

sep 10, 15 3:52

UrlLogger.java

Page 1/1

```

1 package Monitor;
2
3 import java.io.UnsupportedEncodingException;
4 import java.util.concurrent.BlockingQueue;
5 import java.io.FileNotFoundException;
6 import java.io.PrintWriter;
7
8 /**
9  * Created by adrian on 10/09/15.
10 */
11 public class UrlLogger implements Runnable {
12
13     private BlockingQueue<String> urlQueue;
14     private PrintWriter writer;
15
16     public UrlLogger(BlockingQueue<String> urlQueue) {
17         this.urlQueue = urlQueue;
18     }
19
20     public void initializeLogger (String fileName) {
21         try {
22             writer = new PrintWriter(fileName, "UTF-8");
23         } catch (FileNotFoundException e) {
24             e.printStackTrace();
25         } catch (UnsupportedEncodingException e) {
26             e.printStackTrace();
27         }
28     }
29
30     public void run() {
31         while (!Thread.interrupted()) {
32             try {
33                 String url = urlQueue.take();
34                 writer.println(url);
35                 writer.flush();
36             } catch (InterruptedException e) {
37                 e.printStackTrace();
38             }
39         }
40     }
41 }
42
43 }
```

sep 09, 15 21:15

UrlIncreaseMessage.java

Page 1/1

```

1 package Monitor;
2
3 /**
4  * Created by adrian on 09/09/15.
5  */
6 public class UrlIncreaseMessage implements MonitorMessage {
7     public void updateMonitor(MonitorStats monitor) {
8         monitor.increaseAnalyzedUrl();
9     }
10 }
```

sep 10, 15 14:08

ThreadUpdateMessage.java

Page 1/1

```

1 package Monitor;
2
3 import Crawler.ThreadState;
4
5 /**
6  * Created by adrian on 09/09/15.
7  */
8 public class ThreadUpdateMessage implements MonitorMessage {
9
10     private ThreadState state;
11
12     public ThreadUpdateMessage (ThreadState state) {
13         this.state = state;
14     }
15
16     public void updateMonitor(Monitor.MonitorStats monitor) {
17         monitor.setStatus(this.state);
18     }
19 }

```

sep 10, 15 14:08

MonitorWriter.java

Page 1/1

```

1 package Monitor;
2
3 import java.io.*;
4 import java.util.Date;
5
6 /**
7  * Created by adrian on 09/09/15.
8  */
9 public class MonitorWriter implements Runnable {
10
11     private MonitorStats stats;
12
13     private PrintWriter writer;
14     private Integer logInterval;
15
16     private Date flushDate;
17     private Integer flushInterval;
18
19     public MonitorWriter(MonitorStats stats) {
20         this.stats = stats;
21     }
22
23     public void initializeMonitor (String fileName, Integer logInterval, Integer
flushInterval) {
24         flushDate = new Date();
25         this.logInterval = logInterval;
26         this.flushInterval = flushInterval;
27         try {
28             writer = new PrintWriter(fileName, "UTF-8");
29         } catch (FileNotFoundException e) {
30             e.printStackTrace();
31         } catch (UnsupportedEncodingException e) {
32             e.printStackTrace();
33         }
34     }
35
36     public void run() {
37         while (!Thread.interrupted()) {
38             String statsString = stats.getFormattedStats();
39             writer.println(statsString);
40
41             if (((new Date().getTime() - this.flushDate.getTime()) / 1000) > flu
shInterval) {
42                 writer.flush();
43                 this.flushDate = new Date();
44             }
45
46             try {
47                 Thread.sleep((long) logInterval * 1000);
48             } catch (InterruptedException e) {
49                 e.printStackTrace();
50             }
51         }
52     }
53 }

