

MSDF: Memory Statistics Data Format used in system monitoring

La Quoc Nhut Huan^{1,2,3,*}, Nguyen Manh Thin^{1,3}, Nguyen Le Duy Lai^{1,3}, Nguyen Quang Hung^{1,3}, Thoai Nam^{1,2,3}



Use your smartphone to scan this QR code and download this article

ABSTRACT

High performance computing (HPC) system or computing system is very different from ordinary service system. In general, service system only run some specific applications, e.g. web server or mail server to serve as many requests from users as possible while in computing system, users have the permission to run their own applications and isolated with each other. Monitoring technique is the key to ensure system efficiency and users satisfaction, and by combining monitoring together with data analysis, system administrators can solve several operating problems specific to computing system such as resource allocation, application scheduling, abnormal detection, etc. Different from service system while administrators usually prefer system overall information rather than information of each individual user applications in computing system. Since computing system usually contains many applications executed simultaneously, monitoring computing system with traditional approaches would potentially consume a huge amount of storage space and would cost more charge fee if system is deployed in cloud environment.

This article focuses on analyzing monitored memory usage data retrieved from computing program in order to benefit its next resource allocation. Different from traditional approaches with batch processing technique in which collected data is all stored in database before analyzing, we utilized online analysis approaches in which every new coming data is captured, processed, cached in order to transform into useful information, and only allow necessary data be stored in database. We propose Memory Statistics Data Format (MSDF), an on-the-fly processing technique used in monitoring memory usage of computing application for saving storage space while still preserve enough information to solve resource allocation problem. MSDF can help to save more than 95% of storage space while allocation efficiency is always guaranteed depend on the parameter and MSDF can be extended to solve other operating problem or adapted to monitor and analyze other remaining application metrics.

Key words: System Monitoring, Memory Monitoring, Streaming Processing, Online Analysis, Memory Allocation

¹High Performance Computing Laboratory, Faculty of Computer Science and Engineering (HPC Lab), Ho Chi Minh City University of Technology (HCMUT), 268 Ly Thuong Kiet Street, District 10, Ho Chi Minh City, Vietnam

²Advanced Institute of Interdisciplinary Science and Technology (iST), Ho Chi Minh City University of Technology (HCMUT), 268 Ly Thuong Kiet Street, District 10, Ho Chi Minh City, Vietnam

³Vietnam National University Ho Chi Minh City (VNU-HCM), Linh Trung Ward, Thu Duc City, Ho Chi Minh City, Vietnam.

Correspondence

La Quoc Nhut Huan, High Performance Computing Laboratory, Faculty of Computer Science and Engineering (HPC Lab), Ho Chi Minh City University of Technology (HCMUT), 268 Ly Thuong Kiet Street, District 10, Ho Chi Minh City, Vietnam

Advanced Institute of Interdisciplinary Science and Technology (iST), Ho Chi Minh City University of Technology (HCMUT), 268 Ly Thuong Kiet Street, District 10, Ho Chi Minh City, Vietnam

Vietnam National University Ho Chi Minh City (VNU-HCM), Linh Trung Ward, Thu Duc City, Ho Chi Minh City, Vietnam.

Email: huan@hcmut.edu.vn.

1 INTRODUCTION

2 High Performing computing (HPC) system or computing system or computing system in general is critical for scientific research. One of its prominent examples of application is the artificial intelligence training for Smart Village project. Computing system contains a large number of nodes with special network topology and technologies to boost the parallel computing ability as much as possible. Different from ordinary service system such as web or mail server, users in computing system have the permission to access and execute their own programs, i.e. application. system resources are shared between multiple users and is allocated based on user's requirements and allocation policies defined by system administrators. As a consequence, efficiently managing and operating computing system with multiple users and a vast number of applications running simultaneously would cost admins much more effort.

Monitoring computing infrastructure helps admins continuously follow system operation, profile abnormal behavior, etc. and facilitate admins to update their management and operating plan in the future. "We are drowning in data but starving for information"¹, system monitoring often lacks of analysis ability to transform raw data into useful information and knowledge, which leads to the situation where every piece of data collected from any metrics considered potential for later analysis has to be permanently stored. As a result, since monitored metrics are collected at application level, applying traditional monitoring in computing system would potentially consume a huge number of storage space and cost more charge fee if system is deployed in cloud environment. The gap between data saved in storage and actual data in use can be filled with online-analysis which stores only necessary information retrieved from processing raw monitoring data.

Cite this article : Huan L Q N, Thin N M, Lai N L D, Hung N Q, Nam T. **MSDF: Memory Statistics Data Format used in system monitoring.** *Sci. Tech. Dev. J. – Engineering and Technology* 2024; ():1-8.

History

- Received: 25-9-2023
- Accepted: 29-3-2024
- Published Online:

DOI :



Copyright

© VNUHCM Press. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International license.



39 Memory usage which directly reflects system health,
 40 is always listed as one of the most critical metrics.
 41 Information get from analyzing the consumption of
 42 application memory benefits from solving several
 43 problems including resource allocation, application
 44 scheduling, abnormal detection, etc. In this article,
 45 we focus on monitoring memory data which is re-
 46 trieved from our monitoring framework, and analyz-
 47 ing these data to update memory allocation policy
 48 for the next execution. Different from normal ap-
 49 proaches in which data is collected and immediatly
 50 passed to storage, we utilized online-analysis in
 51 every new coming data is being processed, cached,
 52 and later converted into useful information before be-
 53 ing stored in database. We named this method *Mem-*
 54 *ory Statistics Data Format (MSDF)*, for later reference
 55 convenience. MSDF highlights the ability to poten-
 56 tially save a massive amount of storage size when ap-
 57 ply in monitoring and analyzing application memory
 58 while still preserves enough information for alloca-
 59 tion problem.

