

Formula Based Performance Design

-



Gerrit Muller

Embedded Systems Institute

Den Dolech 2 (Laplace Building 0.10) P.O. Box 513, 5600 MB Eindhoven The Netherlands

`gerrit.muller@embeddedsystems.nl`

Abstract

Performance models are mostly simple mathematical formulas. The challenge is to model the performance at an appropriate level. In this presentation we introduce several levels of modeling, labeled zeroth order, second order, et cetera. AS illustration we use the performance of MRI reconstruction.

The complete course ASPTM is owned by Embedded Systems Institute. To teach this course a license from Embedded Systems Institute is required. This material is preliminary course material. The final material and course information can be found at: www.esi.nl/cursus.

All Gaudí documents are available at:
<http://www.gaudisite.nl/>

1 Introduction

We recommend to model performance by using simple, secondary school, mathematical formulas. Frequently designers tend to start using more advanced techniques and formalisms, in an attempt to be accurate. In this paper we discuss an approach using simple mathematical formulas, starting with the most simple formulas and refining these formulas as far as needed.

2 Using n-order formulas

The basis for most performance models are simple mathematical formulas, using secondary school math. The challenge is to keep the models as simple as possible, as discussed in the section about control design. We can express the degree of detail in formulas by the order of the formula. Figure 1 shows such classification.

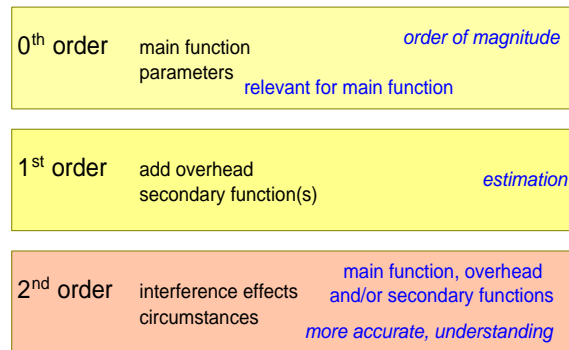


Figure 1: Theory Block 1: n Order Formulas

Figure 2 shows an example of a highly simplified model of the CPU load for image processing. This formula assumes that the CPU load is directly proportional to the number of pixels plus some time to perform the user interface tasks. We call such a formula, where only the main parameter is present, a zeroth order formula.

$$t_{\text{cpu total}} = t_{\text{cpu processing}} + t_{\text{UI}}$$
$$t_{\text{cpu processing}} = n_x * n_y * t_{\text{pixel}}$$

Figure 2: CPU Time Formula Zero Order

It could be that the 0-order formula does not work well enough, for example because overhead is significant. In Figure 3 the biggest overhead contribution is added to the formula, in this example the context switch overhead.

$$t_{\text{cpu total}} = t_{\text{cpu processing}} + t_{\text{UI}} + t_{\text{context switch overhead}}$$

Figure 3: CPU Time Formula First Order

However, in a heavily loaded system may suffer additional loads due to the context switches, the so-called second order effects. In Figure 4 these second order effects are added to the formula. The second order impact may depend on the type of system load. The second order terms might be parameterized to express this relation. For example signal processing loads might cause low penalties, due to high cache efficiency, while control processing might be much more sensitive to these effects.

$$t_{\text{cpu total}} = t_{\text{cpu processing}} + t_{\text{UI}} + t_{\text{context switch overhead}} + t_{\text{stall time due to cache efficiency}} + t_{\text{stall time due to context switching}}$$

signal processing: high efficiency
control processing: low/medium efficiency

Figure 4: CPU Time Formula Second Order

3 Example of n-order formulas in MR reconstruction

The reconstruction of MR images is a processing intensive operation. Fast reconstructions are beneficial for the throughput of MRI scanners and are prerequisite for a number of performance critical applications. Figure 5 shows a simplified block diagram of an MRI scanner, the context of the MR reconstruction. The MR data is digitized in the acquisition subsystem and transferred to the reconstruction subsystem. The reconstructed images are stored in the data base and viewed at the operator or viewing console. All subsystems are controlled by a central host computer.

