

МОДЕЛИ ВЫЧИСЛЕНИЙ ДЛЯ МАСШТАБИРУЕМЫХ КЛАСТЕРНЫХ СИСТЕМ

В работе рассматривается подход к разработке программных компонент для масштабируемых кластерных систем, основанный на использовании моделей вычислений. Предложены модели вычислений, основанные на применении теории последовательных взаимодействующих процессов Ч.Хоара. Показана возможность использования моделей для повышения качества программ и сокращения времени их разработки.

In paper the problem of creation of effective program components for scalable cluster systems based on models of computation is considered. Models of computation based on C.Noare theory of communicating sequential processes are presented. Use models of computation for improving software and reduction design time is showed.

Кластерная архитектура сегодня рассматривается как одна из основных при построении высокопроизводительных компьютерных систем, оптимальных по критерию производительность/стоимость и включает набор вычислительных узлов, соединенных быстродействующей системой связей узлов (ССУ). Узлы кластерной системы (КС) обеспечивают обработку и хранение данных, ССУ обеспечивает быструю передачу больших объемов данных между узлами КС.

Современные КС характеризуются увеличением количества узлов, а также использованием в узлах многоядерных процессоров [1]. Кластерная система с многоядерной архитектурой (КСМА) представляет собой сложную систему, разработка программного обеспечения для которой представляет сложную задачу, связанную с организацией вычислений на двух уровнях: между узлами и внутри узлов. При этом необходимо обеспечить с помощью механизма процессов (поток) поведение каждого узла, а также обеспечить оптимальное по времени и объемам передаваемых данных взаимодействие процессов.

Современная практика разработки приложений для сложных компьютерных систем основывается на использовании моделей вычислений, которые выступают в качестве интерфейса между архитектурой системы и программой, реализующей выбранную модель параллельного (распределенного) программирования (рис.1). Назначение модели вычислений:

- описать архитектуру КС;
- описать поведение процессов;

- описать взаимодействие процессов.

На основе модели выполняются анализ корректности поведения процессов и их взаимодействия еще на этапе проектирования программного продукта. При этом математический аппарат модели позволяет выявить возможные причины возникновения тупиковых ситуаций. Кроме того, удачно выбранная модель может упростить процесс дальнейшего создания написания программы.



Рис.1. Модель вычислений

Кластерные системы являются хорошо масштабируемыми системами с точки зрения увеличения производительности. Кластерная архитектура допускает практически неограниченное наращивание числа узлов. Поэтому проблема разработки программного обеспечения для масштабируемых кластерных систем (МКС) является важной задачей и требует разработки методов и средств, позволяющих повысить качество приложений для МКС.

Анализ последних исследований и публикаций

Рассматриваемая тематика широко представлена в публикациях. В основном, они сводятся к выбору и обоснованию подхода к построению модели параллельных вычислений [2, 4, 5, 6, 9].

Модели вычислений могут базироваться на сетях Петри [9], теории конечных автоматов [5], CSP-теории [2, 3, 6].

Большое число разнообразных моделей вычислений не гарантирует их адекватность множеству практических задач, с которыми приходится сталкиваться при разработке распределенных приложений. Причина заключается в противоречивости требований к этим моделям и важности представления тех или иных особенностей функционирования вычислительных систем.

Нерешенные проблемы

В работе [2] предложено использование теории Ч.Хоара для построения моделей вычислений для кластерных систем. Показано, что математический аппарат CSP теории позволяет описать архитектуру кластерной системы, а также поведение и взаимодействие процессов. Однако в работе не рассмотрены вопросы построения, реализации и использования моделей вычислений на основе CSP теории для масштабируемых систем. Учитывая, что масштабируемость является важной составляющей компьютерной системы, решение проблем, связанных с созданием программного обеспечения для МКС, является актуальной задачей.

Постановка задачи

В работе [4] выделены две особенности распределенных вычислений в МКС, которые определяют требования к их разработке. Это *недетерминизм*, означающий произвольные последовательности формирования компонентов структурированных данных и *неоднородность*, которая связана с произвольной структурированностью данных, которыми обмениваются взаимодействующие процессы.

