

Database Languages and their Compilers

Prof. Dr. Torsten Grust

Database Systems Research Group
U Tübingen

Winter 2010



Copyright © 1995 United Feature Syndicate, Inc.
Redistribution in whole or in part prohibited

2 A Query Calculus for SQL

- ▶ Although SQL comes with a quite limited number of syntactic forms, we can express a **broad diversity of query types**:

- **Projection**

```
SELECT  A, B
FROM    S
```

- **Aggregation**

```
SELECT  MAX(A)
FROM    S
```

- **Quantification**

```
FORALL x IN S : x <> 0
```

- **+ selection, joins, grouping, sorting, ...**

- ▶ The techniques presented in this chapter try to find a **uniform representation** for this diversity found in SQL.

2.1 SQL's Nested Loops Semantics

- ▶ SQL's operation can largely be understood by re-expressing queries via **nested iteration** (deeper nesting occurs with joins, grouping, ...):

■ Projection

```
SELECT  A, B
FROM    S
```

```
c ← ∅;
FOREACH x ∈ S DO
  ⊔ c ← c ⊕ {(x.A, x.B)};
RETURN c;
```

■ Aggregation

```
SELECT  MAX(A)
FROM    S
```

```
c ← -∞;
FOREACH x ∈ S DO
  ⊔ c ← max(c, x.A);
RETURN c;
```

■ Quantification

FORALL x IN S : $x <> 0$

 How would the appropriate loop look like in this case?

- Apparently, it's merely the **initializer** for c and the **accumulating function** that make a difference here:

query type	initializer	accumulator
Projection	\emptyset	\uplus
Aggregation (MAX)	$-\infty$	max
⋮	⋮	⋮
Quantification (FORALL)	<i>true</i>	\wedge
⋮	⋮	⋮

2.2 *fold*

- ▶ It seems worth to try to map the SQL diversity onto a **single loop skeleton**, parameterized by the initializer z and accumulator function \oplus :

 **Definition of loop skeleton *fold*:**

fold (z, \oplus, X) :

$c \leftarrow z$;

FOREACH $x \in X$ DO

└ $c \leftarrow c \oplus x$;

RETURN c ;

 ***fold* is a higher-order function!**

Observe that the \oplus parameter to *fold* is a function itself. If we have $z :: \beta, X :: \text{bag } \alpha$, then $\oplus :: (\beta, \alpha) \rightarrow \beta$, and thus $\text{fold} :: (\beta, (\beta, \alpha) \rightarrow \beta, \text{bag } \alpha) \rightarrow \beta$.

- The simple query types are readily expressed via *fold* if we plug-in suitable definition for z and \oplus :

■ Projection

SELECT	A, B	$fold (\emptyset, \oplus, S)$
FROM	S	with $\oplus (c, x) = c \uplus \{(x.A, x.B)\}$

■ Aggregation

SELECT	$MAX(A)$	$fold (-\infty, \oplus, S)$
FROM	S	with $\oplus (c, x) = \max(c, x.A)$

■ Quantification

FORALL x IN S :	$x \neq 0$	$fold (true, \oplus, S)$
		with $\oplus (c, x) = c \wedge (x \neq 0)$

► This idea extends quite naturally to more complex types of SQL queries and arbitrary levels of query nesting.

■ A join between two tables leads to nested *fold*s of depth 2:

<pre>SELECT x.A, y.B FROM R x, S y</pre>	<pre>fold (∅, ⊕, R) with ⊕ (c, x) = c ⊕ fold (∅, ⊗, S) with ⊗ (c', y) = c' ⊕ {(x.A, y.B)}</pre>
--	---

■ We can also handle **selection** as illustrated below:

<pre>SELECT A FROM S WHERE A > B</pre>	<pre>fold (∅, ⊕, S) with ⊕ (c, x) = c ⊕ (if (x.A > x.B) {x.A} else ∅)</pre>
---	--



The role of \emptyset , $-\infty$, *true*, ...

$$c \uplus \emptyset = c \quad \max(c, -\infty) = c \quad c \wedge \text{true} = c$$

► It turns out that one can embrace **all** SQL constructs via nested cascades of *fold*s.

