

Folding
User guide manual
for version 1.1.0rc

`tools@bsc.es`

October 7, 2015

Contents

Contents	iii
List of Figures	v
List of Tables	1
1 User guide	3
1.1 Quick start guide	3
1.1.1 Decompressing the package	3
1.1.2 Contents of the package	3
1.1.3 Quick run	3
1.1.4 Exploring the results	5
1.1.5 Folding options	8
1.2 Configuration, build and installation	9
1.2.1 Libtools package	9
1.2.2 Folding package	10
2 Generate a trace-file for the Folding	11
2.1 Enabling the sampling mechanism	11
2.2 Collecting the appropriate performance counters	12
2.2.1 Intel Haswell processors	12
2.2.2 Intel SandyBridge processors	13
2.2.3 Intel Nehalem processors	13
2.2.4 IBM Power8 processors (experimental, under revision)	13
2.2.5 IBM Power7 processors	13
2.2.6 IBM Power5 processors	14
2.2.7 ARM v7l / Samsung Exynos 5 processors	14
2.2.8 Other architectures	14
2.3 Capturing the call-stack at sample points	14
3 Tool design	19
3.1 First component: trace-file processing	19
3.2 Second component: applying the folding	20
4 API	27
4.1 Generation of input files for the Folding	27
4.1.1 Usage example	27

List of Figures

1.1	Evolution of graduated instructions for 444.namd.	6
1.2	Evolution of multiple counters for 444.namd.	6
1.3	Instruction mix decomposition for Cluster 1 of Nemo.	7
1.4	Evolution of graduated instructions for 444.namd in Paraver.	8
1.5	Paraver time-line showing the callers for Cluster 1 of Nemo.	8
3.1	Data-flow for the first component of the folding.	19
3.2	Main instances and samples related diagram classes.	21
3.3	Instance selection related diagram classes.	21
3.4	Sample selector related diagram classes.	22
3.5	Performance counter interpolation related diagram classes.	22
3.6	Performance models related diagram classes.	23
3.7	Call-stack processing related diagram classes.	25

List of Tables

1.1	Contents of the folding package.	4
-----	--	---

Chapter 1

User guide

1.1 Quick start guide

The Foldingis a mechanism that provides instantaneous performance metrics, source code references and memory references¹. This mechanism receives a trace-file (currently generated by Extrae- see further details on generating a trace-file for the Foldingin Appendix 2) and generates plots and an additional trace-file depicting the fine evolution of the performance. The Foldinguses information captured through instrumentation and sampling mechanisms and smartly combines them. In this context, the samples are gathered from scattered computing regions into a synthetic region by preserving their relative time within their original region so that the sampled information determines how the performance evolves within the region. Consequently, the folded samples represent the progression in shorter periods of time no matter the monitoring sampling frequency, and also, the longer the runs the more samples get mapped into the synthetic instance. The framework has shown mean differences up to 5% when comparing results obtained sampling frequencies that are two orders of magnitude more frequent.

1.1.1 Decompressing the package

The Foldingpackage is distributed in a .tar.bz2 file that can be uncompressed in the working directory by executing the following command:

```
# tar xvz folding-1.0rc8-x86_64.tar.bz2
```

where folding-1.0rc8-x86_64.tar.bz2 refers to the Foldingpackage as distributed from the BSC web page².

1.1.2 Contents of the package

After decompressing the package, the working directory should be populated with the directories (and corresponding descriptions) as listed in Table 1.1.

1.1.3 Quick run

This section provides examples of two types of execution of the Foldingtool. These examples take benefit of the included sample trace-files from the package. For further information on how to generate trace-files for the Foldingtool, check Appendix 2.

¹This last option is experimental at the moment of writing this document

²<http://www.bsc.es/computer-sciences/performance-tools/downloads>

Table 1.1: Contents of the folding package.

Directory	Contents
bin/	Binary packages
etc/	
extrae-configurations/	Minimal configuration files for Extrae
models/	Configuration files to calculate performance models
basic/	
ibm-power5/	
ibm-power7/	
ibm-power8/	
intel-haswell/	
intel-nehalem/	
intel-sandybridge/	
include/	Header files for the development of 3rd party tools
lib/	Libraries for the folding
share/	Miscellaneous files
cfg/	Configuration files for Paraver
doc/	Documentation
html	
examples/	
folding-writer/	Example on how to generate data for the folding
user-functions/	Sample tracefile with manually instrumented regions
clusters/	Sample tracefile with automatically detected regions

Applied to manually instrumented regions

This first example uses a trace-file from the 444.namd SPEC benchmark that contains manually instrumented information that is located in

```
${FOLDING_HOME}/etc/share/examples/user-functions
```

This trace-file was generated by Extrae and delimiting the main loop using the ExtraeAPI³, more precisely the `Extrae_user_function` which emits events with label `User function` (or event type 60000019). To apply the Foldingprocess to this trace-file, simply execute the following commands:

```
# cd ${FOLDING_HOME}/etc/share/examples/user-functions
# ${FOLDING_HOME}/bin/folding 444.namd.prv "User function"
```

Applied to automatically characterized regions

This example consists of a trace-file for the Nemo application when executed in MareNostrum3. This trace-file contains information regarding automatically characterized regions. This characterization has been done using the Clustering tool⁴. This tool enriches the trace-file by adding events labeled as `Cluster ID` (and event type 90000001) into the trace-file. In this context, these events identify similar computation regions based on the event value. To apply the Foldingprocess to this trace-file, simply execute the following commands:

³Please refer to <http://www.bsc.es/computer-sciences/performance-tools/documentation> for the latest ExtraeUser's Guide.

⁴Please refer to <http://www.bsc.es/computer-sciences/performance-tools/documentation> for the latest documentation with respect to the Clustering tool.

```
# cd ${FOLDING_HOME}/etc/share/examples/user-functions
# ${FOLDING_HOME}/bin/folding \
  nemo.exe.128tasks.chop1.clustered.prv "Cluster ID"
```

This trace-file also contains all the necessary performance counters in order to take benefit of several performance models based on performance counters. Simply add the `-model intel-sandybridge` option to the `Foldingscript` to generate the plots with information of the models instead of providing each performance counter individually. The commands to execute should look like this:

```
# cd ${FOLDING_HOME}/etc/share/examples/user-functions
# ${FOLDING_HOME}/bin/folding -model intel-sandybridge \
  nemo.exe.128tasks.chop1.clustered.prv "Cluster ID"
```

1.1.4 Exploring the results

The `Foldingmechanism` generates two types of output inside a directory named as the trace-file given (without the `.prv` suffix). The first type of results include a set of gnuplot files where each of these represents the evolution of the performance counters within the region. The tool also generates a `Paravertrace`-file with synthetic information derived from the `Foldingmechanism`.

Using gnuplot

With respect to the gnuplot files, the `Foldingmechanism` generates as many files as the combination of analyzed regions (clusters, OpenMP outlined routines, taskified `OmpSs` routines, or manually delimited regions) and the counters gathered during the application execution. The user can easily list the generated gnuplot files calling `ls *.gnuplot` within the directory created. The name of the gnuplot files contain the trace-file prefix, the identification of the region folded, and the performance counter shown. For instance, the example described in Section 1.1.3 generates output files that can be explored by executing the command:

```
# gnuplot -persist \
  444.namd.codeblocks.fused.any.any.any.main.Group_0.PAPI_TOT_INS.\
  gnuplot
```

When executing the aforementioned command, the `gnuplot` command should open a window that resembles that in Figure 1.1⁵. The Figure shows that the application faces six phases that execute at 4,500 MIPS approximately. Most of the code occurs in three code locations (being line 76 the most observed line), and we also observe that phases related to high MIPS are related with the activity in the middle of the code-line plot.

