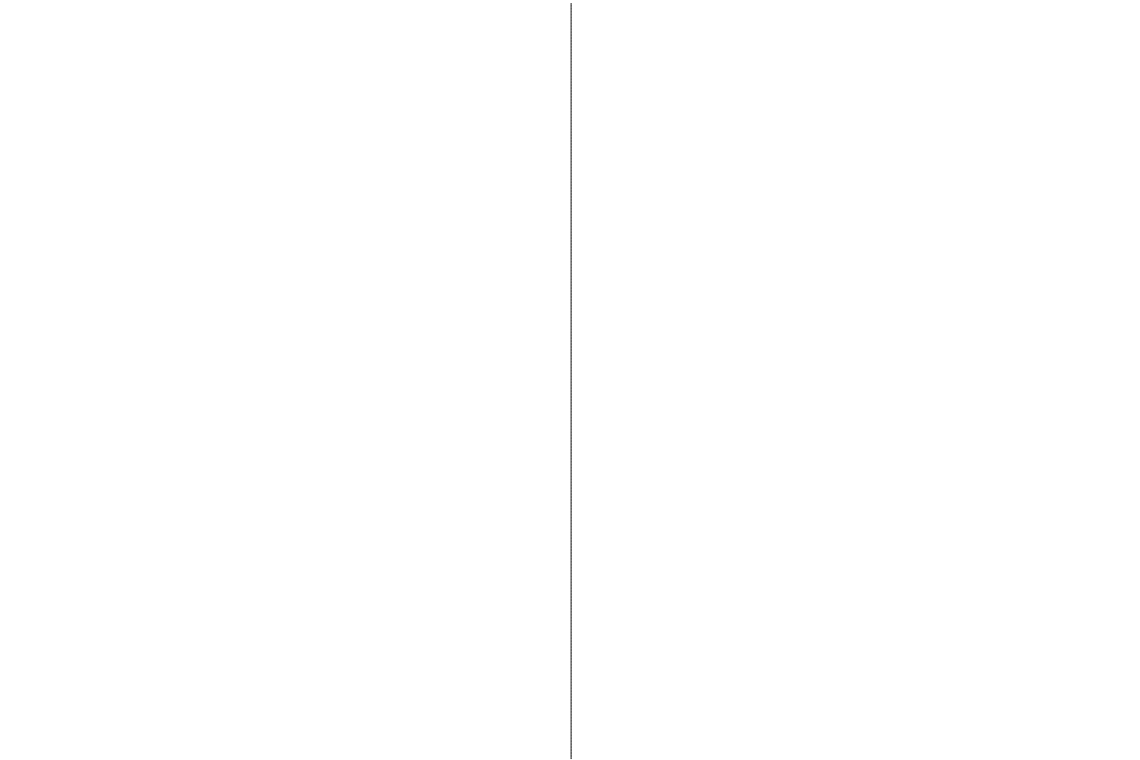


FIGURE A.2 – Arbre linéaire de taille 500 sans labels obtenu avec TikZ.

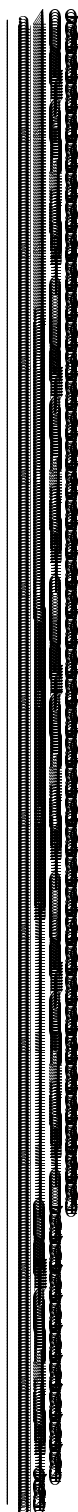


FIGURE A.3 – Arbre linéaire de taille 500 avec labels obtenu avec Asymptote.

