

# Séance 6

## Parsing des documents XML

**Prof. Yassin Aziz REKIK**  
**[Yassin.rekik@he-arc.ch](mailto:Yassin.rekik@he-arc.ch)**

# Manipulation des documents XML



- **Introduction**

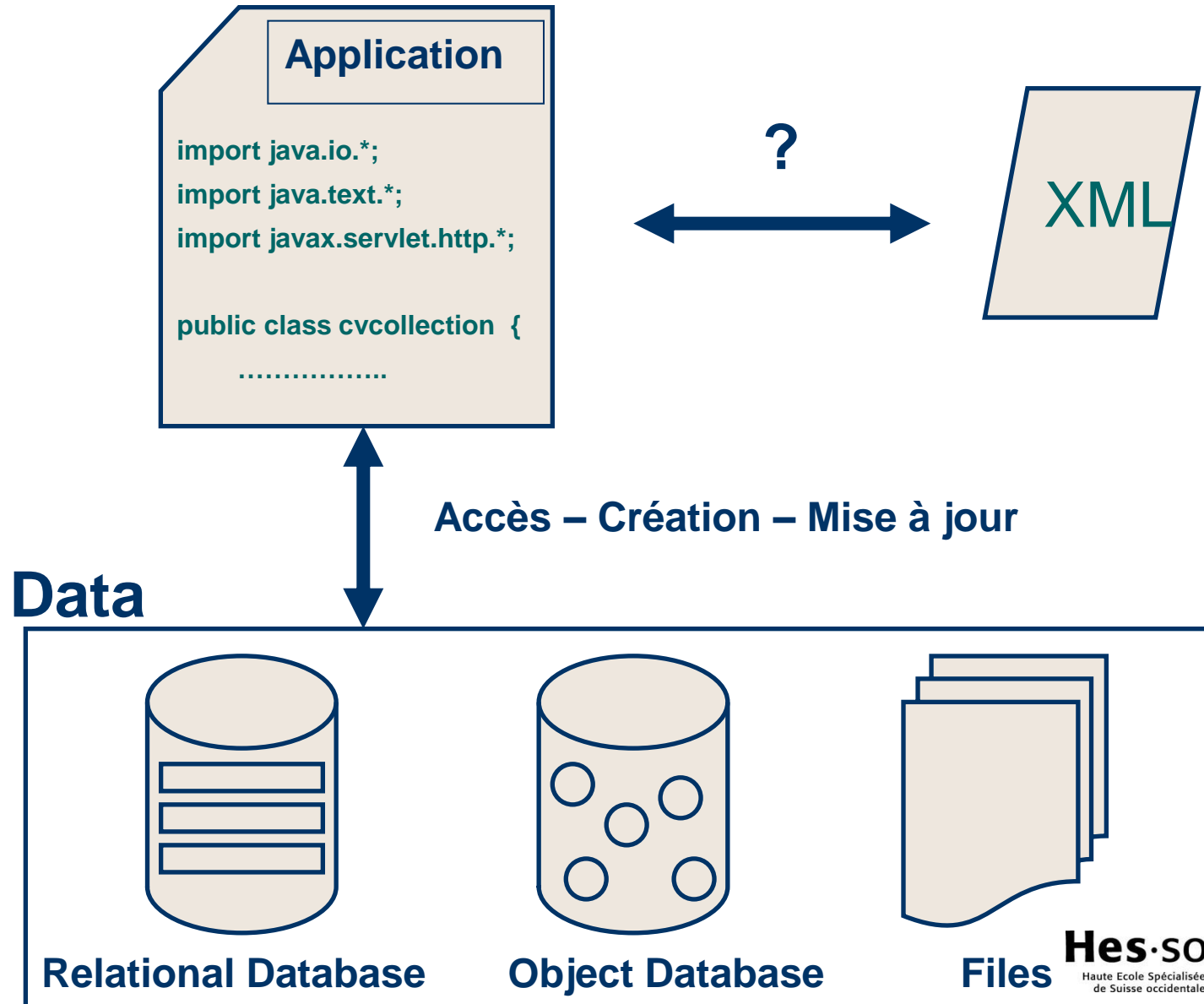
- **Parsing SAX**

- **Exercices simples**

- **Parsing DOM**

- **Exercices simples**

- **Outils**



# Deux approches de traitement



- **Introduction**

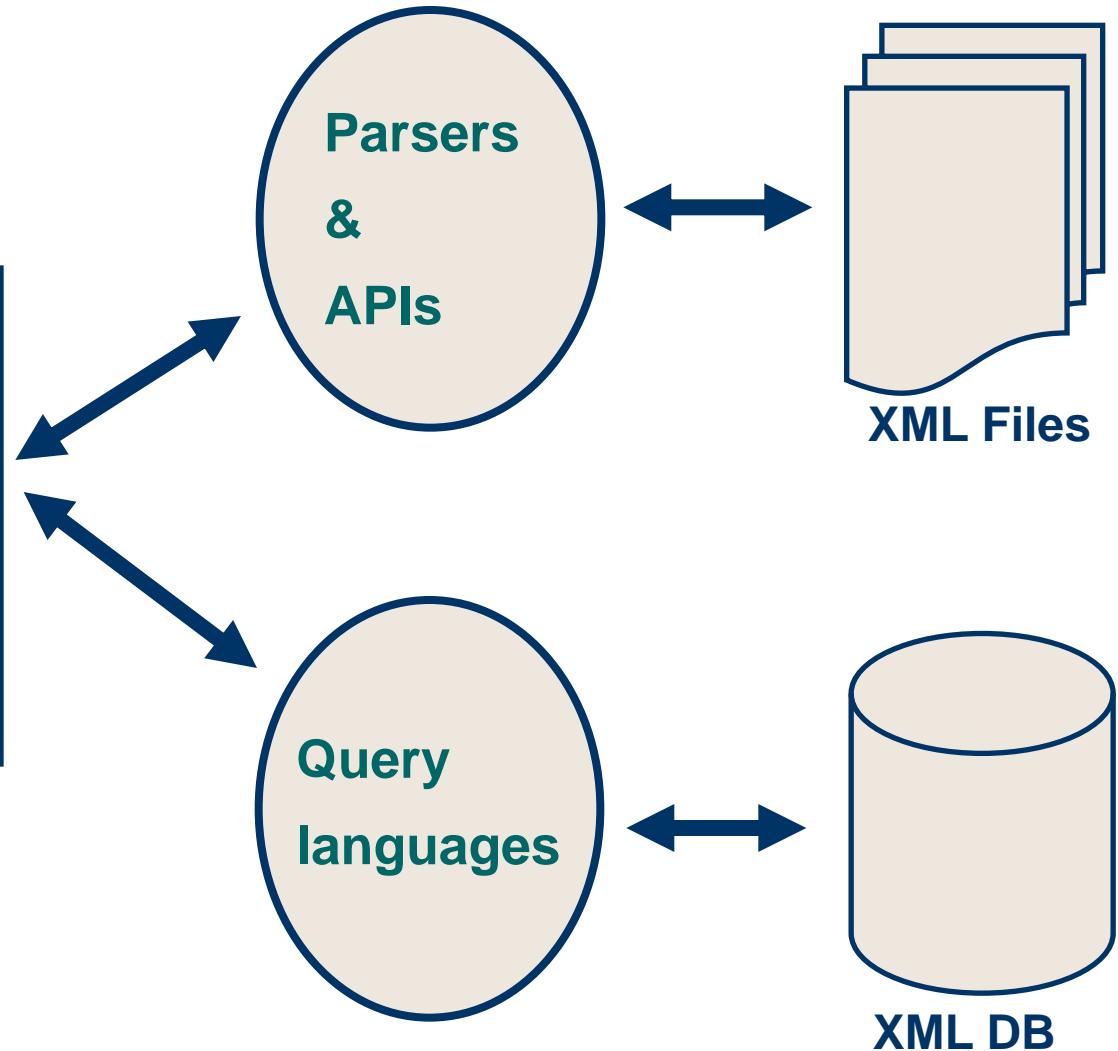
- **Parsing SAX**

- **Exercices simpl**

- **Parsing DOM**

- **Exercices s**

- **Outils**





- **Introduction**

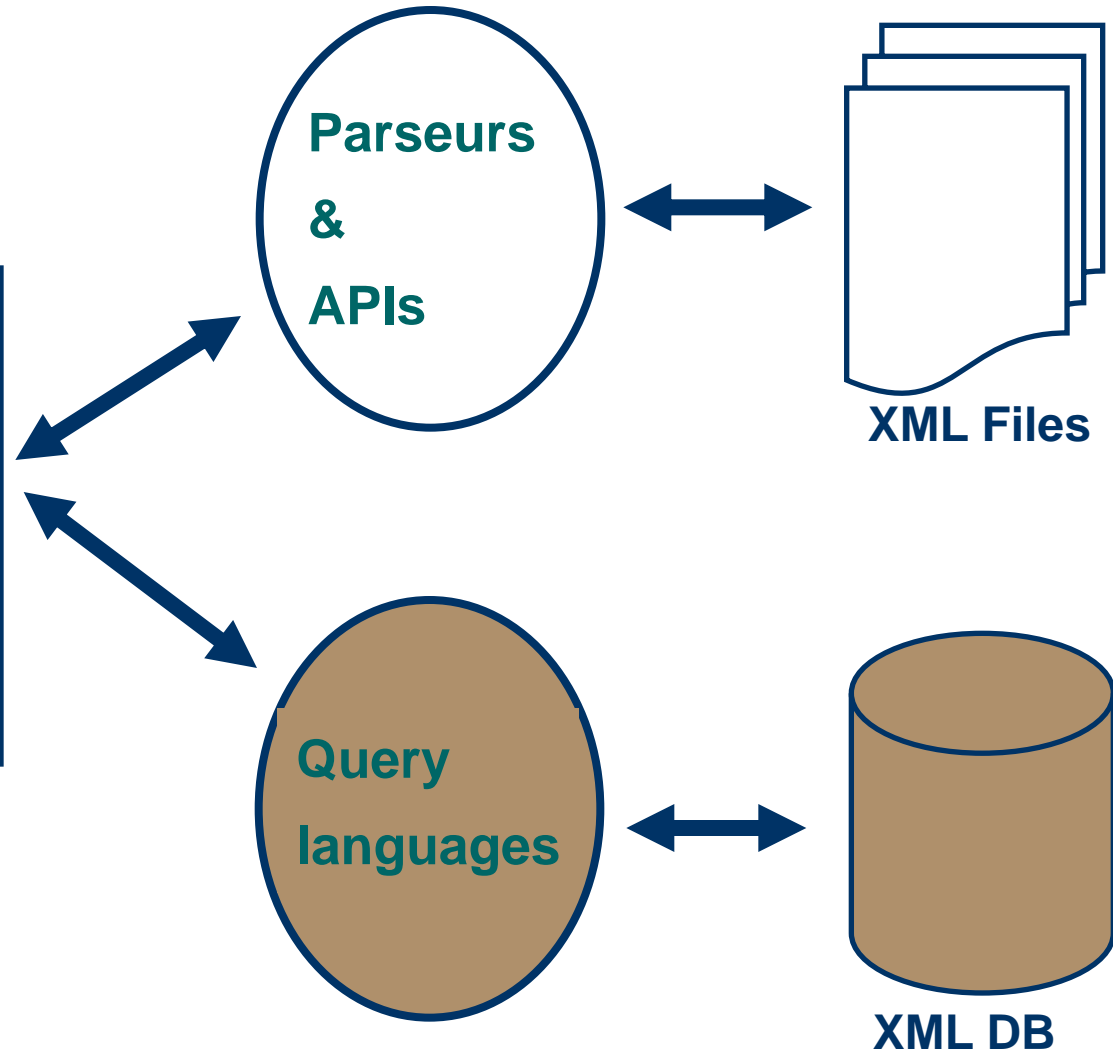
- **Parsing SAX**

- **Exercices simpl**

- **Parsing DOM**

- **Exercices s**

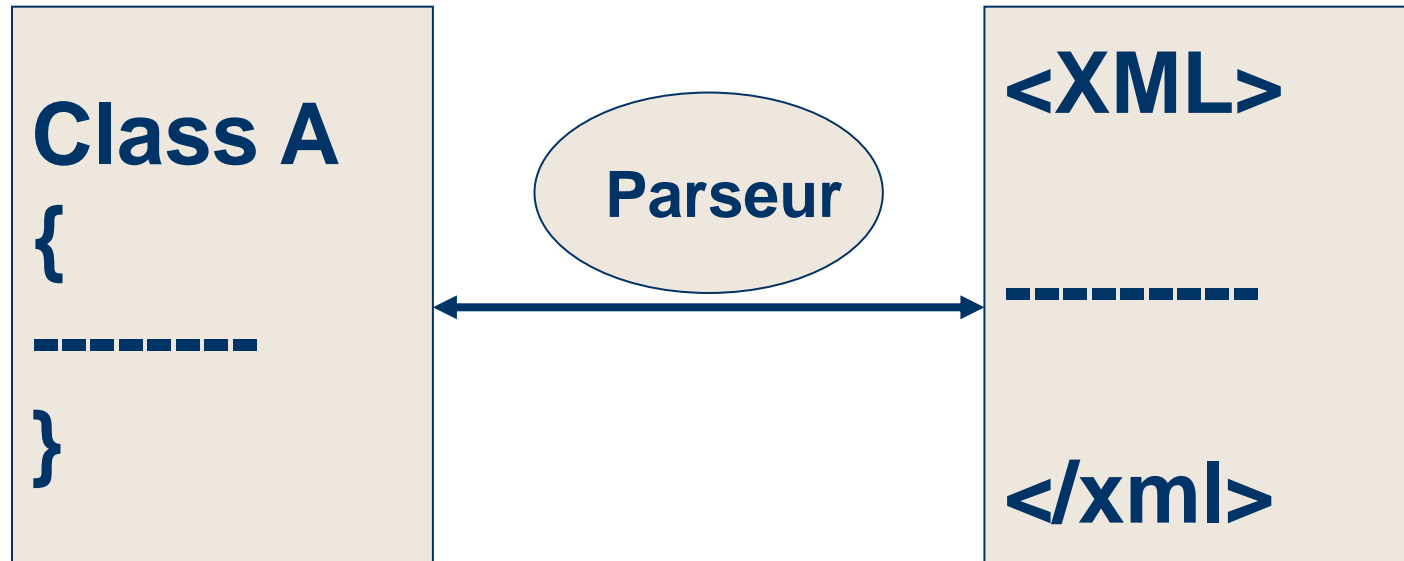
- **Outils**





- **Introduction**
- **Parsing SAX**
- **Exercices simples**
- **Parsing DOM**
- **Exercices simples**
- **Outils**

- **Un parseur XML est un outil (package, librairie, programme, ...) permettant d'accéder, de créer et de mettre à jour des documents XML à travers un code.**





