

File Statistics 2.0a Component Specification

1. Design

File Statistics provides a framework to collect statistic data from a subset of the file system and renders the format in configurable reporting style.

As a framework component the first question comes into mind when designing the component would be ***what aspects of the component need to be pluggable***. It's obvious that the algorithm to extract the statistical data as well as the way to render the data should be typically configurable. A third aspect that is factored out is the filtering logic to specify which subset of the files are processed by certain processor. Though this could be tied to the processor the design would not as flexible that way. Once the three aspects are cleared the strategy design patterns come into play.

Once we have the interfaces the next question would be ***how to glue them together***. The design finally chooses to use a non-sharable manager class because the use of the component involves multiple steps and result is aggregated when each file is processed. (The design does not prevent from using stateless reporting mechanism, like logging, though.) The manager also fulfills the task to instantiate the framework from config files for better configurability and for the use case of the command line interface. Each of the pluggable implementation has contract constructor to configure itself.

Once the framework is settled the rest of the design is to produce plug-ins to fulfill the requirements. The design should be relatively easy to understand.

Changes since version 2.0a

The version 2.0a requires the component to support non-numeric data. The following paragraphs focus on detailing the changes in order to accommodate this requirement.

➤ **Custom Data**

In the first release all data are stored as long. An interface needs to be specified so that non-numeric result can be held and used. This interface comes to be named Statistic. It should define the minimum contract to ease implementation. Normally the data go through three processes – generation, aggregation and render. The minimum operations required are a clone(), an aggregate() as well as a render().

Now that Statistic is specified StatsResult should be changed to hold mapping to Statistic instead of Long. However the original long version of the setStat() method is still supported as a convenience, with the default implementation of LongStatistic.

➤ **Custom Renderer**

Letting the data render itself probably solves the problem with some very simple data (for instance no matter how a long will probably always render itself as its numeric representation) but this does not work for complex data. For example to render the most commonly occurring tokens as plain text we would probably write:

Most commonly occurring tokens are apple, sky which occurred 17 times each.

However as XML it would probably be:

```
<token>apple</token><token>sky</token><occurrence>17</occurrence>
```

Using a consistent format does not make sense in such situations. This is where the renderers come into the play. The most natural point to add this functionality would be AbstractStatsReporting. The renderers are mapped to metrics which provide more flexibility than mapped to the data type. As usually the renderers need to be configurable so that this works with command line interface.

However the render() on Statistic is still supported. This should be quite convenient for simple data (a long probably always renders as the numeric representation). Renderers work as an overridden behavior for the self-render.

➤ **Common Token Analyzer**

The upgrade will not be complete without a real processor that generates custom data. This should be interpreted as both a demonstration and an implementation guide. Over the examples provided by requirement the common token one is chosen because it introduces another tricky point about the process.

Some custom data are not immediately aggregatable. If we only hold the tokens as well as the occurrence we can not aggregate results from two files together, because we do not know if a less occurring word in both files could become the most occurring one. This means for some data we need to hold more data than necessary in order to accurately aggregate.

This is not a problem itself, but since we can hold a lot of entries on a reporting structure. There is one structure for each file, which actually do not need to aggregate at all. This could cost a lot of memory. To avoid this we declare a release() method on the Statistic interface to release unnecessary resources (and only keep those for rendering).

1.1 Design Patterns

- **Strategy Pattern:** used for FileMatcher, FileProcessor and StatsReporting and their respective implementations. This is essential for a framework component in order to provide maximum pluggability. Since version 2.0a Statistic and StatisticRenderer interfaces are added which also implements the strategy pattern.
- **Template Method Pattern:** used for FileNameMatcher, FileCategoryMatcher and FileTypeMatcher to minimize redundancy for extensions.

1.2 Industry Standards

None

1.3 Required Algorithms

1.3.1 Counting Lines in C Style Source Codes

For this release we will include a naive C style source line counter. It provides moderately accurate figure with a relatively simple implementation. This does not necessarily work for some cases.

The design recommends a state machine implementation. Here is the transition table.

