

UDC 681.234

A.V. DERGUNOV

Lobachevsky State University of Nizhny Novgorod

SPECIFICATION AND AUTOMATIC DETECTION OF PERFORMANCE PROBLEMS IN MESSAGE PASSING (MPI) APPLICATIONS

Traditional way to analyze performance of message passing (MPI) applications is via visualization of their execution traces. Several tools were developed to aid this activity; one example of such tool is Jumpshot. However, performance analysis with such tools is a complex task due to large traces and complex interactions between processes. In this paper a new way to analyze performance is proposed by automatic detection of performance problems in message passing applications. Performance problem is defined as a set of actions that inhibit good performance and is specified using tracing and analysis rules.

Keywords: performance problems specification, performance problems detection, message passing applications.

Introduction

Performance properties are very important for parallel applications, so several tools were developed to analyze performance. Most popular performance analysis tools for MPI applications perform visualization of execution traces. Example of such tool is Jumpshot [1] which performs visualization using timeline view. But the task of performance analysis using such tools is quite complex, because trace files usually consists of many events and interaction between events is very complex. Performance analysis also requires expert knowledge about MPI implementation.

This paper proposes a new way to do performance analysis which is based on expert methodology of performance problems descriptions.

1. Model of an MPI application

An MPI application [2] is defined as a set of communicating processes $PR = \{PR_1, \dots, PR_N\}$. Communication is performed using initiating actions in defined sequence:

$$PR_i \Rightarrow a_{i1}, \dots, a_{iM_i}; i = 1, \dots, N.$$

Every action is a call of a function defined by MPI standard $F = \{f_k\}; k = 1, \dots, K$. Each function f_k has input and output arguments:

$$FA_k^{IN} = (fa_{k1}^{IN}, \dots, fa_{kL_k}^{IN})$$

$$FA_k^{OUT} = (fa_{k1}^{OUT}, \dots, fa_{kL_k}^{OUT}).$$

Thus, every action is represented by the calling function and the values of input and output arguments:

$$a_{ij} = \langle f_k, FAval_k^{IN}, FAval_k^{OUT} \rangle;$$

$$FAval_k^{IN} = (av_{k1}^{IN}, \dots, av_{kL_k}^{IN});$$

$$FAval_k^{OUT} = (av_{k1}^{OUT}, \dots, av_{kL_k}^{OUT});$$

$$i = 1, \dots, N; j = 1, \dots, M_i; f_k \in F.$$

Every action a_{ij} has start time t_{ij}^a and duration d_{ij}^a . Loss of performance due to communication is defined as:

$$D^a = \sum_{i=1}^N \sum_{j=1}^{M_i} d_{ij}^a.$$

The task of performance improvement is defined as minimization of this value.

2. Model of a performance problem

Performance problem is defined as a set of actions that inhibit good performance, because the actions are not synchronized. Fig. 1 shows synchronous execution of send and corresponding receive actions which results in good performance. Fig. 2 and fig. 3 show cases when send and receive actions are not synchronized thus producing late sending and late receiving problems correspondingly.

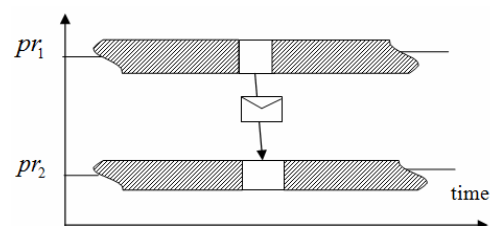


Fig. 1. Synchronous send and receive actions

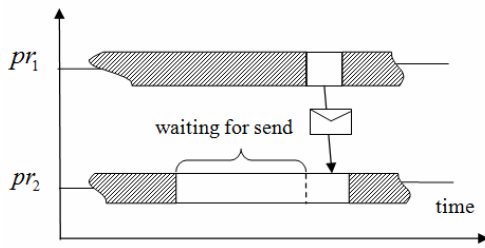


Fig. 2. Problem of late sending

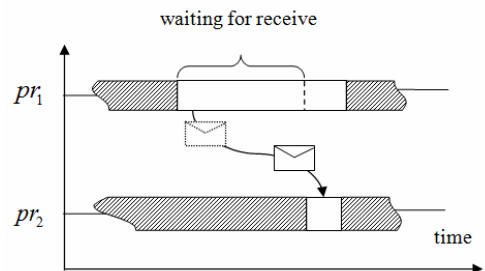


Fig. 3. Problem of late receiving

Formally *performance problem* is defined as:

