

Visualisation d'arbres de grandes tailles ¹

Rapport de PSTL

Érika Baëna Diana Malabard
erika.baena@etu.upmc.fr diana.malabard@etu.upmc.fr

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

Université Pierre et Marie Curie
15 mai 2014

1. Disponible sur: https://github.com/BErika/PSTL_TreeDisplay

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	7
3.2.3	Fonctionnement de l'algorithme de calcul de coordonnées	8
3.2.4	Fonctionnement de la génération du code	10
3.3	Complexité	11
4	Étude de performances	12
4.1	Protocole	12
4.2	Résultats	13
4.2.1	Comparaison des temps d'exécution des différents parsers.	13
4.2.2	Évolution du temps de calcul des coordonnées en fonction du nombre de nœuds de l'arbre.	14
4.2.3	Comparaison des temps d'exécution des différents générateurs.	15
4.2.4	Comparaison des temps d'exécution entre GraphViz et la meilleur combinaison parser/générateur de TreeDisplay.	17
5	Conclusion	18
5.1	Bilan	18
5.2	Pour la suite	18
A	Images obtenues lors de l'étude préliminaire	19
B	Extraits de l'implémentation	27
	Bibliographie	31

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 de manière triviale. Ceci pose problème lorsque l'on souhaite représenter des arbres dont l'ordre des fils est primordial : les arbres de recherche par exemple.

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 intégrer le code à un document L^AT_EX ; Asymptote, une alternative à TikZ ; et NetworkX, pour pouvoir générer automatiquement un pdf ou une image de l'arbre que l'on pourra ultérieurement insérer dans n'importe quel document.

Chapitre 1

État des lieux

Les articles que nous avons étudiés concernaient avant tout la visualisation des arbres binaires de grande taille sans prendre en compte l'ordre des fils. Il en ressort néanmoins plusieurs principes :

1. Les arêtes de l'arbre ne doivent pas s'intersecter.
2. Les nœuds de même profondeur doivent être dessinés sur la même ligne horizontale.
3. Les arbres doivent être dessinés de la manière la plus compacte possible.
4. Un nœud parent doit être centré par rapport à ses fils.
5. Un sous-arbre doit être dessiné de la même façon, peu importe où il est placé dans l'arbre.
6. Les nœuds fils d'un nœud père doivent être espacés de manière homogène.

Un des enjeux soulevés par l'ensemble de ces principes est que les nœuds ne doivent pas se chevaucher, tout comme les arbres en eux-même.

Il existe plusieurs algorithmes permettant de dessiner des arbres de grande taille, mais tous ne respectent pas tous les principes décrits ci-dessus.

Par exemple, l'algorithme le plus simple, l'algorithme de Knuth [2], ne respecte que les deux premiers principes. Cet algorithme décrit déjà une idée de "slot disponible". Mais ce sont Charles Wetherell et Alfred Shannon [6], en 1979, qui introduiront l'utilisation d'un tableau qui associera à chaque profondeur le prochain slot disponible. Ils arrivent ainsi à respecter tous les principes, à l'exception des principes 4 et 5. Ils introduisent malgré tout également l'idée de parcourir l'arbre de bas en haut, plutôt que l'inverse, ceci afin de centrer facilement un père selon ses fils.

Le principal problème est alors : Comment respecter tous ces principes et traiter le chevauchement d'arbres sans perdre en complexité ? C'est l'algorithme The Mods and the Rockers qui répondra à cette question. Au lieu de reparcourir les sous-arbres pour les décaler afin d'éviter tout chevauchement, on raisonne en deux passes de l'arbre. Lors de la première passe, on mémorise un *modifier* qui indiquera le décalage qui devra être appliqué sur chaque sous-arbre lors de la deuxième passe.

Il existe enfin des algorithmes dont le but est surtout d'optimiser les concepts des algorithmes existants. Citons le concept de contour d'arbre, qui permet de ne parcourir que le contour et donc de ne pas rentrer au cœur de l'arbre. Ceci est un gain conséquent puisque nous traitons ici des arbres de grande taille, et donc potentiellement très larges.

Nous allons donc voir à présent comment nous nous sommes inspirées de ces algorithmes afin de généraliser les concepts aux arbres n-aires, et en conservant l'ordre des fils.

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 A.5 et A.6.

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
import networkx as nx
import sys

nbIte = int(sys.argv[1])
longueur_arete = float(sys.argv[2])
if (sys.argv[3] == "true"):
    labels = True
else:
    labels = False
filename="testNetworkx%d" % (nbIte,)
if (labels):
    filename += ".png"
else:
    filename += ".B.png"
G = nx.path_graph(nbIte)
pos={x: (5, x*longueur_arete) for x in G.nodes()}
if (labels):
    labels={x: "label" for x in G.nodes()}
nx.draw(G, pos, labels=labels, node_size=5, with_labels=True)
else:
    nx.draw(G, pos, node_size=5, with_labels=False)
plt.savefig(filename)
```

2.3.1 500 nœuds

De même que précédemment, on génère d'abord un arbre de 500 nœuds. On choisi le format de sortie png. Le résultat est visible aux figures A.7 et A.8.

2.3.2 10000 nœuds

Passons maintenant à 10000 nœuds. Le résultat est aux figures A.9 et A.10.

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 ce 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/MIS2/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, depreciate 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/MIS2/PSTL/Implementation$
```

