

# **An evaluation of Linux Kernel Samepage Merging on simulations**

Sébastien Valat

*CEA, DAM, DIF F-91340 Arpajon, France*

**PhD supervisor** : William Jalby

**CEA,DAM advisor** : Marc Pérache



énergie atomique • énergies alternatives

- **Introduction of KSM**
- **KSM configuration**
- **Testing KSM on HERA**
- **Conclusion**



énergie atomique • énergies alternatives

# KSM : Kernel Samepage Merging



énergie atomique • énergies alternatives

- **Introduced in Linux kernel by version 2.6.32**
- **Goal : reduce memory usage in multi-VM environnement.**
- **Base idea : use the page mapping to merge identical physical pages.**
- **Dynamically detect similar pages with a scanning daemon.**
- **Implemented at kernel level**



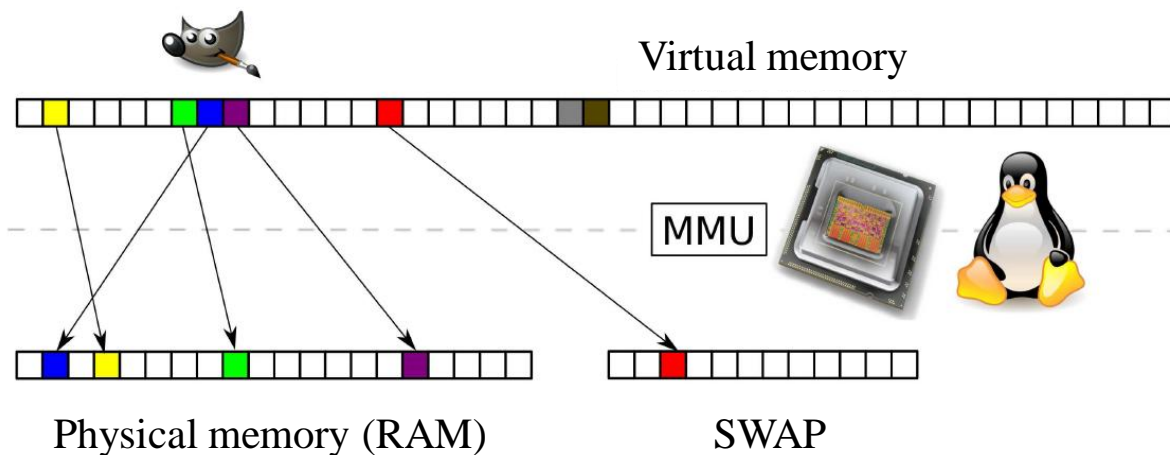
- **Efficient for VM with same of similar operating systems.**
  - Same binary files
  - Same resources
  - Zeroed pages
- **Transparent for the user.**
- **Price : CPU usage for the scanning daemon**
- **Generic implementation : target standard process, not only VM**

