

# Data Warehousing

Eng<sup>a</sup> Informática e Computação

FEUP

27 de Maio de 2002

## Agenda

- *Perspectiva de Negócio*
- *Perspectiva Técnica*
- *Perspectiva do Projecto*
- *Razões para construir um DW*
- *Modelação de Dados*
- *Exploração de Dados*
- *DW como processo*
- *Leitura recomendada*

## Perspectiva de Negócio

### Agenda

Negócio

Técnica

Projecto

### ✓ Sistema Analítico vs Sistemas Operacional

*Análise e Acção - Exemplo de aplicações:*

“Balanced Scorecard” - orçamento

Nível de serviço de Fornecedores - rupturas / negociações

Rentabilidade de Espaço - comparação entre lojas / ...

Preços de Venda, Stocks - Excepção e acção

### ✓ Sistema analítico

*Possibilidade de cruzar informação de diversas fontes*

*Facilidade de acesso aos dados e sua agregação - query Wizzard*

*Extracção de dados para Excel -pivot tables*

## Perspectiva Técnica

### Agenda

Negócio

Técnica

Projecto

### ✓ Repositório de dados

*Modelação de dados - Redundância de dados vs Normalização*

*Modelo em Estrela - Factos e Dimensões*

Performance - bons tempos de resposta ao utilizador

Cruzamento de dados provenientes de diversas fontes

### ✓ Exploração de Dados

*OLAP, ROLAP, MOLAP, Query Ad-hoc, Reporting*

## Perspectiva do Projecto

### Agenda

Negócio

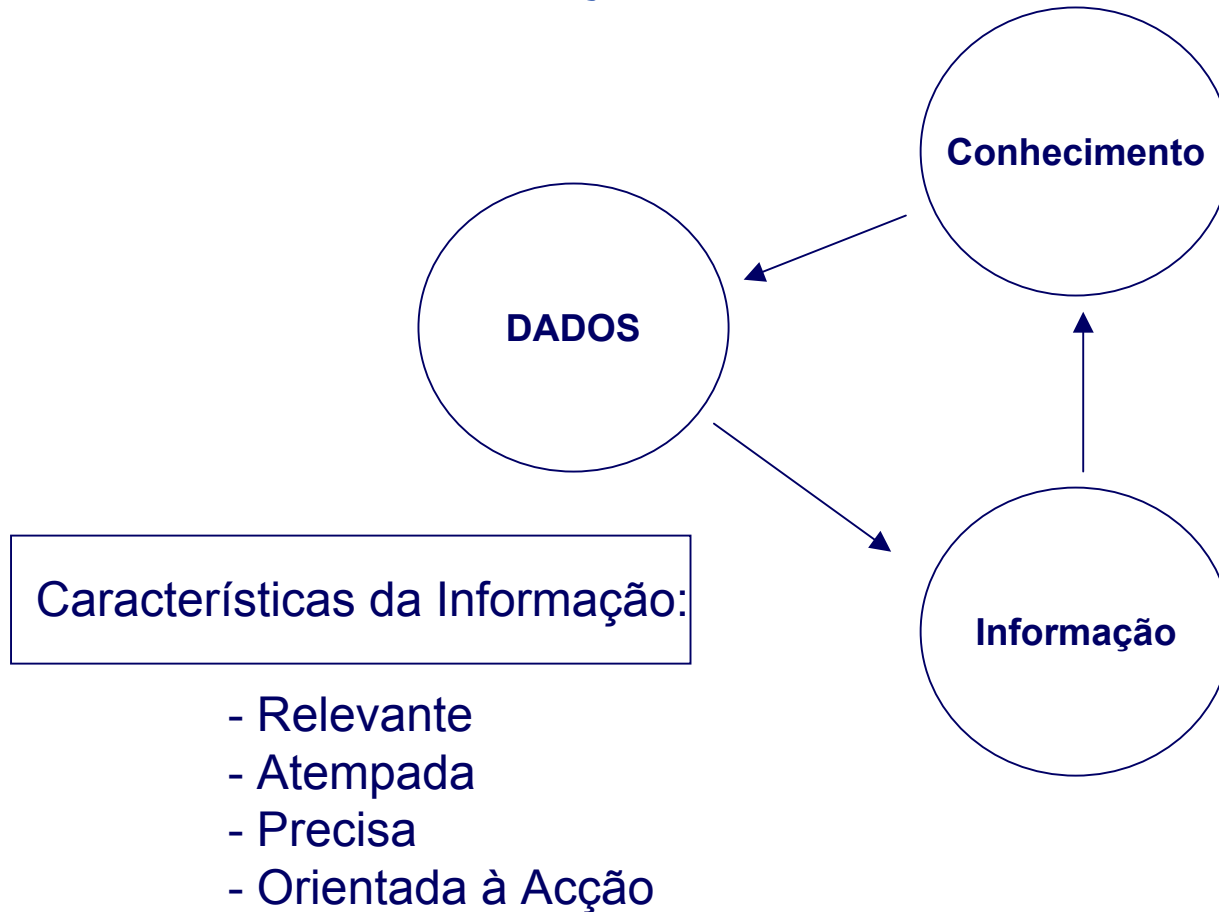
Técnica

Projecto

- ✓ Deve ter início e fim, conhecidos
- ✓ Perspectivas de resultados, conhecidos
- ✓ Comunicação
  - A organização deve saber com o que contar, e quando.*
  - Todos devem saber o que lhes compete no projecto, e quando*
- ✓ Envolvimento do negócio
  - Fases de Análise, Desenho, Testes de aceitação e Implementação*
- ✓ Riscos: qualidade de dados

## Razões para construir um DW

### ✓ Dados vs. Informação



## Razões para construir um DW

✓ Dispersão das fontes de dados

### BANCA

1. B.O.
2. F.O.
3. Call-Center
4. Gestão de Produtos (específicos)

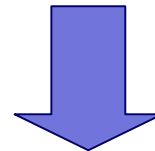
### RETALHO

1. F.O.
2. ERP
3. Entrepasto
4. Financials

### TELECOMUNICAÇÕES

1. Billing
2. Customer Care
3. Engenharia (rede)
4. Financials

- Multiplicidade de Aplicações
- Multiplicidade de Interfaces
- Multiplicidade de Sistemas