3.2 Plus en détails

3.2.1 Structure d'arbre

Nous avons une structure d'arbre n-aire classique avec une liste pour les enfants ainsi que des attributs pour le calcul de coordonnées et la génération finale. En particulier, les attributs `width` et `height` représentent respectivement la plus grande abscisse et la plus grande ordonnée rencontrée dans le sous-arbre. Au niveau de la racine, cela représente la largeur et la hauteur du graphe.

```
1 class Tree(object):
2     "Local representation of tree"
3
4     def __init__(self, x = -1, depth=0, label="", children=None, offset = 0, isRoot
5         = False):
6         self.x = x
7         self.y = depth
8         self.label = label
9         self.offset = offset
10        self.height = None
11        self.width = None
12        if children is None:
13            self.children = list()
14        else:
15            self.children = children
```


3.2.2 Fonctionnement des algorithmes de parsing

Mots bien parenthésés

Le parser de mots bien parenthésés (cf. B.1) respecte la grammaire suivante :

```
ARBRE ::= '(' LABEL NOEUDS ')'  
NOEUDS ::= ARBRE NOEUDS |  $\epsilon$   
LABEL ::= [a-zA-Z1-9]* |  $\epsilon$ 
```

Compte-tenu du fait qu'on veut pouvoir travailler sur des arbres de grande taille, le parser bufferise le fichier d'entrée pour ne pas saturer la mémoire en chargeant tout le fichier dans une variable.

Dot

Le parser dot (cf. B.2) parse une sous-partie du langage dot, à savoir :

```
DOT ::= STRICT GRAPH ID '' SEQINST ''  
STRICT ::= strict |  $\epsilon$   
GRAPH ::= digraph | graph  
SEQINST ::= INST ';' SEQINST |  $\epsilon$   
INST ::= ID '[' label = "LABEL" ']' | ID LINK ID  
LINK ::= -- | ->  
ID ::= [0-9]*  
LABEL ::= [a-zA-Z1-9]* |  $\epsilon$ 
```

XML

Le parser XML (cf. B.3) respecte la grammaire ci-dessous. Comme beaucoup de langages de programmation, Python a une librairie dédiée au XML. Nous l'avons utilisée.

```
XML ::= <?xml version="1.0"?><tree> NOEUDS </tree>  
NOEUDS ::= NOEUD NOEUDS |  $\epsilon$   
NOEUD ::= <node type=TAG id=ID> NOEUDS </noeud> | <leaf type=TAG id=ID />  
TAG ::= " [a-zA-Z1-9]* "  
ID ::= [0-9]*
```

Remarques

On veut pouvoir utiliser les fichiers générés par **Arbogen** [4] qui est un générateur d'arbre. Il nous permet de générer des arbres de grandes tailles pour nos tests. C'est en analysant les fichiers fournis par ce générateur que nous avons établi les différences d'interprétations entre LABEL, ID et TAG.

Un LABEL est une étiquette quelconque sur un nœud, qui peut apparaître à différents endroits de l'arbre. Cela nous assure qu'un mot bien parenthésés ne peut représenter qu'un arbre car seul l'arborescence des parenthèses permet de définir les arcs entre nœuds.

Un ID est un identifiant unique d'un nœud. On remarque alors que la grammaire Dot peut représenter des arbres mais aussi plus généralement des graphes quelconques car les ID sont la seule façon de représenter les arcs entre nœuds.

Un TAG correspond en fait à une partie gauche dans la grammaire donnée en entrée à **Arbogen**. Il n'est pas donc unique dans l'arbre et est interprété dans le parser comme un label. Cependant, la grammaire XML comporte aussi des ID. Si en plus de considérer l'arborescence des balises XML, on considérait aussi les ID, on pourrait aussi représenter des graphes. Ce n'est pas le cas de notre parser XML car on s'intéresse aux arbres.

3.2.3 Fonctionnement de l'algorithme de calcul de coordonnées

On décide que la distance minimale entre deux nœuds est de 1 unité selon l'axe des abscisses et celui des ordonnées.

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.

On veut qu'un père soit centré sur ses fils. On commence donc par placer ses fils s'il en a puis on fait la moyenne de l'abscisse de ses 2 fils extrémaux pour le centrer.

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 nœud a des fils, il se peut qu'en calculant son abscisse en fonction de ses fils on se retrouve avec un nœud qui est trop proche de son frère gauche comme dans l'exemple ci-dessous.

Considérons l'arbre suivant : (root (n11) (n12) (n13 (n21) (n22) (n23))).

	root		n11		n12		n13		n21		n22		n23	
	x	y	x	y	x	y	x	y	x	y	x	y	x	y
Initialement	-1	0	-1	0	-1	0	-1	0	-1	0	-1	0	-1	0
Appel sur root	-1	0	-1	0	-1	0	-1	0	-1	0	-1	0	-1	0
Appel sur n11	-1	0	0	1	-1	0	-1	0	-1	0	-1	0	-1	0
Retour sur root	-1	0	0	1	-1	0	-1	0	-1	0	-1	0	-1	0
Appel sur n12	-1	0	0	1	1	1	-1	0	-1	0	-1	0	-1	0
Retour sur root	-1	0	0	1	1	1	-1	0	-1	0	-1	0	-1	0
Appel sur n13	-1	0	0	1	1	1	-1	0	-1	0	-1	0	-1	0
Appel sur n21	-1	0	0	1	1	1	-1	0	0	2	-1	0	-1	0
Retour sur n13	-1	0	0	1	1	1	-1	0	0	2	-1	0	-1	0
Appel sur n22	-1	0	0	1	1	1	-1	0	0	2	1	2	-1	0
Retour sur n13	-1	0	0	1	1	1	-1	0	0	2	1	2	-1	0
Appel sur n23	-1	0	0	1	1	1	-1	0	0	2	1	2	2	2
Retour sur n13	-1	0	0	1	1	1	-1	0	0	2	1	2	2	2

La position ensuite calculée pour n13 est 1. Or il y a déjà n12 à cet endroit. La première place disponible à profondeur 1 est 2. On va donc devoir effectuer un décalage de 1.

