

DATA MANAGEMENT FOR BUSINESS INTELLIGENCE

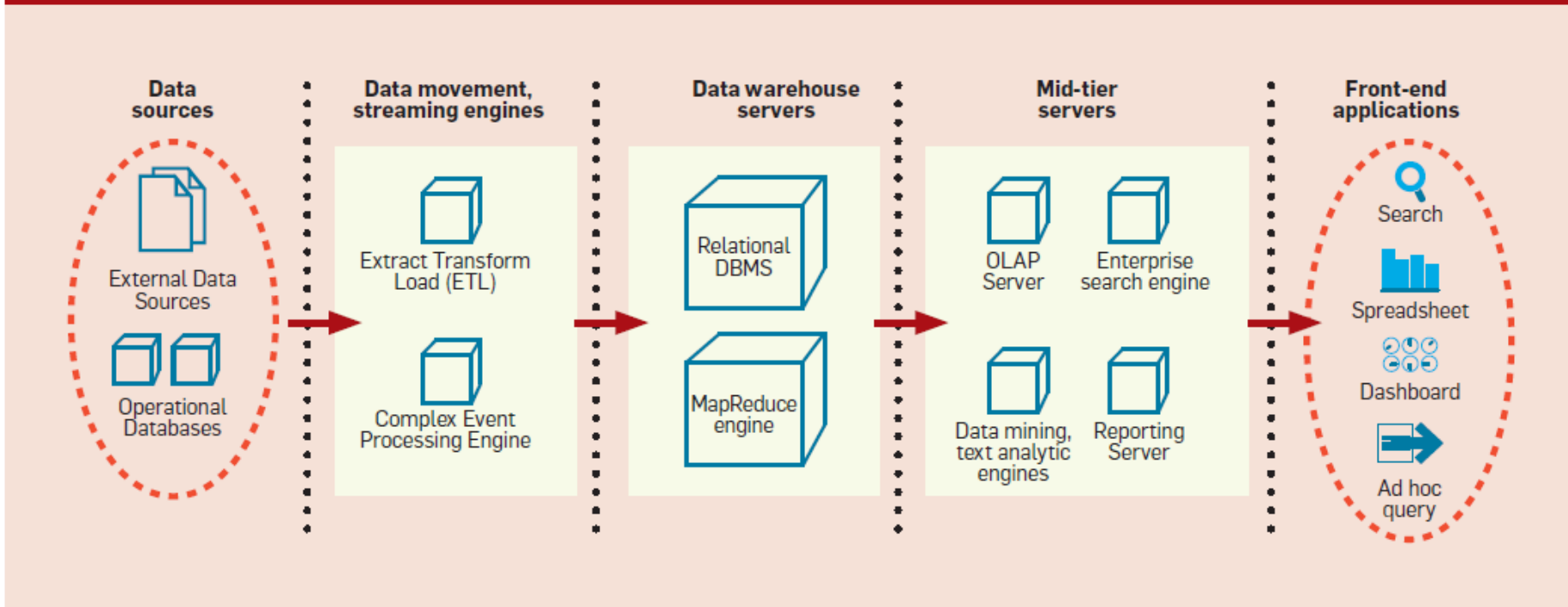
Data Access: Files



BI Architecture

2

Figure 1. Typical business intelligence architecture.



Two issues

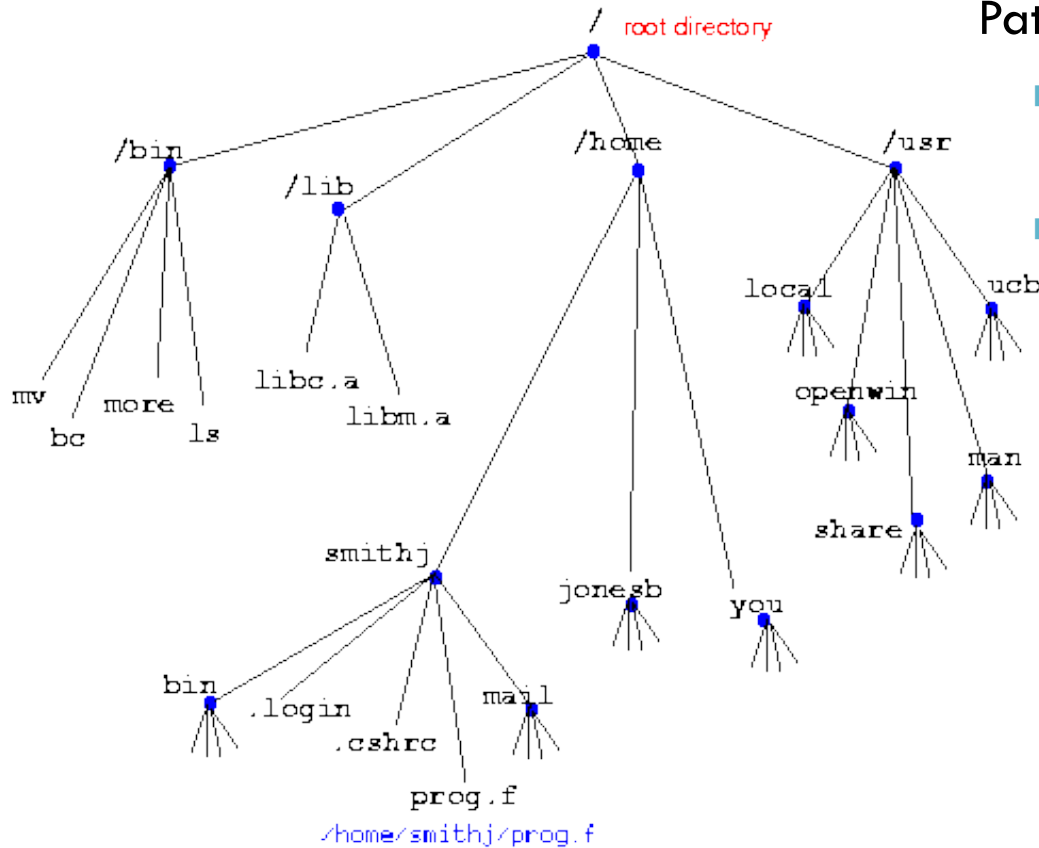
3

- **Where** are my files?
 - Local file systems
 - Distributed file systems
 - Network protocols

- Which **format** is file data in?
 - Text
 - CSV, JSON

Local file system

4



Path of a resource

Windows:

■ C:\Program Files\Office\sample.doc

Linux:

■ /usr/home/r/ruggieri/sample.txt

Local file system

5

A logical abstraction of persistent mass memory

- hierarchical view (tree of directories and files)
- types of resources (file, directory, pipe, link, special)
- resource attributes (owner, rights, hard links)
- services (indexing, journaling)

Sample file system:

