

# ПРОГРАММНЫЙ КОМПЛЕКС ДЛЯ МОДЕЛИРОВАНИЯ И ОПТИМИЗАЦИИ РАСПРЕДЕЛЕННЫХ БАЗ ДАННЫХ КОМПЬЮТЕРНЫХ ИНФОРМАЦИОННЫХ СИСТЕМ

Лаздынь С.В., Телятников А.О.

Донецкий национальный технический университет, кафедра АСУ.  
Alexander.Telyatnikov@gmail.com

## **Abstract**

*Lazdyn S.V. Telyatnikov A.O. Program complex for modeling and optimization of the distributed databases of computer information systems. During the program complex development the new approach based on sharing of objective models and genetic algorithms has been used. It allowed to carry out the analysis of characteristics of the DDB functioning with allowance for dynamics of processes proceeding in it and to find suboptimum schemes of allocating of the data on units of system.*

В настоящее время широкое распространение получают компьютерные информационные системы (КИС), одной из основных частей которых являются распределенные базы данных (РБД). Производительность РБД во многом определяется ее организацией и тем, как данные размещены на узлах КИС. Однако ожидаемый эффект от внедрения тех или иных аппаратных или программных решений, а следовательно, и от капиталовложений, необходимых для их реализации, зачастую оценивается интуитивно и не подкреплен расчетами. Поэтому проектировщикам и администраторам РБД необходимы инструментальные средства для моделирования и оптимизации РБД КИС.

В предыдущих работах [1, 2] для моделирования РБД предлагалось использование аналитических моделей, а для оптимизации – методов математического программирования. Однако применение аналитических моделей не позволяет учесть динамику процессов, протекающих в РБД, а методы математического программирования затруднительно применять для оптимизации систем с большим числом узлов и фрагментов данных. Для оптимизации распределения фрагментов данных в РБД предлагалось использование генетических алгоритмов совместно с аналитической моделью [3]. Однако в данном подходе также не учитывается динамика процессов, протекающих в РБД, в частности, не учитываются такие аспекты как возникновение задержек, вызванных высокой нагрузкой на сеть или на узлы обработки запросов.

Предлагаемый подход основан на совместном использовании объектной модели РБД [4] и модифицированных генетических алгоритмов [5]. На основе этого подхода разработан программный комплекс, предназначенный для моделирования РБД с целью выявления аппаратных и программных компонентов, параметры которых ограничивают повышение производительности системы, а также определения оптимальной или субоптимальной схемы распределения данных по узлам информационной сети.

### **Структура программного комплекса**

Программный комплекс состоит из следующих подсистем: подсистема моделирования РБД, подсистема анализа результатов моделирования, подсистема оптимизации распределения данных. Структура программного комплекса представлена на рис. 1.