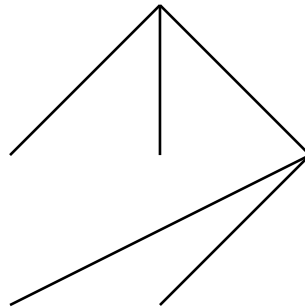


FIGURE 3.1 – Résultat de l'exemple à la fin de l'étape 1.

On compare la position calculée avec la première position disponible et on prend le maximum. 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-arbres dans un second temps pour des raisons de complexité. En effet, si on faisait chaque décalage lorsqu'il se présentait, le décalage serait quadratique car il y a potentiellement un décalage pour chaque ancêtre d'un nœuds.

```

1  def setup (self , depth=0, nexts=None, offset=None):
2      if nexts is None:
3          nexts = defaultdict(lambda:0)
4          if offset is None:
5              offset = defaultdict(lambda:0)
6
7      # L'ordonnée est triviale , c'est la profondeur.
8      self.y = depth
9
10     # On calcule d'abord les coordonnées des enfants.
11     for c in self.children:
12         c.setup(depth+1, nexts, offset)
13
14     # On centre le noeud au milieu de ses enfants.
15     nbChildren = len(self.children)
16     if (nbChildren == 0):
17         place = nexts[depth]
18         self.x = place
19     else:
20         place = (self.children[0].x + self.children[nbChildren-1].x) / 2
21
22     # On calcule l'éventuel décalage engendré.
23     offset[depth] = max(offset[depth], nexts[depth]-place)
24
25     # On applique le décalage de la profondeur.
26     if (nbChildren != 0):
27         self.x = place + offset[depth]
28
29     # On met à jour la prochaine place disponible à cette profondeur.
30     nexts[depth] = self.x + 1

```

Il faut ensuite appliquer récursivement les décalages calculés lors de la première passe. On calcule à la même occasion la hauteur et la largeur de l'image pour pouvoir connaître le ratio de l'image lors de la génération.

```

1  def addOffsets (self , offsum=0):
2      self.x = self.x + offsum
3      offsum = offsum + self.offset
4
5      self.height = self.y
6      self.width = self.x
7
8      for c in self.children:
9          c.addOffsets(offsum)
10         self.height = max (self.height, c.height)
11         self.width = max (self.width, c.width)

```

Le résultat final est le suivant :

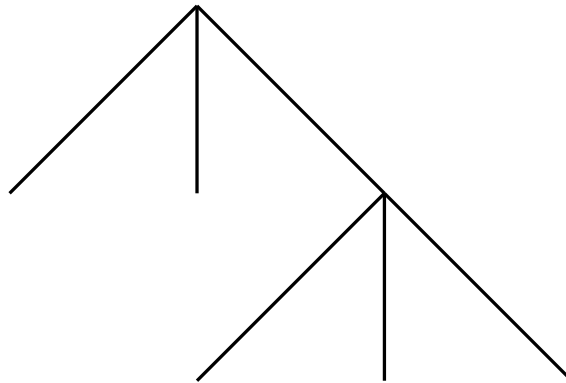


FIGURE 3.2 – Résultat de l'exemple à la fin de l'étape 2.

3.2.4 Fonctionnement de la génération du code

TikZ

```
1 from tree import *
3 def toTikz (t, withLabels = False, fileName = "treeTikz.tex"):
    "Generate an output file which name is fileName with the coordonates of the tree
    t. Node labels are displayed if withLabels is True"
5     f = open (fileName, "w")
    if (withLabels):
6         toTikzRecWith (t, f)
    else:
7         toTikzRecWithout (t, f)
8     f.close()
11
13 def toTikzRecWith (t, f):
    f.write("\\node (a%d) at (%f, -%f) {%s};\\n" % (id(t), t.x, t.y, t.label,))
15     for c in t.children:
        toTikzRecWith (c, f)
17     f.write("\\draw (a%d) — (a%d);\\n" % (id(t), id(c),))
19
21 def toTikzRecWithout (t, f):
    for c in t.children:
        toTikzRecWithout (c, f)
23     f.write("\\draw (%f, -%f) — (%f, -%f);\\n" % (t.x, t.y, c.x, c.y,))
```

Asymptote

```
from tree import *
2
3 def toAsymptote (t, withLabels=False, fileName="treeAsy.tex"):
4     f = open (fileName, "w")
5
6     f.write("\\begin{asy}\\n") #LaTeX file doesn't compile if asy environment is
    written in the LaTeX file ... strange!
7     f.write("size(20cm, 20cm);\\n")
8     if (withLabels): # Avoid to do the test at each iteration, even if
        toAsymptoteRecWith and toAsymptoteRecWithout are really similar
        toAsymptoteRecWith(t, f)
9     else:
        toAsymptoteRecWithout (t, f)
10    f.write("\\end{asy}\\n")
11
12    f.close()
14
16 def toAsymptoteRecWithout (t, f):
    for c in t.children:
18         toAsymptoteRecWithout(c, f)
        f.write("draw((%f, -%f) — (%f, -%f));\\n" % (t.x, t.y, c.x, c.y,))
20
22 def toAsymptoteRecWith (t, f):
    f.write("label(\"%s\", (%f, -%f), E);\\n" % (t.label, t.x, t.y,))
    for c in t.children:
24         toAsymptoteRecWith(c, f)
        f.write("draw((%f, -%f) — (%f, -%f));\\n" % (t.x, t.y, c.x, c.y,))
```

