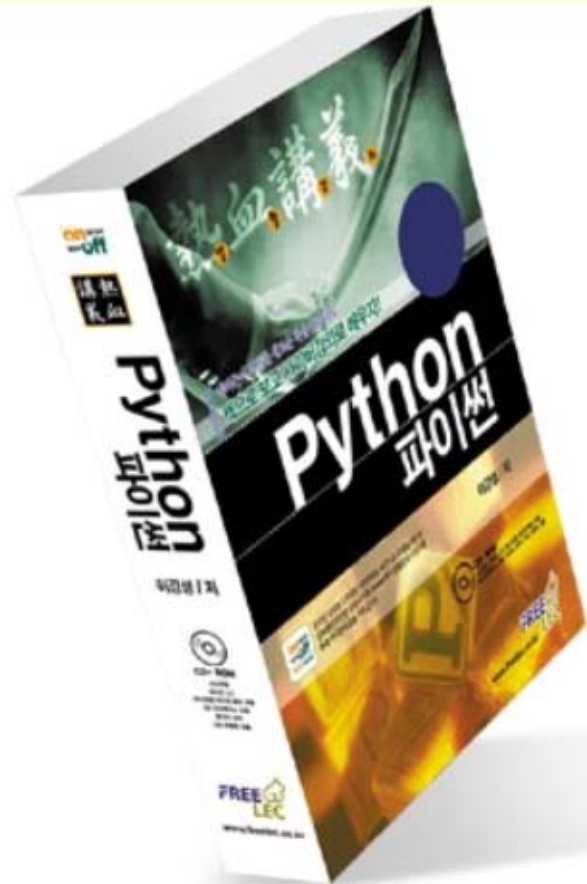熱 血 講 義

프리렉의 열혈강의 시리즈

# Python
## 파이썬

# Python

❖ 21-2

: XML (DOM)

**( gslee@mail.kw.ac.kr )**

❖

1. ..
2. SAX
3. DOM
4. Traversal
5. XPath
6. XSLT
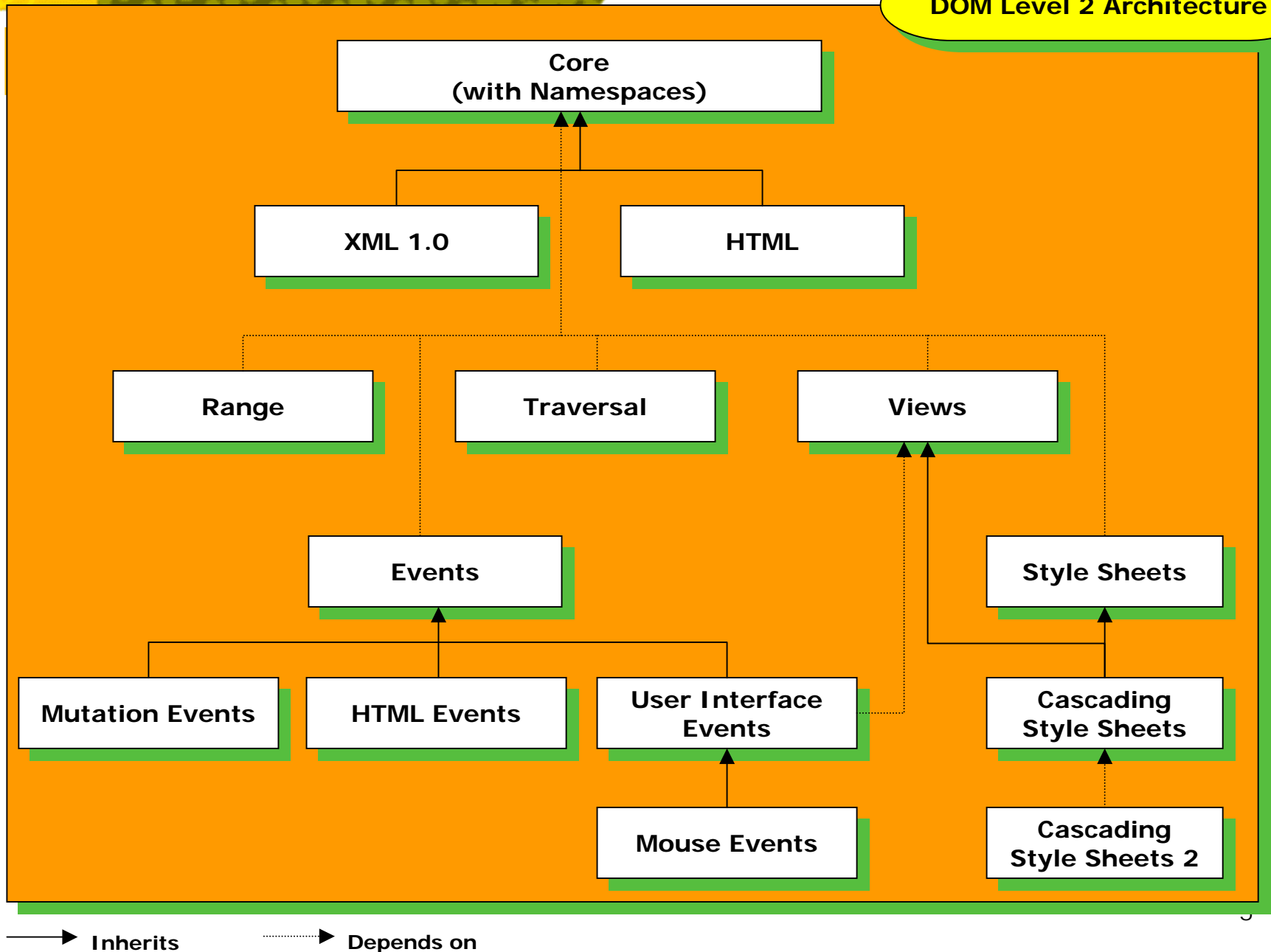
# DOM

