

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

3 Variable Scopes and Type Inference for SQL

- ▶ We have gone a long way to recognize syntactically correct SQL queries but there is still more work to do to further analyze the meaning of queries.
- ▶ This chapter will discuss how we can get control of two **semantic properties** of SQL which go beyond syntax:
 - SQL **variable scopes**, and
 - the **typing** of SQL (sub)queries.
- ▶ MCC will turn out to be an intermediate query representation that allows us to treat both issues relatively easy.

3.1 Variables in SQL

- ▶ From the viewpoint of the query processor, variables in SQL are a language feature of core importance.
- ▶ (Im)proper usage of variables in SQL *can* make a difference. Consider:

```
SELECT  COUNT(*)
FROM    R x
WHERE   FORALL y IN S : x.A > y.A
```

- If we replace tuple variable y by x (e.g., through a user typo), this query computes something completely different.

Static Query Analysis

Trading y for x in the above query has consequences that we can account for **statically**, i.e., before query execution begins.

- ▶ Full analysis of queries with variables seems to be tricky enough; on top of that, SQL allows users to be lax when it comes to variable usage.
- ▶ Consider two relations $X :: \text{bag } (x : \mathbb{N}, y : \mathbb{B})$ and $Y :: \text{bag } (y : \mathbb{S}, z : \mathbb{S})$.

■ Shadowing I:

```
SELECT   $x, \underline{y}$ 
FROM     $X \ x, (\text{SELECT } \underline{y} \text{ FROM } Y)$ 
```

The attribute references \underline{y} are ambiguous: is this y of X (i.e., $x.y$) or y of Y ?
(SQL: \underline{y} of relation Y .)

■ Shadowing II:

```
SELECT   $\underline{x}, y$ 
FROM     $X \ \underline{x}, (\text{SELECT } z \text{ FROM } Y)$ 
```

Usage of tuple variable \underline{x} collides with presence of attribute x of the same name.
(SQL: \underline{x} refers to the attribute of X , not the tuple variable.)

3.1.1 Variable Environments

- ▶ The problem is tractable, though. The minimalistic syntactic structure of MCC and its equally simple **variable scoping rules** come to the rescue.
- We make use of a simple data structure, the **variable environment** E , which we associate with each (sub)expression in MCC:

$$E = \{v_1 \mapsto \tau_1, \dots, v_n \mapsto \tau_n\} .$$

- An entry $v_i \mapsto \tau_i$ in the environment specifies that a variable $v_i :: \tau_i$ **is in scope**.
- If a new variable $v :: \tau$ comes into scope, we **add** $\{v \mapsto \tau\}$ to E :

$$E + \{v \mapsto \tau\} = \{v_1 \mapsto \tau_1, \dots, v_n \mapsto \tau_n, v \mapsto \tau\}$$

- A **lookup** for v in $\{\dots, v \mapsto \tau, \dots, v \mapsto \tau', \dots\}$ returns τ' .
Thus, $E + E' \neq E' + E$.

- Comprehensions interact with the environment they are evaluated in in a simple manner. Let E be the current environment in which we evaluate

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

Then (“*the environment is augmented from left to right*”):

- q_1 is evaluated in E ;
- if the current environment is E and $q_i = v_i \leftarrow e_i$, with $e_i :: \text{bag } \tau_i$, then q_{i+1}, \dots, q_n are evaluated in environment $E + \{v_i \mapsto \tau_i\}$;
- e is evaluated in the environment associated with q_n .

 Which value is computed by the monoid comprehension below?

$$[x \mid x \leftarrow \{1, 2, 3\}, x \leq 2, x \leftarrow \{true, false\}]^{\text{bag}}$$

- ▶ To return to our example of variable shadowing (I), here is how we can make use of environments to resolve the ambiguity:

- Original SQL query:

```

SELECT  x, y
FROM    X x, (SELECT y FROM Y)

```

✍ **Resolve variable reference ambiguities by use of environments:**

Equivalent MCC expression:

$$[(x : x, y : y) \mid x \leftarrow X, R_2 \leftarrow [(y : y) \mid Y_1 \leftarrow Y]^{bag}]^{bag}$$

↑
↑
↑

③
①
②

(Check the environment at points \textcircled{i} , $i = 1 \dots 3$.)