Недетерминизм связан с обеспечением однозначности результатов вычислений, при которой одному входному набору данных соответствует один выходной набор результатов. Неоднородность сообщений связана с тем, что процессы могут принимать (переда-

вать) данные произвольного размера, зависящего от количества узлов в системе. Выбор средств для построения моделей вычислений должен обеспечить решение указанных проблем.

Целью настоящей работы является развитие подхода, предложенного в работе [2], для создания и использования моделей вычислений при разработке программных компонент для масштабируемых кластерных систем с многоядерной архитектурой.

Модели вычислений

При разработке масштабированных распределенных программ необходимо обеспечить возможность статического (динамического) создания (порождения) множества процессов, описание их поведения, запуск на выполнение и организацию взаимодействия.

Существуют различные средства описания процессов в современных языках и библиотеках параллельного программирования [7]. Языки программирования Java, Ада позволяют описать процесс с помощью специальных программных модулей (классов). В языке С# и библиотеке Win32 для этого используются потоковые функции. OpenMP позволяет выделять параллельные участки программы с помощью операторов `parallel/end parallel`. Библиотека MPI создает параллельные процессы как копии основной программы.

Программа для масштабируемой архитектуры должна базироваться на возможности создания (как статически, так и динамически) любого множества процессов. Это возможно только с помощью тиражирования базового процесса, параметризация которого позволит описать все различия между создаваемыми копиями этого процесса. В языках, поддерживающих объектно-ориентированную парадигму, это можно сделать с помощью классов и конструкторов, в языке Ада для этого используется специальный тип с дискриминантом, в MPI – идентификатор задачи, автоматически порождаемый при создании копии процесса.

Организация взаимодействия процессов включает синхронизацию процессов и обмен данными [6]. Модели взаимодействия при этом базируются на общих переменных или на передаче сообщений.

При использовании общих переменных возникают задачи взаимного исключения и

синхронизации, решение которых надо обеспечить для произвольного числа процессов. Для задачи взаимного исключения это не представляет сложности, так как примитивы типа семафоры, мютексы, критические секции, также мониторы обеспечивают контроль общих ресурсов для любого числа процессов, которые их используют. Проблемы могут возникнуть при решении задачи синхронизации, так как здесь имеют место сложные формы синхронизации вида: *процесс* → *группа процессов*, *группа процессов* → *процесс*, *группа процессов* → *группа процессов*. При использовании низкоуровневых примитивов потребуются динамическое создание и использование объектов синхронизации. Эта проблема решается путем использования мониторных объектов синхронизации, которые позволяют эффективно реализовать указанные формы групповой синхронизации.

При использовании механизма отправки сообщений в МКС основная проблема связана с формированием объема передаваемых данных, который может меняться в зависимости от количества динамически создаваемых процессов. Здесь важной является возможность организация обмена данными произвольной размерности. В MPI это обеспечивается функциями `MPI_Gatherv`, `MPI_Scatterv`, в языке Ада – использованием типов данных с уточняемыми границами, что позволяет описать вход задачи для передачи любого объема данных.

Кроме того, эффективность взаимодействия увеличивается, если язык (библиотека) поддерживают коллективные формы взаимодействия, как при передаче сообщений, так и при синхронизации. В этом плане представляют интерес коллективные операции MPI, которые реализованы в виде операций групповой рассылки (сборки) данных `MPI_Bcast`, `MPI_Gather`, `MPI_Scatter` и операций групповой синхронизации `MPI_Barrier`.

Таким образом, можно сформулировать основные требования к моделям вычислений для масштабируемых систем: поддержка недетерминизма для неограниченного числа динамически порождаемых процессов, взаимодействующих путем отправки сообщений произвольного объема и через групповую синхронизацию.