$$pb = \left\langle \begin{matrix} pd, dur, TRRULES, \\ ANRULES, REC(A^{INFO}) \end{matrix} \right\rangle,$$

where *pd* – textual description of the problem; *dur* – duration of the problem; *TRRULES* – trace rules for actions that introduce the problem; *ANRULES* – analysis rules to recognize the problem in sequence of events in trace file; *REC* – recommendations to fix the problem; $A^{INFO} = \{\{f_i, t_i, d_i, pr_i, cs_i\}\}$ – description of actions that introduced the problem (where f_i – the function that was called, t_i – time when the function was called, d_i – duration of the function execution, pr_i – process that initiated the call, cs_i – the call site represented, for example, by source file name and line number in MPI application).

3. Performance Expert system

The task of automatic detection of performance problems in MPI applications turns out to be quite complex and it is hard to solve it using formal methods. Thus, expert methodology was proposed and implemented in Performance Expert system.

The workflow of Performance Expert system is illustrated in fig. 4.

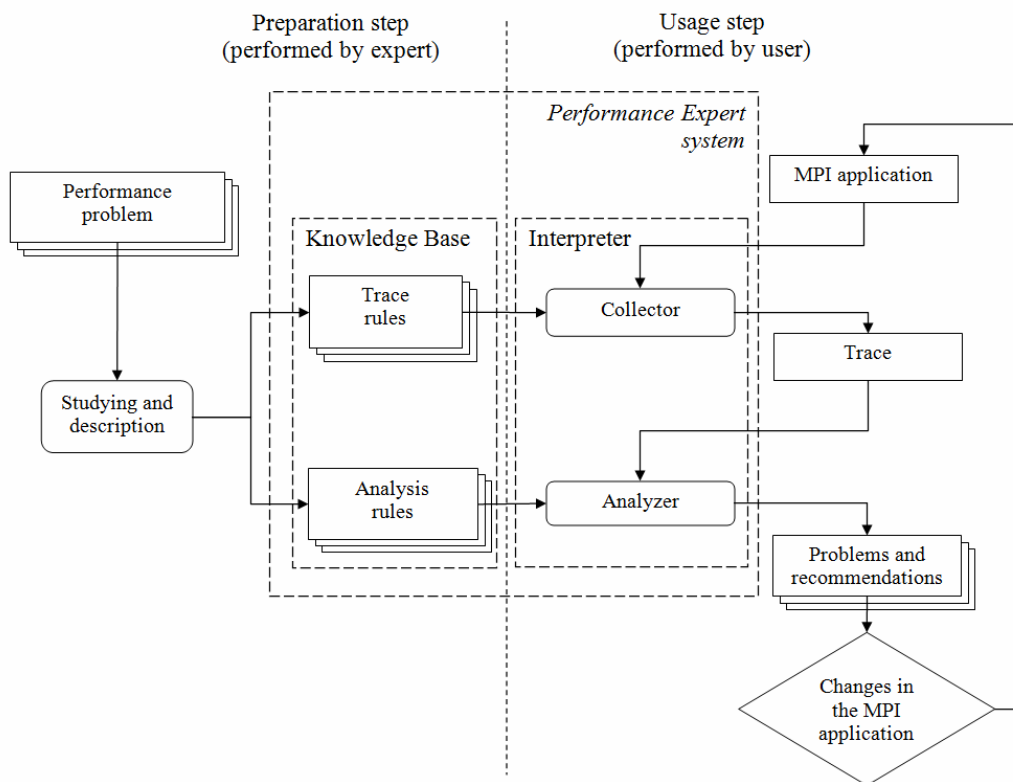


Fig. 4. Performance Expert system

The usage of the system supposes two steps:
Preparation step. This step is performed by expert in performance analysis of MPI applications. The task of the expert is to describe performance problems that she encounters using *trace and analysis rules*. The trace and analysis rules are described using languages de-

scribed later in this paper and they constitute *Knowledge Base* of the system which represents knowledge about all performance problems known so far. At the moment Knowledge Base of Performance Expert contains description of 10 typical performance problems of MPI applications [3].

