

Flamingo Auto-Tuning

∞ User's Guide ∞

Ben Spencer

ben@mistymountain.co.uk

Last updated: 28th September 2011

Introduction

Flamingo is a general-purpose auto-tuning framework for software performance optimisation.

The software automates the process of finding optimal settings for program parameters, even reducing the number of tests needed compared with a brute-force approach. These parameters typically control how a program performs, for example the sub-problem size that a large problem is broken down into. It has been designed to be very general-purpose and makes few assumptions about how it should be used.

This guide aims to show you how to use the tuner and describes each feature in detail.

If you have any comments or questions about the tuner or this guide then please feel free to get in touch.

Contents

| | |
|-----------------------------|----|
| Background..... | 2 |
| Tutorial..... | 2 |
| The Configuration File..... | 2 |
| The Example Programs..... | 6 |
| The Utilities..... | 7 |
| Variable Independence..... | 9 |
| Figures of Merit..... | 10 |
| Data Movement..... | 11 |
| Return Codes..... | 13 |
| Running the Tuner..... | 14 |
| Dependencies..... | 14 |

Background

The auto-tuner is used to find the optimal settings of program parameters. These parameters usually control some form of optimisation, such as breaking a large problem into smaller pieces. If it is unclear exactly which combination of settings to use in a particular situation, the tuner can find them. For example, splitting a problem into tiny pieces makes the pieces easy to solve, but the overheads of splitting the problem and recombining the solutions may become high. The best performance is likely to come from breaking the problem into medium-sized blocks. The tuner can be used to find the optimal setting for parameters such as sub-problem size.

Usually, we want to optimise the running time of a program, so this is the default. However, it is possible to provide a custom 'figure-of-merit' which is used to rank the different tests by some other measurement.

If some parameters are independent of each other, this independence can be exploited to speed up the tuning process.

This guide will go through the features of the tuner in detail, and can be used as a reference.

Tutorial

If you have never used Flamingo before, the tutorial is the best place to start. It will lead you through the process of preparing and tuning an example program, explaining each step. This guide describes each feature in more detail and will be useful as a reference.

The tutorial can be found in `doc/tutorial.pdf`, included with the tuner.

The Configuration File

The tuner is controlled by a configuration file determining what to optimise and how tests should be run. This file is used to set up and run the optimisation.

The file must contain five sections, each beginning with a `[section_name]` header. Within each section, options are set using the syntax `name = value`. Lines that begin with `#` are treated as comments and ignored.

A template configuration file (`examples/template.conf`) is provided, containing some explanation of each option. The example programs also each come with a sample configuration file used to tune them.

All paths and commands used must be relative to the configuration file.

It is normal to create the configuration file in the directory containing the program you are tuning, so all commands are entered into the configuration file exactly as you would type them to compile and run the program by hand.

The [variables] Section

This section defines the program parameters which are to be optimised, and any independence between them. There is only a single option:

`variables` (Required)

This option gives a list of the program parameters which should be optimised. There are two possible formats:

- A simple comma separated list of variable names, for example:
`variables = F00, BAR, BAZ`
- A list describing the independence between variables, as described in the *Variable Independence* section. For example:
`variables = {{F00}}, {BAR, BAZ}}`

The [values] Section

This section gives the possible values that each variable above can take. They are specified as a comma separated list for each variable.

Each variable from the [variables] section must have an entry here. Any entries here which are not mentioned in the variable list will be ignored.

All values are interpreted as text strings, which are used without conversion, for example as compiler flags or the program's arguments.

As an example, if we used the above setting of variables, the [values] section may look like:

```
[values]
F00 = -01, -02, -03, -01 -funroll-loops
BAR = 8, 16, 32, 64
BAZ = Alice, Bob, Charlie
```

The [testing] Section

This section defines how to compile, if needed, and run the tests. Tests can also be removed once their testing is complete, if needed.

All the options in this section specify commands which will be executed by the tuner. The value of any variable from the [variables] section can be substituted into the command which is run, using the syntax `%VAR_NAME%`.

For example, if there is a variable `FOO` being optimised, then the test command might be: `./myTestProgram -op %FOO%`. If for a particular test the variable `FOO` takes the value 3, then to get a score for that test, the command `./myTestProgram -op 3` is executed and timed.