Рассмотрим возможности теории Ч.Хоара с точки зрения реализации сформулированных требований. Теория Ч.Хоара позволяет описать поведение последовательных недетерминированных процессов с произвольным независимым чередованием процессов. Это позволяет описать особенности работы МКС. Поддерживается динамическое порождение процессов без ограничения их количества. Обеспечивается синхронная однопользовательная без буферизации отправка сообщений. При этом рассматривается одна форма, когда взаимодействуют только две задачи, то есть не используются групповые операции обмена. Не рассматривает общие переменные как средства взаимодействия процессов, тем более для групповой синхронизации и взаимного исключения.

Таким образом, теория Ч.Хоара может быть использована для построения моделей вычислений для МКС при ее расширении (доработке) в плане реализации двунаправленных механизмов отправки сообщений произвольной размерности, а также монитороподобных средств для взаимного исключения и синхронизации при описании группового взаимодействия процессов.

Модели вычислений для МКС

Рассмотрим МКС как сложную динамическую систему, структура которой может меняться как путем изменения числа узлов и связей между узлами. В общем случае, любая кластерная система является полностью связанной, так как любые два узла связаны между собой и допускают прямую передачу данных между собой. В этом случае масштабирование реализуется простым добавлением новых узлов. Если архитектура МКС построена путем объединения узлов в группы (подкластеры), то взаимодействие осуществляется между подкластерами и внутри каждого подкластера. При этом возможно использование различных коммуникационных систем для взаимодействия между подкластерами и внутри каждого подкластера.

Рассмотрим использование теории Ч.Хоара для построения моделей вычислений для масштабируемых кластерных систем с многоядерной архитектурой, в узлах которых используются многоядерные процессоры. Масштабируемость КС на структурном уровне связана с возможностью изменения количества узлов и количества ядер

в многоядерных процессорах в узле системы. Теория Ч.Хоара позволяет поведение любой системы рассматривать в терминах объект, абстрагируясь при этом от аппаратной или программной сущности системы построении модели вычислений. Это позволяет упростить переход от описания аппаратной части к программной реализации модели.

При разработке модели вычислений будем исходить из того, что она должна быть двухуровневой и соответствовать двум уровням структуры МКС. Модель *МКС1* описывает поведение системы, связанное с взаимодействием узлов, а модель *МКС2* описывает поведением системы внутри каждого узла.

Модель первого уровня МКС1. Для модели *МКС1* используем модель клиент-сервер, которая включает один объект клиент и P объектов серверов. Параметр P может меняться, тем самым обеспечивая масштабируемость на множестве узлов.

Набор событий, связанных с объектом клиент K можно представить в виде:

$$\alpha K = \{ K.ввод, K.вывод, K.принять.C, K.передать.C \},$$

где события: $K.ввод$ – ввод данных на клиенте, $K.вывод$ – вывод результат на клиенте, $K.принять.C$ – групповой прием данных от всех P серверов C_i , ($i = 1..P$), $K.передать.C$ – групповая передача данных всем серверам C_i . Масштабируемость здесь связана с групповыми передачами и приемами данных между клиентом и всеми серверами. Для реализации групповых операций $K.принять.C$ и $K.передать.C$ их необходимо параметризовать через параметр P . Кроме того, эти операции должны позволять либо передачу любого объема данных, либо допускать настройку на конкретный объем передаваемых данных. В MPI при передаче данных используются функции `Send/Receive`, в которых указывается объем передаваемых данных. В языке Ада - типы с уточняемыми границами (`type Vector is array(integer range <>) of float`), позволяющие передавать объемы данных любой размерности [8].

Набор событий, связанных с i -м объектом сервер (C_i , $i = 1..P$) можно представить в виде

$$\alpha C_i = \{ C_i.принять.K, C_i.передать.K, C_i.счет \}$$

где события: $C_i.принять.K$ – прием данных от клиента, $C_i.передать.K$ - передача данных

клиенту, $C_i.счет$ - счет на сервере. Масштабируемость на серверной стороне не фиксируется.