60 Monitoring framework and MSDF is currently be-
 61 ing developed at high performing computing center
 62 (HPCC) from HCMUT-VNU for our specific
 63 SuperNode-XP system. The following briefly summa-
 64 rizes the contribution of this paper:

- Propose Kafka based monitoring framework which can collect application metrics and perform online-analysis.
- Propose dynamic allocation in terms of memory and extract its efficiency boundary.
- Propose online analysis in memory monitoring to save storage space relating to memory allocation for applications in computing system.

73 The remain of this paper is structured as follow. Sec-
 74 tion 2 the related work that we have surveyed. Section
 75 3 introduces the architecture overview of our mon-
 76 itoring framework and shortly describes its compo-
 77 nents. Section 4 demonstrates how we apply MSDF
 78 in memory monitoring and allocation and Section 5
 79 will discuss the evaluation results. Section 6 further
 80 discusses about analyzing memory data and MSDF.
 81 Finally Section 7 summarizes all of our work and out-
 82 lines the future works.

83 RELATED WORK

84 There are immense of open source and commercial
 85 monitoring tools for computing system with differ-
 86 ent feature, different architecture. For wide range of
 87 used, *Zabbix* provides a stable solution for monitor
 88 host level metrics and hardware through Simple Net-
 89 work Management Protocol (SNMP); *Prometheus* is

a new generation tool with flexible user defined met-
 90 rics through custom exporter and a powerful build-in
 91 query language PromQL; *Nagios* highlights the con-
 92 tinuous real time network monitor with sensitive fail-
 93 ure detection and many other similar softwares can be
 94 listed as *Datadog*, *Icinga*, *SolarWinds*, etc. In general,
 95 these tools are all lack of streaming processing abil-
 96 ity when directly forward collected metrics into long
 97 term storage.

In recent years, *Apache Kafka* always stays as one of
 99 the best distributed data streaming platforms. Kafka
 100 provides both message queue and pub sub server
 101 at the same time with consistency through syn-
 102 chronize and fault tolerant through controller repli-
 103 cated mechanism and can be scaled horizontally by
 104 simply adding more instances to expand the traffic
 105 bandwidth. Besides, Kafka come along with a vast
 106 amount of processing framework including it own
 107 Application Programming Interface (API) in JAVA,
 108 PYTHON, SCALA, etc. which has the lowest over-
 109 head and mostly used in small scale problem; *Apache*
 110 *Spark*², a micro-batch processing framework with
 111 build in Kafka compatible Structured Streaming lib-
 112 rary, mostly used in large scale big data problem;
 113 *Apache Flink*³, a real streaming processing framework
 114 with build in Kafka library also used for large scale
 115 problem.

Analyzing data from system monitoring is not a new
 117 task, there have already been published several re-
 118 searches about this field in recent year. For instance,
 119 analyzing power usage of Central Processing Unit
 120 (CPU) by streaming linear regression using big data
 121 processing framework⁴; reconstructing application
 122 heap from monitoring tool trace file to detect mem-
 123 ory leak by offline-analysis^{5,6}; detecting software ag-
 124 ing based on memory leak investigation at software
 125 runtime⁷. But to the best of our knowledge, this is the
 126 first work which is applied streaming analysis in mon-
 127 itoring of application memory to save storage space
 128 with memory allocation problem case study.

MONITORING FRAMEWORK

As depicted in Figure 1, framework contains four
 131 different layers. Communication Layer positions in
 132 computing nodes and hardware devices, where met-
 133 rics is directly collected. Communication and An-
 134 alyzing Layer position in head nodes, i.e. manage-
 135 ment nodes, where contain central processing logic of
 136 monitoring framework in order to minimize comput-
 137 ing node overhead. Finally Storage Layer positions in
 138 storage nodes where long term data is accumulated for
 139 later purposes such as visualization or offline-analysis.

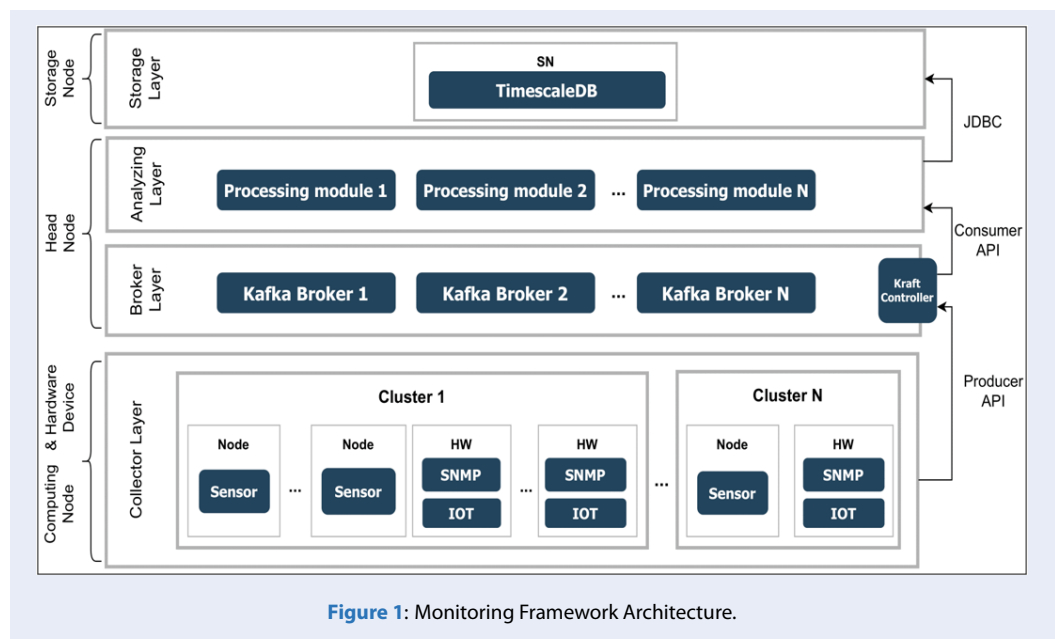


Figure 1: Monitoring Framework Architecture.