This file refers to the user routine `main` (which was manually instrumented) of the trace-file `444.namd.prv` and provides information of the total graduated instructions (`PAPI_TOT_INS`). The user will notice that there are additional files for the different performance counters and they can explore them individually. The `Folding` also generates an additional plot that combines the metrics of all the counters into a single plot. This plot mainly provides information with respect to the MIPS rate (referenced on the right Y-axis), and ratio of the remaining performance

⁵Warning! If the user has problems to open the gnuplot, they should check whether the gnuplot installation is compatible and supports the default terminal. Otherwise, simply select the appropriate terminal (or leave it blank) in the first four lines in the gnuplot script

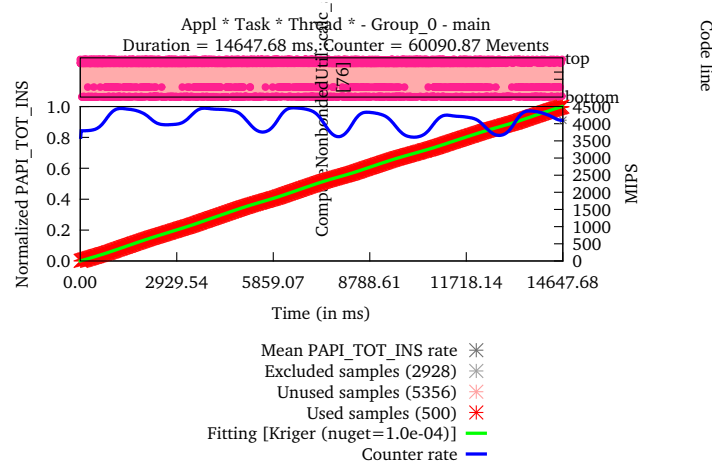


Figure 1.1: Evolution of graduated instructions for 444.namd.

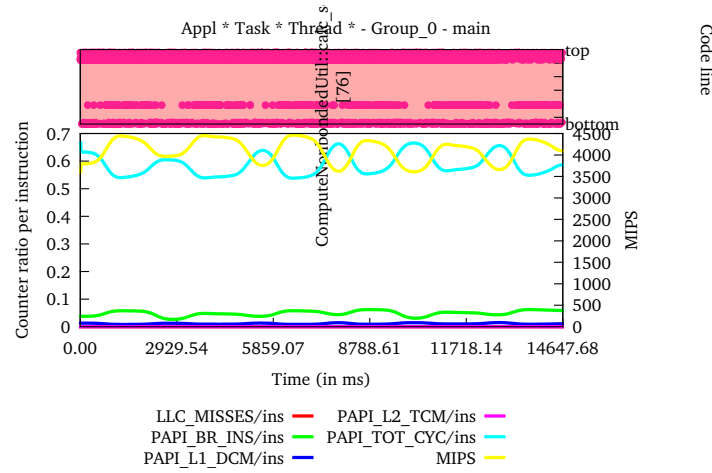


Figure 1.2: Evolution of multiple counters for 444.namd.

counters per instruction (referenced on the left Y-axis). For the particular case of the example from Section 1.1.3, this plot can be explored calling:

```
# gnuplot -persist \
  444.namd.codeblocks.fused.any.any.any.main.\
  Group_0.ratio_per_instruction.gnuplot
```

This command should generate an output combining all the performance counter slopes as shown in Figure 1.2.

The aforementioned instructions also apply to the automatically delimited example described in Section 1.1.3. In this case, the region names are numbered as **Cluster_1** to **Cluster_11**, but they also contain the trace-file prefix and the performance counters to explore them individually. If the user requested the performance models, then additional gnuplot files are created to provide information regarding these models. For the particular case of the Intel SandyBridge model, it generates three models that always generate the MIPS rate and add different metrics:

- *instruction-mix*
Gives insight of the type of instructions executed along the region.

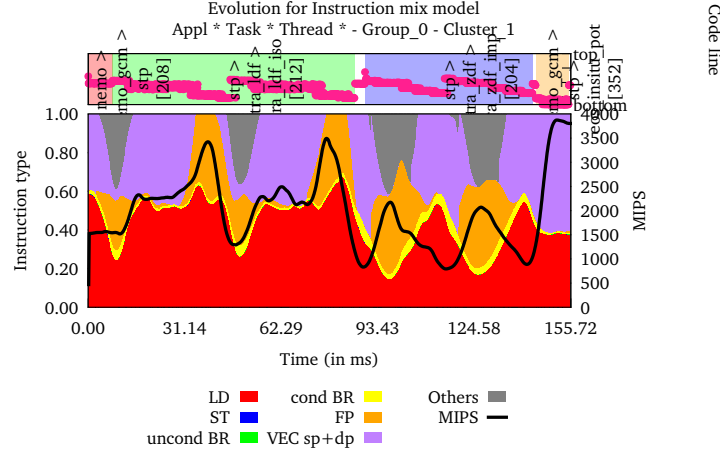


Figure 1.3: Instruction mix decomposition for Cluster 1 of Nemo.

- *architecture-impact*

Provides information regarding to the cache misses at different levels and the branch mispredicts along the region.

- *stall-distribution*

This plot shows information regarding on which components of the processor are stalling the processor pipeline.

For instance, to open the instruction mix for the region labeled as Cluster 1 of the Nemo application executed in Section 1.1.3, the user needs to open the plot invoking the commands below and should obtain a plot similar to Figure 1.3. The reader may see that the application shows two distinctive phases (green and blue) and within each of them there are two repetitions of the same performance.

```
# gnuplot -persist \
nemo.exe.128tasks.chop1.clustered.codeblocks.fused.any.any.any.\
Cluster_1.Group_0.instructionmix.gnuplot
```

The tool also provides a GUI-based tool to explore the plots. the user may invoke a visualizer named `wxfolding-viewer`, by invoking it from the newly created directory such as:

```
# ${FOLDING_HOME}/bin/wxfolding-viewer *.wxfolding
```

Using Paraver

The `Foldingprocess` generates a trace-file with a suffix `.folded.prv` that lets `Paraverto` display some parts of the folded results. The `Foldingpackage` includes several configuration files in the `${FOLDING_HOME}/share/cfg` directory for `Paraverto` help analysing the results. From the configuration files contained in that directory, we outline the following:

- *views/*

- `win_folded_type.cfg`

Generates a time-line that shows in which instances the `Foldingresults` have been

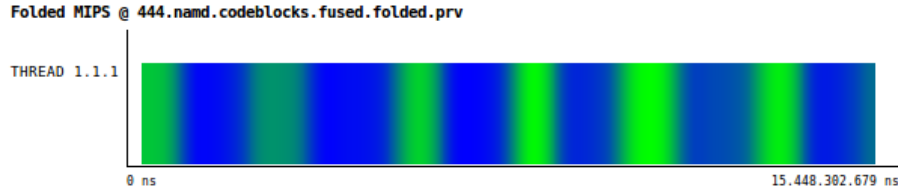


Figure 1.4: Evolution of graduated instructions for 444.namd in Paraver.