```

sep 10, 15 14:08

MonitorStats.java

Page 1/3

```

1 package Monitor;
2
3 import Crawler.ThreadState;
4
5 import java.text.DateFormat;
6 import java.text.SimpleDateFormat;
7 import java.util.*;
8
9 /**
10  * Created by adrian on 09/09/15.
11  */
12 public class MonitorStats {
13
14     private Integer urlsCount;
15
16     private HashMap<String, Integer> fileTypes;
17
18     private HashMap<Integer, ThreadState.Status> urlAnalyzerState;
19     private HashMap<Integer, ThreadState.Status> fileAnalyzerState;
20     private HashMap<Integer, ThreadState.Status> fileDownloaderState;
21
22     public MonitorStats() {
23         urlsCount = 0;
24         fileTypes = new HashMap<String, Integer>();
25         urlAnalyzerState = new HashMap<Integer, ThreadState.Status>();
26         fileAnalyzerState = new HashMap<Integer, ThreadState.Status>();
27         fileDownloaderState = new HashMap<Integer, ThreadState.Status>();
28     }
29
30     public synchronized void initMonitorStats (Integer urlAnalyzers, Integer fileAnalyzers, Integer fileDownloaders) {
31         // I synchronize this block to be consistent, in that every public method should be atomic, but it's not necessary
32         this.initializeStats(fileAnalyzerState, fileAnalyzers);
33         this.initializeStats(fileDownloaderState, fileDownloaders);
34         this.initializeStats(urlAnalyzerState, urlAnalyzers);
35
36         urlsCount = 0;
37     }
38
39     private void initializeStats(HashMap<Integer, ThreadState.Status> map, Integer amount) {
40         for (Integer i = 0; i < amount; ++i) {
41             map.put(i, ThreadState.Status.UNKNOWN);
42         }
43     }
44
45     public synchronized void setStatus(ThreadState state) {
46         HashMap<Integer, ThreadState.Status> map = null;
47         switch (state.getThreadName()) {
48             case URL_ANALYZER:
49                 map = urlAnalyzerState;
50                 break;
51             case HTML_ANALYZER:
52                 map = fileAnalyzerState;
53                 break;
54             case FILE_DOWNLOADER:
55                 map = fileDownloaderState;
56                 break;
57         }
58         map.put(state.getThreadId(), state.getThreadStatus());
59     }
60
61     public synchronized void increaseAnalyzedUrl() {
62         urlsCount++;
63     }

```

sep 10, 15 14:08

MonitorStats.java

Page 2/3

```

64     }
65
66     public synchronized void increaseFileType(String fileType) {
67         if (fileTypes.containsKey(fileType)) {
68             fileTypes.put(fileType, fileTypes.get(fileType) + 1);
69         } else {
70             fileTypes.put(fileType, 1);
71         }
72     }
73
74     public synchronized String getFormattedStats () {
75         String result;
76         DateFormat formatter = new SimpleDateFormat("yyyy/MM/dd HH:mm:ss");
77         String stringDate = formatter.format(new Date());
78
79         result = stringDate + " - Analyzed URLs: " + urlsCount + "\n";
80         result = result + "Downloaded files by type:\n" + this.getFileTypesStats()
81         + "\n";
82         result = result + this.getThreadStats("Analyzing URL", urlAnalyzerState);
83         result = result + this.getThreadStats("Downloading File", fileDownloaderState);
84         result = result + this.getThreadStats("Analyzing File", fileAnalyzerState);
85
86         return result;
87     }
88
89     private String getFileTypesStats () {
90         String filesStats = "";
91         for (Map.Entry pair : fileTypes.entrySet()) {
92             filesStats = filesStats + "\t" + pair.getKey() + ":" + pair.getValue() + "\n";
93         }
94         return filesStats;
95     }
96
97     private String getThreadStats (String threadName, HashMap<Integer, ThreadState.Status> map) {
98         String threadStats;
99         Integer working = 0;
100         Integer blocked = 0;
101         Integer unknown = 0;
102         Integer starting = 0;
103
104         for (Map.Entry pair : map.entrySet()) {
105             ThreadState.Status status = (ThreadState.Status) pair.getValue();
106             switch (status) {
107                 case UNKNOWN:
108                     unknown++;
109                     break;
110                 case STARTING:
111                     starting++;
112                     break;
113                 case BLOCKED:
114                     blocked++;
115                     break;
116                 case WORKING:
117                     working++;
118                     break;
119             }
120         }
121
122         threadStats = "\t" + threadName + "\n\t\t";
123         threadStats = threadStats + "Unknown: " + unknown;
124         threadStats = threadStats + "Starting: " + starting;