Autre

```
from tree import *
2
3 import networkx as nx
4 import sys
5 import matplotlib.pyplot as plt
6
7 node_size=5
8
```

```

#Main function
10 def toNetworkX(treeToDraw, addLabels=False, outputfile="treeNetworkX",
    outputformat="png") :
    if (not(treeToDraw==None)) :
12         print("Generating NetworkX...")
        (G, pos) = toNetworkXrec(treeToDraw)
14         print("Printing...")
        plt.clf()
16         plt.gca().invert_yaxis()
        nx.draw(G, pos, with_labels=False, node_size=node_size)
18         # Draw the labels a little upper, so they don't overlap the nodes
        if (addLabels):
20             node_labels = nx.get_node_attributes(G, 'label')
            pos_labels=dict()
22             for k, v in pos.items():
                pos_labels[k]=(v[0], v[1]-0.1)
24             nx.draw_networkx_labels(G, pos_labels, labels=node_labels, with_labels=True,
                node_size=node_size)
            plt.savefig(outputfile+'.'+outputformat)
26             G.clear()
        else :
28             print("Tree must not be empty.")

30 #Aux rec function
def toNetworkXrec(treeToDraw, pos=None) :
32     if (not(pos)) :
        pos=dict()
34     #create a graph and add the root and its position
    G=nx.Graph()
36     G.add_node(hex(id(treeToDraw)), label=treeToDraw.label)
    pos[hex(id(treeToDraw))] = (treeToDraw.x, treeToDraw.y)
38     #print("noeud courant : ",treeToDraw.label)
    #print("est une feuille ? ",(not(treeToDraw.children)))
40     if (not(treeToDraw.children)) :
        #the current node is a leaf -> return the graph and the associated positions
42         return (G, pos)
    else :
44         for t in treeToDraw.children :
            #the current tree is not a leaf ->
46             #add one edge from the root to each son, make a graph
            #from each son, and combine all the results
            #G.add_edge(treeToDraw.label, t.label)
48             G.add_edge(hex(id(treeToDraw)), hex(id(t)))
            (H, pos2) = toNetworkXrec(t, pos)
50             for (node, data) in H.nodes_iter(data=True):
                G.add_node(node, **data)
52             G.add_edges_from(H.edges())
            pos.update(pos2)
54         return (G, pos)

```

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 $O(n)$ car on parcourt linéairement le fichier.
2. Le calcul des coordonnées est en $O(n)$ car on effectue 2 passes sur l'arbre et lors de chaque passe on visite une et une seule fois chaque nœud.
3. La génération du fichier de sortie est en $O(n)$ car on visite une unique fois chaque nœud pour écrire ses coordonnées dans le fichier de sortie.

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

Étude de performances

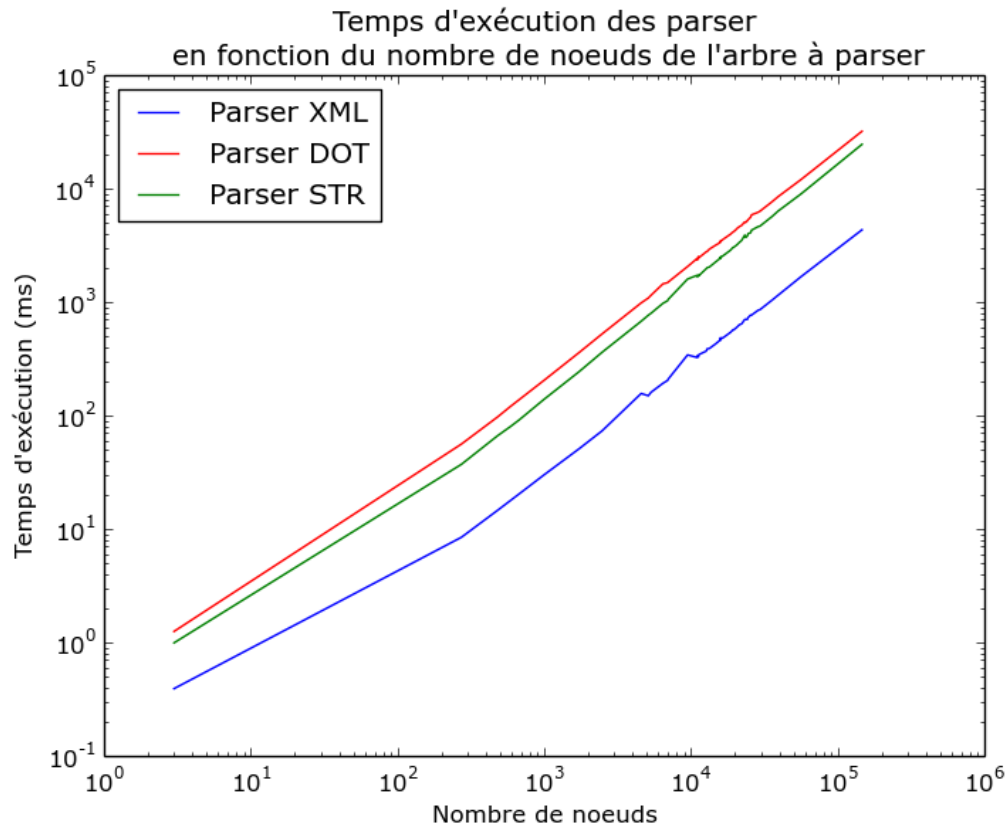
4.1 Protocole

Avant toute chose, il est important de préciser que les tests effectués ci-dessous ont été réalisés sans aucun label. Comme expliqué précédemment, nous pouvons considérer que si la taille des labels est constante, la complexité des algorithmes ne dépendra que du nombre de nœuds compris dans l'arbre. C'est donc l'hypothèse que nous faisons ici.