3.2 The Type of an SQL Query

► **Type checking** concerns query processing in much the same way as the compilation of programming languages.

■ Try to detect “semantic nonsense”, e.g.:

```
SELECT   $x.A$ 
        FROM  (VALUES 1, TRUE, "foo")  $x$ 
        WHERE EXISTS  $z$  IN  $x : x + z$  ;
```

- types help to **prepare query execution** (allocation of buffers for (intermediate) results of correct size and shape);
- static type checking helps to **avoid runtime sanity checks**.

- ▶ Once more we can employ MCC as a concise intermediate representation for queries.
- Instead of type checking SQL directly we type check the simpler MCC. The basic observation is that we have, for any SQL query q :

$$q :: \tau \quad \Leftrightarrow \quad \mathbb{T}(q) :: \tau .$$

3.2.1 Type Inference

- ▶ Inferring the type of an MCC expression is actually a simple undertaking given that we already know the types of its subexpressions. Consider:
 - Given $e_1 :: \mathbb{N}$ and $e_2 :: \mathbb{N}$, we have $(e_1 + e_2) :: \mathbb{N}$;
 - Given $e :: \tau$, we have $[e \mid \dots]^{set} :: set \ \tau$;
 - Given $e :: bag \ \tau$ and $[\dots \mid \dots, v \leftarrow e, p, \dots]^{all}$, p is evaluated in an environment of the form $E + \{v \mapsto \tau\}$;
 - Given an environment of the form $E + \{v \mapsto \tau\}$, we have $v :: \tau$.

- ▶ This type inference process is rooted in trivial type assignments such as

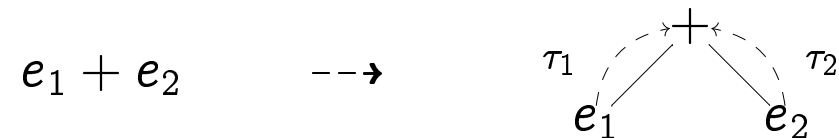
$42 :: \mathbb{N}$ or $"foo" :: \$$.

- ▶ If we want to implement type inference using such a scheme, clearly we have to proceed **bottom-up**.
- ▶ The ultimate goal is to annotate any subexpression (including the overall expression) with its type τ of the following shape:

τ	→	\mathbb{N}	
		$\$$	
		\mathbb{B}	
		<i>bag</i> τ	
		<i>set</i> τ	
		<i>list</i> τ	
		$(l_1 : \tau, \dots, l_n : \tau)$	$(l_i \text{ attribute identifier})$

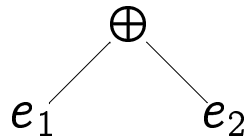
3.3 Semantic Attributes

- ▶ Leaving the variable scoping business aside for a minute, for type inference purposes, it will be sensible to denote an MCC expression by an equivalent expression tree, e.g.:

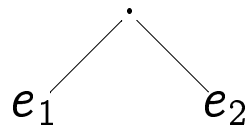


- Inferring types bottom up then simply means to **propagate type information from the leaves up to the root** of the expression tree (i.e., along the dashed arrows).
- ▶ Luckily, MCC is so simple that we can define the structure of expression trees on a single slide.

$e_1 \oplus e_2$



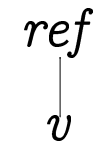
$e_1.e_2$



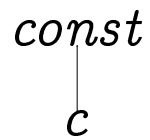
$\neg e$



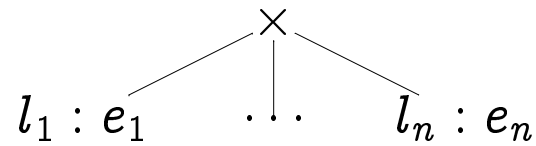
v (variable)



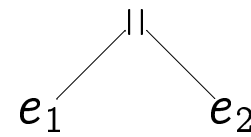
c (constant)