- **Introduction**

- **Parsing SAX**

- **Exercices simples**

- **Parsing DOM**

- **Exercices simples**

- **Outils**

- **Un parseur doit assurer les fonctions suivantes:**

- Analyse lexicale et syntaxique
- Vérification de la conformité du document à la norme
- Substitution des entités
- Validation par rapport à une DTD ou un schéma
- Manipulation de base : modifier, créer, supprimer, ...
- Etc. ....

**Un parseur est plus qu'un compilateur**



- **Introduction**

- **Parsing SAX**

- **Exercices simples**

- **Parsing DOM**

- **Exercices simples**

- **Outils**

- **Deux modèles de parsing XML**

- Modèle événementiel - Event-based model
- Modèle par compilation ou basé arbre - Tree-based model

- **Modele événementiel**

- Standard de facto : SAX (Simple API for XML)

- **Modèle arborescent**

- W3C standard : DOM (Document Object Model)

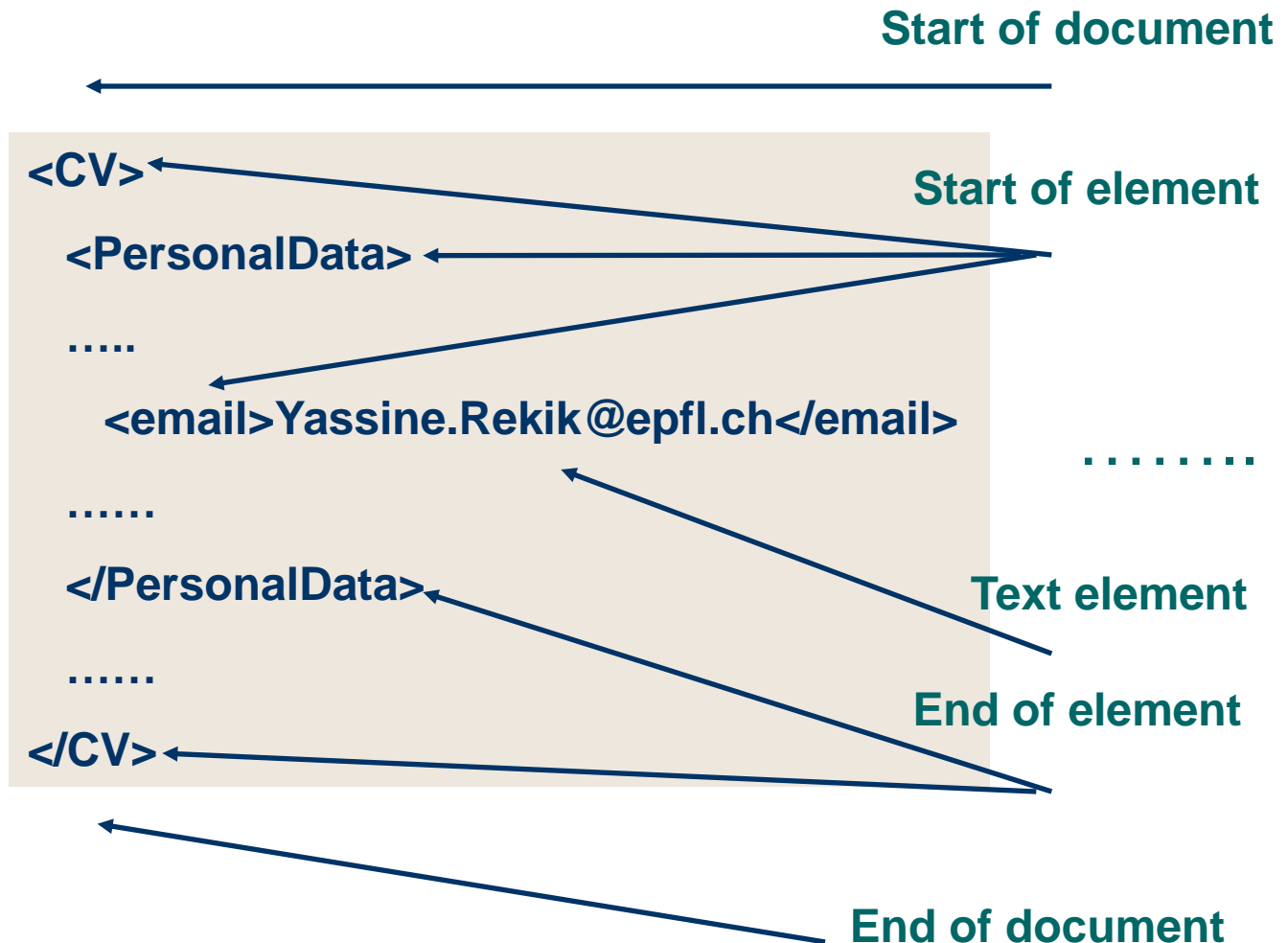
- **Les parseurs SAX et DOM existent pour presque tous les langages de programmation**