- Windows
 - NTFS, FAT32
- Linux
 - EXT2, EXT3, JFS, XFS, REISERFS, FAT32

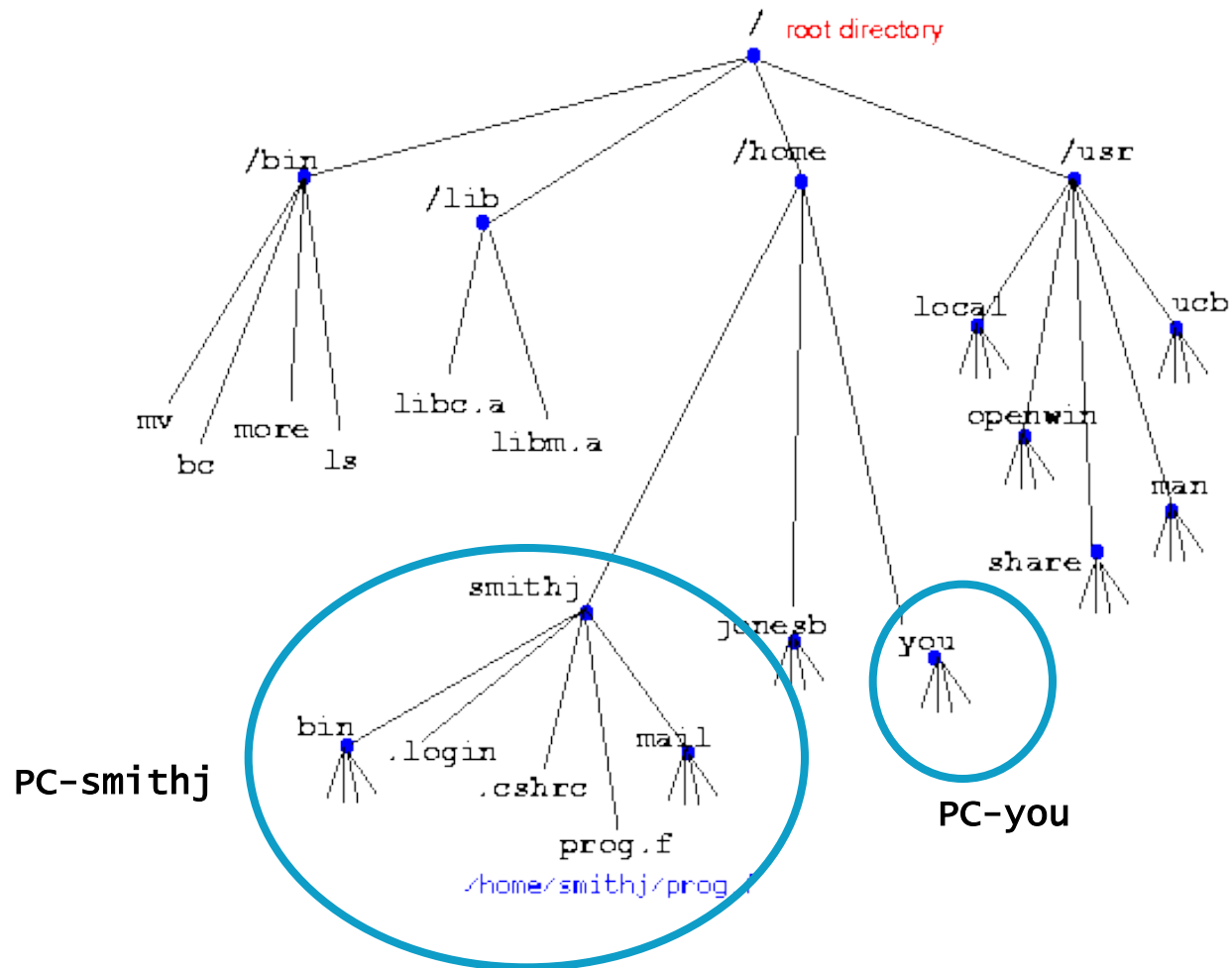
Disk file systems [\[edit\]](#)

Disk file systems are usually block-oriented. Files in a

- ADFS – Acorn's Advanced Disc filing system, such as
- AdvFS - Advanced File System, designed by Digital Equipment Corporation
- AFS (Not to be confused with Andrew File System)
- AFS - Ami File Safe, a commercial file system shipped with AmigaOS
- AofsFS - File System used by the Oberon and A2000
- AthFS - AtheOS File System, a 64-bit journaled file system
- BFS - the Boot File System used on System V releases
- BFS – the Be File System used on BeOS, occasionally on Linux
- Btrfs - is a copy-on-write file system for Linux announced in 2007
- CBMFS – The filesystem used on most Commodore 64 computers
- CMDFS – A filesystem extension added to CBMFS
- CP/M file system — Native filesystem used in the CP/M operating system
- DDFS – Data Domain File System, the data deduplication file system
- DTFS – Desktop File System, featuring file compression
- DOS 3.x - Original floppy operating system and file system
- EAfs – Extended Acer Fast Filesystem, used on Acer Aspire notebooks
- Extent File System (EFS) – an older block filing system
- ext – Extended file system, designed for Linux systems
- ext2 – Second extended file system, designed for Linux systems
- ext3 – A journaled form of ext2.
- ext4 – A follow up for ext3 and also a journaled file system
- ext3cow – A versioning file system form of ext3.
- FAT – File Allocation Table, used on DOS and Microsoft Windows
 - VFAT – Optional layer on Microsoft Windows
 - FATX – A modified version of Microsoft Windows
- FFS (Amiga) – Fast File System, used on Amiga systems
- FFS – Fast File System, used on *BSD systems

Distributed file system

6



Distributed file system

7

Acts as a client for a remote file access protocol

- logical abstraction of remote persistent mass memory

Sample file system:

- Samba (SMB)
or Common Internet File System (CIFS)
- Network File System (NFS)
- Hadoop Distributed File System (HDFS)

Mount/unmount

Distributed file systems [\[edit\]](#)

See also: [Comparison of distributed file system](#)

Distributed file systems are also called network file

- 9P, the Plan 9 from Bell Labs and Inferno distributed file system
- Amazon S3
- Andrew File System (AFS) is scalable and local
- Apple Filing Protocol (AFP) from Apple Inc.. A
- DCE Distributed File System (DCE/DFS) from
- File Access Listener (FAL) is an implementation
- Microsoft Office Groove shared workspace, used
- NetWare Core Protocol (NCP) from Novell is used
- Network File System (NFS) originally from Sun
- OS4000 Linked-OS provides distributed file system:
- **Secure File System (SFS)**
- Self-certifying File System (SFS), a global network
- Server Message Block (SMB) originally from IBM for authentication.

Network protocols

8

- Files accessed through **explicit** request/reply
- A **local copy** has to be made before accessing data
- Resource naming:
 - Uniform Resource Locator (URL)
 - `scheme://user:password@host:port/path`
 - <http://bob:bye@www.host.it:80/home/idx.html>
 - scheme = protocol name (http, https, ftp, file, jdbc, ...)
 - port = TCP/IP port number

HTTP Protocol

9

- HyperText Transfer Protocol
 - URL: <http://user:pwd@www.di.unipi.it>
 - State-less connections
 - Crypted variant: Secure HTTP (HTTPS)
- Windows clients
 - Any browser
 - > wget
 - GNU <http://www.gnu.org/software/wget/>
 - W3C <http://www.w3.org/Library>
- Linux clients
 - Any browser
 - > wget

SCP Protocol

10

- **Secure Copy**
 - `> scp data.zip user@mylinux.unip.it:datacopy.zip`
 - File copy from/to a remote account
 - File paths must be known in advance

- **Client**
 - **command line:**
 - `> scp/pscp > scp2`
 - **Windows GUI**
 - WinSCP <http://winscp.sourceforge.net>
 - SSH Secure Shell
 - **Linux GUI**
 - SCP: default

Two issues

11

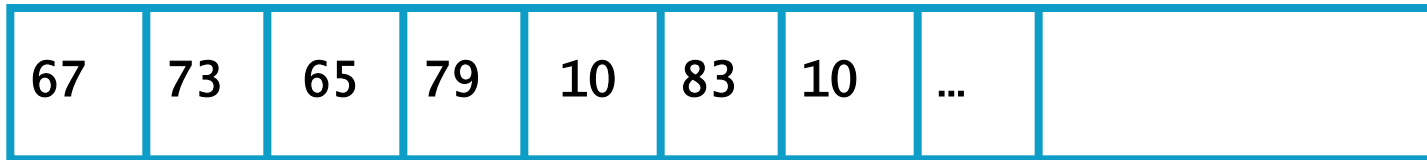
- **Where** are my files?
 - Local file systems
 - Distributed file systems
 - Network protocols

- Which **format** is file data in?
 - Text
 - CSV, ARFF, JSON

What is a file?

