

# An SQL Frontend on top of OCaml for Data Analysis

Yan Dong, Yahui Song, Chin Wei Ngan

National University of Singapore

2<sup>nd</sup> Sep 2022 @ IFL'22 in Copenhagen

### From loops to functional stream processing

```
List<Album> albums = ...;
         List<String> result = new ArrayList<>();
 3
         for (Album album : albums) {
             if (album.getYear() != 1999) {
 6
                 continue;
10
             if (album.getGenre() != Genre.ALTERNATIVE) {
                 continue;
11
12
13
             result.add(album.getName());
14
15
             if (result.size() == 5) {
17
                 break
18
                                                   List<String> result =
19
                                                       albums.stream()
20
         Collections.sort(result);
21
                                                             .limit(5)
```

This code is equivalent to:

```
List<String> result =
    albums.stream()
    .filter(album -> album.getYear == 1999)
    .filter(album -> album.getGenre() == Genre.ALTERNATIVE)
    .limit(5)
    .map(Album::getName)
    .sorted()
    .collect(Collectors.toList());
```

### **Motivation Example - Data Processing**

• Incomes sorted by month in decreasing order in 2021

```
type order = { id: int; price: float; month: string }
\mathbf{2}
   let group f l =
     let eq a b = compare (f a) (f b) = 0 in
4
     List.to_seq l |> Seq.group eq |> Seq.map List.of_seq |> List.of_seq
5
6
   let income by month : order list -> (string * float) list =
      fun orders ->
8
                                                                                     OCaml code written
        orders
9
                                                                                      in a functional manner
        | List.filter (fun o -> o.month >= "2021-01" && o.month <= "2021-12")
10
        |> List.sort (fun a b -> compare a.month b.month)
11
12
        |> group (fun o -> o.month)
        |> List.map (fun l ->
13
           List.fold left
14
                                                                                      With SQL Select Query,
              (fun (m, income) o -> (m, o.price +. income))
15
              ((List.hd 1).month, 0.0) 1)
16
                                                                                      the logic can be rewritten
        | List.sort (fun ( , s1) ( , s2) \rightarrow compare s2 s1)
17
                                                                                      in a more concise way
    let income by month : order list -> (string * float) list =
 1
       fun orders ->
 \mathbf{2}
         SELECT o.month, {sum o.price} FROM o <- orders
 3
         WHERE o.month >= "2021-01" && o.month <= "2021-12"
 4
         GROUP BY o.month ORDER BY {sum o.price} DESC
 5
```

### **Motivation Example - Data Processing**

• Incomes sorted by month in decreasing order in 2021

```
type order = { id: int; price: float; month: string }
                                                                                    modifying general-purposed
\mathbf{2}
   let group f l =
     let eq a b = compare (f a) (f b) = 0 in
                                                                                    programming languages
4
      List.to_seq 1 |> Seq.group eq |> Seq.map List.of_seq |> List.of_seq
5
6
   let income by month : order list -> (string * float) list =
      fun orders ->
8
                                                                                      OCaml code written
        orders
9
                                                                                       in a functional manner
        | List.filter (fun o -> o.month >= "2021-01" && o.month <= "2021-12")
10
        |> List.sort (fun a b -> compare a.month b.month)
11
12
        |> group (fun o -> o.month)
        |> List.map (fun l ->
13
           List.fold left
14
                                                                                       With SQL Select Query,
              (fun (m, income) o -> (m, o.price +. income))
15
              ((List.hd 1).month, 0.0) 1)
16
                                                                                       the logic can be rewritten
        | List.sort (fun ( , s1) ( , s2) \rightarrow compare s2 s1)
17
                                                                                       in a more concise way
    let income by month : order list -> (string * float) list =
 1
       fun orders ->
 \mathbf{2}
```

Implement databases by

```
3 SELECT o.month, {sum o.price} FROM o <- orders
4 WHERE o.month >= "2021-01" && o.month <= "2021-12"</p>
```

```
5 GROUP BY o.month ORDER BY {sum o.price} DESC
```

### We propose: <u>SelectML = OCaml + Select Query</u>

- SelectML is built on top of OCaml
  - Static typed & Functional
- SelectML introduces SQL Select Query onto OCaml
  - Supporting a declarative way for data processing
- Suitable for operating list data
  - Also supporting arrays, sequences, and user-defined data types
- Orthogonal with other OCaml language features
  - Core language, module language, object language, ...

