

Ab initio Session 14

Introduction to Ab Initio



Ab Initio Training

1



- **Concepts of Parallelism**
- **Explanation of Data partitioning**
- **Concept of Repartitioning**

Forms of Parallelism



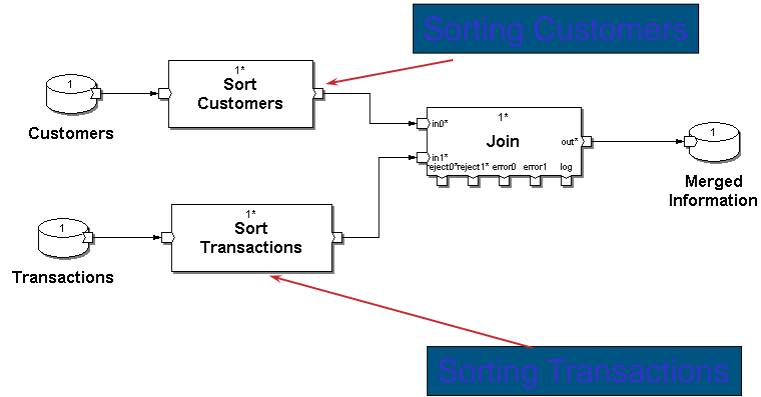
- **Component parallelism**
- **Pipeline parallelism**
- **Data parallelism**

Component parallelism



- A graph with multiple processes running simultaneously on separate data uses component parallelism
- In this two or more components process the records in parallel.

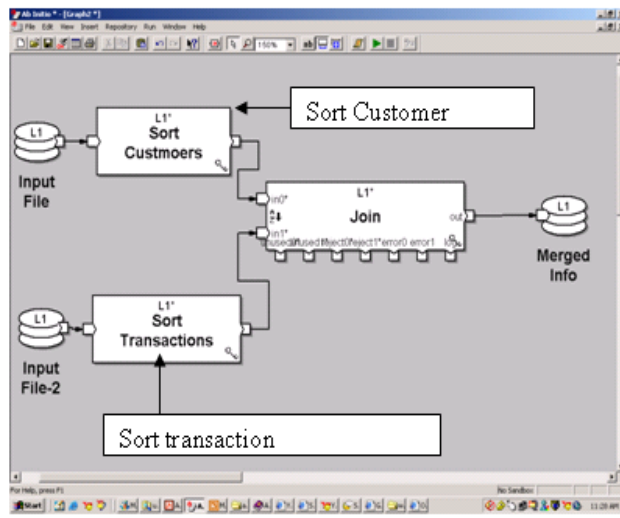
Component Parallelism



CapGemini

Ab Initio Training

5



CapGemini

Ab Initio Training

6

Pipeline Parallelism



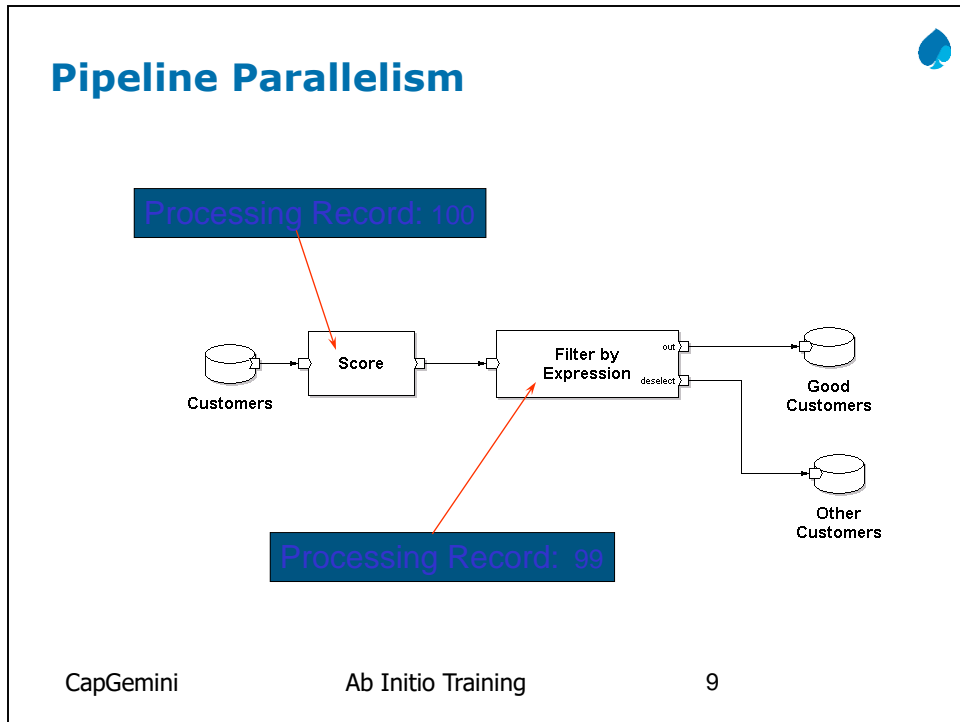
- - A graph with multiple components running simultaneously on the same data uses pipeline parallelism. Each component in the pipeline continuously reads from upstream components, processes data, and writes to downstream components. Since a downstream component can process records previously written by an upstream component, both components can operate in parallel. NOTE: To limit the number of components running simultaneously, set phases in the graph.