```

sep 10, 15 14:08

MonitorStats.java

Page 3/3

```
124     threadStats = threadStats + " Blocked: " + blocked;
125     threadStats = threadStats + " Working: " + working + "\n";
126     return threadStats;
127 }
128 }
```

sep 09, 15 21:28

MonitorMessage.java

Page 1/1

```
1 package Monitor;
2
3 /**
4  * Created by adrian on 09/09/15.
5  */
6 public interface MonitorMessage {
7
8     public void updateMonitor(MonitorStats monitor);
9 }
```

sep 10, 15 14:08

MonitorFetcher.java

Page 1/1

```

1 package Monitor;
2
3 import java.util.concurrent.BlockingQueue;
4
5 /**
6  * Created by adrian on 09/09/15.
7  */
8 public class MonitorFetcher implements Runnable {
9
10     private MonitorStats stats;
11     private BlockingQueue<MonitorMessage> monitorQueue;
12
13     public MonitorFetcher(MonitorStats stats, BlockingQueue<MonitorMessage> moni
14 torQueue) {
15         this.stats = stats;
16         this.monitorQueue = monitorQueue;
17     }
18
19     public void run() {
20         while (!Thread.interrupted()) {
21             try {
22                 MonitorMessage message = monitorQueue.take();
23                 message.updateMonitor(stats);
24             } catch (InterruptedException e) {
25                 e.printStackTrace();
26             }
27         }
28     }
29 }

```

sep 10, 15 14:08

FileTypeUpdateMessage.java

Page 1/1

```

1 package Monitor;
2
3 /**
4  * Created by adrian on 09/09/15.
5  */
6 public class FileTypeUpdateMessage implements MonitorMessage {
7
8     private String fileType;
9
10    public FileTypeUpdateMessage(String type) {
11        this.fileType = type;
12    }
13
14    public void updateMonitor(MonitorStats monitor) {
15        monitor.increaseFileType(this.fileType);
16    }
17 }

```

sep 10, 15 14:08

Launcher.java

Page 1/3

```

1  import Crawler.*;
2  import Monitor.*;
3
4  import java.io.*;
5  import java.net.URI;
6  import java.net.URISyntaxException;
7  import java.util.Properties;
8  import java.util.concurrent.*;
9
10 /**
11  * Created by adrian on 05/09/15.
12  */
13 public class Launcher {
14
15     public static void main(String[] args) {
16
17         Properties prop = new Properties();
18         InputStream input = null;
19
20         try {
21             String configFile = "Config/Config.properties";
22             File f = new File(configFile);
23             if(!f.exists()) {
24                 System.err.println("Config file " + configFile + " missing");
25                 System.err.println("Format:\n" +
26                     "####Config file start####\n" +
27                     "iterations = #\n" +
28                     "fileAnalyzer = #\n" +
29                     "fileDownloader = #\n" +
30                     "urlAnalyzer = #\n" +
31                     "queueSize = #\n" +
32                     "urlLogFile = nameUrlLogFile.log\n" +
33                     "logFile = nameMonitorLogFile.log\n" +
34                     "logInterval = #\n" +
35                     "flushInterval = #\n" +
36                     "####Config file end####");
37                 System.exit(-1);
38             }
39             input = new FileInputStream(configFile);
40
41             prop.load(input);
42
43             int queueSize = Integer.parseInt(prop.getProperty("queueSize"));
44
45             int numberOfFileAnalyzer = Integer.parseInt(prop.getProperty("fileAnalyzer"));
46             int numberOfUrlAnalyzer = Integer.parseInt(prop.getProperty("fileDownloader"));
47             int numberOfFileDownloader = Integer.parseInt(prop.getProperty("urlAnalyzer"));
48
49             // Monitor
50             MonitorStats stats = new MonitorStats();
51             stats.initMonitorStats(numberOfUrlAnalyzer, numberOfFileAnalyzer, numberOfFileDownloader);
52
53             BlockingQueue<MonitorMessage> monitorQueue = new ArrayBlockingQueue<MonitorMessage>(queueSize);
54             MonitorWriter writer = new MonitorWriter(stats);
55             writer.initializeMonitor(prop.getProperty("logFile"), Integer.parseInt(prop.getProperty("logInterval")), Integer.parseInt(prop.getProperty("flushInterval")));
56             (new Thread(writer)).start();
57             MonitorFetcher fetcher = new MonitorFetcher(stats, monitorQueue);
58             (new Thread(fetcher)).start();
59
60             // Url Logger