Usage step. This step is performed by user. She executes her MPI application under *collector*. Collector uses trace rules to produce trace file containing events that correspond to actions executed by MPI application. Then the trace file is examined by *analyzer*. Analyzer detects performance problems in MPI application using analysis rules and generates recommendations for performance improvement. The user is able to follow the recommendations and implement the changes in her MPI application and then repeat usage step until good performance is achieved.

The data flow in Performance Expert system is illustrated in fig. 5:

1. MPI application executes a sequence of ac-

tions. Trace rules are used to describe which events should be generated and saved in the trace file.

2. Analysis of the trace file is performed in two steps:

a) composite events are constructed from simple events of the trace file. Composite event construction rules are used to perform this step;

b) performance problems are identified among the constructed composite events. Performance problem detection rules are used to perform this step. If a problem is identified, recommendations are produced which are related to the exact locations in source code of MPI application (this is done by extracting information about actions that produced the problem).

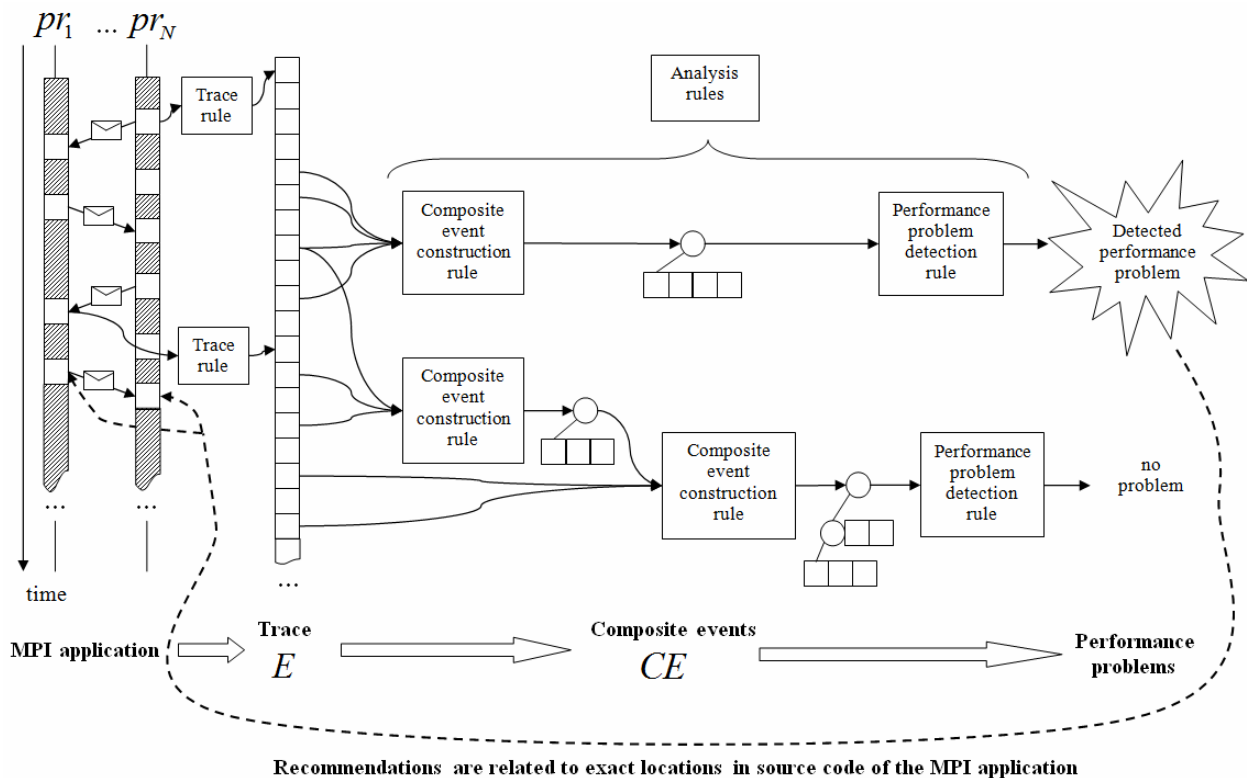


Fig. 5. Data flow in Performance Expert system

4. Tracing model

Simple events (or just events) are stored in trace file generated by collector. Each simple event is represented as:

