

# Do the middle letters of “OLAP” stand for Linear Algebra (“LA”)?

Speaker: Luís A. Bastião Silva

Paper authors: Hugo Daniel Macedo and José Nuno Oliveira

Doctoral Program



# Summary

- + Motivation
- + Goals
- + Background
- + Cross tabulations in LA
- + Higher-dimensional OLAP
- + Conclusion and future work

# Motivation

- Nowadays, companies are creating a huge amount of data
- Big data trend
- They need to access to the information stored in these databases and calculate some metrics
- OLAP (Online Analytical Processing):
  - Summarize huge amount of information
  - Forms of histograms, sub-totals, cross tabulations, roll-up/drill down, data cubes
- Expensive task (computationally)

# Motivation

- Perform data mining and online analytical processing (OLAP) in a efficient way
- OLAP is :
  - Resource-demanding
  - Calls for parallelization
- OLAP operations:
  - Pivot
  - Roll-up
  - Cube

# Related work

- Ng. et al develop a collection of parallel algorithms to data cube construction in low cost PCs (Clustering)
- PARSIMONY: provides a parallel and scalable infrastructure for multidimensional analyses
- There are commercial solutions like Oracle and IBM that also implement their parallel algorithms
- This paper propose a new direction: OLAP and data mining should rely on Linear Algebra

# Cross tabulation

- Provides a **summary** of a data extracted from raw source
- Example:

<b>Model</b>	<b>Year</b>	<b>Color</b>	<b>Sales</b>
Chevy	1990	Red	5
Chevy	1990	Blue	87
Ford	1990	Green	64
Ford	1990	Blue	99
Ford	1991	Red	8
Ford	1991	Blue	7

- How many vehicles sold per colour and model?

# Cross tabulation

- How many vehicles sold per colour and model?
- Selected Color and Model as attributes and Sales as a measure
- Answer is:

Sum of Sales	Model		
Color	Chevy	Ford	Grand Total
Blue	87	106	193
Green		64	64
Red	5	8	13
Grand Total	92	178	270

**In this paper: solve this problem with Linear Algebra.  
But how we can parallelize?**

# OLAP - Cube

- Cross tabulation summaries:
  - Computationally expensive
  - Long time (large datasets)
- OLAP cube compute all dimensions
- Calculate all possible options
- Summarize the table
  - Works like a cache of values
  - Easy to compute and access data in time

	<i>Sales</i>
<i>Chevy 1990 Blue</i>	87
<i>Chevy 1990 Red</i>	5
<i>Ford 1990 Blue</i>	99
<i>Ford 1990 Green</i>	64
<i>Ford 1991 Blue</i>	7
<i>Ford 1991 Red</i>	8
<i>Chevy 1990 ALL</i>	92
<i>Ford 1990 ALL</i>	163
<i>Ford 1991 ALL</i>	15
<i>Chevy ALL Blue</i>	87
<i>Chevy ALL Red</i>	5
<i>Ford ALL Blue</i>	106
<i>Ford ALL Green</i>	64
<i>Ford ALL Red</i>	8
<i>ALL 1990 Blue</i>	186
<i>ALL 1990 Green</i>	64
<i>ALL 1990 Red</i>	5
<i>ALL 1991 Blue</i>	7
<i>ALL 1991 Red</i>	8
<i>Chevy ALL ALL</i>	92
<i>Ford ALL ALL</i>	178
<i>ALL 1990 ALL</i>	255
<i>ALL 1991 ALL</i>	15
<i>ALL ALL Blue</i>	193
<i>ALL ALL Green</i>	64
<i>ALL ALL Red</i>	13
<i>ALL ALL ALL</i>	270



# Cross tabulation – Linear Algebra

- Three matrices:
  - Two associated with dimensions (attributes) – A and B
  - Measure or Metric
- Divide-and-conquer principle, with matrix multiplication:

$$[R|S] \cdot \begin{bmatrix} U \\ V \end{bmatrix} = R \cdot U + S \cdot V$$

- OLAP cross-tabulation can be expressed by:  $t_A \cdot [[T]]_M \cdot t_B^\circ$
- A, B is dimensions and M is the measure

$$t_A(x, r) = \begin{cases} 1 & \text{if } T(A, r) = x \\ 0 & \text{otherwise} \end{cases}$$