# Reminder on memory management

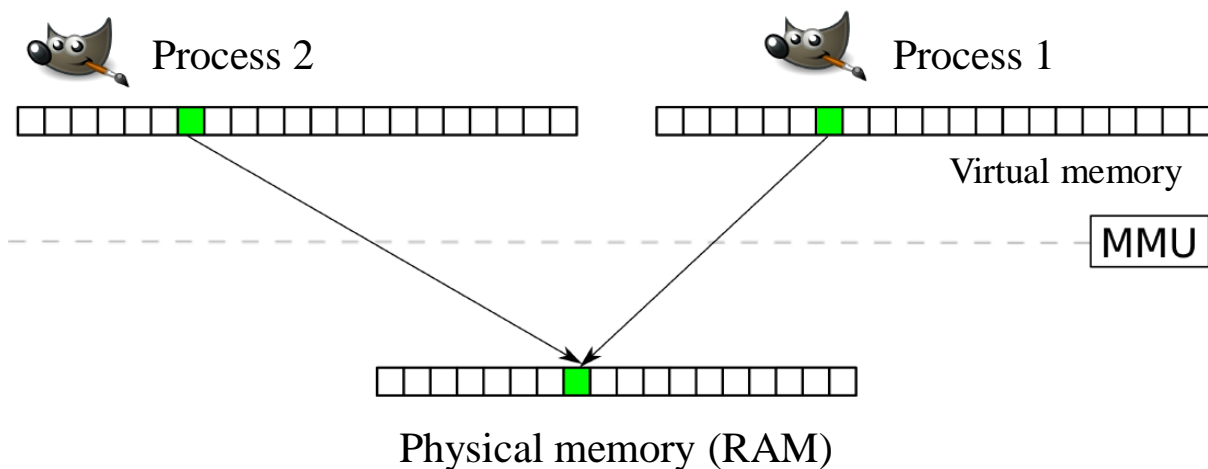


énergie atomique • énergies alternatives

- **Virtual memory :**



- **Shared memory :**

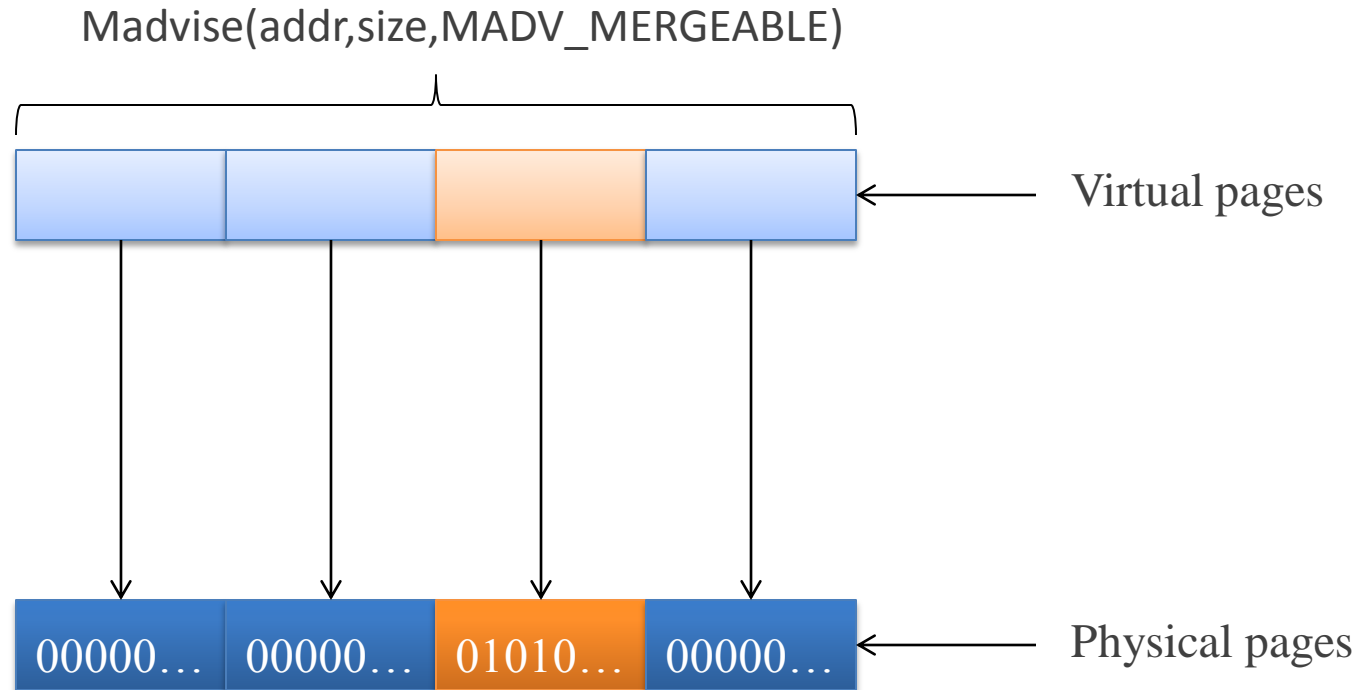


# How to use it



énergie atomique • énergies alternatives

- The user only need to mark the targeted segments.
- Marked with system call : *madvise()*

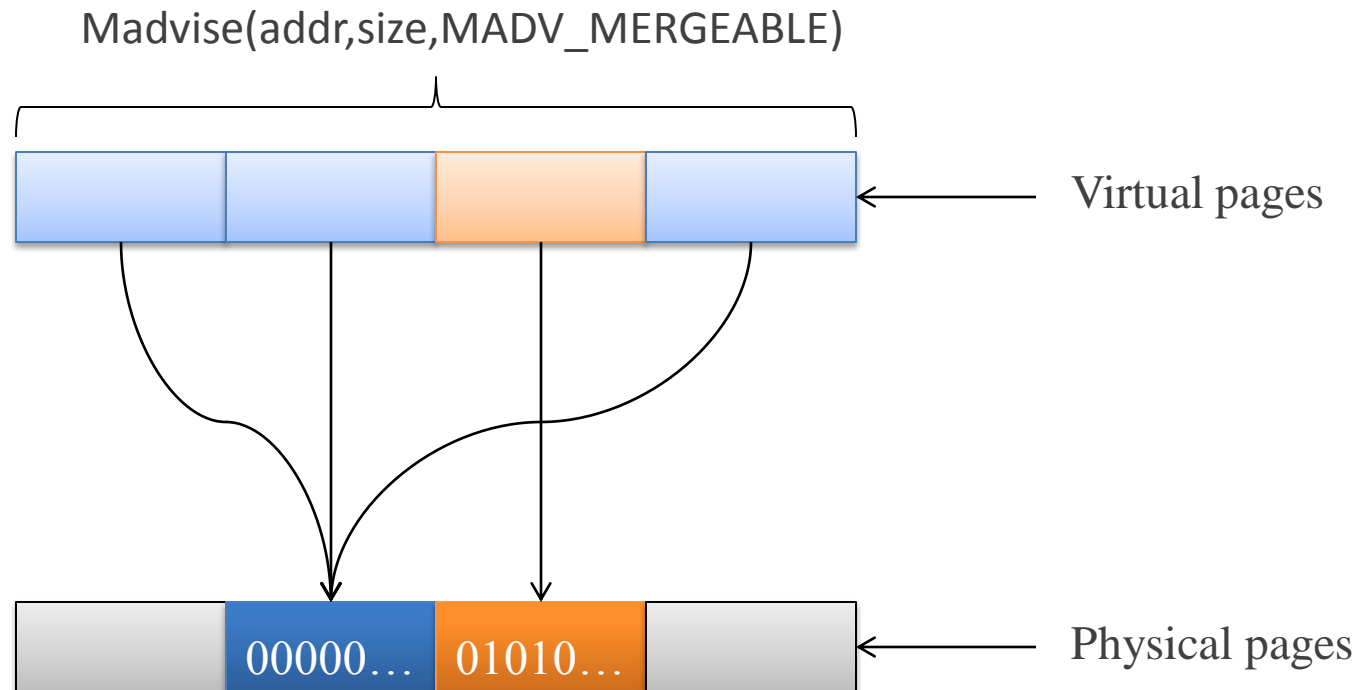


# Page merging



énergie atomique • énergies alternatives

- The daemon (ksmd) scan the pages
- Identical physical pages are transparently merged