There is also one special substitution: `%%ID%%`. This provides a unique id for each test, a counter starting from 1, which can be used in the compilation and testing commands. This can be useful if you want tests to exist at the same time, for example. Compiling tests into an executable called `test_%%ID%%` would result in `test_1`, `test_2`, `test_3` and so on being generated. Each different executable will have a different setting of the parameters being tuned. The 'looping' example in `examples/looping` shows how this works.

`compiler` (Optional)

If specified, this command is executed before testing begins. It can be used to compile a test if a parameter is being changed at compile time (for example a compiler flag or a `#define` statement being overridden by the compiler).

`test` (Required)

To perform a test, this command will be executed and timed by the tuning system. The running time of the command is taken as the score for each test. If a custom figure-of-merit is being used (chosen with the `optimal` option below), then the running time is not measured and instead the score is read from the final line of output from this command.

`cleanup` (Optional)

If specified, this command is executed once the test is complete. It can be used to remove any executables which were compiled.

For some examples of how different types of parameters (`#define` constants, compiler flags, run time arguments, environment variables, ...) can be tuned, see the *Data Movement* section.

The [scoring] Section

This section specifies how tests are scored against each other to choose which is the best.

`repeat` (Optional, defaults to 1, min)

Gives the number of times each test should be repeated. When this option is set, the compilation and cleanup (if present) are still only run once, only the actual test is repeated. The variance, standard deviation and coefficient of variation between the different scores will be displayed at the end of the test.

When a test is repeated, the repeated test scores must be combined into a single overall score for the test. This aggregation can be done by one of the following functions: `min`, `max`, `med` or `avg`.

The number of repetitions and the aggregation function are specified as a comma separated pair. If only the number of repetitions is given, `min` is used as the default aggregate. For example, `repeat = 3, min` is equivalent to `repeat = 3`.

`optimal` (Optional, defaults to `min_time`)

Specifies whether to take the minimum or maximum score as being the best. The settings `min_time` and `max_time` use the running time of the test command as the score for a test. The settings `min` and `max` will use a custom figure-of-merit, which is read from the final line of output. This is described fully in the *Figures of Merit* section.

The [output] Section

Finally, a log can be saved, detailing the testing process:

`log` (Optional)

Specifies the name of a `.csv` file which will be generated. The file will contain details of all the tests performed, which parameter values were used for each and what the individual and overall scores were. This can be used for more detailed analysis after tuning.

Warning: If this file already exists, it will be overwritten.

`script` (Optional)

Specifies the name of a file which will be used to log a 'script' of the tuner's work. This includes which tests are being run and any output from the compilation and testing.

When this option is used, only a summary is shown on screen.

Warning: If this file already exists, it will be overwritten.

The Example Programs

There are some example programs to tune in the `examples` directory. These show how a few different parametrised programs can be tuned, the changes required to the programs and makefiles, and the configuration files used.

`examples/hello/`

A trivial test case, which demonstrates how different parts of the system are connected. It compiles a 'hello world' program, written in C.

Two parameters, named *FOO* and *BAR*, are completely ignored and one, *OPTLEVEL*, controls the compiler optimisation level flag.

The running time of the program is used as a figure of merit, aiming to find the minimum (fastest). There are also example settings in the configuration file to optimise the file size of the generated executable. To try these, simply comment out the existing settings and uncomment these alternatives.

`examples/maths/`

A very simple test case to demonstrate the use of a custom figure of merit. The parameters *X*, *Y* and *Z* are simply summed using the `expr` utility, with no compilation being performed.

`examples/looping/`

This is another fairly trivial C program, which runs a loop as determined by the parameters *XLOOP* and *YLOOP*.

There are settings in the configuration file to measure the running time as measured by the tuner, or by using a custom script, `loop_test.sh`, which uses the `time` utility.

These alternative settings are initially commented out in favour of the tuner's built-in timing. They show how a 'wrapper script' can be used to find a custom figure-of-merit and return it to the tuner.

`examples/matrix/`

Blocked matrix multiplication test case. This is the example used in the tutorial to demonstrate the changes required to the program. The original (not tunable) version is given in `original/`, the modified version after completing the tutorial is given in `modified/` and finally, there is a version which checks the blocked version against the naive version in `comparison/`.