### **System Overview**

- **OCaml frontend** Source Program  $\xrightarrow{Parsing}$  Parsetree  $\xrightarrow{Typing}$  Typedtree  $\xrightarrow{Translating}$  Lambda
- SelectML is implemented by modifying the parsing and typing phases

### **System Overview**

• **OCaml frontend** Source Program  $\xrightarrow{Parsing}$  Parsetree  $\xrightarrow{Typing}$  Typedtree  $\xrightarrow{Translating}$  Lambda

• SelectML is implemented by modifying the parsing and typing phases



✓ Adding the syntax support for SelectML

✓ Conduct type checking before code refracturing

✓ Code refracturing via a formally defined translation schema with a set of query plans

# SelectML Language Design

#### OCaml + Select Expression

- SELECT, DISTINCT, FROM, WHERE, GROUP BY, HAVING, ORDER BY
- Aggregate functions & applications

#### Common features for data analysis

• mapping, filtering, grouping, ordering

Variables	x		
Expressions	e	::=	$\dots \mid q \mid \{e \mid e\}$
Select Expressions	q	::=	SELECT [DISTINCT] $e$ [FROM $s$ ] [WHERE $e$ ]
			$[\texttt{GROUP BY } e] \ [\texttt{HAVING } e] \ [\texttt{ORDER BY } o]$
Source Expressions	s	::=	$x \leftarrow e \mid (x_1, \dots, x_n) \leftarrow e \mid s, s$
Order Expressions	0	::=	$e \; [\texttt{ASC} \mid \texttt{DESC} \mid \texttt{USING} \; e] \mid o, o$

```
1 (* as arguments *)
2 f 123 (SELECT x FROM x <- xs)
3
4 (* as return values *)
5 let g xs = SELECT x FROM x <- xs
6
7 (* as operands *)
8 x :: SELECT y FROM y <- ys</pre>
```

- FROM clause
  - xs AS x changed to ' $x \leftarrow xs'$
  - x::'a, xs::'a list
- GROUP BY clause
  - COUNT(y) changed to {count y}

1 /\* SQL \*/
2 SELECT x, y FROM xs AS x, ys AS y;
1 (\* SelectML \*)

2 **SELECT** x, y **FROM** x <- xs, y <- ys;;

1 /\* SQL \*/
2 SELECT x, COUNT(y) FROM t GROUP BY x;

1 (\* SelectML \*)
2 SELECT x, {count y} FROM (x, y) <- t GROUP BY x;;</pre>

- FROM clause
  - xs AS x changed to ' $x \leftarrow xs'$
  - x::'a, xs::'a list
- GROUP BY clause
  - COUNT(y) changed to {count y}

- 1 /\* SQL \*/
  2 SELECT x, y FROM xs AS x, ys AS y;
- 1 (\* SelectML \*)
- $_2$  Select x, y from x <- xs, y <- ys;;

```
1 /* SQL */
2 SELECT x, COUNT(y) FROM t GROUP BY x;
```

1 (\* SelectML \*)
2 SELECT x, {count y} FROM (x, y) <- t GROUP BY x;;</pre>

#### • WHERE & HAVING clause

- Boolean expressions
- WHERE: before GROUP BY
- HAVING: after GROUP BY

SELECT x FROM x <- xs WHERE f x;; SELECT {h x} FROM x <- xs HAVING g {h x};;

#### • Aggregate functions

• Line 2:

('a, 'b) agg is the built-in type for aggregate functions

• Line 4~9:

Common aggregate functions: AVG, COUNT, SUM, MIN, MAX

1 type ('a, 'b, 'c) aggfunc = 'c \* ('c -> 'a -> 'c) \* ('c -> 'b)
2 type (\_, \_) agg = Agg : ('a, 'b, 'c) aggfunc -> ('a, 'b) agg
3
4
5
6
7
8
1 val mkagg : 'c -> ('c -> 'a -> 'c) -> ('c -> 'b) -> ('a, 'b) agg
2
3
(\* create an aggregate function \*)
4 let count = mkagg 0 (fun n \_ -> n + 1) (fun n -> n);;
9
6 (\* usage inside Select Expression \*)
7 SELECT {count x} FROM x <- [1;2;3];;</pre>

#### • ORDER BY clause

- Ordering key must be comparable expressions
- Ordering directions: ASC, DESC, USING cmp\_func

```
1 let odd_first a b =
2    let x = a mod 2 = 0 in
3    let y = b mod 2 = 0 in
4    compare x y;;
5
6 SELECT x, y
7 FROM (x, y) <- [("a", 2); ("a", 3); ("b", 4); ("b", 5)]
8 ORDER BY x ASC, y USING odd_first;;
9
10 -:(string * int) SelectML.src=[("a",3); ("a",2);("b",5);("b",4)]</pre>
```

