# Ab Initio

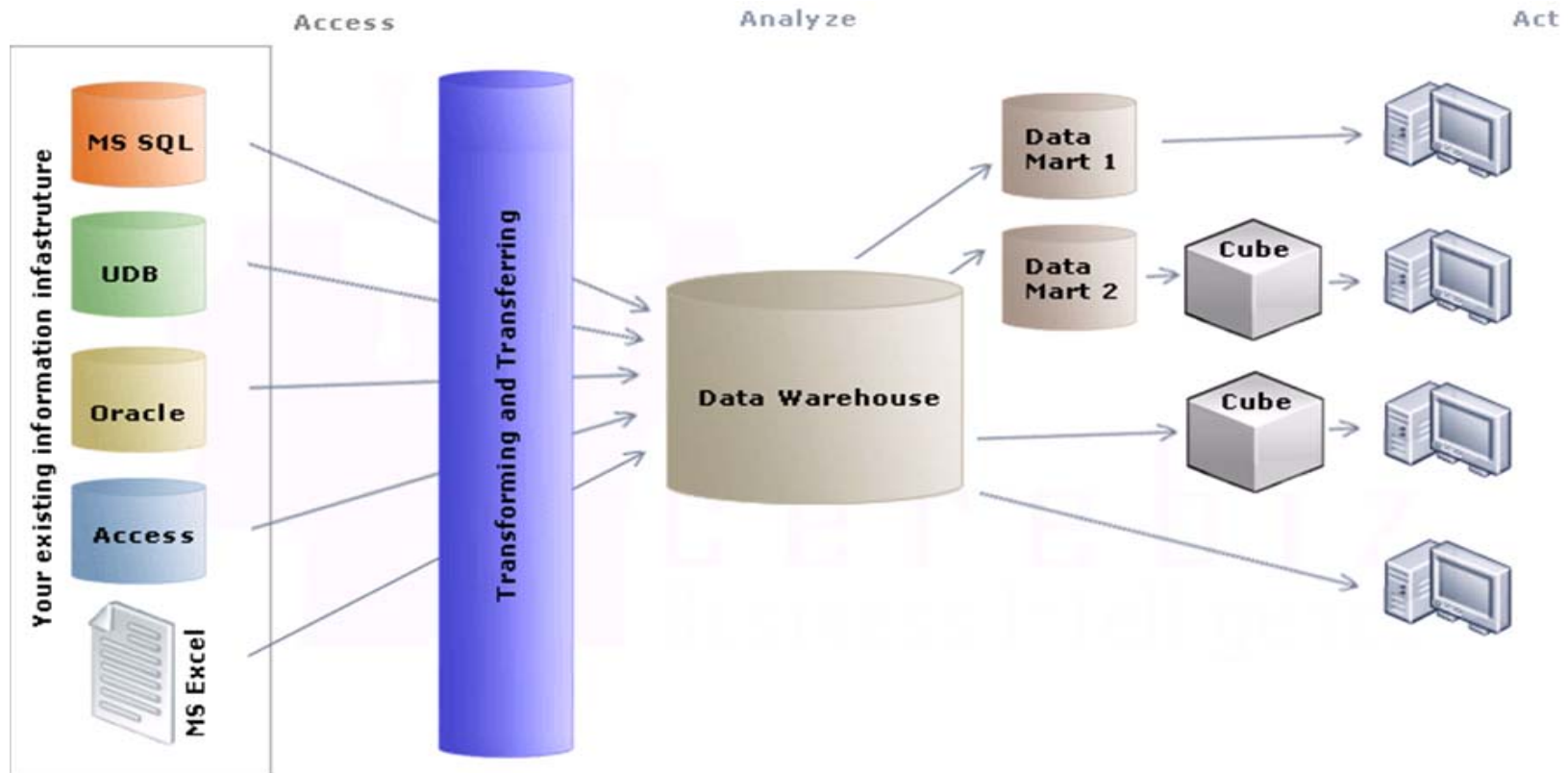## Introduction

- **What is Data warehousing?**

- **Why Data warehousing?**

- **ETL process**

- **Various ETL tools**

- **Introduction about Ab Initio**

- **why Ab Initio**

- **How Unix involved with Ab Initio**

- **GDE window**

- **EME Repository**

- **Sandboxes- User and Standard Sandbox**

- **Ab Initio - Components**

- **Creation of simple graphs**

**TATA**
*Confidential*

**TATA** CONSULTANCY SERVICES
Asia's Largest Global Software & Services Company

# Data warehousing and ETL Process

**Data Warehouse is a collection of 'logical' Data Marts, each of which is designed for a particular line of business i.e. Sales, Marketing( designed to favor/facilitate data analysis and reporting).**

## Example of Source Data

| System Name | Attribute Name | Column Name | Datatype | Values |
|---|---|---|---|---|
| Source System 1 | Customer Application Date | CUSTOMER_APPLICATION_DATE | NUMERIC(8,0) | 11012005 |
| Source System 2 | Customer Application Date | CUST_APPLICATION_DATE | DATE | 11012005 |
| Source System 3 | Application Date | APPLICATION_DATE | DATE | 01NOV2005 |

## Example of Target Data(Data Warehouse)

| Target System | Attribute Name | Column Name | Datatype | Values |
|---|---|---|---|---|
| Record #1 | Customer Application Date | CUSTOMER_APPLICATION_DATE | DATE | 01112005 |
| Record #2 | Customer Application Date | CUSTOMER_APPLICATION_DATE | DATE | 01112005 |
| Record #3 | Customer Application Date | CUSTOMER_APPLICATION_DATE | DATE | 01112005 |

**TATA** CONSULTANCY SERVICES
Asia's Largest Global Software & Services Company

- **Data is first stored temporarily in a 'Staging Table/Area' and is called 'Staging Data' i.e. Data queued for processing.**

- **The processing tool reads the 'Staged Data', performs qualitative processing, filtering, cleansing (As required for the OLAP i.e. reporting/analysis) and finally loads/writes them into Data Warehouse.**

- **All these data flow (both inward and outward) and data processing activities (Extraction from Source System – Transformation of data by cleansing/filtering – Loading into Data Warehouse) are performed using an ETL tool i.e. Ab Initio, Informatica etc.**

- **This entire process is said to be as "ETL process".**

**Extract:**

**The first phase of an ETL process is to extract the data from the source systems. Each separate system may also use a different data organization/ format. Common data source formats are relational databases, and flat files, but other source formats exist. Extraction converts the data into records and columns.**

**Transform: The transform phase applies a series of rules or functions to the extracted data.**

**Examples:**
**Derive a new calculated value     (e.g. sale_amount = qty * unit_price)**
**Summarize multiple rows of data   (e.g. total sales for each region)**