Nous avons généré environ 2200 arbres de tailles croissantes en XML avec **Arbogen**. Sur ces 2200 arbres, nous en gardé aléatoirement 200 qui constituent notre jeu de tests. Chacun de ces arbres au format XML a été converti aux formats dot et str pour exécuter les parsers sur les mêmes arbres. Ensuite, chacune des trois étapes de notre application - à savoir parsing, calcul de coordonnées et génération - a été exécuté 3 fois et la moyenne a été prise comme valeur de référence.

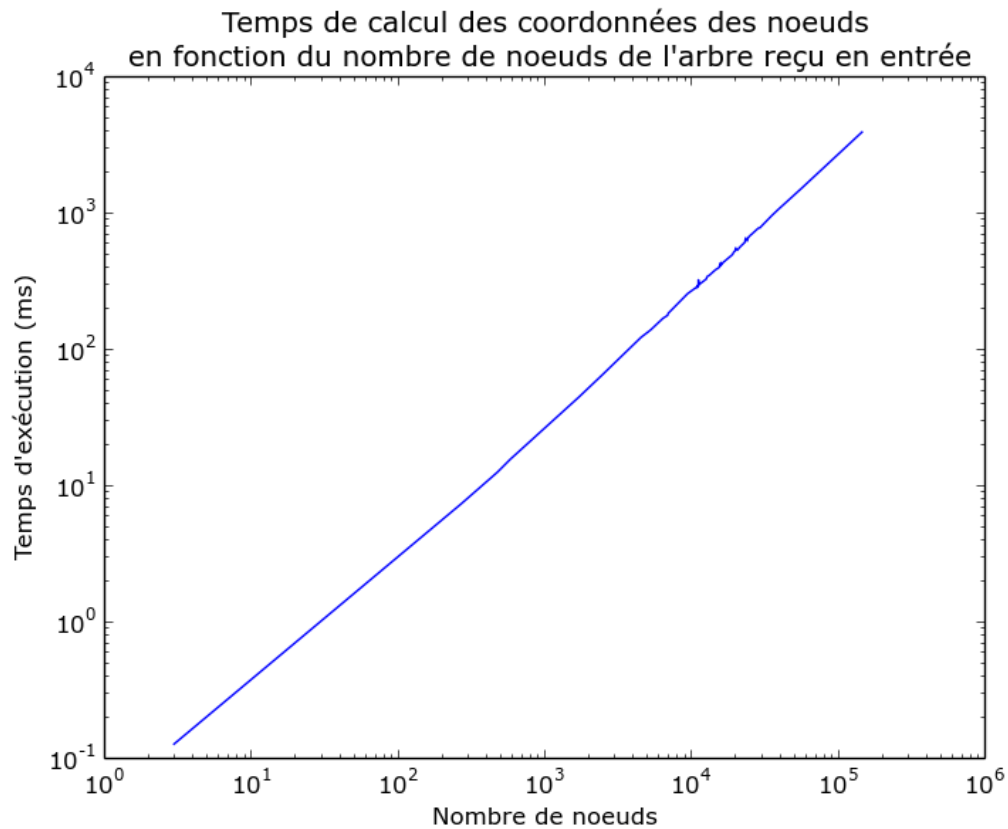
4.2 Résultats

4.2.1 Comparaison des temps d'exécution des différents parsers.



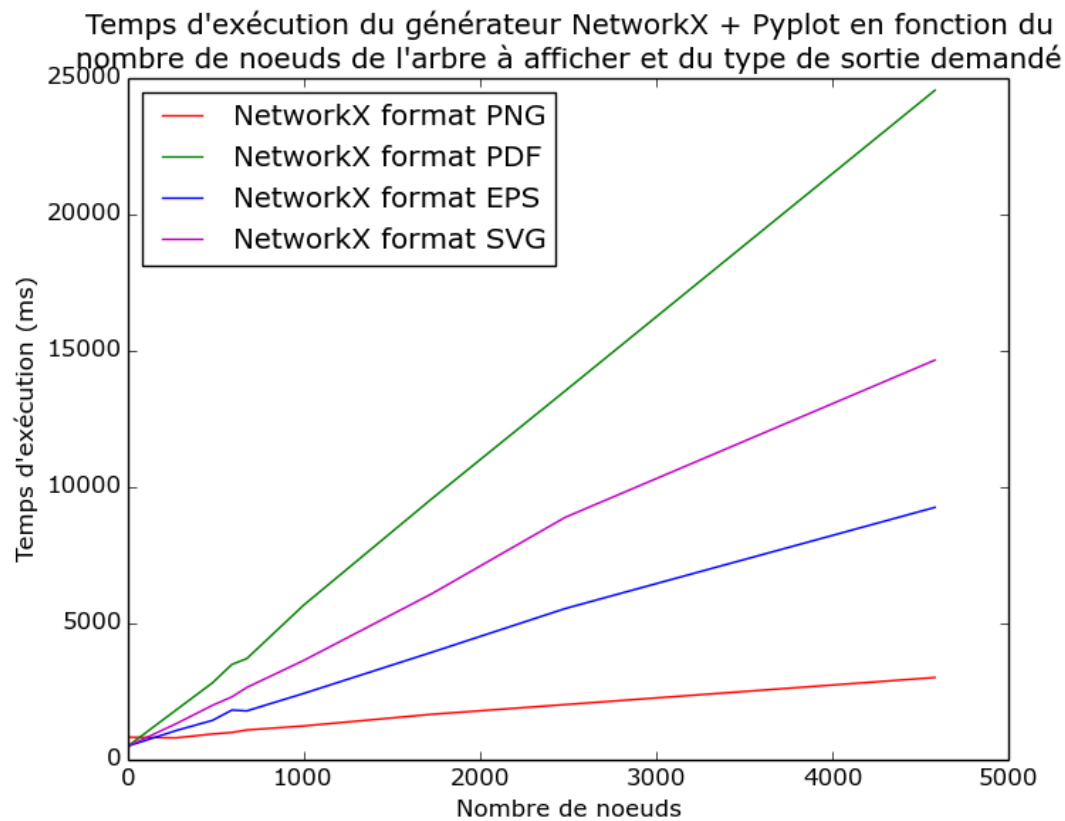
Comme on avait pu le montrer précédemment, on peut observer que les parser ont une complexité linéaire. Ce qu'il ressort également de cette comparaison, c'est que le parser XML est