# ORCHIDEE-MICT

Profiling

# Table of Contents

# Overview

This document tries to understand Orchidee MICT computing time behavior. In the latest version 6.5 it takes a lot of time to compute. Around 8h in 0.5 degrees for 1 year. So it is necessary to understand why It happens. Once the problems are identified it is possible to apply different solutions for each issue.

In order to make such thing possible the code is profiled. Different tools are used. They provide an easy way to identify basic hotspots in the code.
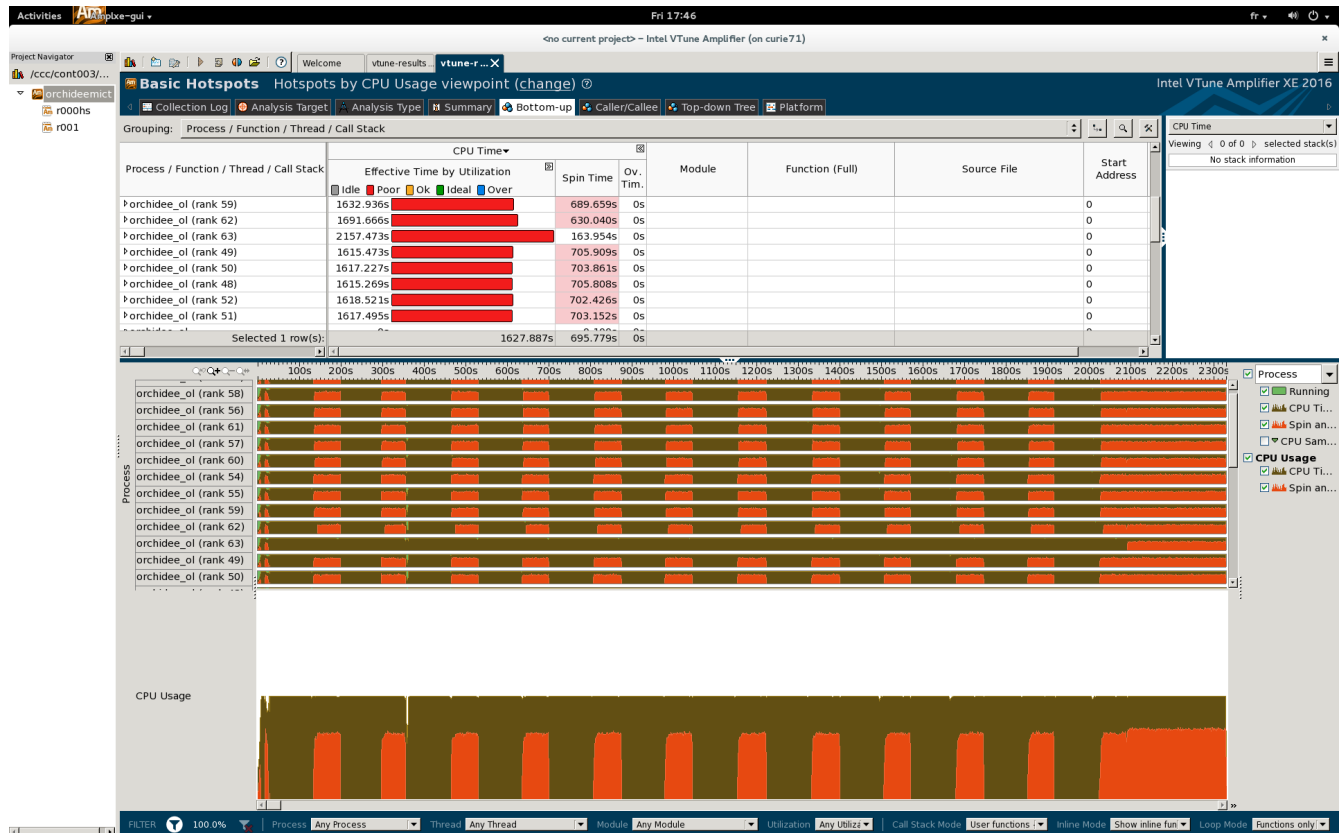
# History file size

## Overview

Orchidee history output files have different file sizes. It might be interesting to understand why. This could lead to major issues or not.

