



Um pouco sobre o quarto paradigma da  
ciência

# Agenda

- Conceitos
  - Quatro paradigmas da ciência
  - SGBDs relacionais
  - Map-Reduce & Hadoop
  - Sistemas analíticos
  - NoSQL
- SciDB
  - Definição de dados e comparações
  - Linguagem de consultas e comparações
  - Arquitetura
- Hands-on
  - Como usar?
  - Alguns exemplos práticos

# Conceitos

...

# Quatro paradigmas da ciência

- Milhares de anos atrás: **EMPÍRICA**
  - Descrição dos fenômenos naturais
- Últimos séculos: **TEÓRICA**
  - Modelos, generalizações etc.
- Últimas décadas: **COMPUTACIONAL**
  - Simulações de fenômenos complexos
- Hoje: **EXPLORAÇÃO DE DADOS (eScience)**
  - Unifica teoria, experimento e simulação

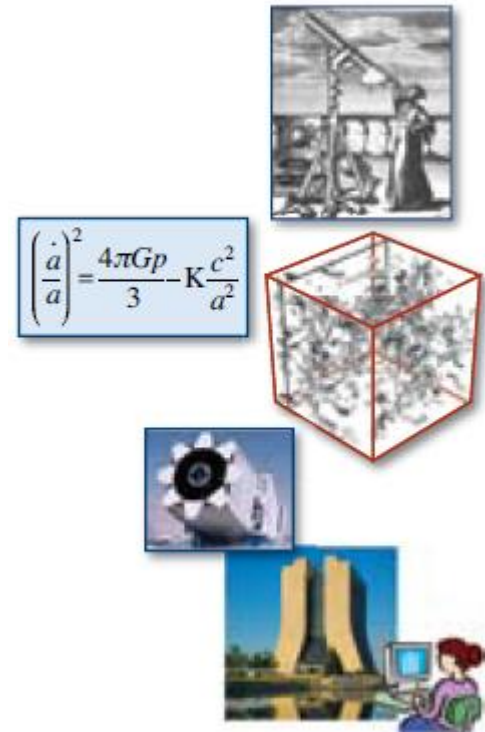
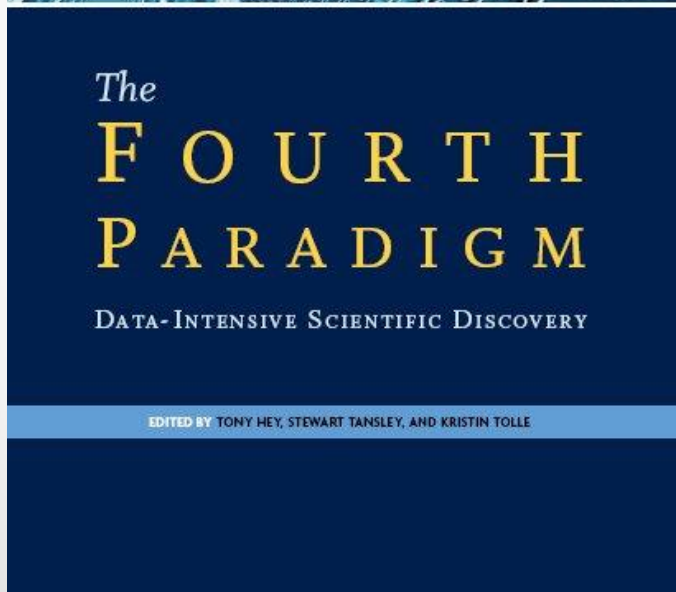


Imagem creditada a [2]

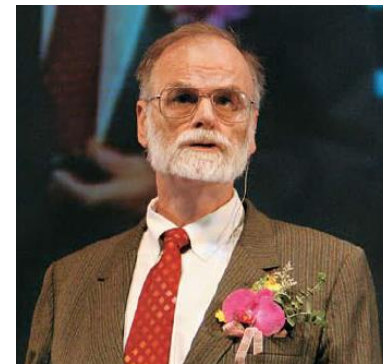
# The Fourth Paradigm: Data-Intensive Scientific Discovery



Publicação da Microsoft Research, 2009.

Disponível na íntegra online (CC):  
<http://research.microsoft.com/en-us/collaboration/fourthparadigm/>

Idealização de Jim Gray (Turing, 1998).



Imagens creditadas a [2]

# Como lidar com os dados?

- Todos os dados científicos online?!