# Cross tabulation – Linear Algebra

	1	2	3	4	5	6
<i>Chevy</i>	1	1	0	0	0	0
<i>Ford</i>	0	0	1	1	1	1

$$|Model| \xleftarrow{t_{Model}} n$$

	1	2	3	4	5	6
<i>Blue</i>	0	1	0	1	0	1
<i>Green</i>	0	0	1	0	0	0
<i>Red</i>	1	0	0	0	1	0

$$|Color| \xleftarrow{t_{Color}} n$$

$$t_{Color} \cdot t_{Model}^{\circ} = \begin{array}{c|cc} & \textit{Chevy} & \textit{Ford} \\ \hline \textit{Blue} & 1 & 2 \\ \textit{Green} & 0 & 1 \\ \textit{Red} & 1 & 1 \end{array}$$

# Cross tabulation – Linear Algebra

$$t_{Color} \cdot [T]_{Sales} \cdot t_{Model}^{\circ} = \begin{array}{r|cc} & \text{Chevy} & \text{Ford} \\ \hline \text{Blue} & 87 & 106 \\ \text{Green} & 0 & 64 \\ \text{Red} & 5 & 8 \end{array}$$

$$ctab_{Color, Model; Sales}(T) = \begin{array}{r|ccc} & \text{Chevy} & \text{Ford} & \text{ALL} \\ \hline \text{Blue} & 87 & 106 & 193 \\ \text{Green} & 0 & 64 & 64 \\ \text{Red} & 5 & 8 & 13 \\ \text{ALL} & 92 & 178 & 270 \end{array}$$

# Rolling-up on functional dependences

- Rolling-up means replacing a dimension by another which is more general in some sense (eg. grouping, classification, containment).

<b>Model</b>	<b>Year</b>	<b>Color</b>	<b>Sales</b>	<b>Month</b>	<b>Season</b>
Chevy	1990	Red	5	March	Spring
Chevy	1990	Blue	87	April	Spring
Ford	1990	Green	64	August	Summer
Ford	1990	Blue	99	October	Autumn
Ford	1991	Red	8	January	Winter
Ford	1991	Blue	7	January	Winter

- Also works for checking functional dependences

# Rolling-up on functional dependences

- Rolling-up means replacing a dimension by another which is more general in some sense (eg. grouping, classification, containment).

$$t_{Season} \cdot t_{Month}^{\circ} =$$

	<i>January</i>	<i>March</i>	<i>April</i>	<i>August</i>	<i>October</i>
<i>Spring</i>	0	1	1	0	0
<i>Summer</i>	0	0	0	1	0
<i>Autumn</i>	0	0	0	0	1
<i>Winter</i>	2	0	0	0	0

- Also works for checking functional dependences

# Rolling-up on functional dependences

- Rolling-up means replacing a dimension by another which is more general in some sense (eg. grouping, classification, containment).

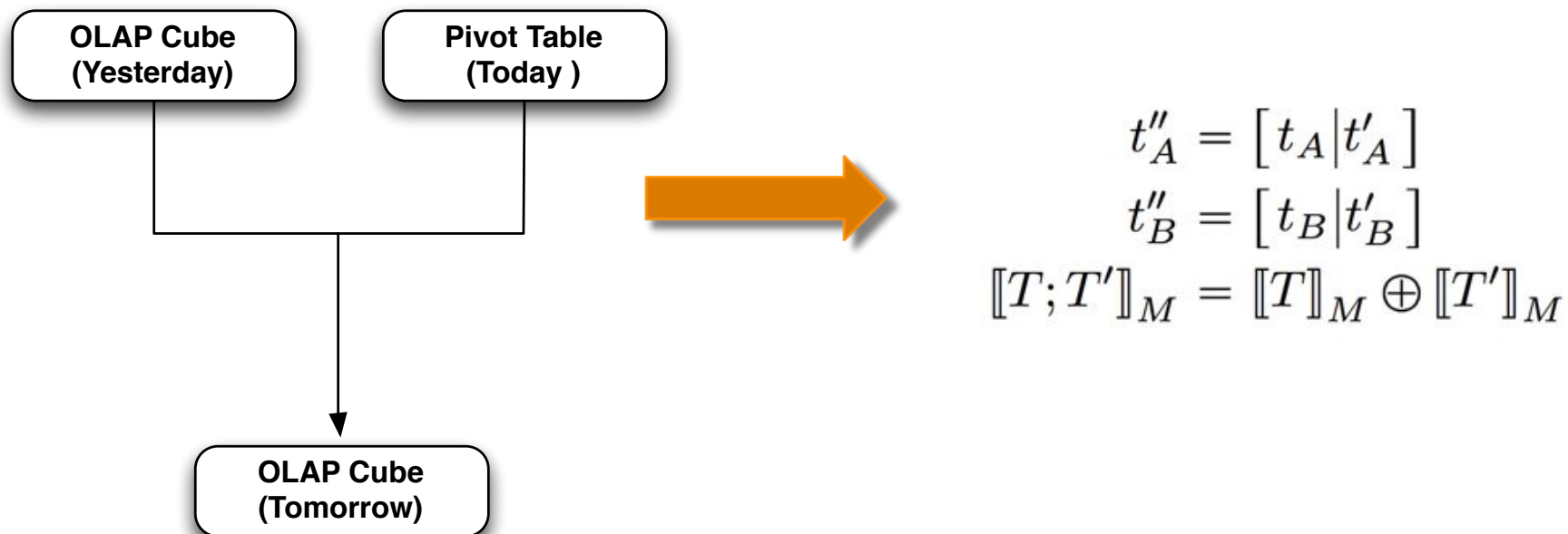
$$(t_{Season \leftarrow Month} \oplus id) \cdot ctab_{Month, Model; Sales}(T) =$$