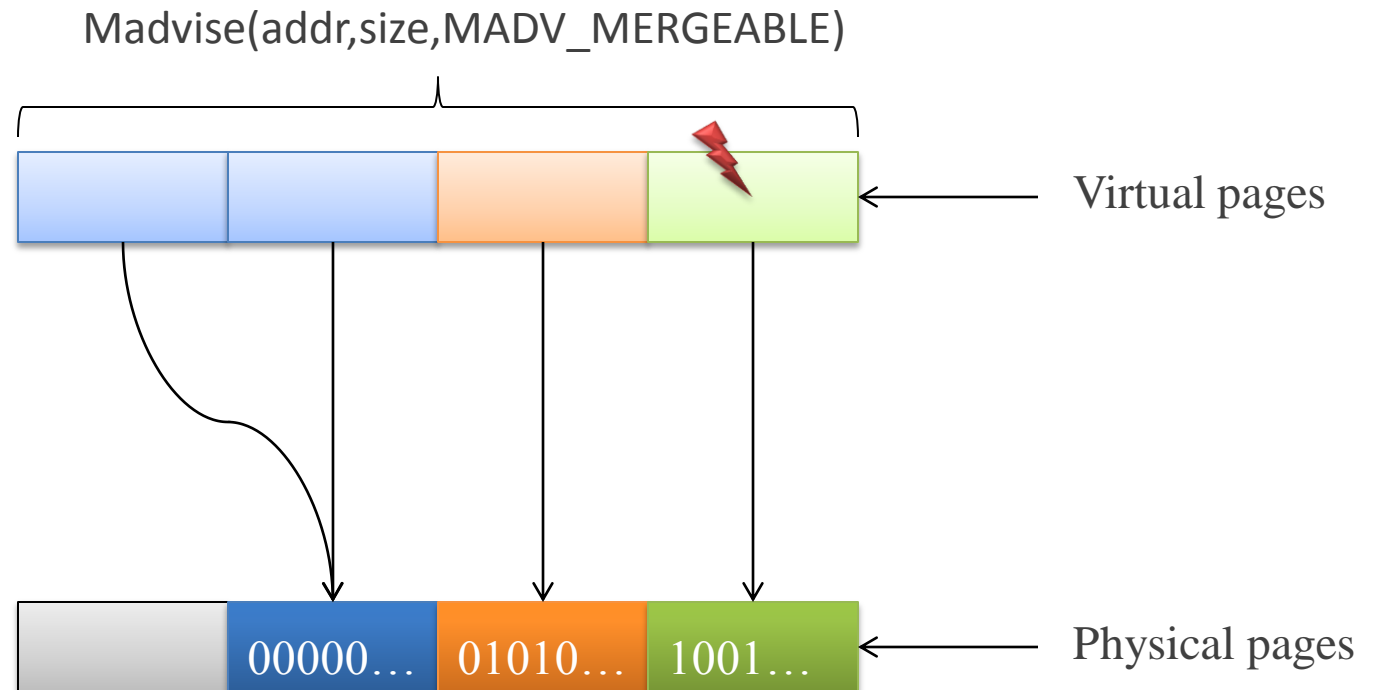


# Modifying pages



énergie atomique • énergies alternatives

- Usage of *Copy On Write* mechanisms
- All steps are transparent for the user



- **Configuration files : /sys/kernel/mm/ksm/\***
- **Parameters :**

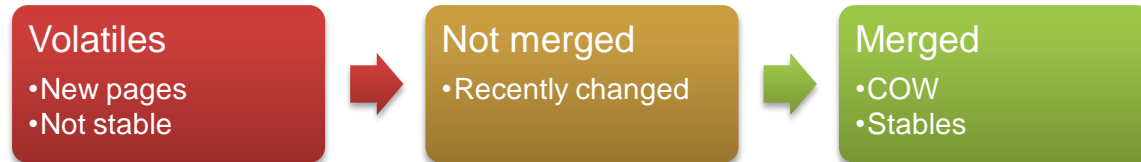
| Name                   | Default value | For our tests |
|------------------------|---------------|---------------|
| <i>Run</i>             | 0             | 1             |
| <i>Sleep_millisecs</i> | 20            | 20            |
| <i>Pages_to_scan</i>   | 100           | 2000          |

- **Access limits :**
  - Root to configure and enable/disable KSM
  - Everybody to use the *advise()* call.
- **Better to select candidate regions**



- **A daemon to automatically setup KSM parameters**
- **Permit to :**
  - Automatically enable/disable KSM on memory threshold
  - Automatically increase/decrease *pages\_to\_scan* parameter depending of the number of pages in queue
- **Available on :**
  - Debian
  - Fedora
  - Redhat

- **Three page class :**



- **Two red-black trees :**

- Stable (merged pages, sorted)
- Not stable (page ordering not fixed)

- **Tree search by comparing pages bits**

- **Checksum only used to detect unstable, but not for page search.**

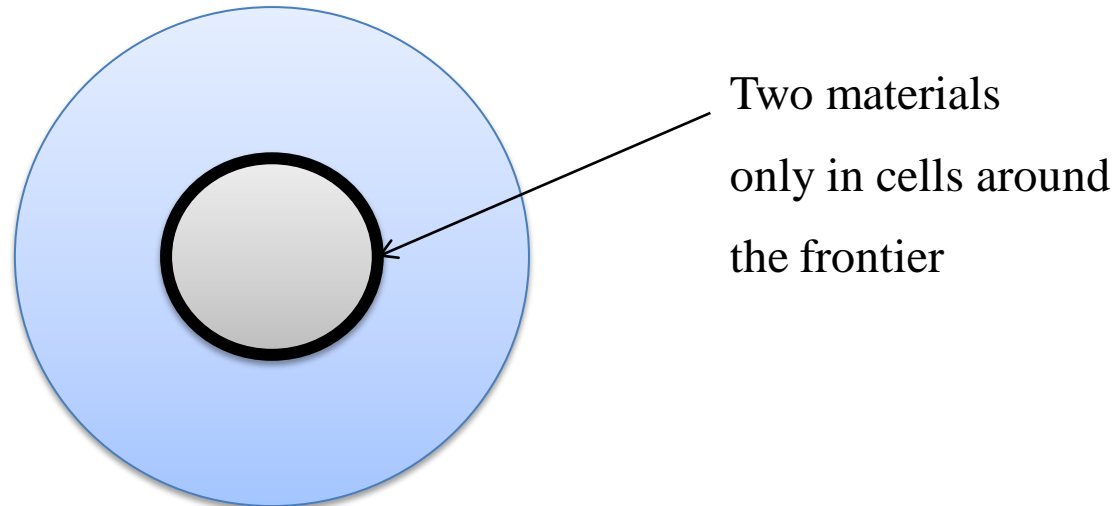


énergie atomique • énergies alternatives

# Evaluation on HERA



- **HERA is a 3D multi-physics, multi-material AMR platform**
- **Multi-materials imply :**
  - Require a cells entry for each materials (density)
  - Many location contain only one material
  - It imply many 0 in memory representation





- **Group blocs of 0 by using indirection at code level**
- **Done by overloading C++ operator []**
  - « Transparent » for the solver
- **Impact :**
  - Memory reduction of roughly 25%
  - Add indirections and loss of compiler optimizations
  - Can double the runtime
- **Is KSM an alternative ?**



énergie atomique • énergies alternatives

# How to measure KSM memory ?





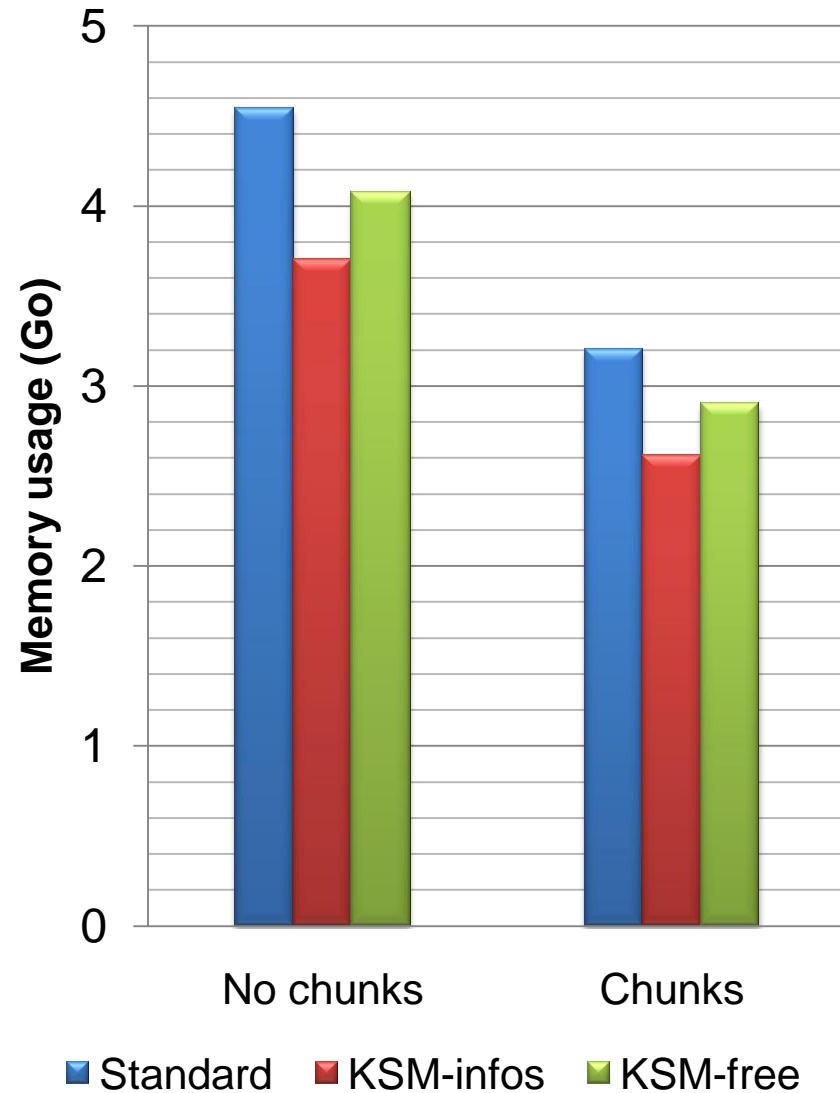
- **Memory measurement, not easy....**
  - How to evaluate merges between MPI process
  - Evaluator significance
- **Possibilities :**
  - ~~RSS of the current process (*Resident Segment Size*)~~  
*=> not updated by KSM*
  - Total free memory
  - KSM provide some stats in `/sys/kernel/mm/ksm/*`
    - Number of volatile pages (new or not stable)
    - Number of unshared pages
    - Number of shared pages

