

МОДЕЛИ ВЫЧИСЛИТЕЛЬНЫХ ПРОЦЕССОВ В КЛАСТЕРНЫХ СИСТЕМАХ С МНОГОЯДЕРНОЙ АРХИТЕКТУРОЙ

В работе рассматриваются подходы к разработке математических моделей, описывающих поведение кластерной системы с многоядерной архитектурой. Проанализированы особенности организации вычислений в таких системах, связанных с двухуровневой параллельной обработкой. Представлена математическая модель кластерной системы, основанная на теории последовательных взаимодействующих процессов.

In paper a problem of creation a model of cluster systems with multicore architecture is considered. An analyze of computing in such systems related with two level of parallelisms is presented. Mathematical model of cluster system based on theory of sequential communicating processes is presented.

1. Введение

Развитие высокопроизводительных компьютерных систем в настоящее время связано с распределенными (кластерными) системам. Последняя редакция рейтинга суперкомпьютеров TOP500 (31-я редакция, июнь 2008 года) включает 61 кластерную систем в первых ста позициях. Количество процессоров в этих кластерных системах варьируется в диапазоне от 1992 до 122400, максимальная производительность – в диапазоне от 18,81 до 1026 Tflop/s, пиковая производительность – в диапазоне от 24,58 до 1375,78 Tflop/s. На сегодня наиболее мощной является кластерная система Roadrunner компании IBM, установленная в Лос-Аламосской национальной лаборатории (США).

Современная кластерная система (КС) включает набор вычислительных узлов, соединенных специальной быстродействующей системой связей, основанной на сетевых технологиях. Узел КС обеспечивает обработку и хранение данных, система связей – передачу данных между узлами КС.

Система связи узлов (ССУ) является важнейшей составляющей КС, так как эффективность КС в значительной мере определяется скоростью обмена данными между узлами системы. Основу всех реализаций систем связей составляют принципы, заимствованные их сетевых технологий, поэтому в КС широко используются технологии типа Ethernet (Giga-bitEthernet), которые на сегодня обеспечивают пропускную способность сети до 10Gbit/sec. Наряду с Ethernet технологией в современных КС используются специализированные системы связей типа Myrinet, Infiniband, cLan, SCI, Quadrics. Анализ кластерных систем, предста-

вленных в первых 100 позициях 31-го рейтинга TOP500 показал, что наиболее используемой коммуникационной технологией является InfiniBand (42 позиция). За ней следует технология GigEthernet (10). В списке представлены также технологии Myrinet (7) и QsNet (2).

Рост производительности КС традиционно связывался с увеличением количества узлов системы. Однако наращивание числа узлов ограничено из-за повышенного требования к ССУ при увеличении узлов и нагрузке на систему связей. Имеется второй путь повышения производительности КС, основанный на увеличении вычислительной мощности непосредственно самого узла. Такой подход обеспечивается использованием в узлах мощных процессоров, также реализацией в узлах КС принципов параллельной обработки за счет применения многопроцессорных систем. Использование в узлах многопроцессорных систем позволяет перейти в КС к двухуровневой параллельной обработке. Первый уровень обеспечивается параллельной обработкой на множестве узлов КС, второй – на множестве процессоров внутри узла.

Появление многоядерных процессоров открывает перспективы увеличения мощности узлов за счет более эффективной реализации второго уровня параллельной обработки в КС. Простая замена в узле обычного процессора на двухядерный процессор позволяет удвоить вычислительную мощность КС без изменения ее структуры. Появление на рынке четырехядерных процессоров и анонсированное в ближайшем будущем появления процессоров с восьмью и более ядрами, открывает возможности резкого увеличения производительности КС при построении систем на основе мно-

гоядерных процессоров. В этой связи следует отметить, что 31-я редакция TOP500 характеризуется широким использованием в компьютерных системах четырехядерных процессоров нового поколения (258 систем).

