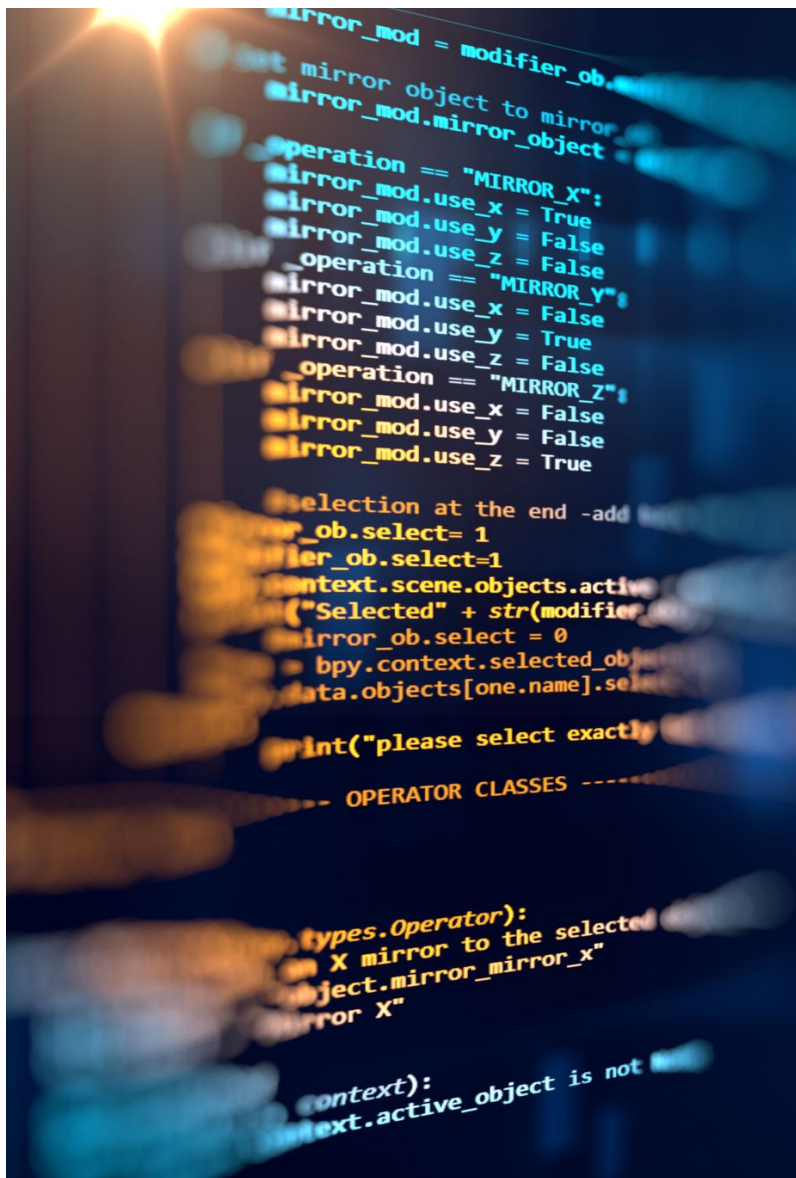# Design Pattern

김희성 20122372
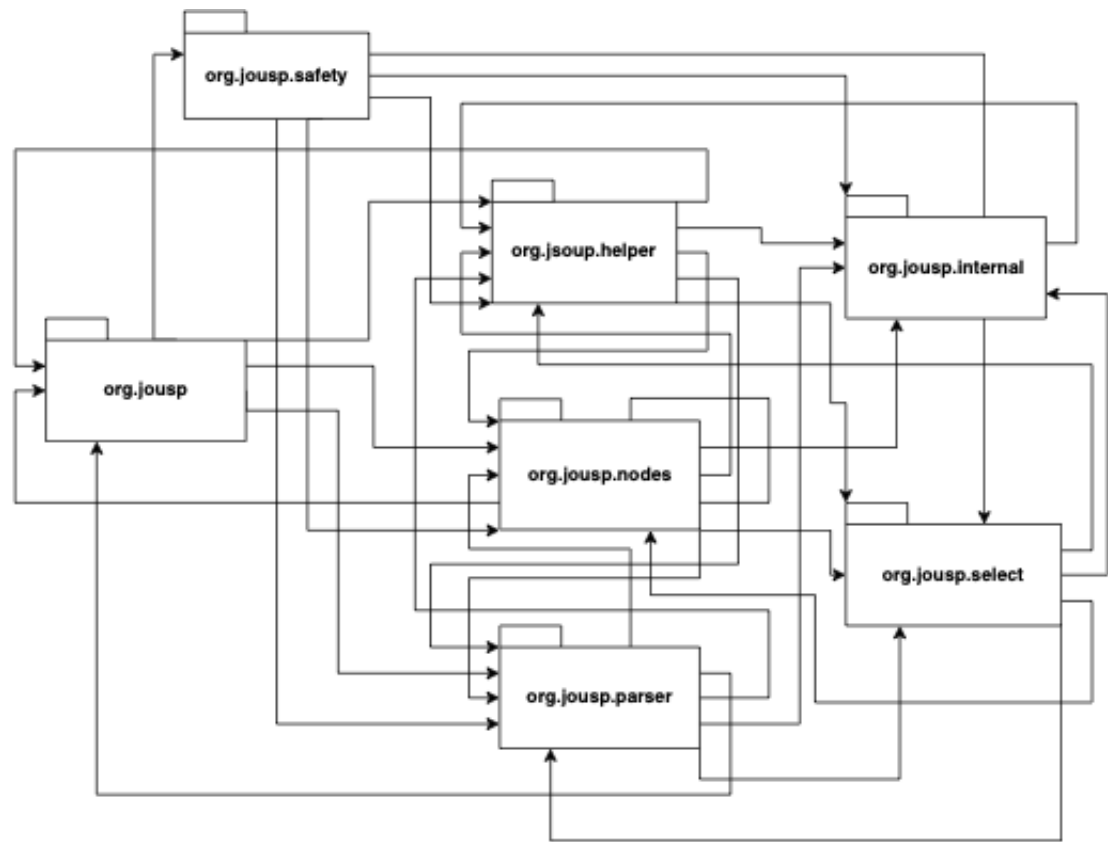김경태 20175119
주광우 20165953

13조

# Jsoup Introduction

Jsoup는 HTML 문서에 있는 데이터를 parse, extract, manipulate하기 위해 사용되는 오픈 소스 자바 라이브러리이다.