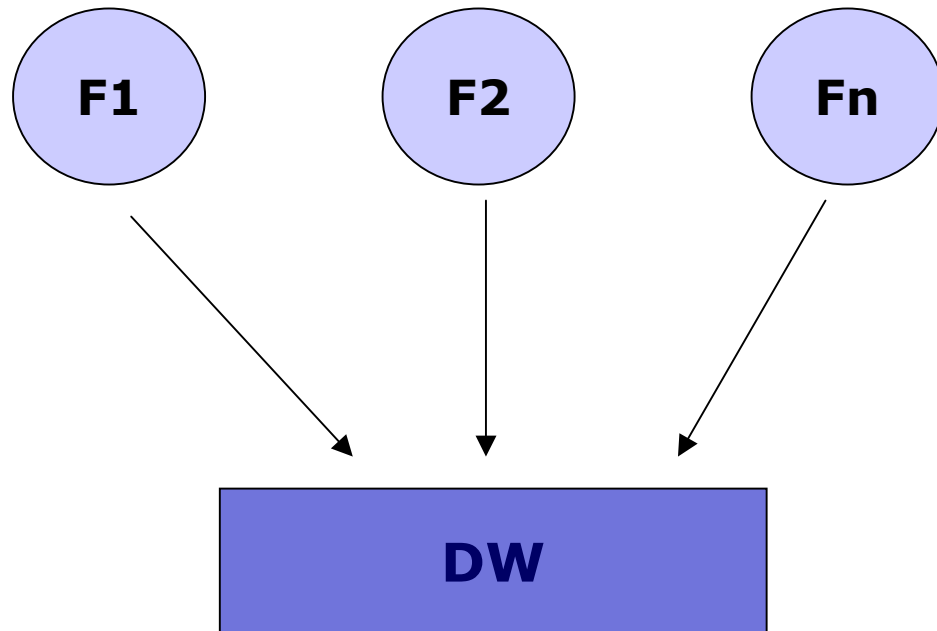


- Visibilidade dificultada
- Orientação operacional (não analítica) das fontes de dados