le plus efficace des trois. Cela est certainement dû au fait que nous utilisons un module interne à Python pour traiter le fichier d'entrée. En effet, les modules internes des langages sont généralement bien optimisés car ils sont pensés pour le langage concerné. De plus, en à peine une dizaine de secondes, nous pouvons parser un arbre contenant un million de nœuds.

4.2.2 Évolution du temps de calcul des coordonnées en fonction du nombre de nœuds de l'arbre.

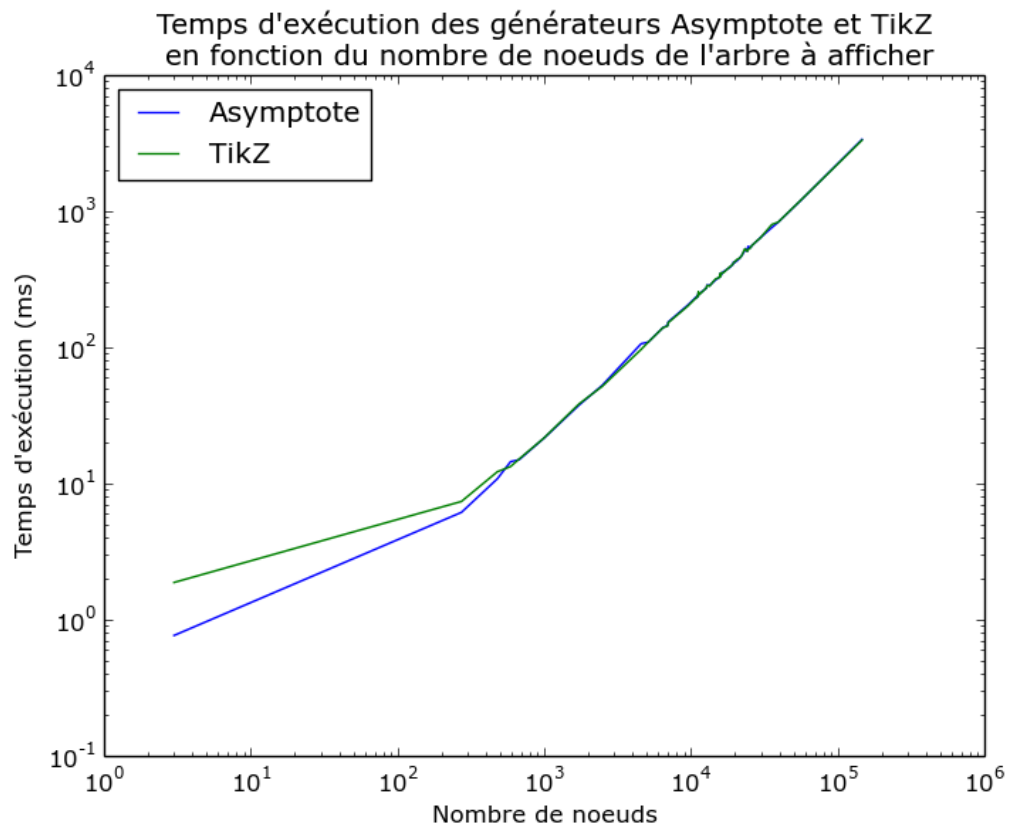


Les tests confirment que notre algorithme est d'une complexité linéaire en temps d'exécution. Nous pouvons également supposer, en observant l'évolution de la courbe, que pour un million de nœuds, nous n'aurions besoin que d'à peine une seconde pour effectuer le calcul de coordonnées, ce qui semble un résultat acceptable.