In Figure 6 a visualization and mathematical formulas are used in combination to model the performance of the MR reconstruction. The visualization shows the processing steps that are performed as reconstruction. Above the arrows it is shown

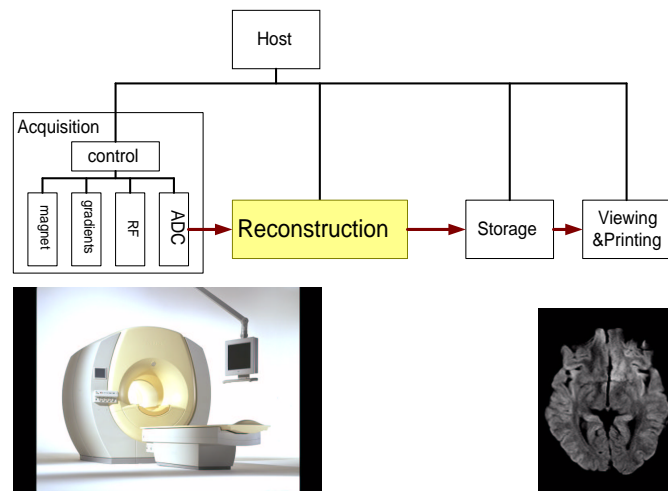


Figure 5: MR Reconstruction Context

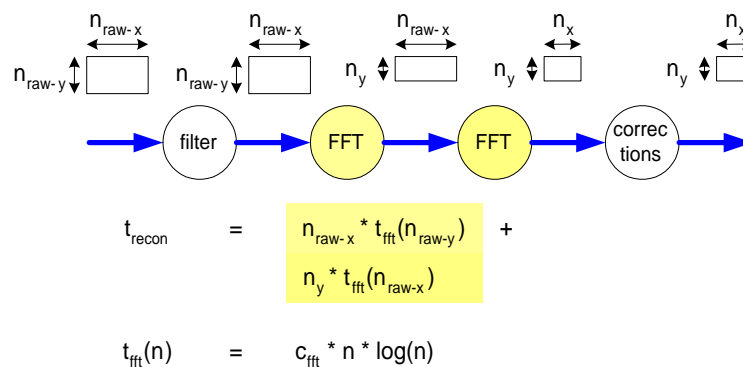


Figure 6: MR Reconstruction Performance Zero Order

what the size of the data matrices is at that phase.

This 0-order model uses the Fast Fourier Transform (FFT) as the dominating term contributing to the performance. Most operations are directly proportional to the matrix size shown above the formulas. The FFT itself is an order $n \log(n)$ term, parameterized with its corresponding load c_{fft} .

Unfortunately the formulas don't tell us much without quantification. Figure 7 provides us with some quantified input based on a FFT micro-benchmark: an FFT on thousand points executes in about 5 msecs (typical performance figures for processing hardware around 1990). The figure takes one typical use case, where a 512*256 raw image is reconstructed on a 256*256 image, to calculate the reconstruction performance. For this use case and assumptions we get 1.2 seconds.

Figure 9 extends the model to also take the non-FFT processing into account.

Typical FFT, 1k points ~ 5 msec
 (scales with $2 * n * \log(n)$)

using:

$$\begin{aligned}
 n_{\text{raw-x}} &= 512 & t_{\text{recon}} &= n_{\text{raw-x}} * t_{\text{fft}}(n_{\text{raw-y}}) + \\
 n_{\text{raw-y}} &= 256 & & n_y * t_{\text{fft}}(n_{\text{raw-x}}) + \\
 n_x &= 256 & & 512 * 1.2 + 256 * 2.4 \\
 n_y &= 256 & & \sim= 1.2 \text{ s}
 \end{aligned}$$

Figure 7: Zero Order Quantitative Example

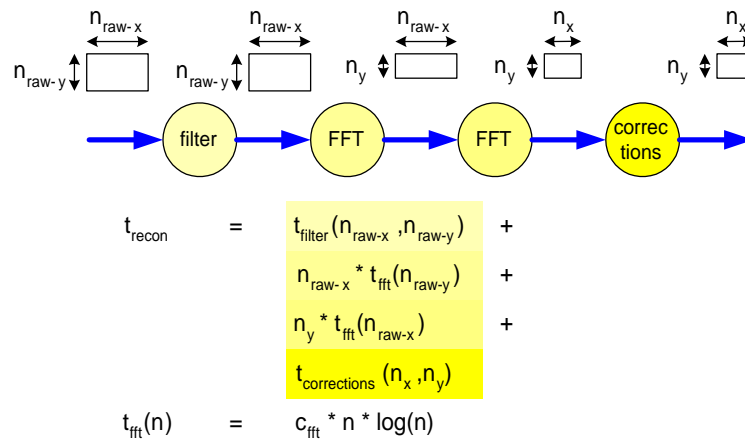


Figure 8: MR Reconstruction Performance First Order

These operations filter the raw data and perform some simple corrections on the image. Both operations are proportional to the number of pixels that is processed.