## Razões para construir um DW

- ✓ Necessidades de Convergência + Navegação

*Visibilidade e Navegação: num único sistema*



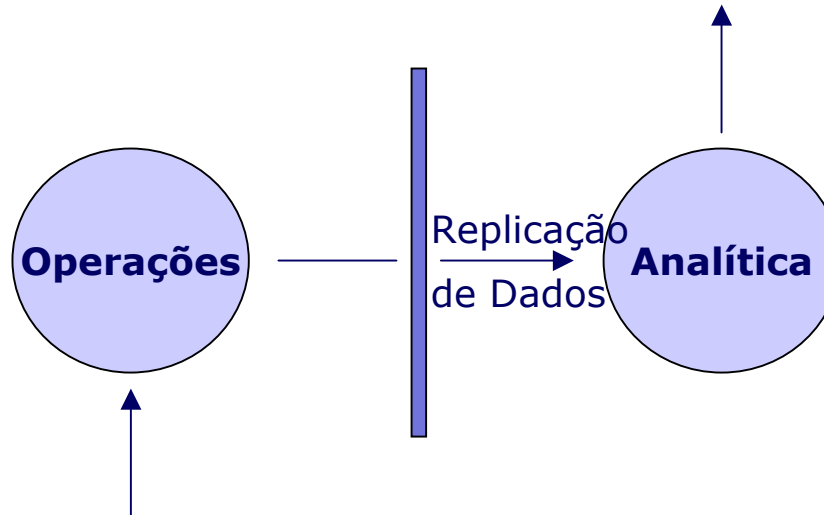


## Razões para construir um DW

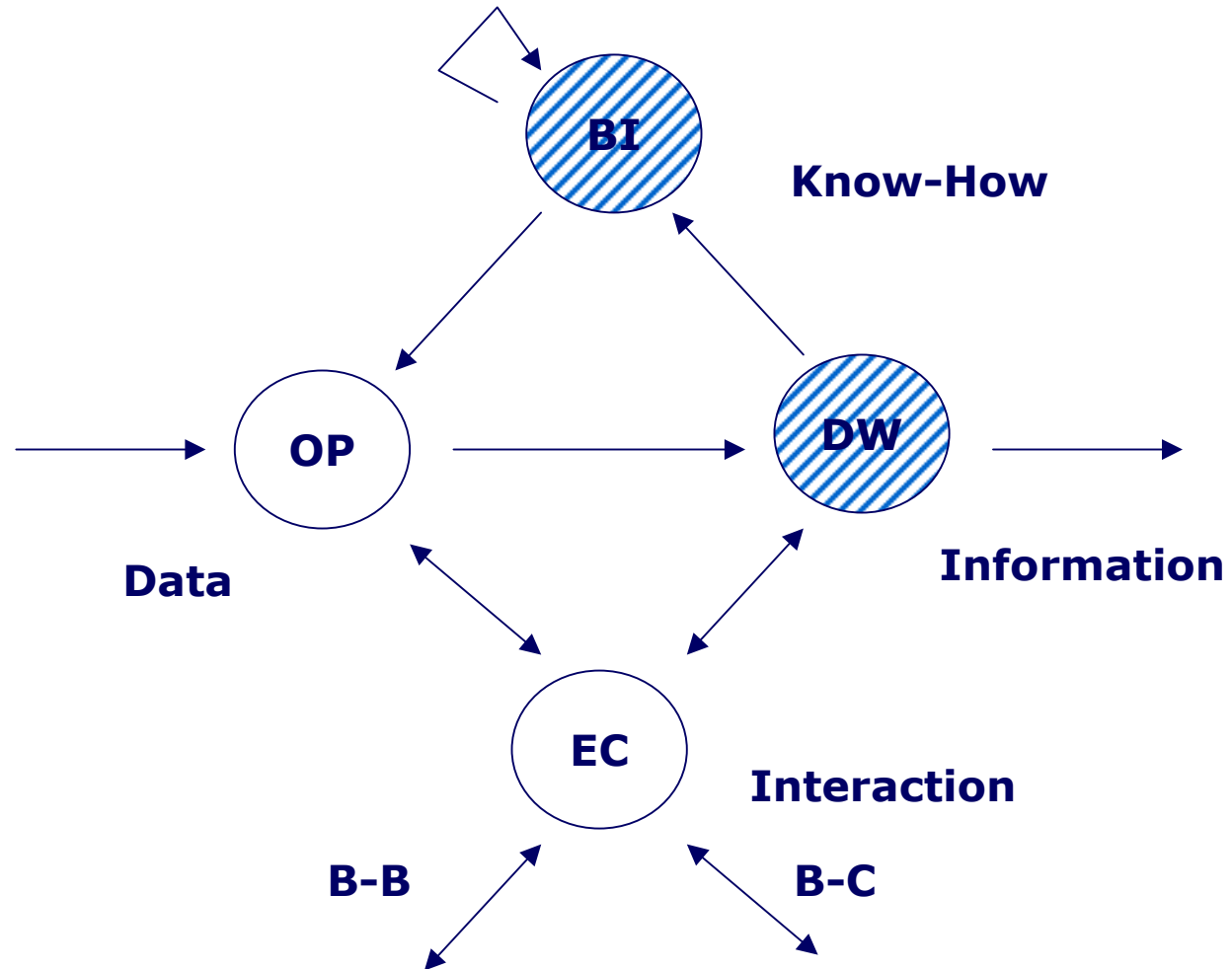
- ✓ Separação Operações / Analítica

*Isolar o impacto das explorações analíticas das operações*

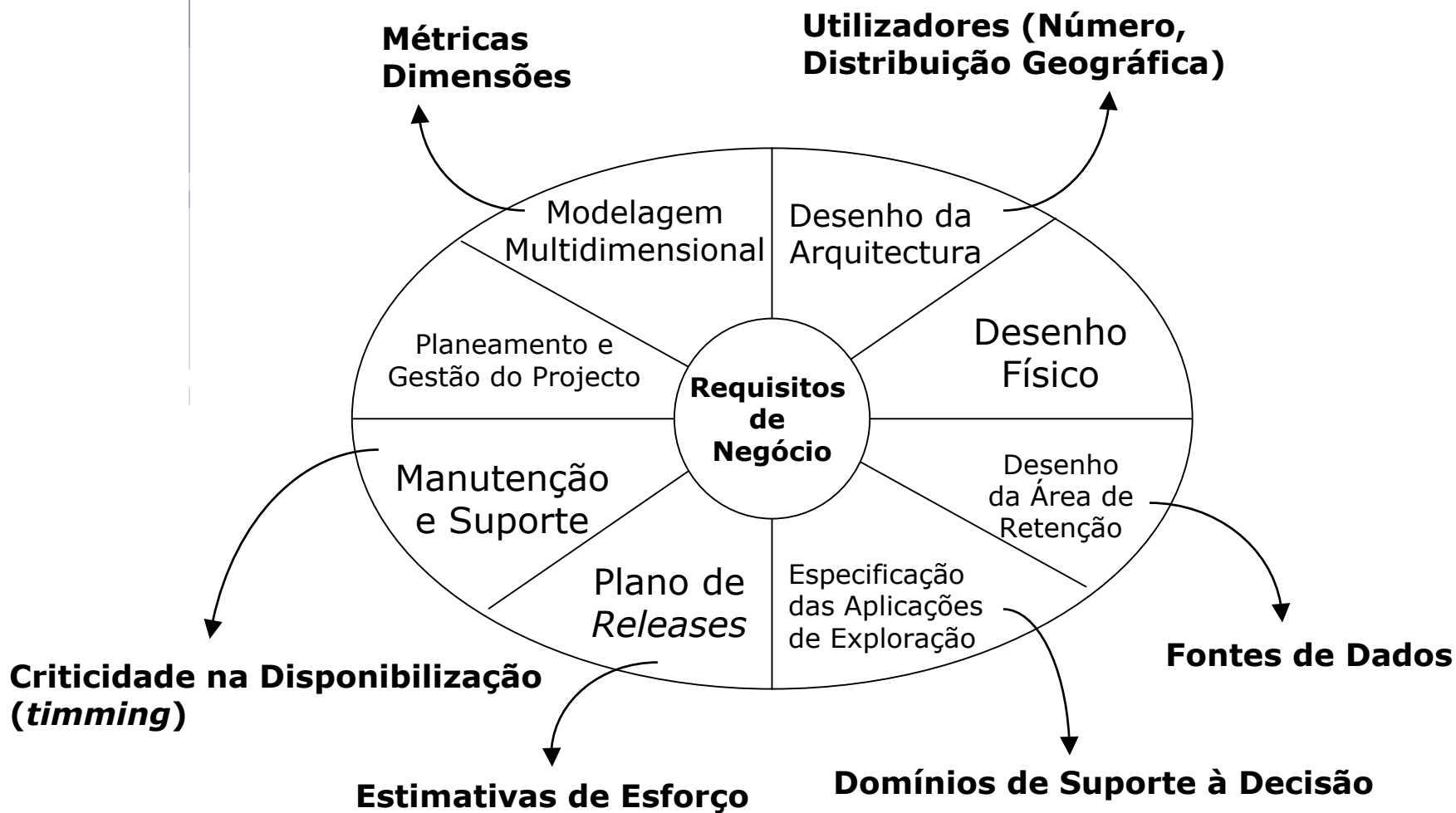
*Nasce a janela nocturna do batch de replicação*