**Load:**
**The load phase loads the data into the data warehouse. Depending on the requirements of the organization, this process ranges widely.**
**Simple-Overwrite old data with new.**
**More complex systems ->Maintenance of history and audit trail of all changes to the data**

| Tool Name | Company Name |
|---|---|
| Informatica | Informatica Corporation |
| DT/studio | Embarcadero technologies |
| Datastage | IBM |
| Abinitio | Abinitio Software corporation |
| Talend | Talend corporation |
| Pentaho | Pentaho corporation |
| Datajunction | Pervasive Software |
| Oracle warehouse builder | Oracle Corporation |
| Microsoft SQl Server Integration | Microsoft |

# Introduction-Ab-Initio

TATA CONSULTANCY SERVICES
Asia's Largest Global Software & Services Company

- **Data processing tool from Ab Initio software corporation (http://www.abinitio.com)**
- **Latin for "from the beginning"**
- **Designed to support largest and most complex business applications**
- **Graphical, intuitive, and "fits the way your business works".**

**Focus:**
  **Moving Data -**
  **Move small and large volumes of data in an efficient manner.**
          **Deal with the complexity associated with business data.**
  **High Performance**
              **Scalable Solutions**
    **Better productivity**
**Usage:**
  - **Data Warehousing**
  - **Batch Processing**
  - **Data Movement**
  - **Data Transformation**

| Graphical Development Environment (GDE) |
|---|

SSH
REXEC
TELNET
DCOM

| Co-operating System (Co>Ops) |
|---|

EME
DB
Conduct>IT
CF

| Product | Functionality |
|---|---|
| GDE | User Interface for creating Graphs and Plans in Ab Initio |
| Data Profiler | Ab Initio Tool for Data Profiling |
| Co>Ops | Server Component for running deployed Ab Initio programs |
| EME | Ab Initio Technical Repository – Part of Co>Ops Install |
| Database | Ab Initio Server Database Components |
| Conduct>IT | Ab Initio Server Component for running Ab Initio Plans |
| Continuous Flow | Ab Initio Server Components for running CF programs |

- All server components are installed by default

- **"AB_HOME"** refers to installation location of Ab Initio
- Various Connectors and Plugins installed in "**AB_HOME/Connectors** & **AB_HOME/plugins"** location

All binaries and library files available in "**AB_HOME/bin** & **AB_HOME/lib"** respectively

User Applications

Development Environments

*GDE*      Shell

Ab Initio

EME

Component Library

User-defined Components

3rd Party Components

The Ab Initio Co>Operating® System

Native Operating System (Unix, Windows, OS/390)

## Host Machine 1

### Unix Shell Script or NT Batch File

- **Supplies parameter values to underlying programs through arguments and environment variables**

- **Controls the flow of data through pipes**

- **Usually generated using the GDE**

### GDE

- **Ability to graphically design batch programs comprising Ab Initio components, connected by pipes**

- **Ability to test run the graphical design and monitor its progress**

- **Ability to generate a shell script or batch file from the graphical design**

### Co>Operating System

**Ab Initio Built-in Component Programs (Partitions, Transforms etc)**

**User Programs**

**Operating System**

**( Unix , Windows NT )**

## Host Machine 2

**Co-Operating System**

**User Programs**

**Operating System**

* **Co>Operating System**

  ❑ **Layered on the top of the operating system.**

  ❑ **Unites a network of computing resources – CPUs, storage disks, programs, datasets into a data-processing system with scalable performance.**

* **GDE**

  ❑ **can talk to the Co-operating system using several protocols like Telnet, Rexec and FTP**

  ❑ **It is GUI for building applications in Ab Initio**

**Graph**

- **is the logical modular unit of an application.**

- **consists of several components that forms the building blocks of an Ab Initio application**

- **Start Script (Host Setup) - Local to the Graph**

- **End Script - Local to the Graph**

**Component**

- **is a program that does a specific type of job and can be controlled by its parameter settings. Ex: Join, Re-format etc**

**Component Organizer**

- **Groups all components under different categories.**

**Setup Command**

- **Ab Initio Host (AIH) file**

- **Builds up the environment to run an Ab Initio application.**

* **Files**

* **Formats**

* **Components**

* **Flows**

* **Layouts**

* **Building with *mp job***

* **Building with *mp run***

# A Sample Graph …

## A Sample Graph ...



Expression Metadata

Record format metadata

Ports

Layout

✴ **A graph, after development, is deployed to the back-end server as a Unix shell script or Windows NT batch file.**

✴ **This becomes the executable to run at the back-end with the help of the Co-operating system.**

✴ **The execution can be done from the GDE itself or manually from the back-end**

✴ **Ab Initio runtime environment is different from the development environment.**

- ❑ **Unix serves as backend for Ab-initio.**

- ❑ **All the graphs/Jobs in Ab-initio can be accessible through Unix(backend)**

- ❑ **Putty connectivity**

**Environment – Quick Overview:**

- ❑ **$AI_RUN,$AI_BIN—run directory, .ksh scripts
  $AI_PLAN, $AI_SERIAL_<LOG/ERROR/SUMMARY>**

- ❑ **$AI_DML—record format files**

- ❑ **$AI_XFR—transform files**

- ❑ **$AI_MP—graphs**

- ❑ **$AI_DB—database config files**

- ❑ **$AI_SERIAL - serial source data, other serial data**

- ❑ **$AI_MFS - Ab Initio multifile directory – in training will also contain partition directories (more about this later!)**

- ❑ **$AI_LOG - A location to place logging files, etc**

**Standard Sandbox**

Project

| V5 | File 1 |
|----|--------|
| V6 | File 2 |
| V3 | File 3 |

• Sandboxes are work areas used to develop, test or run code associated with a given project. Only one version of the code can be held within the sandbox at any time.

• The EME Datastore contains all versions of the code that have been checked into it.

*Check-out*

**User Sandbox**

Project1

| V5 | File 1 |
|----|--------|

*Check-in* →

*Check-out* ←

**Project in EME Datastore**

Project1

| V1 | File 1 |
|----|--------|
| V2 | File 1 |
| V3 | File 1 |
| V4 | File 1 |
| V5 | File 1 |

# The GDE window

The GDE window is made up of several panes and toolbars, as shown in the following figure:



Component Organizer    Sandbox View

GDE main menu
GDE toolbar
Workspace
Drawing canvas
Status bar
Component description pane

**How a job runs**

❑ The execution of an Ab Initio graph is a job.

❑ To run a job, need to invoke a shell script that the GDE generates from a graph.

❑ The script process initiates job processes that control the execution of the programs represented by the graph.