12

- File = sequence of bytes



How bytes are mapped to chars?

13

- Character set = alphabet of characters
- Coding bytes by means of a character set
 - ▣ ASCII, EBCDIC (1 byte per char)
 - ▣ UNICODE (1/2/4 bytes per char)

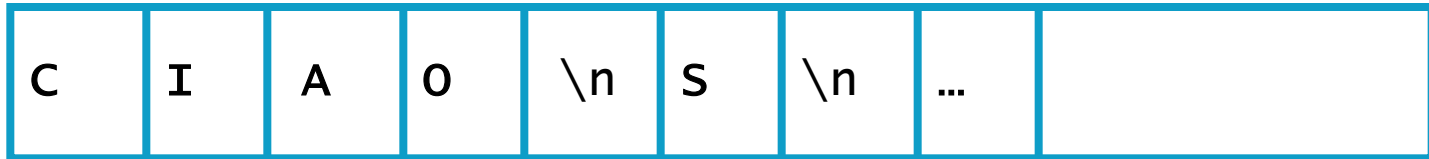
American Standard Code for Information Interchange

CODE	CHAR	CODE	CHAR	CODE	CHAR	CODE	CHAR	CODE	CHAR
0	NUL	26	SUB	52	4	78	N	104	h
1	SOH	27	ESC	53	5	79	O	105	i
2	STX	28	FS	54	6	80	P	106	j
3	ETX	29	GS	55	7	81	Q	107	k
4	EOT	30	RS	56	8	82	R	108	l
5	ENQ	31	US	57	9	83	S	109	m
6	ACK	32	SP	58	:	84	T	110	n
7	BEL	33	!	59	;	85	U	111	o
8	BS	34	"	60	<	86	V	112	p
9	HT	35	#	61	=	87	W	113	q
10	LF	36	\$	62	>	88	X	114	r
11	VT	37	%	63	?	89	Y	115	s
12	FF	38	&	64	@	90	Z	116	t
13	CR	39	'	65	A	91	[117	u
14	SO	40	(66	B	92	\	118	v
15	SI	41)	67	C	93]	119	w
16	DLE	42	*	68	D	94	^	120	x
17	DC1	43	+	69	E	95	_	121	y
18	DC2	44	,	70	F	96	`	122	z
19	DC3	45	-	71	G	97	a	123	{
20	DC4	46	.	72	H	98	b	124	
21	NAK	47	/	73	I	99	c	125	}
22	SYN	48	0	74	J	100	d	126	~
23	ETB	49	1	75	K	101	e	127	DEL
24	CAN	50	2	76	L	102	f		
25	EM	51	3	77	M	103	g		

Text file = file+character set

15

- Text file = sequence di characters



Viewing text files

16

- By a text editor
 - ▣ Emacs, Notepad++, TextPad, GEdit, Vi, etc.
- “Carriage return” character
 - ▣ Start a new line
 - ▣ Coding
 - Unix: 1 char ASCII(0A) ('\n' in Java)
 - Windows: 2 chars ASCII(0D 0A) (“\r\n” in Java)
 - Mac: 1 char ASCII(0D) ('\r' in Java)
 - ▣ Conversions
 - > **dos2unix**
 - > **unix2dos**

Text file = file+character set

17

- Text file = sequence di **lines**

C	I	A	O
S			
...			

Tabular data format

18

Column

Row

Mario	Bianchi	23	Student
Luigi	Rossi	30	Workman
Anna	Verdi	50	Teacher
Rosa	Neri	20	Student

Representing tabular data in text files

19

□ Comma Separated Values (CSV)

- A row per line
- Column values in a line separated by a special character
- Delimiters: comma, tab, space

```
Mario,Bianchi,23,Student  
Luigi,Rossi,30,Workman  
Anna,Verdi,50,Teacher  
Rosa,Neri,20,Student
```

Representing tabular data in text files

20

□ Fixed Length Values (FLV)

- A row per line
- Column values occupy a fixed number of chars
 - Allow for random access to elements
 - Higher disk space requirements

Mario	Bianchi	23	Student
Luigi	Rossi	30	Workman
Anna	Verdi	50	Teacher
Rosa	Neri	20	Student

Quoting

21

- What happens in CSV if a delimiter is part of a value?
 - ▣ Format error
- Solution: **quoting**
 - ▣ Special delimiters for start and end of a value (ex. “ ... “)

Mario Bianchi 23 Student
Luigi Rossi 30 Workman
Anna Verdi 50 Teacher
Rosa Neri 20 Student



“Mario Bianchi” 23 Student
“Luigi Rossi” 30 Workman
“Anna Verdi” 50 Teacher
“Rosa Neri” 20 Student

Missing values

22

- How to represent missing values in CSV or FLV?
 - ▣ A reserved string: “?”, “null”, “”

“Mario Bianchi” 23 Student
“Luigi Rossi” 30 ?
“Anna Verdi” 50 Teacher
“Rosa Neri” ? Student

Meta-data

23

- Describe properties of data
 - ▣ Table name, column name, column type, ...

name	surname	age	occupation
string	string	int	string
Mario	Bianchi	23	Student
Luigi	Rossi	30	Workman
Anna	Verdi	50	Teacher
Rosa	Neri	20	Student

How to represent meta-data in text files?

24

- One or two rows: names and types

name	surname	age	occupation
string	string	int	string



name,surname,age,occupation
string,string,int,string

Meta-data and data in text files

25

- In the same file
 - ▣ Meta-data first (header), then data

name	surname	age	occupation
string	string	int	string
Mario	Bianchi	23	Student
Luigi	Rossi	30	Workman
Anna	Verdi	50	Insegnante
Rosa	Neri	20	Studente



```
name,surname,age,occupation  
string,string,int,string  
Mario,Bianchi,23,Studente  
Luigi,Rossi,30,Operaio  
Anna,Verdi,50,Insegnante  
Rosa,Neri,20,Studente
```

Two issues

26

- **Where** are my files?
 - Local file systems
 - Distributed file systems
 - Network protocols

- Which **format** is file data in?
 - Text
 - CSV, JSON

Data interchange issue

27

- Problem: **data interchange** between applications
 - Proprietary data format do not allow for easy interchange
 - CSV with different delimiters, or column orders
 - Similar limitations of FLV, ARFF, binary data, etc.