# Enquadramento do DW numa Meta-Arquitectura Organizacional

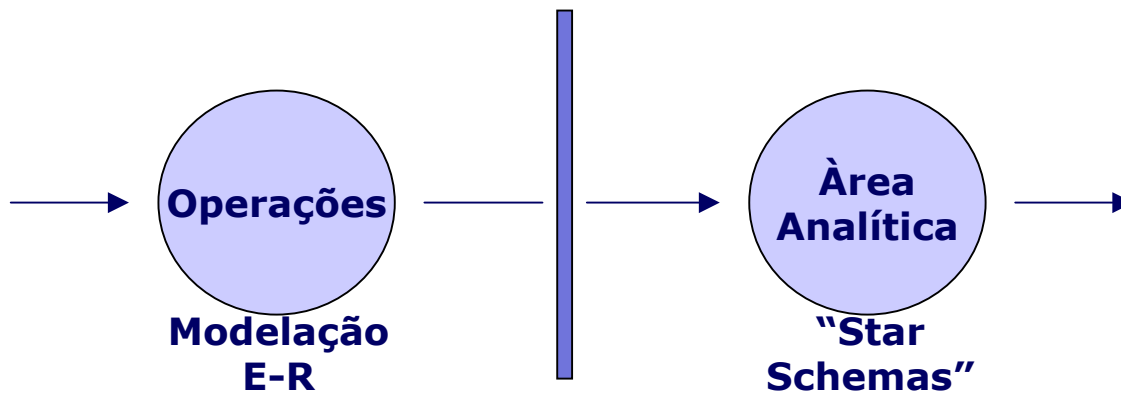


## Levantamento de Requisitos



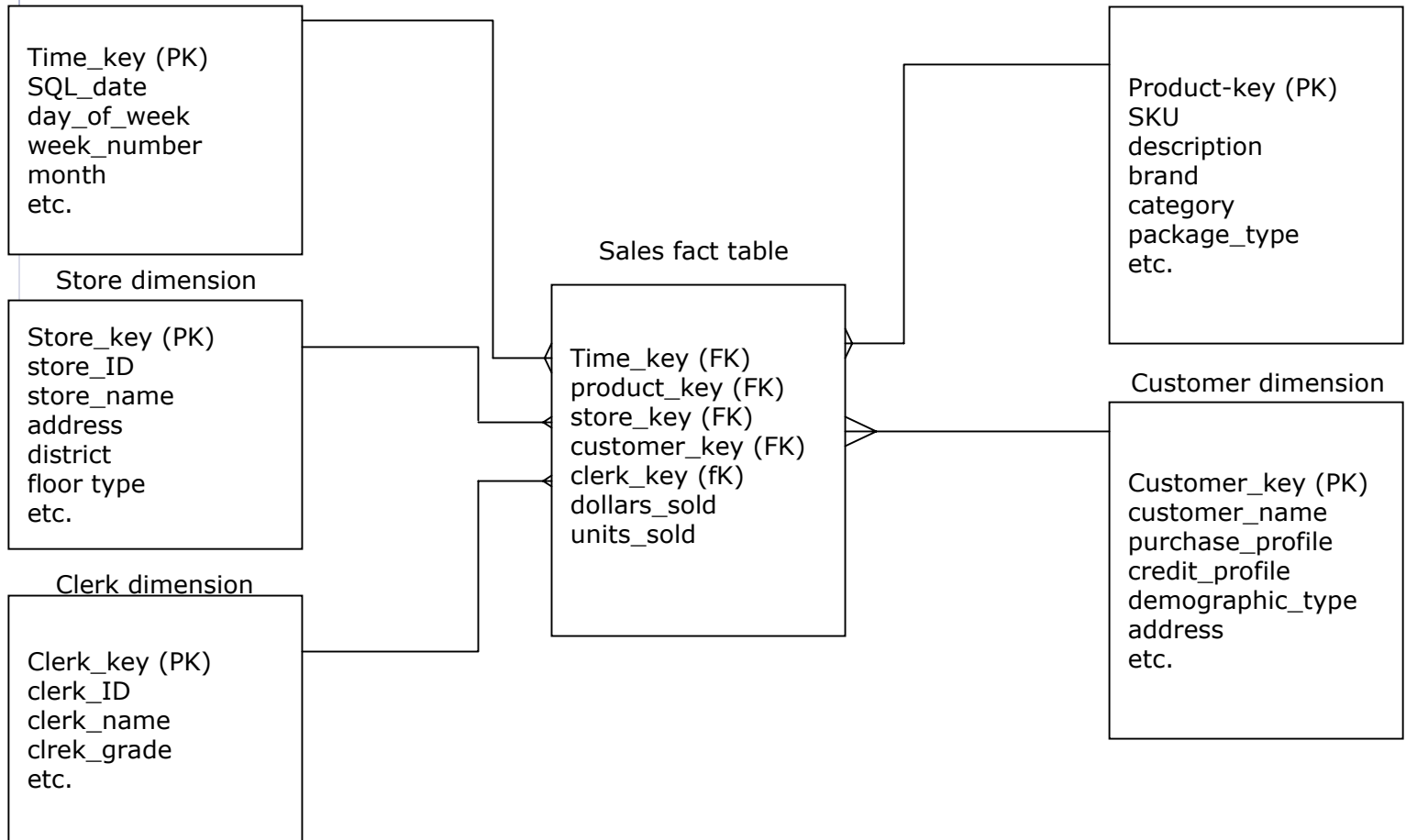
## Modelação de Dados

- ✓ Nas Operações: 3ª Fórmula normal é Lei!  
(CODD)
- ✓ Na área analítica: Emergência do conceito de "Star Schema" (Kimball)



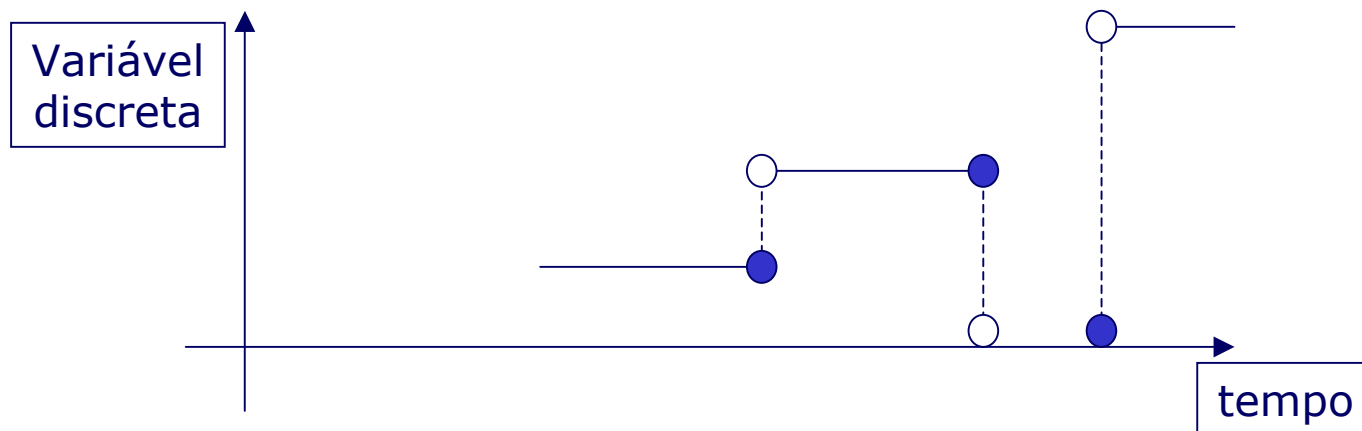
# Modelação de Dados

✓ "Star schema"



## Modelação de Dados

### ✓ Modelação do Tempo



- *As transições de estado não estão normalmente arquivadas (história das alterações) nas bases operacionais; estas só guardam o último valor da variável*
- *Para efeitos de DW, este aspecto é crucial (Análise baseada em históricos, análise de evolução e tendências)*

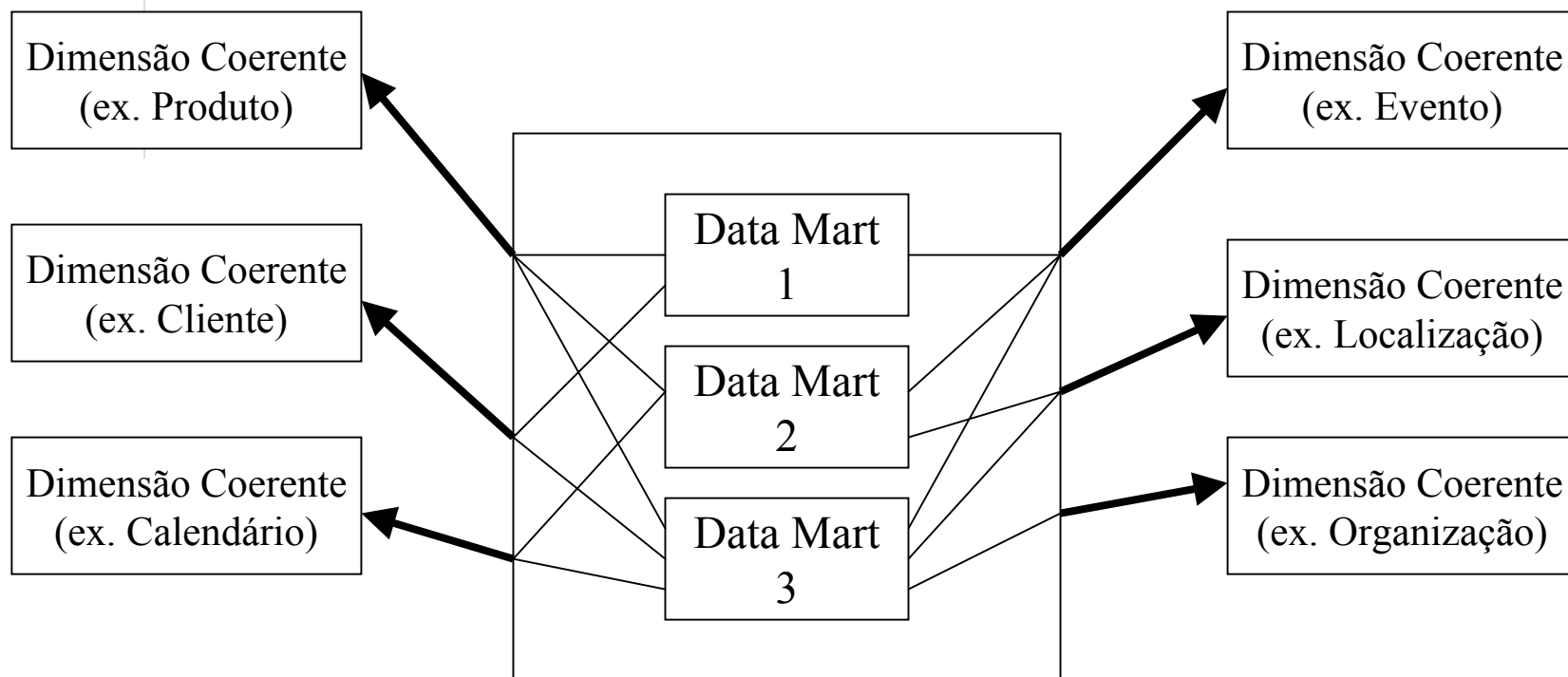
## Modelação de Dados

- ✓ Resolução de Problemas de Incoerência de Dados
  - *Incompletude de Dados*
  - *Eventual incoerência de Dados - DW só pode ter uma versão da verdade!*
  - *Como resolver esta questão?*
    - Definições claras e horizontalmente partilhadas
    - Validação semântica dos conteúdos (Dados) contra as definições