Figure 9 provides the quantifications obtained by micro-benchmarking both operations: 2 msec to process 1k points. Using the same numbers as Figure 7 we get for filtering $512 * 256 * 2/1024ms \approx 0.26s$ and for correction $256 * 256 * 2/1024ms \approx 0.13s$. Both processing steps can not be ignored compared to the FFT operation!

Finally we add bookkeeping and I/O type operations to the formula, see Figure 10. In practice both terms often ruin the performance of well designed processing kernels, mostly by a lack of attention.

Typical FFT, 1k points ~ 5 msec
 (scales with $2 * n * \log(n)$)

Filter 1k points ~ 2 msec
 (scales linearly with n)

Correction ~ 2 msec
 (scales linearly with n)

Figure 9: First Order Quantitative Example

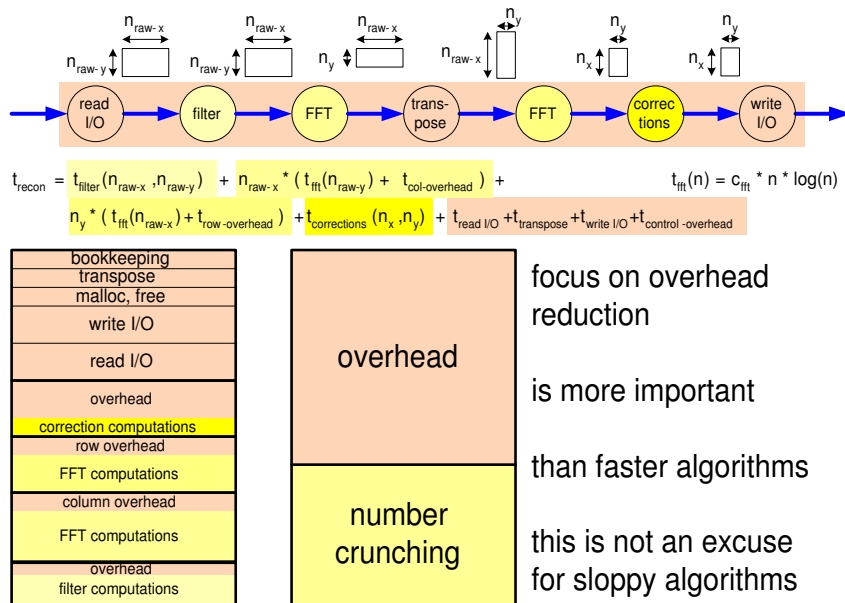


Figure 10: MR Reconstruction Performance Second Order

4 Summary

We have shown that performance can be modeled by starting with the formula for the main function and its parameters. This formula is refined by adding factors that significantly contribute to the (non-)performance.

We used MRI reconstruction as an example of these types of formulas. The formulas are limited to multiplications and logarithms. In this example we have to add quite some factors outside of the main functionality to obtain a usable performance model: simple correction and filter functions, bookkeeping, data-restructuring, and input/output. We also showed that the formulas provide insight, especially the impact of the different parameters, but that actual quantifications also add insight in actual performance numbers and the relevance of the different terms.

5 Acknowledgements

The diagrams are a joined effort of Roland Mathijssen, Teun Hendriks and Gerrit Muller. The approach is based on the EXARCH course created by Ton Kostelijk and Gerrit Muller.

References

- [1] Gerrit Muller. The system architecture homepage. <http://www.gaudisite.nl/index.html>, 1999.

History

Version: 1.0, date: 3 September, 2007 changed by: Gerrit Muller

- Created text by copying from PerformanceEngineering
- changed logo to PHRTreconstructionMRzeroOrder
- changed status to draft

Version: 0, date: 17 February, 2007 changed by: Gerrit Muller

- Created, via refactoring of ASP course, no changelog yet