



# An Advanced Enterprise Information Search and Delivery System

BY YURDAER N. DOGANATA, YOUSSEF DRISSI, TONG-HAING FIN, GREG BROWN, MOON J. KIM, AND LEV KOZAKOV

## Fulfilling IBM's one-Web vision



### ABOUT THE AUTHOR

Dr. Yurdaer Doganata is the manager of the Integrated Networked Services group and the dBlue core research team at the IBM T.J. Watson Research Center. His current focus is on alternative techniques and methods for effective information search and retrieval in unstructured and semistructured environments. Yurdaer holds several patents and research awards and has published numerous papers. He is the technical committee cochair of ISCC '03.

**E-MAIL**  
yurdaer@us.ibm.com

One of the biggest complaints we hear about many company Web sites, from customers and employees alike, is that it's too hard to find what you need. At IBM, with 2.5 million Internet pages and more technical content than any single entity, including the Pentagon, that's no surprise.

A new IBM advanced information search and delivery system for the IBM support site ([www.ibm.com/support](http://www.ibm.com/support)) is expected to solve this problem. Code-named Digital Blue (dBlue), this project is a digital interface to IBM customers. The result of two years of work and five patentable inventions, dBlue is now available to IBM customers.

The team that created dBlue is calling it "the next generation of enterprise information search-and-delivery systems." This is a WebSphere-based technology with breakthroughs in storing, searching, and retrieving information. Customers will be able to search for IBM technical support information using natural language and will receive results that are categorized, prioritized, and personalized. dBlue provides the foundation for a set of user-oriented support services applicable to all IBM support sites worldwide.

Rich Vazzana, vice president of [ibm.com](http://ibm.com) Support and Enablement, took on this project to improve the effectiveness and performance of IBM's Web-enabled post-sales support services. It became the underlying architecture of the "one-Web" vision across multiple IBM Web sites, improving adherence to IBM's company-wide standards and setting the stage for more advanced service offerings. The program will provide customers with IBM support experience, a single IBM support/service portal, toolset, and infrastructure. Hence, cross-IBM "common" support functions will be realized.

"The business goal is to improve goal achievement on

the IBM Internet," said Frank Cummiskey, director of IBM eSupport & Services. "The primary reason that customers visit IBM's support sites is to resolve a technical problem. Today, only about 60% actually achieve their goal. Improving our customers' ability to find what they are looking for, as well as to find value in the information they find, will increase self-service on the Web, saving millions of dollars and increasing customer satisfaction."

### System Architecture

Although dBlue architecture does not depend on the WebSphere software platform, it's the platform of choice of the dBlue architects for its scalability, flexibility, reliability, and high performance required for dynamic Web applications hit by millions of customers every month. In addition to the application server mechanisms, the WebSphere software platform provides reliable communication middle-

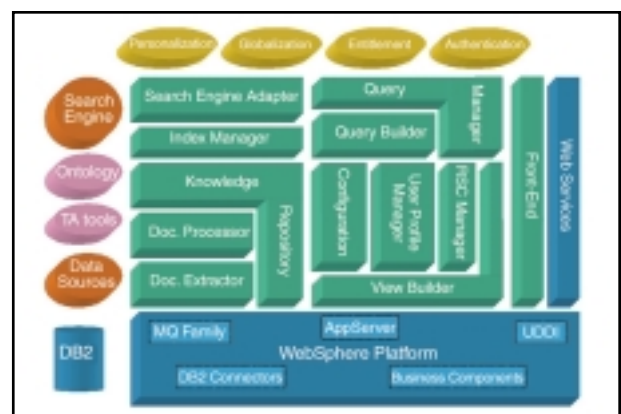


FIG. 1: DBLUE ARCHITECTURE OVERVIEW

ware – the WebSphere MQ family. It also supports DB2 Universal Databases, provides a foundation for Web services, and integrates business components for text analysis and machine translation. The WebSphere Everyplace Suite provides an integrated software platform for extending the reach of business applications, enterprise data, and Internet content into the realm of pervasive computing. All this makes the WebSphere software platform the perfect fundament for the dBlue system. Figure 1 is an overview of the dBlue architecture.

The dBlue architecture connects three important elements from the information search world – information sources, search engines, and end users – on the basis of the WebSphere software platform. This is done through a set of components called “The Knowledge Builder.” Information sources are data sources such as document repositories, DB2 and Lotus Notes databases, Web sites, and so on. Search engines are programs that can index content and enable searching of the indexed data. End users access dBlue through a front-end interface; the current default interface is a Web interface. The content is extracted from information sources using the Document Extractor and mapped to a unified XML Schema; then it’s processed by the Document Processor and stored in the Knowledge Repository.

