

УДК 004.4'244

С.Н. Гриша, О.С. Родичева, Е.О. Роздольский

ПОСТРОЕНИЕ ТЕКСТА SQL ЗАПРОСА ИЗ РЕЛЯЦИОННОГО ДЕРЕВА ЗАПРОСА

Постановка задачи

В статье [1] рассматривался общий подход к автоматизации процесса проектирования систем управления бизнесом. В ней были обоснованы актуальность и необходимость автоматического построения базы данных (БД) и программного обеспечения (ПО) на основе семантического описания предметной области.

При решении задачи автоматического построения БД и ПО необходимо решить ряд технических и теоретических задач. В предложенном подходе аналитик-разработчик описывает систему в терминах показателей. Это означает, что аналитик-разработчик сначала проектирует систему показателей, и после этого для каждого документа системы задает ему запрос на выборку над показателями, используя правила реляционной алгебры. При генерировании БД и ПО специальный алгоритм генерирует БД и при генерировании ПО алгоритм замещает запросы реляционной алгебры над показателями запросами над таблицами, канонизирует полученные выражения и преобразует их в текст SQL. В этой статье решается задача преобразования реляционного дерева запроса над таблицами в текст SQL. Входные данные для задачи – это реляционное дерево запроса (соответствует реляционной алгебре Кода [2]). На основании его необходимо построить SQL запрос стандарта SQL-92.

Обзор публикаций

При изучении публикаций на выбранную тематику было найдено много работ по преобразованию SQL-запросов в реляционное дерево [3], но практически отсутствуют материалы по обратному преобразованию. Есть материалы по преобразованию одной структурированной информации в другую. Так, в [4] описан метод преобразования структурированных документов одной грамматики в другую, причем это должны быть контекстно-независимые грамматики и входные данные – структурированный текст.

© Гриша С.Н., Родичева О.С., Роздольский Е.О., 2009

Также существуют преобразователи xml-документов, такие как XSLT [5]. Но такие преобразователи тоже не подходят для решения поставленной задачи, так как преобразуют из xml в xml и малопригодны для преобразования в SQL.

Изучение трансформационно-порождающих (ТП) грамматик в литературе показало, что большинство исследований направлено на изучение языковых структур и лингвистических преобразований из одного разговорного языка в другой [6]. Применение же термина "ТП грамматика" не в области лингвистики встречалось крайне редко. Несмотря на это, идея применения ТП грамматики в лингвистике подходит и для преобразования других структурированных языков, но только контекстно-независимых. Для контекстно-зависимых грамматик также описаны алгоритмы, но они неточны и требуют конкретного описания для каждого языка.

Алгоритм

Для преобразования реляционного дерева в SQL запрос применяется трансформационная порождающая грамматика, на вход которой подается реляционное дерево S (запускается выполнение продукций грамматики). Рассмотрена была стандартная грамматика SQL стандарта ANSI [7] в расширенной форме Бэкуса-Наура [8].

На применение порождающей грамматики накладываются такие условия:

1. за один шаг должны быть применены одна продукция, т.е. дерево S уменьшится на один узел;
2. продукции должны просматриваться в том порядке, в котором они представлены.

Каждый реляционный оператор порождает или терминал(ы) или появление вложенного подзапроса, причем второе случается, если соответствующая позиция в текущем транслированном SQL-выражении уже занята.

Очевидно, грамматика для преобразования контекстно-зависима [9], так как узлам реляционного дерева соответствуют несколько правил вывода, что зависит от текущего вида полученного выражения SQL. Например, рассмотрим такие два реляционных дерева: $Q_1 = \gamma_A(\sigma_C(T))$ и $Q_2 = \sigma_C(\gamma_A(T))$. В первом случае при преобразовании в SQL запрос необходимо в нем сделать подзапрос, во

втором случае не нужно хотя операторы использованы одинаковые.

