

ArrayQL Integration into Code-Generating Database Systems

Maximilian E. Schüle, Tobias Götz, Alfons Kemper, Thomas Neumann

25th International Conference on Extending Database Technology (EDBT)

29th March-1st April, 2022, Edinburgh, UK



TUM Uhrenturm

Motivation: In-Database Machine Learning

ML2SQL

MLearn: Example 1 ▼

TRANSLATE

Target Language: Python ▼

EXECUTE

```

1 A =
2 [[1.1,0.98,87.3,3],
3 [0.1,3.15,42.05,3.3],
4 [100.5,26.8,10.1,225.1],
5 [1097.5,23000,10.1,24850.1]]
6 X = A[:, 0:2]
7 y = A[:, 3]
8 bias[1,len(X,0)] : 1
9 X = (bias::X.T).T
10 Xt = X.T
11 w = (Xt*X)^(-1) * Xt * y
12 print '%', w

```

```

1 import numpy as np
2 def main():
3     A = np.array([ np.array([ 1.1,0.98,87.3,3]),
4 np.array([0.1,3.15,42.05,3.3]),np.array([ 100.5,26.8,
5 10.1,225.1]),np.array([1097.5,23000,10.1,24850.1])])
6     X = A[:,0:2 +1]
7     y = A[:,3:3 +1]
8     bias = np.full( ( 1,np.size(X,0) ), 1)
9     X = (np.append(bias,X.T, axis=0)).T
10    Xt = X.T
11    w = np.dot( np.dot( np.linalg.matrix_power(
12    (np.dot( Xt,X)), (-1)),Xt),y)
13    print( '{}'.format( w)
14    if __name__ == "__main__": main()

```

```

$ python < main.py
[[ 0.12238861]
 [ 1.9759119 ]
 [ 0.98614977]
 [-0.00300479]]

```

Motivation: In-Database Machine Learning

ML2SQL

MLearn: Example 1 TRANSLATE

```

1 A =
2 [[1.1,0.98,87.3,3],
3 [0.1,3.15,42.05,3.3],
4 [100.5,26.8,10.1,225.1],
5 [1097.5,23000,10.1,24850.1]]
6 X = A[:, 0:2]
7 y = A[:, 3]
8 bias[1,len(X,0)] : 1
9 X = (bias::X.T).T
10 Xt = X.T
11 w = (Xt*X)^(-1) * Xt * y
12 print '%', w

```

```

$ python < main.py
[[ 0.12238861]
 [ 1.9759119 ]
 [ 0.98614977]
 [-0.00300479]]

```

- **ML2SQL**: Translate a universal machine learning language into Python or SQL
- **Problem**: first converted into arrays, expensive
- **Goal**: work directly on arrays
- **Solution**: implement ArrayQL syntax draft presented at XLDB 2012 within a relational, code-generating database system

Structure



ArrayQL Algebra
completed syntax draft
added statements for joins



Umbra Integration
First implementation of ArrayQL
Relational representation



**Relational Array
Representation**

ArrayQL Algebra

- **Data Definition Language:** array creation as new array or out of an existing array

```
CREATE ARRAY m (i INTEGER DIMENSION [1:2], j INTEGER DIMENSION [1:2], v INTEGER);  
CREATE ARRAY m2 FROM SELECT [i], [j], v FROM m;
```

- **Data Query Language**

- SELECT: indices for the dimensions (in brackets), arithmetic expressions as attributes.
- FROM: specifies the source array.
- WHERE: filters each entry by a predicate.
- GROUP BY: addresses the dimensions preserved after an aggregation.

```
SELECT [i], SUM(v)+1 FROM m WHERE v>0 GROUP BY i
```

- **New:** temporary arrays and joins

```
WITH ARRAY temp AS (SELECT [i] as j, SUM(v+1) FROM m[[j] WHERE v>0 GROUP BY j) SELECT * FROM temp;  
SELECT [i] as i, [j] as j, v, v2 FROM m[[i+2,j+2] JOIN m2[[i-2,j-2];
```