State	Input	Action
START	slash(/)	SLASH
	comma(,)	Increment line
	open brace({)	Increment line
	other	START
SLASH	star(*)	BLOCK_COMMENT
	slash(/)	LINE_COMMENT
	other	START
STAR	slash(/)	START
	start(*)	STAR
	other	BLOCK_COMMENT
SQ_BACKSLASH	any	SINGLE_QUOTE
DQ_BACKSLASH	any	DOUBLE_QUOTE
SINGLE_QUOTE	backslash(\)	SQ_BACKSLASH
	single quote(')	START
	other	SINGLE_QUOTE
DOUBLE_QUOTE	backslash(\)	DQ_BACKSLASH
	double quote(")	START

	other	DOUBLE_QUOTE
LINE_COMMENT	any	LINE_COMMENT
	line feed	START
BLOCK_COMMENT	star(*)	STAR
	other	BLOCK_COMMENT

1.3.2 Command Line Interface

The command line interface should be powerful to expose all the functionality and flexibility to the user, which should still be easy to use (meaning most options should have well defined defaults).

Here is what we have got:

[-c config-file] [-n namespace] [-b base-path] [-r] [-o output-file] path [path2 path3 ...]

➤ ***-c config-file***

This specifies the configuration file to use. If the namespace is provided at the same time it will be loaded to the namespace in order to instantiate the component. Otherwise it is loaded to default namespace.

➤ ***-n namespace***

This specifies the namespace to use. If the config file is provided at the same time the config file is loaded to the namespace. Otherwise the namespace should be preloaded. If neither config file nor namespace is provided the component instantiates from the default namespace (assuming preload).

➤ ***-b base-path***

This specifies the base directory for the reporting. If this is provided all path are supposed to be relative to this path and relative paths are used for reporting. If not provided all path are supposed to be either absolute to relative to executing path, with reporting based on absolute path.

➤ ***-r***

This specifies whether the directories are processed in recursive mode. If not provided directories do not go into sub-directories. If no directory is processed this switch does not make any effects.

➤ ***-o output-file***

This specifies the output file which is either absolute or relative to executing path (does not rely on base-path). If not provided reporting is printed to console.

➤ ***path [path2 path3 ...]***

This specifies a list of files or directories to process. If base directory is provided they are relative, otherwise absolute to relative to executing path. Files and/or directories can mix.

1.3.3 Xml Reporting

In order to provide better human readability and easy stylesheet support, the output xml uses dynamic element names for the stats result. The display name of the metric will be used to render the element node (which consequently requires well-formed name to be used but the framework does not validate over that). Refer to the sample output which should be pretty clear in structure. For readability it is advised to use pretty formatting for the output. When multiple filename, directory name or aliases are listed, alphabetical order should be used (case-insensitive).

Since version 2.0a XML reporting needs to resort to custom renderers to render custom data. This could be problem in formatting. We should assume the custom renderer configured will return well formed XML segments which can be encapsulated in an element named after the metric's display name. There is no more pretty formatting recommendation for this, but supporting pretty formatting is a plus.

1.4 Component Class Overview

1.4.1 Package `com.topcoder.file.statistics`

- **FileStatistics:** This is the main class of the user API. It provides the basic API, which includes maintaining a list of file processors, a reporting handler and controls the way to process the file. Files are processed one by one and the results are aggregated into the reporting. If you want to start a new reporting you should either explicitly set the reporting or get the old reporting to reset it. There are two ways to process the file: process it with the first matching processor and process it with each matching processor. Refer to `ExtendedFileStatistics` for advanced usage.
- **ExtendedFileStatistics:** This is the extended version of `FileStatistics` which supports some advanced operations including directory processing and reporting printing. The reason to separate this out is that we may not need to rely on the directory traversing algorithm incorporated here (which is recursive).
- **Main:** The command line entry point of the File Statistics utility. It uses `ExtendedFileStatistics` to fulfill the task. When argument is invalid help message should be printed.
- **FileProcessor (interface):** `FileProcessor` defines the contract to process a file. One file processor defines an algorithm to produce statistical result for a certain class of files against certain metrics. (For instance count line numbers for a text file.) Multiple `FileProcessors` are aggregated into `FileStatistics`, which can run in two modes (process with first matching processor and process with each matching processor). File processors are mapped with an alias. This alias is dynamically provided by the user and the reporting can be grouped against this alias.
- **FileMatcher (interface):** `FileMatcher` defines the contract to match a certain file (in order for an attaching `FileProcessor` to process the file).
- **StatsReporting (interface):** `StatsReporting` interface defines the contract to aggregate and render the statistical report. Result is fed into this interface one by one and the file report is then generated. There is a way to reset (clear) the current aggregated results.
- **StatsResult:** `StatsResult` is container for multiple results, each of which are mapped from a `Metric` to a `Long`. This class wraps some map operations as well as provides a way to clone and aggregate results. [Since version 2.0a this class adds a few methods to support non-numeric Statistic data. The original `getStat\(StatsResult\):long` is not supported any more. Instead a `Statistic` is returned.](#)
- **Metric:** `Metric` represents a measurement of the stats result. A metric is identified by its name. Metrics with the same name will be aggregated. This class overrides `equals()` and `hashCode()` so that it can be used as map keys.
- **Statistic (interface):** This interface represents a single statistic that is generated by the processor. A processor can generate more than one `Statistic` but it's the basic unit that is aggregatable. This interface is added since version 2.0a in order to accommodate non-numerical statistic result. The data itself must handle the logic to aggregate, clone and simple render. Complex rendering is handled by custom `StatisticsRenderer` implementation.
- **LongStatistic:** This class is added in version 2.0a in order to support the original long type statistic data. It is also included as a demonstration of a simple data, while the one included in the Common Token processor is included as a demonstration of a complex data.