- Java , C , C++ , Perl , TCL Tk , VB , ..... , Prolog



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Modèle événementiel**







- Introduction

- Parsing SAX

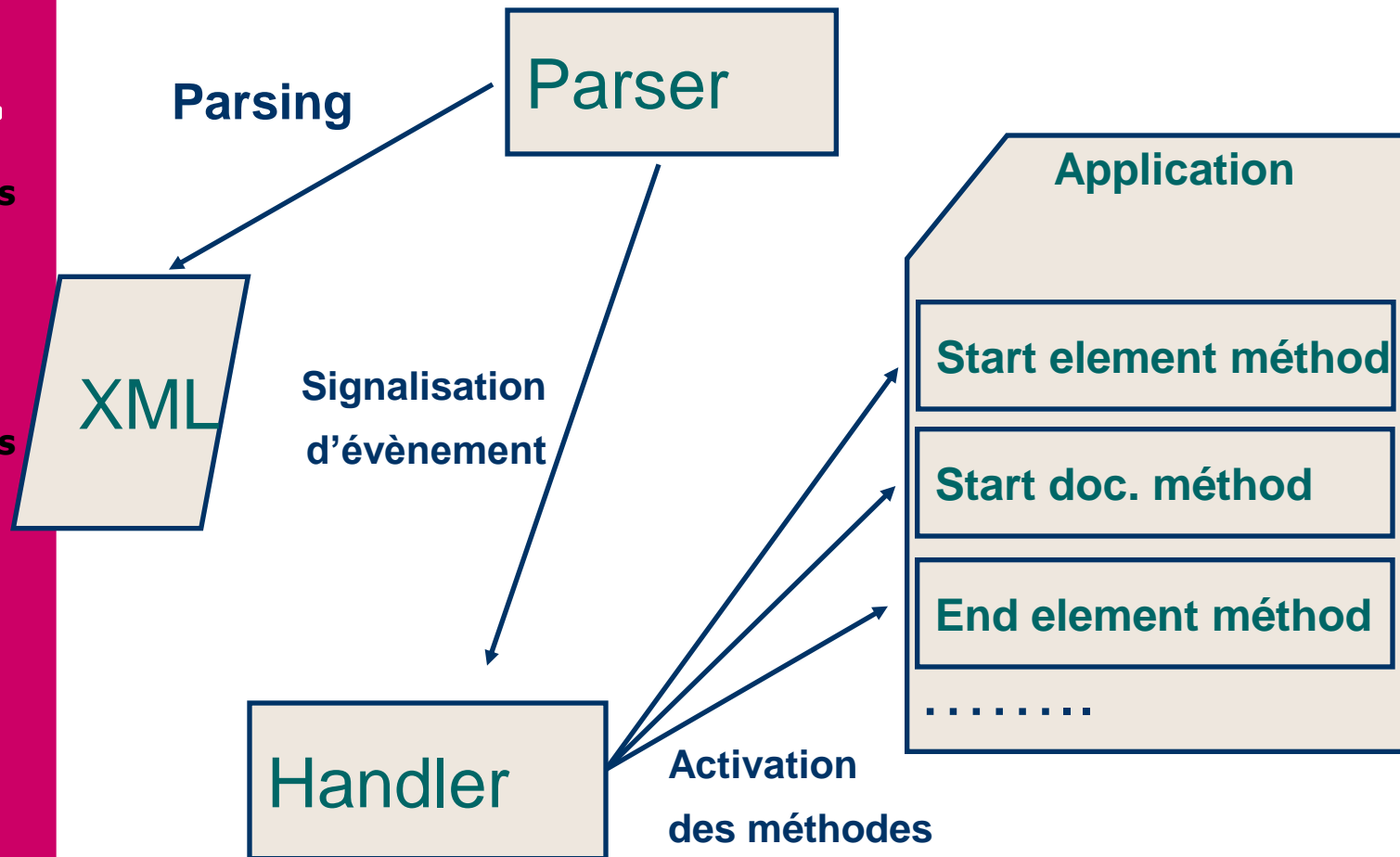
- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

- **Modèle événementiel**





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Applications orientées SAX**

- Décomposition de document
- Translation de document
- Extraction de fragments
- ....

## Accès et manipulation séquentiels



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Versions du SAX**
  - SAX 1 : introduit en Mai 1998
  - SAX 2.0 : introduit en Mai 2000, il supporte :
    - *Les namespaces*
    - *Des filtres de chaînes*
    - *Positionner et interroger des paramètres sur le parsing*



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Implémentation en Java**
  - Extension de la classe HandlerBase

```
Import org.xml.sax.*;  
  
import org.xml.sax.helpers.ParserFactory;  
  
Public class SaxApplication extends HandlerBase  
{  
  
    public static void main(String args[]) {  
        }  
  
}
```



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Implémentation en Java**

- Création du parseur

```
public static void main(args[]) {  
    String parserName =  
        "org.apache.xerces.parsers.SAXParser";  
    try {  
        SaxApplication app = new  
SaxApplication();  
        Parser parser =  
  
        ParserFactory.makeParser(parserName);  
        parser.setDocumentHandler(app);  
        parser.setErrorHandler(app);  
        parser.parse(new InputSource(args[0]));  
    } catch (Throwable t) {  
        // Handle exceptions  
    }  
}
```



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Implémentation en Java**
  - Redéfinir, par surcharge, les évènements voulus

```
Public class SaxApplication extends HandlerBase {  
    public void main (String args[]) {  
        // stuff missing  
    }  
  
    public void startElement(String name,  
AttributeList attrs) {  
        // Process this element  
    }  
}
```



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Implémentation en Java**
- **Surcharger les autres évènements si besoin**
  - startDocument()
  - endDocument()
  - startElement()
  - endElement()
  - Characters
  - error()
  - ...



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Exemples avec les CVs:**
  - Décomposer le documents CVs
  - Compter le nombre de candidats mariés
  - Compter le nombre de candidats femmes
  - Afficher en séquence un CV
  - ...





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Analyse d'un exemple SAX**

**Code SAX**



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

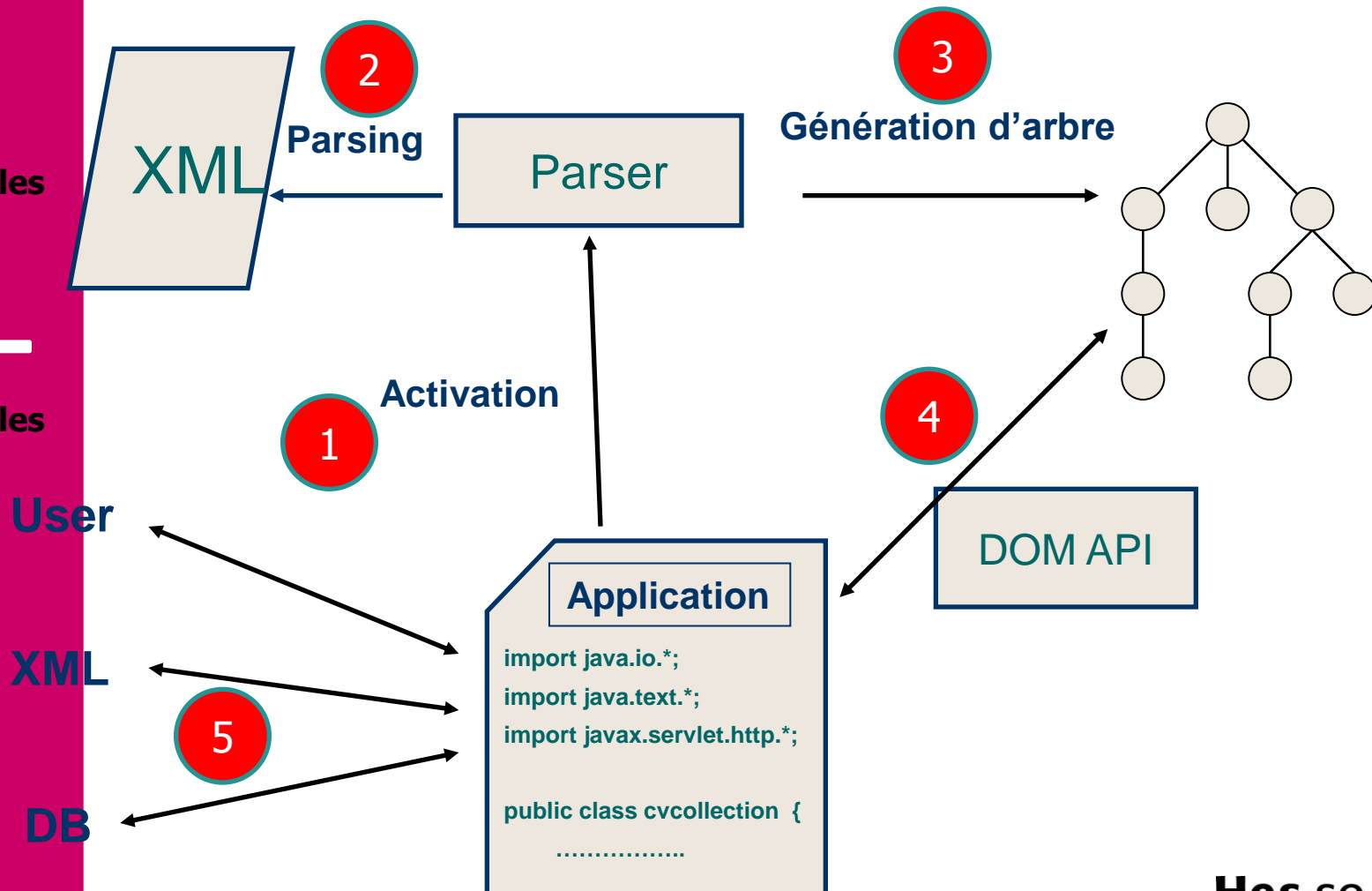
- **Utiliser un parsing SAX pour parser le document «CV.xml » et générer un document « CV.html »**
  - Pour simplifier, vous pouvez générer le html uniquement sur la sortie standard, pas dans un fichier.