MIT License를 따르고 있으며, 따라서 코드 공개 의무가 없다. 또한 크롤링 등의 목적으로 여러 프로젝트에 사용되고 있다.

Jsoup는 WHATWG에서 정한 HTML5 specification에 맞춰 개발하였고, 브라우저와 같이 HTML을 DOM으로 parse 하는 역할을 한다.

사용자들은 Jsoup을 통해 HTML의 특정 요소들을 parsing해서 얻어 낼 수 있다.

Design Overview

org.jousp.safety

org.jsoup.helper

org.jousp.internal

org.jousp

org.jousp.nodes

org.jousp.select

org.jousp.parser

**parser**

Input Stream을 통해 받은 data를 Tokeniser를 통해 토큰으로 구분한다. 이후, 이 토큰으로부터 TreeBuilder를 통해 HTML이나 XML의 Document Tree를 생성하는 역할.

**nodes**

HTML과 XML의 Element를 표현하는 클래스들의 집합체로서 역할.

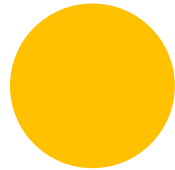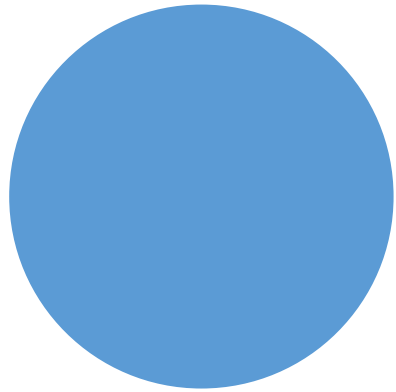각 클래스에는 Document를 생성할 때 필요한 insert, remove 등의 메소드를 제공.

**select**

select 패키지 안에는 CSS select query를 parsing한 후, 요청한 노드를 가져오는 역할.

**safety**

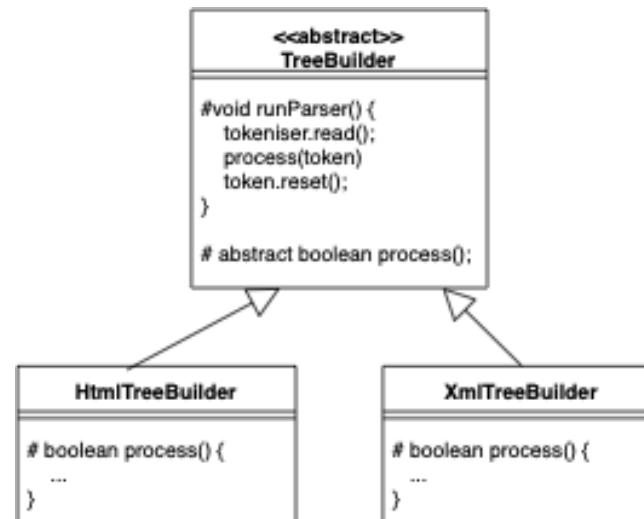safety에는 whitelist 규칙을 정의하는 Whitelist와 HTML 내용중 그 whitelist 규칙에 어긋나는 것을 없애는 Cleaner 클래스가 있음.

Cleaner를 통해 User가 원하는 Element와 attribute를 포함하고 있는 HTML을 제공받을 수 있음. 또한, cross-site scripting attack을 막음.