1.4.2 *Package com.topcoder.file.statistics.processor*

- **FileProcessorBase (abstract):** A base implementation that can be used for FileProcessors. It provides alias and matcher support, including loading them from configuration file. It serves TextLineCounter and CStyleLineCounter as well as can be generically used by future implementations.

1.4.3 *Package com.topcoder.file.statistics.processor.linecounter*

- **TextLineCounter:** TextLineCounter counts lines for text files. It supports a single metric of "Line Count".
- **CStyleLineCounter:** CStyleLineCounter counts lines for C style source files. This implementation can handle C, C++, C# or Java. It only takes braces and comma into account so it would not be very accurate in most cases. Comments and string literals can be handled. It supports a single metric of "Line Count".
- **SimpleStateMachine:** This is a private static inner class to CStyleLineCounter that implements a state machine that can process the source code. Notice this is an implementation recommendation. Developer can choose to improve from this start point. Dropping this class is acceptable as long as the implementation of CStyleLineCounter is stateless (and the clarity and/or efficiency aspects are improved).

1.4.4 *Package com.topcoder.file.statistics.processor.commontoken*

- **CommonTokenAnalyzer:** The CommonTokenProcessor looks for the most commonly occurring tokens in a text file. A token can be either defined with an alphabet or with a list of separators. A token can not span over lines. Since the constructor is a bit confusing otherwise, static factory methods should be used to create instance of this processor.
- **CommonTokenStatistic:** CommonTokenStatistic is created as the custom data generated by the CommonTokenAnalyzer. It also acts as the temporary structure in the processor (to aggregate the tokens).
- **CommonTokenPlainRenderer:** This renderer renders the CommonTokenStatistic as plain text.
- **CommonTokenXmlRenderer:** This renderer renders the CommonTokenStatistic as XML.

1.4.5 *Package com.topcoder.file.statistics.matcher*

- **AnyFileMatcher:** AnyFileMatcher is a trivial implementation that matches any file.
- **FileNameMatcher (abstract):** FileNameMatcher is a template class that handles filename oriented matching logic.
- **ExtensionMatcher:** FileNameMatcher concrete implementation that matches one or more file extensions. The extensions are matched in a case-sensitive manner.
- **RegexMatcher:** FileNameMatcher concrete implementation that matches file names with one or more regular expression.
- **FileCategoryMatcher (abstract):** FileCategoryMatcher is a template class that handles file category (text/binary) oriented matching logic. The purpose to still create subclasses for such minor functionality is about the convenience in configuration.
- **TextFileMatcher:** FileCategoryMatcher concrete implementation that matches text files.
- **BinaryFileMatcher:** FileCategoryMatcher concrete implementation that matches binary files.

- **FileTypeMatcher (abstract):** FileTypeMatcher is a template class that handles file type oriented matching logic. File type is resolved by the Magic Numbers component.
- **TypeNameMatcher:** FileTypeMatcher concrete implementation that matches file types with one or more type names (as configured in Magic Numbers).
- **MimeMatcher:** MimeMatcher concrete implementation that matches file types with one or more mimes (as configured in Magic Numbers).

