# Querying structured and unstructured data - combined approach

**Asen Stefanov, Atanas Palazov**

*Institute of Oceanology, Bulgarian Academy of Sciences, Varna, Bulgaria*

**Objective:** The term "unstructured data" mean different things in different contexts. In the context of relational database systems, it refers to data that can't be stored in rows and columns. Here, unstructured data are understood to include e-mail files, word-processing text documents, PowerPoint presentations, JPEG and GIF image files, and MPEG video files. In the past there were two options to store unstructured data file in the database: either to store the file in the database as a BLOB column, that was good for management but too bad for performance, or to keep the file in the file system and store a link in a database- good for performance, bad for managing the files.

Now a completely new architecture inside the present DBMS for handling file or unstructured data are available. The current DBMS offers many possibilities to store file content such as images, audio, video, PDFs, spreadsheets etc. side by side with other structured information.
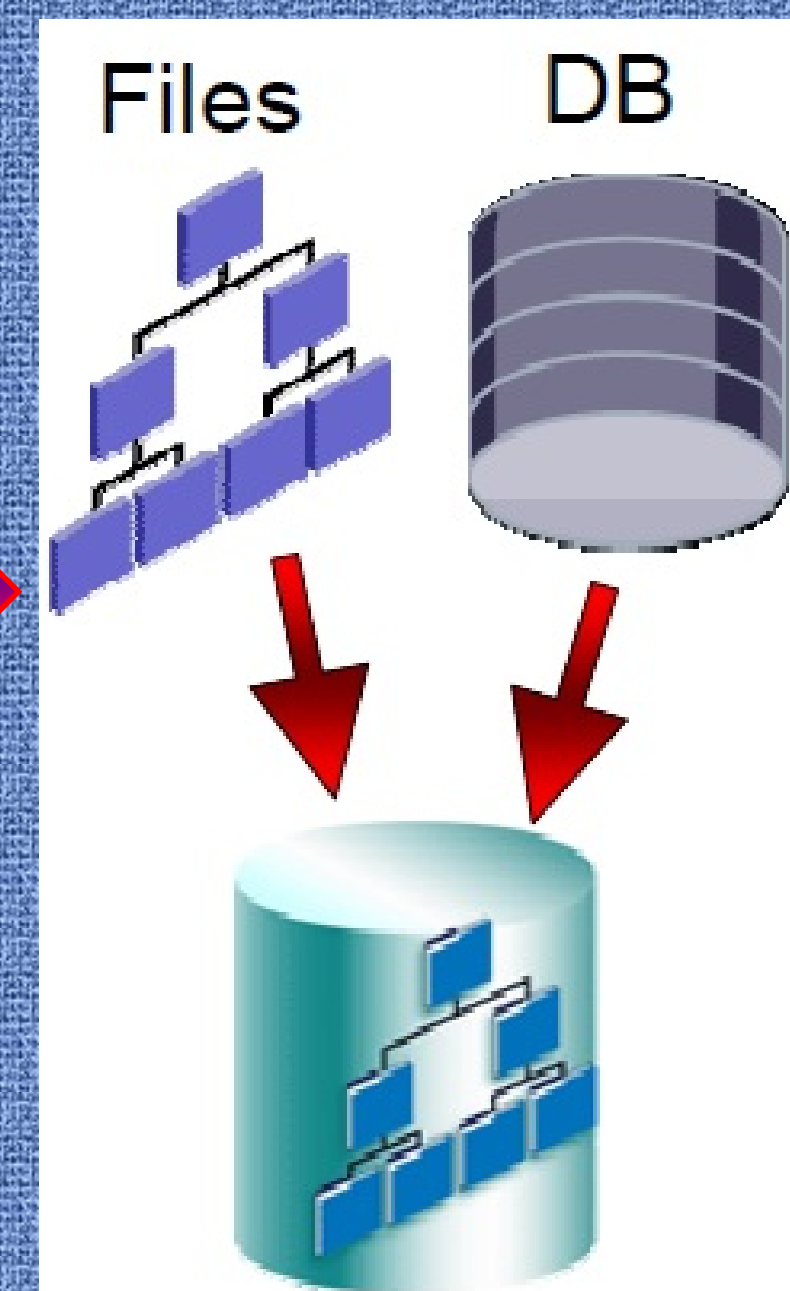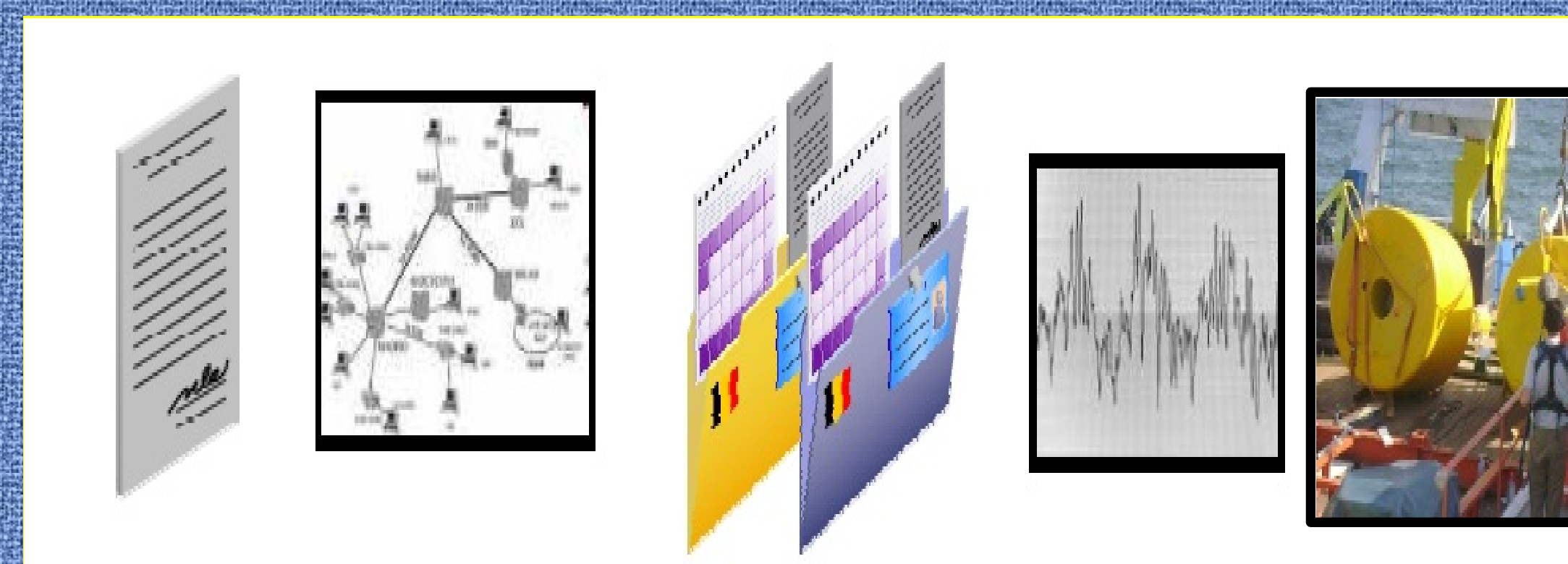


