

# Towards an Efficient Linear Algebra Based OLAP Engine

---

João Afonso - 2017/2018



Useful links:

[RPD](#)

[Dissertation](#)

[Tech Report \(PI\)](#)

[RPD presentation](#)

[PODS\\_PAPER](#)

[Incrementality](#)

[Paper Perform Increm VLDB14](#)



## The main challenges

- Plan and develop LA based OLAP engine, featuring:
  - Low latency (query time)
  - High throughput (queries/time)
  - Horizontal scalability
  - Schema free data access
- To ensure some of this topic we rely on:
  - Columnar data storage and processing
  - Well defined LA properties

## Goals for the meeting

12 Oct

- Study implemented systems (Apache, Amazon, Google, ...)
- First considerations on the system architecture
- Study sparse matrix representations:
  - Read the paper suggested by Prof. Rui Ralha
- Optimise CSC/COO memory usage
- Start learning C++

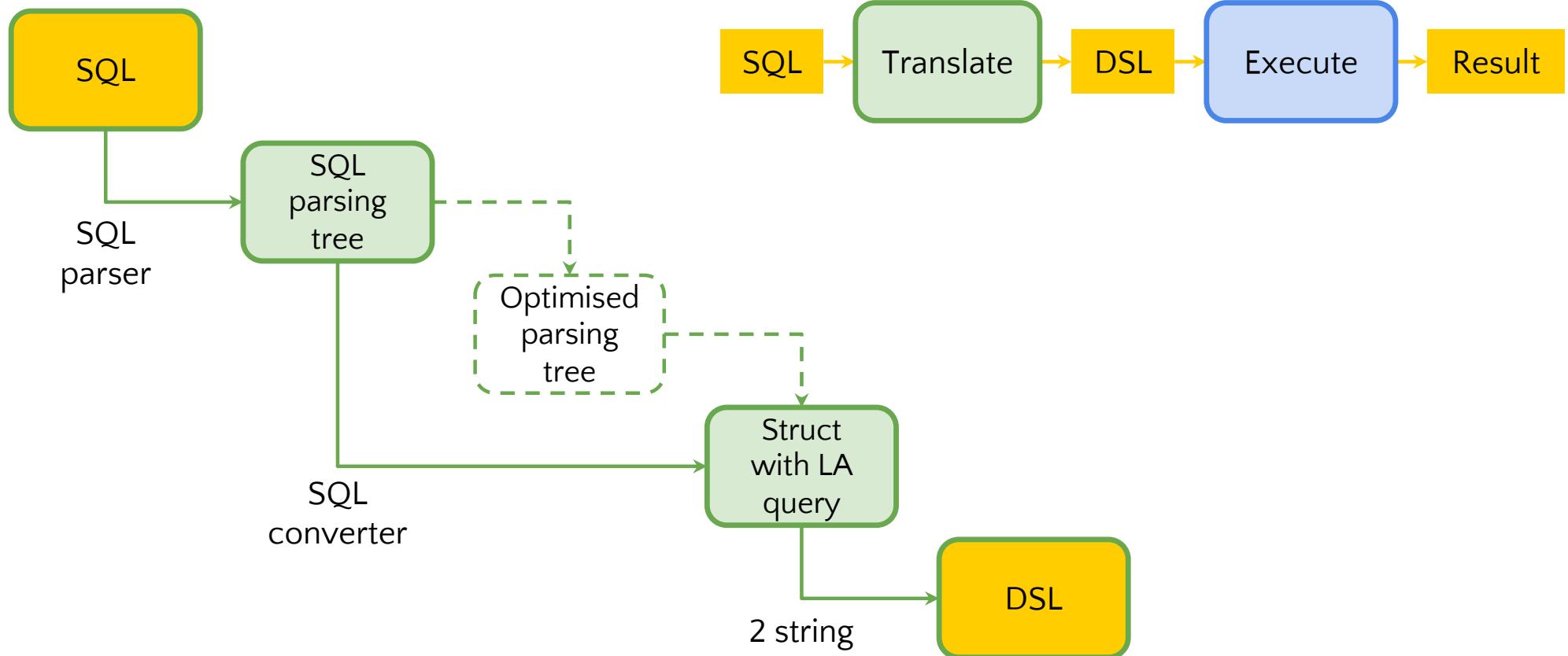


## Implemented systems

- Google provides database solutions considering both **usability**, **efficiency** and **scalability**
- Apache **HBase** is based on Google **BigTable**
- Apache **Drill** is based on Google **BigQuery**
- **EDIT: Impala, Hive + Tez/Pig**
- Both should be considered in the final system design!



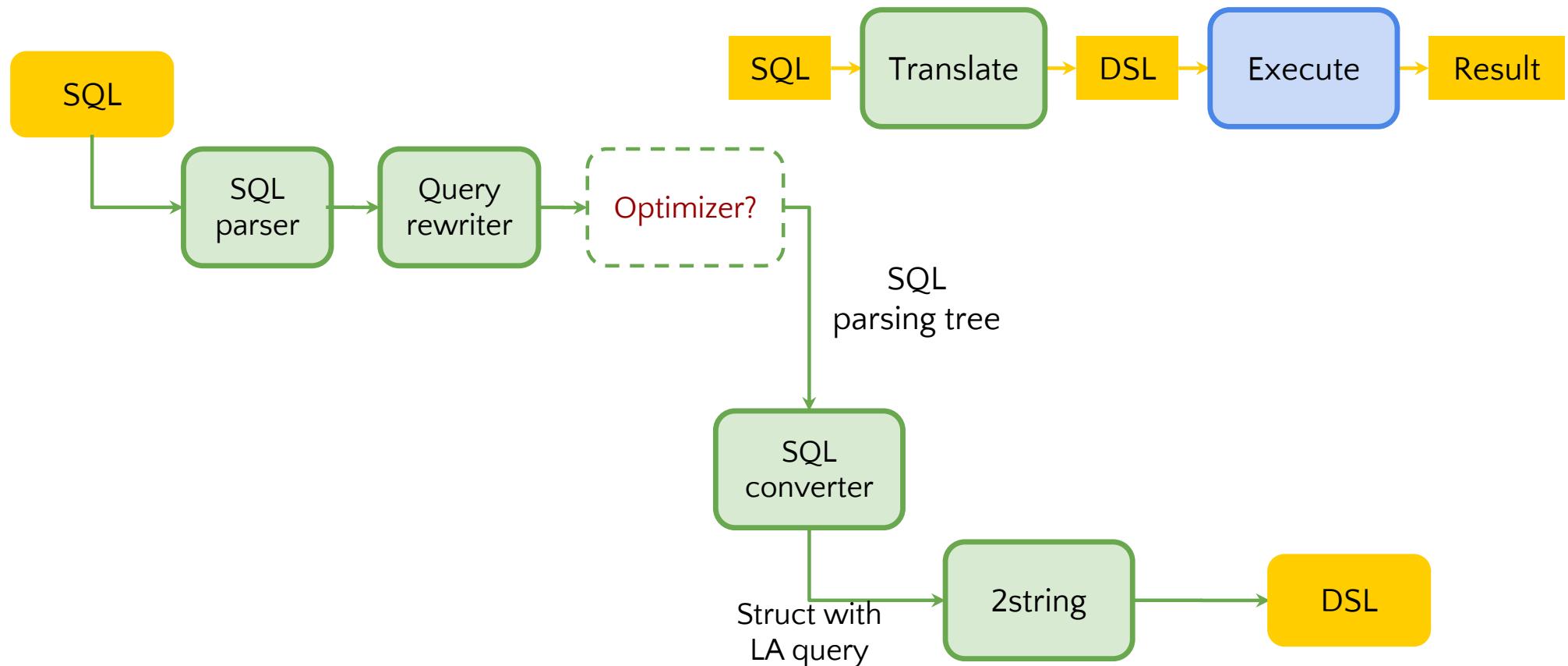
## Possible architecture



Too soon for a functional implementation!

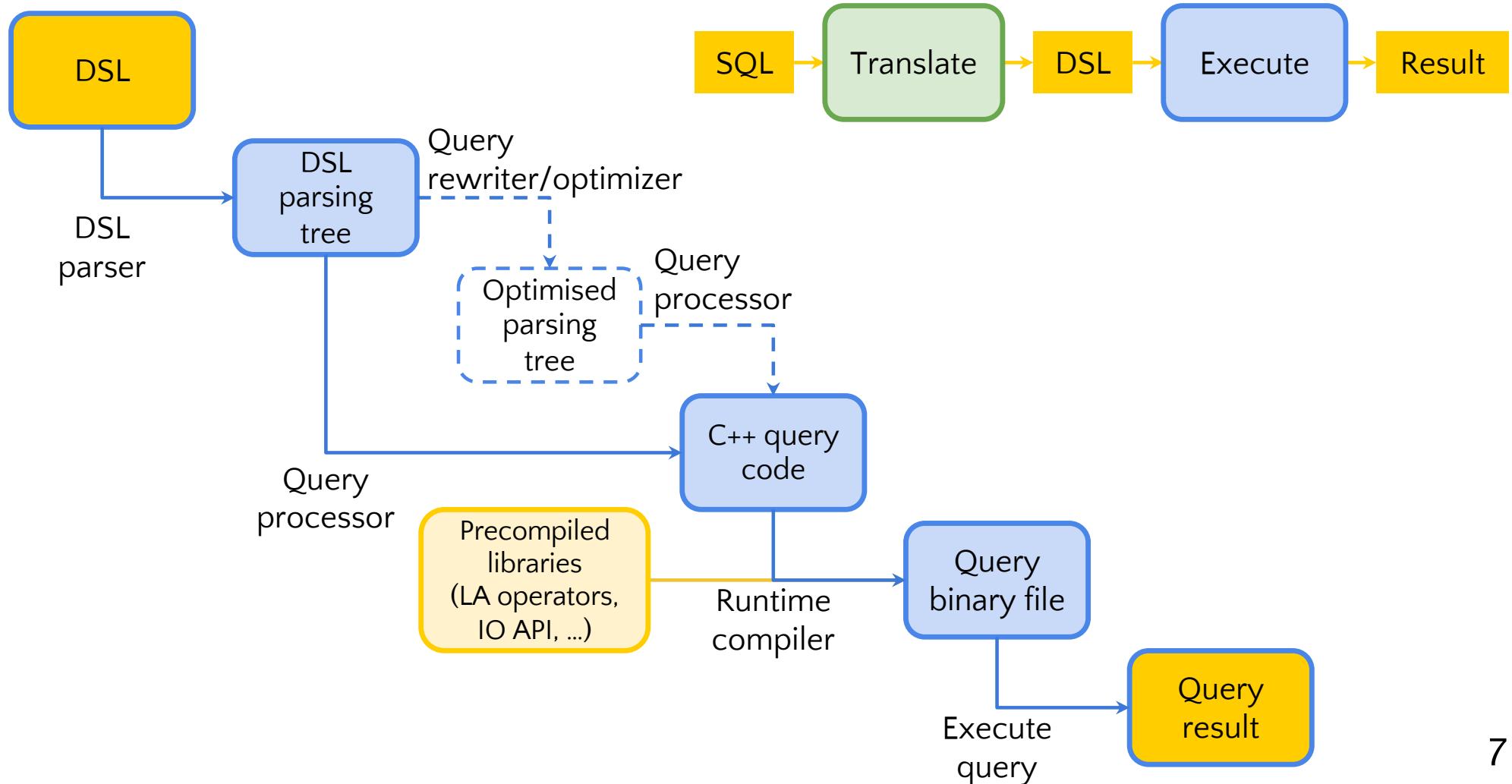


## Possible architecture



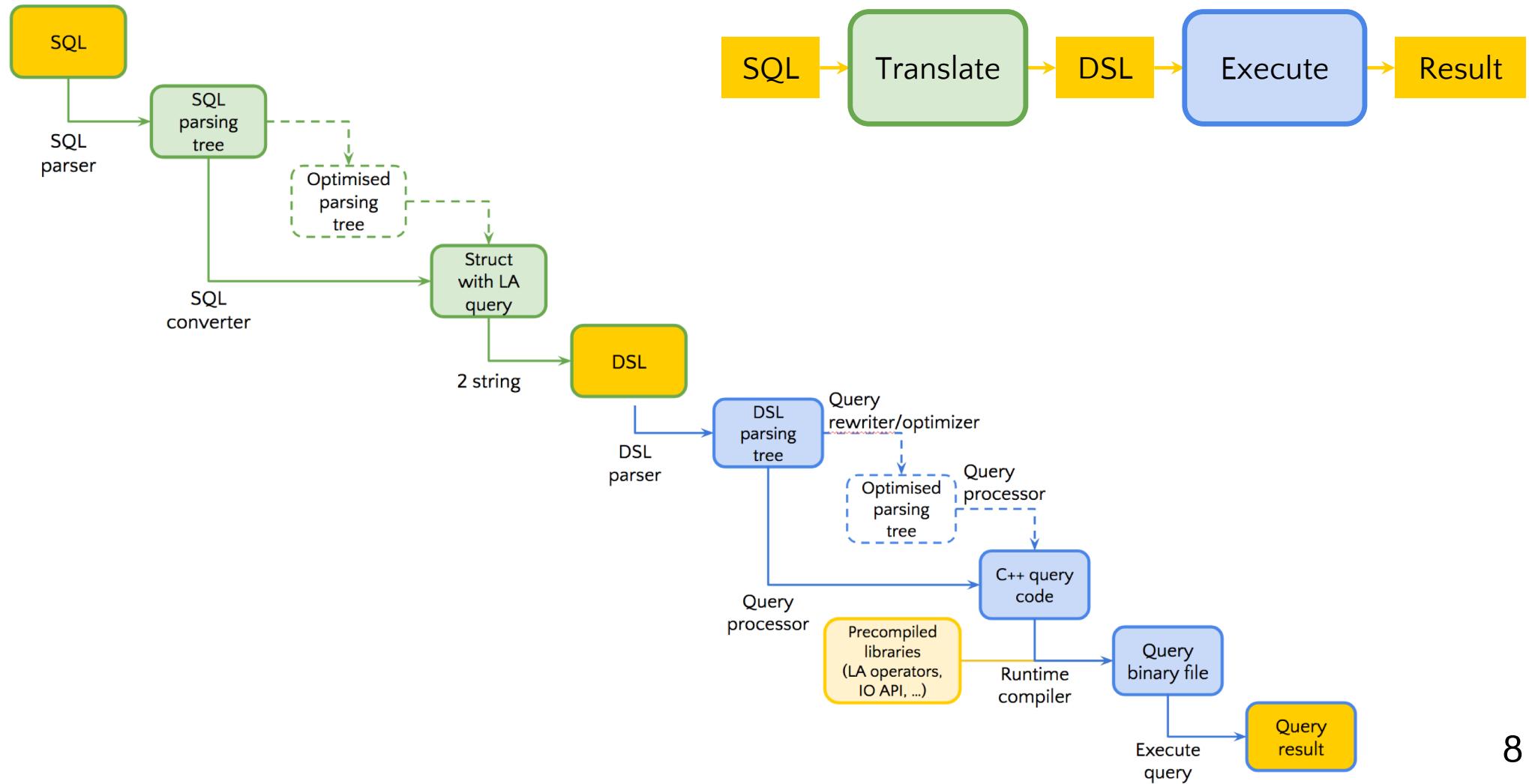


## Possible architecture





# Possible architecture





## Sparse matrix formats

- Previously studied formats: COO, LIL, CSR, CSC
- CSC can compress data when  $\text{nnz} > n$ 
  - We have  $\text{nnz} \leq n$  so we are wasting space
- CSB – Compressed Sparse Blocks:
  - Compressed representation for blocks
  - Discarded by having **dense outer matrix** and **binary search** inside blocks



## CSC memory usage

- CSC format:
  - **Values** [nnz]
  - **Row\_index** [nnz]
  - **Col\_pointer** [n cols]
- Depending on the matrix to be represented, some of them may be redundant:



## CSC/COO memory usage

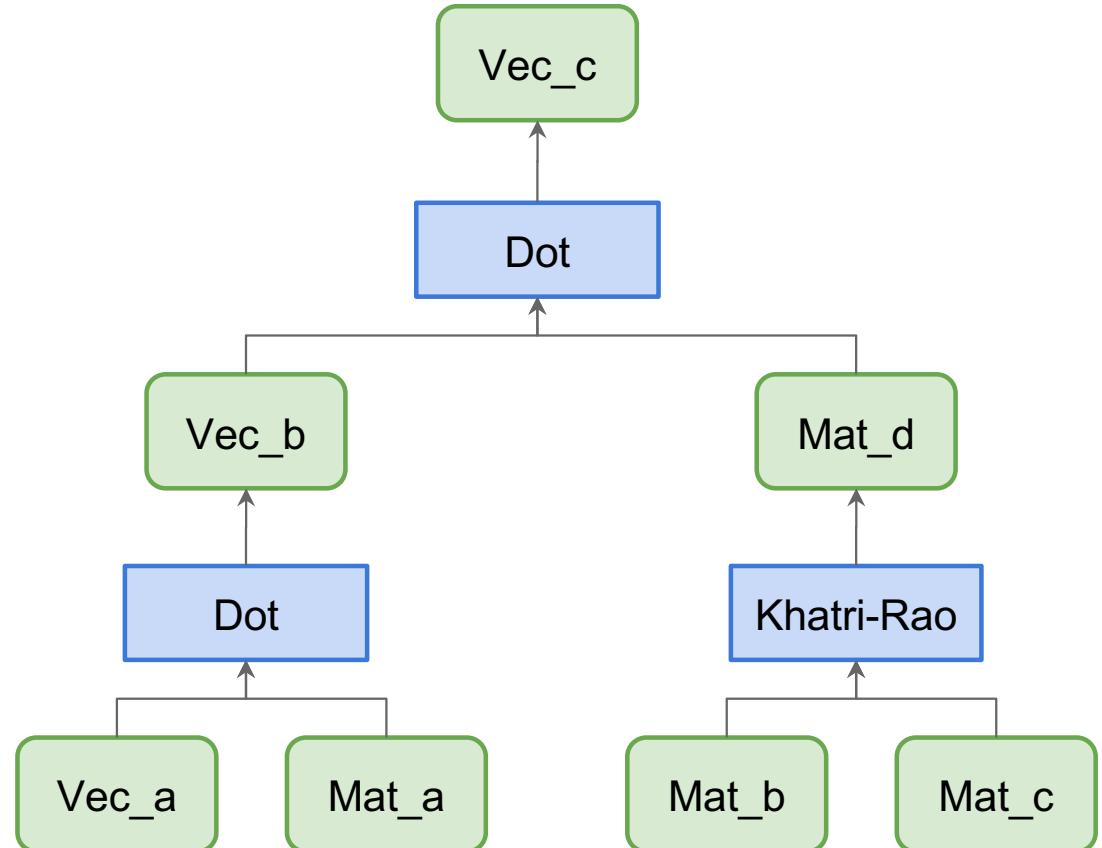
	Values	Rows	Columns
Decimal Vector			
Bitmap			
Filtered Bit Vector			
Decimal Map			
Filtered Bitmap			
Filtered Decimal Vector			
Filtered Decimal Map			

- ◉ Considerations on sparse matrix representations:
  - Case of study
  - Matrix format (data)
  - CSC and COO
  - CSB



## Case of study

- What are we trying to optimise?
  - A parallel flow of data (typically petabytes) through the operation tree (sequential or parallel)
  - Minimize the temporary data storage and waiting times.



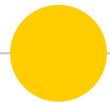


## Matrix format (data)

- Brands: Ford, BMW, Ford, VW, Audi, VW, Ford, Audi;
- Brand  $\neq$  VW (matrix representation):

Ford	1		1				1
BMW		1					
VW							
Audi				1			1

- At most one element per column (**partial** function)*
- The distribution of the non-zero rows is unpredictable



## CSC and COO

Ford	1		1				1
BMW		1					
VW							
Audi					1		1

**4x8**

### ○ CSC representation:

**values:** [1, 1, 1, 1, 1, 1]

**Rows:** [0, 1, 0, 3, 0, 3]

**Cols:** [0, 1, 2, 3, 3, 4, 4, 5, 6]

### ○ COO representation:

**values:** [1, 1, 1, 1, 1, 1]

**Rows:** [0, 1, 0, 3, 0, 3]

**Cols:** [0, 1, 2, 4, 6, 7]

**3x6**

In this ex compacting saves 44% ...



