# Ab Initio Interview Questions

**Q. What do you understand by Ab Initio? / Give a brief introduction of Abinitio.**

Ab Initio, also known as Abinitio, is a tool used to extract, transform, and load data. 'Abinitio' is a Latin word that means 'from the beginning'.

It was named Abinitio because Sheryl Handler and their team started it after the bankruptcy of their previous company. Sheryl Handler was the former CEO of Thinking Machines Corporation, and he decided to start this company as a new beginning when the Thinking Machines Corporation went bankrupt.

It is mainly used for data analysis, data manipulation, batch processing, and graphical user interface (GUI) based parallel processing for businesses.

**Q. What is Ab Initio Software?**

Ab Initio Software is an American multinational private enterprise software corporation headquartered in Lexington, Massachusetts. Ab Initio Software specializes in high-volume data processing applications and enterprise application integration. The Ab Initio software provides several products on a platform for parallel data processing applications.

**Q. What is Abinitio?**

- "Abinitio" is a latin word meaning "from the beginning."
- Abinitio is a tool used to extract, transform, and load data.
- It is also used for data analysis, data manipulation, batch processing, and graphical user interface based parallel processing.
- Ab Initio is popular software that offers data processing applications and enterprise application integration. It provides a single platform for data analysis, complex events, batch processing, data manipulation, quantitative, and qualitative data processing.
- Ab Initio software helps you build large-scale data processing applications and run them in parallel environments.

**Q. Abinitio Applications / used**

Abinitio Software applications are most widely used in Business Intelligence Data Processing Platforms to build most business applications such as operational systems, distributed application integration, complex event processing to data warehousing, and data quality management systems.

**Applications of Ab initio**

- It is used for application development and execution environments.
- It is a tool used in most of the big industries like insurance, banking, logistics, stock market, retail, finance, telecommunication to process complex and enormous volumes of data which gives stability and reliability.
- Ab initio solves the most challenging data processing issues for many of the leading organizations in many fields like finance and insurance.
- ETL for data warehouses, data marts and some operational data sources.
- It is also used for the parallel data cleansing and also for the validation.
- It can be seen in the parallel data transformation and also in the filtering.
- Ab initio ETL tool can also be seen in the high-performance analytics.
- It can be seen in real time and also in the parallel data capture.

**Financial Services:**

Ab initio is widely used in the financial services industry for data integration, processing, and analysis. It can be used to integrate data from various sources, including transactional data, market data, and customer data, allowing financial institutions to make more informed decisions and manage risk more effectively.

**Telecommunications:**

Ab initio is used by many telecommunications companies to process and analyze large amounts of data related to customer usage patterns, call records, and other data sources. This data can be used to improve customer service, optimize network performance, and identify opportunities for new products and services.

**Healthcare:**

Ab initio is used in the healthcare industry for data integration and processing, allowing healthcare providers to combine data from various sources, such as electronic health records, clinical research data, and claims data, to gain insights into patient health and improve the quality of care.

**Retail:**

Ab initio is used in the retail industry for supply chain management, inventory management, and customer analytics. Retailers can use Ab initio to integrate data from various sources, such as point-of-sale systems, customer relationship management systems, and supply chain systems, to gain insights into customer behavior and optimize their operations.

**Government:**

Ab initio is used by many government agencies for data integration and processing, allowing agencies to combine data from various sources and analyze it to identify trends and make informed policy decisions.

**Energy and Utilities:**

Ab initio is used by energy and utility companies to integrate and process data from various sources, such as smart meters, energy usage data, and weather data. This data can be used to optimize energy generation, reduce waste, and improve customer service.

**Manufacturing:**

Ab initio is used by manufacturing companies to integrate data from various sources, such as production line sensors, supply chain systems, and customer relationship management systems. This data can be used to optimize production processes, reduce waste, and improve customer service.

**Transportation:**

Ab initio is used in the transportation industry for data integration and processing, allowing companies to combine data from various sources, such as GPS tracking data, customer relationship management systems, and supply chain systems, to gain insights into customer behavior and optimize their operations.

**Marketing:**

Ab initio is used by marketing departments to integrate data from various sources, such as customer relationship management systems, social media, and web analytics tools. This data can be used to gain insights into customer behavior and optimize marketing campaigns.

**E-commerce:**

Ab initio is used by e-commerce companies to integrate data from various sources, such as customer transactions, website analytics, and supply chain systems. This data can be used to optimize website performance, improve customer service, and drive sales.

**Q. Advantages of Abinitio**

- large amount of data can be handled.
- Development time is less.
- more parallelism (3).
- Scalability (utilize the resources).
- Version control (EME).
- Enterprise management's.

1. **Scalability:** Ab initio can process and analyze large amounts of data, making it well-suited for enterprise-level applications. It can scale horizontally across multiple machines, allowing it to process large datasets quickly.

2. **Graphical Development Environment:** Ab initio uses a graphical development environment (GDE) that allows developers to design and build data integration and processing workflows quickly and efficiently. The GDE offers a drag-and-drop interface that is easy to use and can be customized to meet the needs of the user.

3. **Metadata Management:** Ab initio offers robust metadata management capabilities, allowing developers to track changes to data sources, transformations, and workflows. This makes it easier to manage complex data integration projects and ensures that data is accurately tracked and audited.

4. **Data Quality:** Ab initio includes built-in data quality tools that can be used to clean, validate, and enrich data. This ensures that data is accurate and reliable, leading to more accurate insights and better decision-making.

5.  **Reusability:** Ab initio offers a modular design approach that allows developers to reuse components across different projects. This reduces development time and effort and ensures that workflows are consistent across different applications.

6.  **Platform Independence:** Ab initio can be deployed across a variety of platforms and operating systems, making it easy to integrate into existing IT infrastructure.

7.  **High Performance:** Ab initio has been designed to deliver high performance, even when processing large volumes of data. It is able to make use of parallel processing, which allows it to handle large datasets more efficiently.

8.  **Flexibility:** Ab initio is a very flexible platform that allows users to work with a variety of different data sources and formats. This makes it easy to integrate with existing systems and processes and means that developers don't have to worry about data format issues.

9.  **Collaboration:** Ab initio has been designed to support team-based development, making it easy for multiple developers to work on the same project. This is supported through features like shared components and version control.

10. **Extensibility:** Ab initio is a highly extensible platform that can be customized to meet the specific needs of individual organizations. This is supported through features like custom components and user-defined functions.

11. **Reliability:** Ab initio is a very reliable platform that has been designed to ensure that data is processed correctly and that no data is lost. This is supported through features like error handling and recovery.

12. **Security:** Ab initio is a very secure platform that has been designed to protect data from unauthorized access. This is supported through features like encryption and access control.

**Q. define data warehouse**

A data warehouse is *a large collection of business data used to help an organization make decisions*.

A data warehouse is a large, centralized repository of structured data that is designed to support business intelligence activities, such as data analysis and reporting.
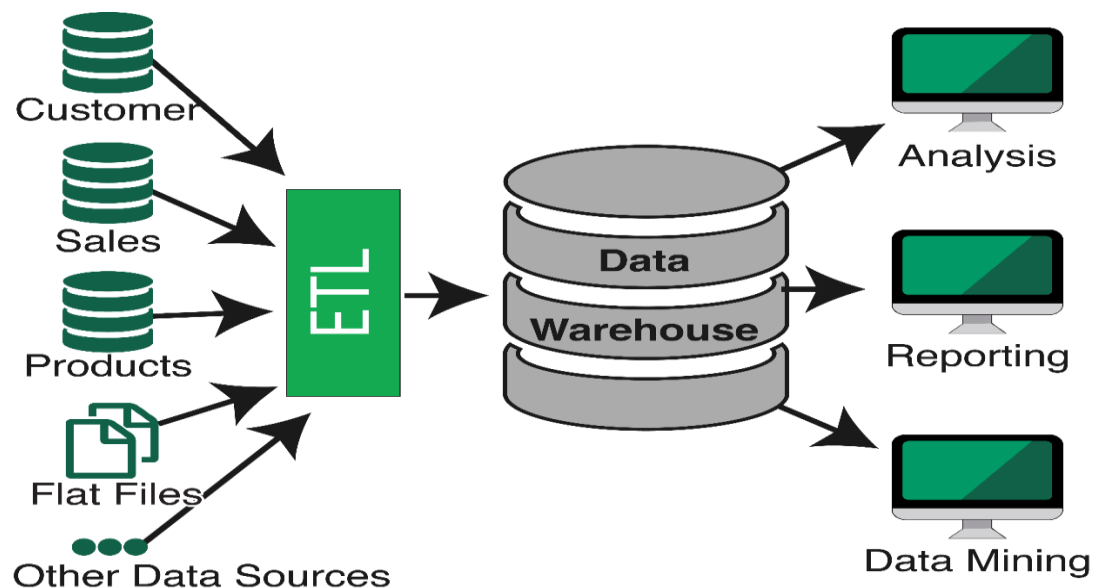
The primary goal of a data warehouse is to provide a unified view of an organization's data that can be used to support decision-making and strategy formation.
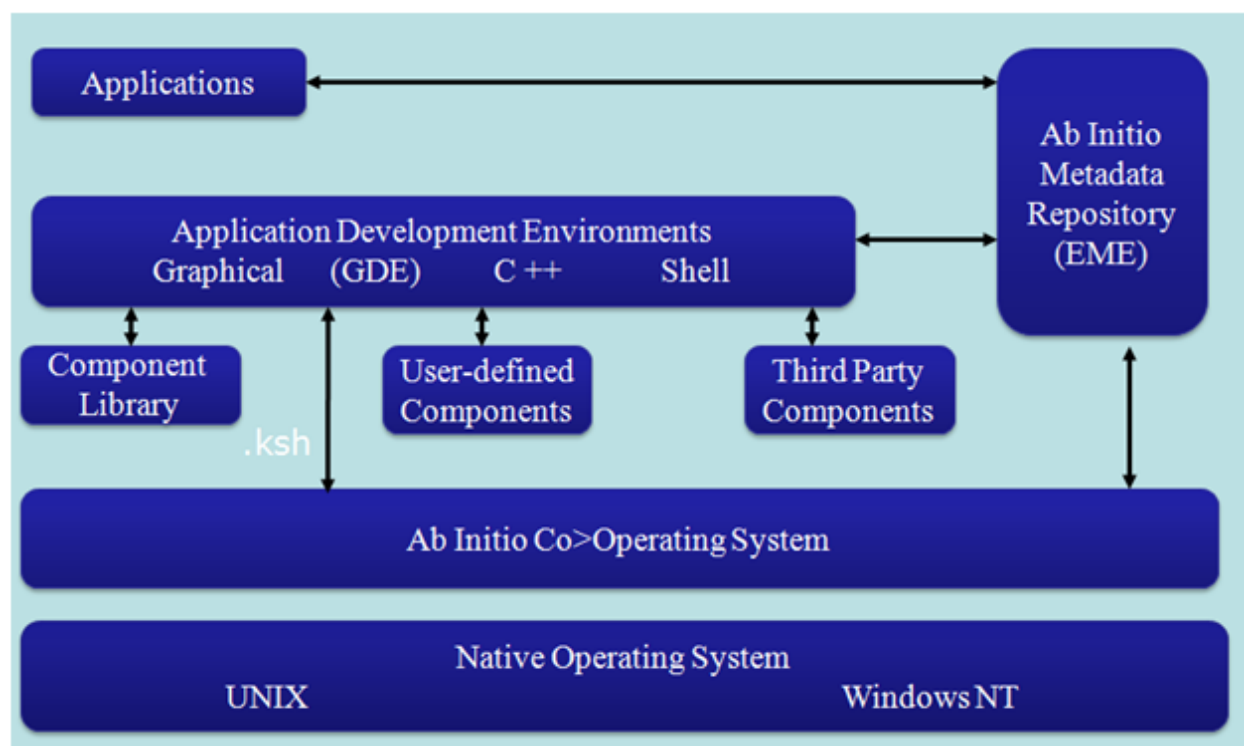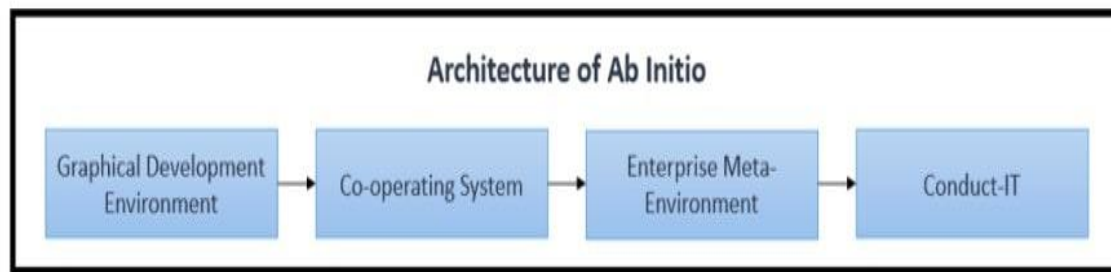
**Q. What is a ETL Tool**

ETL stands for Extract, Transform, Load, and refers to a process in data warehousing that involves extracting data from a variety of sources, transforming it into a format that can be loaded into a target database for analysis and reporting purposes.

An ETL tool is software that automates this process and makes it more efficient. ETL tools typically provide a visual interface that allows users to define the data extraction, transformation, and loading rules and processes, making it easier to manage the flow of data between different systems. The tool can handle a wide range of data formats, including structured data, semi-structured data, and unstructured data.

Some popular ETL tools include Apache NiFi, Talend, Informatica PowerCenter, and Microsoft SQL Server Integration Services (SSIS). These tools are widely used in industries such as finance, healthcare, retail, and e-commerce, among others, to manage large amounts of data and ensure its accuracy and completeness for business intelligence and reporting purposes.
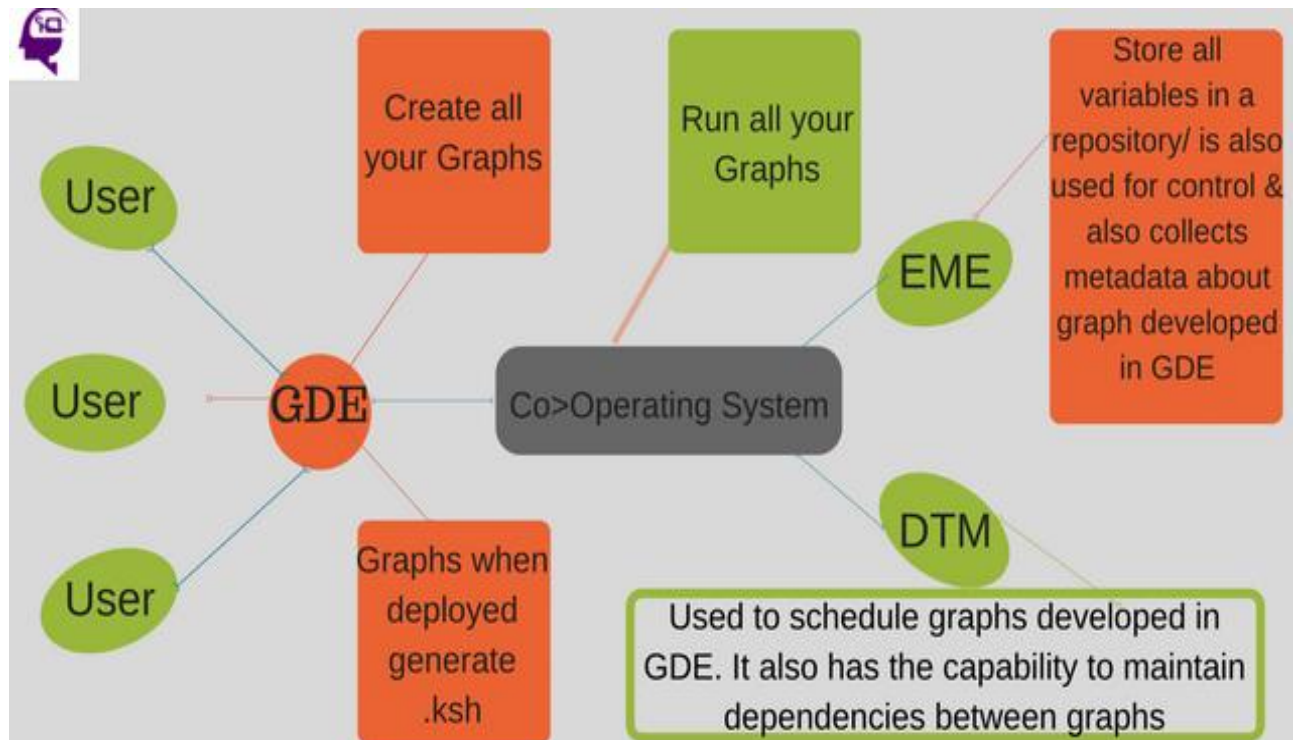
**Q. What is the architecture Of Ab Initio?**



**Ab Initio ETL tool Architecture**

- ➢ Ab-Initio program works with the client-server model.
- ➢ The client is referred to as "Graphical development environment" (GDE).
- ➢ The server is called Co-Operating System".
- ➢ The Co-working system can dwell on a mainframe or UNIX machine.
- ➢ The Ab-Initio code is referred to as graph, which has received .mp extension.
- ➢ The graph from GDE is required to be deployed in corresponding .Ksh version.
- ➢ In a Co-operating approach the corresponding .Ksh in a run to do the required job.