❑ Graph->mp/graph1.mp ; Shell script->run/graph1.ksh

```
0 1 0 -i in0 0 1   -i in1 0 0  {sales.str_nbr; sales.sls_dt; sales.rgstr_nbr; sales.pos_trans_id; sales.crt_ts}
abstbat   13986 13859 10 06:52 ?        00:00:00 /opt/isv/abinitio/ab2.15.8/abinitio-V2-15/bin/unitool sort {pvndr_nbr} 1006632
96
abstbat   13994 13859  7 06:52 ?        00:00:00 /opt/isv/abinitio/ab2.15.8/abinitio-V2-15/bin/unitool sort {pvndr_nbr} 1006632
96
abstbat   13995 13859 17 06:52 ?        00:00:00 /opt/isv/abinitio/ab2.15.8/abinitio-V2-15/bin/unitool reformat_transform 0 0.0
00000 !is_blank(pvndr_nbr) 0 -x out0 Reformat_Filter_out_blank_PNVDR_7_PKG.reformat
abpobat   14475  1654  9 06:52 ?        00:00:00 /opt/isv/abinitio/ab2.15.8/abinitio-V2-15/bin/unitool logger /opt/hd/eg/admin/
wh/td_dist/log/./run_sql_header_update_dc_inb_case_1916407.log Start End
abpobat   14497 18644 61 06:52 ?        00:00:00 /opt/isv/abinitio/abinitio-V3-1-2/bin/unitool sort {join_key_1; join_key_2} 10
0663296
abpobat   14498 18627 64 06:52 ?        00:00:00 /opt/isv/abinitio/abinitio-V3-1-2/bin/unitool sort {join_key_1; join_key_2} 10
0663296
abpobat   14499 18646 56 06:52 ?        00:00:00 [unitool] <defunct>
abpobat   14501 18625 60 06:52 ?        00:00:00 /opt/isv/abinitio/abinitio-V3-1-2/bin/unitool sort {join_key_1; join_key_2} 10
0663296
abpobat   14502 18630 55 06:52 ?        00:00:00 /opt/isv/abinitio/abinitio-V3-1-2/bin/unitool sort {join_key_1; join_key_2} 10
0663296
abpobat   14503 18628 46 06:52 ?        00:00:00 /opt/isv/abinitio/abinitio-V3-1-2/bin/unitool sort {join_key_1; join_key_2} 10
0663296
bissupt   14542 13743  0 06:52 pts/12   00:00:00 grep unitool
abstbat   17204 16505  7 06:17 ?        00:02:29 /opt/isv/abinitio/abinitio-V3-1-2/bin/unitool broadcast
```

**You can invoke the script in two ways:**

❑      **From the GDE**

❑      **From a command line**

❑      **To invoke the script from the GDE, click the Run button or choose Run > Start from the GDE menu bar.**

❑      **To invoke the script through command line,**

❑      **For bin script: ksh scriptname.ksh in bin path.**

❑      **To run a graph from backend: $AI_RUN Graphname.ksh parameters(if needed) in run path.**

# Components - Overview

Component Organizer

- Datasets
  - MVS
  - Block-Compressed Lookup Template
  - Input File
  - Input Table
  - Intermediate File
  - Lookup File
  - Lookup Template
  - Output File
  - Output Table
  - Read Input File Slices
  - Read Multiple Files
  - Read Shared
  - Write Multiple Files
- Departition
- Deprecated
- Examples
- FTP
- Miscellaneous
- Partition
- Sort
- Transform
  - Aggregate
  - Combine
  - Dedup Sorted
  - Denormalize Sorted

Transform components modify or manipulate data records by using one or more transform functions.

Path: I:\Program Files\Ab Initio\Ab Initio GDE 1.15.10.1\Components\Transform\

**TATA CONSULTANCY SERVICES**
Asia's Largest Global Software & Services Company

**An application script**

Shown below is the entire script (**classify-customers.mp**) needed to implement the customer-classification program we have been examining:

```
#! /bin/ksh
set -e

mp job classify-customers


#   Component           Name                Arguments
#   ---------           ------              ---------
mp ifile   customers            file:customers
mp ofile   telesales            file:telesales
mp ofile   mailing-list         mfile:/mfs/mlist

mp hash-partition partition      "zipcode" -layout customers
mp classifier classify           -layout mailing-list
mp householder household         -layout mailing-list
mp gather gather        -layout telesales


#   Define record format
mp metadata             format              -file customers.dml


#   Flow type       Name From Port          To Port             Data format
#   ---------       ---- ---- ----          -------             ----------
mp straight-flow f1   customers.read partition.in     -metadata format
mp fan-out-flow  f2   partition.out  classify.input   -metadata format
mp straight-flow f3   classify.mail  household.in     -metadata format
mp straight-flow f4   household.out  mailing-list.write -metadata format
mp straight-flow f5   classify.phone gather.input     -metadata format
mp fan-in-flow   f6   gather.output  telesales.write -metadata format

# Run the Application
mp run
```

TATA CONSULTANCY SERVICES
Asia's Largest Global Software & Services Company

Double click on a component to bring up its Properties Page

Click on the Ports Tab to view the 'Port(s)' Properties

* **DML**
  * **Ab Initio stores metadata in the form of record formats.**
  * **Metadata can be embedded within a component or can be stored external to the graph in a file with a ".dml" extension.**
* **XFR**
  * **Data can be transformed with the help of transform functions.**
  * **Transform functions can be embedded within a component or can be stored external to the graph in a file with a ".xfr" extension.**

**DML Syntax**

✳ **Record types begin with *record* and end with *end***

✳ **Fields are declared: *data_type*(*len*) *field_name*;**

✳ **Field names consist of letters(a...z,A...Z),digits(0...9) and underscores(_) and are *Case sensitive***

✳ **Keywords/Reserved words are *record*, *end*, *date*....**

**Some of the Data Types available**

✳ **String**

✳ **Decimal**

✳ **Integer**

　✧ **Storing Data in binary form**

✳ **Date and Datetime**

✳ **EBCDIC and ASCII records**

✳ **Null in Ab Initio - Non-existence of column values.**

File data:
id|name|domain
1234|bh|retail
12|vars|telecom
154|abcdefghij|retai|

File Data Model:
decimal id("|") not null;
string name("|") not null;
string domain("\n") null;

DML

record

decimal("|") id;
string("|") name;
String("\n") domain;
end

# DML format created for a data

```
0345John    Smith
0212Sam     Spade
0322Elvis   Jones
0492Sue     West
0121Mary    Forth
0221Bill    Black
```

Record Format Editor - port read of transform Simple - simple.dml

File   Edit   View   Help

Fields | Functions

| Field Name | Data Type | Length |
|---|---|---|
| Top-level Record | record | |
| id | decimal | 4 |
| first_name | string | 6 |
| last_name | string | 6 |
| | | |

## DML creation



**Field name**     **Field type**     **Field length**

View… Attributes.

