

Optimisation Cycles for Modular Supercomputing and Energy Efficiency

Alexander Geiß, Technical University of Darmstadt Mathieu Stoffel, Eviden

16.01.2024

The Research leading to these results has received funding from the European Commission's FP7, H2020, and EuroHPC Programmes, under Grant Agreements n° 287530, 610476, 754304, and 955606



A lot of tools in *DEEP-SEA*





A. Geiß, M. Stoffel – Optimisation Cycles for Modular Supercomputing and Energy Efficiency, 16.01.2024

User perspective





3



Optimisation Cycles

... represent the typical workflow of an application developer implementing and optimising an HPC-code



Optimisation Cycles

- Identify the software components that have to play together
 - Optimally support the developer
 - Efficiently utilise the resources of an MSA system
- Define the relevant interfaces between the connected components
 - Describe the information flow





5



Optimising for Modular Supercomputing

Alexander Geiß, Technical University of Darmstadt





III DEEP-SEA

DEEP Software for Exascale Architectures

- Better manage and program compute and memory heterogeneity
- Targets easier programming for Modular Supercomputers
- Continuation of the DEEP projects series

Modular Supercomputing

- Cost-efficient scaling
- Effective resource-sharing
- Composability of heterogeneous resources



Which module(s) should run your application?



- You probably have a good idea
 - Guided by intuition and expectations
 - e.g., the Booster is expected to be faster than the Cluster

- Are these expectations fulfilled?
 - Hard to answer when parameters change
 - e.g., problem size, number of nodes
- The same issue exists on individual CPU-GPU nodes



8

Application Mapping Optimisation Cycle







MSA-related Optimisation Cycle





Optimised execution configuration

Application Mapping OC (simplified)







Cluster Nodes

Instrumentation and Measurement

- Use Score-P
 - Repeat measurement 5 times
 - Measure floating point operation count and DRAM accesses once

Booster Nodes

- Use the Extra-Prof CPU-GPU profiler
 - Repeat measurement 5 times

Measurement grid



Example: Quicksilver mini-app

Solves a simplified dynamic monte-carlo particle transport problem



12

Performance Modelling: Extra-P



More about performance modelling with Extra-P:

- <u>https://youtu.be/Cv2YRCMWqBM</u>
- <u>https://github.com/extra-p/extrap</u>





13



Comparison







Interactive Comparison and Exploration







Interactive Comparison and Exploration





A. Geiß, M. Stoffel – Optimisation Cycles for Modular Supercomputing and Energy Efficiency, 16.01.2024



Interactive Comparison and Exploration





A. Geiß, M. Stoffel – Optimisation Cycles for Modular Supercomputing and Energy Efficiency, 16.01.2024

Determining the Expectation



• We formulate the performance expectation in terms of the roofline model



PEP **Determining the Expectation** Projects We formulate of the roofline model Compare With Other Experiment \times Projection Empirical Roofline Graph (Results.cn.deep. Its.esb.deep.fz-juelich.de.01) Target experiment O CPU GPU 10000 visits Metrics to project 6985.4 GFLOPs/sec B time (FP64 Maxi mum) 2970.0 Ø Ø 1000 sec Roofline model information 179 A830. Memory bandwidth Peak performance DRAM GELOPS Load ERT JSON file... CPU 197,17 GBytes/s 1977,83 GFlops/s \$ 6985,40 GFlops/s Load ERT JSON file... 100 779,22 GBytes/s GPU 50% Use arithmetic intensity (operation intensity) to improve accuracy Floating point operations (double precision) PAPI_DP_OPS Ŧ 10 UNC_M_CAS_COUNT:ALL Number of memory transfers 0.01 0.1 10 100 + FLOPs / Byte Bytes per memory transfer 8 Next > < Back Cancel

19

Comparison Against Roofline Expectation







Use Case: PATMOS





Solves the neutron transport equations to simulate evolution of physical quantities for complex systems

- Cross-sections computation represents 60% to 90% of total runtime
 - Porting cross section computation to GPU
 - Offload batch-size particles at a time



What is a good batch-size?







A. Geiß, M. Stoffel – Optimisation Cycles for Modular Supercomputing and Energy Efficiency, 16.01.2024



Optimizing for Energy Efficiency



Mathieu Stoffel, Eviden



The topic of energy-efficiency in HPC



- The most important metric still is performance
- Hardware tends to be more energy-efficient
- Kernels/applications not using the full computational power of the system have room for improvement of their energy-efficiency
- Resource management is still far from optimal regarding energyefficiency

Opportunities for optimizing the energy-efficiency of HPC applications at runtime



What is Bull Dynamic Power Optimizer (BDPO)?



- A lightweight tool to be run in parallel of an HPC application
- It aims at improving the energy-efficiency associated with the execution of an HPC application
- It is agnostic of the target HPC application
- It features two compatible modes: profiling and optimizing
- lifecycle is fully integrated with Slurm (and, Its standard hence, ParaStation)



How does BDPO work?