# How to observe the KSM memory usage ?



énergie atomique • énergies alternatives

- RSS not updated
  - What about SLURM ?
- KSM infos are biased and overestimate gains
- Free not limited to the process





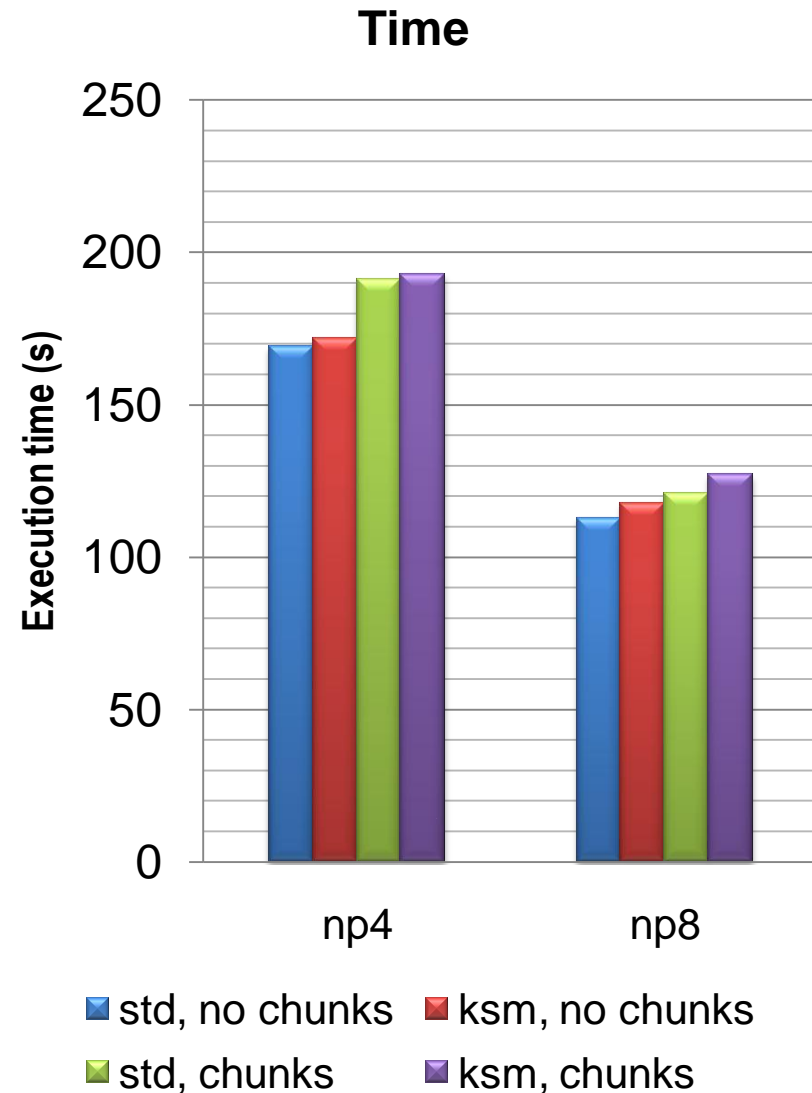
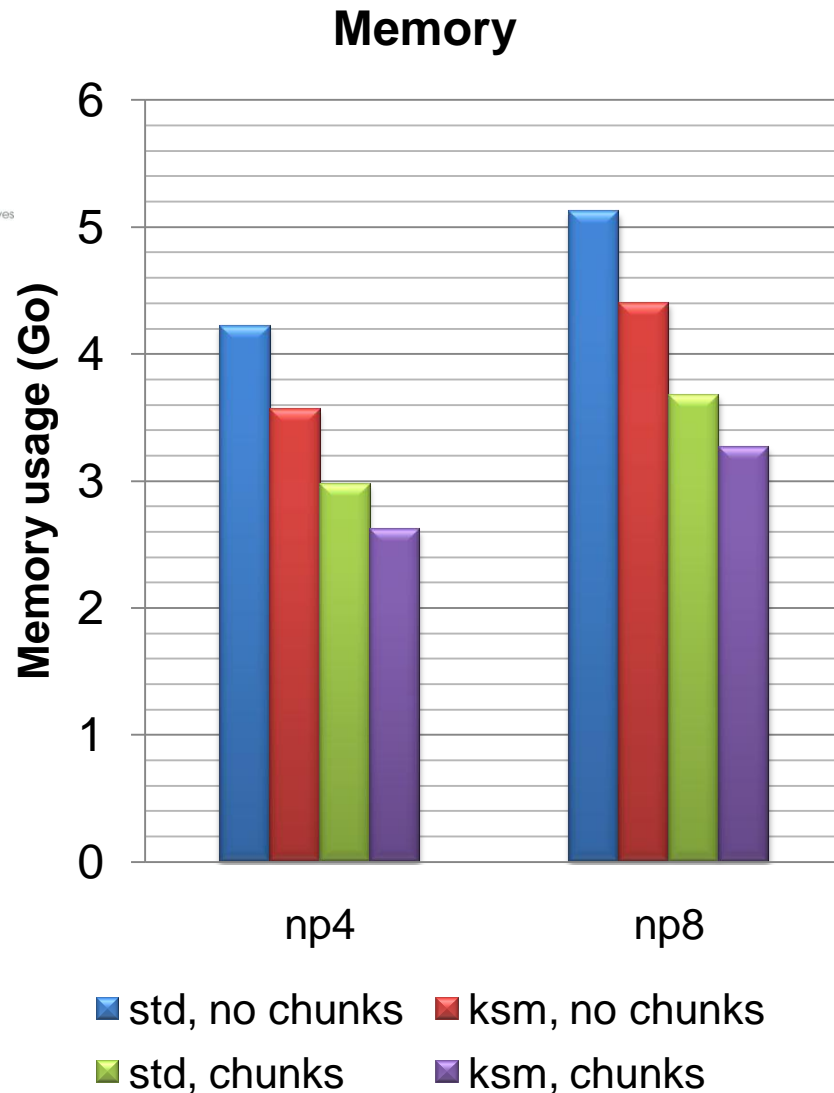
énergie atomique • énergies alternatives