SQL запросы для деревьев Q_1 и Q_2 будут такими:

$\text{select } (Q_1) = \text{select } A \text{ from } (\text{select } * \text{ from } T \text{ where } C) \text{ group by } A / A_{agg}$

$\text{select } (Q_2) = \text{select } A \text{ from } T \text{ where } C \text{ group by } A / A_{agg}$

где A_{agg} - агрегирующие аргументы оператора группировки.

Таким образом, терминальные символы вместе с нетерминалом будут определять продукцию в трансформационной грамматике. Для упрощения грамматики (уменьшения количества правил) добавим предварительную обработку дерева для его упорядочивания. Сделать это можно после этапа канонизации, добавив в правила канонизации правила упорядочивания структуры дерева.

Определим необходимый порядок для структуры дерева. Пусть это будет такой порядок для операторов с одним аргументом: $\delta, \tau, \pi, \sigma, \gamma$. Операторы с двумя аргументами (операторы соединения и объединения) тут нельзя рассматривать так это относится к оптимизации запроса, над чем работает алгоритм канонизации. Алгоритм канонизации выражений меняет дерево запроса так, чтобы привести его к оптимальной структуре запроса. Изменения же порядка в дереве необходимо только для упрощения процесса построения текста SQL запроса и не влияет на быстродействие запроса.

Для получения нужного порядка необходимо применить такие правила как показано в таблице Таблица 1. По строкам показан первый оператор, по колонкам – оператор-аргумент. В соответствующих ячейках на пересечении элементов стоит:

– прочерк в том случае, если такой порядок следования операторов корректен, и его не нужно менять;

– " $\rightarrow Q$ " в том случае, если соответствующее изменение порядка не может быть сделано (или из-за несоответствия аргументов или из-за несохранения результирующего набора данных);

– формула – это правая часть продукции, которую необходимо применить к соответствующим операторам, чтобы упорядочить дерево; левая часть продукции образуется из значения строки и колонки.

Таблица 1

Правила упорядочивания дерева

	(δ)	(π_L)	(σ_C)	(γ_A)	(τ_M)
δ	$\rightarrow \delta(*)$	—	—	—	—
π_L	$\rightarrow \delta(\pi_L(*))$	$\rightarrow \pi_{L(L)}(*)$	—	—	$\rightarrow Q$
σ_C	$\rightarrow \delta(\sigma_C(*))$	$\rightarrow \pi_{L'}(\sigma_C(*))$	$\rightarrow \sigma_{C \text{ and } C'}(*)$	—	$\rightarrow \tau_M(\sigma_C(*))$
γ_A	$\rightarrow Q$	$\rightarrow Q$	$\rightarrow Q$	$\rightarrow Q$	$\rightarrow Q$
τ_M	$\rightarrow \delta(\tau_M(*))$	—	—	—	$\rightarrow \tau_M(*)$

Рассмотрим подробнее таблицу. В ячейках по диагонали (когда оператор и его аргумент одинаковы), происходит слияние для всех операторов кроме группировки γ_A с изменением аргументов. Если аргумент оператора группировки не объединение, соединение или таблица, то этот аргумент порождает вложенный подзапрос. Все аргументы оператора селекции кроме группировки поднимаются вверх по дереву запроса. Оператор группировки Оператор удаления дубликатов всегда необходимо перемещать по дереву вверх. Это возможно во всех случаях кроме группировки, так как это повлияло бы на результирующий набор данных. Оператор проекции можно поднимать вверх по дереву только если он является аргументом селекции и сливать с оператором проекции. Оператор сортировки поднимается вверх по дереву, если он является аргументом селекции, в остальных случаях или порядок правильный, или порядок нельзя менять, так как это приведет к изменению набора данных.