## CSC and COO

### ○ CSC:

- Can be iterated through the non-zero values
- Can't be iterated row-wise
- $\Theta(1)$  column access cost

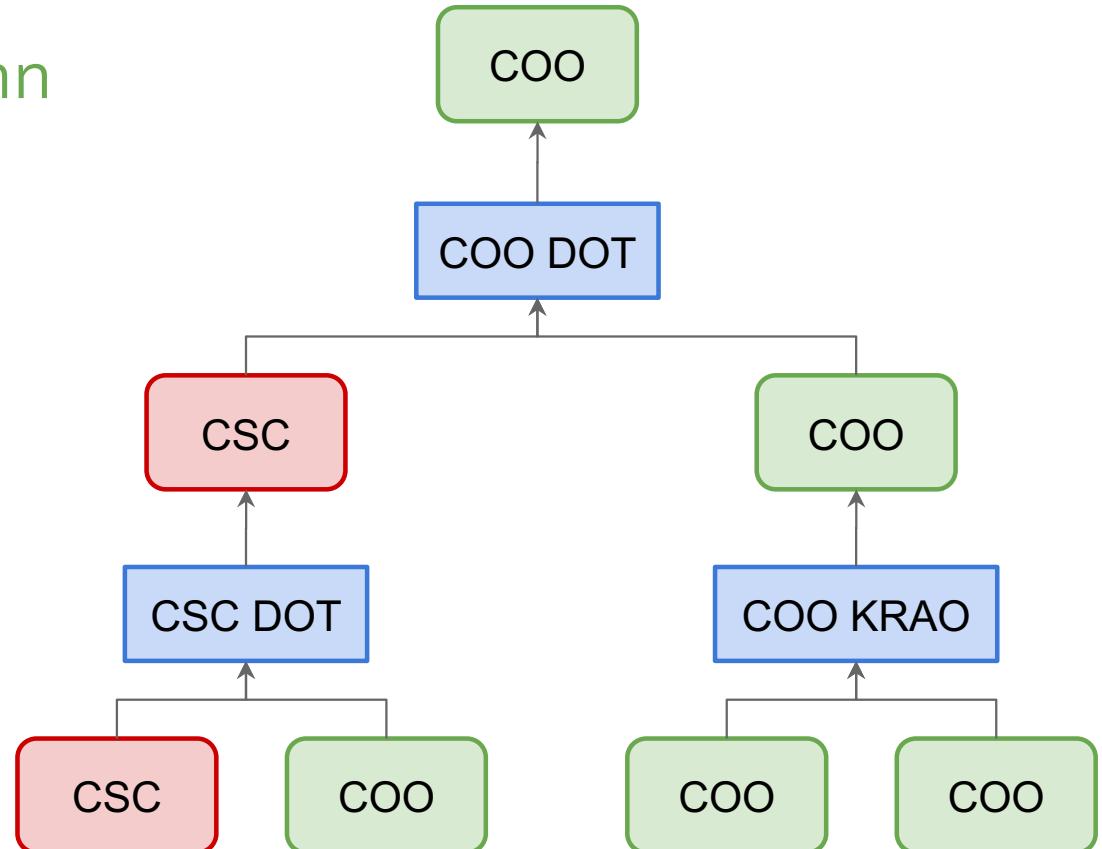
### ○ COO:

- Can be iterated through the non-zero values
- Can't be iterated row-wise
- $\Theta(\log(n))$  column access cost (binary search)



# CSC and COO

- $\Theta(1)$  left matrices column access cost
  - Compressed column pointer array on COO matrices
  - Two different matrix formats
  - Two versions of each operator





**CSB**

---

- CSB:
  - Dense outside matrix pointing to COO blocks
- There is any safe way to define a maximum size (in bytes) for:
  - each individual block?
  - the outer matrix?
- How can the maximum number of columns can be calculated?
- [LINK TO PAPER](#)

## Meeting remarks

- Matrix format
- Data stream operations
  - Dot product dependencies
- Parallel dot product
  - Data ring
- Bibliography



## Matrix format

- ◉ CSB is not suited to our study case
- ◉ CSC/COO as a possible approach
  - Khatri-Rao and Hadamard can use COO both sides
- ◉ Professor Rui Ralha will further investigate an adequate representation for functional sparse matrices



## Stream dot product

- Consider  $y = x \cdot A$
- Divide work by columns “i” of A:

$$\begin{array}{|c|c|} \hline y1 & \dots \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline x1 & x2 & x3 & x4 \\ \hline \end{array} .$$

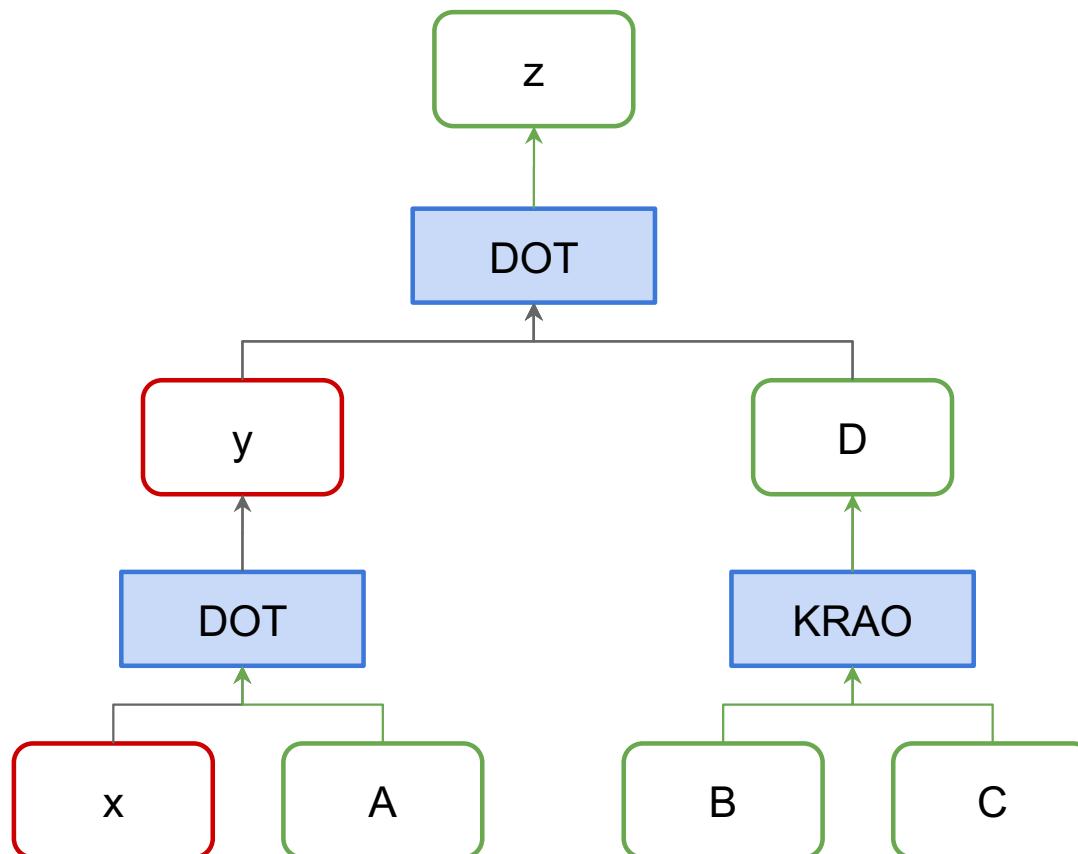
A11	...
A21	...
A31	...
A41	...

- $y[i] = x \cdot A[i]$  ; **x must be in memory!**
- Same for  $C[i] = A \cdot B[i]$  ; A must be in memory!
- EDIT: This is a broadcast JOIN**

NB: Cf. fusion-- law (8) of “A Linear Algebra ... to OLAP” (JNO)



## Stream operation tree



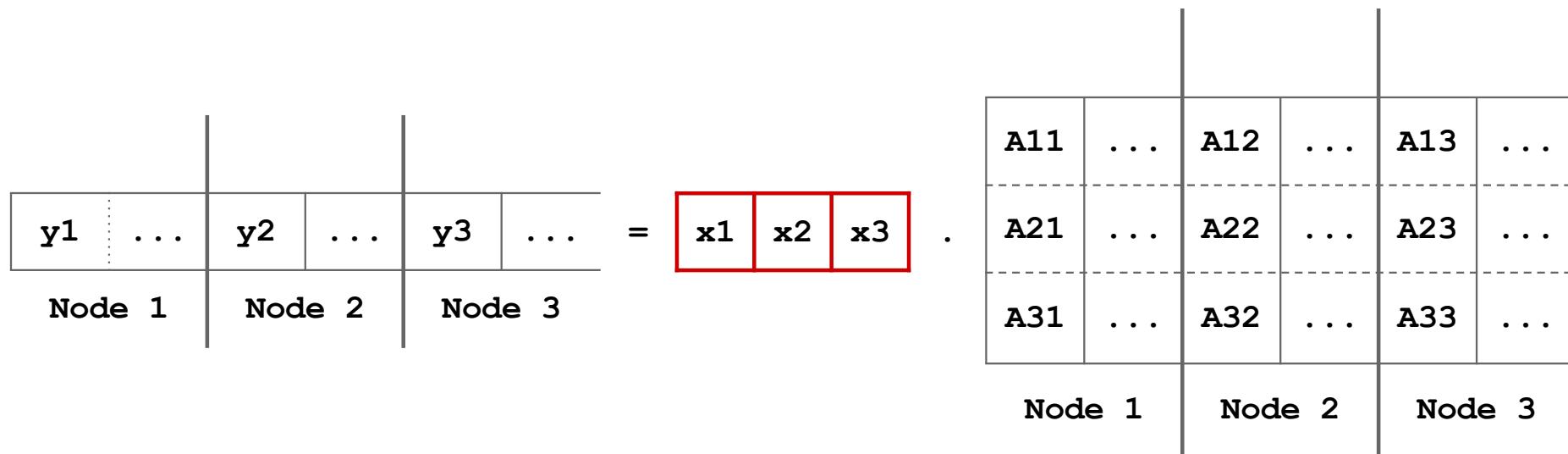


## Parallel stream dot product

**y**: 1  $\leftarrow \infty$

**x**: 1  $\leftarrow 3$

**A**: 3  $\leftarrow \infty$



What if **x** does not fit in memory?



## Parallel stream dot product

$$\begin{array}{c|c|c|c} & & & \\ \hline C1 & \dots & C2 & \dots & C3 & \dots & \dots \\ \hline \end{array} = \begin{array}{c|c|c} A1 & A2 & A3 \\ \hline \end{array} \cdot \begin{array}{c|c|c} & & \\ \hline \end{array}$$

Node 1      Node 2      Node 3

$$\begin{array}{c|c|c|c|c|c} & & & & & \\ \hline B11 & \dots & B12 & \dots & B13 & \dots \\ \hline \dots & & \dots & & \dots & \\ \hline B21 & \dots & B22 & \dots & B23 & \dots \\ \hline \dots & & \dots & & \dots & \\ \hline B31 & \dots & B32 & \dots & B33 & \dots \\ \hline \dots & & \dots & & \dots & \\ \hline \end{array}$$

Node 1      Node 2      Node 3

$$y_1 = x_1 \cdot A_{11} + x_2 \cdot A_{21} + x_3 \cdot A_{31}$$

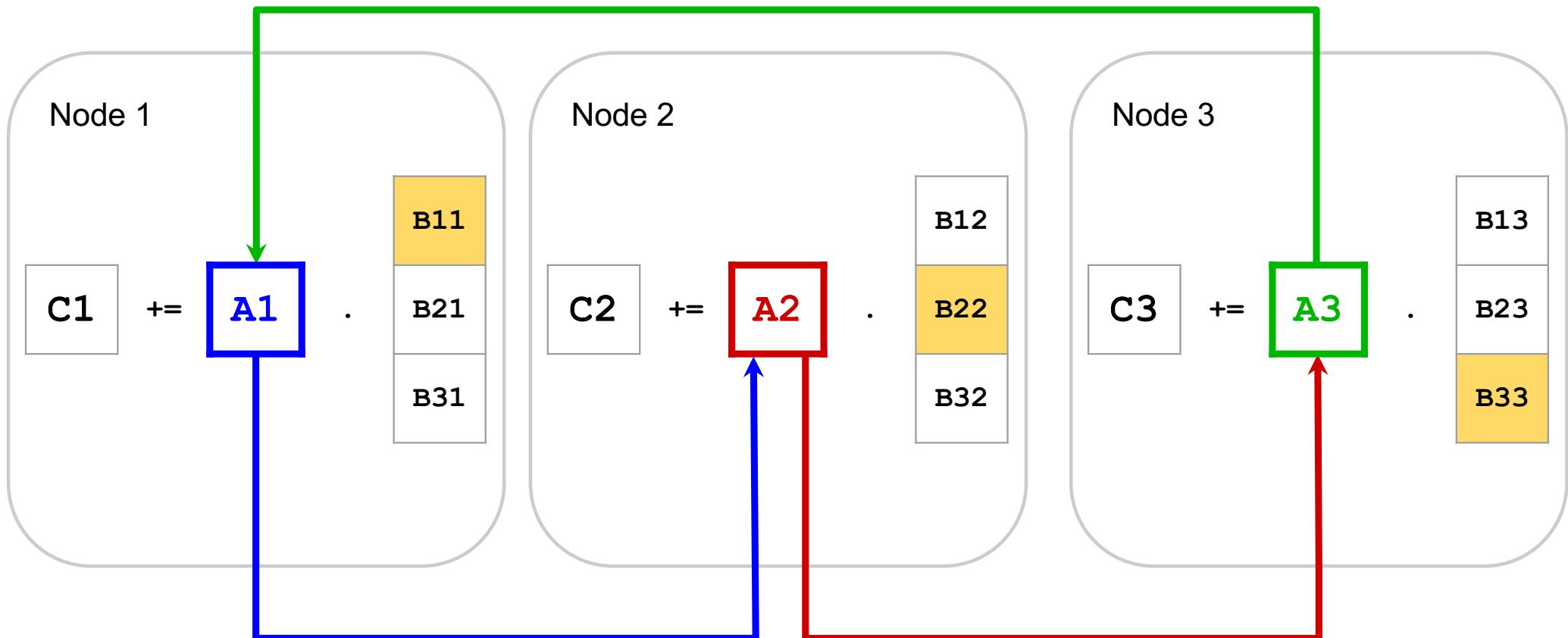
$$y_2 = x_1 \cdot A_{12} + x_2 \cdot A_{22} + x_3 \cdot A_{32}$$

$$y_3 = x_1 \cdot A_{13} + x_2 \cdot A_{23} + x_3 \cdot A_{33}$$

**Scatter x across all nodes!**

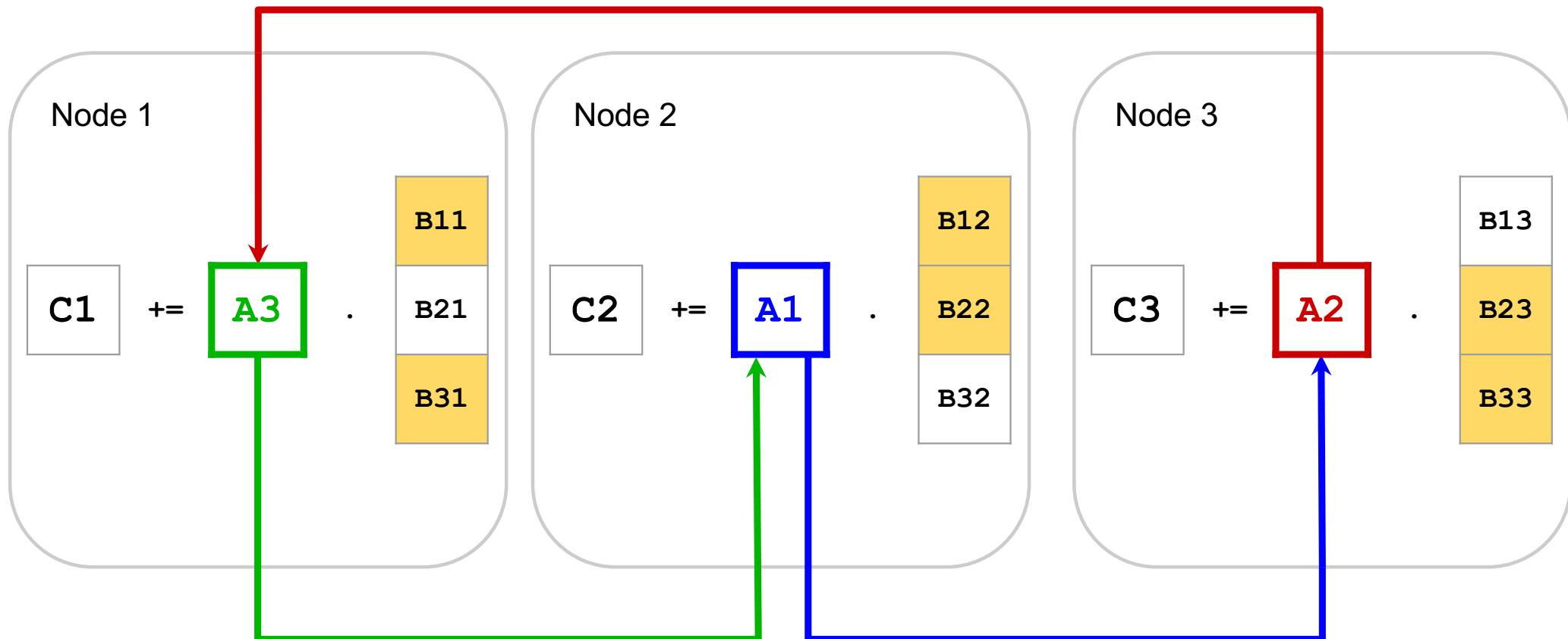


## Parallel stream dot product



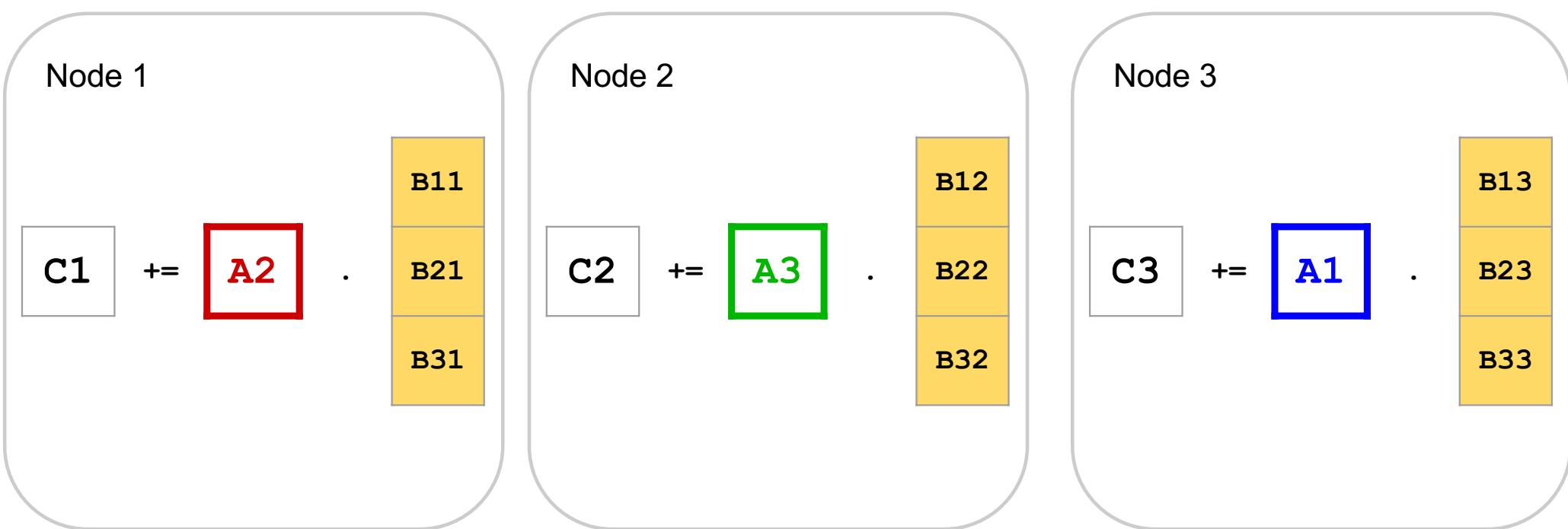


## Parallel stream dot product





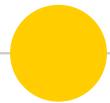
## Parallel stream dot product





## Parallel stream dot product

- Problem – we either have:
  - Both matrices in memory = large memory usage
  - Huge number of rotations = communication overhead.
- Possible solution:
  - Instead of rotating the left matrix, each node creates a thread to send data on demand (only the necessary values are passed).



## Parallel stream dot product

### ◉ Other possibility

- Rotate small blocks / elements of the right side matrix instead of all the left one
- Problems:
  - **C** must be reduced (MPI all to all + multi-merge)

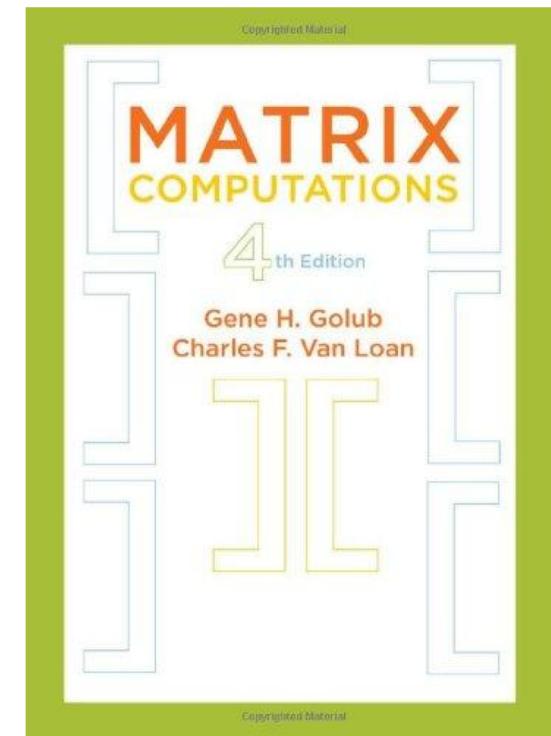
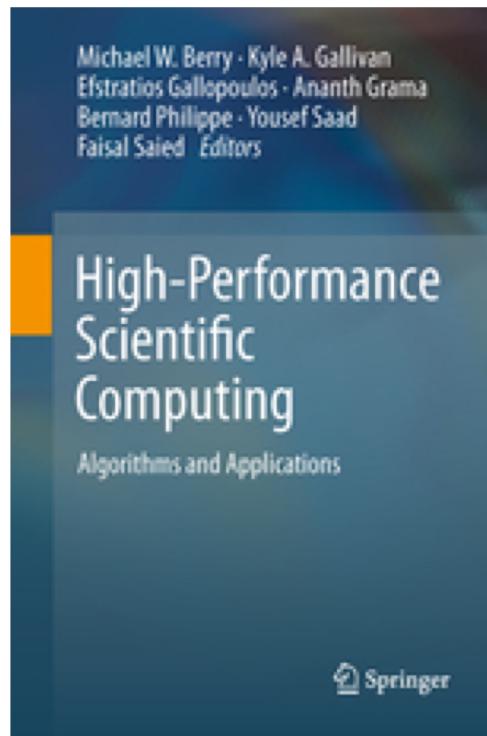


## Parallel stream dot product

- If the result of dot product is passed to a *bang()*, it can benefit from the horizontal partition!
  - Does it happen frequently?
- The *bang* also has dependencies that must be considered to build the operations tree
  - To be studied !!