# Existing Patterns in Jsoup

# TreeBuilder

Template
Method
Pattern



```
<<abstract>>
TreeBuilder

#void runParser() {
    tokeniser.read();
    process(token)
    token.reset();
}

# abstract boolean process();
```

```
HtmlTreeBuilder

# boolean process() {
    ...
}
```

```
XmlTreeBuilder

# boolean process() {
    ...
}
```

```java
abstract class TreeBuilder {
// other codes
    Document parse(Reader input, String baseUri, Parser parser) {
        initialiseParse(input, baseUri, parser);
        runParser();
        return doc;
    }

    protected void runParser() {
        while (true) {
            Token token = tokeniser.read();
            process(token);
            token.reset();

            if (token.type == Token.TokenType.EOF)
                break;
        }
    }

    protected abstract boolean process(Token token);
}
```
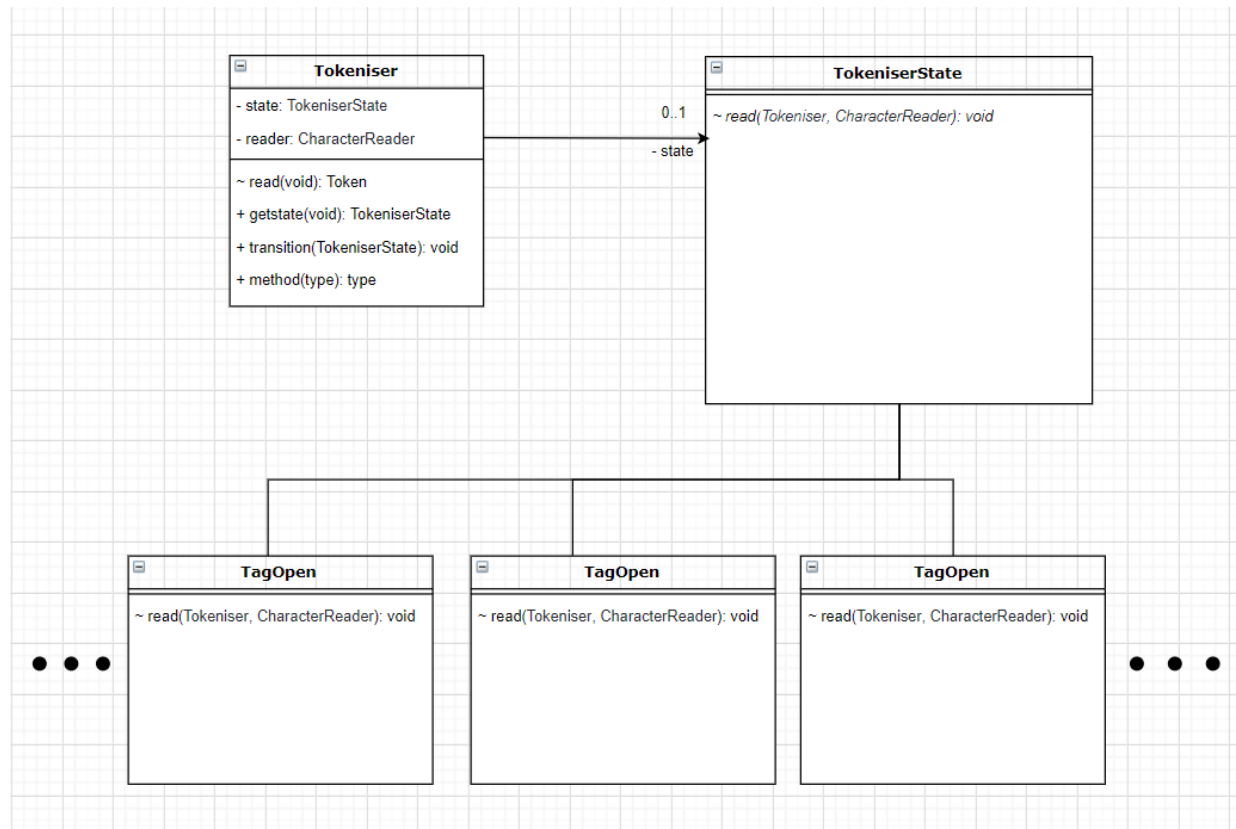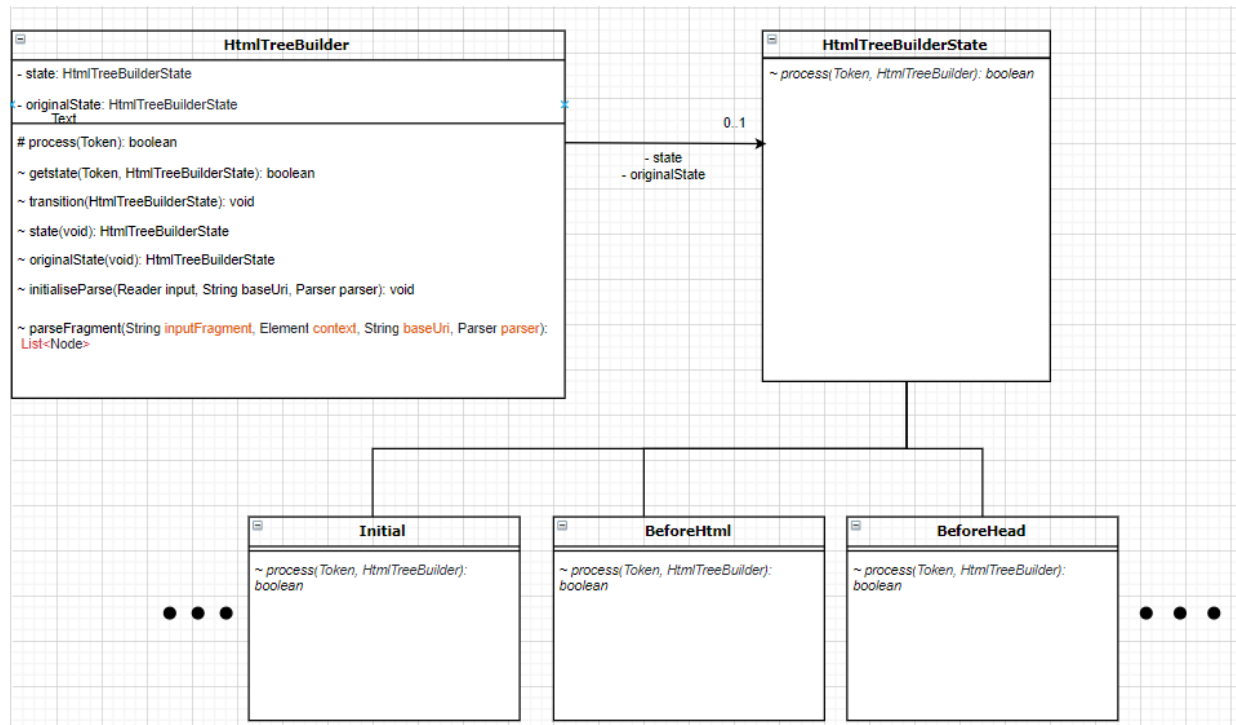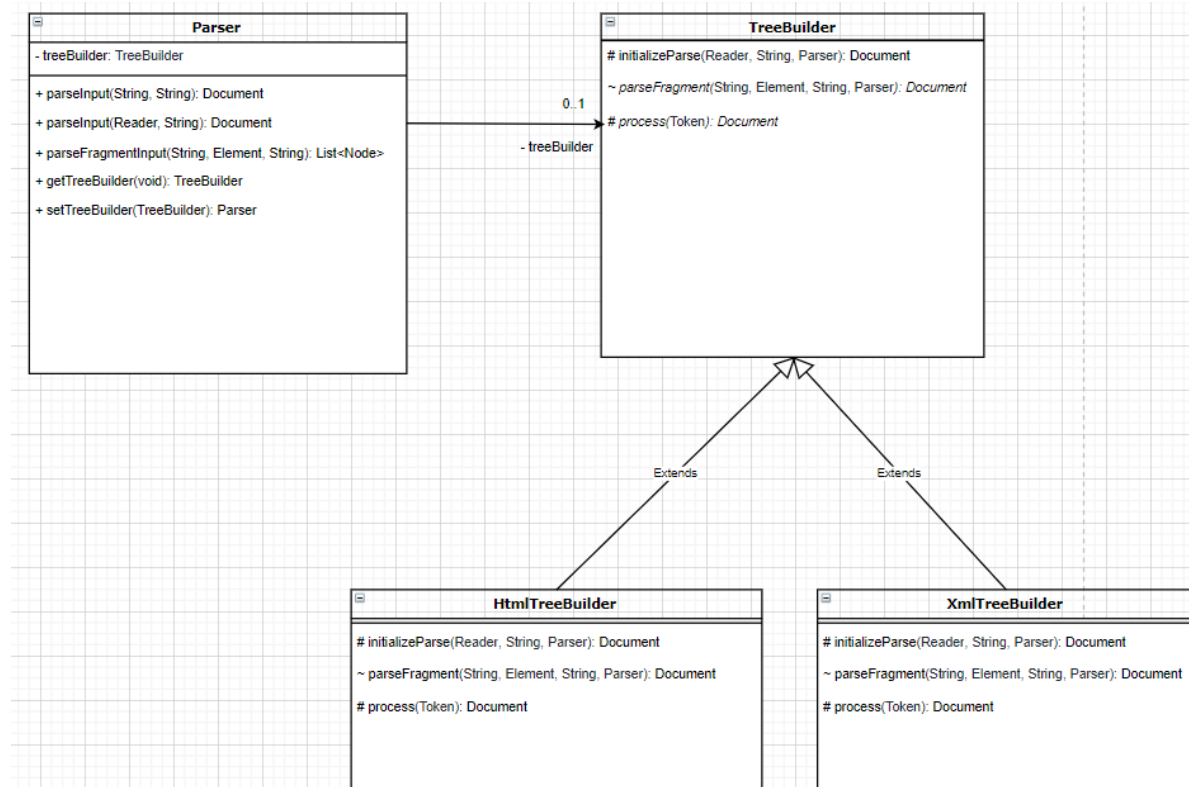
# TokeniserState