Таким образом, получаем следующий набор продукций грамматики $G_0\{T_0, N_0, P_0, S_0\}$ для упорядочивания дерева:

$$T_0 = U \setminus N_0, N = \{\pi_L, \sigma_C, \gamma_A, \delta, \tau_M\}$$

P_0 :

$$\delta(\delta(*)) \rightarrow \delta(*)$$

$$\pi_L(\delta(*)) \rightarrow \delta(\pi_L(*))$$

$$\pi_L(\pi_L(*)) \rightarrow \pi_{L(L)}(*)$$

$$\sigma_C(\delta(*)) \rightarrow \delta(\sigma_C(*))$$

$$\sigma_C(\pi_L(*)) \rightarrow \pi_{L'}(\sigma_C(*))$$

$$\sigma_C(\sigma_C(*)) \rightarrow \sigma_{C \text{ and } C'}(*)$$

$$\sigma_C(\tau_M(*)) \rightarrow \tau_M(\sigma_C(*))$$

$$\tau_M(\delta(*)) \rightarrow \delta(\tau_M(*))$$

$$\tau_M(\tau_{M'}(*)) \rightarrow \tau_M(*)$$

Этот набор продукций необходимо применить к реляционному дереву перед применением порождающей грамматики. Тогда порождающая грамматика упрощается так и будет обрабатывать операторы по мере их появления. Таким образом, получаем грамматику $G(T, N, P, S)$ такую:

Терминалы грамматики G :

$$T = \{select, from, where, left, right, full, natural, join, distinct, order by, group by, *, dual, t_1, \dots, t_n, l_1, \dots, l_m\}$$

где t_1, \dots, t_n - названия таблиц, l_1, \dots, l_m - названия всех аргументов реляционных операторов кроме таблиц каждого дерева:
 $t_1, \dots, t_n, l_1, \dots, l_m \notin T \setminus \{t_1, \dots, t_n, l_1, \dots, l_m\}$

Нетерминалы грамматики G :

$$N = \{\delta, \tau, \pi, \sigma, \gamma, \times, \triangleright \triangleleft, \triangleright \triangleleft_C, \triangleright \triangleleft_{\dot{C}}, \triangleright \triangleleft_{\dot{L}}, \triangleright \triangleleft_{\dot{R}}, \triangleright \triangleleft_{\dot{C}}, \cup\}$$

Продукции грамматики:

$$\varepsilon \rightarrow \varepsilon S$$

$$\varepsilon S \rightarrow S(Q)$$

$$(t_i)(Q) \rightarrow t_i$$

$$S(Q) \rightarrow (\text{select } S(\text{from}))$$

$$\delta(S)(\text{from}) \rightarrow \text{distinct } S(\text{ from})$$

$$\tau_i(S)(\text{from}) \rightarrow S(\text{from}) \text{ order by } l_i$$

$$\pi_i(S)(\text{from}) \rightarrow S(l_i \text{ from})$$

$$\sigma_i(S)(l_j \text{ from}) \rightarrow S(l_j \text{ from where } l_i)$$

$$\sigma_i(S)(\text{from}) \rightarrow S(* \text{ from where } l_i)$$

$$\gamma_{l_i \rightarrow l_i'}(S)(\text{from}) \rightarrow l_i' \text{ from } S(Q) \text{ group by } l_i$$

$$\gamma_{l_i \rightarrow l_i'}(S)(l_j \text{ from}) \rightarrow l_j \text{ from } (\gamma_{l_i \rightarrow l_i'}(S)(Q))$$

$$\gamma_{l_i \rightarrow l_i'}(S)(l_j \text{ from where } l_k) \rightarrow l_j \text{ from } (\gamma_{l_i \rightarrow l_i'}(S)(Q)) \text{ where } l_k$$

$$\gamma_{l_i \rightarrow l_i'}(S)(* \text{ from where } l_j) \rightarrow l_i' \text{ from } (S) \text{ where } l_j \text{ group by } l_i$$

$$(t_i)(\text{ from}) \rightarrow * \text{ from } t_i$$

$$(t_i)(l_i \text{ from}) \rightarrow l_i \text{ from } t_i$$

$$(t_i)(l_j \text{ from where } l_i) \rightarrow l_j \text{ from } t_i \text{ where } l_i$$