Length can be delimiter string

Date format goes here

Field Type drop-down

**Ab Initio**

* **User-defined function producing one or more output from one or more input**

* **Associated with transform components**

* **Rules that computes expression from input values and local variable and assigns the result to output objects**

* **Syntax**

  ◇ **Functions :**

**output-records** *: :* **function-name (input-records)** **=**

*begin*

**assignments**

*End;*

**Assignments :**

**Direct Mapping without any transformation: out.\* :: in.\***

```
> ls
ab-initio-cmd-CPNTTAA3-50506579.ksh   abtalendstudy_POC_0.2   dd_signs_output   inp_filter.dat
abtalendstudy_POC_0.1                 dd_signs_lookup         file_chk.sh       xml
[ 2.15.7.8  ] bxm8666@cpliad50.homedepot.com:/home/bxm8666
> touch outp.dat
[ 2.15.7.8  ] bxm8666@cpliad50.homedepot.com:/home/bxm8666
> cat inp_filter.dat
318922|Bh|Retail
675790|ha|Retail
343434|vi|Telecom
998989|Mo|Banking
454556|va|Telecom
232345|vn|Banking
[ 2.15.7.8  ] bxm8666@cpliad50.homedepot.com:/home/bxm8666
```

**TATA CONSULTANCY SERVICES**
Asia's Largest Global Software & Services Company

✴**Serial or Multifiles**

✴**Held in main memory**

✴**Searching and Retrieval is key-based and faster as compared to files stored on disks**

✴**associates key values with corresponding data values to index records and retrieve them**

✴**Lookup parameters**

   ✧ **Key**

   ✧ **Record Format**

❑ **Filter by Expression**

❑ **Reformat**

❑ **Redefine Format**

❑ **Sort**

❑ **Join**

❑ **Replicate**

❑ **Dedup**

❑ **Rollup**

- **Reads record from *input* port**

- **Evaluate the *select_expr***

- **If result is true, record written to *out* port**

- **If result is false, record written to *deselect* port**

**TATA CONSULTANCY SERVICES**

Asia's Largest Global Software & Services Company

✳ **REJECT**

  ✧ **Input records that caused error**

✳ **ERROR**

  ✧ **Associated error message**

✳ **LOG**

  ✧ **Logging records**

1. **Reads record from *input* port**

2. **Record passes as argument to *transform function or xfr***

3. **Records written to *out* ports, if the function returns a success status**

4. **Records written to *reject* ports, if the function returns a failure status**

5. **Parameters of Reformat Component**

   ✴ **Count**

   ✴ **Transform (Xfr) Function**

   ✴ **Reject-Threshold**

   ◇ **Abort**

   ◇ **Never Abort**

   ◇ **Use Limit & Ramp**

   ⇨ **Limit – Number of errors to tolerate**

   ⇨ **Ramp – Scale of errors to tolerate per Input**

A drop-down menu specifying the number of errors to tolerate.

# Reformatted output

**Sort Component**

✸ **Reads records from input port, sorts them by key, writes result to output port**

✸ **Parameters**

    ✧ **Key**

    ✧ **Max-core**

**Keys**

**A key identifies a field or set of fields to organize a dataset**

    ✧ **Single Field: *employee_number***

    ✧ **Multiple field or Composite key: *(last_name; first_name)***

    ✧ **Modifiers: *employee_number descending***

*Max-core: Maximum memory usage in bytes*

1. **Reads records from multiple input ports**

2. **Operates on records with matching keys using a multi-input transform function**

3. **Writes result to the output port**

*PORTS*

- **in**
- **out**
- **unused**
- **reject (optional)**
- **error (optional)**
- **log (optional)**

*PARAMETERS*

- **count**
- **key**
- **override key**
- **transform**
- **limit**
- **ramp**

**Rollup evaluates a group of input records that have the same key, and then generates records that either summarize each group or select certain information from each group.**

**Parameters:**

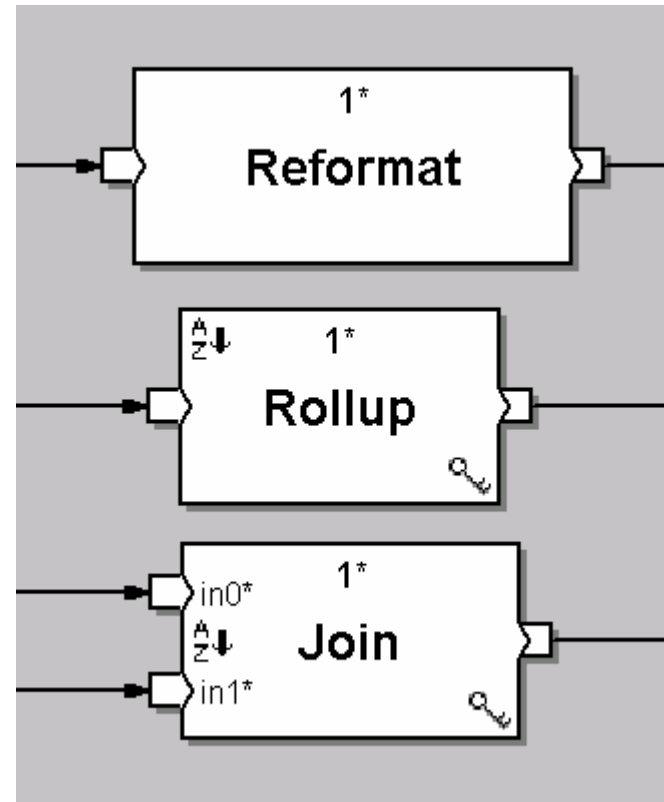| | |
|---|---|
| **check-sort,sorted input** | **limit,Ramp** |
| **logging** | **log_group** |
| **log_input** | **log_intermediate** |
| **log_output** | **grouped-input** |
| **error_group** | **key** |
| **key-method** | **major-key** |
| **log_reject** | **max-core** |

- **The following aggregation functions are predefined and are only available in the rollup component:**

  - ❑ **avg**

  - ❑ **max**

  - ❑ **min**

  - ❑ **count**

  - ❑ **first**

  - ❑ **Product**

  - ❑ **last**

  - ❑ **Sum**

- **Multi-stage Transform – initialize,iterate,finalize,use of variables**

Note the use of an aggregation function in the expression

**In these components the record format metadata does not change from input to output**

**In these components the record format metadata typically changes (goes through a transformation) from input to output**