State Pattern
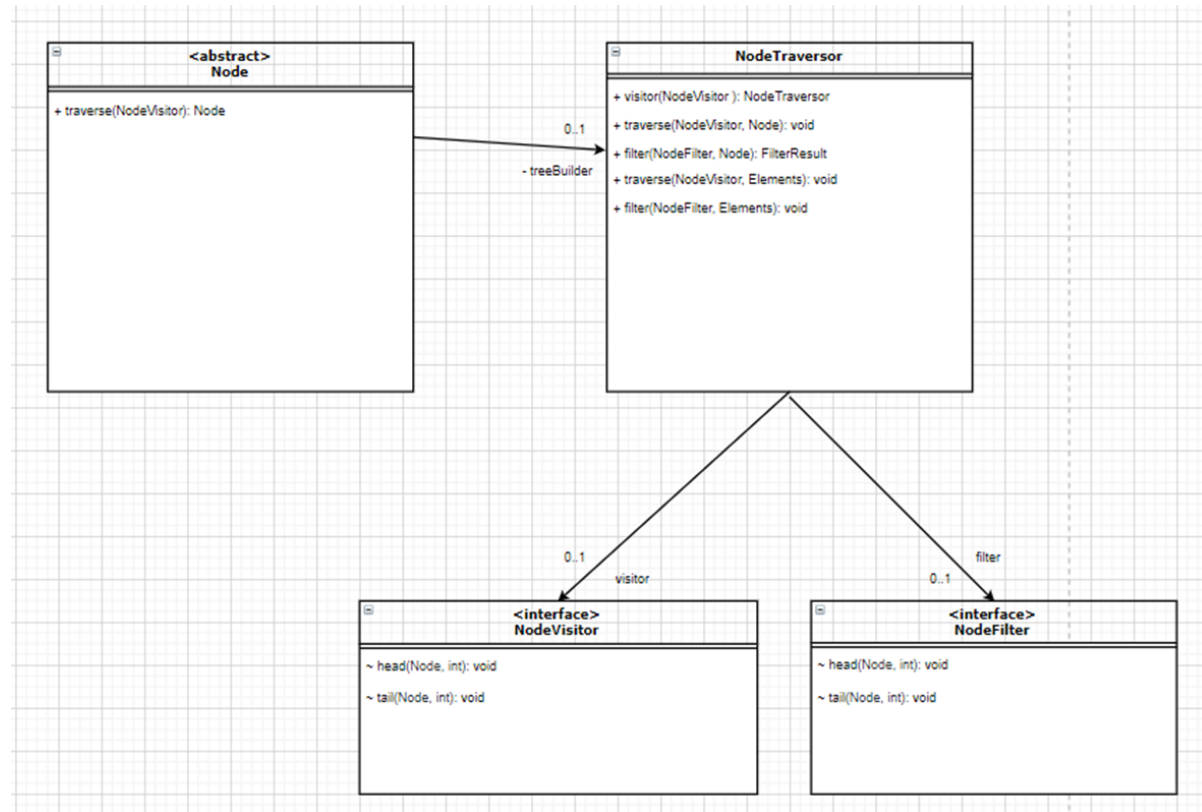
# HtmlTreeBuilderState



State Pattern

# Parser



Strategy Pattern

**Parser**

- treeBuilder: TreeBuilder

+ parseInput(String, String): Document
+ parseInput(Reader, String): Document
+ parseFragmentInput(String, Element, String): List<Node>
+ getTreeBuilder(void): TreeBuilder
+ setTreeBuilder(TreeBuilder): Parser

**TreeBuilder**

# initializeParse(Reader, String, Parser): Document

~ *parseFragment*(String, Element, String, Parser*): Document*

# *process*(Token*): Document*

0..1

- treeBuilder

**HtmlTreeBuilder**

# initializeParse(Reader, String, Parser): Document

~ parseFragment(String, Element, String, Parser): Document

# process(Token): Document

Extends

**XmlTreeBuilder**

# initializeParse(Reader, String, Parser): Document

~ parseFragment(String, Element, String, Parser): Document

# process(Token): Document

Extends

# NodeTraversor



**Strategy Pattern**

****
**Node**

+ traverse(NodeVisitor): Node

**NodeTraversor**

+ visitor(NodeVisitor ): NodeTraversor
+ traverse(NodeVisitor, Node): void
+ filter(NodeFilter, Node): FilterResult
+ traverse(NodeVisitor, Elements): void
+ filter(NodeFilter, Elements): void