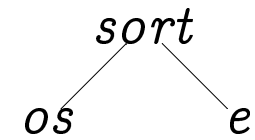
$(l_1 : e_1, \dots, l_n : e_n)$



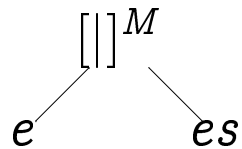
$e_1 \parallel e_2$



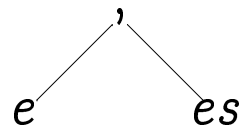
$sort\ os\ e$



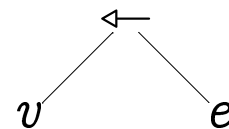
$[e \mid es]^M$



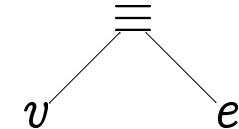
e, es



$v \leftarrow e$



$v \equiv e$



✎ What would be the shape of the expression tree for the comprehension below?

$$[(x : x, y : y) \mid x \leftarrow X, R_2 \leftarrow [(y : y) \mid y \leftarrow Y]^{bag}]^{bag}$$

- ▶ We can now organize our type inference process by simply **attaching an attribute named *type* to each node** of the expression tree.
- In accordance to our plan to proceed **bottom up**, we can refer to the *type* attribute of child nodes when we try to determine the *type* for a current node (*type* is said to be **synthesized**), e.g.:

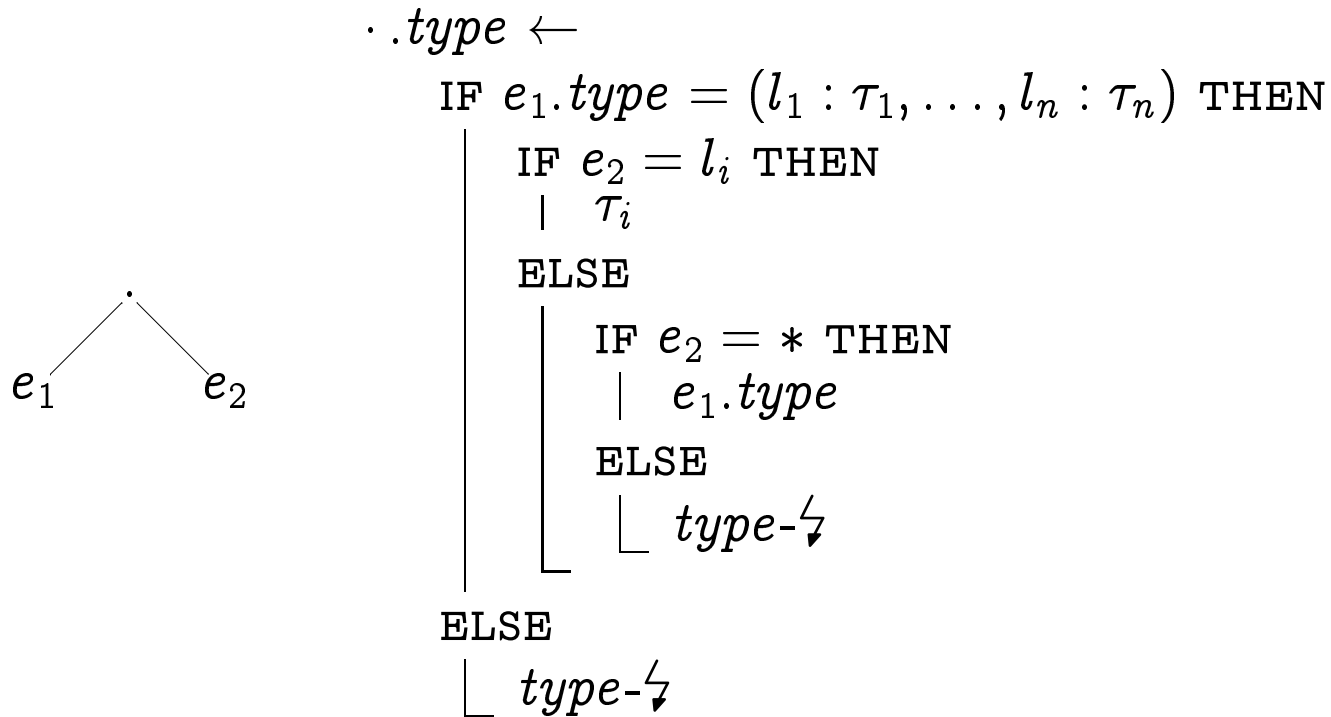
```

      ¬.type ←
      IF e.type = B THEN
      | B
      ELSE
      | type-↯

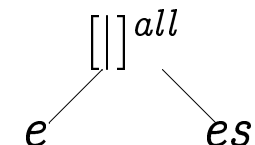
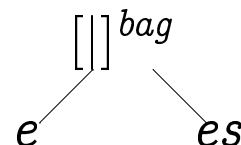
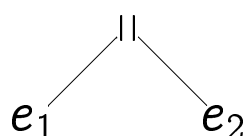
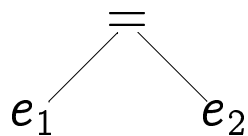
```