# Results with HERA

# HERA - 4 or 8 MPI process – 6 materials



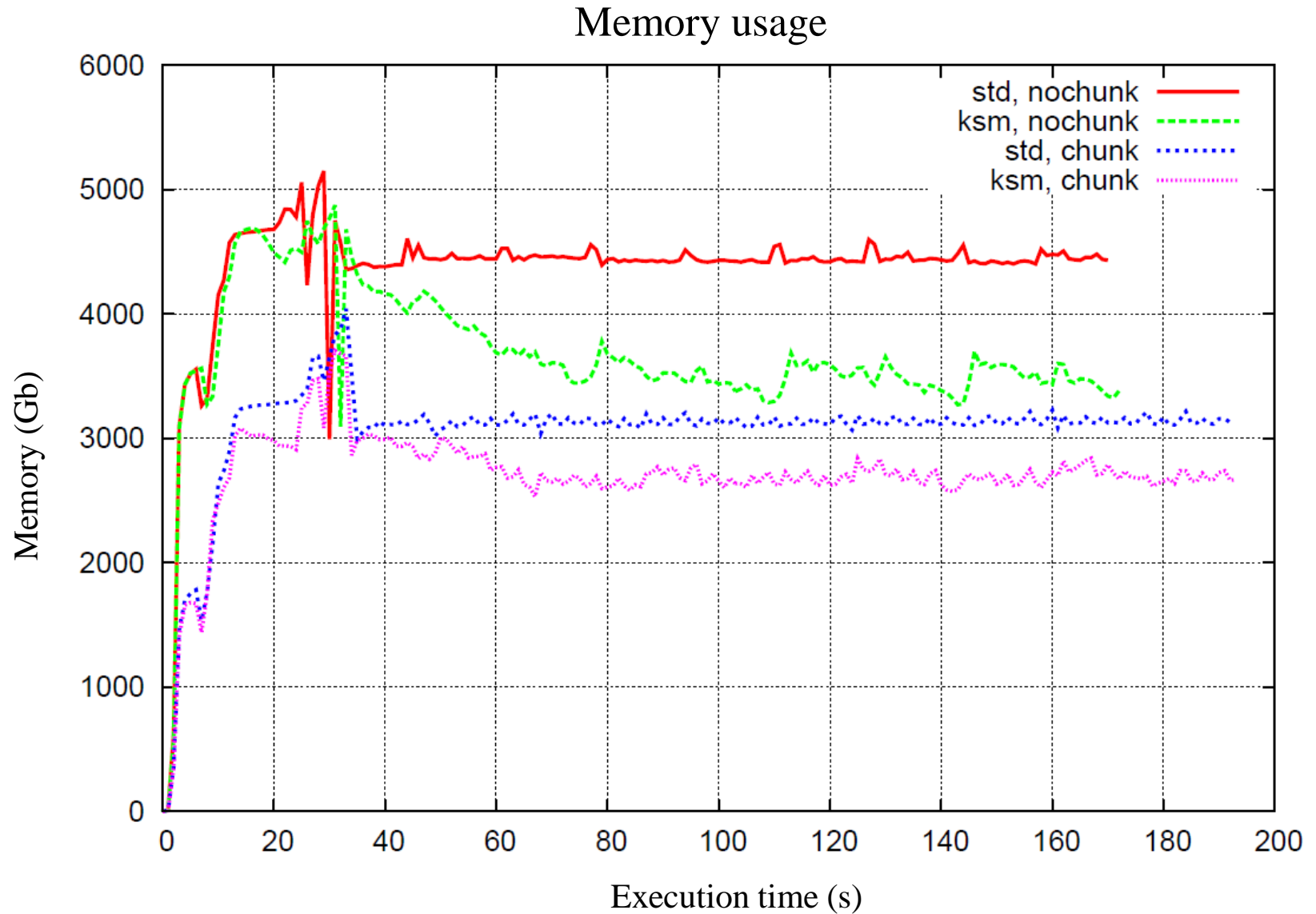
énergie atomique • énergies alternatives