Поведение объекта клиент, используя алфавит αK , можно описать как процесс

$$K = (K.ввод \rightarrow K.передать.C \rightarrow K.принять.C \rightarrow K.вывод).$$

Для множества объектов клиентов C_i их поведение описывается следующим процессом

$$C_i = (C_i.принять.K \rightarrow C_i.счет \rightarrow C_i.передать.K).$$

Модель *МКС1* можно описать как параллельную комбинацию поведения процессов K и C_i :

$$\begin{aligned} \text{СЕРВЕРЫ} &= (C_1 \parallel C_2 \parallel \dots \parallel C_p), \\ \text{КЛИЕНТ} &= (K), \end{aligned}$$

$$\text{МКС1} = (\text{КЛИЕНТ} \parallel \text{СЕРВЕРЫ}).$$

Взаимодействие процессов в модели МКС1. Взаимодействие процессов в модели *МКС1* осуществляется на основе посылки сообщений. Для описания взаимодействия процессов используем объект коммуникации (OK), предложенный в работе [2]. Взаимодействие через OK реализуется как обращение к объекту (вызов OK и прием вызова OK):

$$\alpha OK = \{ OK.передать, OK.принять \}.$$

Процесс, который описывает поведение объекта коммуникации, задается рекурсивным соотношением:

$$OK = ((OK.принять \mid OK.передать) \rightarrow OK).$$

OK является параметризованным объектом $OK(D, N, B, \text{Б})$, где параметр D определяет тип и объем передаваемых данных, параметр N - направление передаваемых данных, параметр B - синхронное или асинхронное взаимодействие, параметр Б – буферизированное или не буферизованное взаимодействие. Для реализации масштабированности параметр D должен допускать возможность указания объема передаваемых данных ($D.X$), зависящего от P .

В теории Ч.Хоара поведение объекта OK описывается процессом, который является подчиненным для процессов, которые через него взаимодействуют. Корректность такого взаимодействия определяется следующими правилами:

Правило 1 определяет необходимое условие взаимодействия процессов через OK относительно алфавитов процессов: “Если A и

B процессы, Z - объект коммуникации, то необходимым условием взаимодействия A и B через Z есть выполнение соотношения $\alpha Z \subseteq (\alpha A \cup \alpha B)$ “.

Правило 2 определяет необходимое условие взаимодействия процессов через OK относительно поведения процессов: “Если A и B процессы, Z - объект коммуникации, то необходимым условием взаимодействия A и B через Z есть следующие действия в процессах $Z = ((Z.принять \mid Z.передать) \rightarrow Z)$, $A = ((Z.передать) \rightarrow A)$, $B = ((Z.принять) \rightarrow B)$, $(Z \parallel (A \parallel B))$. При этом $A.Z.D.X = B.Z.D.X$ ”.

Реализация объекта коммуникации для взаимодействия узлов МКС можно выполнить различными способами. С учетом сетевой технологии, которая определяет аппаратное взаимодействие узлов, это могут быть сокет, а также удаленные методы [6]. При использовании удаленных методов для операции умножения вектора на матрицу $A = B * MC$ в МКС, включающей P серверов, каждый сервер содержит удаленную процедуру $Счет()$, реализующую вычисления части результата $A_n = B * MC_n$, где $N = N/P$.

Описание удаленной процедуры на стороне сервера:

```
type Vector is array(integer
                        range <>) of
float;
type Matrix is array(integer
                       range <>) of Vector(1..N);
procedure Счет(B: in Vector; MC: in
               Matrix; A: out
               Vector);
```

Удаленный вызов на стороне клиента для i -го сервера:
`Сервер.Счет(B(1..N), MC(h*(i-1)+1..i*N), A(h*(i-1)+1..i*N));`

Модель второго уровня МКС2. Для модели $MKS2$, описывающей поведение узла, используем модель мастер-рабочий. Объект мастер отвечает за связь с клиентским узлом МКС и здесь связываются первый и второй уровни системы. Объект мастер (B) принимает исходные данные от клиента, распределяет их в узле по всем K объектам рабочий, собирает результаты счета в узле и отправляет результат клиенту. Объект рабочий (W) при-

нимает данные от объекта мастер, выполняет счет и возвращает результат мастеру.