- **DOM(Document Object Model)**
  - ➢ **W3C** **API**
  - ➢ **API**
  - ➢ **binding (Python, Java, C++, VB..)**
  - ➢ **1 → 2 → 3(Draft)**

**DOM Level 2 Architecture**

```
                        Core
                  (with Namespaces)

              XML 1.0              HTML

        Range          Traversal          Views

                       Events                        Style Sheets

Mutation Events    HTML Events    User Interface          Cascading
                                    Events               Style Sheets

                                   Mouse Events          Cascading
                                                        Style Sheets 2
```

⟶ Inherits        ┈┈▶ Depends on

# 파이썬 (Python)

**DOM Level 3 Architecture**

**Core**
**(with Namespaces, XML Base)**

**XML 1.0**

**HTML**

**Range**

**Traversal**

**Document Editing**

**Load/Save**

**Views**

**XPath**

**Events**
**(with Events Group)**

**Style Sheets**

**Mutation Events**

**HTML Events**

**User Interface Events**

**Cascading Style Sheets**

**Text Events**

**Mouse Events**

**Cascading Style Sheets 2**

→ **Inherits**      ┈┈► **Depends on**

# DOM

- 
  - ➢ **minidom**
  - ➢ **pulldom**
- **PyXML**
  - ➢ **http://pyxml.sourceforge.net/**
  - ➢ **4DOM, javadom (for Jython)**
- **4Suite**
  - ➢ **cDomlette**

# Load and Save

- ## DOM Core        2        Load and Save                binding
  .

- ## minidom

```
# dom01.py
from xml.dom.minidom import *

dom1 = parse('sample01.xml')  #              DOM

f = open('sample01.xml')
dom2 = parse(f)        #              DOM

dom3 = parseString('<value><int>1</int></value>')
```

# Load and Save

## • 4DOM Load

```
# dom03.py
from xml.dom.ext.reader.Sax2 import *

f = open('sample01.xml')
reader = Reader()  # Reader           .
dom1 = reader.fromStream(f)
dom2 = reader.fromUri('sample01.xml')
dom3 = reader.fromString('<value><int>1</int></value>')
```

