

Dynamic Trace-based Sampling Algorithm for Memory Usage Tracking of Enterprise Applications

Houssem Daoud
Ecole Polytechnique Montreal
Montreal, Quebec h3t 1j4
houssem.daoud@polymtl.ca

Naser Ezzati-jivan
Ecole Polytechnique Montreal
Montreal, Quebec h3t 1j4
n.ezzati@polymtl.ca

Michel R. Dagenais
Ecole Polytechnique Montreal
Montreal, Quebec h3t 1j4
michel.dagenais@polymtl.ca

Abstract—Excessive memory usages in software applications has become an increasingly frequent issue. A high degree of parallelism and difficulty of monitoring for the developer can quickly lead to memory deficiency or can increase the time needed for possible garbage collections. There are several solutions introduced to monitor memory usages of softwares however they are not usually efficient or scalable. Massif is one of the most common memory monitoring tools but it easily doubles the execution time of the under-investigating application. Execution tracing introduces low overhead means of runtime memory and resource usages monitoring, however the huge size of the collected trace data is a challenge for later analysis. In this paper, we propose a dynamic trace-based sampling algorithm to collect and analyse run time information and metrics of memory usages. In the reality, it is implemented as a kernel module which gathers memory usages data from operating system structures when only a predefined condition is set or a threshold is passed. The thresholds and conditions are preset but they can be changed dynamically based on the application behavior. We tested our solutions to monitor several applications and our evaluation results show that the proposed method generates smaller trace data and reduces the time needed for analysis without losing the precision.

I. INTRODUCTION

Analysing software systems is becoming difficult and complex when we know that the systems that are hosting them are getting complex in hardware, include multiple nodes in parallel, each containing a large number of cores as well as parallel co-processors for graphics, signal processing and networking. When virtualization comes to standards for servers it becomes even more tedious which complicates understanding of their architecture, functionalities and possible performance misbehaviors.

Dynamic analysis through execution tracing is a solution for run time analysis of such systems. Tracing works by hooking in the different places of the software and collecting runtime data while the software is running. It usually gathers the valuable information about system execution and can be used in software comprehension and finding problems and misbehavior.

Although tracing in general is a great solution to analyse runtime behavior of systems, it may present some challenges: the trace size is usually huge (magnitude of gigabytes for only a few second of tracing). Huge size of tracing data means more required storage space, long analysis time and maybe lost opportunities in the timely detection of sensitive problems. The

techniques to alleviate this problem can be well received by the community.

Our goal is basically to analyse memory usages of the systems from the kernel point of view. Many state-of-the-art tools perform virtual memory monitoring from the userspace [10], [3]. Their methods mostly require the instrumentation of the memory allocator of the programming language or the use of a pre-loaded library that overrides original memory manipulation functions. Those solutions are not portable and can only target a specific system.

We propose a generic trace-based architecture to monitor and analyse memory usages of any applications. The challenges we face and aim to solve in this trace-based dynamic memory usage analysis method are as the following:

- The high frequency of memory operations makes the trace file huge
- It is not usually possible to reduce trace size by just targeting a single process using basic trace filtering techniques, since the actual physical memory releasing is done out of the process context
- Tracing can contribute in polluting the memory of the system

This paper provides a dynamic sampling technique to reduce the amount of trace data while it is generating. The technique works by hooking on some memory related functions (i.e., in their tracepoints) in the operating system kernel and listen to the events. Then, instead of generating events for each occurrence of the function, it only samples when a predefined condition is true, e.g., a predefined but changeable threshold is passed or a particular time duration is elapsed. The proposed method overallly:

- Instruments the Kernel to get the required information
- Provides filtering and aggregation mechanisms based on some thresholds to reduce the frequency of events.
- Generates metrics and visualization from the trace file

This can lead to a greater efficiency when we use it to trace high frequent operations like memory allocation, in which the frequency of the function call (and therefore the potential trace generate rate) is really huge (10k+ per just a millisecond). As we will see later in the paper, this approach reduces the size of trace and analysis time while giving the same analysis output and precision.

The contributions of the paper lies in the proposal of an efficient architecture for an in-kernel trace sampling and aggregation method. The solution is tested to track high frequent kernel memory allocation functions and confirms the efficiency and usefulness of the method.