## Designing parallel solutions



Click on the image to open the book

## Goals for the meeting

19 Oct

- Implement DSL parser
  - Name the DSL?
- Meeting with professor Rui Ralha
  - See slides above



## DSL Parser

### ○ Current state

- Lex regular expressions are done
- Grammar is written
- Lost C functional version ( C++ conversion is not trivial )

### ○ To do

- C++ version of the parser
- Create list of operations
- Create operations dependency tree based on the list
- Reorder operations based on the previous concepts



## DSL Name

- [File extension website](#)
- Suggestions:
  - LAQ - 'LA Query' language
  - LAS - 'LA Scripting' language

## Goals for the meeting

26 Oct

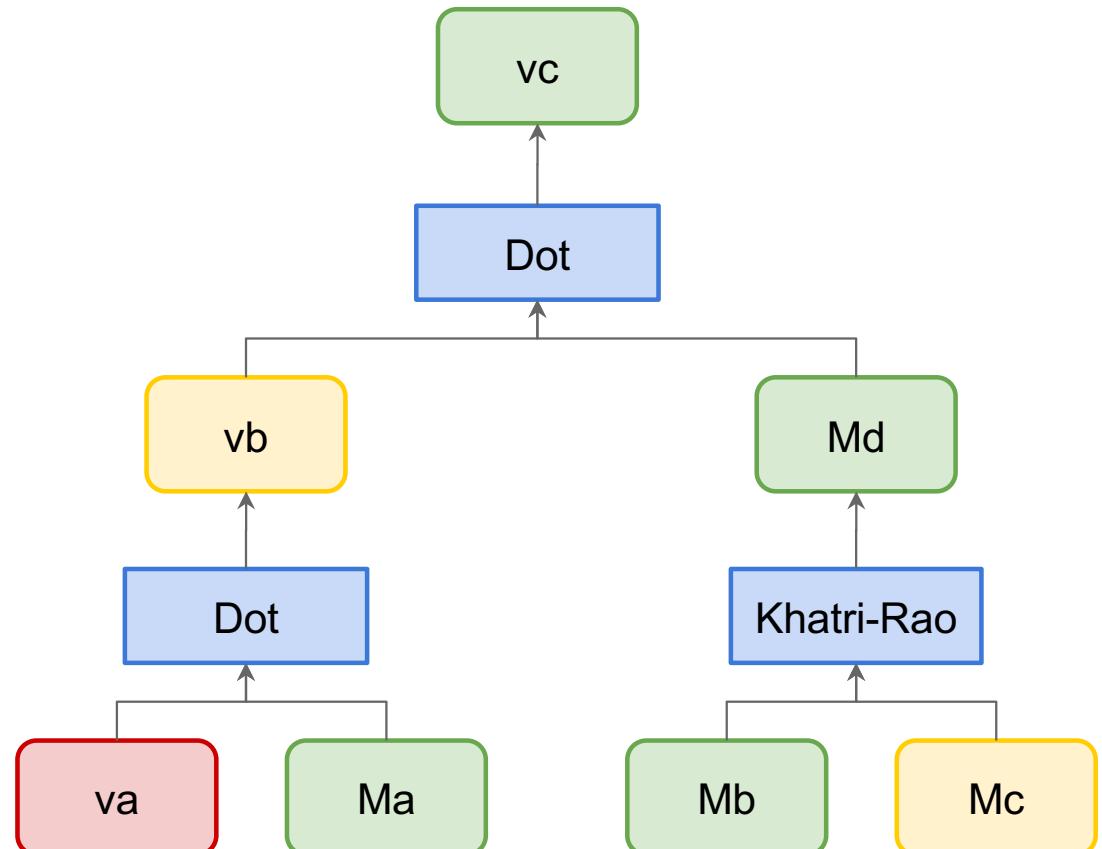
- Problems of a block approach
- Matrix Transpose
- Matrix Vectorization
- Topics to cover in the dissertation's state of the art



## Problems of a block approach

- $va$  in-memory
- $vb[i] = va \cdot Ma[i]$
- How can we calculate  $vc$  streaming both  $vb$  and  $Md$ ?  
 $vc[i] += vb[j] \cdot Md[i][j]$   
 $Md$  must be iterated row-wise

Let's calculate  $Md$  row-wise:



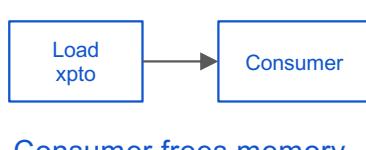
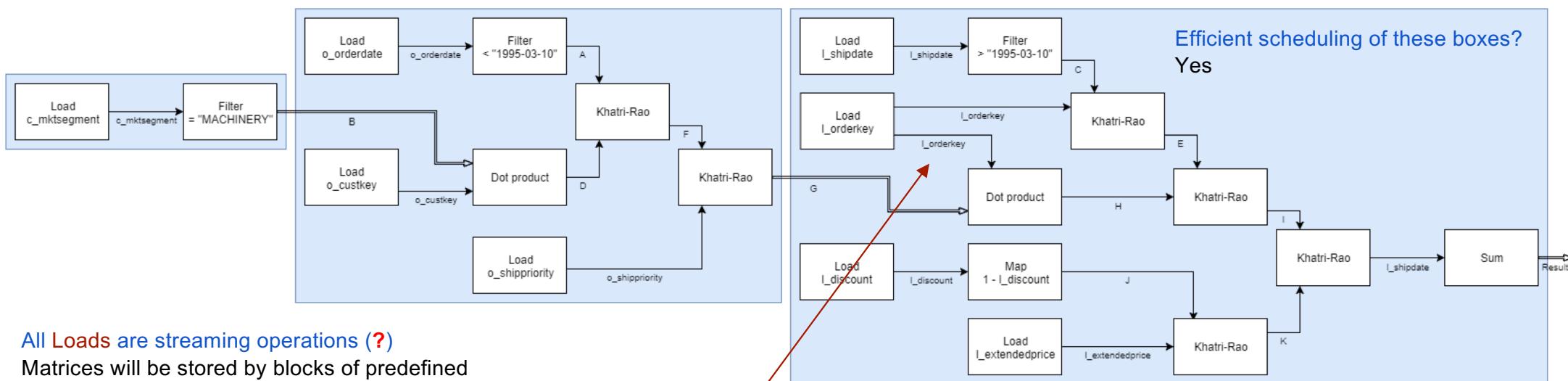
# Vectorized image

```

A = filter_to_matrix( o_orderdate, <'1995-03-10' )
B = filter_to_vector( c_mktsegment, ='MACHINERY' )
C = filter_to_vector( l_shipdate, >'1995-03-10' )
D = dot( B, o_custkey )
E = kraq( l_orderkey, C )
F = kraq( A, D )
G = kraq( F, o_shippriority )
H = dot( G, l_orderkey )
I = kraq( E, H )
J = map( \x -> 1 - x, l_discount )
K = hadamard( l_extendedprice, J )
L = kraq( I, K )
M = bang( sum, L )

```

Full matrix data  
Streaming data

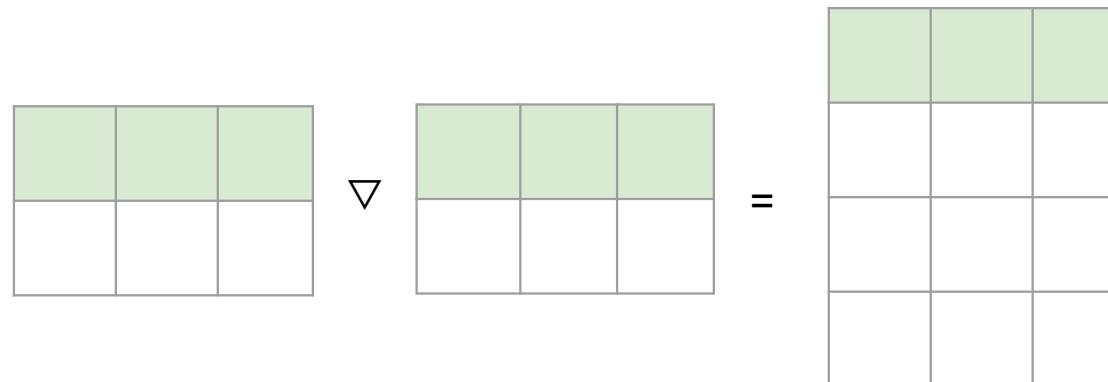


Only last consumer frees memory  
TRUE, and this forks also occur after other operations (not only loads)  
E.g.: TPC-H Query 14 after map operation (G)



## Problems of a block approach

- Option 1) Dense outer matrix grows



- After some KR<sub>s</sub>, dense matrix may not fit in memory (even on disk?)



## Problems of a block approach

- Option 2) Sparse inner blocks grow
  - Keeping the same NNZ

$$\begin{array}{|c|c|c|} \hline \text{light green} & \text{light green} & \text{light green} \\ \hline \text{white} & \text{white} & \text{white} \\ \hline \end{array} \quad \nabla \quad \begin{array}{|c|c|c|} \hline \text{pink} & \text{pink} & \text{pink} \\ \hline \text{pink} & \text{pink} & \text{pink} \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline \text{light green} & \text{light green} & \text{light green} \\ \hline \text{white} & \text{white} & \text{white} \\ \hline \end{array}$$

- We need the full matrix B to compute any row of C
  - The memory optimization is lost!



## Matrix Transpose

1			
	1		
			1

COO:

c: 0 1 3

r: 0 1 0

1	
	1
1	

COO:

c: 0 0 1

r: 0 3 1

- 1st - Swap column and row pointers:  $\Theta(1)$
- 2nd - Sort both arrays (based on columns):
  - $\Theta(\text{NNZ} * \log(\text{NNZ}))$  - depending on sorting algorithm
- Comparing to the other operations –  $\Theta(\text{NNZ})$ 
  - For a million records: **≈20 times slower**
  - For an american trillion records: **≈40 times slower**
  - Can be worst depending on the parallelization strategy

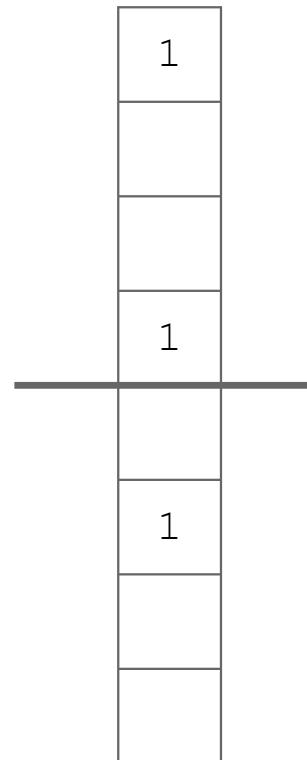
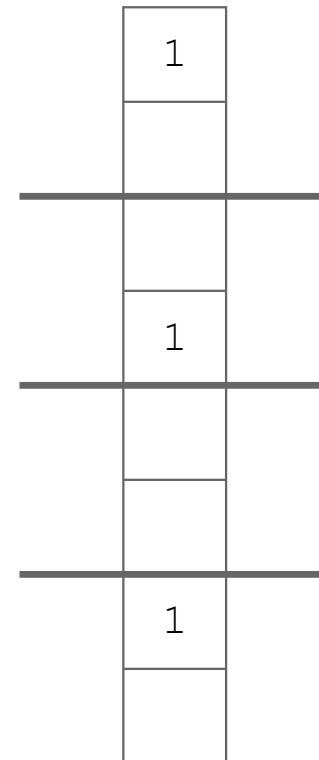


## Matrix Vectorization

1				1
	1			

### COO:

- c: 0 1 3
- r: 0 1 0



○  $R \leftarrow C \Rightarrow C \times R \leftarrow 1$

- $nr = c * R + r$
- $nr: 0 \ 3 \ 6$

○  $R \leftarrow C \Rightarrow R \times C \leftarrow 1$

- $nr = r * C + c$
- $nr: 0 \ 5 \ 3$

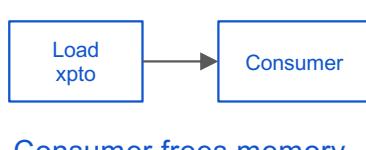
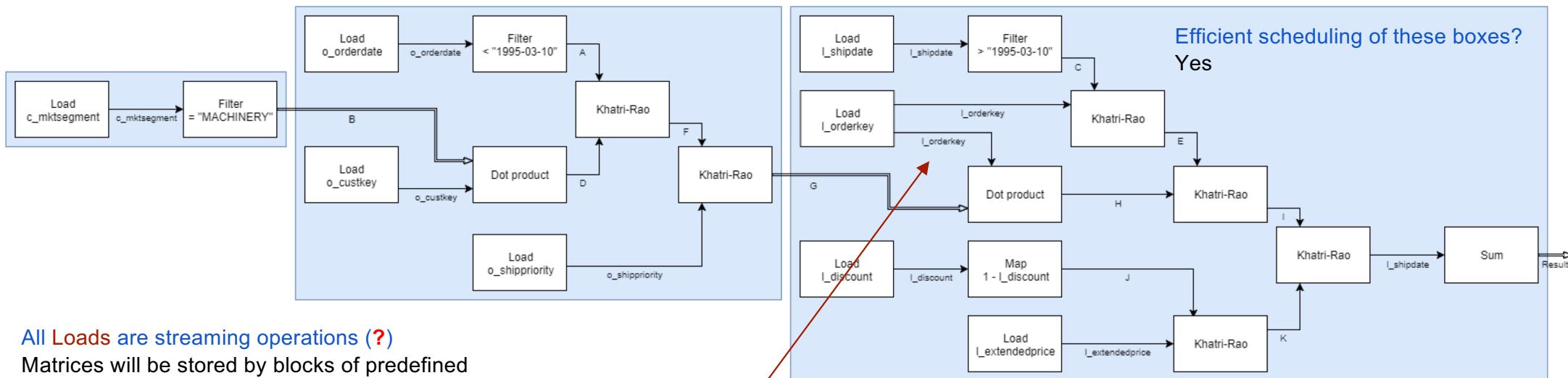
# Vectorized image

```

A = filter_to_matrix( o_orderdate, <'1995-03-10' )
B = filter_to_vector( c_mktsegment, ='MACHINERY' )
C = filter_to_vector( l_shipdate, >'1995-03-10' )
D = dot( B, o_custkey )
E = kraq( l_orderkey, C )
F = kraq( A, D )
G = kraq( F, o_shippriority )
H = dot( G, l_orderkey )
I = kraq( E, H )
J = map( \x -> 1 - x, l_discount )
K = hadamard( l_extendedprice, J )
L = kraq( I, K )
M = bang( sum, L )

```

Full matrix data  
Streaming data



Only last consumer frees memory  
TRUE, and this forks also occur after other operations (not only loads)  
E.g.: TPC-H Query 14 after map operation (G)



## Matrix Vectorization

- Consider a generic COO matrix:  $A \times B \leftarrow C \times D$ 
  - The following vectorizations are trivial and cost  $\Theta(\text{NNZ})$ :
    - $B \leftarrow C \times D \times A$
    - $1 \leftarrow C \times D \times A \times B$
    - $D \times A \times B \leftarrow C$
    - $C \times D \times A \times B \leftarrow 1$
  - Any other needs sorting, thus costing  $\Theta(\text{NNZ} * \log(\text{NNZ}))$



## Vectorization analysis

From	To	Orientation	Read Cost	Write Cost	Memory Used
COO	COO	Vertical	$\Theta(\text{NNZ})$	$\Theta(\text{NNZ})$	NNZ
		Horizontal	$\Theta(\text{NNZ})$	$\Theta(\text{NNZ})$	NNZ
	CSC	Vertical	$\Theta(\text{NNZ})$	$\Theta(\text{NNZ})$	NNZ
		Horizontal	$\Theta(\text{NNZ})$	$\Theta(R \times C)$	$R \times C$
CSC	COO	Vertical	$\Theta(C)$	$\Theta(\text{NNZ})$	NNZ
		Horizontal	$\Theta(C)$	$\Theta(\text{NNZ})$	NNZ
	CSC	Vertical	$\Theta(C)$	$\Theta(\text{NNZ})$	NNZ
		Horizontal	$\Theta(C)$	$\Theta(R \times C)$	$R \times C$



## **State of the art**

---

- **Relational DBMS**
- **Business Intelligence**
  - OLAP
- **Big Data / Business Analytics**
  - OLTP / Stream processing
- **Market solutions**
  - Relational: PostgreSQL
  - OLAP: Apache kylin, Hive?
  - Hybrid: Apache Drill, Impala?
  - Cloud: Google BigQuery, Amazon RedShift
- **The Typed LA approach**
  - Foundation of the previous report
- **Benchmarks**
  - TPC-H, papers

## Goals for the meeting

2 Nov

- Matrix transpose via swapping & vectorization
- DSL
- DSL parser



## Matrix transpose via swapping & vectorization

- $T = \text{swap} \cdot \text{vec}(M)$
- $\text{vec}(M)$  can be calculated in  $\Theta(\text{NNZ})$  [cf slide 40]
- $\text{swap} = \text{snd} \triangleright \text{fst}$ 
  - $\text{row} = (\text{col} \% R) * C + (\text{col} / R)$
  - $\text{col} = (\text{row} \% C) * R + (\text{row} / C)$
- $\text{swap}$  does not need to be created!
- For each elem in  $\text{vec}(M)$ 
  - $\text{col\_swap} = \text{row\_elem}$
  - $\text{row\_swap} = (\text{col\_swap \% R}) * C + (\text{col\_swap} / R)$
  - **$T.\text{insert}(\text{col}=\text{row\_swap})$  Too costly: search+shift!**



## DSL new operator

- Renamed from LAS to LAQ due to name conflict with [LASSO](#) (.las)
- Added **return ( v1, v2, ..., vn )** statement
  - The return is always the root node of the operators tree!
  - Otherwise multiple endpoints were allowed
    - E.g.: **Select v1, v2, ..., vn From ...**



## DSL parser

- Slow development while learning C++
- Bison and flex completed!
- Operation tree constructed using adjacency lists (graphs)
  - $\Theta(\approx 1)$  search cost
  - Tree structure using “return” as root node
- What to do next?
  1. Implement operators to convert the DSL in C++
  2. Implement optimizations on the operators tree

## Goals for the meeting

16 Nov

- Optimized Khatri–Rao between blocks
- From block to stream level
  - Who manages matrix blocks? EDIT: Considering doing it in the main function
- Data load
  - Concurrent queue? EDIT: Is protobuf thread-safe?



## Block level Khatri-Rao

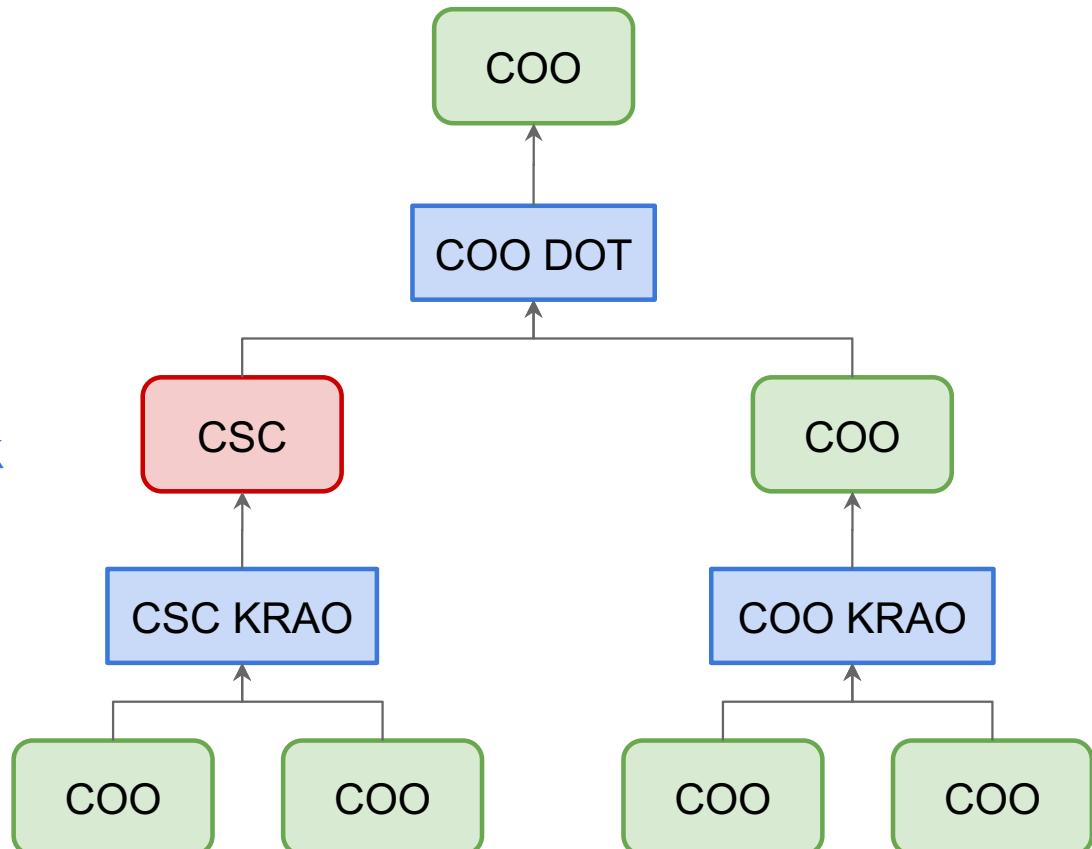
	Values	Rows	Columns
Decimal Vector			
Bitmap			
Filtered Bit Vector			
Decimal Map			
Filtered Bitmap			
Filtered Decimal Vector			
Filtered Decimal Map			

7 types of matrices:  $7 \times 7 = 49$  possible operations



## Block level Khatri-Rao

- CSC or COO output
- $C = A \downarrow B$ 
  - **Possibly destroy A and/or B to optimize operations**
  - **EDIT: Too soon to do it with protobuf, wait for benchmark results**
- **$49 * 2 * 2 = 196$  combinations**
- Some functions can be reused:
  - **Total: 89**
  - **EDIT: maybe 61**



## Goals for the meeting

23-11

- Write some example queries in C++ format
  - Look for a default syntax
- Matrix data serialization module
  - Choose a library/framework
  - Define the data structure
- Rewrite operations (Khatri–Rao) using the lib
- Define unit tests for the parser and the krao operator
- Design data io API
- Design bang: sum/count/avg functions



## Data serialization

- Select a serialization API
  - Based on: [C++ serializers benchmark](#), TensorFlow & Apache Calcite:
    - Google Protocol Buffer
- Code modified to use protobuf generated classes
- Functional version tested on local machine
  - Protobuf (and ALL dependencies) installed on SeARCH
    - Still with linking issues
    - Possibly replace Makefile with CMake
      - Short available time to learn CMake: **is it worth the time?**



## Matrix protobuf specification

```
1 syntax = "proto3";
2
3 package engine;
4 import "src/block.proto";
5 import "src/label_block.proto";
6
7 option cc_enable_arenas = true;
8
9 message matrix {
10     // matrix identifiers
11     string database_path = 1;
12     string table = 2;
13     string column = 3;
14
15     // matrix dimensions
16     int32 nnz = 4;
17     int32 nrows = 5;
18     int32 ncols = 6;
19
20     // matrix content
21     map<int32, block> blocks = 7;
22
23     // matrix labels
24     map<int32, label_block> labels = 8;
25 }
26
```



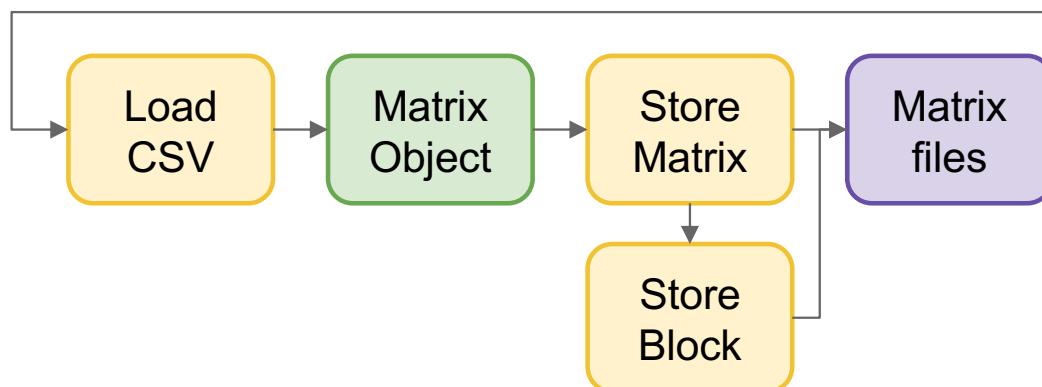
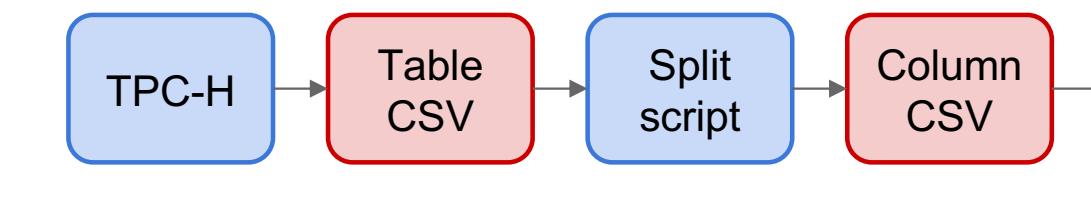
## Disk data structure

- database\_path/
  - tableA/
    - column1/
      - blocks/
        - 1.dat, 2.dat, ...
      - matrix.dat
      - labels/
        - 1.dat, 2.dat, ...
    - column2/
    - ...
  - tableB/
  - ...