**The Ab Initio is a business intelligence software containing 6 data processing products:**

**GDE & Co-operating system is a Client-Server Architectures**

1. Co>Operating system (Co>Op v2.14, 2.15..)

2. The Component Library

3. Graphical Development Environment (GDE v1.14, 1.15…)

4. Enterprise Meta>Environment (EME v3.0…)

5. Data Profiler

6. Conduct>IT

**Explanation: -**

1. **Graphical Development Environment (GDE**): This is a visual interface for designing data processing flows. Users can drag and drop components and connect them together to create complex data processing pipelines.

2. **Co>Operating System (Co>OS):** This is the runtime environment that executes Ab initio graphs. It provides a distributed processing architecture that allows for scalable and high-performance data processing.

3. **Ab initio Enterprise Meta>Environment (EME):** This is a central repository for storing and managing metadata about Ab initio projects. It allows users to manage components, track changes, and collaborate with other team members.

4. **Ab initio Continuous>Flows (ACF):** This is a tool for creating real-time data processing flows. It allows users to create data pipelines that can handle high volumes of data in near real-time.

## 1. Co>Operating System:

- Co>Operating system is an engine for you graph to Run. It is not visible, it installed on your OS.
- Co>Operating system is nothing like distributed operating system which runs as a back end for GDE.
- Co>operating system is a server for Abinitio.
- The Co>Operating System® product is a platform for building and deploying large, sophisticated business data processing applications.
- These applications can reliably process large amounts of data by running scalably on many CPUs within a single server and across networks of servers or containers (computing clouds).
- The Co>Operating System product provides complete portability of applications across Unix, Linux, Hadoop, Windows, and mainframe (z/OS) servers.
- Co>Operating System, which your system administrator installs on a host UNIX or Windows NT Server, as well as on processing nodes. (The host is also referred to as the control node.)

Co>Operating System **operates on the top of the operating system and works as a base for all Ab Initio processes**. It can run on operating systems like Windows, Linux, Solaris, AIX, HP-UX, and z/OS. It provides Ab Initio extensions through which ETL processes can be controlled.

Ab Initio Co>operating system is a foundation for all Ab Initio functions and provides a base for all Ab Initio processes. It runs on a variety of process environments like AIX, HP-UX, Solaris, Linux, z/OS, and windows.

**It presents following points:**

- Control and run Ab Initio graphs and controls the ETL methods
- Provides Ab Initio extensions to the operating system
- ETL processes monitoring and debugging
- Metadata administration by interacting with EME

**Co>Operating system** is the foundation for all the Ab initio applications. It provides a general engine for the integration of all kinds of data processing and communication between all the ETL tools within the given platform.

It runs on Windows, Linux, Unix, OS/390, zOS on the mainframe. It enables distributed and parallel execution, platform independent data transport, establishing checkpoints, and process monitoring.

It is the foundation for all the Ab initio functions and provides a base for all the Ab initio processes. It controls and runs the Ab initio graphs and controls all of the ETL methods.

It provides the extensions to the operating systems powered by Ab initio. It also helps for debugging the information.

**Co>Operating System** is installed on the operating system.

CoOperating System It operates on top of the operating system, and the ab initio provides this and it the base for all Ab Initio processes.

Air commands are one of the features that can be installed on different operating systems like UNIX, Linux, IBM etc.

This co-op is installed on particular o.s platform that is called native o.s .coming to the eme, it's just as a repository in Informatica, its hold the metadata, transformations, dbconfig files source and targets information.

**These are the following features that it provides,**

- Manages and runs Ab Initio graphs and control the ETL processes
- Providing the extensions
- ETL processes monitoring and debugging
- Metadata management and interaction with the EME.

## 2. Graphical Development Environment:

GDE is a graphical application for developers which is used for designing and running Ab Initio graphs.

You are going to create a business application or Graph.

- Graph is a combination of dataset component and program.
- Component connected by a flow.

The Graphical Development Environment™ (GDE®) product allows for extremely rapid development of complex and enterprise-wide applications. The graphical nature of the product encourages a clearer understanding of these applications by large and diverse groups of people.

**GDE** It's a designing component and used to run the ab initio graphs.

Graphs are formed by the components (predefined or user-defined) and flows and the parameters. It provides the ETL process in Ab Initio that is represented by graphs.

Ability to run, debug the process logs jobs and trace execution logs

GDE provides an intuitive graphical interface for enhancing and executing applications. Possible without any difficulty drag-and-drop add-ons from the library on to a canvas, configure them and fasten them into flowcharts. Ab Initio GDC graph compilation system outcome within the new release of a UNIX shell script which may be executed on a machine where GDE is not installed.

**Graphical Development Environment (GDE),** which you install on your PC

(Client node) and configure to communicate with the host (control node).

**Graphical Development Environment (GDE)**

- Is the graphical application for designing and running Ab Initio graphs.
- Users can access the EME metadata through the GDE web browser or Co>Operating system command line.
- Coming to gde its is an end-user environment where we can develop the graphs (mapping just like in Informatica) designer uses the gde and designs the graphs and save to the eme or sandbox it is at the user side. Where eme is at the server-side.

**Graphical Development Environment (GDE):**

Ab initio Graphical Development Environment (GDE) gives an intuitive graphical interface for editing and executing the applications. It is very easy to drag components from the library onto a canvas, to configure them, and connect them via flowcharts. Here, the code development tool is developer friendly.

### 3. Enterprise Meta>Environment:

EME is repository and environment to retailer and managing metadata. It has the potential to retailer both business and technical metadata.

EME metadata can also be accessed from the Ab Initio GDE, internet browser or Ab Initio Co>Op command line.

The Enterprise Meta>Environment® (EME®) product is an enterprise data catalog and metadata management solution. It provides broad capabilities for all metadata created and used by the Co>Operating System product and a large number of third-party products. The EME product versions and stores schema-level metadata, business metadata, data transformation rules, applications (graphs), and operational statistics. The EME product uses technical metadata to automatically derive dependency and lineage information. It also includes a set of connectors for integrating information from a wide range of third-party design and reporting tools.

Ab initio Enterprise Meta>Environment is a data store that allows tracking changes in the developed graphs and metadata used in their development. It also offers some tools such as independent analysis, metadata management, statistical analysis, and version controlling. It has the potential to retailer both the business and technical metadata. It can also be accessed from the Ab initio Graphical Development Environment (GDE) or from any internet browser.

#### 4. Component Library:

The Ab Initio Aspect Library is a reusable program module for sorting, data transformation, and excessive-space database loading and unloading.

The Ab initio Component Library is a reusable software module for sorting, data transformation, and high-speed database loading and also unloading. It is a flexible and extensible tool which adapts at the runtime to the old formats of the given data entered.

The Ab initio Component Library is a reusable software module for sorting, data transformation, and high-speed database loading and also unloading.

It is a flexible and extensible tool which adapts at the runtime to the old formats of the given data entered.

It allows creation and the incorporation of new components obtained from any program that gives permission for the integration and reuse of external legacy codes and the search engines.

#### 5. Data Profiler:

Data Profiler is a analytical software that can specify data variety, scope, distribution and variance. It runs in a graphic environment on top of the Co>Op.

I've never used abinitio data profiler, can some one throw some light on how it will help in data quality, as I understand data profiler provides statistical information. Internally it uses abinitio graph. So why we use this at all. Can't we get all those things by using various components.

Ab initio Data Profiler is an analytical application that helps us to specify the given data range, scope, distribution, variance, and the quality. It also runs in a graphical environment on the top of the Co>Operating System.

#### 6. Conduct>It:

It is powerful GUI-based parallel processing tool for ETL data management and analysis.

Conduct>it's a high-volume knowledge processing systems developing tool. It allows combining graphs from GDE with customized scripts and programs from other vendors.

Ab Initio provides both Graphical and command line interface to Conduct>It.

The Conduct>It® product provides sophisticated job automation and monitoring for Ab Initio applications and other programs, linking them together into a single execution environment. The Conduct>It product can launch applications that run in parallel across distributed sets of Unix,

Linux, Windows, and z/OS servers, and containers. It also provides facilities for job alerting, job recovery, and resource management.

Ab initio Conduct>It is a high-volume data processing systems developing tool for developing the data processing in a big manner. It enables us to combine graphs from Graphical Development Environment (GDE) with custom scripts and programs from other vendors. It helps for monitoring and providing the operational statistics for the execution of graphs. It also provides both the graphical and command line interface to the Conduct>It.

**Q. Explain the relation between EME, GDE, and Co-operating system.**

**Enterprise Meta>Environment (EME)**

Is nothing but a repository for storing and managing metadata.

Enterprise Meta-Environment (EME) It's an environment for storage and also metadata management (Both business and technical metadata).

The metadata is accessed from the graphical development environment and also the web browser or the cooperating command line. It is an ab initio repository for any placeholders.

Eme is said as enterprise metadataenv, gde as graphical development env and the co-operating system can be said as abinitio server relation b/w this co-op, eme and gde are as follows operating system is the abinitio server.

**EME:**

- EME stands for Enterprise Metadata Environment
- It is a repository to AbInitio. It holds transformations, database configuration files, metadata, and target information

**GDE:**

- GDE – Graphical Development Environment
- It is an end user environment. Graphs are developed in this environment
- It provides GUI for editing and executing AbInitio programs

**Co-operative System:**

- Co-operative system is the server of AbInitio.
- It is installed on a specific OS platform known as Native OS.
- All generated graphs in GDE are later deployed and executed in co-operative system.

**Q. what is the role of Co-operating system in Abinitio?**

- Manage and run Abinitio graph and control the ETL processes
- Provide Ab initio extensions to the operating system
- ETL processes monitoring and debugging
- Meta-data management and interaction with the EME

Co>Operating System operates on the top of the operating system and works as a base for all Ab Initio processes.

It can run on operating systems like Windows, Linux, Solaris, AIX, HP-UX, and z/OS.

It provides Ab Initio extensions through which ETL processes can be controlled. It manages metadata by interacting with EME, manages, and runs Ab Initio graphs.

**Q. Explain what is SANDBOX?**

- A SANDBOX is referred for the collection of graphs and related files that are saved in a single directory tree and behaves as a group for the purposes of navigation, version control, and migration.

- A sandbox is like a directory that contains a collection of Ab Initio graphs and related files. It is local to any user and is a replica of the Project in the EME. It will be helpful for version control, migration, and navigation.

- A SANDBOX is referred for the compilation of graphs and connected files that are saved in a particular directory tree and behaves as a set for version control, navigation, and relocation.

**BASIC TERMS of In Ab Initio**

**Q. What is Dataset?**

In Simple terms dataset is a file. A file can be main frame file or any fixed or delimited

files. There are various **types of datasets**

➢ FIDXED EBCDIC
➢ FIXED ASCII
➢ DELIMITED ASCII
➢ SAS dataset etc...

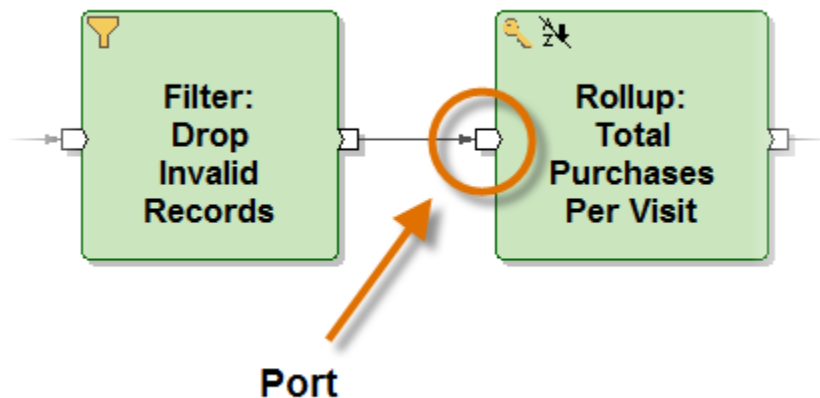You can also think of dataset as a table in database world.

**Q. What is Component?**

- Component is Ab Initio Program.
- There are various components like SELECT, FILTER, SORT, JOIN, MERGE, DEDUP, ROLLUP, SCAN, USER DEFINED etc

**Q. What is Port?**

- A port is a connection point for the input or output to a component.
- The place where a flow connects to a component is called a port**.**



Port

**Q. What is Flow?**

- A flow carries a stream of data between components in a graph. Flows connect components via ports.
- Ab Initio supplies four kinds of flows with different patterns: straight, fan-in, fan-out, and all-to-all. We will discuss various kinds of flows as we go through this training.

**Q. What is Graph?**

A graph is a diagram that defines the various processing stages of a task and the streams of data as they move from one stage to another. Visually, stages are represented by components and streams are represented by flows. The collection of components and flows comprise an Ab Initio graph.

**Q. What is Field?**

A filed is Equivalent to column of a table in Database world. Filed is also called variable

**Q. What is Key?**

Key is used many places in Ab Initio Development, We use key to sort the data, Join two files, Rollup data etc...

**Components**

**Q. List the commonly used components in an Ab Initio graph?**

The commonly used components in an Ab Initio graph are,

- Input file/output file
- Lookup file
- Input table/output table
- Join
- Sort
- Partition
- Partition by key
- Gather
- Reformat

- Concatenate
- Rollup
- Normalize
- Scan
- Dedup
- Filter
- Replicate
- Merge

**Q. What are the components or functions available in ab initio?**

| Component | Purpose |
|---|---|
| Dedup | To remove duplicates |
| Join | To join multiple input dataset based on a common key value. |
| Sort | This component reorders the data. It takes the collation order and dumps data to memory. <br><br> Sort the elements ascending order or descending order <br><br> Ex:- Score(Ace) |
| Filter | Any conditional related removal of data. <br><br> Ex:- mark>75 |
| Replicate | This is component is mainly for parallelism as an additional copy of data is useful while any other nodes go unavailable. |
| merge | This component is to combine multiple input data. |

## 1. Sort Component

The Sort component reads records from its input port, sorts them by one or more key fields, then writes the records in sorted order to its output port.

**DML Example: -**

```
record
    decimal("|") ID;
    string("|") Fname;
    string("|") Lname;
    decimal("|") Mobile;
    string("\n")  City;
end;
```

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CRARE | 85 |

**Sort**

| | fname | lname | score |
|---|---|---|---|
| 1 | JAMES | LOWEK | 51 |
| 2 | TERI | ANSELL | 64 |
| 3 | JAMES | NEVIUS | 64 |
| 4 | CHERRY | JORGE | 64 |
| 5 | ALLAN | RISO | 65 |
| 6 | LOIS | DAVIAU | 66 |
| 7 | TERRY | GILZOW | 67 |
| 8 | LINH | TORIAN | 67 |
| 9 | SCOTT | FORGUE | 69 |
| 10 | DAPHNE | HECT | 69 |

key: score

**Q. Explain what is Sort Component in Abinitio?**

The Sort Component in Abinitio re-orders the data. It comprises of two parameters "Key" and "Max-core".

**Key:**

➢ It is one of the parameters for sort component which determines the collation order
➢ It represents the collation order.

**Max-core:**

➢ This parameter controls how often the sort component dumps data from memory to disk.
➢ It defines how often the sort component should dump data from memory to disk.

