

平成30年度

筑波大学情報学群情報科学類

卒業研究論文

題目

グループ化機能を持つ統合言語クエリの設計と
効率的実装のためのプログラム変換

主専攻 知能情報メディア主専攻

著者 大倉 瑠維

指導教員 亀山幸義, 海野広志

要 旨

本研究の目的は、SQL のグループ化に直接対応する統合言語クエリを設計することであり、設計した構文の意味論や書き換え規則を追加することである。

SQL におけるグループ化は、グループ化前後でデータの型を変化させてしまうだけでなく、サブクエリのようなネストされた構造を用いることでしか表現できないクエリが存在するため、統合言語クエリにおけるグループ化の拡張は困難とされてきた。本研究では、従来の **for** 項にグループ化機能を付与した **gfor** 項といった構文を新たに提案し、束縛変数の追加や条件式の組み込み等により、一つの **gfor** 項に対して単一のグループ化機能つき SQL クエリに対応するように設計した。グループ化の拡張に伴い新たに追加した書き換え規則による書き換えについては、複数のテーブルに対する処理の仕方や、グループ化キーを増やす処理等、従来とは少し異なる視点で書き換え規則を追加することによって、書き換えを可能にした。また、集約関数について、代数的に使えるものみに制限することで、書き換えによる集約関数の正規化を可能にした。

目次

第 1 章	はじめに	1
1.1	背景	1
1.2	本研究の目的と構成	1
第 2 章	前提知識	3
2.1	正規化	3
2.2	グループ化	6
第 3 章	統合言語クエリ Que Λ	7
3.1	Que Λ の型付け規則	8
3.2	Que Λ の意味論	9
3.3	SQL 変換	9
第 4 章	グループ化機能の拡張	11
4.1	型付け規則の拡張	12
4.2	グループ化機能を追加した Que Λ の意味論	13
4.2.1	集約関数	13
4.2.2	操作的意味論の拡張	14
4.3	拡張する正規形と正規化規則	17
4.3.1	正規形	17
4.3.2	正規化規則	18
4.4	SQL 変換の拡張	19
4.5	正規化、SQL 変換の例	21
第 5 章	結論	23
5.1	まとめ	23
5.2	今後の課題	23
	謝辞	24
	参考文献	25
付録 A	正規化規則	26

目次

2.1	データテーブル	3
2.2	図 2.1 のデータ構造	4
2.3	クエリ Q_1, Q_2 とそれぞれに対応する SQL クエリ	4
2.4	クエリ Q_3 の正規化と SQL 変換	5
3.1	Que Λ の構文	7
3.2	Que Λ の型付け規則	8
3.3	Que Λ の値、評価文脈、操作的意味論	9
3.4	SQL サブ言語	10
3.5	SQL 変換	10
4.1	追加したグループ化機能の構文	11
4.2	グループ化を行う SQL クエリとそれに対応する Que Λ	12
4.3	グループ化機能を追加した Que Λ の型付け規則	13
4.4	グループ化に対する新たな補助関数	15
4.5	グループ化機能を追加した Que Λ における評価文脈, 操作的意味論の定義	16
4.6	グループ化機能を追加した Que Λ の正規形	17
4.7	追加する正規化規則	19
4.8	グループ化機能を追加した SQL サブ言語	20
4.9	SQL 変換	20
A.1	正規化規則 (Stage 1)	27
A.2	正規化規則 (Stage 2)	28

表目次

4.1	for 項と gfor 項の関係	18
-----	--	----

第 1 章

はじめに

1.1 背景

データベースを活用するシステムでは、データベース問い合わせ言語とプログラム言語の 2 つの言語を扱う必要がある。しかし、2 つの言語を扱うことは、コードの複雑化やセキュリティ上の問題により様々な負担をプログラマに強いることになるため、20 年以上前から問題とされてきた。データベース問い合わせ言語である SQL では、抽象や合成、入れ子のデータ構造等を直接表現できないが、静的型付き言語のようなプログラム言語ではこれらの表現を可能とし、データベース問い合わせ言語とプログラム言語を 1 つの言語として統合することでプログラマの負担を減らした効率の良いシステムが開発できると期待されている。1 つに統合された言語を統合言語クエリと呼び、代表的なものとして Microsoft LINQ 等が挙げられる [9]。LINQ 等の従来の統合言語クエリでは、様々なデータ構造を持つシステムに対して効率よくアクセスできるようになったが、複数のクエリ項の合成を許すことでクエリ雪崩 [1] という問題を発生させた。クエリ雪崩とは、コード変換によって多数のクエリを生成してしまうという問題であり、生成したクエリの数だけデータベースへアクセスするという点から実行効率を大幅に下げてしまう。

この問題に対し、Cheney らはクエリ項を直接実行せずに、生成したコードを書き換え規則に基づいて書き換えを行うことで単一の SQL クエリを生成し、クエリ雪崩を防ぐことに成功した [2]。一方で、書き換え規則において型付けが保存されるかは手動で証明されている。そのため、書き換え規則の追加や構文の拡張をする際には手動で証明しなければいけないため、拡張は困難である。

そこで、Suzuki らは、終タグレス法 [3] に基づき、型安全性をホスト言語の型付けに帰着させることで、書き換え規則の型付け保存性を自動的に検査されるようにした [4]。これにより、まだ拡張されていない構文に対して拡張可能であることを提唱した。しかし、SQL でのグループ化において、OCaml や Haskell のような関数型言語で表現しようとするグループ化前後でデータの型が保存されない [5] だけでなく、既存の構文の拡張では意味論や書き換え規則の等価性を表すことができないため、統合言語クエリのグループ化に対する拡張は困難とされてきた。

1.2 本研究の目的と構成

本研究では、SQL でのグループ化に対して新たな統合言語クエリの構文を設計し、その構文に対する型付け規則、意味論、書き換え規則等を提案する。変換の対象となる言語は Suzuki らの提案した Que Λ とし、先行研究と

同様にクエリの抽象や合成が可能で、クエリ雪崩を起こさない。本研究の貢献や従来の研究と異なる点としては以下の点が挙げられる。

- グループ化機能付き SQL の構文に直接対応するようなグループ化の構文を提案し、追加した書き換え規則を用いて書き換えを行うことで、単一の SQL クエリを生成することを可能にした。
- 複数のテーブルに対する解釈を変更し、従来の構文では書けなかった書き換え規則の記述を可能にした。
- ネストされた構造を一部許すことで、SQL においてサブクエリを使うことでしか表すことのできない構文に直接対応できるようにした。

本論文の構成は次の通りである。まず、第 2 章で前提知識として、コードの正規化とグループ化について概説する。第 3 章においては、本研究での対象言語 QueA の型付け規則や意味論、SQL 変換について記述する。第 4 章では、本研究において新たに提案するグループ化の構文 **gfor** をはじめ、グループ化に関する構文全般の型付け規則や意味論や正規形等について記述する。また、同時にグループ化において不可欠となる集約関数の働きについても記述する。第 5 章では、本研究のまとめと今後の課題について述べる。

第 2 章

前提知識

本章では、前提知識として、クエリの正規化の必要性とグループ化の難点について例題を用いて概説する。

2.1 正規化

統合言語クエリにおけるクエリ雪崩の問題を解決するために、Cheney ら [2] はクエリをコードとして生成した後、正規化規則によってコードを書き換えることで単一の SQL クエリを生成することに成功した。この書き換えを正規化と呼び、正規化の必要性を図 2.1 のテーブルを用いて述べる。

en_test テーブル				class テーブル	
date	sid	name	results	sid	class
2018-10-01	1	Sato	30	1	A
2018-10-01	2	Suzuki	50	2	B
2018-10-02	1	Sato	40	3	A
2018-10-02	2	Suzuki	45	4	B
2018-10-02	3	Takahashi	60	5	C
2018-10-03	3	Takahashi	85		

図 2.1 データテーブル

en_test テーブルは全生徒の英語の小テスト情報として、テスト受験日 (date), 生徒 ID(sid), 生徒名 (name), テスト結果 (results) を持つ。また、class テーブルは生徒のクラス情報として、生徒 ID(sid), クラス (class) を持つ。

データベース上のテーブルのデータ構造は、*Int*, *String* の基本的な型を持つレコードのバッグとして表現できる。例えば、図 2.1 の en_test テーブルの型を *EnTest* とすると、型 *EnTest* は以下のように表せる。なお、SQL には *Date* 型という日付の型を持つデータ型が存在するが、ここでは簡単のために日付は *String* 型で表すこととする。

$$EnTest = \langle date : String, sid : Int, name : String, results : Int \rangle$$

この型に基づき、図 2.1 のデータ構造は、図 2.2 のように記述できる。

```

edata : Bag EnTest =
[ < date = "2018-10-01", sid = 1, name = "Sato", results = 30 >;
  < date = "2018-10-01", sid = 2, name = "Suzuki", results = 50 >;
  < date = "2018-10-02", sid = 1, name = "Sato", results = 40 >;
  < date = "2018-10-02", sid = 2, name = "Suzuki", results = 45 >;
  < date = "2018-10-02", sid = 3, name = "Takahashi", results = 60 >;
  < date = "2018-10-03", sid = 3, name = "Takahashi", results = 85 > ]

cdata : Bag Class =
[ < sid = 1, class = "A" >;
  < sid = 2, class = "B" >;
  < sid = 3, class = "A" >;
  < sid = 4, class = "B" >;
  < sid = 5, class = "C" > ]