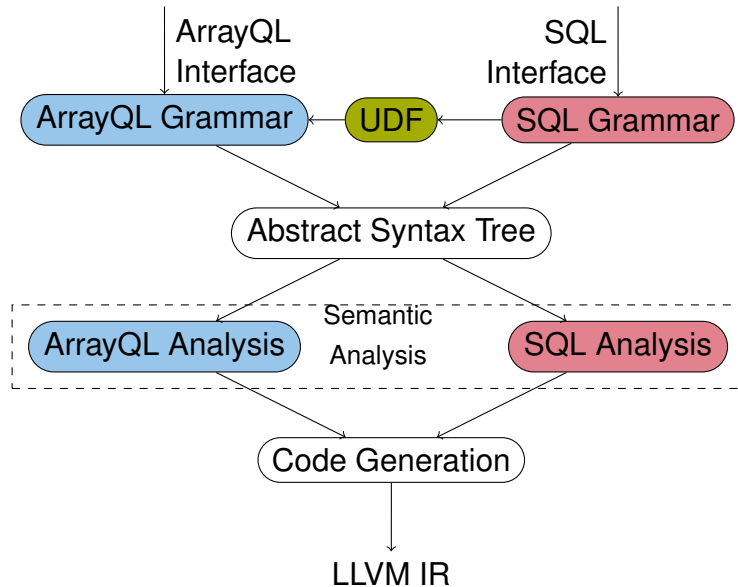
- **New: Data Modification Language**

```
UPDATE ARRAY m [[1][1] (select [[2][1],v FROM m);  
UPDATE ARRAY m [[1][1:2] (values (1),(2));
```

Umbra Integration



Umbra Integration

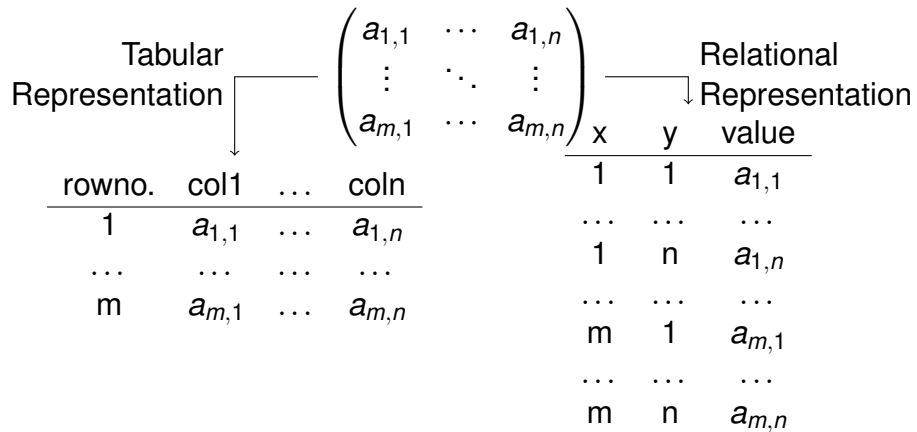


- Umbra: a code-generating database system (producer-consumer concept)
- SQL → abstract syntax tree → semantic analysis → operator plan
- Each user-defined function: a separate grammar file/semantic analysis
- Code compiled together
- Only schema known during compile-time → relational representation needed

Relational Array Representation



Relational Array Representation



```
CREATE ARRAY m (
  i INTEGER DIMENSION [1:2],
  j INTEGER DIMENSION [3:4],
  v INTEGER);
```

→

x	y	v
1	3	NULL
2	4	NULL

- **Goal:** Uniform representation (access from SQL and an array query language)
- Either coordinate list (*relational representation*) or tables with row order (*tabular representation*)
- Relational representation: saves memory (sparse arrays) and primitive data types are sufficient even for more than two dimension
- n -dimensional array with m values per cell
→ table with $n + m$ attributes
- Attributes for the indices are unique (the primary key)
- ArrayQL: *bounding box*, *validity map* and *content*
 - creation: insert tuples defining the bounding box
 - an entry within the bounding box is valid: if it exists and at least one attribute is not NULL.