```

sep 10, 15 14:08

Launcher.java

Page 2/3

```

61     BlockingQueue<String> urlLoggerQueue = new ArrayBlockingQueue<String>(queueSize);
62     UrlLogger urlLogger = new UrlLogger(urlLoggerQueue);
63     urlLogger.initializeLogger(prop.getProperty("urlLogFile"));
64     new Thread(urlLogger).start();
65
66     // Crawler Queues
67     BlockingQueue<UrlMessage> urlToAnalyzeQueue = new ArrayBlockingQueue<UrlMessage>(queueSize);
68     BlockingQueue<Pair<String, UrlMessage>> fileToAnalyzedQueue = new ArrayBlockingQueue<Pair<String, UrlMessage>>(queueSize);
69     BlockingQueue<UrlMessage> urlToDownloadQueue = new ArrayBlockingQueue<UrlMessage>(queueSize);
70
71     ConcurrentHashMap<String, Boolean> hashMap = new ConcurrentHashMap<String, Boolean>();
72
73     // Crawler workers
74     for (int i = 0; i < numberOfFileAnalyzer; ++i) {
75         (new Thread(new FileAnalyzer(i, fileToAnalyzedQueue, urlToAnalyzeQueue, monitorQueue))).start();
76     }
77
78     int maxIterations = Integer.parseInt(prop.getProperty("iterations"));
79     for (int i = 0; i < numberOfUrlAnalyzer; ++i) {
80         (new Thread(new FileDownloader(i, maxIterations, urlToDownloadQueue, fileToAnalyzedQueue, monitorQueue))).start();
81     }
82
83     for (int i = 0; i < numberOfFileDownloader; ++i) {
84         (new Thread(new UrlAnalyzer(i, urlToAnalyzeQueue, urlToDownloadQueue, hashMap, monitorQueue, urlLoggerQueue))).start();
85     }
86
87     while(!Thread.interrupted()) {
88         BufferedReader br = new BufferedReader(new InputStreamReader(System.in));
89         System.out.print("Enter Url:");
90         String src = br.readLine();
91
92         try {
93             URI u = new URI(src);
94             if (u.isAbsolute()) {
95                 urlToAnalyzeQueue.put(new UrlMessage(0, src));
96             }
97         } catch (URISyntaxException e) {
98             e.printStackTrace();
99         }
100     }
101
102     System.out.println("Finishing web crawler");
103
104     } catch (IOException ex) {
105         ex.printStackTrace();
106     } catch (InterruptedException e) {
107         e.printStackTrace();
108     } finally {
109         if (input != null) {
110             try {
111                 input.close();
112             } catch (IOException e) {
113                 e.printStackTrace();
114             }
115         }
116     }
117 }

```

sep 10, 15 14:08

Launcher.java

Page 3/3

```
118     }  
119 }  
120 }
```

sep 10, 15 2:19

UrlMessage.java

Page 1/1

```
1 package Crawler;  
2  
3 /**  
4  * Created by adrian on 10/09/15.  
5  */  
6 public class UrlMessage {  
7  
8     private Integer iteration;  
9     private String url;  
10  
11     public UrlMessage (Integer iteration, String url) {  
12         this.iteration = iteration;  
13         this.url = url;  
14     }  
15  
16     public Integer getIteration() {  
17         return iteration;  
18     }  
19  
20     public String getUrl() {  
21         return url;  
22     }  
23 }
```