In the remaining of the paper, we firstly investigate the related work. Then, after talking about the motivations of the work, we propose our architecture and conclude with the use cases and evaluations.

II. RELATED WORK

Tracing is a dynamic analysis method that collects trace logs of a runtime software systems executions. A trace log can present a function call, a system call, a signal, a network packet, etc. Unlike debugging which is a step by step procedure to go through the program execution to get its current state or discover a problem, tracing is more a background process which collects runtime data while the execution is going on and stores them in a disk or a network storage, to be analysed later [9]. The impact of tracing must be as low as possible to preserve the normal behavior of the software.

LTTng [6] is a low impact open source Linux tracing tool developed at DORSAL Lab¹ to provide tracing capabilities for Linux kernel and user-space applications. Kernel tracing can be performed dynamically using Kprobes or statically using the TRACE_EVENT macro. Trace generated by LTTng can be used to analyse run-time behavior of systems. For instance, trace-driven tools are proposed in the literature for disk block analysis [5], virtual machine analysis [1], [2], userspace level applications (e.g., Chromium browser) [15], [12], web applications [4] and live stream analysis [8]. Since LTTng is a low-impact tracing [7] and provides data at multiple levels of the system, it has been used in this paper to collect the tracing data about the memory usages.

Griswold et al. [10] instrumented memory allocation and liberation operations of the *Icon* programming language. The instrumentation is done using macros, which is low performance and not optimized of multi-threaded applications. The trace generated is used to present the memory usage as a 2D graph where each object allocated is shown as a rectangle proportional to the size of the allocation.

GCspy [14] a memory monitoring framework developed by Printezis et al. This tool has a client-server architecture: the collection of data is done on the server side and the visualization on the client side. The data collection requires the use of an instrumented memory allocator. The authors started by instrumenting the Java Virtual Machine (JVM) to track memory allocation in the Java programming language.

Cheadle et al. [3] extended GCspy to support dlmalloc, the C memory allocator used by the Glibc library. They also proposed different optimization to GCspy to reduce the frequency of communication between the client and the server as well as an automatic problem detector. GCspy uses animated images

to show the state of the heap memory throughout time. This visualization is not appropriate with some applications since the refresh rate of the animation can not handle high frequency memory events.

Jurenz et al. [11] extended VampirTrace, a performance analysis tool, to provide a detailed memory analysis based on execution traces. Instrumentation is done using shared library pre-loading. Original memory manipulation symbols like malloc, free, calloc, etc. are overridden with new functions that contain the required tracepoints. The limitation of this method is that it targets only the processes that are started with the pre-loaded library. There is no way to trace an already running program or to trace all the programs running on the system at the same time. Tracecompass also provides a memory analysis view that uses library pre-loading to collect Lttng-UST traces.

Massif is a heap memory profiler provided with Valgrind [13]. It hooks into the application loading code using LD_PRELOAD and shows its heap memory usages and can be used to optimize and reduce the memory usages of programs. Although Massif is a useful tool to monitor memory usages but it easily doubles the execution time of the application. We address an scalable and efficient solution in our paper to monitor memory usages of applications from operating system point of view.

III. MOTIVATION

Tracking the memory usage is important to know the processes that are demanding and consuming more memory and therefore to detect performance problems. The operating system uses complex mechanism to manage the physical memory. When a process allocates memory, a new virtual memory space is created and assigned to the process. The real physical memory allocation is done afterward when the memory is actually used by the process. The opposite operation is even more tricky. When a process asks kernel to release (free) the memory, kernel releases it from its virtual memory space. However, the main physical release is happens only when another process needs to get more physical memory.

Many tools use sampling to track the memory usage of processes using the */proc* filesystem. Most of those tools use a sampling frequency of 1 Hz, which gives a low precision result. Using a fixed sampling frequency is not a good solution since the the same frequency value can be very slow for some processes and fast for others, depending on the memory usage pattern.

Tracing is another way to get information from the kernel. If we trace memory operations like allocation and liberation, we can track in a very precise way the memory usage of processes. The problem is that those events are high frequency events and it is very likely to have lost events, which prevent us from having precise computations. Even if we are able to collect all events using big tracing buffers, the trace file will be huge and very difficult to read and analyse.