Целью настоящей работы является анализ особенностей организации вычислительных процессов в кластерных системах с многоядерной архитектурой (КСМА) и построение математических моделей, позволяющих формализовать эти особенности, которые позволят в полной мере реализовать все преимущества структурной организации КСМА.

2. Модели вычислительных процессов в КСМА

Рассмотрим КСМА как систему, поведение которой можно описать с помощью алгебры процессов, предложенной в теории взаимодействующих последовательных процессов (CSP) Ч.Хоара [1]. Выбор теории Ч.Хоара объясняется тем, что она позволяет описать любой объект с помощью процессов. Модель, основанная на понятии процесс, в свою очередь может быть эффективно реализована с помощью современных языков параллельного программирования, базирующихся на механизме потоков (легких процессов).

При разработке модели будем исходить из того, что КСМА является сложной системой и имеет два уровня. Первый уровень – система узлов, второй уровень – система ядер внутри узла. Тогда модель КСМА имеет вид иерархической структуры, состоящей из моделей двух уровней (КСМА1 и КСМА2). Модель первого уровня описывает поведение системы, связанные с узлами, а модель второго уровня – с поведением системы внутри узла.

2.1 Модель КСМА1. Для модели первого уровня КСМА1 используем модель клиент-сервер, которая включает объект – сервер и множество объектов – клиентов.

В теории Т.Хоара процесс определяет поведение объекта и может быть описан с помощью ограниченного набора событий, выбранного в качестве его алфавита.

Набор событий, связанных с объектом сервер (алфавит сервера) C можно представить в виде

$$\alpha C = \{ \text{ввод}C, \text{вывод}C, \text{принять}SK_i, \text{передать}SK_i, \text{счет}C \}$$

где события: $\text{ввод}C$ – ввод данных на сервере, $\text{вывод}C$ – вывод результат на сервере, $\text{принять}SK_i$ – прием сервером данных от клиента, $\text{передать}SK_i$ – передача данных от сервера клиенту K_i , $\text{счет}C$ – счет на сервере.

Набор событий, связанных с i -м объектом клиентом (K_i) можно представить в виде

$$\alpha K_i = \{ \text{принять}K_iC, \text{передать}K_iC, \text{счет}K_i, \text{принять}K_iK_j, \text{передать}K_iK_j \}$$

где события: $\text{принять}K_iC$ – прием данных от сервера, $\text{передать}K_iC$ – передача данных серверу, $\text{счет}K_i$ – счет на клиенте, $\text{принять}K_iK_j$ – прием данных от клиента K_j , $\text{передать}K_iK_j$ – передача данных клиенту K_j .

Поведение объекта сервера используя алфавит αC , можно описать как процесс

$$C = (\text{ввод}C \rightarrow \text{передать}SK_i \rightarrow \text{счет}C \rightarrow \text{принять}SK_i \rightarrow \text{вывод}C)$$

Для множества объектов клиентов K_i их поведение (клиентский процесс) описывается следующим образом

$$K_i = (\text{принять}K_iC \rightarrow \text{передать}K_iK_j \rightarrow \text{счет}K_i \rightarrow \text{принять}K_iK_j \rightarrow \text{передать}K_iC)$$

Поведение КСМА1 можно описать как параллельную комбинацию поведения процессов C и K_i :

$$\text{КЛИЕНТЫ} = (K1 || K2 || \dots || KP), \text{СЕРВЕР} = (C) \\ \text{КСМА1} = (\text{СЕРВЕР} || \text{КЛИЕНТЫ})$$

Взаимодействие процессов в модели КСМА1

Взаимодействие узлов в КСМА1 осуществляется с помощью посылки сообщений. Будем рассматривать взаимодействие как событие. Взаимодействие состоит в передаче сообщения и является событием, требующим участия двух процессов.