*Fig.1 Unstructured data - new approach*

**Implementation:** Today in the IO-BAS a large amount of oceanographic information is believed to be unstructured. As a result, the summary reports often lack key project data that might be relevant to decision making. "For a complete picture of the project, you have to look at both structured and unstructured data,"
IO-BAS has hands-on some experience in the still relatively new domain of unstructured data processing, and has helped the scientists to combine structured and unstructured solutions as follow:

➢ Initially the sources and file types of unstructured information are identified – ship's model characteristics, project documentation, contracts, cost sheets and etc.
➢ The various files and documents are bulk loaded, updated as well as managed in SQL in DBMS, and access is provided for standard applications as if they were stored in the file system. There are mostly text heavy, but can contain also scanned documents and images, presentations and etc.
➢ To access the unstructured part of data base the common network share for local users was implemented and an FTP portal was installed for Internet users.
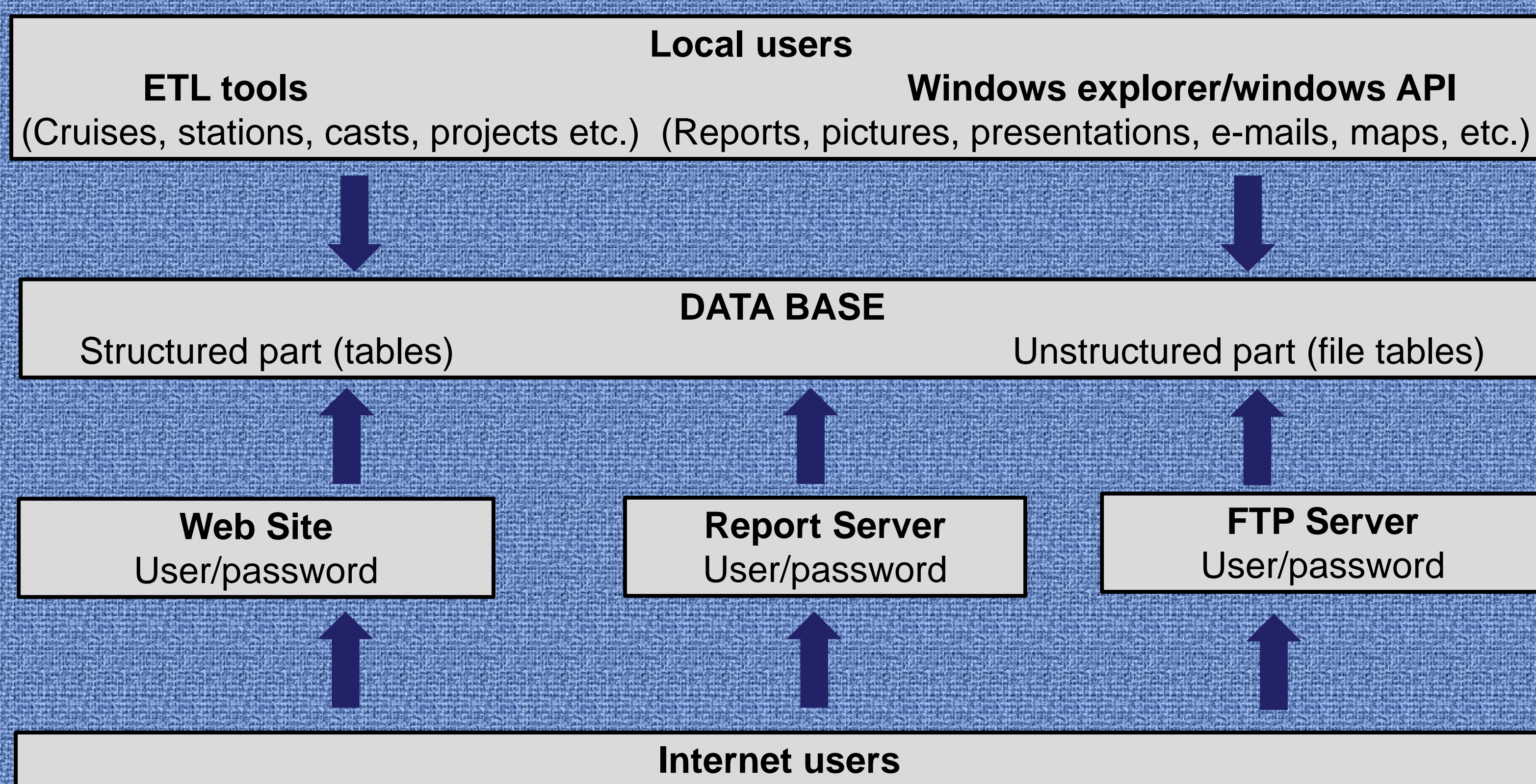


*Fig.2 Data flow search - query- open*

➢ To implement a combined approach in queries a set of drill-down reports from the structured parts of the database are developed. These reports are used as search filters over unstructured part of the database.

➢ A quantum improvement in search efficiency is gained by adding sufficient *meta information to the files.* This approach, often called "property selection", allows users to locate and retrieve information by taking advantage of available meta data by filtering and sorting on known meta data fields such as author, keywords, version etc.

➢ each level of the hierarchy has its own security



➢ Unstructured data are organized in a hierarchical directory structure, referred to as a taxonomy. A taxonomy operates like a directory on a PC by providing a convenient, intuitive way to navigate and access information

*Fig.3 Case cruise – project – documents search*

**Technical prerequisites:**
➢ MS SQL 2012 Express (free edition) with file stream option and reporting services enabled
➢ Windows OS
➢ IIS 7

**Conclusion:**
It is necessary to read, integrate, and precondition the unstructured data before it can be used for the purposes of analytics in the structured environment. The benefits of combining the structured and unstructured data are:
➢ All information can be part of a database transaction, freeing the application from the complexity of guaranteeing atomicity, read consistency and other backup and recovery requirements
➢ The unstructured data can be analyzed
➢ The unstructured text can be accessed by direct or indirect searches
➢ The unstructured data can be linked to structured data and composite queries can be created
➢ At a certain point the saved in data base unstructured information can be transferred, for instance, into an established oceanographic relationship management system