```

図 2.2 図 2.1 のデータ構造

ここで、2つのクエリを合成することで起きるクエリ雪崩と、その解決方法であるコードの正規化を説明するために、2つのクエリ (Q_1, Q_2) を用意する。クエリ Q_1 は、指定したクラスに該当する生徒を返す操作をするクエリであり、クエリ Q_2 は、指定された生徒 ID に該当し、かつ指定した日付の生徒の英語の点数を返す操作をするクエリである。内包表記を用いて記述したクエリとそれに対応する SQL クエリを図 2.3 に示す。

<p>(クエリ Q_1)</p> <pre> Q1 = λclass. for (c ← table("class")) where (c.class = class) yield c </pre>	<pre> SELECT c.* FROM class AS c WHERE c.class = class </pre>
<p>(クエリ Q_2)</p> <pre> Q2 = λc. λdate. for (e ← table("en_test")) where (e.sid = c.sid ∧ e.date = date) yield { date = e.date, name = e.name, results = e.results, class = c.class } </pre>	<pre> SELECT e.date AS date, e.name AS name, e.results AS results, c.class AS class FROM en_test AS e WHERE e.sid = c.sid AND e.date = date </pre>

図 2.3 クエリ Q_1, Q_2 とそれぞれに対応する SQL クエリ

内包表記 $\mathbf{for} (x \leftarrow M) N$ は、データベース上のテーブルに対応するバッグ M から要素を一つずつ取り出し、変数 x に束縛した環境の下で項 N を評価する。 \mathbf{table} 構成子はテーブル名を受け取り、そのテーブルをレコードのバッグとして返す働きをする。 $\mathbf{where} L M$ は、 L が \mathbf{true} であれば、 M を評価した結果を返し、それ以外であれば空のバッグを返す。 $\mathbf{yield} M$ は、 M を評価した結果を単一の要素とするバッグを返す。

ここで、2つのクエリを合成する関数 *compose* を用意する。クエリ Q_1 を実行して得られた N 件のデータに対して、クエリ Q_2 を N 回実行するような合成関数を仮に用意すると、 $N+1$ 件のクエリが発行されてしまうため、クエリ雪崩を起こす。このクエリ雪崩の現象を回避しつつ、クエリを合成する関数を以下に示す。

$$\begin{aligned} \text{compose} &= \lambda q. \lambda r. \lambda x. \lambda z. \\ &\quad \text{for } (y \leftarrow q x) r y z \end{aligned}$$

関数 *compose* を用いてクエリ Q_1, Q_2 を合成したクエリを Q_3 とする。しかし、 Q_3 はそのままでは直接 SQL に変換することができない。そのため、正規化を行うことで単一の SQL クエリに変換する。クエリ Q_3 と正規化後の Q_3 を図 2.4 に示す。

```

Q3 = λx. λz. compose Q1 Q2 x z
    = λx. λz.
      for (y ← (λclass.
                for (c ← table("class"))
                where (c.class = class)
                yield c) x)
        (λc. λdate.
          for (e ← table("en_test"))
          where (e.sid = c.sid ∧ e.date = date)
          yield { date = e.date,
                  name = e.name,
                  results = e.results,
                  class = c.class }) y z

```

↓ 正規化

```

Q3 = for (c ← table("class"))
     for (e ← table("en_test"))
     where (c.class = "A" ∧ e.sid = c.sid ∧ e.date = "2018-10-02")
     yield { date = e.date,
            name = e.name,
            results = e.results,
            class = c.class }

```

↓ SQL 変換

```

SELECT e.date AS date, e.name AS name,
         e.results AS results, c.class AS class
FROM en_test AS e, class AS c
WHERE c.class = "A" AND e.sid = c.sid AND e.date = "2018-10-02"

```

図 2.4 クエリ Q_3 の正規化と SQL 変換

このように、正規化を行うことで、クエリ雪崩を起こすことなく、単一の SQL クエリを生成することができる。また、正規化を前提とすることで、中間のデータ構造をネストすることができ、より見通しの良いクエリを記述することができる。

2.2 グループ化

SQLにおけるグループ化とは、データベース上のテーブルにおいて同列内で同じ値を持つデータをまとめることを言い、**GROUP BY** 句を用いる。前節の図 2.1 に対して、**date** でグループ化し、その日の英語の最高点を返す SQL クエリは以下のように書くことができる。

```
SELECT e.date AS date, MAX(e.results) AS max
FROM en_test AS e
GROUP BY e.date
```

一方で、上記のようなグループ化をする SQL クエリに直接対応する統合言語クエリは我々の知る限り存在しない。その理由として、以下の点が挙げられる。

1. 既存言語での構文の拡張が困難
2. 正規化できないクエリの組み合わせが存在

1 について、既存の統合言語クエリの構文において単純にグループ化機能のみを拡張すると、複数のテーブルに対してグループ化する場所次第で意味が異なってしまう。そのため、意味論の記述が困難となり、既存の構文にグループ化機能のみをそのまま拡張することは困難である。

2 について、例えばグループ化したクエリをさらにグループ化するという操作に対して、SQL では必ずサブクエリを用いて表現しなければならない。サブクエリとは、クエリ内部に入れ子になっているクエリのことであり、SQL でサブクエリを使用するときにはサブクエリによって作成された仮のテーブルを用いて、外側のクエリを評価する。ただ、サブクエリを用いることは実行効率を遅らせるとされているため、Cheney らの先行研究ではサブクエリ付きのクエリは、正規化によってフラットな構造のクエリに変換される。本研究ではグループ化機能によって表されるクエリのすべてを表現するため、グループ化の操作のみサブクエリを許す必要がある。また、項 **for** に単純にグループ化機能を付与した構文を **group** とすると、**for** と **group** の組み合わせに対して既存の書き換え規則の書き換えではどれも正規化できないため、異なるアプローチで書き換え規則を追加しなければならない。これらの観点からグループ化は困難とされ、本研究ではこれらの問題に対して、新たな提案方式を提唱し、グループ化機能を持つ統合言語クエリを設計する。

第 3 章

統合言語クエリ Que Λ

本研究では、Suzuki らが提案した Que Λ を対象言語とした変換を定義する。Que Λ は Cooper の言語 [6] からエフェクトを除いたものであり、Cheney らの T-LINQ からコード化の機能を除いたものに似ている。Que Λ の構文を図 3.1 に示す。

Syntax

Base types	$O ::= Int \mid Bool \mid String$	
Types	$A, B ::= O$	
	$A \rightarrow B$	(function type)
	$\mathbf{Bag} A$	(bag)
	$\langle l : \vec{A} \rangle$	(record)
Type environment	$\Gamma ::= \phi$	(empty environment)
	$\Gamma, x : A$	(term variable binding)
Primitives	\oplus	
Table	s, t	
Field names	l	
Terms	$L, M, N ::= \lambda x. M$	(λ -abstraction)
	$M N$	(application)
	$\oplus(\vec{M})$	(primitive)
	x	(variable)
	c	(constant)
	$\mathbf{for} (x \leftarrow M) N$	(comprehension)
	$\mathbf{where} L M$	(test)
	$\mathbf{yield} M$	(singleton)
	$[]$	(empty)
	$M \uplus N$	(bag union)
	$\mathbf{exists} M$	(not empty)
	$\mathbf{table}(t)$	(table reference)
	$\langle l = \vec{M} \rangle$	(record constructor)
	$L.l$	(projection)

図 3.1 Que Λ の構文

型は、 $Int, Bool, String$ からなる基本型 O と、関数型やバッグ、レコードからなる A, B がある。型環境は Γ で表し、拡張する場合はカンマの右側に新たな変数を加える。 \oplus はプリミティブな演算子を表し、テーブルやレコードラベルとして t, l がある。項はバッグ内包表記 **for** $(x \leftarrow M) N$ 、条件分岐 **where** $L M$ 、単一要素のバッグを返す **yield** M だけでなく、 λ 抽象や関数適用がある。また、バッグが空かどうかを判定するものは **exists** M で表し、データベースのテーブルを参照するものは **table**(t) として表す。レコードは $\langle \overrightarrow{l = \overrightarrow{M}} \rangle$ 、射影は $L.l$ と表す。

