

# CS257 Advanced Computer Architecture

## Coursework Assignment

Term 2, 2024/25

### Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Submission</b>	<b>2</b>
<b>3</b>	<b>Introduction to ACACGS</b>	<b>3</b>
<b>4</b>	<b>Compiling and Running the Code</b>	<b>3</b>
4.1	Visualisation Generation . . . . .	4
<b>5</b>	<b>Hardware Details</b>	<b>6</b>
<b>6</b>	<b>How will my code be tested for performance?</b>	<b>7</b>
<b>7</b>	<b>Rules</b>	<b>7</b>
<b>8</b>	<b>Where do I start?</b>	<b>7</b>
<b>9</b>	<b>Instructions for Submission</b>	<b>7</b>
<b>10</b>	<b>Support</b>	<b>7</b>

# 1 Introduction

The purpose of this coursework is to give you some hands-on experience in code optimisation. By the time you read this, you will have encountered a variety of code optimisation techniques including loop unrolling and vectorisation.

## 2 Submission

Your submission will consist of two parts:

### 1. Optimised Code (50%)

A piece of C code based on the initial implementation provided. This C code will be assessed with respect to your selection and understanding of optimisations, functional correctness, i.e., producing the right answer, and execution speed. Please consider adding comments to your code to increase clarity and explain/justify certain design choices.

### 2. Written Report (50%)

A report (maximum 4 pages, excluding references) detailing your design and implementation decisions. Your report will be evaluated with respect to your understanding of code optimisation techniques and the optimisations you attempted. This means that your report should explain:

- (a) which optimisations you did and did not use;
- (b) why your chosen optimisations improve performance; and
- (c) how your chosen optimisations affect floating-point correctness.

The report should demonstrate your understanding of the subject and your ability to critically evaluate different approaches to the problem as well as characterise them in terms of pros and cons. Please consider using references to support your statements. Your report will also be evaluated in terms of clarity and effectiveness of communication using appropriate language and style.

Given that you may apply many different optimisations, a sensible approach is to build your solution incrementally, saving each partial solution and documenting the impact of each optimisation you make. This means that it is in your interest to attempt as many different optimisations or combinations of optimisations as you can. Consider using graphs or tables to illustrate the effectiveness of each optimisation.

You may discuss optimisation techniques with others but you are not allowed to collaborate on solutions to this assignment. Please remember that the University takes all forms of plagiarism seriously.

**Report Format** Please format the report using a sensible and readable font, e.g. Times New Roman or a similar serif font, at a size of at least 11pt, with at least 0.8 inches margins. The overall formatting should be clean and professional.

**Use of Generative AI tools (GAITs)** Please read carefully the Department policy regarding the use of GAITs available in the Student Handbook.

If you have utilized GAITs for any aspect of this work, including but not limited to brainstorming, summarizing information, or checking grammar this should be clearly acknowledged using a disclaimer like the following:

*I certify that this work is my own original work, except where otherwise explicitly cited. I understand that the use of AI tools is subject to the guidelines outlined in the course handbook. I acknowledge that I am ultimately responsible for the accuracy, originality, and integrity of this work.*

### 3 Introduction to ACACGS

ACACGS is a conjugate gradient proxy application for a 3D mesh. The simulation will execute for either a fixed number of timesteps or alternatively until the residual value falls below a given threshold. This is done for a given mesh size, which is passed in at runtime through command-line arguments.

In this proxy application, a force is applied to each edge boundary of the cuboid, which is then propagated throughout the mesh. As each time step passes, the force is dissipated within the mesh, until the amount of residual is significantly small that the simulation stops (as there are no more calculations to perform), or a set number of time steps have passed.

In addition to providing numeric solutions, the code can also generate visuals which depict the pressure within the mesh throughout the simulation run. Creating the visualisations relies on two optional packages, Silo and VisIt, which are available on the DCS systems <sup>1</sup>. Please note that the visualisations are solely intended to assist students in understanding the problem and the behavior of the simulation and therefore not marked.