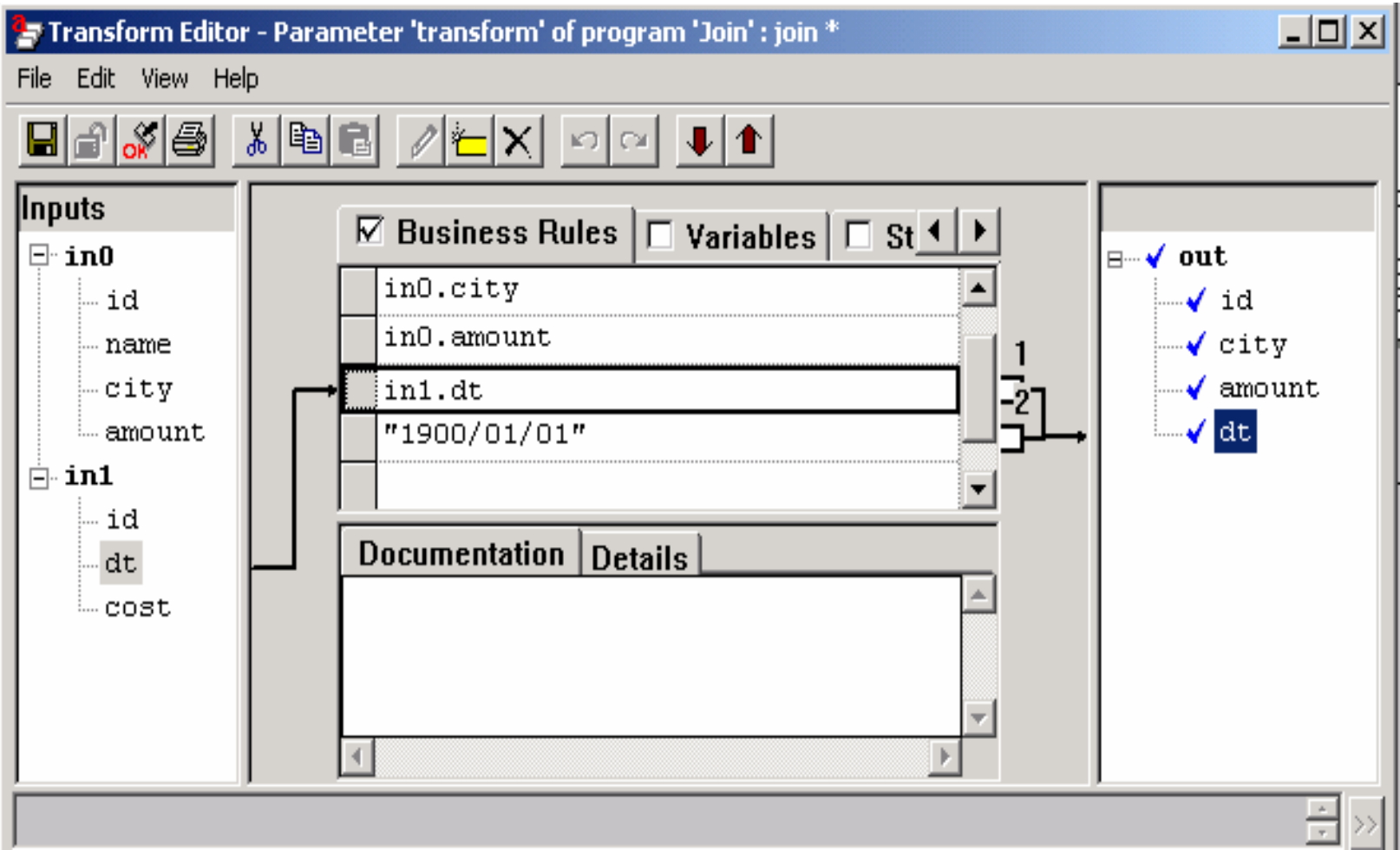The *Priority* is the order of evaluation of rules in a transform function.
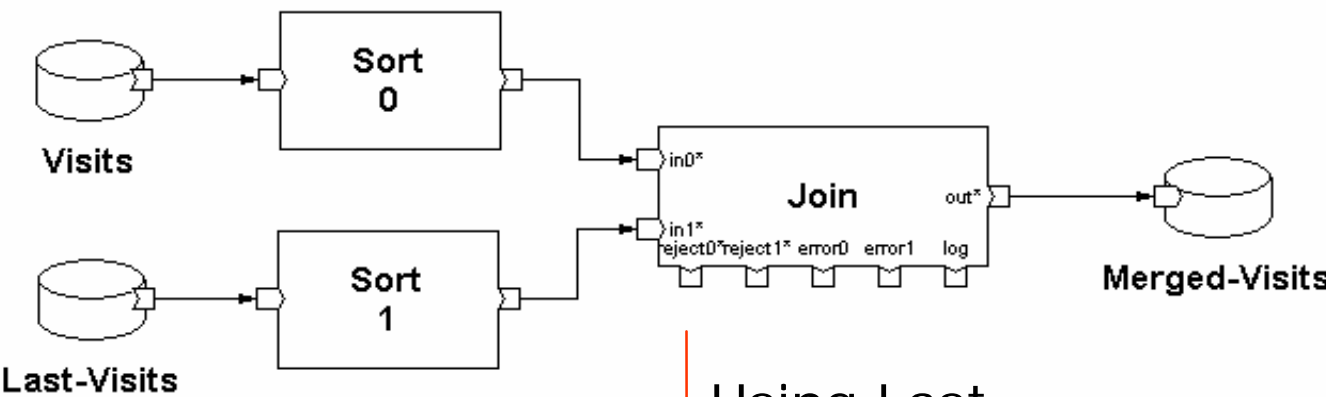
**An example**

A join component may have a transform function with    prioritized rules as
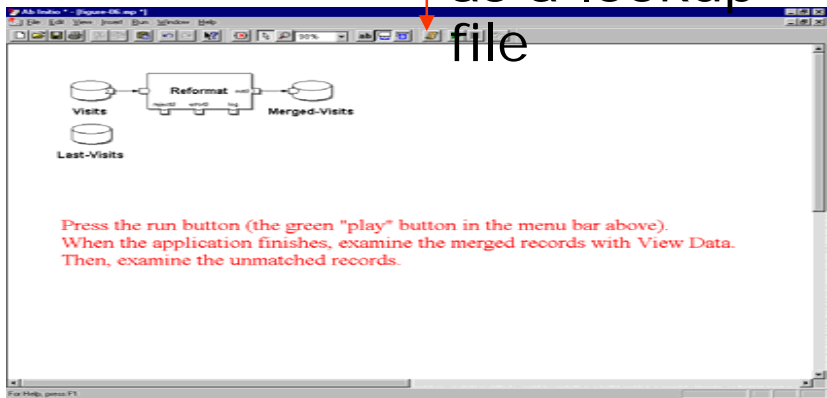
out.ssn :1: in1.ssn;

out.ssn :2: in2.ssn;

out.ssn :3: "999999999";

**TATA** CONSULTANCY SERVICES
Asia's Largest Global Software & Services Company

Using Last-
Visits
as a lookup
file

Input 0 record format:

```
record
  decimal(4) id;
  string(6) name;
  string(8) city;
  decimal(3) amount;
end
```
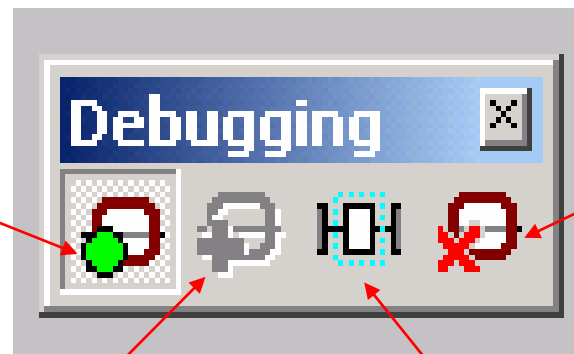
Output record format:

```
record
  decimal(4) id;
  string(8) city;
  decimal(3) amount;
  date("YYYY/MM/DD") dt;
end
```

Transform function:

```
out :: lookup_info(in) =
begin
  out.id      : : in.id;
  out.city    : : in.city;
  out.amount  : : in.amount;
  out.dt      :1 : lookup("Last-Visits", in.id).dt;
  out.dt      :2 : "1900/01/01";
end;
```