Для описания взаимодействия процессов введем понятие объекта коммуникации (OK). Алфавит OK включает следующие события:

$$\alpha OK = \{ \text{запись}OK, \text{чтение}OK \}$$

Процесс, описывающий поведение объекта OK :

$$OK = ((\text{запись}OK \rightarrow \text{чтение}OK | \text{чтение}OK \rightarrow \text{запись}OK) \rightarrow OK)$$

Составляющей OK является спецификация C , который определяет тип и объем передаваемых данных ($OK.C$).

В модели КСМА1 взаимодействие объектов связано с событиями принять и передать . При этом взаимодействие должно быть связано с OK , куда один объект передаст данные, а другой объект произведет считывание данных. Дополним события принять и передать параметром, определяющим используемый

объект коммуникации: *передатьСК(ОК)*, *принятьСК(ОК)*. Процесс *ОК* можно рассматривать как сопроцесс [1].

2.2 Модель КСМА2. Узел будем рассматривать как систему, включающую набор объектов, связанных с действиями, выполняемыми в многоядерном процессоре. При организации процессов внутри узла используем модель взаимодействия хост-рабочий [2]. Хост объект отвечает за связь с остальными узлами КСМА, и, в частности, – с серверным узлом КСМА. Хост объект (*X*) принимает исходные данные от сервера КСМА, распределяет их в узле по объектам – рабочим, собирает результаты счета в узле и отправляет результат в серверный узел КСМА. Объект рабочий (*P*) принимает данные от хост объекта, выполняет счет и возвращает результат хосту.

Набор событий, связанных с объектом хост *i*-го узла (*X_i*), можно представить в виде

$$\alpha X_i = \{ \text{принять} X_i, \text{передать} X_i P_{ij}, \text{счет} X_i, \text{принять} X_i P_{ij}, \text{передать} X_i \}$$

где события: *принять**X_i* – прием данных от другого узла, *передать**X_iP_{ij}*, – передать данные в узле объекта рабочий *P_iP_{ij}*, *счет**X_i*, – счет, *принять**X_iP_{ij}* – принять данные в узле от объекта рабочий *P_{ij}*, *передать**X_i* – передать данные в другой узел.

Набор событий, связанных с объектом *j*-ий рабочий *i*-го узла (*P_{ij}*) можно представить в виде

$$\alpha P_{ij} = \{ \text{принять} P_{ij} X_i, \text{передать} P_{ij} X_i, \text{счет} P_{ij}, \text{принять} P_i P_{ik}, \text{передать} P_{ij} P_{ik} \}$$

где события: *принять**P_{ij}X_i* – прием данных от хост объекта, *передать**P_{ij}X_i* – передача данных хост объекту, *счет**P_{ij}* – счет в хост объекте, *принять**P_iP_{ik}* – прием данных от объекта рабочего *P_{ik}*, *передать**P_{ij}P_{ik}* – передача данных объекту рабочему *P_{ik}*.

Описание поведения объектов хост и рабочий будет зависеть от системы связей ядер в многоядерном процессоре, которая будет определять форму взаимодействия объектов хост и рабочий. Увеличение количества ядер ставит проблему организации связей между ядрами. Используемая в 2-4 ядрах шинная архитектура станет тормозом для эффективной связи ядер и неизбежным будет переход на распределенную архитектуру. Так, например, связывание ядер с помощью топологии типа решетка позволит отказаться централизованных ресурсов, а также решить проблему

надежности при выходе из строя ядер или связей.

Если связь ядер в многоядерном процессоре реализована через распределенную архитектуру с помощью, например, топологии звезда, где все ядра связаны через хост ядро, то поведение объектов внутри узла можно описать следующим образом.