9

# Load and Save

- ## 4DOM          Save

```
# dom04.py
from xml.dom.ext.reader.Sax2 import Reader
from xml.dom.ext import PrettyPrint


reader = Reader()  # Reader            .
dom = reader.fromUri('sample01.xml')  # URI      .
PrettyPrint(dom)         # dom                  .
```

```
f = open('test.xml', 'w')
PrettyPrint(dom, f)      # dom                  .
f.close()
```
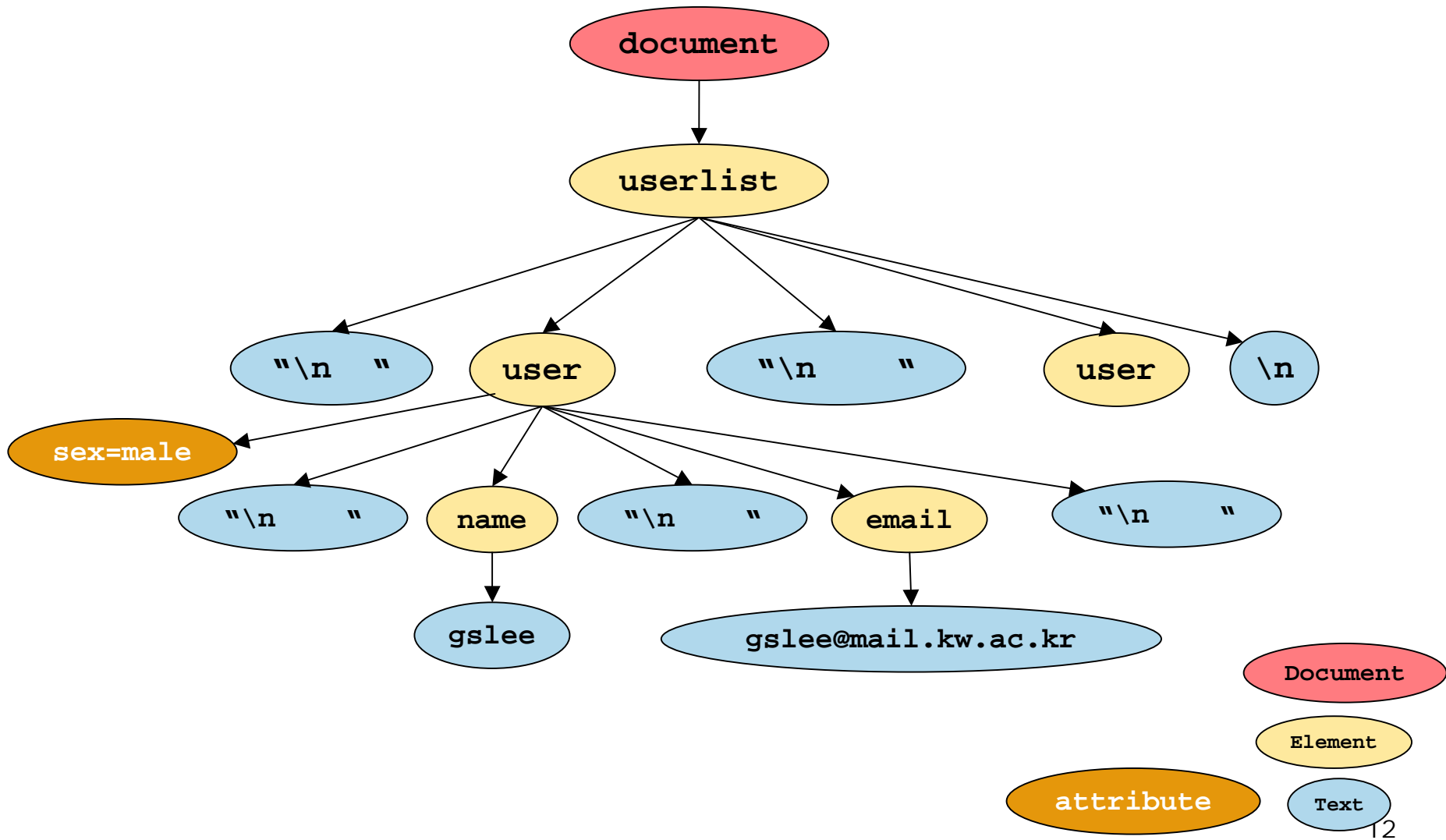
```
from StringIO import StringIO


strio = StringIO()    #                         .
PrettyPrint(dom, strio) #
xmlstr = strio.getvalue() #                      .
```

10

# DOM

- ## sample02.xml

```xml
<?xml version="1.0" ?>
<userlist>
    <user sex="male">
        <name>gslee</name>
        <email>gslee@mail.kw.ac.kr</email>
    </user>
    <user sex="female">
        <name>spam</name>
        <email>spam@mail.kw.ac.kr</email>
    </user>
</userlist>
```
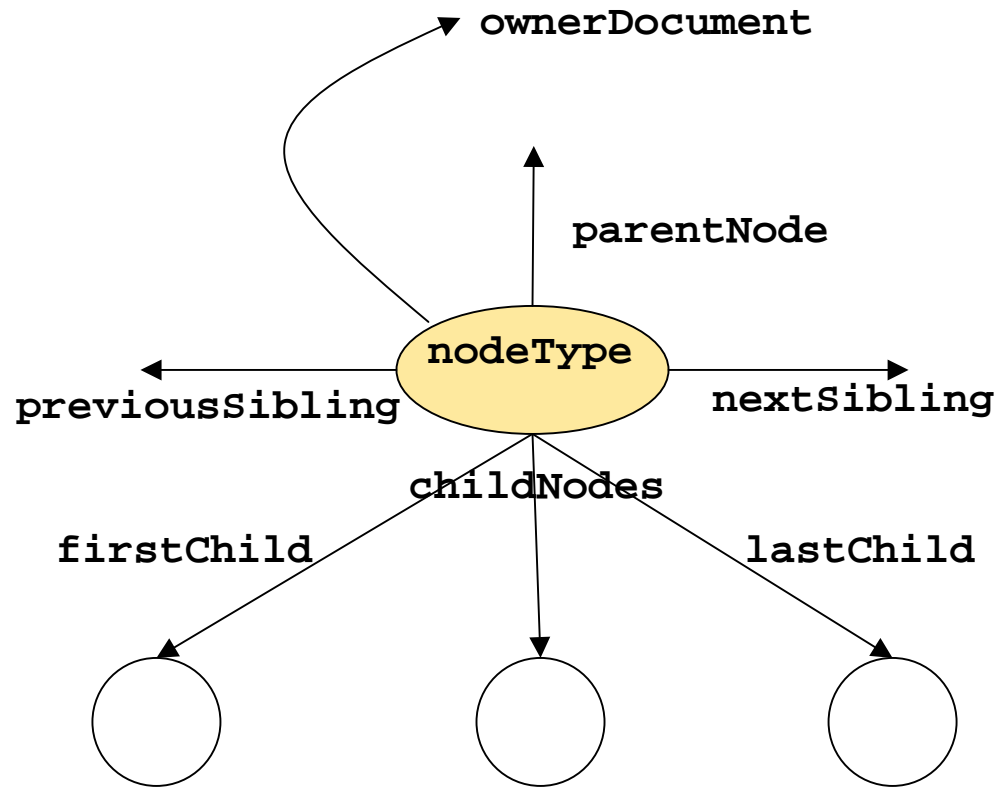
# DOM