- Solution:
 - definition of an interchange format...
 - ... marking data elements with their meaning ...
 - ... so that any other party can easily interpret them.

DATA MANAGEMENT FOR BUSINESS INTELLIGENCE

Data Access: Relational Data Bases

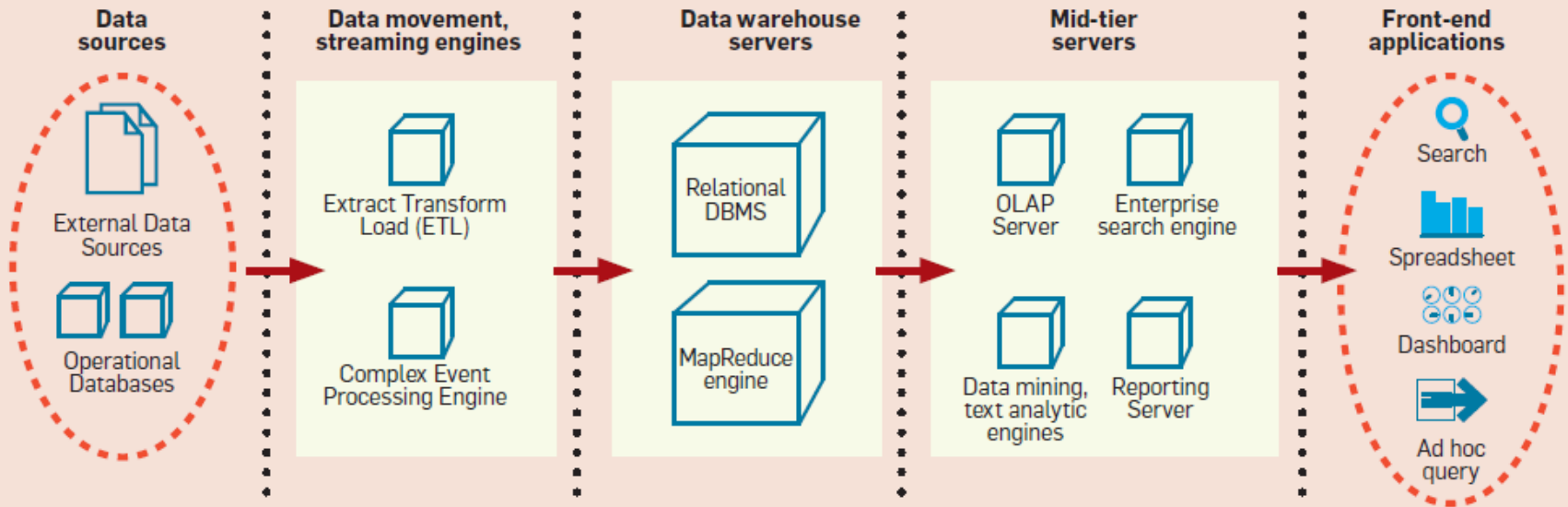
Computer Science Department, University of Pisa



BI Architecture

29

Figure 1. Typical business intelligence architecture.



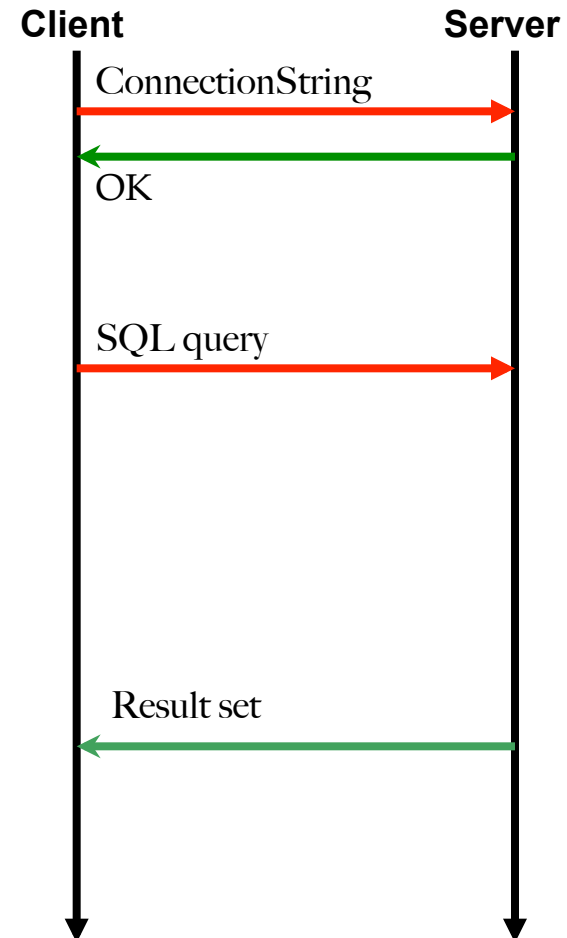
Connecting to a RDBMS

30

- **Connection protocol**
 - locate the RDBMS server
 - open a connection
 - user authentication

- **Querying**
 - query SQL
 - SELECT
 - UPDATE/INSERT/CREATE
 - stored procedures
 - prepared query SQL

- **Scan Result set**
 - scan row by row
 - access result meta-data



Connection Standards

31

- ODBC - Open DataBase Connectivity
 - ▣ Windows: [odbc](#) Linux: [unixodbc](#), [iodbc](#)
 - ▣ Tabular Data

- JDBC
 - ▣ Java APIs for tabular data

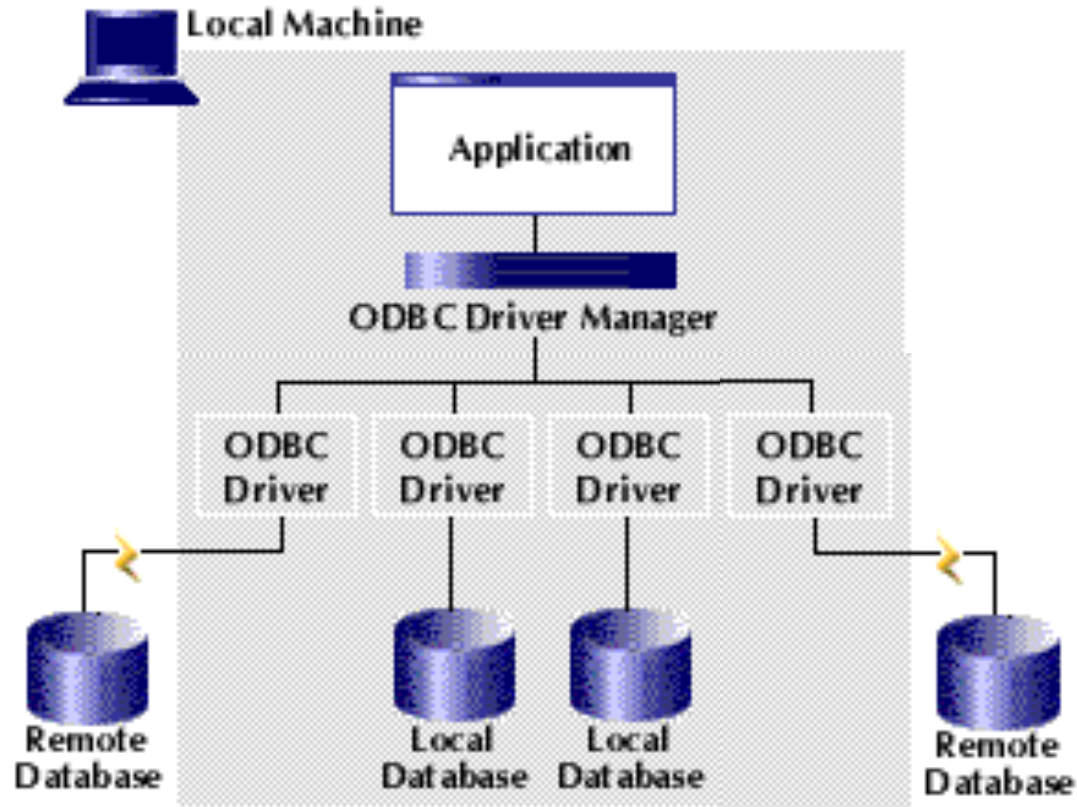
- OLE DB (Microsoft)
 - ▣ Tabular data, XML, multi-dimensional data