**2. Filter by Expression Component**

**"Filter By Expression" is used to filter the input data according to the given condition.**

**"Filter By Expression" is a component in the Ab Initio software used to filter records from an input data stream based on specified conditions.**

It evaluates a given expression for each record in the input stream and passes only those records that meet the specified condition to the output stream. This component is often used in data integration and ETL (Extract, Transform, Load) processes to extract only the records that are relevant for further processing, and discard those that do not meet the specified criteria.

The Filter by Expression component reads records from its input port and evaluates a selection expression for each one. If the result is true, it writes that record to its output port.

**Example: -**

Suppose we have a flat file containing sales data, with the following columns: date, product_name, price, and quantity. We want to filter out all the records where the price is less than $10.

We start by reading in the input file using a "Read File" component, and passing the output to the "Filter By Expression" component.

In the "Filter By Expression" component, we specify the condition as price >= 10. This will evaluate to true for all records where the price is greater than or equal to $10.

We connect the output of the "Filter By Expression" component to a "Write File" component, which will write the filtered records to a new file.

Read File  →  Filter By Expression  →  Write File

price >= 10

When we run the graph, the "Filter By Expression" component will only pass records where the price is greater than or equal to $10 to the output file. This allows us to filter out irrelevant records and focus on the data that meets our criteria.

**DML Example: -**

```
record
    string("|") Fname;
    string("|") Lname;
    decimal("|") Score;
end;
```

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CPARE | 85 |

**Filter by Expression**

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | ANNIE | MARCIL | 95 |
| 4 | MARTHA | CRAPE | 85 |
| 5 | NATHAN | TOMEY | 92 |
| 6 | TAMI | HOLTON | 86 |
| 7 | LYDA | ARNOW | 91 |
| 8 | GLEN | KITKA | 83 |
| 9 | JAMES | EHRMAN | 100 |
| 10 | PAUL | MARULLO | 88 |

select_expr: **score > 75**

### 3. Dedup Component

**"Dedup" component is used to eliminate duplicate records from an input data stream.** The Dedup component compares each incoming record with the previous records and filters out any records that are exact duplicates of previously processed records.

**Dedup component Example**

Suppose we have a flat file containing sales data, with the following columns: date, product_name, price, and quantity. We want to eliminate duplicate records based on the date and product_name columns.

We start by reading in the input file using a "Read File" component, and passing the output to the "Dedup" component.

In the "Dedup" component, we specify the key as date and product_name. This means that the component will compare each record based on these two fields to determine if it is a duplicate.

We connect the output of the "Dedup" component to a "Write File" component, which will write the deduplicated records to a new file.
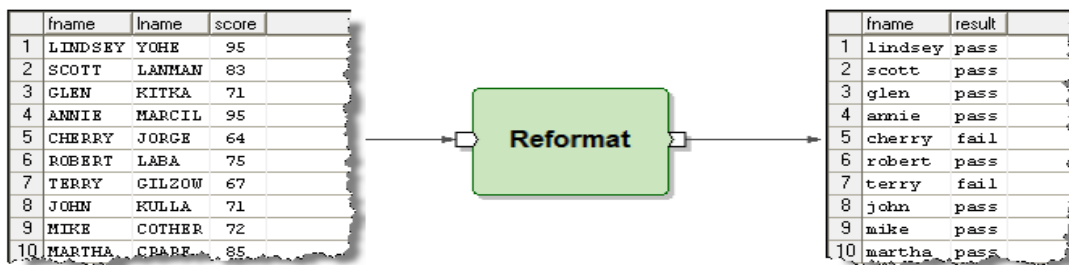

Read File  →  Dedup  →  Write File

           date, product_name

When we run the graph, the "Dedup" component will remove any records that are exact duplicates of previously processed records based on the date and product_name fields. This allows us to eliminate redundant data and focus on the unique records in the dataset.

## 4. Reformat Component

**The "Reformat" component in Ab Initio is a powerful tool used for transforming data in a variety of ways.** It allows developers to specify a set of rules that the component will follow to modify the data, which can include things like converting data types, reformatting strings, and applying various types of mathematical operations.

The Reformat component reads records from its input port and, for each, applies a set of rules to create a new record, which it writes to its output port.



| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CRABE | 85 |

**Reformat**

| | fname | result |
|---|---|---|
| 1 | lindsey | pass |
| 2 | scott | pass |
| 3 | glen | pass |
| 4 | annie | pass |
| 5 | cherry | fail |
| 6 | robert | pass |
| 7 | terry | fail |
| 8 | john | pass |
| 9 | mike | pass |
| 10 | martha | pass |

```
transform:
    string_downcase(fname)                    → fname
    if (score > 70) "pass" else "fail"  → result
```

**DML Eaxample:-**

```
record
    decimal("|") ID;
    string("|") Fname;
    string("|") Lname;
    decimal("|") Mobile;
    string("\n")  City;
end;
```

**Input table** data = ID, Fname, Lname, Mobile, City

**Output table** after the use Reformat=ID, Mobile, City, Full_Name

**transforming Exapresion:-**

```
out :: reformat(in) =
 begin
     out.ID :: in.ID;
     out.Mobile :: in.Mobile;
     out.City :: in.City;
     out.Full_Name :: in.Fname + in.Lname;
end;
```

### 5. Rollup component

**Rollup component is used to perform aggregate calculations on a set of input records based on one or more grouping fields. The Rollup component can be used to calculate various aggregate functions, such as sum, count, maximum, minimum, average, and so on.**

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.

The Rollup component aggregates collections of records grouped by a key. It applies a set of rules to produce an output record for each group.

**Example: -**

Assume we have a file containing transaction data with the following fields: Date, Customer ID, Transaction ID, Item ID, Quantity, and Price. We want to calculate the total revenue and average price for each customer across all transactions.
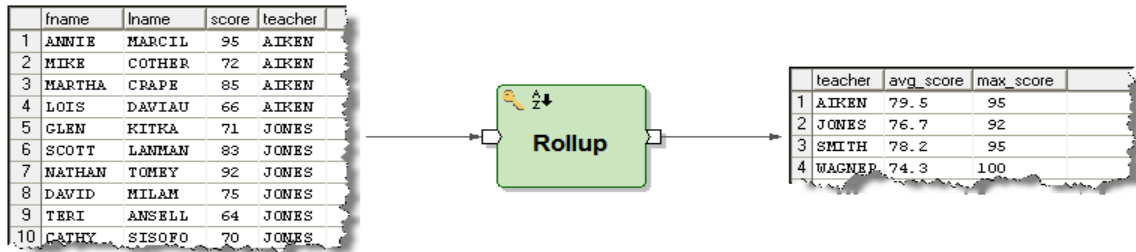
To achieve this, we can use the Rollup component as follows:

1. Use a Sort component to sort the input records by Customer ID and Transaction ID.
2. Connect the output of the Sort component to the input of the Rollup component.
3. In the Rollup component, select Customer ID as the grouping field and select the Sum and Average aggregate functions for the Quantity and Price fields, respectively.
4. Connect the output of the Rollup component to a Write component to write the aggregated data to a file.

When the graph is executed, the Rollup component will group the input records by Customer ID and calculate the total revenue and average price for each customer. The output will be a file

containing a single record for each customer, with the following fields: Customer ID, Total Revenue, and Average Price. This can be useful for analyzing customer behavior and identifying high-value customers.

| | fname | lname | score | teacher |
|---|---|---|---|---|
| 1 | ANNIE | MARCIL | 95 | AIKEN |
| 2 | MIKE | COTHER | 72 | AIKEN |
| 3 | MARTHA | CRAPE | 85 | AIKEN |
| 4 | LOIS | DAVIAU | 66 | AIKEN |
| 5 | GLEN | KITKA | 71 | JONES |
| 6 | SCOTT | LANMAN | 83 | JONES |
| 7 | NATHAN | TOMEY | 92 | JONES |
| 8 | DAVID | MILAM | 75 | JONES |
| 9 | TERI | ANSELL | 64 | JONES |
| 10 | CATHY | SISOFO | 70 | JONES |

**Rollup**

| | teacher | avg_score | max_score |
|---|---|---|---|
| 1 | AIKEN | 79.5 | 95 |
| 2 | JONES | 76.7 | 92 |
| 3 | SMITH | 78.2 | 95 |
| 4 | WAGNER | 74.3 | 100 |

transform:

teacher → teacher

avg(score) → avg_score

max(score) → max_score

## Q. What is Rollup Component?

- Roll-up component enables the users to group the records on certain field values. It is a multiple stage function and consists initialize 2 and Rollup 3.

- A rollup component is used to group the records based on certain field values. It is a multi-stage transform function which contains functions like initialize, rollup and finalize.

- The roll-up component facilitates users to collect or group the records on certain field values. It is called for each of the records in the group and consists of initializing 2 and Rollup 3.

**Q. Have you used the rollup component? Describe how?**

If the user wants to group, the records on particular field values then rollup is the best way to do that. Rollup is a multi-stage transform function and it contains the following mandatory functions.

1. initialize

2. rollup

3. finalize

Also, need to declare one temporary variable if you want to get counts of a particular group.

For each of the group, first, it does call the initialize function once, followed by rollup function calls for each of the records in the group and finally calls the finalize function once at the end of last rollup call.

\** Rollup Component has been used for grouping the records which are based on Field Values. It also is a Multi-Stage Transforming function that contains functions such as:
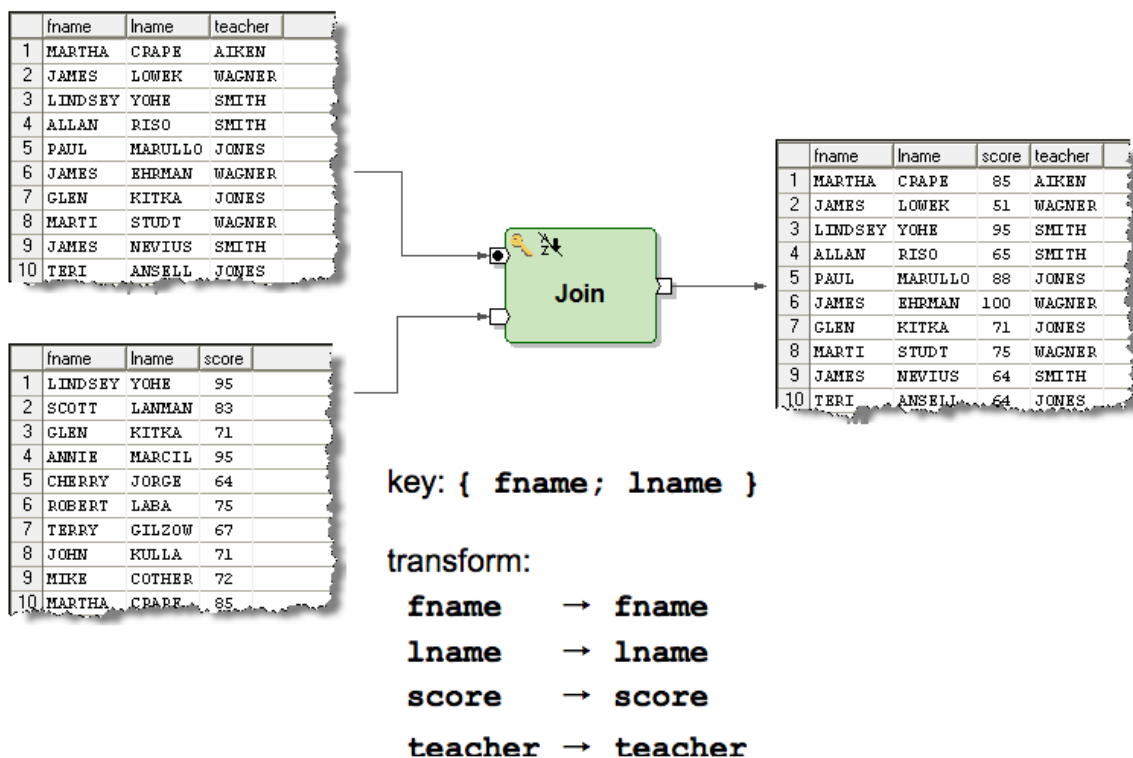
Initialize

Rollup

Finalize

## 6. Join component

**Join reads data from two or more input ports, combines records with matching keys according to the transform you specify, and sends the transformed records to the output port. Additional ports allow you to collect rejected and unused records.**

join component is used to combine data from multiple input flows based on a matching condition. The join component is typically used when you have two or more input data sources that need to be merged together based on a common key.
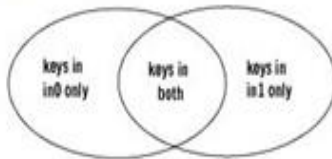
The Join component combines the data arriving on its input ports. Each pair of input records with matching keys is combined using rules to produce an output record.



| | fname | lname | teacher |
|---|---|---|---|
| 1 | MARTHA | CRAPE | AIKEN |
| 2 | JAMES | LOWEK | WAGNER |
| 3 | LINDSEY | YOHE | SMITH |
| 4 | ALLAN | RISO | SMITH |
| 5 | PAUL | MARULLO | JONES |
| 6 | JAMES | EHRMAN | WAGNER |
| 7 | GLEN | KITKA | JONES |
| 8 | MARTI | STUDT | WAGNER |
| 9 | JAMES | NEVIUS | SMITH |
| 10 | TERI | ANSELL | JONES |

| | fname | lname | score | teacher |
|---|---|---|---|---|
| 1 | MARTHA | CRAPE | 85 | AIKEN |
| 2 | JAMES | LOWEK | 51 | WAGNER |
| 3 | LINDSEY | YOHE | 95 | SMITH |
| 4 | ALLAN | RISO | 65 | SMITH |
| 5 | PAUL | MARULLO | 88 | JONES |
| 6 | JAMES | EHRMAN | 100 | WAGNER |
| 7 | GLEN | KITKA | 71 | JONES |
| 8 | MARTI | STUDT | 75 | WAGNER |
| 9 | JAMES | NEVIUS | 64 | SMITH |
| 10 | TERI | ANSELL | 64 | JONES |

**Join**

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CRAPE | 85 |

key: { fname; lname }

transform:
```
fname    → fname
lname    → lname
score    → score
teacher  → teacher
```

# Examples of join types

This section gives examples of various join types, involving different settings for the **join-type**, **record-required**$n$, and **dedup**$n$ parameters.

The two intersecting ovals in the diagrams that follow represent the key values in the records on the two ports — **in0** and **in1** — that are the inputs to JOIN:
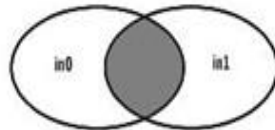


For each possible setting of **join-type** — or (if **join-type** is **Explicit**) combination of settings for **record-required**$n$ — the shaded region of each of the following diagrams represents the inputs for which JOIN calls the transform. JOIN ignores the records that have key values represented by the white regions, and consequently those records go to the **unused**$n$ port.



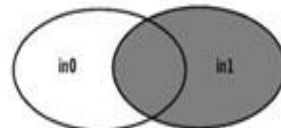JOIN calls transform          JOIN does not call transform
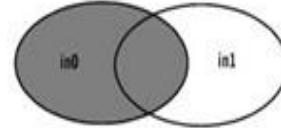
**Case 1**
Inner Join join-type



**Case 2**
Full Outer Join join-type



**Case 3a**
Explicit join-type

record-required0: False
record-required1: True



**Case 3b**
Explicit join-type

record-required0: True
record-required1: False



## Types of joins

1. Inner Join
2. Left Outer Join
3. Right Outer Join
4. Full Outer Join

**1. Inner Join:**

**An inner join returns only the matching records from two or more data sources.** To perform an inner join in Ab Initio, you can use the "join" component with the "inner" join type. Here is an example:

Suppose you have two data sets: "Customers" and "Orders". Both data sets have a "CustomerID" column that can be used to match records between the two data sets.

To perform an inner join between these two data sets, you can use the following Ab Initio graph:

|---read Customers---|      |---read Orders---|

            join(inner)

              |

            write

In this graph, the "read Customers" and "read Orders" components read the data from the two input data sets. The "join" component performs an inner join using the "CustomerID" column as the key. Finally, the joined data is written to an output data set using the "write" component.