- 

  ➢ **.nodeType**

```
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
```
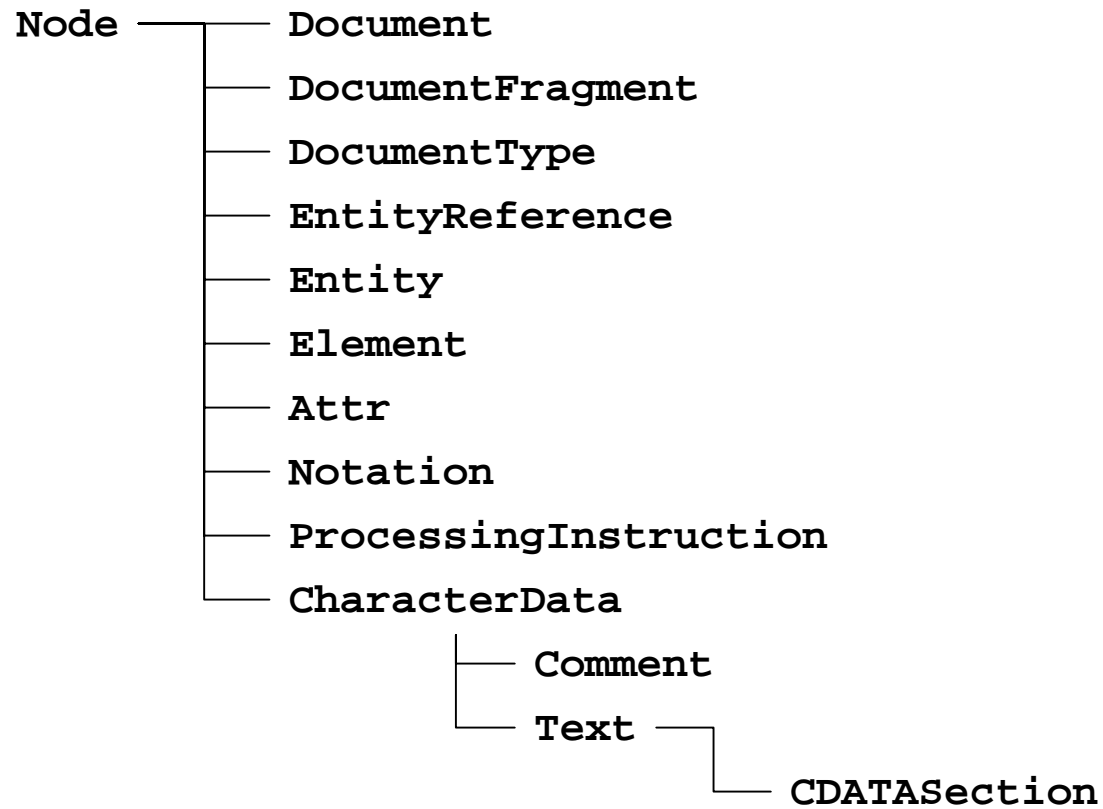
```
>>> dom.nodeType == dom.DOCUMENT_NODE

1
```

ownerDocument

parentNode

nodeType

previousSibling

nextSibling

childNodes

firstChild

lastChild

```
Node ───┬─── Document
        ├─── DocumentFragment
        ├─── DocumentType
        ├─── EntityReference
        ├─── Entity
        ├─── Element
        ├─── Attr
        ├─── Notation
        ├─── ProcessingInstruction
        └─── CharacterData
                  ├─── Comment
                  └─── Text ───┐
                               └─── CDATASection
```