Поведение хост объекта *X_i* (процесс хоста) используя алфавит αX_i , можно описать как $X_i = \{ \text{принять} X_i \rightarrow \text{передать} X_i P_{ij} \rightarrow \text{счет} X_i \rightarrow \text{принять} X_i P_{ij} \rightarrow \text{передать} X_i \}$

Поведение объекта рабочий *P_{ij}*: $P_{ij} = (\text{принять} P_{ij} X_i \rightarrow \text{счет} P_{ij}, \rightarrow \text{передать} P_{ij} X_i)$

Модель второго уровня для *i*-го клиентского узла можно описать как параллельную комбинацию поведения процессов *X_i* и *P_{ij}*:

$$\text{ХОСТ}_i = (X_i), \text{РАБОЧИЕ}_i = (P_i)$$

$$\text{УЗЕЛ}_i = (X_i \parallel \text{РАБОЧИЕ}_i)$$

Взаимодействие процессов в модели КСМА2.

Организация взаимодействия процессов внутри узла зависит от выбранной системы связей ядер. При прямой связи ядер (распределенная архитектура) используется модель, основанная на послышке сообщений, для связи через память (шинная архитектура) – модель, основанная на общих переменных.

При использовании послышки сообщений справедлива схема взаимодействия через объекты коммуникации (*ОК*), рассмотренная для модели первого уровня КСМА1.

Рассмотрим подробнее организацию взаимодействия процессов при использовании модели, основанной на общих переменных. Данная модель связана с решением двух задач: задачи взаимного исключения и задачи синхронизации [4]. Введем понятие объекта синхронизации (*ОС*), который будет использоваться процессами при взаимодействии через общую память. Выделим три вида объектов синхронизации. Первый вид *ОС(А)* связан с неделимыми (атомик) операциями на переменными. Второй вид *ОС(S)* реализует концепцию низкоуровневых средств синхронизации, используемых в примитивах типа семафор, мютекс, событие. Третий вид *ОС(M)* реализует концепцию мониторов.

Объект синхронизации типа А

$$\alpha A = \{ \text{чтение} A, \text{запись} A, \text{счет} A \}$$

где *чтение**A* – чтение неделимой переменной, *запись**A* – запись в неделимую переменную,

счетA – счет с использованием неделимой переменной.

Объект синхронизации типа S

$\alpha S = \{ \text{установить}S, \text{проверить}S, \text{сигнал}S \}$
 где *установитьS* – установка начального не-сигнального состояния *S*, *проверитьS* – проверка состояния *S*; блокирование процесса при занятом (несигнальном) состоянии – *S*, *сигналS* – установка *S* в сигнальное состояние (при создании *OC* принимает сигнальное состояние).

Объект синхронизации типа M

Набор событий, связанных с объектом-монитором *M* можно представить в виде

$$\alpha M = \{ \text{запись}M, \text{чтение}M, \text{счет}M, \text{ждать}M, \text{сигнал}M \}$$

где события: *записьM* – запись данных в монитор, *чтениеM* – чтение (копирование) данных из монитора, *счетM* – счет в мониторе, *ждатьM* – ожидание сигнала, *сигналM* – сигнал.

Использование объектов синхронизации в задачах взаимного исключения и синхронизации.

Объект A. Обращение процессов к общему ресурсу связано со следующими событиями в объектах *X* и *P*:

$$X: (\text{запись}A), \quad P: (\text{чтение}A)$$

Объект S. Обращение процессов к общему ресурсу связано со следующими событиями в объектах *X* и *P*:

$$X: (\text{проверить}S \rightarrow \text{счет}X \rightarrow \text{сигнал}S) \\ P: (\text{проверить}S \rightarrow \text{счет}P \rightarrow \text{сигнал}S)$$

Синхронизация связано со следующими событиями (хост посылает сигнал рабочему):

$$X: (\text{сигнал}S), \quad P: (\text{проверить}S)$$

Взаимное исключение. Объект M. Обращение процессов к общему ресурсу связано со следующими событиями в описании объектов *X* и *P*:

$$X: (\text{запись}M, \text{чтение}M, \text{счет}M, \text{ввод}M, \text{вывод}M) \\ P: (\text{чтение}M, \text{запись}M, \text{счет}M)$$

Синхронизация связано со следующими событиями в описании объектов *X* и *P*:

$$X: (\text{ждать}M, \text{сигнал}M), P: (\text{сигнал}M, \text{ждать}M)$$

3. Применение моделей

Предложенные модели позволяют применить математический аппарат теории алгебры процессов Ч.Хоара для анализа вычислительных процессов в КСМА.

