

Unit 4.

Distributed and Multimedia IR.

- **Distributed IR:**

- Distributed Computing is the application of multiple computers connected by a network to solve a single problem
- Distributed Computing system can be viewed as a MIMD parallel processor with a relatively slow inter-processor communication channel and freedom to employ a heterogeneous collection of processor in the system.
- Distributed IR is the practice of distributing the processes and tasks involved in IR across multiple computers or nodes in network.
- This approach often used to improve efficiency and scalability of large-scale IR system by dividing the workload and processing tasks in parallel.
- Distributed IR commonly applied in search engines where indexing, query processing and result ranking can be distributed to optimize performance.
- Computational model is very similar to the MIMD parallel processing model but difference here is that:
 1. Sub-tasks runs on different computers and the communication between the subtasks is performed using network protocol such as TCP/IP rather than the shared memory based inter-process communication.

tion mechanism.

2. In a distributed system it is more common to employ a procedure for selecting a subset of the distributed servers for processing a particular request rather than broadcasting every request to every server in the system.

* Architecture of Distributed IR

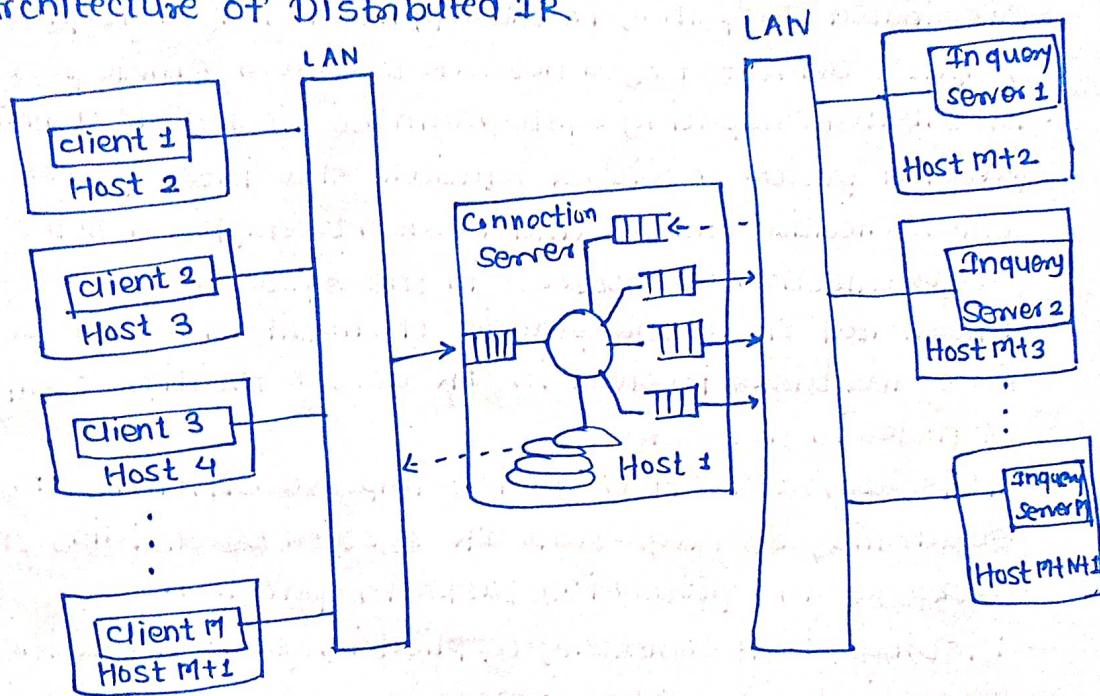


Fig. Architecture of Distributed IR.

- At a minimum, the protocol should allow a client to:
 - Obtain information about search servers, e.g. a list of databases available for searching at the server and possibly statistics associated with the databases.
 - Submit a search request for one or more databases using a well defined query language.