#### • ORDER BY clause

- Ordering key must be comparable expressions
- Ordering directions: ASC, DESC, USING cmp\_func

```
1 let odd_first a b =
2    let x = a mod 2 = 0 in
3    let y = b mod 2 = 0 in
4    compare x y;;
5
6 SELECT x, y
7 FROM (x, y) <- [("a", 2); ("a", 3); ("b", 4); ("b", 5)]
8 ORDER BY x ASC, y USING odd_first;;
9
10 -:(string * int) SelectML.src=[("a",3); ("a",2);("b",5);("b",4)]</pre>
```

#### • SELECT clause

- DISTINCT for deduplication
- Type of Select Expressions
  - Assume the type of SELECT clause is 'a
  - Line 1~8: general cases
    - type 'a list
  - Line 10~16: returns exactly one row
    - type 'a

```
1 SELECT x FROM x <- [1;1;2;2;3;3];;</pre>
2 - : int list = [1; 1; 2; 2; 3; 3]
3
4 SELECT DISTINCT x FROM x <- [1;1;2;2;3;3];;
5 - : int list = [1; 2; 3]
6
7 SELECT y, x FROM (x, y) <- [("a", 1); ("b", 2)];;
8 - : (int * string) list = [(1, "a"); (2, "b")]
10 (* without FROM, WHERE, and HAVING *)
11 SELECT x, y;;
12 SELECT x, y GROUP BY z;;
13 SELECT x, y ORDER BY z;;
14
  (* aggregation without GROUP BY, WHERE, and HAVING *)
15
16 SELECT {count x} FROM x <- t;;
```

### **Query Plans and Translation Schema**



#### **Query Plans and Translation Schema**

- Translation Schema: Query plans → Typed OCaml code
  - Semantics of the query plans are captured by primitives listed in module **SelectML**



### **Translation Schema in Detail**

- Query plans are translated to primitives defined in Stdlib module **SelectML**
- Primitives defined on line 1~10 are used by intermediate computations
- Primitives defined on line 12~14 are used for casting type

```
1 type 'a t = 'a list
_2 val one : 'a t -> 'a
3 val singleton : 'a -> 'a t
4 val product : ('a -> 'b -> 'c) -> 'a t -> 'b t -> 'c t
5 val map : ('a -> 'b) -> 'a t -> 'b t
6 val filter : ('a -> bool) -> 'a t -> 'a t
7 val sort : ('a -> 'a -> int) -> 'a t -> 'a t
8 val unique : 'a t -> 'a t
9 val group_all : ('a, 'b) agg -> 'a t -> 'b
10 val group : ('a -> 'c) -> ('a, 'b) agg -> 'a t -> 'b t
11
12 type 'a src = 'a list
13 val input : 'a src -> 'a t
14 val output : 'a t -> 'a src
```