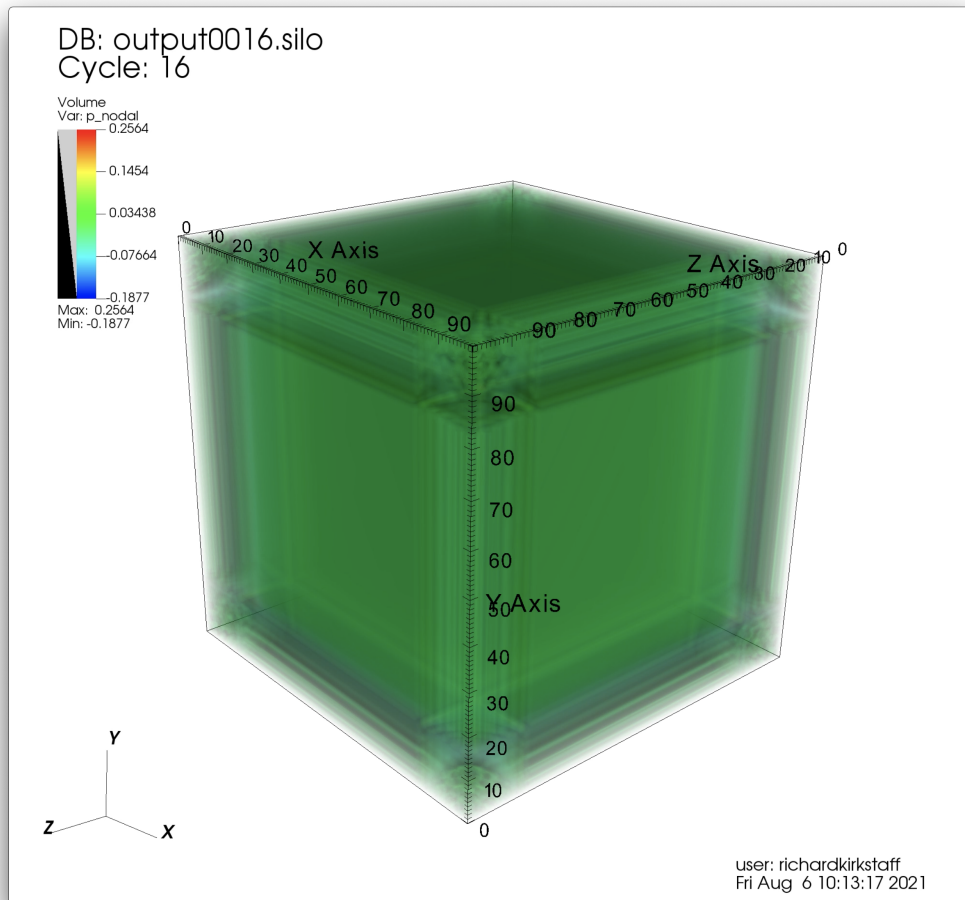


Figure 1: Pressure Matrix Visualisation

### 4 Compiling and Running the Code

You may download the ACACGS Code from the DCS machines under: `/modules/cs257/ACACGS` The code includes a **Makefile** file to build the program. You can compile all of the code using the command `make`. You should not modify the **Makefile** file, but examining it may prove helpful in some situations.

<sup>1</sup>You can use VisIt on the DCS labs machines, but due to compatibility issues, it cannot run on the remote systems - see FAQ on moodle for further details on this

Typing `make` to build the code, will create an executable named `acacgs` in the directory. To clean up the directory, you can run `make clean`.