# Node IDL

- ## IDL(interface definition language)

```
interface Node {
  //
  const unsigned short      ELEMENT_NODE                    = 1;
  const unsigned short      ATTRIBUTE_NODE                  = 2;
  const unsigned short      TEXT_NODE                       = 3;
  const unsigned short      CDATA_SECTION_NODE              = 4;
  const unsigned short      ENTITY_REFERENCE_NODE           = 5;
  const unsigned short      ENTITY_NODE                     = 6;
  const unsigned short      PROCESSING_INSTRUCTION_NODE     = 7;
  const unsigned short      COMMENT_NODE                    = 8;
  const unsigned short      DOCUMENT_NODE                   = 9;
  const unsigned short      DOCUMENT_TYPE_NODE              = 10;
  const unsigned short      DOCUMENT_FRAGMENT_NODE          = 11;
  const unsigned short      NOTATION_NODE                   = 12;
```

# Node IDL

```
//
 readonly attribute DOMString        nodeName;    //
         attribute DOMString         nodeValue;   //
 readonly attribute unsigned short   nodeType;    //


 //
 readonly attribute Node             parentNode;  //
 readonly attribute NodeList         childNodes;  //
 readonly attribute Node             firstChild;  //
 readonly attribute Node             lastChild;   //
 readonly attribute Node             previousSibling; //
 readonly attribute Node             nextSibling;     //
 readonly attribute NamedNodeMap     attributes;
 readonly attribute Document         ownerDocument;   //
```

# Node IDL

```
//
 Node      insertBefore(in Node newChild, in Node refChild)  //
 Node      replaceChild(in Node newChild, in Node oldChild)  //
 Node      removeChild(in Node oldChild)  //
 Node      appendChild(in Node newChild)  //
 boolean   hasChildNodes();                  //
 Node      cloneNode(in boolean deep);     //
 void      normalize();                      //


 //
 boolean   isSupported(in DOMString feature, in DOMString version); //


 boolean   hasAttributes();                      //


 //
 readonly attribute DOMString        namespaceURI; //          URI
         attribute DOMString        prefix;       //          prefix
 readonly attribute DOMString        localName;
};
```

# Text

- 

```
>>> def getText(node):
        L = []
        for childNode in node.childNodes:
            if childNode.nodeType == childNode.TEXT_NODE:
                L.append(childNode.data)
        return ''.join(L)

>>> getText(dom.childNodes[1])
u' \n      \n      \n'
```

# Text

- 

```
# dom07.py
from xml.dom.ext.reader.Sax2 import Reader

def getAllText(node):
    s = ''
    for node in node.childNodes:
        if node.nodeType == node.TEXT_NODE:
            s += node.nodeValue
        elif node.nodeType == node.ELEMENT_NODE:
            s += getAllText(node)
    return s

if __name__ == '__main__':
    reader = Reader()
    dom = reader.fromUri('sample02.xml')
    print getAllText(dom)
```

- ## getElementsByTagName(tagName)
  - ### tagName
- ## getElementsByTagNameNS(namespaceURI, localName)
  - ### ...

```
>>> dom.getElementsByTagName('user')   # user                    .
<NodeList at 1925280: [<Element Node at 18c33a0: Name='user' with 1
attributes and 5 children>, ...
>>> dom.getElementsByTagNameNS(None, 'user')
<NodeList at 1928f40: [<Element Node at 18c33a0: Name='user' with 1
attributes and 5 children>, ...

>>> print [node.nodeName for node in doc.getElementsByTagNameNS('*', '*')]

[u'userlist', u'user', u'name', u'email', u'user', u'name', u'email']
```

- 

  - ➢ **getAttribute(attrName)**

    - ▪

  - ➢ **getAttributeNS(NmSpace, attrName)**

- **attributes
  NameNodeMap**

```
>>> attr = elem.attributes
>>> attr[(None, 'sex')]     #
<Attribute Node at 18c5a20: Name="sex", Value="male">
>>> for key, value in attr.items():
        print key, value


(None, u'sex') <Attribute Node at 18c5a20: Name="sex", Value="male">
>>> for key, value in attr.items():    #      (      ,    )
        print key, value.value


(None, u'sex') male
>>> for value in attr.values():        #
        print value.name, value.value


sex male
```

# DOM Core -

| | |
|---|---|
| createElement(tagname) | . |
| createElementNS(namespaceURI, qualifiedName) | createElement |
| createTextNode(data) | . |
| createComment(data) | . |
| createDocumentFragment() | DocumentFragment . |
| createCDATASection(data) | CDATA |
| createEntityReference(name) | . |
| createProcessingInstruction(target, data) | ProcssingInstruction |

| | |
|---|---|
| **appendChild(newChild)** | .<br>**newChild** |
| **insertBefore(newChild, refChild)** | **(refChild)** .<br>**newChild** |
| **removeChild(oldChild)** | . **oldChild**<br>. |
| **replaceChild(newChild, oldChild)** | . **oldChild**<br>. |
| **cloneNode(deep)** | . **deep**<br>. . |
| **normalize()** | .<br>. |