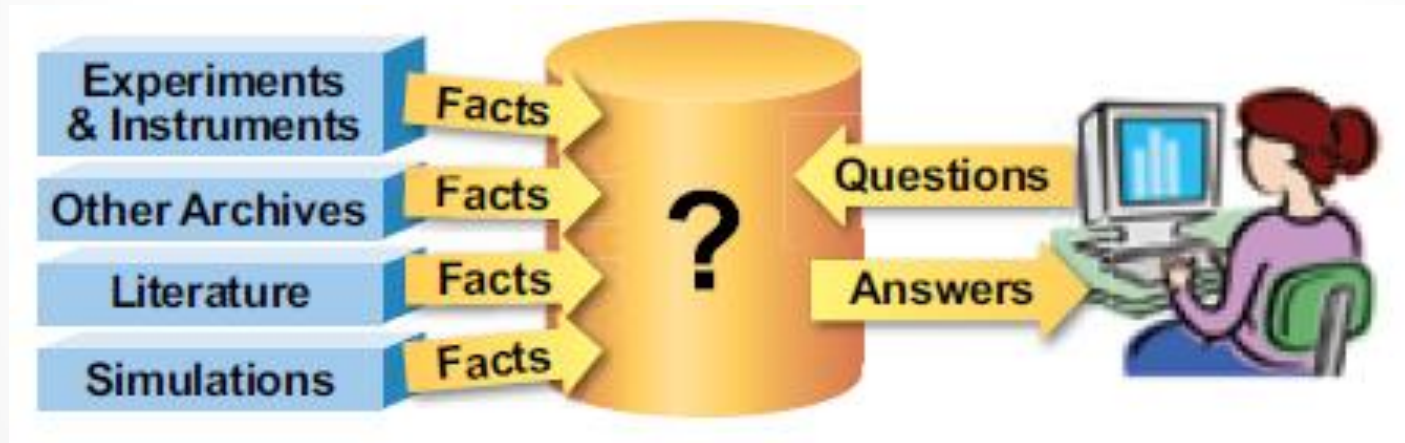
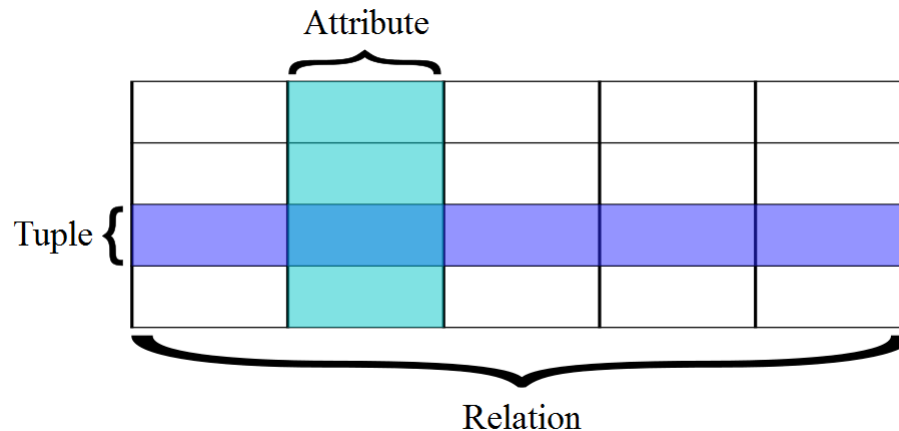


Imagem creditada a [2]

# SGBDs relacionais

	Route No.	Miles	Activity
Record 1	I-95	12	Overlay
Record 2	I-495	05	Patching
Record 3	SR-301	33	Crack seal