■ One last example:

```
SELECT  COUNT(*)
FROM    R x
WHERE   EXISTS y IN S : x.A = y.B
```

$fold(0, \oplus, fold(\emptyset, \otimes, R))$
with $\oplus(c, x) = c + 1$
 $\otimes(c, x) = c \uplus (if (fold(false, \odot, S)) \{x\} else \emptyset)$
with $\odot(c, y) = c \vee (x.A = y.B)$

 Give the *fold* equivalent for the SQL query shown here:

```
SELECT  SUM(x)
FROM    R x, (SELECT y
              FROM  S y
              WHERE x.A = y.B)
```


- ▶ The resulting *fold* programs look rather forbidding, though. Explicit introduction of accumulators, variables, and their nested scopes make this representation hardly readable for a human.

2.3 The Monoid Comprehension Calculus

- ▶ There is, however, a way to eat the cake and have it, too:

We will now introduce a different representation for SQL queries, the **Monoid Comprehension Calculus (MCC)** [1, 2], with the following characteristics:

- ① SQL queries can be mapped to MCC via a straightforward **parse tree matching** mechanism;
- ② MCC is **human readable** and amenable to simple **analysis** techniques (variable scoping, type checking);
- ③ MCC has a **formal semantics**, completely based on *fold*.

- The core of MCC are **comprehension** expressions:

$$[e \mid q_1, \dots, q_n]^M$$

- the **qualifiers** q_i are either **generators** $v_i \leftarrow e_i$ (v_i variable) or **predicates**;
 - e , the **head**, is an arbitrary expression;
 - a variable bound in a generator q_i is visible in qualifiers q_{i+1}, \dots, q_n and in e ;
 - M denotes a monoid.
- **Example** (the bag of all x in X :: $bag \mathbb{N}$ so that $x \neq 42$):

$$[x \mid x \leftarrow X, x \neq 42]^{bag}$$

(here: head x , qualifiers $q_1 : x \leftarrow X, q_2 : x \neq 42$).

All bindings for x that pass the predicate are accumulated using the given monoid's binary operation (here: \oplus) to form the result.

- ▶ A comprehension over a monoid M can obviously express **joins**, **selection**, and **projection** over collections ($M \in \{bag, set, list\}$).
- ▶ Assume monoid $all = (true, \wedge)$, then we can express universal **quantification** in MCC, e.g. (with $X :: bag\ \alpha, Y :: bag\ \beta$ and α, β comparable):

$$[x > y \mid x \leftarrow X, y \leftarrow Y]^{all}$$

- ▶ Analogously, comprehensions may denote different types of **aggregations** as they are common in SQL (using monoids $sum = (0, +)$, $min = (\infty, \min), \dots$).

COUNT in MCC

How could we express the SQL COUNT aggregation, as in

```
SELECT  COUNT(*)
FROM    R x
WHERE   x > 0
```

in MCC?

- ▶ As variable binding and scoping in MCC and SQL coincide, it's rather simple to map SQL queries to equivalent MCC expressions.

- SQL formulation:

```
SELECT  COUNT(*)
FROM    R x
WHERE   EXISTS y IN S : x.A = y.B
```

MCC equivalent ($some = (false, \vee)$):

$$[1 \mid x \leftarrow R, [x.A = y.B \mid y \leftarrow S]^{some}]^{sum}$$

 How does the MCC equivalent look like here?

```
SELECT  DISTINCT x, y
FROM    R x, S y
WHERE   x.A = y.A AND FORALL z IN T : x.A > z.A
```

- In preparation for the development of a systematic mapping, let us list all monoids we will need to express SQL constructs in MCC (the purpose of $lift^M$ will become clear in a minute):

M	carrier	$lift^M$	z^M	\oplus^M
<i>bag</i>	<i>bag</i> α	$\{\cdot\}$	\emptyset	\uplus
<i>set</i>	<i>set</i> α	$\{\cdot\}$	\emptyset	\cup
<i>list</i>	<i>list</i> α	$[\cdot]$	$[\]$	$++$
<i>all</i>	\mathbb{B}	<i>id</i>	<i>true</i>	\wedge
<i>some</i>	\mathbb{B}	<i>id</i>	<i>false</i>	\vee
<i>sum</i>	\mathbb{N}	<i>id</i>	0	$+$
<i>max</i>	α	<i>id</i>	$-\infty$	max
<i>min</i>	α	<i>id</i>	∞	min

N.B.

- $id(x) = x$, $\{\cdot\}(x) = \{x\}$, $[\cdot](x) = [x]$;
- $++$ is list concatenation, e.g. $[1, 2, 3]++[4] = [1, 2, 3, 4]$.