Набор событий, связанных с объектом мастер i -го узла (B_i), можно представить в виде

$\alpha B_i = \{B_i.принять.W_{ij}, B_i.передать.W_{ij}, B_i.счет\}$, где события: $B_i.передать.W_{ij}$ – передать данные в узле объекту рабочий W_{ij} , $B_i.счет$ – счет, $B_i.принять.W_{ij}$ – принять данные в узле от объекта рабочий W_{ij} , ($j = 1..K$, K – количество ядер в многоядерном процессоре).

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

$\alpha W_{ij} = \{W_{ij}.принять.B_i, W_{ij}.передать.B_i, W_{ij}.принять.W_{ik}, W_{ij}.передать.W_{ik}\}$, где события: $W_{ij}.принять.B_i$ – прием данных от мастер объекта, $W_{ij}.передать.B_i$ – передача данных мастер объекту, $W_{ij}.счет$ – счет в объекте рабочий, $W_{ij}.принять.W_{ik}$ – прием данных от объекта рабочий W_{ik} , $W_{ij}.передать.W_{ik}$ – передача данных объекту рабочий W_{ik} .

Описание поведения объектов мастер и рабочий зависит от системы связей ядер в многоядерном процессоре, которая будет определять форму взаимодействия объектов мастер и рабочий. В общем случае для полностью связанной архитектуры узла, когда все ядра связаны между собой, поведение объектов в узле на основе модели послышки сообщений можно представить следующим образом.

Поведение объекта мастер B_i (процесс мастера) используя алфавит αB_i , можно описать как $B_i = (B_i.передать.W_{ij} \rightarrow B_i.счет \rightarrow B_i.принять.W_{ij})$.

Поведение объекта рабочий W_{ij} :
 $W_{ij} = (W_{ij}.принять.B_i \rightarrow W_{ij}.счет \rightarrow W_{ij}.передать.B_i)$.

Модель второго уровня $MKS2$ для i -го клиентского узла можно описать как параллельную комбинацию поведения процессов B_i и W_{ij} :

$$МАСТЕР_i = (B_i), \quad РАБОЧИЕ_i = (W_{ij}), \\ УЗЕЛ_i = (B_i \parallel РАБОЧИЕ_i).$$

Взаимодействие процессов в модели МКС2. При рассмотрении вопросов масштабирования внутри узлов МКС, построенных с использованием многоядерных процессоров, взаимодействие процессов осуществля-

ется на основе общей памяти и модели общих переменных. Проблема масштабирования в узле МКС характерна для SMP – систем и связана с наличием общей памяти и когерентностью кэш-памяти.

В работе [2] введено понятие объекта синхронизации (*ОС*) для организации взаимодействия через общие переменные. Для реализации групповой синхронизации будем использовать объект синхронизации типа *M*, который реализует концепцию мониторов. При этом необходимо дополнить *ОС(M)* средствами, обеспечивающими эффективную реализацию групповой синхронизации.

Набор событий, связанных с параметризованным объектом *ОС(M)* для синхронизации можно представить в виде

$$aM(k) = \{ M.ждать(k), M.сигнал, \\ M.несигнал, M.установить(k) \},$$

где события: *M.ждать(k)* – ожидание сигнала от *k* процессов, *M.сигнал* – сигнал о событии, *M.несигнал* – установка *ОС(M)* в не-сигнальное состояние, *M.установить* – настроить объект на синхронизацию с *k* процессами.

В задаче групповой синхронизации, когда один процесс (*A*) ждет сигналов от *s* процессов ($E_1 .. E_s$), необходимо поведение *ОС(M)* описать как $M = ((M.несигнал) \rightarrow M.установить(s) \rightarrow ((M.ждать(s) \mid M.сигнал) \rightarrow M))$, ожидание сигналов в процесс *A* описать как *M.ждать(s)*, а посылку сигналов в процессах $E_1 .. E_s$ реализовать событием *M.сигнал*.

Использование моделей вычислений

После построения модели вычислений она может быть использована для анализа корректности взаимодействия процессов, который выполняется на основе проверки предложенных правил. Это позволит выявить и предотвратить тупиковые ситуации, оптимизировать объемы передаваемой информации, выбрать более эффективные средства коммуникации и синхронизации.