3.1 Que Λ の型付け規則

前節における Que Λ の構文の型付け規則を図 3.2 に表す。

$\Gamma \vdash M : A$		
$\frac{\text{CONST} \quad \Sigma(c) = O}{\Gamma \vdash c : O}$	$\frac{\text{OP} \quad \Sigma(\oplus) = O_1 \times \dots \times O_n \rightarrow O \quad \Gamma \vdash M_i : O_i \text{ (for each } 1 \leq i \leq n\text{)}}{\Gamma \vdash \oplus(\overrightarrow{M}) : O}$	
$\frac{\text{VAR} \quad x : A \in \Gamma}{\Gamma \vdash x : A}$	$\frac{\text{ABS} \quad \Gamma, x : A \vdash N : B}{\Gamma \vdash \lambda x. N : A \rightarrow B}$	$\frac{\text{APP} \quad \Gamma \vdash L : A \rightarrow B \quad \Gamma \vdash M : A}{\Gamma \vdash L M : B}$
$\frac{\text{RECORD} \quad \Gamma \vdash M_i : A_i \text{ (for each } 1 \leq i \leq n\text{)}}{\Gamma \vdash \langle \overrightarrow{l = \overrightarrow{M}} \rangle : \langle \overrightarrow{l : \overrightarrow{A}} \rangle}$	$\frac{\text{PROJECT} \quad \Gamma \vdash L : \langle \overrightarrow{l : \overrightarrow{A}} \rangle}{\Gamma \vdash L.l_i : A}$	$\frac{\text{SINGLETON} \quad \Gamma \vdash M : A}{\Gamma \vdash \mathbf{yield} M : \text{Bag } A}$
$\frac{\text{EMPTY}}{\Gamma \vdash [] : \text{Bag } A}$	$\frac{\text{UNIONALL} \quad \Gamma \vdash M : \text{Bag } A \quad \Gamma \vdash N : \text{Bag } A}{\Gamma \vdash M \uplus N : \text{Bag } A}$	$\frac{\text{TABLE} \quad \Sigma(t) = \text{Bag } \langle \overrightarrow{l : \overrightarrow{O}} \rangle}{\Gamma \vdash \mathbf{table}(t) : \text{Bag } \langle \overrightarrow{l : \overrightarrow{O}} \rangle}$
$\frac{\text{FOR} \quad \Gamma \vdash M : \text{Bag } A \quad \Gamma, x : A \vdash N : \text{Bag } B}{\Gamma \vdash \mathbf{for}(x \leftarrow M) N : \text{Bag } B}$	$\frac{\text{WHERE} \quad \Gamma \vdash L : Bool \quad \Gamma \vdash M : \text{Bag } A}{\Gamma \vdash \mathbf{where} L M : \text{Bag } A}$	$\frac{\text{EXISTS} \quad \Gamma \vdash M : \text{Bag } A}{\Gamma \vdash \mathbf{exists} M : Bool}$

図 3.2 Que Λ の型付け規則

CONST 規則は定数項 c の型が基本型 O であるならば、その定数項は型 O を持つという意味である。OP 規則はプリミティブ演算子の型がそれぞれ異なる型 $O_1 \times \dots \times O_n$ からある型 O への関数型であり、かつ型環境 Γ の下で n 個の項がそれぞれプリミティブ演算子の引数の型に対応する型を持つならば、 $\oplus(\overrightarrow{M})$ は型 O を持つという意味である。他にも、FOR 規則は M が型 $\text{Bag } A$ を持ち、型環境 Γ と型 A を持つ変数 x の下で N が型 $\text{Bag } B$ を持つならば、**for** $(x \leftarrow M) N$ は型 $\text{Bag } B$ を持つという意味である。

3.2 Que Λ の意味論

次に、Que Λ の値、評価文脈、操作的意味論を図 3.3 に定義する。

Value and Evaluation Contexts

Value	$V, W ::= c \mid \lambda x.M \mid \langle l = \vec{V} \rangle \mid [V]$
Evaluation Contexts	$\mathcal{E} ::= [] \mid \oplus(\vec{V}, \mathcal{E}, \vec{M}) \mid \mathcal{E} M \mid V \mathcal{E}$ $\mid \langle l = \vec{V}, l' = \mathcal{E}, l'' = \vec{M} \rangle \mid \mathcal{E}.l \mid \mathbf{yield} \mathcal{E}$ $\mid \mathcal{E} \uplus M \mid V \uplus \mathcal{E} \mid \mathbf{for} (x \leftarrow \mathcal{E}) N \mid \mathbf{exists} \mathcal{E}$ $\mid \mathbf{where} \mathcal{E} M$

Operational Semantics

$\oplus(\vec{V})$	\longrightarrow	$\sigma(\oplus, \vec{V})$	(E-OP)
$(\lambda x.N) V$	\longrightarrow	$N[x := V]$	(E-ABS- β)
$\langle l = \vec{V} \rangle.l_i$	\longrightarrow	V_i	(E-RECORD- β)
where true M	\longrightarrow	M	(E-WHERETRUE)
where false M	\longrightarrow	$[]$	(E-WHEREFALSE)
for $(x \leftarrow \mathbf{yield} V) M$	\rightsquigarrow	$M[x := V]$	(E-FORYIELD)
for $(x \leftarrow []) N$	\longrightarrow	$[]$	(E-FOREMPTY ₁)
for $(x \leftarrow L \uplus M) N$	\longrightarrow	$\mathbf{for} (x \leftarrow L) N \uplus \mathbf{for} (x \leftarrow M) N$	(E-FORUNION ₁)
exists $[]$	\longrightarrow	false	(E-FOREXISTFALSE)
exists $[\vec{V}]$	\longrightarrow	true ($ \vec{V} > 0$)	(E-EXISTTRUE)
$\frac{M \longrightarrow N}{\mathcal{E}[M] \longrightarrow \mathcal{E}[N]}$		(E-CONTEXT)	

図 3.3 Que Λ の値、評価文脈、操作的意味論

定数、 λ 抽象、レコード、レコードのバッグが値となる。 $[V]$ は、**yield** $V_1 \uplus \dots \uplus \mathbf{yield} V_n \uplus []$ の略記とし、左側の項から評価される。評価文脈では、図 3.1 で定義した Que Λ の構文に対してプリミティブ演算子、関数適用、バッグの結合が左側から評価される。操作的意味論では、 σ は各々のプリミティブ演算子に対する解釈であり、 σ によって意味論をパラメータ化している。E-FORUNION₁ は、レコードのバッグ同士を \uplus で結合してから評価を進める操作に対し、SQL クエリに直接変換するために、それぞれのレコードのバッグの評価をしてから \uplus によってレコードのバッグを結合する操作に書き換えている。これは正規化規則としても定義される。

3.3 SQL 変換

次に、Que Λ によって書かれたクエリを SQL に変換する対象となる SQL サブ言語と変換規則を図 3.4、図 3.5 に示す。SQL クエリでは、基本的にそれぞれの問い合わせ式を記述する **SELECT** 句、テーブルとその束縛を記述する **FROM** 句、条件式を記述する **WHERE** 句からなる。Que Λ において **for** $(x \leftarrow M) N$ を連続して書いたもの

は複数のテーブルに対する操作であるため、SQL クエリでは **FROM** 句にて並べて記述され、そこで複数のテーブルの内部結合が行われる。また、 $\text{Que}\Lambda$ におけるプリミティブ演算は SQL クエリでは \oplus_{sql} と書くこととし、変数 x のレコード全てに対応するものは、SQL において要素すべてという意味である $x.*$ と書く。図 3.5 での *where* の条件式は一部の項に対して事前に SQL 変換を行った結果を表しており、その結果を用いてクエリを変換する。

$$\begin{aligned}
P, Q & := P \text{ UNION ALL } Q \mid S \\
S & := \text{SELECT } \vec{s} \text{ FROM } t \text{ AS } x \text{ WHERE } e \\
s & := e \text{ AS } l \mid x.* \\
e & := c \mid x.l \mid e \wedge e' \mid \neg e \mid \text{EXISTS}(Q) \mid \oplus_{sql}(\vec{e})
\end{aligned}$$

図 3.4 SQL サブ言語

$$\begin{aligned}
S[U_1 \uplus U_2] & = S[U_1] \text{ UNION ALL } S[U_2] \\
S[[]] & = \text{SELECT } \overrightarrow{\text{null}} \text{ AS } \vec{l} \text{ FROM } \phi \text{ WHERE FALSE} \\
S[\text{for } (x \leftarrow \text{table}(s)) F] & = \text{SELECT } \overrightarrow{e} \text{ AS } \vec{l} \text{ FROM } s \text{ AS } x, t \text{ AS } \vec{y} \text{ WHERE } B \\
& \quad \text{where } S[F] = (\text{SELECT } \overrightarrow{e} \text{ AS } \vec{l} \text{ FROM } t \text{ AS } \vec{y} \text{ WHERE } B) \\
S[\text{where } B Z] & = \text{SELECT } \overrightarrow{e} \text{ AS } \vec{l} \text{ FROM } \vec{t} \text{ WHERE } S[B] \wedge B' \\
& \quad \text{where } S[Z] = (\text{SELECT } \overrightarrow{e} \text{ AS } \vec{l} \text{ FROM } \vec{t} \text{ WHERE } B') \\
S[\text{table}(s)] & = \text{SELECT } \overrightarrow{s.l} \text{ AS } \vec{l} \text{ FROM } s \text{ WHERE TRUE} \\
S[\text{yield } R] & = \text{SELECT } S[R] \text{ FROM } \phi \text{ WHERE TRUE} \\
S[\langle \vec{l} = \vec{B} \rangle] & = \overrightarrow{S[B]} \text{ AS } \vec{l} \\
S[\text{exists } U] & = \text{EXISTS } (S[U]) \\
S[\oplus(\vec{B})] & = \oplus_{sql}(S[\vec{B}]) \\
S[x.l] & = x.l \\
S[x] & = x.* \\
S[c] & = c
\end{aligned}$$