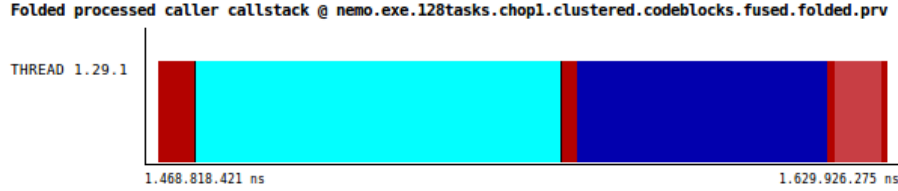


Figure 1.5: Paraver time-line showing the callers for Cluster 1 of Nemo.

integrated. This helps correlating the original trace-file and its contents with the folded trace-file. Notice that only one instance per type (where type refers to function, cluster, *etc...*) is folded.

- `win_folded_mips.cfg`
Generates a time-line showing a signal of the MIPS rate within the folded instances. See in Figure 1.4 a time-line depicting the MIPS rate achieved in the example trace-file for 444.namd.
- `win_folded_processed_call-stack_caller.cfg`
Generates a time-line showing the most time-dominant routines as they have been executed within the folded instances. Figure 1.5 shows a time-line depicting the called routines for the Nemo example.
- `win_folded_processed_call-stack_callerline.cfg`
Generates a time-line showing the most time-dominant source code references (as pair of line and file) as they have been executed within the folded instances.
- *histograms/*
 - `3dh_folded_mips_per_caller.cfg`
Generates a histogram that shows the achieved MIPS rate depending on the caller (columns) for a particular region folded.
 - `3dh_folded_mips_per_callerline.cfg`
Generates a histogram that shows the achieved MIPS rate depending on the source code references (pair of line and file in columns) for a particular region folded.

1.1.5 Folding options

These are the available options for the folding command:

- `-model M` (*optional*)
Requests to use the performance model M when generating plots. Currently, available models include: basic, ibm-power5, ibm-power7, intel-haswell, intel-sandybridge and intel-nehalem.

- **-source S** (*optional*)
Indicates where the source code of the application is located. This is used to detect routines and loops when using the piece-wise linear regression method through the R/strucchange packages.
- **-region R** (*optional*)
Applies the folding only to the instance type named R. Please, notice that region names will not contain special characters (spaces, tabs, ...) because these are swaped by under-scores, so if the folding should apply to **Cluster 2** use **-region Cluster_2**. The user can provide multiple of these flags like: **-region R1 ... -region RN**.
- **trace.prv** (*mandatory*)
The Paraver trace-file to which apply the folding. Should be the penultimate parameter.
- **InstanceSeparator** (*mandatory*)
The label or value of the event type to separate instances within the trace-file. Should be the last parameter. Cannot be given at the same time with **SemanticSeparator**.
- **SemanticSeparator** (*mandatory*)
.csv file generated from Paraver to separate instances within the trace-file. Should be the last parameter. Cannot be given at the same time with **InstanceSeparator**.

1.2 Configuration, build and installation

This section describes how to build and install the Foldingpackage. The Foldingpackage (and its dependencies) requires the Boost library (only the headers suffice), a C compiler, a Fortran compiler and a C++ compiler that supports the C++ 2011 specification (such as g++ version 4.8). This package optionally uses the **strucchange** package from the R statistical application (and may execute in parallel if the **doParallel** is available) to use the piece-wise linear regression interpolation mechanism. Additionally, the Foldingpackage requires the libtools package to be installed first. This package helps on the parsing of Paraver trace-files and can be downloaded from the BSC download web page.

1.2.1 Libtools package

This package is included within the folding package and needs to be installed first. This package requires the boost header files⁶. If the boost header files are located in the system's default, simply run the following command:

```
# ./configure --prefix=/home/harald/aplic/libtools/1.0 \
&& make && make install
```

where **--prefix** indicates the destination folder for this package.

If the boost header files are located elsewhere in the system, run the following command:

```
# ./configure --prefix=/home/harald/aplic/libtools/1.0 \
--with-boost=/path/to/boost \
&& make && make install
```

⁶The libtools package has been successfully tested against version from 1.48 to 1.54.

1.2.2 Folding package

The most basic configuration for the Foldingpackage honors the following commands:

```
# ./configure --with-libtools=$HOME/aplic/libtools/1.0 \  
--prefix=$HOME/aplic/folding/1.0rc8 && \  
make && make install
```

where `--with-libtools` refers to the location of the libtools package installed in Section 1.2.1 and `--prefix` indicates where to install the Foldingtool. If the compilation and installation succeed, the contents of the target installation should look like as the contents defined in Section 1.1.2.

The Foldingtool supports several compilation flags that modify the behavior or enable additional functionalities of the tool. The following list groups the flags according to the behavior they enable.

- `--with-clustering-suite=<DIR>`
The Foldingtool relies on the similarity between the folded instances in order to generate its results. By default, the Foldingtool includes two mechanisms to reduce the noise that appear from using instances with significant different behavior. However, this flag allows using the BSC Clustering suite as a third alternative in order to reduce the noise.
- `--with-R=<DIR1>`, `--with-cube=<DIR2>`, `--with-clang=<DIR3>`
Enables the usage of piece-wise linear regressions on top of the strucchange package⁷ from the R statistical application⁸. This functionality requires the clang compiler⁹ and can generate input files for the Cube3 performance analysis package¹⁰.
- `--with-boost=<DIR>`
This flag lets the Foldingto use a given Boost installation package.
- `--enable-gui`
The results of the Foldingtool is a set of gnuplot files that have to explored manually. If this flag is given at the configure step, the Foldingpackage would include a GUI written in Python that helps exploring all the results from the tool.
- `--enable-callstack-analysis`
This flag enables the call-stack analysis of the segments captured during the measurement step. Enabling this option results in gnuplot files that depict the performance progression collocated with the source code progression.
- `--enable-reference-analysis`
This flag enables the memory references analysis of the references captured during the measurement step (currently, through the `perf` system tool). Enabling this option results in gnuplot files that depict the performance progression collocated with the memory address space and the sampled references.

⁷<http://cran.r-project.org/web/packages/strucchange/index.html>

⁸<http://www.r-project.org>

⁹<http://clang.llvm.org>

¹⁰<http://www.scalasca.org/software/cube-3.x/download.html>

Chapter 2

Generate a trace-file for the Folding

This chapter covers the minimum and necessary steps so as to configure **Extrae**¹ in order to use its resulting trace-files for the Folding process. There are three requirements when monitoring an application with **Extrae** in order to take the most benefit from the **Foldingtool**. First, it is necessary to enable the sampling mechanism in addition to the instrumentation mechanism (see Section 2.1). Second, it is convenient to collect the appropriate performance counters for the underlying processor (see Section 2.2). Finally, **Extrae** needs to capture a segment of the call-stack in order to allow the **Foldingtool** to provide information regarding the progression of the executed routines. The forthcoming sections provide information on how to enable these functionalities through the XML tags for the **Extrae** configuration file.