- Receive Search results in a well defined format.
- Retrieve items identified in the search results.
- Example of DIR: (Google search Engine)
- Federated search

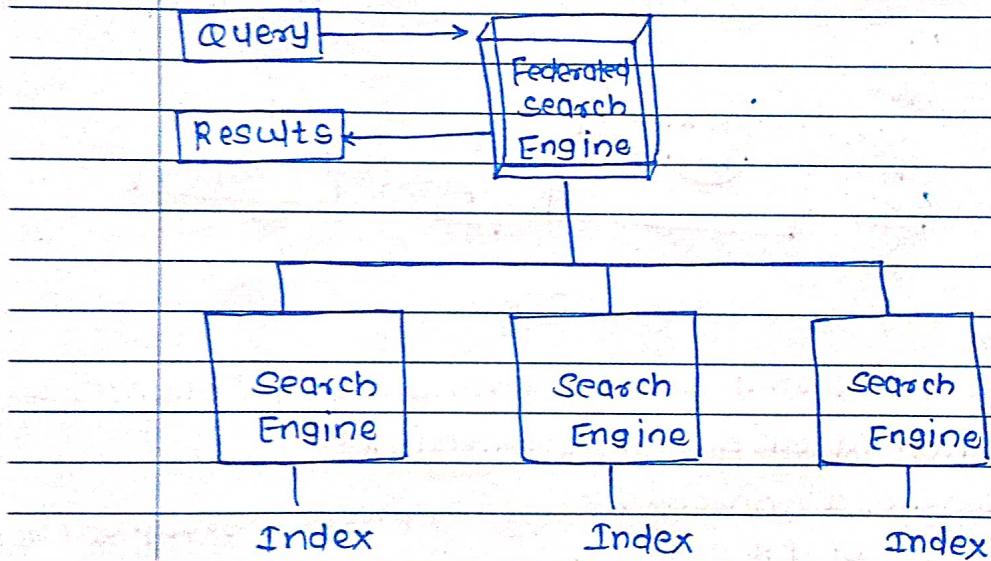


Fig. Example of DIR
(Google search Engine)

- Federated search is another name of DIR.
- It has its own collection of data.
- Instead of crawling a site, Federated search systems send a user query to a resource's search tools.
- Users can search multiple data sources simultaneously using Federated search by sending a single query.
- The Federator gathers result from one or more search

engines and then present all the results in the single user terminal.

- Architecture of Distributed IR

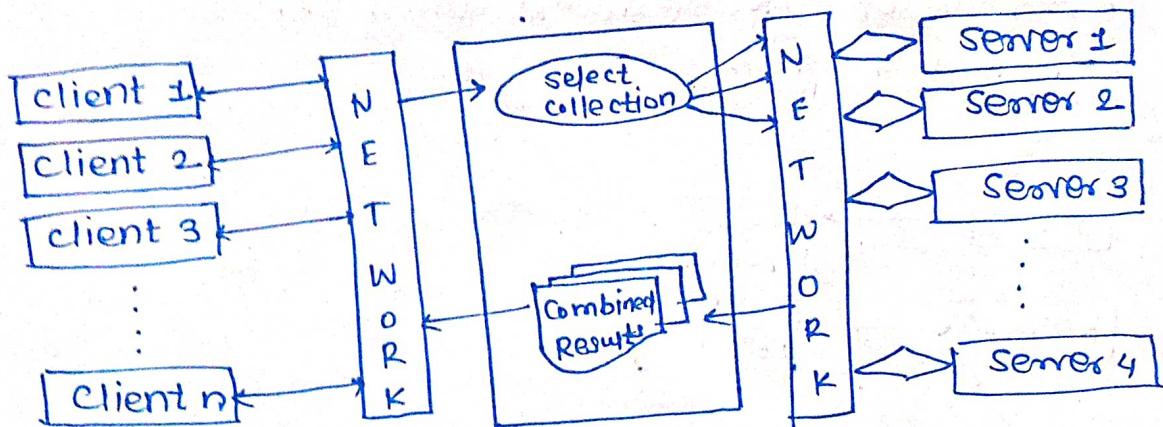


Fig. Architecture of DIR.

- DIR conceptual model that enables several users (clients) to continuously search various document collections.
- In DIR architecture, a connection server connects a group of clients (users) to a group of IR system.
- Distributed system typically consist of a set of server processes, each running on a separate processor node, and a designated broker process for accepting client requests, distributing the requests to the servers, collecting intermediate results from the servers, and combining the intermediate into a final result for the client.
- In DIR system, prototype that enables several users to simultaneously search various documents search.

- All communications betn clients and (servers) IR systems are handled by connection servers.
- Additionally, it is simply for us to modify the architecture and assess its success. Our findings demonstrate that our architecture can support high number of clients.

• Collection Partitioning:

* What exactly a collection?