When a user accesses the system and submits a search query, the Query Manager, along with all the submitted parameters, processes this query. The Query Builder then collects the query and parameters submitted by the user, along with information coming from the user’s profile and the system configuration, to build a standard Query object. The Query object is submitted to the search engine through the Search Engine Adapter. The search results flow back to the user through the Search Engine Adapter, the Search Query Manager, and the View Builder. The View Builder uses the Remote Site Customization component and data to construct a personalized view of the search hit list. When the user requests a view of a specific document, this request is processed by the View Builder, which accesses the Knowledge Repository to get the document content and builds a coherent document view.

Enabled by the WebSphere software platform, dBlue introduces various innovative solutions in the areas of information search and delivery. In dBlue:

- Content is indexed using the concept of virtual URLs.
- Search results and documents are rendered by employing dynamic layout features.
- Keyword and navigational search are combined for effective searching.

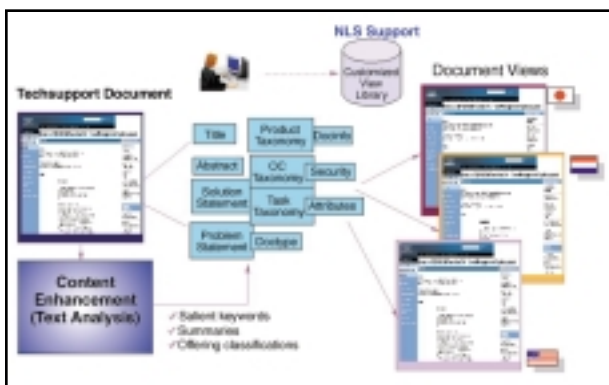


FIG. 2: CREATING MULTIPLE VIEWS FROM THE SAME CONTENT

- Search results and indexing are improved by using text analysis technologies.
- Architecture is enabled for globalization and dual language search.

### Virtual URLs and Dynamic Layout

dBlue is a search system, but it doesn’t depend on a particular search engine. The technical content to be indexed can be pushed to any search engine using the concept of virtual URLs. Until now, search systems have had to crawl content off a particular address where it’s stored. Hence, the documents are replicated redundantly for the purpose of indexing the same information in a different context. With virtual URLs, documents to be indexed are built on the fly from building blocks, eliminating the need for replication and crawling. In other words, the virtual URLs aren’t associated with any physically stored documents. This motivates another breakthrough in content storage. In the back end, the documents are broken down into components, such as title, problem, solution, reference, and category, allowing for true knowledge mining and the building of multiple views of the same content. Extracting the documents from their original sources and creating components based on unified XML Schema for technical documents accomplishes this, giving users a great deal of flexibility and allowing them to receive a wider range of information.

In a typical search system, the documents are stored and retrieved with a layout defined by the content providers. In this case the layout is static and cannot be changed to meet customers’ needs. dBlue solves this problem by introducing the concept of dynamic layout for creating multiple views from the same content (see Figure 2).

The component-based storage system invented by the dBlue team decomposes documents into data elements without breaking the ties to their original documents. When customers request information in a specific layout, components are analyzed to ensure that they have all the necessary elements for a specific document, which is then built dynamically. This gives the flexibility to separate user experience from the content-generation process and also enables rapid localization and internationalization of the pages.

### OC Taxonomy

One of the first challenges was to institute a consistent structure for content creation, since the huge amount of support content that already existed was not suitable for search. In order to structure the content and organize the content-creation process, the unified XML Schema for technical documents was created. This schema incorporates content components, such as title, abstract, problem statement, and solution statement, along with multiple attributes, keywords, references, and attachments.

The second step in organizing the content was creation of the content repository schema that allows storage of both unstructured and structured data. This schema contains more than 30 DB2 tables that provide storage for the document content, along with all associated information, and supports a variety of queries. Then, of course, both existing and new content had to be migrated to this structure. The content migration pipe is powered by the WebSphere MQ family of communication middleware. The documents extracted from their original repositories were converted to XML format based on the unified XML Schema and transferred to the new storage. All document attachments were



#### ABOUT THE AUTHOR

Youssef Drissi is an advisory software engineer at the IBM T.J. Watson Research Center in Hawthorne, New York. He holds several patents and is the author of several publications. Youssef is a member of the dBlue project architecture and research teams. His current work involves research, architecture, and development of next-generation unstructured information and knowledge management systems.

#### E-MAIL

youssefd@us.ibm.com



#### ABOUT THE AUTHOR

Dr. Tong-Haing Fin is a senior software engineer at the IBM T.J. Watson Research Center. He holds a number of patents and is the author of several publications. He is a member of the dBlue project architecture and research teams. His current work involves research and system integration of text analysis and knowledge management systems.

