

Modéliser en XMLSchema et stocker des données au format XML c'est bien ; pouvoir les manipuler, c'est mieux ! Dans les chapitres qui suivront celui-ci, nous apprendrons à écrire des feuilles de transformation XSLT et à lire et écrire des documents XML. Dans les deux cas, il est nécessaire de naviguer dans le document, c'est-à-dire de sélectionner des éléments ou attributs particuliers, ou même des listes d'éléments, de passer d'un élément à un autre, ... de la même manière que l'on navigue dans l'arborescence des fichiers d'un disque dur.

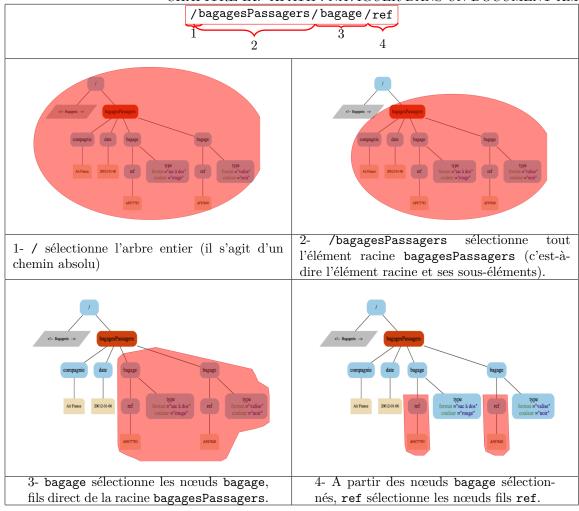
Le langage XPath offre un moyen d'identifier un ensemble de nœuds dans un document XML. Toutes les applications (écrites en Java, Javascript, C++, C#...) ayant besoin de récupérer un fragment de document XML peuvent utiliser ce langage. Nous verrons dans des chapitres ultérieurs de nombreux exemples d'usage avec des technologies utilisant XPath (XSLT, DOM) ; d'autres existent comme XQuery, etc.

XPath est fondé sur la représentation d'un document XML sous forme d'arbre, comme illustré dans les exemples du chapitre précédent: II.7 et III.2.

Une expression XPath (encore appelée chemin XPath ou encore *location path*) est **une expression non XML**, avec une syntaxe compacte. Elle s'évalue sur un nœud courant dans un contexte donné.

Il s'agit d'un ensemble d'étapes de navigation (*location steps*) séparées par des '/'. Par exemple, l'expression /bagagesPassagers/bagage/ref appliquée sur le document XML des figures II.7 et III.2 donne le résultat expliqué dans le tableau suivant<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>voir également l'animation proposée sur la plateforme de cours.



Une expression XPath désigne un ou plusieurs *chemins* dans l'arbre à partir du nœud courant. Elle a pour résultat soit un ensemble de nœuds (*node-set*), soit une valeur (numérique, bouléenne ou alphanumérique).

Remarque: comme vous le verrez en TD et TP, debogguer une application dépendant d'une expression xpath peut s'avérer compliqué. Fort heureusement, il existe de nombreux sites sur Internet permettant de tester des expressions xpath sur des documents XML valides, par exemple <a href="http://www.utilities-online.info/xpath">http://www.utilities-online.info/xpath</a>.

# 1 Node-set

Un nœud est un élément $^2$  + ses attributs.

A partir d'un nœud courant, une étape de chemin XPath (*location step*) renvoie **un ensemble de nœuds**, c'est-à-dire un *node-set*. Chaque nœud de ce *node-set* devient à son tour le nœud courant pour l'évaluation du *location-step* suivant.

On commence par la racine dénotée par /.

# 2 Les axes de recherche

Un axe est un ensemble de nœuds relativement à un autre nœud. On distingue les axes verticaux et les axes horizontaux.