2.1 Enabling the sampling mechanism

Extrae is an instrumentation package that by default collects information from different parallel runtimes, including but not limited to: MPI, OpenMP, pthreads, CUDA and OpenCL (and even combinations of them). **Extrae** can be configured so that it also uses sampling mechanisms to capture performance metrics on a periodic basis. There are currently two alternatives to enable sampling in **Extrae**: using alarm signals and using performance counters. For the sake of simplicity, this document only covers the alarm-based sampling. However, if the reader would like to enable the sampling using the performance counters they must look at section 4.9 in the **Extrae** User's Manual for more details.

Listing 2.1: Enable default time-based sampling in **Extrae**.

```
1 <sampling enabled="yes" type="default" period="50m" variability="10m"/>
```

The XML statements in Listing 2.1 need to be included in the **Extrae** configuration file. These statements indicate **Extrae** that sampling is enabled (**enabled="yes"**). They also tell **Extrae** to capture samples every 50 milliseconds (ms) with a random variability of 10 ms, that means that samples will be randomly collected with a periodicity of 50 ± 10 ms. With respect to **type**, it determines which timer domain is used (see **man 2 setitimer** or **man 3p setitimer** for further information on time domains). Available options are: **real** (which is also the **default** value, **virtual** and **prof** (which use the **SIGALRM**, **SIGVTALRM** and **SIGPROF** respectively). The default timing accumulates real time, but only issues samples at master thread. To let all the threads collect samples, the **type** must be set to either **virtual** or **prof**.

¹Please refer to <http://www.bsc.es/computer-sciences/performance-tools/documentation> for the latest **Extrae** User's Guide.

Additionally, the **Foldingmechanism** is able to combine several performance models and generate summarized results that simplify understanding the behavior of the node-level performance. Since these performance models are heavily-tighted with the performance counters available on each processor architecture and family, the following sections provide **ExtraeXML** configuration files ready to use on several architectures. Since each architecture has different characteristics, the user may need to tune the XML presented there to make sure that all the list performance counters are gathered appropriately.

2.2 Collecting the appropriate performance counters

The **Foldingmechanism** provides, among other type of information, the progression of performance metrics along a delimited region through instrumentation points. These performance metrics include the progression of performance counters of every performance counter by default. To generate these kind of reports, **Extrae** must collect the performance counters during the application execution and this is achieved by defining counter sets into the **<counters>** section of the **Extrae** configuration file (see Section 4.19 of the **ExtraeUser's** guide for more information).

There has been research that has developed some performance models based on performance counters ratios among performance counters in order to ease the analysis of the reports. Each of these performance models aims at providing insight of different aspects of the application and system during the execution. Since the availability of the performance counters changes from processor to processor (even in the same processor family), the following sections describe the performance counters that are meant to be collected in order to calculate these performance models. These sections include the minimal **<counters>** sections to be added in a previously existing **Extrae** configuration file, but the **Foldingpackage** also includes full **Extrae** configuration files in `${FOLDING_HOME}/etc/extrae-configurations`.

2.2.1 Intel Haswell processors

Listing 2.2: Counter definition sets for the **Extrae** configuration file when used on Intel Haswell processors.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_L1_DCM , PAPI_L2_DCM
4   </set>
5   <set enabled="yes" domain="all" changeat-time="500000us">
6     PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_L3_TCM , RESOURCE_STALLS:SB , RESOURCE_STALLS:ROB
7   </set>
8   <set enabled="yes" domain="all" changeat-time="500000us">
9     PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_SR_INS , PAPI_BR_CN , PAPI_BR_UCN
10  </set>
11  <set enabled="yes" domain="all" changeat-time="500000us">
12    PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_BR_MSP , PAPI_LD_INS
13  </set>
14  <set enabled="yes" domain="all" changeat-time="500000us">
15    PAPI_TOT_INS , PAPI_TOT_CYC , RESOURCE_STALLS , RESOURCE_STALLS:RS
16  </set>
17 </cpu>

```

The listing 2.2 indicates **Extrae** to arrange five performance counter sets with performance

counters that are available on Intel Haswell processors. The collection of these performance counters allows the Foldingto apply the models contained in the $\${FOLDING_HOME}}/etc/models/intel-sandybr$ that include: instruction mix, architecture impact and stall distribution. Unfortunately, the PMU of the Intel Haswell processors do not count neither floating point nor vector instructions.

2.2.2 Intel SandyBridge processors

Listing 2.3: Counter definition sets for the Extrae configuration file when used on Intel Sandy-Bridge procesors.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_L1_DCM , PAPI_L2_DCM , PAPI_L3_TCM , PAPI_BR_MSP ,
4     PAPI_BR_UCN , PAPI_BR_CN
5   </set>
6   <set enabled="yes" domain="all" changeat-time="500000us">
7     PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_VEC_SP , PAPI_LD_INS , PAPI_SR_INS ,
8     RESOURCE_STALLS , RESOURCE_STALLS:RS , RESOURCE_STALLS:ROB
9   </set>
10  <set enabled="yes" domain="all" changeat-time="500000us">
11    PAPI_TOT_INS , PAPI_TOT_CYC , PAPI_VEC_DP , PAPI_FP_INS , RESOURCE_STALLS:SB ,
12    RESOURCE_STALLS:LB
13  </set>
14 </cpu>

```

The listing 2.4 indicates Extraeto configure five performance counter sets with performance counters that are available on Intel SandyBridge processors. The collection of these performance counters allows the Foldingto apply the models contained in the $\${FOLDING_HOME}}/etc/models/intel-sandybr$ that include: instruction mix, architecture impact and stall distribution.

2.2.3 Intel Nehalem processors

The listing 2.3 indicates Extraeto prepare three performance counter sets with performance counters that are available on Intel Nehalem processors. The collection of these performance counters allows the Foldingto apply the models contained in the $\${FOLDING_HOME}}/etc/models/intel-nehalem$ that include: instruction mix, architecture impact and stall distribution.

2.2.4 IBM Power8 processors (experimental, under revision)

The listing 2.5 indicates Extraeto arrange six performance counter sets with performance counters that are available on IBM Power8 (and similar) processors. The collection of these performance counters allows the Foldingto calculate the CPIStack model for the IBM Power8 processor which is contained in $\${FOLDING_HOME}}/etc/models/ibm-power8$.

2.2.5 IBM Power7 processors

The listing 2.6 indicates Extraeto prepare six performance counter sets with performance counters that are available on IBM Power7 (and similar) processors. The collection of these performance counters allows the Foldingto calculate the CPIStack model for the IBM Power7 processor which is contained in $\${FOLDING_HOME}}/etc/models/ibm-power7$.