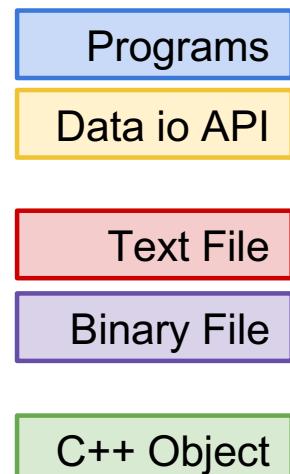
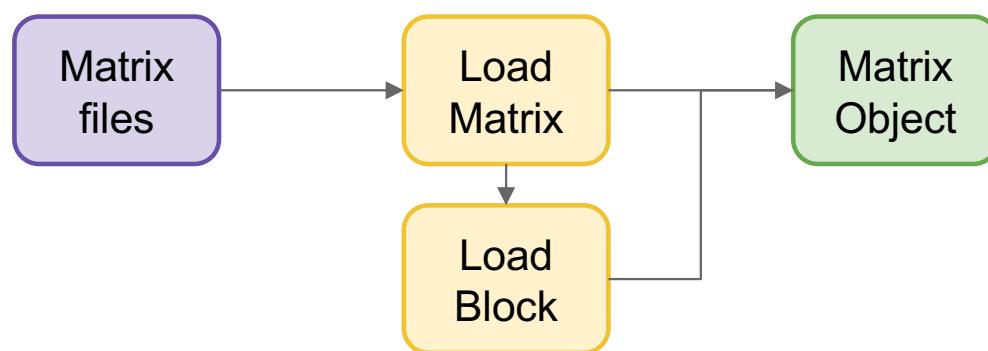


## Data io API

Import CSV /  
store data:



Load Data:





## Current status

Needed to  
run TPC\_H  
query 6

	Design	Implementation	Test
Khatri-Rao (krao)			
Hadamard (krao)			
Load			
Bang			
Filter			
Map			
Dot product			

- Queries 3, 6 execution plan (streaming)
- Query rewriter
- Label normalization
- Load API
  - Considerations on database schemas
- MySQL as example
  - More user friendly source code than PostgreSQL
- In-memory databases
  - Performance of non-volatile RAM? E.g. Intel Optane

# TPC-H

## Query 3

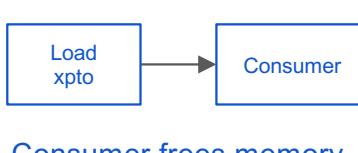
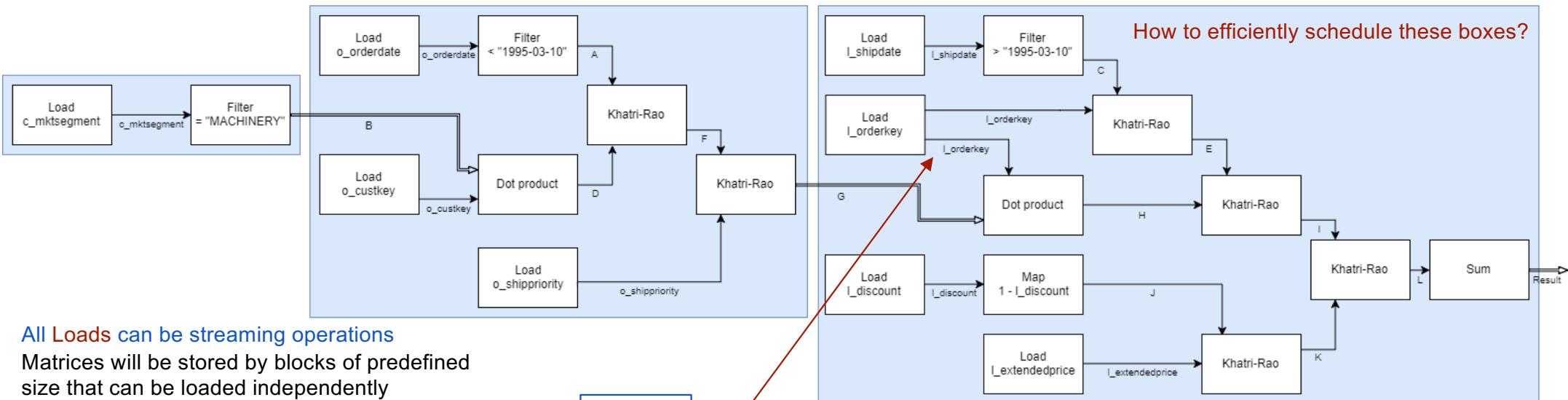
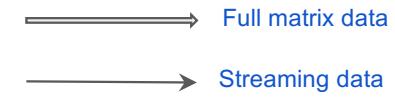
Streaming approach

```

A = filter_to_matrix( o_orderdate, <'1995-03-10' )
B = filter_to_vector( c_mktsegment, ='MACHINERY' )
C = filter_to_vector( l_shipdate, >'1995-03-10' )
D = dot( B, o_custkey )
E = krao( l_orderkey, C )
F = krao( A, D )
G = krao( F, o_shippriority )
H = dot( G, l_orderkey )
I = krao( E, H )
J = map( \x -> 1 - x, l_discount )
K = hadamard( l_extendedprice, J )
L = krao( I, K )
M = bang( sum, L )

```

SVG



Only last consumer frees memory  
**TRUE**, and this forks also occur after other operations (not only loads)  
E.g.: TPC-H Query 14 after map operation (G)

# TPC-H Query 6

Streaming approach  
(Unoptimized)

```

A = filter( l.shipdate >= "1994-01-01" AND l.shipdate <= "1995-01-01" )
B = filter( l.discount >= 0.05 AND l.discount <= 0.07 )
C = filter( l.quantity < 24 )
D = hadamard( A, B )
E = hadamard( C, D )
F = map( l.extendedprice * l.discount )
G = hadamard( E, F )
H = sum( G ) |
return ( H )

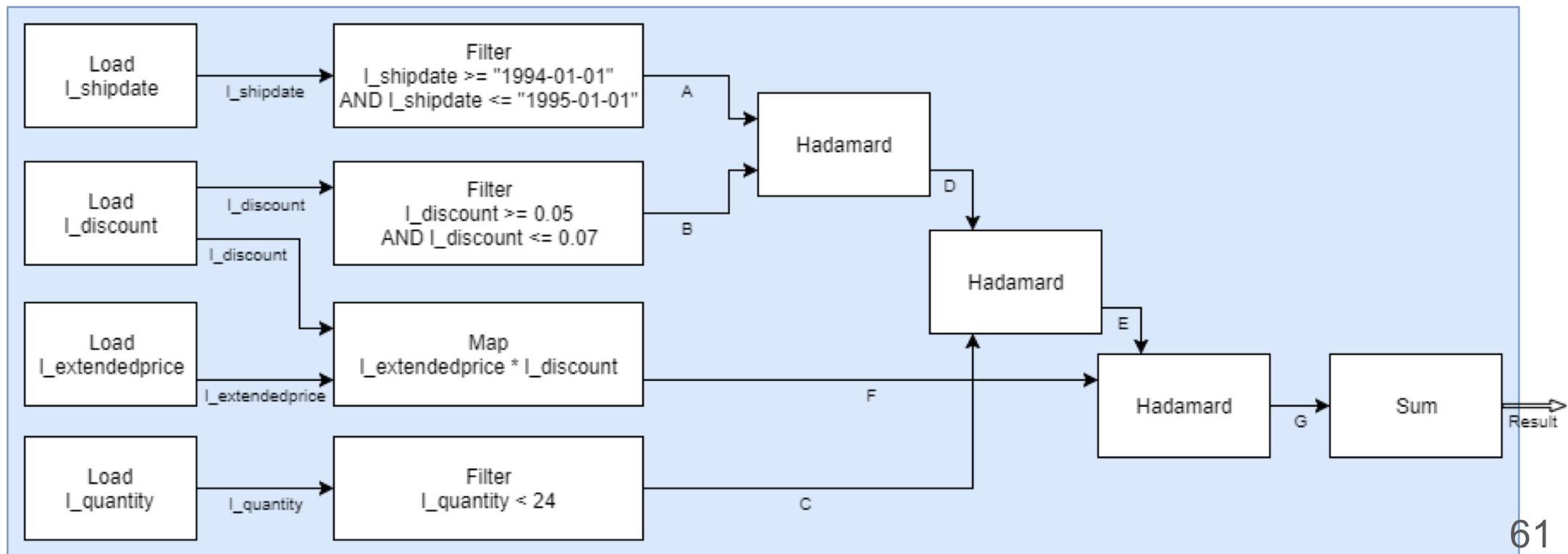
```

```

A = filter_to_vector( l_shipdate, >='1995-03-10', <'1995-03-10' )
B = filter_to_vector( l_discount, >='0.49', <='0.51' )
C = filter_to_vector( l_quantity, <='3' )
D = hadamard( A , B )
E = hadamard( C , D )
F = hadamard( l_extendedprice , l_discount )
G = hadamard( E , F )
H = bang( sum, G )

```

Before...



# TPC-H

## Query 6

Streaming approach  
(Optimized)

Optimized how? More parallelism,  
shorter pipeline  
=> does not guarantee speedup!  
Can this be automated? Yes, but  
LA properties must be considered  
Which ones?

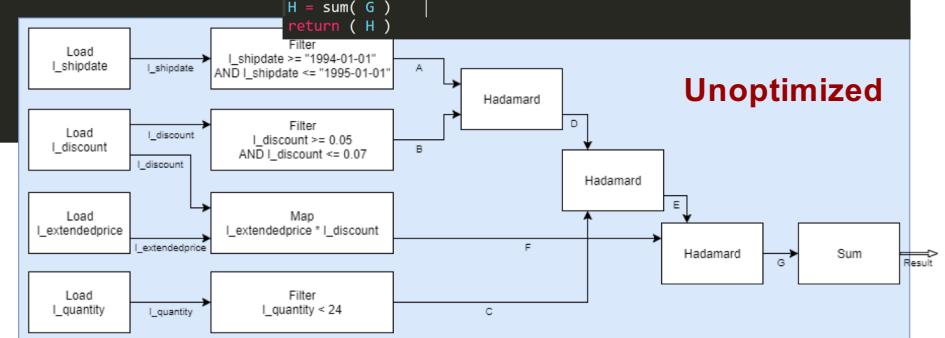
E.g. associative? in Hadamard,  
filter and map  
( $A \times B \times C = A \times (B \times C)$ )

```

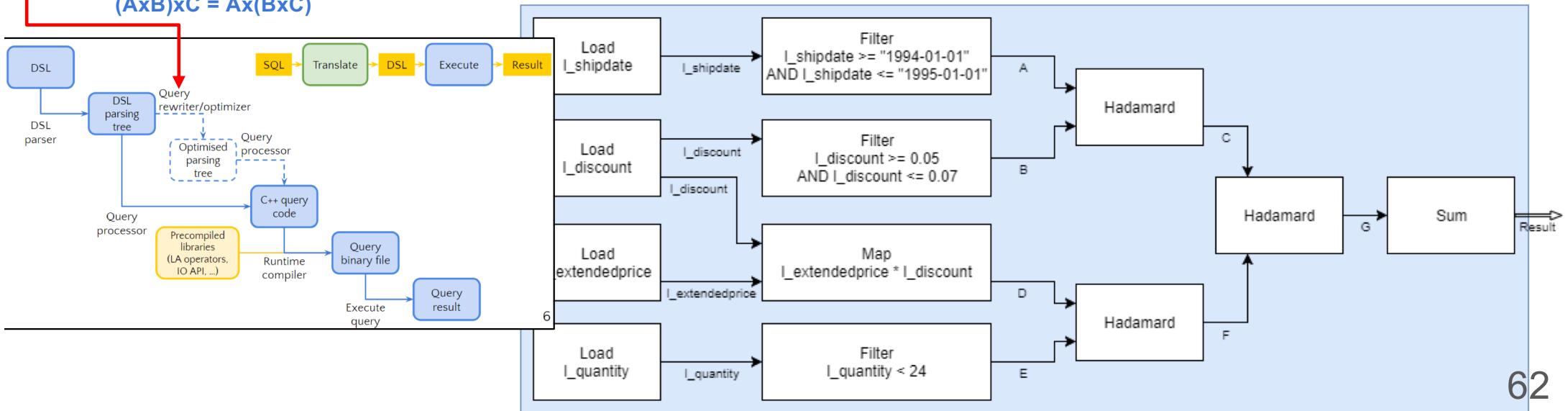
A = filter( l.shipdate >= "1994-01-01" AND l.shipdate <= "1995-01-01" )
B = filter( l.discount >= 0.05 AND l.discount <= 0.07 )
C = hadamard( A, B )
D = map( l.extendedprice * l.discount )
E = filter( l.quantity < 24 )
F = hadamard( D, E )
G = hadamard( C, F )
H = sum( G )
return ( H )
    
```

```

A = filter( l.shipdate >= "1994-01-01" AND l.shipdate <= "1995-01-01" )
B = filter( l.discount >= 0.05 AND l.discount <= 0.07 )
C = filter( l.quantity < 24 )
D = hadamard( A, B )
E = hadamard( C, D )
F = map( l.extendedprice * l.discount )
G = hadamard( E, F )
H = sum( G )
return ( H )
    
```



Unoptimized



62



## Query rewriter

- Before proceeding to query translation, rewrite them to avoid
  - **SELECT FROM/WHERE/IN SELECT ? see slides 63-64!**
  - CASE statements
  - Complex aggregations
    - E.g. AVG + SUM => SUM + COUNT
  - (INNER / OUTER / LEFT / RIGHT / ... ) JOIN statements
    - Replace with WHERE clause(s)
- Example in the next slides...



## Query rewriter (TPC-H query 1)

```
select
    l_returnflag,
    l_linenumber,
    sum(l_quantity) as sum_qty,
    sum(l_extendedprice) as sum_base_price,
    sum(l_extendedprice * (1 - l_discount)) as sum_disc_price,
    sum(l_extendedprice * (1 - l_discount) * (1 + l_tax)) as sum_charge,
    avg(l_quantity) as avg_qty,
    avg(l_extendedprice) as avg_price,
    avg(l_discount) as avg_disc,
    count(*) as count_order
```

- L\_quantity     L\_extendedprice     L\_discount     L\_tax
  - Sum
  - Avg
  - Count
- L\_quantity     L\_extendedprice     L\_discount     L\_tax
  - Sum
  - Avg
  - Count
- L\_quantity     L\_extendedprice     L\_discount     L\_tax
  - Sum
  - Avg
  - Count
- L\_quantity     L\_extendedprice     L\_discount     L\_tax
  - Sum
  - Avg
  - Count



## Query rewriter (TPC-H query 1)

- l\_extendedprice \* (1 - l\_discount) is calculated multiple times

```
select
    sum(l_extendedprice * (1 - l_discount)) as sum_disc_price,          A
    sum(l_extendedprice * (1 - l_discount) * (1 + l_tax)) as sum_charge,
from
    lineitem
```



```
select
    sum(tmp) as sum_disc_price,                                     B
    sum(tmp * (1+l_tax)) as sum_charge
from
    lineitem,
    select
        l_extendedprice * (1 - l_discount) as tmp
    from
        lineitem
```

Considering the translation algorithm:

- is it easier to process a SELECT FROM SELECT (B)
- or to identify similar aggregations, creating tmp variables in the DSL to avoid repeated calculus (A) ?

## ○ Let:

- `count(*) = count(l_quantity)`
- `Avg = sum / count`

```
select ...  
        sum(l_quantity) as sum_qty,  
        sum(l_extendedprice) as sum_base_price,  
        ...  
        avg(l_quantity) as avg_qty,  
        avg(l_extendedprice) as avg_price,  
        avg(l_discount) as avg_disc,  
        count(*) as count_order
```

A



```
select ...  
        sum_qty/count_order as avg_qty,  
        sum_base_price/count_order as avg_price,  
        sum(l_discount)/count_order as avg_disc  
from lineitem,  
     select ...  
        sum(l_quantity) as sum_qty,  
        sum(l_extendedprice) as sum_base_price,  
        ...  
        count(l_quantity) as count_order
```

B

Again:

- is it easier to process a `SELECT FROM SELECT` (B)
- or to dive AVG into SUM/COUNT avoiding repeated calculus (A) ?



## Summary

---

- SQL must be processed before the translation takes place
- Current examples can become deprecated, depending on both **rewriting** and **translation** algorithms
  - Complete DSL engine before using SQL



## Label normalization

PKey (A) : a, b, c, d

A	#0	#1	#2	#3
a	1			
b		1		
c			1	
d				1

FKey (B) : b, c, b, a

B	#0	#1	#2	#3
b	1		1	
c		1		
a				1

Transforming matrix B in conformity with the labels of A

```
new_row = labels_A(labels_B(old_row))
```

B	#0	#1	#2	#3
a				1
b	1		1	
c			1	
d				



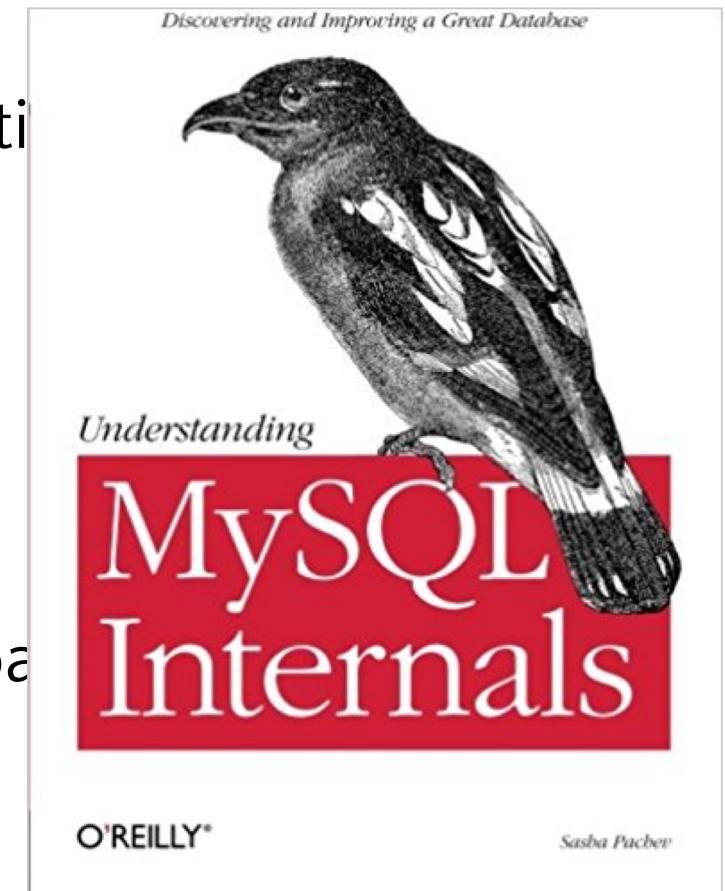
## Load API

- Import data from CSVs to bitmaps
  - Database with a defined schema
    - Load foreign keys respecting pk labels
      - No overhead; drops functionalities
  - Schema free database
    - Load the needed CSVs for each query execution
      - High performance overhead
    - Store matrices with labels extracted from their own data; map fk matrices into pk labels
      - Limited overhead



## DBMS design patterns

- Reading about MySQL implementation
  - Trying to understand
    - Log based consistency
    - Storage mechanisms
      - Database schemas and rules
      - Table data
- **NOTE:** MySQL has no support for parallel execution of a single query





## In-memory databases

- In-memory databases keep all data on RAM
  - Not only during query processing
- Scalability in data size is ensured by the addition of computing nodes, each one storing part of the database
- Persistence is granted by non-volatile RAM
  - Intel Optane [benchmark](#) shows similar performance to SSD disks, although outscaling in read throughput of small files (4k)

## Goals for the meeting

14-12

- Storage format
- Operations
  - Map
  - Filter
- Data creation path



## Storage format

- Constraints:

- Redundant data should not be stored
- Values should only be loaded in filter/map/bang operations
- Allow streamed operations

- Solution:

- Store Dimensions in bitmaps
- Store Measures in dense arrays
  - Otherwise map can't be streamed



## Operations

- Map

- Can only be applied to measures
- Simply iterate over blocks

- Filter (Measure)

- Same as map, but returns a sparse bit vector

- Filter (Dimension)

- Iterate through the labels instead of values
- **A Dot Product is needed to complete the operation**
  - Query 6 won't be fully streamed: 2 filters on dimensions



## Current status

Needed to  
run TPC\_H  
query 6

	Design	Implementation	Test
Khatri-Rao (krao)			
Hadamard (krao)			
<b>Data creation</b>			
Bang			
Filter			
Map			
<b>Dot product</b>			



## Data creation path

- Create

- Database
  - Table (with all columns)

- Insert

- Value into column block

- Store

- Block
  - Database metadata

- Load

- Database metadata
  - Block

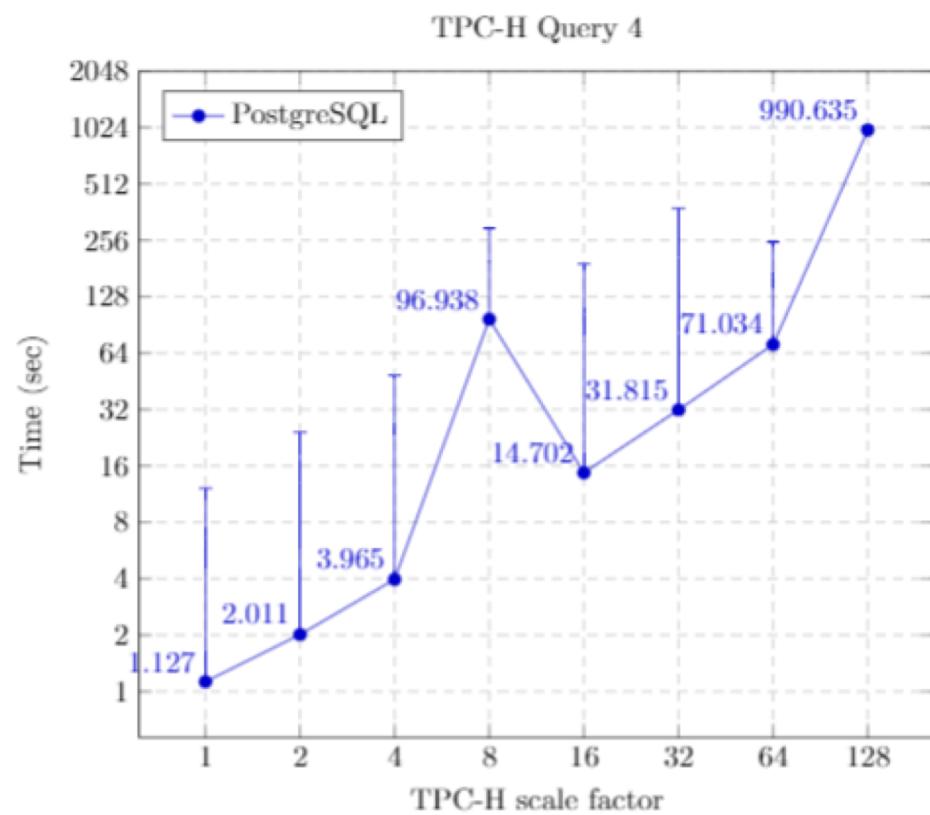
## Goals for the meeting

23-01

- Benchmarks for the paper



## Variable times due to caching issues:



Mark: best of 2

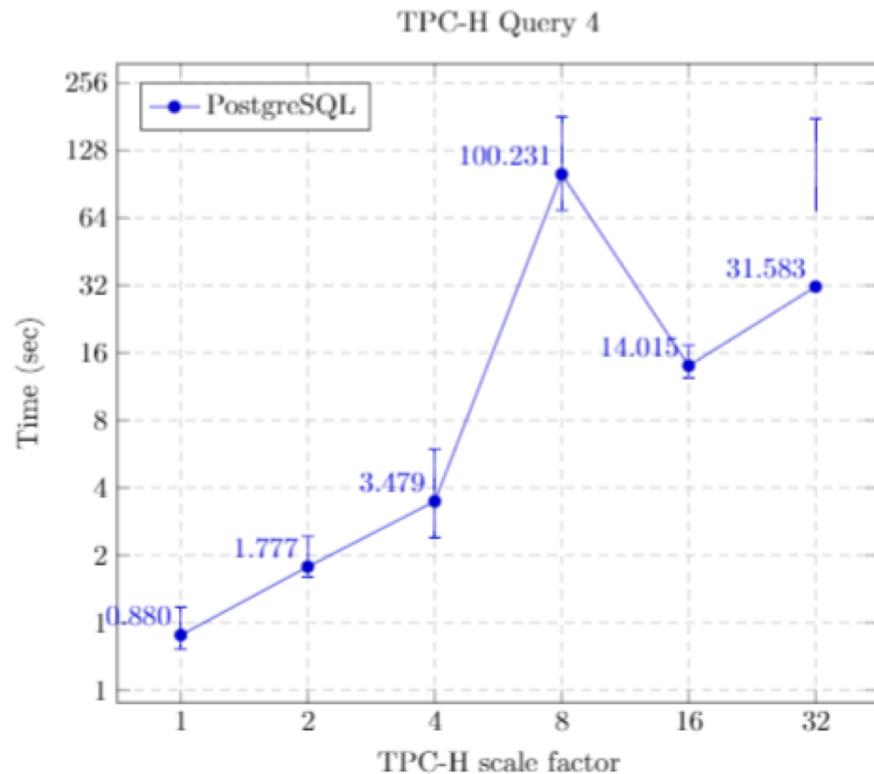
Bar center: Average

Bar size: Standard deviation



# Postgres

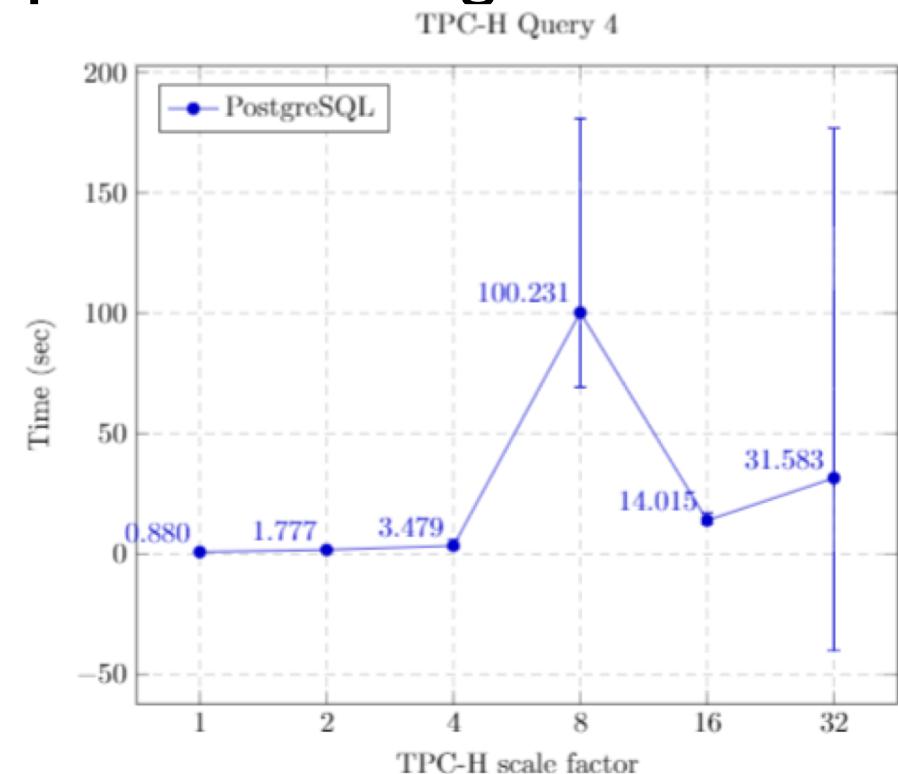
Sometimes SD>AVG, thus unplotable in logarithmic



