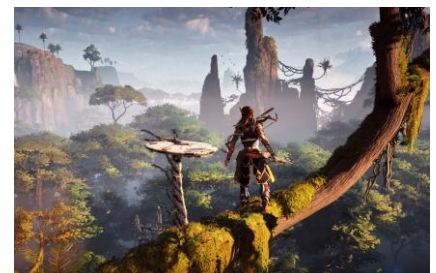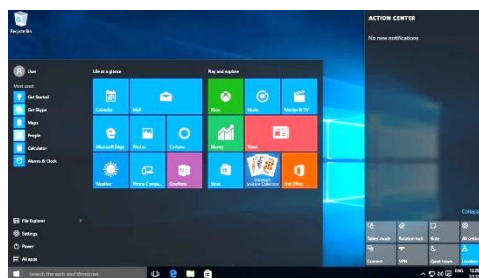# #0 - Profiling

**Author: Jacco Bikker**

## TL;DR

These documents contain an overview of many of the topics discussed in the lectures. Often, additional links to relevant materials are provided. These documents will ultimately provide a full overview of the theory of the course. For the 2019/2020 edition, this is not yet expected to be the case.

In this '#0' a global overview of the optimization process is provided. We also briefly discuss the topic of profiling.

## Raison d'être  /ˌreɪzõ ˈdɛtrə/

There are many reasons to optimize software. The *goal* of the optimization affects the optimization process, and, in particular, the 'termination criterion', i.e. when to stop. Some software simply must run as fast as possible, in which case we start with low hanging fruit (as indicated by a profiling session) and proceed with profiling / optimization cycles until time runs out. In other cases, software should run on a CPU that is as simple as possible, e.g., to save on hardware or energy cost. Perhaps we want to optimize for memory usage, or application response time. Software that is already fast enough also exist. And in some cases, optimization effort is more expensive than new hardware. And finally: sometimes we should refrain from optimization to minimize the risk of breaking something that works.



If faster hardware is not an option, we need to make better use of the resources we have. Often this is a matter of improving algorithms, which is the focus of many research areas. In this course however, we focus on the '$C$' in '$O(N)$'; on how a clever algorithm uses the hardware. Besides the obvious (multithreading and GPGPU) we dive into the not-so-obvious (vectorization, data locality, low level hardware details) to obtain a substantial and often predictable speedup.

## Consistent Approach

Software optimization happens when functionality is complete. This often means that optimization is squeezed in just before the deadline, probably with a tight budget. An important question is: *what can we do in X days?* Answering that question, and delivering what we promise, requires a consistent and reproduceable process, as well as some intuition. E.g.:

- Single-threaded software can run up to five times faster on a quad-core CPU with hyperthreading.
- Code that is clearly suitable for data parallelism can run up to eight times faster on AVX hardware.
- Most software makes poor use of the memory system. Improving this can double application performance.

For many programs, it is safe to estimate that a 25x speedup can be achieved in a modest amount of time. If the GPU can be successfully deployed, this can be even more.
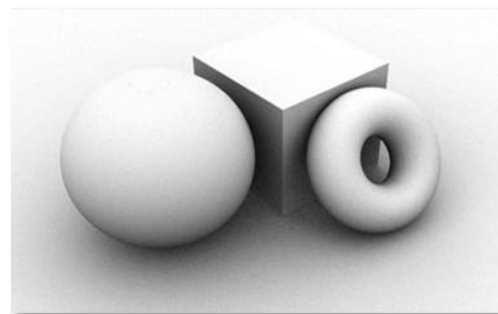
Do note that optimization tends to reduce maintainability of the software. Also note that optimization may break functionality.

I mentioned a consistent and reproduceable process. This is roughly:

0. Determine optimization requirements.
1. **Profile** the program to determine hotspots.
2. Determine scalability of the hotspots.
3. Improve scalability, if possible.
4. **Profile** the improved code.
5. Parallelize / vectorize the code and/or use GPGPU.
6. **Profile** the improved code.
7. Apply low-level optimizations to hotspots.
8. Repeat step 6 and 7 until time runs out.
9. Report.

A repeating term here is 'profiling'. Without profiling, the effort lacks a clear course and becomes essentially random. It is very likely that you spend your efforts on the wrong code lines. You may break things and slow down sections of the program.