Relational Array Representation: Mapping to Relation Algebra

Operator	Input	Output	Relational Algebra
apply	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }, f \in (\mathbb{R} \rightarrow \mathbb{R})$	$\mathbb{R}^{ i_1 \times \dots \times i_n }$	$\pi_{i_1, \dots, i_n, f(v)}(\mathbf{a})$
combine	$\mathbf{a}, \mathbf{b} \in \mathbb{R}^{ i_1 \times \dots \times i_n }$	$\mathbb{R}^{ i_1 \times \dots \times i_n }$	$\mathbf{a} \bowtie_{\mathbf{a}.i_1 = \mathbf{b}.i_1 \wedge \dots \wedge \mathbf{a}.i_n = \mathbf{b}.i_n}(\mathbf{b})$
i. dim. join	$\mathbf{a}, \mathbf{b} \in \mathbb{R}^{ i_1 \times \dots \times i_n }$	$(\mathbb{R}, \mathbb{R})^{ i_1 \times \dots \times i_n }$	$\mathbf{a} \bowtie_{\mathbf{a}.i_1 = \mathbf{b}.i_1 \wedge \dots \wedge \mathbf{a}.i_n = \mathbf{b}.i_n}(\mathbf{b})$
fill	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }$	$\mathbb{R}^{ i_1 \times \dots \times i_n }$	$\dots \mathbf{0}_{ a.i_1 , \dots, a.i_n } \dots$
filter	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }, \rho \in (\mathbb{R} \rightarrow \mathbb{B})$	$\mathbb{R}^{ i_1 \times \dots \times i_n }$	$\sigma_{\rho(v)}(\mathbf{a})$
rebox	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }, i_1^u, i_1^l, \dots, i_n^u, i_n^l \in \mathbb{I}$	$\mathbb{R}^{i_1^u - i_1^l \times \dots \times i_n^u - i_n^l}$	$\sigma_{i_1^l \leq i_1 \leq i_1^u \wedge \dots \wedge i_n^l \leq i_n \leq i_n^u}(\mathbf{a})$
reduce	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }, f \in (\mathbb{R}^{ i_n } \rightarrow \mathbb{R})$	$\mathbb{R}^{ i_1 \times \dots \times i_{n-1} }$	$\gamma_{i_1, \dots, i_n, f(v)}(\mathbf{a})$
rename	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }$	$\mathbb{R}^{ i_1 \times \dots \times i_n }$	$\rho(\mathbf{a})$
shift	$\mathbf{a} \in \mathbb{R}^{ i_1 \times \dots \times i_n }, i_1^l, \dots, i_n^l \in \mathbb{I}$	$\mathbb{R}^{ i_1 \times \dots \times i_n }$	$\pi_{i_1 + i_1^l, \dots, i_n + i_n^l, v}(\mathbf{a})$

- map nine ArrayQL operators to SQL operators considering the underlying schema
- one relation $\mathbf{a} \subseteq \mathbb{I}^n \times \mathbb{R}^m$ with schema $sch(\mathbf{a}) = \{\underline{i_1}, \dots, \underline{i_n}, r_1, \dots, r_m\}$ represents one n -dimensional array $\mathbf{a} \in (\mathbb{R}^m)^{|i_1| \times \dots \times |i_n|}$ with m attributes of domain \mathbb{R} as content.
- coordinates $(i_1, \dots, i_n) \subseteq \mathbb{I}^n$ form the primary key and delimit the bounding box.
- validity map of an array \mathbf{a} : set of indices $d_a \subseteq \mathbb{I}^n$ of valid entries.

