

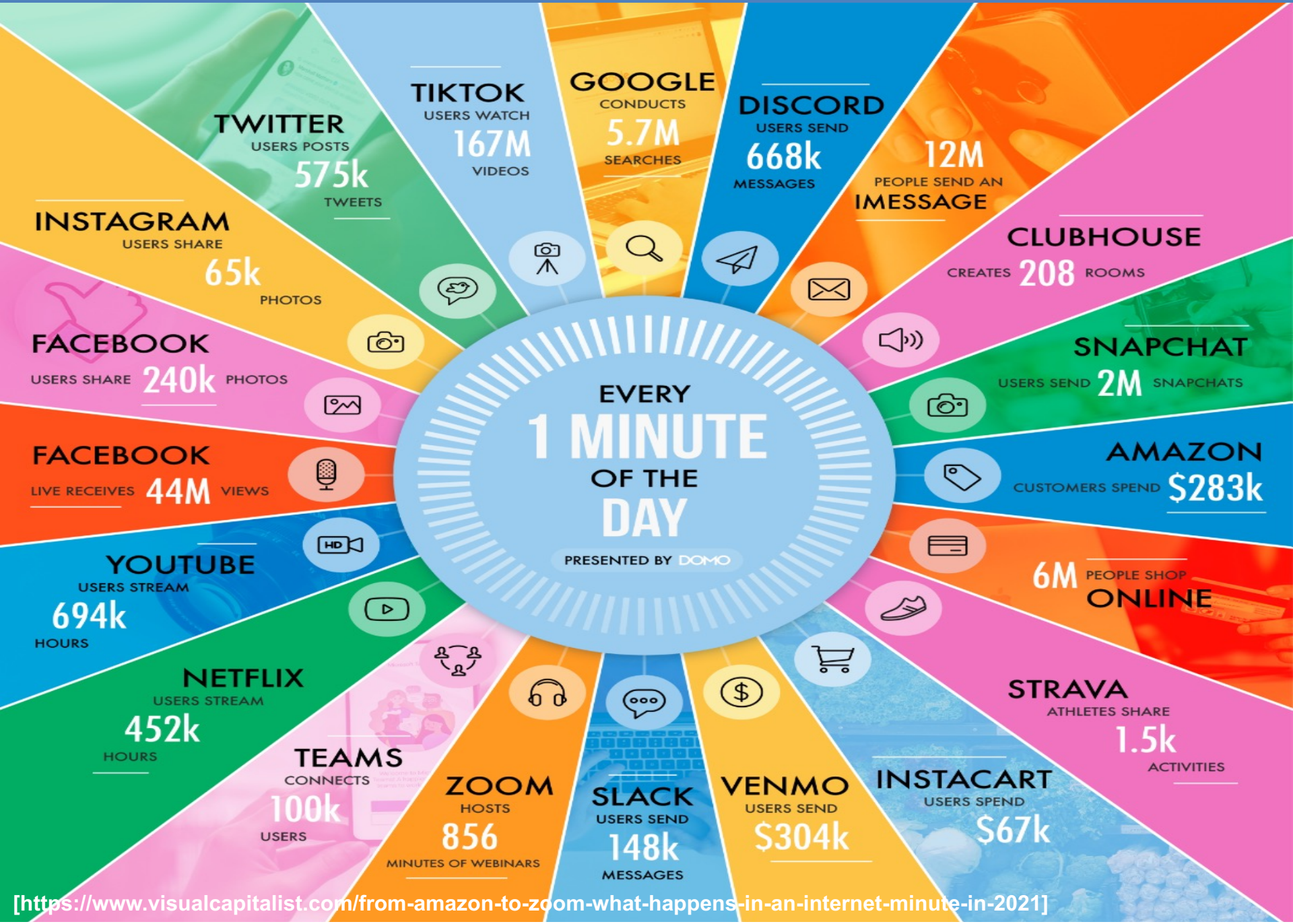
Sparse Tensor Algebra Compilation using Equality Saturation

Amir Shaikhha

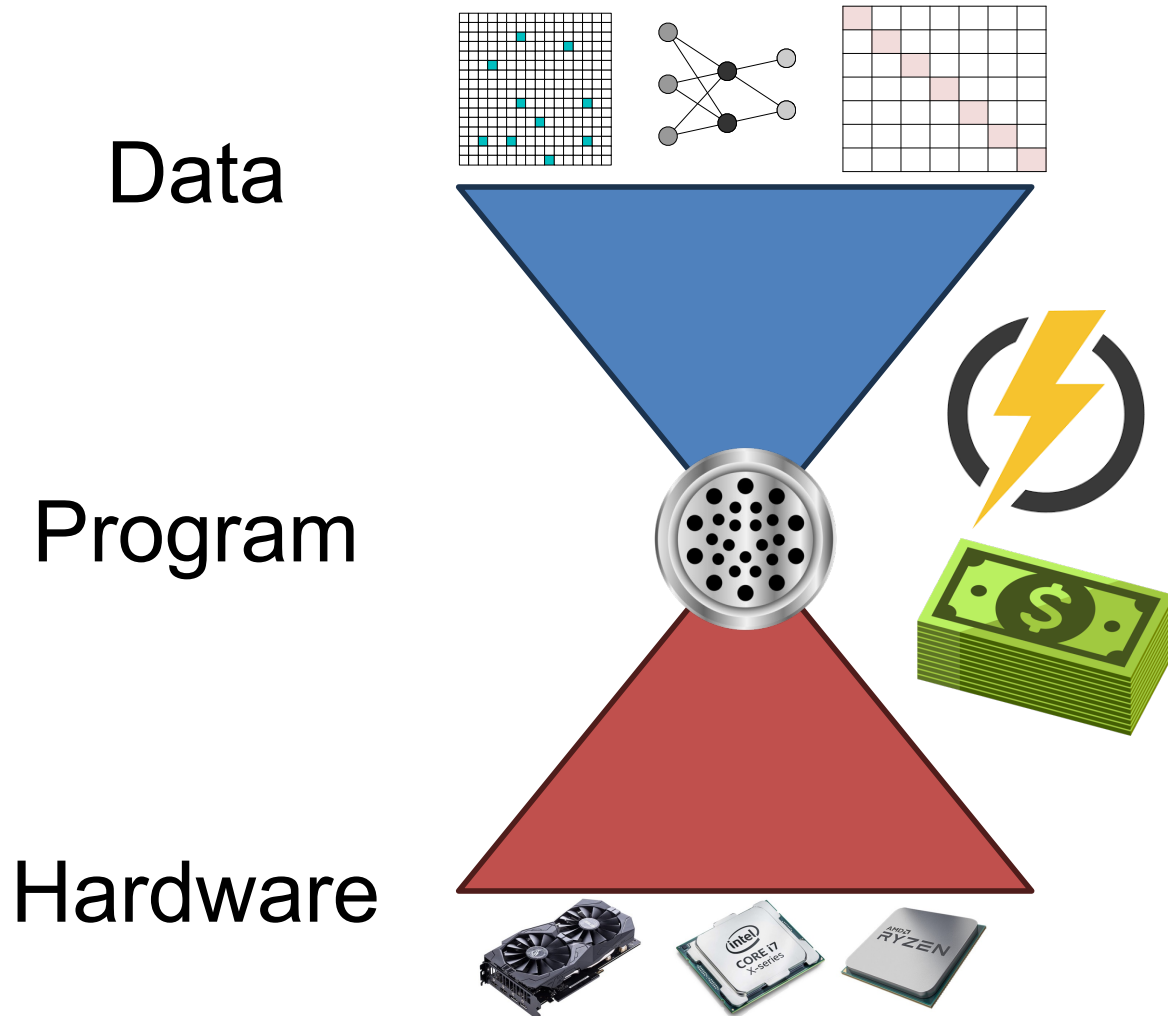
joint work with Mathieu Huot, Shideh Hashemian,
Dan Olteanu, Jaclyn Smith, Dan Suciu, Max Schleich