図 3.5 SQL 変換

第 4 章

グループ化機能の拡張

本研究で新たに提案する、グループ化機能を拡張した QueA の構文を図 4.1 に示す。

Grouping key	k	
Terms	$L, M, N ::= \dots$	
	$\mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow M}) L) N$	(grouping)
	$\mathbf{having} L M$	(grouped test)
	$\mathbf{gyield} M$	(grouped singleton)
	$\odot(M)$	(aggregation function)

図 4.1 追加したグループ化機能の構文

$\mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow M}) L) N$ は、グループ化の列を指定するキー \vec{k} に基づいて、グループ化を行う構成子である。 L の部分には **where** 項での条件式と同等の働きをするものが入り、グループ化をする機能を除けば、テーブルから要素を取り出し評価を進めるという点で **for** 項と同じ働きをする。評価の手順としては、 M から取り出した要素を y に束縛し、その環境の下で L が **true** のもののみ \vec{k} によってグループ化する。そして、グループ分けされたバッグのレコードを変数 z に束縛し、その環境の下で項 N を評価していく。ここで、**for** 項は複数のテーブルを並べる際、**for** 項を連続して並べることで表現していたが、**gfor** 項で同じように表現しようとする、正しく意味付けがされない。そこで、複数のテーブルを記述する際は $(y \leftarrow M)$ 部分で並べ、列として表現することで記述可能とした。項 N を評価した結果はレコードのバッグであり、それが **gfor** 項の結果となる。なお、**gfor** 項での変数 z は、SQL において登場しない変数となっているので SQL に変換するときには適切なレコードの変数として変換する。**having** $L M$ は、SQL での **HAVING** 句と同じ働きをし、基本的には **where** 項と同じ働きをするが、 L には集約関数によって集約された列のみを許し、集約されていない列は次節の型付け規則によって制約され、記述できないようになっている。**gyield** M も M には集約された列のみを許し、それ以外は **yield** 項と同じ働きをする。 \odot は集約関数を表し、SQL における値の平均値を算出する **AVG**、値の最大値を算出する **MAX** 関数や最小値を算出する **MIN** 関数等がある。

ここで、グループ化を行う SQL クエリと、それに対応するグループ化機能付き QueA のコードを図 4.2 に示す。

```
(SQL)
SELECT e.date AS date, c.class AS class, MAX(e.results) AS max
FROM en_test AS e, class AS c
WHERE e.sid = c.sid
GROUP BY c.class, e.date
```

```
(QueΛ)
gfor (z ← ((date, class) ((e ← table("en_test")), (c ← table("class")))) (e.sid = c.sid))
gyield ⟨date = z.date,
        class = z.class,
        max = max(z.results)⟩
```

図 4.2 グループ化を行う SQL クエリとそれに対応する QueΛ

図 4.2 では、2.1 節におけるテーブルに対してグループ化の処理を行っているが、ここで示すように、**gfor** の提案によって SQL クエリの **GROUP BY** だけでなく、テーブルが複数の時やクエリ中に **WHERE** 構文が存在する時にもサブクエリを使わずに直接対応するクエリの記述を書くことに成功した。

4.1 型付け規則の拡張

本節では新たに追加する **gfor** 等の型付け規則について記述していく。型付け規則を図 4.3 に示す。SQL において、**GROUP BY** 句でグループ化したクエリに対する **SELECT** 句や **HAVING** 句では、集約された列のみを許すという制約がある。そのため、新たに追加する型付け規則では、既存の型判断 $\Gamma \vdash M : A$ に加え、集約された項に対する型判断 $\Gamma \vdash_a M : A$ 、グループ化した後に記述できる構文に対する $\Gamma \vdash_G M : A$ を追加する。 Σ は定数 c 、プリミティブ演算子 \oplus 、集約関数 \odot からそれらの型への写像である。グループ化の手順としては、**gfor** ($z \leftarrow \vec{k} (\overrightarrow{y \leftarrow \vec{M}}) L$) N にて、最初に GFOR 規則が適用される。 \vec{k} は、 $k_i \in \vec{l}$ より、型 A を持つレコードのフィールド l_i のいずれかが指定される。項 M はレコードのバッグであり、その要素を変数 y に束縛し、その環境の下で項 L を評価する。項 L が **true** であれば、 \vec{k} によってグループ化し、変数 z に束縛する。 z はグループ化するためのキーと、同一のキーを持つ複数のレコードのバッグ $\langle \vec{k}, \text{Bag } A \rangle$ を持つバッグを型として持つ。そして、この環境の下で項 N を評価する。 N 項は型判断 $\Gamma \vdash_G M : A$ によって導出され、**HAVING**、**SINGLETONG**、**EMPTYG** のいずれかの規則が適用され、**having** $L M$ 、**gyield** M 、 $[\]$ のいずれかを記述することができる。

$$\boxed{\Gamma \vdash M : A}$$

GFOR

$$\frac{\Gamma \vdash M : \text{Bag } \langle \vec{l} : \vec{A} \rangle \quad \Gamma, y : \langle \vec{l} : \vec{A} \rangle \vdash L : \text{Bool} \quad \Gamma, z : \langle \vec{k}, \text{Bag } \langle \vec{l} : \vec{A} \rangle \rangle \vdash_G N : \text{Bag } B \quad k_i \in \vec{l} \text{ (for each } 1 \leq i \leq n)}{\Gamma \vdash \mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow M) L) N : \text{Bag } B}$$

$$\boxed{\Gamma \vdash_a M : A}$$

CONSTA

$$\frac{\Sigma(c) = O}{\Gamma \vdash_a c : O}$$

OPA

$$\frac{\Sigma(\oplus) = O_1 \times \dots \times O_n \rightarrow O \quad \Gamma \vdash_a M_i : O_i \text{ (for each } 1 \leq i \leq n)}{\Gamma \vdash_a \oplus(\vec{M}) : O}$$

RECORDA

$$\frac{\Gamma \vdash_a M_i : A_i \text{ (for each } 1 \leq i \leq n)}{\Gamma \vdash_a \langle \vec{l} = \vec{M} \rangle : \langle \vec{l} : \vec{A} \rangle}$$

PROJECTA

$$\frac{\Gamma \vdash_a L : \langle \vec{l} : \vec{A} \rangle}{\Gamma \vdash_a L.l_i : A_i}$$

PROJECTK

$$\frac{z : \langle \vec{k}, \text{Bag } \langle \vec{l} : \vec{A} \rangle \rangle \in \Gamma}{\Gamma \vdash_a z.k_i : A_i}$$

AGGREGATION

$$\frac{\Sigma(\odot) = O_1 \rightarrow O_2 \quad \Gamma \vdash M : O_1}{\Gamma \vdash_a \odot(M) : O_2}$$

$$\boxed{\Gamma \vdash_G M : A}$$

HAVING

$$\frac{\Gamma \vdash_a L : \text{Bool} \quad \Gamma \vdash_G M : \text{Bag } A}{\Gamma \vdash_G \mathbf{having} L M : \text{Bag } A}$$

SINGLETONG

$$\frac{\Gamma \vdash_a M : A}{\Gamma \vdash_G \mathbf{yield} M : \text{Bag } A}$$

EMPTYG

$$\frac{}{\Gamma \vdash_G [] : \text{Bag } A}$$

図 4.3 グループ化機能を追加した Que Λ の型付け規則

4.2 グループ化機能を追加した Que Λ の意味論

本節では、グループ化で必要不可欠な集約関数の動きについて記述した後、前節での型付け規則に基づいた **gfor** 等の評価文脈、操作的意味論について記述する。

4.2.1 集約関数

SQL において、**GROUP BY** 句を用いて最終的に返せる評価結果の一つが、集約関数を用いて集約された列である。集約関数には、前にも述べたように SQL における **AVG** 関数、**MAX** 関数、**MIN** 関数等が存在し、Que Λ では **avg**、**max**、**min** 等として記述する。本研究において、これらの集約関数の多くはそのまま表現することができるが、一部で例外が存在する。例えば、以下のようなクエリに対する場合である。

$$\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow M \uplus N) L) (\mathbf{avg}(S))$$

このクエリは、それぞれのテーブルを内部結合したテーブルに対して、 \vec{k} でグループ化した上で、 S の平均値

を求めるという操作である。このクエリを SQL で直接表現するには、以下のように \boxplus を外側に出し、最後にそれぞれのクエリを結合させる形にしなければならない。

$$\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow M) L) S' \boxplus \mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow N) L) S'$$