1.4.6 Package *com.topcoder.file.statistics.reporting*

- **AbstractStatsReporting (abstract):** This class provides a StatsReporting implementation base. It aggregates results from each processor and classifies them based on single file, directory or alias. It also provides overall stats. It defers the actual rendering logic to concrete implementations (BasicStatsReporting and XmlStatsReporting). [During version 2.0 this class is changed to support custom renderers to render the Statistic data. A mapping from metric to StatisticRenderer is maintained and configurable through configuration manager.](#)
- **StatsCollection:** StatsCollection is a generic container for mapping from String keys to StatsResult instances. It is used in AbstractStatsReporting to hold file, directory and alias classified mappings. The class provides basic manipulation methods for the mapping.
- **BasicStatsReporting:** BasicStatsReporting is a reporting implementation that renders user-friendly reports. [Since version 2.0a the generateReport\(\) method is refactored but the API does not change.](#)
- **XmlStatsReporting:** XmlStatsReporting is a reporting implementation that renders XML reports. Refer to component specification for formatting of the XML. It is capable of using a stylesheet to transform the XML into user- oriented format (HTML, CSV, etc.) [Since version 2.0a the generateReport\(\) method is refactored but the API does not change.](#)
- **StatisticRenderer (interface):** This interface is added since version 2.0a to render custom (complex) Statistic. Implementation only needs to implement a single method. If implementation needs to be configured through configuration file it is required to have a public no-arg constructor.

1.5 Component Exception Definitions

1.5.1 Custom Exceptions

- **FileStatisticsException:** Framework exception which provides the extension base for all exceptions. It does not subclass Base Exception since JDK 1.3 is not supported.
- **ConfigurationException:** Used to cover configuration related exceptions for all the constructors with a namespace or with a property.
- **FileMatchingException:** Encapsulates errors from the FileMatcher interface to indicate a dependency error. If the file operation fails IOException is used instead.
- **FileProcessingException:** Encapsulates errors from the FileProcessor interface to indicate an implementation specific error. If the file operation fails IOException is used instead.
- **StatsReportingException:** Encapsulates errors from the StatsReporting interface to indicate an implementation specific error.
- **StatsAggregationException:** [New exception added since version 2.0a to represent the error in Statistic data aggregation. Because the data is now custom instead of a simple Long we can not predict the behavior any more. Custom implementation of the Statistic can throw this exception if the data can not be aggregatable \(since the type does not match, or the data is released\).](#)

- **NoSuchRendererException:** New exception added since version 2.0a to indicate a Statistic can not render within the reporting class since the data does not render itself and there is no renderer configured at the respective metric.
- **StatsRenderingException:** New exception added since version 2.0a when the custom renderers reports the Statistic can not be rendered by it. This could either be a type mismatch or some other error.

1.5.2 System Exceptions

- **NullPointerException:** Used wherever null argument is used while not acceptable.
- **IllegalArgumentException:** Used wherever empty String argument is used while not acceptable. Normally an empty String is checked with trimmed result.
- **IOException:** Propagated from file operations.

1.6 Thread Safety

This component is not thread-safe. The current FileMatcher and FileProcessor implementations are thread-safe, which the rest part of the components is not thread-safe. Different thread should instantiate separate FileStatistics and StatsReporting to fulfill concurrent tasks (but the FileMatcher and FileProcessor implementations can potentially be shared). This will not cause problem with the command line interface because everything is run in a single thread in a stand-alone JVM.

2. Environment Requirements

2.1 Environment

JDK 1.4 (in order to use regular expression)

2.2 TopCoder Software Components

- **Configuration Manager 2.1.4:** used to support component configuration.
- **Magic Numbers 1.0:** used to examine file types. Notice this component also has the ability to distinguish text and binary files but unfortunately it is not exposed on API.
- **Command Line Utility 1.0:** used to process command line arguments.

2.3 Third Party Components

- **Apache Xalan 2.6.0:** used to apply XSLT against reporting XML.

3. Installation and Configuration

3.1 Package Name

com.topcoder.file.statistics
com.topcoder.file.statistics.processor
com.topcoder.file.statistics.processor.linecounter
[com.topcoder.file.statistics.processor.common.token](#)
com.topcoder.file.statistics.matcher
com.topcoder.file.statistics.reporting