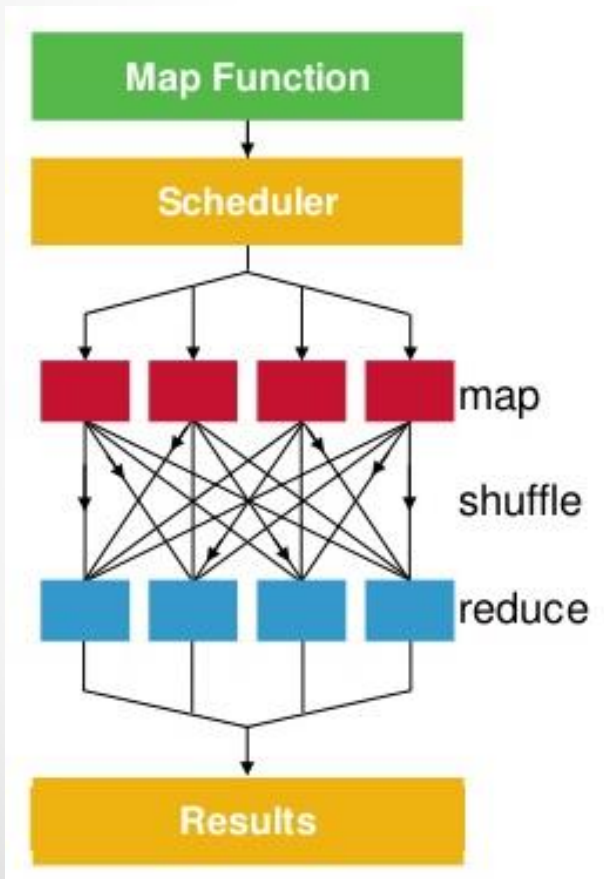
*Flat files* → redundância



Edgar Codd (IBM) → 12 regras

# Map-Reduce & Hadoop

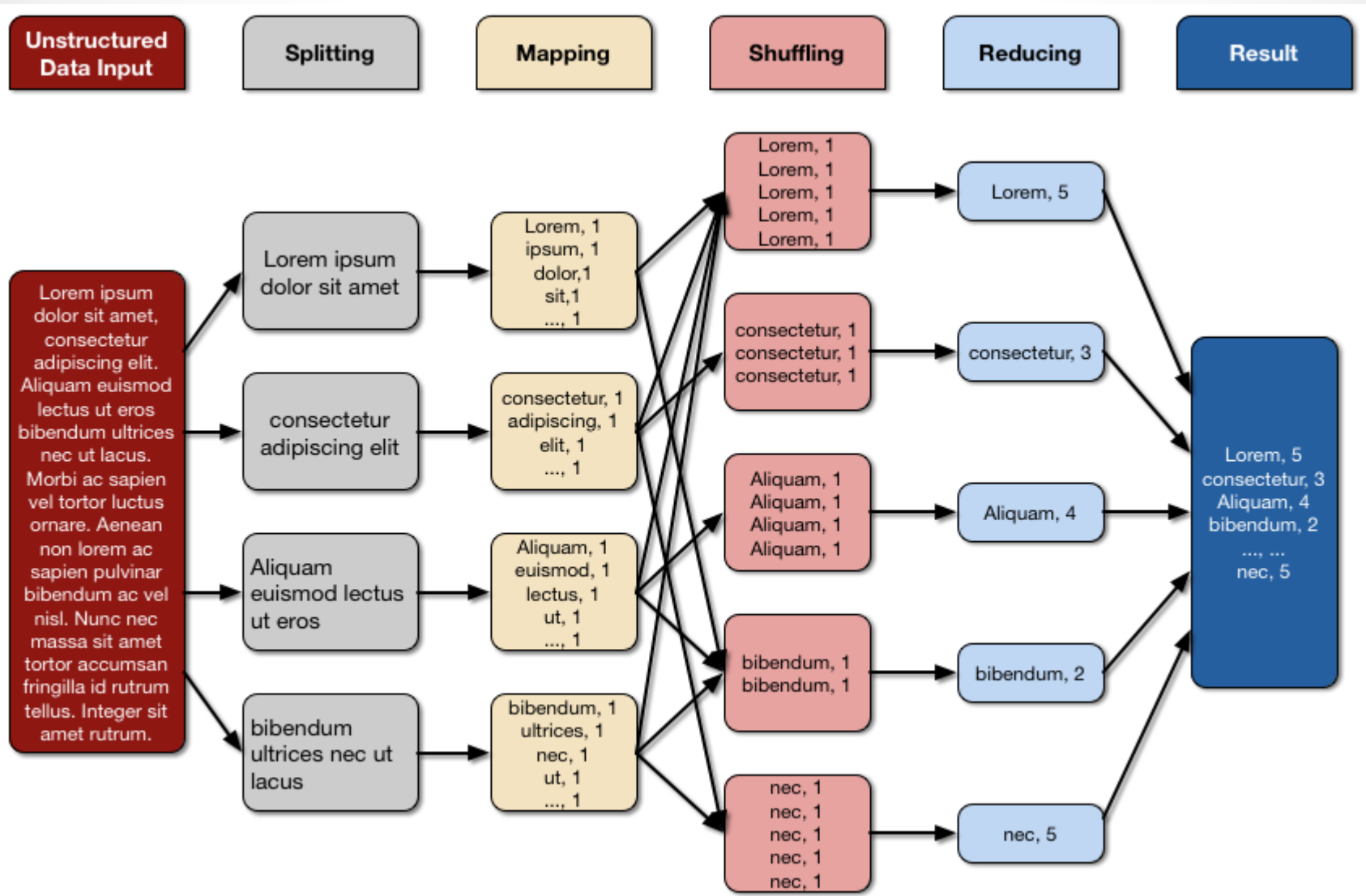
- Framework de armazenamento e processamento distribuído de grandes volumes de dados em clusters.