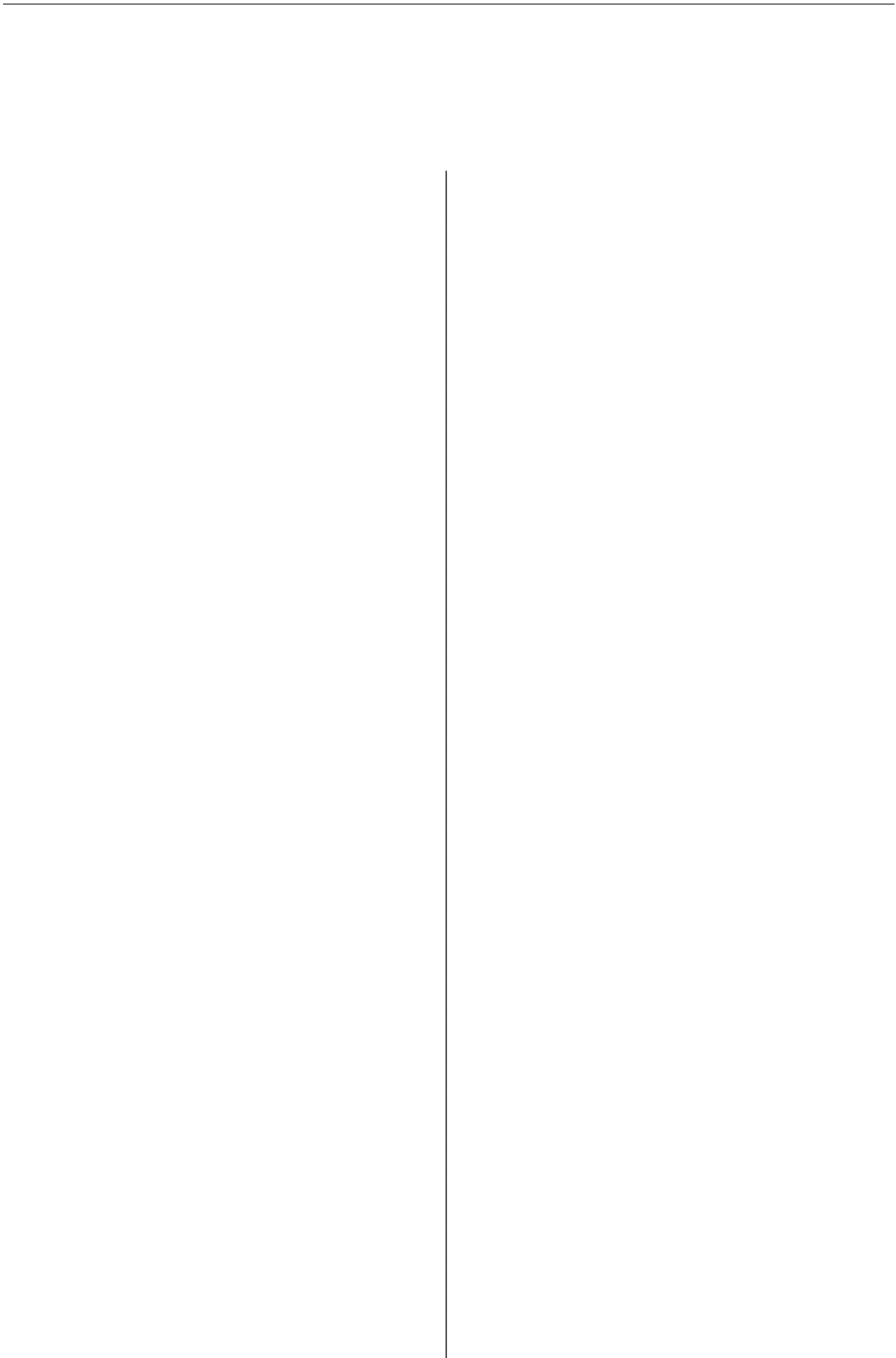


FIGURE A.4 – Arbre linéaire de taille 500 sans labels obtenu avec Asymptote.



FIGURE A.5 – Arbre linéaire de taille 10000 avec labels obtenu avec Asymptote.