0..1

- treeBuilder

0..1          visitor

filter

0..1

**&lt;interface&gt;**
**NodeVisitor**

~ head(Node, int): void

~ tail(Node, int): void

**&lt;interface&gt;**
**NodeFilter**

~ head(Node, int): void

~ tail(Node, int): void

# Tag

### Singleton Pattern



| Tag |
| --- |
| - tags: static final Map<String, Tag> |
| - Tag(String)<br>+ getName(void): String<br>+ valueOf(String, ParseSettings): Tag<br>+ valueOf(String): Tag |

```java
public class Tag {
    private static final Map<String, Tag> tags = new HashMap<>(); // map

    private String tagName;

    private Tag(String tagName) {
        this.tagName = tagName;
        normalName = Normalizer.lowerCase(tagName);
    }

    public static Tag valueOf(String tagName, ParseSettings settings) {
        Validate.notNull(tagName);
        Tag tag = tags.get(tagName);

        if (tag == null) {
            tagName = settings.normalizeTag(tagName);
            Validate.notEmpty(tagName);
            tag = tags.get(tagName);

            if (tag == null) {
                // not defined: create default; go anywhere, do anything!
                tag = new Tag(tagName);
                tag.isBlock = false;
            }
        }
        return tag;
    }
}
```
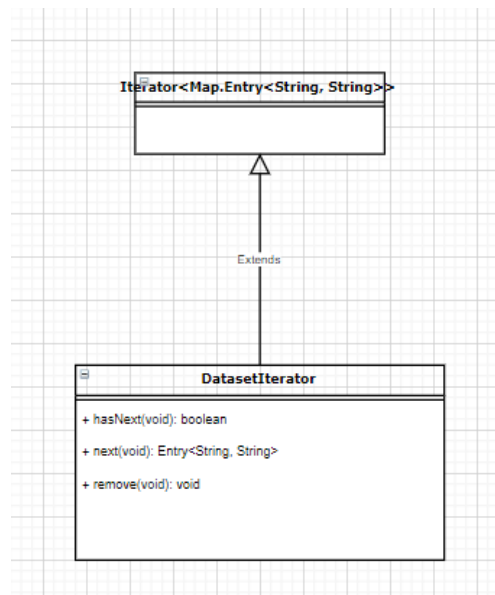
# Tag

**Singleton Pattern**

```java
private static final String[] blockTags = {
    "html", "head", "body", "frameset", "script", "noscript", "style", "meta",
    "noframes", "section", "nav", "aside", "hgroup", "header", "footer", "p",
    "ul", "ol", "pre", "div", "blockquote", "hr", "address", "figure", "figcapt
    "del", "dl", "dt", "dd", "li", "table", "caption", "thead", "tfoot", "tbody
    "td", "video", "audio", "canvas", "details", "menu", "plaintext", "template
    "svg", "math", "center"
};
...
```

```java
static {
    // creates
    for (String tagName : blockTags) {
        Tag tag = new Tag(tagName);
        register(tag);
    }
    for (String tagName : inlineTags) {
        Tag tag = new Tag(tagName);
        tag.isBlock = false;
        tag.formatAsBlock = false;
        register(tag);
    }

    // mods:
    for (String tagName : emptyTags) {
        Tag tag = tags.get(tagName);
        Validate.notNull(tag);
        tag.canContainInline = false;
        tag.empty = true;
    }
    ...

}

private static void register(Tag tag) {
    tags.put(tag.tagName, tag);
}
```

# DatesetIterator

Iterator
Pattern



```
Iterator<Map.Entry<String, String>>

            △
            │ Extends
            │

        DatasetIterator

+ hasNext(void): boolean

+ next(void): Entry<String, String>

+ remove(void): void
```

```java
private class DatasetIterator implements Iterator<Map.Entry<String, String>> {
    private Iterator<Attribute> attrIter = attributes.iterator();
    private Attribute attr;
    public boolean hasNext() {
        while (attrIter.hasNext()) {
            attr = attrIter.next();
            if (attr.isDataAttribute()) return true;
        }
        return false;
    }

    public Entry<String, String> next() {
        return new Attribute(attr.getKey().substring(dataPrefix.length()), attr.getValue());
    }

    public void remove() {
        attributes.remove(attr.getKey());
    }
}
```

## Façade Pattern

## Jsoup

```java
public class Jsoup {
    private Jsoup() {}

    Parse HTML into a Document. The parser will make a sensible, balanced document tree out of any HTML.
    public static Document parse(String html, String baseUri) {

    Parse HTML into a Document, using the provided Parser. You can provide an alternate parser, such as a simple XML
    public static Document parse(String html, String baseUri, Parser parser) {

    Parse HTML into a Document. As no base URI is specified, absolute URL detection relies on the HTML including a
    public static Document parse(String html) {

    * Creates a new {@link Connection} to a URL. Use to fetch and parse a HTML page.
    public static Connection connect(String url) {

    Parse the contents of a file as HTML.
    public static Document parse(File in, String charsetName, String baseUri) throws IOException {

    Parse the contents of a file as HTML. The location of the file is used as the base URI to qualify relative URLs.
    public static Document parse(File in, String charsetName) throws IOException {

    Read an input stream, and parse it to a Document.
    public static Document parse(InputStream in, String charsetName, String baseUri) throws IOException {

    Read an input stream, and parse it to a Document. You can provide an alternate parser, such as a simple XML
    public static Document parse(InputStream in, String charsetName, String baseUri, Parser parser) throws IOException {

    Parse a fragment of HTML, with the assumption that it forms the {@code body} of the HTML.
    public static Document parseBodyFragment(String bodyHtml, String baseUri) {

    Parse a fragment of HTML, with the assumption that it forms the {@code body} of the HTML.
    public static Document parseBodyFragment(String bodyHtml) {

    Fetch a URL, and parse it as HTML. Provided for compatibility; in most cases use {@link #connect(String)} instead.
    public static Document parse(URL url, int timeoutMillis) throws IOException {

    Get safe HTML from untrusted input HTML, by parsing input HTML and filtering it through a white-list of permitted
    public static String clean(String bodyHtml, String baseUri, Whitelist whitelist) {

    Get safe HTML from untrusted input HTML, by parsing input HTML and filtering it through a white-list of permitted
    public static String clean(String bodyHtml, Whitelist whitelist) {

    * Get safe HTML from untrusted input HTML, by parsing input HTML and filtering it through a white-list of
    public static String clean(String bodyHtml, String baseUri, Whitelist whitelist, Document.OutputSettings outputSettings) {

    Test if the input body HTML has only tags and attributes allowed by the Whitelist. Useful for form validation.
    public static boolean isValid(String bodyHtml, Whitelist whitelist) {

}
```
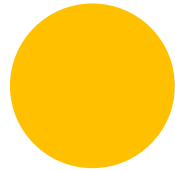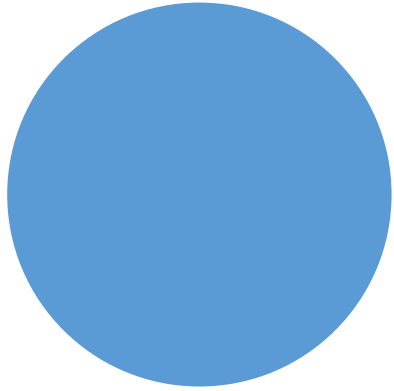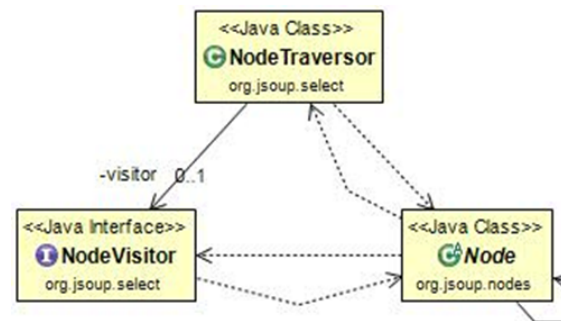
Improvement