しかし、この変換において等価なクエリを表現しようとする、それぞれのクエリで合計値とデータの数を求め、2つのテーブルの合計値の和をデータの数の和で割ることでは平均値を求めることはできない。すなわち、変換の前後で異なる集約関数を使わなければならない。そこで本研究では、この変換を可能にするために $norm(\odot(S))$ といったマクロを用意し、正規化の前に必要に応じて集約関数を書き換える。しかし、中央値を求めるような集約関数は全データを必要とし、 $norm$ を用いても表すことができないため、本研究ではこのような統計的な集約関数を制限して、代数的な集約関数のみを使えるようにする。

4.2.2 操作的意味論の拡張

意味論において、グループ化に対する解釈を必要とするため、レコードのバッグに対し、グループ化直後のレコードのバッグを $[\langle \text{key} = k.., \text{body} = [x..] \rangle]$ として内部的に用意し、以下のような補助関数 ($sub_group, group, sub_eval_group, eval_group$) を与える。これらに対する定義を図 4.4 に示す。

sub_group

グループ化するためのキー (\vec{k})、レコード (x)、グループ化した途中段階のレコードのバッグ (N) を引数にとる。 N の要素に対して、 $x.k$ が一致すれば、グループ化後のバッグの $body$ 部分にレコード x を加え、一致しなければ新たなレコードとしてバッグに加える。

group

グループ化するためのキー (\vec{k})、レコードのバッグとその要素を束縛した変数 ($(y \leftarrow V)$)、条件式 (L) を引数にとる。 L が **true** であれば要素 x を sub_group に渡し、グループ化する。最終的にすべての要素に対してグループ化したバッグが返される。

sub_eval_group

グループ化したバッグ (G)、グループ化後評価する項 (M) を引数にとる。グループ化したバッグのそれぞれの要素に対して項 M を評価し、集約関数等を用いたフラットな構造のレコードのバッグが返される。

eval_group

グループ化するためのキー (\vec{k})、レコードのバッグとその要素を束縛した変数 ($(y \leftarrow V)$)、条件式 (L)、グループ化後評価する項 (M) を引数にとる。 $\vec{k}, (x \leftarrow V), L$ を $group$ に渡し、グループ化したバッグの要素に対して項 M を評価する。

```

sub_group( $\vec{k}$ , v, N) =
  N = [] のとき [ $\langle \text{key} = v.\vec{k}, \text{body} = [v] \rangle$ ]
  N = n :: N' のとき if n.key = v. $\vec{k}$  then [ $\langle \text{key} = n.\text{key}, \text{body} = v :: n.\text{body} \rangle$ ] @ N'
    else n :: sub_group( $\vec{k}$ , v, N')

group( $\vec{k}$ , (y ← V), L) =
  V = [] のとき []
  V = v :: V' のとき let gr = group( $\vec{k}$ , (y ← V'), L) in
    if L[v/y] then sub_group( $\vec{k}$ , v, gr)
    else gr

sub_eval_group(G, M) =
  G = [] のとき []
  G = g :: G' のとき M @ sub_eval_group(G', M)

eval_group( $\vec{k}$ , (y ← V), L, M) =
  let g = group( $\vec{k}$ , (y ← V), L) in sub_eval_group(g, M)

```

図 4.4 グループ化に対する新たな補助関数

新たに追加する評価文脈、操作的意味論を図 4.5 に示す。評価文脈には、集約関数、**gfor**、**having**、**gyield** が追加され、**gfor** においては、テーブルの要素を取り出してから、条件分岐をするという評価順序となっている。操作的意味論には、集約関数、**having**、**gfor** の意味論が追加されている。**gfor** に関して E-GFOR で記述されている *sub_eval_group*, *group*, *sub_group* については、図 4.4 のプログラム形式で書かれている補助関数を意味論の形式で定義し直したものである。集約関数は最初に前節の *norm* によって正規化され、その環境の下で **gfor** 項の意味論が定義されている。また、ここでの $M[g.\text{body}/z]$ では、項 M において z が存在すれば、 $g.\text{body}$ を代入する操作を行っている。例えば、以下のような全クラスのテスト結果の平均点を求める M と、それぞれのテスト結果、クラス情報を持つ $g.\text{body}$ に対し、 $M[g.\text{body}/z]$ では項 M の z に $g.\text{body}$ を代入し、全てのテストの点数の合計点をクラス数で割ることで平均点を求める操作をしている。ここでは集約関数が正規化されているため、平均点を求める操作が **sum** 関数と **count** 関数で表されている。

```

M = [ $\langle \text{avg} = \text{sum}(z.\text{results})/\text{count}(z.\text{class}) \rangle$ ]
g.body = [ $\langle \text{results} = 40, \text{class} = \text{"A"} \rangle; \langle \text{results} = 60, \text{class} = \text{"B"} \rangle$ ]
M[g.body/z] = [ $\langle \text{avg} = \text{sum}([40; 60])/\text{count}([\text{"A"}; \text{"B"}]) \rangle$ ]
              = [ $\langle \text{avg} = 50 \rangle$ ]

```

Value and Evaluation Contexts

Evaluation Context $\mathcal{E} ::= \dots$

| $\odot(\vec{V}, \mathcal{E}, \vec{M})$ | **gfor** $(z \leftarrow \vec{k} ((\overline{y_1 \leftarrow \vec{V}}), (y'_1 \leftarrow \mathcal{E}), (\overline{y''_1 \leftarrow \vec{L}})) M) N$

| **gfor** $(z \leftarrow \vec{k} (\overline{y_1 \leftarrow \vec{V}}) \mathcal{E}) N$ | **having** $\mathcal{E} M$ | **gyield** \mathcal{E}

Operational Semantics

...

$\odot(\vec{V})$	\longrightarrow	$\sigma(\odot, \vec{V})$	(E-AGGREGATION)
having true M	\longrightarrow	M	(E-HAVINGTRUE)
having false M	\longrightarrow	$[]$	(E-HAVINGFALSE)
gfor $(z \leftarrow \vec{k} (y \leftarrow []) L) N$	\longrightarrow	$[]$	(E-GFOREMPTY ₁)
gfor $(z \leftarrow \vec{k} (y \leftarrow M) L) N$	\longrightarrow	$sub_eval_group(z \leftarrow group(\vec{k}, (y \leftarrow V), L), M)$	(E-GFOR)

① $sub_eval_group((z \leftarrow G), M) \rightarrow []$ (if $G = []$)

$sub_eval_group((z \leftarrow G), M) \rightarrow$

$M[g.body/z] :: sub_eval_group((z \leftarrow G'), M)$ (if $G = g :: G'$)

② $group(\vec{k}, (y \leftarrow V), L) \rightarrow []$ (if $V = []$)

$group(\vec{k}, (y \leftarrow V), L) \rightarrow sub_group(\vec{k}, v, group(k, (y \leftarrow V), L))$

(if $V = v :: V' \wedge L[v/y] \rightarrow \mathbf{true}$)

$group(\vec{k}, (y \leftarrow V), L) \rightarrow group(\vec{k}, (y \leftarrow V'), L)$

(if $V = v :: V' \wedge L[v/y] \rightarrow \mathbf{false}$)

③ $sub_group(\vec{k}, v, N) \rightarrow [\langle \text{key} = v.\vec{k}, \text{body} = [v] \rangle]$ (if $N = []$)

$sub_group(\vec{k}, v, N) \rightarrow [\langle \text{key} = n.\text{key}, \text{body} = v :: n.\text{body} \rangle] @ N'$

(if $N = n :: N' \wedge n.\text{key} = v.\vec{k}$)

$sub_group(\vec{k}, v, N) \rightarrow n :: sub_group(\vec{k}, v, N')$

(if $N = n :: N' \wedge n.\text{key} \neq v.\vec{k}$)

図 4.5 グループ化機能を追加した QueA における評価文脈, 操作的意味論の定義

4.3 拡張する正規形と正規化規則

グループ化機能を拡張するために $\text{Que}\Lambda$ にいくつかの構成子を新たに追加した。本節では $\text{Que}\Lambda$ によって書かれたグループ化のクエリを SQL に変換するための正規形と正規化規則を記述する。

4.3.1 正規形

SQL の構文と一対一に対応する統合言語クエリの構文を正規形と呼び、既存の正規形と新たに追加する正規形を図 4.6 に示す。

Queries	U	$::=$	$U_1 \uplus U_2 \mid [] \mid F \mid F_g$
Comprehension	F	$::=$	for ($x \leftarrow H$) $F \mid Z$
	F_g	$::=$	gfor ($z \leftarrow \overrightarrow{K_g} (\overrightarrow{y_1 \leftarrow H}) B$) Z_g
Body	H	$::=$	$F_g \mid \text{table}(t)$
	Z	$::=$	where $B Z \mid \text{yield } R$
Record	Z_g	$::=$	having $B_g Z_g \mid \text{gyield } R_g$
	R	$::=$	$\langle \overrightarrow{l = B} \rangle \mid x$
Primitives	R_g	$::=$	$\langle \overrightarrow{l = B_g} \rangle$
	K_g	$::=$	k
Aggregation	B	$::=$	exists $U \mid \oplus(\overrightarrow{B}) \mid x.l \mid y.l \mid c \mid z.l$
	B_g	$::=$	$\odot(\overrightarrow{B}) \mid c \mid z.k \mid \oplus(\overrightarrow{B_g})$