One idea to limit the trace size is to filter the tracing events to include just one process, but this solution is not possible

¹Distributed open reliable systems analysis lab (DORSAL) <http://www.dorsal.polymtl.ca>

because, as mentioned earlier, the physical memory liberation doesn't always occur in the context of the target process. It may happen in the context of kernel threads or in the context of another process that reclaims the memory.

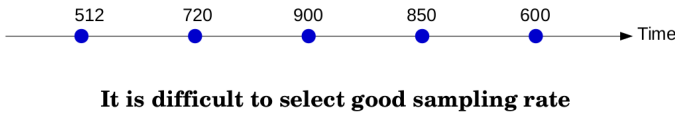


Fig. 1. Sampling

In this paper, we combine the benefits of both approaches, tracing and sampling, by proposing a dynamic trace-based sampling method. The sampling rate is defined based on the frequency of the related trace events.

IV. ARCHITECTURE

Memory management of modern operating systems is much more complex than before. Each process has a contiguous virtual address space in which he allocates required memory objects. A physical memory page is associated to a virtual one by the page fault handler only when the process actually accesses it.

In this paper, we provide a tool to monitor virtual and physical memory usage using a combination of tracing and sampling techniques. The proposed architecture is shown in Figure 2.

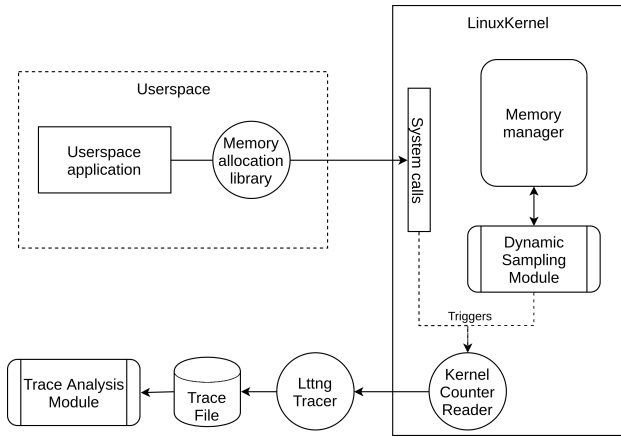


Fig. 2. Architecture

A. Virtual memory monitoring

The method tracks the virtual memory usage from the kernel space. It traces the different system calls related to memory allocation and release and uses them as triggers for the Kernel Counters Reader.

Memory-related functions like `malloc()`, `calloc()`, `realloc()`, and `free()` use system calls to interact with the operating system modules where the memory management is accomplished. Table I shows the mapping between library functions and system calls.

TABLE I
MAPPING BETWEEN MEMORY FUNCTIONS AND SYSTEM CALLS

	size \leq MMAP_THRESHOLD	size $>$ MMAP_THRESHOLD
Malloc calloc realloc	sbrk	mmap
free	None, or sbrk(negative) depending on M_TRIM_THRESHOLD	munmap

The behavior of the allocator differs depending on the size of the allocation. Small allocations are done using `sbrk` system calls, which increases the size of an existing virtual space region. In this case, releasing the memory allocated doesn't automatically reduce the size of the virtual space. In contrast, big allocations are done using `mmap` and released right away by `munmap` after the memory is freed by the applications (Figure 3)

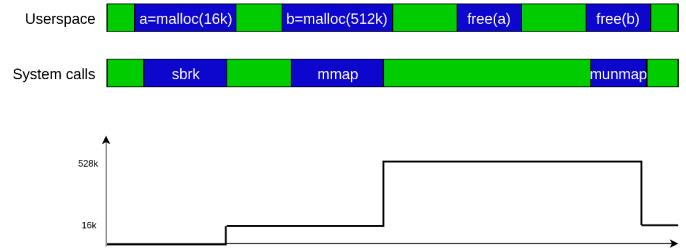


Fig. 3. Virtual memory growth after allocation and release operations

The virtual memory of a specific process also grows when a shared library is loaded or when a shared segment is mapped into the address space using `shmat`.

The system calls cited in Table I are used as triggers for the Kernel Counter Reader which reads the exact value of virtual memory usage from the `mm_struct` of the concerned process and generates an Ltng trace event.

B. Physical memory monitoring: Dynamic sampling algorithm

Physical memory monitoring presents a big challenge compared to virtual memory: the memory manager of the operating system generates a huge number of events and recording all of them in a trace file is almost impossible.