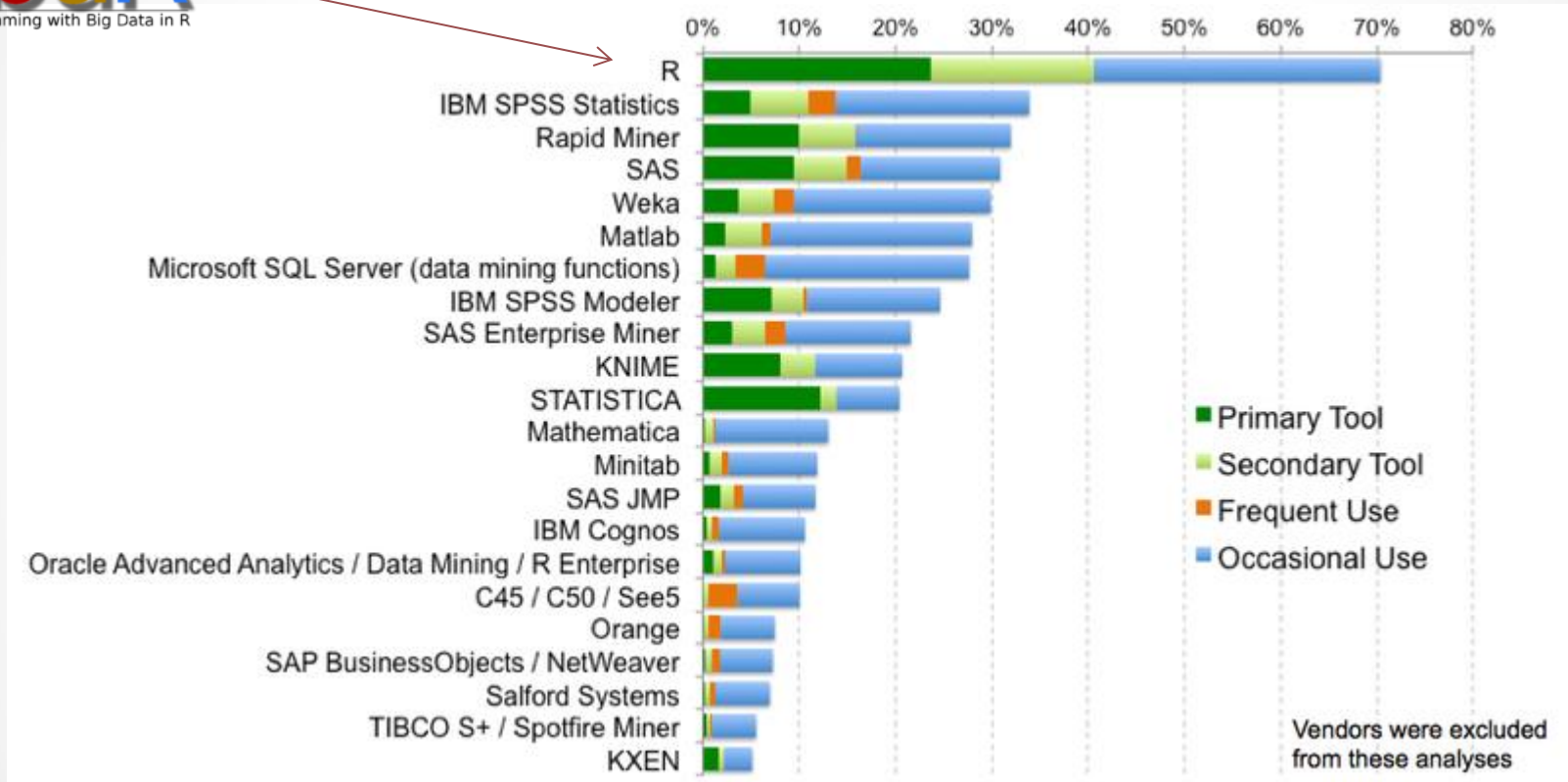
- HDFS (Hadoop Distributed File System)
- HBase (roda em HDFS, é um BigTable)
- Pig (MapReduce em “alto nível”)
- Hive (DataWarehouse)
- R
- ...