図 4.6 グループ化機能を追加した $\text{Que}\Lambda$ の正規形

グループ化をすることが先行研究と異なる点は、ネストされた構文に対して、コードを最適化することができない場合があるという点である。これまでの先行研究では、**for** 項のテーブル部分がネストされた構造になっていれば、正規化規則によりネストがないフラットな構造にコードを書き換えられたのに対し、例えば、グループ化したものをさらにグループ化するという操作においては、ネストされた構造と同じ意味を持つフラットな構造な SQL クエリは存在しない。そのため、ネストされたものが最小のクエリとなり、グループ化に対し、ネストされた構造を正規形として許さなければならない。図 4.6 では、**gfor** 項と **table** 項を内包表記 H にまとめ、**for** 項、**gfor** 項のテーブル部分に書けるようにした。**gfor** 項を用いて複数のテーブルをグループ化するクエリについては、テーブルの要素を変数に束縛する操作を $\overrightarrow{y \leftarrow H}$ 部分に列として並べて表現するため、グループ化付き SQL クエリに対応する **gfor** 項は常に一つとなる。本研究では、**for** 項と **gfor** 項をどのように対応付けするかが肝になっている。そこで、**for** 項と **gfor** 項の関係をわかりやすく記述したものを表 4.1 に示す。**for** 項同士の組み合わせは先行研究 [2] における正規形と正規化規則を用いて表現することができる。**gfor** 項と **for** 項が連続して書かれるクエリは、グループ化した途中段階のテーブルのそれぞれの要素に対して、**for** 項を評価するという意味になる。グループ化した途中段階のテーブルに対する操作というのは、前節で記述したような **key** と **body** を持つレコードのバッグに対しての操作ということであり、型が合わないだけでなく、記述されたクエリに対して適切な意味付

けがされないことからそういったクエリを記述することはできない。**gfor** 項と **gfor** 項が連続して書かれるクエリに対しても同様の理由で適切な意味付けをすることができない。そのため、**gfor** のボディー部分に記述することができる項は **having** 項と **gyield** 項の 2 つのみであり、これを Z_g として表記する。**having** 項の条件式、もしくは **gyield** 項のボディーとして書けるものは定数、グループ化したキー、集約関数を用いて評価された値であり、これらを R_g として表記する。本研究では表 4.1 内の下の二つに対して注目し、次項で述べる正規化規則を用いることで正規化することに成功した。

表 4.1 **for** 項と **gfor** 項の関係

for 項と gfor 項の組み合わせ	意味
for ($x \leftarrow \text{for} \dots$)	(先行研究での) 正規化規則により書き換え可
for ($x \leftarrow M$) for ...	正規形
gfor ($z \leftarrow \dots$) for ...	適切な意味付けがされない
gfor ($z \leftarrow \dots$) gfor ...	適切な意味付けがされない
for ($x \leftarrow \text{gfor} \dots$)	サブクエリとして記述
gfor ($z \leftarrow (\vec{k} (y \leftarrow \text{gfor} \dots))$)	サブクエリとして記述
gfor ($z \leftarrow (\vec{k} (y \leftarrow \text{for} \dots))$)	(本研究での) 正規化規則により書き換え可
for ($x \leftarrow M$) gfor ...	(本研究での) 正規化規則により書き換え可

4.3.2 正規化規則

新たに追加する正規化規則を図 4.7 に示す。参考のため、全体の正規化規則は付録 A に記載する。

Stage 1 では、ネストされている中間のデータ構造に対して、グループ化せずに書かれているクエリがあれば、フラットなデータ構造に書き換える操作をしている。GFORFOR 規則、GFORWHERE 規則、GFORYIELD 規則では、ネストされている **for** 項、**where** 項、**yield** 項に対して、**gfor** 項のテーブル、条件式、ボディーの一部として書き換えることでフラットなデータ構造にしている。本研究では、**gfor** 項のテーブル部分で、レコードを束縛する処理を列として並べて記述したことにより、複数のテーブルに対する処理を可能にした。この処理によって GFORFOR 規則の書き換えは可能になっており、条件式を **gfor** 項の中に埋め込んだことで、GFORWHERE 規則の書き換えを可能にした。

Stage 2 では、クエリを意味的に同じ SQL クエリに直接変換できるように並べかえる操作を行なっている。FORFOR 規則は、**for** 項のレコードそれぞれに対して **gfor** 項の要素をグループ化するという意味であり、これを SQL に直接変換できるように書き換えると、2 つのテーブルそれぞれの要素に対し、 \vec{k}, x のレコード全体でグループ化するというものと同義である。故に一つの **gfor** 項にまとめることができ、同じ意味を持つ SQL クエリに直接変換できる。しかしこのとき、レコード y とレコード x のフィールド名に同一のものが存在していれば、判別するためにフィールド名を書き換えなければならない。そのため、正規化によってフィールド名を書き換える関数 f_{re} を以下のように用意する。

$$f_{re} x = \langle \vec{l}^i = x.\vec{l}^i \rangle$$

関数 f_{re} ではレコード x を引数にとり、レコードのそれぞれのフィールドにダッシュを付け、そのレコードを

返す操作をしている。この操作により同じフィールド名に対しての判別を可能にしている。また、関数 $FN(m)$ ではレコード m のフィールド名を取り出す操作をしている。WHEREFOR 規則において、グループ化したバグの要素に対し、条件式で制御するというのは、グループ化後に **having** を使うことと同義であるため、図 4.7 のように書き換えることができる。

(Stage 1)

$$\begin{aligned}
\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow \mathbf{for} (x \leftarrow M) N) L) O &\rightsquigarrow \\
&\mathbf{gfor} (z \leftarrow \vec{k} ((x \leftarrow M), (y \leftarrow N)) L) O && \text{(GFORFOR)} \\
\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow \mathbf{where} M N) L) O &\rightsquigarrow \\
&\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow N) (L \wedge M)) O && \text{(GFORWHERE)} \\
\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow \mathbf{yield} M) L) N &\rightsquigarrow \mathbf{where} L (N[y := M]) && \text{(GFORYIELD)} \\
\mathbf{gfor} (z \leftarrow \vec{k} (y \leftarrow []) L) N &\rightsquigarrow [] && \text{(GFOREMPTY}_1\text{)} \\
\mathbf{having} \mathbf{true} M &\rightsquigarrow M && \text{(HAVINGTRUE)} \\
\mathbf{having} \mathbf{false} M &\rightsquigarrow [] && \text{(HAVINGFALSE)}
\end{aligned}$$

(Stage 2)

$$\begin{aligned}
\mathbf{for} (x \leftarrow M) \mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow N}) L) O &\leftrightarrow \\
&\mathbf{gfor} (z \leftarrow (\vec{k}, \vec{l}') ((x \leftarrow M), (\overrightarrow{y \leftarrow N})) L) O && \\
&\quad (\mathit{where} m' = \mathit{fre}(m :: M), l' = FN(m')) && \text{(FORGFOR)} \\
\mathbf{where} L \mathbf{gfor} (\vec{z} \leftarrow \vec{k} (\overrightarrow{y \leftarrow N}) M) O &\leftrightarrow \\
&\mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow N}) M) \mathbf{having} L O && \text{(WHEREGFOR)} \\
\mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow M}) L) [] &\leftrightarrow [] && \text{(GFOREMPTY}_2\text{)} \\
&\mathbf{having} L [] &\leftrightarrow [] && \text{(HAVINGEMPTY)} \\
\mathbf{having} L (\mathbf{having} M N) &\leftrightarrow \mathbf{having} (L \wedge M) N && \text{(HAVINGHAVING)}
\end{aligned}$$

図 4.7 追加する正規化規則

4.4 SQL 変換の拡張

図 4.6 の正規形に対応する SQL サブ言語と、SQL へ変換するための変換規則を図 4.8、図 4.9 に示す。SQL における **GROUP BY** 句, **HAVING** 句を追加しただけでなく、グループ化のキーを表す $y.k$ や集約関数を表す \odot_{sql} を追加した。Que Λ において、**gfor** 項における変数 z は、元のレコードをグループ化した中間的なレコードであり、変数 z に直接対応する SQL クエリは存在しないため、 z が扱うフィールド名によって変数を変更して変換している。

$$\begin{aligned}
P, Q &:= P \text{ UNION ALL } Q \mid S \\
S &:= \text{SELECT } \vec{s} \text{ FROM } t \text{ AS } x \text{ WHERE } e \\
T &:= \text{SELECT } \vec{u} \text{ FROM } t \text{ AS } y \text{ WHERE } e \text{ GROUP BY } \vec{k} \text{ HAVING } f \\
s &:= e \text{ AS } l \mid x.* \\
u &:= f \text{ AS } l \\
e &:= c \mid x.l \mid e \wedge e' \mid \neg e \mid \text{EXISTS}(Q) \mid \oplus_{sql}(\vec{e}) \mid y.l \\
f &:= c \mid y.k \mid f \wedge f' \mid \neg f \mid \oplus_{sql}(\vec{f}) \mid \odot_{sql}(\vec{e})
\end{aligned}$$