Sampling can be a good solution for this case, but choosing a sampling rate is not an easy task. Some processes demand and access the memory very frequently during the execution and others access it less frequently. A low sampling rate gives a bad precision, but on the other hand, a high sampling rate generates a huge amount of useless data for inactive processes.

In Algorithm 1, the method that dynamically changes the sampling rate depending on processes activity is displayed.

It is a 2D sampling algorithm that uses the time and the memory variability to choose the appropriate moment of get memory usage information from the Kernel data structures. An event is generated if the timer finishes or before that

Input:

Sampling rate
 Variability threshold

```

//Main thread
startTimer(rate)
if Timer tick then
  processes ← ListSystemProcesses()
  for process in processes do
    trace_memory(process)
    variability[process] ← 0
  end
  restartTimer()
end

//Kprobe hook
if memory page allocated/released then
  process ← getCurrentProcess()
  variability[process] += PAGE_SIZE * direction
  if variability[process] exceeds the threshold then
    trace_memory(process)
  end
end

```

Algorithm 1: Dynamic sampling Algorithm

if the memory variability of a process exceeds a certain threshold as shown in Figure 4. The blue points represent events timer events, and the red points are the events caused by the threshold.

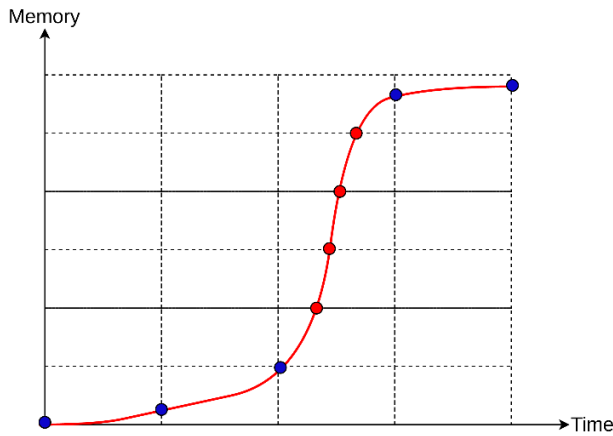


Fig. 4. (Time, Space) Sampling

Memory variability is computed by hooking on `kmem_mm_page_alloc` and `kmem_mm_page_free` events which occurs when a physical page is allocated or released. The Kernel keeps information about physical memory usage in the `mm_struct` data structure (RSS: Resident set size) This counter is adjusted each time a physical is inserted or removed from the page table of the process.

The proposed algorithm is implemented as a Kernel module and is configurable through the `proc` file system (sampling rate, variability threshold).

Lock-free data structures are used to provide a good scalability:

- RCU Hashmap is used to hold process information
- Memory variability is defined as `atomic_long` to avoid using heavy synchronization mechanisms.

V. EVALUATION

In this section, we evaluate the performance and the usefulness of our tool. We performed benchmarking with a synthetic workload and then with real applications. We compared our method with Massif, another state-of-the-art tool, to confirm the correctness of our analysis.

A. Performance

The performance tests are executed on an Intel i7-4790 CPU @ 3.60GHz with 6 GB of main memory and an Intel SSD 530 Series 240 GB hard disk, running Linux Kernel version 4.4. The traces are collected using LTTng 2.8.

The following cases are used for benchmarking:

- *No tracing*: the program runs without any tracing mechanism
- *Dynamic sampling*: We used our dynamic sampling module with a sampling period of 10 ms and a memory variability threshold of 10 MB.
- *LTTng all memory events*: We traced all memory allocation and release events. We used the tracepoints `kmem_mm_page_alloc` and `kmem_mm_page_free`.
- *Massif*: We used Massif, a widely used memory monitoring tool.

A program is developed to generate a memory access workload. It allocates, accesses and frees a memory buffer of a certain size, and repeats the operation until it reaches 20GB of workload. The execution time of this program with the different configurations is presented in Table II and Figure 5. The same benchmark is also performed with real applications: Firefox, a widely-used web browser, and Totem, the default movie player of GNOME desktop. The results are reported in Table III.