**2. Left Outer Join:**

**A left outer join returns all the records from the left data source and matching records from the right data source.** To perform a left outer join in Ab Initio, you can use the "join" component with the "left outer" join type. Here is an example:

Suppose you have two data sets: "Customers" and "Orders". Both data sets have a "CustomerID" column that can be used to match records between the two data sets.

To perform a left outer join between these two data sets, you can use the following Ab Initio graph:

|---read Customers---|      |---read Orders---|

            join(left outer)

              |

            write

In this graph, the "read Customers" and "read Orders" components read the data from the two input data sets. The "join" component performs a left outer join using the "CustomerID"

column as the key. Finally, the joined data is written to an output data set using the "write" component.

## 3. Right Outer Join:

**A right outer join returns all the records from the right data source and matching records from the left data source.** To perform a right outer join in Ab Initio, you can use the "join" component with the "right outer" join type. Here is an example:

Suppose you have two data sets: "Customers" and "Orders". Both data sets have a "CustomerID" column that can be used to match records between the two data sets.

To perform a right outer join between these two data sets, you can use the following Ab Initio graph:

```
|---read Customers---|      |---read Orders---|

          join(right outer)

                |

            write
```

In this graph, the "read Customers" and "read Orders" components read the data from the two input data sets. The "join" component performs a right outer join using the "CustomerID" column as the key. Finally, the joined data is written to an output data set using the "write" component.

## 4. Full Outer Join:

**A full outer join returns all the records from both data sources, along with matching records.** To perform a full outer join in Ab Initio, you can use the "join" component with the "full outer" join type.
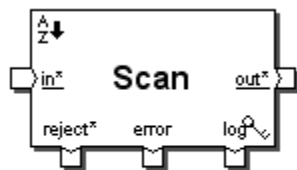
### Example:

Suppose you have two data sets: "Customers" and "Orders". Both data sets have a "CustomerID" column that can be used to match records between the two data sets.

### 7. Scan Component

**Scan generates a series of cumulative summary records for groups of data records.**

For every input record, Scan generates an output record that consists of a running cumulative summary for the group to which the input record belongs, up to and including the current record.

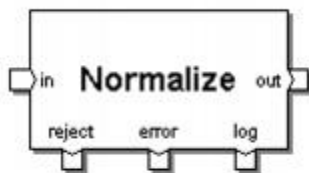**For example,** the output records might include successive year-to-date totals for groups of records.



### 8. Normalize Component

**Normalize generates multiple output records from each of its input records. You can directly specify the number of output records for each input.**

Normalize component is also one of the multistage components in ab initio.



Normalize generates multiple output records from each of its input records. You can directly specify the number of output records for each input record, or you can make the number of output records dependent on a calculation.

In contrast, to consolidate groups of related records into a single record with a vector field for each group — the inverse of NORMALIZE — you would use the accumulation function of the ROLLUP component.

**Q. What is the difference between scan and rollup?**

- Rollup is used to aggregate data.
- Scan is for successive total and is used to produce summary(cumulative) of data.

**9. Replicate component**

**It replicates the data for a particular partition and send it out to multiple out ports of the component but maintains the partition integrity**.

**Example: -** Your incoming flow to replicate has a data parallelism level of 2. with one partition having 10 recs & other one having 20 recs.

replicate component is used to duplicate data and distribute it to multiple downstream components. The replicate component is typically used in situations where the same data needs to be processed by multiple components in parallel.
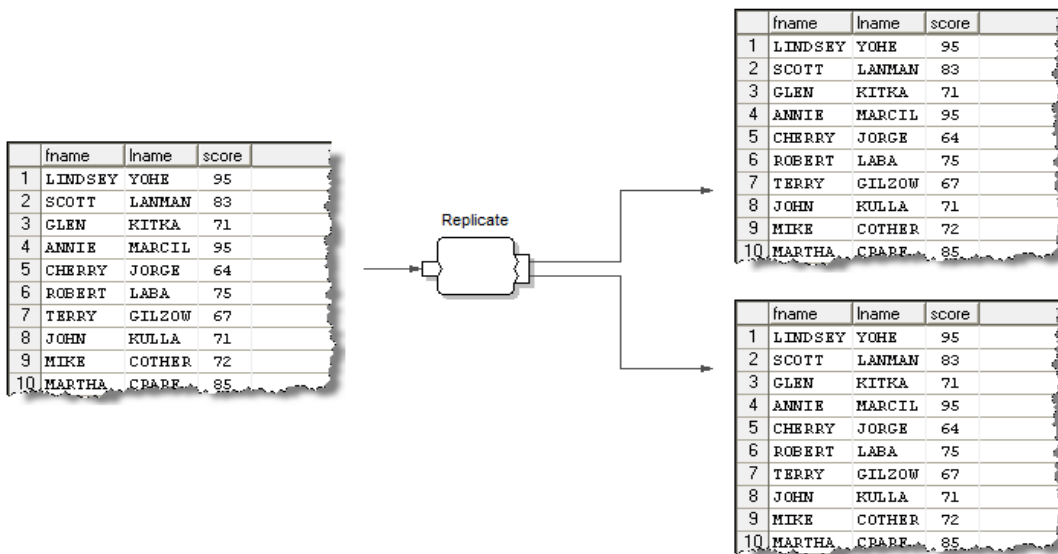
When a replicate component receives an input record, it creates multiple copies of the record and sends each copy to a different downstream component. The number of copies created is determined by the replication factor property of the component. This property specifies the number of copies that should be created for each input record.

The replicate component is often used in conjunction with a gather component, which combines the output from multiple downstream components back into a single stream. The gather component ensures that all copies of the input record are processed and combined into a single output record.

One of the benefits of using the replicate component is that it can improve the overall performance of a data integration solution by distributing the processing load across multiple components. This can help to reduce the time it takes to process large volumes of data, which is often a key requirement in data-intensive applications.

The Replicate component simply copies every record arriving on its input port to every flow connected to its output port.

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CPARE | 85 |

Replicate

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CPARE | 85 |

| | fname | lname | score |
|---|---|---|---|
| 1 | LINDSEY | YOHE | 95 |
| 2 | SCOTT | LANMAN | 83 |
| 3 | GLEN | KITKA | 71 |
| 4 | ANNIE | MARCIL | 95 |
| 5 | CHERRY | JORGE | 64 |
| 6 | ROBERT | LABA | 75 |
| 7 | TERRY | GILZOW | 67 |
| 8 | JOHN | KULLA | 71 |
| 9 | MIKE | COTHER | 72 |
| 10 | MARTHA | CPARE | 85 |

**Q. Mention what dedup-component and replicate component does?**

**Dedup component:**

- ➢ It is used to remove duplicate records.
- ➢ It is used to remove duplicate records from the flow based on a specified key.

**Replicate component:**

- ➢ It combines the data records from the inputs into one flow and writes a copy of that flow to each of its output ports.
- ➢ It is used to combine input records from multiple sources into one flow and write a copy of that flow to output ports.

**Q. What is the usage of dedup component and replicate component in Abinitio?**

In Abinitio, the dedup component is used to eliminate duplicate records. On the other hand, the replicate component combines the data records from the inputs into one run and writes a copy of that run to each of its output ports.

<u>**Parallelisms**</u>

**Q. What Parallelisms Abinitio Support?**

**1. Pipeline Parallelism:**

- Data is passed from one component to another component. Data is worked on both of the components.

- In this the records are processed in pipeline.

**2. Data Parallelism :**

- Same data is parallelly worked in a single application.
- Data is processed at the different servers at the same time.

**3. Component Parallelism:**

- Different data is worked parallelly in a single application.

- In this two or more components process the records in.

**4. Multi-File Parallelism:**

- Processing multiple input or output files in parallel on different processing nodes.

**Q. What are the different types of parallelism used in Abinitio?**

**Parallelism is the concurrent execution of multiple components, and it's a good way to increase performance.**
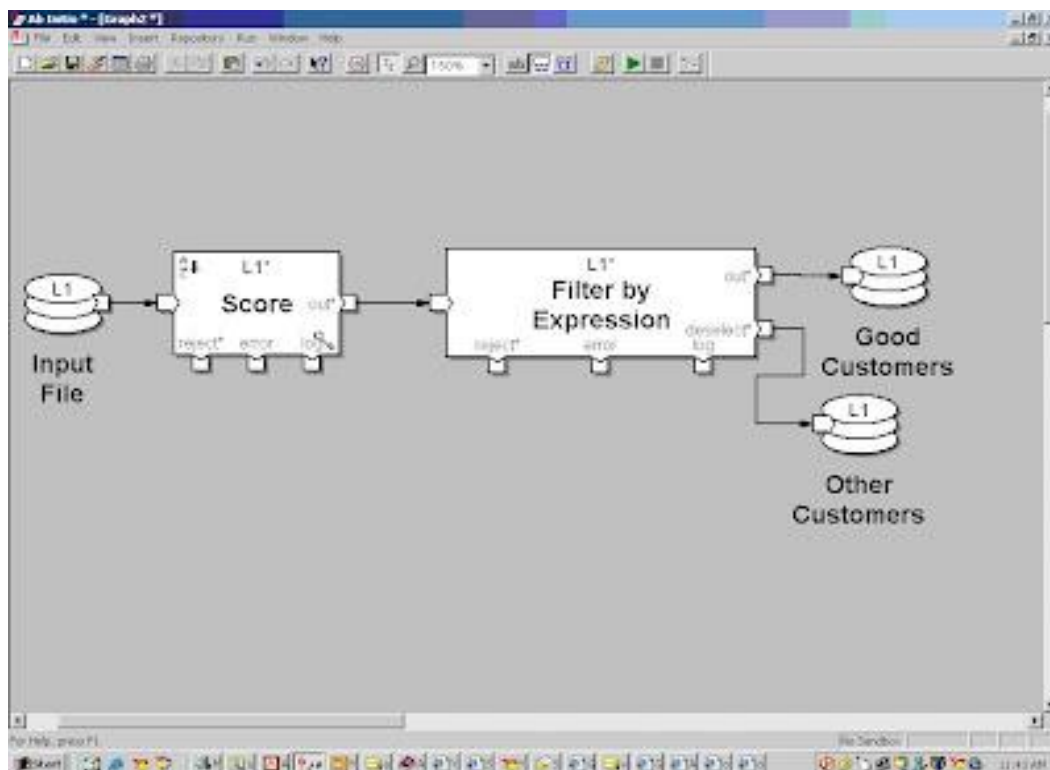
**1. Pipeline parallelism:** (By default parallelism)

occurs when connected sequence of program components executes simulteneously on same branch of the graph.

A graph that deals with multiple components executing simultaneously on the same data uses pipeline parallelism. Each component in the pipeline read continuously from the upstream components, processes data and writes to downstream components.  Both components can operate in parallel.

It is used by the graph that deals with multiple components executing simultaneously on the same data. An application with multiple components but running on the same dataset. This uses pipeline parallelism.

An application with multiple components running simultaneously on the same data uses pipeline parallelism. Each component in the pipeline continuously reads from upstream components, processes data, and writes to downstream components. Since a downstream component can process records previously written by an upstream component, both components can operate in parallel.

### 2. Component parallelism:

The component parallelism is used by a graph with multiple processes executing simultaneously on separate data.
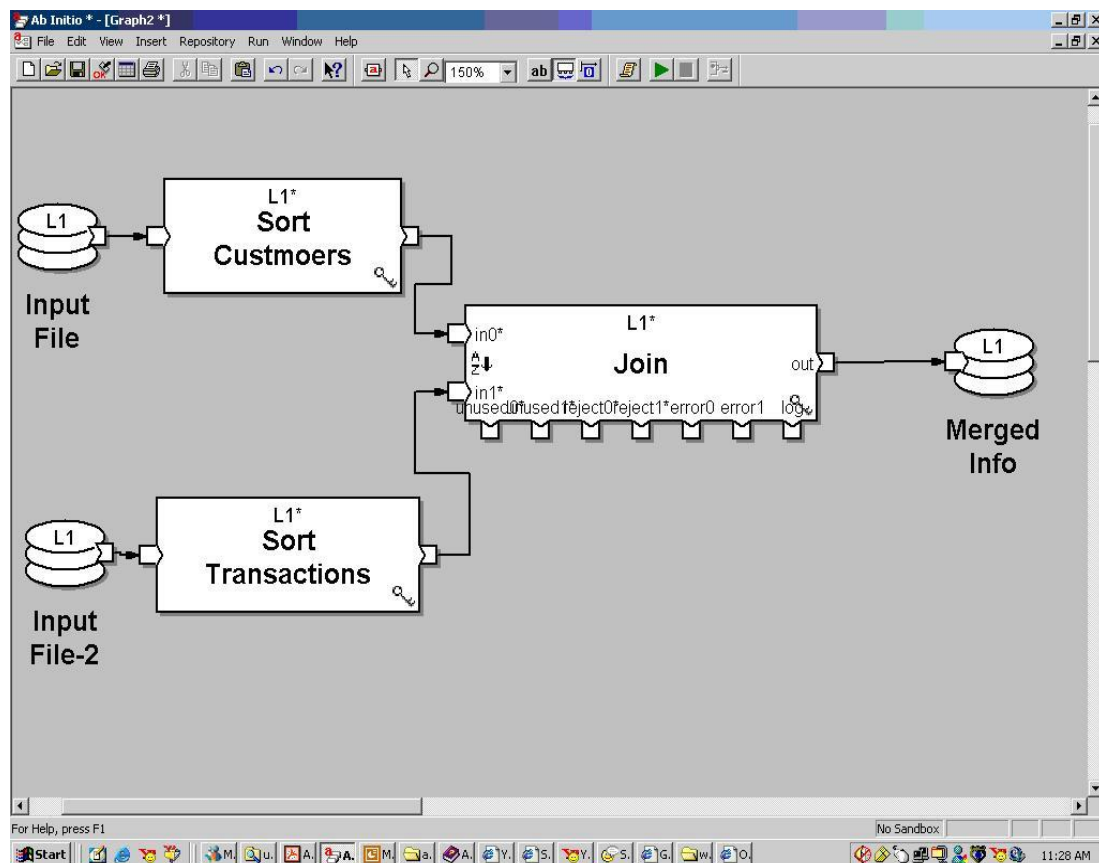
A graph with multiple processes executing simultaneously on separate data uses parallelism

**It is used by the graph that has multiple processes executing simultaneously on separate data.**

An application that has multiple components running on the system simultaneously. But the data are separate. This is achieved through component level parallel processing.

occurs when program components execute simuteneously on diff branches of the graph.

An application with multiple processes running on separate data uses component-level parallelism.
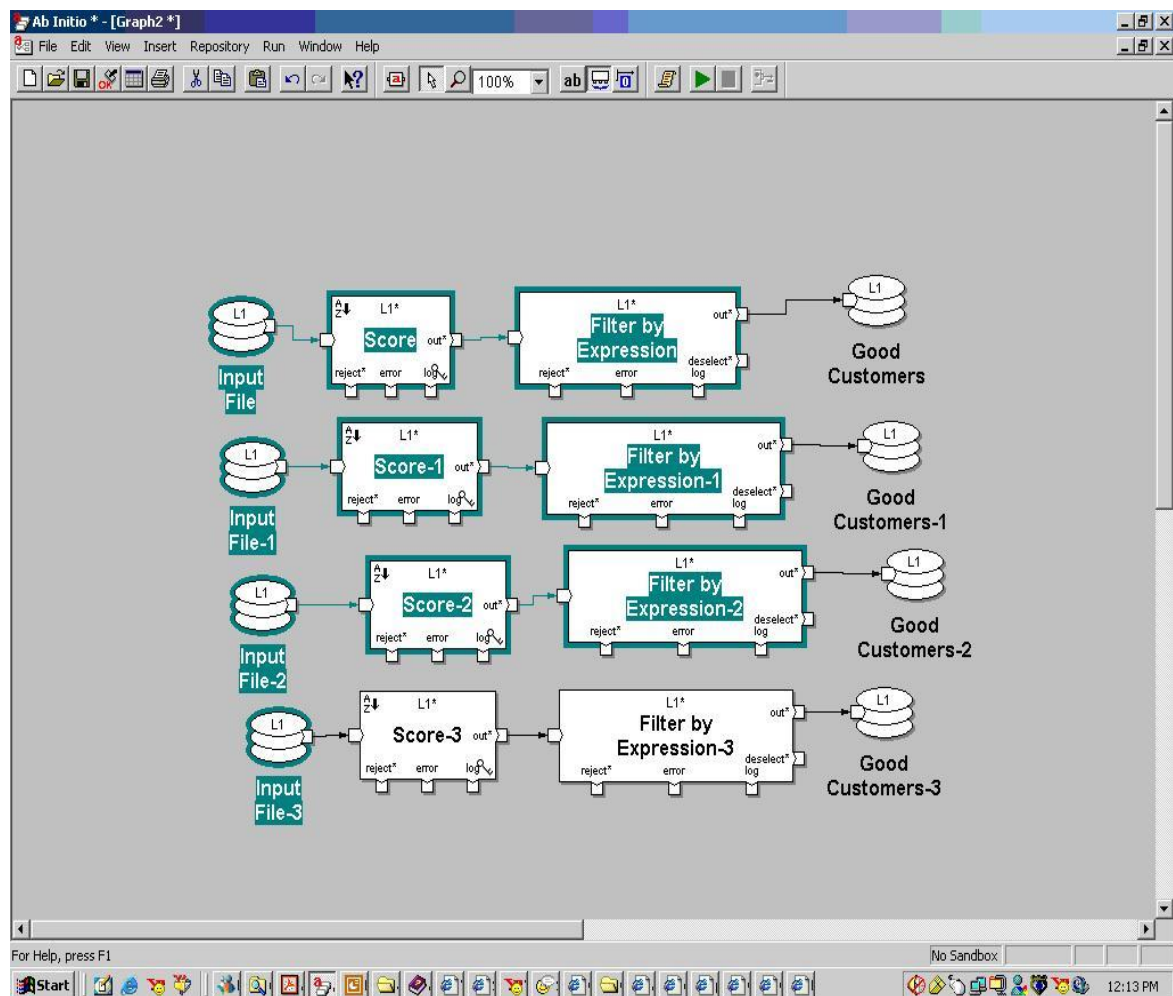
### 3. Data parallelism:

The data parallelism is used by a graph that works with data divided into segments and operates on each segment respectively.