► Earlier we have said that comprehensions have a semantics based on the *fold* loop skeleton.

■ The three equations below reveal comprehensions to be mere syntactic sugar for (nested) *fold*s:

$$[e \mid]^M = \text{lift}^M(e) \quad (1)$$

$$[e \mid v_1 \leftarrow e_1, q]^M = \text{fold}(z^M, \otimes, e_1) \quad (2)$$

with $\otimes(c, v_1) = c \oplus^M [e \mid q]^M$

$$[e \mid p, q]^M = \text{if } (p) [e \mid q]^M \text{ else } z^M \quad (3)$$

✎ Derive the *fold* form of this SQL query via MCC (SQL \rightarrow MCC \rightarrow *fold*):

```
SELECT  COUNT(*)
FROM    R x
WHERE   EXISTS y IN S : x.A = y.B
```

2.4 Mapping SQL to MCC

- ▶ In what follows we will develop a detailed mapping, say \mathbb{T} , for the complete subset of SQL constructs to MCC.
- ▶ For any given SQL query q , $\mathbb{T}(q)$ will be an equivalent MCC expression.
 - \mathbb{T} will be a **syntax-directed mapping**, i.e., for each syntactic construct of SQL we will define which MCC expression \mathbb{T} shall emit.
 - As the SQL parser provides us with a parse tree for query q , \mathbb{T} will essentially be a **tree pattern matcher**, e.g.

$$\mathbb{T} \left(\begin{array}{c} \text{expr} \\ \swarrow \quad \downarrow \quad \searrow \\ \text{expr} \quad \langle + \rangle \quad \text{expr} \\ \downarrow \quad \quad \downarrow \\ e_1 \quad \quad \quad e_2 \end{array} \right) = \mathbb{T}(e_1) + \mathbb{T}(e_2)$$

- To save space, we will abbreviate the above and write

$$\mathbb{T}(e_1 + e_2) = \mathbb{T}(e_1) + \mathbb{T}(e_2) .$$

- ▶ We shall make use of a small number of simple **helper functions** which will be introduced as soon as we need them. The two most important of these are:
 - `rv` (derive MCC variable name from SQL fragment and/or index i):

$$\begin{aligned}
 \text{rv}(i, e \text{ AS } \langle \text{ID}, v \rangle) &= v \\
 \text{rv}(i, \langle \text{ID}, v \rangle) &= v \\
 \text{rv}(i, e) &= R_i
 \end{aligned}$$

We will use `rv` to generate unique MCC variable names in situations of lax SQL variable usage:

$$\begin{array}{ccccccc}
 \text{SELECT} & \dots & & & & & \\
 \text{FROM} & \underbrace{T}_{\text{rv}(1,T)=T_1}, & \underbrace{S \text{ AS } s}_{\text{rv}(2,S \text{ AS } s)=s}, & \underbrace{T}_{\text{rv}(3,T)=T_3}, & \underbrace{(\text{SELECT} \dots \text{FROM})}_{\text{rv}(4,\text{SELECT} \dots)=R_4}
 \end{array}$$

- `col` (derive column name from SQL fragment and/or index i):

$$\begin{aligned}
 \text{col}(i, e \text{ AS } \langle \text{ID}, c \rangle) &= c \\
 \text{col}(i, \langle \text{ID}, c \rangle) &= c \\
 \text{col}(i, \langle \text{ID}, v \rangle . \langle \text{ID}, c \rangle) &= v \cdot c \\
 \text{col}(i, e) &= X_i
 \end{aligned}$$

(a) **Basic SELECT-FROM block:**

$$\begin{aligned} \mathbb{T} \left(\begin{array}{ll} \text{SELECT} & p_1, \dots, p_n \\ \text{FROM} & e_1, \dots, e_m \end{array} \right) = \\ \mathbb{T} \left(\begin{array}{ll} \text{SELECT} & p_1, \dots, p_n \\ \text{FROM} & e_1, \dots, e_m \\ \text{WHERE} & \text{TRUE} \end{array} \right) \end{aligned}$$