`examples/matlab/`

Demonstrates the use of run-time parameters to determine the optimum level of 'strip-mining' vectorisation in a MATLAB program.

`examples/laplace3d/`

Uses compile-time parameters to tune the block size used by a CUDA GPU program.

The Utilities

The tuner comes with several utilities which can be used to analyse or visualise the results of the tuning process. They use the `.csv` log files generated by the tuner with the `log` option.

`utilities/output_gnuplot.py`

This script converts a `.csv` log file into a `gnuplot .plt` file. This `.plt` file can be used with the `gnuplot` plotting program to produce a detailed graph of the testing process. If required, the `.plt` file can be modified by hand to alter how the graph looks.

The following options are available:

`-h` or `--help`

Outputs some usage information and exits.

`-r SCORE` or `--reference=SCORE`

Plots a reference score for comparison with the tuner's results. This can be useful to compare the tuning scores to some fixed benchmark, such as the running time of the program with the default, non-tuned parameter values.

If `mylog.csv` is the tuning log and `myplot.plt` is the plot file to generate (which will be overwritten if it exists):

```
./output_gnuplot.py [-h] [-r SCORE] mylog.csv myplot.plt
```

Some instructions for generating a graph are given at the top of the `.plt` file generated.

For more information about `gnuplot`, see www.gnuplot.info.

`utilities/output_screen.py`

This script reads a `.csv` log file and produces a graph displayed on the screen, which can be saved if needed. The `matplotlib` python library is required, which may not be installed by default. Information about `matplotlib` can be found at matplotlib.sourceforge.net.

The following options are available:

`-h` or `--help`

Outputs some usage information and exits.

`-r SCORE` or `--reference=SCORE`

Plots a reference score for comparison with the tuner's results.

`-s` or `--stddev`

Plots the standard deviation of multiple test repetitions.

If `mylog.csv` is the tuning log:

```
./output_screen.py [-h] [-r SCORE] [-s] mylog.csv
```

utilities/csv_plot.m

This is a MATLAB program which can be used to display a graph of the testing process.

To use, modify the file as needed and use MATLAB to generate a graph.

utilities/exhaustive_results.py

This script can be used to reconstruct a full set of exhaustive test results from the results of the tests made by the tuner. Because the tuner only ignores tests which are made redundant by the variable independence, the results of these ignored tests can be predicted by looking at the tests which *were* run.

If `results.csv` is the tuning log, `tune.conf` is the configuration file used for the tuning (containing the variable independence information), and `table.csv` is a new results file to be generated (and overwritten if it already exists):

```
./exhaustive_results.py results.csv tune.conf table.csv
```

utilities/node_importance.py and parameter_importance.py

These scripts can analyse a results file, either directly from the tuner, or a reconstructed file from `exhaustive_results.py`, to display the importance of the different nodes in the variable tree or parameters, respectively.

The utilities show which nodes/parameters had the most effect on the running time of the tests. The 'importance' values returned are the average variation in score caused by changing that node/parameter.

If `results.csv` is a results log and `tune.conf` is the configuration file used for tuning:

```
./node_importance.py results.csv tune.conf
```

If `results.csv` is a results log (preferably a complete one, either from a tuning run with no independence or one that was constructed by using `exhaustive_results.py`):

```
./parameter_importance results.csv
```


Variable Independence

We say two variables are independent when they can be optimised separately. That is, variables A and B are independent if when we optimise A with B held at some fixed value, then that optimal value of A will still be optimal for *any* setting of B (and vice-versa).

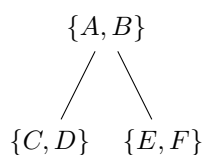
Independent variables are written as separate ‘sets’ of variables. The example we just saw would be written as $\{A, B\}$ if they were dependent and $\{\{A\}, \{B\}\}$ if they were independent.

More complex independences can also be written. As an example, suppose A and B control the operation of the entire program (maybe they are compiler flags), C and D control one aspect of the program and E and F another aspect. The aspect controlled by C and D is not related to the aspect controlled by E and F . This independence would be written like this:

$$\{A, B, \{C, D\}, \{E, F\}\}$$

This shows the two sub-lists $\{C, D\}$ and $\{E, F\}$ are independent of each other and that the variables A and B ‘dominate’ those sub-lists. However, the sub-list $\{C, D\}$ shows that C and D are dependent on each other and must be tuned together. Similarly, A and B are dependent on each other.