141 *Collector Layer.* The lowest layer responsible for collecting monitoring metrics. Since load from computing system mostly from scientific programs executed by system users, monitoring framework must be able to collect metrics from each user application. Sensor from monitoring framework, which is a daemon program installed at every computing node, is responsible for collecting resource usage corresponding to each user and application such as CPU utilization, memory usage, network traffic, etc. Besides, Sensor also considers to collect other hardware metrics such as network devices which exports through SNMP protocol; Uninterruptible Power Supply (UPS) devices which exports through SNMP; temperature, humidity sensor through Message Queuing Telemetry Transport (MQTT) protocol.

157 Theoretically, each cluster in computing system has different role, different behavior hence Sensor must be able to configure to collect only applicable metrics with appropriate interval. Additionally, sensors also can perform simple preprocessing step without costing too much computing resources if necessary. These data will eventually be pushed up to Communication Layer through Kafka Producer API.

165 *Communication Layer.* Because of a large number of applications executed at the same time, the communication between Analyzing Layer and data source become extremely complicated. Apache Kafka takes role as an intermediate data broker to provide a reliable transmission channel. Kafka implements the message queue and pub sub server where Collector

172 Layer take role as publishers push data to Kafka under specific topics and Analyzing Layer takes role as consumers fetch data from subscribed topics. Furthermore, Kafka provides a powerful streaming processing API for Analyzing Layer to conveniently perform online processing and analysis. Kafka is coordinated by Kraft Controller which acts as a gateway to receive all request from both producers and consumers and dispatch them to appropriate broker server.

181 Different type of metrics can be organized under different topic follow Linux hierarchical file system structure. For example `<organization name>/<cluster name>/<node id>` contains metrics of a specific computing node at host level such as CPU load, free disk, etc.; `<organization name>/<cluster name>/<node id>/<user id>` contains metrics of any applications executed under that user id in a specific computing node. Moreover, independent type of metrics can be put under different topic partitions to leverage parallel processing.

192 *Analyzing Layer.* This layer is responsible for handling streaming data from Communication Layer, each Processing Module will have different operations based on topic-name and partition-id. In relatively small system, Processing Module uses simple API which is JAVA Kafka Consumer to process streaming data and Java Database Connectivity (JDBC) driver to ingress into storage. But in large scale computing system with a vast number of clusters, users and applications where monitoring metrics from Collector Layer are considered big data, streaming processing with Apache Spark is seem to be more suitable.

204 Streaming processing ability can be utilized to create
 205 new metrics in more generalized level, for example
 206 aggregating each individual node health status to get
 207 the overall cluster health; or to immediately response
 208 when errors occur, for example auto shut down all
 209 infrastructure when detected power cut. Moreover,
 210 performing online-analysis in system monitoring can
 211 significantly reduce the total amount of data store and
 212 computing effort. We will further discuss how to ap-
 213 ply online-analysis in case of memory monitoring in
 214 Section 4.

215 *Storage Layer.* This layer contains database for long
 216 term storing where data can later be used in offline-
 217 analysis phrase, e.g applying machine learning algo-
 218 rithm to predict future trend to generate the next up-
 219 grade plan. Naturally, monitoring data is time series
 220 but because of Extract-Load-Transform (ETL) pro-
 221 cess in Analyzing Layer, storage should also support
 222 data warehouse along with time series database. We
 223 use TimescaleDB extension of PostgreSQL utilizes
 224 for time series data but still maintaining relational
 225 data model.

226 MSDF METHOD

227 Memory Allocation

228 Computing resources is shared between multiple pro-
 229 grams executed by different users. An effective alloca-
 230 tion strategy can result in more applications execute at
 231 the same time and hence boosting system efficiency.
 232 Allocation mechanism is usually classified as static
 233 and dynamic allocating. With naive static allocation,
 234 memory is given for application based on the maxi-
 235 mum usage which is collected from execution in the
 236 past. This strategy advantages can be listed as simple
 237 implementation, can prevent application from crash-
 238 ing and memory overflow, but is ineffective because of
 239 memory wasting. As in Figure 2, memory usage usu-
 240 ally will fluctuate within a certain range of value for
 241 a specific time before significant change happen; and in
 242 most of the time, memory usage is far from its maxi-
 243 mum value. Other applications may have more stable
 244 memory lines however allocating with static policy is
 245 not effective with this type of application and com-
 246 puting system in general. Thus, dynamic allocation
 247 should be utilized to only give application the most
 248 appropriate memory at specific period of time in ex-
 249 ecution to increase allocation efficiency.

250 Allocation from serverless computing⁸ is rated as one
 251 the best among dynamic allocation mechanism, in
 252 which application continuously requests the amount
 253 of memory needed and allocator will then give exactly
 254 that amount of memory. Nevertheless, applications

255 in computing system are mostly located inside con-
 256 tainer environment, although we can dynamically ad-
 257 just the environment resources, we can not contin-
 258 uously update its configuration due to technology lim-
 259 ited. Inspired from serverless computing, suppose the
 260 next execution is nearly similar to previous one, appli-
 261 cation runtime can be split into multiple continuous
 262 segments and each with different allocation of mem-
 263 ory. Memory given in each segment should be around
 264 the maximum value recorded in history correspond-
 265 ing to that segment. Recall from Figure 2, theoretic-
 266 ally data points in each segment should be stable and
 267 as close to the maximum value of that segment as pos-
 268 sible to maximize allocation efficiency. Thus, any free
 269 data points, i.e have not yet belonged to any segments
 270 before significant change happen, should be grouped
 271 into the same segment.