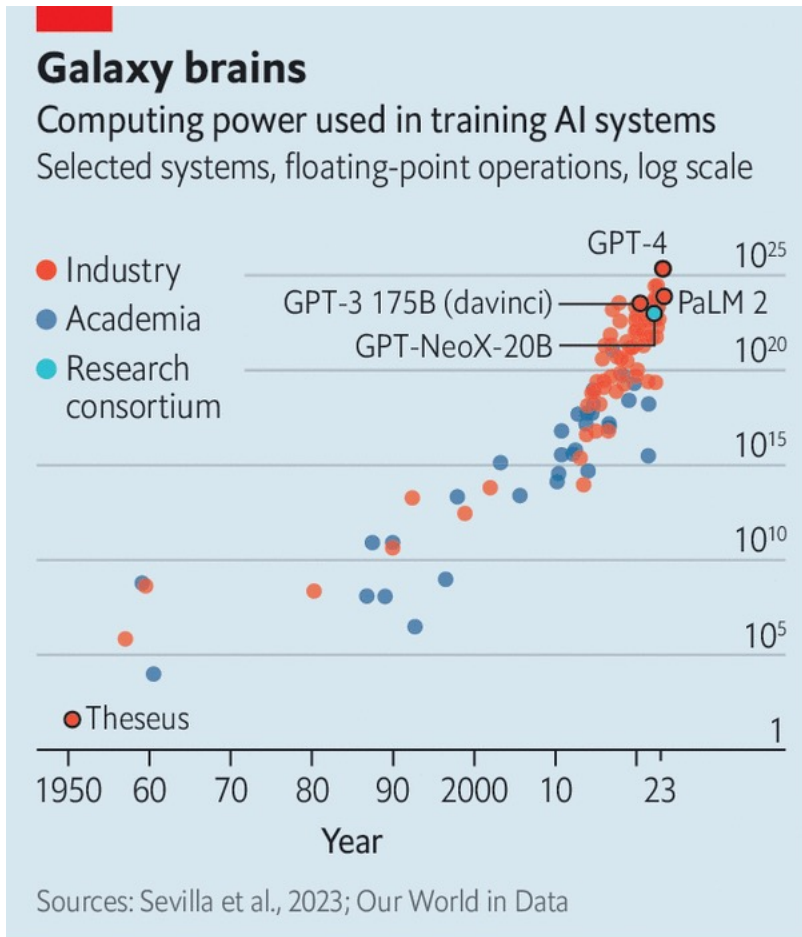
2024/4/18



Data Processing



Planet and Economic Crisis



The Economist



GPT-4

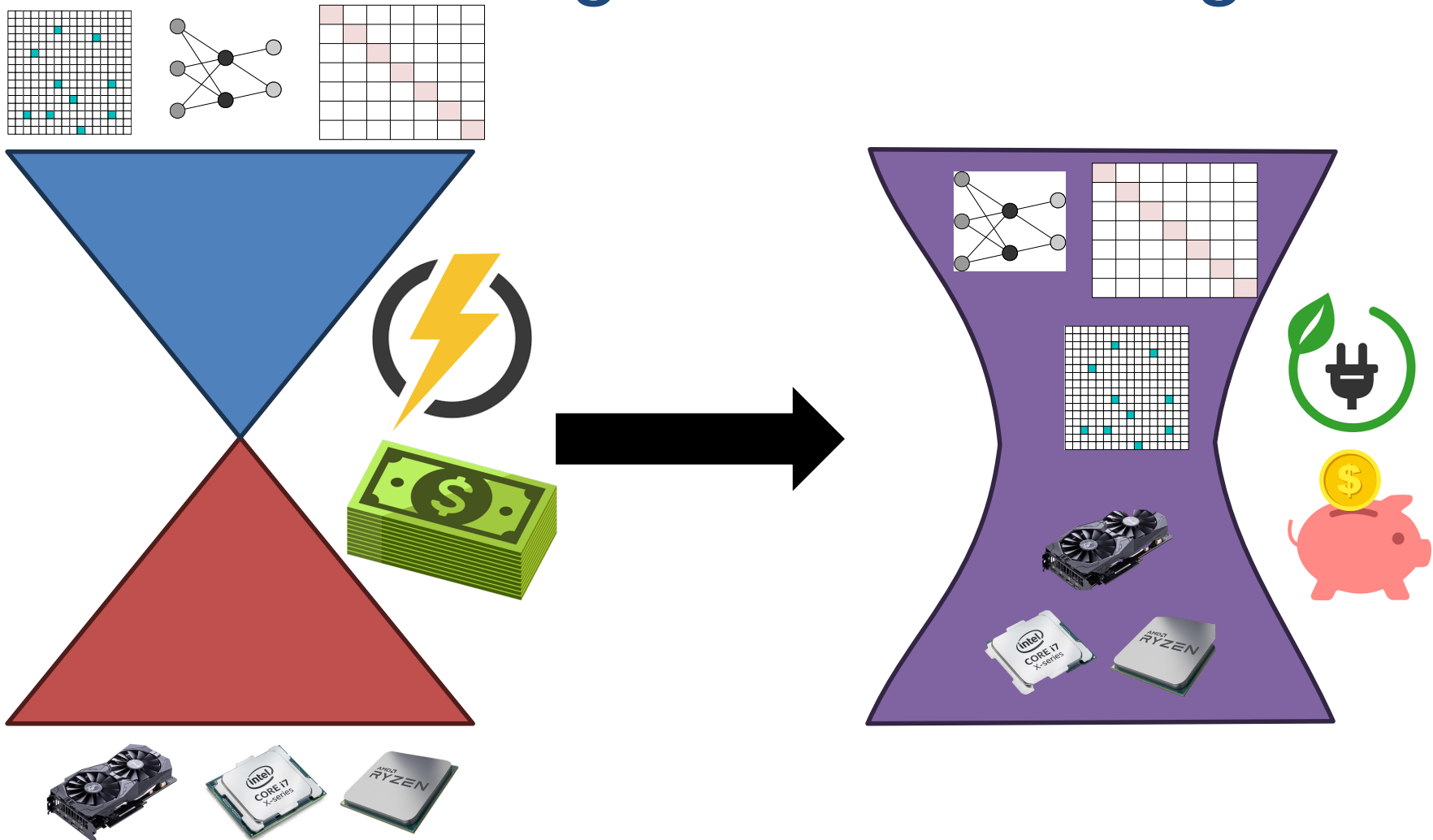


5 years of 3K
UK households



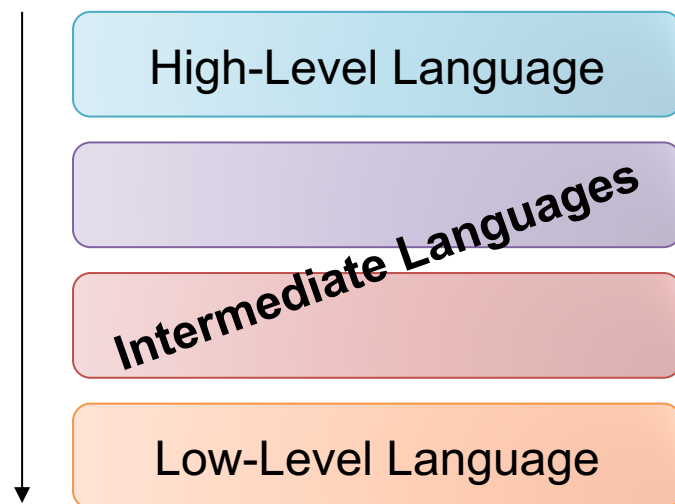
80M £

Revolutionalizing Data Processing



Language Design

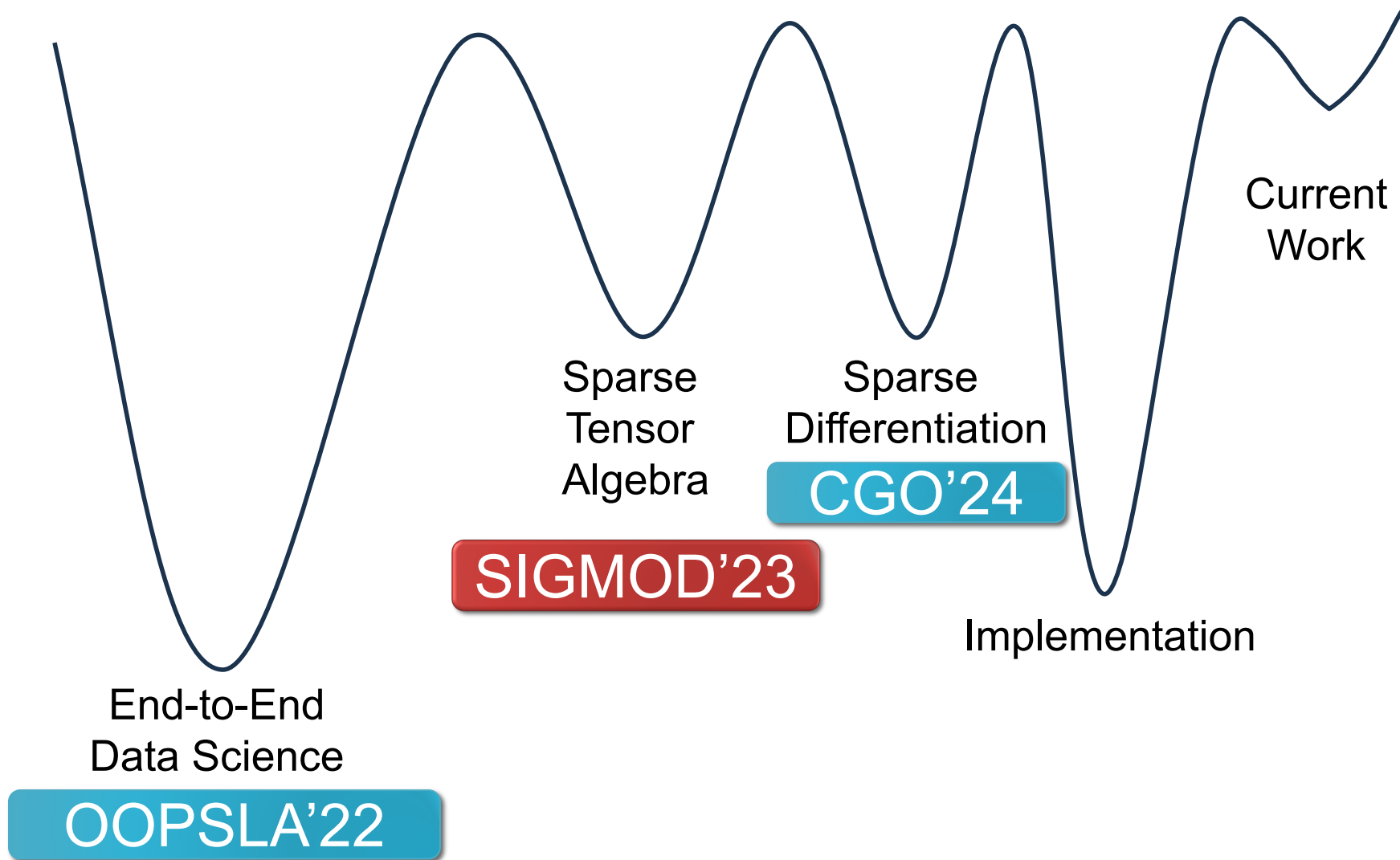
- Domain-Specific Languages (DSLs)
 - Languages for a particular domain
- Exploit
 - Domain knowledge
 - Algebraic structure
 - Structure of data
 - Algorithmic knowledge
 - Specialized data-structures



Databases

Programming Languages

Outline



END-TO-END DATA SCIENCE

Functional Collection Programming with Semi-ring Dictionaries

AMIR SHAIKHHA, University of Edinburgh, United Kingdom

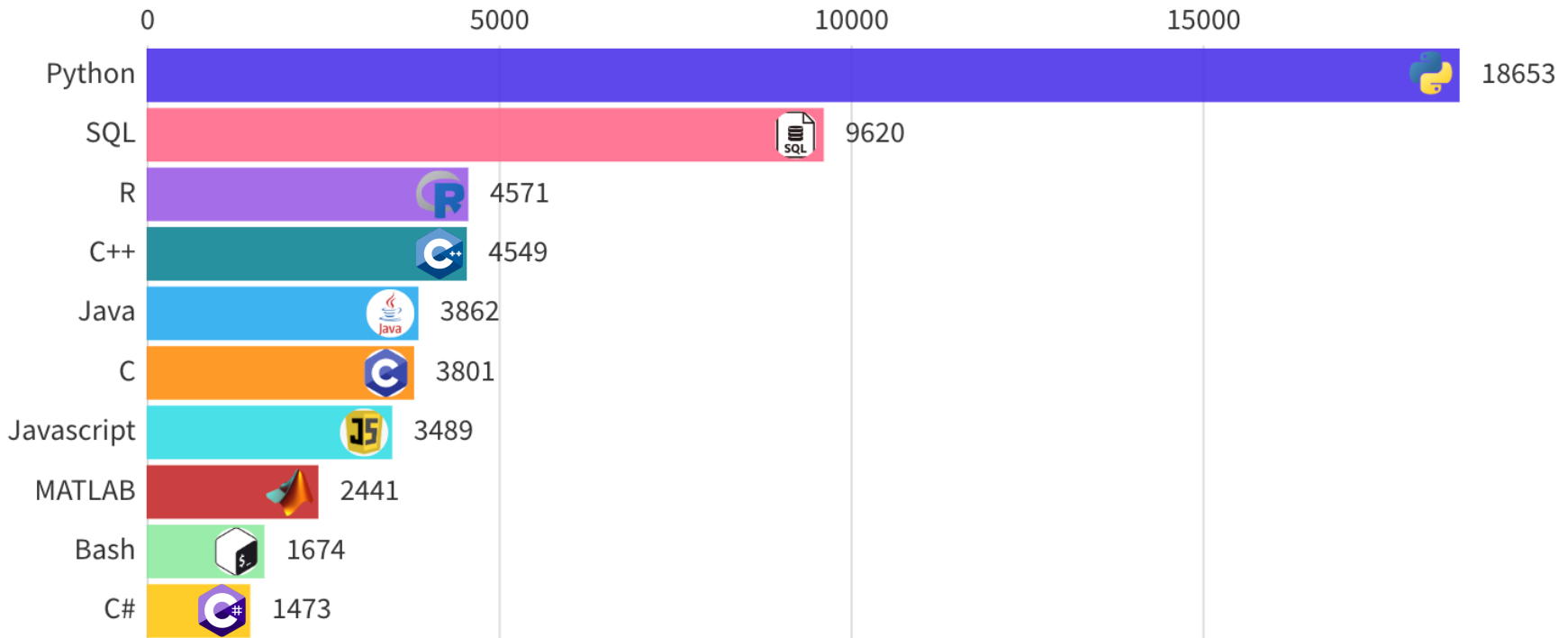
MATHIEU HUOT, University of Oxford, United Kingdom

JACLYN SMITH, University of Oxford, United Kingdom

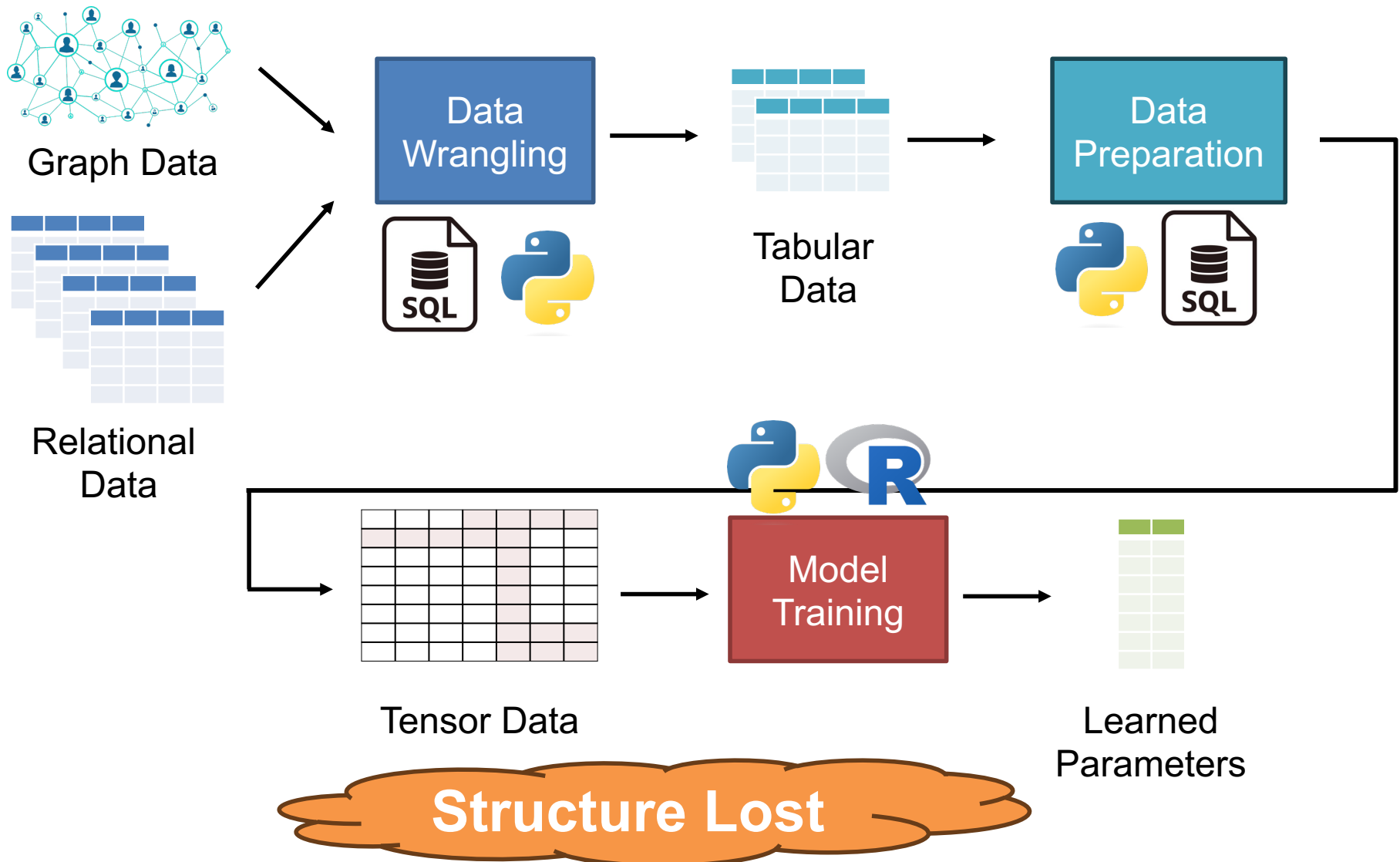
DAN OLTEANU, University of Zurich, Switzerland

OOPSLA'22

Data Science PLs



End-to-End Data Science



Data Science Workloads

DB Workloads

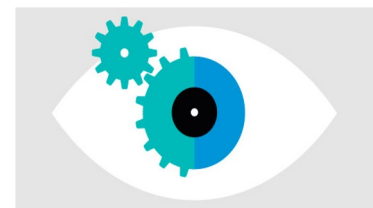


Data Warehouses (OLAP)