- One standalone instance of BDPO per compute node one main control loop
- Monitoring of HW CPU-centric performance metrics through PMC (retired Instructions Per reference Cycle – IPC)
- Detection of phases with low computation intensity
- Enforcement of Dynamic Voltage Frequency Scaling (DVFS) accordingly to the phase detection step





26



How does BDPO work?



user@frontend-hpc \$> srun --bdpo=yes -N 180 --ntasks-per-node 96 ./myapp

4 user@frontend-hpc \$> sbatch --bdpo=yes ./script.sbatch

1	#!/bin/bash
2	
3	#SBATCHjob-name=\$NAME
4	#SBATCHťime=1:10:00
5	#SBATCHexclusive
6	#SBATCHpartition=\$PARTITION
7	#SBATCHnodes=\$NODES
8	#SBATCHntasks=\$NTASKS
9	#SBATCHntasks-per-node=\$NTASKS_PER_NODE
10	#SBATCHcpus-per-task=\$CPUS_PER_TASK
11	#SBATCHbdpo=yes
12	
13	# Export DPO configuration:
14	export BDP0_SAMPLING_MS=100
15	<pre>export BDP0_PFM_INSTRUCTIONS_PER_CYCLE_PROFILING=off</pre>
16	<pre>export BDP0_PFM_INSTRUCTIONS_PER_CYCLE_TRIGGERING=on</pre>
17	<pre>export BDP0_PFM_INSTRUCTIONS_PER_CYCLE_TRIGGERING_THRESHOLD=2.0</pre>
18	<pre>export BDP0_PFM_INSTRUCTIONS_PER_CYCLE_THRESHOLD_CROSSING_UP_ACTION=set_cpufreq:3201000</pre>
19	<pre>export BDP0_PFM_INSTRUCTIONS_PER_CYCLE_THRESHOLD_CROSSING_DOWN_ACTION=set_cpufreq:2400000</pre>
20	export BDP0_COMM_OPTIMIZATION_ENABLED=off
21	<pre>export BDP0_ENTER_MPI_FREQUENCY_ACTION=set_cpufreq:min</pre>
22	<pre>export BDP0_EXIT_MPI_FREQUENCY_ACTION=set_cpufreq:resume</pre>
23	export BDP0_0PTIMIZE_MPI_CLOCK_MODULATION=off
24	<pre>export BDP0_ENTER_MPI_CLOCK_MODULATION_ACTION=set_clock_modulation:87.5</pre>
25	<pre>export BDP0_EXIT_MPI_CLOCK_MODULATION_ACTION=set_clock_modulation:max</pre>
26	

Optimization cycle of BDPO





































Results with ecTrans (part of the IFS) on DEEP



- Each point stems from statistics on 10 executions
- "performance" represents executions with the performance governor and acts as the reference point

Hatfield S. et. al. Mapping a coupled Earth-System simulator onto the Modular Supercomputing Architecture, 20th ECMWF Workshop on High Performance Computing in Meteorology



9000





Hatfield S. et. al. Mapping a coupled Earth-System simulator onto the Modular Supercomputing Architecture, 20th ECMWF Workshop on High Performance Computing in Meteorology

Results with ecTrans (part of the IFS) on DEEP



- BDPO:
 - 16% decrease of the energy consumption
 - 3% increase of the execution time
 - → better energy-efficiency
- No fixed CPU frequency outperforms BDPO





TSMP & BDPO: a few words about co-conception



- TSMP is constituted of three applications: COSMO (modelisation of the atmosphere), ParFlow (modelisation of the ocean), and CLM (coupler between COSMO and ParFlow)
- CLM is a good candidate for BDPO OC. However, the three applications are started in the same Slurm job, on shared nodes, through the multiprog option
- Co-conception: support of the multiprog option in BDPO, and design of a TSMP job where applications do not share nodes (evaluation phase in progress)







A. Geiß, M. Stoffel – Optimisation Cycles for Modular Supercomputing and Energy Efficiency, 16.01.2024



Summary





The Application Mapping OC



Your application

Module 2

Module

Data Analyti

Clust

MSA

Module 4

Module 6 ulti-tier Storage System

> Module 5 Quantum

- Mapping applications onto modules can cause surprises
 - Also, when offloading on heterogeneous nodes
 - Performance might differ from expectations

- The Application Mapping OC helps you to check if expectations are fulfilled
 - Integrated set of tools
 - Systematic approach
 - Flexible, when needed



Two takeaways about BDPO OC

- There are opportunities for optimization of the energy-efficiency of the executions of HPC applications at runtime:
 - → Very good results with ecTrans on 16 nodes:
 - 16% decrease of the energy-to-solution
 - \circ 3% increase of the time-to-solution
- Usability is important:
 - A lot of partners tried BDPO OC on their applications: ecTrans (IFS), TSMP, PATMOS, xPic, GROMACS, ...
 - When asked why: BDPO works out-of-the-box and is easy to use









The DEEP Projects have received funding from the European Commission's FP7, H2020, and EuroHPC Programmes, under Grant Agreements n° 287530, 610476, 754304, and 955606.

The EuroHPC Joint Undertaking (JU) receives support from the European Union's Horizon 2020 research and innovation programme and Germany, France, Spain, Greece, Belgium, Sweden, United Kingdom, Switzerland.