**DW = Oportunidade de definir “Metadata”  
(Dados sobre os Dados) Organizacional**

## Dimensões e Factos Coerentes

- ✓ Um DW com dimensões e factos coerentes são como um "bus interface" que permite a adição sucessiva de novos data marts





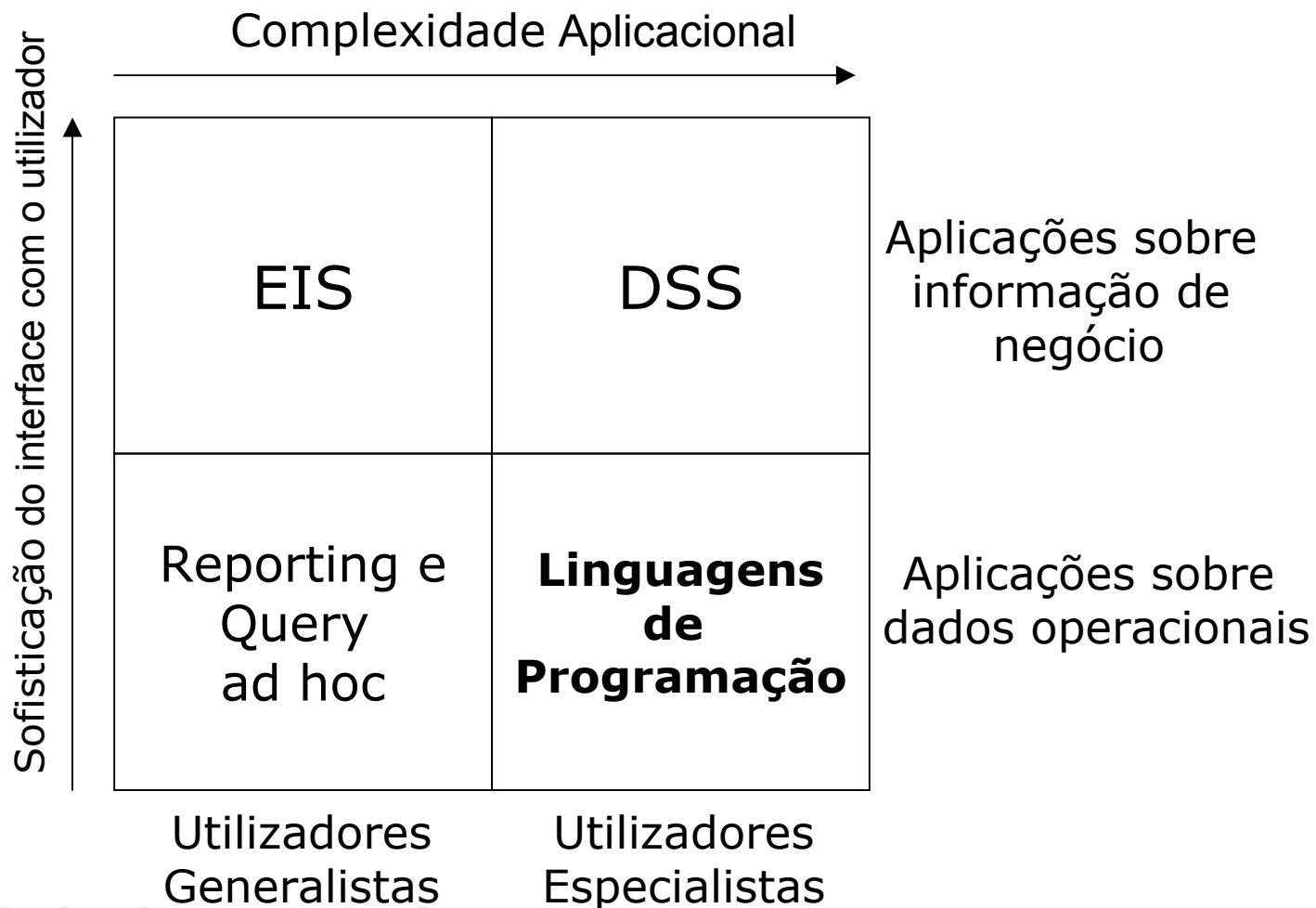
## Dimensões Coerentes\*

\* “conformed dimensions”

- ✓ Dimensão coerente é uma dimensão que tem o mesmo significado qualquer que seja tabela de factos com a qual possa ser ligada
  - *Uma dimensão coerente é partilhada pelos diversos data marts que a referenciam*
  - *Exemplos: cliente, fornecedor, produto, tempo (calendário)*
  - *Do ponto de vista da consistência é uma das vantagens dos modelos ER aplicadas na modelagem multidimensional*
  - *Tornam possível a mesma interpretação do conceito e respectivos atributos ao longo dos diferentes data marts*
  - *Potenciam o cruzamento de informação de diferentes data marts*
  - *Representa 80% do esforço de modelagem*