## Simulation description

Orchidee MICT revision 3526 (interpolation, on going development):
- First year
- 64 cores
- 1 year
- 1 degree
- IOIPSL library
- Monthly history
- Total time: ~2350 seconds

## Profiling (Intel vTune)
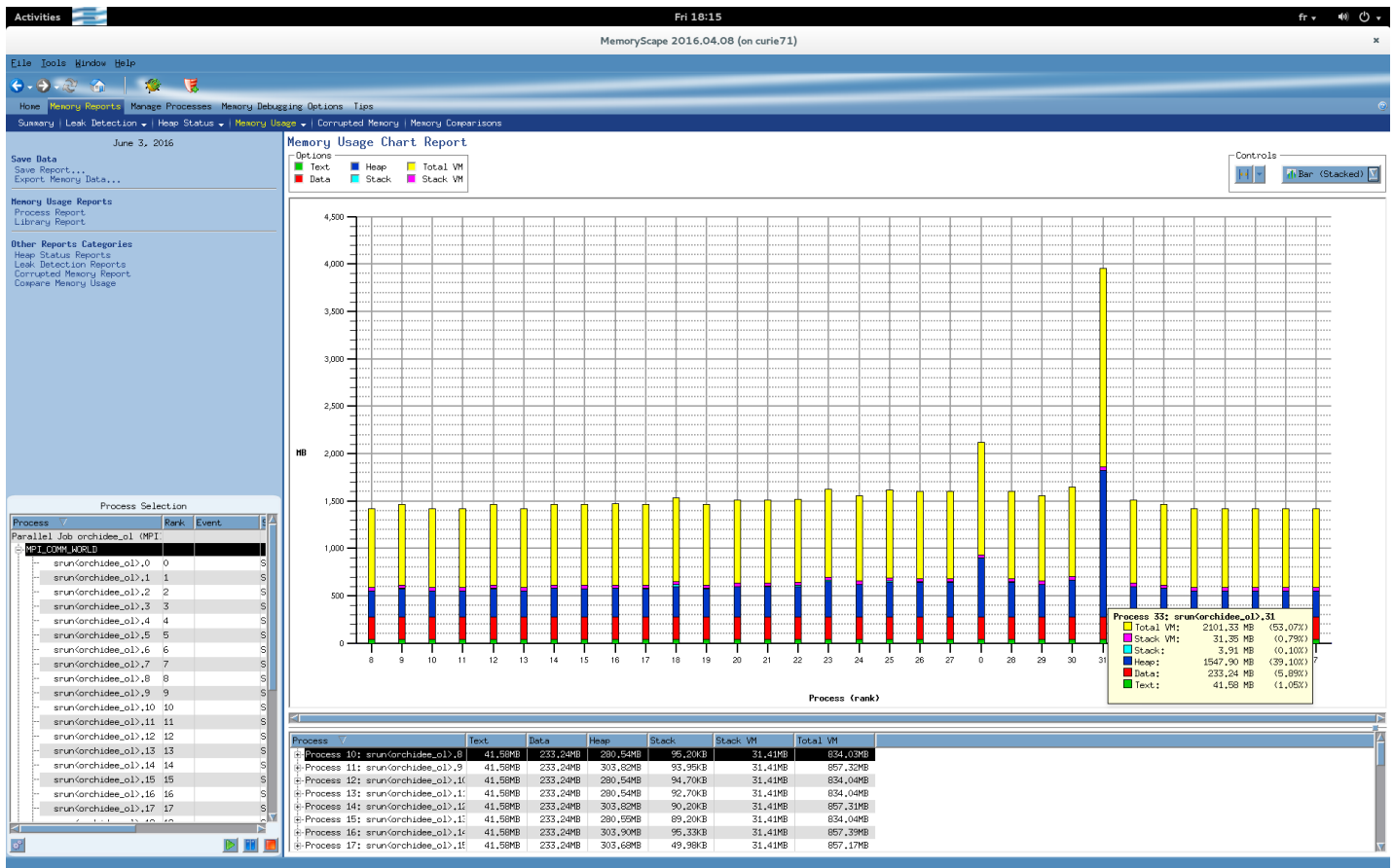
Red color: the processor is waiting. It is doing nothing.
Middle red colors (12 months): at the end of every month history is written.
Latest red part: restart files are written.
Latest processor is always unbalanced. It is slower than the rest. Due to MPI blocking calls they need to wait.