Mark: 3-best 10 (5%)

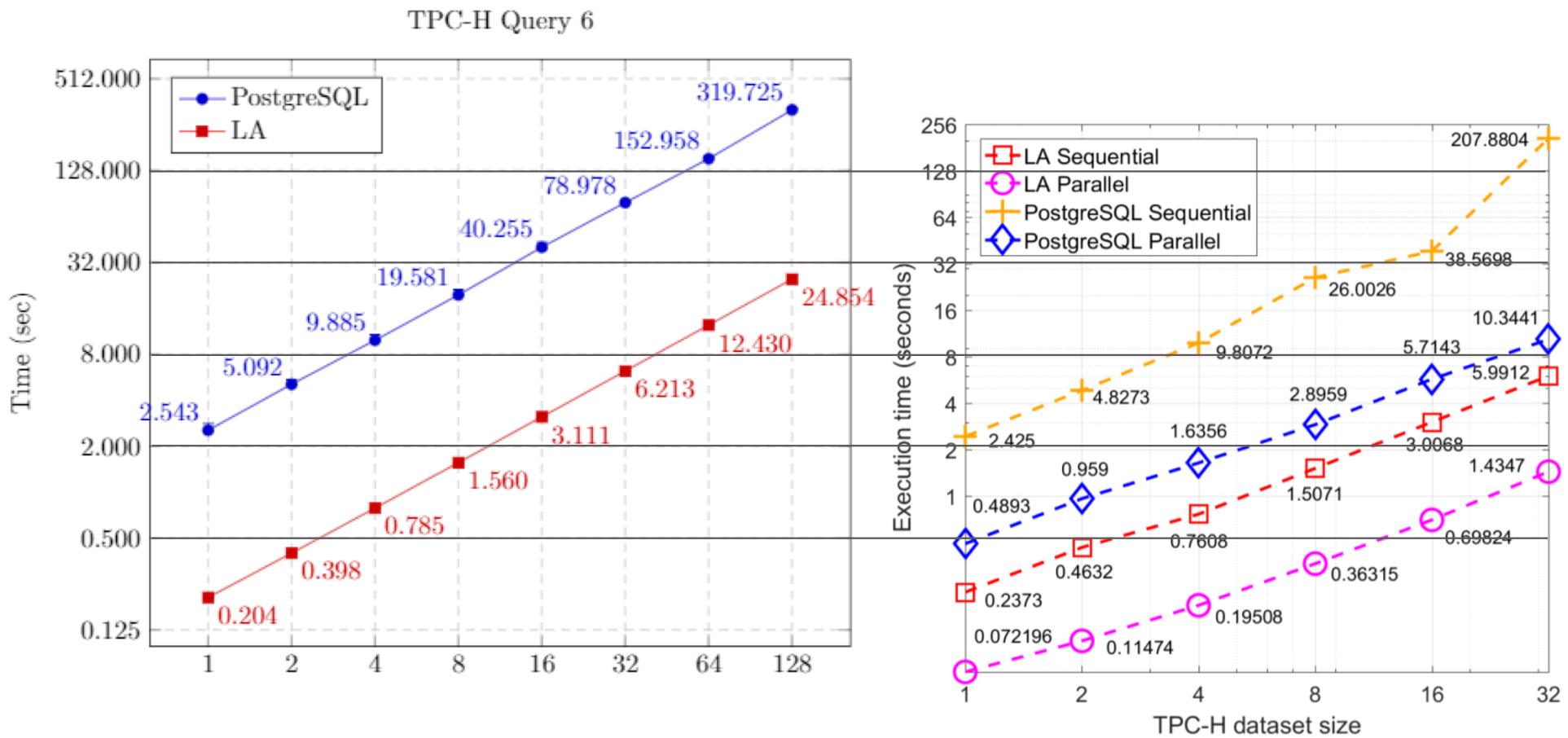
Bar center: AVG

Bar size: SD





## Query 6



- Query 4: semi-join statement



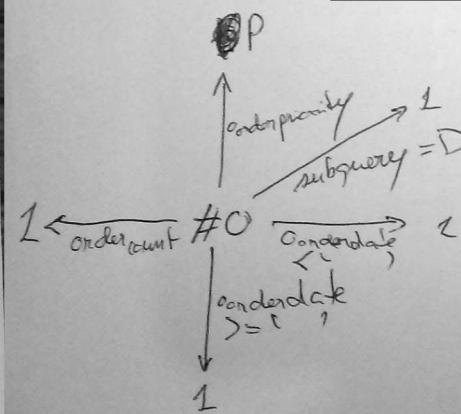
## Query 4

# Slower than Postgres

Query A

```
#l orderkey > #o
1 \ l_c < l_n
A = hadamard (l_commitdate <; l_receiptdate) 1 ← #l
B = krao (A, l_orderkey) #o ← #l
C = avl_exists (B) AVL
NOTE = min OR max OR avg
D = avl_to_vector (C) 1 ← #o boolean
```

**Bottleneck**



```
select
    o_orderpriority,
    count(*) as order_count
from
    orders
where
    o_orderdate >= date ':1'
    and o_orderdate < date ':1' + interval '3' month
    and exists (
        select
            *
        from
            lineitem
        where
            l_orderkey = o_orderkey
            and l_commitdate < l_receiptdate
    )
group by
    o_orderpriority
```

$$\begin{aligned}
 E &= \text{filter}(\text{orderdate} < \text{l_receiptdate}) \quad 1 \leftarrow \#o \\
 F &= \text{filter}(\text{orderdate} \geq \text{l_commitdate}) \quad 1 \leftarrow \#o \\
 G &= \text{hadamard}(E, F) \quad 1 \leftarrow \#o \\
 H &= \text{hadamard}(G, D) \quad 1 \leftarrow \#o \\
 I &= \text{krao}(\text{o_orderpriority}, H) \quad P \leftarrow \#o \\
 J &= \text{avl\_count}(I) \quad \text{AVL} \\
 K &= \text{avl\_to\_vec}(J) \quad 1 \leftarrow P
 \end{aligned}$$



## Query 4

A = filter (l.commitdate < l.recipdate)

B = knac (A, l.orderkey)

C = filter (o.orderdate < 1993...)

D = filter (o.orderdate >= 1993...)

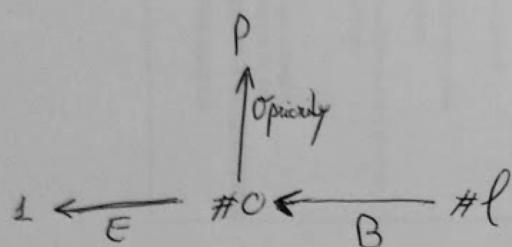
E = hadamard (C, D)

$\perp \leftarrow \#l$   
 $\#o \leftarrow \#l$   
 $\#l \leftarrow \#c$   
 $\perp \leftarrow \#c$   
 $\perp \leftarrow \#c$

**Wrong result!**

DSL error!

(G) is a join, the query requires a semi-join: EXISTS



F = knac (o.orderpriority, E)

G = dot (F, B)

H = count (G)

$P \leftarrow \#o$   
 $P \leftarrow \#l$   
 $P \leftarrow 1$



## Query 4

$A = \text{filter}(\ell_{\text{commitdate}} < \ell_{\text{recipdate}})$

$B = \text{Krao}(A, \ell_{\text{orderkey}})$

$C = \text{filter}(\ell_{\text{orderdate}} < 1993\dots)$

$D = \text{filter}(\ell_{\text{orderdate}} \geq 1993\dots)$

$E = \text{hadamard}(c, \delta)$

$F = \text{Krao}(\ell_{\text{orderpriority}}, E)$

$G = \text{dot}(F, B)$

$H = \text{Krao}(B, G)$

$I = \text{exists}(H)$

avg, min, max have same result

$\bar{J} = \text{unvect}(I)$

$K = \text{count}(\bar{J})$

$1 \leftarrow \#l$

$\#0 \leftarrow \#l$

$1 \leftarrow \#0$

$1 \leftarrow \#0$

$P \leftarrow \#0$

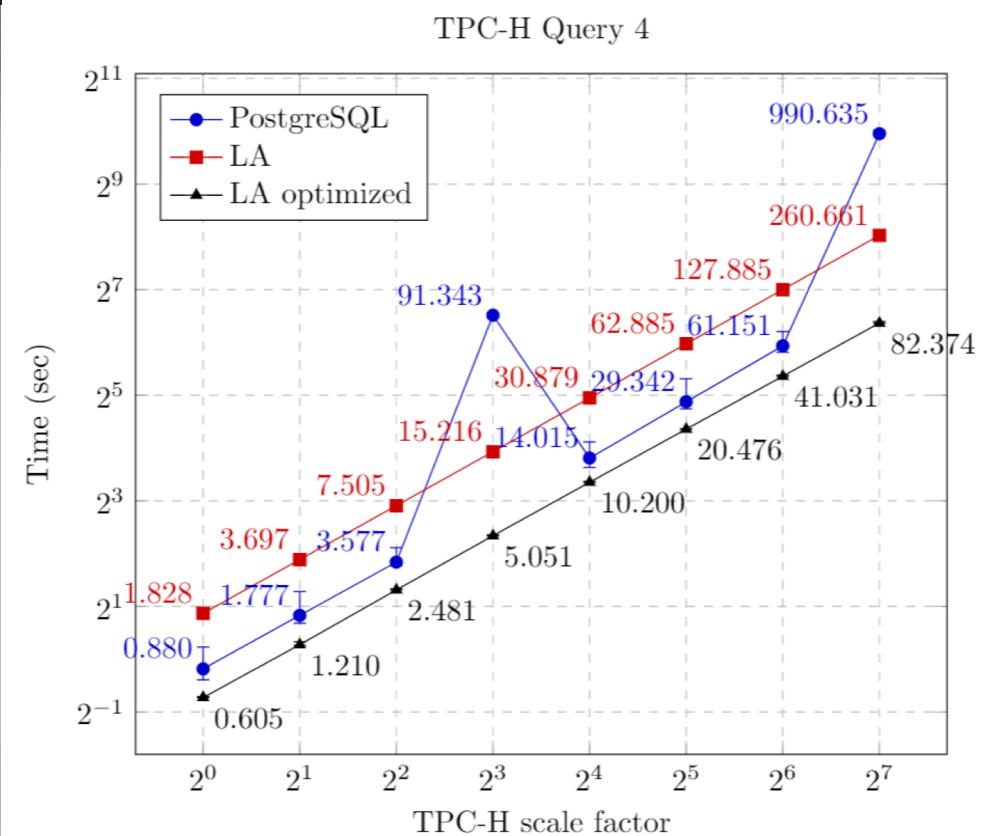
$P \leftarrow \#l$

$\#0 \times P \leftarrow \#l$

$\#0 \times P \leftarrow 1$

$P \leftarrow \#0$

$P \leftarrow 1$



## MySQL configuration file: my.cnf

```
[client]
port = 3306
socket = /home/laolap/new_libs/mysql/mysql/mysql.sock

[mysql]
no_auto_rehash
default_character_set = utf8

[mysql_safe]
user = mysql
log-error = /home/laolap/new_libs/mysql/mysql/laolap_error.log

[mysql]
port = 3306
max_open_files = 2048
table_open_cache = 1024
explicit_defaults_for_timestamp = 1
socket = /home/laolap/new_libs/mysql/mysql/mysql.sock
basedir = /home/laolap/new_libs/mysql/mysql
datadir = /home/laolap/new_libs/mysql/mysql/data
tmpdir = /home/laolap/new_libs/mysql/mysql/tmp

default_storage_engine = MEMORY
disable_partition_engine_check = true

query_cache_type = 1
query_cache_size = 1M
```

- Redundant filters
  - Measures
  - Dimensions
  - Future further optimization
- Further optimize aggregations
- Future Khatri–Rao/Hadamard optimization



# Redundant filters (measures)

## Data

Facts	
Price	Quantity
5.1	9
2.3	8
3.5	9
4.7	7
4.2	5

## SQL query (part of)

```
WHERE price < 2.3  
AND quant > 5
```

## DSL query (1)

```
A = filter( f.price < 2.3 )  
B = filter( f.quant > 5 )  
C = hadamard( A, B )
```

$$\mathbf{A} = \begin{matrix} - & 1 & - & - & - \end{matrix}$$

$$\mathbf{B} = \begin{matrix} 1 & 1 & 1 & 1 & - \end{matrix}$$

$$\mathbf{C} = \begin{matrix} - & 1 & - & - & - \end{matrix}$$

5+5+4 iterations

Desired: 5+1

## DSL query (2)

```
A = filter( f.price < 2.3  
AND f.quant > 5 )
```

## Solutions:

- Iterate both matrices using C && clause (stops on false)
- Iterate **price** and then iterate the resultant vector, filtering by **quantity**



# Redundant filters (dimensions)

Data

Car	
Color	Brand
Black	Ford
Blue	Ford
Black	Ford
White	Ford
White	Ferrari

SQL query (part of)

```
WHERE color = 'Blue'  
      AND brand = 'Ford'
```

DSL query (1)

```
A = filter( c.color = 'Blue' )  
B = filter( c.brand = 'Ford' )  
C = hadamard( A, B )
```

Black	Blue	White
-	1	-

1	-	1	-	-
-	1	-	-	-
-	-	-	1	1

$$A = \begin{bmatrix} - & 1 & - & - & - \end{bmatrix}$$

DSL query (2)

```
A = filter( c.color = 'Blue'  
          AND c.brand = 'Ford' )
```

$$\begin{array}{c|c} \text{Ford} & \text{Ferrari} \\ \hline 1 & - \end{array} \bullet \begin{array}{c|c|c|c|c|c} 1 & 1 & 1 & 1 & - \\ - & - & - & - & - & 1 \end{array}$$

$$B = \begin{bmatrix} 1 & 1 & 1 & 1 & - \end{bmatrix}$$

$$C = AxB = \begin{bmatrix} - & 1 & - & - & - \end{bmatrix}$$

Next slide implements DSL (2)



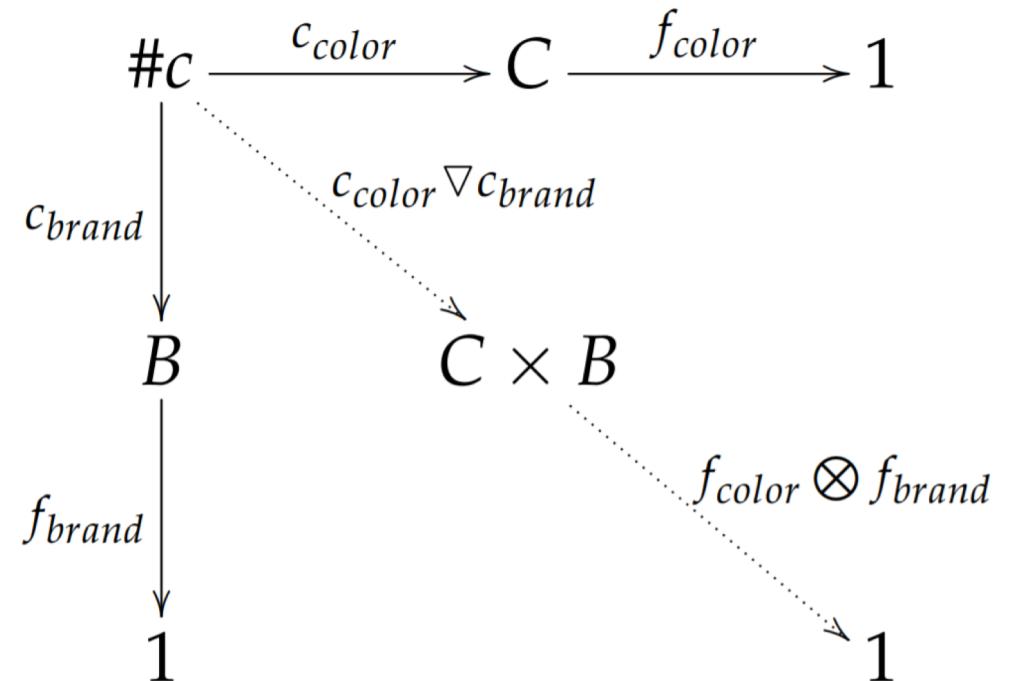
## Redundant filters (dimensions)

DSL query (2)

```
A = filter( c.color = 'Blue'  
          AND c.brand = 'Ford')
```

Is implemented using:

```
A = filter( c.color.labels = 'Blue' )  
B = filter( c.brand.labels = 'Ford' )  
C = kron( A, B )  
D = krao ( c.color, c.brand )  
E = dot ( C, D )
```



**Problem:** Khatri-Rao in (D) is always heavy (2 unfiltered matrices)

**Solution:** Next Slide



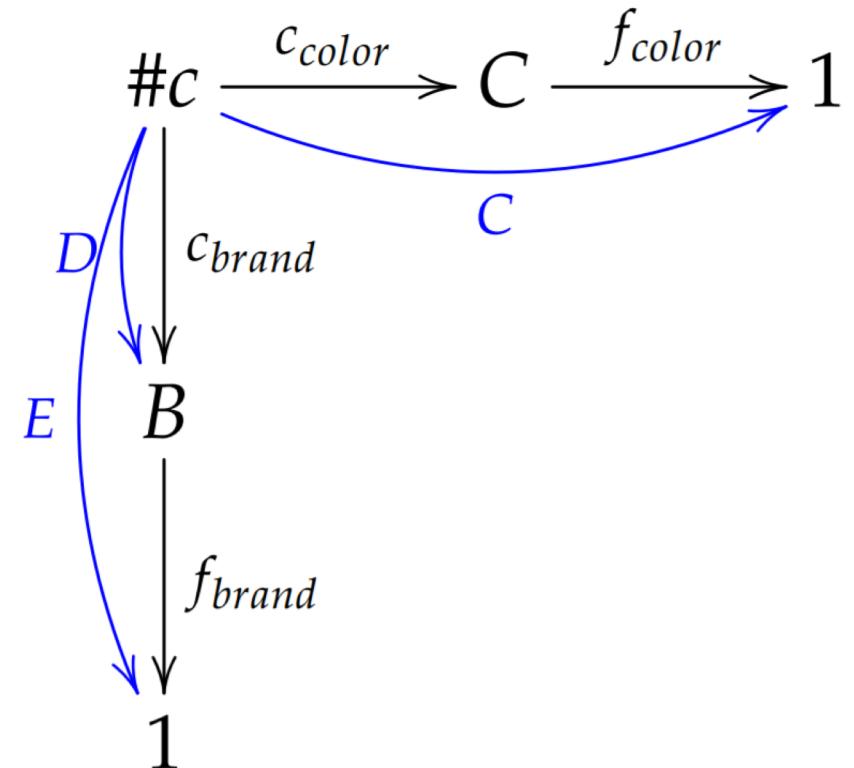
## Redundant filters (dimensions)

DSL query (2)

```
E = filter( c.color = 'Blue'  
          AND c.brand = 'Ford')
```

Is implemented using something like:

```
f_color = filter( c.color.labels = 'Blue' )  
f_brand = filter( c.brand.labels = 'Ford' )  
C = dot ( f_color, c.color )  
D = krao ( C, c.brand )  
E = dot ( f_brand, D )
```



**Restriction:** No filter can have attributes from distinct tables!

**Problem:** The SQL converter doesn't know which attributes are dimensions/measures

**Solution:** Expand the operation (as above) in DSL-runtime

**Result:** Highly complex filter operation, **must be studied** to check if is always solvable



## Redundant filters (dimensions)

Future optimizations :

- $\text{filter}_B = \text{predicate}_B . B$
- $\text{predicate}_A . ( A \setminus \text{filter}_B )$ 
  - Iterate the non-zeros of  $\text{filter}_B$  and get the column
  - Get the correspondent row in  $A$
  - Check its existence in the columns of  $\text{predicate}_A$
- 2 operations in a single run
  - Needs more loaded data! May not be faster...