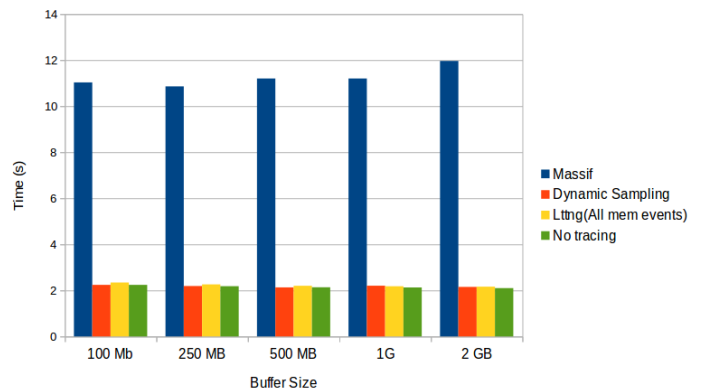


Fig. 5. Tracing impact on execution time

The results show that the overhead of Massif is very high compared to other cases. It is 5x slower with the benchmarking

TABLE II
EXECUTION TIME IN SECONDS IN FUNCTION OF MALLOC BUFFER SIZE WITH DIFFERENT TRACING MECHANISMS

Malloc Buffer Size / Tracing Mechanism	Massif	Dynamic Sampling	LTTng (All memory events)	No tracing
100 Mb	11.03	2.24	2.34	2.239
250 MB	10.86	2.19	2.26	2.183
500 MB	11.2	2.13	2.2	2.134
1G	11.2	2.203	2.18	2.126
2 GB	11.96	2.15	2.16	2.097

TABLE III
EXECUTION TIME IN SECONDS OF SOME APPLICATIONS WITH DIFFERENT TRACING MECHANISMS

Application / Tracing Mechanism	Massif	Dynamic Sampling	LTTng (All memory events)	No tracing
Firefox	51	2.509	2.59	2.51
Totem (10 seconds video)	28	10.641	10.645	10.752
Benchmark application (Buffer size = 500 Mb)	11.2	2.13	2.2	2.134

program, 3x slower with Totem and 20x slower with Firefox. In contrast, the overhead of the two other cases is almost negligible. It doesn't exceed 1% in all cases.

It is expected that the Dynamic sampling algorithm and *LTTng all memory events* are similar in terms of execution time because we are tracing the same kernel functions in both cases, the difference is in the number of events generated. Table IV shows that with a sampling period of 10 ms and a memory variability threshold of 10 MB, we can reduce the size of trace by a factor between 3 and 7 for a normal workload.

An interesting phenomenon happens when the buffer size is more than 4 GB. The operating system goes into a state of thrashing. Memory pages start to be moved between the main memory and the swap space, which creates a huge memory activity. Tracing all memory activity at this moment is very inefficient and somehow impossible since the number of events is very high. The Dynamic Sampling Mechanism is able to handle this case by filtering the unnecessary events at the kernel side.

B. Correctness

In this section, we use our tool to monitor the virtual and the physical memory usage of different applications and we use massif, Valgind memory profiler, to make sure that the result is the same with both tools.

At first, We traced our program using the dynamic benchmarking mechanism and we used TraceCompass² to show the results graphically. The output of our tool (Figure 6) corresponds perfectly to logic behind the code. The program allocates 500 Mb, access the allocated memory and the frees it.

We can see that the virtual memory, shown in red, is allocated when malloc() is called. The physical memory, shown in blue, is allocated when the memory pages are accessed using memset(). The virtual and physical memory are released during the free() function call.

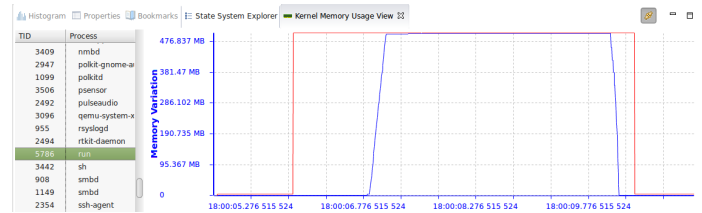


Fig. 6. Virtual and physical memory usage monitoring

Figures 7 and 8 show that both tools give similar memory usage graphs for Firefox. Totem memory usage is also plotted by both tools in Figures 9 and 10 which displays the same output for both approaches.