$$e = \langle f, et, EPval, t, d, pr, cs \rangle,$$

where:

- $f \in F$ – function that was called;
- et – type of event which describes event parameters
- $EP = (ep_1, \dots, ep_K)$;
- $EPval = (v_1, \dots, v_K)$ – values of event parameters;
- t – event occurrence time;
- d – event duration;

pr – process which generated event;

cs – call site.

Trace rule is represented as:

$$a \xrightarrow{trule} e;$$

$$a = \langle f, FAval^{IN}, FAval^{OUT} \rangle;$$

$$e = \langle f, et, EPval, t, d, pr, cs \rangle;$$

$$trule = \langle f, et, EPtemp \rangle,$$

where:

$f \in F$ – function of the rule, a call to that function produces described event;

et – type of event which is generated as a result of the function call;

$EPtemp = \{\langle ep_i, Expr_i, kind_i \rangle\}$ – description of generated event parameters and the way to calculate their values, where:

ep_i – event parameter;

$Expr_i : FAval^{IN} \times FAval^{OUT} \rightarrow v_i$ – function to calculate value of event parameter (using argument values of the called function);

$kind_i \in \{in, out\}$ – kind of that parameter (in – calculated before function call; out – after the call).

An XML-based language was developed to declare trace rules. Trace rules are used by collector generator to produce wrappers of original functions. Function wrappers are used by PIN system [4] for dynamic instrumentation of MPI applications (PIN is a third-party component, a framework for dynamic instrumentation).

5. Analysis model

Composite event is represented as:

$$ce = \langle cet, CEPval, ME \rangle,$$

where:

cet – type of composite event which describes parameters $CEP = (cep_1, \dots, cep_K)$;

$CEPval = (cev_1, \dots, cev_K)$ – values of composite event parameters;

$ME = \{e_i\} \cup \{ce_j\}; i = 1, \dots, N; j = 1, \dots, M$ – set of simple and composite events that are members of this composite event.

Composite event construction rule is:

$$\left. \begin{array}{l} E \xrightarrow{ET} E^R \\ CE \xrightarrow{CET} CE^R \end{array} \right\} \xrightarrow{cerule} ce = \langle cet, CEPval, ME \rangle;$$

$$cerule = \langle ET, CET, \sigma, cet, CEPtemp, ET^S \rangle,$$

where:

$ET = \{et_i\}$ – set of simple event types to select a subset of relevant simple events:

$$E^R = \{\langle f_i, et_i, EPval_i, t_i, d_i, pr_i, cs_i \rangle\} \subseteq E;$$

$CET = \{cet_j\}$ – set of composite event types to select a subset of relevant composite events:

$$CE^R = \{\langle cet_j, CEPval_j, ME_j \rangle\} \subseteq CE.$$

Let us denote:

$$EPV_i = \langle f_i, EPval_i, t_i, d_i, pr_i, cs_i \rangle;$$

$$P = EPV_1 \times \dots \times EPV_N \times CEPval_1 \times \dots \times CEPval_M;$$

$\sigma : P \rightarrow \{1, 0\}$ – Boolean condition for constructing the composite event;

cet – type of the composite event to compose;

$CEPtemp = \{\langle cep_k, Expr_k \rangle\}$ – description of the composite event parameters and the way to calculate their values, where:

cep_k – parameter of the composite event;

$Expr_k : P \rightarrow cev_k$ – function to calculate parameter value;

If rule conditions are satisfied (i.e., subsets of relevant events ET^R and CET^R of the specified types exist and Boolean condition for their parameters σ is satisfied), then composite event ce of type cet is constructed, where:

$CEPval = (cev_1, \dots, cev_K)$ – parameter values;

$ME = E^R \cup CE^R$ – set of member events.