4.2.3 Comparaison des temps d'exécution des différents générateurs.

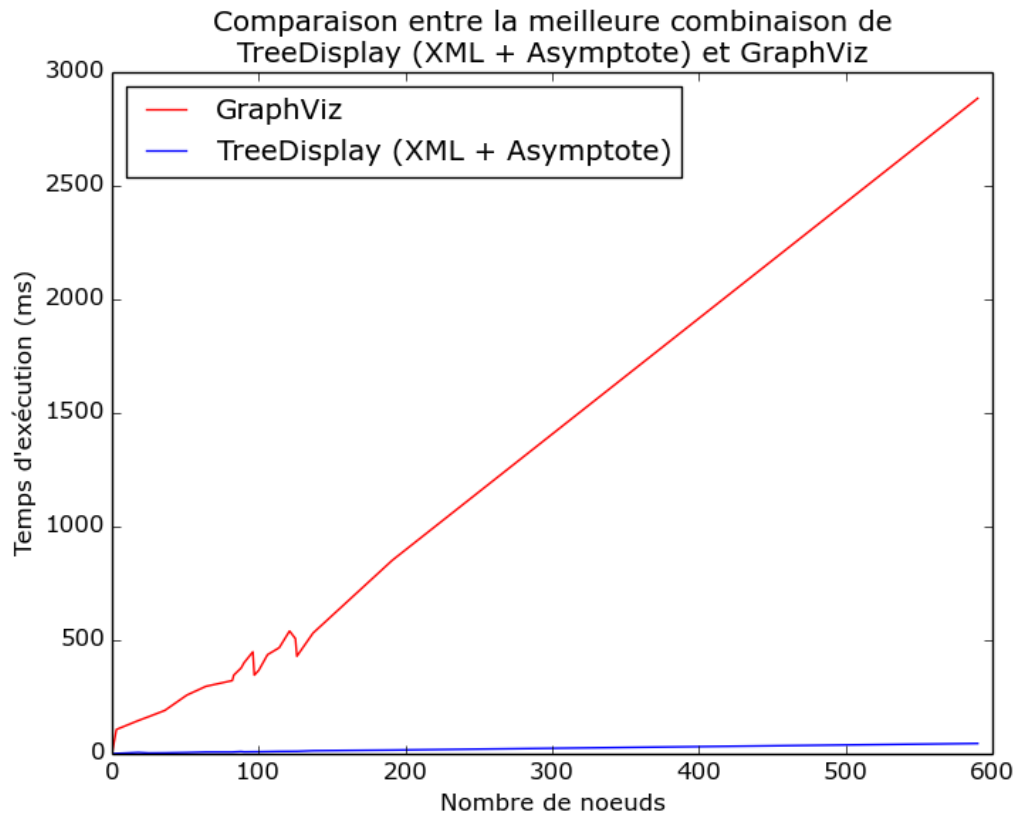


Nous comparons tout d'abord les différentes sorties gérées par NetworkX. Nous observons que nous sommes toujours en complexité linéaire, mais que le temps d'exécution est bien supérieur à ce qui a été observé pour les autres générateurs, avec plusieurs secondes pour seulement 5000 noeuds. Cependant, la complexité de ce générateur vient en grande partie du fait qu'il génère une image via Matplotlib. Ceci ralentit grandement l'exécution.



Contrairement à NetworkX, nous générons ici un fichier qui doit être recompilé. Ceci explique certainement en partie le fait que ce générateur soit plus efficace que le générateur NetworkX. Cependant, nous restons toujours en complexité linéaire.

4.2.4 Comparaison des temps d'exécution entre GraphViz et la meilleur combinaison parser/générateur de TreeDisplay.



Nous comparons maintenant TreeDisplay dans sa meilleure combinaison possible de parser et de générateur de sortie, avec le visualiseur que nous avons cité en début de ce rapport : GraphViz. Il en ressort que TreeDisplay est bien plus puissant que GraphViz. Cependant, GraphViz gère des graphes de manière générale et fournit une image en sortie, et nous avons généré un fichier .tex à compiler. Ceci peut expliquer en partie la différence de performance.

Chapitre 5

Conclusion

5.1 Bilan

L'application `treeDisplay` permet de calculer les coordonnées des nœuds d'un arbre en temps linéaire par rapport à ce nombre de nœuds. Cependant, l'affichage qui en résulte n'est pas aussi élégant que GraphViz. Nous ne pouvons par exemple pas voir la structure de l'arbre en détail. Mais puisque nous travaillons majoritairement avec des arbres de grande taille, les détails nous intéressent peu. On veut avant tout observer des tendances.

Nous avons également observé que la partie de l'application qui prend le plus de temps à l'exécution est le parsing et la génération de sortie. La clé de la complexité de cette application était donc de lire les fichiers d'entrée et de générer les fichiers de sortie de manière optimale, de façon à utiliser le moins de ressources possible.

5.2 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.

L'algorithme de calcul de coordonnées implémenté respecte un certain nombre de contraintes que nous avons appelés *principes*. On peut étudier un autre algorithme qui va respecter d'autres contraintes comme par exemple la minimisation de la largeur.

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.

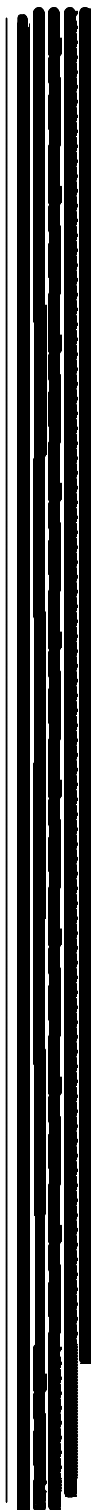


FIGURE A.5 – Arbre linéaire de taille 10000 avec labels obtenu avec Asymptote.