A graph that works with data divided into segments and operates on each segments respectively, uses data parallelism.

It is used by the graph that works with data divided into segments, which operates on each segment, respectively.

Data is split into segments and runs the operations simultaneously. This kind of process is achieved using data parallelism.

An application that deals with data divided into segments and operates on each segment simultaneously uses data parallelism. Nearly all commercial data processing tasks can use data parallelism.

## Partitions

**Q. What is the different way to achieve the partitions?**

| Partitions | Description |
|---|---|
| Expression | Data split according to the data manipulation language. |
| Key | Grouping the data by specific keys |
| Load balance | Dynamic load balancing |
| Percentage | Segregate the data where the output size is on the fractions of 100 |
| Range | Split the data evenly based on a key and a range among the nodes |
| Round robin | Distributing the data evenly in block size across the output partitions. |

1. **Partition by Round-Robin:** Distributing data evenly, in block size chunks, across the output partitions
2. **Partition by Range:** You can divide data evenly among nodes, based on a set of partitioning ranges and key
3. **Partition by Percentage:** Distribution data, so the output is proportional to fractions of 100
4. **Partition by Load balance:** Dynamic load balancing
5. **Partition by Expression:** Data dividing according to a DML expression
6. **Partition by Key:** Data grouping by a key

**Q. What are the types of partition components in Ab Initio?**

In Abinitio, partition is the process of dividing data sets into multiple sets for further processing.

**1. Partition by Round-Robin:**

- ➢ The Round-Robin Partition is used for distributing data evenly, in block size chunks, across the output partitions.
- ➢ Distributing data evenly, in block size chunks, across the output partitions.
- ➢ It distributes the data evenly in block size chunks.

**Round Robin Partitioning:**

In this method, the data is partitioned in a round-robin fashion across the available partitions. This method ensures that the data is evenly distributed across all partitions.

**Round Robin Partitioning Example:**

Suppose we have a large product inventory dataset with multiple columns, and we want to partition it using round-robin partitioning. We can use the Partition component in Ab Initio with the following parameters:

- Partition method: Round-robin partitioning
- Number of partitions: 5 (for example)

This will create five partitions of the data, with each partition containing approximately the same number of records.

**2. Partition by Range:**

- ➢ The Partition by Range facilitates users to divide data evenly among nodes, according to the set of partitioning ranges and keys.
- ➢ You can divide data evenly among nodes, based on a set of partitioning ranges and key.
- ➢ Based on a set of partitioning ranges and key, it divides data evenly among nodes.

Partition by Range is a technique used to split a large dataset into smaller, more manageable parts called partitions. This is achieved by partitioning the data based on the values in one or more key columns.

When using Partition by Range in Ab Initio, the data is first sorted based on the key columns, and then partitioned based on a range of values in those columns. The range can be specified using different techniques, such as specifying a minimum and maximum value for each partition or specifying the number of partitions required.

The main advantage of using Partition by Range in Ab Initio is that it allows for better parallel processing of large datasets. Each partition can be processed independently, allowing for more efficient use of system resources and faster processing times.

Ab Initio provides several components for partitioning data by range, such as Partition by Key and Partition by Expression. These components allow users to define the partitioning strategy and specify the key or expression to be used for partitioning.

In addition to Partition by Range, Ab Initio also supports other partitioning techniques, such as Partition by Round Robin and Partition by Hash, which can be used in combination to achieve optimal performance and scalability for large-scale data processing.

**Example** of how Partition by Range might be used in Ab Initio:

Suppose we have a large dataset containing sales transactions, with columns for the transaction date, the product sold, and the sales amount. We want to partition this dataset by year to improve performance when querying the data.

To do this, we could use the Partition by Key component in Ab Initio, specifying the transaction date column as the key. We would then sort the data by the transaction date, and define a range of values for each partition based on the year of the transaction.

**Example** we might specify the following partition ranges:

- Partition 1: Transaction year <= 2010
- Partition 2: Transaction year > 2010 and <= 2011
- Partition 3: Transaction year > 2011 and <= 2012
- Partition 4: Transaction year > 2012 and <= 2013
- Partition 5: Transaction year > 2013

Ab Initio would then create separate physical storage structures for each partition, allowing us to query the data more efficiently by only accessing the relevant partitions.

When querying the data, we could use the Partition by Key component again to ensure that our query only accesses the relevant partitions. For example, if we wanted to retrieve all sales transactions from 2012, we could specify the relevant partition range in the query to ensure that only the relevant data is accessed.

### 3. Partition by Percentage:

> - The Partition by Percentage is used to distribute data in a way that the output is proportional to fractions of 100.
> - Distribution data, so the output is proportional to fractions of 100.
> - It distributes data in such a way that the output is proportional to fractions of 100.

Partition by Percentage is a technique used in Ab Initio to split a dataset into partitions based on the percentage of total data in each partition. This technique is often used when the data has an uneven distribution and cannot be evenly partitioned by key or range.

When using Partition by Percentage, Ab Initio randomly assigns each record in the dataset to a partition based on a percentage value. The percentage value for each partition is defined by the user, and the sum of all percentage values must equal 100%.

Ab Initio provides several components for partitioning data by percentage, such as Partition by Percentage and Partition by Round Robin with Percentage. These components allow users to specify the number of partitions and the percentage of data to be assigned to each partition.

The main advantage of using Partition by Percentage in Ab Initio is that it allows for better distribution of data across partitions, even when the data is highly skewed. This can help improve performance when processing large datasets by balancing the workload across different nodes in a cluster.

**Example** of how Partition by Percentage might be used in Ab Initio:

Suppose we have a dataset of website user clicks, with columns for the user ID, page URL, click date, and click count. We want to partition this dataset into 10 partitions, with each partition containing roughly the same percentage of total click count.

To do this, we could use the Partition by Percentage component in Ab Initio, specifying the click count column as the percentage key. We would then define the percentage values for each partition based on the total click count, such that each partition contains approximately 10% of the total click count.

**For example,** if the total click count is 1 million, we would define the following percentage values for each partition:

- Partition 1: 10%
- Partition 2: 10%
- Partition 3: 10%
- Partition 4: 10%

- Partition 5: 10%
- Partition 6: 10%
- Partition 7: 10%
- Partition 8: 10%
- Partition 9: 10%
- Partition 10: 10%

Ab Initio would then randomly assign each click record in the dataset to a partition based on the percentage value, creating 10 partitions with roughly the same percentage of total click count. This allows us to distribute the workload evenly across different nodes in a cluster, improving performance when processing large datasets.

When querying the data, we could use the Partition by Percentage component again to ensure that our query only accesses the relevant partitions. For example, if we wanted to retrieve all click data for a particular user, we could specify the relevant partition ranges in the query to ensure that only the relevant data is accessed.

### 4. Partition by Expression:

> The Partition by Expression is used to divide data according to a DML expression.
> Data dividing according to a DML expression.
> It divides data based on a DML expression.

Partition by Expression is a technique used in Ab Initio to partition data based on the result of a user-defined expression. This technique is useful when the data cannot be evenly partitioned by key, range, percentage, or load balance.

When using Partition by Expression, Ab Initio evaluates an expression for each record in the dataset and assigns the record to a partition based on the result of the expression. The expression can be any valid Ab Initio expression, and can reference any columns in the dataset.

Ab Initio provides several components for partitioning data by expression, such as Partition by Expression and Partition by Key and Expression. These components allow users to define the expression to be evaluated and the number of partitions to create, and Ab Initio automatically assigns each record to a partition based on the result of the expression.

The main advantage of using Partition by Expression in Ab Initio is that it allows users to partition data based on complex or custom criteria. For example, users can partition data based on the values of multiple columns, or based on a combination of mathematical or logical expressions.

**Example** of how Partition by Expression might be used in Ab Initio:

Suppose we have a dataset of customer orders, with columns for the order ID, customer ID, order date, and order amount. We want to partition this dataset based on the order date and customer ID.

We can use the Partition by Key and Expression component in Ab Initio to achieve this. First, we can use the "Multifile Partition" component to split the data into multiple files, as this makes it easier to process the data in parallel.

Next, we can use the Partition by Key and Expression component to partition the data based on the order date and customer ID. We can define an expression that concatenates the order date and customer ID columns, and then hash this value to create the partition key. For example, the expression might look like: concat(order_date, customer_id)

We can specify the number of partitions to create, and Ab Initio will automatically assign each record to a partition based on the result of this expression. In this way, the data will be evenly distributed across partitions based on the combination of order date and customer ID.

Once the data is partitioned, we can process each partition in parallel across multiple nodes in the computing cluster. Each node will process a subset of the data, and the workload will be balanced across all nodes, ensuring optimal resource utilization.

Finally, we can use the "Gather" component in Ab Initio to collect the results from each node and combine them into a single output file. This allows us to process the data efficiently in parallel while ensuring that the results are combined correctly.

**5. Partition by Key:**

- ➢ The Partition by Key is used to group data by a key.
- ➢ Data grouping by a key.
- ➢ It partitions data based on a key.

**Partitioning by Key Range**

This method is similar to hash partitioning, but instead of applying a hash function, the data is partitioned based on a range of key values. The range is usually determined based on the minimum and maximum values of the key column.

**Partitioning by Key Range Example:**

Suppose we have a large customer order dataset with multiple columns, and we want to partition it based on the order date column using key range partitioning. We can use the Partition component in Ab Initio with the following parameters:

- Partition method: Key range partitioning
- Key column: Order date
- Number of partitions: 12 (for example)
- Key range: [01-Jan-2023, 01-Feb-2023), [01-Feb-2023, 01-Mar-2023), ..., [01-Dec-2023, 01-Jan-2024)

This will create 12 partitions of the data based on the order date column, with each partition containing data for a specific month.

**6. Partition by Load balance:**

> ➤ The Partition by Load balance is used for dynamic load balancing.
> ➤ Dynamic load balancing.
> ➤ It distributes data based on dynamic load balancing.

Partition by Load Balance is a technique used in Ab Initio to distribute data evenly across multiple nodes in a computing cluster. This technique is often used when the data cannot be evenly partitioned by key, range or percentage.

When using Partition by Load Balance, Ab Initio assigns each record in the dataset to a partition based on the workload of each partition. The workload of each partition is determined by the number of records in the partition and the processing time required to process the records.

Ab Initio provides several components for partitioning data by load balance, such as Partition by Load Balance and Partition by Round Robin with Load Balance. These components allow users to define the maximum number of records or the maximum processing time that each partition can handle, and Ab Initio automatically balances the workload across partitions based on these values.

The main advantage of using Partition by Load Balance in Ab Initio is that it allows for optimal resource utilization by balancing the workload across different nodes in a cluster. This can help improve performance when processing large datasets by optimizing the processing time and resource usage on each node.

**Example** of how Partition by Load Balance can be used in Ab Initio.

Suppose we have a dataset of sales transactions, with columns for transaction ID, customer ID, product ID, transaction date, and transaction amount. The dataset is very large, and we want to process it in parallel across multiple nodes in a computing cluster to improve performance.

We want to partition this dataset into 8 partitions, with each partition having an approximately equal workload in terms of the number of records and processing time. To do this, we can use the Partition by Load Balance component in Ab Initio.

We can start by using the "Multifile Partition" component in Ab Initio to split the data into multiple files. This will make it easier to process the data in parallel.

Next, we can use the "Partition by Load Balance" component to evenly distribute the data across 8 partitions. We can specify a maximum number of records or maximum processing time for each partition, and Ab Initio will automatically balance the workload across partitions based on these values.Once the data is partitioned, we can process each partition in parallel across multiple nodes in the computing cluster. Each node will process a subset of the data, and the workload will be balanced across all nodes, ensuring optimal resource utilization.

Finally, we can use the "Gather" component in Ab Initio to collect the results from each node and combine them into a single output file. This allows us to process the data efficiently in parallel while ensuring that the results are combined correctly.

In this way, we can use Partition by Load Balance in Ab Initio to process large datasets in parallel across multiple nodes in a computing cluster, improving performance and resource utilization.

### 7. **Hash Partitioning:**

In this method, the data is partitioned based on a hash function applied to a particular key column. The key column is usually selected based on the distribution of the data to ensure a balanced distribution across partitions.

### **Hash Partitioning Example:**

Suppose we have a large sales transaction dataset with multiple columns, and we want to partition it based on the customer ID column using hash partitioning. We can use the Partition component in Ab Initio with the following parameters:

- Partition method: Hash partitioning
- Key column: Customer ID
- Number of partitions: 8 (for example)

This will create eight partitions of the data based on the hash function applied to the customer ID column.

### 8. Broadcast Partitioning:

This method broadcasts the data to all available partitions. This method is useful for small data sets that can be efficiently processed by all available resources.

### Broadcast Partitioning Example:

Suppose we have a small reference dataset that we want to use for joining with a large fact dataset. We can use the Broadcast component in Ab Initio with the following parameters:

- Input data: Reference dataset
- Broadcast to: All available partitions of the fact dataset

This will broadcast the reference dataset to all available partitions of the fact dataset, allowing the join to be performed efficiently.


### 9. Hybrid Partitioning:

This method combines two or more partitioning methods to achieve a more efficient data distribution. For example, a combination of hash partitioning and key range partitioning can be used to achieve a balanced data distribution and faster query performance.

### Hybrid Partitioning Example:

Suppose we have a large customer feedback dataset with multiple columns, and we want to partition it based on the feedback date column using a combination of hash partitioning and key range partitioning. We can use the Partition component in Ab Initio with the following parameters:

- Partition method: Hybrid partitioning
- Key column: Feedback date
- Number of partitions: 8 (for example)
- Hash partitioning: Based on customer ID column
- Key range partitioning: Based on feedback date column, with the key ranges defined as in the key range partitioning example

This will create eight partitions of the data, with each partition containing data for a specific month and a balanced distribution of data across the partitions based on the hash function applied to the customer ID column.

**Q. What is the difference between partitioning with key and round-robin?**

**Partition by key:**

In this, we have to specify the key based on which the partition will occur. It results in well-balanced data due to the key-based partitions. It is useful for key- dependent, parallelism.

Partitioning with key involves dividing the data into partitions based on the values of a specific key column. Each record in the dataset is assigned to a partition based on the value of this key column. This technique is useful when the data can be evenly divided into partitions based on the key column, as it allows for efficient data retrieval and processing based on the partition key.

Partitioning with key assigns records to partitions based on the value of a specific key column.