sep 10, 15 14:08

UrlAnalyzer.java

Page 1/2

```

1 package Crawler;
2
3 import Monitor.MonitorMessage;
4 import Monitor.ThreadUpdateMessage;
5 import Monitor.UrlIncreaseMessage;
6
7 import java.util.concurrent.BlockingQueue;
8 import java.util.concurrent.ConcurrentHashMap;
9
10 /**
11  * Created by adrian on 05/09/15.
12  */
13
14 public class UrlAnalyzer implements Runnable {
15
16     private Integer threadId;
17     private BlockingQueue<MonitorMessage> monitorQueue;
18     private BlockingQueue<String> urlLoggerQueue;
19
20     private BlockingQueue<UrlMessage> urlToAnalyzeQueue;
21     private BlockingQueue<UrlMessage> urlToDownloadQueue;
22
23     private ConcurrentHashMap<String,Boolean> analyzedUrls;
24
25     public UrlAnalyzer(Integer threadId, BlockingQueue<UrlMessage> analyzeQueue,
26 BlockingQueue<UrlMessage> downloadQueue, ConcurrentHashMap<String,Boolean> map,
27 BlockingQueue<MonitorMessage> monitorQueue, BlockingQueue<String> urlLoggerQueue) {
28         this.threadId = threadId;
29         this.monitorQueue = monitorQueue;
30         this.urlLoggerQueue = urlLoggerQueue;
31         this.urlToAnalyzeQueue = analyzeQueue;
32         this.urlToDownloadQueue = downloadQueue;
33         this.analyzedUrls = map;
34     }
35
36     public void run() {
37
38         try {
39             this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Thread
40 State.Type.URL_ANALYZER, this.threadId, ThreadState.Status.STARTING)));
41         } catch (InterruptedException e) {
42             e.printStackTrace();
43         }
44
45         while (!Thread.interrupted()) {
46             try {
47                 this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Thread
48 State.Type.URL_ANALYZER, this.threadId, ThreadState.Status.BLOCKED)));
49                 UrlMessage url = urlToAnalyzeQueue.take();
50                 this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Thread
51 State.Type.URL_ANALYZER, this.threadId, ThreadState.Status.WORKING)));
52
53                 System.out.println("Analyzing url: "+url.getUrl());
54                 if (analyzedUrls.putIfAbsent(url.getUrl(), true) == null) {
55                     System.out.println("Url NOT downloaded " + url.getUrl());
56                     this.urlLoggerQueue.put(url.getUrl());
57                     urlToDownloadQueue.put(url);
58                 } else {
59                     // Url already downloaded
60                     System.out.println("Url already downloaded " + url.getUrl());
61                 }
62
63                 this.monitorQueue.put(new UrlIncreaseMessage());
64             }
65         }
66     }
67 }

```

sep 10, 15 14:08

UrlAnalyzer.java

Page 2/2

```

61         } catch (InterruptedException e) {
62             e.printStackTrace();
63         }
64     }
65 }
66 }
67 }

```


sep 10, 15 1:03

ThreadState.java

Page 1/1

```

1 package Crawler;
2
3 /**
4  * Created by adrian on 08/09/15.
5  */
6
7 public class ThreadState {
8
9     private Type threadName;
10    private Integer threadId;
11    private Status threadStatus;
12
13    public enum Status {
14        UNKNOWN,
15        STARTING,
16        BLOCKED,
17        WORKING
18    }
19
20    public enum Type {
21        URL_ANALYZER,
22        HTML_ANALYZER,
23        FILE_DOWNLOADER
24    }
25
26    public ThreadState(Type name, int id, Status status) {
27        threadName = name;
28        threadId = id;
29        threadStatus = status;
30    }
31
32    public Type getThreadName() {
33        return threadName;
34    }
35
36    public int getThreadId() {
37        return threadId;
38    }
39
40    public Status getThreadStatus() {
41        return threadStatus;
42    }
43 }