- **DOMImplementation**

  - **createDocument(namespaceURI, qualifiedName, docType)**

```
>>> from xml.dom import implementation

>>> root = implementation.createDocument(None, None, None)

>>> root

<XML Document at 17a4020>
```

- **createElement(tagName)**
- **createElementNS(namespaceURI, qualifiedName)**

```
>>> rootElem = root.createElementNS(None, 'userlist')
>>> root.appendChild(rootElem)
```

```
>>> user = root.createElementNS(None, 'user')
>>> name = root.createElementNS(None, 'name')
>>> email = root.createElementNS(None, 'email')
>>> user.appendChild(name)
<Element Node at 1873050: Name='name' with 0 attributes and 0
children>
>>> user.appendChild(email)
<Element Node at 187bb60: Name='email' with 0 attributes and
0 children>
>>> rootElem.appendChild(user)
<Element Node at 186c9c0: Name='user' with 0 attributes and 2
children>
>>> PrettyPrint(root)
<?xml version='1.0' encoding='UTF-8'?>
<userlist>
  <user>
    <name/>
    <email/>
  </user>
</userlist>
```

```
>>> name.appendChild(root.createTextNode('gslee'))
<Text Node at 1899de0: 'gslee'>
>>>
email.appendChild(root.createTextNode('gslee@some.where'))
<Text Node at 18a6860: 'gslee@some.where'>
>>> PrettyPrint(root)
<?xml version='1.0' encoding='UTF-8'?>
<userlist>
  <user>
    <name>gslee</name>
    <email>gslee@some.where</email>
  </user>
</userlist>
```

# data

```
>>> email.firstChild.data = 'gslee@mail.gwu.ac.kr'

>>> PrettyPrint(root)

<?xml version='1.0' encoding='UTF-8'?>

<userlist>

  <user>

    <name>gslee</name>

    <email>gslee@mail.gwu.ac.kr</email>

  </user>

</userlist>
```

# setAttributeNS

```
>>> user.setAttributeNS(None, 'sex', 'female')

>>> PrettyPrint(root)

<?xml version='1.0' encoding='UTF-8'?>

<userlist>

  <user sex='female'>

    <name>gslee</name>

    <email>gslee@mail.gwu.ac.kr</email>

  </user>

</userlist>
```

31

## removeAttributeNS

```
>>> user.setAttributeNS(None, 'age', '40')  #

>>> user.removeAttributeNS(None, 'age')     #

>>> PrettyPrint(root)

<?xml version='1.0' encoding='UTF-8'?>

<userlist>

  <user sex='male'>

    <name>gslee</name>

    <email>gslee@mail.gwu.ac.kr</email>

  </user>

</userlist>
```

## removeChild(oldChild)

```
# dom08.py
from xml.dom.ext.reader.Sax2 import Reader
from xml.dom.ext import PrettyPrint

reader = Reader()
doc = reader.fromUri('sample02.xml')

def removeNode(name, node):
    for childNode in node.childNodes:
        if childNode.nodeName == name:
            node.removeChild(childNode)
        else:
            removeNode(name, childNode)

removeNode('email', doc)
PrettyPrint(doc)
```

## ● **replaceChild(newChild, oldChild)**

```
# dom09.py
from xml.dom.ext.reader.Sax2 import Reader
from xml.dom.ext import PrettyPrint

reader = Reader()
doc = reader.fromUri('sample02.xml')

email = doc.getElementsByTagName('email')[0]
oldtext = email.firstChild    #                         .

email.replaceChild(doc.createTextNode('gslee@pymail.net'),
oldtext)  #            .

PrettyPrint(doc)
```

34

## removeChild  appendChild insertBefore

```python
# dom10.py
from xml.dom.ext.reader.Sax2 import Reader
from xml.dom.ext import PrettyPrint

reader = Reader()
doc = reader.fromUri('sample02.xml')

users = doc.getElementsByTagName('user')
parent = users[0].parentNode
parent.removeChild(users[0])
parent.appendChild(users[0])

PrettyPrint(doc)
```

➢ **cloneNode(deep)**

```
# dom11.py
from xml.dom.ext.reader.Sax2 import Reader
from xml.dom.ext import PrettyPrint

reader = Reader()
doc = reader.fromUri('sample02.xml')

firstUser = doc.getElementsByTagNameNS(None, 'user')[0]
clonedUser = firstUser.cloneNode(1)

doc.documentElement.appendChild(clonedUser)
PrettyPrint(doc)
```