# Exemplo de Map-Reduce

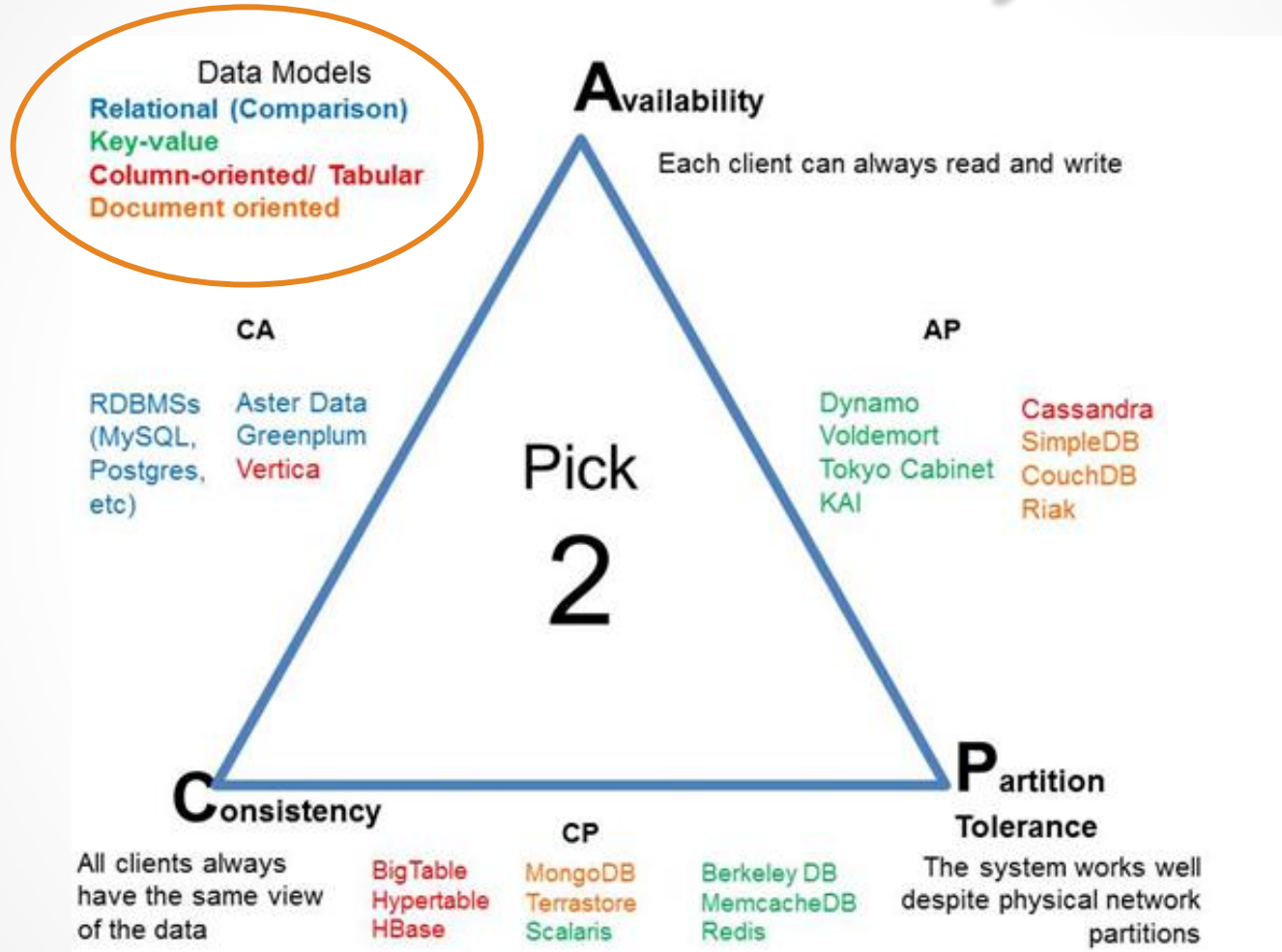


# Sistemas analíticos



Fonte: <http://blog.revolutionanalytics.com/popularity/> (2013)

# NoSQL = Not Only SQL



*CAP Theorem*, de Eric Brewer (Berkeley)

SciDB  
...

# Introdução ao SciDB

- Open source (sob GPL-3)
- Duas versões: comunitária e Paradigm4
- Contribuições: EUA, Rússia, Índia, Europa
- SGBD visando a paralelismo massivo
- Plataforma de armazenamento
- Plataforma analítica
- Modelo de dados baseado em **Arrays**