272 MSDF Approach

273 Allocation strategy from Section 4.1 only requires
 274 the maximum memory usage. With normal ap-
 275 proach, monitoring framework permanently stored
 276 every memory data of application collecting at dif-
 277 ferent timestamp, allocator queries and traverses
 278 through all of that data to compute the maximum
 279 value at each different segment. Based on sensor in-
 280 terval and application run time, each execution could
 281 end up thousand to million of records in database
 282 which will largely cost storage size and computing ef-
 283 fort. Instead with applying online processing in mon-
 284 itoring, when new data is coming, MSDF can calcu-
 285 late segment statistic value immediately and store only
 286 necessary data.

287 In particular, MSDF calculates the max, min, mean,
 288 and amount of data points of each segment. Suppose
 289 the new coming data at n^{th} offset is M_n , and current
 290 max, min are MAX and MIN , new values can be easily
 291 calculated with below formula:

$$MAX = Greater(M_n, MAX). \\ MIN = Lesser(M_n, MIN).$$

292 Calculating mean (\bar{M}_n) from previous mean $n-1$ value
 293 is a bit more complicated by following the formula be-
 294 low:

$$\bar{M}_n = \frac{\bar{M}_{n-1} * (n-1) + M_n}{n}$$

295 We defined $\epsilon \in (0,1]$ is the Threshold parameter and
 296 can be configured, a significant change is considered
 297 to happen when new coming data exceeds the given
 298 ϵ threshold:

$$\frac{|\bar{M}_{n-1} - M_n|}{\bar{M}_{n-1}} > \epsilon. \quad (1)$$

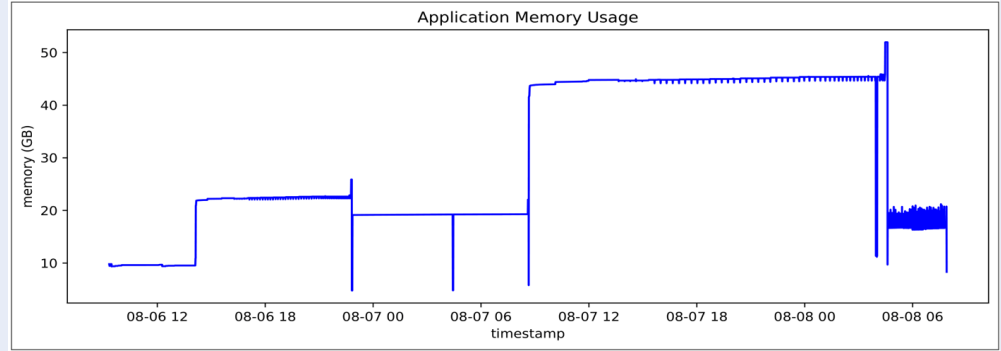


Figure 2: Efficiency and records stored in Storage Layer associate with different thresholds.

299 Then MSDF store the information related to that seg-
 300 ment into database, reset $n=0$, move to the next seg-
 301 ment and calculate new statistics value. By leveraging
 302 streaming processing technique, the final memory in-
 303 formation of application execution saved in database
 304 is only statistics values including mean, max, min and
 305 number of data points corresponding to each seg-
 306 ment.

307 Allocation Efficiency

308 We define a metric to estimate the efficiency of mem-
 309 ory allocation corresponding to each threshold ϵ .
 310 Suppose $f(t)$ is the memory usage function of appli-
 311 cation at time t , segment begins at t_a and ends at t_c ,
 312 and the total memory used by application in this seg-
 313 ment is defined as:

$$MemoryUsage = \int_{t_a}^{t_c} f(t) dt. \quad (2)$$

314 In fact, since the $f(t)$ function is unknown, we
 315 only have the set of discrete data points, where
 316 M_i indicated memory used at a specific time t which
 317 is collected from monitoring application. Suppose
 318 segment contains k data points with i indicated the
 319 sample interval of monitoring sensor. Integral (2) is
 320 now being calculated by discrete rectangle method:

$$\begin{aligned} MemoryUsage &= \sum_{i=1}^k M_i * \Delta t_i \\ \Leftrightarrow MemoryUsage &= I * \sum_{i=1}^k M_i. \end{aligned} \quad (3)$$

321 As mentioned before, the memory allocated to each
 322 segment is approximate to MAX value of that seg-
 323 ment. Applying (3), the expected memory consump-
 324 tion is:

$$\begin{aligned} SegmentEfficiency &= \frac{I * \sum_{i=1}^k M_i}{k * I * MAX} \\ \Leftrightarrow SegmentEfficiency &= \frac{\sum_{i=1}^k M_i}{k * MAX} * 1 \\ \Leftrightarrow SegmentEfficiency &= \frac{\bar{M}}{MAX}. \end{aligned} \quad (5)$$

With (5) is the allocation efficiency of each segment, 325
 suppose we have n segments, T_i is the i^{th} segment 326
 length and K_i is the number of data points in the i^{th} 327
 segment, the total efficiency corresponding to the de- 328
 fined threshold is: 329

$$\begin{aligned} TotalEfficiency &= \frac{\sum_{i=1}^n \frac{\bar{M}_i}{MAX} * T_i}{\sum_{i=1}^n T_i} \\ \Leftrightarrow TotalEfficiency &= \frac{\sum_{i=1}^n \bar{M}_i * K_i}{\sum_{i=1}^n K_i}. \end{aligned} \quad (6)$$

Efficiency Boundary

Efficiency boundary of each segment can guarantee 331
 the quality of whole allocation solution. Naturally, 332
 the upper bound reaches 1 and represented the most 333
 ideal scenario when application used exactly the same 334
 amount of given memory. And lower bound repre- 335
 sents the worst situation that allocation strategy can 336
 be encountered. Thus this subsection will mainly fo- 337
 cus on finding the lower bound of the solution from 338
 Section 4.1. 339