図 4.8 グループ化機能を追加した SQL サブ言語

$$\begin{aligned}
&\dots \\
\mathbb{S}[\text{gfor } (z \leftarrow \vec{K}_g (y \leftarrow \vec{H}) B) Z_g] &= \text{SELECT } \vec{d} \text{ AS } \vec{l} \text{ FROM } (\overline{S[H]}) \text{ AS } \vec{y}, t \text{ AS } \vec{h} \text{ WHERE } S[B] \text{ AND } B' \\
&\quad \text{GROUP BY } S[\vec{K}_g], \vec{k} \text{ HAVING } B_g \\
\text{where } S[Z_g] &= (\text{SELECT } \vec{d} \text{ AS } \vec{l} \text{ FROM } t \text{ AS } \vec{h} \text{ WHERE } B' \text{ GROUP BY } \vec{k} \text{ HAVING } B_g) \\
\mathbb{S}[\text{having } B_g Z_g] &= \text{SELECT } \vec{d} \text{ AS } \vec{l} \text{ FROM } \vec{t} \text{ WHERE } B \\
&\quad \text{GROUP BY } \vec{k} \text{ HAVING } S[B_g] \text{ AND } B'_g \\
\text{where } S[Z_g] &= (\text{SELECT } \vec{d} \text{ AS } \vec{l} \text{ FROM } \vec{t} \text{ WHERE } B \text{ GROUP BY } \vec{k} \text{ HAVING } B'_g) \\
\mathbb{S}[\text{yield } R_g] &= \text{SELECT } S[R_g] \text{ FROM } \phi \text{ WHERE TRUE} \\
&\quad \text{GROUP BY } \phi \text{ HAVING TRUE} \\
\mathbb{S}[\langle \vec{l} = \vec{B}_g \rangle] &= \overline{S[B_g]} \text{ AS } \vec{l} \\
\mathbb{S}[\oplus(\vec{B}_g)] &= \oplus_{sql}(S[\vec{B}_g]) \\
\mathbb{S}[\odot(\vec{B})] &= \odot_{sql}(S[\vec{B}]) \\
\mathbb{S}[y.l] &= y.l \\
\mathbb{S}[z.l] &= y.l \text{ where } l = FN(y) \\
\mathbb{S}[z.k] &= y.k \\
\mathbb{S}[\vec{k}] &= \vec{k}
\end{aligned}$$

図 4.9 SQL 変換

4.5 正規化、SQL 変換の例

前節までの提案に対し、正規化、SQL 変換がどのように行われるかを例題を見ながら概説していく。ここでは、2.1 節のテーブルに対し、2つのクエリ (Q_1, Q_4) を用意し、*compose* 関数によって2つのクエリを合成し、正規化によって単一の SQL クエリに変換するという操作を行う。クエリ Q_1 は、2.1 節におけるクエリと同じものであり、class テーブルから与えられたクラスを取り出す操作をするクエリである。クエリ Q_4 は、en_test テーブルから Q_4 の生徒と同一の生徒を取り出し、その生徒の英語の平均点を求める操作をするクエリである。新たなクエリ Q_4 を以下に示す。

```
 $Q_4 = \lambda c.$   
  gfor ( $z \leftarrow ((\text{sid}, \text{class}) (e \leftarrow \text{table}(\text{"en\_test"})) (e.\text{sid} = c.\text{sid}))$ )  
  gyield  $\langle \text{sid} = z.\text{sid},$   
            $\text{class} = c.\text{class},$   
            $\text{avg} = \text{avg}(z.\text{results}) \rangle$ 
```

compose 関数を用いて2つのクエリ Q_1, Q_4 を合成したクエリを Q_5 とし、展開したものを以下に示す。

```
 $Q_5 = \lambda x.$   
  for ( $y \leftarrow (\lambda \text{class}.$   
    for ( $c \leftarrow \text{table}(\text{"class"})$ )  
    where ( $c.\text{class} = \text{class}$ )  
    yield  $c$ )  $x$ )  
  ( $\lambda c.$   
  gfor ( $z \leftarrow ((\text{sid}, \text{class}) (e \leftarrow \text{table}(\text{"en\_test"})) (e.\text{sid} = c.\text{sid}))$ )  
  gyield  $\langle \text{sid} = z.\text{sid},$   
          $\text{class} = c.\text{class},$   
          $\text{avg} = \text{avg}(z.\text{results}) \rangle$ )  $y$ 
```

まず、λ 抽象を簡約化するために、 Q_5 に “A” を適用する。

```
 $Q_5 = \text{for}$  ( $y \leftarrow (\text{for}$  ( $c \leftarrow \text{table}(\text{"class"})$ )  
  where ( $c.\text{class} = \text{"A"}$ )  
  yield  $c$ )  
  ( $\lambda c.$   
  gfor ( $z \leftarrow ((\text{sid}, \text{class}) (e \leftarrow \text{table}(\text{"en\_test"})) (e.\text{sid} = c.\text{sid}))$ )  
  gyield  $\langle \text{sid} = z.\text{sid},$   
          $\text{class} = c.\text{class},$   
          $\text{avg} = \text{avg}(z.\text{results}) \rangle$ )  $y$ 
```

次に、Stage 1 における FORFOR 規則、FORWHERE 規則、FORYIELD 規則を適用して、フラットなデータ構造に書き換える。

```
 $Q_5 = \text{for}$  ( $c \leftarrow \text{table}(\text{"class"})$ )  
  where ( $c.\text{class} = \text{"A"}$ )  
  gfor ( $z \leftarrow ((\text{sid}, \text{class}) (e \leftarrow \text{table}(\text{"en\_test"})) (e.\text{sid} = c.\text{sid}))$ )  
  gyield  $\langle \text{sid} = z.\text{sid},$   
          $\text{class} = c.\text{class},$   
          $\text{avg} = \text{avg}(z.\text{results}) \rangle$ 
```

次に、SQL に変換するために **Stage 2** の正規化規則を用いて正規化を行っていく。
 まず、WHEREGFOR 規則を適用する。

$$Q_5 = \text{for } (c \leftarrow \text{table}(\text{"class"})) \\
 \text{gfor } (z \leftarrow ((\text{sid}, \text{class}) (e \leftarrow \text{table}(\text{"en_test"})) (e.\text{sid} = c_1.\text{sid})) \\
 \text{having } (z.\text{class} = \text{"A"}) \\
 \text{gyield } \langle \text{sid} = z.\text{sid}, \\
 \text{class} = c.\text{class}, \\
 \text{avg} = \text{avg}(z.\text{results}) \rangle$$

次に、FORGFOR 規則を適用する。

$$Q_5 = \text{gfor } (z \leftarrow ((\text{sid}, \text{class}, \text{sid}', \text{class}') ((e \leftarrow \text{table}(\text{"en_test"})), (c \leftarrow \text{table}(\text{"class"})))) (e.\text{sid} = c.\text{sid})) \\
 \text{having } (z.\text{class} = \text{"A"}) \\
 \text{gyield } \langle \text{sid} = z.\text{sid}, \\
 \text{class} = z.\text{class}, \\
 \text{avg} = \text{avg}(z.\text{results}) \rangle$$

上記のように正規項が得られたので、これを SQL クエリに変換していく。

SQL 変換では、 $\text{gfor } (z \leftarrow (\vec{k} (\overrightarrow{y \leftarrow M}) L) N$ 項において、先に項 N から変換される。よって、最初に **gyield** 項を変換する。

```
SELECT e.sid AS sid, c.class AS class, AVG(e.results) AS avg
FROM  $\phi$ 
WHERE TRUE
GROUP BY  $\phi$ 
HAVING TRUE
```

次に、上記の結果を持って **having** 項を変換する。

```
SELECT e.sid AS sid, c.class AS class, AVG(e.results) AS avg
FROM  $\phi$ 
WHERE TRUE
GROUP BY  $\phi$ 
HAVING TRUE AND c.class = "A"
```

最後に、上記の結果を持って **gfor** 項を変換する。

```
SELECT e.sid AS sid, c.class AS class, AVG(e.results) AS avg
FROM en_test AS e, class AS c
WHERE TRUE AND e.sid = c.sid
GROUP BY e.sid, c.class, c
HAVING TRUE AND c.class = "A"
```

よって最終的な SQL クエリとして上記のクエリが得られ、グループ化機能を持つ $Q_{ue\Lambda}$ のクエリを正規化によって単一の SQL クエリに変換することができた。