Listing 2.4: Counter definition sets for the Extrae configuration file when used on Intel Nehalem processors.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM,PAPI_L2_DCM,PAPI_L3_TCM
4   </set>
5   <set enabled="yes" domain="all" changeat-time="500000us">
6     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_BR_MSP,PAPI_BR_UCN,PAPI_BR_CN,RESOURCE_STALLS
7   </set>
8   <set enabled="yes" domain="all" changeat-time="500000us">
9     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_VEC_DP,PAPI_VEC_SP,PAPI_FP_INS
10  </set>
11  <set enabled="yes" domain="all" changeat-time="500000us">
12    PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_LD_INS,PAPI_SR_INS
13  </set>
14  <set enabled="yes" domain="all" changeat-time="500000us">
15    PAPI_TOT_INS,PAPI_TOT_CYC,RESOURCE_STALLS:LOAD,RESOURCE_STALLS:STORE,
16    RESOURCE_STALLS:ROB_FULL,RESOURCE_STALLS:RS_FULL
17  </set>
18 </cpu>

```

2.2.6 IBM Power5 processors

The listing 2.7 indicates Extrae to configure six performance counter sets with performance counters that are available on IBM Power5 (and similar) processors. The collection of these performance counters allows the Folding to calculate the CPIStack model for the IBM Power5 processor which is contained in `#{FOLDING_HOME}/etc/models/ibm-power5`.

2.2.7 ARM v7l / Samsung Exynos 5 processors

The listing 2.8 indicates Extrae to setup three performance counter sets with counters available in Samsung Exynos 5 processors (based on ARM v7l). The collection of these counters allows the Folding to generate instruction decomposition and architecture impact similar to Intel processors. The model is contained in `#{FOLDING_HOME}/etc/models/samsung-exynos5-armv7l`.

2.2.8 Other architectures

The previous definitions of counter sets included performance counters that are available on the specific stated machines. Since these performance counters may not be available on all the systems, the package also provides a group of counter sets that may be available on a variety of systems. Listing 2.9 defines three Extrae counter sets that may be available on many systems (caveat here, not all systems may provide them). With the use of these counter sets, the Folding can apply the models contained in the `#{FOLDING_HOME}/etc/models/basic` that include: instruction mix and architecture impact.

2.3 Capturing the call-stack at sample points

By default, the sampling mechanism captures the performance counters indicated in the `counters` section and the Program Counter interrupted at the sample point. The Folding provides the instantaneous progression of the routines that last at least a minimum given duration. To enable

Listing 2.5: Counter definition sets for the Extrae configuration file when used on IBM Power8 processors.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL,PM_CMPLU_STALL_DCACHE_MISS,
4     PM_CMPLU_STALL_THRD,PM_GRP_CMPL
5   </set>
6   <set enabled="yes" domain="all" changeat-time="500000us">
7     PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_BRU,PM_GCT_NOSLOT_CYC,
8     PM_CMPLU_STALL_FXU
9   </set>
10  <set enabled="yes" domain="all" changeat-time="500000us">
11    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_SCALAR,PM_CMPLU_STALL_LSU
12  </set>
13  <set enabled="yes" domain="all" changeat-time="500000us">
14    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_STORE,PM_CMPLU_STALL_DMISS_L3MISS
15  </set>
16  <set enabled="yes" domain="all" changeat-time="500000us">
17    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_VECTOR,PM_CMPLU_STALL_REJECT
18  </set>
19  <set enabled="yes" domain="all" changeat-time="500000us">
20    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_DMISS_L2L3
21  </set>
22 </cpu>

```

this type of analysis, it is necessary to instruct Extrae to capture a portion of the call-stack during its execution. Listing 2.10 shows how to enable the collection of the call-stack at the sample points in the Extrae configuration file. The mandatory lines to capture the call-stack at sample points are lines 1 and 4. Line 1 indicates that this section must be processed and Line 4 tells Extrae to capture levels 1 to 5 from the call-stack (where 1 refers to the level below to the top of the call-stack).

Listing 2.6: Counter definition sets for the Extrae configuration file when used on IBM Power7 processors.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL,PM_CMPLU_STALL_DCACHE_MISS,
4     PM_CMPLU_STALL_THRD,PM_GRP_CMPL
5   </set>
6   <set enabled="yes" domain="all" changeat-time="500000us">
7     PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_DFU,PM_CMPLU_STALL_IFU,
8     PM_GCT_NOSLOT_CYC
9   </set>
10  <set enabled="yes" domain="all" changeat-time="500000us">
11    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_FXU,PM_CMPLU_STALL_SCALAR
12  </set>
13  <set enabled="yes" domain="all" changeat-time="500000us">
14    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_LSU
15  </set>
16  <set enabled="yes" domain="all" changeat-time="500000us">
17    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_STORE
18  </set>
19  <set enabled="yes" domain="all" changeat-time="500000us">
20    PM_RUN_INST_CMPL,PM_RUN_CYC,PM_CMPLU_STALL_VECTOR
21  </set>
22 </cpu>

```

Listing 2.7: Counter definition sets for the Extrae configuration file when used on IBM Power5 processors.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PM_INST_CMPL,PM_CYC,PM_GCT_EMPTY_CYC,PM_LSU_LMQ_SRQ_EMPTY_CYC,
4     PM_HV_CYC,PM_1PLUS_PPC_CMPL,PM_GRP_CMPL,PM_TB_BIT_TRANS
5   </set>
6   <set enabled="yes" domain="all" changeat-time="500000us">
7     PM_INST_CMPL,PM_CYC,PM_FLUSH_BR_MPRED,PM_BR_MPRED_TA,
8     PM_GCT_EMPTY_IC_MISS,PM_GCT_EMPTY_BR_MPRED,PM_L1_WRITE_CYC
9   </set>
10  <set enabled="yes" domain="all" changeat-time="500000us">
11    PM_INST_CMPL,PM_CYC,PM_LSU_FLUSH,PM_FLUSH_LSU_BR_MPRED,PM_CMPLU_STALL_LSU,
12    PM_CMPLU_STALL_ERAT_MISS
13  </set>
14  <set enabled="yes" domain="all" changeat-time="500000us">
15    PM_INST_CMPL,PM_CYC,PM_GCT_EMPTY_SRQ_FULL,PM_FXU_FIN,PM_FPU_FIN,
16    PM_CMPLU_STALL_FXU,PM_FXU_BUSY,PM_CMPLU_STALL_DIV
17  </set>
18  <set enabled="yes" domain="all" changeat-time="500000us">
19    PM_INST_CMPL,PM_CYC,PM_IOPS_CMPL,PM_CMPLU_STALL_FDIV,PM_FPU_FSQRT,
20    PM_CMPLU_STALL_FPU,PM_FPU_FDIV,PM_FPU_FMA
21  </set>
22  <set enabled="yes" domain="all" changeat-time="500000us">
23    PM_INST_CMPL,PM_CYC,PM_CMPLU_STALL_OTHER,PM_CMPLU_STALL_DCACHE_MISS,
24    PM_LSU_DERAT_MISS,PM_CMPLU_STALL_REJECT,PM_LD_MISS_L1,PM_LD_REF_L1
25  </set>
26 </cpu>

```

Listing 2.8: Basic counter definition sets for ARMv7l / Samsung Exynos5.

```

1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_FP_INS,PAPI_LD_INS,PAPI_SR_INS,PAPI_VEC_INS
4   </set>
5   <set enabled="yes" domain="all" changeat-time="500000us">
6     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_BR_INS,PAPI_BR_MSP,PAPI_L2_DCM
7   </set>
8   <set enabled="yes" domain="all" changeat-time="500000us">
9     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM
10  </set>
11 </cpu>

```

Listing 2.9: Basic counter definition sets for other processors not stated before.

```
1 <cpu enabled="yes" starting-set-distribution="cyclic">
2   <set enabled="yes" domain="all" changeat-time="500000us">
3     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_L1_DCM,PAPI_L2_DCM,PAPI_L3_TCM
4   </set>
5   <set enabled="yes" domain="all" changeat-time="500000us">
6     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_BR_CN,PAPI_BR_UCN,PAPI_LD_INS,PAPI_SR_INS
7   </set>
8   <set enabled="yes" domain="all" changeat-time="500000us">
9     PAPI_TOT_INS,PAPI_TOT_CYC,PAPI_VEC_SP,PAPI_VEC_DP,PAPI_FP_INS,PAPI_BR_MSP
10  </set>
11 </cpu>
```