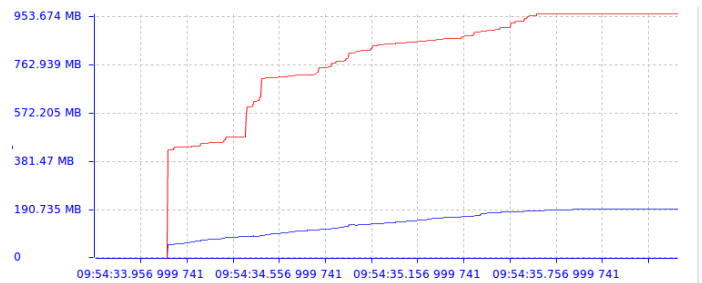


Fig. 7. Firefox memory usage at startup using LTTng

Our tools also provides other important advantages compared to Massif: it provides physical memory usage and the analysis covers all the processes running on the system at the same time, not only one targeted process.

VI. CONCLUSION

In this paper, a framework to collect memory usages information of enterprise application is proposed. It includes a dynamic sampling algorithm to gather runtime information from operating system kernel. The method checks if a certain time is passed or a threshold is reached then it gathers the information from kernel data structures and generates trace

²<http://www.tracecompass.org>

TABLE IV
NUMBERS OF EVENTS IN THE TRACE FILE GENERATED BY LTTng (ALL MEMORY EVENTS) AND THE DYNAMIC SAMPLING MECHANISM

malloc size / tracing mechanism	LTTng (All memory events)	Dynamic Sampling	Reduction factor
5 GB	6367016	150044	31.83
4 GB	4364234	58049	75.18
2 GB	102387	29686	3.45
1G	101200	28317	3.57
500 MB	124644	28425	4.39
250 MB	173558	32834	5.29
100 Mb	300623	40635	7.4

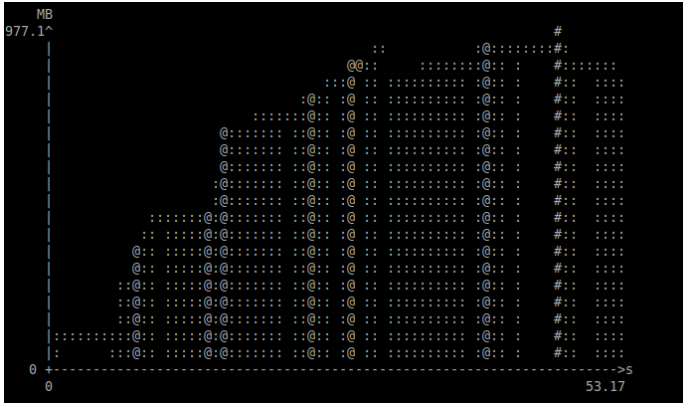


Fig. 8. Firefox memory usage at startup using Massif

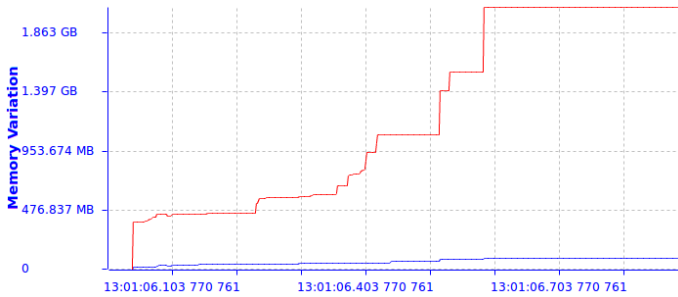


Fig. 9. Totem memory usage to play a video using LTTng

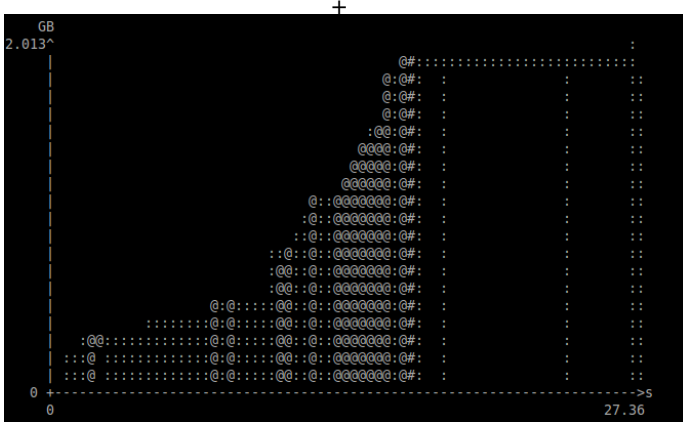


Fig. 10. Totem memory usage to play a video using Massif

events to be processed and analysed later. The thresholds are dynamic and can be decreased or increased based on application behavior and memory usages pattern (e.g., the rate of memory allocation calls).