## Exploração de Dados

### ✓ Tipos de Ferramentas vs Tipos de Utilização



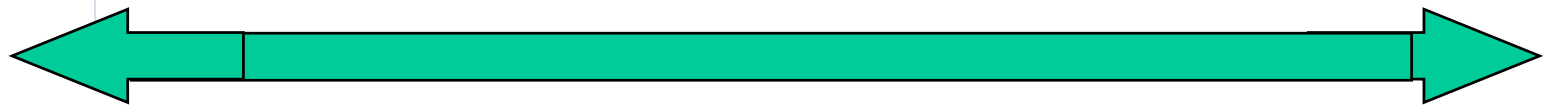
## Exploração de Dados

- ✓ Papel do Analista de Negócio vs Processo de Análise dos Dados

**Conduzido  
pelo Analista**

**Assistido  
pelo Analista**

**Conduzido  
pelos Dados**



**Processamento  
Informação**

- *Query*
- *Reports*

**Processamento  
Analítico**

- OLAP Relacional
- OLAP Multidimensional

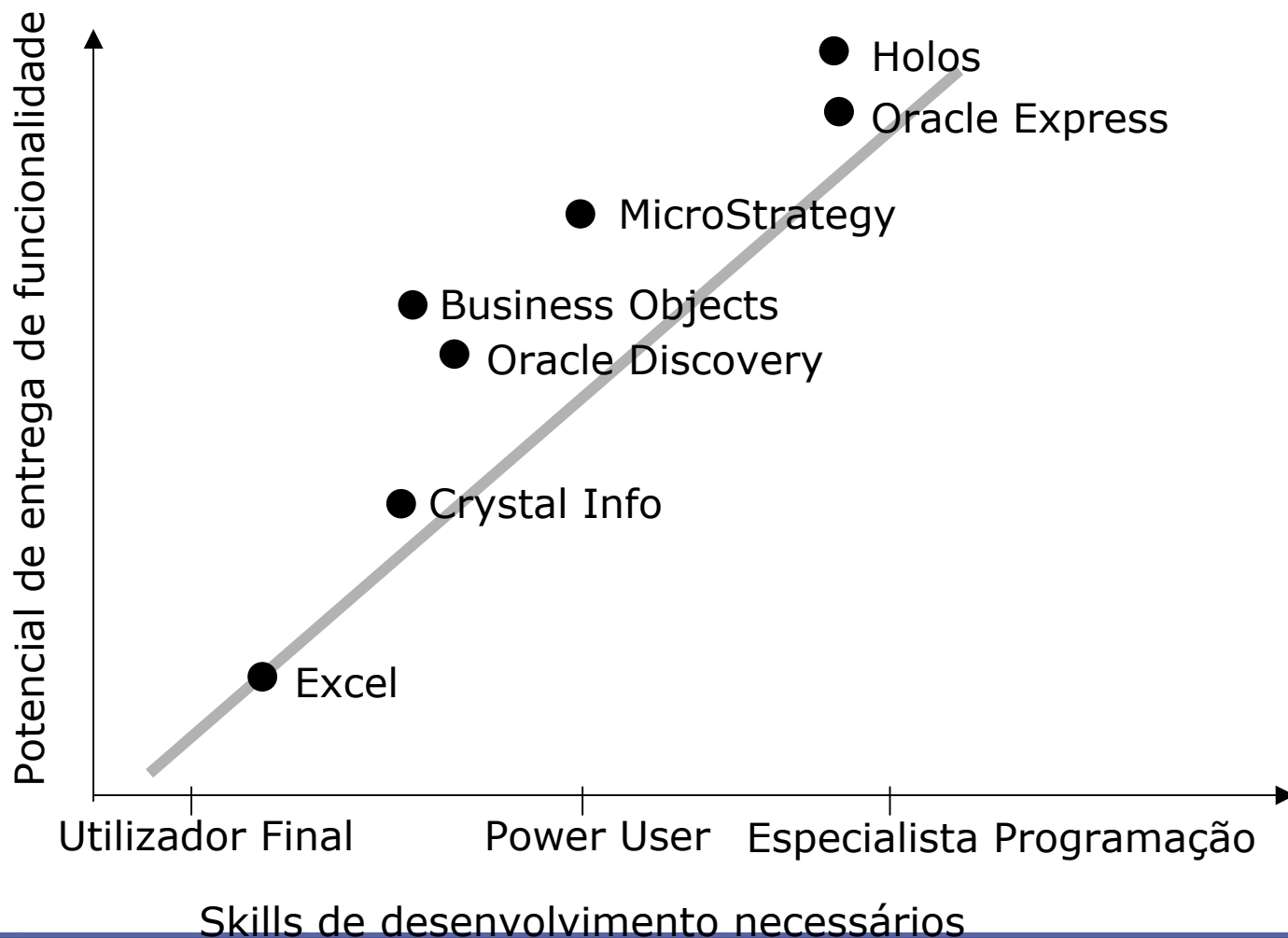
**Métodos  
Estatísticos**

**Sistemas  
Periciais**

***Knowledge  
Discovery***

## Exploração de Dados

### ✓ Facilidade Desenvolvimento vs Potencial Funcionalidade



## Exploração de Dados

### ✓ Algumas Características Importantes

#### ***Interface***

- *Conjunto completo de funcionalidades disponíveis no interface com o utilizador (todas ou quase todas as disponíveis nas ferramentas de reporting e query ad hoc)*
- *Capacidade de gerir alertas e gerar gráficos sofisticados*
- *Sistema de navegação fácil e intuitivo. Capacidade do utilizador criar as suas próprias "bookmarks"*
- *Suportar modelos de previsão*

## Exploração de Dados

- ✓ **Base de Dados Multidimensional (MOLAP)**
  - *A monitorização de alto nível sobre o comportamento do negócio deve ser suportada por instrumentos eficientes*
  - *A transformação dos dados operacionais em informação de negócio introduz forte sumarização (agregações de agregações)*
  - *A representação multidimensional introduz vantagens de performance no acesso aos dados relativamente ao relacional.*
  - *Algumas características importantes:*
    - 4GL potente para programação do carregamento e actualização dos cubos de dados
    - Metacubos ou cubos virtuais
    - Algoritmos adequados à dimensão e esparsidade dos cubos - cubos pequenos podem residir em memória principal

# Exploração de Dados - Exemplos

✓ OLAP (MOLAP: Seagate Holo, Oracle Express, ...)

cm_vendas	Agosto	Setembro	Outubro	Novembro	Dezembro	Acum. Dezembr
Vendas Brutas	39.212.320	35.333.820	36.944.380	37.844.720	59.690.550	431.784.600
Vendas Líquidas	34.836.980	31.409.050	32.673.020	33.551.340	51.731.500	303.351.500
Mrg Teórica	4.249.320	3.960.830	4.277.730	4.495.756	5.623.372	47.231.150
Consumíveis	60.666	72.581	55.699	65.124	99.122	707.693
Armação	630.828	666.769	619.847	822.026	965.839	7.287.314
Transportes	345.016	351.157	308.820	413.593	447.651	3.776.055
Custos Logísticos	975.845	1.017.926	928.667	1.235.419	1.333.490	11.063.370
Mrg c/ Logística	3.212.617	2.870.323	3.293.367	3.195.013	4.210.762	35.460.090
Ajust. Conferência	11.880	13.149	8.953	14.404	65.438	239.014
RDC Deb. Reac. Shop	21.558	38.959	19.438	49.717	140.326	594.651
RDC Deb. Comp. Mrg		3.776	1.888		169.757	178.524
Débitos	21.558	42.735	15.326	49.717	310.083	773.175
Prov Quebras	292.962	274.060	284.748	307.024	498.899	3.352.939
Mrg Comercial	2.929.534	2.625.850	3.014.991	2.923.223	3.966.509	32.641.310
Ajust. + CCC s/ compras	3.265.515	3.296.201	2.533.196	3.820.470	4.256.514	36.263.400
Ajust. Rappel no stock	15.407	-202.073	111.059	-554.800	587.059	-526.573
Rappel Internacional						
Rappel + CCC	3.277.921	3.094.128	2.644.255	3.265.670	4.843.573	35.736.930
RDC Comercial	1.490.703	1.309.785	1.215.860	1.374.688	2.348.151	17.169.150
Benefício Comercial s/s	7.685.751	7.231.936	6.764.047	6.117.780	10.571.170	86.072.820
Benefício Comercial	7.698.150	7.029.763	6.675.106	7.562.901	11.150.230	85.546.250
RDC de Publicidade						

TOTAL DIR.COMERCIAL		TOTAL INSGIAS		TOTAL INSGIAS	
TOTAL INSGIAS		TOTAL INSGIAS		TOTAL INSGIAS	
Real		R/H		R/H	
Mrg Teórica		Mrg Teórica		Mrg Teórica	
Mrg Comercial		Mrg Comercial		Mrg Comercial	
Rappel + CCC		Rappel + CCC		Rappel + CCC	
RDC Comercial		RDC Comercial		RDC Comercial	
Benefício Comercial		Benefício Comercial		Benefício Comercial	
Contribuição Comercial		Contribuição Comercial		Contribuição Comercial	