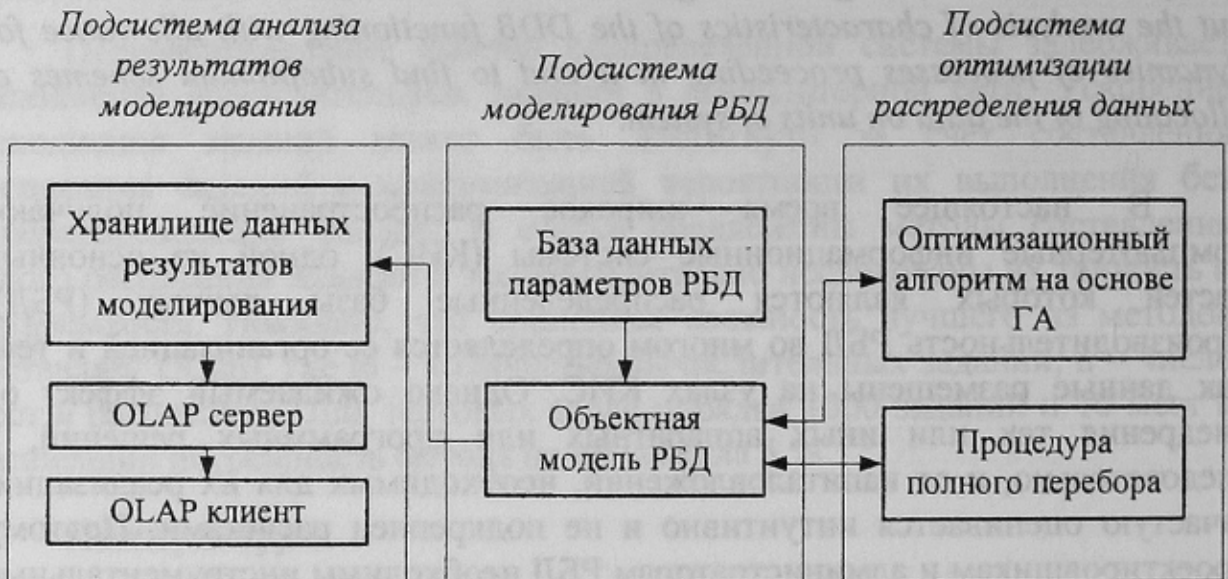


Рисунок 1 - Структура программного комплекса.

Программный комплекс разработан на языке C++ с использованием среды разработки Borland C++ Builder.

Для ввода параметров РБД, а также для работы с моделью и алгоритмом оптимизации был создан графический пользовательский интерфейс, который включает набор экранных форм.

### **Подсистема моделирования РБД**

Подсистема моделирования представляет собой программную реализацию объектной модели РБД [4]. Данная объектная модель построена как система взаимодействующих объектов ее типовых компонентов, выделенных в результате системного анализа. Типовые

компоненты РБД: узел, канал передачи данных, приложение, запрос, таблица РБД. Для каждого из перечисленных компонентов разработаны классы объектов.

Исходными данными для моделирования РБД являются:

1. Множество узлов  $N_i$  ( $i = 1, 2, K, n$ ), где  $n$  – количество узлов РБД.
2. Множество каналов связи  $C_i$  ( $i = 1, 2, K, k$ ), где  $k$  – количество каналов связи.
3. Множество фрагментов данных  $D_i$  ( $i = 1, 2, K, m$ ), где  $m$  – количество фрагментов данных.
4. Схема распределения данных по узлам сети.

Фрагменты данных

		1	2	...	m
Узлы	1	$a_{11}$	$a_{12}$	...	$a_{1m}$
	2	$a_{21}$	$a_{22}$	...	$a_{2m}$
	⋮				
	n	$a_{n1}$	$a_{n2}$	...	$a_{nm}$

$$a_{ij} = \begin{cases} 1, & \text{если на } i\text{-том узле хранится копия } j\text{-того фрагмента;} \\ 0, & \text{иначе.} \end{cases}$$

Рисунок 2 - Схема распределения данных.

Схема распределения данных по узлам сети представляет собой двумерный массив размером  $N \times M$ , количество строк которого равно количеству узлов, а количество столбцов – количеству фрагментов данных (рис. 2). В РБД должна присутствовать хотя бы одна копия каждого набора данных, поэтому в каждом столбце данного массива должна быть хотя бы одна единица.

5. Множество приложений  $A_i$  ( $i = 1, 2, K, N_A$ ), где  $N_A$  – количество приложений, инициирующих запросы на чтение и обновление данных в РБД.

6. Множество запросов  $Q_i$  ( $i = 1, 2, K, N_Q$ ), где  $N_Q$  – количество запросов.

7. Время моделирования  $T$ .

Для хранения параметров моделируемой РБД была создана база данных. В качестве СУБД используется Microsoft Access. Перечень таблиц представлен в таблице 1.

Таблица 1 Перечень таблиц БД моделирования

№	Название	Комментарии
1	RBD	Информация о моделируемой РБД
2	Nodes	Информация об узлах РБД
3	Applications	Информация о приложениях
4	NodeApps	Размещение приложений
5	Queryes	Информация о запросах и обновлениях
6	QueryTypes	Типы запросов
7	AppQueries	Состав приложений
8	Datasets	Информация о фрагментах данных
9	DatasetLocation	Информация о размещении данных
10	Stat	Результаты моделирования

В процессе моделирования РБД происходит имитация процессов выполнения запросов и распространения обновлений. При этом для каждого запроса вычисляются моменты наступления событий:

- постановка в очередь передачи;
- начало передачи;
- завершение передачи;
- постановка в очередь обработки;
- начало обработки;
- завершение обработки.