## Document

### importNode(importedNode, deep)

```
# dom12.py
from xml.dom.ext.reader.Sax2 import Reader
from xml.dom.ext import PrettyPrint

reader = Reader()
doc = reader.fromUri('sample02.xml')      #
doc2 = reader.fromString('<userlist/>')  #

node = doc2.importNode(doc.documentElement.childNodes[1], 1)
doc2.documentElement.appendChild(node)

PrettyPrint(doc2)
```
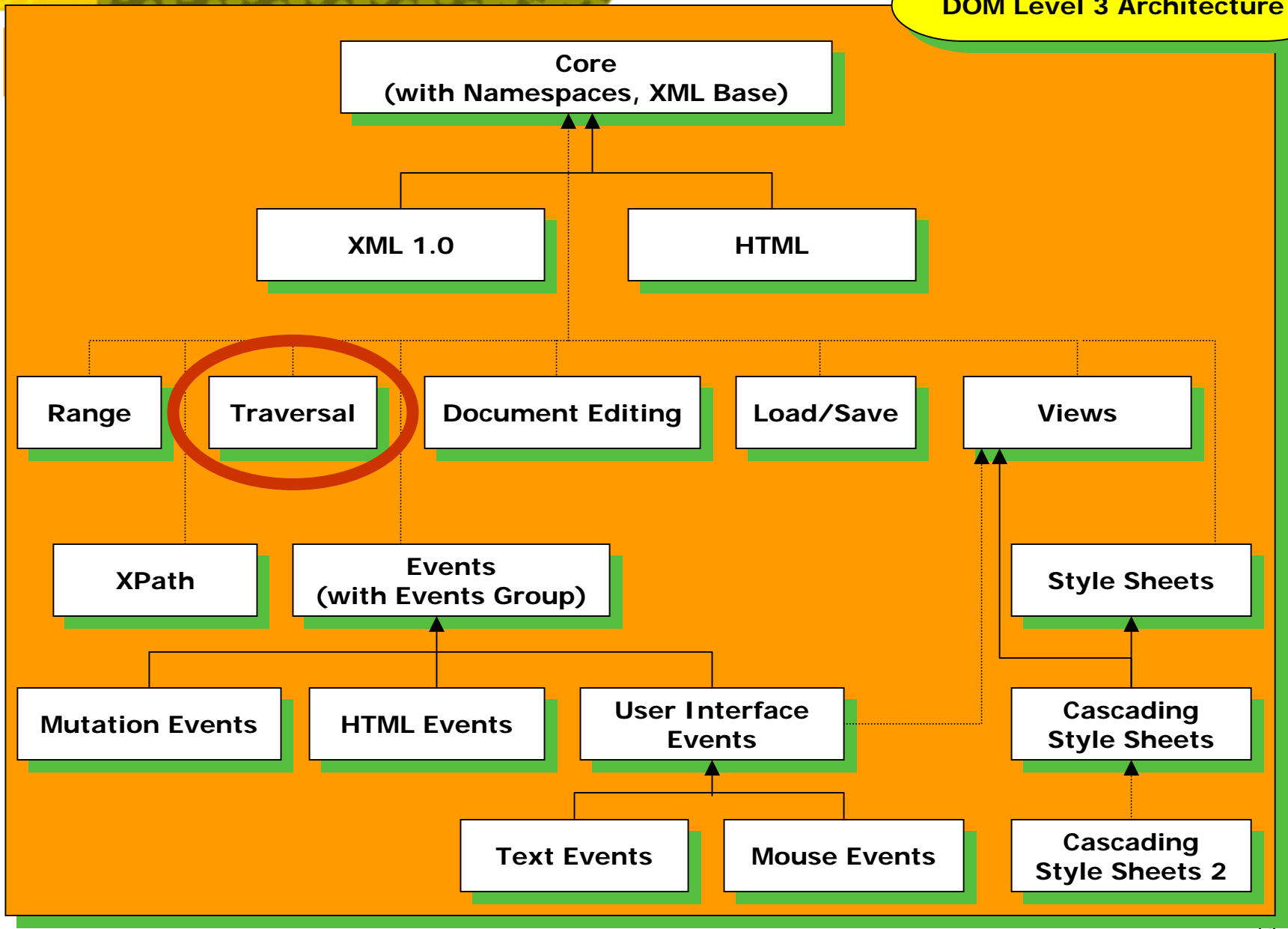
# DOM Traversal