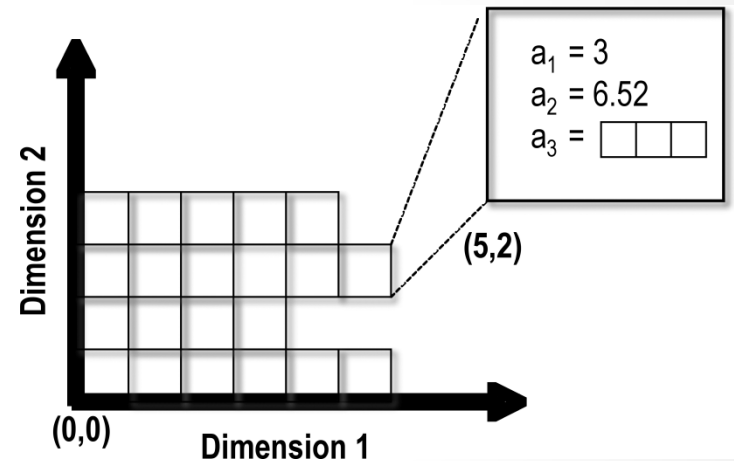


Imagem creditada a [3]

# Por que arrays?

- Fundamentos formais da Matemática
  - A álgebra de arrays é fechada quanto às operações usuais
  - Provas formais de corretude
  - Modelagem matemática oferece várias aplicações
- Experiência na Ciência da Computação
  - APL (*A Programming Language*), 1964



- Muito trabalho em algoritmos para computar e armazenar eficientemente arrays
- Área de pesquisa e desenvolvimento
  - BLAS, LAPack, ScaLAPack – ferramentas open source
  - NAG – Numerical Algorithmics Group
  - R, SAS etc. – pacotes eficientes de álgebra linear

# Definição de dados

```
CREATE ARRAY Simple_Array <
  v1 : double,
  v2 : int64,
  v3 : string >
[ I = 0:*, 5, 0, J = 0:9, 5, 0 ];
```

Attributes  
v1, v2, v3

Dimensions  
I, J

Dimension size  
\* indicates unbounded

Chunk  
size

Chunk  
overlap

*Chunks* estão relacionados a  
armazenamento físico e  
lógico (veremos adiante!)

Imagem creditada a [3]

# Interfaces para programação

**Array  
Functional  
Language  
AFL**

```
aggregate(  
    filter ( Simple_Array,  
            v3 = 'Odd' ),  
    I,  
    avg ( Simple_Array.v1 )  
);
```

**Array  
Query  
Language  
AQL**

```
SELECT avg ( S.v1 )  
  FROM Simple_Array S  
 WHERE S.v3 = 'Odd'  
 GROUP BY S.I;
```

Imagem creditada a [3]



# Capacidade analítica

- Sistema de armazenamento baseados em versões (sem sobrescrita)
  - UPDATE é um append
  - Qualquer estado do sistema pode ser reconstruído num momento específico
- Proveniência
  - Log de todas as consultas para reconstruir como os resultados foram derivados
- Incerteza e metodologia estatística
  - Tipos com margens de erro
  - Funções para testes estatísticos e análise

# Especificação – SciDB

- Stonebraker. M, et al. “Requirements for Science Data Bases and SciDB”, CIDR 2009.
- *The Conference on Innovative Data Systems Research (CIDR) was started in 2002 by Michael Stonebraker, Jim Gray, and David DeWitt to provide the database community with a venue for presenting innovative data systems architectures, as well as a prestigious publication opportunity.*

<http://www-db.cs.wisc.edu/cidr/>

# SciDB vs SGBDs relacionais

- Acelera o acesso aos dados num sistema distribuído
- Armazenamento mais eficiente conforme o número de atributos e dimensões aumentam
- Funções matemáticas são executadas diretamente no formato de armazenamento nativo

**SciDB**

32.5	46.3	81.7	53.6
90.9	35.4	35.9	86.3
42.1	35.7	35.3	45.9
96.7	41.3	89.9	27.6

**16 cells**

**Relational Database**

I	J	value
0	0	32.5
1	0	90.9
2	0	42.1
3	0	96.7
0	1	46.3
1	1	35.4
2	1	35.7
3	1	41.3
0	2	81.7
1	2	35.9
2	2	35.3
3	2	89.9
0	3	53.6
1	3	86.3
2	3	45.9
3	3	27.6

**48 cells**

# SQL vs SciDB

```
CREATE TABLE USNOB (  
    Locn point                not null,  
    -- RA double, DECL double  
    ...  
    B_mag double precision    not null,  
    ...  
    R_mag double precision    not null  
    ...  
);
```

SQL

```
CREATE ARRAY USNOB <  
    B_Mag : double,  
    ...  
    R_Mag : double  
>  
[ RA(double)=*,72000,720, DECL(double)=*,36000,360 ];
```