- If the flow of control reaches *type-↯* we have detected a **type error** and abort the type inference process.

■ Or, quite similarly:



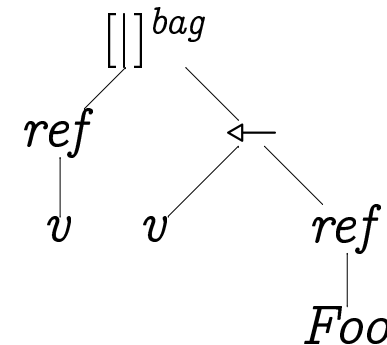
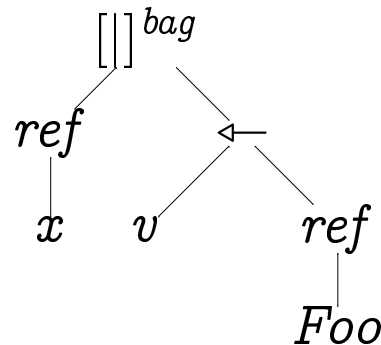
 Devise the type inference rules for the MCC expression trees below:



► We need to be more careful as soon as we come across *ref* nodes.

■ Consider the two comprehensions and their respective expression trees below (let *Foo* :: *bag* (*x* : \mathbb{N}) denote a relation residing in the DBMS):