Rule parameter $ET^S \subseteq ET$ describes event types which are common for the constructed composite event and other event types.

Performance problem detection rules is:

$$ce \xrightarrow{pbrule} pb;$$

$$ce = \langle cet, CEPval, ME \rangle;$$

$$pb = \langle pd, dur, REC(A^{INFO}) \rangle;$$

$$pbrule = \langle cet, \phi, pd, RECtemp, L_{dur} \rangle,$$

where:

cet – type of the composite event which may represent performance problem;

$\phi : CEPval \rightarrow \{1, 0\}$ – Boolean condition of the problem occurrence;

pd – textual description of the problem;

$RECtemp : CEPval \rightarrow REC$ – recommendation

template (A^{INFO} is the description of actions that produced the problem and it is generated by extracting information from all simple events contained in the composite event recursively);

$L_{dur} : CEPval \rightarrow dur$ – function to calculate duration of the problem.

Languages were developed to declare composite event construction rules and performance problem detection rules. Performance Expert system uses CLIPS [5] expert system tool to perform the analysis, so internally these rules are converted into CLIPS rules.

6. Experiment

To investigate the implemented system an experiment was conducted to analyze and improve performance of MPI application which models heart activity [6]. Cells of heart comprise $N \times N$ lattice and each cell is

connected to the nearest neighbors. Each cell is described by differential equations and the application performs numerical integration.

The lattice is split into smaller parts $M \times M$ (where $M < N$) which are distributed among processes. Each process performs the following actions in cycle:

1. Calculate values in lattice points.
2. Exchange values on borders with neighboring processes using `MPI_Sendrecv`.

After numerical integration is done, the calculated data is sent to the main process using `MPI_Gather`.

Analysis of this application by Performance Expert system in automatic mode revealed the following performance problems:

- *Late sending*. Cumulative duration of such problems is 28,02% of total execution time.
- *Early receive for "many-to-one" operation*. Cumulative duration is 24,16%.

The second problem relates to application implementation details (main processor initiates call to `MPI_Gather` which is waiting until data is sent by other processes). So, it was decided to fix the first problem. Performance Expert system provided recommendation to use non-blocking receive operations, so the cycle in application was rewritten:

1. Calculate values in boarder lattice points.
2. Send values on borders to neighboring processes using `MPI_Send`.
3. Initiate receive of values on borders from neighboring processes using `MPI_Irecv` in non-blocking mode.
4. Calculate values in the rest lattice points.
5. Wait until receiving of values is finished using `MPI_Wait`.

Fig. 6 shows achieved performance improvements. The best improvement (1,62x) was achieved for 101 processes.

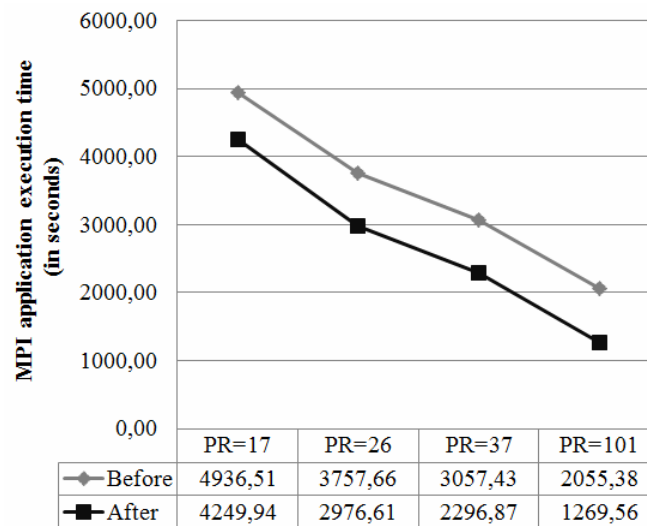


Fig. 6. Performance improvement results (PR denotes processes number)

Data was collected on cluster of Lobachevsky State University of Nizhny Novgorod (it consists of dual-core Intel Xeon 5150 2.66 GHz cores, 4 GB memory, Gigabit Ethernet, Windows Server 2008 x64, Microsoft implementation of MPI library).