SciDB

# SQL vs SciDB

```
SELECT U1.Locn, COUNT(*)  
  FROM USNOB AS U1, USNOB AS U2  
 WHERE box(U1.Locn, U1.Locn) &&  
        box(point(U1.Locn[0] - 0.001,  
                  U1.Locn[1] + 0.001),  
            point(U1.Locn[0] + 0.001,  
                  U1.Locn[1] - 0.001))  
 GROUP BY U1.Locn;
```

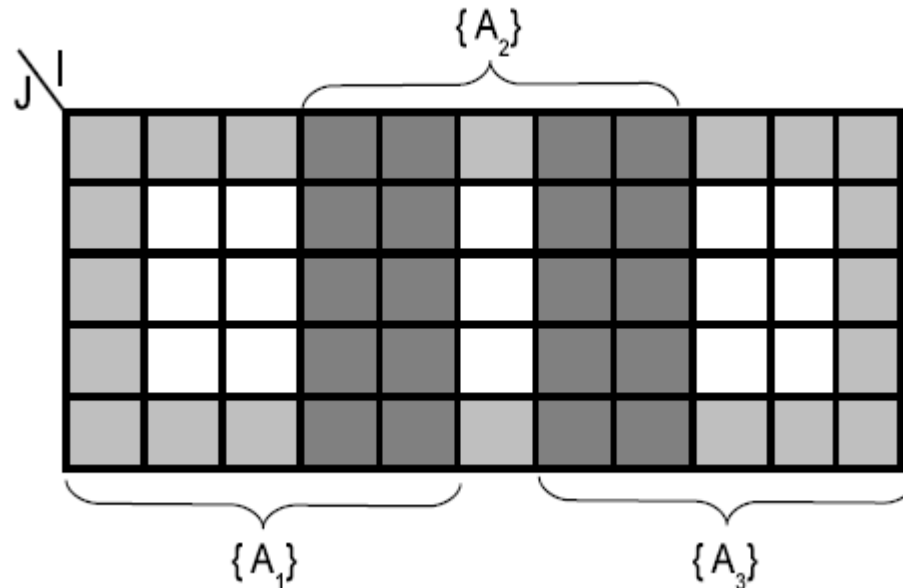
**SQL**

```
window (  
    USNOB,  
    0.001, 0.001,  
    count(*)  
);
```

**SciDB**

Imagem creditada a [3]