# Memory (Totalview)

Note: this is a test with 32 processors. Randomly stop at the beginning to show off the memory usage:



The latest processors uses more memory:
- The heap (allocate/deallocate) is much higher
-  VM is a consequence of the Heap size

Let's take a deeper look. Process 33 (the latest) uses more memory.



Unfold the latest as well as another processors to compare the memory. Modules mathelp and histcom are using most of the memory. They belong to IOIPSL library. It manages input/output.

Again, a more detailed screenshot:

Histdef is responsible of the memory allocation.

# Output File Size

Stomate output file history vs its size.



The latest processor is several times bigger the others.

# Ncview

Using ncview it shows how big is the area managed by the latest proc.  Let's compare it with the latest (63) and another random processor output.



Latest proc (63) has much more data than the other (e.g 24).

# Conclusion

Orchidee selects for each processor the same number of land pixels. But the outputs have different sizes because they represent different world areas (managed by IOIPSL). A direct consequence of is an increase time writing data to netcdf file. The rest of the processors must wait for the latest due to MPI calls.

# XIOS 1 Profiling

To deactivate IOIPSL history it is necessary to change in run.def:
- XIOS_ORCHIDEE_OK=y
- WRITE_STEP=0
- STOMATE_HIST_DT=0

Simulation:
- 64 cores (+ XIOS1 1core ),
- 1Y
- 05DEG
- Simulation time: 2h40 (9600 seconds)
- yearly output
- /ccc/work/cont003/dsm/p529jorn/experiments/SECHSTOM.DGVM.10336_vtune_xios