# Memory usage over time (8 process)



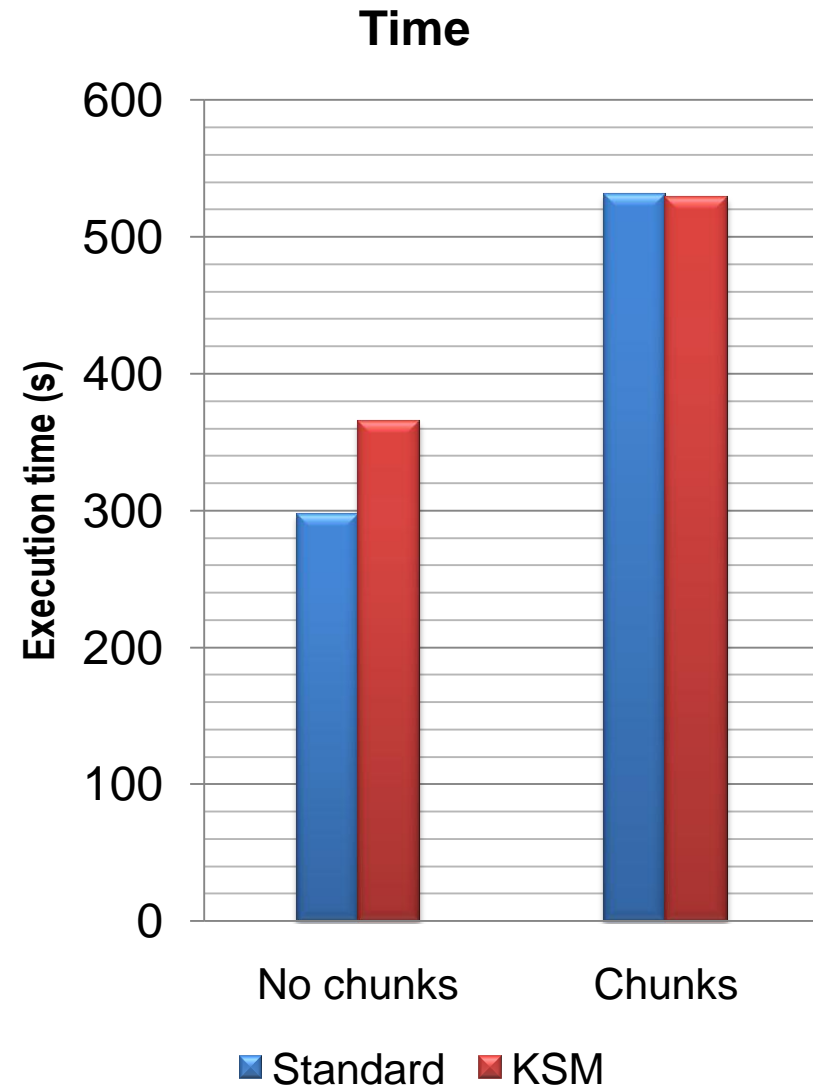
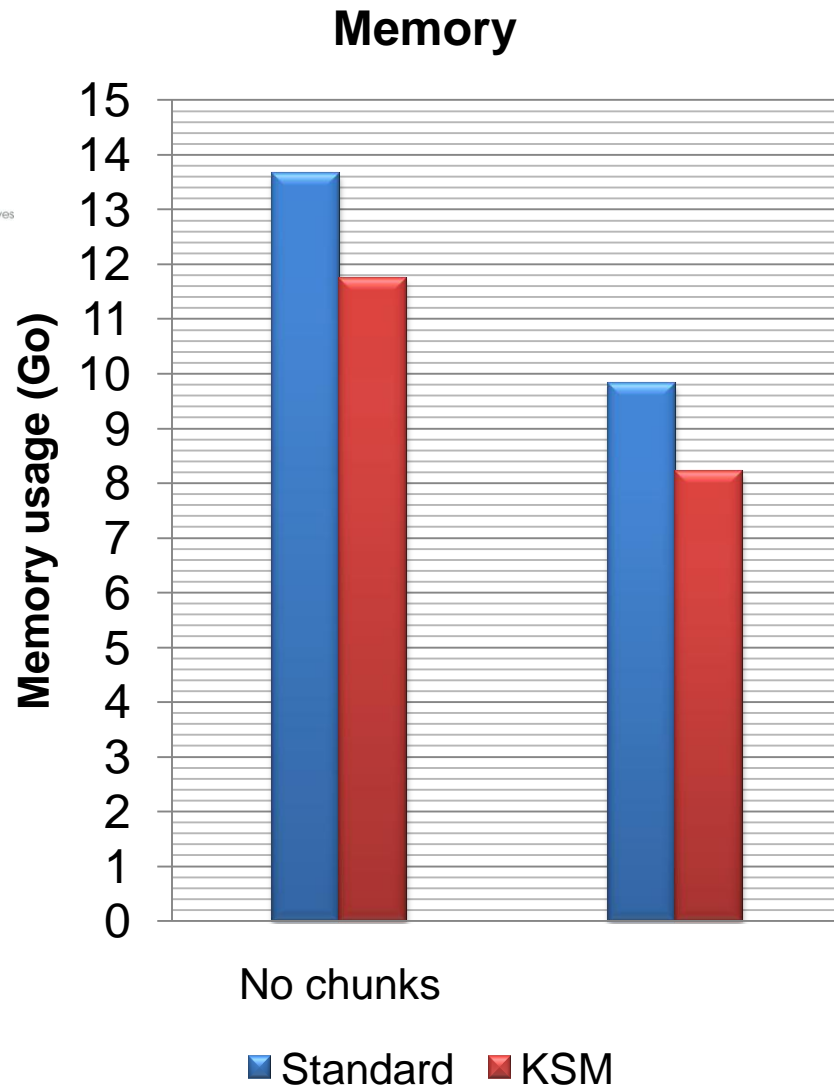
énergie atomique • énergies alternatives



# HERA - 8 MPI process – bigger case



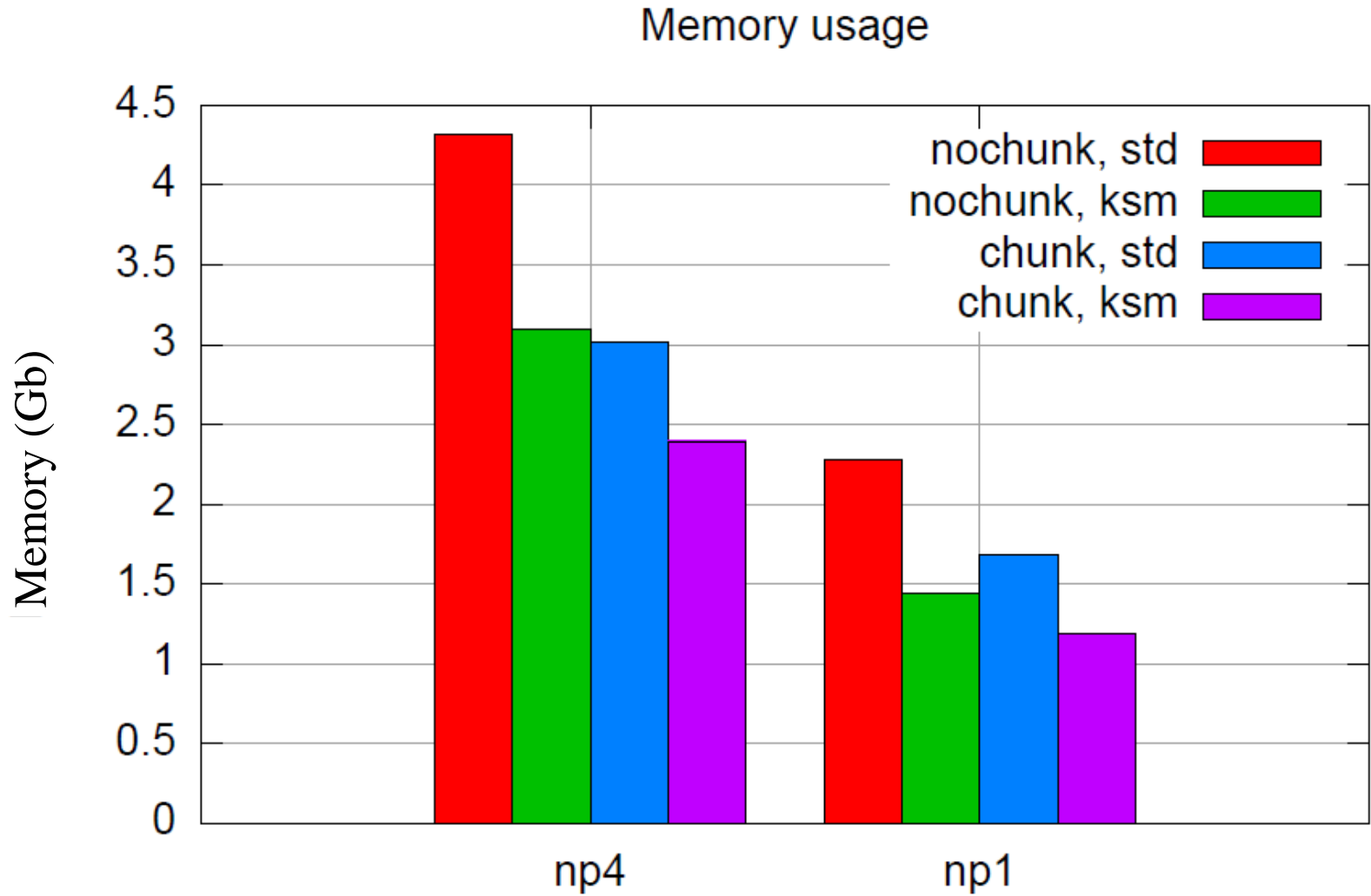
énergie atomique • énergies alternatives