```

sep 10, 15 14:08

Pair.java

Page 1/1

```

1 package Crawler;
2
3 /**
4  * Created by adrian on 06/09/15.
5  */
6 public class Pair<A, B> {
7     private A first;
8     private B second;
9
10    public Pair(A first, B second) {
11        super();
12        this.first = first;
13        this.second = second;
14    }
15
16    public A getFirst() {
17        return first;
18    }
19
20    public B getSecond() {
21        return second;
22    }
23 }

```

sep 10, 15 14:08

FileDownloader.java

Page 1/3

```

1 package Crawler;
2
3 import Monitor.FileTypeUpdateMessage;
4 import Monitor.MonitorMessage;
5 import Monitor.ThreadUpdateMessage;
6
7 import java.io.*;
8 import java.net.*;
9 import java.nio.file.*;
10 import java.nio.file.Files;
11 import java.util.concurrent.*;
12
13 public class FileDownloader implements Runnable {
14
15     private Integer threadId;
16     private Integer iterations;
17     private BlockingQueue<MonitorMessage> monitorQueue;
18     private BlockingQueue<UrlMessage> urlToDownloadQueue;
19     private BlockingQueue<Pair<String, UrlMessage>> fileToAnalyzeQueue;
20
21     public FileDownloader(Integer threadId, Integer iterations, BlockingQueue<Ur
lMessage> downloadQueue, BlockingQueue<Pair<String, UrlMessage>> fileAnalyzeQueu
e, BlockingQueue<MonitorMessage> monitorQueue) {
22         this.threadId = threadId;
23         this.iterations = iterations;
24         this.monitorQueue = monitorQueue;
25         this.urlToDownloadQueue = downloadQueue;
26         this.fileToAnalyzeQueue = fileAnalyzeQueue;
27     }
28
29     public void run() {
30         try {
31             this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Thread
State.Type.FILE_DOWNLOADER, this.threadId, ThreadState.Status.STARTING)));
32         } catch (InterruptedException e) {
33             e.printStackTrace();
34         }
35
36         while (!Thread.interrupted()) {
37             try {
38                 this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Th
readState.Type.FILE_DOWNLOADER, this.threadId, ThreadState.Status.BLOCKED)));
39                 UrlMessage url = urlToDownloadQueue.take();
40                 this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Th
readState.Type.FILE_DOWNLOADER, this.threadId, ThreadState.Status.WORKING)));
41                 downloadFile(url, "Downloads");
42             } catch (InterruptedException e) {
43                 e.printStackTrace();
44             }
45         }
46
47         private void downloadFile(UrlMessage message, String destinationDir) throws
InterruptedException {
48             String localFileName;
49             int slashIndex = message.getUrl().lastIndexOf('/');
50             // We should remode tha last / from the url because we use the string ne
xt to the slash to use as file name
51             if (slashIndex == message.getUrl().length() - 1) {
52                 String noLastSlashUrl = message.getUrl().substring(0,message.getUrl(
).length() - 1);
53                 slashIndex = noLastSlashUrl.lastIndexOf('/');
54                 localFileName = noLastSlashUrl.substring(slashIndex + 1);
55             } else {
56                 localFileName = message.getUrl().substring(slashIndex + 1);
57             }
58         }