(b) **Basic SELECT-FROM-WHERE block:**

$$\begin{aligned} \mathbb{T} \left(\begin{array}{ll} \text{SELECT} & p_1, \dots, p_n \\ \text{FROM} & e_1, \dots, e_m \\ \text{WHERE} & q \end{array} \right) = \\ [(\text{col}(1, p_1) : \mathbb{T}(p_1), \\ \dots, \\ \text{col}(n, p_n) : \mathbb{T}(p_n)) \mid \\ \text{rv}(1, e_1) \leftarrow \mathbb{T}(e_1), \dots, \text{rv}(m, e_m) \leftarrow \mathbb{T}(e_m), \mathbb{T}(q), \\ * \equiv \text{rv}(1, e_1).* \parallel \dots \parallel \text{rv}(m, e_m).*]^{bag} \end{aligned}$$

(c) **Basic SELECT-FROM-WHERE block with DISTINCT:**

$$\begin{aligned} \mathbb{T} \left(\begin{array}{l} \text{SELECT DISTINCT } p_1, \dots, p_n \\ \text{FROM } e_1, \dots, e_m \\ \text{WHERE } q \end{array} \right) = \\ [(\text{col}(1, p_1) : \mathbb{T}(p_1), \\ \dots, \\ \text{col}(n, p_n) : \mathbb{T}(p_n)) \mid \\ \text{rv}(1, e_1) \leftarrow \mathbb{T}(e_1), \dots, \text{rv}(m, e_m) \leftarrow \mathbb{T}(e_m), \mathbb{T}(q), \\ * \equiv \text{rv}(1, e_1).* \parallel \dots \parallel \text{rv}(m, e_m).*]^{\text{set}} \end{aligned}$$

- In (b) and (c) we have used operator \parallel to glue together two tuples:

$$(a_1 : e_1, \dots, a_n : e_n) \parallel (b_1 : e'_1, \dots, b_m : e'_m) = (a_1 : e_1, \dots, a_n : e_n, b_1 : e'_1, \dots, b_m : e'_m)$$

- We use $v \equiv e'$ to denote local named definitions in comprehensions:

$$[e \mid q_1, \dots, v \equiv e', q_{i+1}, \dots, q_n]^M = [e[e'/v] \mid q_1, \dots, q_{i+1}[e'/v], \dots, q_n[e'/v]]^M$$

(d) **VALUES clause:**

$$\begin{aligned} \mathbb{T}(\text{VALUES } e_1, \dots, e_n) = \\ \{ \mathbb{T}(e_1) \} \uplus \dots \uplus \{ \mathbb{T}(e_n) \} \end{aligned}$$

(e) **FORALL quantifier:**

$$\begin{aligned} \mathbb{T}(\text{FORALL } v \text{ IN } e_1 : e_2) = \\ [\mathbb{T}(e_2) \mid v \leftarrow \mathbb{T}(e_1)]^{all} \end{aligned}$$

(f) **EXISTS quantifier:**

$$\begin{aligned} \mathbb{T}(\text{EXISTS } v \text{ IN } e_1 : e_2) = \\ [\mathbb{T}(e_2) \mid v \leftarrow \mathbb{T}(e_1)]^{some} \end{aligned}$$

(g) “False” aggregate (aggregate over a column reference in a grouping, $agg \in \{\text{SUM}, \text{MAX}, \text{MIN}\}$):

$$\begin{aligned} \mathbb{T}(\text{COUNT}(a)) &= \\ & [1 \mid G \leftarrow group]^{sum} \\ \mathbb{T}(agg(a)) &= \\ & [G.col(1, a) \mid G \leftarrow group]^{agg(agg)} \end{aligned}$$

agg is the straightforward mapping from SQL aggregates to MCC monoids:

$$agg(\text{SUM}) = sum \quad agg(\text{MAX}) = max \quad agg(\text{MIN}) = min \quad .$$

These cases of \mathbb{T} refer to the identifier *group* which represents the current group of a partitioned relation. *group* is introduced by the translation rules for GROUP BY, see (j).

(h) Aggregation:

$$\begin{aligned} \mathbb{T} \left(\begin{array}{l} \text{SELECT } agg_1(a_1), \dots, agg_n(a_n) \\ \text{FROM } e_1, \dots, e_m \\ \text{WHERE } q \end{array} \right) = \\ [(\text{col}(1, agg_1(a_1)) : \mathbb{T}(agg_1(a_1)), \\ \dots, \\ \text{col}(n, agg_n(a_n)) : \mathbb{T}(agg_n(a_n))) \mid \\ \text{group} \equiv [(\text{col}(1, a_1) : \mathbb{T}(a_1), \dots, \text{col}(n, a_n) : \mathbb{T}(a_n)) \mid \\ \text{rv}(1, e_1) \leftarrow \mathbb{T}(e_1), \dots, \text{rv}(m, e_m) \leftarrow \mathbb{T}(e_m), \mathbb{T}(q)]^{bag}]^{bag} \end{aligned}$$