LA Workloads



Machine Learning



Computer Vision

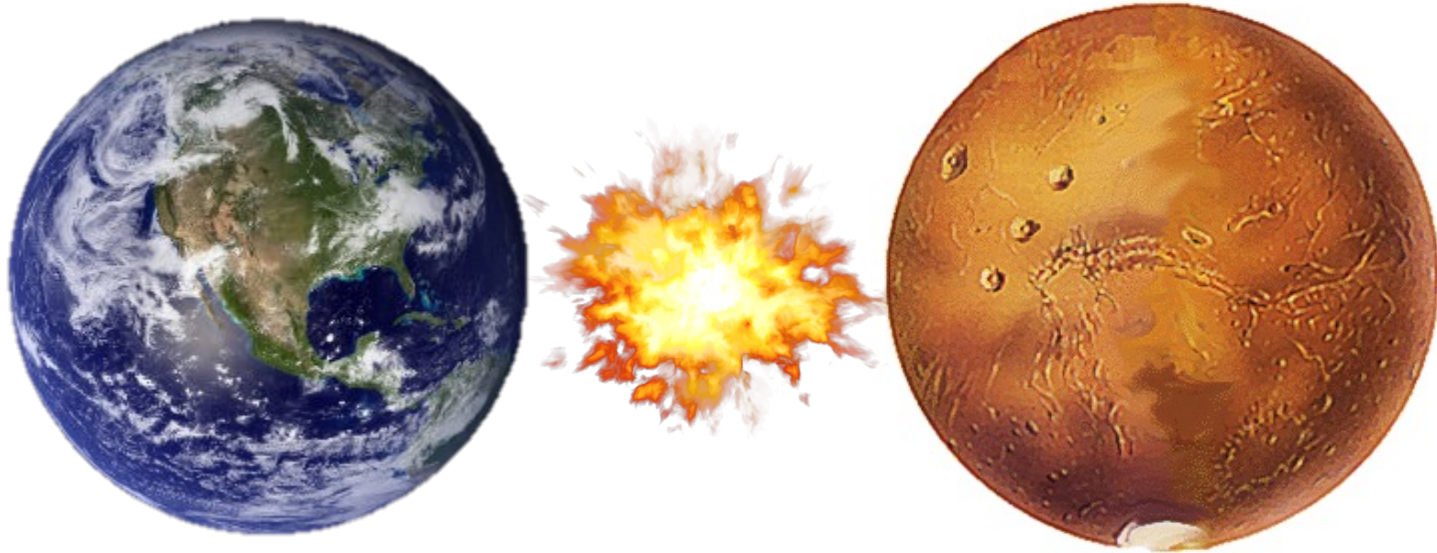


Scientific Computing



Graph Processing

Data Science Workloads



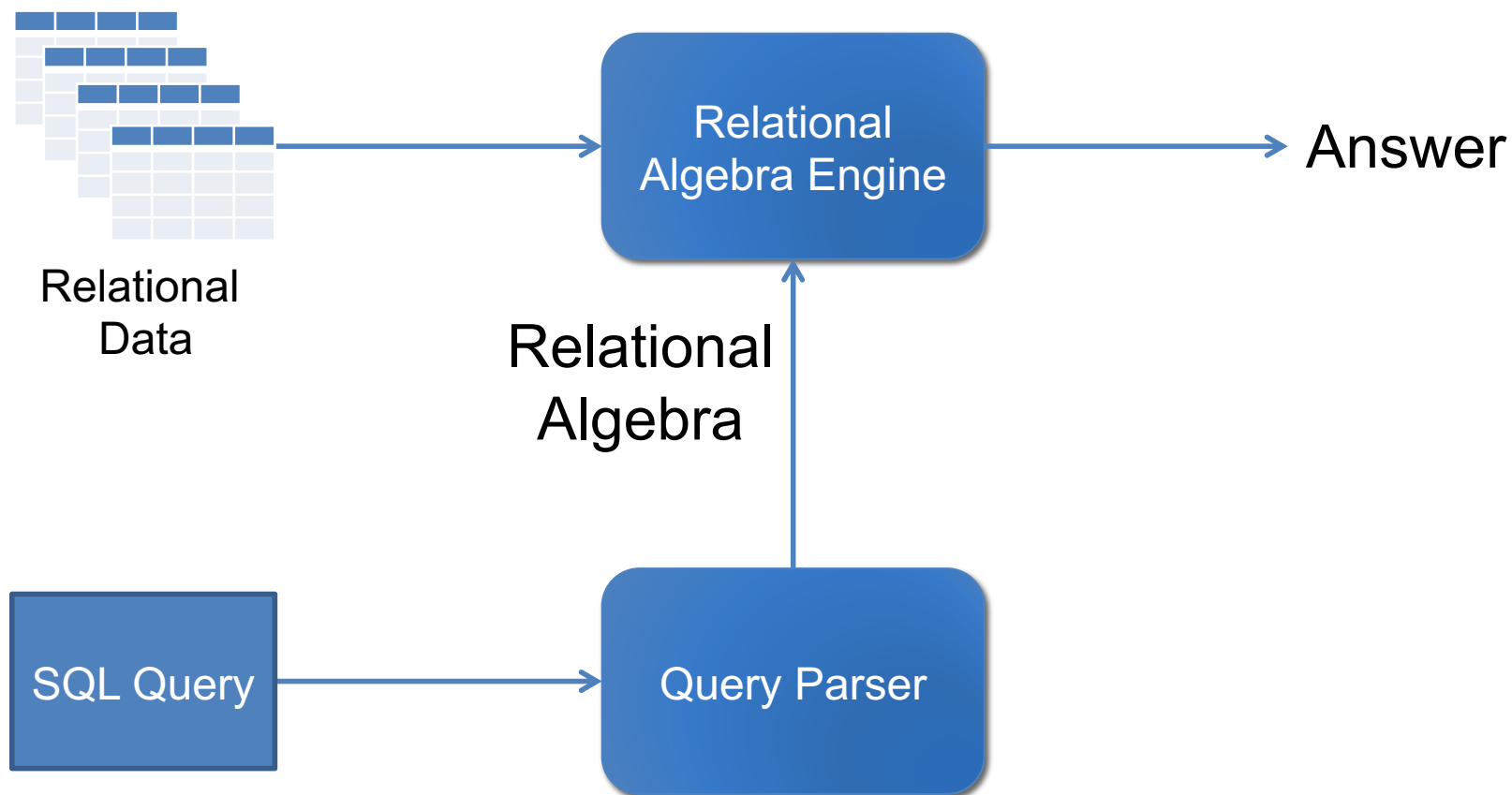
DB Workloads

Relational Algebra
Nested Relational Algebra
RDBMS, Pandas DataFrame

LA Workloads

Linear Algebra
Tensor Algebra
TensorFlow, PyTorch, scipy

Relational DB



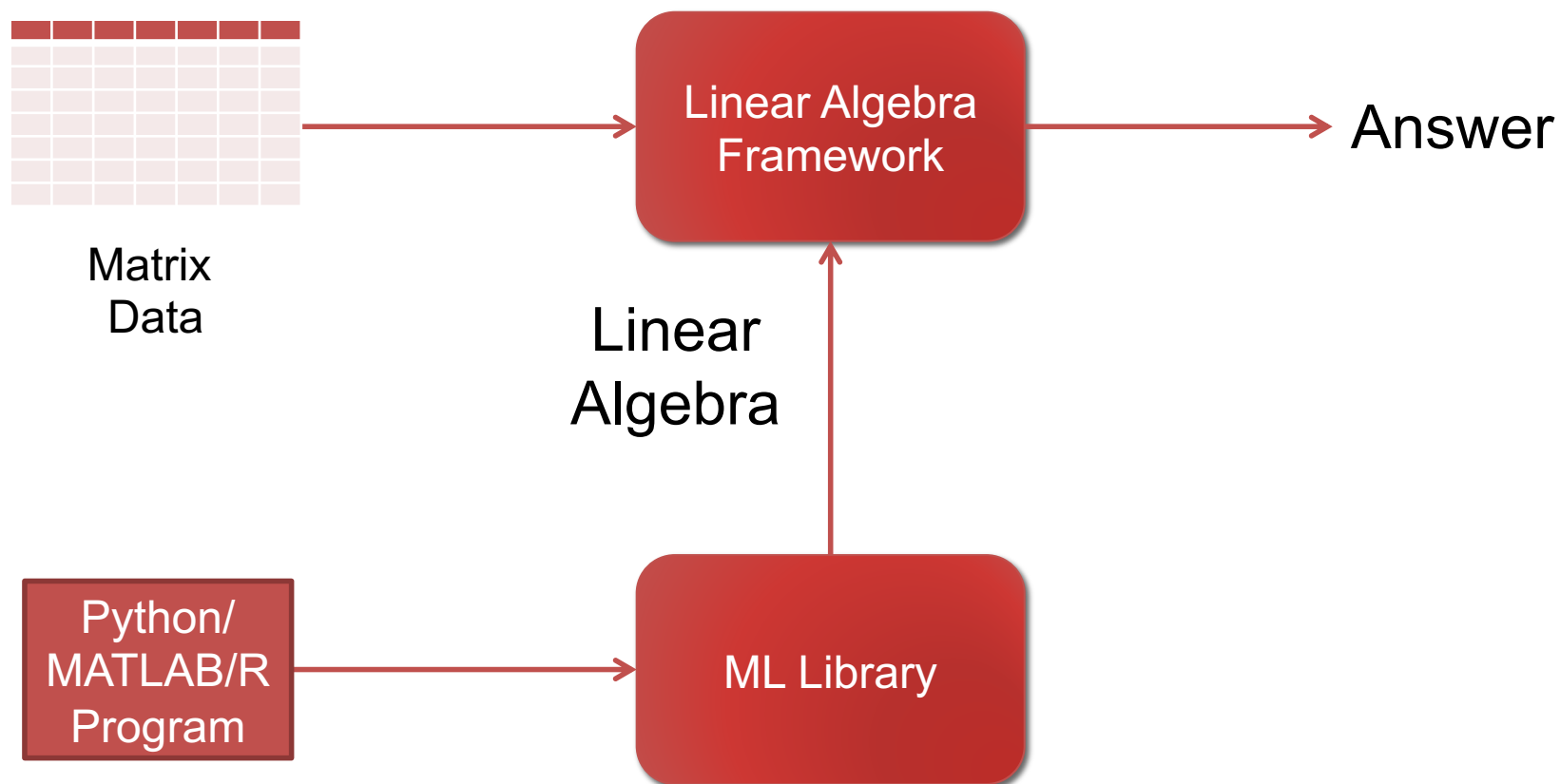
Relational Algebra

- An algebra for relational databases
- Selection (σ)
 - Filters out all tuples that do not satisfy a predicate
- Projection (π)
 - Filters out unnecessary columns of a relation
- Join (\bowtie)
 - Combines the tuples of two relations
 - A complex operator
- Group-By Aggregation (Γ)
 - Partitions data and aggregates!
 - Another complex operator

Relational Algebra Optimizations

- $\sigma_{c_1}(\sigma_{c_2}(R)) = \sigma_{c_2}(\sigma_{c_1}(R))$
- $\sigma_{c_1 \wedge \dots \wedge c_n}(R) = \sigma_{c_1}(\dots(\sigma_{c_n}(R))\dots)$
- $\pi_{a_1}(R) = \pi_{a_1}(\dots(\pi_{a_n}(R))\dots)$
- $R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$
- $R \bowtie S = S \bowtie R$
- $\sigma_{c_1 \wedge \dots \wedge c_n}(R \bowtie S) = \sigma_{c_1 \wedge \dots \wedge c_k}(R) \bowtie \sigma_{c_{p+1} \wedge \dots \wedge c_n}(S)$
- ...

ML Frameworks



Linear Algebra Optimizations

- $M1 + M2 = M2 + M1$
- $M1 + 0 = 0 + M1 = M1$
- $M \times I = I \times M = M$
- $M \times 0 = 0 \times M = 0$
- $M1 \times (M2 \times M3) = (M1 \times M2) \times M3$
- $M1 \times (M2 + M3) = M1 \times M2 + M1 \times M3$
- ...

Can we have a
unified environment?

Similarity of Optimizations

$$Q(a, d) = \Gamma_{a,d}^{\#} R_1(a, b) \bowtie R_2(b, c) \bowtie R_3(c, d)$$

$$N(i, l) = \sum_{j,k} M_1(i, j) \cdot M_2(j, k) \cdot M_3(k, l)$$



$$Q'(a, c) = \Gamma_{a,c}^{\#} R_1(a, b) \bowtie R_2(b, c) \qquad Q(a, d) = \Gamma_{a,d}^{\#} Q'(a, c) \bowtie R_3(c, d)$$

$$N'(i, k) = \sum_j M_1(i, j) \cdot M_2(j, k) \qquad N(i, k) = \sum_k N'(i, k) \cdot M_3(k, l)$$

Pushing aggregates past joins

Matrix chain ordering

SDQL



Semi-Ring Dictionary Query Language

Semi-Ring

$\langle R, 0, 1, +, \times \rangle$

$\forall a, b, c \in R$

- $a+0=a$
- $a+b=b+a$
- $(a + b) + c = a + (b + c)$
- $a \times 1 = 1 \times a = a$
- $a \times 0 = 0 \times a = 0$
- $(a \times b) \times c = a \times (b \times c)$
- $a \times (b + c) = (a \times b) + (a \times c)$
- $(a + b) \times c = (a \times c) + (b \times c)$

Factorization

Semi-Ring Examples

 $\langle \mathbb{R}, 0, 1, +, \times \rangle$ $\langle \mathbb{N}, 0, 1, +, \times \rangle$  $\langle \{\text{false}, \text{true}\},$
 $\text{false}, \text{true}, \vee, \wedge \rangle$  $\langle \mathbb{R} \cup \{+\infty\},$
 $+\infty, 0, \min, + \rangle$ 

Semi-Ring Dictionaries

One collection to rule them all

`{ key -> value }` 

`Relation[T] = { T -> Bool }` *(no duplicates)*

`Relation[T] = { T -> Int }` *(with duplicates)*

`Vector[T] = { Int -> T }`

`Matrix[T] = { (Int, Int) -> T }`

Database Relations (Bag Semantics)

Relation $R(A,B)$

A	B
a_1	b_1
a_1	b_1
a_2	b_1
a_2	b_1
a_2	b_2

A	B	\rightarrow	$R(A, B)$
a_1	b_1	\rightarrow	2
a_2	b_1	\rightarrow	2
a_2	b_2	\rightarrow	1

{ tuple \rightarrow multiplicity }

Linear Algebra (Matrix)

Matrix M

	0	1	2
0	m_1	0	0
1	0	0	m_2
2	0	0	0
3	0	m_3	0

row	col	→	$M_{row,col}$
0	0	→	m_1
1	2	→	m_2
3	1	→	m_3

{ index -> value }

SDQL Examples

SDQL

```
sum(<key, val> in R)
  f(key, val)
```

```
sum(<key, val> in R)
  { g(key) -> f(val) }
```

C++

```
double res = 0;
for(auto&e : R) {
    res += f(e.key, e.val)
}
```

```
dict<K,V> res = dict<K,V>();
for(auto&e : R) {
    res[g(e.key)] += f(e.val)
}
```

Aggregations over Relations (Bag)

```
SELECT COUNT (*) FROM R
```

```
sum(<key, val> in R) val
```

```
SELECT SUM(A) FROM R
```

```
sum(<key, val> in R) key.A * val
```

```
SELECT B, SUM(A) FROM R GROUP BY B
```

```
sum(<key, val> in R) { key.B -> key.A * val }
```

Relational Algebra to SDQL

$$\begin{aligned}
 \llbracket \sigma_p(R) \rrbracket &= \text{sum}(x \leftarrow \llbracket R \rrbracket) \text{if}(p(x.\text{key}))\{ x.\text{key} \} \text{ else } \{ \} \\
 \llbracket \pi_f(R) \rrbracket &= \text{sum}(x \leftarrow \llbracket R \rrbracket) \{ f(x.\text{key}) \} \\
 \llbracket R \cup S \rrbracket &= \llbracket R \rrbracket + \llbracket S \rrbracket \\
 \llbracket R \cap S \rrbracket &= \text{sum}(x \leftarrow \llbracket R \rrbracket) \text{if}(\llbracket S \rrbracket(x.\text{key}))\{ x.\text{key} \} \text{ else } \{ \} \\
 \llbracket R - S \rrbracket &= \text{sum}(x \leftarrow \llbracket R \rrbracket) \text{if}(\llbracket S \rrbracket(x.\text{key}))\{ \} \text{ else } \{ x.\text{key} \} \\
 \llbracket R \times S \rrbracket &= \text{sum}(x \leftarrow \llbracket R \rrbracket) \text{sum}(y \leftarrow \llbracket S \rrbracket) \\
 &\quad \{ \text{concat}(x.\text{key}, y.\text{key}) \} \\
 \llbracket R \bowtie_{\theta} S \rrbracket &= \llbracket \sigma_{\theta}(R \times S) \rrbracket \\
 \llbracket \Gamma_{\theta;f}(e) \rrbracket &= \text{sum}(x \leftarrow \llbracket e \rrbracket) x.\text{val} * \llbracket f \rrbracket(x.\text{key}) \\
 \llbracket \Gamma_{g;f}(e) \rrbracket &= \text{let tmp} = \text{sum}(x \leftarrow \llbracket e \rrbracket) \{ \llbracket g \rrbracket(x.\text{key}) \rightarrow x.\text{val} * \llbracket f \rrbracket(x.\text{key}) \} \\
 &\quad \text{in sum}(x \leftarrow \text{tmp}) \{ \langle \text{key}=x.\text{key}, \text{val}=x.\text{val} \rangle \rightarrow 1 \}
 \end{aligned}$$

Vector Operations

$$V1 + V2$$
$$v1 + v2$$
$$V1 .* V2$$

```
sum(<key, val> in V1) { key -> val * V2(key) }
```

$$V1 \cdot V2$$

```
sum(<key, val> in V1) val * V2(key)
```

Linear Algebra to SDQL

$$\begin{aligned}
 \llbracket V_1 + V_2 \rrbracket &= \llbracket V_1 \rrbracket + \llbracket V_2 \rrbracket \\
 \llbracket a \cdot V \rrbracket &= \llbracket a \rrbracket * \llbracket V \rrbracket \\
 \llbracket V_1 \circ V_2 \rrbracket &= \text{sum}(x \text{ in } \llbracket V_1 \rrbracket) \{ x.\text{key} \rightarrow x.\text{val} * \llbracket V_2 \rrbracket(x.\text{key}) \} \\
 \llbracket V_1 \cdot V_2 \rrbracket &= \text{sum}(x \text{ in } \llbracket V_1 \rrbracket) x.\text{val} * \llbracket V_2 \rrbracket(x.\text{key}) \\
 \llbracket \sum_{a \in V} a \rrbracket &= \text{sum}(x \text{ in } \llbracket V \rrbracket) x.\text{val} \\
 \\
 \llbracket M_1^T \rrbracket &= \text{sum}(\text{row in } \llbracket M_1 \rrbracket) \text{sum}(x \text{ in row.val}) \\
 &\quad \{ x.\text{key} \rightarrow \{ \text{row.key} \rightarrow x.\text{val} \} \} \\
 \llbracket M_1 \circ M_2 \rrbracket &= \text{sum}(\text{row in } \llbracket M_1 \rrbracket) \{ \text{row.key} \rightarrow \\
 &\quad \text{sum}(x \text{ in row.val}) \{ x.\text{key} \rightarrow \\
 &\quad \quad x.\text{val} * \llbracket M_2 \rrbracket(\text{row.key})(x.\text{key}) \} \} \\
 \llbracket M_1 \times M_2 \rrbracket &= \text{sum}(\text{row in } \llbracket M_1 \rrbracket) \{ \text{row.key} \rightarrow \\
 &\quad \text{sum}(x \text{ in row.val}) \text{sum}(y \text{ in } \llbracket M_2 \rrbracket(x.\text{key})) \\
 &\quad \{ y.\text{key} \rightarrow x.\text{val} * y.\text{val} \} \} \\
 \llbracket M \cdot V \rrbracket &= \text{sum}(\text{row in } \llbracket M \rrbracket) \{ \text{row.key} \rightarrow \\
 &\quad \text{sum}(x \text{ in row.val}) x.\text{val} * \llbracket V \rrbracket(x.\text{key}) \} \\
 \llbracket \text{Trace}(M) \rrbracket &= \text{sum}(\text{row in } \llbracket M \rrbracket) \text{row.val}(r.\text{key})
 \end{aligned}$$

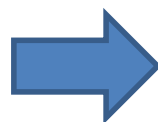
Loop Optimizations

- Vertical Loop Fusion
- Horizontal Loop Fusion
- Loop-Invariant Code Motion (Hoisting)
- Loop Factorization
- Loop Memoization

Loop Memoization & Hoisting

```

sum(<r,r_v> in R)
  sum(<s,s_v> in S)
    if(jkR(r)==jkS(s)) then
      { concat(r,s)->r_v*s_v }
  
```



```

sum(<r,r_v> in R)
  let Sp = sum(<s,s_v> in S)
    { jkS(s) -> {s->s_v} } in
  sum(<s,s_v> in Sp(jkR(r)))
    { concat(r,s)->r_v*s_v }
  
```



```

let Sp = sum(<s,s_v> in S)
  { jkS(s) -> {s->s_v} } in
sum(<r,r_v> in R)
  sum(<s,s_v> in Sp(jkR(r)))
    { concat(r,s)->r_v*s_v }
  