❑ The GDE has a built in debugger capability

❑ To enable the Debugger, Debugger:Enable Debugger

❑ The Debugger Toolbar

**Enable Debugger**

**Remove All Watchers**

**Add Watcher File**

**Isolate Components**

- **Data transformation in multiple stages following several sets of rules**

- **Each set of rule form one transform function**

- **Information is passed across stages by temporary variables**

- **Stages include initialization, iteration, finalization and more**

- **Few multistage components are aggregate, rollup, scan**

**Aggregate/Rollup/Scan**

- **Generates summary records for group of input records**

## *Input Table

- unloads data records from a database into an Ab Initio graph
- Source : DB table or SQL statement to SELECT from table

## *Output Table

- loads data records into a database
- Destination : DB table or SQL statement to INSERT into table

## *Update Table

- executes UPDATE or INSERT statements in embedded SQL format to modify a DB table

**\* Join with DB**

**\* Truncate Table**
- **Deletes all the rows in a specified DB table**

**\* Run SQL**
- **Executes SQL statements in a DB**

**Ab Initio built-in functions are DML expressions that**

- ✧ **can manipulate strings, dates, and numbers**
- ✧ **access system properties**

**Function categories**

- ✧ **Date functions : now(), today(), date_to_int(), ..**
- ✧ **Inquiry and error functions: is_defined(), is_valid(), force_error(), ..**
- ✧ **Lookup functions: lookup(), lookup_local(), ..**
- ✧ **Math functions: ceiling(), floor(), ..**
- ✧ **Miscellaneous functions:decimal_round(), hash_value(), ..**
- ✧ **String functions: string_substring(), is_blank(), ..**

| Name | Description |
|------|-------------|
| Normalize | ✴Generates multiple data records from each input data record<br>✴Separate a data record with a vector field into several individual records, each containing one element of the vector. |
| Denormalize Sorted | ✴Consolidates groups of related data records into a single output record with a vector field for each group<br>✴Requires Grouped Input |
| Validate Records | Separates valid data records from invalid data records |
| Check Order | Tests whether data records are sorted according to a key-specifier. |
| Compare Records | Compares data records from two flows one by one |
| Generate Records | Generates a specified number of data records with fields of specified lengths and types. |
| Gather Logs | Collects the output from the log ports of components for analysis of a graph after execution |
| Sample | Selects a specified number of data records at random from one or multiple input flows |

❑ **Mechanism by which some or all constituents of an application –**

**datasets and processing modules are replicated into a number of**
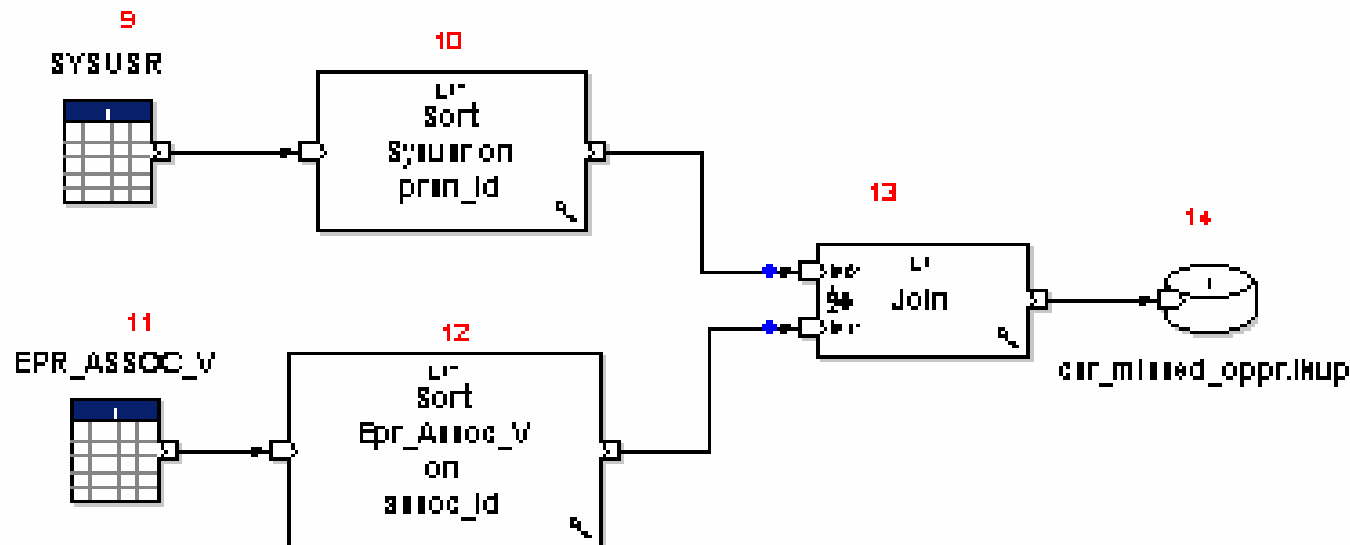
**partitions, each spawning a process.**

❑ **This makes the Ab initio to process considerable huge volume (in**

**millions) of records with an optimum usage of hardware available.**

**The power of Ab Initio lies in the fact that it can process data in parallel**
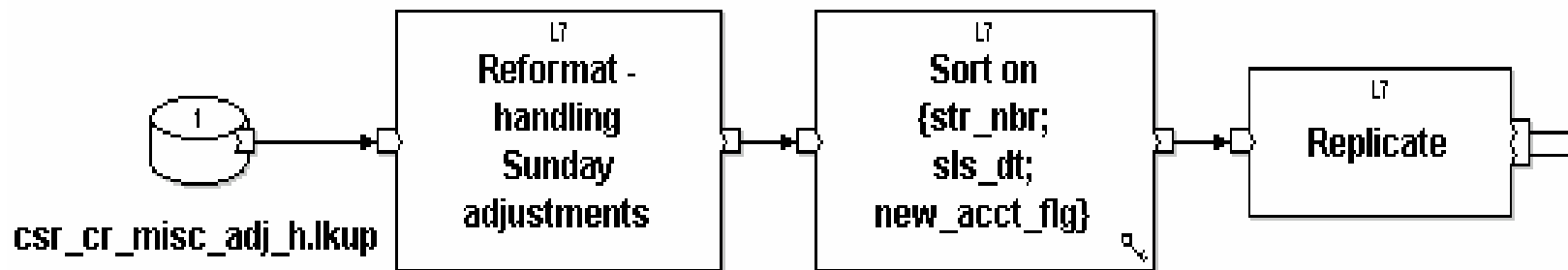
**runtime environment**