Рассмотрим анализ взаимодействия процессов на возможность возникновения тупиковой ситуации при синхронизации двух процессов через объект синхронизации типа *OC*. Пусть *P* и *Q* – процессы, *D* – объект синхронизации типа *OC(S)*. Процессы синхронизируются по событию (ввод данных) в процессе *P* (процесс *Q* ждет сигнала от процесса *P*), после чего процессы начинают счет.

Алфавиты процессов будут иметь следующий вид:

$$\alpha P = \{ \text{ввод}, \text{счет}, \text{установить}D, \text{сигнал}D \} \\ \alpha Q = \{ \text{счет}, \text{проверить}D, \text{установить}D \} \\ \alpha D = \{ \text{проверить}D, \text{установить}D, \text{сигнал}D \}$$

Поведение процессов:

$$P = (\text{установить}D \rightarrow \text{ввод} \rightarrow \text{сигнал}D \rightarrow \text{счет}) \\ Q = (\text{проверить}D \rightarrow \text{счет})$$

Процесс *D* следует рассматривать как подчиненный процесс по отношению к процессам *P* и *Q*. При этом $\alpha D \subseteq (\alpha P \cup \alpha Q)$. В комбинации $(P \parallel Q \parallel D)$ каждое действие *D* может произойти, если это позволяют процессы *P* и *Q*.

Тупиковая ситуация при синхронизации процессов *P* и *Q* может возникнуть в двух случаях: а) объект *D* не установлен изначально в несигнальное состояние; б) процесс *P* не посылает сигнал о событии, которое ожидает процесс *Q*. Эти действия в процессах связаны с событиями *проверитьD*, *сигналD*, *установитьD*. Для борьбы с тупиком необходимо выполнить анализ алфавитов процессов и поведения процессов. Необходимое условие отсутствия тупика – наличие этих событий в алфавитах процессов:

$$\{ \text{проверить}D, \text{сигнал}D, \text{установить}D \} \subseteq (\alpha P \cup \alpha Q)$$

Достаточные условия – наличие события *сигналD* в алфавите процесса *P* и события *проверитьD* в алфавите процесса *Q*:

$$\{ \text{сигнал}D \} \subseteq (\alpha P), \quad \{ \text{проверить}D \} \subseteq (\alpha Q)$$

Для процессов – необходимо присутствие события *сигналD*, события *проверитьD* и *установитьD* при описании поведения процессов *P*, *Q* и *D* соответственно:

$$(\text{проверить}D \rightarrow Q) \parallel (\text{сигнал}D \rightarrow P) \parallel (\text{установить}D \rightarrow D)$$

Выводы

Предложены математические модели, основанные на теории взаимодействующих процессов Ч.Хоара, позволяющие описать орга-

низацию вычислений в кластерных системах с многоядерной архитектурой. Модели отображают двухуровневую организацию параллельных процессов в КСМА, а также систему взаимодействия процессов на основе посылки сообщений или общих переменных. Модель может быть использована для анализа парал-

лельных процессов в КСМА, оптимизации взаимодействия процессов, связанной с передачей данных, а также выявления тупиковых ситуаций, возникающих при взаимодействии процессов в КСМА.

Список литературы

1. Хоар Ч. Взаимодействующие последовательные процессы: Пер с англ. – М.: Мир, 1989.
2. El - Rewini H, Lewis T., Distributed and Parallel Computing. – Manning Pub. Co., NY, 1998.
3. Burns A., Welling A. Real-Time Systems and Programming Languages, Third Edition, Addison-Wesley, NY, 2001.
4. Жуков І., Корочкін О. Паралельні та розподілені обчислення. – Корнійчук, К., 2006.