Второй важной стороной моделей вычислений является возможность их применения при разработке распределенных приложений. Использование моделей позволяет упростить построение алгоритмов процессов и их последующую программную реализацию. Построение алгоритмов процессов осуществляется в терминах событий, которые были определены и использовались для объектов

моделей. Сами алгоритмы полностью соответствуют поведению объектов в модели вычислений. При реализации модели следует учитывать тот факт, что большинство современных языков параллельного программирования реализуют концепцию взаимодействующих последовательных процессов, составляющих основу теории Ч.Хоара. Это упрощает реализацию процессов для объектов с помощью процессов (потоков) в программе.

Выбор и реализация объектов коммуникации и синхронизации также упрощается, так как они разработаны с учетом и анализом различных существующих механизмов взаимодействия процессов, таких как семафоры, мютексы, критические секции, замки, мониторы, рандеву [7].

Авторами были проведены экспериментальные исследования эффективности использования предложенных моделей вычислений, построенных на основе теории Ч.Хоара, для разработки приложений для реальной масштабируемой кластерной системы, оснащенной четырехядерными процессорами. Модели вычислений позволили описать особенности архитектуры КСМА, включающей переменное число узлов. Использование моделей вычислений для разработки распределенных приложений позволило сократить время их программирования и отладки. Использовались языки программирования C, Java, C#, Ада, библиотеки Win32, MPI. Разработанный пакет программ эффективно выполнялся в МКС с любым числом узлов, обеспечивая высокоскоростную передачу данных между узлами и высокий коэффициент загрузки многоядерных процессоров в узлах.

Выводы

Предложен подход к построению моделей вычислений для масштабируемых кластерных систем с многоядерной архитектурой, основанный на расширении теории взаимодействующих последовательных процессов путем введения параметризованных объектов синхронизации и коммуникации. Модели позволяют описать архитектуру кластерной системы, поведение процессов и их взаимодействие с учетом особенностей масштабирования системы. Математический аппарат CSP теории позволил сформулировать ряд требований для оптимального взаимо-

действия процессов, выполнить анализ корректности взаимодействия процессов на обоих уровнях кластерной системы при использовании модели, основанной на общих переменных, и модели, основанной посылке сообщений.

Использование предложенных моделей при проектировании программных компонент МКС позволит повысить качество программ и сократить время их разработки.

Список литературы

1. TOP500 Supercomputing sites. <http://www.top500.org>
2. Садег Раггу. Модели вычислительных процессов в кластерных системах с многоядерной архитектурой // Вісник НТУУ «КПІ». Сер. Інформатика, управління та обчислювальна техніка. № 48. – 2008. С.58-62.
3. Хоар Ч. Взаимодействующие последовательные процессы: Пер. с англ. – М.: Мир, 1989. – 264 с.
4. Топорков В.В. Модели распределенных вычислений. – М.: ФИЗМАТЛИТ, 2004. – 320 с.
5. Любченко В.С. Автоматная модель параллельных вычислений / В.С.Любченко // Высокопроизводительные параллельные вычисления на кластерных системах : 3-я междунар. конф. 2-4 окт. 2006 г. : труды конф. – М. - 2006. – С. 1359 – 1374.
6. Березовский А. Организация вычислений в масштабируемых кластерных системах // Комп'ютерні системи та мережеві технології: друга міжнар.наук.- техн.конф., CSNT'2009, 10-12 черв. 2009 р.: збірник наукових праць. – 2009, вип. 2(24).
7. Жуков І.А., Корочкін О.В. Паралельні та розподілені обчислення. Навч. посібник. – К.: Корнійчук, 2005. – 226 с.
8. Корочкин А.В. Ада95. Введение в программирование. – К.: Корнійчук, 1998. – 240 с.
9. Каляев И.А. Реконфигурируемые мультиконвейерные вычислительные структуры / И.А.Каляев, И.И.Левин, Е.А.Семерников, В.И.Шмойлов. – М.: ид-во ЮНЦ РАН, 2008. – 320 с.

Поступила в редакцию 14.12.2009