# On slower Hera (-finstrument-function)



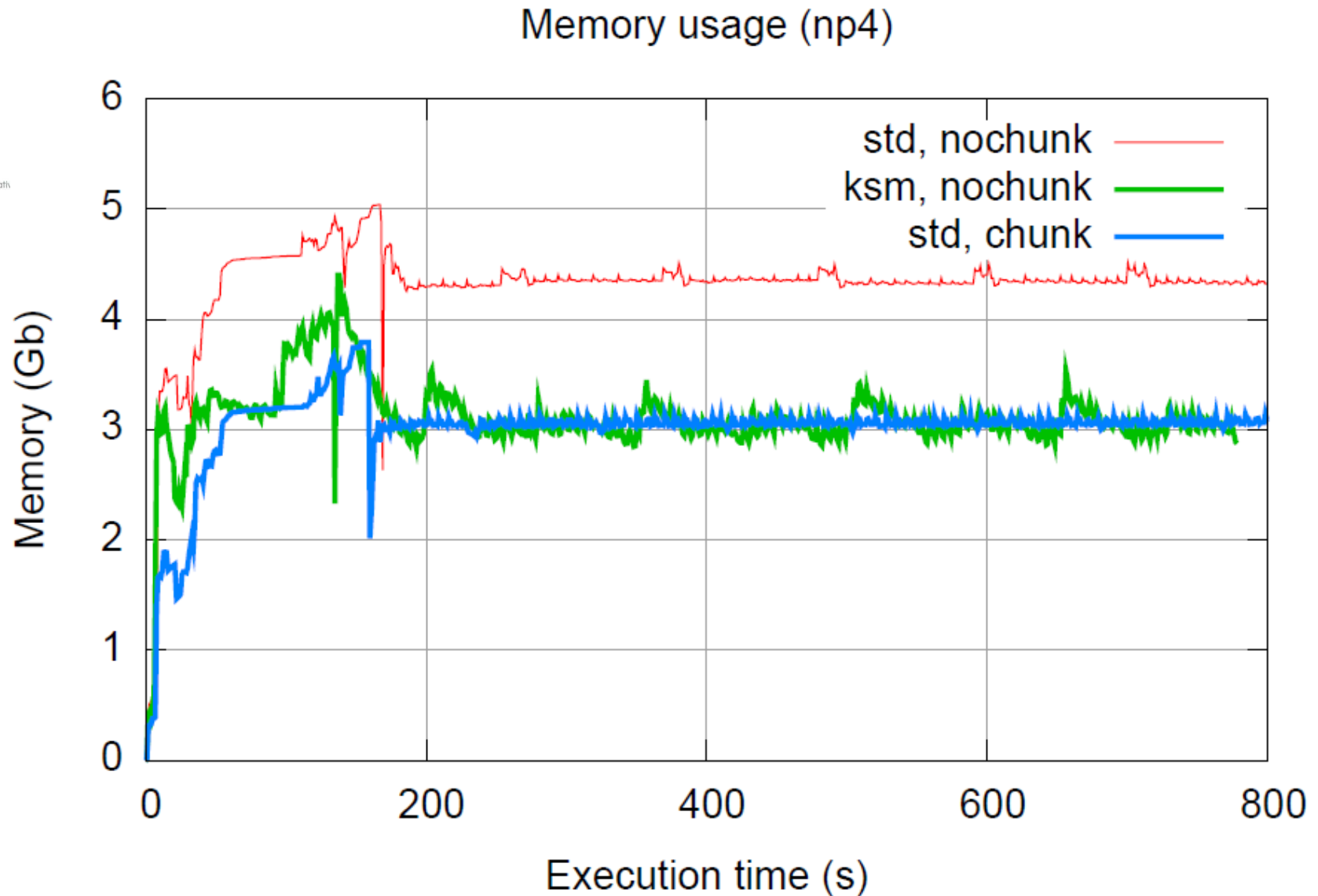
énergie atomique • énergies alternatives



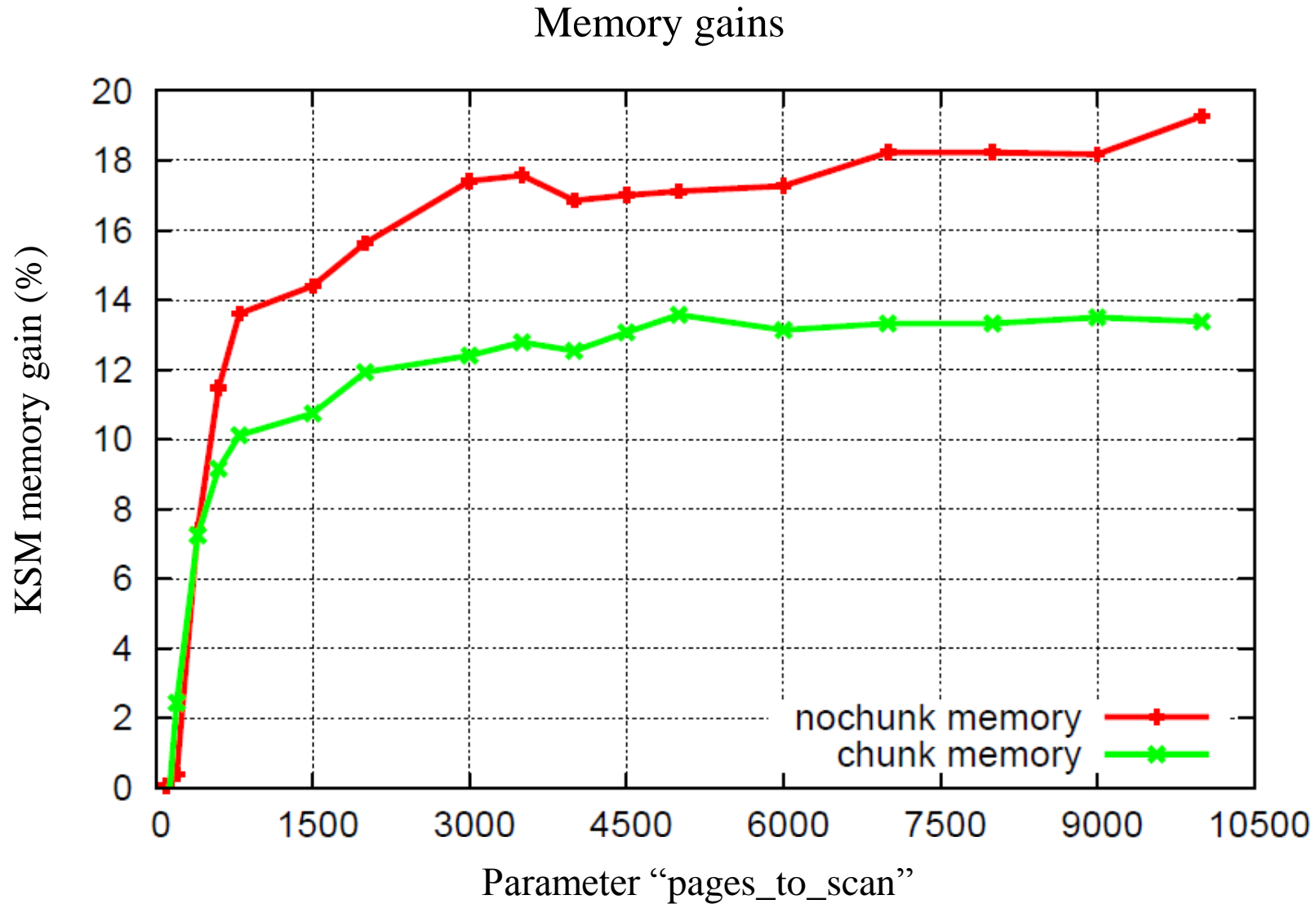
# On slower Hera (-finstrument-function)