- **DOM          2**

-

  ➢ **NodeIterator**

    ▪ 1

  ➢ **TreeWalker**

    ▪

  ➢ **NodeFilter**

    ▪

**DOM Level 3 Architecture**

**Core**
**(with Namespaces, XML Base)**

**XML 1.0**　　　　**HTML**

**Range**　　**Traversal**　　**Document Editing**　　**Load/Save**　　**Views**

**XPath**　　**Events**
**(with Events Group)**　　**Style Sheets**

**Mutation Events**　　**HTML Events**　　**User Interface Events**　　**Cascading Style Sheets**

**Text Events**　　**Mouse Events**　　**Cascading Style Sheets 2**

→ **Inherits**　　⋯▸ **Depends on**

# Node Iterator

- 

  - ➢ **createNodeIterator(root, whatToShow, filter, entityReferenceExpansion)**

- 

  - ➢ **nextNode()**
    - ▪

  - ➢ **previousNode()**
    - ▪

  - ➢ **detach()**
    - ▪

# Node Iterator

- :

```python
# dom13.py
from xml.dom.NodeFilter import NodeFilter
from xml.dom.ext.reader.Sax2 import Reader

reader = Reader()
doc = reader.fromUri('sample02.xml')

iterator = doc.createNodeIterator(doc, NodeFilter.SHOW_ELEMENT,
None, 0)
elem = iterator.nextNode()
while elem:
    print elem.nodeName
    elem = iterator.nextNode()
```

# Node Iterator

- 
  - ## sample02.xml

```xml
<?xml version="1.0" ?>
<userlist>
    <user sex="male">
        <name>gslee</name>
        <email>gslee@mail.gwu.ac.kr</email>
    </user>
    <user sex="female">
        <name>spam</name>
        <email>spam@mail.gwu.ac.kr</email>
    </user>
</userlist>
```

# Node Iterator

- 

```
userlist

user

name

email

user

name

email
```

# Node filter

- 
  - ➢ **from xml.dom.NodeFilter import NodeFilter**

- 
  - ➢ **iterator = doc.createNodeIterator(doc, NodeFilter.SHOW_ELEMENT, None, 0)**

# Node filter

- 

SHOW_ALL

SHOW_ELEMENT

SHOW_ATTRIBUTE

SHOW_TEXT

SHOW_CDATA_SECTION

SHOW_ENTITY_REFERENCE

SHOW_ENTITY

SHOW_PROCESSING_INSTRUCTION

SHOW_COMMENT

SHOW_DOCUMENT

SHOW_DOCUMENT_TYPE

SHOW_DOCUMENT_FRAGMENT

SHOW_NOTATION

# Node filter

●　　　：

```
# dom14.py
from xml.dom.NodeFilter import NodeFilter
from xml.dom.ext.reader.Sax2 import Reader

reader = Reader()
doc = reader.fromUri('sample02.xml')

iterator = doc.createNodeIterator(doc,
NodeFilter.SHOW_TEXT, None, 0)
elem = iterator.nextNode()
while elem:
    print elem.data,
    elem = iterator.nextNode()
```

# Node filter

- 

  - ➢ **NodeFilter**

  - ➢ **acceptNode**

    - ▪ **FILTER_ACCEPT,**

    - ▪ **FILTER_REJECT**

```
class emailFilter(NodeFilter):
    def acceptNode(self, node):
        if node.localName == 'email':
            return self.FILTER_ACCEPT
        else:
            return self.FILTER_REJECT
```

# Node filter

- 

```
# dom15.py
from xml.dom.NodeFilter import NodeFilter
from xml.dom.ext.reader.Sax2 import Reader

class emailFilter(NodeFilter):
    def acceptNode(self, node):
        if node.localName == 'email':
            return self.FILTER_ACCEPT
        else:
            return self.FILTER_REJECT

reader = Reader()
doc = reader.fromUri('sample02.xml')

iterator = doc.createNodeIterator(doc,
NodeFilter.SHOW_ELEMENT, emailFilter(), 0)
elem = iterator.nextNode()
while elem:
    print elem.firstChild.data
    elem = iterator.nextNode()
```

48

# Tree walker

- **Node Iterator**

  - ➤ **nextNode(), previousNode()**
  - ➤ **parentNode()**
  - ➤ **previousSibling(), nextSibling()**
  - ➤ **firstChild() lastChild()**

- **: Preorder**

# Tree walker

```python
# dom17.py
from xml.dom.NodeFilter import NodeFilter
from xml.dom.ext.reader.Sax2 import Reader

def doSomething(node):
    if node.nodeName in ('email', 'name'):
        print node.nodeName, ':', node.firstChild.data
    else:
        print node.nodeName

def processMe(tw):
    n = tw.currentNode
    doSomething(n)
    child = tw.firstChild()
    while child:
        processMe(tw)
        child = tw.nextSibling()
    tw.currentNode = n

reader = Reader()
doc = reader.fromUri('sample02.xml')

tw = doc.createTreeWalker(doc, NodeFilter.SHOW_ELEMENT, None, 1)
processMe(tw)
```

50