ArrayQL for Linear Algebra

$$(a_{ij}) = A \in \mathbb{R}^{2 \times 3}$$

$$(a_{ij}) = A$$

```
CREATE ARRAY A (
i INTEGER DIMENSION [1 : 2],
j INTEGER DIMENSION [1 : 3],
a FLOAT)

SELECT [i], [j], A.a FROM A
```

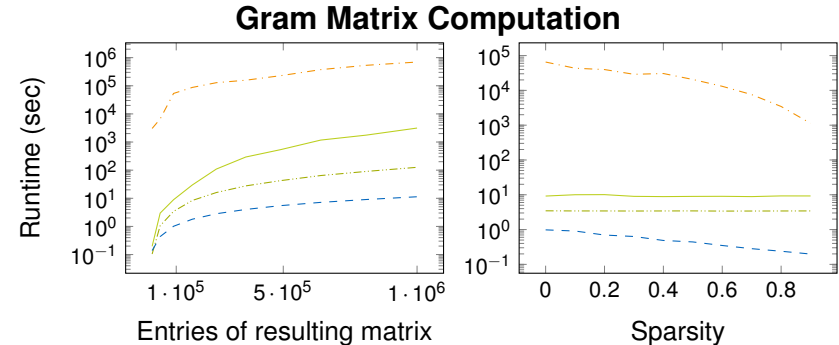
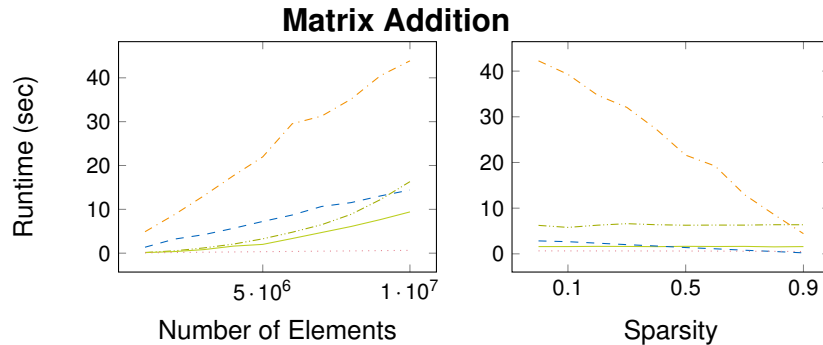
- Bounding box defines the size of a matrix, the attributes determine the field of which each entry is a member
- Scalar operations: part of the `select`-clause; matrix operations belong to the `from`-clause
- sparse matrices: values of non-existing entries are assumed to be zero (keyword `filled` behind `select`)
- Short-cuts for matrix operations as part of the `from` clause
- Linear regression: $\vec{w} = (A^T A)^{-1} A^T \vec{y}$

```
SELECT FILLED [i], [j], * FROM ((a^T * a)^-1 * a)^T * y
```

Evaluation



Evaluation: Matrix Algebra

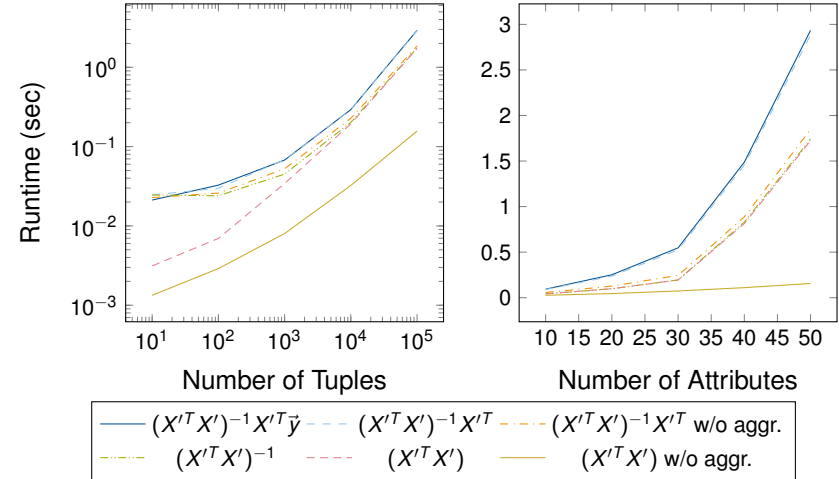
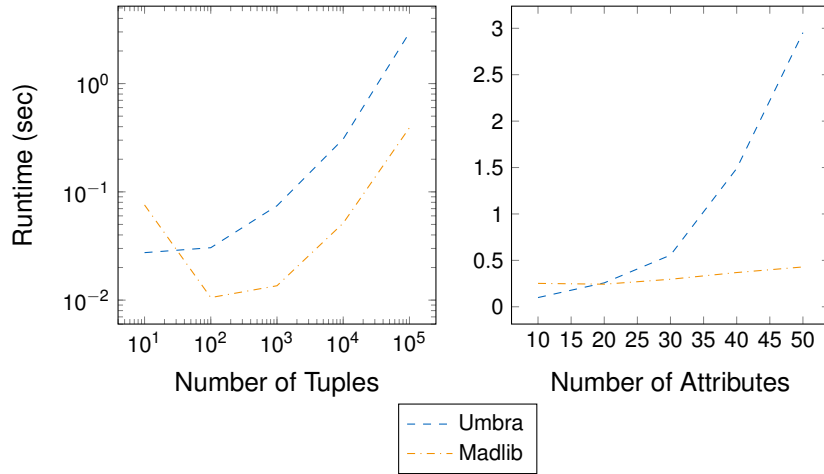