- DIR, when there are multiple collections, searching.
- Local contexts, for instance, partitioning a sizable collection.
- Wide area settings, such as company network or internet.
- Divide huge collections among processors to speed up processing - due to governmental or regulatory requirements.
- Heterogenous environments, many IR systems, networks with hundreds or thousands of collections and collections that are indexed on the web.
- Economic expenses of conducting a full site search.
- Economic costs of conducting a full network search.

Collection Partitioning

partitioning of collections in a decentralized system	partitioning of collections in a centralized system.
---	--

1. Partitioning of collections of in a Decentralized System:-

- Distributed document collections in a system with independently running, heterogeneous search engines will be made and kept separately.
- There is no centralized authority over the document splitting procedure.
- It's feasible that each search engine has a specific area of expertise.

2. Partitioning collection in a centralized system:-

- When the collection is small enough to fit on a single search server but high availability and query processing throughput are required, collection partitioning in a centralized system is appropriate.
- The goal of the broker is to direct inquiries to the search servers and balance the loads on the servers, which is how the parallelism in the system is being utilized through multitasking.
- Every query is broadcast to every search server by the broker, who then compiles the results for the user.
- Explicit semantic partitioning of the documents, which may be done either manually or automatically.
- The articles may already have been organized into collections that make sense semantically, like groupings by technical discipline.

● Source Selection:

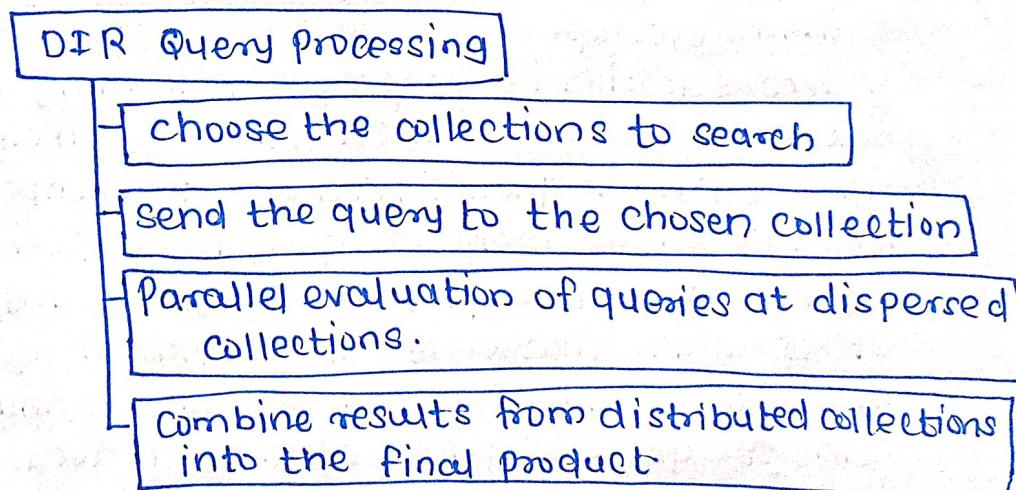
- The process of choosing which of the distributed document collections is most likely to contain document relevant to the present query and as a result should receive the query for processing is known as Source Selection.
- One technique is to simply broadcast the query to all collections, assuming that every collection has an equal likelihood of having relevant items.
- This approach is suitable when are randomly sorted into groups or when there is a significant amount of semantic overlap.
- Source selection is the process of deciding which scattered document collection is most likely to contain the current query's pertinent documents (and therefore should receive the query for processing).
- It is suitable to employ the straightforward method of assuming that every collection has an equal possibility of containing relevant documents and broadcasting the query to all of them when documents are randomly partitioned or when there is sufficient semantic overlap between the collections.
- The collections can also be ranked based on how likely they are to succeed.
- This is feasible if texts are grouped into groupings that have semantic importance, otherwise, searching through each collections individually would be excessively expensive.

- The basic strategy is to generate a collection vector for each collection and compare it to the query vector treating each collection as if it were a single, enormous document.

- Query Processing:

A DIR system's query processing

1. choose the collections to search.
2. send the query to the chosen collections.
3. Parallel evaluation of queries at dispersed collections.
4. Combine results from the distributed collections into final product.



.Fig. DIR Query Processing.

- Merging the Results:

There are various situations for merging the results:

- The final result set is the union of the result sets if the query is boolean and the search servers return Boolean result sets.