# Parseurs et API DOM



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

## Modèle basé arbre



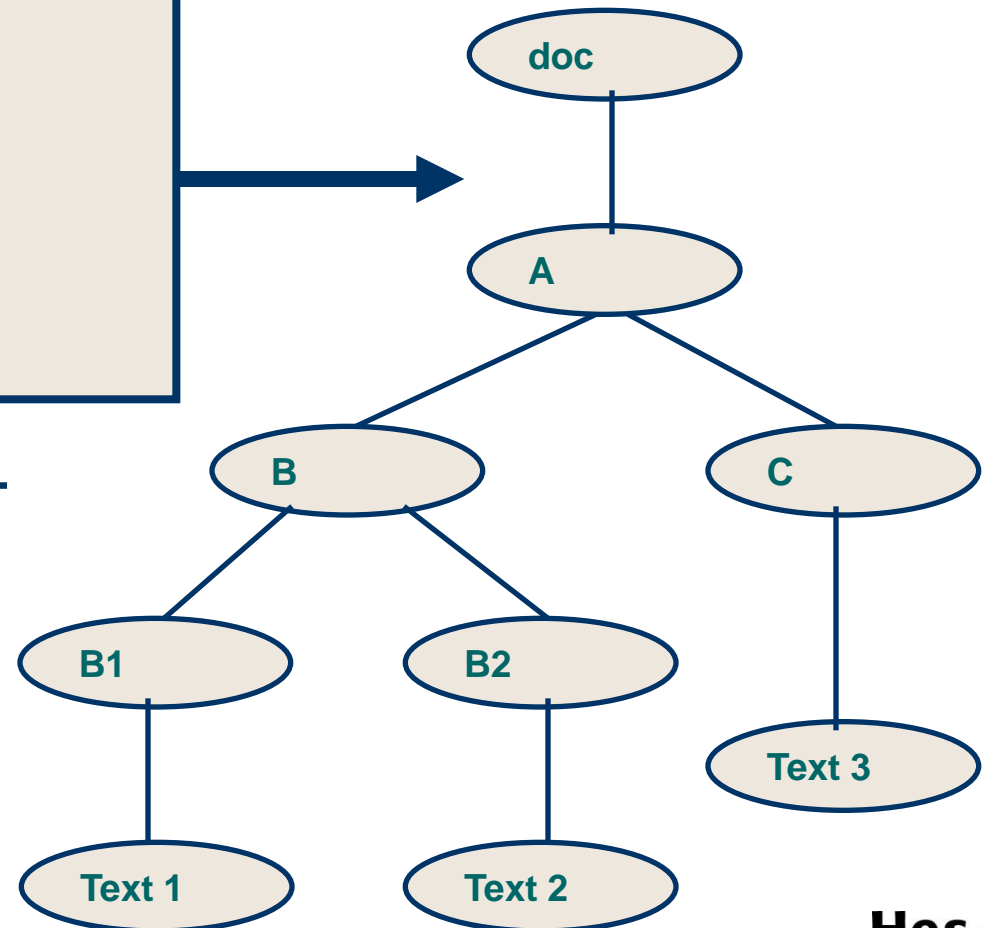


- **Génération de l'arbre DOM**

## Arbre DOM

```
<A>  
  <B>  
    <B1>Text 1</B1>  
    <B2>Text 2</B2>  
  </B>  
  <C>text 3</C>  
</A>
```

## Document XML



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils



- Introduction

- Parsing SAX

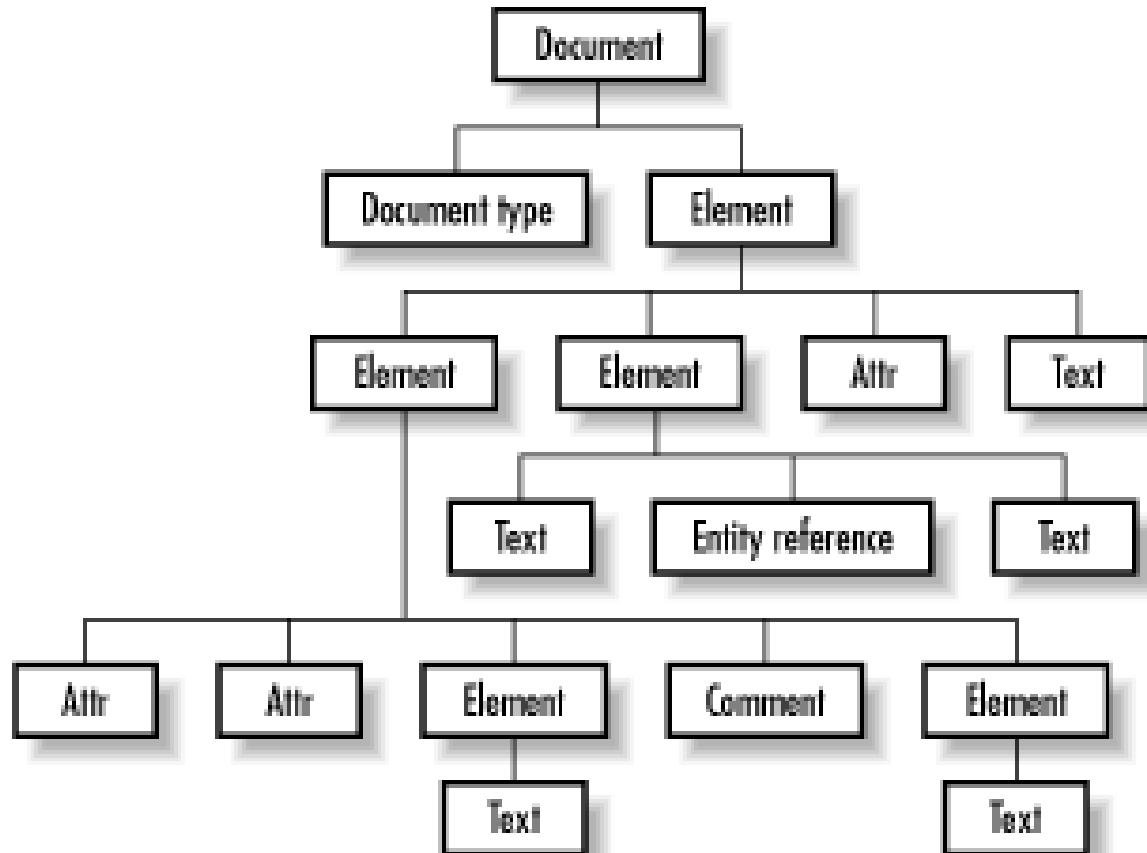
- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