Listing 2.10: Collect call-stack information at sample points.

```
1 <callers enabled="yes">
2   <mpi enabled="yes">1-3</mpi>
3   <pacx enabled="no">1-3</pacx>
4   <sampling enabled="yes">1-5</sampling>
5 </callers>
```

Chapter 3

Tool design

While the end user executes a single command to apply the **Foldingtool**, this command hides two major components that are executed sequentially and all the outputs are generated into a newly created directory with the name of the input trace-file. The first component processes a user-given trace-file that contains instrumented and sampled data and generates a textual file that contains sequences of instances and samples. The second component takes these sequences of instances and samples, then applies the contouring algorithm, any performance model, and the call-stack processing, and, finally, it generates the output results. Both components are grouped together within the **folding.sh** appearing to the user that the **Foldings** simply consists of a single tool. The tool package contains additional components that may be capture the interest of the user.

3.1 First component: trace-file processing

The first component is divided into three phases that are executed one after another with the user-given trace-file and each of these parse the given trace-file and generates another trace-file that will be used in the subsequent phase as depicted in Figure 3.1. Each of these phases are built in a similar fashion. They parse the input trace-file and keep in memory information regarding the thread state, and eventually, add information to the output. The phases are:

1. **codeblocks** (found in **src/codeblocks**)
This phase attributes to each sample information regarding to the loop / code region that it belongs to according to the application source code.
2. **fuse** (found in **src/fuse**)
This phase compacts the trace-file and ensures that the resulting trace-file is well formed.

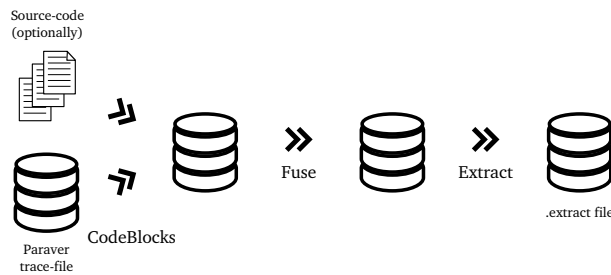


Figure 3.1: Data-flow for the first component of the folding.

3. `extract` (found in `src/extract`)

This is the final phase and extracts information regarding the instances and samples within the trace-file.

The output of this component is a set of files containing information relative to the application. The most notable output is the `.extract` file, which contains the sequence of instances and their samples. For instance, Listing 3.1 shows the contents of the `.extract` file generated using the provided example to demonstrate the API facility. This listing contains information regarding one instance of the `FunctionA` region. The instance starts at timestamp 1,000 ns and lasts 4,500 ns, and it executes up to 2,500 instructions (`PAPI_TOT_INS`) and takes 5,000 cycles (`PAPI_TOT_CYC`) to complete. This instance has two samples associated that occurred at timestamps 2,000 and 4,000, and each of those provides information regarding the aforementioned performance counters.

Listing 3.1: Output example for the `extract` phase of the Foldingmechanism.

1	I	2	2	2	FunctionA	1000	4500	2	PAPI_TOT_CYC	5000	PAPI_TOT_INS	2500
2	S	2000	1000	2	PAPI_TOT_CYC	2000	PAPI_TOT_INS	1000	0	0		
3	S	4000	3000	2	PAPI_TOT_CYC	4000	PAPI_TOT_INS	2000	2	0	1	2 3 1 3 4 5 0

3.2 Second component: applying the folding

The main objective of this component relies on processing the instances and samples extracted and generate the output results. These results include the temporal evolution of the performance counters, any models requested by the user, the source code references and memory references progression, and the results are written in gnuplot and `Paravertrace` files. This section gives a summarized view of the folding work-flow by depicting the most notable class diagrams found in the application source code.

Figure 3.2 shows a portion of the classes that are most important within this tool. The classes *Instance* and *Sample* refer to the instances and samples as-is, without any further processing and as generated by the `extract` tool, in which each *Instance* contains a set of *Sample*, and every *Instance* belongs to an *InstanceContainer*.

After reading every *Instance*, the folding may apply a clustering algorithm (see Figure 3.3) according to the duration of each instance in order to reduce the difference between folded *Instance*. Currently, there are three alternatives regarding the grouping.

- *InstanceSeparatorNone* groups all instances into a single group.
- *InstanceSeparatorAuto* automatically groups the instances according to their duration. The grouping partitions the time-space interval defined by the shortest and longest instances and looks for group of nearby instances.
- *InstanceSeparatorDBSCAN* groups the instances according to a DBSCAN algorithm applied to the duration of the instances. The DBSCAN algorithm groups together instances that are closely packed together (instances with many nearby neighbors) in terms of time and marks as outliers those instances that lie alone in low-density regions. This grouping uses the ClusteringSuite implementation from the BSC performance tools¹.

¹See <http://www.bsc.es/computer-sciences/performance-tools/downloads>.

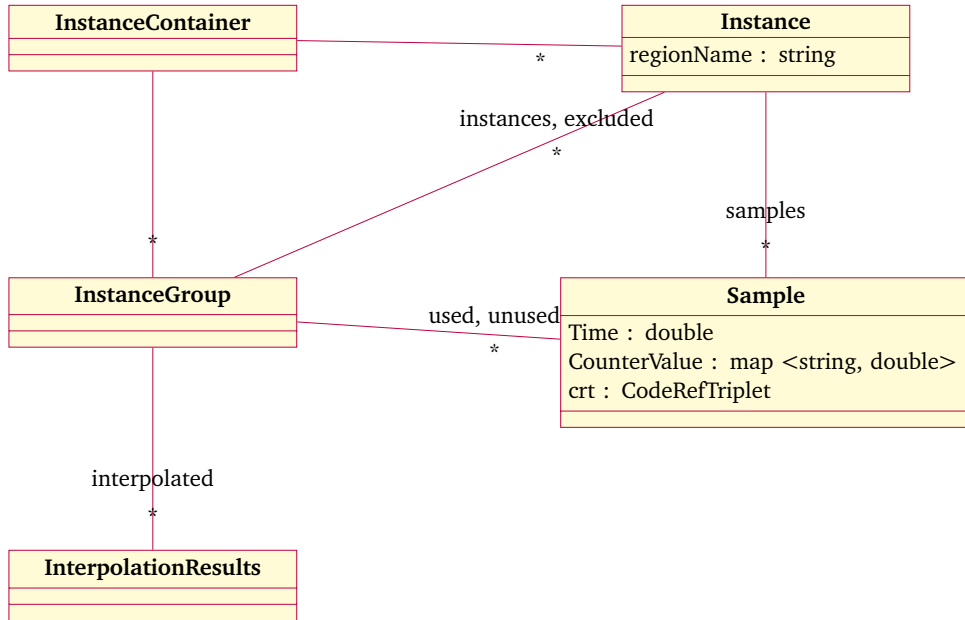


Figure 3.2: Main instances and samples related diagram classes.