```

Nested Loop Join -> Hash Join

Uniform Optimization

- Vertical Loop Fusion

Pipeline Query Engine

Deforestation, Pull/Push Arrays

- Horizontal Loop Fusion

Multi-aggregate Operator

Horizontal Fusion

- Loop Factorization + Memoization

Hash Join, Group Join

Matrix chain ordering

Data Layouts

- Relations
 - Row/Columnar layout
 - Standard Dictionary
 - Factorized

- Tensors
 - Dense (Row/Col Major)
 - COO
 - Compressed

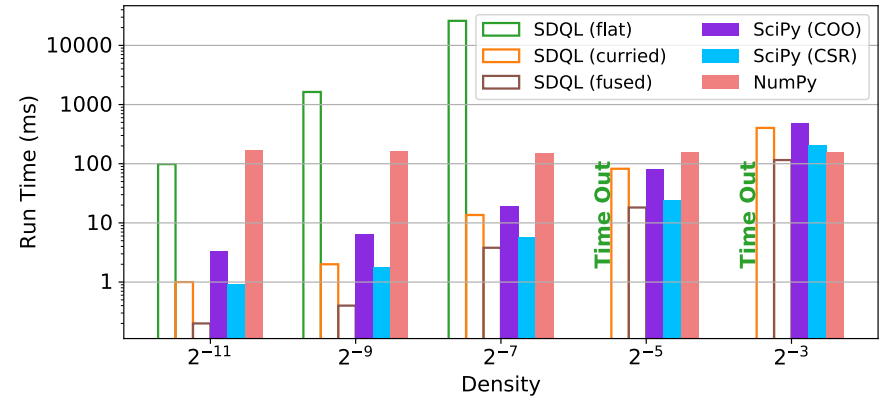
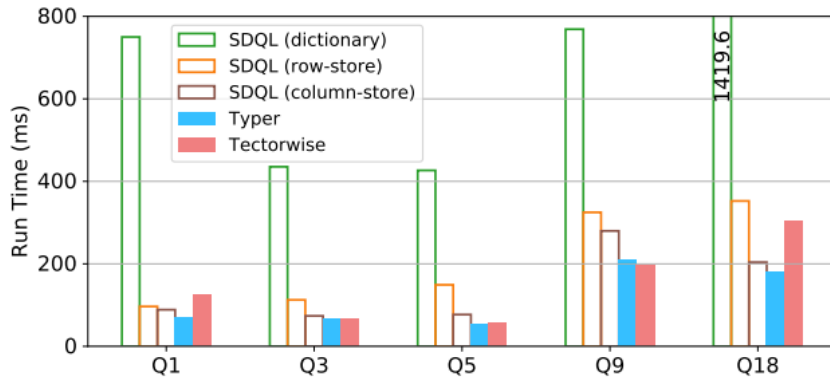
Dictionary		Factorized			Row		Columnar							
$\langle A=a_1, B=b_1 \rangle$	1	a_1	b_1	1	0	$\langle A=a_1, B=b_1 \rangle$	$\langle A=$	0	a_1	,	$B=$	0	b_1	\rangle
$\langle A=a_1, B=b_2 \rangle$	1		b_2	1	1	$\langle A=a_1, B=b_2 \rangle$		1	a_1			1	b_2	
$\langle A=a_2, B=b_3 \rangle$	1	a_2	b_3	1	2	$\langle A=a_2, B=b_3 \rangle$		2	a_2			2	b_3	

Semi-Ring Dictionaries

One collection to rule them all

- Relations
 - `Bag{T}` = `Dict{T, Int}`
 - `Set{T}` = `Dict{T, Bool}`
- Nested Relations
 - `Bag{Bag{T}}` = `Dict{Dict{T, Int}, Int}`
 - `Set{Set{T}}` = `Dict{Dict{T, Bool}, Bool}`
- Tensors
 - `SparseVector{T}` = `Dict{Int, T}`
 - `SparseMatrixCOO{T}` = `Dict{(Int, Int), T}`
 - `SparseMatrixTrie{T}` = `Dict{Int, Dict{Int, T}}`
 - `DenseVector{T}` = `Dict{DInt, T}`
 - `DenseMatrix{T}` = `Dict{DInt, Dict{DInt, T}}`

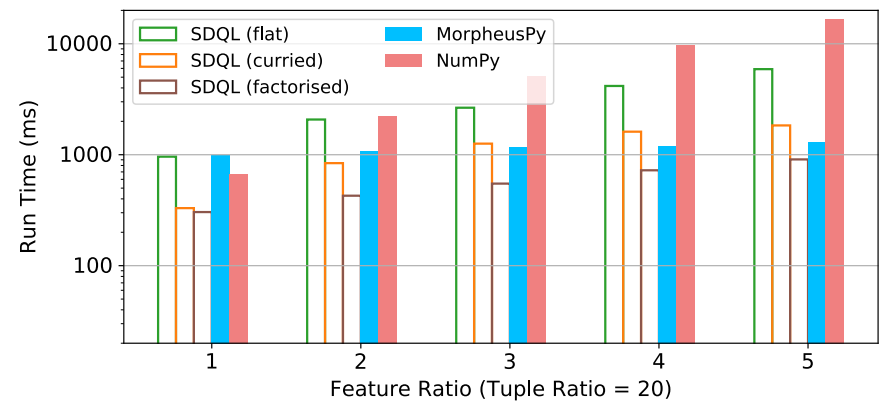
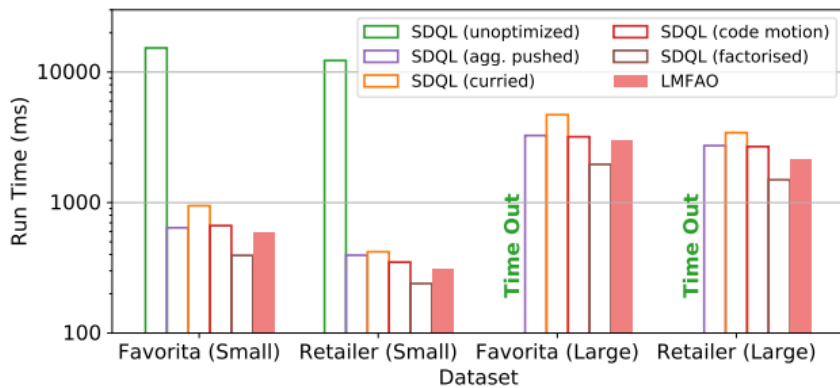
Experimental Results



DB

LA

DB + LA



Competitive with (or better than) specialized systems

SPARSE TENSOR ALGEBRA

Optimizing Tensor Programs on Flexible Storage

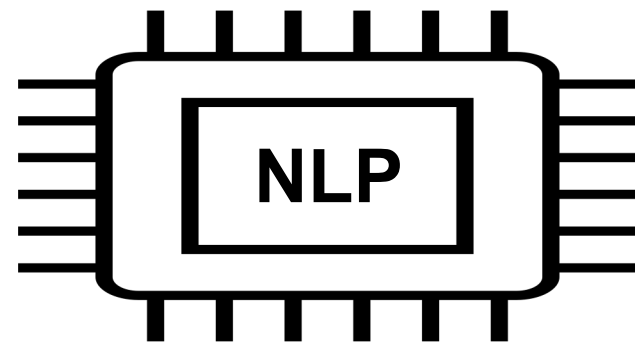
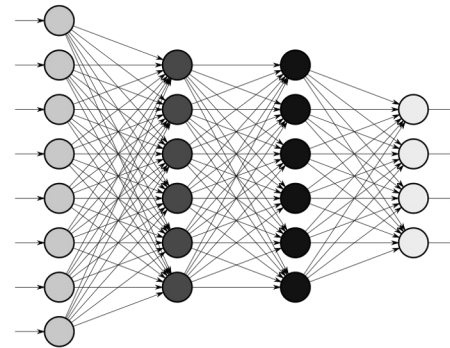
MAXIMILIAN SCHLEICH, RelationalAI, USA

AMIR SHAIKHHA, University of Edinburgh, United Kingdom

DAN SUCIU, University of Washington, USA

SIGMOD'23

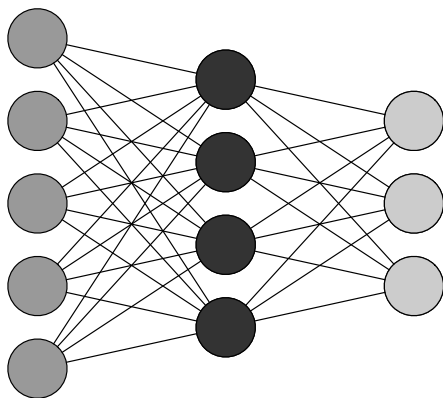
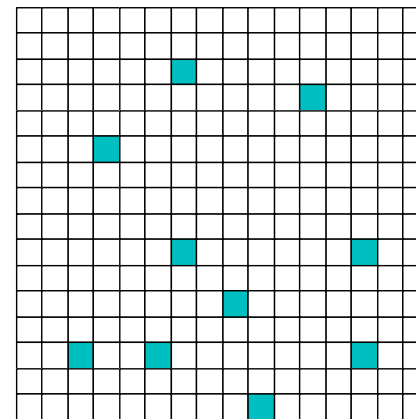
Tensors



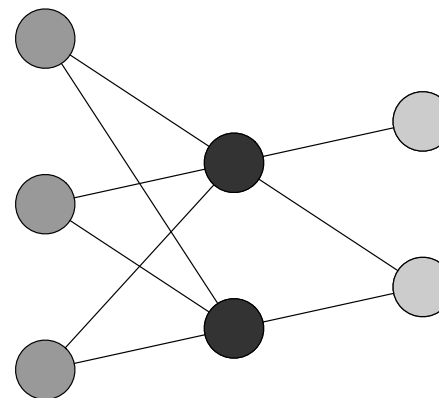
Sparse Tensors



Adjacency
Matrix



Sparsification



Sparse Matrix as Relation

A

	0	1	2	3	4	5
0	5	1				
1	7	3				
2						
3	8			4	9	

$$C_{ik} = \sum_j A_{ij} B_{jk}$$

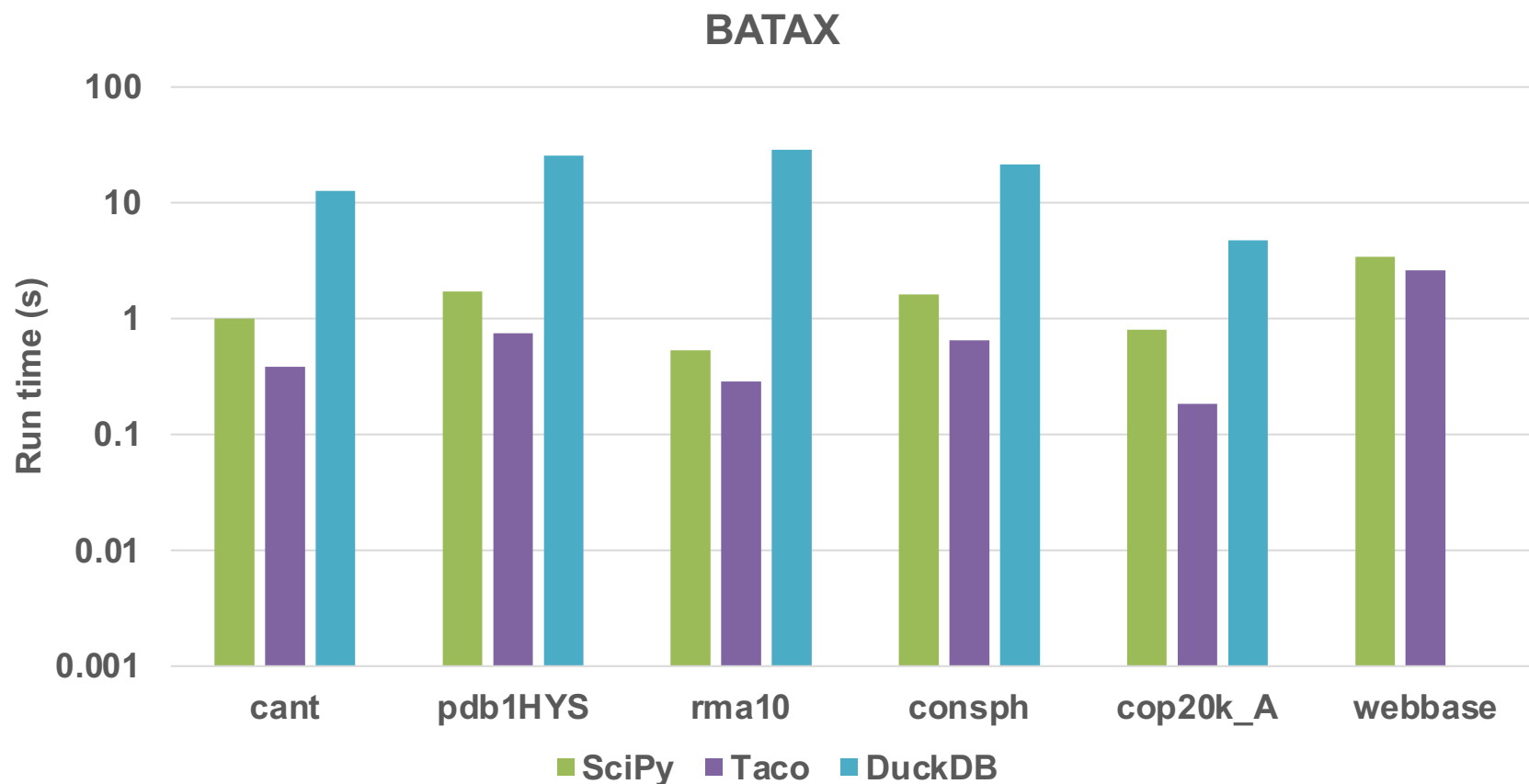
Row	Col	Val
0	0	5
0	1	1
1	0	7
1	1	3
3	0	8
3	3	4
3	4	9

```

SELECT A.row, B.col,
       SUM(A.val*B.val)
FROM A, B
WHERE A.col = B.row
GROUP BY A.row, B.col

```


Are Database Engines Competitive?