To run the code, you need to provide the three dimensions for the mesh as three parameters to the executable. For example to execute the provided code on a small 10x10x10 mesh you would enter `./acacgs 10 10 10`. On my system the output for the code is below. This information is also stored in a file, which is named after the wallclock date and time of when the program was first executed (for example, `2025_01_26_12_00_00.txt`).

```
===== Final Statistics =====
Executable name:      ./acacgs
Dimensions:           10 10 10
Number of iterations: 149
Final residual:       2.226719e-92

=== Time ===
Total:                1.126600e-02 seconds
ddot Kernel:          8.390000e-04 seconds
waxpby Kernel:        1.087000e-03 seconds
sparsemv Kernel:      9.123000e-03 seconds

=== FLOP ===
Total:                9.536000e+06 floating point operations
ddot Kernel:          5.960000e+05 floating point operations
waxpby Kernel:        8.940000e+05 floating point operations
sparsemv Kernel:      8.046000e+06 floating point operations

=== MFLOP/s ===
Total:                8.464406e+02 MFLOP/s
ddot Kernel:          7.103695e+02 MFLOP/s
waxpby Kernel:        8.224471e+02 MFLOP/s
sparsemv Kernel:      8.819467e+02 MFLOP/s

Difference between computed and exact = 1.110223e-15
```

You will find more detailed instructions to build the code in the `README.md` file, including flags to turn on verbose mode, which will output details for each timestep in the simulation, and flags for enabling visualisation.

## 4.1 Visualisation Generation

To enable visualisation outputs, you must build your code using `make SILO=1`. This will then compile your code in a way which produces files suitable for visualisation in *VisIt*. If you are working remotely and want to visualise the coursework, it will be quicker and easier for you to copy the files to your local machine, then utilise *VisIt* on the local machine to visualise the cuboid - as often the screen sharing system crashes. Before you make the program, make sure you load the SILO module (`module load cs257-silo`).

When the program is ran with visualisations, each timestep will produce a SILO file within a directory named after the wallclock date and time (for example: `2025_01_26_12_00_00`). In this directory will be a collection of `.silo` files, each named `outputXXXX.silo`, where XXXX represents the timestep it relates to.

Once the program has finished, these can be utilised in *VisIt*. To do so, load the *VisIt* module (`module load cs257-visit`) and open *VisIt* using the command `visit`. From here, you will get 2 windows. The smaller, skinner one is the control window and is used to manage everything that will be displayed. The larger window is the display window. In the control window, select *Open*, and navigate to the directory with the SILO files. You should then be able to select these SILO files.

Now that the SILO files have been loaded, we can now draw some given variables. To do this, click on the *Add* and select a mode and a variable that should be viewed. One of the nicest ones to use is **Volume** and either `x_nodal` or `p_nodal`. When you have finished adding elements, click on **Draw**. This will generate an image in the display window, that can be dragged around so that the cuboid can be viewed from different angles. The control window has a play button, which will run through each timestep.

Table 1: Visualisation Data File Sizes

x	y	z	Cells	Approximate Data Size
10	10	10	1000	4MB
25	25	25	15,625	39MB
50	50	50	125,000	301MB
100	100	100	1,000,000	2.4GB
200	200	200	8,000,000	19.3GB

**Visualisations are nice to have, but for performance purposes we turn them off as they write a significant amount of data to disk.**

There is the potential to go significantly over your DCS disk quota with large meshes. I recommend that you do not exceed 30x30x30 for producing visualisations on the DCS machines. If you are developing your solution on your personal machine then you may wish to produce larger visualisations.

## 5 Hardware Details

On a Linux system, you can read the processor information using the command `cat /proc/cpuinfo` or `lscpu`. This will provide full details on the CPU in the machine, including the CPU model, number of cores, the clock frequency and supported extensions. I strongly recommend taking a look at this on your development machine.