# Armazenamento no SciDB

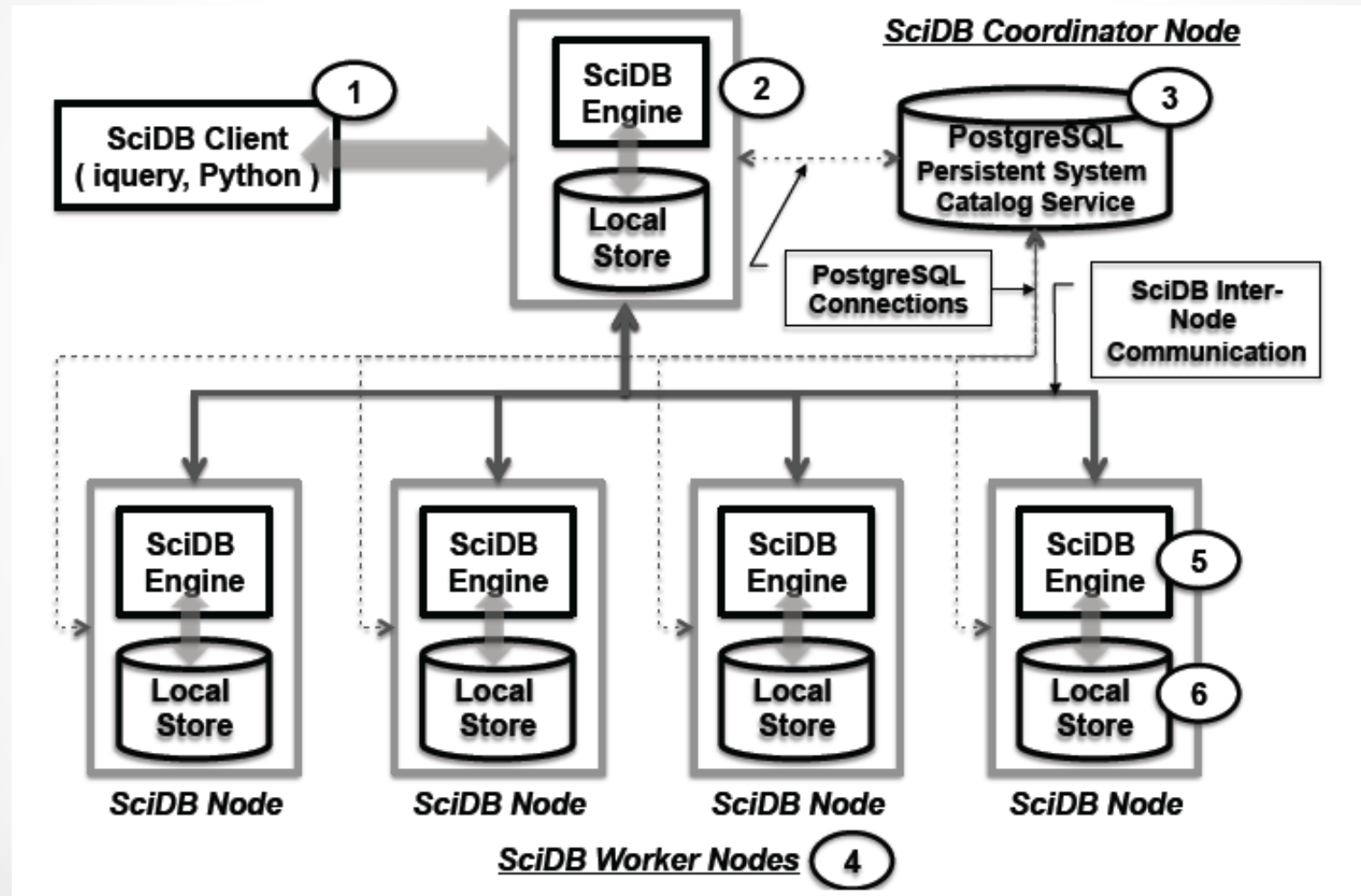


*Arrays e Chunks:* Um *array* 11 x 5 decomposto em três *chunks* 5 x 5 que se sobrepõem. Qualquer janela 3 x 3 desse *array* requer somente a consulta aos dados de um único *chunk* e a operação pode ser paralelizada de três maneiras. As regiões mais escuras do *array* mostram os dados que foram replicados entre os *chunks*, isto é, um particular atributo está armazenado em mais de um *chunk*. As regiões mais claras denotam células não essenciais no *array*, no sentido que qualquer algoritmo que considere uma janela 3 x 3 não conseguirá computar resultados para essas células.

# Chunks

- *Chunks* correspondem a blocos que possuem tamanho lógico fixo e tamanho físico variável;
- Cada *chunk* é armazenado num *container* (arquivo) em disco que pode ser acessado de forma eficiente;
- O tamanho de um *chunk* costuma ser definido a partir de uma média de megabytes, de tal sorte que o custo de realizar buscas é mascarado pela quantidade de dados devolvidos;
- Exemplo: considere um *chunk* que tenha entre 10 e 20 MB de dados. Se os dados forem exclusivamente números de ponto flutuante, o *chunk* conterá entre 500.000 e 1 milhão de elementos (supondo que sejam necessários 8 bytes para cada número de ponto flutuante).
- Leia mais em [7].

# Arquitetura Shared Nothing





# SciDB não é Map-Reduce



# Hands-on

...

# Como usar?

- Para conhecer: AML (Amazon Machine Image)
  - Conta na Amazon AWS: <http://aws.amazon.com/>
  - Distro Ubuntu 12.04
  - SciDB, SciDB-R e SciDB-Py
- Para usar de verdade: instalação tradicional
- Ambos possuem passo-a-passo no fórum:  
<http://scidb.org/forum/viewtopic.php?f=14&t=1431>
- Para os “Paradigm4 employees”... Quick-start VMs! :)
  - <http://downloads.paradigm4.com/>

Alguns exemplos práticos

...

# Referências

- [1] Site oficial  
<http://www.scidb.org/>
- [2] The Fourth Paradigm: Data-Intensive Scientific Discovery  
<http://research.microsoft.com/en-us/collaboration/fourthparadigm/>
- [3] SciDB is not Hadoop  
<http://cdn.oreillystatic.com/en/assets/1/event/63/Big%20Data%20and%20Big%20Analytics%20SciDB%20is%20not%20Hadoop%20Presentation.pdf>
- [4] eScience Institute – University of Washington  
<http://escience.washington.edu/get-help-now/introduction-scidb>
- [5] New England Database Summit, 2013  
<http://db.csail.mit.edu/nedbdays13/slides/lewis.pdf>
- [6] Overview of SciDB  
<http://www.cs.cmu.edu/~pavlo/courses/fall2013/static/papers/sigmod691-brown.pdf>
- [7] Chunk size selection  
[http://scidb.org/HTMLmanual/13.3/scidb\\_ug/ch04s05s02.html](http://scidb.org/HTMLmanual/13.3/scidb_ug/ch04s05s02.html)