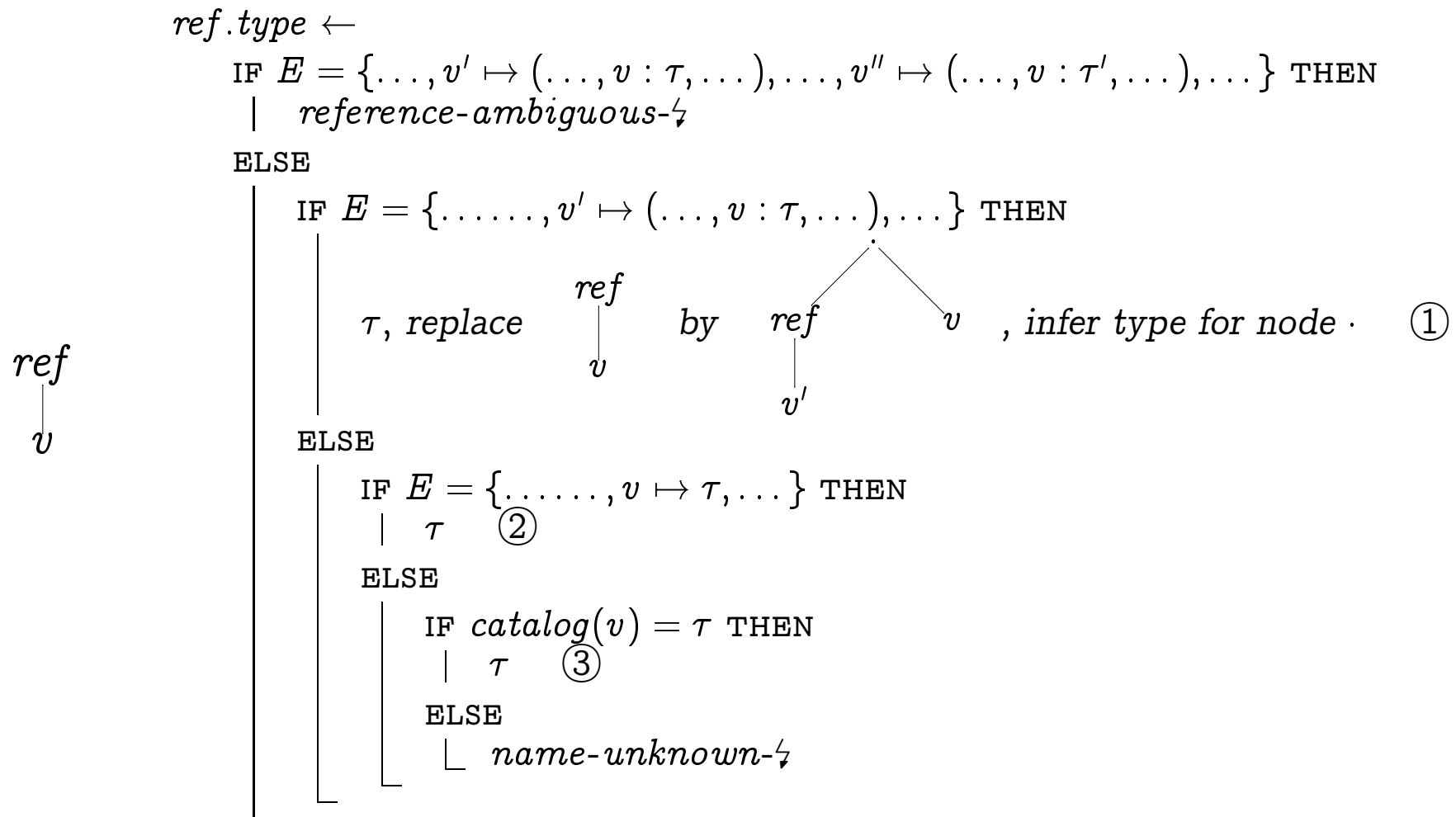
$$[x \mid v \leftarrow Foo]^{bag}$$

$$[v \mid v \leftarrow Foo]^{bag}$$


- ① $\underset{x}{ref}$: reference to a **column name** of a variable in scope;
- ② $\underset{v}{ref}$: reference to a **variable** in scope;
- ③ $\underset{Foo}{ref}$: reference to a **relation registered in the DBMS catalog**.

► At this stage in query analysis we are not yet able to distinguish the three cases.