$$\begin{aligned}
 (t_i)(* \text{ from where } l_i) &\rightarrow * \text{ from } t_i \text{ where } l_i \\
 (t_i) &\rightarrow t_i \\
 \times(S_1, S_2)(l_i \text{ from}) &\rightarrow l_i \text{ from } \times(S_1, S_2) \\
 \times(S_1, S_2)(l_j \text{ from where } l_i) &\rightarrow l_j \text{ from } \times(S_1, S_2) \text{ where } l_i \\
 \times(S_1, S_2)(* \text{ from where } l_i) &\rightarrow * \text{ from } \times(S_1, S_2) \text{ where } l_i \\
 \times(S_1, S_2)(\text{from}) &\rightarrow \text{from } \times(S_1, S_2) \\
 (\times(S_1, S_2)) &\rightarrow (S_1)(Q), (S_2)(Q) \\
 \triangleright\triangleleft(S_1, S_2)(l_i \text{ from}) &\rightarrow l_i \text{ from } \triangleright\triangleleft(S_1, S_2) \\
 \triangleright\triangleleft(S_1, S_2)(l_j \text{ from where } l_i) &\rightarrow l_j \text{ from } \triangleright\triangleleft(S_1, S_2) \text{ where } l_i \\
 \triangleright\triangleleft(S_1, S_2)(* \text{ from where } l_i) &\rightarrow * \text{ from } \triangleright\triangleleft(S_1, S_2) \text{ where } l_i \\
 \triangleright\triangleleft(S_1, S_2)(\text{from}) &\rightarrow \text{from } \triangleright\triangleleft(S_1, S_2) \\
 (\triangleright\triangleleft(S_1, S_2)) &\rightarrow (S_1)(Q) \text{ natural join } (S_2)(Q) \\
 \triangleright\triangleleft_{l_i}(S_1, S_2)(l_j \text{ from}) &\rightarrow l_j \text{ from } \triangleright\triangleleft_{l_i}(S_1, S_2) \\
 \triangleright\triangleleft_{l_i}(S_1, S_2)(l_j \text{ from where } l_k) &\rightarrow l_j \text{ from } \triangleright\triangleleft_{l_i}(S_1, S_2) \text{ where } l_k \\
 \triangleright\triangleleft_{l_i}(S_1, S_2)(* \text{ from where } l_j) &\rightarrow * \text{ from } \triangleright\triangleleft_{l_i}(S_1, S_2) \text{ where } l_j \\
 \triangleright\triangleleft_{l_i}(S_1, S_2)(\text{from}) &\rightarrow \text{from } \triangleright\triangleleft_{l_i}(S_1, S_2) \\
 (\triangleright\triangleleft_{l_i}(S_1, S_2)) &\rightarrow (S_1)(Q) \text{ join } (S_2)(Q) \text{ on } l_i \\
 \overset{\circ}{\triangleright\triangleleft}(S_1, S_2)(l_j \text{ from}) &\rightarrow l_j \text{ from } \overset{\circ}{\triangleright\triangleleft}(S_1, S_2) \\
 \overset{\circ}{\triangleright\triangleleft}(S_1, S_2)(l_j \text{ from where } l_k) &\rightarrow l_j \text{ from } \overset{\circ}{\triangleright\triangleleft}(S_1, S_2) \text{ where } l_k \\
 \overset{\circ}{\triangleright\triangleleft}(S_1, S_2)(* \text{ from where } l_j) &\rightarrow * \text{ from } \overset{\circ}{\triangleright\triangleleft}(S_1, S_2) \text{ where } l_j \\
 \overset{\circ}{\triangleright\triangleleft}(S_1, S_2)(\text{from}) &\rightarrow \text{from } \overset{\circ}{\triangleright\triangleleft}(S_1, S_2) \\
 (\overset{\circ}{\triangleright\triangleleft}(S_1, S_2)) &\rightarrow (S_1)(Q) \text{ full join } (S_2)(Q) \\
 \overset{\circ}{\triangleright\triangleleft}_X(S_1, S_2)(l_j \text{ from}) &\rightarrow l_j \text{ from } \overset{\circ}{\triangleright\triangleleft}_X(S_1, S_2) \\
 \overset{\circ}{\triangleright\triangleleft}_X(S_1, S_2)(l_j \text{ from where } l_k) &\rightarrow l_j \text{ from } \overset{\circ}{\triangleright\triangleleft}_X(S_1, S_2) \text{ where } l_k \\
 \overset{\circ}{\triangleright\triangleleft}_X(S_1, S_2)(* \text{ from where } l_j) &\rightarrow * \text{ from } \overset{\circ}{\triangleright\triangleleft}_X(S_1, S_2) \text{ where } l_j \\
 \overset{\circ}{\triangleright\triangleleft}(S_1, S_2)(\text{from}) &\rightarrow \text{from } \overset{\circ}{\triangleright\triangleleft}(S_1, S_2)
 \end{aligned}$$