We have tested our method against some real world applications like Firefox and Totem video player and the results prove that the performance cost of the proposed approach is neglectable while the precision is preserved.

The proposed solution is used to analyse memory usages, However, the architecture is generic enough to be used in any other resource usages metrics. It can be actually used for other high frequent tracing events within the operating system kernel, like network usages, disk I/O, etc. Extending the proposed method to support other kind of metrics and using other high frequent events will be investigated as a future work.

VII. ACKNOWLEDGMENT

We would like to gratefully acknowledge the Natural Sciences and Engineering Research Council of Canada (NSERC), Ericsson, and EffciOS for funding this project.

REFERENCES

- [1] C. Biancheri, N. Ezzati-Jivan, and M. R. Dagenais. Multilayer virtualized systems analysis with kernel tracing. In *2016 IEEE 4th International Conference on Future Internet of Things and Cloud Workshops (FiCloudW)*, pages 1–6, Aug 2016.
- [2] Cédric Biancheri and Michel R. Dagenais. Fine-grained multilayer virtualized systems analysis. *J. Cloud Comput.*, 5(1):69:1–69:14, December 2016.
- [3] Andrew M Cheadle, AJ Field, JW Ayres, Neil Dunn, Richard A Hayden, and J Nystrom-Persson. Visualising dynamic memory allocators. In *Proceedings of the 5th international symposium on Memory management*, pages 115–125. ACM, 2006.
- [4] Housseem Daoud and Michel R. Dagenais. Multi-level diagnosis of performance problems in distributed systems. *Journal of Systems and Software*, pages 1–14, 2017.
- [5] Housseem Daoud and Michel R. Dagenais. Recovering disk storage metrics from low-level trace events. *SOFTWAREPRACTICE AND EXPERIENCE*, pages 1–14, 2017.
- [6] Mathieu Desnoyers. *Low-impact operating system tracing*. PhD thesis, École Polytechnique de Montréal, 2009.
- [7] Mathieu Desnoyers and Michel R Dagenais. Lockless multi-core high-throughput buffering scheme for kernel tracing. *ACM SIGOPS Operating Systems Review*, 46(3):65–81, 2012.
- [8] Naser Ezzati-Jivan and Michel R Dagenais. Cube data model for multilevel statistics computation of live execution traces. *accepted in Concurrency and Computation: Practice and Experience*, 2014.
- [9] Naser Ezzati-Jivan and Michel R Dagenais. Multiscale abstraction and visualization of large trace data: A survey. *submitted to The VLDB Journal*, 2014.

- [10] Ralph E. Griswold and Gregg M. Townsend. *The visualization of dynamic memory management in the icon programming language*. University of Arizona, Department of Computer Science, 1989.
- [11] Matthias Jurenz, Ronny Brendel, Andreas Knüpfer, Matthias Müller, and Wolfgang E Nagel. Memory allocation tracing with vampirtrace. In *International Conference on Computational Science*, pages 839–846. Springer, 2007.
- [12] K. Kouame, N. Ezzati-Jivan, and M. R. Dagenais. A flexible data-driven approach for execution trace filtering. In *2015 IEEE International Congress on Big Data*, pages 698–703, June 2015.
- [13] Nicholas Nethercote and Julian Seward. Valgrind: A framework for heavyweight dynamic binary instrumentation. *SIGPLAN Not.*, 42(6):89–100, June 2007.
- [14] Tony Printezis and Richard Jones. *GCspy: an adaptable heap visualization framework*, volume 37. ACM, 2002.
- [15] Florian Wininger, Naser Ezzati-Jivan, and Michel R. Dagenais. A declarative framework for stateful analysis of execution traces. *Software Quality Journal*, pages 1–29, 2016. URL: <http://dx.doi.org/10.1007/s11219-016-9311-0>.