- **Arbre DOM**





- Introduction

- Parsing SAX

- Exercices simples

- Parsing DOM

- Exercices simples

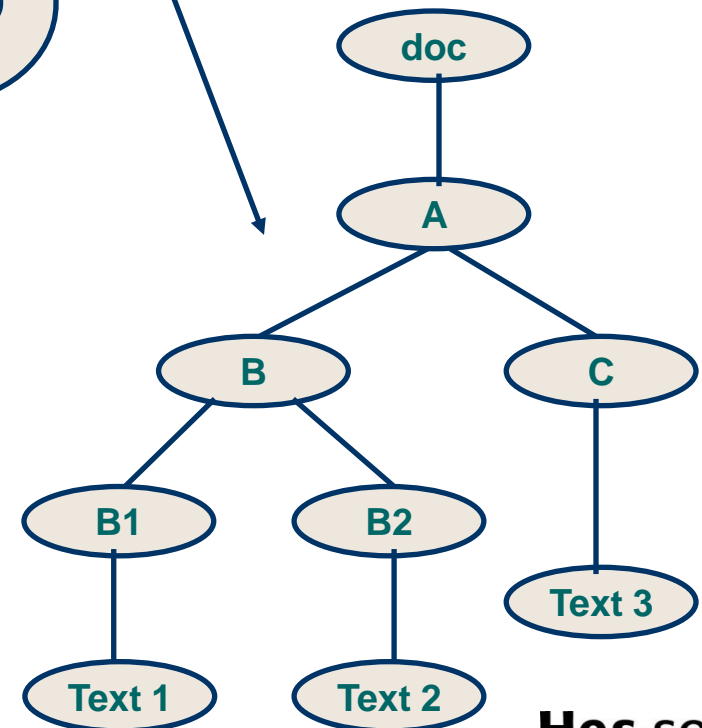
- Outils

- **API DOM**

Application

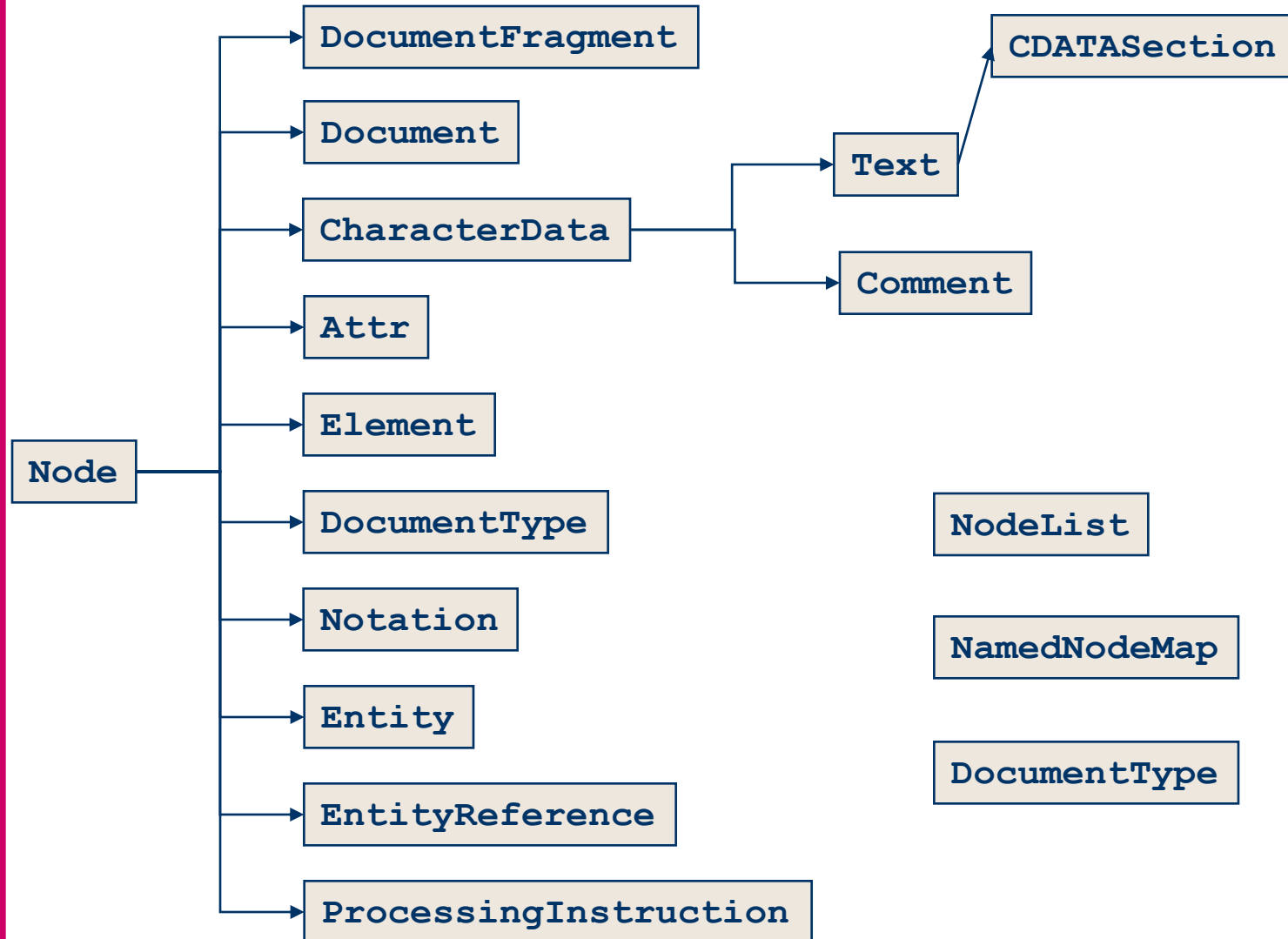
DOM API

Créer un document  
Créer des éléments  
Créer des attributs  
Accès aux éléments  
Accès au textes  
Accès aux attributs  
Supprimer des éléments  
Supprimer des attributs  
Changer des valeurs  
Changer les attributs  
.....  
.....





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils





- Introduction
- Parsing SAX
- Exercices simples
- **Parsing DOM**
- Exercices simples
- Outils

- **Le DOM définit plusieurs interfaces Java**
  - Node : Le type de base de tous les objets DOM
  - Element : Représente les éléments du document
  - Attr : Représente les attributs
  - Text : Le contenu des attributs ou des terminaux
  - Document : Représente le document XML en entier



# DOM : Interface Node



- Introduction

- Parsing SAX

- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

- **Objet de base du DOM**

- **Les nodes décrivent:**

Elements

Attributes

Text

Comments

CDATA sections

Entity declarations

Entity references

Notation declarations

Entire documents

Processing instructions

- **Node collections**

- NodeList, NamedNodeMap, DocumentFragment



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Trois catégories de méthodes:**
  - Node characteristics
    - *name, type, value*
  - Accès relatif et situation dans l'arbre
    - *parents, siblings, children, ancestors, descendants*
  - Modification de la Node
    - *Edit, delete, re-arrange child nodes*

# DOM : Les méthodes Node



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

```
short      getNodeType () ;
```

```
String     getNodeName () ;
```

```
String     getNodeValue ()  
                                     throws DOMException;
```

```
void setNodeValue (String value)  
                                     throws DOMException;
```

```
boolean    hasChildNodes () ;
```

```
NamedNodeMap getAttributes () ;
```

```
Document   getOwnerDocument () ;
```

# DOM : Les types de Node



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

ELEMENT\_NODE = 1

ATTRIBUTE\_NODE = 2

TEXT\_NODE = 3

CDATA\_SECTION\_NODE = 4

ENTITY\_REFERENCE\_NODE = 5

ENTITY\_NODE = 6

PROCESSING\_INSTRUCTION\_NODE = 7

COMMENT\_NODE = 8

DOCUMENT\_NODE = 9

DOCUMENT\_TYPE\_NODE = 10

DOCUMENT\_FRAGMENT\_NODE = 11

NOTATION\_NODE = 12

```
if (myNode.getNodeType() == Node.ELEMENT_NODE) {  
    //process node  
    ...  
}
```

# DOM : Les méthodes de navigation

- Introduction

- Parsing SAX

- Exercices simples

- Parsing DOM

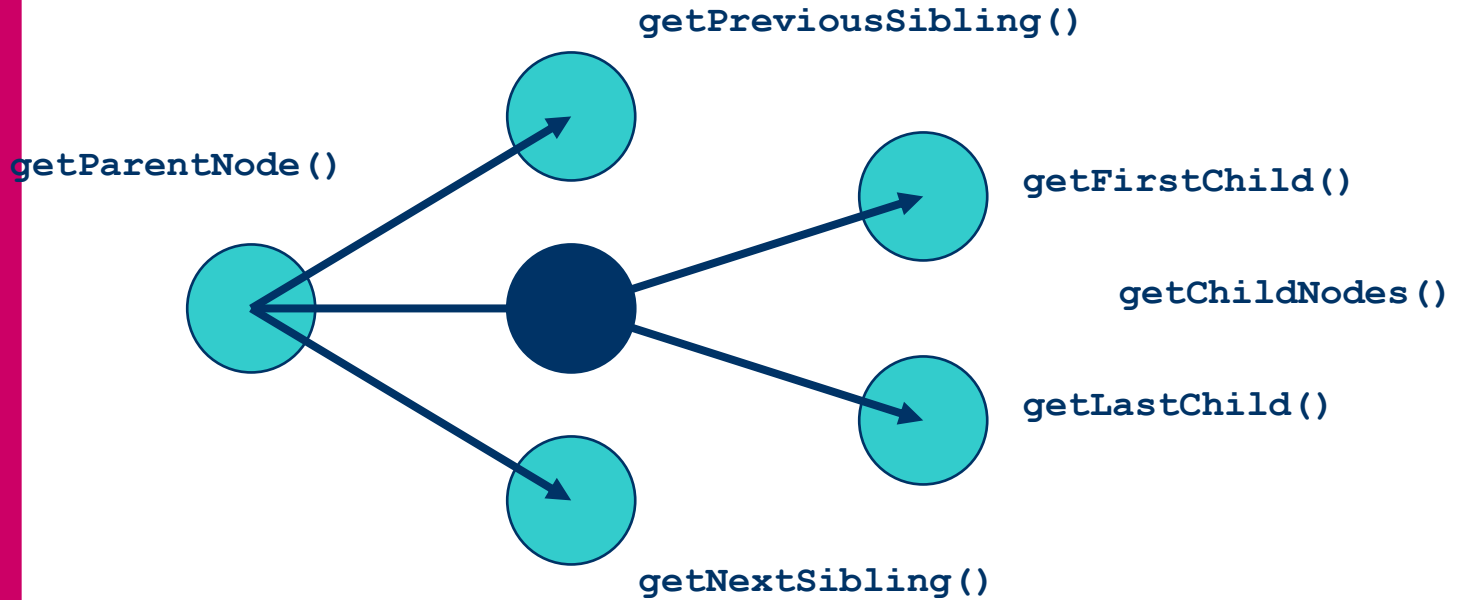
- Exercices simples

- Outils

- Chaque Node possède un emplacement spécifique dans l'arbre DOM
- L'interface Node spécifie des méthodes pour accéder aux environs d'un node :
  - Node      `getFirstChild();`
  - Node      `getLastChild();`
  - Node      `getNextSibling();`
  - Node      `getPreviousSibling();`
  - Node      `getParentNode();`
  - NodeList   `getChildNodes();`

# DOM : Les méthodes de navigation

- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils



```
Node parent = myNode.getParentNode();  
if (myNode.hasChildren()) {  
    NodeList children = myNode.getChildNodes();  
}
```

# DOM : Les méthodes de manipulation

- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Les nodes fils dans un arbre DOM peuvent être manipulés: ajouter, éditer, supprimer, déplacer, copier, etc.**

```
Node removeChild(Node old)
    throws DOMException;
```

```
Node insertBefore(Node new, Node ref)
    throws DOMException;
```

```
Node appendChild(Node new)
    throws DOMException;
```

```
Node replaceChild(Node new, Node old)
    throws DOMException;
```

```
Node cloneNode(boolean deep) ;
```