Add functions

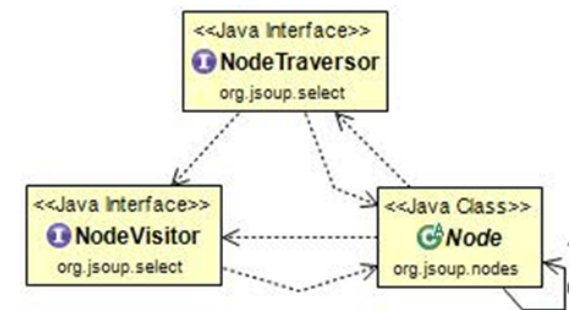Purpose : Find the last element

Approach : Add Node Traversal Algorithm

# Add Visitor Pattern



```
public Node traverse(NodeVisitor nodeVisitor) {
    Validate.notNull(nodeVisitor);
    NodeTraversor.traverse(nodeVisitor, this);
    return this;
}
```
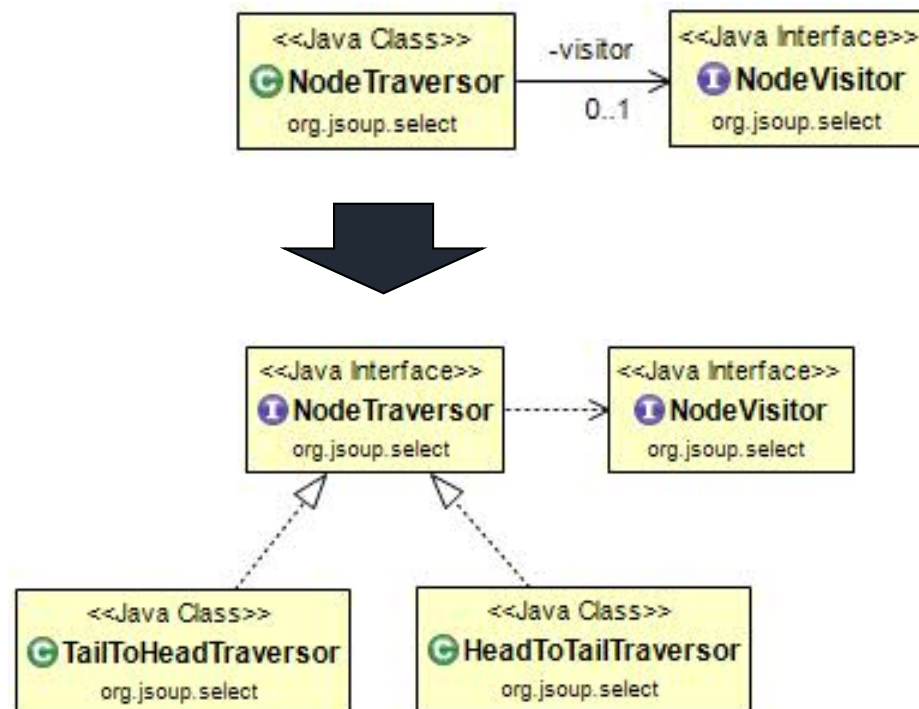
```
public Node traverse(NodeVisitor nodeVisitor) {
    Validate.notNull(nodeVisitor);
    NodeTraversor nodeTraversor
                = new HeadToTailTraversor();
    nodeTraversor.traverse(nodeVisitor, this);
    return this;
}

public Node traverse(   NodeTraversor nodeTraversor,
                        NodeVisitor nodeVisitor) {
    Validate.notNull(nodeTraversor);
    Validate.notNull(nodeVisitor);
    nodeTraversor.traverse(nodeVisitor, this);
    return this;
}
```

Add
Node Traversal
Algorithm

# Add Concrete Class



Add
Node Traversal
Algorithm

# Unit Test for Concrete Class of NodeTraversoer

```
10  public class TraversorTest {
11      // Note: NodeTraversor.traverse(new NodeVisitor) is tested in
12      // ElementsTest#traverse()
13
15⊕     public void HeadToTailFilterVisit() {...
34
36⊕     public void HeadToTailFilterSkipChildren() {...
56
58⊕     public void HeadToTailFilterSkipEntirely() {...
80
82⊕     public void HeadToTailFilterRemove() {...
100
102⊕    public void HeadToTailFilterStop() {...
122
124⊕    public void TailToHeadFilterVisit() {...
143
145⊕    public void TailToHeadFilterSkipChildren() {...
165
167⊕    public void TailToHeadFilterSkipEntirely() {...
189
191⊕    public void TailToHeadFilterRemove() {...
209
211⊕    public void TailToHeadFilterStop() {...
231  }
```

Add
Unit Test

## Add Collector.findLast

Add
Method

```java
public static Element findLast(Evaluator eval, Element root) {
    LastFinder finder = new LastFinder(root, eval);
    NodeTraversor nodeTraversor = new TailToHeadTraversor();
    nodeTraversor.filter(finder, root);
    return finder.match;
}

private static class LastFinder implements NodeFilter {
    private final Element root;
    private Element match = null;
    private final Evaluator eval;

    LastFinder(Element root, Evaluator eval) {
        this.root = root;
        this.eval = eval;
    }
```

# Add Selector.selectLast

Add Method

```java
public static Element selectFirst(String cssQuery, Element root) {
    Validate.notEmpty(cssQuery);
    return Collector.findFirst(QueryParser.parse(cssQuery), root);
}
```

⬇

```java
public static Element selectFirst(String cssQuery, Element root) {
    Validate.notEmpty(cssQuery);
    return Collector.findFirst(QueryParser.parse(cssQuery), root);
}

public static Element selectLast(String cssQuery, Element root) {
    Validate.notEmpty(cssQuery);
    return Collector.findLast(QueryParser.parse(cssQuery), root);
}
```