$$\begin{aligned}
& (\triangleright \triangleleft_X(S_1, S_2)) \rightarrow (S_1)(Q) \text{ full join } (S_2)(Q) \text{ on } l_i \\
& \cup(S_1, S_2)(l_j \text{ from}) \rightarrow l_j \text{ from } \cup(S_1, S_2) \\
& \cup(S_1, S_2)(l_j \text{ from where } l_k) \rightarrow l_j \text{ from } \cup(S_1, S_2) \text{ where } l_k \\
& \cup(S_1, S_2)(* \text{ from where } l_j) \rightarrow * \text{ from } \cup(S_1, S_2) \text{ where } l_j \\
& \cup(S_1, S_2)(\text{from}) \rightarrow \text{from } \cup(S_1, S_2) \\
& \cup(S_1, S_2) \rightarrow (S_1)(Q) \text{ union } ((S_2)(Q))
\end{aligned}$$

Граф переходов автомата, соответствующего грамматике G , показан на рисунке Рисунок 1. Для уменьшения количества узлов, соответствующих операторам соединения и объединения сделана такая замена:

$$\otimes = \{\times, \triangleright \triangleleft, \triangleright \triangleleft_{l_i}, \triangleright \triangleleft_{l_i}^{\circ}, \triangleright \triangleleft_X^{\circ}, \cup\}$$

Таким образом, получили алгоритм преобразования реляционного дерева в SQL запрос:

1. подготовить дерево: построить соответствие между всеми конечными аргументами, кроме таблиц, всех операторов и обозначениями l_1, \dots, l_m .

2. применить упорядочивающую грамматику G_o к реляционному дереву S ;

3. применить порождающую грамматику G к упорядоченному дереву – результату применения продукций G_o к реляционному выражению S и получить терминальную цепочку;

4. в полученной терминальной цепочке заменить все $l_i, i = \overline{1, m}$ на соответствующие им значения аргументов.

Оба описанных алгоритма являются модификациями одного и того же подхода – применения трансформационной порождающей контекстно-зависимой грамматики к реляционному дереву запроса.