第 5 章

結論

5.1 まとめ

本研究では、対象言語 `Que Λ` に対して、`gfor` や `having`、`gyield` などの構文を追加したことで、グループ化機能を持つ統合言語クエリを設計した。それぞれのレコードを変数に束縛する処理を列として並べ、複数のテーブルを表現するようにならただけでなく、グループ化の前に条件分岐を組み込むことで、グループ化機能付き SQL に直接対応する統合言語クエリの記述を可能にした。また、SQL における、グループ化の `UNION ALL` について、代数的に使える集約関数しか表現できないことを発見し、`AVG` 関数のような集約関数を書き換えることで正規化において意味を変えない等価なクエリとして表現することを可能にした。

本研究における正規化規則は、先行研究とは異なるアプローチにより追加されている。先行研究における正規化規則ではネストされた構造のクエリや SQL に直接変換できないクエリをただ並べ替えるような正規化規則を追加することで正規化を可能としてきた。しかし、本研究では並べ替えるだけでは適切な正規化がされず、4.3.2 における `GFORFOR` 規則や `FORGFOR` 規則のように複数のテーブルの処理やグループ化キーを増やす等、先行研究とは異なるやり方で正規化規則を追加することで正規化を可能とした。

5.2 今後の課題

本研究ではグループ化機能を持つ統合言語クエリを設計したが、処理系を用いた正規化規則等の型付け保存性はチェックされていない。これらの型付け保存性を証明するために、先行研究では OCaml での終タグレス法による実装を行っていたが [4]、グループ化後に生成されるレコードのバグや集約関数の組み込みなどグループ化によって型が変わってしまうことを考えると、終タグレス法では表しきれないと考えられる。よって、別のアプローチを用いて本研究で提案した統合言語クエリの正しさを証明することが今後の課題としてあげられる。

謝辞

本研究を進めるにあたり、様々なご指導をしてくださった亀山幸義教授および海野広志准教授に深く感謝いたします。また、多くのご助言を下されたプログラム論理研究室の皆様感謝を申し上げます。

参考文献

- [1] Torsten Grust, Jan Rittinger, and Tom Schreiber. Avalanche-safe linq compilation. In *Proceedings of the VLDB Endowment*, Volume 3 Issue 1-2, pp. 162–172, September 2010.
- [2] James Cheney, Sam Lindley, and Philip Wadler. A practical theory of language-integrated query. In *Proceedings of 18th ACM SIGPLAN international conference on Functional programming*, ICFP 2013, pp. 36–51, New York, NY, USA, 2013. ACM.
- [3] Jacques Carette, Oleg Kiselyov, and Chung-chieh Shan. Finally tagless, partially evaluated: Tagless staged interpreters for simpler typed languages. In *Journal of Functional Programming*, Volume 19 Issue 5, pp. 509–543, New York, NY, USA, September 2009.
- [4] Kenichi Suzuki, Oleg Kiselyov, and Yuki-yoshi Kameyama. Finally, safely-extensible and efficient language-integrated query. In *Proceedings of the 2016 ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation*, PEPM '16, pp. 37–48, New York, NY, USA, 2016. ACM.
- [5] Simon Jones, Peyton and Philip Wadler. Comprehensive comprehensions comprehensions with 'order by' and 'group by'. In *Proceedings of the ACM SIGPLAN workshop on Haskell workshop*, Haskell '07, pp. 61–72, New York, NY, USA, 2007. ACM.
- [6] Ezra Cooper. The script-writer's dream: How to write great sql in your own language, and be sure it will succeed. In *Proceedings of the 12th International Symposium on Database Programming Languages*, DBPL '09, pp. 36–51, Berlin, Heidelberg, 2009. Springer-Verla.
- [7] Tatsuya Katsusima and Oleg Kiselyov. Sound and efficient language-integrated query. In *Asian Symposium on Programming Languages and Systems*, APLAS 2017, pp. 364–383, LNCS 10695, 2017. Springer.
- [8] James Cheney, Sam Lindley, and Philip Wadler. Query shredding: Efficient relational evaluation of queries over nested multisets. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*, SIGMOD '14, pp. 1027–1038, New York, NY, USA, 2014. ACM.
- [9] Microsoft corporation. the linq project: .net language integrated query. White paper, September 2005.

付録 A

正規化規則

付録では、Suzuki らが対象言語を QueA として Cheney らの正規化規則を記述した正規化規則に加え、本研究で追加した正規化規則を図 A.1、図 A.2 に示す。なお、ネストされた構造のクエリをフラットな構造のクエリへ変換する規則を **Stage 1** として記述し、 $M \rightsquigarrow N$ と表記する。また、SQL に直接変換できるように並べ替える規則を **Stage 2** として記述し、 $M \leftrightarrow N$ と表記する。

Stage 1

$(\lambda x.N) M$	\rightsquigarrow	$N[x := M]$	(ABS- β)
$\overrightarrow{\langle l = \overline{M} \rangle}.l_i$	\rightsquigarrow	M_i	(RECORD- β)
for $(x \leftarrow \mathbf{yield} M) N$	\rightsquigarrow	$N[x := M]$	(FORYIELD)
for $(x \leftarrow \mathbf{for} (y \leftarrow L) M) N$	\rightsquigarrow		
for $(y \leftarrow L) \mathbf{for} (x \leftarrow M) N$	\rightsquigarrow	(if $y \notin FV(N)$)	(FORFOR)
for $(x \leftarrow \mathbf{where} L M) N$	\rightsquigarrow	where $L (\mathbf{for} (x \leftarrow M) N)$	(FORWHERE)
for $(x \leftarrow []) N$	\rightsquigarrow	$[]$	(FOREMPTY ₁)
for $(x \leftarrow M_1 \uplus M_2) N$	\rightsquigarrow		
		for $(x \leftarrow M_1) N \uplus \mathbf{for} (x \leftarrow M_2) N$	(FORUNIONALL ₁)
where true M	\rightsquigarrow	M	(WHERETRUE)
where false M	\rightsquigarrow	$[]$	(WHEREFALSE)
gfor $(z \leftarrow \overrightarrow{k} (y \leftarrow \mathbf{for} (x \leftarrow M) N) L) O$	\rightsquigarrow		
		gfor $(z \leftarrow \overrightarrow{k} ((y \leftarrow N), (x \leftarrow M)) L) O$	(GFORFOR)
gfor $(z \leftarrow \overrightarrow{k} (y \leftarrow \mathbf{where} M N) L) O$	\rightsquigarrow		
		gfor $(z \leftarrow \overrightarrow{k} (y \leftarrow N) (L \wedge M)) O$	(GFORWHERE)
gfor $(z \leftarrow \overrightarrow{k} (y \leftarrow \mathbf{yield} M) L) N$	\rightsquigarrow	$N[y := M]$	(GFORYIELD)
gfor $(z \leftarrow \overrightarrow{k} (y \leftarrow []) L) N$	\rightsquigarrow	$[]$	(GFOREMPTY ₁)
having true M	\rightsquigarrow	M	(HAVINGTRUE)
having false M	\rightsquigarrow	$[]$	(HAVINGFALSE)

図 A.1 正規化規則 (Stage 1)

Stage 2

$$\begin{aligned}
\mathbf{for} (x \leftarrow M) (N_1 \uplus N_2) &\leftrightarrow \mathbf{for} (x \leftarrow M) N_1 \uplus \mathbf{for} (x \leftarrow M) N_2 && (\text{FORUNIONALL}_2) \\
\mathbf{for} (x \leftarrow M) [] &\leftrightarrow [] && (\text{FOREMPTY}_2) \\
\mathbf{where} L (M \uplus N) &\leftrightarrow (\mathbf{where} L M) \uplus (\mathbf{where} L N) && (\text{WHEREUNION}) \\
\mathbf{where} L [] &\leftrightarrow [] && (\text{WHEREEMPTY}) \\
\mathbf{where} L (\mathbf{where} M N) &\leftrightarrow \mathbf{where} (L \wedge M) N && (\text{WHEREWHERE}) \\
\mathbf{where} L (\mathbf{for} (x \leftarrow M) N) &\leftrightarrow \mathbf{for} (x \leftarrow M) (\mathbf{where} L N) && (\text{WHEREFOR}) \\
\mathbf{for} (x \leftarrow M) \mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow \vec{N}}) L) O &\leftrightarrow \mathbf{gfor} (z \leftarrow (\vec{k}, x_1) ((\overrightarrow{y \leftarrow \vec{N}}), (x \leftarrow M))) L) O && (\text{FORGFOR}) \\
\mathbf{where} L \mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow \vec{N}}) M) O &\leftrightarrow \mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow \vec{N}}) M) \mathbf{having} L O && (\text{WHEREGFOR}) \\
\mathbf{gfor} (z \leftarrow \vec{k} (\overrightarrow{y \leftarrow \vec{M}}) L) [] &\leftrightarrow [] && (\text{GFOREMPTY}_2) \\
\mathbf{having} L [] &\leftrightarrow [] && (\text{HAVINGEMPTY}) \\
\mathbf{having} L (\mathbf{having} M N) &\leftrightarrow \mathbf{having} (L \wedge M) N && (\text{HAVINGHAVING})
\end{aligned}$$

図 A.2 正規化規則 (Stage 2)