## Optimize aggregations

First approach by BFG:

### Performance Results (2GB dataset)

The data below shows the execution time for each query step  
≈ 95% of the time is occupied by the last step

Execution Time: 33.830901

NNZ: 22767

0.152911	0.013123	0.597610	0.003269	0.083978	0.040801
0.121182	0.042822	0.090620	0.014628	0.009637	0.062796
0.049104	0.084366	0.111407	0.111351	32.228314	

[Previous performance](#)  
by Filipe

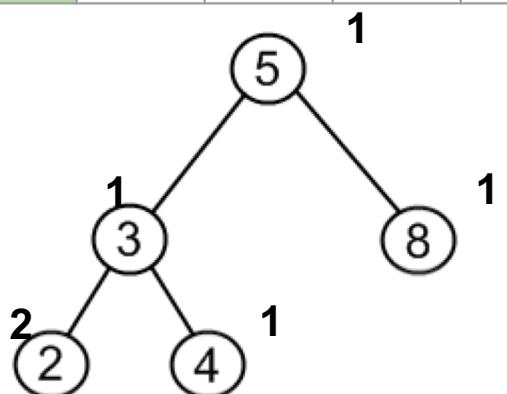
**Aggregations  
were the  
bottleneck!**

# Sparse Matrix-Vector Multiplication

```
for each non-zero in vector B :  
    col <- B_row_ind  
    if A_col_ptr[col+1] > A_col_ptr[col] :  
        temp[A_row_ind[A_col_ptr[col]]] +=  
B_values[curr_nnz]  
  
for each row in matrix A :  
    if temp[i] != 0.0 :  
        result_row_ind <- curr_row  
        result_values <- temp[curr_r  
nnz++  
  
col_ptr[0] = 0  
col_ptr[1] = nnz
```

Basically F.m<sup>o</sup>  
(using a complex  
algorithm...)

	#0	#1	#2	#3	#4	#5
1: A						
2: B		1		1		
3: C					1	
4: D				1		
5: D						1
6: E						
7: F						
8: G	1					



AVL avoids  
iterating rows

NNZ.log(NNZ)  
complexity!

### ◎ Query 3

- ~ 3 600 000 000 empty rows
- ~ 0 000 011 000 rows with data

### ◎ Problem

- `for (i=0; i<n_rows; i++) ...`
- Explained by the other group



## Optimize aggregations

**Assuming:** F is a (filtered) sparse Functional matrix  
m is an unfiltered measure vector

```
for i := 0 to f_nnz
    col = F_cols[i]
    row = F_rows[i]
    val = m_vals[col]
    res_vals[row] += val
```

Streamable since N is used transposed

$$M \cdot N = [ A \mid B ] \cdot \left[ \begin{matrix} C \\ D \end{matrix} \right] = A \cdot C + B \cdot D$$

**Doesn't work: result vector too long (use AVL instead)**  
Still avoids a khatri-Rao!

0	0	0	0	0
0	0	1	0	0
0	0	0	0	1

2	3	4	5	6
---	---	---	---	---

0	4	6
---	---	---



## Khatri-Rao/Hadamard optimization

- COO used less memory, but [discovered during paper writing]
- Sometimes iterating the non-zeros is **NOT** the best case
  - Remember that this approach is similar to the merge on the merge sort...
- Eg.:
  - $M :: i \leftarrow k, k \text{ non-zeros}$
  - $N :: j \leftarrow k, 1 \text{ non-zero}$
- Iterating M & N like in the merge, we have **~k+1 iterations**
- Solution:
  - Encode **M in CSC**: same memory for k non-zeros, slightly more for  $NNZ < k$  if  $NNZ \sim k$
  - Iterate the few (single) non-zeros of **N**
  - Directly add  $M_{row[col]} * M_{size} + N_{row[col]}$  to the result: **1 iteration**

## Goals for the meeting

15-03

- VLDB repository
- Coding
- DSL optimization



## VLDB repository

- Completed ([LINK](#))
  - Add in README: contact address
- Some parts verified
- Never tested from end to end
  - If it fails to run, at least serves as an installation and configuration guide



## DSL optimization

```
A = filter( l.shipdate >= "1994-01-01" AND l.shipdate < "1995-01-01" )
B = filter( l.discount >= 0.05 AND l.discount <= 0.07 )
C = hadamard( A, B )
D = filter( l.quantity < 24 )
E = lift( l.extendedprice * l.discount )
F = hadamard( D, E )
G = hadamard( C, F )
H = sum( G )
return ( H )
```

DSL provided by the LCC students:

1 attribute per filter operation

Optimized DSL (operation A still has an hidden dot product):

```
A = filter( l.shipdate >= "1994-01-01" AND l.shipdate < "1995-01-01" )
B = filter( A, l.discount >= 0.05 AND l.discount <= 0.07 AND l.quantity < 24 )
C = lift( B, l.extendedprice * l.discount )
D = sum( C )
return ( D )
```



## Coding

- Not all matrices can be encoded in COO
  - Reimplemented everything using CSC
  - COO may be explored later when fine tuning...
- Pbs using serialization of Google's *Protocol Buffers* (for disk accesses)
  - Must use the classes they generate, but not recommended...
  - Reimplemented using [Cereal](#) (faster, simpler, but with larger outputs)
- Current status:
  - Kernel implemented to run query 6
    - All possibilities of operators are easy to code, but development too long

## Goals for the meeting

05-04

- Query 6 streamed version
  - Profiling
  - GPU



## Query 6 profiling (serial streaming)

```
[laolap@compute-431-7 dbms]$ perf stat -e cpu-clock,page-faults,context-switches,cpu-migrations engine/bin/q6
```

```
Performance counter stats for 'engine/bin/q6':  
 3144.982266      cpu-clock (msec)  
 24203            page-faults  
 23778            context-switches  
      0            cpu-migrations  
 5.680358410 seconds time elapsed
```

Analysys developed as learned in  
*Engenharia de Sistemas da Computação*

**Perf** can introduce some overhead, but  
is a powerful tool to find hotspots (no  
changes to the source code are needed)

**TBC: Ongoing issues with perf on Search...**

NOTE: Consider compressing files on disk to make the  
engine less IO bound



## Queries on GPU



Jose, 9:30 AM

Viva - Como disse, não poderei estar na reunião de amanhã (nem por kstype) pois estou de viagem. Mas, aqui, surgiu na conversa (quando falei da nossa LAOLAP a colegas daqui) surgiu a ideia de correr os scripts em GPUs, uma ideia antiga; qual é a vossa perspectiva sobre isto, agora que se tem muitos mais experiência nestas coisas?



Possible reasons to reject this suggestion:

1. The key bottleneck is the data access in memory, and we are talking about BIG DATA; the GPU memory is not shared with the multicore memory, and data transfers between memories is through PCIe. If we manage to overcome this...
2. ...

João:

1. True, I'm working to complete query 6 and then make the profile of the streamed version. The GPU could make part of the pipeline if there are computation heavy LAQ(DSL) operators
2. Despite load time, COO operations cannot be vectorized and directly ported to GPU, but some operations over CSC can

## Goals for the meeting

12-04

19-04

- VLDB paper review
  - Incremental querying
  - Testbed environment
- DSL to C++ conversion
- Parallel streaming simple test



## VLDB paper review (1)

- Validate incremental querying
  - 1) Step by step incremental Q3 (paper & pen)
    - INSERT on **Lineitem**, **Orders** and **Customer**
  - 2) A possible controller must store:
    - The query result
    - DB Schema with counters
      - to get the delta
    - Query parsing tree (or some kind of key eg: MD5)
      - to compare with new queries
- Testbed environment: see testbed for dissertation?



## VLDB paper review (2)

- Comparison with columnar format (Hyper, ...)
- Compare SQL coverage in DB Toaster with ours
- Clearly show that other LA+OLAP works in the literature are not in similar areas
- What experimental results should be presented?
  - Data in memory vs. in disk (Sugg for the paper: only mem!)
  - Incremental querying



## Querying in C++ (Q6)

```
6 #include <iostream>
7 #include <vector>
8 #include "include/block.hpp"
9 #include "include/database.hpp"
10 #include "include/dot.hpp"
11 #include "include/filter.hpp"
12 #include "include/fold.hpp"
13 #include "include/krao.hpp"
14 #include "include/lift.hpp"
15 #include "include/matrix.hpp"
16 #include "include/types.hpp"
```

- 1 - Load precompiled libraries
- Always include **vector** and **iostream**
  - Always include **block**, **matrix**, **database** and **types**
  - Include the operators as needed
    - **Simplification:** include all headers

DONE

```

18 inline bool filter_a(std::vector<engine::Literal> args) {
19     return args[0] >= "1994-01-01" && args[0] < "1995-01-01";
20 }
21
22 inline bool filter_b(std::vector<engine::Decimal> args) {
23     return args[0] >= 0.05 && args[0] <= 0.07;
24 }
25
26 inline bool filter_c(std::vector<engine::Decimal> args) {
27     return args[0] < 24;
28 }
29
30 inline engine::Decimal lift_f(std::vector<engine::Decimal> args) {
31     return args[0] * args[1];
32 }

```

## 2 - Define the necessary predicates

**DONE**

- For each statement in the parsing tree:
  - If the operation is a **filter** or **lift**:
    - Check if the (first) attribute is a **measure** or **dimension**
    - Create the predicate using the parsed **expression**

```
34 ▼ int main() {  
35 ▼   engine::Database db(  
36       "data/la",  
37       "TPCH_1");  
--
```

### 3 - Select the database

DONE

- The **LAQ driver** receives **engine::Database** as argument
  - get() the arguments from the provided DB

```
39 ▼   engine::Bitmap *lineitem_shipdate =  
40 ▼     new engine::Bitmap(db.data_path,  
41           db.database_name,  
42           "lineitem",  
43           "shipdate");  
44 ▼   engine::DecimalVector *lineitem_discount =  
45 ▼     new engine::DecimalVector(db.data_path,  
46           db.database_name,  
47           "lineitem",  
48           "discount");  
49 ▼   engine::DecimalVector *lineitem_quantity =  
50 ▼     new engine::DecimalVector(db.data_path,  
51           db.database_name,  
52           "lineitem",  
53           "quantity");  
54 ▼   engine::DecimalVector *lineitem_extendedprice =  
55 ▼     new engine::DecimalVector(db.data_path,  
56           db.database_name,  
57           "lineitem",  
58           "extendedprice");
```

DONE

### 4 - Load the attributes metadata

- Iterate the list of **database vars**
  - Dimensions are Bitmaps
  - Measures are DecimalVectors

```

engine::FilteredBitVector *a_pred =
    new engine::FilteredBitVector(lineitem_shipdate->nLabelBlocks);
engine::FilteredBitVector *a =
    new engine::FilteredBitVector(lineitem_shipdate->nBlocks);
engine::FilteredBitVector *b =
    new engine::FilteredBitVector(lineitem_discount->nBlocks);
engine::FilteredBitVector *c =
    new engine::FilteredBitVector(lineitem_quantity->nBlocks);
engine::FilteredBitVector *d =
    new engine::FilteredBitVector(a->nBlocks);
engine::FilteredBitVector *e =
    new engine::FilteredBitVector(c->nBlocks);
engine::DecimalVector *f =
    new engine::DecimalVector(lineitem_extendedprice->nBlocks);
engine::FilteredDecimalVector *g =
    new engine::FilteredDecimalVector(e->nBlocks);
engine::Decimal *h =
    new engine::Decimal();

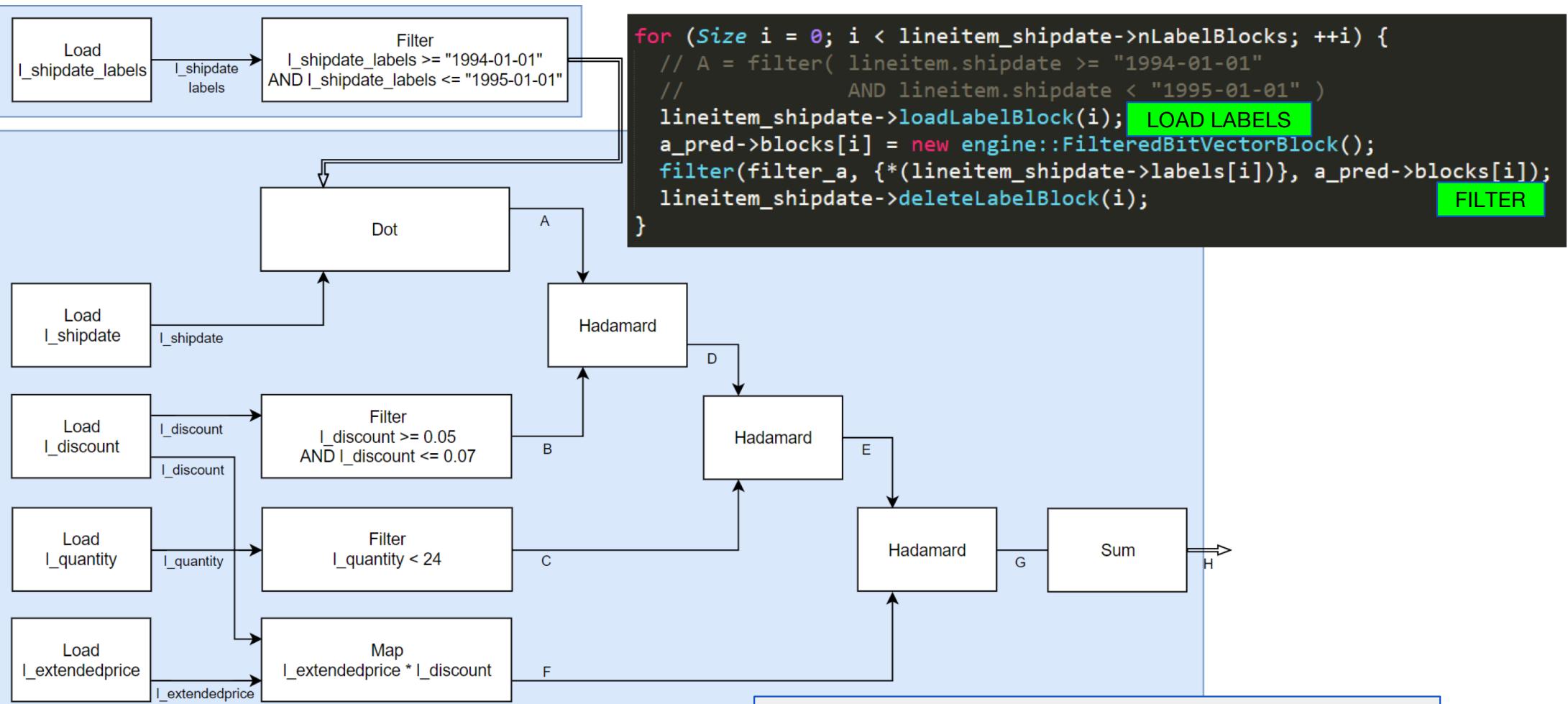
```

## 5 - Declare **DONE** temporary matrices

- Define a map that dictates the type of the output considering the operation and the input type(s)

**INCOMPLETE**

- Recursively declare all nodes of the operation tree
- The argument is the number of blocks of the first argument in the operation



## 6 - Build the for loops

- 2 Blue boxes = 2 for loops
- Inside each loop all operations follow any sequence that respects dependencies



## 6 - some considerations

- The value in the for loop stop condition
  - Is the number of blocks of the first variable in the first operation
- Attributes may be reused, thus we need a list like this:

Variable	I_discount	I_shipdate	...	A	B	...
Type	DecimalVector	Bitmap	...	FilteredBitVector	FilteredBitVector	...
Total Uses	2	1	...	1	1	...
Actual Uses	<runtime>	<runtime>	...	<runtime>	<runtime>	...

- So if **Actual Uses = 0** and **Variable ∈ database attributes**
  - A variable block must be loaded
- If **Actual Uses = Total Uses** a variable block must be deleted



## 6 - some considerations

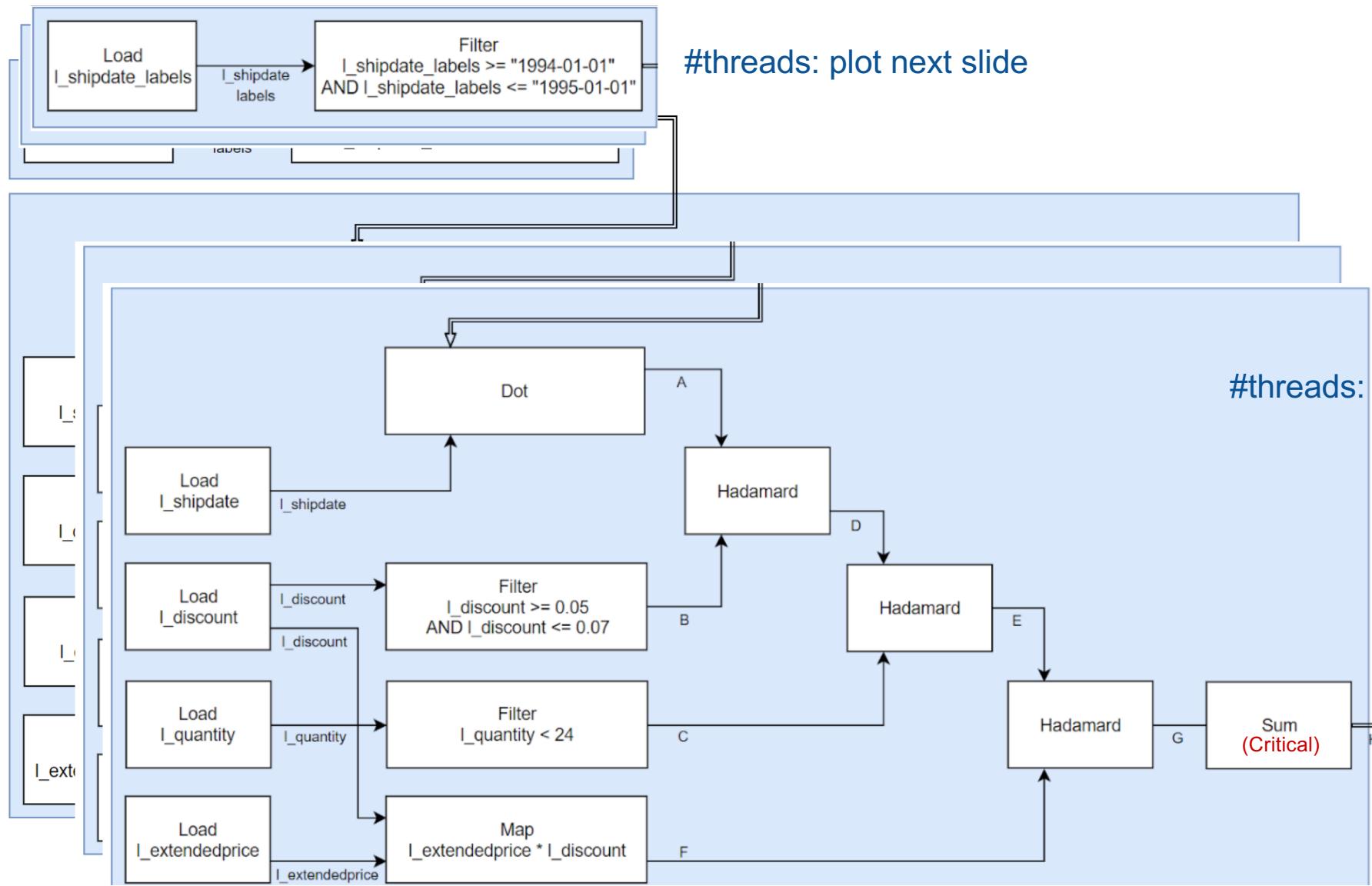
- When the operation is a **Dot**, **Filter on dimensions** or **Fold**
  - The previous loop must be closed and appended to a list
  - Some variables may have to be deleted (**strong arrows**)
    - Delete blocks inside streams and full matrices between for loops
  - A new loop must be created
- Special cases (unsolvable with left recursion -> solve in parent node)
  - Dot (orders\_orderdate, lineitem\_orderkey)
    - lineitem\_orderkey loaded in stream
    - orders\_orderdate loaded in a separate cycle (strong arrow)
  - Filter on dimensions
    - Must be replaced by a filter on labels and a dot product

```
139     std::cout << (*h) << std::endl;  
140  
141     return 0;  
142 }
```

## 7 - Print the result

- **No algorithm (yet) for multiple outputs**
- **The final classes must implement <<**
  - In this case: double has << implemented by default

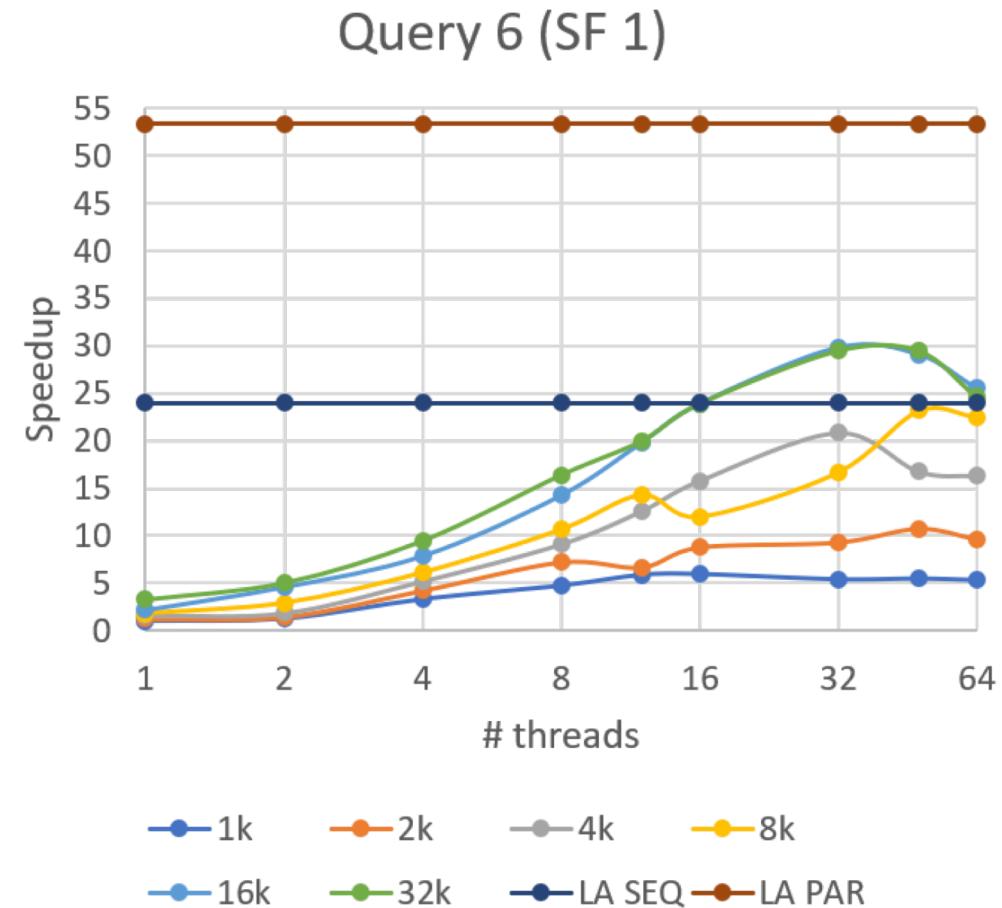
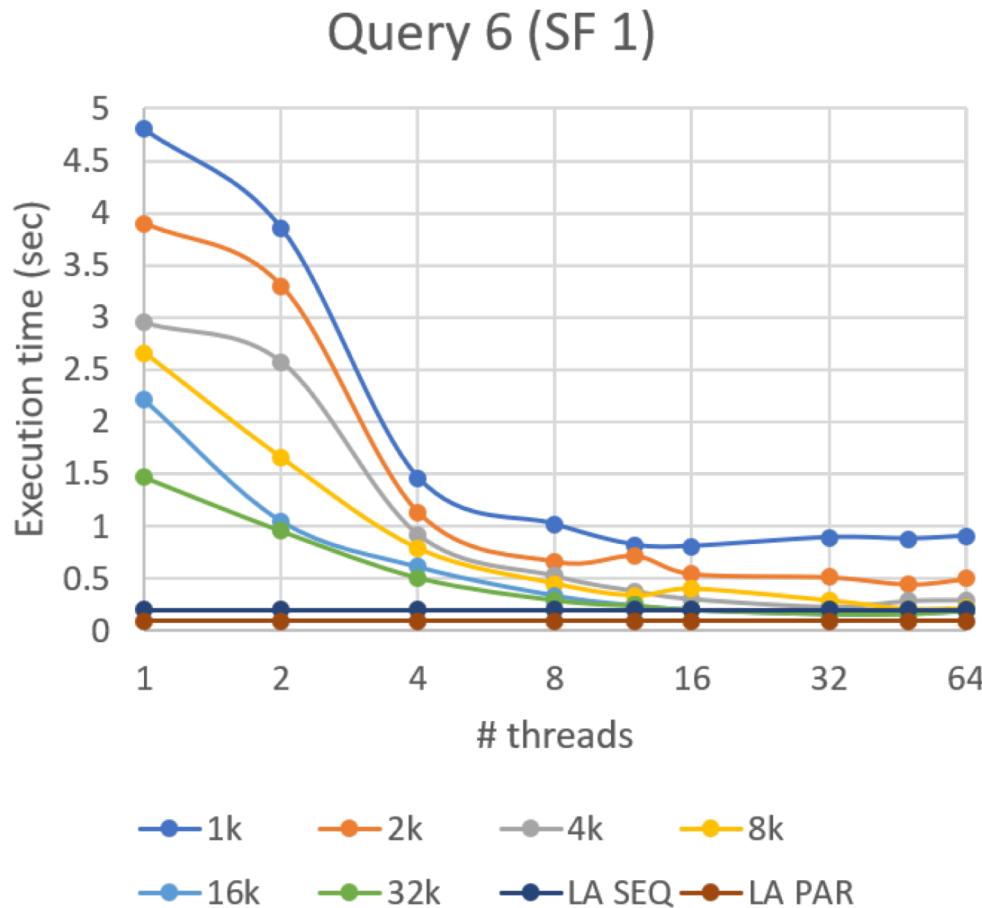
## 8 - Close main function **DONE**