No, because they do not support optimized storage formats for sparse tensors

Sparse Storage Formats



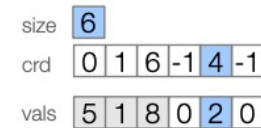
(a) An 8-vector



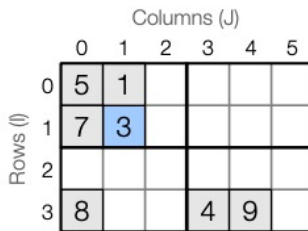
(b) Dense array



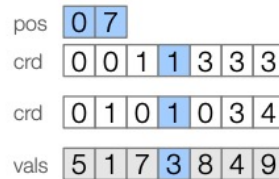
(c) Sparse vector



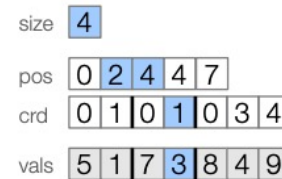
(d) Hash map



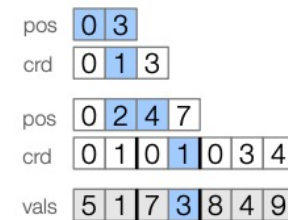
(e) A 4x6 matrix



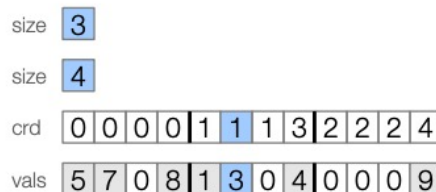
(f) COO



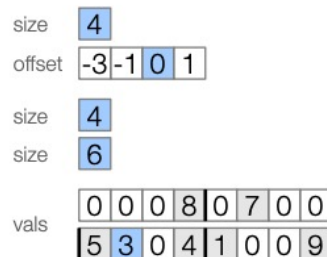
(g) CSR



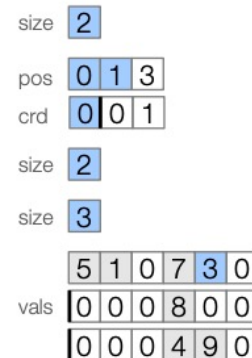
(h) DCSR



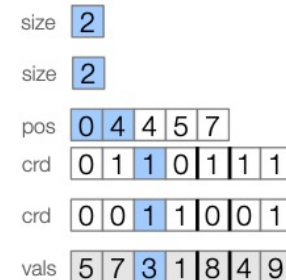
(i) ELL



(j) DIA

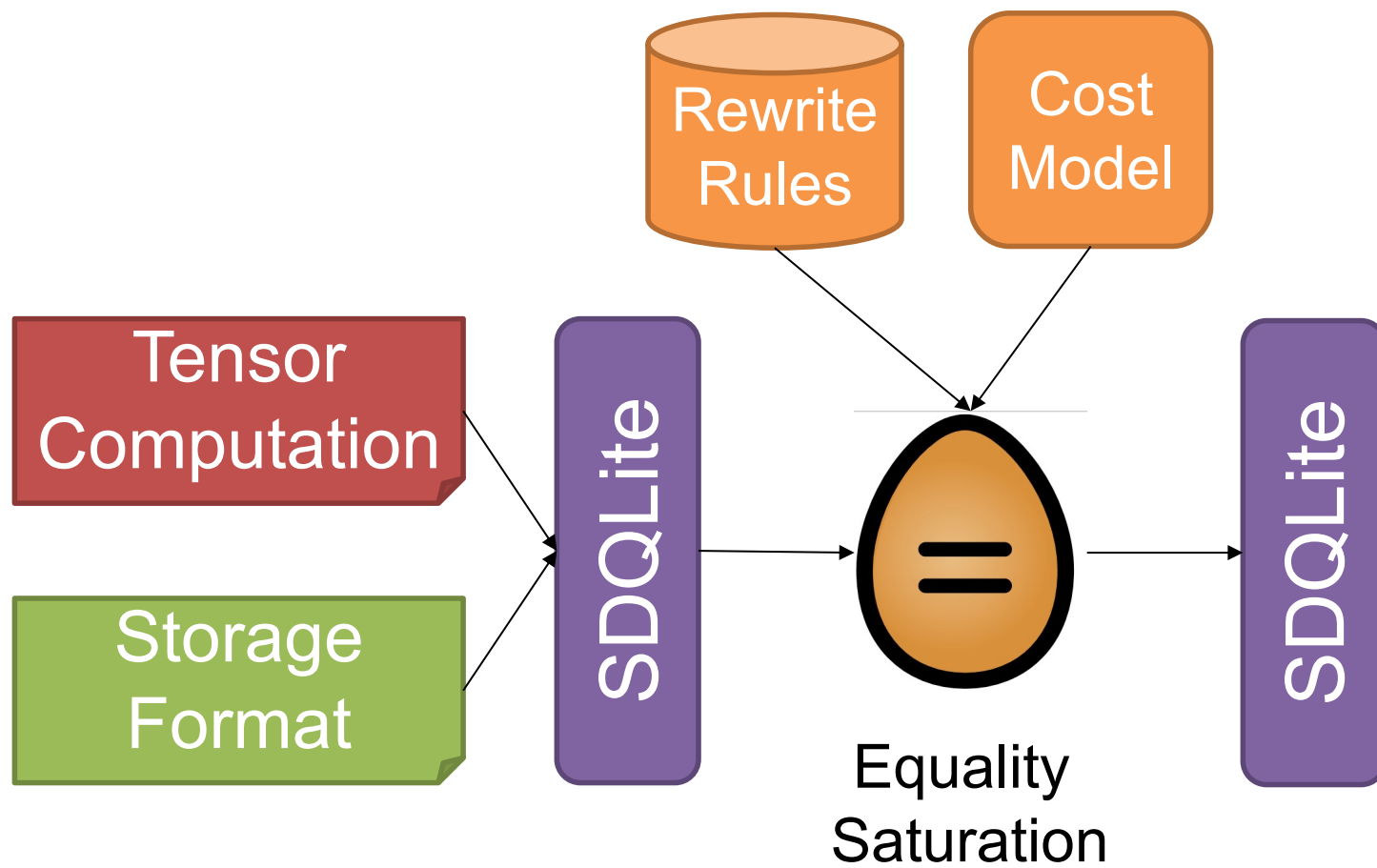


(k) BCSR

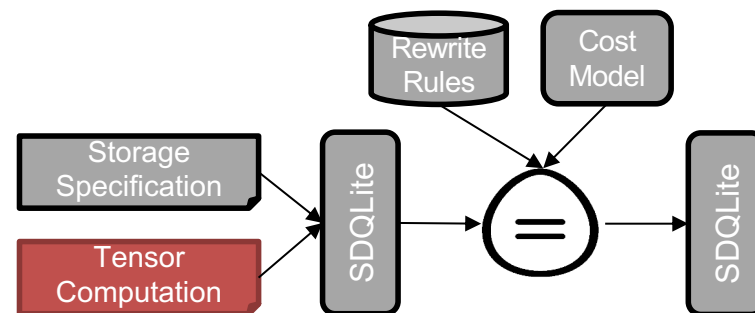


(l) CSB

Flexible Storage Formats in SDQLite



Tensor Computation

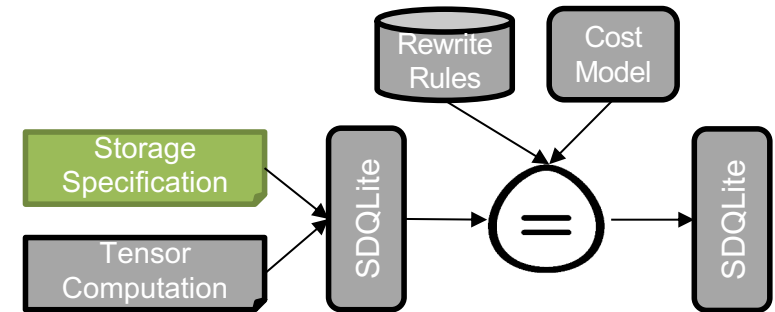


$$C_{ik} = \sum_j A_{ij} B_{jk}$$

```

sum (<(i,j), A_v> in A,
     <(j,k), B_v> in B)
{ (i,k) -> A_v*B_v }
  
```

Storage Specification



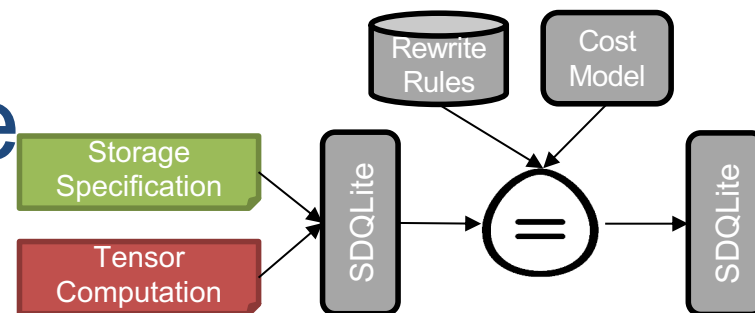
size	4						
pos	0	2	4	4	4	7	
cols	0	1	0	1	0	3	4
vals	5	1	7	3	8	4	9

	0	1	2	3	4	5
0	5	1				
1	7	3				
2						
3	8			4	9	

```

sum(<r, _> in 0:size)
  sum(<i, c> in cols(pos(r):pos(r+1)))
    { (r, c) -> vals(i) }
  
```

Computation+Storage



```

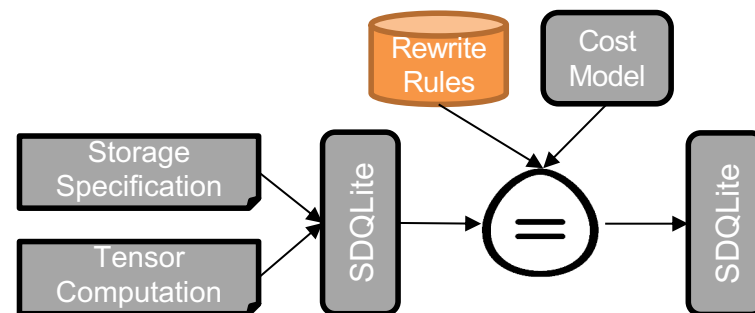
let A =
  sum(<r,_> in 0:A_size)
    sum(<i,c> in A_cols(A_pos(r):A_pos(r+1)))
      { (r,c) -> A_vals(i) } in
let B =
  sum(<r,_> in 0:B_size)
    sum(<i,c> in B_cols(B_pos(r):B_pos(r+1)))
      { (r,c) -> B_vals(i) } in
sum(<(i,j),A_v> in A, <(j,k),B_v> in B)
  { (i,k) -> A_v*B_v }

```

Inefficient

Rewrite Rules

- We have 44 rewrite rules
- Loop Fusion



let A =

```
sum(<k1,v1> in D)
  { k1 -> f(k1,v1) } in
```

```
sum(<k2,v2> in A)
  g(k2,v2)
```



```
sum(<k1,v1> in D)
  let v2 = f(k1,v1) in
  g(k1,v2)
```

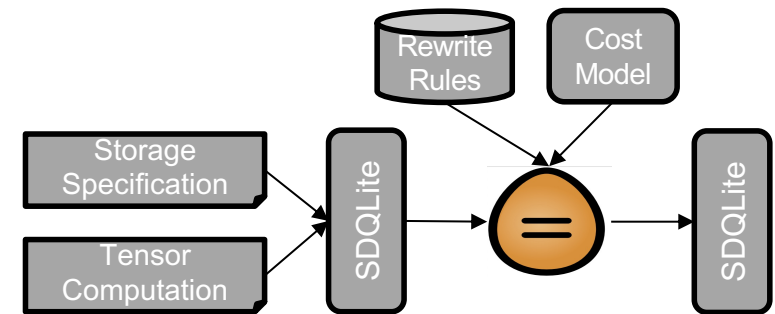
- Factorization

```
sum(<k,v> in A)
  e * f(k,v)
```

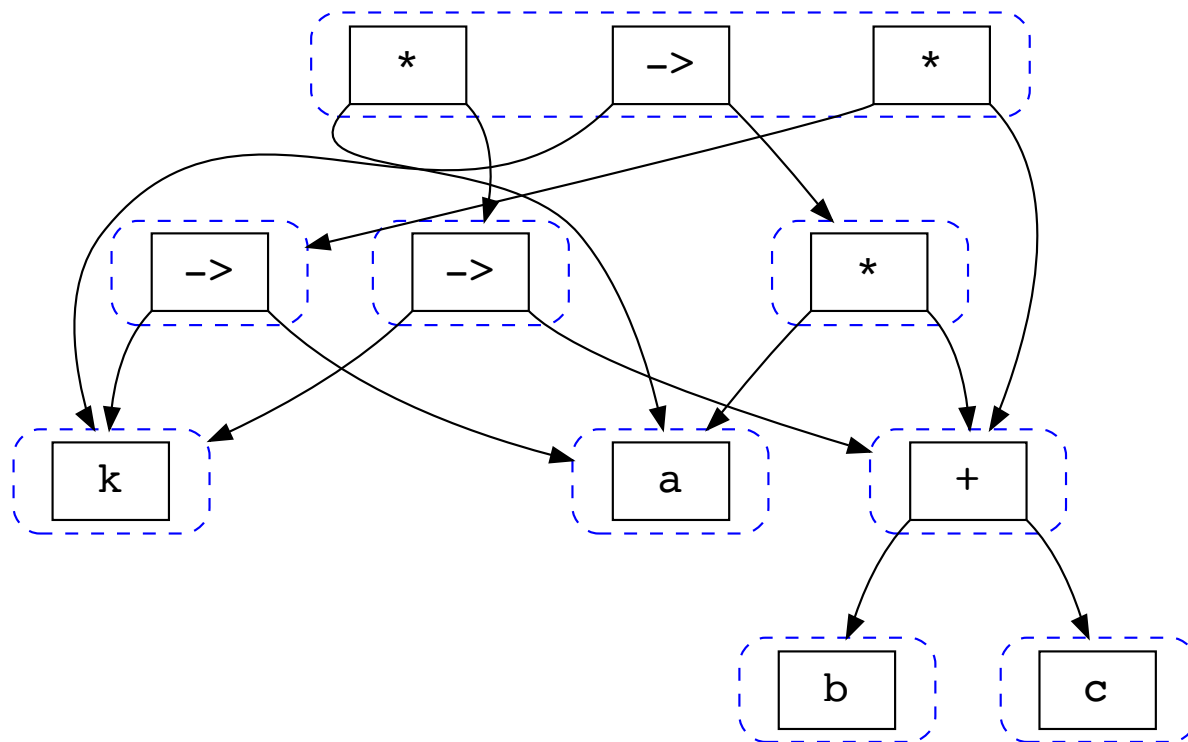


```
e * sum(<k,v> in A)
      f(k,v)
```

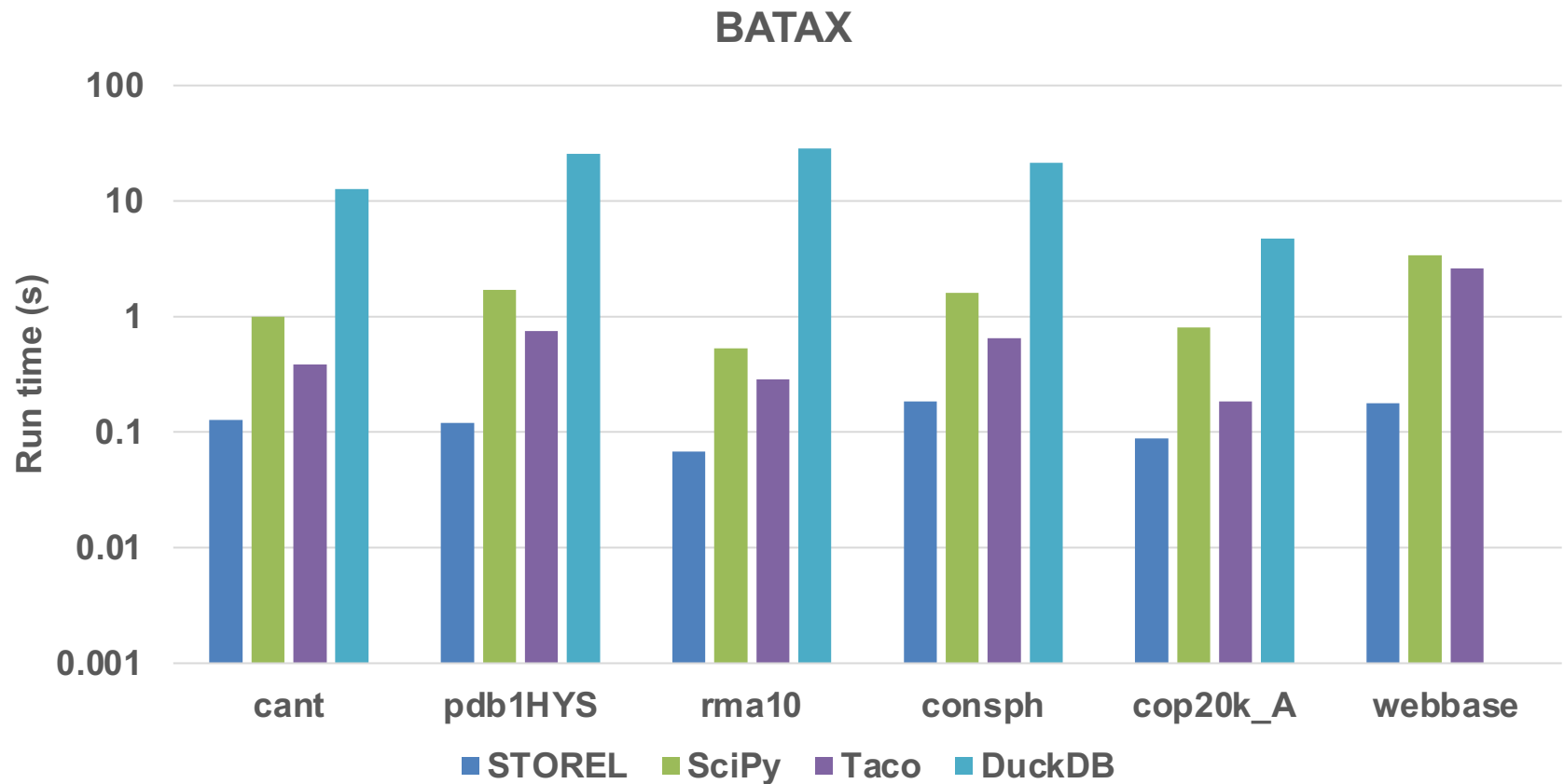
Equality Saturation



- E-graph: Compressed representation of Search Space



Performance Results



Optimizations + Compressed Storage

SPARSE DIFFERENTIATION

A Tensor Algebra Compiler for Sparse Differentiation

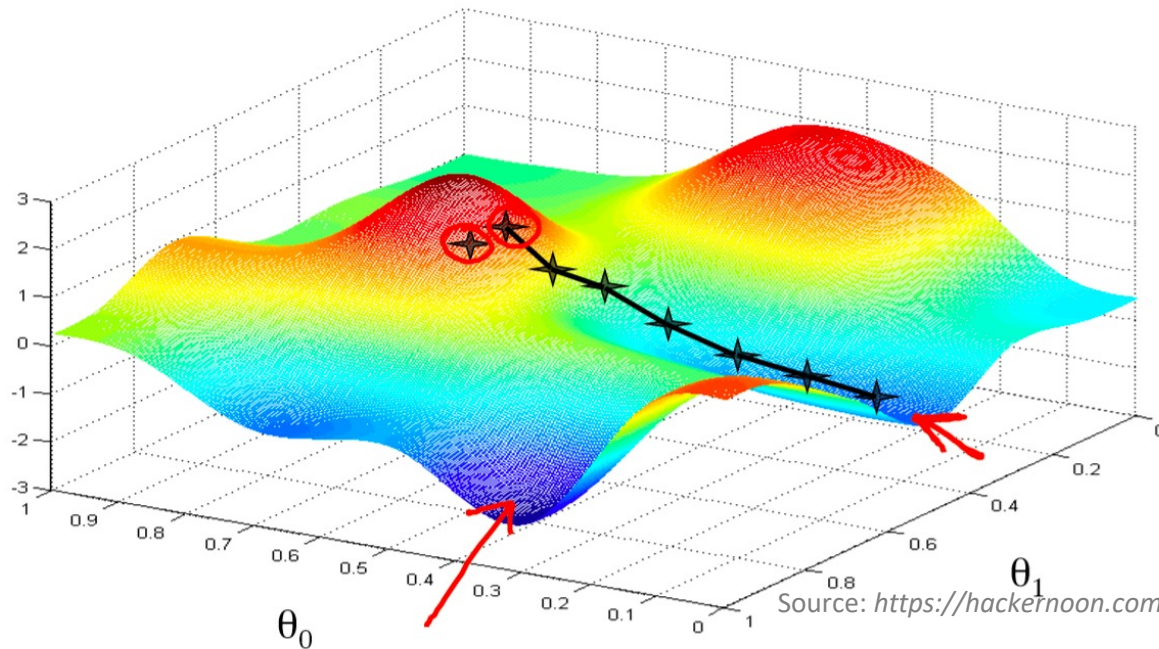
Amir Shaikhha
School of Informatics
University of Edinburgh
Edinburgh, United Kingdom
amir.shaikhha@ed.ac.uk