For the purposes of assessment, your code will be run on a DCS machine with 4 cores. The output from `lscpu` can be seen below:

```
Architecture:          x86_64
CPU op-mode(s):        32-bit, 64-bit
Byte Order:            Little Endian
CPU(s):                4
On-line CPU(s) list:   0-3
Thread(s) per core:    1
Core(s) per socket:    4
Socket(s):              1
NUMA node(s):          1
Vendor ID:              GenuineIntel
CPU family:             6
Model:                  158
Model name:             Intel(R) Core(TM) i5-7500 CPU @ 3.40GHz
Stepping:               9
CPU MHz:                3400.000
CPU max MHz:            3800.0000
CPU min MHz:            800.0000
BogoMIPS:               6816.00
Virtualization:         VT-x
L1d cache:              32K
L1i cache:              32K
L2 cache:               256K
L3 cache:               6144K
NUMA node0 CPU(s):     0-3
Flags:                  fpu vme de pse tsc msr pae mce cx8 apic sep mtrr pge mca cmov pat pse36
                        clflush dts acpi mmx fxsr sse sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm
                        constant_tsc art arch_perfmon pebs bts rep_good nopl xtopology nonstop_tsc cpuid
                        aperfmperf tsc_known_freq pni pclmulqdq dtes64 monitor ds_cpl vmx smx est tm2 ssse3
                        sdbg fma cx16 xtpr pdcm pcid sse4_1 sse4_2 x2apic movbe popcnt tsc_deadline_timer aes
                        xsave avx f16c rdrand lahf_lm abm 3dnowprefetch cpuid_fault invpcid_single pti ssbd
                        ibrs ibpb stibp tpr_shadow vnmi flexpriority ept vpid ept_ad fsgsbase tsc_adjust bmi1
                        avx2 smep bmi2 erms invpcid mpx rdseed adx smap clflushopt intel_pt xsaveopt xsavec
                        xgetbv1 xsaves dtherm ida arat pln pts hwp hwp_notify hwp_act_window hwp_epp md_clear
                        flush_lld arch_capabilities
```

Machines matching this specification are available in the `cs257` queue of the Batch Compute System in the Department (referred to as `kudu` in the labs). You will learn how to use this system during the lab sessions, so there will be time to get used to it.

## 6 How will my code be tested for performance?

Your submission will be tested on the DCS batch system on a range of input sizes to evaluate how robust your performance improvements are. It is recommended that you try testing your solution on inputs that are not cubes to see if there are any weaknesses in your optimisation strategies. The 7-pt stencil option will **not** be used for testing your code.

Your code will be executed five times for each problem size on the target hardware. The highest and lowest runtimes will be discarded, and the mean of the three remaining values will be taken as your runtime for that problem size.

## 7 Rules

Your submitted solution **must**:

- Compile on the DCS workstations.

Your submitted solution **must not**:

- Alter the Makefile or add or edit any compiler flags - for example adding the directive `"#pragma GCC optimize("O3")"` is not allowed;
- Use instruction sets not supported by the DCS machines;
- Require additional hardware e.g., GPUs;
- Add relaxed math options to the compile line, e.g., `-ffast-math`. Note: Manual use of approximate math functions is acceptable.

## 8 Where do I start?

This can seem like a daunting project, but we can break it down into a number of steps.

1. Compile and run the code as provided. This is a quick easy check to make sure your environment is setup correctly.
2. Read the code. Start in `main.c` and follow it through. The functions are well documented with Doxygen comments. Don't panic - you are not expected to understand the physics in the code.
3. Measure the runtime of the code for reference purposes.
4. Figure out where the most intensive sections of code are.
5. Develop a small optimisation.
6. Run the code and review the impact of your changes.
7. Repeat steps 5 and 6 until you have exhausted your performance ideas.

## 9 Instructions for Submission

Your solution should be submitted using Tabula. Please ensure that your code works on DCS machines prior to submission.

**Submission Deadline:** Thursday 6<sup>th</sup> March 2025 @ 12 Noon

**Files Required:** A single file named `coursework.zip` which should contain all of your code at the top-level (i.e. no subdirectories) and the report file as a PDF. All files should be submitted through Tabula.

## 10 Support

Support can be found from one of your Teaching Assistants during the labs, or by submitting your question on this form or getting in touch with the module organiser via email.