(i) GROUP BY:

$$\begin{aligned} \mathbb{T} \left(\begin{array}{l} \text{SELECT } p_1, \dots, p_n, agg_1(a_1), \dots, agg_k(a_k) \\ \text{FROM } e_1, \dots, e_m \\ \text{GROUP BY } g_1, \dots, g_r \end{array} \right) = \\ \mathbb{T} \left(\begin{array}{l} \text{SELECT } p_1, \dots, p_n, agg_1(a_1), \dots, agg_k(a_k) \\ \text{FROM } e_1, \dots, e_m \\ \text{WHERE TRUE} \\ \text{GROUP BY } g_1, \dots, g_r \\ \text{HAVING TRUE} \end{array} \right) \end{aligned}$$

(j) GROUP BY constrained by HAVING:

$$\mathbb{T} \left(\begin{array}{l} \text{SELECT } p_1, \dots, p_n, \text{agg}_1(a_1), \dots, \text{agg}_k(a_k) \\ \text{FROM } e_1, \dots, e_m \\ \text{WHERE } q_1 \\ \text{GROUP BY } g_1, \dots, g_r \\ \text{HAVING } q_2 \end{array} \right) =$$

$$\left[\begin{array}{l} (\text{col}(1, p_1) : \mathbb{T}(p_1), \dots, \\ \text{col}(n, p_n) : \mathbb{T}(p_n), \\ \text{col}(n+1, \text{agg}_1(a_1)) : \mathbb{T}(\text{agg}_1(a_1)), \dots, \\ \text{col}(n+k, \text{agg}_k(a_k)) : \mathbb{T}(\text{agg}_k(a_k))) \mid \\ P \leftarrow [(\text{col}(1, p_1) : \mathbb{T}(p_1), \dots, \\ \text{col}(1, p_n) : \mathbb{T}(p_n), \\ \text{group} : [(\text{col}(1, a_1) : \mathbb{T}(a_1), \dots, \text{col}(k, a_k) : \mathbb{T}(a_k)) \mid \\ \text{rv}(1, e_1) \leftarrow \mathbb{T}(e_1), \dots, \text{rv}(m, e_m) \leftarrow \mathbb{T}(e_m), \mathbb{T}(q_1), \\ G.\text{col}(1, g_1) = \mathbb{T}(g_1), \dots, G.\text{col}(r, g_r) = \mathbb{T}(g_r)]^{bag}) \mid \\ G \leftarrow [(\text{col}(1, g_1) : \mathbb{T}(g_1), \dots, \\ \text{col}(r, g_r) : \mathbb{T}(g_r)) \mid \\ \text{rv}(1, e_1) \leftarrow \mathbb{T}(e_1), \dots, \text{rv}(m, e_m) \leftarrow \mathbb{T}(e_m)]^{set}]^{bag}, \\ \mathbb{T}(q_2)]^{bag} \end{array} \right.$$

(k) ORDER BY:

$$\mathbb{T}(e \text{ ORDER BY } a_1, \dots, a_n) = \text{sort}(\text{ord}(a_1), \dots, \text{ord}(a_n)) \mathbb{T}(e)$$

with

$$\text{ord}(a \text{ ASC}) = (\mathbb{T}(a), \leq) \quad \text{ord}(a \text{ DESC}) = (\mathbb{T}(a), \geq) \quad .$$

Let *sort* denote a function that can sort its input according to the criteria given in (\dots) , producing an output of type *list* α if $\mathbb{T}(e) :: \text{bag } \alpha$ or $\mathbb{T}(e) :: \text{set } \alpha$.

(l) UNION, EXCEPT, INTERSECT:

$$\begin{aligned} \mathbb{T}(e_1 \text{ UNION } e_2) &= \mathbb{T}(e_1) \uplus \mathbb{T}(e_2) \\ \mathbb{T}(e_1 \text{ EXCEPT } e_2) &= \mathbb{T}(e_1) \setminus^+ \mathbb{T}(e_2) \\ \mathbb{T}(e_1 \text{ INTERSECT } e_2) &= \mathbb{T}(e_1) \uplus \mathbb{T}(e_2) \end{aligned}$$

$\uplus, \setminus^+, \uplus$ denote the respective bag operators.