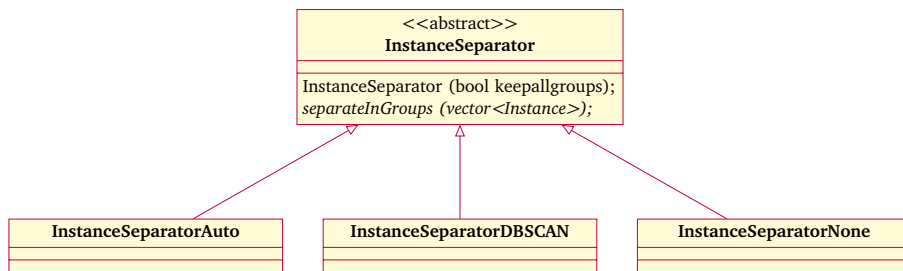


Figure 3.3: Instance selection related diagram classes.

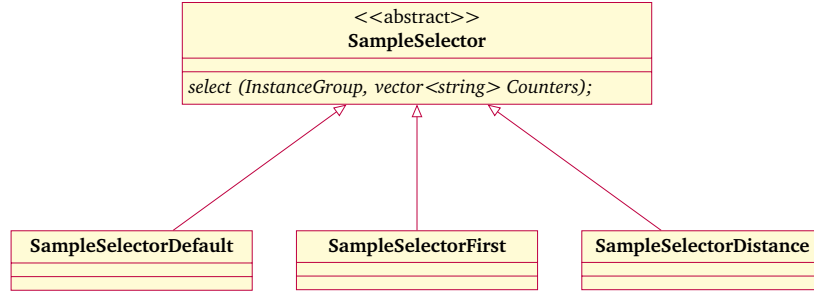


Figure 3.4: Sample selector related diagram classes.

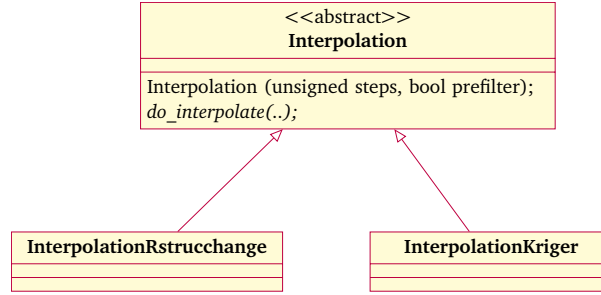


Figure 3.5: Performance counter interpolation related diagram classes.

This grouping begets the *InstanceGroup* objects which contains references to those *Instance* that belong to that particular group. Then, the folding removes the outliers to each *Instance* within every *InstanceGroup* and store the outliers and the remaining in the *excluded* and *instances* associations, respectively.

Since the complexity of the contouring algorithms depends on the number of points to connect, and therefore the number of samples to fold, the **Foldingtool** supports limiting the number of samples given to these algorithms. Figure 3.4 depicts the class diagram of the available *SampleSelector* mechanisms to limit the number of samples.

- *SampleSelectorDefault* the *select* method returns all of the samples within the *InstanceGroup*. This is useful when the user does not impose any limit to the number of samples to be folded.
- *SampleSelectorFirst* receives a threshold (N) in the class constructor. Then, the *select* method tags the first N samples for the processing while the rest are marked as *unused*.
- *SampleSelectorDistance* receives a threshold (N) in the class constructor. Then, the *select* method tags N samples that are equidistant within the *Instance* duration, while the rest are *unused*.

Then the **Folding** repeatedly applies the contouring algorithm to the *used* samples among the different *InstanceGroup* objects. The contouring algorithm applies to each performance counter individually, and as of writing this document, there are two approaches that honor the *Interpolation* super-class virtual method (mainly *do_interpolate*):

- *InterpolationKriger* uses the self-provided contouring algorithm based on the Kriging mechanism to implement the *do_interpolate*.

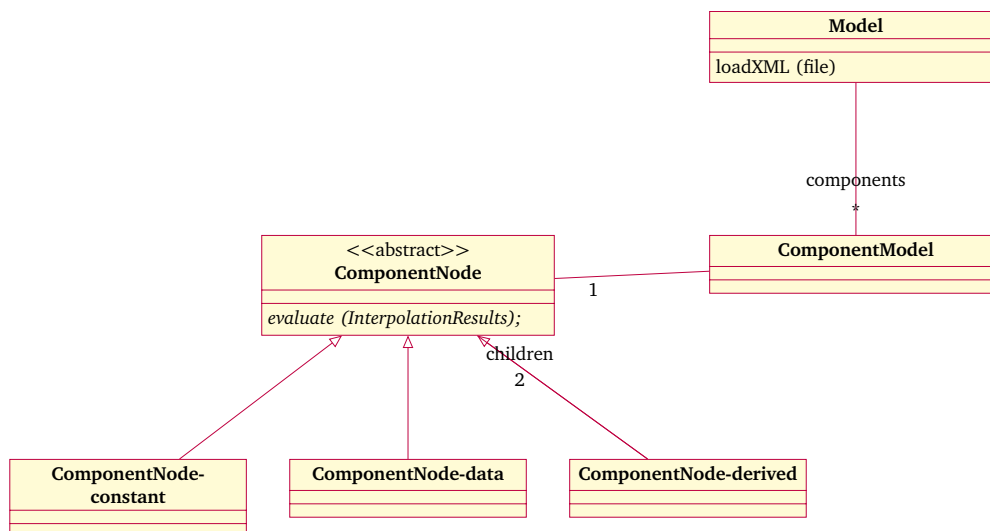


Figure 3.6: Performance models related diagram classes.

- *InterpolationRstrucchange* employs the *strucchange* package² from the R statistical package³ to use piece-wise linear regressions to the folded samples. Additionally, this package may benefit from parallel environments if the *doParallel* package⁴ is available on the system.

The interpolation results are stored, per performance counter, into *InterpolationResults* objects that are associated by *InstanceGroup* by the attribute *interpolated* (as depicted in Figure 3.2). The *interpolated* attribute is implemented as a hash function indexed by the performance counter, so that the interpolation results can be fetched easily.

The *Foldingallows* defining performance models based on performance counters using XML files (see Listing 3.2 for exemplification purposes and `#{FOLDING_HOME}/etc/models` for more detailed examples). Within every XML there may be one or several components (in the last Listing these are: `l1_dcm_ratio`, `l2_dcm_ratio` and `mips`) that will be later represented in the resulting gnuplot using the selected colors and Y-axis (left [y1] or right [y2]). Each component may refer to the instantaneous value of a certain performance counter (as in the `mips` component), a constant value or the operation (addition, subtraction, multiplication and division) between two other values (as in `l1_dcm_ratio` and `l2_dcm_ratio` components). The *Folding* implements the performance models based on performance counters employing the diagram classes show in Figure 3.6. The XML model files are loaded into the *Model* class and each of them may contain multiple components (*ComponentModel*). The *ComponentModel* implements the definition of the component on top of the *ComponentNode* derived sub-classes. These sub-classes allow referencing constant values (*ComponentNode_constant*), interpolated results from a specific performance counter (*ComponentNode_data*) and operation between other two *ComponentNode* objects).