# **Translation Schema in Deta**

- Query plans are translated to primitives defined in Sto
- Primitives defined on line 1~10 are used by intermed
- Primitives defined on line 12~14 are used for casting

```
1 type 'a t = 'a list
2 val one : 'a t -> 'a
3 val singleton : 'a -> 'a t
4 val product : ('a -> 'b -> 'c) -> 'a t -> 'b t ->
5 val map : ('a -> 'b) -> 'a t -> 'b t
6 val filter : ('a -> bool) -> 'a t -> 'a t
7 val sort : ('a -> 'a -> int) -> 'a t -> 'a t
8 val unique : 'a t -> 'a t
9 val group_all : ('a, 'b) agg -> 'a t -> 'b
10 val group : ('a -> 'c) -> ('a, 'b) agg -> 'a t ->
                                                     27
11
12 type 'a src = 'a list
13 val input : 'a src -> 'a t
14 val output : 'a t -> 'a src
```

```
1 [\mathcal{E}] \implies \text{SelectML.singleton} ()
5 \llbracket \mathcal{P}_1 \times \mathcal{P}_2 \rrbracket \Longrightarrow  SelectML.product
        (\operatorname{fun} \llbracket \phi(\mathcal{P}_1) \rrbracket \llbracket \phi(\mathcal{P}_2) \rrbracket \rightarrow \llbracket \phi(\mathcal{P}_1) \rrbracket + \llbracket \phi(\mathcal{P}_2) \rrbracket) \llbracket \mathcal{P}_1 \rrbracket \llbracket \mathcal{P}_2 \rrbracket
 8 \llbracket \sigma_e(\mathcal{P}) \rrbracket \Rightarrow \llbracket \mathcal{P} \rrbracket \mid >  SelectML.filter (fun \llbracket \phi(\mathcal{P}) \rrbracket \rightarrow e)
|_{10} [ [\Pi_{e/\chi}(\mathcal{P}) ] ] \Rightarrow [ [\mathcal{P} ] ] |> SelectML.map (fun [ [\phi(\mathcal{P}) ] ] \rightarrow e)
12 \llbracket e S_{e_f}(\mathcal{P}) \rrbracket \Rightarrow \llbracket \mathcal{P} \rrbracket > (\text{let cmp} = e_f \text{ and key } \llbracket \phi(\mathcal{P}) \rrbracket = e \text{ in }
              SelectML.sort (fun a b -> cmp (key a) (key b)))
15 \llbracket \mathcal{U}(\mathcal{P}) \rrbracket \implies \llbracket \mathcal{P} \rrbracket > \mathsf{SelectML.unique}
    \llbracket e \mathcal{G}_{e_1,\ldots,e_n}(\mathcal{P}) \rrbracket \Rightarrow \llbracket \mathcal{P} \rrbracket \mid > \mathsf{SelectML.group}
18 (fun \llbracket \phi(\mathcal{P}) \rrbracket \rightarrow e) \llbracket combine(e_1, \dots, e_n) \rrbracket
20 \llbracket \mathcal{G}_{e_1,\ldots,e_n}(\mathcal{P}) \rrbracket \implies \llbracket \mathcal{P} \rrbracket \mid > \texttt{SelectML}.group\_all \llbracket combine(e_1,\ldots,e_n) \rrbracket
              > SelectML.singleton
\begin{bmatrix} combine(e_1,\ldots,e_n) \end{bmatrix} =>
         let Agg (init1, update1, final1) = e_1 in
         let Agg (initn, updaten, finaln) = e_n in
         Agg ((init1, ..., initn),
              (fun (acc1, ..., accn) (x1, ..., xn) ->
                  (update1 acc1 x1, ..., updaten accn xn)),
              (fun (acc1, ..., accn) -> (final1 acc1, ..., finaln accn)))
\mathcal{P} with exactly one row \mathbb{P} \Rightarrow \mathbb{P} \mathbb{P} > \mathsf{SelectML}.one
\mathcal{A}  [[\mathcal{P} at the outermost]] => [[\mathcal{P}]] |> SelectML.output
```

#### Implementation

- Core implementation: ~ 1000 LOC in OCaml
- Test for validation: ~ 60 testcases, manually marked with the expected outputs
- Flexible Customisations
  - ✓ Customise the semantics of Select Expression: rewrite the WHERE implementation to keep the falsas
  - ✓ Change the input and output type of Select Expression: from 'a list to 'a array
  - ✓ Change the intermediate type of Select Expression: from list to array.
  - ✓ Generalization with Functors: to deal with both list type and array type.

https://github.com/dyzsr/ocaml-selectml

#### Implementation



Figure 36: Generalising the Type and Operations

- Core implementation: ~ 1000 LOC in OCaml
- Test for validation: ~ 60 testcases, manually marked with the expected outputs

```
    Flexible Customisations
```

 $\checkmark$  Customise the semantics of Select

 $\checkmark$  Change the input and output type

✓ Change the intermediate type of S

✓ Generalization with Functors: to  $d_{11}^{10}$ 

```
https://github.com/dyzsr/o
```

```
1 module F (SelectML : SelectMLType) = struct
    let run xs ys = SELECT x, y FROM x <- xs, y <- ys WHERE x < y;;
2
3 end;;
5 module F : (* REPL output *)
    functor (SelectML : SelectMLType) ->
      sig
        val run : 'a SelectML.src -> 'a SelectML.src ->
          ('a * 'a) SelectML.src
      end
  (* Supplying primitives with the list implementation *)
13 let open F (ListImpl) in run ...
  (* Supplying primitives with the array implementation *)
16 let open F (ArrayImpl) in run ...
```