- Several strategies are available, ranging from straight-forward to sophisticated and accurate, if the inquiry incorporates free-text ranking.
- The earliest method is to use round robin interleaving to merge the ranking result lists.
 1. 1st document from the first list
 2. 1st document from the second list.

N: 1st document from Nth list.

N+1: 2nd document from the first list.

- Because hits from irrelevant collections are assigned the same status as hits from highly relevant collections, this method is likely to produce low-quality results.
- It is better to combine the result lists depending on relevance score.
- Re-ranking is necessary if the document collections are semantically divided or are managed by separate parties.
- We risk getting the wrong findings if adequate global term statistics are not employed to calculate the document ratings.

You can calculate a collection's weight using formula:

$$W = 1 + |C| \cdot (S - S^{\bar{}}) / S^{\bar{}}$$

where, |C| - no. of collections you searched.

S - Collection score , $S^{\bar{}}$ - mean of those scores.

• Issues in DIR

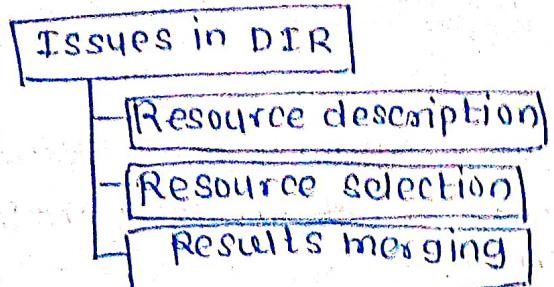


Fig. Issues in DIR

- The scattered location and control of information in a wide area computer network are reflected in that multi-database model of DIR.
- However, it is also more complicated than the IR model using single database necessitating the solution to a number of additional issues:
 1. Resource description:
 - Each text databases contents must be explained.
 2. Resource Selection:
 - choosing which database(s) to search require consideration of an information demand and list of resource descriptions.
 3. Merging results:
 - combining ranked lists that each database returned to create a single cohesive ranked list.
 - During a 5 year period, addresses many of problems that arises in DIR system.
 1. Difficulty of briefly explaining the contents of each database or resource that is available.
 2. Proposes a scoring system for database based on how well they will likely meet information needs.
 3. issue of combining the results produced by various search systems is covered. Examines the process by which a DIR system in a multi-party setting. obtains resources descriptions.

Multimedia TR

Introduction

Multimedia Information Retrieval (MIR) involves the retrieval of various forms of multimedia content such as images, audio, video, and text from large datasets. It combines techniques from information retrieval, computer vision, audio processing and natural language processing to enable efficient and effective searching, organizing, and retrieving of multimedia data.

* Data Modelling

Data Modelling in multimedia information retrieval (MIR) refers to the process of creating structured representations of multimedia data in a way that allows efficient storage, organization, and retrieval of the content.

In MIR, data modelling involves:

- 1) Feature Extraction : Extracting meaningful features from multimedia data. For images, this could be color, texture, or shape features

2) Feature Representation

Representing these extracted features as numerical values or vectors, which can be used for indexing and retrieval.

3) Indexing : Storing these feature representation in a way that allows for efficient retrieval. Indexing structures may include databases, trees, or hash tables.

4) Semantic Information:

Incorporating semantic information if needed. This could involve associating metadata, tags, or annotations with the multimedia data.

5) Query processing : Developing algorithms & methods for comparing query features with the stored representation to retrieve relevant multimedia content.

* Query Language

- In Multimedia Information Retrieval (MIR), a query language is a formal way for users to express their information needs when searching for multimedia content within a database.
- A query language in MIR typically provides ways to express:

i) Content-based queries : users can describe the visual or auditory characteristics they are seeking in the multimedia content. For instance, they might specify

colors, shapes, or audio patterns.

- 2) Semantic Queries : users can express their information needs using terms or concepts related to the content they're looking for.
- 3) Spatial and Temporal Constraints : In the case of videos or images, users might specify spatial regions of interest or time intervals for their search.
- 4) Combination of Criteria: users can often combine different criteria to refine their search, such as specifying both visual & textual requirements.

* Background - Spatial Access Method

- The Background-Spatial Access Method (BSAM) is an indexing technique used in Multi-media Information Retrieval (MIR) to efficiently retrieve multimedia content, particularly images, based on their spatial characteristics within the image.