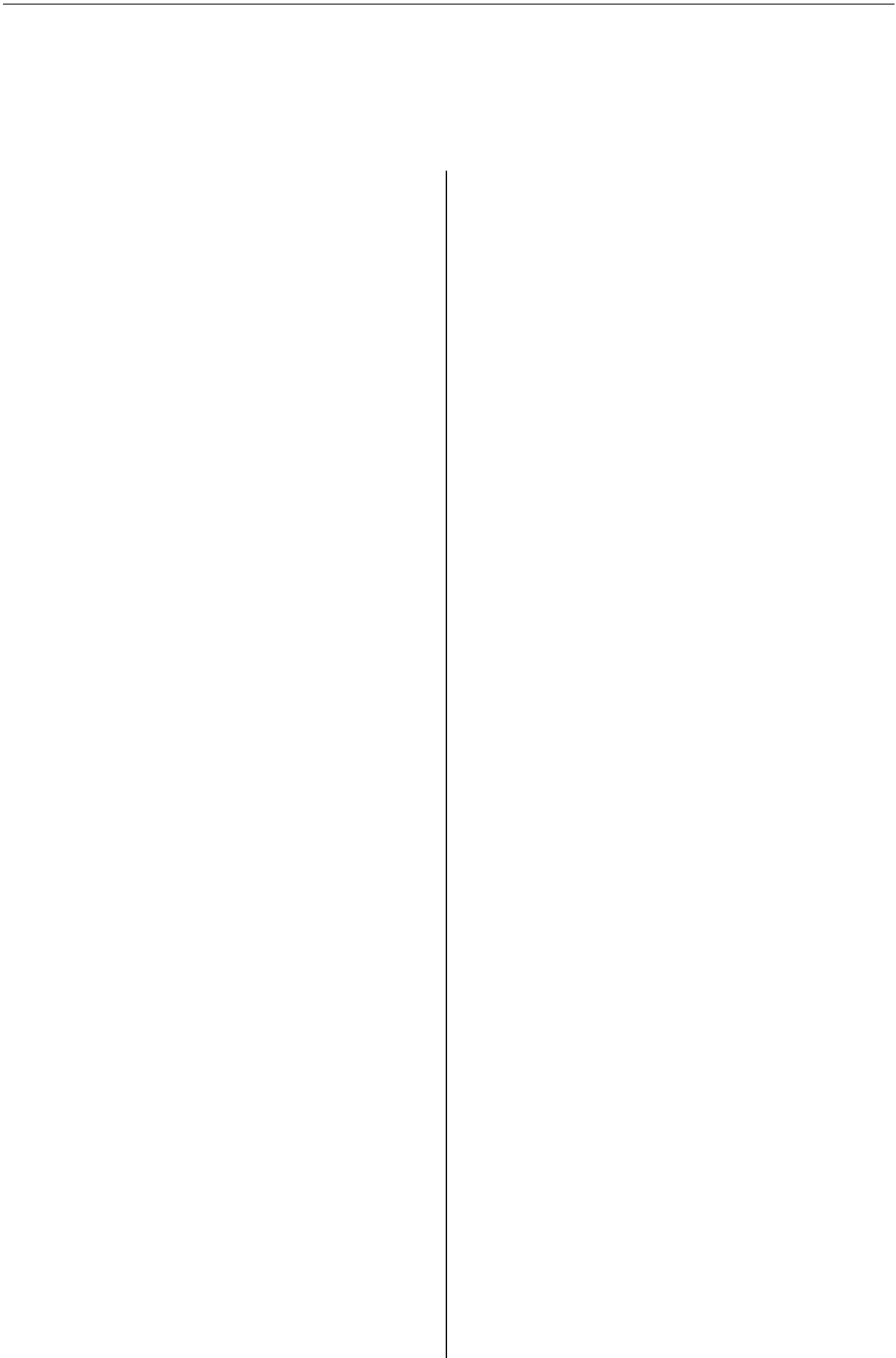


FIGURE A.6 – Arbre linéaire de taille 10000 sans labels obtenu avec Asymptote.