### **Possible More Features**

#### 1. Joins ⇔ Cartesian product + filter: may reduce memory consumption

```
val join : ('a -> 'b -> 'c option) -> 'a t -> 'b t -> 'c t
1
  val join eq : ('a -> 'd) -> ('b -> 'd) -> 'a t -> 'b t -> 'c t
  (* when hash key can be determined *)
1
  SELECT ... FROM x <- xs JOIN y <- ys ON x = y;;
    (* translation *)
3
   SelectML.join_eq (fun x \rightarrow x) (fun y \rightarrow y) xs ys;;
\mathbf{4}
\mathbf{5}
  (* when hash key cannot be determined *)
6
  SELECT ... FROM x <- xs JOIN y <- ys ON f x y;;
     (* translation *)
8
   SelectML.join (fun x y -> if f x y then Some (x, y) else None) xs ys;;
9
```

#### 2. Window Functions

• They are like aggregate functions, but without causing rows to become grouped into a single output

<ul> <li>Additional window fur</li> </ul>	• 1 2	INSERT INTO t VALUES (1,1),(1,2),(1,3),(2,2),(2,3) SELECT x, COUNT(y) OVER (PARTITION BY x) FROM t;				3),(2,2),(2,3),( 3Y x) FROM t;	(3,3);	
x=1, count=3	x=1, count=3	$\begin{vmatrix} 3 \\ 4 \\ 5 \end{vmatrix}$	SELECT x, COUNT(y) OVER (ORDER BY x) FROM t;					
<pre>x=1, count=3 x=1, count=3</pre>	x=1, count=3		/* PARTITION BY */		*/ /* 01	RDER BY *,	/	
	x=1, count=3	6	x   co	unt	x	$\mathtt{count}$		
x=2, count=2	x=2, count=5	7	+		+			
x=2, count=2	x=2, count=5	8	1	3	1	3		
x=3, count=1	x=3, count=6	9	1	3	1	3		
		10	1	3	1	3		
) Window Frames of <b>PARTITION BY</b>	(b) Window Frames of ORDER BY	11	2	2	2	5		
		12	2	2	2	5		
		13	3	1	3	6		

### **Possible Optimizations**

#### 1. Rule-based optimisation

- Eliminate unnecessary plans  $\Pi_{(c_1,c_2)}(\Pi_{(x,y)/(c_1,c_2)}(\mathcal{P})) \Longrightarrow \Pi_{(x,y)}(\mathcal{P})$
- Push down filters  $\sigma_{x<1\land y>1}(\mathcal{D}_x(\mathrm{xs})\times\mathcal{D}_y(\mathrm{ys})) \Longrightarrow \sigma_{x<1}(\mathcal{D}_x(\mathrm{xs}))\times\sigma_{y>1}(\mathcal{D}_y(\mathrm{ys}))$

#### 2. Cost-based optimisation

• Select the optimal algorithm for a query plan, based on the runtime information of data

#### 3. Indexes

- Optimising table scan
- May require modifications

on the OCaml type system

$1 \\ 2 \\ 3 \\ 4$	<pre>/* SQL */ SELECT x FROM xs AS x WHERE x BETWEEN lb AND ub SELECT x FROM xs AS x, ys AS y WHERE x &lt; 1 SELECT x FROM xs AS x, ys AS y WHERE x &lt; y</pre>
16	<pre>let f (module XS : (int * int) list ORDER BY 0) =</pre>
17	SELECT x, y FROM $(x, y) \le (module XS)$ WHERE lb $\le x \&\& x \le ub;;$
18	
19	<pre>type order = { id: int; price: float; month: string }</pre>
20	
21	let f' (module XS : order list ORDER BY price) =
22	SELECT x FROM x <- (module XS) WHERE 1b <=, x && x <=, ub::

#### **Create Database**

- From REPL (Read–eval–print loop) to Database
  - Behaviors
    - #CREATE introduces a top-level variable xs of type (int \* int) table
    - #INSERT inserts two new rows (1,2), (3,4) to xs
    - #UPDATE updates xs by setting increasing the first column by 1



### Conclusions

- SelectML = OCaml + Select Query
- Language Design: Select Expression & Aggregate functions
- Semantics: Typing & Planning rules & Translation schema
- Implementation: Different Flexible Customisations
- More possible features and future work
  - Joins, window functions, optimisations, indexes
  - Connection to database via libraries such as 'caqti' with 'postgresql' database

#### Conclusions

- SelectML = OCaml + Select Query
- Language Design: Select Expression & Aggregate functions
- Semantics: Typing & Planning rules & Translation schema
- Implementation: Different Flexible Customisations

### Thanks!

- More possible features and future work
  - Joins, window functions, optimisations, indexes
  - Connection to database via libraries such as 'caqti' with 'postgresql' database