

Set 4 - Amdahl's Law, Vectorization and ISPC

Issued: October 19, 2018

Hand in (optional): October 26, 2018 23:59

Question 1: Amdahl's Law

Gene Amdahl's Law expresses the effectiveness of improving the performance of a particular part of a system for which the size of the system (associated work) *does not* change. We denote these particular parts of the system as f_1 and f_p which relate to a *serial* and *improvable* fraction, respectively. Note that $f_p = 1 - f_1$. Our goal is to improve the performance of part f_p by a factor of p . If our system without optimization runs for a time t_1 , we expect the improved time for the part f_p to be $f_p t_1 / p$. The execution time of the improved system then is

$$t_p = f_1 t_1 + \frac{f_p t_1}{p} = \left(f_1 + \frac{f_p}{p} \right) t_1.$$

Because the size of the system does not change, we can equate the expected *speedup* for a factor of p as $S_p = t_1 / t_p$, or

$$S_p = \frac{1}{f_1 + \frac{f_p}{p}}. \quad (1)$$

Note that the formulation is very general. We can map this formalism to a computer code by denoting that f_p corresponds to a parallel fraction of the code that can be improved using p parallel elements. For example, if f_p corresponds to 70% for a particular code and we want to improve this part by using $p = 4$ cores, Amdahl's Law tells us that we can expect a net speedup for the whole system of $2.1\times$. Even though we have improved a rather large fraction of the code, we only get a 50% efficiency for our improvement.

- a) Suppose you live in the beautiful city of [Luzern¹](#) in the heart of Switzerland and you are studying Computational Science at ETH Zürich. The distance between the two cities is 55 km and you are required to take a 90 minutes bus ride every day (in any direction). While on the bus today, you read in the news that the Swiss government is planning a new 36 km train route from Rotkreuz (which is one of the bus stops) to Zürich with new trains that can average 70 km/h.
- i) Assume that you can switch from the bus to the train (and vice versa) without waiting time. What is the net speedup you will get from this announcement?

¹http://d2eohwa6gpdg50.cloudfront.net/content/uploads/2018/01/12081536/29_Winter_Chapel_Bridge_Rathaussteg.jpg

- ii) How much faster could you get to ETH if the trains were traveling at the speed of light?
- b) You are working on a science project and you are responsible for the simulation code. Your advisor asks you to work on the code performance to run your simulations $8\times$ faster. From a previous analysis you know that you can potentially improve 90% of the code and that it does not benefit much from multiple hardware threads on the compute cores (hyper threading). Your group currently has access to a single node on the Euler cluster, where each node is a NUMA architecture with *two sockets* and a 12 core CPU on each socket.
- i) How many cores do you need to reach your performance goal?
 - ii) If you need to buy new nodes, you can only buy them in the configuration described above. How many nodes do you need and what is the maximum speedup you can get from your groups investment?
 - iii) Would it make sense to buy more nodes?

Question 2: Manual Vectorization of Reduction Operator

We are interested in optimizing the performance of the following compute kernel

$$r = \sum_{i=1}^n a_i, \quad (2)$$

where a_i are the elements of a vector $\mathbf{a} \in \mathbb{R}^n$ and $r \in \mathbb{R}$ is the result of reducing \mathbf{a} by summing up its elements. We aim at improving the performance of this kernel by exploiting the data level parallelism (DLP) of \mathbf{a} using the SIMD capabilities available on the CPU. We will utilize the streaming SIMD extensions (SSE) to *manually* vectorize the kernel shown in Equation (2). We want to test our implementation for single precision data (32bit) and double precision data (64bit).

- a) Implement vectorized code for the reduction kernel in Equation (2). A baseline version `gold_red` has already been implemented and can be used as a reference. You are asked to work on the items marked with "TODO" in the source file `vectorized_reduction.cpp` located in the directory with the same name inside the skeleton code directory. A guide with a possible workflow can be found in the `Readme.html` file within the source directory. You may use your browser to open the file. Furthermore, the [Intel intrinsics guide²](https://software.intel.com/sites/landingpage/IntrinsicsGuide/) is a useful reference for this task.
- b) In the previous subquestion we were concerned with the vectorization of our compute kernel in Equation (2), that is, the mapping of the execution flow to SIMD lanes for a single core. We can further exploit thread level parallelism (TLP) by mapping the SIMD execution flow to multiple cores using OpenMP. In this task you are going to extend the benchmarking routine in `vectorized_reduction.cpp` with TLP using the OpenMP framework. You can find further hints for this task in the `Readme.html` file and the comments in the code.
- c) You can measure the performance of your vectorized code by running

```
make measurement
```