Conclusion

A new way to analyze performance is proposed by automatic detection of performance problems in message passing applications. Experiment conducted on a real MPI application shows validity of this approach.

References

1. *Toward Scalable Performance Visualization with Jumpshot* [Текст] / O. Zaki [et al.] // *High-Performance*

Computing Applications. – 1999. – Vol. 13, № 2. – P. 277 – 288.

2. *MPI: A Message-Passing Interface Standard* [Электронный ресурс]. – Режим доступа: <http://www.mpi-forum.org/docs/mpi-11-html/mpi-report.html>. – 28.08.2011 г.

3. Дергунов, А.В. Автоматизация выявления причин потери производительности MPI программ на экзафлопсных и других больших суперкомпьютерах [Текст] / А.В. Дергунов // *Научный сервис в сети Интернет: экзафлопное будущее: труды Междунар. суперкомпьютерной конф.* – М.: Изд-во МГУ, 2011. – С. 491 – 496.

4. Pin: Building Customized Program Analysis Tools with Dynamic Instrumentation [Текст] / С.-К. Luk [et al.] // *Proceedings of the 2005 ACM SIGPLAN conference on Programming language design and implementation, Chicago, IL, 2005*. – P. 190 – 200.

5. Джарратано, Дж. Экспертные системы: принципы разработки и программирование / Дж. Джарратано, Г. Райли; пер. с англ. К. А. Птицына. – 4-е изд. – М.: Вильямс, 2007. – 1152 с.

6. Моделирование сердечной активности / Г.В. Осипов [и др.] // Суперкомпьютерные технологии в науке, образовании и промышленности. – М.: Изд-во Моск. ун-та, 2010. – С. 35 – 40.

Поступила в редакцию 17.02.2012

Рецензент: д-р техн. наук, ст. наук. співр. В.М. Опанасенко, Інститут кібернетики ім. В.М. Глушкова, Київ, Україна.

ОПИСАНИЕ И АВТОМАТИЧЕСКОЕ ВЫЯВЛЕНИЕ ПРОБЛЕМ ПРОИЗВОДИТЕЛЬНОСТИ В MPI ПРИЛОЖЕНИЯХ

А.В. Дергунов

Традиционно для анализа производительности MPI приложений используются программные средства для визуализации трассы их работы. Для выполнения этой задачи разработано несколько инструментов, примером является Jumpshot. Но анализ производительности с использованием таких инструментов является сложной задачей из-за больших размеров трасс и сложных взаимодействий процессов. В работе предложен новый подход к анализу производительности с использованием автоматического выявления проблем производительности MPI приложений. Под проблемой производительности понимается множество действий, которые негативно влияют на производительность программы. Проблемы производительности описываются с использованием правил трассировки и анализа.

Ключевые слова: описание и выявление проблем производительности, MPI приложения.

ОПИС І АВТОМАТИЧНЕ ВИЯВЛЕННЯ ПРОБЛЕМ ПРОДУКТИВНОСТІ В MPI ДОДАТКАХ

А.В. Дергунов

Традиційно для аналізу продуктивності MPI додатків використовуються програмні засоби для візуалізації траси їх роботи. Для виконання цього завдання було розроблено декілька інструментів, наприклад Jumpshot. Однак аналіз продуктивності з використанням цих інструментів є складним завданням через великий розмір трас та складних взаємодій процесів. У роботі запропонований новий підхід до аналізу продуктивності MPI додатків. Під проблемою продуктивності розуміється певні дії, котрі негативно впливають на продуктивність програми. Проблеми продуктивності описуються з використанням правил трасування та аналізу.

Ключові слова: опис і виявлення проблем продуктивності, MPI додатки

Дергунов Антон Владимирович – аспирант кафедры математического обеспечения ЭВМ Нижегородского государственного университета им. Н.И. Лобачевского, Нижний Новгород, Россия, e-mail: anton.dergunov@gmail.com.