```

sep 10, 15 14:08

FileDownloader.java

Page 2/3

```

59
60
61         if (localFileName.length() > 0) {
62             File dir = new File(destinationDir);
63
64             // We create the download directory if it doesn't exists
65             if (!dir.exists()) {
66                 dir.mkdirs();
67             }
68
69             OutputStream outStream = null;
70             URLConnection connection;
71             InputStream inStream = null;
72             try {
73                 URL url;
74                 byte[] buf;
75                 int bytesRead;
76                 url = new URL(message.getUrl());
77                 connection = url.openConnection();
78
79                 outStream = new BufferedOutputStream(new FileOutputStream(destin
ationDir + "/" + localFileName));
80                 inStream = connection.getInputStream();
81
82                 buf = new byte[1024];
83                 while ((bytesRead = inStream.read(buf)) != -1) {
84                     outStream.write(buf, 0, bytesRead);
85                 }
86
87                 String fileType = Files.probeContentType(FileSystems.getDefault(
).getPath(destinationDir, localFileName));
88                 if ((message.getIteration() < this.iterations-1) ^ fileType.equ
als("text/html")) {
89                     fileToAnalyzeQueue.put(new Pair<String, UrlMessage>(destinat
ionDir + "/" + localFileName, new UrlMessage(message.getIteration()+1,message.ge
tUrl())));
90                 }
91
92                 this.monitorQueue.put(new FileTypeUpdateMessage(fileType));
93
94             } catch (MalformedURLException e) {
95                 e.printStackTrace();
96             } catch (FileNotFoundException e) {
97                 e.printStackTrace();
98             } catch (IOException e) {
99                 e.printStackTrace();
100             } finally {
101
102                 if (inStream != null) {
103                     try {
104                         inStream.close();
105                     } catch (IOException e) {
106                         e.printStackTrace();
107                     }
108                 }
109
110                 if (outStream != null) {
111                     try {
112                         outStream.close();
113                     } catch (IOException e) {
114                         e.printStackTrace();
115                     }
116                 }
117             }
118         }
119     }

```

sep 10, 15 14:08

FileDownloader.java

Page 3/3

120 }

sep 10, 15 14:04

FileAnalyzer.java

Page 1/2

```

1  package Crawler;
2
3  import Monitor.MonitorMessage;
4  import Monitor.ThreadUpdateMessage;
5
6  import javax.swing.text.AttributeSet;
7  import javax.swing.text.html.HTML;
8  import javax.swing.text.html.HTMLDocument;
9  import javax.swing.text.html.HTMLEditorKit;
10 import javax.swing.text.html.parser.ParserDelegator;
11
12 import java.io.*;
13 import java.net.URI;
14 import java.net.URISyntaxException;
15
16 import java.util.concurrent.BlockingQueue;
17
18
19 /**
20  * Created by adrian on 04/09/15.
21  */
22 public class FileAnalyzer implements Runnable {
23
24     private Integer threadId;
25     private BlockingQueue<MonitorMessage> monitorQueue;
26     private BlockingQueue<UrlMessage> urlToAnalyzeQueue;
27     private BlockingQueue<Pair<String, UrlMessage>> fileToAnalyzedQueue;
28
29     public FileAnalyzer(Integer threadId, BlockingQueue<Pair<String, UrlMessage>
> fileQueue, BlockingQueue<UrlMessage> urlQueue, BlockingQueue<MonitorMessage> m
onitorQueue) {
30         this.threadId = threadId;
31         this.monitorQueue = monitorQueue;
32         this.urlToAnalyzeQueue = urlQueue;
33         this.fileToAnalyzedQueue = fileQueue;
34     }
35
36     public void run() {
37         try {
38             this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Thread
State.Type.HTML_ANALYZER, this.threadId, ThreadState.Status.STARTING)));
39         } catch (InterruptedException e) {
40             e.printStackTrace();
41         }
42
43         while (!Thread.interrupted()) {
44             try {
45                 this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Th
readState.Type.HTML_ANALYZER, this.threadId, ThreadState.Status.BLOCKED)));
46                 Pair<String, UrlMessage> file = fileToAnalyzedQueue.take();
47                 this.monitorQueue.put(new ThreadUpdateMessage(new ThreadState(Th
readState.Type.HTML_ANALYZER, this.threadId, ThreadState.Status.WORKING)));
48                 System.out.println("Analyzing file: "+file.getFirst());
49                 analyzeFile(file.getFirst(), file.getSecond());
50             } catch (InterruptedException e) {
51                 e.printStackTrace();
52             }
53         }
54     }
55
56     public void analyzeFile (String fileName, UrlMessage url) throws Interrupted
Exception {
57         InputStream is = null;
58         try {
59             is = new FileInputStream(fileName);
60         } catch (FileNotFoundException e) {