on Euler. The job will create two pdf files, one for 32bit precision and another for 64bit precision results, each with speedup plots for a small n and another with a large n .
 - i) What is the maximum speedup you expect for 32bit precision and 64bit precision data?
 - ii) Study the speedup plots in the generated pdf files and clearly explain the reason for differences you may observe for a small vector size n and a large vector size n (if there are any).

²<https://software.intel.com/sites/landingpage/IntrinsicsGuide/>

Question 3: Intel SPMD Program Compiler (ISPC)

Manual vectorization can be cumbersome and requires the use of macros to easily switch between different floating point precision. ISPC is a compiler that helps the programmer to avoid such difficulties. As a result, you will be able to write optimized code faster, often with satisfying up to excellent results. This depends on the complexity of the code you want to optimize and yourself. It is very easy to write vectorized code that performs *worse* than the baseline version (true for manual optimizations as well as with ISPC). In this exercise we will utilize the single program multiple data (SPMD) programming model used by ISPC to map program instances to the SIMD lanes on the hardware. We will study ISPC together with the general matrix multiplication kernel (GEMM)

$$C = \alpha AB + \beta C, \quad (3)$$

where $A \in \mathbb{R}^{p \times r}$, $B \in \mathbb{R}^{r \times q}$ and $C \in \mathbb{R}^{p \times q}$ are matrices. For this exercise, we consider the scalars $\alpha = 1$ and $\beta = 0$.

ISPC is used to optimize a performance critical kernel (usually defined in a function, for example) and then compile optimized machine code for that kernel only. In order to use it in our main application code, we need to *link* to the optimized code at compile time. The application code for this exercise is contained in the file `gemm.cpp` inside the `ispc_gemm` directory in the skeleton codes folder. We want to target SSE2 and AVX2 instruction sets, which are both supported by the CPUs on the euler nodes. You are asked to complete the items marked with "TODO" in the skeleton codes. Have a look at the `Readme.html` file inside the source directory for a suggestion of how to solve the exercise and further tips. The ISPC [documentation](#)³ is a helpful resource for this task. Your code should compile for 32bit precision and 64bit precision data (single precision and double precision).

- a) Start with a baseline implementation of the kernel in Equation (3) in the application code `gemm.cpp`. You can compile and test your code with

```
make debug=true gemm_serial
```

You may omit the debug flag if you do not need debugging symbols in your code. Your baseline GEMM implementation should return a norm of truth of 255.966 for double precision data and 256.211 for single precision data.

- b) To get started with the ISPC compiler, you need to install it. You can install it with

```
make install_ispc_linux_x86_64
```

This works on Euler or other 64bit Linux distributions. See also the `Readme.html` file. Windows and MacOS binaries can be downloaded [here](#)⁴. All explanations given in this exercise were tested with the Linux binary of ISPC.

- c) Complete the ISPC related `Makefile` flags and implement the ISPC code for the kernel of Equation (3) in the file `gemm.ispc`. We target optimized kernels for SSE2 and AVX2 instruction sets. You can compile and link to ISPC code with the following (default) target

```
make debug=true gemm
```

³<https://ispc.github.io/ispc.html>

⁴<https://ispc.github.io/downloads.html>

You may omit the debug flag if you do not need debugging symbols in your code. You can submit a job on Euler to test your code using

```
make job
```

For convenience, you may want to work with an interactive node on Euler to omit the latency associated with submitting jobs to the queue.

- d) What are the speedups you expect for your SSE2 and AVX2 optimized kernels? Report two numbers for each optimization, one for single precision data and one for double precision data. If your ISPC code does not reach these expectations, please state the reason for this behavior.
- e) If there is a non-zero error associated with your optimized kernel, explain the reason for this error.
- f) Do you observe differences in errors generated by the SSE2 and AVX2 instruction sets? If so, why is that?