**Partition by round-robin:**

In this, distributing data evenly in block size chunks, the records are partitioned in a sequential way across the output partition. It is not key-based, and results are well-balanced data, especially with a block size of 1. It is useful for record independent parallelism.

round-robin involves dividing the data into partitions in a sequential manner. Each record is assigned to a partition in a circular fashion, with the first record going to the first partition, the second record going to the second partition, and so on, until all partitions have been assigned a record. This technique is useful when the data cannot be evenly divided based on a key column, as it allows for an even distribution of data across partitions.

partitioning with round-robin assigns records to partitions in a sequential and circular manner.

The choice of partitioning technique depends on the nature of the data and the desired parallel processing strategy.

**Q. Explain what is de-partition in Abinitio?**

**Q. What do you understand by de-partition in Abinitio?**

De-partition is used to read data from multiple flows or operations and re-join data records from different flows. There are Several de-partition components are available in Abinitio, such as Gather, Merge, Interleave, Concatenation, etc.

De-partition helps in allowing us to read from multiple flow and operations, is used for re-joining data records from different flows. De-partition consists of components that include:

**Q. List out the file extensions used in Abinitio?**

1. .mp: This file extension is used to store Abinitio graph or graph components.
2. .mpc: This file extension is used to specify a custom component or program.
3. .mdc: This file extension is used to specify data-set or custom data-set components.
4. .dml: This file extension is used to specify data manipulation language file or record type definition.
5. .xfr: This file extension is used to specify transform function files.
6. .dat: This file extension is used to specify data files (multifile or serial file).

1. .mp: It stores Ab initio graph or **graph component**
2. .mpc: **Custom component or program**
3. .mdc: Dataset or custom **data-set component**
    i. dataset template files or custom dataset components
4. .dml: **Data manipulation language** file or record type definition
    i. record format files
5. .xfr: Transform **function file**
6. .dat: **Data file** (multifile or serial file)
7. .dbc - **database table** files
8. .ksh - **shell scripting file**

1. **.dat:** This extension is used for data files in Ab Initio, which can be in various formats such as delimited, fixed width, or binary.

2. **.dml:** This extension is used for Data Manipulation Language files in Ab Initio. These files describe the structure of the input and output data files.

3.  **.mp:** This extension is used for Ab Initio graph files, which define the flow of data and processing steps for a particular job.

4.  **.roll:** This extension is used for checkpoint files in Ab Initio. These files store the state of a particular job so that it can be resumed if it is interrupted.

5.  **.db:** This extension is used for Ab Initio database files, which store metadata about data files and processing jobs.

6.  **.log:** This extension is used for log files in Ab Initio, which contain information about the execution of a particular job or process.

7.  **.csv:** This extension is commonly used for comma-separated value files in Ab Initio, which are a type of data file that uses commas to separate values in a table
8.  **.txt:** This extension is used for plain text files in Ab Initio, which can be used for storing data or for specifying configuration parameters for a job.

**Q. air commands used in Abintio?**

**1. air object Is<EME path for the object-/Projects/edf/..> :** It is used to see the listings of objects in a directory inside the project

**2. air object rm<EME path for the object-/Projects/edf/..> :** It is used to remove an object from the repository

**3. air object versions-verbose<EME path for the object-/Projects/edf/..> :** It gives the version history of the object.

1.  **m_dbaccess:** This command is used to perform database operations such as querying or updating data in a database.

2.  **m_dump:** This command is used to dump the contents of a data file in a human-readable format for troubleshooting and debugging.

3.  **m_gather_stats:** This command is used to collect performance statistics for a particular job or component.

4. **m_hdr**: This command is used to display the header information of a data file, such as the file format, field lengths, and data types.

5. **m_mkidx:** This command is used to create an index file for a data file, which can improve the performance of queries and searches.

6. **m_partition:** This command is used to partition data files for parallel processing.

7. **m_proj:** This command is used to select specific fields or columns from a data file.

8. **m_replicate:** This command is used to replicate data to multiple targets, such as databases or data warehouses.

9. **m_sort:** This command is used to sort data files based on one or more fields.

10. **m_transform:** This command is used to transform data files by applying business rules, data validation, or data enrichment.

**Q. What are the air commands used in Ab Initio?**

- Air object ls
- Air object rm
- Air project modify
- Air object versions -verbose
- Air lock show -user
- Air sandbox status

**Q. In how many segments the Abinitio EME can be segregated?**
**Q. Explain how Abinitio EME is segregated?**

- Data Integration Portion
- User Interface (Access to the meta-data information)

The Ab Initio EME is logically segregated into data integration portion and user interface to access metadata information.

**Q. Mention how can you connect EME to Abinitio Server?**

1. Login to EME web interface-
   http://serverhost:[serverport]/abinitio to connect EME to Abinitio Server.
2. Set AB_AIR_ROOT
3. We can connect to the EME data store through GDE.
4. We can also use air-command to connect EME to Abinitio Server.

**Q. How to connect EME to Ab Initio Server?**

EME can be connected to the Ab Initio server in several ways. The following are the ways to connect to EME.

1. Set AB_AIR_ROOT
2. Air commands
3. EME web interface at http://serverhost:[serverport]/abinitio
4. Connect to EME data store through GDE

**For connecting Eme with Ab Initio, there are many ways such as:**

1. We can login to EME Web Interfaces.
2. We can connect to EME Data Store.
3. By going through Air Command.
4. By setting AB_AIR_ROOT.

**Q. What are the types of output formats that we can get after processing data?**

- Charts
- Tables
- Vectors
- Plain Text files
- Maps
- Raw files
- Image files

**Q. Explain what does dependency analysis mean in Abinitio?**

Dependency Analysis is used for retrieving the maximum from the existing data and by using next_in_sequence/reformat we can also generate the sequence. The data comes from dependency analysis, the applications can depend on and produce this data, etc.

In Ab initio, dependency analysis is a process through which the EME examines a project entirely and traces how data is transferred and transformed- from component-to-component, field-by-field, within and between graphs.

The dependency analysis is a process through which the EME analyzes the project for the dependencies within and between the graphs. It examines the entire project and tracks how data is being transformed and transferred from component to component and field by field. The steps involved in dependency analysis are translation and analysis.

Dependency analysis is a procedure through which the EME examines a project completely and traces how data is transferred and distorted from field-by-field, component-to-component, within and among graphs.

Dependency analysis will answer the questions regarding data lineage. That is where does the data come from, what applications produce and depend on this data etc.

We can retrieve the maximum (surrogate key) from the existing data, the by using the scan or next_in_sequence/reformat we can generate a further sequence for new records. (E Learning pottal).

In Abinitio, dependency analysis is a process through which the EME examines a project entirely and traces how data is transferred and transformed- from component-to-component, field-by-field, within and between graphs.

**Q. What is a local lookup?**

- Local lookup file has records which can be placed in main memory.
- They use transform function for retrieving records much faster than retrieving from the disk.

**Q. Mention what the difference between "Look-up" file and "Look is up" in Abinitio?**

**Lookup file**

Lookup file defines one or more serial file (Flat Files); it is a physical file where the data for the Look-up is stored.

A lookup file represents one or more serial files, also known as Flat files. It is a physical file where lookup data is stored, which is small enough to be held in memory.

**Look-up**

While Look-up is the component of abinitio graph, where we can save data and retrieve it by using a key parameter.

Lookup is a component of the Ab Initio graph where data resides along with a key parameter. The data can be retrieved using this key parameter.

**Q. What is a local lookup?**

The local lookup function will be used before the lookup function call when the lookup file is a multifile and partitioned/sorted on a particular key.

It will be local to a partition, depending on the key. The data records in the lookup file can be loaded into memory. This way, the transform function retrieves records faster than retrieving from disk.

**Q. What Is A Local Lookup?**

- Local lookup file has records which can be placed in main memory
- They use transform function for retrieving records much faster than retrieving from the disk.

**Q. Define a local lookup?**

Local lookup file has documentations which can be located in major memory, and they are also used to change function for retrieving records much earlier than retrieving from the disk.

**Q. What do you understand by the "lookup" file in Abinitio?**

In Abinitio, the lookup file is used to define one or more serial files (also known as flat files). It is a physical file that stores the data for the Lookup. It is a two-dimensional table of data that has been stored in a disk file. It stores the name and display format for each column of data depending on the file format.

**Q. What is a ramp limit?**

- A limit is an integer parameter which represents a number of reject events
- Ramp parameter contain a real number representing a rate of reject events of certain processed records
- The formula is – No. of bad records allowed = limit + no. of records x ramp
- A ramp is a percentage value from 0 to 1.
- These two provides the threshold value of bad records.

**Q. What is the relation between limit and ramp?**

- Limit and ramp are used to set the reject tolerance of a graph.
- Limit is the number of rejects and ramp is rate of rejection.
- The formula for rejects tolerance is, **limit + (ramp*no_of_records_processed)**

**Q. What is a ramp limit?**

- ➢ A limit is an integer parameter which represents several reject events
- ➢ Ramp parameter contains a real number representing a rate of reject events of certain processed records.
- ➢ The formula is – No. of bad records allowed = limit + no. of records x ramp.
- ➢ A ramp is a percentage value from 0 to 1.
- ➢ These two provides the threshold value of bad records.

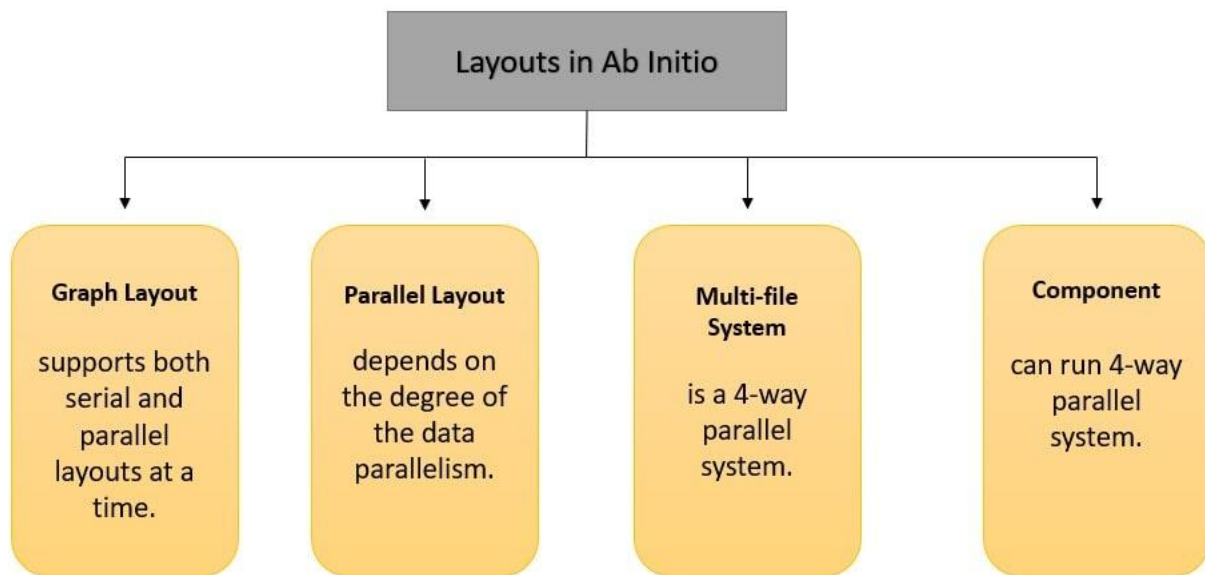**Q. Define ramp limit in AB Initio?**

Generally, the ramp is referred to the percentage value ranging from 0 to 1. To represent the number of reject events, the limit parameter possesses an integer called a ramp limit.

**Q. What do you mean by ramp limit?**

An edge is an integer parameter which represents a numeral of reject events. Ramp parameter encloses a real number on behalf of a rate of decline events of certain processed records; it has a percentage value from 0 to 1.

**Q. What are the kinds of layouts that Ab Initio supports?**

A layout defines which component should run where. Ab Initio has two kinds of layouts.



**Serial layout -** the level of parallelism is 1.

**Parallel layout -** the level of parallelism depends on the data partition.

**Q. What kind of layouts does Ab Initio support?**

- Supports serial and parallel layouts.
- A graph layout supports both serial and parallel layouts at a time.
- A multi-file system is a 4-way parallel system.
- A component in a graph system can run 4-way parallel system.

**Q. Mention the kind of layouts does Abinitio supports?**

There are successive and parallel layouts supported by Abinitio. A graph can have both at the same time. The parallel one depends on the amount of data parallelism. If the multi-file system is 4-way parallel, then a part in a graph can run four-way parallel if the layout is distinct such as it's same as the degree of parallelism.

**Q. How to improve the performance of a graph?**

- Use lookup instead of join and merge components.
- When we have to join two files and don't want duplicates, use a union function instead of a duplicate remover.
- Minimize the use of sort components.
- Reduce the use of regular expression functions in the transfer functions.
- Don't use broadcast as partitioner for large datasets.
- Use only the fields that are required in sort, reformat, and join components.

**Q. How do you improve the performance of a graph?**

- ➢ Reduce the usage of multiple components on certain phases.
- ➢ Use a refined and well-defined value of max core values for sort and join components.
- ➢ Minimize the use of regular expression functions like re_index in the transfer functions
- ➢ Minimize sorted join component and, if possible, replace them with in-memory join/hash join.
- ➢ Use only required fields in the sort, reformat, join components.
- ➢ Using Phase or the flow buffering during the cases of merge or sorted joins.
- ➢ Use hash join if the two sets of input are small; else better to choose the sorted join for the huge input size.
- ➢ For large dataset better not use broadcast as partitioned
- ➢ Reduce the number of sort components while processing.
- ➢ Avoid repartitioning of data unnecessarily.

**Q. Explain how you can run a graph infinitely in Ab initio?**

To execute graph infinitely, the graph end script should call the .ksh file of the graph. Therefore, if the graph name is abc.mp then in the end script of the graph it should call to abc.ksh. This will run the graph for infinitely.

Calling the .ksh file of the graph at the end script runs the graph infinitely. If the graph name is xyz.mp, then the end script of the graph should make a call to xyz.ksh.

**Q. Is it possible to run a graph infinitely in Ab Initio? If yes, how?**

Yes, it is possible to run a graph infinitely in Ab Initio. To do so, the graph end script should call the .ksh file of the graph. After that, if the graph name is xyz.mp then in the end script of the graph, it should call to xyz.ksh. By following the above steps, we can run the graph for infinitely.

**Q. What do you mean by the overflow errors?**

The overflow errors are the errors that occur when the computer cannot process the bulk data. While processing data, overflow errors occur if the bulky calculations exceed the range of memory provided to them.

Overflow errors are the errors raised when dealing with processing bulky sets of data. While processing data, bulky calculations might not fit the memory allocated for them. And when a character of more than 8-bits is stored, an overflow error is raised.

While processing data, bulky calculations are often there, and it is not always necessary that they fit the memory allocated for them. In case a character of more than 8-bits is stored there, these errors result simply.

**Q. Mention what information does a .dbc file extension provides to connect to the database?**

- The .dbc extension provides the GDE with the information to connect with the database are.
- Name and version number of the data-base to which you want to connect
- Name of the computer on which the data-base instance or server to which you want to connect runs, or on which the database remote access software is installed
- Name of the server, database instance or provider to which you want to link

**Q. What information does a .dbc file extension provide to connect to the database?**

- ➢ The .dbc file provides the below information to connect to a database.
- ➢ Name and version number of the database
- ➢ Name of the system on which the database server or instance is running.
- ➢ Name of the server or database instance

**Q. What are is_valid and is_define used for?**

- **is_valid -** It is used to test if a value is valid or not. If the expression is a valid data item, the value will be 1. If the expression is not a valid data item, the value will be 0.
- **is_define -** It is used to test if an expression is not NULL. If the expression is non-NULL, the value will be 1. The value will be 0 otherwise.

**Q. Mention what is the syntax for m_dump in Abinitio?**

The syntax for m_dump in Abinitio is used to view the data in multifile from unix prompt. The command for m_dump includes

**m_dump a.dml a.dat:** This command will print the data as it manifested from GDE when we view data in formatted text

**m_dump a.dml a.dat>b.dat:** The output is re-directed in b.dat and will act as a serial file.b.dat that can be referred when it is required.