# MICT vs CROP

## Overview
Compare MICT vs MICT merged CROP disabled. Time performance and its outputs.

## Simulation Description
1Year
2 Degree

## MICT
/ccc/work/cont003/dsm/p529jorn/experiments/SECHSTOM.DGVM.10336_interpol_improve

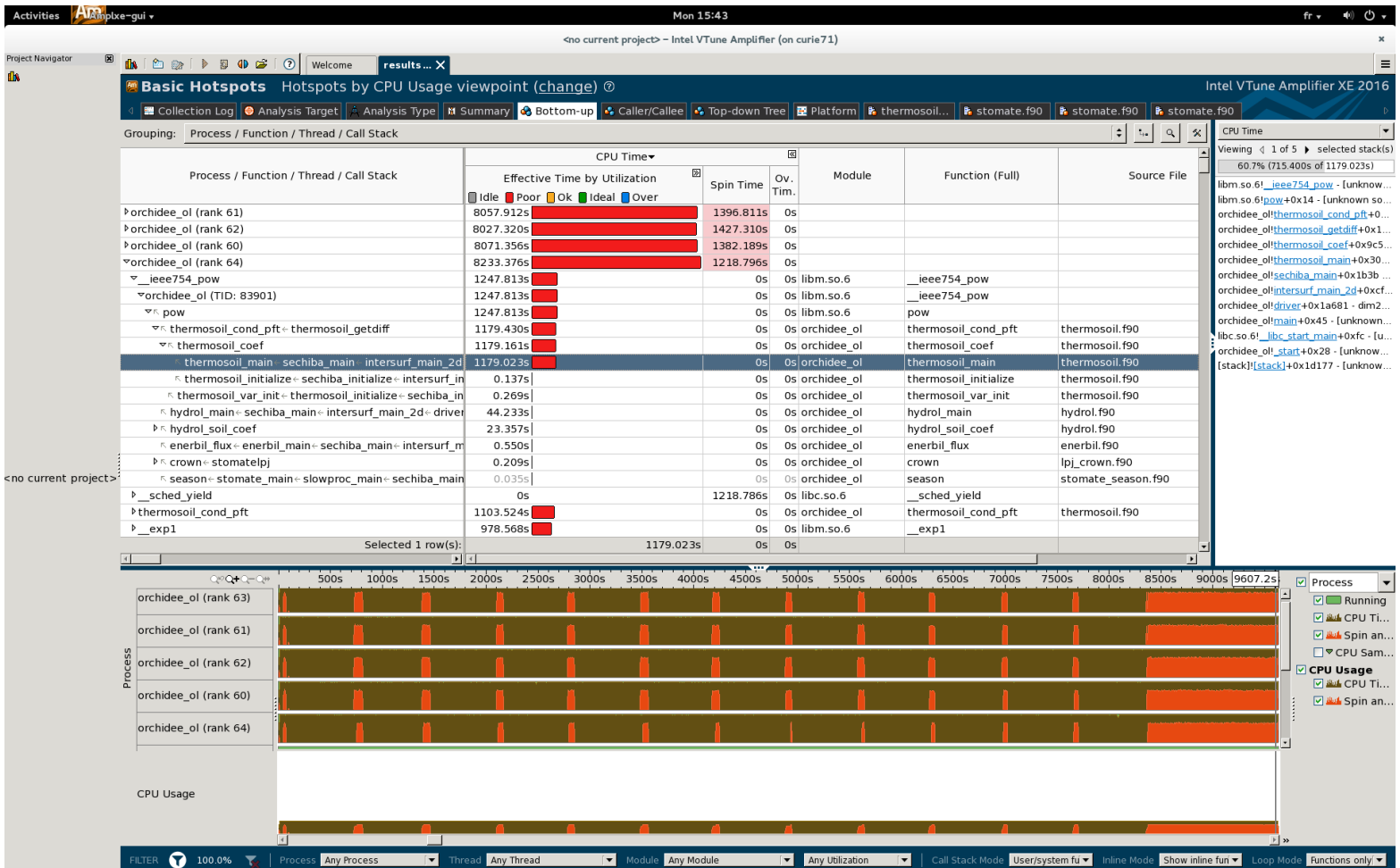| JobID | JobName | Ntasks | Ncpus | Nnodes | Layout | Elapsed | Ratio | CPusage | Eff | State |
|-------|---------|--------|-------|--------|--------|---------|-------|---------|-----|-------|
| 4996984 | M7_test | - | 32 | 2 | - | 00:12:55 | 100 | - | - | - |
| 4996984.0 | orchidee_ol | 32 | 32 | 2 | BBlock | 00:12:54 | 99.8 | 00:12:39 | 98.0 | COMPLETED |

774s

## MICT-CROP
/ccc/work/cont003/dsm/p529jorn/experiments/SECHSTOM.DGVM.10336

| JobID | JobName | Ntasks | Ncpus | Nnodes | Layout | Elapsed | Ratio | CPusage | Eff | State |
|-------|---------|--------|-------|--------|--------|---------|-------|---------|-----|-------|
| 5001560 | M65_test | - | 32 | 2 | - | 00:16:39 | 100 | - | - | - |
| 5001560.0 | orchidee_ol | 32 | 32 | 2 | BBlock | 00:16:38 | 99.8 | 00:16:22 | 98.3 | COMPLETED |

998s

20% slower