- [ADO](#) (Microsoft)
 - ▣ Object-oriented API on top of OLE DB

- [ADO.NET](#)
 - ▣ Evolution of ADO in the .NET framework

ODBC Open DataBase Connectivity

32



DATA MANAGEMENT FOR BUSINESS INTELLIGENCE

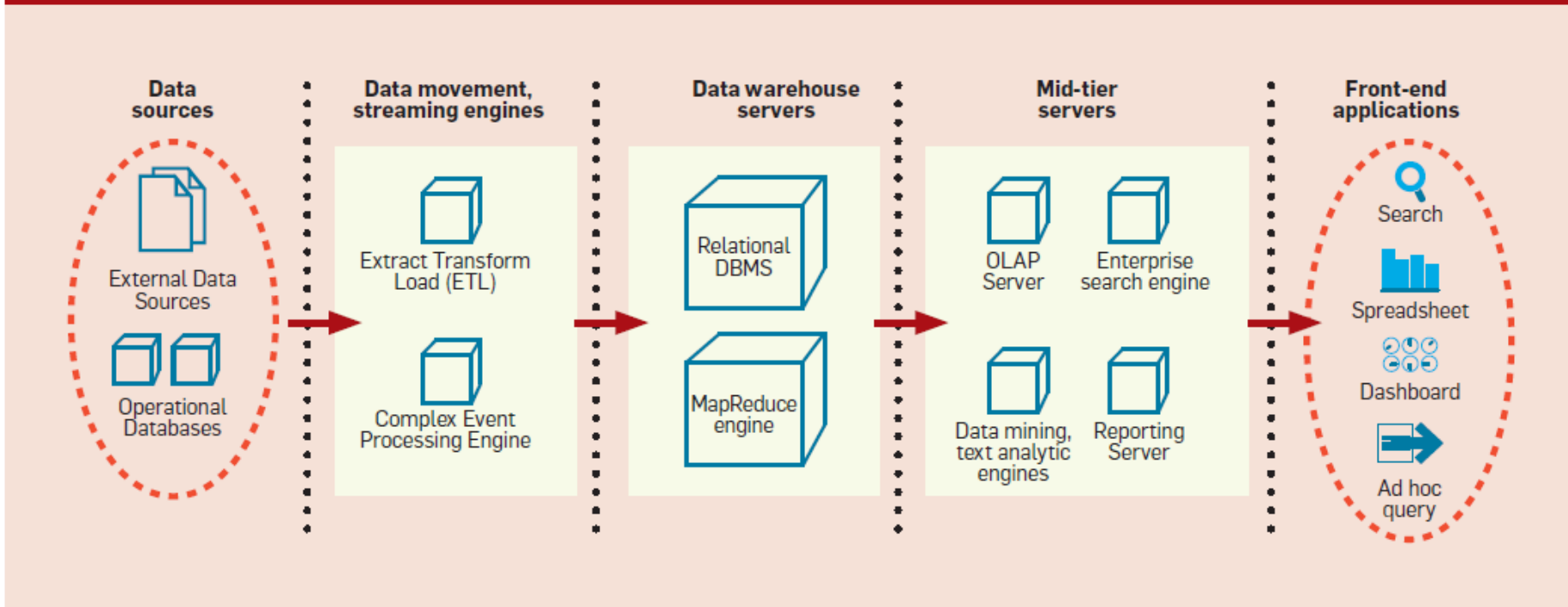
ETL – Extract, Transform and Load

Computer Science Department, University of Pisa

BI Architecture

34

Figure 1. Typical business intelligence architecture.



Extract, Transform and Load

35

ETL (extract transform and load) is the process of extracting, transforming and loading data from heterogeneous sources in a data base/warehouse.

- ▣ Typically supported by (visual) tools.

No.	List of ETL Tools	Version	ETL Vendors
1.	Oracle Warehouse Builder (OWB)	11gR1	Oracle
2.	Data Services	XI 3.2	SAP Business Objects new!
3.	IBM Information Server (Datastage)	9.1	IBM
4.	SAS Data Integration Studio	4.21	SAS Institute new!
5.	PowerCenter	9.0	Informatica
6.	Elixir Repertoire	7.2.2	Elixir
7.	Data Migrator	7.7	Information Builders new!
8.	SQL Server Integration Services	10	Microsoft
9.	Talend Open Studio & Integration Suite	4.0	Talend
10.	DataFlow Manager	6.5	Pitney Bowes Business Insight
11.	Data Integrator	9.2	Pervasive
12.	Open Text Integration Center	7.1	Open Text
13.	Transformation Manager	4.1.4	ETL Solutions Ltd.
14.	Data Manager/Decision Stream	8.2	IBM (Cognos)
15.	Clover ETL	2.9.2	Javlin
16.	Centerprise	5.0	Astera new!
17.	DB2 Warehouse Edition	9.1	IBM
18.	Pentaho Data Integration	4.1	Pentaho
19.	Adeptia Integration Suite	5.1	Adeptia

ETL tasks

36

- **Extract:** access data sources
 - ▣ Local, distributed, file format, connectivity standards

- **Transform:** data manipulation for quality improvment
 - ▣ Selecting data
 - remove unnecessary, duplicated, corrupted, out of limits (ex., age=999) rows and columns, sampling, dimensionality reduction
 - ▣ Missing data
 - fill with default, average, filter out
 - ▣ Coding and normalizing
 - to resolve format (ex., CSV, ARFF), measurement units (ex., meters vs inches), codes (ex., person id), times and dates, min-max norm, ...
 - ▣ Attribute Splitting/merging
 - of attributes (ex., address vs street+city+country)

ETL tasks

37

- Managing surrogate key & Slowly changing dimensions
 - generation and lookup
- Aggregating data
 - At a different granularity. Ex., grain “orders” (id, qty, price) vs grain “customer” (id, no. orders, amount), discretization into bins, ...
- Deriving calculated attributes
 - Ex., margin = sales – costs
- Resolving inconsistencies – record linkage
 - Ex., Dip. Informatica Via Buonarroti 2 is (?) Dip. Informatica Largo B. Pontecorvo 3
- Data merging-purging
 - from two or more sources (ex., sales database, stock database)