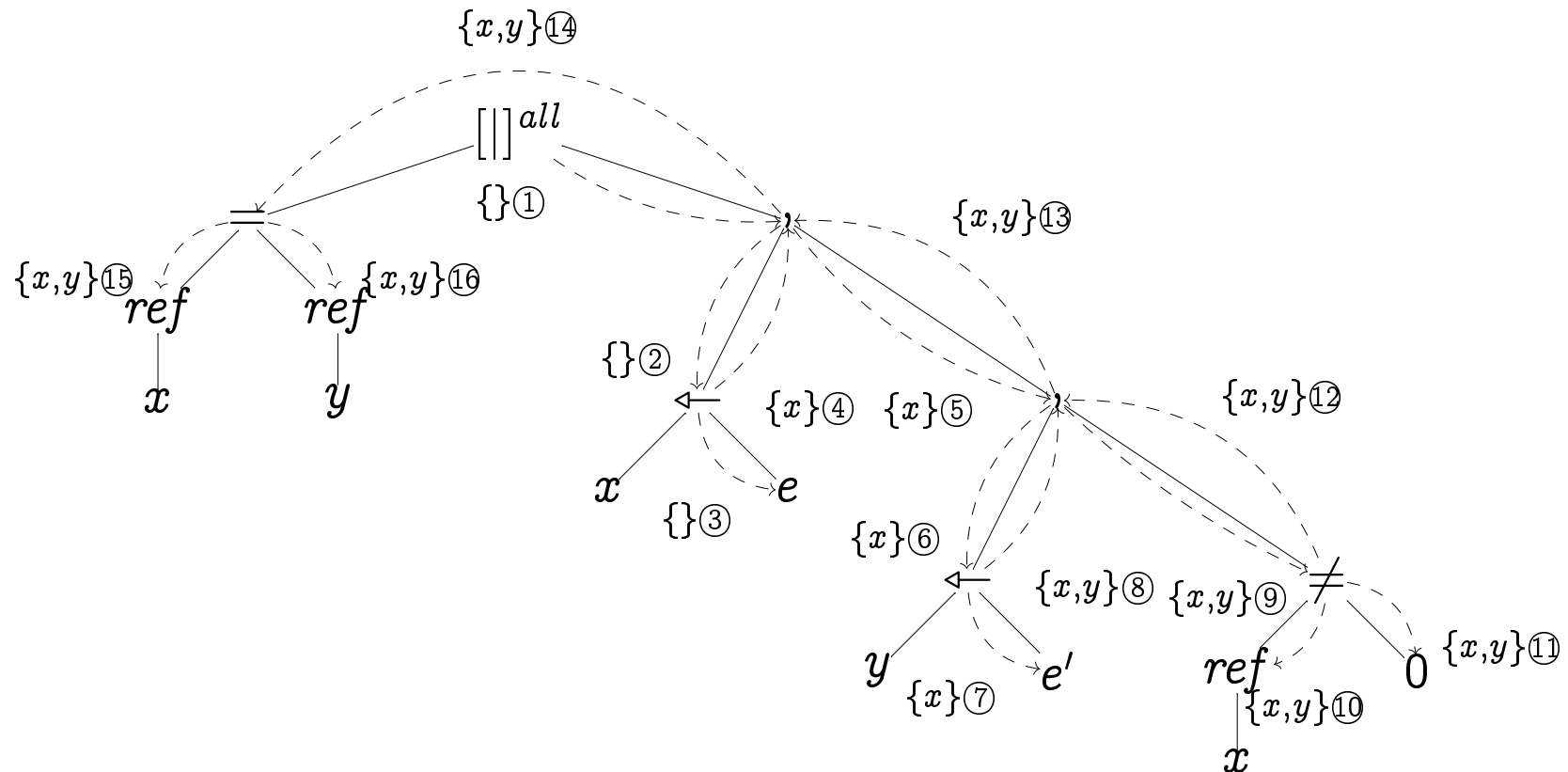
■ What we need is access to the current environment E during type inference:



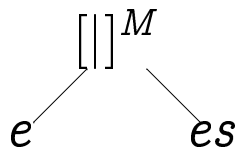
► To make the current variable environment accessible let us attach another semantic attribute, *env*, to each node in the MCC expression trees.

■ Here is an example of how the value of *env* has to be propagated:

$$[x = y \mid x \leftarrow e, y \leftarrow e', x \neq 0]^{all}$$



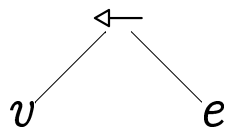
- ▶ Note that—for almost all node types—the *env* attribute is propagated **top-down** during this process. Attribute *env* is said to be **inherited**.
- ▶ An exception to this are the node types $[[]]^M$, \leftarrow , and $(,)$ which represent the comprehension construct in MCC:



```

es.env ←  $[[ ]]^M$ .env ;
do environment propagation in es (may update es.env) ;
e.env ← es.env ;

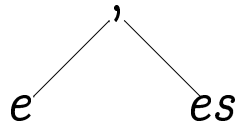
```



```

e.env ←  $\leftarrow$ .env ;
do environment propagation in e ;
v.type ← IF e.type = bag  $\tau$  THEN  $\tau$  ELSE type- $\leftarrow$ ;
 $\leftarrow$ .env ←  $\leftarrow$ .env + {v  $\mapsto$  v.type} ;

```



```
e.env ← ,.env ;  
do environment propagation in e (may update e.env) ;  
es.env ← e.env ;  
do environment propagation in es (may update es.env) ;  
,.env ← es.env ;
```

- ▶ We could have done without explicit instructions like *do environment propagation in ...*, if we adopt a simple rule:

The occurrence of a semantic attribute (of a child node) on the right-hand side of an assignment (\leftarrow) implies that his node must have been visited prior to the assignment.

- ▶ Traversing expression (or parse) trees to annotate nodes with semantic attributes is a powerful and well-understood technique in Computer Science: see **attribute grammars**.

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
2.4	Mapping SQL to MCC	81

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