In each segment, suppose M_1 is the first data point 340
 of segment and $\bar{M}_1 = M_1$. Maximum value of seg- 341
 ment reaches its highest threshold when data points in 342
 segment progressively increase by a largest allowable 343
 value between any data points. From (1), we have: 344

$$\begin{aligned} \bar{M}_{n-1} - M_n &\geq -\epsilon \bar{M}_{n-1} \\ \Leftrightarrow M_n &\leq \bar{M}_{n-1} * (1 + \epsilon). \end{aligned} \quad (7)$$

Equal sign from (7) occurs in any data points, for all 345
 $k > 1$, mean value of the first k data points in segment 346
 can be calculated as below: 347

$$\begin{aligned} \bar{M}_k &= \frac{\bar{M}_{k-1} * (k-1) + \bar{M}_{k-1} * (1 + \epsilon)}{k} \\ \Leftrightarrow \bar{M}_k &= \bar{M}_{k-1} * (1 + \frac{\epsilon}{k}). \quad (9) \\ \square \\ (9) \Rightarrow \bar{M}_2 &= M_1 * (1 + \frac{\epsilon}{2}) \quad (10) \\ (9), (10) \Rightarrow \bar{M}_3 &< M_1 * (1 + \frac{\epsilon}{2})^{k-1}. \quad (12) \end{aligned}$$

348 In contrast, Minimum value of segment reaches it
 349 lowest threshold when data points in segment pro-
 350 gressively decrease by a largest allowable value be-
 351 tween any data points. From (1), we have:

$$\begin{aligned} \overline{M_{n-1}} - M_n &\leq \varepsilon \overline{M_{n-1}} \\ \Leftrightarrow M_n &\geq \overline{M_{n-1}} * (1 - \varepsilon). \end{aligned} \quad (9)$$

352 And similarly, equal sign from (9) occurs in any data
 353 points. The mean value of first k data points in seg-
 354 ment is:

$$\begin{aligned} \overline{M}_k &= \overline{M_{k-1}} * \left(1 - \frac{\varepsilon}{k}\right) \\ \Rightarrow M_k &> M_1 * \left(1 - \frac{\varepsilon}{k}\right)^{k-1}. \end{aligned} \quad (10)$$

355 In order to figure out the allocation lower bound, the
 356 situation when efficiency become worst must be first
 357 determined. There are three different cases of mem-
 358 ory consumption:

359 A: Memory is progressively increases within thresh-
 360 olds.

361 B: Memory is progressively decreases within thresh-
 362 olds.

363 C: Memory is randomly changes within thresholds.

364 From (7) and (8), suppose segment has k data points,
 365 the efficiency of case A is:

$$\begin{aligned} \frac{\overline{M}}{MAX} &= \frac{M_1 * \prod_{n=2}^k \left(1 + \frac{\varepsilon}{n}\right)}{M_1 * \prod_{n=2}^{k-1} \left(1 + \frac{\varepsilon}{n}\right) * (1 + \varepsilon)} \\ \Leftrightarrow \frac{\overline{M}}{MAX} &= \frac{1 + \frac{\varepsilon}{k}}{1 + \varepsilon}. \end{aligned} \quad (11)$$

366 From (9) and (10), suppose segment has k data points,
 367 the efficiency of case B is:

$$\begin{aligned} \frac{\overline{M}}{MAX} &= \frac{M_1 * \prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right)}{M_1} \\ \frac{\overline{M}}{MAX} &= \prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right). \end{aligned} \quad (12)$$

368 For all $k > 2$ and $\varepsilon \in (0,1]$, (11) $>$ (12) by using inves-
 369 tigating function approach. Therefore, the efficiency
 370 from case B is worse than case A.

$$371 \quad k = 2 \Rightarrow 1 + \frac{\varepsilon}{2} > \left(1 - \frac{\varepsilon}{2}\right) (1 + \varepsilon)$$

$$372 \quad k = 3 \Rightarrow 1 + \frac{\varepsilon}{3} > \left(1 - \frac{\varepsilon}{2}\right) \left(1 - \frac{\varepsilon}{3}\right) (1 + \varepsilon)$$

$$373 \quad k > 3 \Rightarrow \frac{1 + \frac{\varepsilon}{k}}{1 + \varepsilon} > \left(1 - \frac{\varepsilon}{2}\right) \left(1 - \frac{\varepsilon}{3}\right) \left(1 - \frac{\varepsilon}{4}\right) (1 + \varepsilon) \Rightarrow$$

$$374 \quad \frac{1 + \frac{\varepsilon}{k}}{1 + \varepsilon} > \prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right).$$

375 The efficiency from case C is equal to the case when all
 376 data points in case C is sorted in gradually decreasing
 377 order. In this case, we can consider each data point
 378 is changed α time compared to mean value of previ-
 379 ous data points. Additionally, the difference between

MAX and MIN in case C is smaller than the difference
 in case B. Hence, we have:

$$\left\{ \begin{array}{l} M_n = \overline{M_{n-1}} * (1 - \alpha_n) \\ \alpha_n < \varepsilon, \forall n \end{array} \right\} \quad (13)$$

From (13), following similar step from case B, the
 mean value of first k data points in segment is:

$$\overline{M}_k = M_1 * \prod_{n=2}^k \left(1 - \frac{\alpha_n}{n}\right). \quad (14)$$

From (13) and (14), suppose segment has k data
 points, the efficiency of case C is:

$$\begin{aligned} \frac{\overline{M}}{MAX} &= \frac{M_1 * \prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right)}{M_1} \\ \Leftrightarrow \frac{\overline{M}}{MAX} &= \prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right) \end{aligned} \quad (15)$$

The efficiency in case B is also worse then case C be-
 cause of (12), (13) and (15):

$$\begin{aligned} 1 - \frac{\alpha_n}{n} &> 1 - \frac{\varepsilon}{n} \quad \forall n \\ \Leftrightarrow \prod_{n=2}^k \left(1 - \frac{\alpha_n}{n}\right) &> \prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right). \end{aligned} \quad (16)$$

Thus it is confident to say the worst efficiency belongs
 to case C when memory consumption is progressively
 decreased. The boundary of segment efficiency is:

$$\prod_{n=2}^k \left(1 - \frac{\varepsilon}{n}\right) < Efficiency < 1$$

When $k \rightarrow +\infty$, lower bound will go toward 0 but in
 face, k is always a limited number and my depend on
 program type, program execution time, sensor col-
 lecting interval, etc. Table 1 showed the lower bound
 of different ε with different k. Decreasing ε can po-
 tentially lead to decrease of knumber in all segments,
 the smaller ε value, the more efficiency could be guar-
 teed.