Эти значения записываются в базу данных для дальнейшей аналитической обработки. Также по этим значениям могут вычисляться различные параметры функционирования РБД, в том числе критерий оценки эффективности РБД, который будет описан ниже.

### **Подсистема анализа результатов моделирования**

Для того чтобы результаты, полученные при моделировании, способствовали принятию решений относительно конфигурации РБД и распределения данных по узлам, они должны быть представлены в удобной, информативной форме. Для этого существуют различные технологии аналитической обработки данных. Комплексный взгляд на собранную в процессе моделирования информацию, ее обобщение и многомерный анализ позволяют реализовать системы оперативной аналитической обработки данных (Online Analytical Processing, OLAP). В основе концепции OLAP лежит принцип многомерного представления данных.

Данные, полученные в результате моделирования, представляются в виде многомерных кубов, осями которых служат измерения. На пересечении этих осей в соответствующей ячейке находятся значения показателей.

Для анализа характеристик работы РБД представляет интерес информация о времени выполнения пользовательских запросов и распространения обновлений, загруженности узлов и каналов передачи данных. Поэтому анализ функционирования РБД можно разбить на 3 составляющие:

- анализ времени выполнения запросов и распространения обновлений;
- анализ загруженности узлов РБД;
- анализ загруженности каналов передачи данных.

Для анализа результатов моделирования РБД было создано 3 многомерных OLAP куба.

Первый куб предназначен для анализа времени выполнения запросов и распространения обновлений. Данный куб состоит из таблицы фактов "QueryProcessing", в которой хранятся значения показателя "Длительность (Duration)", и трех таблиц измерений (рис. 3).

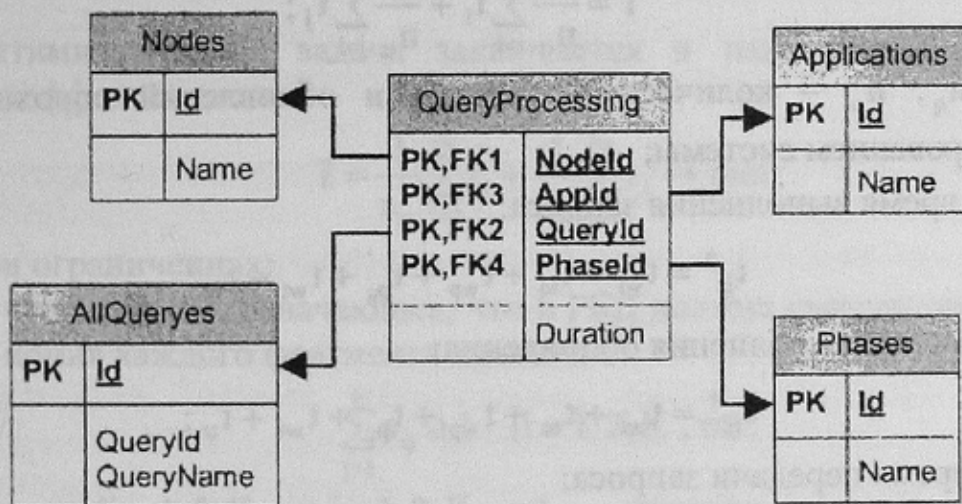


Рисунок 3 - Схема OLAP куба "Выполнение запросов".

Измерение "Узел (Nodes)" содержит перечень всех узлов РБД. Предназначено для анализа времени выполнения запросов и распространения обновлений по узлам.

Измерение "Приложения (Applications)" содержит названия приложений, инициирующих запросы и обновления.

Измерение "Все запросы (AllQueryes)" содержит название и уникальный номер каждого запроса, инициированного при моделировании.

Измерение “Этапы (Phases)” содержит этапы выполнения запросов и обновлений: а) ожидание передачи; б) передача; в) ожидание обработки; г) обработка; д) ожидание передачи ответа; е) передача ответа.