	<i>Chevy</i>	<i>Ford</i>	<i>ALL</i>
<i>Spring</i>	92	0	92
<i>Summer</i>	0	64	64
<i>Autumn</i>	0	99	99
<i>Winter</i>	0	15	15
<i>ALL</i>	92	178	270

- Also works for checking functional dependences

# Incremental construction

- Cross tabulations defined by Linear Algebra is amenable to incremental constructions



- Advantage: is not necessary to build all the CUBE every single day!

# Higher dimensionality - OLAP

- + Consider n-dimensions: aggregate, group-by, cross tabulations and cube
- + Generalization based on Khatri-Rao product
  - + Works like a Cartesian product
- + Khatri-Rao product:

$$s = [5 \ 87 \ 64 \ 99 \ 8 \ 7]$$

$$6 \xleftarrow{s \odot id} 6 = \begin{bmatrix} 5 & 0 & 0 & 0 & 0 & 0 \\ 0 & 87 & 0 & 0 & 0 & 0 \\ 0 & 0 & 64 & 0 & 0 & 0 \\ 0 & 0 & 0 & 99 & 0 & 0 \\ 0 & 0 & 0 & 0 & 8 & 0 \\ 0 & 0 & 0 & 0 & 0 & 7 \end{bmatrix}$$



# Higher-dimensional OLAP

- All dimensions
- Whole dimension part
- Raw-data table
- The Khatri-Roa of:
  - tModel and tColor

	1	2	3	4	5	6
<i>Chevy</i>	1	1	0	0	0	0
<i>Ford</i>	0	0	1	1	1	1

$$|Model| \xleftarrow{t_{Model}} n$$

	1	2	3	4	5	6
<i>Blue</i>	0	1	0	1	0	1
<i>Green</i>	0	0	1	0	0	0
<i>Red</i>	1	0	0	0	1	0

$$|Color| \xleftarrow{t_{Color}} n$$

		1	2	3	4	5	6
<i>Chevy</i>	<i>Blue</i>	0	1	0	0	0	0
<i>Chevy</i>	<i>Green</i>	0	0	0	0	0	0
<i>Chevy</i>	<i>Red</i>	1	0	0	0	0	0
<i>Ford</i>	<i>Blue</i>	0	0	0	1	0	1
<i>Ford</i>	<i>Green</i>	0	0	1	0	0	0
<i>Ford</i>	<i>Red</i>	0	0	0	0	1	0

# Higher-dimensional OLAP

- All dimensions
- Whole dimension part
- Raw-data table

$$t_{Model \times Year \times Color}$$

$$= t_{Model} \odot t_{Year} \odot t_{Color}$$

$$=$$

			1	2	3	4	5	6
<i>Chevy</i>	1990	<i>Blue</i>	0	1	0	0	0	0
<i>Chevy</i>	1990	<i>Green</i>	0	0	0	0	0	0
<i>Chevy</i>	1990	<i>Red</i>	1	0	0	0	0	0
<i>Chevy</i>	1991	<i>Blue</i>	0	0	0	0	0	0
<i>Chevy</i>	1991	<i>Green</i>	0	0	0	0	0	0
<i>Chevy</i>	1991	<i>Red</i>	0	0	0	0	0	0
<i>Ford</i>	1990	<i>Blue</i>	0	0	0	1	0	0
<i>Ford</i>	1990	<i>Green</i>	0	0	1	0	0	0
<i>Ford</i>	1990	<i>Red</i>	0	0	0	0	0	0
<i>Ford</i>	1991	<i>Blue</i>	0	0	0	0	0	1
<i>Ford</i>	1991	<i>Green</i>	0	0	0	0	0	0
<i>Ford</i>	1991	<i>Red</i>	0	0	0	0	1	0