- BSAM is designed to work with content-based image retrieval systems, where users can specify spatial relationships between objects within an image.

- In BSAM, images are divided into multiple regions, and each region is described using relevant visual features. These features might include color histograms, texture descriptors, or other attributes that capture the characteristics of the image content. The spatial relationships between regions are also recorded.
- The key idea behind BSAM is to create an index that allows for efficient retrieval of images based on their spatial arrangement. This is particularly useful when users want to find images where certain objects or features are located in specific position relative to each other.

* Generic Multimedia Indexing Approach

- The Generic Multimedia Indexing Approach (GMIA) is a framework or methodology used in Multimedia Information Retrieval (MIR) to provide a unified way of indexing & retrieving various types of multimedia content such as images, videos, and text.
- The goal of GMIA is to create a common indexing structure & techniques that can be applied across different types of multimedia, making it easier to manage and search multimedia data.

GMIA typically involves the following steps:

- 1) Feature Extraction
- 2) Feature Representation
- 3) Indexing
- 4) Query Processing
- 5) Semantic Information

- The aim of Generic Multimedia Indexing Approach is to provide a standardized approach that can handle different types of media while sharing common indexing and retrieval techniques. This helps reduce the complexity of managing multimedia data and enables more efficient search and retrieval in multimedia databases.

* One Dimensional Time Series

- A one dimensional time series refers to a sequence of data points collected over time, where each data point corresponds to a single numerical value of feature at a specific time instance.

- One-dimensional time series are commonly used to represent temporal patterns & variations in multimedia data, such as audio signals, video frames, or other time-dependent information.

For example:

• Heart Rate: A person's heart rate measured at regular intervals can be represented as a one-dimensional time series, showing variations in heart rate over time.

One-dimensional time series in MIR are used for various purposes:

- 1) Pattern Recognition: Detecting recurring patterns or trends within the time series data.
- 2) Similarity Search: finding other time series that are similar in behaviour to a given time series.
- 3) Event Detection: Identifying specific events or transitions in the time series.
- 4) Classification: Assigning a category or label to a time series based on its behaviour.
- 5) Anomaly Detection: Detecting unusual or unexpected behaviour in the time series.

* Two-Dimensional color Images

- Two-dimensional color Images are a type of multimedia data that represent visual content using two dimensions (width & height) & include information about color.

In a two dimensional color image?

- 1) Two dimensional : The image is structured with a width and a height, forming a grid of pixels. Each pixel is located at a specific position within the grid.
 - 2) Color Information : Each pixel in the image is associated with color information. This color information can be represented using various color models such as RGB (Red, Green, Blue), CMYK (Cyan, Magenta, Yellow, Black), or HSV (Hue, Saturation, Value).
- Two dimensional color images are commonly used in various applications within MIR.
 - o Image Retrieval : Two-dimensional color images are often used as queries or target content for retrieval systems, where users search for images similar to a given query image.
 - o Content Analysis : Color information is valuable for content analysis tasks, such as identifying dominant colors, detecting objects, or recognizing patterns in images.

- **Visual Search & Image-based search engines** use color features to help users find images that match their visual preferences.
- **Object Recognition**: color features can contribute to identifying and categorizing objects within images.

TRENDS AND RESEARCH ISSUE

* Trends:

- ① One of the most difficult and rapidly expanding study fields is Content extraction, indexing, and retrieval of multimedia data.
- ② Specifically, we require Scalable browsing algorithms that enable access to very large multimedia databases, reliable methods for indexing / retrieving and compressing multimedia data, and semantic visual interfaces that integrate the aforementioned elements into Unified multimedia browsing and retrieval systems.
- ③ The performance of Content-based retrieval is improved for retrieval of multimedia by combining several integrated media types.

* Research Issues:

Research Issues:

→ Analysis of multimedia input.

→ output generation of multimedia

→ multimedia Collaboration

→ Interfaces for agents

A) Analysis of multimedia Input -

- ① In domains like inter-media segmentation, partial input parsing and interpretation, and partial multimedia reference resolution, there are still many research issues to be solved.

- ② It is necessary to create and test new interactive tools (such as force, and facial expression detectors) to open up fresh possibilities, like the detection and tracking of human emotional states.

B) Output Generation for multimedia -

- ① There are still significant unanswered problems surrounding the best ways to choose effective media, content and modality allocation (such as language, non-speech audio, or gesture to guide attention). (e.g. realizing language as visual, text or aural speech.)