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



FIGURE A.7 – Arbre linéaire de taille 500 avec labels obtenu avec NetworkX et Maptplotlib.



FIGURE A.8 – Arbre linéaire de taille 500 sans labels obtenu avec NetworkX et Maptplotlib.



FIGURE A.9 – Arbre linéaire de taille 10000 avec labels obtenu avec NetworkX et Matplotlib.



FIGURE A.10 – Arbre linéaire de taille 10000 sans labels obtenu avec NetworkX et Matplotlib.

Annexe B

Extraits de l'implémentation

```
from tree import *
2 import re

4 nbCharAlire = 1024

6 def strParser (src):
    "Parses a string file into the intern tree structure"
    8
    f=open(src, 'r')
    10 s=f.read(nbCharAlire)
    #(T, f) = parse(f, f.read(1))
    12 T=parse(f, s)[0]
    f.close()

    14 return T

16 def parse (f, s, currentIndex=0):
    18
    isWhiteSpace=re.compile("[\s]")
    20
    (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
    22
    if s[currentIndex]!='(':
    24     raise Exception ("The tree is ill-formed : A tree should start with a '('")

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

    30
    (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
    32 if (s==[]):
        raise Exception ("The tree is ill-formed : End of file reached")
    34 c = s[currentIndex]

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

    42
    (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
    44 c = s[currentIndex]

    46 listChildren = []
    while (s!=[] and c == '('):
    48     (child, s, currentIndex) = parse (f, s, currentIndex)
        listChildren.append(child)
    50     #print("fils ",label,"ajouté")
        (f, s, currentIndex) = skipSpaces(f, s, currentIndex)
    52     c = s[currentIndex]
        #print("c après avoir ajouté le fils",label," : ",c)
    54
```

```

(f, s, currentIndex) = skipSpaces(f, s, currentIndex)
56 c = s[currentIndex]
if (c == ' '):
58     (f, s, currentIndex) = incrementIndex(f, s, currentIndex)
    return (Tree(label = label, children=listChildren), s, currentIndex)
60 else:
    raise Exception ("The tree is ill-formed : The tree",label,"should end with a
        ' '")
62
def incrementIndex(f, s, currentIndex):
64     currentIndex+=1
    l=len(s)
66     if (currentIndex>=l):
        #print("fin du buffer -> on re-remplit")
68         s=f.read(nbCharAlire)
        currentIndex=0
70     if (s==[]):
        raise Exception ("The tree is ill-formed : End of file reached")
72     #print("currentIndex dans incrementIndex :",currentIndex)
    return (f, s, currentIndex)
74
def skipSpaces(f, s, currentIndex) :
76     isWhiteSpace=re.compile("[\s]")
    while (s!=[] and isWhiteSpace.match(s[currentIndex])):
78         (f, s, currentIndex) = incrementIndex(f, s, currentIndex)
    return (f, s, currentIndex)

```

Listing B.1 – ../Src/strParser.py

```

from tree import *
2 from pyparsing import Word, alphas
import re
4
def dotParser (src):
6     "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]
    while (toto[k] == " "):
18        k += 1

20    if (toto[k] == 's'): #read "strict"
        if re.compile('strict').search(toto):
22            if (re.search(u'strict', toto).start() != 0):
                raise Exception ("Dot syntaxe error")
            else:
24                k = re.search(u'strict', toto).end()
        else:
26            raise Exception ("Dot syntaxe error")

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

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

```

```

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)

```

Listing B.2 – ../Src/dotParser.py

```

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

```

Listing B.3 – ../Src/xmlParser.py

Bibliographie

- [1] Christoph BUCHHEIM, Michael JÜNGER et Sebastian LEIPERT : Improving walker's algorithm to run in linear time. *In Graph Drawing*, 2002. <http://dirk.jivas.de/papers/buchheim02improving.pdf>.
- [2] D. E. KNUTH : Optimum binary search trees. *Acta Informatica*, 1(1):14–25, 1971. <http://wenku.baidu.com/view/a4feec5e312b3169a451a42f>.
- [3] Kim MARRIOTT et Peter J. STUCKEY : Np-completeness of minimal width unordered tree layout. *Journal of Graph Algorithms and Applications*, 8(3):295–312, 2004. <http://www.emis.de/journals/JGAA/accepted/2004/MarriottStuckey2004.8.3.pdf>.
- [4] Frédéric PESCHANSKI : Arbogen, 2014. <https://github.com/fredokun/arbogen>.
- [5] Edward M. REINGOLD et John S. TILFORD : Tidier drawing of trees. *IEEE Transactions on Software Engineering*, 7(2):223–227, mars 1981. http://scholar.google.fr/scholar_url?url=fr&q=http://emr.cs.iit.edu/~reingold/tidier-drawings.pdf&sa=X&scisig=AAGBfm26weQoxIBL-KjosMM1a-_bR1Isnw&oi=scholar&ei=QsZsU_j900nX0QWL24HwAw&ved=OCC8QgAMoADAA.
- [6] Charles WETHERELL et Alfred SHANNON : Tidy drawing of trees. *IEEE Transactions on Software Engineering*, 5(5):514–519, septembre 1979. http://scholar.google.fr/scholar_url?url=fr&q=http://poincare.matf.bg.ac.rs/~tijana/geometrija/seminarski/tree_drawing.pdf&sa=X&scisig=AAGBfm34LGKfBgU7qNXhH-xoqu4zpmbfdA&oi=scholar&ei=QsZsU_j900nX0QWL24HwAw&ved=OCDAQgAMoATAA.