**Q. How to add default rules in the transformer?**

Go to component properties, navigate to the parameter tab page, and double click on the transform parameter.

The transform editor page will open. Click on the edit menu and select the 'Add Default Rules' option from the drop-down. You can choose from Match names and Wildcard options.

- Double click on the transform parameter in the parameter tab page in component properties
- Click on Edit menu in Transform editor
- Select Add Default Rules from the dropdown list box.
- It shows Match Names and Wildcard options. Select either of them.

**Q. What is a surrogate key?**

The system generated unique sequential number is called a surrogate key. It acts as a primary key.

**Q. How will you handle it if DML is changing dynamically?**

We have a lot of ways through which we can handle dynamically changing DML. Some of the methods are,

- Use a conditional DM
- Call the vector functionality while calling the DMLs
- Use the MULTI REFORMAT component

**Q. What is Data Encoding?**

Data Encoding is used for keeping many cases confidential, it also is used for making sure if the information remains in the form by which no one else other than the receiver and the sender will understand.

**Q. How data is processed and what are the fundamentals of this approach?**

There are certain activities which require the collection of the data and the best thing is processing largely depends on the same in many cases. The fact is data needs to be stored and analyzed before it is processed. This task depends on some major factors are they are

1. Collection of Data

2. Presentation

3. Final Outcomes

4. Analysis

5. Sorting

These are also regarded as the fundamentals that can be trusted to keep up the pace in this matter.

**Q. What are the factors on which storage of data depends?**

It depends on the sorting and filtering. In addition to this, it largely depends on the software one uses.

**Q. What do you mean by data sorting?**

It is not always necessary that data remains in a well-defined sequence. It is always a random collection of objects. Sorting is nothing but arranging the data items in desired sets or sequence.

**Q. Explain data sorting?**

It is not always essential that data remains in a well-defined series. It is always a chance compilation of objects. Sorting is nil but arranging the data items in preferred sets or progression.

**Q. When running a stored procedure definition script how would you guarantee the definition could be rolled back in the event of problems?**

There are quite a few factors that determine the approach such as what type of version control is used, what is the size of the change, what is the impact of the change, is it a new procedure or replacing an existing and so on.

If it is new, then just drop the wrong one.

if it is a replacement then how big is the change and what will be the possible impact, depending upon you can have the entire database backed up or just create a script for your original procedure before messing it up or you just do ed and change the file back to original and reapply. you may rename the old procedure as old and then work on new and so on.

few issues to keep in mind are synonyms, dependencies, grants, any job calling the procedure at the time of change and so on. In nutshell, the scenario can be varied and the solution also can be varied.

**Q. What is AB_LOCAL expression where do you use it in ab-initio?**

ablocal_expr is a parameter of the table component of Ab Initio.LOCAL() is replaced by the contents of ablocal_expr. Which we can make use in parallel unloads. There are two forms of AB_LOCAL() construct, one with no arguments and one with a single argument as a table name(driving table).

The use of AB_LOCAL() construct is in Some complex SQL statements contain grammar that is not recognized by the Ab Initio parser when unloading in parallel. You can use the LOCAL() construct in this case to prevent the Input Table component from parsing the SQL (it will get passed through to the database). It also specifies which table to use for the parallel clause.

**Q. How to get DML using Utilities in UNIX?**

If your source is a COBOL copybook, then we have a command in Unix which generates the required in Ab Initio.

**Q. State the first_defined function with an example?**
This function is similar to the function NVL() in Oracle database
It performs the first values which are not null among other values available in the function and assigns to the variable.
**Example:** A set of variables, say v1,v2,v3,v4,v5,v6 are assigned with NULL.
Another variable num is assigned with value 340 (num=340)
num = first_defined(NULL, v1,v2,v3,v4,v5,v6,NUM)
The result of num is 340.

**Q.  How does the force_error function work?**

If any mentioned conditions are not met, the force_error it forces an error. It will be useful when you want to stop the execution of a graph if it doesn't meet the set condition.

It will send the records to the reject port and error message to the error port.

**Q. What is a multifile system?**
Multifile is a set of directories on different nodes in a cluster. They possess an identical directory structure. The multifile system leads to better performance as it is parallel processing where the data resides on multiple disks.

It is created with the control partition on one node and data partitions on the other nodes to distribute the processing in order to improve the performance.

**Q. What is data encoding in Abinitio?**

Data needs to be kept classified in many cases, and it can be complete through this approach. It ensures information scraps in a form which no one else than the sender and the receiver can appreciate.
Data needs to be kept confidential in many cases and it can be done through this approach. It simply makes sure of information remains in a form which no one else than the sender and the receiver can understand.

In Abinitio, data encoding is an approach that is used to keep data confidential. In this approach, we ensure that the information remains in a form that cannot be understood by someone else other than the sender and the receiver.

**Q. Define the Evaluation of parameters order?**

- ➢ The host setup script will be executed first
- ➢ All Common parameters, that is, included, are evaluated
- ➢ All Sandbox parameters are evaluated
- ➢ The project script – project-start.ksh is executed
- ➢ All form parameters are evaluated
- ➢ Graph parameters are evaluated
- ➢ The Start Script of the graph is executed

**Q. Define local and formal parameter?**

Both are graph level parameters, but in local you need to initialize the worth at the time of announcement whereas globe no need to initialize the data it will produce at the time of operation the graph for that parameter.

**Q. Name the operations that support avoiding duplicate record?**

- Performing aggregation

- Utilizing the Rollup part

- Using Dedup class

**Q. What is the use of aggregation when we have rollup as we know rollup component in abinitio is used to summarize a group of data record? Then where we will use aggregation?**
Aggregation and Rollup both can summarize the data, but rollup is much more convenient to use. To understand how a particular summarization being rollup is much more explanatory compared to aggregate.

Rollup can do some other functionality like input and output filtering of records. Aggregate and rollup perform the same action, rollup display intermediate result in main memory, Aggregate does not support intermediate result.

**Q. Describe the elements you would review to ensure multiple scheduled batch jobs do not collide with each other?**

Because every job depends upon another job for example if your first job result is successfully then another job **will execute otherwise your job doesn't work.**

**Q. How to create a repository in abinitio for stand-alone system(LOCAL NT)?**

If you are trying to install the Ab -Initio on a stand-alone machine, then it is not necessary to create the repository, While installing It creates automatically for you under abinitio folder ( where you installing the Ab-Initio) If you are still not clear please ask your Question on the same portal.

**Q. Describe the process steps you would perform when defragmenting a data table. Does this table contain mission-critical data?**

**1**) We can move the table in the same or other tablespace and rebuild all the indexes on the table. alter table move this activity reclaims the defragmented space in the table analyze table table_name compute statistics to capture the updated statistics.

**2**) Reorg could be done by taking a dump of the table, truncate the table and import the dump back into the table.

**Q. Why might you create a stored procedure with the with recompile option?**

Recompile is useful when the tables referenced by the stored proc undergoes a lot of modification/deletion/addition of data. Due to the heavy modification activity, the execute plan becomes outdated and hence the stored proc performance goes down.

If we create the stored proc with recompile option, the sql server won't cache a plan for this stored proc and it will be recompiled every time it is run.

**Q. Explain PDL with an example?**

To make a graph behave dynamically, PDL is used

- Suppose there is a need to have a dynamic field that is to be added to a predefined DML while executing the graph
- Then a graph level parameter can be defined
- Utilize this parameter while embedding the DML in the output port.

**For Example:**

define a parameter named field with a value "string(" | ") name;"

Use ${mystring} at the time of embedding the DML in out port.

Use $substitution as an interpretation option.

**Q. What is BROADCASTING and REPLICATE?**

**Broadcast –** Takes data from multiple inputs, combines it and sends it to all the output ports. Eg – You have 2 incoming flows (This can be data parallelism or component parallelism) on Broadcast component, one with 10 records & other with 20 records. Then on all the outgoing flows (it can be any number of flows) will have $10 + 20 = 30$ records.

**Replicate: -** It replicates the data for a particular partition and send it out to multiple outports of the component but maintains the partition integrity.

**Example: -** Your incoming flow to replicate has a data parallelism level of 2. with one partition having 10 recs & other one having 20 recs. Now suppose you have 3 output flows from replicate. Then each flow will have 2 data partitions with 10 & 20 records respectively.

**Q. What are the benefits of data processing according to you?**

Well, processing of data derives a very large number of benefits. Users can put separate many factors that matter to them. In addition to this, with the help of this approach, one can easily keep up the pace simply by deriving data into different structures from an unstructured format.

In addition to this, the processing is useful in eliminating various bugs that are often associated with the data and cause problems in a later section.

It is because of no other reason than this, data processing has a wide application in several tasks.

**Q. What exactly do you understand with the term data processing and businesses can trust this approach?**

Processing is a procedure that simply covert the data from a useless form into a useful one without making a lot of efforts. However, the same may vary depending on factors such as the size of data and its format.

A sequence of operations is generally carried out to perform this task and depending on the type of data, this sequence could be automatic or manual. Because in the present scenario, most of the devices that perform this task are PC's automatic approach is more popular than ever before.

Users are free to obtain data in forms such as a table, vectors, images, graphs, charts and so on. This is the best things that business owners can simply enjoy.

**Q. Describe how you would ensure that database object definition (Tables, Indices, Constraints, Triggers, Users, Logins, Connection Options, and Server Options, etc) are consistent and repeatable between multiple database instances (i.e.: a test and production copy of a database) ?**

- Take an entire database backup and restore it in a different instance.
- Take statistics of all valid and invalid objects and match.
- Periodically refresh.

**Q. What is local and formal parameter?**

Two are graph level parameters but in local you need to initialize the value at the time of declaration whereas globe no need to initialize the data it will print at the time of running the graph for that parameter.

**Q. What Is Max Core of A Component?**

- MAX CORE is the space consumed by a component that is used for calculations
- Each component has different MAX COREs
- Component performances will be influenced by the MAX CORE's contribution
- The process may slow down / fasten if a wrong MAX CORE is set

**Q. How can we find old recovery files on Ab Initio production boxes?**

We can find the recovery files by using the following command: find . -name \*.rec

**Q. What Is the Function That Transfers A String Into A Decimal?**

Use decimal cast with the size in the transform() function, when the size of the string and decimal is same.

**Ex**: If the source field is defined as string(8).

- The destination is defined as decimal(8)
- Let us assume the field name is salary.
- The function is out.field :: (decimal(8)) in salary
- If the size of the destination field is lesser that the input then string_substring() function can be used.

**Ex:** Say the destination field is decimal(5) then use…

- out.field :: (decimal(5))string_lrtrim(string_substring(in.field,1,5))
- The ' lrtrim ' function is used to remove leading and trailing spaces in the string

## Q. Explain Pdl with An Example?

To make a graph behave dynamically, PDL is used

- Suppose there is a need to have a dynamic field that is to be added to a predefined DML while executing the graph
- Then a graph level parameter can be defined
- Utilize this parameter while embedding the DML in output port.

**Example: -** define a parameter named myfield with a value "string(" | "") name;"

- Use ${mystring} at the time of embedding the dml in out port.
- Use $substitution as an interpretation option

## Q. State the Working Process of Decimal_strip Function.

- A decimal strip takes the decimal values out of the data.
- It trims any leading zeros
- The result is a valid decimal number

**Ex:**

decimal_strip("-0184o") := "-184"

decimal_strip("oxyas97abc") := "97"

decimal_strip("+$78ab=-&^&%cdw") := "78"

decimal_strip("Honda") "0"

## Q. State the First_defined Function with An Example.

- This function is similar to the function NVL() in Oracle database
- It performs the first values which are not null among other values available in the function and assigns to the variable

**Example:**

A set of variables, say v1,v2,v3,v4,v5,v6 are assigned with NULL.

Another variable num is assigned with value 340 (num=340)

num = first_defined(NULL, v1,v2,v3,v4,v5,v6,NUM)

The result of num is 340

**Q. What Are the Operations That Support Avoiding Duplicate Record?**

Duplicate records can be avoided by using the following:

- Using Dedup sort
- Performing aggregation
- Utilizing the Rollup component

**Q. What Is A Deadlock And How It Occurs?**

- A graphical / program hand is known as deadlock.
- The progression of a program would be stopped when a dead lock occurs.
- Data flow pattern likely causes a deadlock
- If a graph flows diverge and converge in a single phase, it is potential for a deadlock
- A component might wait for the records to arrive on one flow during the flow converge, even though the unread data accumulates on others.
- In GDE version 1.8, the occurrence of a dead lock is very rare.

**Q. What is the component that can be used to lower a file in size?**

The components that can be used are Deflate and Compress.

**Q. Can a graph be infinitely run? If yes, how?**

Yes, it can, by calling the .ksh file in the end of the script.

**Q. What meaning has lock in Ab Initio?**

A graph must be locked in order to give permission only to the developer for editing it or the objects that belong to the graph. When other users try to lock the graph, they will be warned on the monitor that another user has already made a lock. This is a protection mechanism for developers.

**Q. What possible errors we can receive from a graph execution?**

The compilation errors in Abinition can appear when for example the output DML is not as the input DML.

**Q. What role has the XFR function?**

The purpose of this function is usually to store mappings, and this is very useful because the rewriting of code takes time and the XFR spares us of this effort.

**Q. What is the Difference between DML Expression and XFR Expression?**

DML represents the format of the metadata.
XFR represent the transform functions. which will contain business rules.

**Q. How much memory do we need for a graph?**

Some calculations lead to approximately eight MB plus the MAX_CORE and the phase file size.

**Q. How the term Standard Environment can be explained?**

The Standard Environment is the type of environment with more than one project that is private and another public one.

**Q. State the difference between checkpoint and phase?**

**CHECKPOINT:**

- The break of the procedure will be sustained following the checkpoint
- Data from the checkpoint is gathered and carry on executing after correction.
- When a graph fails in the centre of the procedure, a recovery point is shaped known as a Checkpoint.

**PHASE:**

- All the phases will sprint one by one.
- The middle file will be deleted.
- If a graph is shaped with phases, each phase is assigned to some division of memory one after another.

**Q. What is the difference between check point and phase in Ab Initio?**

| Check point | Phase |
|---|---|
| A check point is a recovery point that is created when a graph fails in the middle of the process. | A graph consists of phases. If a graph is created with phases, each phase is assigned to some part of the memory. |
| The rest of the process will be continued after the check point. | All the phases run one by one. |
| Data from the check point is fetched and continues to execute after correction. | In phase, the intermediate file will be deleted. |

**CREATE command**

This command is used to create a new database, table, view, or index in a database management system (DBMS). It allows you to define the structure of the object being created, such as the columns in a table.

**Example**

```
CREATE TABLE users (

  id INT PRIMARY KEY,

  name VARCHAR(50) NOT NULL,

  email VARCHAR(100) UNIQUE

);
```

This command creates a new table named "users" with three columns: "id", "name", and "email". The "id" column is set as the primary key, the "name" column is set to not allow NULL values, and the "email" column is set as a unique key.

**SELECT command**

This command is used to retrieve data from one or more tables in a database. It allows you to specify the columns you want to retrieve and filter the results using conditions.

**Example**

```
SELECT name, email

FROM users

WHERE id = 1;
```

This command retrieves the "name" and "email" columns from the "users" table where the "id" column is equal to 1.

### ALTER command

This command is used to modify the structure of a database object. It can be used to add or remove columns from a table, change the data type of a column, or rename a table.

**Example**

ALTER TABLE users

ADD COLUMN phone VARCHAR(20);

This command adds a new column named "phone" to the "users" table.


### DELETE command

The DELETE command is used to **remove rows from a table** in a database.

**Example: -**

DELETE FROM users WHERE id = 1;

This command removes the row(s) from the "users" table where the "id" column is equal to 1. If the condition is not specified, the DELETE command will remove all rows from the table. It's important to use this command with caution and double-check the conditions before executing it, as it can permanently delete data from the database. It's also worth noting that the DELETE command only removes the data from the table but doesn't free up the space on disk. To free up the space on disk, you may need to run a vacuum command or rebuild the table.


### DROP command

This command is **used to delete a database**, table, view, or index from a DBMS. It permanently removes the object and its data from the database.

It is DML command, generally used to delete a record, clusters or tables. Rollback command can be performed, to retrieve the earlier deleted things. To make deleted things permanently, "commit" command should be used.