(m) **Binary/unary operators** ($op \in \{+, *, \text{AND}, \text{OR}, =, <>, <=, >=, <, >, \text{NOT}\}$):

$$\begin{aligned}\mathbb{T}(e_1 \text{ op } e_2) &= \mathbb{T}(e_1) \text{ op}(op) \mathbb{T}(e_2) \\ \mathbb{T}(op \ e) &= \text{op}(op) \mathbb{T}(e)\end{aligned}$$

with

$$\begin{array}{llll}\text{op}(+) = + & \text{op}(\ast) = \ast & \text{op}(\text{AND}) = \wedge & \text{op}(\text{OR}) = \vee \\ \text{op}(=) = = & \text{op}(<>) = \neq & \text{op}(<=) = \leq & \text{op}(>=) = \geq \\ \text{op}(<) = < & \text{op}(>) = > & \text{op}(\text{NOT}) = \neg & .\end{array}$$

(n) **Parenthesized queries, ignoring tuple variables:**

$$\begin{aligned}\mathbb{T}((e)) &= \mathbb{T}(e) \\ \mathbb{T}(e \text{ AS } v) &= \mathbb{T}(e)\end{aligned}$$

(o) **Basic query elements** (constants, column references, ...):

$$\mathbb{T}(e) = e$$

2.4.1 Applying \mathbb{T} to Translate SQL Queries

- ▶ The definition of \mathbb{T} is **compositional** in the sense that
 - ① we can apply \mathbb{T} to a complex SQL query q ,
 - ② then use \mathbb{T} to translate the **subqueries** of q ,
 - ③ and finally assemble (compose) the resulting MCC parts to form the overall translation result (an MCC expression).
- ▶ This simplifies the implementation of \mathbb{T} and also helps us to easily track the progress of \mathbb{T} . Consider:

$$\mathbb{T} \left(\begin{array}{l} \text{SELECT} \quad \text{COUNT}(\ast) \\ \quad \text{FROM} \quad R \text{ AS } x \\ \quad \text{WHERE} \quad \text{EXISTS } y \text{ IN } S : x.A = y.B \end{array} \right)$$

$$\begin{aligned}
& \stackrel{=}{=} \left[\left(\text{col}(1, \text{COUNT}(*)) : \mathbb{T}(\text{COUNT}(*)) \right) \mid \right. \\
& \stackrel{(h)}{=} \left. \text{group} \equiv \mathbb{T} \left(\begin{array}{l} \text{SELECT} \quad * \\ \text{FROM} \quad R \text{ AS } x \\ \text{WHERE} \quad \text{EXISTS } y \text{ IN } S : x.A = y.B \end{array} \right) \right]^{bag} \\
& \stackrel{=}{=} \left[\left(X_1 : [1 \mid G \leftarrow \text{group}]^{sum} \right) \mid \right. \\
& \stackrel{(g)}{=} \left. \text{group} \equiv \mathbb{T} \left(\begin{array}{l} \text{SELECT} \quad * \\ \text{FROM} \quad R \text{ AS } x \\ \text{WHERE} \quad \text{EXISTS } y \text{ IN } S : x.A = y.B \end{array} \right) \right]^{bag} \\
& \stackrel{=}{=} \left[\left(X_1 : [1 \mid G \leftarrow \text{group}]^{sum} \right) \mid \right. \\
& \stackrel{(b)}{=} \left. \text{group} \equiv \left[\left(\text{col}(1, *) : \mathbb{T}(*) \right) \mid \text{rv}(1, R \text{ AS } x) \leftarrow \mathbb{T}(R \text{ AS } x), \right. \right. \\
& \qquad \qquad \qquad \left. \mathbb{T}(\text{EXISTS } y \text{ IN } S : x.A = y.B), \right. \\
& \qquad \qquad \qquad \left. * \equiv \text{rv}(1, R \text{ AS } x).* \right]^{bag} \Big]^{bag} \\
& \stackrel{=}{=} \left[\left(X_1 : [1 \mid G \leftarrow \text{group}]^{sum} \right) \mid \right. \\
& \stackrel{(f)}{=} \left. \text{group} \equiv \left[\left(\text{col}(1, *) : \mathbb{T}(*) \right) \mid x \leftarrow R, \right. \right. \\
& \qquad \qquad \qquad \left. \left[\mathbb{T}(x.A = y.B) \mid y \leftarrow \mathbb{T}(S) \right]^{some}, \right. \\
& \qquad \qquad \qquad \left. * \equiv x.* \right]^{bag} \Big]^{bag}
\end{aligned}$$