- - - RMA Optimiser
 — RMA Runtime
 - - - Umbra
 - - - Madlib Matrices
 - - - Madlib Arrays

- - - RMA Optimiser
 — RMA Runtime
 - - - Umbra
 - - - Madlib Matrices

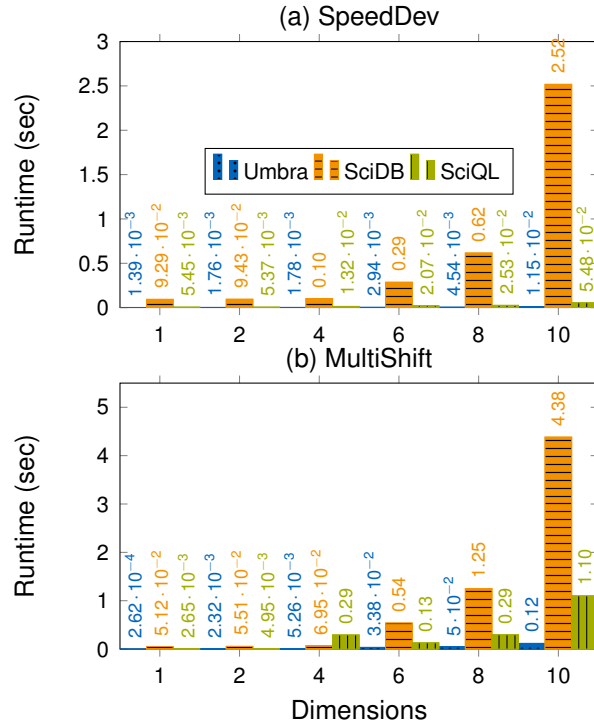
- **System:** Ubuntu 20.04 LTS, six Intel Core i7-3930K CPUs at 3.20GHz, 64 GB main-memory
- Varying either the number of elements in a dense array or the sparsity
- Varying the sparsity: gram matrix XX^T on a resulting matrix with 90000 entries, matrix addition $X + X$ with 10^6 elements
- With increasing size, ArrayQL computes the matrix sum faster than RMA
- MADlib matrices and Umbra benefit from sparse matrices (relational representation)
- Relational representation: comparable performance to existing database extensions for linear algebra

Evaluation: Linear Regression



- **System:** Ubuntu 20.04 LTS, six Intel Core i7-3930K CPUs at 3.20GHz, 64 GB main-memory
- synthetic dataset: ArrayQL using linear algebra only vs. MADlib (dedicated table function)
- Most of the time: aggregation part of each matrix product (summation)
- ability to solve numerical problems when suitable table functions are available