This notation essentially describes a ‘tree’ of variables, where higher nodes dominate their subtrees and sibling nodes are independent:



Variable lists can be nested in this way as much as needed to describe the independences in the parameters being optimised.

Optimising Independent Variables

To optimise the tree shown above, the tuner would do the following: For each possible valuation of A and B , test all possible valuations of C and D with some arbitrary fixed setting of E and F . Then fix C and D at some setting and test every possible valuation of E and F . This gives the optimal values for C , D , E and F for that particular setting of A and B . Once these ‘sub-optimums’ have been found for each valuation of A and B , the one which results in the best score is returned.

This method works because sibling subtrees are independent, so the optimal value of C and D which was found at an arbitrary setting of E and F will also be optimal when combined with the optimal setting of E and F (and vice-versa). So the optimisation algorithm avoids many of the redundant tests which a brute-force search would have tried, but can still guarantee to return an optimal valuation as long as the stated independence holds.

Figures of Merit

By default (when using `optimal = min_time`), each test is run and its execution time is measured. This timing is used to rank the tests and choose the best. Sometimes, it is more useful to be able to choose a custom figure of merit. This can be because the part of the program being optimised is not the longest running part of the program, or because you wish to measure some other property, such as memory or network usage.

When the `optimal` option is set to `min` or `max`, the auto-tuner reads the output from the program and interprets the last line of output as the figure of merit. You are free to make any measurements necessary within the test itself (given by the `test` option). The score for the test (a floating point or integer number with no units or other text) must be output as the final line of the commands output.

The 'maths' example program in `examples/maths` is a demonstration of using a custom figure of merit.

Data Movement

Because arbitrary commands can be used for compilation, execution and cleaning, it is possible to use the parameters in any way you could from a command prompt. However, it is not always obvious how the parameter values can be 'passed' through the tool-chain to where they are needed. The following list of tips may be helpful. Although they are fairly linux/make/gcc specific, the general ideas are still applicable to any build tools.

To perform a 'dry-run' with only one test to check these settings, simply set a single possible value for each parameter.

Tuner \implies Shell Command

To pass the value of a parameter to a shell command, simply use the % substitution. If you are tuning a parameter named *FOO* (in the variable list) then `%FOO%` will substitute the current value of *FOO* being tested into the command to execute.

There is also a special substitution, `%%ID%%`, which provides a unique ID for each test. See the section *The Configuration File* for details.

This can be used to set compiler flags, program arguments, and so on. For example, `hello.conf` contains the following lines:

```
examples/hello/hello.conf
49 compile = gcc %OPTLEVEL% -o bin/test_%%ID%% hello.c
...
56 test = ./bin/test_%%ID%% %FOO% %BAR%
...
66 clean = rm ./bin/test_%%ID%%
```

This will cause the following commands to be executed by the tuner:

```
gcc -O0 -o bin/test_1 hello.c
Timed: ./bin/test_1 1 1
Timed: ./bin/test_1 1 1
Timed: ./bin/test_1 1 1
rm ./bin/test_1
gcc -O0 -o bin/test_2 hello.c
Timed: ./bin/test_2 2 1
Timed: ./bin/test_2 2 1
Timed: ./bin/test_2 2 1
rm ./bin/test_2
...
```

The commands you enter are executed directly, as if typed into a command prompt. You are not limited to only running a compiler, or just your test program. It can be useful to use a `Makefile` for compilation or a shell script for measuring and returning a custom figure of merit.

The shell used is `/bin/sh`. This provides all the power and flexibility of being able to use any shell command. However, it also means that if you type `'rm -rf /'` as the command to run a test, that is exactly what will be executed. Be careful.

Tuner ⇒ Makefile

Parameters can be passed to the `make` command simply by appending `NAME=VALE` to the call to `make`, which will allow them to be used within a `Makefile`. However, changing the parameters does not trigger a recompilation of the affected files (because there has been no change to the source code), so this must be forced with the `-B` option to `make`.