With respect to the analysis of the call-stack, the *Foldingtool* has implemented this analysis through the *CallstackProcessor* related-classes that receives a set of *Sample* objects to explore. Currently, the unique implementation available relies on aligning the call-stacks from the given samples and then exploring the call-stack frames at a given level whether consecutive samples refer to the same routine. If the number of samples surpasses a given threshold, then applies it

²<http://cran.r-project.org/web/packages/strucchange/index.html>

³<http://www.r-project.org>

⁴<http://cran.r-project.org/package=doParallel>

Listing 3.2: Folding example model that generates the L1D and L2D misses per instruction, in addition to the MIPS rate

```

1 <?xml version='1.0'?>
2
3 <model name="sample" title-name="Sample_model"
4     y1="Ratios" y2="MIPS" y1-stacked="no">
5
6     <component name="l1_dcm_ratio" title-name="L1_DCM" where="y1"
7         color="red">
8         <operation type='/'>
9             <value> PAPI_L1_DCM </value>
10            <value> PAPI_TOT_INS </value>
11        </operation>
12    </component>
13
14    <component name="l2_dcm_ratio" title-name="L2_DCM" where="y1"
15        color="blue">
16        <operation type='/'>
17            <value> PAPI_L2_DCM </value>
18            <value> PAPI_TOT_INS </value>
19        </operation>
20    </component>
21
22    <component name="mips" title-name="MIPS" where="y2" color="black">
23        <value> PAPI_TOT_INS </value>
24    </component>
25
26 </model>

```

recursively to the next level until no more levels are available or the number of samples do not surpass the threshold.

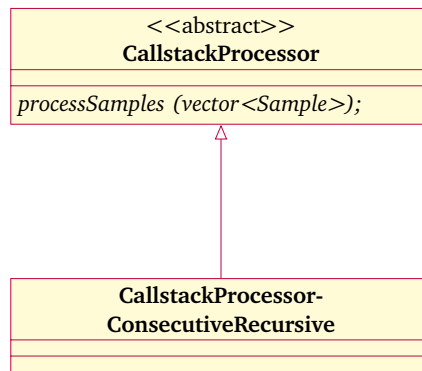


Figure 3.7: Call-stack processing related diagram classes.

Chapter 4

API

This section covers the public API available in the Foldingpackage. This API is meant to allow the Foldingtool to interact with other performance analysis tools in addition to Extrae.

4.1 Generation of input files for the Folding

4.1.1 Usage example

The directory `$FOLDING_HOME/share/examples/folding-writer` contains an example that shows how to generate an input file for the folding from a programatically point of view. The example can be compiled using the following command:

```
# cd $FOLDING_HOME/share/examples/folding-writer
# make
```

The Listing 4.1 shows the example provided in the distributed/installed package. This example demonstrates how to programatically create an `.extract` file for the `interpolate` binary of the Foldingpackage.

Listing 4.1: Example of generating an input file for the Foldingmechanism.

```
1 /*****\
2 *              ANALYSIS PERFORMANCE TOOLS *
3 *              Folding *
4 *      Instrumentation package for parallel applications *
5 *****/
6 *      --- This library is free software; you can redistribute it and/or *
7 *      /  _ _ modify it under the terms of the GNU LGPL as published *
8 *      /  /  _ _ by the Free Software Foundation; either version 2.1 *
9 *      /  /  /  _ _ of the License, or (at your option) any later version. *
10 * ( ( ( B S C ) *
11 * \ \ \ _ _ _ _ / This library is distributed in hope that it will be *
12 * \ \ \ _ _ useful but WITHOUT ANY WARRANTY; without even the *
13 * \ _ _ implied warranty of MERCHANTABILITY or FITNESS FOR A *
14 *      PARTICULAR PURPOSE. See the GNU LGPL for more details. *
15 * *
16 * You should have received a copy of the GNU Lesser General Public License *
17 * along with this library; if not, write to the Free Software Foundation, *
18 * Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301 USA *
19 * The GNU LEsser General Public License is contained in the file COPYING. *
20 *      ----- *
21 *      Barcelona Supercomputing Center - Centro Nacional de Supercomputacion *
22 \*****/
```

23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74

```
#include "folding-writer.H"
#include <fstream>

using namespace std;

int main (int argc, char *argv[])
{
    string nameRegion = "FunctionA";
    unsigned long long startRegion = 1000;
    unsigned long long durationRegion = 4500;

    /* NOTE:: Counters are given in deltas from their previous read, not as absolute values */
    map<string, unsigned long long> c1;
    c1["PAPI_TOT_INS"] = 1000;
    c1["PAPI_TOT_CYC"] = 2000;
    map<unsigned, CodeRefTriplet> crt1;
    Sample *s1 = new Sample (2000, 2000-startRegion, c1, crt1);

    map<string, unsigned long long> c2;
    c2["PAPI_TOT_INS"] = 1000;
    c2["PAPI_TOT_CYC"] = 2000;
    map<unsigned, CodeRefTriplet> crt2;
    CodeRefTriplet codeinfo_l0 (1,2,3);
    CodeRefTriplet codeinfo_l1 (3,4,5);
    crt2[0] = codeinfo_l0;
    crt2[1] = codeinfo_l1;
    Sample *s2 = new Sample (4000, 4000-startRegion, c2, crt2);

    /* Last sample typically coincides with end of region -- see durationRegion,
       Folding::Write won't emit in a S entry */
    map<string, unsigned long long> c3;
    c3["PAPI_TOT_INS"] = 500;
    c3["PAPI_TOT_CYC"] = 1000;
    map<unsigned, CodeRefTriplet> crt3;
    Sample *s3 = new Sample (4500, 4500-startRegion, c3, crt3);

    vector<Sample*> vs;
    vs.push_back (s1);
    vs.push_back (s2);
    vs.push_back (s3);

    ofstream f("output.extract");
    if (f.is_open())
    {
        FoldingWriter::Write (f, nameRegion, 1, 1, 1, startRegion,
                               durationRegion, vs);
        f.close();
    }

    return 0;
}
```

The given example considers that the region `FunctionA` has been identified somehow by the underlying monitoring mechanism, starts at 1,000 ns and lasts 4,500 ns (lines 31-33). Within this period of time, three samples have occurred (*s1-s3*, created in lines 40, 50 and 58, respectively). Samples contain performance counter information and source code references. The performance counter information is given in a relative manner, thus each sample contains the difference from the previous sample (or starting point). For instance, sample *s1* captured information from two

performance counters (`PAPI_TOT_INS` and `PAPI_TOT_CYC`) that counted 1,000 and 2,000 events since the start of the region at time-stamp 2,000 ns (lines 36-40). The second sample (*s2*) does not only contain information from performance counters, but also contains a call-stack segment referencing two call-stack frames. The first frame (`codeinfo_10`) refers to the routine coded as 1, which has source code information coded as 2, and AST-block information coded as 3 (line 46). The same applies to second frame (`codeinfo_11`) - (line 47). These frames are mapped into depths 0 and 1 (where 0 refers to the top of the call-stack) in lines 48 and 49, and then the sample is built using the performance counter information and the call-stack information in line 50. Finally, the last sample (*s3*) only accounted 500 and 1,000 events for the `PAPI_TOT_INS` and `PAPI_TOT_CYC` performance counters respectively, but did not capture any source code reference (lines 54-58). This last sample should coincide with the end of the region (`FunctionA`), and may not be necessarily information captured from a sample point, but from an instrumentation point that indicates the end of the region. All these samples are packed together in a STL vector container (lines 60-63), and then the `FoldingWriter::Write` static method dumps all the information from the samples using the given output stream (lines 65-71).