RESULT AND EVALUATION

MSDF proposes a way to store monitoring memory
 data and retain only necessary information in order
 to save storage space. In case of memory allocation
 problem, ε value indicates the trade off between allo-
 cation efficiency and storage saving. We define some
 metrics to clarify MSDF efficiency and the trade off
 between these two factors with different ε value. Ef-
 ficiency score showed the efficiency of memory al-
 location could potentially achieve with segment in-
 formation corresponding to the ε value. Number of
 blocks indicates the disk block used by storage to save
 monitoring data, suppose the field size of all type of
 monitoring data in database is all equal to exactly on

Table 1: Efficiency lower bound of allocation strategy corresponding to different ϵ and different k .

ϵ	1	0.8	0.5	0.35	0.2	0.1
K = 1000	0	0	0.04	0.1	0.27	0.52
K = 10000	0	0	0.01	0.04	0.17	0.41

413 block. And finally, storage metrics compares the stor- 414
 414 age space of MSDF approach to normal approaches 415
 415 such as Zabbix or Prometheus.

416 As mentioned before, to our knowledge, there have 417
 417 not been any works which is similar to ours. MSDF 418
 418 is evaluated on four different computing programs 419
 419 from civil engineering research at our SuperNode- 420
 420 XP system with 10 seconds collector sensor interval. 421
 421 Each program used VASP library and executed in total 422
 422 total 196.6 hours on computing node with 48 cores, 96 423
 423 threads and 256 GB memory configuration. And normal 424
 424 approaches which store all monitoring data is set 425
 425 as the baseline for comparison purpose.

426 Efficient scores from Table 2 confirmed the validity of 427
 427 Table 1 where efficient scores are all greater than the 428
 428 lower bound in the same ϵ value. When the threshold 429
 429 value is decreased, the efficiency increased but number 430
 430 of blocks, i.e storage size also increased. Because 431
 431 lower threshold means more strict in allowing value 432
 432 change to be happen, and hence will be split to more 433
 433 segments which cost more storage size. But the fluctuation 434
 434 of data points in each segment will become stable 435
 435 thus increasing the overall efficiency score.

436 When ϵ goes toward 0, efficiency goes toward 1 and 437
 437 number of records reaches total raw records in which 438
 438 mean $\$=\$ \max \$=\$ \min \$=\$$ value. Based on the ef- 439
 439 ficiency score, storage saving and the statistics value 440
 440 at each segment, ϵ value can be reconfigured to find 441
 441 the best trade off between efficiency and storage saving. 442
 442 Different type of applications may yield more or 443
 443 less optimistic result, however with VASP programs 444
 444 above, MSDF is able to save 99% storage when allocation 445
 445 effectiveness reach more than 80%.

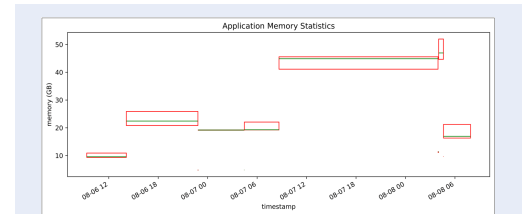
446 **DISCUSSION**

447 Allocation strategy from Section 4 suggest that at each 448
 448 segment, program should be allocated to the maximum 449
 449 memory usage of that segment. In fact, at each segment, 450
 450 resources must be allocated before, but the maximum value 451
 451 can not be found until reaching the end of that segment. 452
 452 Fortunately, applications in computing system usually belong 453
 453 to parameter-sweep class⁹, i.e program executes each time 454
 454 in the same behavior but with different input. Thus applica- 455
 455 tions in computing system can be assumed that memory or re- 456
 456 sources between its different executions do not vary 457
 457

458 much, so it is feasible to apply history information in- 459
 459 cluding segment and segment maximum to the next 460
 460 execution.

461 The core idea of MSDF is to group together any contin- 462
 462 uous and stable data points. MSDF accepts new 463
 463 coming data changing below certain threshold compared 464
 464 to previous data points. In case of applications which 465
 465 memory usage is gradually increased or decreased within 466
 466 the allowable threshold, MSDF eventually will have only 467
 467 one segment with low efficiency allocation. As a consequence, 468
 468 MSDF should not be applied in applications with resource 469
 469 usage gradually changed behavior. 470

471 Additionally, since the final data saved in storage of 472
 472 each application executions is only segment information, 473
 473 These value can be directly visualized as shown in Figure 3 474
 474 without being recomputed. The more closer between line 475
 475 and upper rectangle boundary indicated the more efficiency 476
 476 of allocation in the corresponding segment. Moreover, by 477
 477 visualizing different application executions and stacking 478
 478 these graphs together, system-operating questions such as 479
 479 whether these applications are able to executed simultane- 480
 480 ously can be easily answered. In general MSDF can be utilized 481
 481 to use in scheduling problem as well. 482



483 **Figure 3:** Statistics visualization corresponding to 484
 484 Figure 2. In each segment, the line represented 485
 485 mean value, rectangle represented min and max 486
 486 boundary. Lines without rectangle boundary rep- 487
 487 represents the mean=min=max situation. 488

489 **CONCLUSION**

489 In this paper, we first introduced our monitoring 490
 490 framework architecture and briefly detailed its com- 491
 491 ponents. To sum up, monitoring framework can 492
 492 collect metrics at application level, utilize Apache 493
 493

Table 2: Efficiency and records stored in storage associate with different thresholds.