ETL tasks

38

□ Load

□ Data staging area

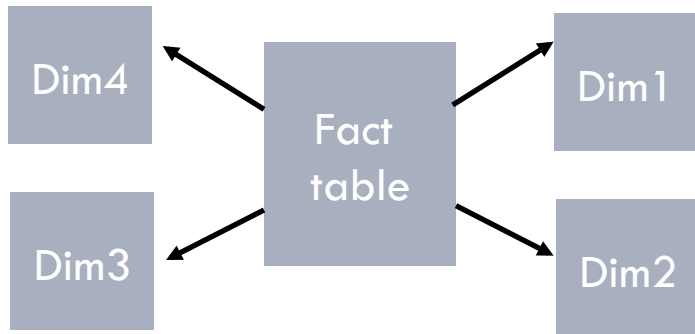
- Area containing intermediate, temporary, partially processed data

□ Types of loading:

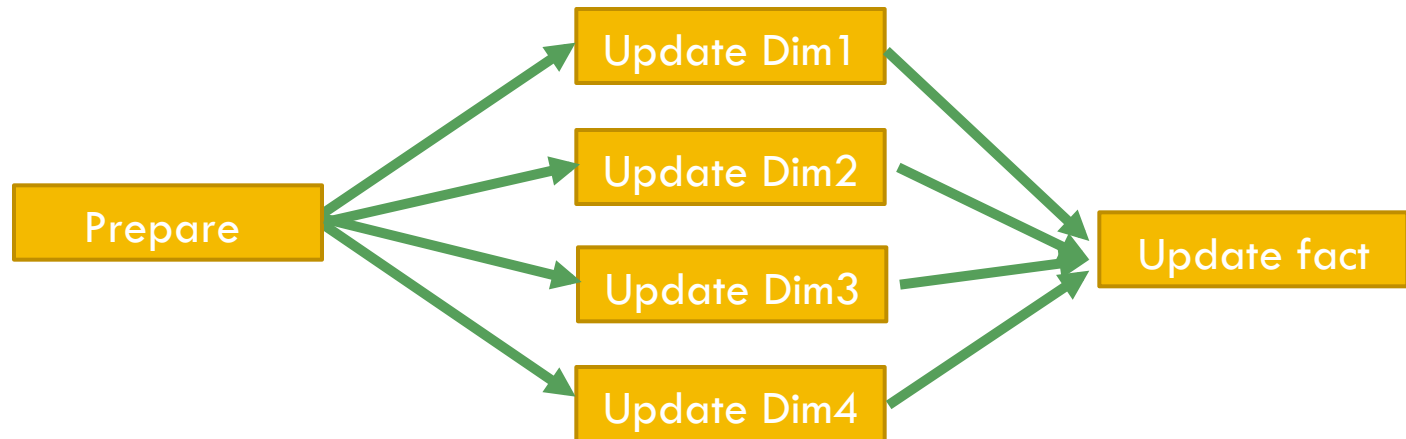
- Initial load (of the datawarehouse)
- Incremental load
 - Types of updates: append, destructive merge, constructive merge
- Full refresh

ETL process for DW

39



Control Flow



BUSINESS INTELLIGENCE

SSIS - SQL Server Integration Services



Background

41

- **SSIS** is a tool for ETL
 - ▣ It can be used independently from SQL Server
 - ▣ Formerly called Data Transformation Services (in SQL Server 2000)

- Docs and samples
 - ▣ Tutorial from Books on Line
 - <http://msdn.microsoft.com/en-us/library/ms141026.aspx>
 - ▣ CodePlex samples
 - <http://www.codeplex.com/SqlServerSamples#ssis>
 - ▣ On-line community
 - <http://sqlis.com>

Developing SSIS projects

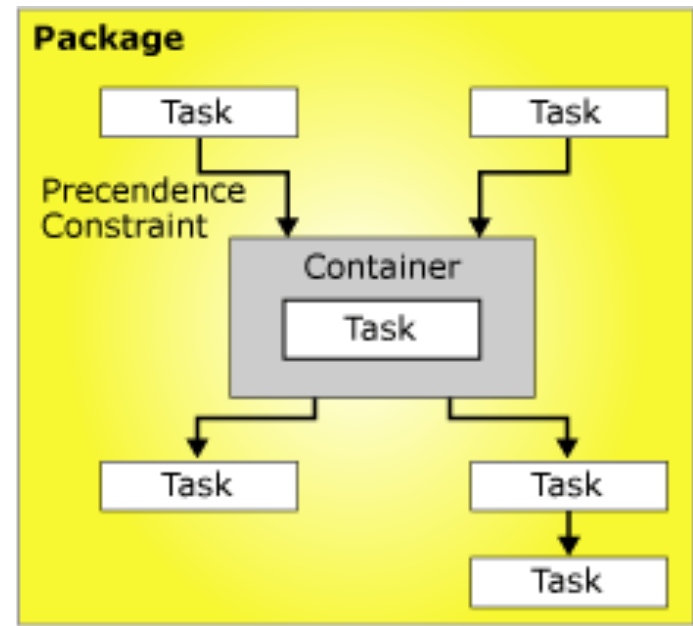
42

- Developer framework
 - ▣ Integrated within SSDT/BIDS
 - Solution = collection of projects
 - Project = developer project (C++, C#, IS, ...)
- Demo
 - ▣ File → New Project → Integration Services
 - ▣ Panels: solution explorer, server explorer, others
 - ▣ SSIS packages (.dtsx extension)
 - Panels: control flow, data flow

Control flow / Jobs

43

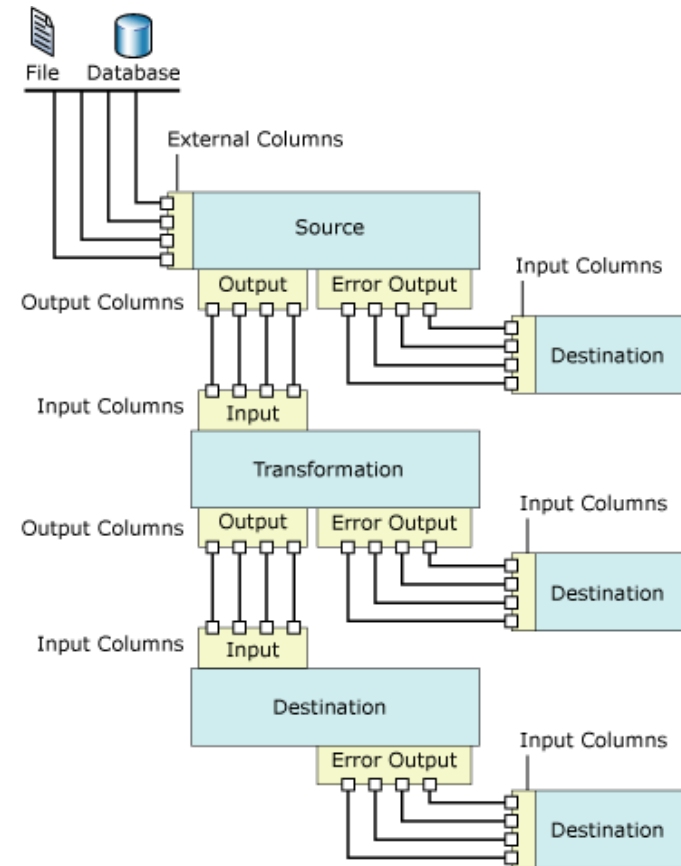
- **Tasks, Containers & Precedence**
 - **Tasks**
 - ETL tasks (list in the Toolbox panel)
 - **Container**
 - Iteration
 - **Precedence**
 - Arrows connecting tasks specify precedence type