## Parallel Streaming

- Testing block sizes (#elements); relative to worst case (longest, 1k\_elem)



## Goals for the meeting

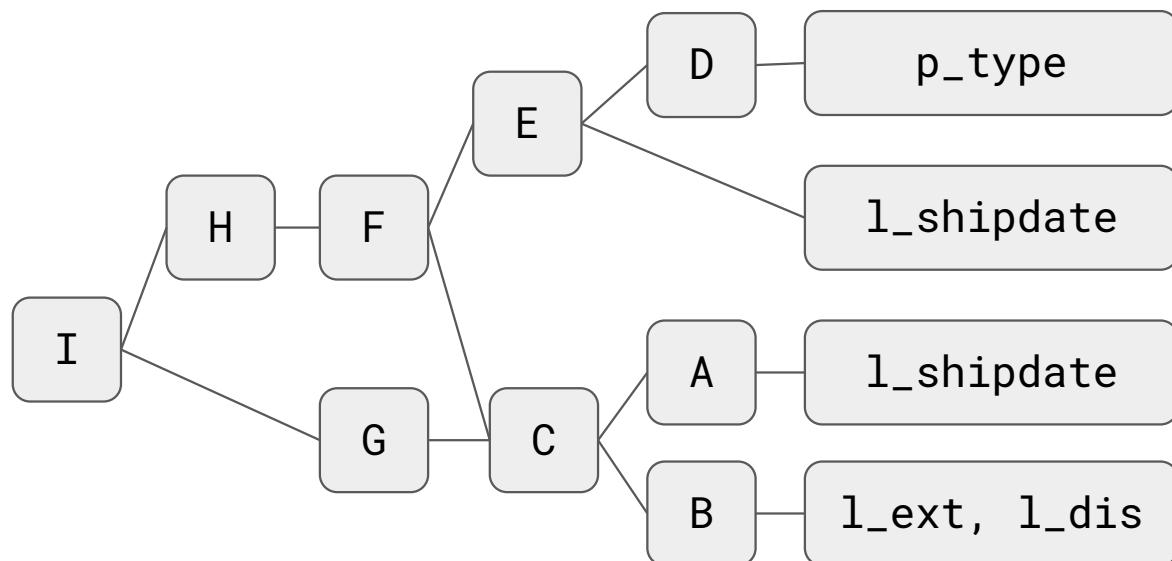
10-05

- DSL to C++ conversion (also see 12 April)
- Dissertation topics
- Define work until July



## Q14 LAQ

```
5 A = filter( lineitem.shipdate >= "1995-09-01" AND lineitem.shipdate < "1995-10-01") // 1 <-- #l
6 B = lift( lineitem.extendedprice * (1 - lineitem.discount) ) // 1 <-- #l
7 C = hadamard( A, B ) // 1 <-- #l
8 D = filter( match( part.type , ".+PROMO.+" ) ) // 1 <-- #p
9 E = dot( D, lineitem.partkey ) // 1 <-- #l
10 F = hadamard( C, E ) // 1 <-- #l
11 G = sum( C ) // 1 <-- 1
12 H = sum( F ) // 1 <-- 1
13 I = lift( 100.00 * H / G ) // 1 <-- 1
14 return ( I ) // RESULT
```



## Goals for the meeting

23-05

- Paper deadlines
- Time planning
  - Paper
  - Dissertation



## Paper submission info (PODS'19)

Submitted papers should be at most twelve pages.

Additional details may be included in an appendix. These will be read at the submission time (online appendices are not allowed). Appendices will be read at the discretion of the program committee.

FIRST SUBMISSION CYCLE:

**June 15, 2018:** Abstract submission

**June 22, 2018:** Paper submission

**August 31, 2018:** First notification

**September 28, 2018:** Revised submission

**November 2, 2018:** Final notification

SECOND SUBMISSION CYCLE:

**December 14, 2018:** Abstract submission

**December 21, 2018:** Paper submission

**March 8, 2019:** Final notification

All deadlines end at 11:59pm AoE.

**Another submission without  
appendix page limit that has the  
appendix before bibliography**

## Goals for the meeting

30-05

- Incremental Query 3 script
- Column store engines
  - MonetDB - install on Search



# Incremental Q3 script

```
δA = dot( filter( [c_mktseg.labels | δc_mktseg.labels] ), δc_mktseg )
δB = dot( [A | δA], δo_custkey )
δC = dot( filter( [o_orderdate.labels | δo_orderdate.labels] ), δo_orderdate )
δD = kraq( δB, δC )
δE = kraq( δD, δo_orderdate )
δF = kraq( δE, δo_shippriority )
δG = dot( [F | δF], δl_orderkey )
δH = dot( filter( [l_shipdate.labels | δl_shipdate.labels] ), δl_shipdate )
δI = kraq( δG, δH )
δJ = kraq( δl_orderkey, δI )
δK = lift( δl_extendedprice * (1 - δl_discount) )
δL = kraq( δJ, δK )
δM = sum( δL )
M' = add(M, δM)
return( M' )
```



## MonetDB - install on Search

### ○ MonetDB ./configure

- Error on OpenSSL. **Too old version?**

```
checking for openssl... no
configure: error: OpenSSL library not found but required for MonetDB5
[laolap@compute-662-5 build]$ openssl version -v -b -o
OpenSSL 1.0.0-fips 29 Mar 2010
built on: Thu Aug 23 04:57:41 UTC 2012
options: bn(64,64) md2(int) rc4(16x,int) des(idx,cisc,16,int) blowfish(idx)
```

When a component is disabled, it is usually due to a missing dependency. A dependency is often a library that the component in question relies on for its functionality. The following two dependencies are known trouble makers on most systems that aren't fairly recent: `openssl`, `libxml2` and `pcre` or `libpcre`. The former is usually available, but a not recent enough version. The latter two often lack the development headers on systems. You have to install missing dependencies either through the native package management system of your host system, or by compiling and installing them manually. Since this in general differs from package to package, we cannot provide generalized help on this topic. Make sure you install the `xxx-dev` versions on binary distributions of the necessary packages as well!



## MonetDB - install on Search

- OpenSSL ./configure

- Error on Perl modules

```
Can't locate Pod/Simple.pm in @INC (@INC contains: . /usr/local/lib64/perl5 /usr/lib64/perl5 /usr/share/perl5 .) at /usr/share/perl5/Pod/Text.pm line 33.  
BEGIN failed--compilation aborted at /usr/share/perl5/Pod/Text.pm line 33.  
Compilation failed in require at /usr/share/perl5/Pod/Usage.pm line 452.  
BEGIN failed--compilation aborted at /usr/share/perl5/Pod/Usage.pm line 459.  
Compilation failed in require at configdata.pm line 17370.  
BEGIN failed--compilation aborted at configdata.pm line 17370.  
Compilation failed in require.  
BEGIN failed--compilation aborted.
```



## MonetDB - install on Search

- Install Perl 5.26.2 from source – **OK**
- Install OpenSSL 1.1 from source- **OK**
- Update system paths – **OK**
- Install MonetDB from source – **OpenSSL not found**
- `.../configure --with-openssl=no ...` – **OK**
  - Does it affect any important feature?

```
[laolap@compute-781-2 monetdb_opt]$ bin/mclient -d laolap  
server requires unknown hash 'SHA512'
```

# End of meetings' slides



$$(M \otimes N) \cdot (P \nabla Q) = (M \cdot P) \nabla (N \cdot Q)$$

JNO: found this title while reviewing a paper,

Golub, G., Van Loan, C.: *Matrix Computations*. Johns Hopkins University Press, fourth edn. (2013),

Cf. state-of-the-art etc if not yet there.

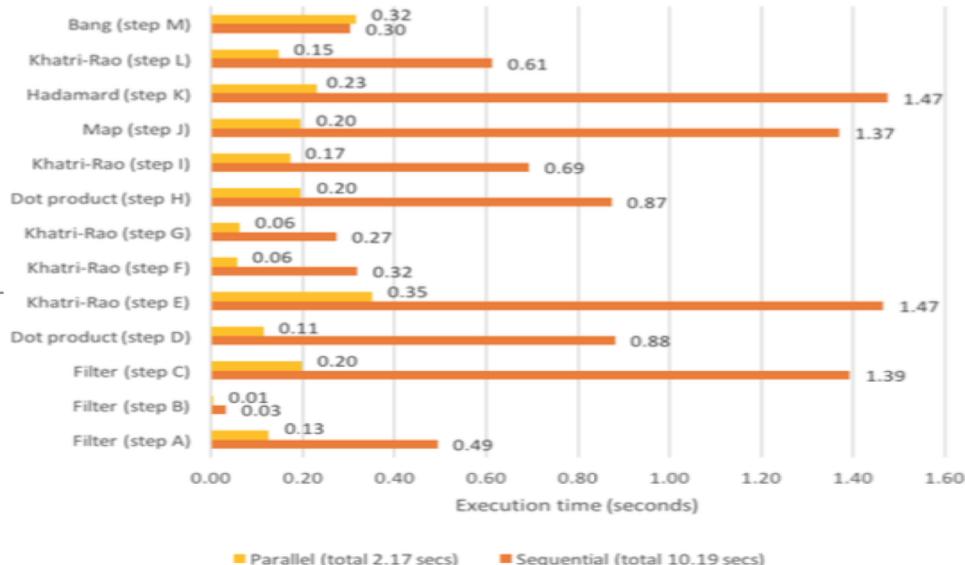
# TPC-H

## Query 3

```

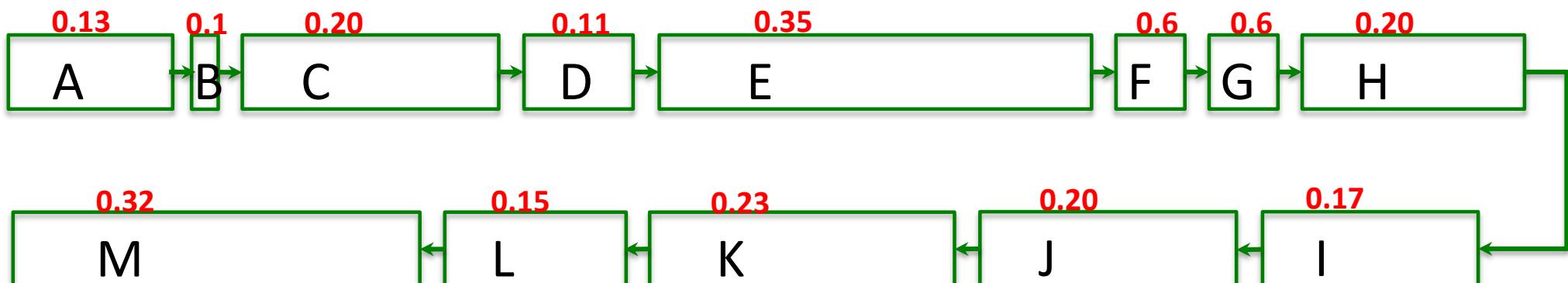
A = filter_to_matrix( o_orderdate, <'1995-03-10' )
B = filter_to_vector( c_mktsegment, ='MACHINERY' )
C = filter_to_vector( l_shipdate, >'1995-03-10' )
D = dot( B, o_custkey )
E = krao( l_orderkey, C )
F = krao( A, D )
G = krao( F, o_shippriority )
H = dot( G, l_orderkey )
I = krao( E, H )
J = map( \x -> 1 - x, l_discount )
K = hadamard( l_extendedprice, J )
L = krao( I, K )
M = bang( sum, L )