For example, `looping.conf` contains the following (Note that the argument `EXEC_NAME` is a new variable for the `Makefile` only, it is not one of the variables being tuned):

```
examples/looping/looping.conf
50 compile = make -f MakeLoop -B EXEC_NAME=./bin/loop_%%ID%% XLOOP=%XLOOP% \
    YLOOP=%YLOOP% OPTLEVEL=%OPTLEVEL%
```

These variables can now be used within the `Makefile` as if they were environment variables, with the syntax `$(NAME)`.

Makefile ⇒ Compiler

Once the parameters have been passed into the `Makefile` as described above, they can be used in the commands to be executed by `make`. For example, the `Makefile` for the ‘looping’ example contains the following, which passes the parameters to the compiler:

```
examples/looping/MakeLoop
4 gcc $(OPTLEVEL) -o $(EXEC_NAME) loop.c -D XLOOP=$(XLOOP) -D YLOOP=$(YLOOP)
```

Compiler ⇒ Program Code

As shown above, the ‘looping’ `Makefile` uses the `-D` option to `gcc` to set the parameter values for `XLOOP` and `YLOOP`. This flag sets them as if they had been set with a `#define` statement in the program itself, so they can be used as constants within the program source code.

Tuner ⇒ Environment Variable

Some programs may be controlled by a certain environment variable, which you want to tune. For example, programs using OpenMP for parallelism can have the number of parallel threads controlled with the environment variable `OMP_NUM_THREADS`. To tune this variable, the `env` command can be used to run the program to be tested in an environment where `OMP_NUM_THREADS` has the value we want:

```
test = env OMP_NUM_THREADS=%OMP_NUM_THREADS% ./program
```

Note that `env` is not a feature of the tuner, simply a standard linux command. This demonstrates the power of being able to run any tuning command you need—the tuner doesn’t need to support changing environment variables directly.

Command Sequencing

The commands given in the `compile`, `test` and `clean` options are passed directly to the shell, so features of the shell, such as using `;` to run one command after another, can be used. For example, another possible solution to the above problem of setting an environment variable would be to export it, then run the command which depends on it:

```
test = export OMP_NUM_THREADS=%OMP_NUM_THREADS%; ./program
```

Return Codes

The tuner will check the return code of all the commands which are executed (`compile`, `test` and `clean`) to see if they failed. Any non-zero return code is considered a failure. This check is used to discount tests which either do not compile or do not run with one particular setting of the parameters.

This is part of the tuning process, as it discounts tests which will not run in your environment which may still run in a different situation. For example, if the number of parallel threads is being tuned, valid possible values on one machine may fail to compile on machines where there are fewer processing cores.

If this causes problems, for example if your program routinely returns a non-zero code, then the `test` command could be run from a 'wrapper' shell script which always succeeds, or the shell's sequencing operator (`;`) could be used to set the return code afterwards, the following will always return with code 0 (success), whatever the behaviour of `./program`:

```
test = ./program; true
```

If you do this, the system will not be able to detect failed tests, so they will still be timed (and will often fail quickly, leading to very low running times), so be careful of erroneous results.

Running the Tuner

The main script which runs the tuner is `autotune`, in the top-level `Autotuning` directory. If the tuner is extracted into your home directory, it can be run from anywhere in the system with the command:

```
~/Autotuning/autotune [configuration file]
```

While the tuner is running tests it can be interrupted with `Ctrl+C`. If the `log` option was used in the configuration file, then a partial log of the tests completed so far will be saved.

Dependencies

The tuner is written in **Python**, so you will need to have the Python interpreter installed. At least **version 2.5** is required, and Python 3000 is not supported.

So far, the tuner has been tested under linux, but is designed to work on any platform (the most compelling reason for choosing Python was this portability and flexibility). On linux, the tuner will be directly executable with `autotune`, but on windows the python file itself will need to be given as an argument to the python interpreter: `python Autotuning\tuner\tune.py`.

The utility `output_screen.py` (which plots graphs of the testing process on screen) requires the `matplotlib` Python module, which may not be installed with a standard python distribution. If it is not installed, you will not be able to use this utility, but the rest of the system will still work as normal. To get `matplotlib`, see matplotlib.sourceforge.net.