**Análise de Quadrantes por Categorias**  
R/O Vendas Brutas vs R/O Mrg Teórica %

Utilizador: ltopol\_a  
Data: 04-FEB-2000  
Versão: 2.0

Critério: 100,6% (R/O Vnd Brt), -3,3% (R-O Mrg%)

	Vnd Brt	Mrg%
DC Não Operacional #DCNO	100	54,3
DC Perceções #DCP	101	-2,4
DC Textil #DCT	100,6	-0,3

+ (R-O) Mrg%

	Vnd Brt	Mrg%
DC Alimentar #DCA	110,3	-2,8

+ R/O Vnd Brt

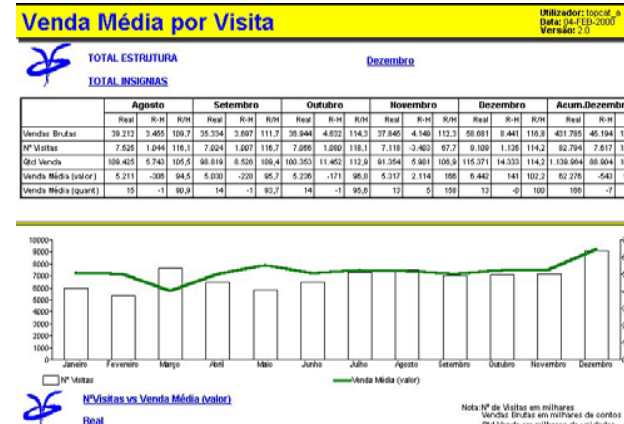
	Vnd Brt	Mrg%
DC Bazar Leveiro #DCBL	90,7	-5,4

+ R/O Mrg%

	Vnd Brt	Mrg%
DC Bazar Pesado #DCBP	137,2	-3,7

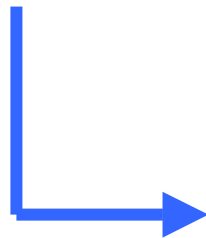
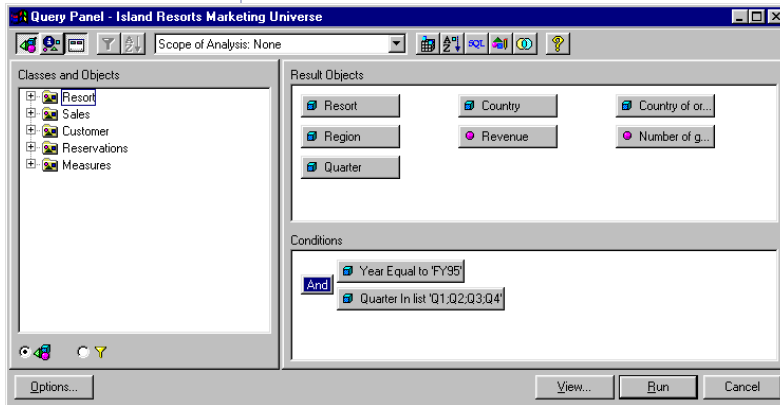
TOTAL DIR.COMERCIAL  
TOTAL INSGIAS  
Real vs. Ocorrêdo  
Dezembro

TOTAL DIR.COMERCIAL	Vnd Brt	Mrg%
	108,6	-3,3



# Exploração de Dados - Exemplos

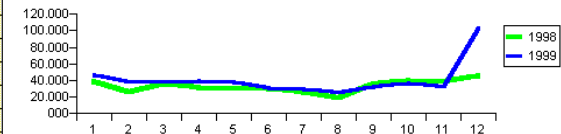
- ✓ Ad-Hoc Query (MicroStrategy, Brio, Business Objects, Oracle Discovery, ...)



## Análise Evolutiva de Vendas

BASE MOLH.MAIZENA EXPRES.250G CX (002.000.022)

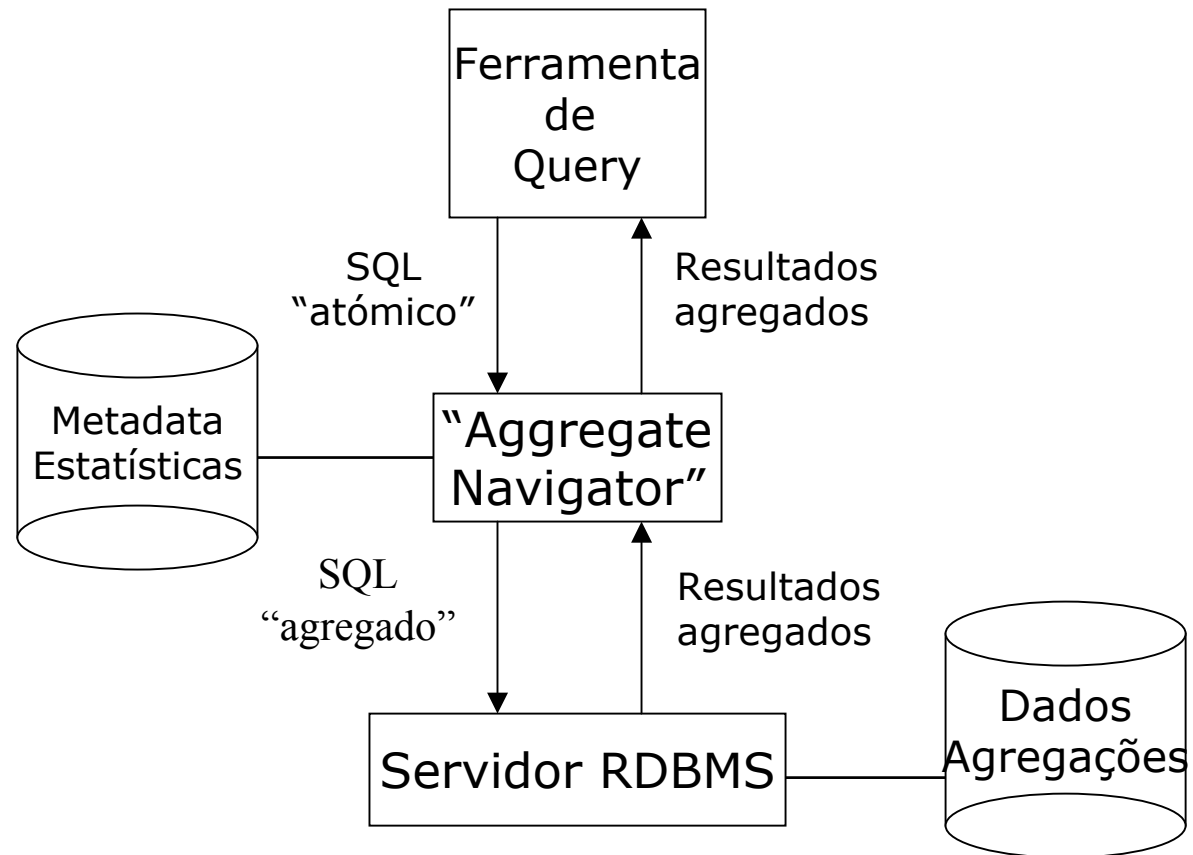
Ano	Mês	VND BRT
1998	1	39.004
1998	2	24.864
1998	3	35.616
1998	4	31.584
1998	5	30.240
1998	6	29.904
1998	7	25.200
1998	8	17.808
1998	9	35.280
1998	10	39.984
1998	11	37.968
1998	12	45.696
1999	1	47.376
1999	2	37.632
1999	3	38.976
1999	4	40.320
1999	5	38.304
1999	6	31.248
1999	7	29.376
1999	8	28.224
1999	9	38.400
1999	10	38.400
1999	11	38.400
1999	12	38.400