## 2.a Axes verticaux

Parmi les axes verticaux, les axes avant permettent de descendre dans l'arbre:

 $<sup>^2</sup>$ Rappel: un élément est soit une balise ouvrante + un contenu + une balise fermante, soit une balise ouvrante/fermante.

self désigne le nœud courant. Par exemple /bagagesPassagers/bagage/type/self::type[format="sac à dos"] désigne le seul nœud type dont l'attribut format a pour valeur sac à dos.

child désigne les fils du nœud courant (c'est l'axe par défaut). L'expression bagage/child::type est équivalente à l'expression bagage/type.

descendant désigne les descendants (les fils, et les fils des fils des fils des fils des fils, etc.) du nœud courant. par exemple l'expression /bagagesPassagers/descendant::ref désigne les 2 éléments ref de l'arbre.

descendant-or-self désigne les descendants du nœud courant ainsi que le nœud courant lui-même.

Les axes arrière permettent de remonter dans l'arbre:

parent désigne le nœud parent du nœud courant

ancestor désigne les ascendants du nœud courant. Par exemple si le nœud courant est le premier élément ref, l'élément ancestor::bagagesPassagers/compagnie fait référence à l'élément compagnie.

ancestor-or-self désigne les ascendants du nœud courant ainsi que le nœud courant lui-même

### 2.b Axes horizontaux

following-sibling désigne les nœuds frères (i.e. issus du même parent) placés après le nœud courant. Par exemple, /bagagesPassagers/date/following-sibling::node désigne les 2 éléments bagage

preceding-sibling désigne les nœuds frères placés avant le nœud courant.

Par exemple, /bagagesPassagers/date/preceding-sibling::node désigne l'élément compagnie

## 2.c Axe pour les attributs

attribute désigne les attributs du nœud courant

### 2.d Abréviation des Axes

Afin d'éviter une trop grande lourdeur, les simplifications suivantes sont autorisées:

| abréviation | Axe   |
|-------------|---|
| (rien)      | child::                                       |
| @           | attribute::                                   |
|             | self::node()                                  |
| //          | desendant-or-self::node()/                    |
| //X         | <pre>desendant-or-self::node()/child::X</pre> |
|             | <pre>parent::node()</pre>                     |
| /X          | <pre>parent::node/child::X</pre>              |

# 2.e Remarque sur les Chemins

On peut faire un parallèle entre les chemins XPath et les chemins sur les disques durs des ordinateurs. Par exemple, sur l'arborescence d'un ordinateur, / désigne la racine du disque. De même, en XPath, /bagagesPassagers désigne la racine de l'arbre.

De même, sur un disque dur, /cours/XML désigne le fichier XML dans le répertoires cours à la racine du disque ; le chemin /bagagesPassagers/date désigne l'élément date, fils de l'élément racine bagagesPassagers.

Dans les 2 cas, si le chemin commence par / il s'agit d'un chemin absolu (qui part de la racine), sinon, il s'agit d'un chemin relatif (qui part de l'endroit/du nœud courant).

Attention cependant, cette comparaison a des limites, car dans le cas d'un chemin XPath, un chemin absolu (ou relatif) peut sélectionner plusieurs nœuds !

# 86 **3** Sélectionner des parties de l'arbre

#### Les filtres 3.a

- :: (nom) le filtre nom sélectionne les nœuds de l'axe qui portent ce nom.
  - Par exemple, /bagagesPassagers/child::bagage sélectionne les 2 nœuds bagage fils de bagagesPassagers.
- \* le filtre \* sélectionne les nœuds de l'axe qui ont un nom.

Par exemple, /bagagesPassagers/\*/type/attribute::\* sélectionne les 4 attributs des 2 éléments type de l'arbre.

text() le filtre text() sélectionne les nœuds de type texte.

Par exemple, /bagagesPassagers/date/text() renvoie "2012-01-06".

comment() le filtre comment() renvoie les nœuds de type commentaire de l'axe.

Par exemple, /comment() renvoie "Bagagerie".