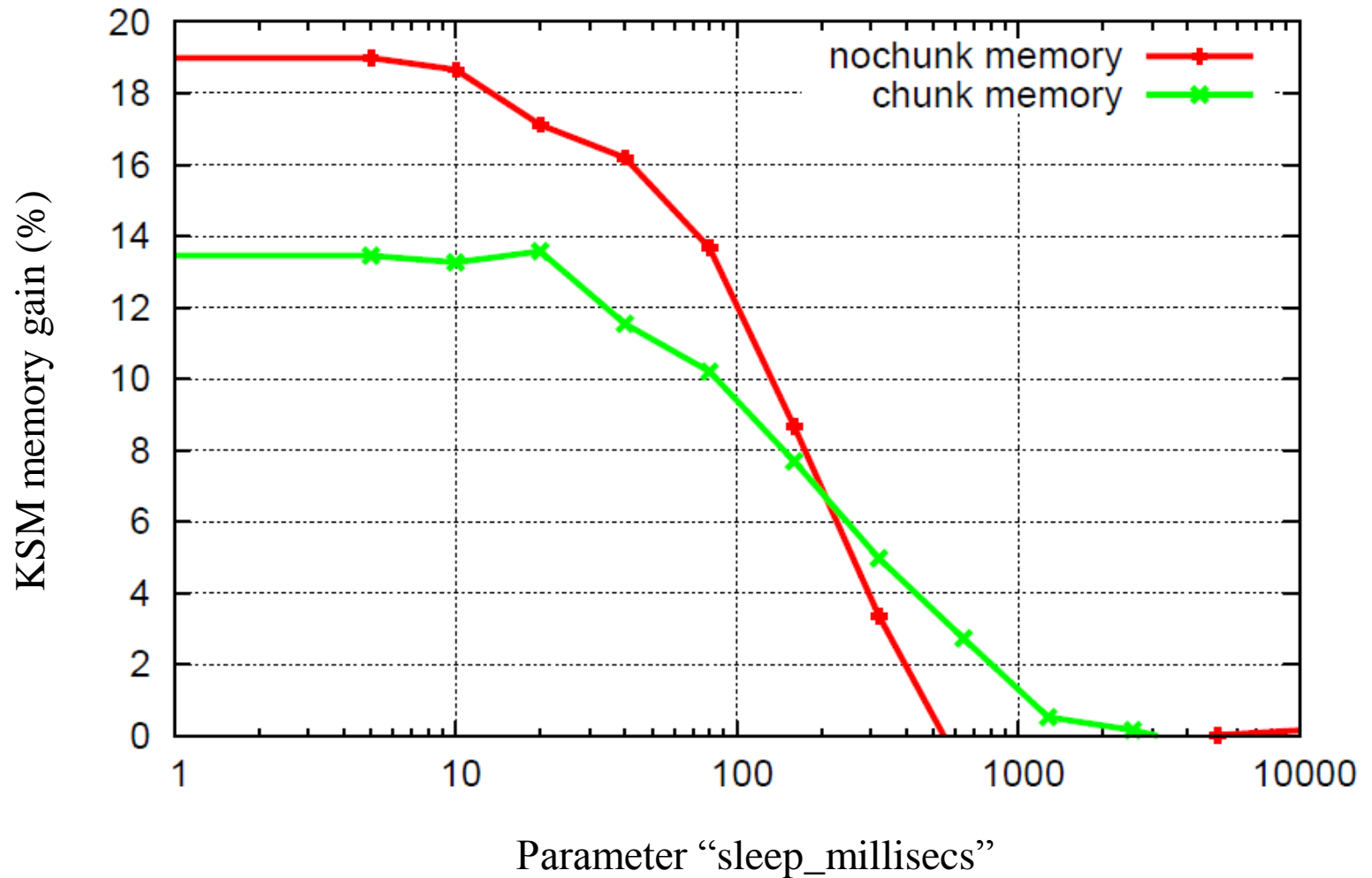
énergie atomique • énergies alternatives







## Memory gains





énergie atomique • énergies alternatives

# Conclusions



énergie atomique • énergies alternatives

- **Hardware limitations :**

- Require large similar segments (multiple of 4Kb)

- **Implementation limitations :**

- Non support of NUMA architecture
- Only one scanning thread
- Parameters are global and limited to root
- RSS isn't updated by KSM
- No way to do scan on demand or control KSM from process



- **Initially designed for VM but usable for applications**
- **KSM can reduce memory usage of some applications**
  - ~12-16% on Hera with or without chunks
  - Max performance degradation observed 23%, but mostly around of a few percent.
- **Provide a transparent way to implement indirections**
- **Get gains even with 4Kb limitation.**
- **Interesting, but need some improvement for usage in HPC**



énergie atomique • énergies alternatives

**THANKS**

- **[1] Andrea Arcangeli, Izik Eidus, and Chris Wright.  
Increasing memory density by using KSM. In Proceedings of  
the Linux Symposium, pages 19–28, July 2009.**



énergie atomique • énergies alternatives

# BACKUP