Mathieu Huot
Department of Computer Science
University of Oxford
Oxford, United Kingdom
mathieu.huot@stx.ox.ac.uk

Shideh Hashemian
School of Informatics
University of Edinburgh
Edinburgh, United Kingdom
s.hashemian@sms.ed.ac.uk

CGO'24

Gradient Descent Based Algorithms



Repeat until converges {

$$\theta_i = \theta_i - \alpha \frac{\partial}{\partial \theta} f(\theta_i)$$

}

Derivative of a function

Automatic Differentiation (AD)

- Differentiable Programming
- A systematic approach for computing the derivative of a function
- Functions represented as programs

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad \longrightarrow \quad \frac{\partial f}{\partial x}: \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$$

AD in a Nutshell

$$\mathcal{D}_{\mathcal{T}}[\text{Num}] = \text{Num} \times \text{Num}$$

Original

Tangent

Dual
Number

$$\mathcal{D}_{\mathcal{T}}[\text{Array}\langle M \rangle] = \text{Array}\langle \mathcal{D}_{\mathcal{T}}[M] \rangle$$

$$\mathcal{D}_{\mathcal{T}}[T_1 \Rightarrow T_2] = \mathcal{D}_{\mathcal{T}}[T_1] \Rightarrow \mathcal{D}_{\mathcal{T}}[T_2]$$

$$\mathcal{D}_{\mathcal{T}}[M_1 \times M_2] = \mathcal{D}_{\mathcal{T}}[M_1] \times \mathcal{D}_{\mathcal{T}}[M_2]$$

Sparse Differentiation

- Tensor Libraries: TensorFlow, PyTorch

```
term = lambda V1: tensordot(V1, V2, 1)
dTerm = jacobian(term, V1)
```

- Function

- Support

- Lack of

- Cryptic

- Not expressive enough



Challenges for Sparse Differentiation

- Differentiated functions require control-flow constructs / User-Defined Functions
 - Beyond tensor algebra
- Sparse tensor programs are complicated
 - Imperative loopy code (co-iteration)
 - Data formats (CSR, CSC, CSF)

Imperative Loops in Sparse TA

```

1 for (int i = 0; i < m; i++) {
2
3
4
5   for (int j = 0; j < n; j++) {
6     int pB2 = i * n + j;
7     int pA2 = i * n + j;
8
9
10    for (int k = 0; k < p; k++) {
11      int pB3 = pB2 * p + k;
12
13
14
15
16
17
18      A[pA2] += B[pB3] * c[k];
19
20
21
22    }
23  }
24 }

```

Fig. 1. $A_{ij} = \sum_k B_{ijk} c_k$

```

for (int pB1 = B1_pos[0];
    pB1 < B1_pos[1];
    pB1++) {
  int i = B1_idx[pB1];
  for (int pB2 = B2_pos[pB1];
      pB2 < B2_pos[pB1+1];
      pB2++) {
    int j = B2_idx[pB2];
    int pA2 = i * n + j;
    for (int pB3 = B3_pos[pB2];
        pB3 < B3_pos[pB2+1];
        pB3++) {
      int k = B3_idx[pB3];
      A[pA2] += B[pB3] * c[k];
    }
  }
}

```

Fig. 2. $A_{ij} = \sum_k B_{ijk} c_k$ (sparse B)

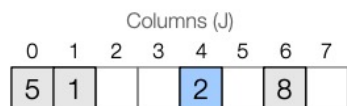
```

for (int pB1 = B1_pos[0];
    pB1 < B1_pos[1];
    pB1++) {
  int i = B1_idx[pB1];
  for (int pB2 = B2_pos[pB1];
      pB2 < B2_pos[pB1+1];
      pB2++) {
    int j = B2_idx[pB2];
    int pA2 = i * n + j;
    int pB3 = B3_pos[pB2];
    int pc1 = c1_pos[0];
    while (pB3 < B3_pos[pB2+1] &&
          pc1 < c1_pos[1]) {
      int kB = B3_idx[pB3];
      int kc = c1_idx[pc1];
      int k = min(kB, kc);
      if (kB == k && kc == k) {
        A[pA2] += B[pB3] * c[pc1];
      }
      if (kB == k) pB3++;
      if (kc == k) pc1++;
    }
  }
}

```

Fig. 3. $A_{ij} = \sum_k B_{ijk} c_k$ (sparse B, c)

Sparse Storage Formats



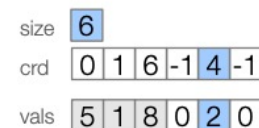
(a) An 8-vector



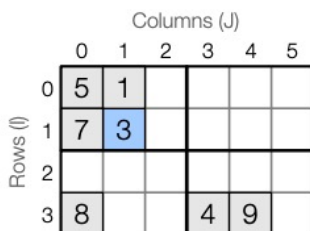
(b) Dense array



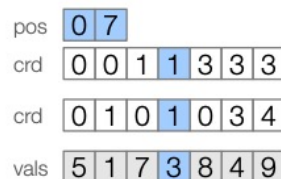
(c) Sparse vector



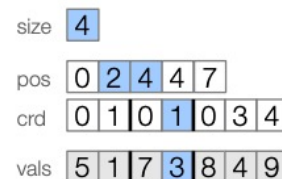
(d) Hash map



(e) A 4x6 matrix



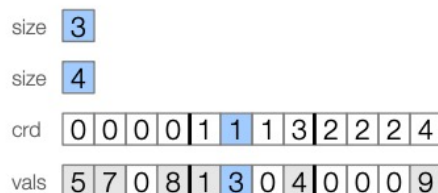
(f) COO



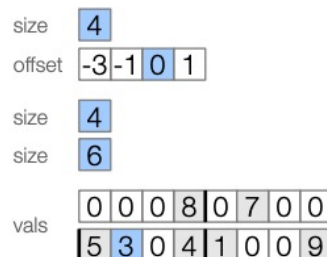
(g) CSR



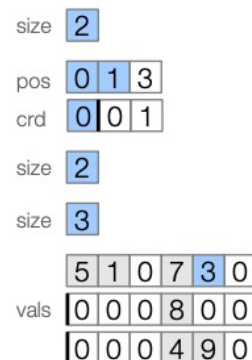
(h) DCSR



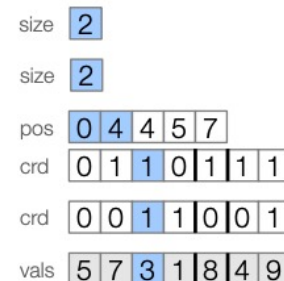
(i) ELL



(j) DIA



(k) BCSR



(l) CSB

Our solution: Separation of Concerns

- Logical SDQL
 - Tensor algebra + more
 - AD friendly
- Physical SDQL
 - Sparse storage formats
 - Efficient



Logical SDQL

$$\sum_i A_i B_i$$

`sum (<i , a> in A) a * B(i)`

AD in Logical SDQL

$$\frac{\partial \sum_i A_i B_i}{\partial B}$$

gradient

`(sum(<i, a> in A) a * B(i)) B`

AD in Logical SDQL (cont.)

$\mathcal{D}_\tau[\mathbb{T}]$ Tensorized FAD on Types

$$\begin{aligned} \mathcal{D}_\tau[\mathbb{D}] &= \mathbb{D} \otimes \tau \\ \mathcal{D}_\tau[\mathbf{bool}] &= \mathbf{real} \\ \mathcal{D}_\tau[\mathbf{int}] &= \mathbf{real} \end{aligned}$$

$\mathcal{D}_\tau[\Gamma]$ Tensorized FAD on Context

$$\begin{aligned} \mathcal{D}_\tau[\emptyset] &= \emptyset \\ \mathcal{D}_\tau[\Gamma, x:\mathbb{T}] &= \mathcal{D}_\tau[\Gamma], x:\mathbb{T}, x':\mathcal{D}_\tau[\mathbb{T}] \end{aligned}$$

$\mathcal{D}_\tau[e]$ Tensorized FAD on Expressions

– Invariant: If $\Gamma \vdash e : \mathbb{T}$, then $\mathcal{D}_\tau[\Gamma] \vdash \mathcal{D}_\tau[e] : \mathcal{D}_\tau[\mathbb{T}]$

$$\begin{aligned} \mathcal{D}_\tau[\mathbf{sum}(\langle k, v \rangle \mathbf{in} e1) e2] &= \mathbf{sum}(\langle k, v \rangle \mathbf{in} e1) \mathbf{let} \langle k', v' \rangle = \langle 0, \mathcal{D}_\tau[e1(k)] \rangle \mathbf{in} \mathcal{D}_\tau[e2] \\ \mathcal{D}_\tau[\mathbf{let} x = e1 \mathbf{in} e2] &= \mathbf{let} \langle x, x' \rangle = \langle e1, \mathcal{D}_\tau[e1] \rangle \mathbf{in} \mathcal{D}_\tau[e2] \\ \mathcal{D}_\tau[\mathbf{if} e1 \mathbf{then} e2] &= \mathbf{if} e1 \mathbf{then} \mathcal{D}_\tau[e2] & \mathcal{D}_\tau[\{ e1 \rightarrow e2 \}] &= \{ e1 \rightarrow \mathcal{D}_\tau[e2] \} \\ \mathcal{D}_\tau[e1 * e2] &= e1 * \mathcal{D}_\tau[e2] + \mathcal{D}_\tau[e1] *^T[\tau] e2 & \mathcal{D}_\tau[e1(e2)] &= \mathcal{D}_\tau[e1](e2) \\ \mathcal{D}_\tau[e1 + e2] &= \mathcal{D}_\tau[e1] + \mathcal{D}_\tau[e2] & \mathcal{D}_\tau[\mathbf{uop}(e)] &= \mathbf{uop}'(e) * \mathcal{D}_\tau[e] \\ \mathcal{D}_\tau[x] &= x' & \mathcal{D}_\tau[r] &= \mathbf{zero}[\tau] & \mathcal{D}_\tau[n] &= \mathcal{D}_\tau[\mathbf{false}] = \mathcal{D}_\tau[\mathbf{true}] = 0 \end{aligned}$$

$$\begin{aligned} e1 *^T[\mathbf{real}] e2 &\triangleq e1 * e2 \\ e1 *^T[\mathbf{tensor} n] e2 &\triangleq \mathbf{sum}(\langle i_1, r_2 \rangle \mathbf{in} e1) \dots \mathbf{sum}(\langle i_m, v \rangle \mathbf{in} r_m) \\ \mathbf{if} e1 : \mathbf{tensor} (m + n) &\quad \{ i_1 \rightarrow \dots \{ i_m \rightarrow 1 \} \dots \} * e2 * v \end{aligned}$$

AD in Logical SDQL (cont.)

```
gradient (sum(<i, a> in A) a * B(i)) B
```



Tensorized AD

```
let A' = {} in
let B' = sum(<i, _> in B) {i -> {i -> 1}} in
sum(<i, a> in A)
  let <i', a'> = <{}, A'(i)> in
  a * B'(i) + a' * B(i)
```



Optimizations

```
sum(<i, a> in A) { i -> a }
```

Optimizations in Logical SDQL

```

let A' = {} in
let B' = sum(<i, _> in B) {i -> {i -> 1}} in
sum(<i, a> in A)
  let <i', a'> = <{}, A'(i)> in
  a * B'(i) + a' * B(i)

```



```

let B' = sum(<i, _> in B) {i -> {i -> 1}} in
sum(<i, a> in A)
  let <i', a'> = <{}, {}> in
  a * B'(i) + a' * B(i)

```

Optimizations in Logical SDQL (cont.)

```
let B' = sum(<i, _> in B) {i -> {i -> 1}} in
sum(<i, a> in A)
```

```
let <i', a'> = <{}, {}> in
a * B'(i) + a' * B(i)
```



```
let B' = sum(<i, _> in B) {i -> {i -> 1}} in
sum(<i, a> in A)
a * B'(i) + {} * B(i)
```


Optimizations in Logical SDQL (cont.)

```
let B' = sum(<i, _> in B) {i -> {i -> 1}} in  
sum(<i, a> in A)  
  a * B'(i)
```



```
sum(<i, a> in A)  
  a * {i -> 1}
```