#### 3.b Les critères

On peut ajouter des critères de sélection en ajoutant des crochets aux location-steps. Chemin[critère] ne sélectionne que les nœuds du Chemin qui correspondent à la sélection par critère. On peut lire ce chemin xpath ainsi : "renvoie les noeuds du Chemin tel que ce critère est vrai".

### Critères d'attributs

On peut sélectionner tous les nœuds de l'arbre ayant/n'ayant pas un certain attribut:

//\*[@format] sélectionne tous le nœuds qui ont un attribut format.

//\*[not(@format)] sélectionne tous les nœuds de l'arbre qui n'ont pas d'attribut format.

//bagage/\*[@\*] sélectionne tous les nœuds fils de bagage qui ont un attribut (quel qu'il soit)

//bagage/\*[not(@\*)] sélectionne tous les nœuds fils de bagage qui n'ont aucun attribut

//type[@couleur='noir'] sélectionne tous les noœuds type qui ont un attribut couleur de valeur noir

# Critères d'ordre

/bagagesPassagers/bagage[1] sélectionne le premier nœud bagage.

//bagage[2] sélectionne le deuxième nœud bagage.

/bagagesPassagers/bagage[last()] sélectionne le dernier nœud bagage.

# 3.c Expressions arithmétiques

On peut également utiliser des expressions arithmétiques pour les sélections:

| Opérations                |                                       |  |
|---------------------------|---------------------------------------|--|
| +                         | addition                              |  |
| _                         | soustraction                          |  |
| *                         | multiplication                        |  |
| div                       | division (attention, ce n'est pas /)  |  |
| mod                       | modulo (reste de la division entière) |  |
| =                         | égalité                               |  |
| !=                        | inégalité                             |  |
| Opérateurs Booléens       |                                       |  |
| and                       | et puis                               |  |
| or                        | ou bien                               |  |
| not()                     | négation                              |  |
| Opérateurs de comparaison |                                       |  |
| <                         | strictement inférieur à               |  |
| <=                        | inférieur ou égal à                   |  |
| >                         | strictement supérieur à               |  |
| >=                        | supérieur ou égal à                   |  |

En XPath, une expression est vraie s'il a au moins un nœud sélectionné. Par exemple, l'expression //type[@couleur='noir'] sélectionne tous les nœuds type dont l'attribut couleur a la valeur noir. L'expression //type[@couleur!='noir'] sélectionne tous les nœuds type dont l'attribut couleur n'a pas la valeur noir. Ces deux expressions sont *vraies*. En effet, il existe un nœud dont l'attribut couleur a la valeur noir et un nœud dont l'attribut couleur n'a pas la valeur noir.

### 3.d Fonctions

## Fonctions sur les chaînes de caractères

| <pre>string-length(s: string) : int  concat(s1: string, s2: string,): string</pre> | renvoie la longueur (en nombre de caractères) de la chaîne de caractères s  Exemple string-length('bonjour') renvoie 7.  renvoie la chaîne composée de s1, s2, à la suite  Exemple: concat('bonjour', 'toi') renvoie la chaîne 'bonjourtoi'. |
|--|--|
| contains(s1: string, s2: string): boolean  | renvoie vrai si la chaîne s1 contient la chaîne s2 Exemple: contains('bonjour', 'njo') renvoie true.   |
| starts-with(s1: string, s2: string): boolean                                       | renvoie vrai si la chaîne s1 commence par la chaîne s2 Exemple: starts-with('bonjour', 'njo') renvoie false.   |
| <pre>substring(s: string, d:int, [, f:int]): string</pre>                          | renvoie la sous-chaîne de s à partir de l'indice d (éventuellement jusqu'à l'indice f)  Exemple: substring('bonjour', 3) renvoie 'jour'.   |

### Fonctions sur les node-set

| count([node-set]) : int  | nombre de nœuds [contenu dans le node-set]. |
|--------------------------|---|
| last(): int              | taille du node-set courant.                 |
|                          |   |
| name([node-set]): string | Nom du nœud courant [1er du node-set].      |
| position(): int          | Position du nœud courant.                   |
| position(): int          | 1 Osition du nœdd courant.                  |
|                          |   |

# 4 Ecrire correctement du XPath

Vous connaissez la syntaxe XPath ? C'est bien ! Vous découvrirez vite cependant que l'écriture de chemins complexes en XPath peut vite devenir un casse-tête sans fin.