C) multimedia Collaboration -

- ① As our world grows more networked, finding more efficient ways of human-human computer-mediated contact becomes more and more crucial.

D) Interfaces for Agents -

- ① In educational settings, video games, and customer service software, agents are present.
- ② Important user interface requirements include the development of lifelike behaviors and the provision of speaking and gesturing agent displays.

Multimedia Data Support in Commercial DBMS -

- ① A multimedia database is a collection of related multimedia data that may include text, graphics (sketches, drawings), photos, animations, video, audio, and other elements.
- ② It may also contain large amounts of multi-source multimedia data. The term "multimedia database management system" refers to the framework that controls various multimedia data types that can be provided, stored, and used in various ways.

Content of multimedia Database Management System -

- ① media data - the real data used to portray an object
- ② media format data - details concerning the format of media data after it has been acquired, processed, and encoded, such as sampling rate, resolution, encoding technique, etc.
- ③ media keyword data - keywords that describe how data is generated. It also goes by the name of "content descriptive data". Example: the recording's date, time, and location.
- ④ media feature data - information that is depending on the content, such as the distribution of colors, types of texture, and shapes.
- ⑤ There are still many challenges to multimedia databases, some of which are:-

Challenges of multimedia databases

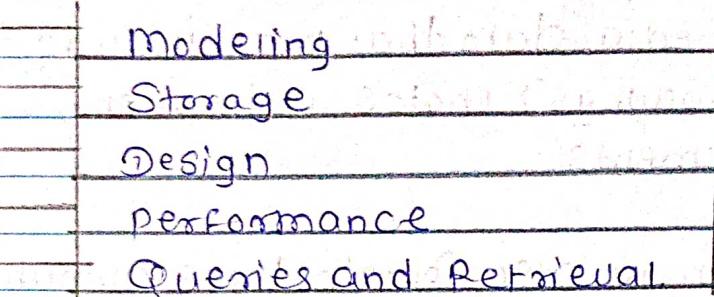


Fig - Challenges OF Multimedia Databases

- ① Modeling— By enhancing database Retrieval methods rather than information retrieval techniques, this Field of Study is Specialized and warrants Special attention.
- ② Design— Because multimedia databases use a range of Formats, Including JPEG, GIF, PNG, and MPEG, which are difficult to Convert between one another the Conceptual, logical And physical design.
- ③ Storage— The Representation, Compression, mapping to device hierarchies, archiving, and buffering during input-output operations Create challenges when Storing multimedia databases on any Conventional disc.
- ④ performance— In a situation where Video playback or audio-video Synchronization is involved, physical Constraints take precedence.

⑤ Queries and Retrieval -

For multimedia material, such as photos, video and audio, through the use of queries, data can be accessed, but this raises a number of challenges that need to be addressed, including effective query formulation, query execution and query optimization.

JMP MULTOS - (The multimedia filing System)

- ① MULTOS (MULTimedia OFFICE Server) is a multimedia document server with advanced document retrieval capabilities, developed in the context of an ESPRIT project in the area of office systems. MULTOS is based on a Client/Server architecture.
- ② The objective of MULTOS was to provide a practical and affordable system for storing and retrieving multimedia content in an office setting.
- ③ A conceptual model providing a semantic-oriented description of documents forms the basis of much of the processing within the MULTOS system.
- ④ The content-oriented processing, categorization, and retrieval algorithms used by MULTOS are all built on top of this semantic document model.
- ⑤ Document modeling, document type handling, query language design and query processing optimization, image data modeling and retrieval, knowledge-based document classification, document management,

The main goals were to,

- ① Create a System For the Filing and Retrieval OF multimedia documents that is effective and affordable Integrate Optical Storage media for the filing OF numerous documents.
- ② Implement methods For Retrieving text and attribute data based on Content and look into methods for Retrieving picture data based on Content.
- ③ Implement document management Features like access Control , security , integrity , and Version Support.
- ④ There are two different kinds of document servers that Correspond to two separate sets of documents that have various retrieval requirements.
- ⑤ while the document on the Archive Server are reliable and read far less frequently than those on the dynamic Server , the dynamic Server allows for document updates.
- ⑥ magnetic storage is used to store documents unproduced or Collected in the Client environment
- ⑦ The Classification enables portions of the document to be linked to Conceptual elements for use in later document retrieval.