# DOM : Les méthodes de manipulation

- Introduction

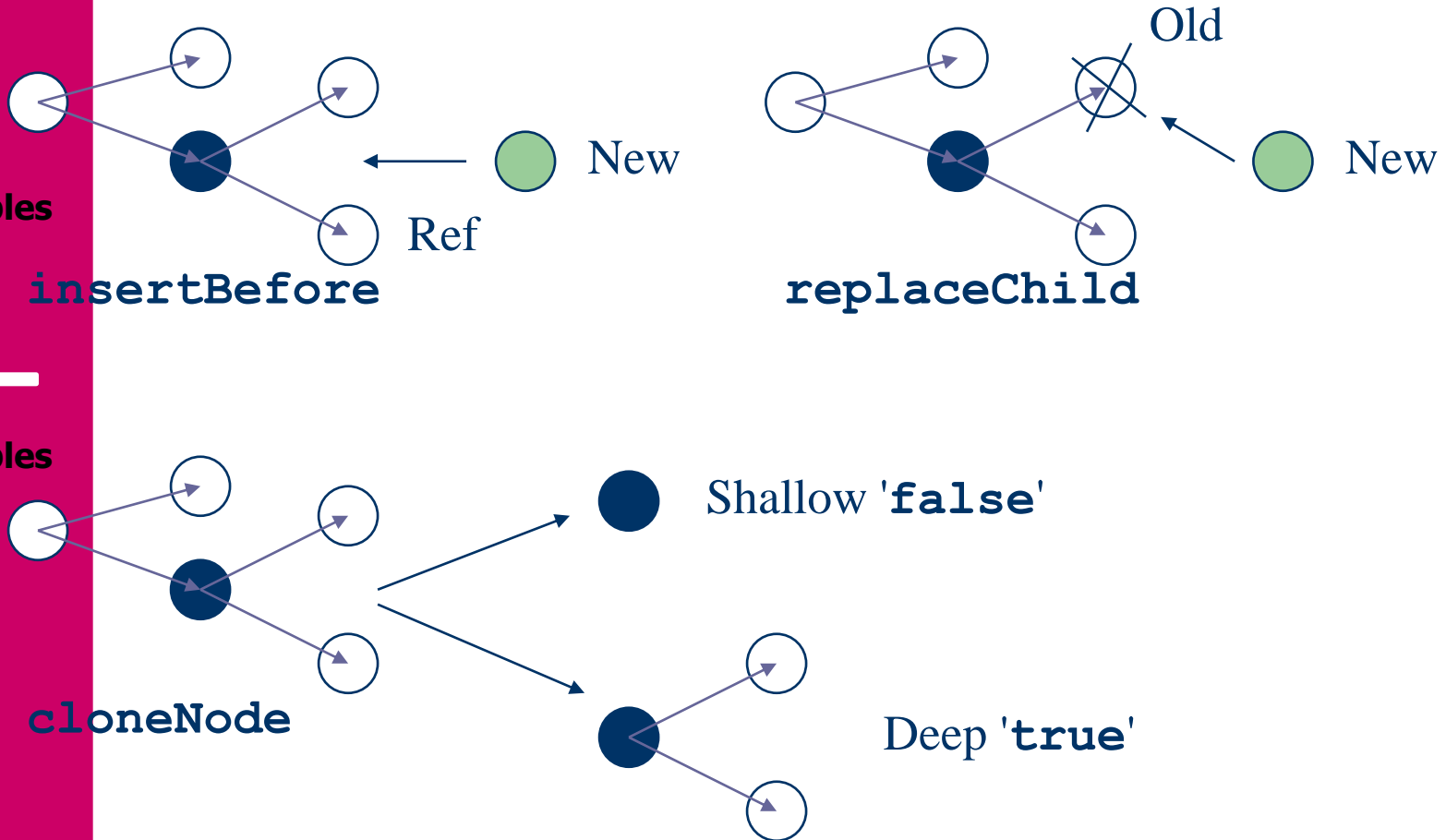
- Parsing SAX

- Exercices simples

- Parsing DOM

- Exercices simples

- Outils





# DOM : Interface Document



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Réprésente le document entier : Root élément**
- **Méthodes**

```
//Information from DOCTYPE - See 'DocumentType'  
DocumentType      getDocumentType() ;
```

```
//Returns reference to root node element  
Element           getDocumentElement() ;
```

```
//Searches for all occurrences of 'tagName' in nodes  
NodeList          getElementsByTagName(String tagName) ;
```



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Méthodes de création de nodes**

```
Element createElement(String tagName)
                        throws DOMException;

DocumentFragment createDocumentFragment();

Text createTextNode(String data);

Comment createComment(String data);

CDATASection createCDATASection(String data)
                        throws DOMException;

ProcessingInstruction createProcessingInstruction(
    String target, String data) throws DOMException;

Attr createAttribute(String name) throws DOMException;

EntityReference createEntityReference(String name)
                throws DOMException;
```



- Introduction

- Parsing SAX

- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

- **Deux catégories de méthodes**

- Méthodes générales sur les éléments

```
String    getTagName() ;  
NodeList getElementsByTagName() ;  
void      normalize() ;
```

- Méthodes pour la gestion des attributs

```
String    getAttribute(String name) ;  
void      setAttribute(String name, String value)  
                                throws DOMException;  
void      removeAttribute(String name)  
                                throws DOMException;  
Attr      getAttributeNode(String name) ;  
void      setAttributeNode(Attr new)  
                                throws DOMException;  
void      removeAttributeNode(Attr old)  
                                throws DOMException;
```



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Interface sur les nodes Attributes**

```
//Get name of attribute  
String getName();
```

```
//Get value of attribute  
String getValue();
```

```
//Change value of attribute  
void setValue(String value);
```

```
//if 'true' - attribute defined in element,  
//else in DTD  
boolean getSpecified();
```



- Introduction

- Parsing SAX

- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

- Permet de manipuler les éléments textuel
- Utilisée sur les élément text et comment

```
String getData()  
    throws DOMException;  
void    setData(String data)  
    throws DOMException;  
int     getLength();  
void    appendData(String data)  
    throws DOMException;  
String  substringData(int offset, int length)  
    throws DOMException;  
void    insertData(int offset, String data)  
    throws DOMException;  
void    deleteData(int offser, int length)  
    throws DOMException;  
void    replaceData(int offset, int length,  
String data)  
    throws DOMException;
```



- Introduction

- Parsing SAX

- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

- **Représente le contenu textuel de** `Element` **ou** `Attr`
  - Généralement, se sont des fils de ces éléments
- **Toujours des objets terminaux**
- **Méthode ajouter à l'interface Character Data**
  - `Text splitText(int offset) throws DOMException`
- **Factory method dans** `Document` **pour la création**
- **La méthode** `normalize()` **sur** `Element` **concatène les** `Text` **objets de ce dernier**