Pour éviter au maximum de se retrouver dans une situation inextricable, il FAUT que vous travailliez méthodiquement, à savoir :

# 4.a Formuler correctement le problème

Prenons tout de suite un exemple avec le document XML suivant qu stocke les données de la bagagerie d'un aeroport pour les bagages hors format :

```
<?xml version="1.0" encoding="UTF-8"?>
   <bagagerie</pre>
2
       xmlns:xsi='http://www.w3.org/2001/XMLSchema-instance'
3
       xmlns='http://www.timc.fr/nicolas.glade/bagages'
4
       xsi:schemaLocation='http://www.timc.fr/nicolas.glade/bagages bagagerieHF.
5
           xsd'>
       <passagers>
6
           <passager nom="Gator" prenom="Ali" id="AG18300427"></passager>
           <passager nom="Nancière" prenom="Sophie" id="SF23890311"></passager>
           <passager nom="Kass" prenom="Oscar" id="OK17401023"></passager>
       </passagers>
10
       <baggages>
11
           <bagge idPassager="OK17401023">
12
                <ref>AF677793</ref>
13
                <type format="autre" couleur="blanc"/>
14
                <poids>12.1</poids>
15
                <description>skis</description>
16
           </bagage>
17
           <bagge idPassager="OK17401023">
                <ref>AF677793</ref>
19
                <type format="autre" couleur="bleu"/>
20
                <poids>8.7</poids>
21
                <description>skis</description>
22
           </bagage>
23
           <bagge idPassager="AG18300427">
24
                <ref>AF48717</ref>
25
                <type format="autre" couleur="marron"/>
26
27
                <poids>6.4</poids>
                <description>tam-tam</description>
           </bagage>
           <bagge idPassager="SF23890311">
30
                <ref>AF67840</ref>
31
                <type format="autre" couleur="blanc"/>
32
                <poids>10.4</poids>
33
                <description>skis</description>
34
           </bagage>
35
       </bagages>
36
   </bagagerie>
37
```

et le schéma associé :

```
<?xml version="1.0"?>
   <schema version="1.0"
2
           xmlns="http://www.w3.org/2001/XMLSchema"
3
           xmlns:bg="http://www.timc.fr/nicolas.glade/bagages"
4
           targetNamespace="http://www.timc.fr/nicolas.glade/bagages"
5
           elementFormDefault="qualified">
6
       <include schemaLocation="bagages.xsd"/>
7
       <element name="bagagerie" type="bg:Bagagerie"/>
9
10
       <complexType name="Bagagerie">
11
```

```
12
            <sequence>
                <element name="passagers">
13
                     <complexType>
14
                         <sequence>
15
                              <element name="passager" type="bg:Passager" maxOccurs="</pre>
16
                                  unbounded"/>
                          </sequence>
17
                     </complexType>
18
                </element>
19
                <element name="bagages">
20
                     <complexType>
21
                         <sequence>
22
                              <element name="bagage" type="bg:BagageHF" maxOccurs="</pre>
23
                                  unbounded"/>
                          </sequence>
24
                     </complexType>
25
                </element>
26
            </sequence>
27
        </complexType>
28
        <complexType name="Passager">
30
            <attribute name="nom" type="string" use="required"/>
31
            <attribute name="prenom" type="string" use="required"/>
32
            <attribute name="id" type="bg:IDPassager" use="required"/>
33
        </complexType>
34
35
36
        <complexType name="BagageHF">
37
            <complexContent>
38
                <extension base="bg:Bagage">
39
                     <sequence>
40
                         <element name="description" type="string"/>
41
                     </sequence>
                     <attribute name="idPassager" type="bg:IDPassager"/>
43
                </extension>
44
            </complexContent>
45
        </complexType>
46
47
        <simpleType name="IDPassager">
48
            <restriction base="string">
49
                <pattern value="[A-Z]{2}[0-9]{8}"/>
50
            </restriction>
51
        </simpleType>
52
53
   </schema>
```

Supposons que l'on souhaite connaître le nom de famille de tous les passagers dont les bagages hors format sont des skis.