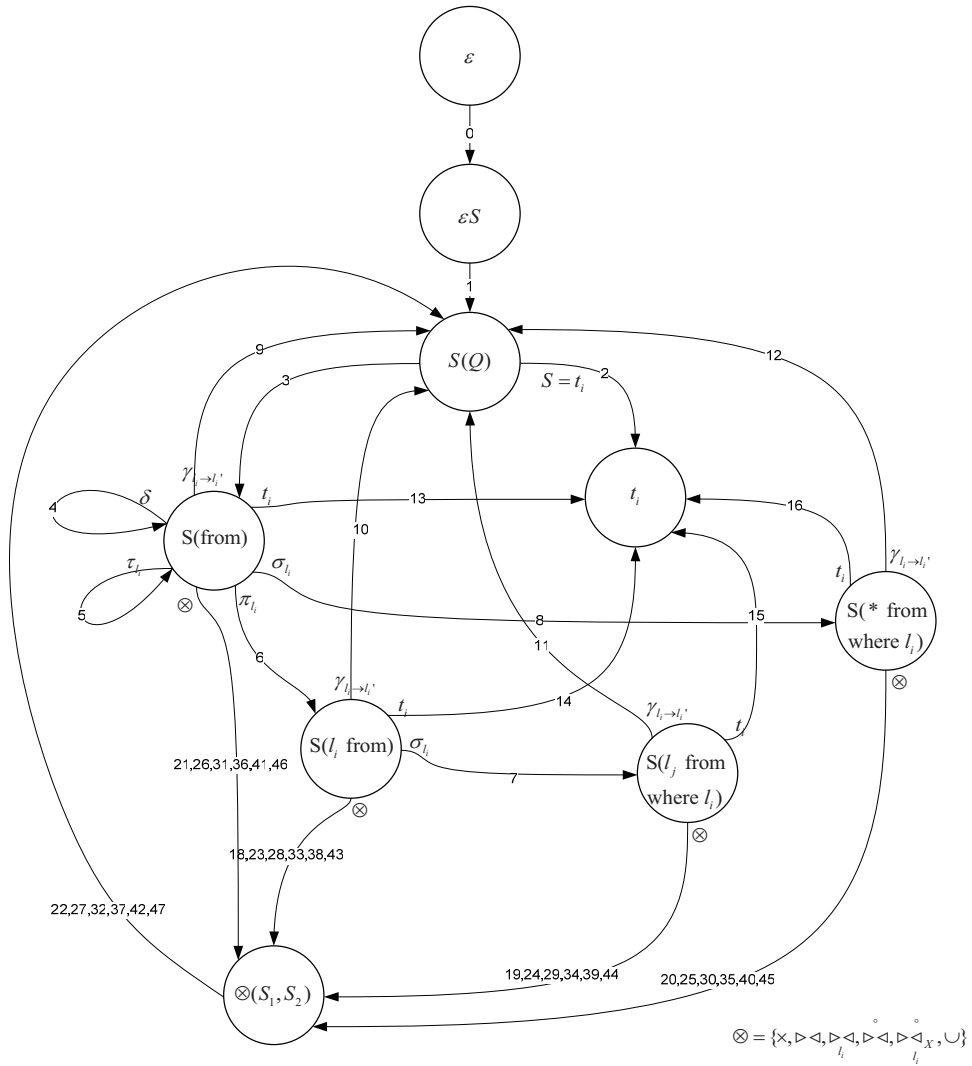


Рисунок 1 – Граф переходов автомата

Примеры преобразования

Рассмотрим пример преобразования дерева $S = \sigma_C(\pi_L(\times(t_1, t_2)))$

Применим алгоритм к заданному примеру.

Шаг 1. Упорядочивание дерева с помощью G_O .

Таблица 2

Упорядочивание дерева G_O

1) $\sigma_C(\pi_L(\times(t_1, t_2))) \rightarrow \pi_L(\sigma_C(\times(t_1, t_2)))$	Правило 5
--	-----------

Шаг 2. Применение грамматики G

Таблица 3

Применение грамматики G

1) $\varepsilon \rightarrow \varepsilon \pi_L(\sigma_C(\times(t_1, t_2)))$	Правило 0
2) $\varepsilon \pi_L(\sigma_C(\times(t_1, t_2))) \rightarrow \pi_L(\sigma_C(\times(t_1, t_2)))(Q)$	Правило 1
3) $\pi_L(\sigma_C(\times(t_1, t_2)))(Q) \rightarrow (\text{select } \pi_L(\sigma_C(\times(t_1, t_2)))(\text{from}))$	Правило 3
4) $(\text{select } \pi_L(\sigma_C(\times(t_1, t_2)))(\text{from})) \rightarrow (\text{select } \sigma_C(\times(t_1, t_2))(l_1 \text{ from}))$	Правило 6, $L = l_1$
5) $(\text{select } \sigma_C(\times(t_1, t_2))(l_1 \text{ from})) \rightarrow (\text{select } \times(t_1, t_2)(l_1 \text{ from where } l_2))$	Правило 7, $C = l_2$
6) $(\text{select } \times(t_1, t_2)(l_1 \text{ from where } l_2)) \rightarrow (\text{select } l_1 \text{ from } \times(t_1, t_2) \text{ where } l_2)$	Правило 19
7) $(\text{select } l_1 \text{ from } \times(t_1, t_2) \text{ where } l_2) \rightarrow (\text{select } l_1 \text{ from } t_1(Q), t_2(Q) \text{ where } l_2)$	Правило 22
8) $(\text{select } l_1 \text{ from } t_1(Q), t_2(Q) \text{ where } l_2) \rightarrow (\text{select } l_1 \text{ from } t_1, t_2 \text{ where } l_2)$	Правило 2