Threshold	Efficiency	Number of Blocks	Storage Saving
1	0.04	60	99.91%
0.8	0.28	204	99.71%
0.5	0.67	400	99.43%
0.35	0.79	476	99.32%
0.2	0.85	660	99.06%
0.1	0.91	1136	98.39%
Normal Approaches	~1	70693	100.00%

489 Kafka as a data broker to organize system hierar-
 490 chy structure under different topic and also leverage
 491 Kafka streaming processing ability to perform online-
 492 analysis. As a use case, we demonstrated how to ap-
 493 ply online-analysis in monitoring memory for alloca-
 494 tion problem through MSDF. MSDF showed the trade
 495 off between allocation efficiency and storage saving
 496 based on the threshold ϵ value. In conclusion, apply-
 497 ing MSDF in monitoring and analyzing memory usage
 498 of computing application can save a huge storage
 499 capacity while still ensure allocation efficiency.

500 In future, since setting threshold in MSDF can affect
 501 both allocation efficiency and storage saving, we aim
 502 to fine tuning MSDF ϵ value in order to get the best
 503 trade off between the two factors, and in advanced
 504 providing MSDF ability to auto update that threshold
 505 value at application runtime. Moreover, also in mem-
 506 ory monitoring area, we planned to adapt MSDF to
 507 solve application scheduling problem as well.

508 **ACKNOWLEDGMENT**

509 This research was conducted within 58/20-
 510 DTDL.CN-DP Smart Village project sponsored
 511 by the Ministry of Science and Technology of
 512 Vietnam.

513 We gratefully acknowledge the valuable support and
 514 resources provided by HPC Lab at HCMUT, VNU-
 515 HCM.

516 **CONFLICT OF INTEREST**

517 The authors confirm that there is not any conflict of
 518 interest related to the content reported in this paper.

519 **AUTHORS CONTRIBUTION**

520 **La Quoc Nhut Huan:** first author, writing & editing,
 521 investigation, formal analysis, provide solution.

522 **Nguyen Manh Thin:** supervision, validation, func-
 523 tion testing, resources providing.

524 **Nguyen Quang Hung:** solution advising, reviewing,
 525 methodology.

Nguyen Le Duy Lai: solution advising, reviewing, 526
 methodology. 527

Thoai Nam: funding acquisition, supervision, con- 528
 ceptualization, instruction. 529

REFERENCES 530

- 531 Choenni S, Bakker R, Blok HE, de Laat R. Supporting tech- 532
 nologies for knowledge management. In: Knowledge Man- 533
 agement and Management Learning: Extending the Hor- 534
 izons of Knowledge-based Management. Springer; 2005. p. 89- 535
 112; Available from: https://doi.org/10.1007/0-387-25846-9_6.
- 536 Shoro AG, Soomro TR. Big data analysis: Ap spark perspective. 537
 Global Journal of Computer Science and Technology: Csoft- 538
 ware & Data Engineering. 2015;15(1):7-14; 539
- 540 Javed MH, Lu X, Panda DK. Characterization of big data stream 541
 processing pipeline: a case study using flink and kafka. In: Pro- 542
 ceedings of the Fourth IEEE/ACM International Conference on 543
 Big Data Computing, Applications and Technologies. 2017. p. 544
 1-10; Available from: <https://doi.org/10.1145/3148055.3148068>.
- 545 Beneventi F, Bartolini A, Cavazzoni C, Benini L. Continuous 546
 learning of hpc infrastructure models using big data analyt- 547
 ics and in-memory processing tools. In: Design, Automation & 548
 Test in Europe Conference & Exhibition (DATE), 2017. IEEE; 2017. 549
 p. 1038-1043; Available from: <https://doi.org/10.23919/DATE.2017.7927143>.
- 550 Wenginger M, Lengauer P, Mössenböck H. User-centered offline 551
 analysis of memory monitoring data. In: Proceedings of the 8th 552
 ACM/SPEC on International Conference on Performance Engi- 553
 neering. 2017. p. 357-360; Available from: <https://doi.org/10.1145/3030207.3030236>.
- 554 Wenginger M, Mössenböck H. User-defined classification and 555
 multi-level grouping of objects in memory monitoring. In: 556
 Proceedings of the 2018 ACM/SPEC International Conference 557
 on Performance Engineering. 2018. p. 115-126; Available from: 558
<https://doi.org/10.1145/3184407.3184412>.
- 559 Matias R, Costa BE, Macedo A. Monitoring memory-related soft- 560
 ware aging: An exploratory study. In: 2012 IEEE 23rd Interna- 561
 tional Symposium on Software Reliability Engineering Work- 562
 shops. IEEE; 2012. p. 247-252; Available from: <https://doi.org/10.1109/ISSREW.2012.90>.
- 563 Li Z, Guo L, Cheng J, Chen Q, He B, Guo M. The serverless com- 564
 puting survey: A technical primer for design architecture. ACM 565
 Computing Surveys (CSUR). 2022;54(10s):1-34; Available from: 566
<https://doi.org/10.1145/3508360>.
- 567 Chirigati F, Silva V, Ogasawara E, de Oliveira D, Dias J, Porto 568
 F, Valduriez P, Mattoso M. Evaluating parameter sweep work- 569
 flows in high performance computing. In: Proceedings of the 570
 1st ACM SIGMOD Workshop on Scalable Workflow Execution 571
 Engines and Technologies. 2012. p. 1-10; Available from: <https://doi.org/10.1145/2443416.2443418>. 572
 573
 574