Остальные OLAP кубы, предназначенные для анализа загруженности узлов РБД и каналов передачи данных, выполнены аналогично.

### **Подсистема оптимизации распределения данных**

Подсистема оптимизации распределения данных состоит из двух блоков. Первый, представляет собой программную реализацию оптимизационного алгоритма на основе генетических алгоритмов (ГА), позволяющего находить субоптимальные решения [5]. Второй – программную реализацию процедуры полного перебора, позволяющий находить глобальный оптимум задачи, но имеющий большую вычислительную сложность.

В качестве критерия эффективности РБД (целевой функции) предложено использовать среднее время выполнения запросов и распространения обновлений, порожденных функционированием системы в течении времени моделирования:

$$T = \frac{1}{n_q} \sum_{i=1}^{n_q} t_i + \frac{1}{n_u} \sum_{j=1}^{n_u} t_j;$$

где  $n_q$ ,  $n_u$  – количество запросов и обновлений порожденных функционированием системы;

$t_{ij}^q$  – время выполнения запроса:

$$t_{ij}^q = t_{wt} + t_{tq} + t_{wp} + t_{p_{ij}} + t_{wt} + t_{ta};$$

$t_{ij}^u$  – время выполнения обновления:

$$t_{ij}^u = t_{wt} + t_{tu} + t_{wp} + t_{p_{ij}} + t_{wt} + t_{tr};$$

$t_{tq}$  – время передачи запроса:

$$t_{tq} = \frac{V_q}{\min(B_l)};$$

$t_{ta}$  – время передачи ответа на запрос:

$$t_{ta} = \frac{V_a}{\min(B_l)};$$

$t_{p_{ij}}$  – время обработки запроса:

$$t_{p_{ij}} = \frac{K_j}{P_i} \cdot 60 (c);$$

$t_{tu}$  – время передачи обновления:

$$t_{tu} = \frac{V_u}{\min(B_l)};$$

$t_{tr}$  – время передачи сообщения о завершении обновления:

$$t_{tr} = \frac{V_r}{\min(B_l)};$$

$t_{wt}$  – время ожидания передачи;

$t_{wp}$  – время ожидания обработки;

$V_q$  – объем запроса;

$V_a$  – объем ответа на запрос;

$P_i$  – производительность  $i$ -того узла передачи данных (tpmC);

$K_j$  – количество транзакций ТРС-С соответствующее  $j$ -тому запросу;

$V_u$  – объем обновления;

$V_r$  – объем сообщения о завершении обновления;

$B_l$  – пропускная способность  $l$ -го канала передачи данных.

Оптимизационная задача заключается в нахождении минимума критерия эффективности:

$$T = \frac{1}{n_q} \sum_{i=1}^{n_q} t_i + \frac{1}{n_u} \sum_{j=1}^{n_u} t_j \rightarrow \min;$$

при ограничениях:

1) ограничение, означающее, что в РБД должна присутствовать хотя бы одна копия каждого фрагмента данных:

$$\sum_{j=1}^n a_{ij} \geq 1 \quad (i = 1, 2, K, m);$$

где  $a_{ij}$  ( $i = 1, 2, K, n$ ;  $j = 1, 2, K, m$ ) – признак наличия фрагмента данных на узле РБД, определяемый по формуле:

$$a_{ij} = \begin{cases} 1, & \text{если } j\text{-тый фрагмент данных находится на } i\text{-том узле;} \\ 0, & \text{в противном случае;} \end{cases}$$

2) ограничение, указывающее на то, что суммарный объем данных, хранящихся на узле, не должен превышать общее дисковое пространство данного узла:

$$\sum_{j=1}^m V_j \cdot a_{ij} \leq L_i \quad (i = 1, 2, K, n);$$

где  $V_j$  – объем  $j$ -того фрагмента данных;

$L_i$  – общее дисковое пространство  $i$ -того узла.

Для оптимизации распределения данных по узлам сети предложен новый подход, основанный на совместном использовании объектной модели РБД и аппарата ГА.