```

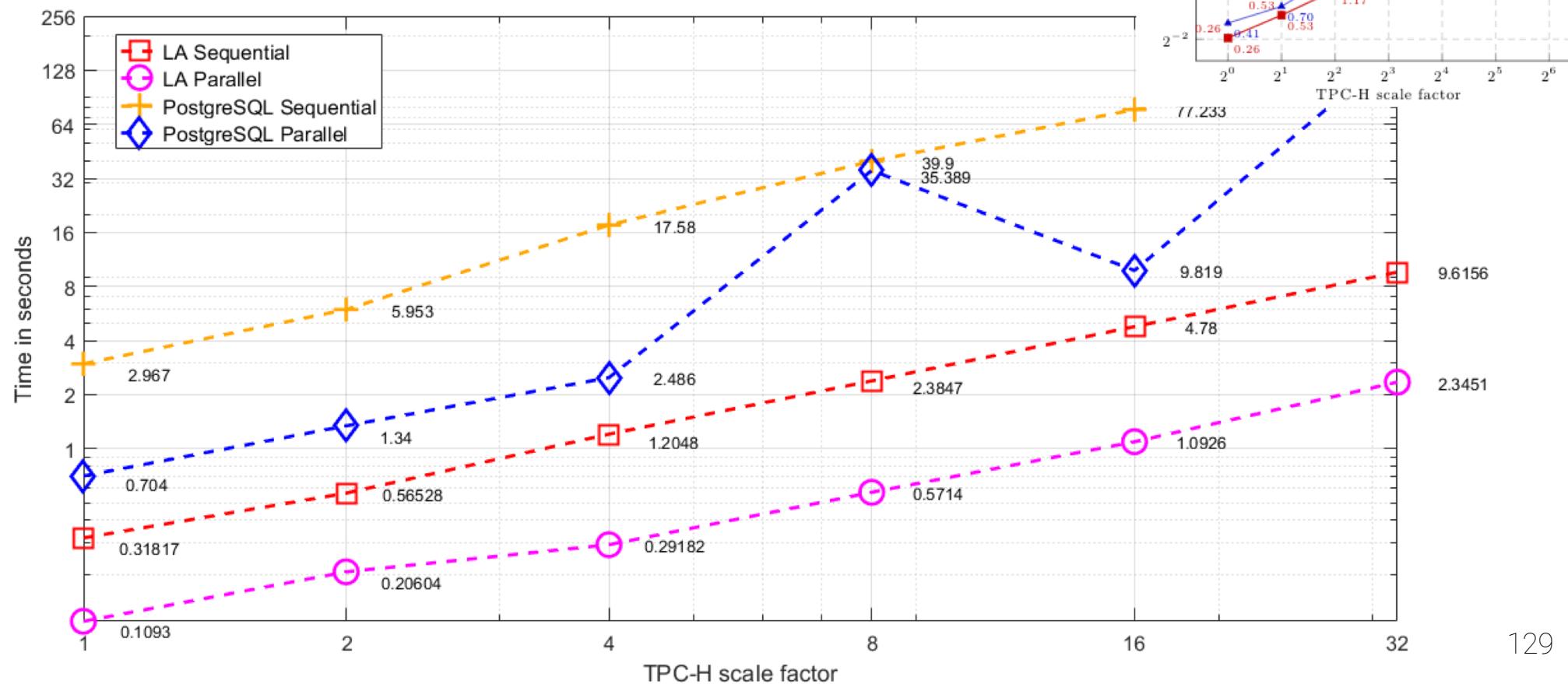


### Proposed parallel implementation

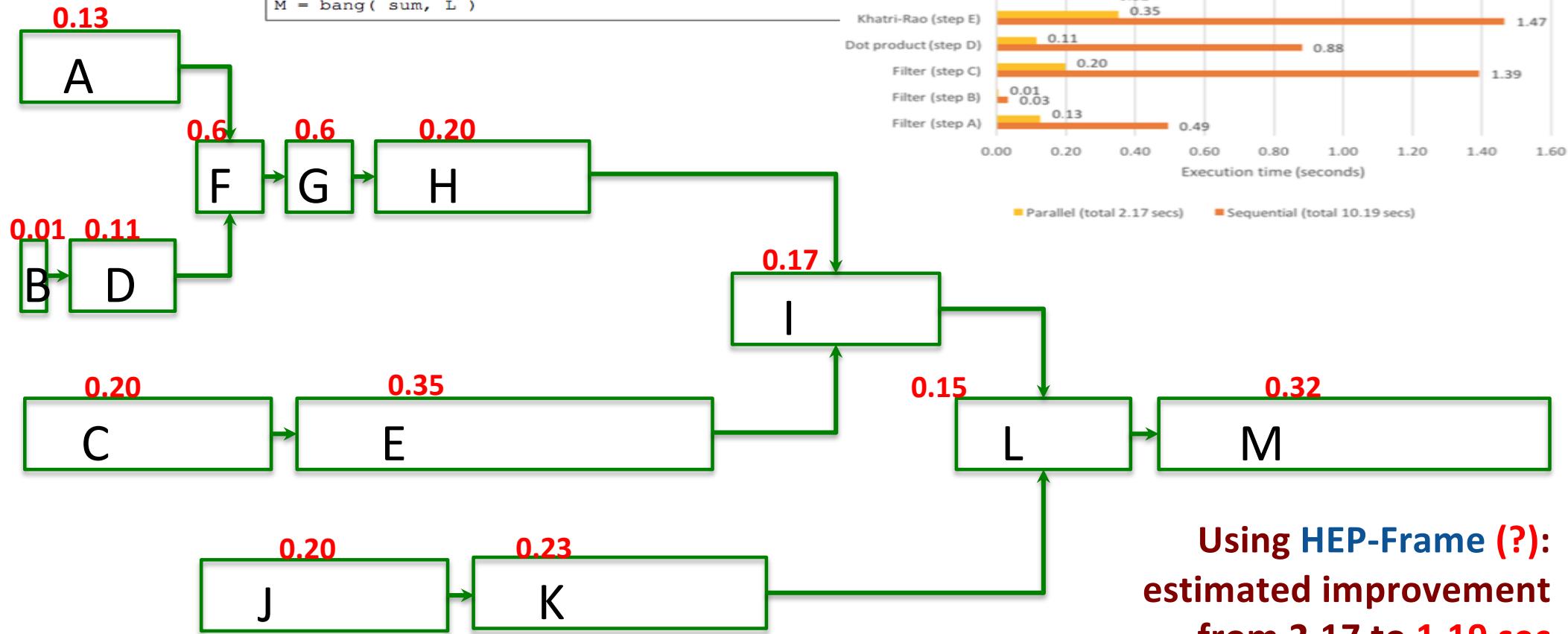
Total = **2.17 sec**



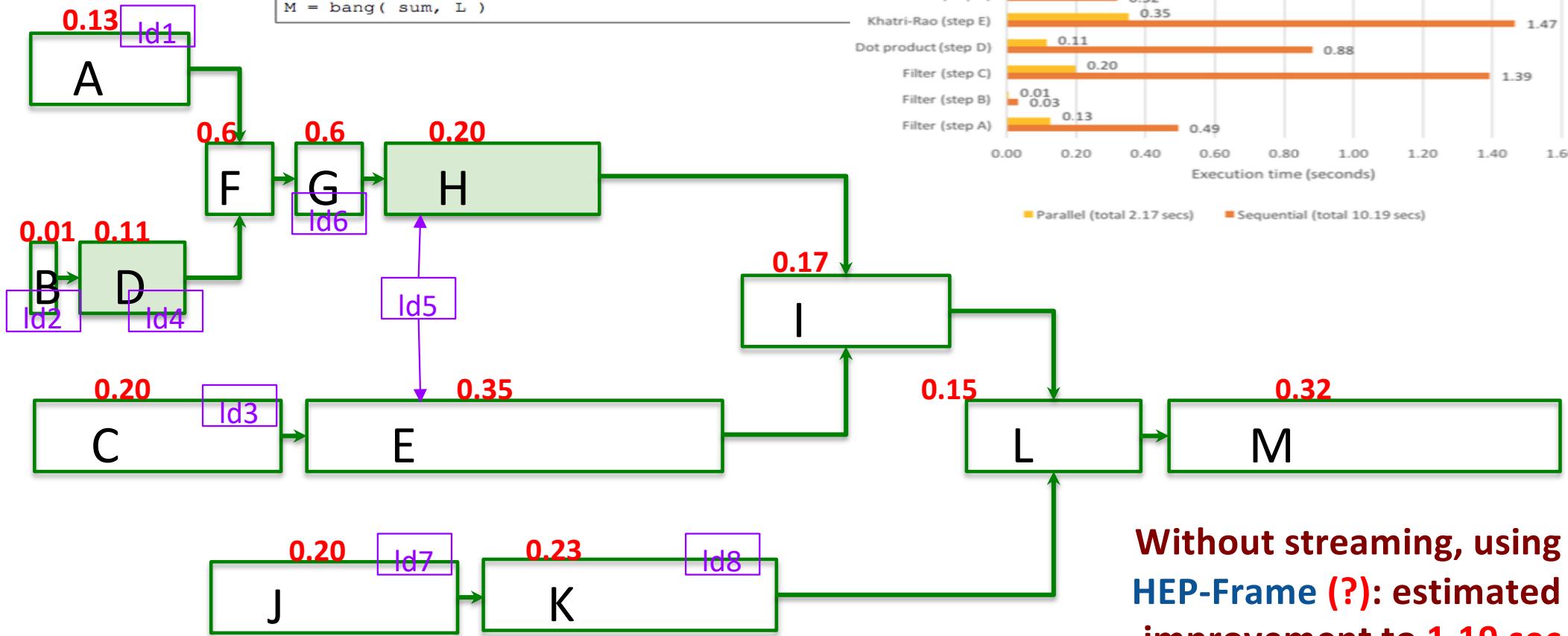
# Query 3 performance result



# TPC-H Query 3



# TPC-H Query 3



# TPC-H

## Query 3

Streaming approach

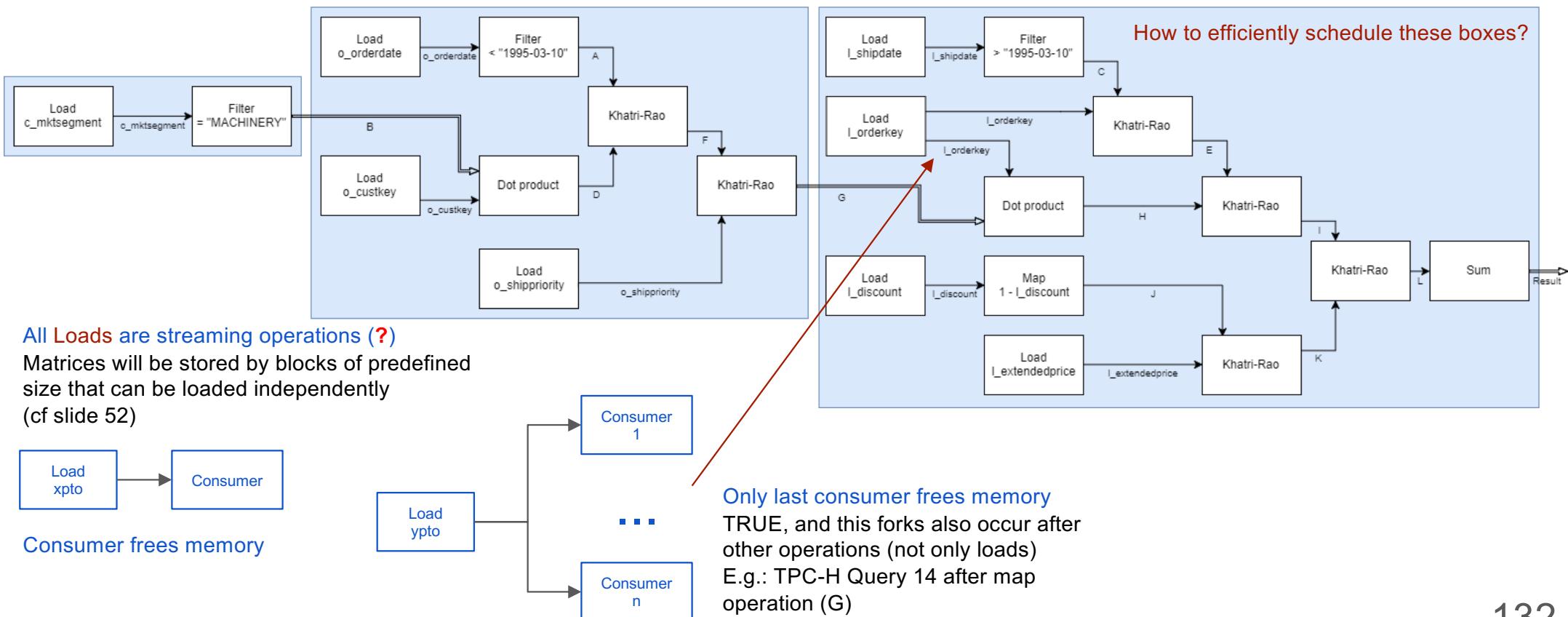
```

A = filter_to_matrix( o_orderdate, <'1995-03-10' )
B = filter_to_vector( c_mktsegment, ='MACHINERY' )
C = filter_to_vector( l_shipdate, >'1995-03-10' )
D = dot( B, o_custkey )
E = krao( l_orderkey, C )
F = krao( A, D )
G = krao( F, o_shippriority )
H = dot( G, l_orderkey )
I = krao( E, H )
J = map( \x -> 1 - x, l_discount )
K = hadamard( l_extendedprice, J )
L = krao( I, K )
M = bang( sum, L )

```

## Vectorized image

→ Full matrix data  
→ Streaming data



# TPC-H

## Query 6

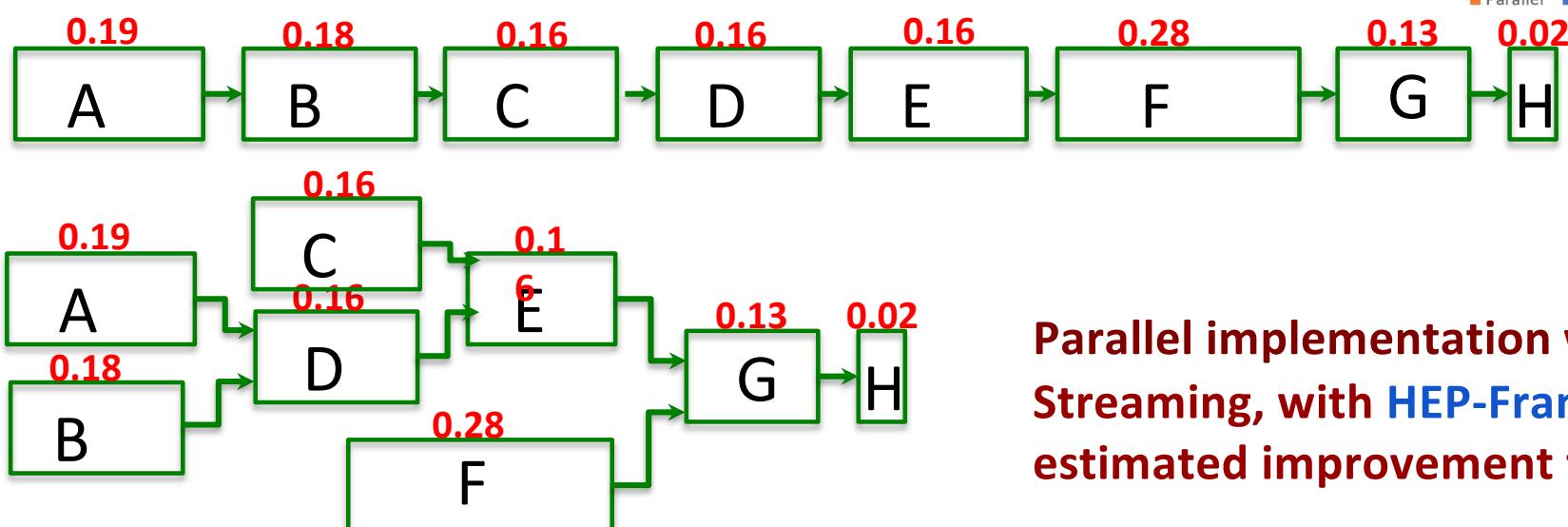
```

SELECT      sum(l_extendedprice * l_discount) as revenue
FROM        lineitem
WHERE       l_shipdate >= date '1995-03-10'
            AND l_shipdate < date '1995-03-10' + interval '1' year
            AND l_discount between 0.05 - 0.01 AND 0.05 + 0.01
            AND l_quantity < 3;

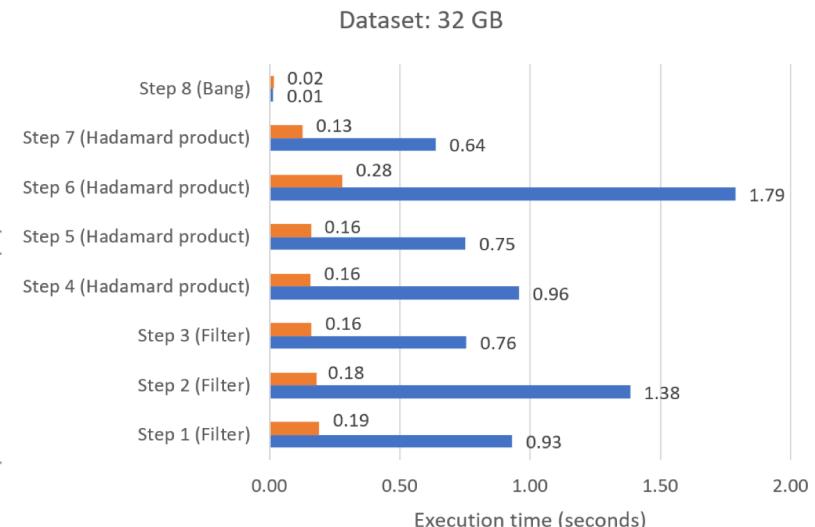
A = filter_to_vector( l_shipdate, >='1995-03-10', <'1995-03-10' )
B = filter_to_vector( l_discount, >='0.49', <='0.51' )
C = filter_to_vector( l_quantity, <='3' )
D = hadamard( A , B )
E = hadamard( C , D )
F = hadamard( l_extendedprice , l_discount )
G = hadamard( E , F )
H = bang( sum, G )

```

Parallel implementation w/out HEP-Frame => 1.28 sec



Query 6 Profiling Step-by-Step



Parallel implementation without Streaming, with HEP-Frame, estimated improvement to 0.66 sec

# TPC-H

## Query 6

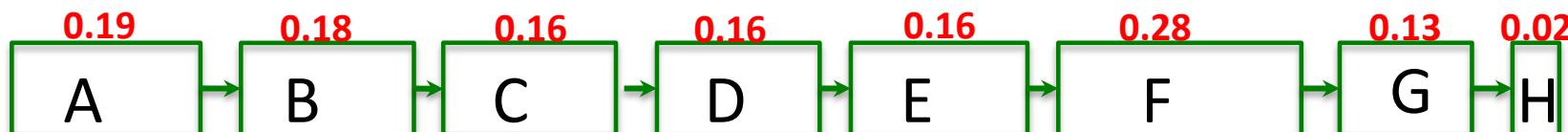
```

SELECT
    sum(l_extendedprice * l_discount) as revenue
FROM
    lineitem
WHERE
    l_shipdate >= date '1995-03-10'
    and l_shipdate < date '1995-03-10' + interval '1' year
    and l_discount between 0.05 - 0.01 and 0.05 + 0.01
    and l_quantity < 3;
  
```

```

A = filter_to_vector( l_shipdate, >='1995-03-10', <'1995-03-10')
B = filter_to_vector( l_discount, >='0.49', <='0.51' )
C = filter_to_vector( l_quantity, <='3' )
D = hadamard( A , B )
E = hadamard( C , D )
F = hadamard( l_extendedprice , l_discount )
G = hadamard( E , F )
H = bang( sum, G )
  
```

**Sequential implementation for scale  $2^5 \Rightarrow 6.21$  sec**



Cluster node 781.1

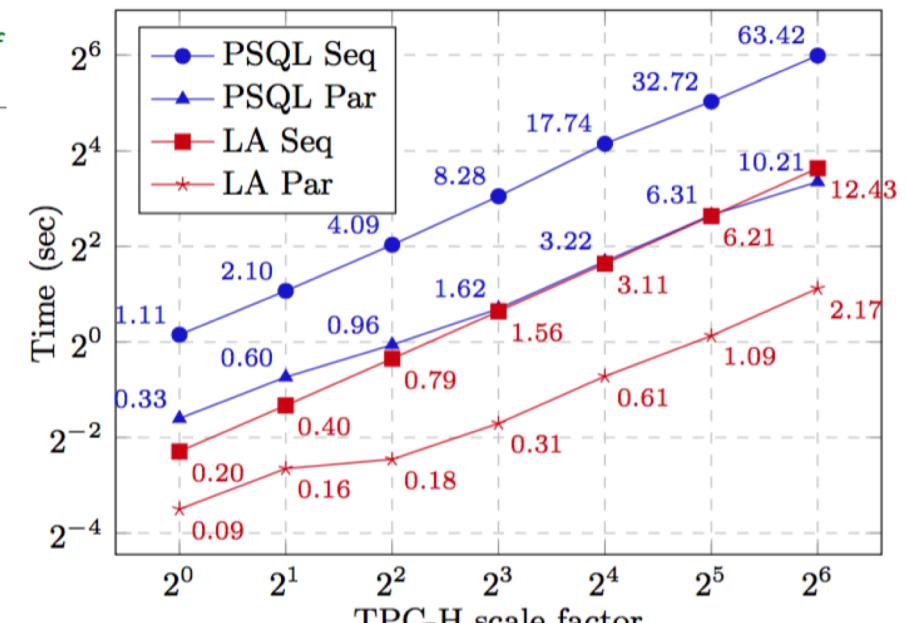
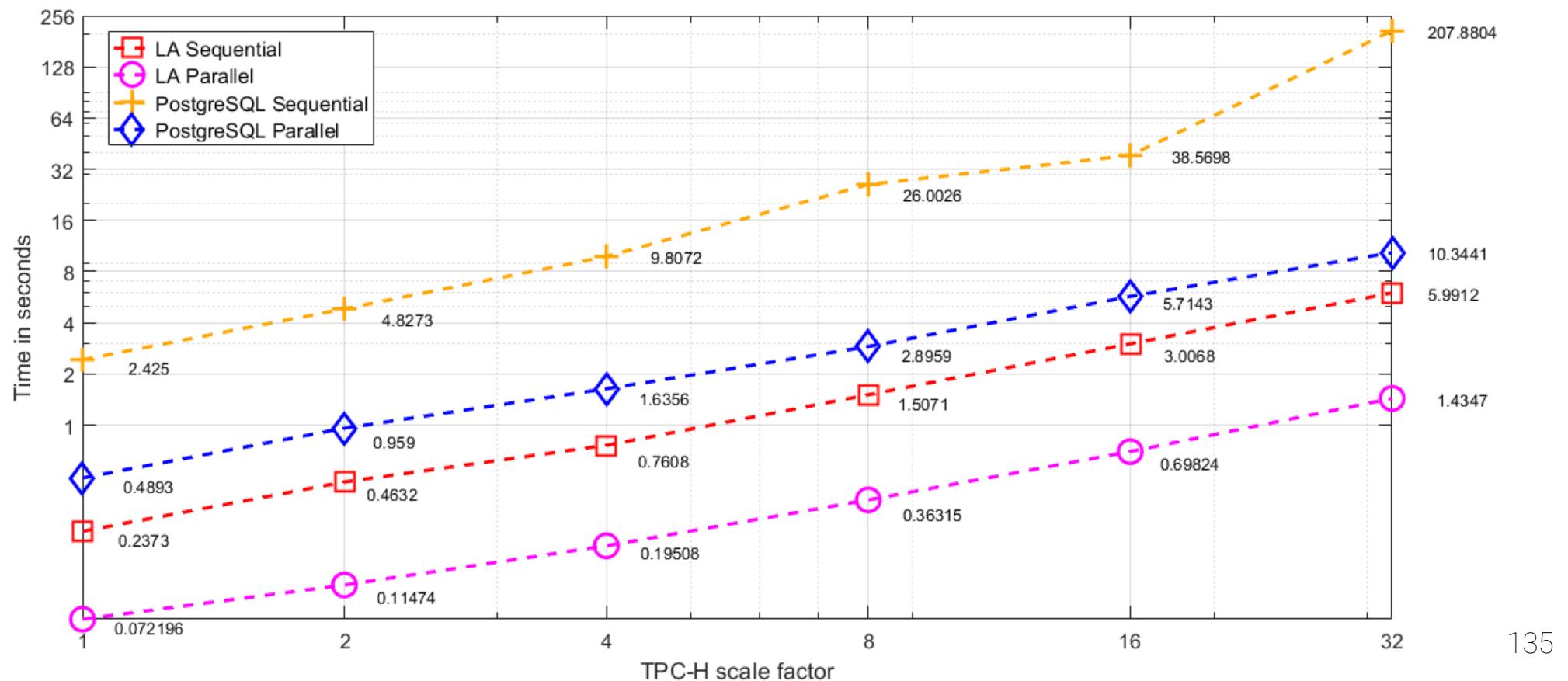


Figure 7: Sequential and Parallel Query 6

Needs to be updated...

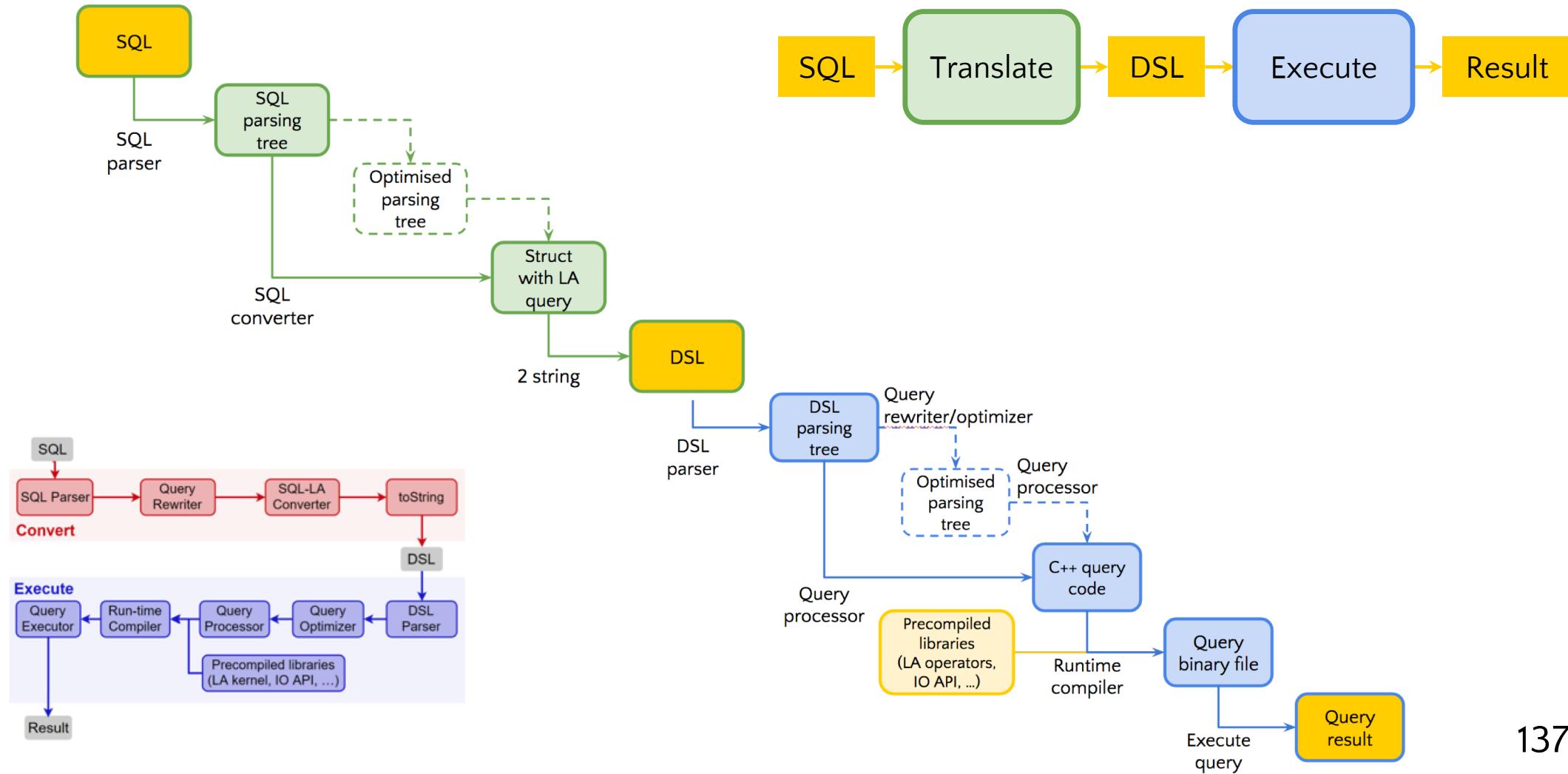
# Query 6 performance results







## Possible architecture (initial proposal)





# RPD: Contents

## CONTENTS

<b>1 INTRODUCTION</b>	<b>1</b>	
1.1 Challenges & Goals	2	
1.2 Dissertation Outline	2	
<b>2 THE RELATIONAL MODEL AND OLAP</b>	<b>3</b>	
2.1 Data Storage Origins	3	
2.1.1 Data Processing Machines	4	
2.1.2 The First Databases	4	
2.2 The Relational Model	5	
2.2.1 Relational Algebra	6	
2.2.2 SQL	8	
2.3 Database System Implementation	10	
2.3.1 Query Compilation	10	
2.3.2 Query Execution	13	
2.4 Data Warehousing	15	
2.4.1 Data Modeling	15	
<b>3 TYPED LINEAR ALGEBRA FOR OLAP</b>	<b>17</b>	
3.1 Data Representation	17	
3.1.1 Dense Vectors	17	
3.1.2 Sparse Matrices	18	
3.2 Type Diagrams	18	
3.3 Linear Algebra Operations	19	
3.3.1 Dot Product	20	
3.3.2 Khatri-Rao Product	21	
3.3.3 Hadamard-Schur Product	21	
3.3.4 Attribute Filter	22	
3.3.5 Matrix Aggregation	23	
3.3.6 Elementwise Map Operator	23	
3.4 Conversion Example	24	
<b>4 A TLA-DB ENGINE FOR RELATIONAL SQL</b>	<b>29</b>	
4.1 SQL Driver	29	
4.1.1 SQL Parser	29	
4.1.2 Query Rewriter	30	
4.1.3 SQL Converter	31	
4.1.4 DSL toString	31	
4.2 DSL Engine	31	
4.2.1 DSL Parser	31	
4.2.2 Query Optimizer	31	
4.2.3 Query Processor	32	
4.2.4 Run-time Compiler	32	
4.2.5 Query Execution	32	
4.3 Engine Libraries	32	
4.3.1 Matrix Representation	32	
4.3.2 IO API	34	
4.3.3 LA Operators	34	
<b>5 VALIDATION AND PERFORMANCE RESULTS</b>	<b>35</b>	
5.1 TPC Benchmark H	35	
5.2 Preliminary Results	37	



# Dissertation: Contents

1. Introduction
  - a. Challenges & Goals
  - b. Dissertation Outline
2. The relational model and OLAP
  - a. **Data Storage Origins**
  - b. The Relational Model
  - c. Database System Implementation
  - d. **Data Warehousing**
  - e. **OLAP**
    - i. **ROLAP** (current framework)
    - ii. **MOLAP** (future work)
  - f. **Distributed Databases**  
(to introduce distributed LA)
1. Typed linear algebra for OLAP
  - a. Linear Algebraic Encoding of Data
  - b. Type Diagrams
  - c. Linear Algebra Query language
    - i. Linear Algebra Operations
  - d. From SQL to LAQ
  - e. **Cubing in Linear Algebra**
  - f. **LA OLAP operations**  
(dicing, slicing, roll up, drill down, ...)
4. A TLA DB engine for relational OLAP
  - a. Matrix representation
  - b. LAQ operations
  - c. “Streaming” approach
  - d. Controller
  - e. SQL Driver
  - f. LAQ Engine
5. Validation and performance results
  - a. TPC Benchmark H
  - b. In-memory
6. **Framework extension**
  - a. **Incremental querying**  
(introduce this in the state of the art?)
  - a. **Distributed querying**  
(already studied distributed LAQ operations)
    - a. **Future work**
4. Conclusion

**Sections to be deleted**  
**Possible (new) sections**  
**Comments**



## Dissertation topics - FUTURE WORK

---

- Refer data cubes theory in the state of the art ?
- Better thread-level scheduler? (HEP-Frame?)
- Better controller
  - When flush data to disk? ACID vs performance
  - Multi-user/query system
    - Study what should be cached
    - Study data share among queries executed in parallel
    - Study resource allocation
  - JDBC & console APIs
  - Multi-node system
    - Replicate / distribute data
    - Reduce communications on querying & minimize data transfers