# Higher-dimensional OLAP

- All dimensions
- Whole dimension part
- Raw-data table

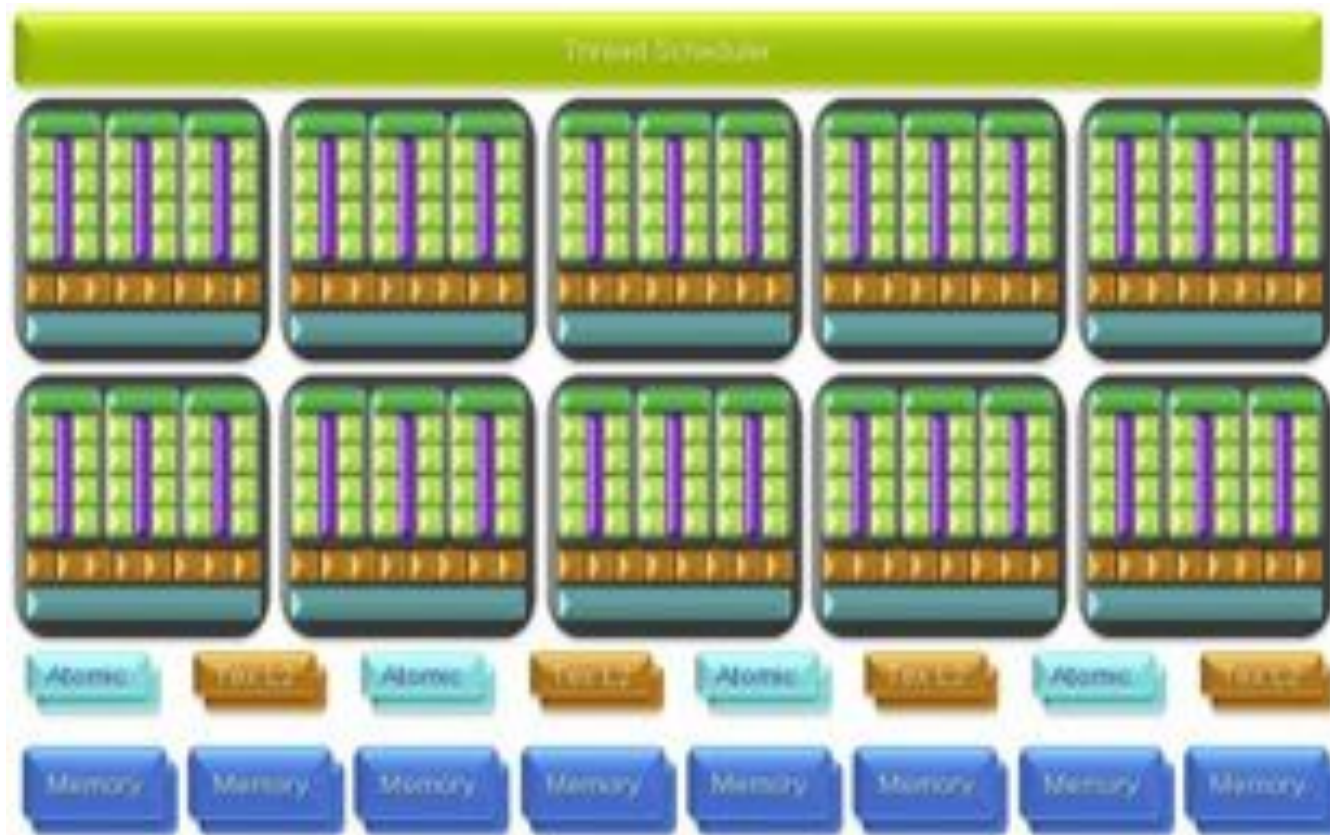
	<i>Sales</i>
<i>Chevy 1990 Blue</i>	87
<i>Chevy 1990 Green</i>	0
<i>Chevy 1990 Red</i>	5
<i>Chevy 1991 Blue</i>	0
<i>Chevy 1991 Green</i>	0
<i>Chevy 1991 Red</i>	0
<i>Ford 1990 Blue</i>	99
<i>Ford 1990 Green</i>	64
<i>Ford 1990 Red</i>	0
<i>Ford 1991 Blue</i>	7
<i>Ford 1991 Green</i>	0
<i>Ford 1991 Red</i>	8

$t_{Model \times Year \times Color} \cdot [T]_{Sales} \cdot !^{\circ} =$

# Conclusion and future work

- OLAP computationally problematic
- Parallelization is already possible, but not with linear algebra
- Encoding OLAP in concepts of Linear Algebra – formal method
- Rely on theory of parallel sparse matrix/matrix multiplication to achieve parallelism
- Cross tabulation is incremental
- Future:
  - Extending LA for other OLAP features
  - Implement in Multi-core and GPU and replace the OpenOffice/ LibreOffice pivot table calculator

# Future work (GPGPU)



Questions?