MSDF: Định dạng thống kê dữ liệu cho bộ nhớ ứng dụng trong giám sát hệ thống

La Quốc Nhựt Huân^{1,2,3,*}, Nguyễn Mạnh Thìn^{1,3}, Nguyễn Lê Duy Lai^{1,3}, Nguyễn Quang Hùng^{1,3}, Thoại Nam^{1,2,3}



Use your smartphone to scan this QR code and download this article

TÓM TẮT

Hệ thống máy tính hiệu năng cao (HPC) hoặc hệ thống tính toán có sự khác biệt nhất định với hệ thống dịch vụ thông thường. Nhìn chung, hệ thống dịch vụ chỉ chạy một số ứng dụng cụ thể, ví dụ như máy chủ web hoặc máy chủ mail và phục vụ cùng lúc nhiều người dùng nhất có thể trong khi với hệ thống tính toán, người dùng trong hệ thống có quyền chạy các ứng dụng của riêng họ và hoàn toàn cô lập với người dùng khác. Kỹ thuật giám sát là chìa khóa để đảm bảo hiệu quả sử dụng hệ thống và sự hài lòng của người dùng, bằng cách kết hợp kỹ thuật giám sát cùng với phân tích dữ liệu, quản trị viên có thể giải quyết một số bài toán vận hành cụ thể như phân bổ tài nguyên, lập lịch ứng dụng, phát hiện bất thường, v.v. Khác với hệ thống dịch vụ trong khi các quản trị viên thường sẽ giám sát những thông tin tổng quát của hệ thống trong khi với hệ thống tính toán sẽ cần giám sát thông tin của từng ứng dụng khởi chạy bởi từng người dùng. Do hệ thống tính toán thường đồng thời thực thi rất nhiều ứng dụng nên việc giám sát bằng các phương pháp truyền thống sẽ tiêu tốn một lượng lớn dung lượng lưu trữ và khiến ta chi trả nhiều phí hơn nếu hệ thống được triển khai trên môi trường điện toán đám mây.

Bài viết này tập trung vào việc phân tích dữ liệu sử dụng bộ nhớ của chương trình tính toán nhằm giải quyết bài toán phân bổ tài nguyên cho lần khởi chạy tiếp theo của ứng dụng đó. Khác với các phương pháp truyền thống trong đó tất cả dữ liệu được giám sát thu thập sẽ được lưu trữ trong cơ sở dữ liệu trước khi phân tích, chúng tôi sử dụng các phương pháp phân tích trực tuyến trong đó mọi dữ liệu mới sẽ được thu thập, xử lý, lưu trữ trong bộ nhớ đệm để chuyển đổi thành thông tin hữu ích và chỉ cho phép dữ liệu cần thiết được ghi xuống đĩa cứng. Chúng tôi đề xuất Định Dạng Thống Kê Dữ Liệu Cho Bộ Nhớ (MSDF), một kỹ thuật xử lý trực tuyến được sử dụng trong giám sát bộ nhớ sử dụng của ứng dụng nhằm tiết kiệm dung lượng lưu trữ trong đĩa cứng trong khi vẫn lưu giữ đủ thông tin để giải quyết bài toán phân bổ tài nguyên cho ứng dụng. MSDF có thể giúp tiết kiệm hơn 95% dung lượng lưu trữ trong khi luôn đảm bảo hiệu quả phân bổ tài nguyên tùy thuộc vào tham số ϵ và MSDF có thể được mở rộng để giải quyết thêm nhiều bài toán vận hành khác hoặc tinh chỉnh để thích ứng trong việc giám sát và phân tích các thông số khác của ứng dụng.

Từ khóa: giám sát hệ thống, giám sát bộ nhớ, xử lý dòng dữ liệu, phân tích trực tuyến (online), phân bổ tài nguyên bộ nhớ

¹Phòng thí nghiệm tính toán hiệu năng cao, khoa Khoa học và Kỹ thuật máy tính (HPC Lab), Trường Đại học Bách Khoa Thành Phố Hồ Chí Minh (HCMUT), 268 Lý Thường Kiệt, Quận 10, Thành Phố Hồ Chí Minh, Việt Nam,

²Viện Khoa học và Công nghệ tiến tiến liên ngành (iST), Trường Đại học Bách Khoa Thành Phố Hồ Chí Minh (HCMUT), 268 Lý Thường Kiệt, Quận 10, Thành Phố Hồ Chí Minh, Việt Nam,

³Đại học Quốc gia Thành Phố Hồ Chí Minh (VNU-HCM), Phường Linh Trung, Thành Phố Thủ Đức, Thành Phố Hồ Chí Minh, Việt Nam.

Liên hệ

La Quốc Nhựt Huân, Phòng thí nghiệm tính toán hiệu năng cao, khoa Khoa học và Kỹ thuật máy tính (HPC Lab), Trường Đại học Bách Khoa Thành Phố Hồ Chí Minh (HCMUT), 268 Lý Thường Kiệt, Quận 10, Thành Phố Hồ Chí Minh, Việt Nam,

Viện Khoa học và Công nghệ tiến tiến liên ngành (iST), Trường Đại học Bách Khoa Thành Phố Hồ Chí Minh (HCMUT), 268 Lý Thường Kiệt, Quận 10, Thành Phố Hồ Chí Minh, Việt Nam,

Đại học Quốc gia Thành Phố Hồ Chí Minh (VNU-HCM), Phường Linh Trung, Thành Phố Thủ Đức, Thành Phố Hồ Chí Minh, Việt Nam.

Email: huan@hcmut.edu.vn.

Trích dẫn bài báo này: Huân L Q N, Thìn N M, Lai N L D, Hùng N Q, Nam T. **MSDF: Định dạng thống kê dữ liệu cho bộ nhớ ứng dụng trong giám sát hệ thống.** *Sci. Tech. Dev. J. - Eng. Tech.* 2024; ():1-1.