- Introduction

- Parsing SAX

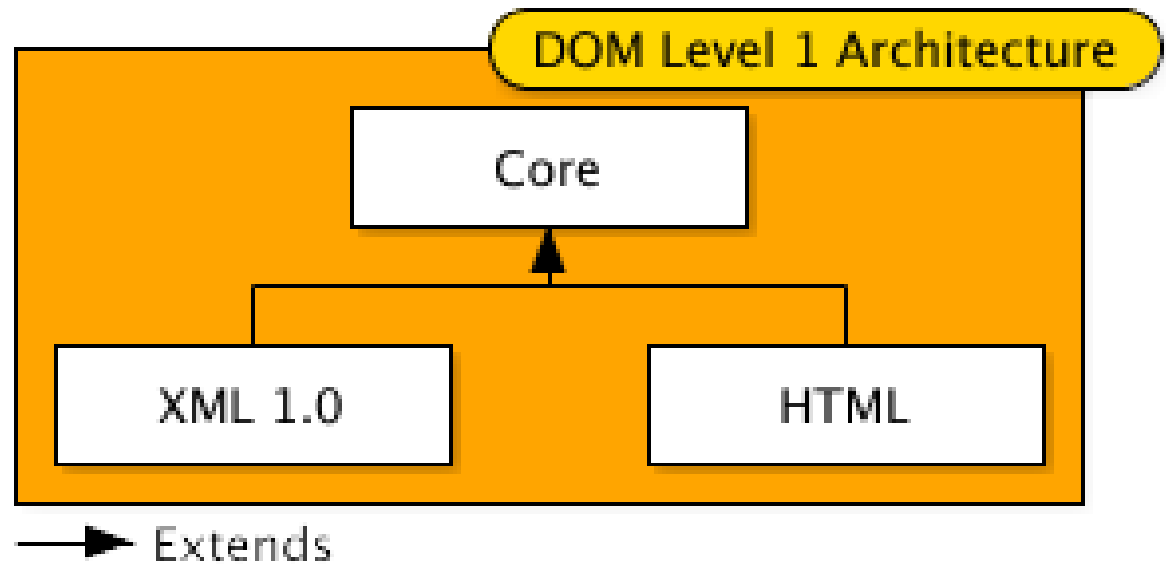
- Exercices simples

- Parsing DOM

- Exercices simples

- Outils

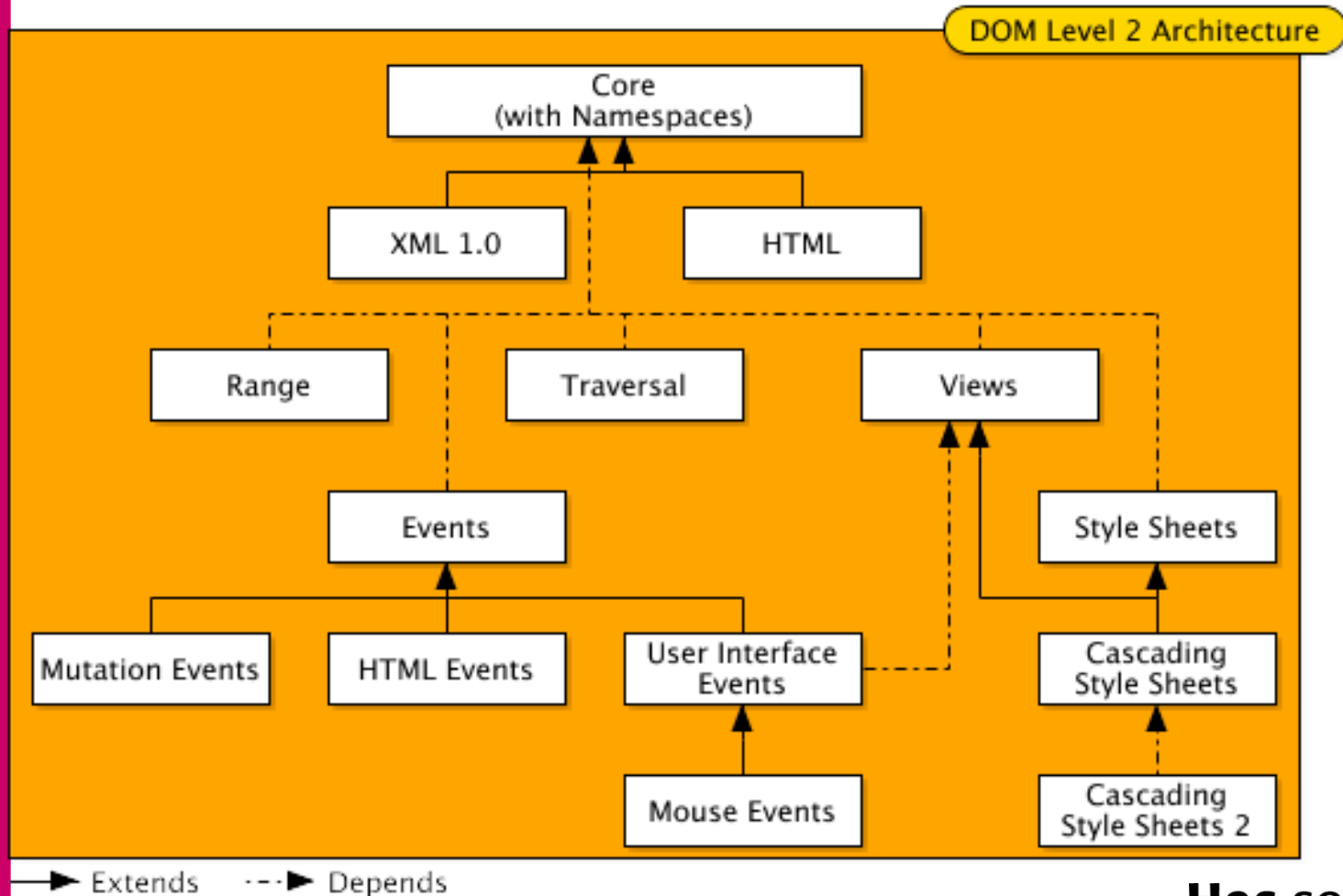
- Evolution du DOM : Level 1





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- Evolution du DOM : Level 2

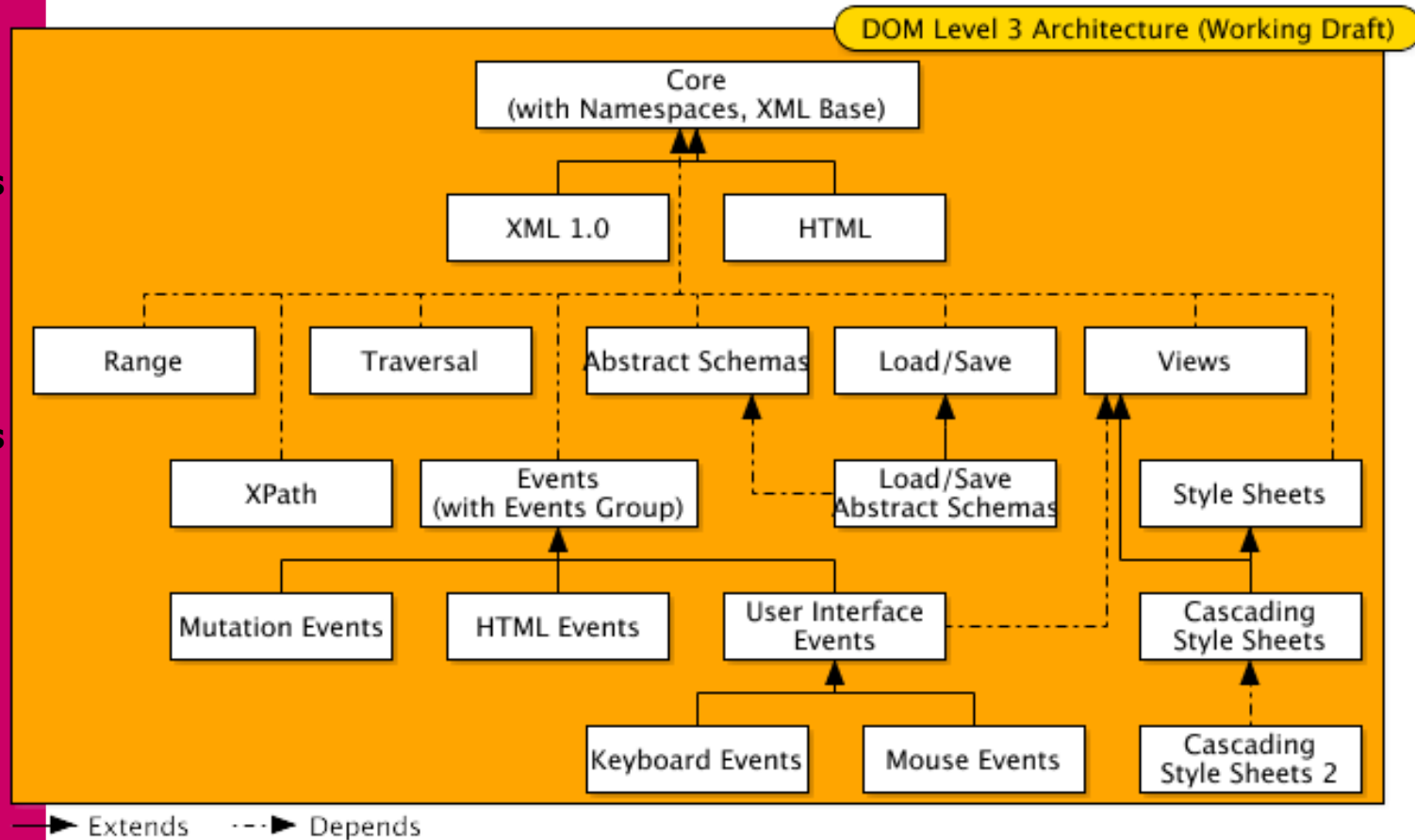






- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- Evolution du DOM : Level 3





- Introduction

- Parsing SAX

- Exercices simples

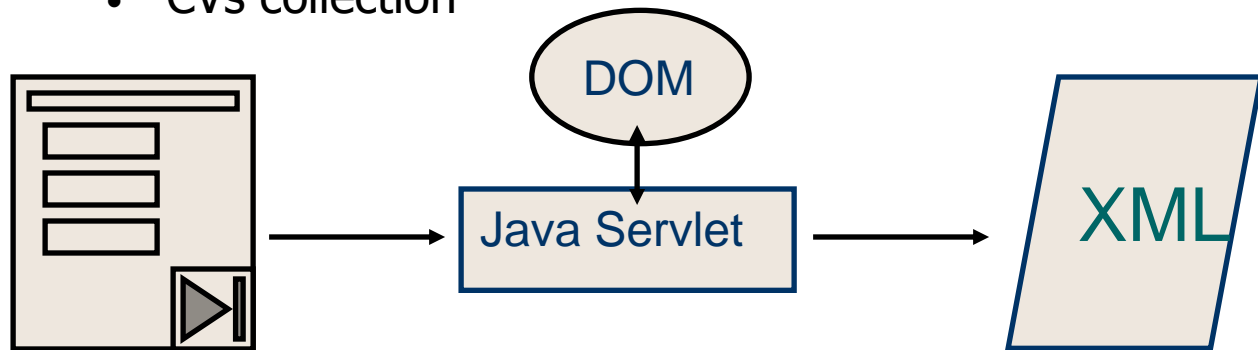
- Parsing DOM

- Exercices simples

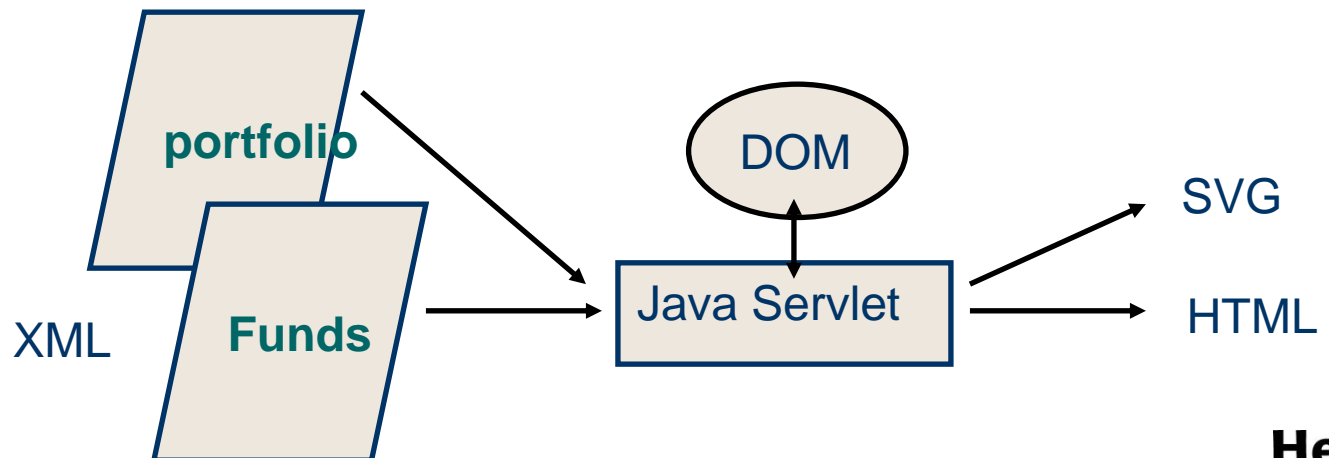
- Outils

- Exemples DOM

- CVs collection



- Portfolio Management





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Analyse d'exemple de code DOM**

**Code DOM**



- Introduction
  - Parsing SAX
  - Exercices simples
  - Parsing DOM
  - Exercices simples
  - Outils
- **Utiliser un parsing DOM pour parser le document «CV.xml » et générer un document « CV.html »**
    - Pour simplifier, vous pouvez générer le html uniquement sur la sortie standard, pas dans un fichier.