Annexe B

Extraits de l'implémentation

```
1  Éë
from tree import *
3  import re

5  nbCharAlire = 1024

7  def strParser (src):
    "Parses a string file into the intern tree structure"
9
    f=open(src , 'r')
11   s=f.read(nbCharAlire)
    #(T, f) = parse(f, f.read(1))
13   T=parse(f, s)[0]
    f.close()

15   return T

17  def parse (f, s, currentIndex=0):
19
    isWhiteSpace=re.compile("[\s]")

21
    (f, s, currentIndex) = skipSpaces(f, s, currentIndex)

23
    if s[currentIndex] != '(':
25         raise Exception ("The tree is ill-formed : A tree should start with a '('")

27
    (f, s, currentIndex) = incrementIndex(f, s, currentIndex)
    size = len(s)
29
    c = s[currentIndex]
    label=""

31
    (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
33
    if (s==[]):
        raise Exception ("The tree is ill-formed : End of file reached")
35
    c = s[currentIndex]

37
    while (s!=[] and c!='(' and c != ')'):
        #print("currentIndex=",currentIndex)
39         if (not isWhiteSpace.match(c)):
            label = label + c
41         (f, s, currentIndex) = incrementIndex(f, s, currentIndex)
            c = s[currentIndex]

43
    (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
45
    c = s[currentIndex]

47
    listChildren = []
    while (s!=[] and c == '('):
49         (child, s, currentIndex) = parse (f, s, currentIndex)
            listChildren.append(child)
51         #print("fils ",label,"éajout")
            (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
53         c = s[currentIndex]
        #print("c éaprs avoir éajout le fils",label," : ",c)
```

```

55 (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
57 c = s[currentIndex]
   if (c == ')'):
59     (f, s, currentIndex) = incrementIndex(f, s, currentIndex)
   return (Tree(label = label, children=listChildren), s, currentIndex)
61 else:
   raise Exception ("The tree is ill-formed : The tree",label,"should end with a
   ')")
63
def incrementIndex(f, s, currentIndex):
65     currentIndex+=1
   l=len(s)
67   if (currentIndex>=l):
   #print("fin du buffer -> on re-remplit")
69     s=f.read(nbCharAlire)
   currentIndex=0
71   if (s==[]):
   raise Exception ("The tree is ill-formed : End of file reached")
73 #print("currentIndex dans incrementIndex :",currentIndex)
   return (f, s, currentIndex)
75
def skipSpaces(f, s, currentIndex) :
77   isWhiteSpace=re.compile("[\s]")
   while (s!=[] and isWhiteSpace.match(s[currentIndex])):
79     (f, s, currentIndex) = incrementIndex(f, s, currentIndex)
   return (f, s, currentIndex)

```

```

Éë
2 from tree import *
  from pyparsing import Word, alphas
4 import re

6 def dotParser (src):
   "Parses a dot file into the intern tree structure"
8
   f= open(src)
10
   s = f.read()
12
   (toto, i) = getWord (s, 0, '{') # [strict] (graph | digraph) [ID]
14
   k=0
16
   # [strict]
18   while (toto[k] == " "):
   k += 1
20
   if (toto[k] == 's'): #read "strict"
22     if re.compile('strict').search(toto):
   if (re.search(u'strict', toto).start() != 0):
24       raise Exception ("Dot syntaxe error")
   else:
26     k = re.search(u'strict', toto).end()
   else:
28     raise Exception ("Dot syntaxe error")

30 # (graph|digraph)
   while (toto[k] == " "):
32     k += 1

34   if (toto[k] == 'g'): #read "graph"
   if re.compile('graph').search(toto):
36     if (re.search(u'graph', toto).start() != 0):
   raise Exception ("Dot syntaxe error")
38     else:
   k = re.search(u'graph', toto).end()
40     else:
   raise Exception ("Dot syntaxe error")
42   elif (toto[k] == 'd'): #read "digraph"
   if re.compile('digraph').search(toto):
44     if (re.search(u'digraph', toto).start() != 0):
   raise Exception ("Dot syntaxe error")