3.2 Configuration Parameters

Parameter	Description	Values
processors	Multiple sub-properties each specify a FileProcessor instance to be aggregated in the FileStatistics.	Property container required
reporting	Specifying a StatsReporting instance to be used with FileStatistics.	Property container optional, XmlStatsReporting is instantiated if not specified
matchall	Flag to indicate whether file is processed with each matching processor or the first matching processor.	“yes”/“true”/“on” is translated to true, other values translated to false. optional, false if not specified
matcher	Specifying a FileMatcher instance to be used with the associating FileProcessor.	Property container optional, AnyFileMatcher is instantiated if not specified
classname	Used within containers to specify the classname to instantiate.	Fully qualified class name. required
alias	Used in processor container to specify the processor alias.	Non-empty name. required
renderers	Used in AbstractStatsReporting to hold a list of StatisticRenderers.	Property container optional
metric	Used in AbstractStatsReporting to specify the metric name for the associated renderer.	Non-empty name. required
stylesheet	Used in XmlStatsReporting to apply on the raw XML generated.	Valid file path. optional, no XSLT applied if not specified.
extensions	Used in ExtensionMatcher to specify the file extensions to match.	Multi-valued. Empty is allowed. required
patterns	Used in RegexMatcher to specify the regular expressions to match filenames.	Multi-valued. Valid regular expression. required
types	Used in TypeNameMatcher to specify the file type names to match (which are configurable in Magic Numbers).	Multi-valued. Non-empty. required
mimes	Used in MimeMatcher to specify the file mimes to match (which are configurable in Magic Numbers).	Multi-valued. Non-empty. required
alphabet	Used in CommonTokenAnalyzer to define the alphabet for the tokens. Alphabet is the complement set to separators.	Non-empty string, contains no duplicates, carriage return or new line. Either alphabet or separator is required.
separator	Used in CommonTokenAnalyzer to define the separators for the tokens.	Non-empty string, contains no duplicates. Either alphabet or separator is required.
casesensitive	Used in CommonTokenAnalyzer to specify whether the tokens are compared case sensitively.	“yes”/“true”/“on” is translated to true, other values translated to false. optional, false if not specified

3.3 Dependencies Configuration

Magic Numbers require configuration.

4. Usage Notes

4.1 Required steps to test the component

- Extract the component distribution.
- Follow [Dependencies Configuration](#).
- Execute 'ant test' within the directory that the distribution was extracted to.

4.2 Required steps to use the component

Follow configuration instructions.

4.3 Demo

4.3.1 *Configure File Statistics*

```
// create a file statistics instance
FileStatistics statistics = new FileStatistics();

// create file matchers
FileMatcher matcher1 = new AnyFileMatcher();
FileMatcher matcher2 = new ExtensionMatcher("txt");
FileMatcher matcher3 = new RegexMatcher(Pattern.compile("pattern"));
FileMatcher matcher4 = new TextFileMatcher();
FileMatcher matcher5 = new BinaryFileMatcher();
FileMatcher matcher6 = new TypeNameMatcher("Java");
FileMatcher matcher7 = new MimeMatcher("text/plain");

// create file processors
FileProcessor processor1 = new TextLineCounter("Text File", matcher1);
FileProcessor processor2 = new CStyleLineCounter("Source File", matcher6);

// create stats reporting implementation
StatsReporting reporting1 = new BasicStatsReporting();
StatsReporting reporting2 = new XmlStatisReporting("stylesheet/xml_to_html.xml");

// manipulate file processors
statistics.addFileProcessor(processor1);
statistics.addFileProcessor(processor2);
FileProcessor removed = statistics.removeFileProcessor("Source File");
FileProcessor query = statistics.getFileProcessor("Text File");
List processors = statistics.getAllFileProcessors();
statistics.clearAllFileProcessors();

// manipulate reporting
StatsReporting original = statistics.getReporting();
statistics.setReporting(reporting2);

// manipulate match all flag
if (!statistics.isMatchAll()) statistics.setMatchAll(true);
```

4.3.2 *Configure Custom Renderers*

```
// create some custom renderers
StatisticRenderer renderer1 = new CommonTokenPlainRenderer();
StatisticRenderer renderer2 = new CommonTokenXmlRenderer();

// manipulate the renderers
reporting.addRenderer(metric1, renderer1);
reporting.addRenderer(metric2, renderer2);
```

```
StatisticRenderer removed = reporting.removeRenderer(metric);
StatisticRenderer query = reporting.getRenderer(metric2);
Map renderers = reporting.getAllRenderers();
reporting.clearAllRenderers();
```

4.3.3 *Process File*

```
// process single file
statistics.processFile(new File("test_files/a/p.txt"));
statistics.processFile(new File("test_files/b/q.java"));
// process directory
statistics.processDirectory(new File("test_files/a"));
statistics.processDirectory(new File("test_files"), true);
```

4.3.4 *Generate Report*

```
// obtain report
String report = statistics.getStatsReporting().generateReport();
// output report
statistics.printReport();
statistics.printReport(new File("test_files/output.html"));
// reset reporting
statistics.getStatsReporting().reset(new File("test_files"));
```

4.3.5 *Command Line Interface*

```
// specify configuration file
-c conf/FileStatistics.xml a.txt b.java
// specify namespace
-n com.topcoder.file.statistics a.txt b.java
// specify base directory
-b test_files a.txt b.java
// recursively process sub-directory
-r test_files
// specify output file
-o test_files/output.html a.txt b.java
```

5. **Future Enhancements**

Provide useful file matchers and file processors.