**Example**

**DROP TABLE users;**

This command deletes the "users" table from the database.

**TRUNCATE command**

This command is **used to remove all rows from a table but leaves the table structure intact**. It is a faster way to delete all data from a table compared to using the DELETE command.

It is a DDL command, used to delete tables or clusters. Since it is a DDL command hence it is auto commit and Rollback can't be performed. It is faster than delete.

**Example**

TRUNCATE TABLE users;

This command removes all data from the "users" table, but leaves the table structure intact.


**Q. What are the primary keys and foreign keys?**


In RDBMS the relationship between the two tables is represented as Primary key and foreign key relationship. Whereas the primary key table is the parent table, and the foreign key table is the child table.


**PRIMARY KEY**

This is a constraint that is used to uniquely identify each row in a table. It is a column or set of columns that are used as the primary key for the table.

**Example**

```
CREATE TABLE orders (
  id INT PRIMARY KEY,
  user_id INT,
  total DECIMAL(10,2),
  FOREIGN KEY (user_id) REFERENCES users(id)
);
```


This command creates a new table named "orders" with three columns: "id", "user_id", and "total". The "id" column is set as the primary key and the "user_id" column is set as a foreign key that references the "id" column in the "users" table.

**FOREIGN KEY**

This is a constraint that is used to link two tables together in a database. It is a column or set of columns in one table that refers to the primary key of another table, creating a relationship between the two tables. It is used to maintain data integrity by ensuring that the values in the foreign key column(s) match the values in the primary key column(s) they are referring to.

**Example**

ALTER TABLE orders

ADD FOREIGN KEY (product_id)

REFERENCES products(id);

This command adds a new foreign key to the "orders" table that references the "id" column in the "products" table.

**Q. Describe the Grant/Revoke DDL facility and how it is implemented?**

- GRANT means **permissions** for example GRANT CREATE TABLE, CREATE VIEW AND MANY MORE.
- REVOKE means **cancel the grant (permissions).** So, Grant or Revoke both commands depend upon D.B.A.

**Q. When using multiple DML statements to perform a single unit of work, is it preferable to use implicit or explicit transactions, and why?**

Because implicit is using for internal processing and explicit is using for using open data required.

**Q. Explain an outer join?**

An outer join is used when one wants to select all the records from a port – whether it has satisfied the join criteria or not.

### Linux Commands

**PWD Command: -**

PWD stands for Print Working Directory

This command is used to show the current working directory.

**For example**  /home/foobar

**ls Command: -**

ls a will list all files including hidden files (files with names beginning with a dot).

**ls –ltr   (ls -lrt) Command:-**

This command will list you all files according to the order of time in which they were created. Here "ltr" stands for l- long listing, t- time, r- recursive. The list displayed contains a file name, file permissions, owner of the file, group, date and time of file creation and links.

**ls -lart Command: -**

this Command is used to display hidden file or directories.

**ls command: -**

The ls command is one of the most commonly used commands in daily Linux/UNIX operations.

**ls -l Command: -**

The "ls -l" option **displays the contents of the current directory in a long listing format, one per line**. The line begin with the file or directory permission, owner and group name, file size, created/modified date and time, file/folder name as some of the attributes.

**mkdir command**

mkdir command in Linux **allows the user to create Multiple directories** .

Just **type "mkdir <dir name> , in place of <dir name> type the name of new directory, you want to create and then press enter**.

**Example:** mkdir created. by using. (dot) you can create hidden directories**.**

**For example**:-  mkdir .abdc


**touch command**

touch filename. We can use the touch command to create new files by giving file names as the input. We can create multiple files

**Example:-** touch a.dat b.dat c.txt by using. (dot) you can create hidden files.

**Example:-** touch .abdc


**cat Command**

Using the cat command, you can quickly create a file and put text into it. To do that, use the > redirect operator to redirect the text in the file.
cat > filename.txt
The file is created, and you can begin populating it with text. To add multiple lines of text just press **Enter** at the end of each line.  Once you're done, hit **CTRL+D** to exit the file.



**Cp command**
'cp' means copy. 'cp' command is used to copy a file or a directory.
copy a file into the same directory
**syntax** cp <existing file name> <new file name>

**cat command**

The 'cat' command can be used to display the content of a file.

**Syntax:** cat <fileName>

| Option | Function |
| --- | --- |
| cat > [fileName] | To create a file. |
| cat [oldfile] > [newfile] | To copy content from older to new file. |
| cat [file1 file2 and so on] > [new file name] | To concatenate contents of multiple files into one. |
| cat -n/cat -b [fileName] | To display line numbers. |
| cat -e [fileName] | To display $ character at the end of each line. |
| cat [fileName] <<EOF | Used as page end marker. |

**wc Command**

Linux wc command helps in counting the lines, words, and characters in a file. It displays the number of lines, number of characters, and the number of words in a file. Mostly, it is used with pipes for counting operation.

**Syntax:**

wc [OPTION]... [FILE]...

 wc [OPTION]... --files0-from=F

**rm Command**

The Delete File 'rm' means remove. This command is used to remove a file. The command line doesn't have a recycle bin or trash unlike other GUI's to recover the files. Hence, be very much careful while using this command. Once you have deleted a file, it is removed permanently.

**Syntax:**

rm <filename>

Example:

rm myfile1

**head Command**

The 'head' command displays the starting content of a file. By default, it displays starting 10 lines of any file.

**Syntax:** head <file name>

**Example:**

head jtp.txt

**tail Command**

Linux tail command is used to display the last ten lines of one or more files. Its main purpose is to read the error message. By default, it displays the last ten lines of a file. Additionally, it is used to monitor the file changes in real-time. It is a complementary command of the head command.

**Syntax:** tail <file name>

1. **ls:** Used to list the files and directories in a directory.
   **Example:** "ls /home/user" will list all the files and directories in the "/home/user" directory.

2. **cd:** Used to change the current working directory.
   **Example:** "cd /home/user" will change the current working directory to "/home/user".

3. **pwd:** Used to print the current working directory.
   **Example:** "pwd" will display the absolute path of the current working directory.

4. **cp:** Used to copy files and directories.
   **Example:** "cp file.txt /home/user/backup" will copy the file "file.txt" to the "/home/user/backup" directory.

5. **mv:** Used to move or rename files and directories.
   **Example:** "mv file.txt /home/user/backup" will move the file "file.txt" to the "/home/user/backup" directory.

6. **rm:** Used to delete files and directories.
   **Example:** "rm file.txt" will delete the file "file.txt".

7. **mkdir:** Used to create a new directory.

    **Example:** "mkdir /home/user/new_dir" will create a new directory called "new_dir" in the "/home/user" directory.

8. **rmdir:** Used to remove an empty directory.

    **Example:** "rmdir /home/user/empty_dir" will remove the directory "empty_dir" if it is empty.

9. **cat:** Used to display the contents of a file.

    **Example:** "cat file.txt" will display the contents of the file "file.txt".

10. **less:** Used to view the contents of a file page by page.

    **Example:** "less file.txt" will display the contents of the file "file.txt" page by page.

11. **head:** Used to display the first few lines of a file.

    **Example:** "head -n 5 file.txt" will display the first 5 lines of the file "file.txt".

12. **tail:** Used to display the last few lines of a file.

    **Example:** "tail -n 5 file.txt" will display the last 5 lines of the file "file.txt".

13. **grep**: Used to search for a pattern in a file.

    **Example:** "grep 'error' file.txt" will display all the lines in the file "file.txt" that contain the word "error".

14. **find:** Used to search for files and directories.

    **Example:** "find /home/user -name '*.txt'" will search for all ".txt" files in the "/home/user" directory and its subdirectories.

15. **chmod:** Used to change the permissions of files and directories.

    **Example:** "chmod 755 file.txt" will give the owner full permissions, and all others read and execute permissions for the file "file.txt".

1. **cd:** This command is used to **change the current directory.**

   **For example**, "cd /home/user/" will change the current directory to the user's home directory.

2. **ls:** This command is used to **list the files** in the current directory.

   **For example,** "ls -l" will list the files in the current directory in a long format.

3. **chmod:** This command is used to **change the permissions of a file or directory**.

   **For example,** "chmod 755 file.txt" will give read, write, and execute permissions to the owner of the file and read and execute permissions to all other users.

4. **rm:** This command is used to **remove a file or directory**.

   **For example,** "rm file.txt" will remove the file named "file.txt".

5. **mkdir:** This command is used to **create a new directory**.

   **For example,** "mkdir mydirectory" will create a new directory named "mydirectory".

6. **rmdir**: This command is used **to remove a directory**.

   **For example,** "rmdir mydirectory" will remove the directory named "mydirectory".

7. **tar**: This command is used **to create, manipulate, and extract files from a tar archive**. **For example,** "tar -cvf archive.tar file1 file2" will create a new tar archive containing files "file1" and "file2".

8. **scp**: This command is **used to securely copy files between hosts**.

   **For example,** "scp file.txt [user@remote.host](mailto:user@remote.host):/remote/directory/" will copy the file "file.txt" to the remote host at directory "/remote/directory/".

9. **ssh**: This command is used to connect to a remote host over a secure shell.

   **For example,** "ssh [user@remote.host](mailto:user@remote.host)" will establish a secure shell connection to the remote host as the user "user".

10. **cat**: This command is **used to display the contents of a file**.

    **For example,** "cat file.txt" will display the contents of the file "file.txt" on the console.

11. **tail**: This command is used to **display the last few lines of a file**.

    **For example,** "tail -n 10 file.txt" will display the last 10 lines of the file "file.txt".

12. **head**: This command is used to **display the first few lines of a file**.

    **For example,** "head -n 5 file.txt" will display the first 5 lines of the file "file.txt".

13. **grep**: This command **is used to search for a pattern** in a file.

    **For example,** "grep 'pattern' file.txt" will search for the pattern "pattern" in the file "file.txt".

14. **ps**: This command is used to **display the currently running processes.**

    **For example,** "ps aux" will display all the running processes along with their details.

15. **kill**: This command is **used to terminate a running process**.

    **For example,** "kill -9 PID" will terminate the process with the process ID "PID".

16. **top**: This command is used to **display the system resource usage and the currently running processes in real-time**.

    **For example,** running "top" will display the resource usage of the system and the processes that are currently running.

17. **ifconfig**: This command is used to display the network interface configuration.

    **For example,** "ifconfig eth0" will display the configuration details of the network interface "eth0".

18. **netstat**: This command is used to display the network connections and their status.

    **For example,** "netstat -an" will display all the active network connections and their status.

19. **df**: This command is used to **display the disk space usage**.

    **For example,** "df -h" will display the disk space usage in a human-readable format.

20. **du**: This command is used to **display the disk usage of a file or directory**.

    **For example,** "du -sh /path/to/directory" will display the disk usage of the directory in a human-readable format.

21. **find**: This command is used to search for files and directories based on various criteria. **For example,** "find /path/to/search -name 'filename'" will search for the file named "filename" in the directory "/path/to/search" and its subdirectories.

22. **tar**: This command is also used to **extract files from a tar archive**.

    **For example,** "tar -xvf archive.tar" will extract the files from the tar archive "archive.tar".

23. **ping**: This command is **used to test the connectivity to a network host**.

    **For example,** "ping google.com" will test the connectivity to the website "google.com".

24. **ifup/ifdown:** These commands are used to bring up or take down a network interface. **For example, "**ifup eth0" will bring up the network interface "eth0".

25. **uname:** This command is used to **display the system information**.

    **For example,** "uname -

26. **echo**: This command is used to **display text on the console or to redirect text to a file**. **For example, "**echo 'Hello, World!'" will display the text "Hello, World!" on the console.

27. **wc**: This command is **used to count the number of lines,** words, and characters in a file. **For example,** "wc -l file.txt" will count the number of lines in the file "file.txt".

28. **sort**: This command is **used to sort the contents of a file**.

    **For example,** "sort file.txt" will sort the lines in the file "file.txt" in ascending order.

29. **uniq**: This command is **used to remove duplicate lines from a file.**

    **For example**, "uniq file.txt" will remove the duplicate lines in the file "file.txt".

30. **diff**: This command is used to compare the contents of two files**.**

    **For example**, "diff file1.txt file2.txt" will **display the differences between the files** "file1.txt" and "file2.txt".

31. **tar**: This command can also be **used to create a compressed archive of files**.

    **For example,** "tar -czvf archive.tar.gz file1 file2" will create a compressed tar archive "archive.tar.gz" containing files "file1" and "file2".

32. **chown**: This command is used to **change the owner of a file or directory**.

    **For example,** "chown user:group file.txt" will change the owner of the file "file.txt" to the user "user" and the group "group".

33. **top**: This command can be used to monitor the system resource usage and the currently running processes in real-time.

34. **chmod**: This command is **used to change the file permission modes.**

    **For example,** "chmod 755 file.txt" will give the owner of the file full permission and everyone else read and execute permission.

35. **scp**: This command is used to copy files securely between hosts over the network.

    **For example,** "scp file.txt user@hostname:/path/to/destination" will copy the file "file.txt" to the remote host "hostname" to the specified path.

36. **ssh**: This command is used to securely connect to a remote host over the network**.**

   **For example,** "ssh user@hostname" will connect to the remote host "hostname" as the user "user".

37. **sed**: This command is used to perform text transformations on a file or stream.

   **For example,** "sed 's/old_text/new_text/g' file.txt" will replace all occurrences of "old_text" with "new_text" in the file "file.txt".

38. **awk**: This command is used to process and manipulate text data.

   **For example, "**awk '{print $1,$3}' file.txt" will print the first and third fields of each line in the file "file.txt".

39. **crontab**: This command is used to schedule tasks to run at specific times**.**

   **For example**, "crontab -e" will open the crontab editor where you can schedule tasks to run at specific intervals.

40. **history**: This command **is used to display the command history of the current session. For example, "**history | tail -n 10" will display the last 10 commands that were executed in the current session.

41. **cut**: This command **is used to extract specific columns from a file.**

   **For example, "**cut -d',' -f1,3 file.csv" will extract the first and third columns of the comma-separated file "file.csv".

42. **grep**: This command is **used to search for patterns in a file**.

   **For example,** "grep 'pattern' file.txt" will search for the pattern "pattern" in the file "file.txt".

43. **diff**: This command is **used to compare two files line by line.**

   **For example**, "diff file1.txt file2.txt" will display the differences between the files "file1.txt" and "file2.txt".

44. **kill**: This command is used to terminate a running process.

   **For example**, "kill PID" will terminate the process with the specified PID.

45. **ps**: This command is used to display the currently running processes.

   **For example**, "ps -ef" will display all the currently running processes.

46. **netstat**: This command is used to display network connections and statistics.

   **For example,** "netstat -an" will display all the current network connections.

47. **wget**: This command is used to download files from the internet.

   **For example**, "wget https://example.com/file.txt" will download the file "file.txt" from the website "example.com".

50. **curl**: This command is used to transfer data from or to a server. It can be used to download files or upload data to a server.

   **For example,** "curl -O https://example.com/file.txt" will download the file "file.txt" from the website "example.com".

51. **awk**: This command is used for text processing and data extraction. It can be used to extract specific columns from a file or to perform calculations on the data.

   **For example,** "awk '{print $1,$3}' file.txt" will print the first and third columns of each line in the file "file.txt".

52. **sed**: This command is used to perform text transformations on a file or stream. It can be used to find and replace text or to delete lines from a file.

   **For example,** "sed 's/old_text/new_text/g' file.txt" will replace all occurrences of "old_text" with "new_text" in the file "file.txt".

53. **tee**: This command is used to redirect output to a file and also to the console.

   **For example,** "ls | tee file.txt" will list the files in the current directory and also redirect the output to the file "file.txt".

54. **tailf**: This command is used to display the last few lines of a file and also to follow the file as it grows.

   **For example,** "tailf file.txt" will display the last few lines of the file "file.txt" and also keep following the file as new data is added to it.

55. **ping**: This command is used to test the connectivity to a remote host by sending packets to the host and measuring the response time**.**

   **For example**, "ping google.com" will test the connectivity to the website "google.com".

56. **du**: This command is used to display the disk usage of files and directories.

   **For example,** "du -sh /path/to/directory" will display the total disk usage of the directory at the specified path.

57. **df**: This command is used to display the disk space usage of the file system.

   **For example,** "df -h" will display the disk space usage in a human-readable format.