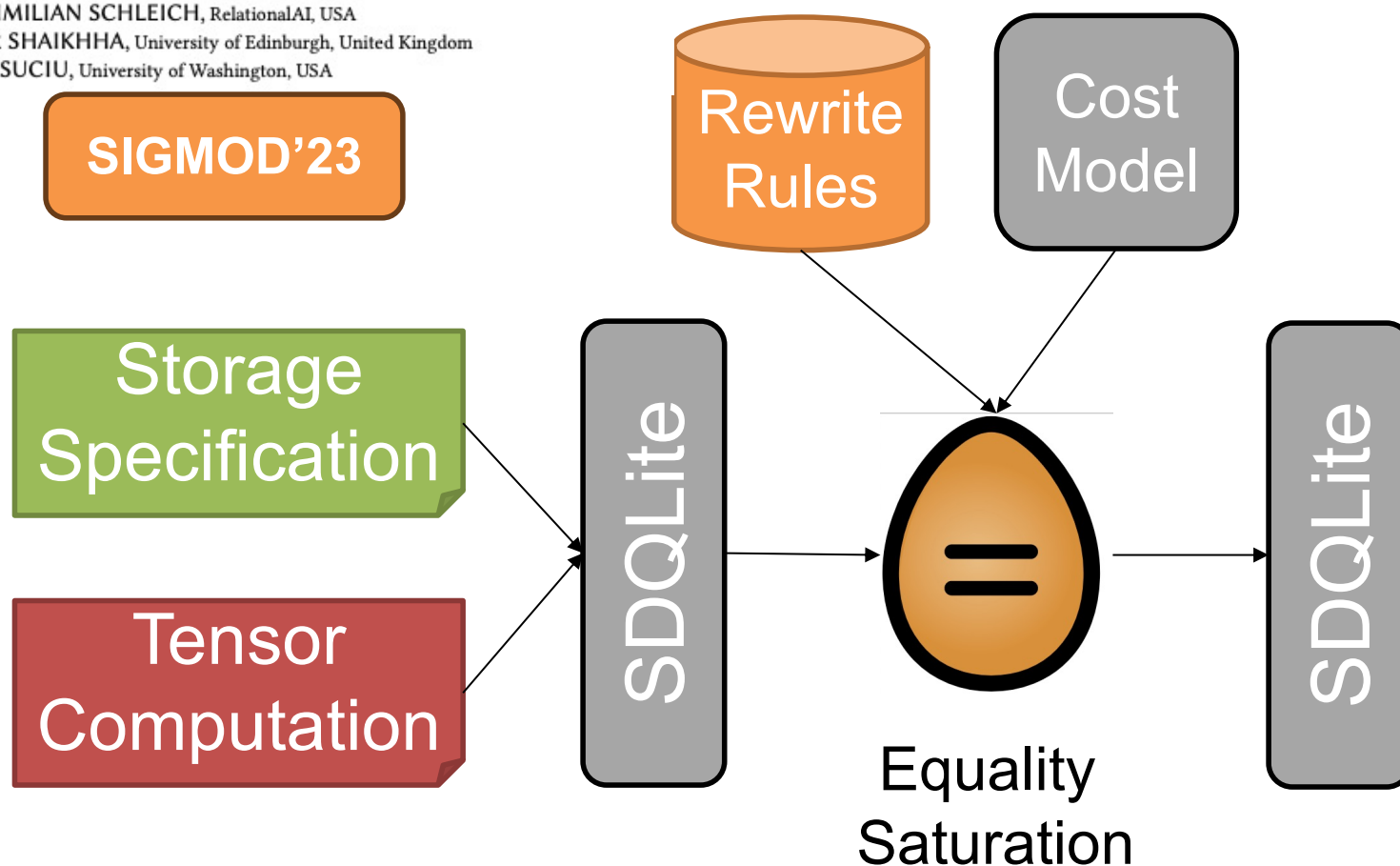
Physical SDQL == SDQLite

Optimizing Tensor Programs on Flexible Storage

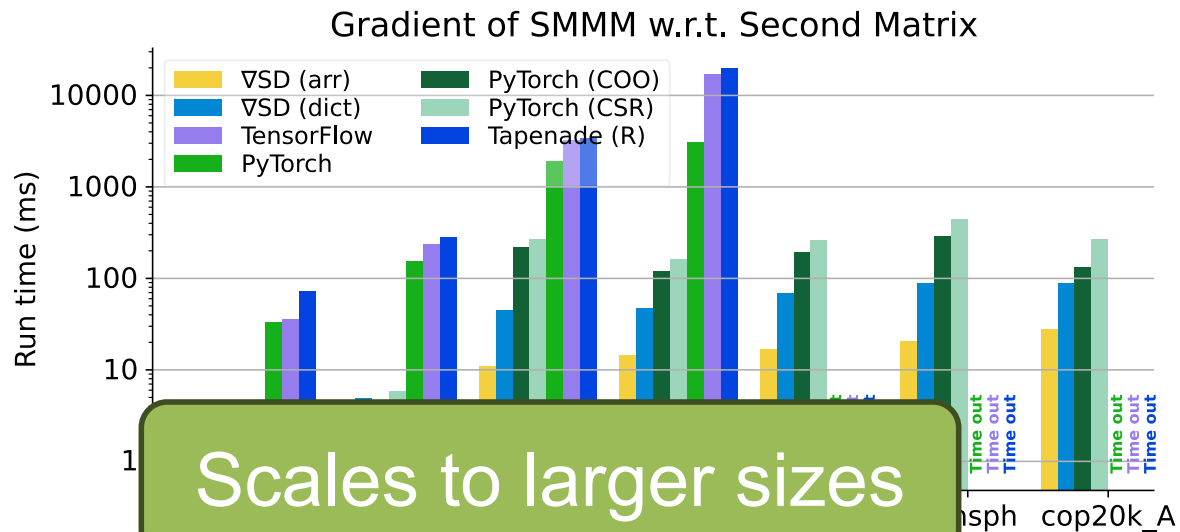
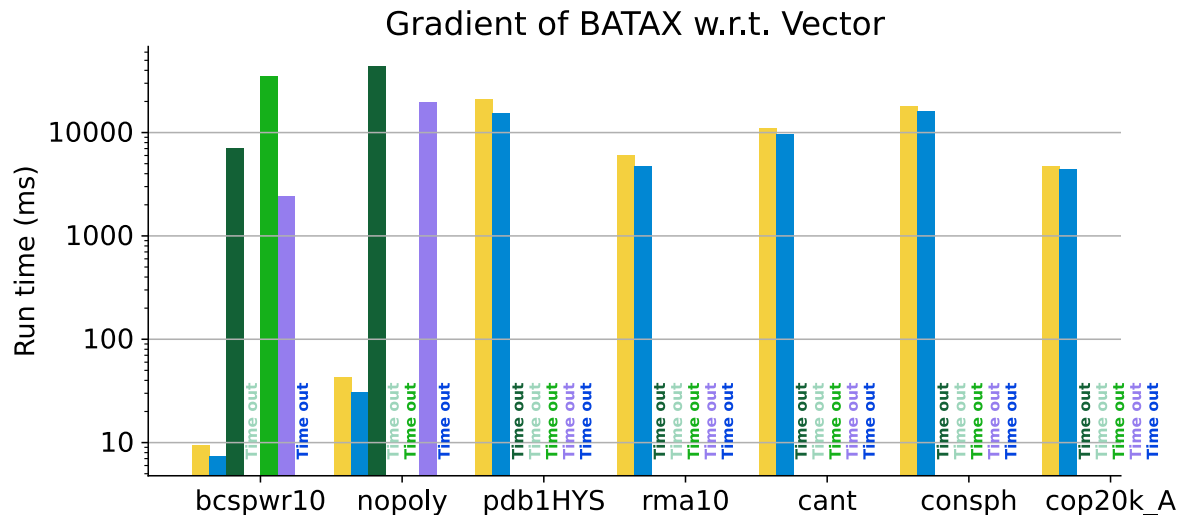
MAXIMILIAN SCHLEICH, RelationalAI, USA

AMIR SHAIKHHA, University of Edinburgh, United Kingdom

DAN SUCIU, University of Washington, USA

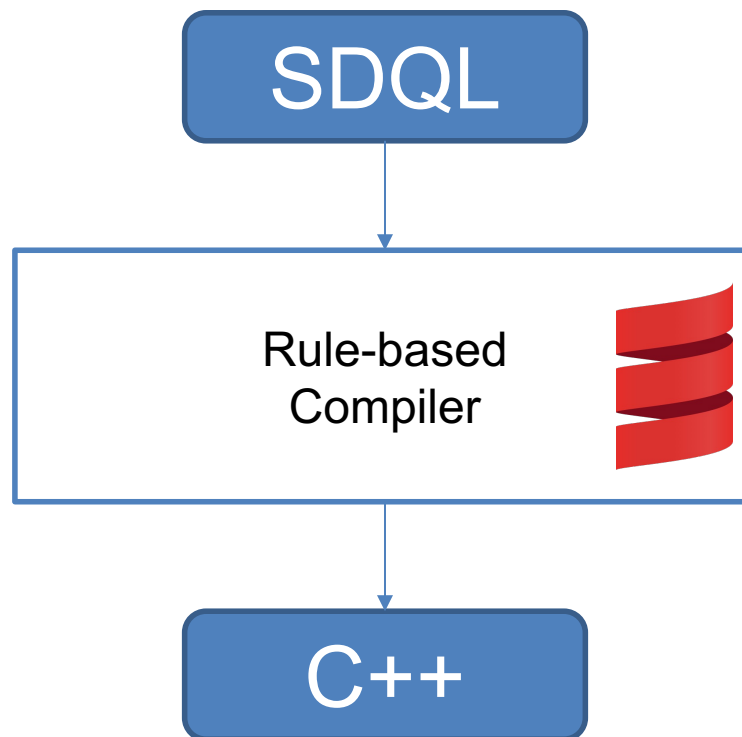


Performance Results

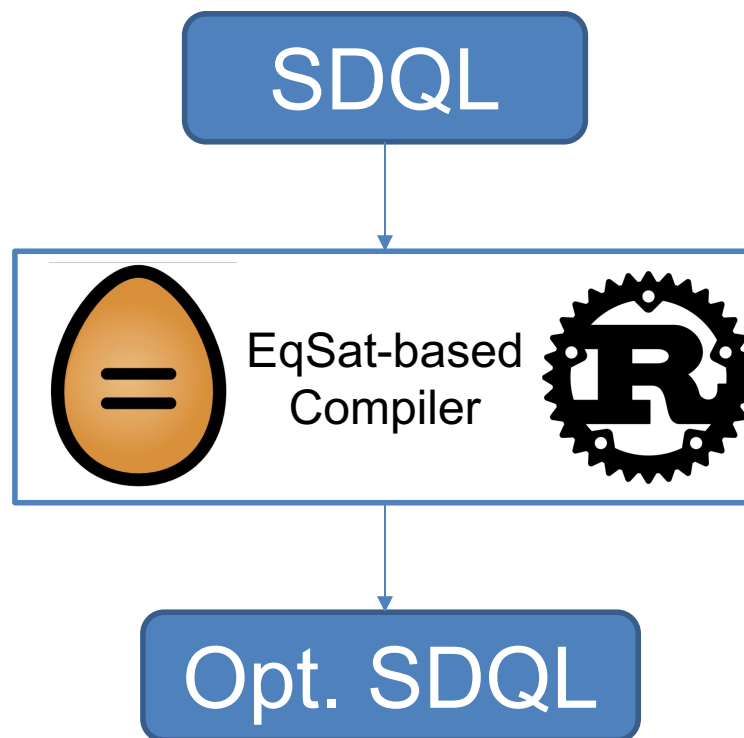


IMPLEMENTATION

Compiler v1 [OOPSLA'22]



Towards Eq-Sat-Based Compiler



Challenges with Equality Saturation

- Scalability
 - Variable binding
- Scalability
 - Associativity/Commutativity
- Scalability
 - One directional rules
- Scalability
 - Search space is BIG!
- Analysis
 - Cardinality estimation
 - Cost estimation

Variable binding

- De bruijn indexing

Sketch-Guided Equality Saturation

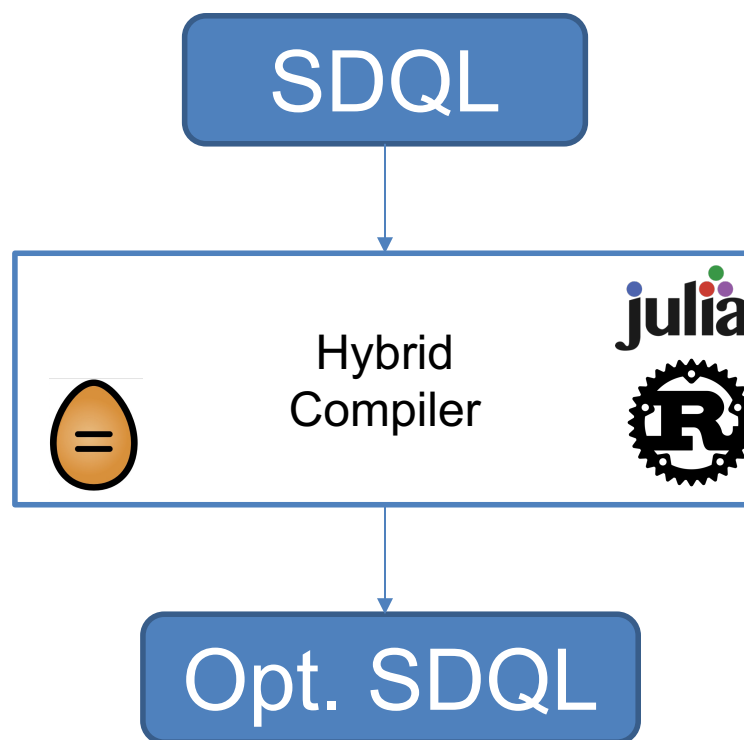
Scaling Equality Saturation to Complex Optimizations in Languages with Bindings

Thomas Köhler
University of Glasgow
Scotland, UK
thomas.koehler@thok.eu

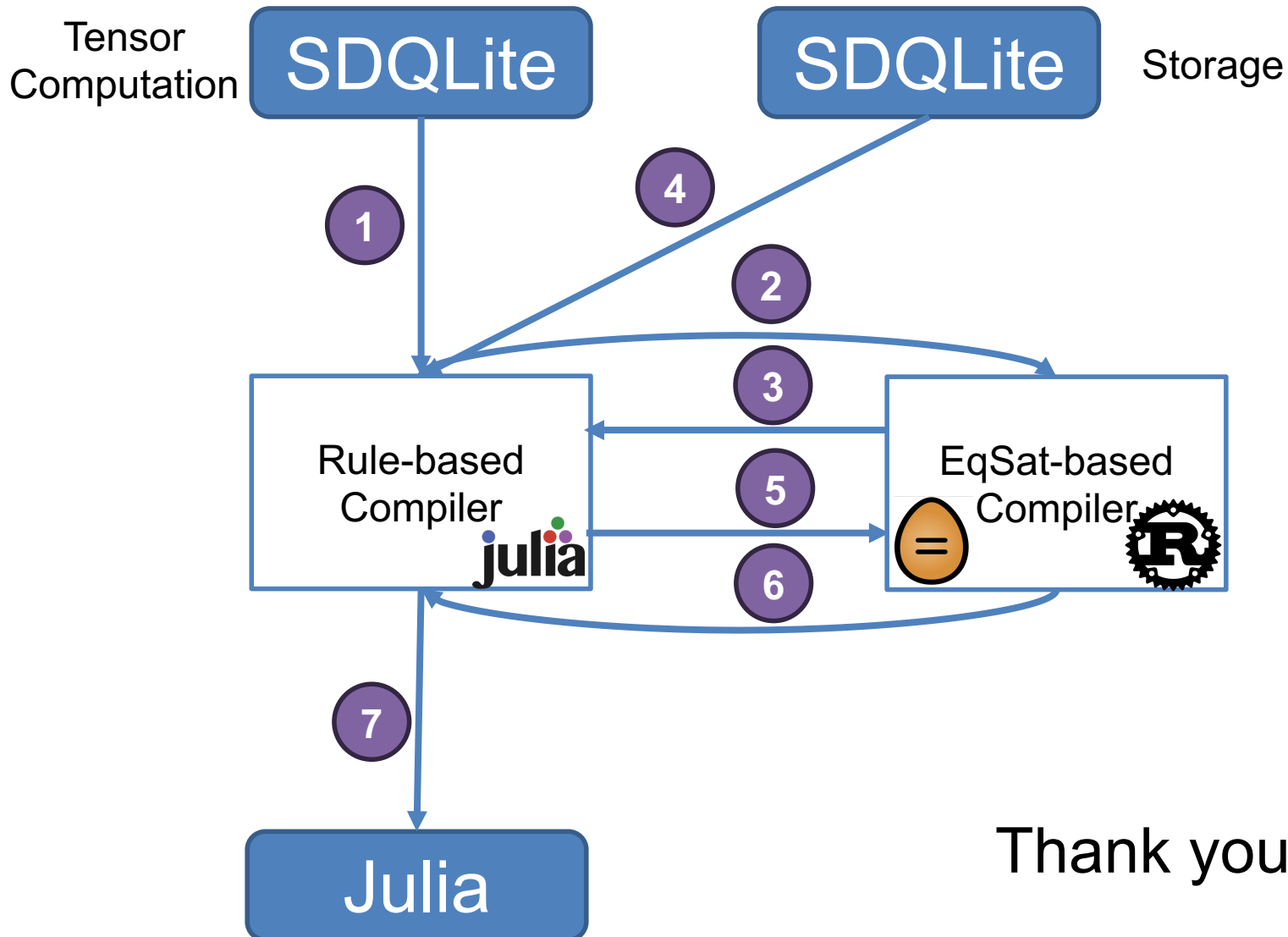
Phil Trinder
University of Glasgow
Scotland, UK
Phil.Trinder@glasgow.ac.uk

Michel Steuwer
University of Edinburgh
Scotland, UK
michel.steuwer@ed.ac.uk

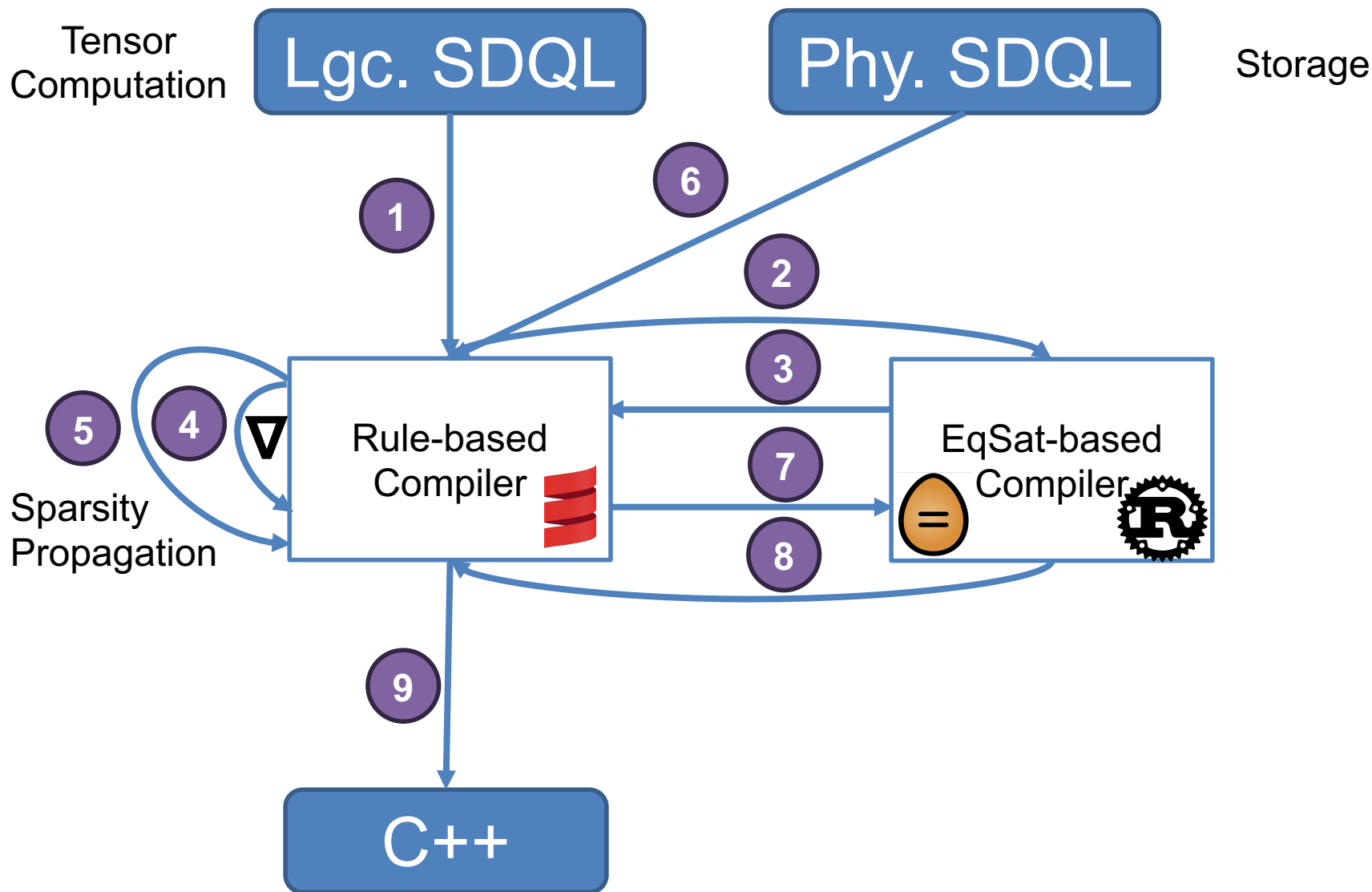
Hybrid Compiler



Compiler v2 [SIGMOD'23]