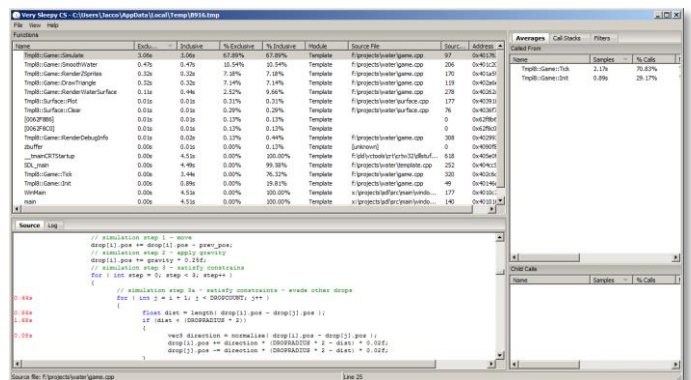*Optimized code can be hard to read.*

# Profiling

Profiling can be as simple as starting a stopwatch when our code starts, and stopping it when the program completes. This can be refined by using a stopwatch per function, or even per line of code.

Luckily, we don't need a timer around each line of code for this. A *profiler* is a program that interrupts running software at fixed intervals. During such an interruption, it probes the *program counter* (PC) of the CPU, which stores the address of the instruction that is about to be executed by the CPU. By storing the PC many times, we get a histogram of where the CPU spends its time.

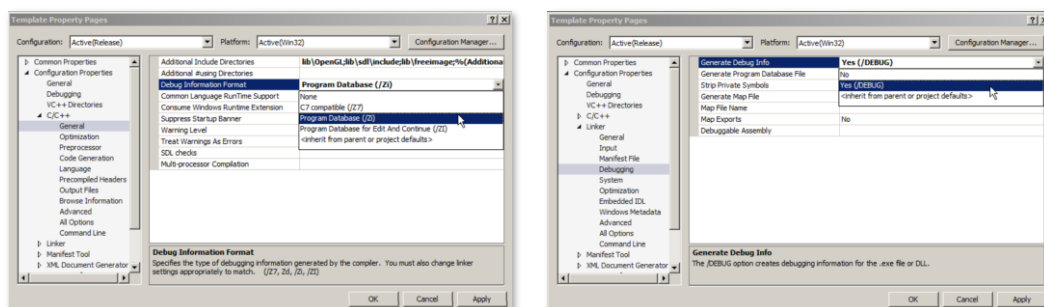A popular profiler for Windows is Very-Sleepy. The profiler of Visual Studio is a good alternative.

Be aware that profiling influences your measurements: interrupting the running program and storing profiling information takes time and changes the contents of caches. Also be aware that modern CPUs and GPUs do not run at fixed clock rates: a hot CPU may slow down. This is especially true on laptops.
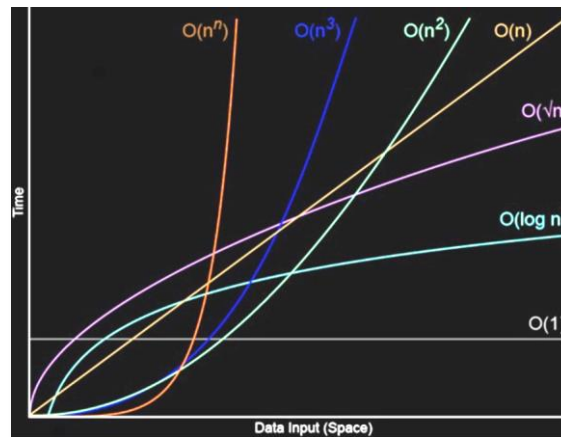


*VerySleepy*

**PRACTICAL NOTES** – In Microsoft Visual Studio you typically chose between a 'Debug' and 'Release' configuration. Note that these are vastly different. You will want to profile the 'Release' version, but this poses a problem: this configuration by default lacks debug information, which is needed for linking instruction addresses to source lines.

Luckily, you can enable this information in a Release build as well. Right-click on the project in the solution explorer and select 'Properties'. Then, for the Release configuration, enable 'Debug Information Format' (C/C++ - General) and tell the linker to generate debug info (Linker – Debugging).

## Scalability

A typical goal of an optimization effort is to make the program fast enough for a larger data set. If this is indeed the case, we may want to measure more than just the performance of the original code. By comparing the performance for different workloads, we get an impression of the *scalability* of the code.



A function in the program that takes little time for 100 input elements may become a bottleneck at 1000 elements. We need to make sure to optimize the software for the *desired* workload, which may not be the current workload.

## Temporally Varying

For real-time applications such as games it is important that the performance of individual frames does not diverge too much from the average performance. To get real-time data on application performance it is sometimes beneficial to add a custom profiler to an application. Many game engines provide this.



*Profiling data in CryEngine.*

## The End

Up next in #1: Low Level Optimization, in which we try to measure the cost of a single line of code.

This material is part of the Optimization & Vectorization course of the MGT master program at the Utrecht University in the Netherlands. More information at:

http://www.cs.uu.nl/docs/vakken/mov.