```

sep 10, 15 14:04

FileAnalyzer.java

Page 2/2

```

61         e.printStackTrace();
62     }
63
64     InputStreamReader isr = new InputStreamReader(is);
65     BufferedReader br = new BufferedReader(isr);
66
67     HTMLToolkit htmlKit = new HTMLToolkit();
68     HTMLDocument htmlDoc = (HTMLDocument) htmlKit.createDefaultDocument();
69     HTMLToolkit.Parser parser = new ParserDelegator();
70     HTMLToolkit.ParserCallback callback = htmlDoc.getReader(0);
71
72     try {
73         parser.parse(br, callback, true);
74     } catch (IOException e) {
75         e.printStackTrace();
76     }
77
78     getResource(url, htmlDoc, HTML.Tag.IMG, HTML.Attribute.SRC);
79     getResource(url, htmlDoc, HTML.Tag.A, HTML.Attribute.HREF);
80     getResource(url, htmlDoc, HTML.Tag.LINK, HTML.Attribute.HREF);
81     getResource(url, htmlDoc, HTML.Tag.SCRIPT, HTML.Attribute.SRC);
82 }
83
84 private void getResource(UrlMessage url, HTMLDocument htmlDoc, HTML.Tag tag,
85 HTML.Attribute attribute) throws InterruptedException {
86     for (HTMLDocument.Iterator iterator = htmlDoc.getIterator(tag); iterator
87 .isValid(); iterator.next()) {
88         AttributeSet attributes = iterator.getAttributes();
89         String src = (String) attributes.getAttribute(attribute);
90
91         if (src != null) {
92             URI u = null;
93             try {
94                 u = new URI(src);
95                 String finalUrl;
96                 if (u.isAbsolute()) {
97                     finalUrl = src;
98                 } else {
99                     if (src.startsWith("//")) {
100                         int index = url.getUrl().indexOf(":");
101                         String protocol = url.getUrl().substring(0, index);
102                         finalUrl = protocol + "://" + src;
103                     } else {
104                         finalUrl = u.resolve(url.getUrl()).normalize().toString();
105                     }
106                 }
107                 urlToAnalyzeQueue.put(new UrlMessage(url.getIteration(), finalUrl));
108             } catch (URISyntaxException e) {
109                 e.printStackTrace();
110             }
111         }
112     }
113 }
114 }

```

sep 10, 15 14:24

Table of Content

Page 1/1

1	Table of Contents																		
2	1 <i>UrlLogger.java</i>	sheets	1 to	1 (1) pages	1- 1	44 lines													
3	2 <i>UrlIncreaseMessage.java</i>	sheets	1 to	1 (1) pages	2- 2	11 lines													
4	3 <i>ThreadUpdateMessage.java</i>	sheets	2 to	2 (1) pages	3- 3	20 lines													
5	4 <i>MonitorWriter.java</i> ..	sheets	2 to	2 (1) pages	4- 4	54 lines													
6	5 <i>MonitorStats.java</i> ...	sheets	3 to	4 (2) pages	5- 7	129 lines													
7	6 <i>MonitorMessage.java</i> .	sheets	4 to	4 (1) pages	8- 8	10 lines													
8	7 <i>MonitorFetcher.java</i> .	sheets	5 to	5 (1) pages	9- 9	30 lines													
9	8 <i>FileTypeUpdateMessage.java</i>	sheets	5 to	5 (1) pages	10- 10	18 lines													
10	9 <i>Launcher.java</i>	sheets	6 to	7 (2) pages	11- 13	121 lines													
11	10 <i>UrlMessage.java</i>	sheets	7 to	7 (1) pages	14- 14	24 lines													
12	11 <i>UrlAnalyzer.java</i>	sheets	8 to	8 (1) pages	15- 16	68 lines													
13	12 <i>ThreadState.java</i>	sheets	9 to	9 (1) pages	17- 17	44 lines													
14	13 <i>Pair.java</i>	sheets	9 to	9 (1) pages	18- 18	24 lines													
15	14 <i>FileDownloader.java</i> .	sheets	10 to	11 (2) pages	19- 21	121 lines													
16	15 <i>FileAnalyzer.java</i> ...	sheets	11 to	12 (2) pages	22- 23	115 lines													