Таким образом, получили

$$S = \sigma_C(\pi_L(\times(t_1, t_2))) \rightarrow (\text{select } l_1 \text{ from } t_1, t_2 \text{ where } l_2)$$

Осталось заменить $l_i, i = \overline{1, 2}$ на соответствующие значения аргументов реляционных операторов и текст SQL запроса готов.

Программные реализация преобразования реляционного дерева в текст SQL запроса

Текст SQL запроса можно генерировать из реляционного дерева, представленного в виде:

1. xml-файла;
2. объектной программной модели;
3. таблицы базы данных.

Все три варианта программного представления реляционного дерева взаимозаменяемы, и при наличии одного из них это представление легко преобразовать в другой.

В зависимости от вида входных данных можно применить такие программные подходы к реализации генеративной грамматики:

1. из xml-описания дерева с использованием сгенерированного с помощью генератора интерпретатора (например модификации Coco/R или ANTLR [10]) парсера xml;

2. также с использованием генератора интерпретатора (например, модификации Coco/R или ANTLR), но при этом алгоритм сразу преобразует входной xml-файл в sql запрос, используя порождающую грамматику;

3. из объектной программной модели: для этого необходимо осуществить обход дерева объектов в глубину с использованием описанных правил грамматики.

В нашей реализации выбран метод обхода объектного дерева, так как реляционное дерево представлено в программной реализации системы в виде объектной модели.

Выводы

В статье рассмотрена и решена задача преобразования реляционного дерева в текст SQL запроса. Сделан обзор существующих публикаций на данную тематику. Показано алгоритм преобразования с использованием порождающей трансформационной грамматики. Показаны примеры применения грамматики к конкретным примерам реляционных деревьев и приведено описание программной реализации алгоритма.

Эти результаты являются частью технологии интеллектуализированного генерирования базы данных и программного обеспечения "КИТ-XXI", основные концепции которой описаны в [1].

ЛИТЕРАТУРА

1. Гриша С.М., Родічева О.С., Приліпко Д.І. Технологічно інтелектуалізовані інформаційні системи для управління бізнесом (ІСУБ) на сонові алгебри показників // Вісник Національного Технічного Університету «ХПІ» – 2008. – №5. – С.123-133.
2. Codd, E.F. *The Relational Model For Database Management Version 2*. Reading, Mass.: Addison-Wesley, 1990. ISBN 0-201-14192-2.
3. Гарсиа-Молина Г., Ульман Д., Уидом Д. Системы баз данных. Полный курс. – М.: Вильямс, 2002. – 1088 с.

4. Eila Kuikka, Martti Penttonen, Transformation of structured documents with the use of grammar. Electronic publishing, VOL. 6(4), 373–383 (december 1993).
5. М. Kay. XSL Transformations (XSLT) Version 2.0. W3C Candidate Recommendation 3th November 2005. November 1999.
6. Касевич В. Б. Элементы общей лингвистики. М.: Наука, 1977.
7. Электронный ресурс: <http://savage.net.au/SQL/sql-92.bnf>
8. Электронный ресурс: ISO/IEC 14977 : 1996(E)
9. Chomsky, Noam (1956). "Three models for the description of language". IRE Transactions on Information Theory (2): 113–124
10. Электронный ресурс: The Catalog of Compiler Construction Tools: <http://catalog.compilertools.net/>

Получено 05.02.2009г.