## Explorações de Dados

- ✓ ROLAP => "Aggregate awareness"



## Sumário: Razões para construir um DW

- ✓ Convergência / Visibilidade num ponto único
- ✓ Fácil navegabilidade
- ✓ Separação Operacional / Analítica
  - *Não impacto de queries nas operações*
  - *Modelação diferente*
- ✓ Resolução de problemas de coerência de definições e de dados
- ✓ Sustentação de novas aplicações de Business Intelligence

## Processos base de um DW

- ✓ Extracção de Dados
- ✓ Transformação:  
*Limpeza, Reformatação, Combinação*
- ✓ Carregamento
- ✓ Controlo de Qualidade
- ✓ Publicação; Actualização/Refresh; Interrogação
- ✓ Auditoria
- ✓ Back-Up e Recuperação

## A Função de DW Manager

- ✓ Definição da Função
  - *"Guardião do Templo"*
  - *"Editor" responsável pela qualidade dos dados publicados*
  - *Responsável pela Metadata Organizacional*
  - *Facilitador da priorização do desenvolvimento de todas as aplicações de Suporte à Decisão*
  - *Responsável pela Publicação, nos timings acordados, das novas versões dos dados*

## A Função de DW Manager

- ✓ Valências exigidas pela Função
  - *Conhecimento do Negócio (em particular, das necessidades dos "Knowledge Workers")*
  - *Capacidade de Comunicação e Facilitação / Geração de Consenso*
  - *Capacidade Organizativa e Disciplina de Entrega (c/Controlo de Qualidade)*
  - *Conhecimentos Técnicos específicos de DW / DSS / BI*
  - *Resiliência (DW é um Processo!)*
  - *Capacidade de Modelagem Avançada*

## Dimensão de um DW

Qual o tamanho da  
"Piscina" ?

**Dados**

- **Números de Tabelas**
- **Volume de Dados**
- **Número de Fontes operacionais (instâncias)**
- **Número de Data-Marts (instâncias)**

**Processos**

- **Número de Processos distintos do Batch**
- **Número de Queries por dia  
(atendidos / não atendidos)**

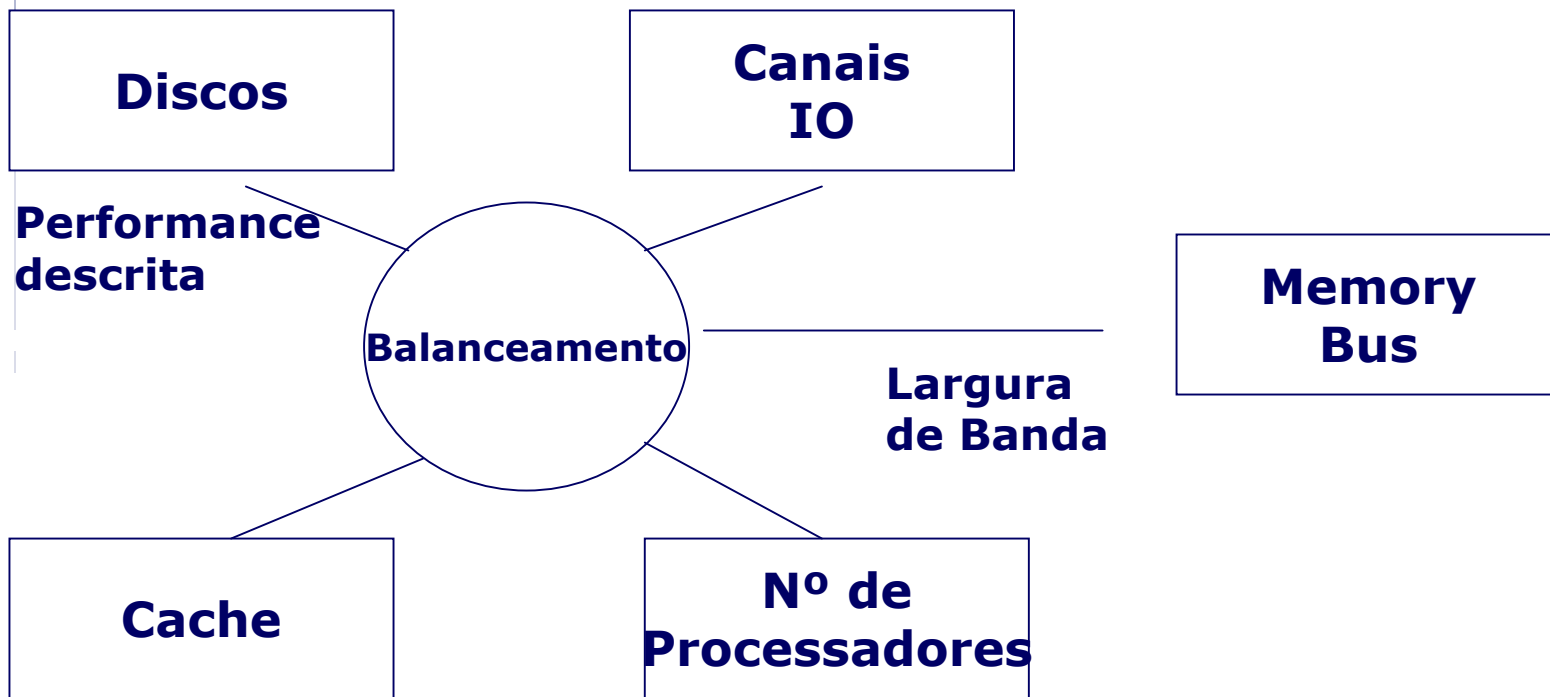


**IMPACTO NA INFRAESTRUTURA DE SERVIDORES!**

## Problemas Técnicos que se levantam

- ✓ Distribuição da computação SMP ou MPP?
- ✓ Como diminuir ao máximo a contenção entre processos?
- ✓ Como facilitar a paralelização implícita do software?
- ✓ Como garantir Back-Up's íntegros e prontos a suportar uma recuperação?
- ✓ Como implementar "Disaster Recovery"?
- ✓ Como garantir que as operações pouco extensas de delete/update não têm grande impacto de performance?
- ✓ Como promover o balanceamento dinâmico de carga entre processadores?

# Balanceamento da Arquitectura Tecnológica como um todo





## **Pensar Arquitectura Tecnológica DW ... a tempo**

**A escalabilidade não acontece:  
Arquitectura-se!**

- ✓ DW: A explosão de utilização é sempre mais rápida do que o previsto
- ✓ Atenção à "Procura Escondida" ( reprimida )
- ✓ Necessidade de revisão da Arquitectura Tecnológica é normal

## Leituras recomendadas

- ✓ “The Data Warehouse, Lifecycle Toolkit”, Ralph Kimball