```

```

46         else:
47             k = re.search(u'digraph', toto).end()
48         else:
49             raise Exception ("Dot syntaxe error")
50     else:
51         raise Exception ("Dot syntaxe error")
52
53     # [ID]
54     while (k<i and toto[k] == " "):
55         k += 1
56
57
58     i+= 1 #skip the {
59     print (s[i])
60
61     dico = {}
62     roots = set()
63     children = set()
64
65     try:
66         while (True):
67             (seq, i) = getWord (s, i, ';' )
68             seq = re.sub(r'\s', "", seq)
69             j=0
70             if (s[i-1] == ' '):
71                 # A node should be matched
72
73                 (num, j) = getWord(seq, j, '[' )
74                 j += 1 # skip the '['
75
76                 (etiquette, j) = getWord (seq, j, '=')
77                 j += 1 # skip the '='
78
79                 (label, j) = getWord (seq, j, ' ')
80                 j += 1 # skip the ' '
81
82                 if (num == "" or etiquette != "label" or label[0] != '\\' or
83                     label[len(label)-1] != '\\'):
84                     raise Exception ("Tree is ill-formed")
85
86                 dico[num] = Tree (label=label[1:len(label)-1])
87                 roots.add(dico[num])
88
89     else:
90         # An arrow should be matched
91
92         (start, j) = getWord (seq, j, '-')
93         j+=2 #skip — or ->
94
95         end = seq[j:]
96
97         # Check if start and end are in dico
98         t = dico[start]
99         c = dico[end]
100         if (t == None or c == None):
101             raise Exception ("Tree is ill-formed")
102
103         # If start not in children add to roots, add end to children
104         #if t not in children:
105         # roots.add (t)
106
107         if c in roots:
108             roots.remove (c)
109
110         children.add (c)
111         #print (t, c, t.children, c.children)
112
113         #Add child to node
114         t.children.append(c)
115
116         i+=2 # skip the ; and go to the following item
117 except IndexError:

```

```

118     seq = s[i:]
119     seq = re.sub(r'\s', "", seq)
120     if (seq == "{}"):
121         # end of tree
122
123         #verify that there is only one root and found it for returning
124         if len(roots) != 1:
125             raise Exception ("Tree is ill-formed")
126
127         f.close()
128         for x in roots:
129             return x
130         #return roots.get()
131
132     #Exception indice out of array should have been raised
133     f.close()
134     raise Exception ("Tree is ill-formed")
135
136 def getWord (s, start, end):
137     "Return a substring of s that starts at indice start and ends with character end
138     excluded. Return also the indice of caracer end."
139     res = ""
140     j = start
141     while (s[j] != end):
142         res = res + s[j]
143         j+=1
144     return (res, j)

```

```

1  Éë
2  import xml.etree.ElementTree as etree
3  from tree import *
4
5  def xmlParser (src):
6      "Parses a xml file into the intern tree structure"
7
8      tree = etree.parse(src)
9      root = tree.getroot()[0] #getRoot returns the node with the "tree" tag,
10                             #we want the "node" or "leaf" inside that node
11
12      return parse(root)
13
14 # Given an xml tree parsed by ElementTree, returns the corresponding Tree object
15 def parse(xmltree):
16
17     if(xmltree.tag=='leaf'):
18         # Case 1 : the tree is a leaf
19         # Try to get the id
20         l=xmltree.get('id')
21         if (l):
22             return Tree(label=l)
23         else:
24             return Tree()
25     else:
26         # Case 2 : the tree is a node with children
27         # Get the children
28         children=[]
29         for child in xmltree:
30             children.append(parse(child))
31         # Try to get the id of the node
32         l=xmltree.get('id')
33         if (l):
34             return Tree(label=l, children=children)
35         else:
36             return Tree(children=children)

```