Схема распределения фрагментов данных по узлам РБД кодируется в виде хромосомы ГА. Популяция ГА представляет собой набор некоторых точек пространства поиска. Начальная популяция генерируется случайным образом. В процессе оптимизации с помощью операторов ГА генерируются хромосомы, то есть схемы распределения данных. Полученные схемы являются исходной информацией для объектной модели, с помощью которой вычисляются оценки критерия эффективности РБД. Эти оценки, в свою очередь, являются значениями функции приспособленности ГА для данного варианта решения. То есть, другими словами, предлагается использовать разработанную объектную модель РБД для вычисления функции приспособленности ГА.

Исходной информацией для алгоритма оптимизации является:

- количество поколений ГА;
- размер популяции (количество особей в популяции);
- размер мультихромосомы особи (количество хромосом в одной мультихромосоме);
- размер хромосомы (количество ген в хромосоме);
- вероятность применения оператора рекомбинации к хромосомам двух особей;
- вероятность применения оператора скрещивания к хромосомам двух особей;
- вероятность применения оператора мутации.

Разработанный алгоритм оптимизации реализует случайный направленный поиск. На каждой итерации алгоритма рассматривается некоторое множество точек пространства поиска, представленное популяцией особей ГА. Затем, с помощью операторов ГА, генерируется новое множество точек пространства поиска. При этом каждое последующее множество решений, то есть каждая последующая популяция, будет содержать решения в целом лучше предыдущих. После того, как выполнено определенное число итераций алгоритма (рассмотрено заданное количество поколений ГА), процесс оптимизации заканчивается. Лучшее из полученных решений является субоптимальным.

Разработанный программный комплекс прошел экспериментальную проверку. В качестве объекта экспериментальных исследований выбрана РБД КИС ЗАО ПО “Киев-Континти”. Моделирование данной РБД позволило выявить аппаратные элементы, параметры которых ограничивают повышение производительности системы. С помощью подсистемы оптимизации распределения данных определена оптимальная и субоптимальная схемы распределения данных по узлам КИС. На основе



полученных результатов были разработаны рекомендации по повышению эффективности РБД КИС ЗАО ПО "Киев-Конти" с минимальными материальными затратами на модернизацию.

Проведенные экспериментальные исследования показали работоспособность и адекватность разработанного программного комплекса для моделирования и оптимизации РБД.

### **Заключение**

Разработанное инструментальное средство может быть использовано как на этапе проектирования КИС с РБД, так и при эксплуатации уже существующих систем.

Для вновь проектируемых КИС с РБД разработанный программный комплекс ускоряет процесс создания эффективной структуры РБД, позволяет найти оптимальную схему распределения данных по узлам сети.

На этапе эксплуатации уже существующих КИС с РБД данное инструментальное средство позволяет увеличить производительность системы путем выявления и устранения так называемых «узких мест», а также путем оптимизации распределения данных по узлам сети.

В дальнейших исследованиях планируется проведение дополнительных вычислительных экспериментов с использованием разработанного программного комплекса по моделированию и оптимизации РБД различной размерности. Дополнительные эксперименты необходимы для определения зависимости параметров оптимизационного алгоритма, при которых происходит наилучшее приближение субоптимальных решений, полученных с его помощью, к глобальному минимуму, полученному путем полного перебора, от размерности РБД.

### **Литература**

1. Галкин В.Е. Методы оптимальной организации распределенной информационной системы. // Автоматизация и современные технологии. – 2004. – № 4. – С. 13 – 17.
2. Цегелик Г.Г. Системы распределенных баз данных. – Львов.: Свит, 1990. – 168 с.
3. Corcoran A.L., Hale J. A genetic algorithm for fragment allocation in a distributed database system // Proceedings of the 1994 ACM symposium on Applied computing. – Phoenix: ACM Press, 1994. – P. 247 – 250.
4. Телятников А.О. Разработка объектной модели распределенной базы данных // Наукові праці ДонНТУ. Випуск 74. – Донецьк: ДонНТУ, 2004. – С. 192 – 200.
5. Лаздынь С.В., Телятников А.О. Оптимизация распределенных баз данных с использованием генетических алгоритмов // Вестник ХГТУ № 1(19) 2004. – С. 236 – 239.