Evaluation: Taxi Data

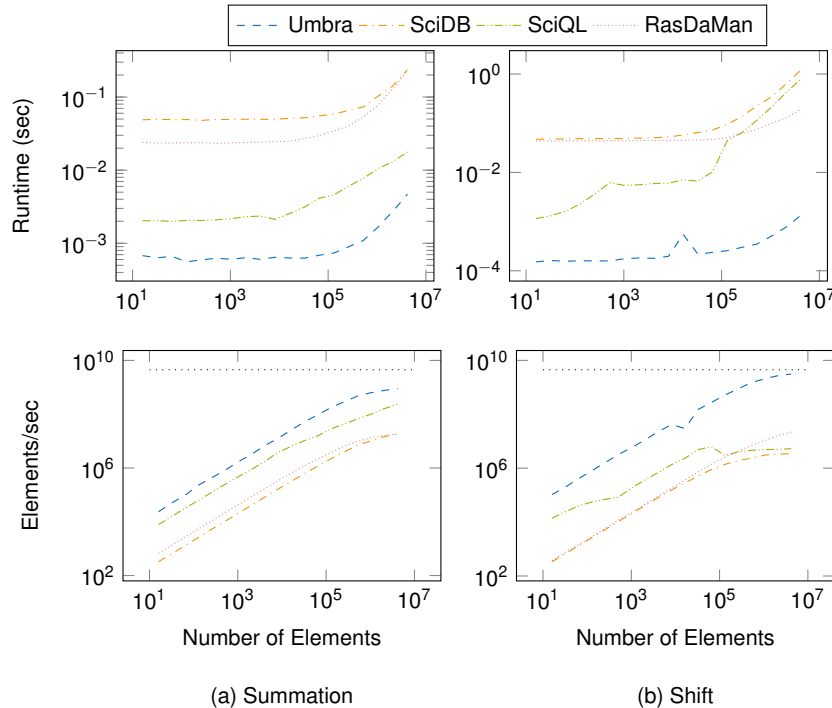


```
SELECT MAX(3600.0*(tmp3.v3-tmp1.v)) as speed
FROM (SELECT pickup_day, AVG(trip_distance/total_duration) as v
FROM taxiData WHERE total_duration<>0 GROUP BY pickup_day)
as tmp1, (SELECT AVG(trip_distance/total_duration) as v3
FROM taxiData WHERE total_duration <> 0) as tmp3;
```

```
SELECT [pickup_day] as a, ..., * FROM taxiData[a+1, ...]
```

- **System:** Ubuntu 20.04 LTS, six Intel Core i7-3930K CPUs at 3.20GHz, 64 GB main-memory
- impact of involved dimensions on aggregations
- *SpeedDev*: maximum deviation of the average speed per day compared to the average speed of the whole dataset
- *MultiShift*: shifts all array's dimensions.
- scales linearly with an increasing number of dimensions.

Evaluation: Random Data



- **System:** Ubuntu 20.04 LTS, six Intel Core i7-3930K CPUs at 3.20GHz, 64 GB main-memory
- two-dimensional arrays with synthetic data and an increasing number of elements.
- summation and shifting the indices
- Umbra the fastest system when performing aggregates
- shifting slows down the performance as all indices have to be changed
- upper constant lines: the maximum throughput of $4.5 \cdot 10^9$ elements per second

Conclusion

- First ArrayQL integration (as another query interface and addressable inside SQL)
- Integration into a code-generating database system
- Introduced a relational array model and an ArrayQL operator translation to relational algebra
- In summary, ArrayQL in Umbra benefits from sparse matrices as well as the performance of an in-memory database system
- Future: ArrayQL to SQL translator

```
CREATE FUNCTION exampletable() RETURNS TABLE (x INT, y INT, v INT) LANGUAGE 'arrayql' AS
'SELECT [x], [y], [v] FROM [m]';
CREATE FUNCTION exampleattribute() RETURNS INT[] [] LANGUAGE 'arrayql' AS
'SELECT [x], [y], [v] FROM [m]';
```

Listing 1: ArrayQL as part of a user-defined function returns either an SQL table or an SQL array.

Thank you for your attention!