Data flow / Transformations

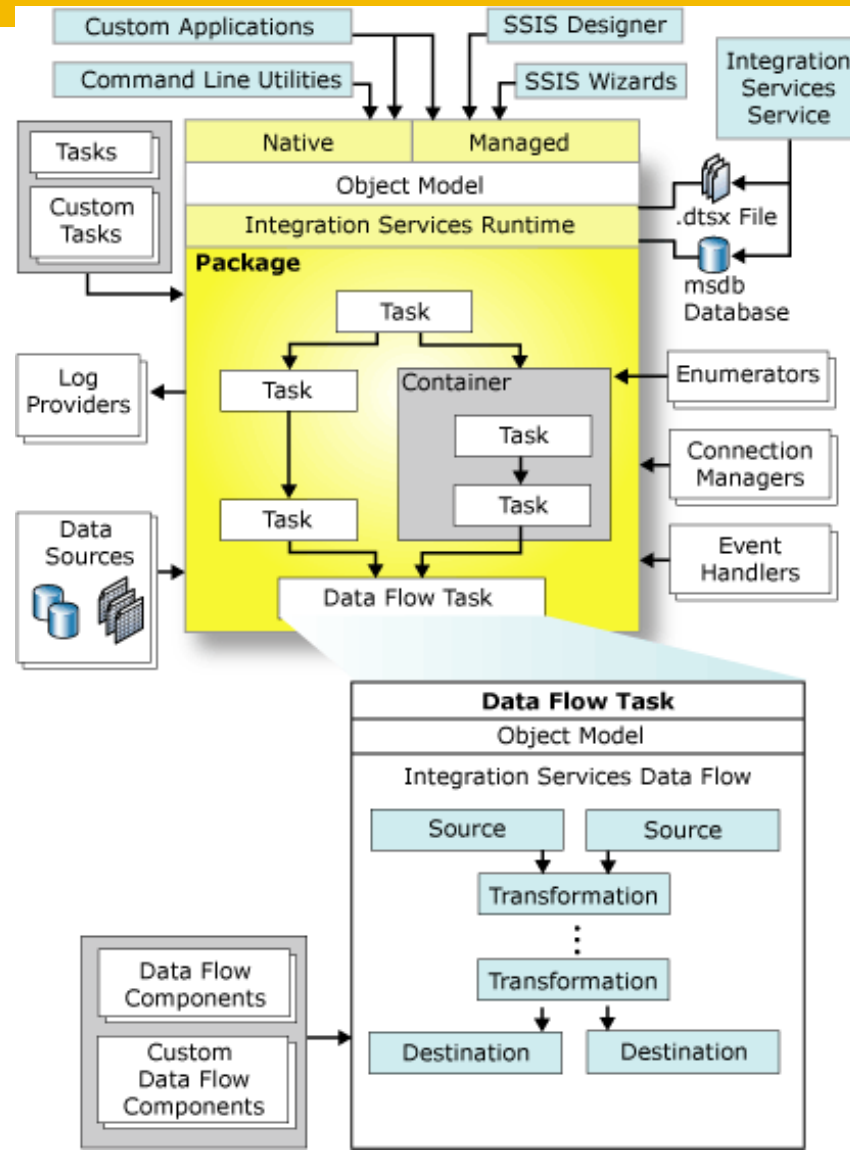
44

- Special tasks
- Define pipelines of data flows from sources to destination
 - ▣ Data flow sources
 - ▣ Data flow transformation
 - ▣ Data destination
 - ▣ Toolbox panel for list



SSIS projects structure

45



SSIS data types

46

- SSIS defines a set of reference data types
 - ▣ As seen for connectivity standards (ODBC, JDBC, OLE DB)
 - ▣ <http://msdn.microsoft.com/en-us/library/ms141036.aspx>
- Data type from sources are mapped into SSIS types
- SSIS transformations works on SSIS types
- SSIS types are mapped to destination data types

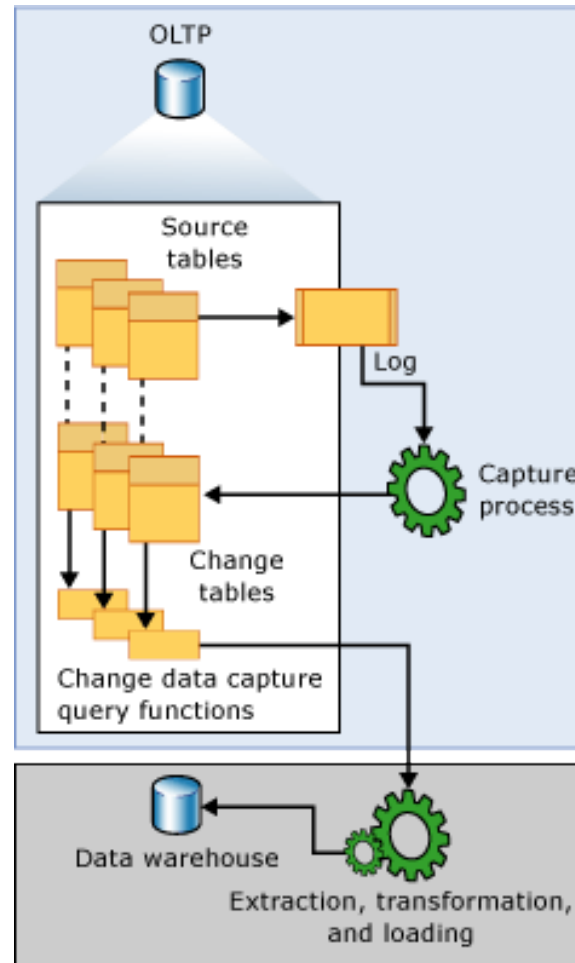
Debug, deployment, scheduling

47

- Debug
 - ▣ Data viewers
- Deployment
 - ▣ Save project on file
 - ▣ Save project on remote SSIS server
 - Project->Deploy
 - ▣ Load project from remote SSIS server
 - File->Add new project->Integration Services Import Project Wizard
- Launch
 - ▣ Local run
 - From Visual Studio
 - From command line: dtexec
 - From explorer: double click on .dtsx files
 - ▣ Remote run on SSIS servers
 - On demand / scheduled

Change data capture

48



BUSINESS INTELLIGENCE LABORATORY

ETL Demo: Pipeline, Sampling and Surrogate Keys



Pipeline

50

- Consider the Foodmart sales database
- Design an ETL project for writing to a CSV file the list of products ordered descending by gain
 - ▣ Gain of a single sale is defined as $(\text{store_sales} - \text{store_cost}) * \text{unit_sales}$
 - ▣ Avg gain of a product is the sum of gains of sales of the product divided by the total units_sales sold
- Do not use views or queries! Do all work in ETL.