$$\begin{aligned}
& \stackrel{=}{(m)} \quad [(X_1 : [1 \mid G \leftarrow group]^{sum}) \mid \\
& \quad \quad \quad group \equiv [(col(1, *) : \mathbb{T}(*)) \mid x \leftarrow R, \\
& \quad \quad \quad \quad \quad \quad [x.A = y.B \mid y \leftarrow \mathbb{T}(S)]^{some}, \\
& \quad \quad \quad \quad \quad \quad * \equiv x.*]^{bag}]^{bag} \\
& \\
& \stackrel{=}{(\equiv)} \quad [(X_1 : [1 \mid G \leftarrow [(X_1 : x.*) \mid x \leftarrow R, \\
& \quad \quad \quad [x.A = y.B \mid y \leftarrow S]^{some}]^{bag}]^{sum}) \mid]^{bag}
\end{aligned}$$

- Note that, due to MCC semantics equation (1), this comprehension computes a result of type $(bag \ \mathbb{N})$ which is what we expect of the original SQL query.

✎ Use \top to compile this SQL query into an MCC expression
(with $X :: \text{bag } (x : \mathbb{B}, y : \mathbb{N})$):

```
SELECT   $x$ , SUM( $y$ )  
FROM     $X$   
GROUP BY  $x$ 
```

- ▶ You will notice that you have to concern yourself with
 - variable use and scoping in monoid comprehensions,
 - comprehension semantics (see slide 80),
 - and typing of subexpressions.

2.5 Recommended Reading

- [1] L. FEGARAS AND D. MAIER, *Towards an Effective Calculus for Object Query Languages*, in Proc. of the ACM SIGMOD Int'l Conference on Management of Data, May 1995.

- [2] T. GRUST, J. KRÖGER, D. GLUCHE, A. HEUER, AND M. H. SCHOLL, *Query Evaluation in CROQUE—Calculus and Algebra Coincide*, in Proc. of the 15th British National Conference on Databases (BNCOD15), C. Small, P. Douglas, R. Johnson, P. King, and N. Martin, eds., no. 1271 in Lecture Notes in Computer Science (LNCS), London, Birkbeck College, July 1997, Springer Verlag, pp. 84–100.

Contents

0	Introduction to Query Compilation	1
0.1	Welcome!	1
0.2	Administrativa	3
0.3	Some Remarks on these Slides	4
0.4	Relational Databases and SQL	6
0.5	A Guided Tour through an SQL Query Processor	14
0.5.1	Character Streams and Tokens	15
0.5.2	Identify Valid SQL Syntax	16
0.5.3	Resolve the Meaning of Variables and Identifiers	17
0.5.4	Check Types and Schemas	18

0.5.5	How Could Query Evaluation Look Like?	19
0.5.6	One Query, Many Programs	20
0.5.7	Query Operators	21
0.5.8	Query Operator Trees and Rewriting	22
0.5.9	The Cheaper, the Better: Query Cost Models	23
0.5.10	Executing Operator Trees on a Machine	24
0.6	Recommended Reading	27
1	Query Parsing	28
1.1	The Scanner (Lexer)	28
1.1.1	Regular Expressions	33
1.1.2	Regular Expressions for SQL	35
1.1.3	<u>lex – A Scanner Generator</u>	<u>39</u>

1.2	Syntax Analysis (Parsing)	42
1.2.1	Context-Free Grammars	44
1.2.2	Derivations and Parse Trees	46
1.2.3	Ambiguous Grammars	50
1.2.4	Predictive Parsing and Recursive Descent	54
1.2.5	yacc – A Parser Generator	63
1.3	Recommended Reading	67
2	A Query Calculus for SQL	68
2.1	SQL’s Nested Loops Semantics	69
2.2	<i>fold</i>	71
2.3	The Monoid Comprehension Calculus	75
<u>2.4</u>	<u>Mapping SQL to MCC</u>	<u>81</u>

2.4.1	Applying \mathbb{T} to Translate SQL Queries	91
2.5	Recommended Reading	95
3	Variable Scopes and Type Inference for SQL	96
3.1	Variables in SQL	97
3.1.1	Variable Environments	99
3.2	The Type of an SQL Query	102
3.2.1	Type Inference	103
3.3	Semantic Attributes	105