❑ **Types of Parallelism**

● **Component Parallelism**

● **Pipeline Parallelism**

● **Data Parallelism**

Component Parallelism is achieved when different instances of same component run on separate data sets. Component parallelism scales to the number of branches of a graph — the more branches a graph has, the greater the component parallelism. <u>If a graph has only one branch, component parallelism cannot occur</u>.
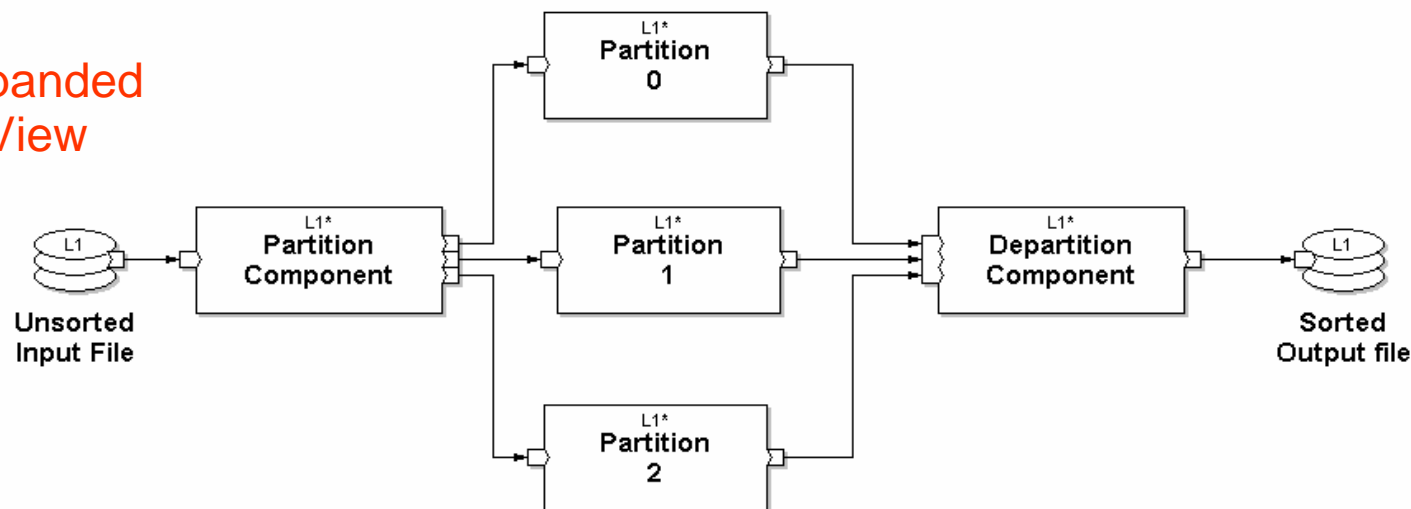
**Pipeline parallelism occurs when several connected program components on the same branch of a graph execute simultaneously. In this kind the two processing stages of the graph run concurrently.**
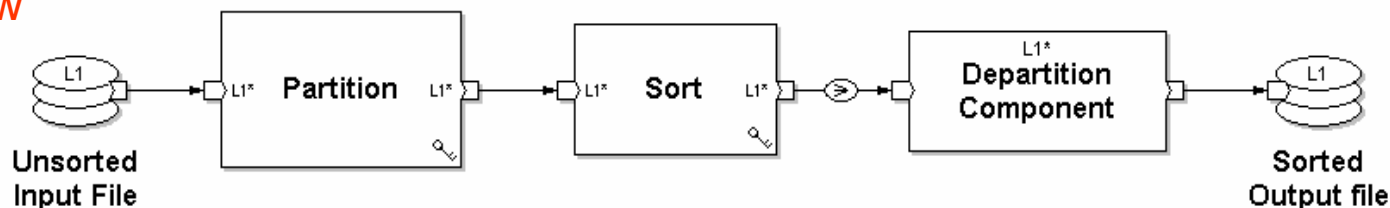
**When data is divided into segments or *partitions* and multiple instances of program components run simultaneously on each *partition***

Expanded View

Linear View

**Multifiles**

✳**A global view of a set of ordinary files called *partitions* usually located on different disks or systems**

✳**Ab Initio provides shell level utilities called *"m_ commands"* for handling multifiles (copy, delete,move etc.)**