SQL SOLUTION

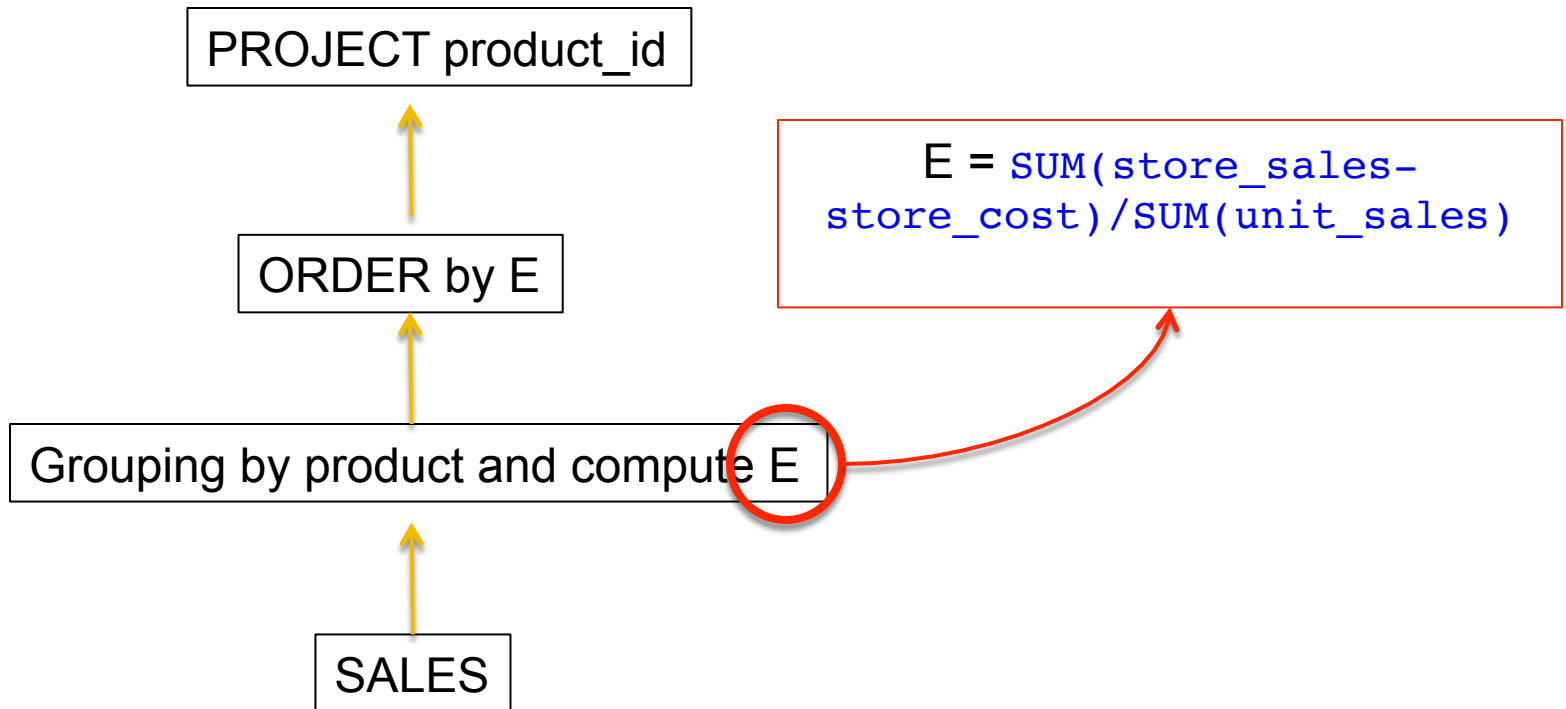
51

```
SELECT product_id
FROM Sales
GROUP BY product_id
ORDER BY SUM(store_sales-store_cost)/
         SUM(unit_sales)
```

... and what about adding Product_name?

BASIC IDEA OF SISS SOLUTION

52



Stratified subsampling

53

- Consider the census table in the *MasterBigData* db
- Design an ETL project for writing to a CSV a random sampling of 30% stratified by sex
 - ▣ 30% of males plus 30% of females
- Do not use views or queries! Do all work in ETL.

BUSINESS INTELLIGENCE LABORATORY

Lab exercise on ETL: SCD

SCD: background

55

□ Slowly Changing Dimensions

- Datawarehouse dimensions members updates
- Three types:
 - Type 1: overwrite previous value
 - Type 2: keep all previous values
 - Type 3: keep last N previous values ($N \sim 1, 2, 3$)
- Each attribute of the dimension can have its own type
 - Type 1: name, surname, ...
 - Type 2: address, ...

SCD: input and output tables

56

- Database FoodMart in SQL Server
- Input
 - ▣ table `customer`
- Output in Lbi database
 - ▣ create a table `customer_dim`
 - columns
 - `surrogate_key` (PK), `customer_id`, `customer_name`, `address`, `date_start`, `date_end`
 - with
 - `surrogate_key` being a surrogate key, `customer_name` including name and surname, `address` made of `address1-city-zip-province-country`, `date_start` and `date_end` are dates

Preliminary step

57

- Develop a SSIS package that adds to `customer_dim` the customers in `customer` that are not already in it

SCD: type 1 updates

58

- Overwrite previous value
- Changes on the input table **customer**
 - On 10/3/2007
 - 231, Mario Rosi, Via XXV Aprile Pisa
 - On 12/3/2007
 - 231, Mario Rossi, Via XXV Aprile Pisa
 - Surname has been corrected

SCD: type 1 updates

59

- The DW `customer_dim` table looks as:
 - ▣ On 10/3/2007, and up to 12/3/2007

surrogate_key, customer_id, name, address, date_start, date_end
874, 231, Mario Rosi, Via XXV Aprile Pisa, 10/3/2007, NULL

- ▣ On 12/3/2007

surrogate_key, customer_id, name, address, date_start, date_end
874, 231, Mario Rossi, Via XXV Aprile Pisa, 10/3/2007, NULL

SCD: type 2 updates

60

- Keep all previous values
- Changes on the input table **customer**
 - On 12/3/2007
 - 231, Mario Rossi, Via XXV Aprile Pisa
 - On 25/9/2008
 - 231, Mario Rossi, Via Risorgimento Pisa
 - Customer has changed his address

SCD: type 2 updates

61

- The DW `customer_dim` table looks as:
 - ▣ On 12/3/2007, and up to 25/9/2008

surrogate_key, customer_id, name, address, date_start, date_end
874, 231, Mario Rossi, Via XXV Aprile Pisa, 10/3/2007, NULL

- ▣ On 25/9/2008

surrogate_key, customer_id, name, address, date_start, date_end
874, 231, Mario Rossi, Via XXV Aprile Pisa, 10/3/2007, 25/9/2008
987, 231, Mario Rossi, Via Risorgimento Pisa, 25/9/2008, NULL

Lab exercise

62

- Design a SSIS project to update `customer_dim` starting from `customer` as follows:
 - ▣ Customers in `customer` that are not in `customer_dim` are added to it
 - ▣ Updates of `customer_name` are of Type 1
 - ▣ Updates of `address` are of Type 2

Sales during travels

63

- A sale in *sales_fact* was done during a travel if the store of the sale was not in the city of residence of the customer. Develop a SSIS package which produces a CSV file with a row for every customer with:
 - the customer full name
 - the total sales to the customer
 - the ratio of sales done during travels

Sales in weekends of previous month

64

- For a given customer and month, the frequency of purchases in weekends (FPW) is the number of distinct weekend days (Saturdays or Sundays) of the **previous** month in which the customer made a purchase. Develop a SSIS package which produces a CSV file with a row for every customer and month with:
 - ▣ the customer full name
 - ▣ the month and year
 - ▣ the customer FPW