# Add Element.selectLast

**Add Method**

```java
public Element selectFirst(String cssQuery) {
    return Selector.selectFirst(cssQuery, this);
}
```
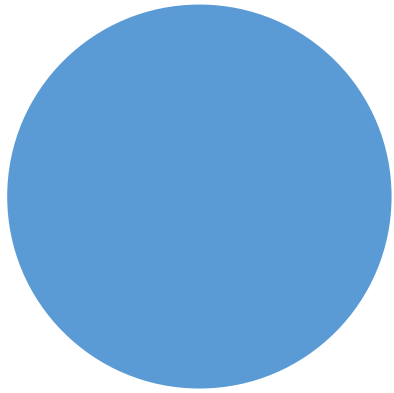
```java
public Element selectFirst(String cssQuery) {
    return Selector.selectFirst(cssQuery, this);
}

public Element selectLast(String cssQuery) {
    return Selector.selectLast(cssQuery, this);
}
```

# Unit Test for selectLast method

## Add Unit Test

```java
@Test public void selectLast() {
    String html = "<p>One<p>Two<p>Three";
    Document doc = Jsoup.parse(html);
    assertEquals("Three", doc.selectLast("p").text());
}

@Test public void selectLastWithAnd() {
    String html = "<p>One<p class=foo>Two<p>Three";
    Document doc = Jsoup.parse(html);
    assertEquals("Two", doc.selectLast("p.foo").text());
}

@Test public void selectLastWithOr() {
    String html = "<p>One<p>Two<p>Three<div>Four";
    Document doc = Jsoup.parse(html);
    assertEquals("Four", doc.selectLast("p, div").text());
}
```
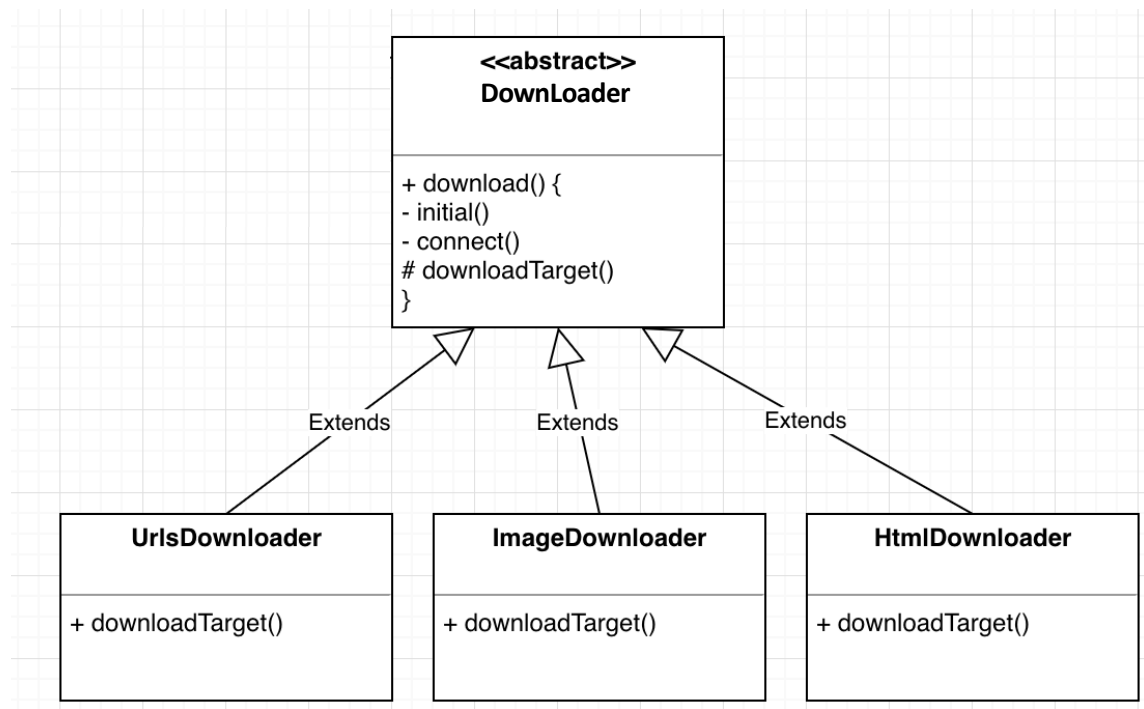
Extension

# Add package

**Purpose :** Download contents of page
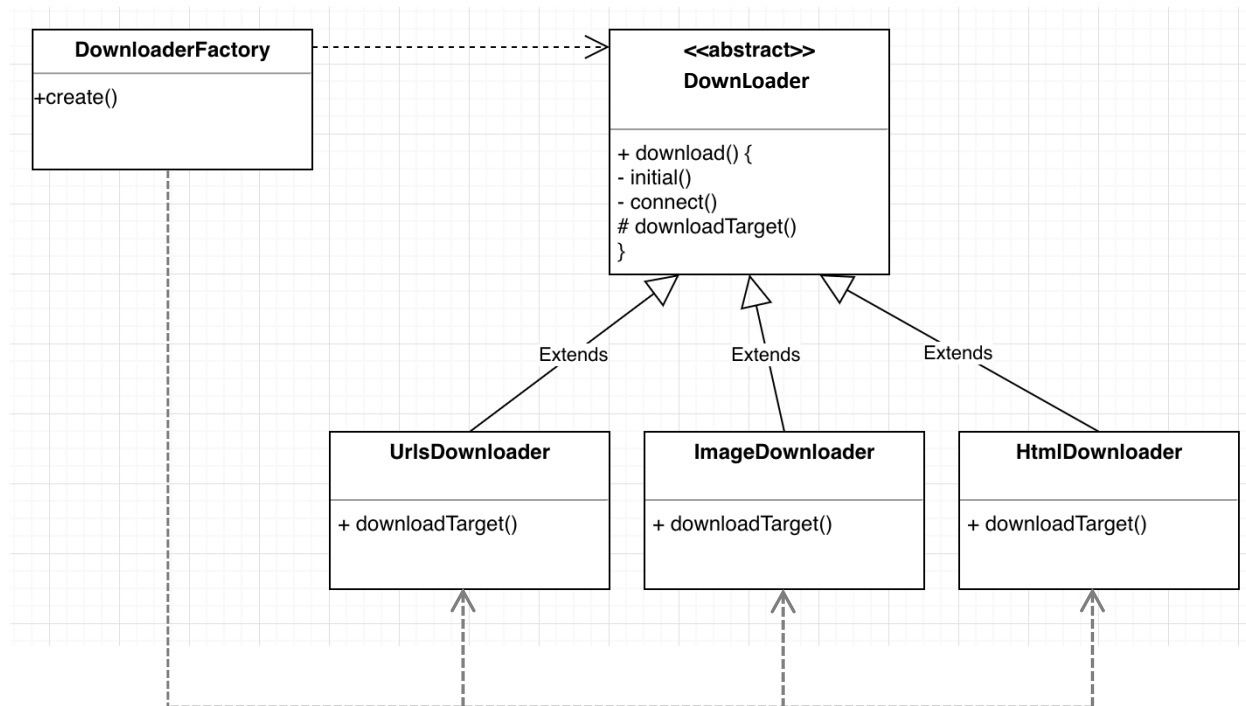
**Approach :** Add Content Downloader classes