✳**Multifiles reside on Multidirectories**

✳**Each is represented using URL notation with *"mfile"* as the protocol part:**

⇨  **mfile://pluto.us.com/usr/ed/mfs1/new.dat**
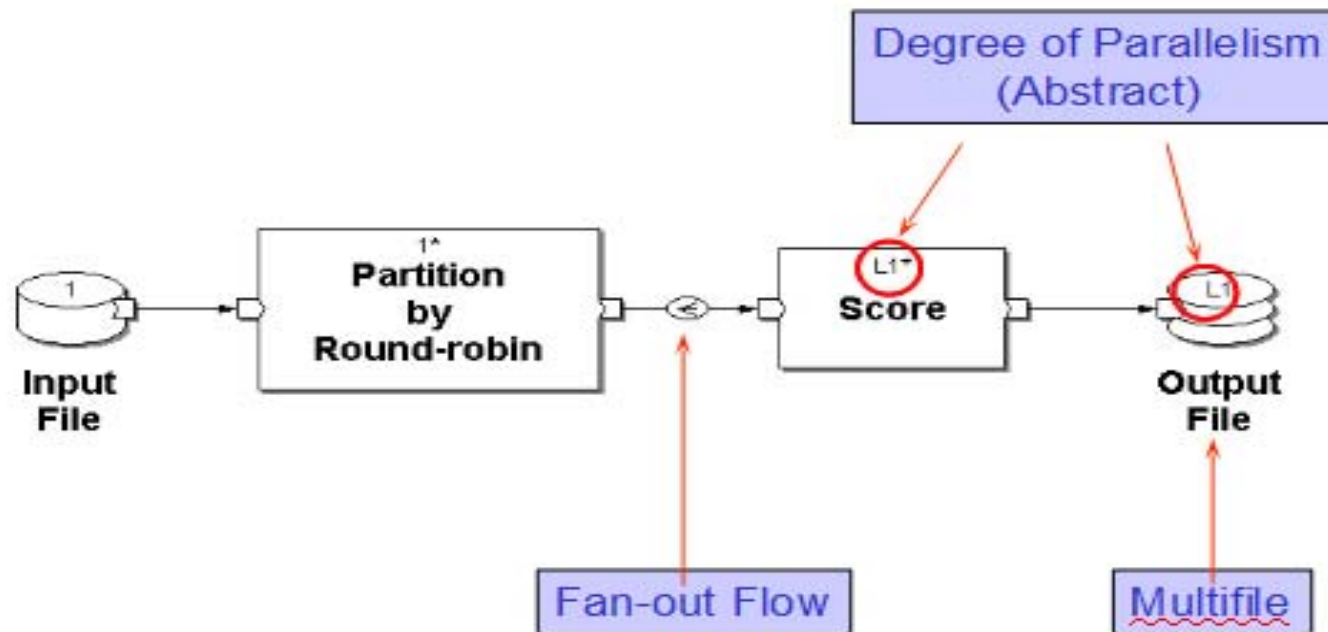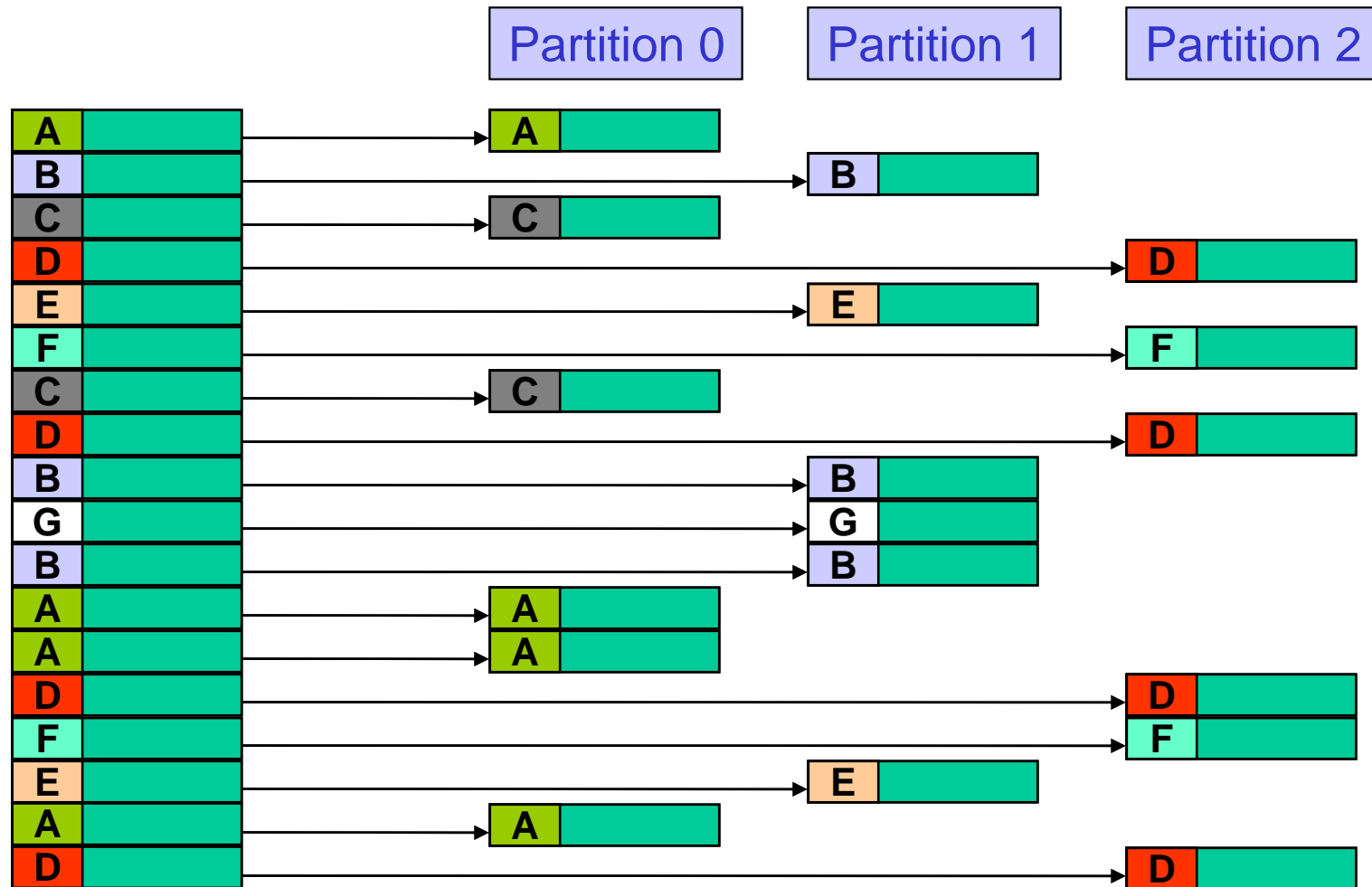
**Data can be partitioned using**

- **Partition by Round-robin**

- **Partition by Key**

- **Broadcast**

- **Partition by Expression**

- **Partition by Range**

- **Partition by Percentage**

- **Partition by Load Balance**

**Writes records to each partition evenly**
**Block-size records go into one partition before**
**moving on to the next.**

Partition 0      Partition 1      Partition 2

A hash code computed using the key determines which partition a record will be written on, meaning that records with the same key value will go to the same partition

✴**Gather**

  ✧ **Reads data records from the flows connected to the input port**

  ✧ **Combines the records arbitrarily and writes to the output**

✴**Concatenate**

  ✧ **Concatenate appends multiple flow partitions of data records one after another**

✴**Merge**

  ✧ **Combines data records from multiple flow partitions that have been sorted on a key**

  ✧ **Maintains the sort order**

**Phasing:**
- Breaking an application into phases limits the contention for
    Main memory.
    Processor(s).
- Breaking an application into phases cost
    Disk space.

**Checkpoint - Purpose:**
- ✧ Provide same functionality as phase
- ✧ Additional: Provide restart capability

**How does it work ?**
- ✧ At job start, output datasets are copied to temporary files (in .WORK-serial or .WORK-parallel directories)
- ✧ At checkpoint completion, intermediate datasets and job state are stored in temporary files
- ✧ Recovery information is stored in host and vnode directories represented by AB_WORK_DIR defined in the Ab Initio environment

- **Directory dedicated to Co>Ops**
- **Should have enough free space; Cannot be NFS or NAS mounted**
- **Holds Storage of Internal Log Files (used in recovery of Ab Initio Graph)**
- **Used when components are connected via name pipes**
- **Sub-directories of  AB_WORK_DIR**
  - **host – Holds Control Node Recovery Files**
  - **vnode – Holds Processing Node Recovery Files**
  - **data – Holds files for Layouts**
  - **cache – Holds Cache Files needed by remote components**
- **Important logging information in "host" and "vnode" directories**
- **Usually does not have data files.**
- **Components with host layouts  or database layouts, data written to "data" subdirectory**
- **AB_WORK_DIR fill up leads to non-recovery of Ab Initio Jobs.**

## A sample log file ..

✹**Reading the Log : CPU**
- ✧ CPU time:  total processing for component
- ✧ Status:  [ Running : Finished ]
- ✧ Skew: among CPU times of each partition
- ✧ Vertex:  component

✹**Reading the Log : DATA**
- ✧ Data bytes:  # processed
- ✧ Records:  # processed
- ✧ Status:  [ unopened : opened : closed ]
- ✧ Skew:  among data bytes in partitions
- ✧ Flow: link between components
  - ⇨ data tracking info is displayed on flows in GDE
- ✧ Vertex:  component
- ✧ Port: of component

✹**Interpreting the log**

- ✧ Compute data bytes/sec through component, in each partition
- ✧ Look for serialization: effective CPU = (cpu time)/(elapsed time)
- ✧ compare open vs. closed partitions:serialized when some partitions remain open long after others have closed    data skew
- ✧ Deadlock:no change in record counts over couple of intervals

- **Avoid Sorts as it is consuming more memory.**

- **Avoid components like Join with DB(hitting db for each and every record) .**

- **Use Lookups.**

- **Use In-memory Join/Rollup.**

- **Assign Driving Port of Join correctly.**

- **Filtering un-required data before processing.**

- **Phasing.**

# THANK YOU

**TATA CONSULTANCY SERVICES**
Asia's Largest Global Software & Services Company