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Applications basées sur le Web : croissance**
- **Manipulation des documents XML via le Web**
  - Besoin courant
- **Types de manipulations**
  - Accès et publication : consommation
  - Création et mise à jour : production



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Architecture 1 : statique**

XML



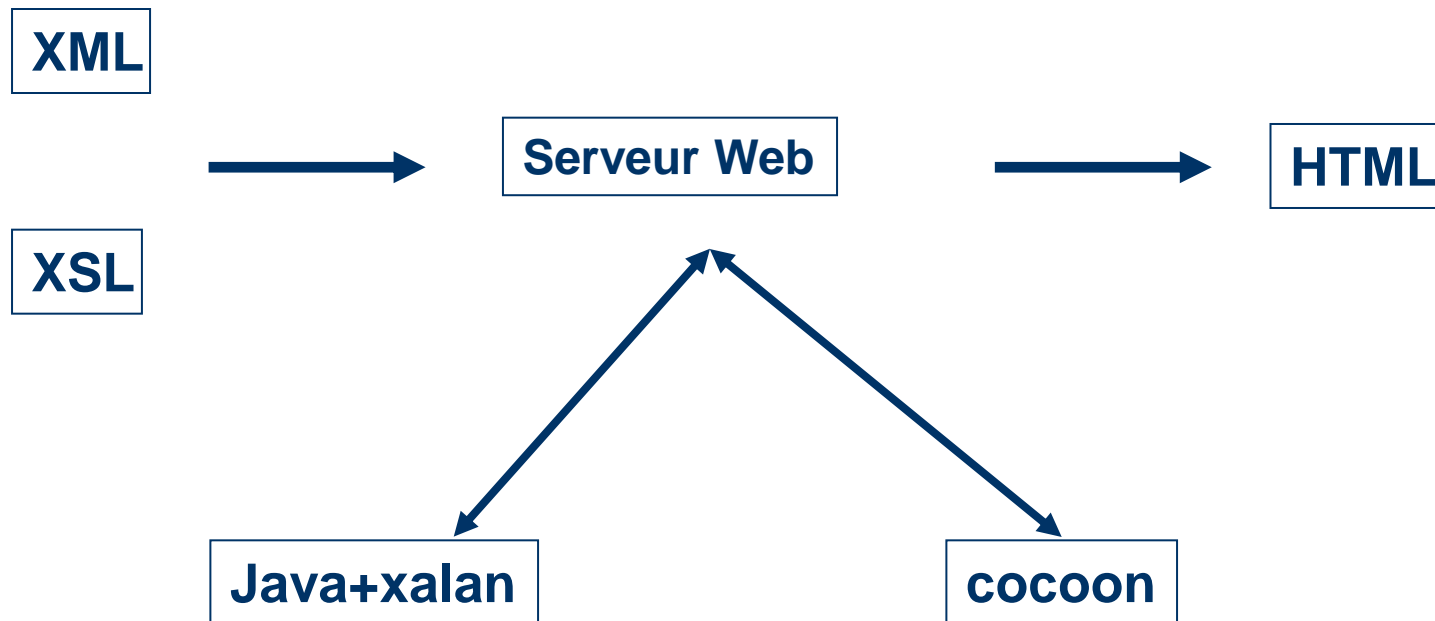
Navigateurs

CSS



- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

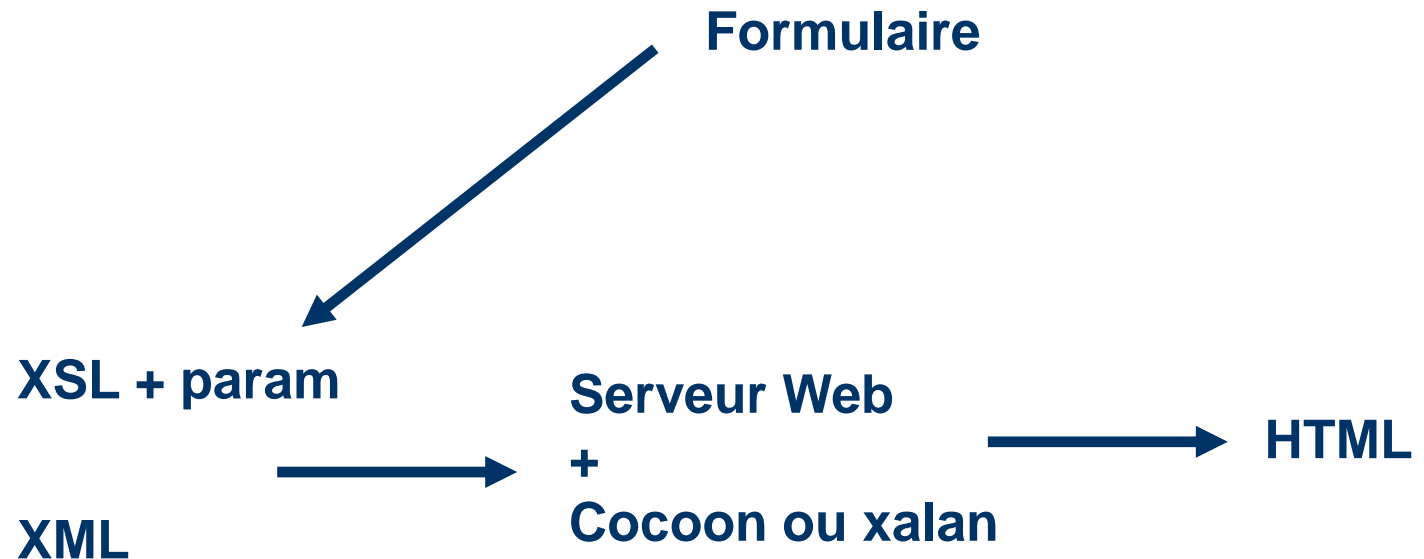
- **Architecture 2 : statique**





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Architecture 3 : paramétrable**







- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

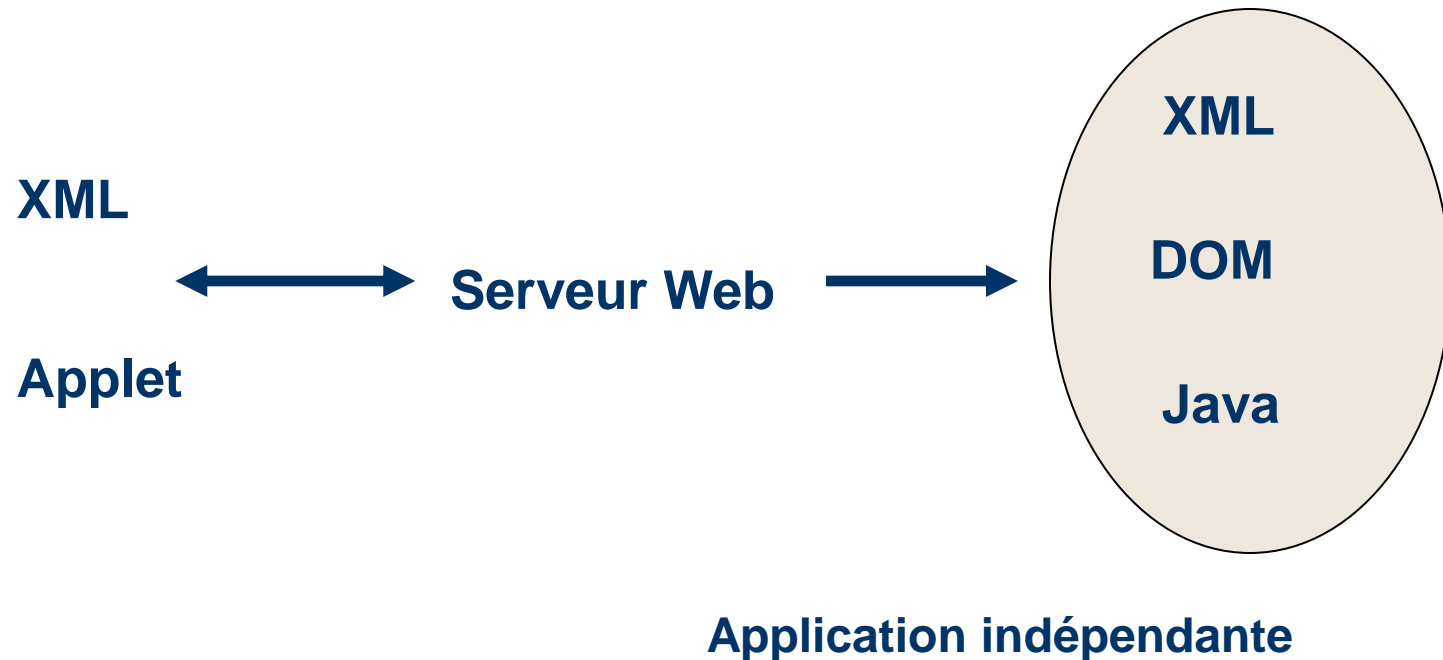
- **Architecture 4 : dynamique**





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Architecture 5 : XML coté client**





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Architecture 1 : Application serveur**





- Introduction
- Parsing SAX
- Exercices simples
- Parsing DOM
- Exercices simples
- Outils

- **Architecture 2 : XML crée coté client**