# Add Template Method Pattern

Add
Contents Downloader
Class

```
           <<abstract>>
            DownLoader
  ┌────────────────────────────┐
  │ + download() {             │
  │ - initial()                │
  │ - connect()                │
  │ # downloadTarget()         │
  │ }                          │
  └────────────────────────────┘
         △     △     △
   Extends   Extends   Extends

 UrlsDownloader  ImageDownloader  HtmlDownloader
 + downloadTarget()  + downloadTarget()  + downloadTarget()
```

# Add Façade Pattern with Simple Factory

**Add Contents Downloader Class**

```
DownloaderFactory
-----------------
+create()
```

```
<<abstract>>
DownLoader
-----------------
+ download() {
- initial()
- connect()
# downloadTarget()
}
```

Extends          Extends          Extends

```
UrlsDownloader
-----------------
+ downloadTarget()
```

```
ImageDownloader
-----------------
+ downloadTarget()
```

```
HtmlDownloader
-----------------
+ downloadTarget()
```

# Add Unit Test for DownloaderFacotry

Add
Unit Test

```java
public class DownloaderFactoryTest {

    @Test
    public void testCreateUrlDownloader() {
        DownloaderFactory df = new DownloaderFactory();
        assertEquals(true, df.create("UrlsDownloader") instanceof UrlsDownloader);
    }

    @Test
    public void testCreateImageDownloader() {
        DownloaderFactory df = new DownloaderFactory();
        assertEquals(true, df.create("ImageDownloader") instanceof ImageDownloader);
    }

    @Test
    public void testCreateHtmlDownloader() {
        DownloaderFactory df = new DownloaderFactory();
        assertEquals(true, df.create("HtmlDownloader") instanceof HtmlDownloader);
    }

}
```

# Add Unit Test for UrlsDownloader

Add
Unit Test

```java
public class DownloaderTest {

    String testUrl =
            "https://raw.githubusercontent.com/JsoupMaster/jsoup/master/src/main/javadoc/overview.html";
    String storePath = System.getProperty("user.dir") + "/download/";

    @Test
    public void testUrlsDownloader() {
        File file = new File(storePath);
        file.mkdir();

        Downloader downloader = new UrlsDownloader();
        downloader.download(testUrl, storePath + "urls.txt");

        try {
            InputStream in = new FileInputStream(storePath + "urls.txt");
            try {
                String urls = new String(in.readAllBytes());
                assertTrue(urls.equalsIgnoreCase(
                        "http://whatwg.org/html\n" +
                        "http://jonathanhedley.com/\n" +
                        "https://jsoup.org/\n"));

            } catch (IOException e) {
                fail("file is not readable.");
            }
        } catch (FileNotFoundException e) {
            fail("file is not found.");
        }

        file = new File(storePath + "urls.txt");
        file.delete();

        file = new File(storePath);
        file.delete();
    }
}
```
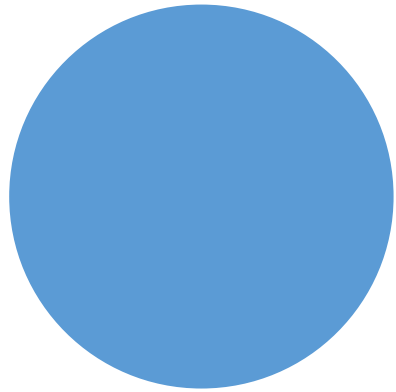
# DepthOneCrawler
# (ImageDownloader, HtmlDownloader)

**Add Example**
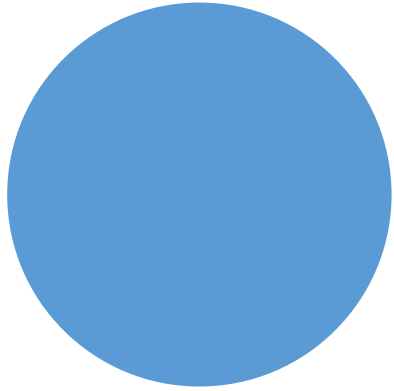
| | | |
|---|---|---|
| 31px-Commons-logo.svg.png | 35px-Mediawiki-logo.png | 35px-Wikibooks-logo.svg.png |
| 35px-Wikime....svg.png | 35px-Wikiquote-logo.svg.png | 35px-Wikisource-logo.svg.png |
| 35px-Wikispe....svg.png | 35px-Wikivoy....svg.png | 35px-Wiktionary-logo-v2.svg.png |
| 41px-Wikivers...svg.png | 47px-Wikidata-logo.svg.png | 51px-Wikinews-logo.svg.png |
| 120px-ADIL_A...d%29.jpg | 120px-Nofretet...eum.jpg | 120px-Sanjay_...0%99.jpg |

2005_ACC_Championship...me.html
2019_Hong_Kong_protests.html
2019_Iraqi_protests.html
2019_Uruguayan_general_...ion.html
Adil_Abdul-Mahdi.html
Atlantic_hurricane.html
Battle_of_the_Falkland_Islands.html
Chesapeake_Bay.html
Cold_front.html
daily-article-l.html
December_2019.html
Encyclopedia.html
English_language.html
Euler%27s_theorem_in_ge...etry.html
File/190309....jpg.html
File/ADIL_A...).jpg.html
File/Hurrica...r.gif.html
Free_content.html

# Improvement that Failed

# Tag

## Singleton Pattern

```java
public static Tag valueOf(String tagName, ParseSettings settings) {
    Validate.notNull(tagName);

    Tag tag = tags.get(tagName);

    if (tag == null) {
        tagName = settings.normalizeTag(tagName);
        Validate.notEmpty(tagName);
        tag = tags.get(tagName);

        if (tag == null) {
        // zgy: Handle multithreading
            synchronized(Tag.class) {
                tag = tags.get(tagName);
                if (tag == null) {
                    tag = new Tag(tagName);
                    tag.isBlock = false;
                    register(tag);
                }
            }
        }
    }
    return tag;
}
```

# Unit Test

# Test Result