D'abord, on commence par reformuler cette phrase sous cette forme: [Ce qu'on veut] [tel que] [la ou les conditions]. Cela donne dans notre cas: [Le nom de famille des passagers] [tel que] [les bagages sont des skis].

Mais si l'on regarde mieux, on voit que ce que l'on doit faire, c'est relier l'identifiant des bagages qui sont des skis, à l'identifiant passager. Donc on peut reformuler la condition ainsi : [les identifiants des bagages qui sont des skis]  $\Rightarrow$  [les identifiants des bagages] [tel que] [ce sont des skis]

On voit donc qu'on va devoir mettre en relation un 1er chemin (les noms des passagers) avec un deuxième chemin (les identifiants des bagages)

## 4.b Exprimer ce que l'on veut et ce que sont les conditions

Que veut-on finalement ?  $\Rightarrow$  les noms des passagers ... donc on commence par écrire le chemin XPath qui donne les noms de tous les passagers : //passager/@nom. Comme attendu, on obtient la liste de tous les noms du document : { Gator, Nancière, Kass}

Pour les bagages, que veut-on ?  $\Rightarrow$  les identifiants des bagages ... donc on écrit cela sous la forme d'un autre chemin XPath : //bagage/@idPassager . Cette fois, on obtient la liste de tous les identifiants passagers dans les bagages : { OK17401023, OK17401023, AG18300427, SF23890311}. Vous remarquerez que l'on peut obtenir plusieurs fois le même identifiant.

Ensuite on se pose la question de comment réunir les deux. Ce qui lie les bagages aux passagers, ce sont les identifiants des passagers. Mais avant de réunir les deux chemins, on va tester des choses simples, par exemple des identifiants fixes. Reprenons les noms des passagers et ajoutons une sélection liée à l'identifiant : je veux les noms des passagers dont l'identifiant est X : //passager[@id='AG18300427']/@nom ce qui se lit les noms des passagers tel que l'identifiant passager est AG18300427. On obtient la liste contenant 1 seul nom (mais il s'agit bien d'une liste) : {Gator}

De même, pour les bagages, récupérons les identifiants des bagages qui sont des skis : //bagage[contains(./description/text(),'skis')]/@idPassager . On obtient la liste d'identifiant suivante : { OK17401023, OK17401023, SF23890311 }

Ensuite seulement, on peut relier les 2 chemins :

//passager[@id=//bagage[contains(./description/text(),'skis')]/@idPassager]/@nom, ce qui renvoie bien la liste de noms attendus { Nancière, Kass} qui sont les deux seules personnes à partir au ski en avion.

# 4.c Procéder par petites étapes qu'on teste

De la même manière qu'on n'écrit pas un programme entier avant de le tester, on n'écrit pas une ligne complexe de XPath sans procéder par petites étapes. Chacune des étapes écrite au dessus doit être testée individuellement soit avec un outil en ligne (par exemple <a href="http://xpather.com/">http://xpather.com/</a> qui marche bien) soit avec un programme appelant en Java ou autre.

Autrement dit, vous devez tester dans l'ordre les expressions suivantes :

- //passager/@nom
- //bagage/@idPassager
- //passager[@id='AG18300427']/@nom
- //bagage[contains(./description/text(),'skis')]/@idPassager
- //passager[@id=//bagage[contains(./description/text(),'skis')]/@idPassager]/@nom