Pipeline Parallelism-cont.



- 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.

NOTE: To limit the number of components running simultaneously, set phases in the graph.



Pipeline Parallelism-cont.



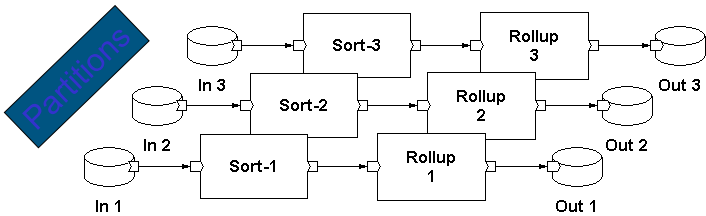
- In this the records are processed in pipeline, i.e. the components do not have to wait for all the records to be processed. The records that got processed are passed to next component in pipeline.

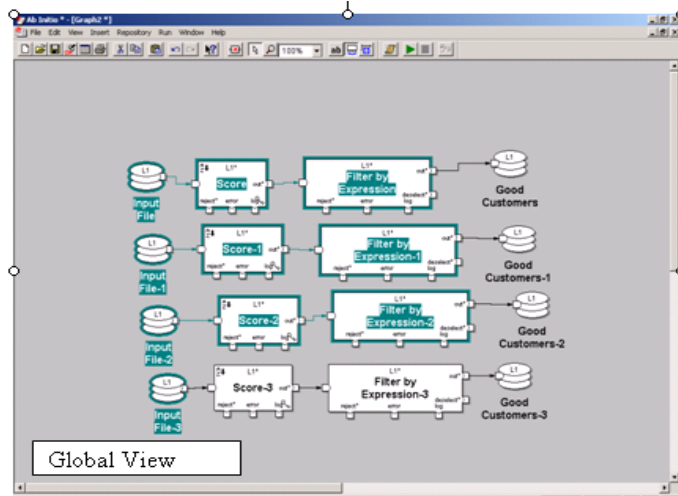
Data Parallelism



- A graph 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. To support this form of parallelism, Ab Initio provides Partition components to segment data, and Departition components to merge segmented data back together.
- Partitioning is an example of data parallelism.

Data Parallelism





CapGemini

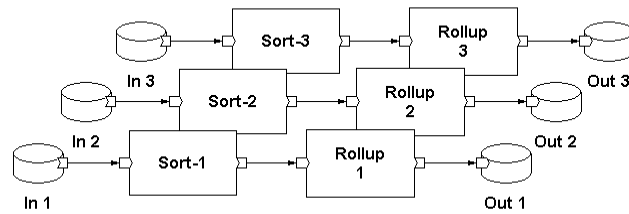
Ab Initio Training

13

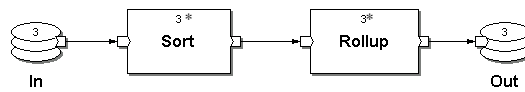
Two Ways of Looking at Data Parallelism



Expanded View:



Global View:



CapGemini

Ab Initio Training

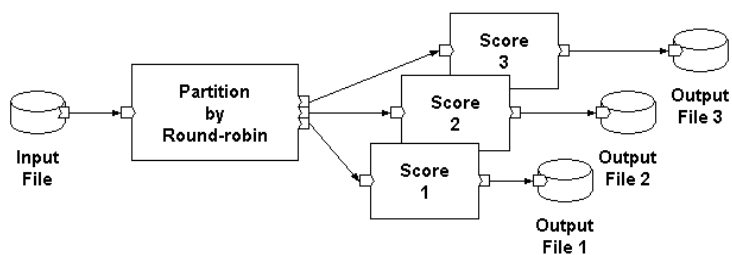
14

Data Parallelism



- Scales with data.
- Requires *data partitioning*.
- Dependent upon the application, different partitioning methods are available.

A Data Parallel Application: The Expanded View

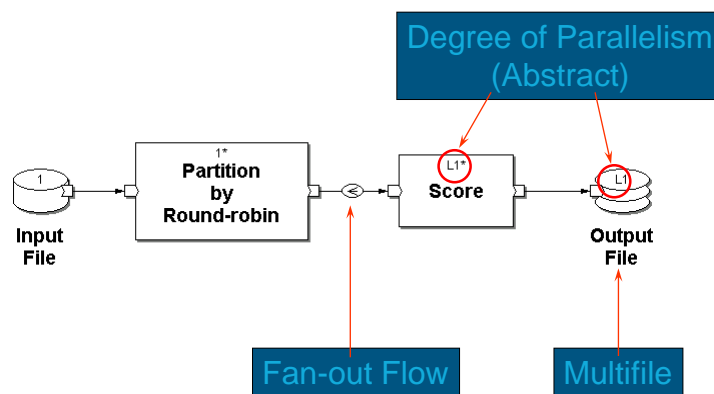


CapGemini

Ab Initio Training

16

A Data Parallel Application: The Global View



CapGemini

Ab Initio Training

17



Thank You