Compiler v3 [CGO'24]



Challenges with EqSat

- Scalability
 - DB/PL/ML ideas
 - For Reverse-mode AD, limited use of EqSat
- Analysis

Better Together: Unifying Datalog and Equality Saturation

YIHONG ZHANG, University of Washington, USA

YISU REMY WANG, University of Washington, USA

OLIVER FLATT, University of Washington, USA

DAVID CAO, UC San Diego, USA

PHILIP ZUCKER, Draper Laboratory, USA

ELI ROSENTHAL, Google, USA

ZACHARY TATLOCK, University of Washington, USA

MAX WILLSEY, University of Washington, USA

Latent Idiom Recognition for a Minimalist Functional Array Language using Equality Saturation

Jonathan Van der Cruysse

McGill University

Montreal, Quebec, Canada

jonathan.vandercruysse@mail.mcgill.ca

Christophe Dubach

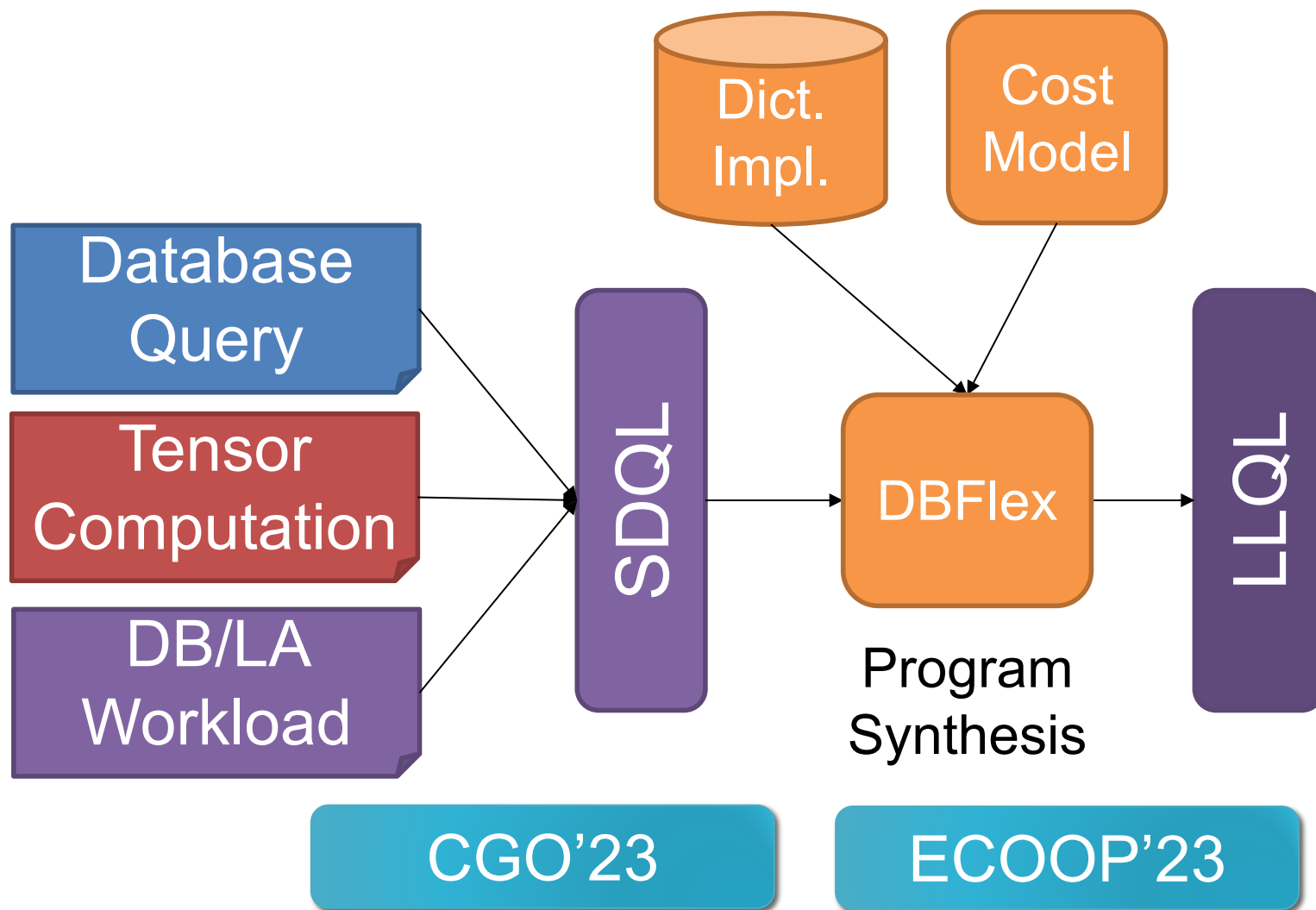
McGill University & Mila

Montreal, Quebec, Canada

christophe.dubach@mcgill.ca

CURRENT WORK

ML for Dictionary Tuning



Optimizing Nested Recursive Queries

AMIR SHAIKHHA, University of Edinburgh, United Kingdom

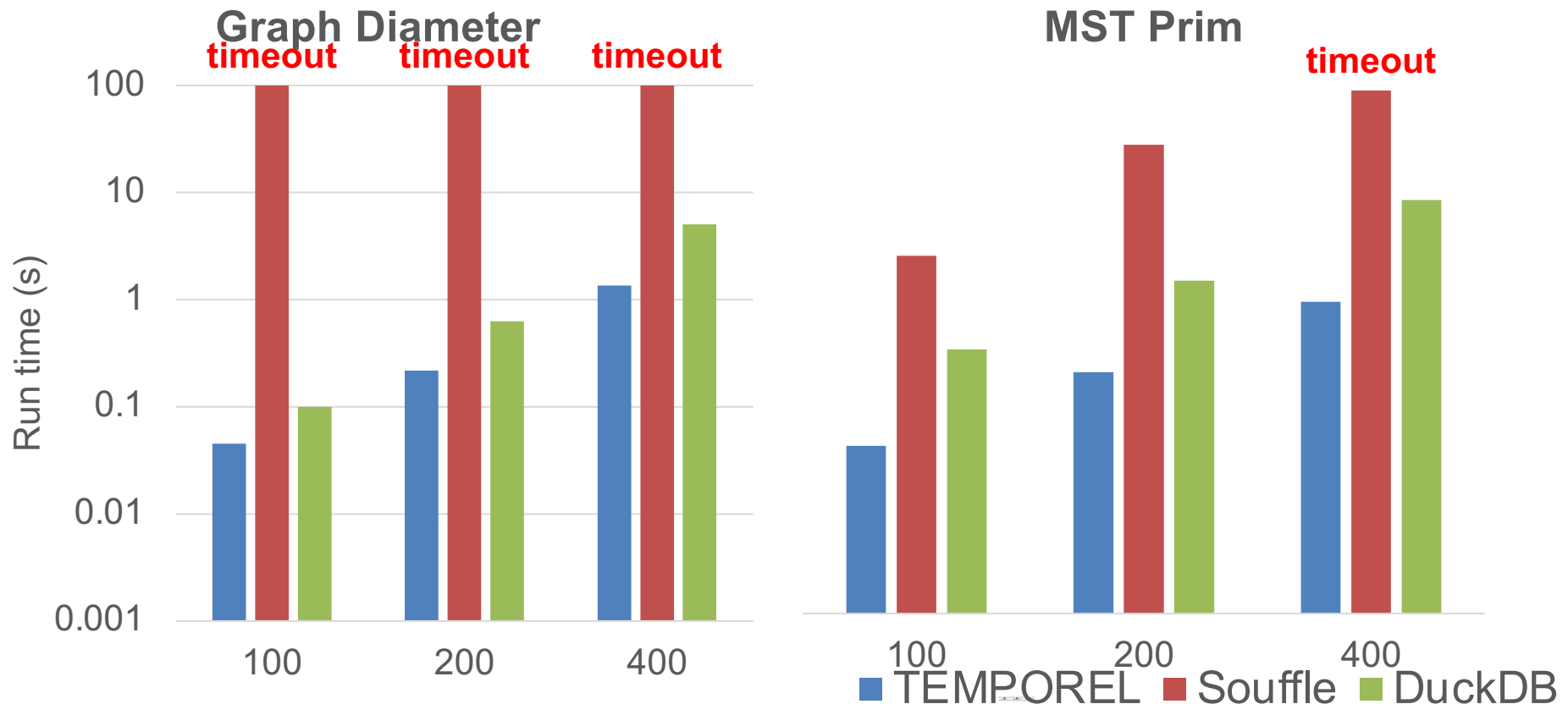
DAN SUCIU, University of Washington, USA

MAXIMILIAN SCHLEICH, RelationalAI, USA

HUNG NGO, RelationalAI, USA

Recursive Queries

- Datalog \rightarrow TempoDL \rightarrow SDQL + Recursion



SIGMOD'24

Advanced Joins in SDQL

Worst-case Optimal Join Algorithms

Hung Q. Ngo
University at Buffalo, SUNY
hungngo@buffalo.edu

Ely Porat
Bar-Ilan University
porately@cs.biu.ac.il

Christopher Ré
University of
Wisconsin–Madison
chrisre@cs.wisc.edu

Atri Rudra
University at Buffalo, SUNY
atri@buffalo.edu

Leapfrog Triejoin: A Simple, Worst-Case Optimal Join Algorithm

Todd L. Veldhuizen
LogicBlox Inc.
Two Midtown Plaza
1349 West Peachtree Street NW
Suite 1880, Atlanta GA 30309
tveldhui@{logicblox.com,acm.org}

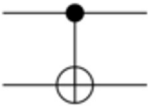
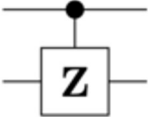
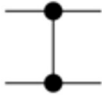

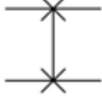
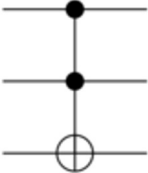
Free Join: Unifying Worst-Case Optimal and Traditional Joins

YISU REMY WANG, University of Washington, USA

MAX WILLSEY, University of Washington, USA

DAN SUCIU, University of Washington, USA

Quantum Simulation in SDQL

Operator	Gate(s)	Matrix
Controlled Not (CNOT, CX)		$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$
Controlled Z (CZ)	 	$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & -1 \end{bmatrix}$
SWAP	 	$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$
Toffoli (CCNOT, CCX, TOFF)		$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$

Conclusion

- Data science is a double-edged sword!
- Radical rethinking of data science pipelines
- Language design
 - SDQL (Semi-ring Dictionary)
 - TempoDL (Low-Level Datalog) [SIGMOD'24]
 - TondIR (Python to SQL) [ICDE'24]
 - STUR (Structured Tensor Algebra) [OOPSLA'23]
 - BTL (Probabilistic Language) [CC'23]
- Leverage structure
- Optimize across pipeline

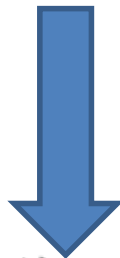
Thank you

BACKUP

Vertical Loop Fusion

<pre>let y=sum(x in e1) {x.key->f1(x.val)} sum(x in y){x.key->f2(x.val)}</pre>	\rightsquigarrow	<pre>sum(x in e1) { x.key -> f2(f1(x.val)) }</pre>
--	--------------------	---

```
let At = sum(row in A) sum(x in row.val) { x.key -> {row.key -> x.val } }
sum(row in At) { row.key ->
  sum(x in row.val) sum(y in A(x.key))
  { y.key -> x.val * y.val } }
```

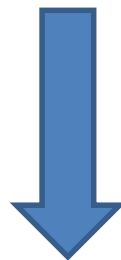


```
sum(row in A)
  sum(x in row.val) { x.key ->
    sum(y in row.val) { y.key ->
      x.val * y.val } }
```

Horizontal Loop Fusion

<pre>let y1=sum(x in e1)f1(x) let y2=sum(x in e1)f2(x) f3(y1, y2)</pre>	\rightsquigarrow	<pre>let tmp = sum(x in e1) <y1 = f1(x), y2 = f2(x)> f3(tmp.y1, tmp.y2)</pre>
---	--------------------	---

```
let Rsum = sum(r in R) r.key.A * r.val in
let Rcount = sum(r in R) r.val in
Rsum / Rcount
```



```
let RsumRcount = sum(r in R) < Rsum = r.key.A * r.val, Rcount = r.val > in
RsumRcount.Rsum / RsumRcount.Rcount
```

Loop Factorization

- Scalars

`sum(x in NR) sum(y in x.key.C) x.key.A * x.val * y.key.D * y.val`



`sum(x in NR) x.key.A * x.val * sum(y in x.key.C) y.key.D * y.val`

- Dictionaries

`sum(x in NR) sum(y in x.key.C) { x.key.B -> x.key.A * x.val * y.key.D * y.val }`



`sum(x in NR) { x.key.B -> x.key.A * x.val * sum(y <- x.key.C) y.key.D * y.val }`

Other Approaches

	Expressiveness					Data Representation					Specialization				
	Relational Algebra	Nested Rel. Calc.	Group-by Aggregates	Efficient Equi-Joins	Linear Algebra	Set & Bag	Dense Array	Sparse Tensor	Dictionary	Semi-rings	Loop Fusion	Loop Hoisting	Loop Memoization	Code Generation	Vectorization
SDQL (This Paper)	●	●	●	●	●	●	●	●	●	●	●	●	●	●	○
Query Compilers (HyPer)	●	○	●	●	○	●	●	○	●	○	●	●	○	●	○
Vectorized Query Engines (Vectorwise)	●	○	●	●	○	●	●	○	●	○	●	●	○	○	●
Monad Calculus (NRC ⁺)	●	●	○	○	○	●	○	○	○	○	●	●	○	○	○
Monoid Comprehension	●	●	●	○	○	●	○	○	○	○	●	●	○	○	○
Monad Calc. + Agg. (Kleisli, Trance)	●	●	●	○	○	●	○	○	○	○	●	●	○	●	○
Lang. Integrated Queries (LINQ, CompComp)	●	●	●	○	●	●	○	○	○	○	●	●	○	○	○
Functional Lists (Generalized Stream Fusion)	●	●	●	○	●	●	○	○	○	○	●	●	○	●	●
Functional APL (Futhark, SAC)	○	○	○	○	●	○	●	○	○	○	●	●	○	●	●
Dense LA Library (NumPy)	○	○	○	○	●	○	●	○	○	○	○	○	○	○	●
Dense LA DSL (Lift, Halide, LGen)	○	○	○	○	●	○	●	○	○	○	●	●	○	●	●
Sparse LA Library (SPLATT, SciPy)	○	○	○	○	●	○	●	○	○	○	○	○	○	○	○
Sparse LA DSL (TACO)	○	○	○	○	●	○	●	○	○	○	○	○	○	●	○
DB/LA by casting to LA (Morpheus)	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○
DB/LA by casting to DB (LMFAO)	●	○	●	●	○	●	●	○	○	●	○	○	○	●	○
DB/LA by new DSL (IFAQ)	●	○	●	●	●	●	○	●	●	●	○	○	○	●	○