#### E-MAIL

thfin@us.ibm.com



### ABOUT THE AUTHOR

Greg Brown is the dBlue solutions architect and team lead for OneWeb Infrastructure and Integration for the IBM.com e-Support & Service Delivery team. Greg holds several patents related to dBlue and multiple IBM.com awards.

### E-MAIL

browng@us.ibm.com



### ABOUT THE AUTHOR

Dr. Moon J. Kim is an IBM senior technical staff member with responsibility for the development of the e-Support advanced Web system. Moon also developed many large system solutions, such as S/390 and MPP and was involved in the development of the network systems that later called the broadband high-speed access system, including HFC and FSN. Moon is an IBM Master Inventor who holds 10 patents, has published 10 invention technical papers, and has filed 12.

### E-MAIL

moonkim@us.ibm.com

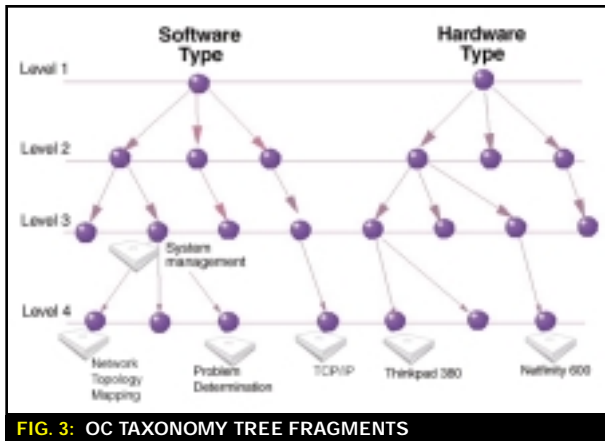


FIG. 3: OC TAXONOMY TREE FRAGMENTS

encoded using "Base64" encoding and incorporated in XML objects. To eliminate unnecessary XML parsing, the transportation was done in a binary format.

Another challenge was determining how to store and dynamically retrieve this information in a scalable and flexible way. The team adopted a categorization scheme based on IBM product offerings, called offering classification (OC). The common library classification can be used, but for IBM technical support all contents are associated with IBM products. With the OC taxonomy attached to the content, the content can easily be shown where it belongs. Figure 3 shows a fragment of the OC taxonomy tree with sample documents that may be found under certain leaves.

Having OC taxonomy information attached to the documents made it possible to combine a keyword with the navigational search. This way, users can narrow down search results with single click.

### Combining Keyword with Navigational Search

The way the system is architected allows combining keyword search with navigational search. Based on a topic or a document type, users can narrow down search findings with a single click. This increases the chances of finding the requested information when the user query isn't specific enough to narrow down the search results on the first attempt. The categorized results are returned with the option of filtering the results based on IBM's product offerings and the document types.

Although combining keyword and navigational search refines the search results, it doesn't improve relevancy or precision/recall rates. The following sections discuss some text-analysis techniques used to improve precision/recall.

### Content Enhancement for Search Improvement

The quality of full-text search depends mainly on query terms and on how documents are indexed by the search engine. The search results contain the documents that are indexed against the query terms and scored based on certain statistical criteria. In many real-life situations, the relevant documents can't be found or may not appear at the top of the search results because they are scored low or they don't contain the terms exactly as in the query. This is common when users choose variations of the query terms, including inflections, misspellings, abbreviations, and so

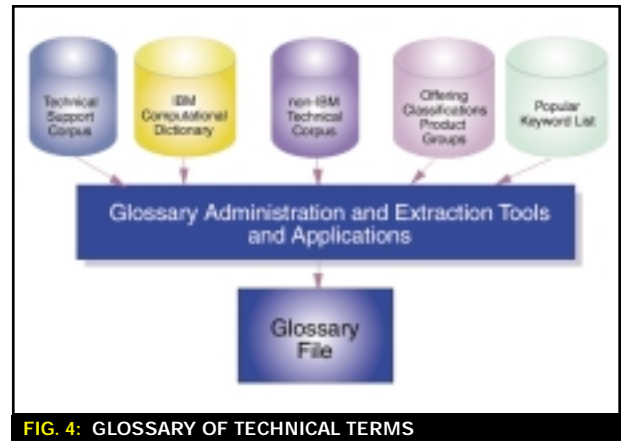


FIG. 4: GLOSSARY OF TECHNICAL TERMS

on. To improve the user experience, dBlue uses text analysis tools developed by IBM Research to enhance the contents of documents. This process is started by extracting terms from a large collection of documents in the IBM technical support domain to create a domain-specific glossary. The terms in the glossary can consist of canonical form, variant form (inflection, abbreviation, misspelling, etc.), synonym, term definition, statistical data, and other information. This initial glossary is enhanced by eliminating irrelevant terms and reranking terms using special dictionaries and algorithms. The process of generating and enhancing the glossary is semi-automatic, using glossary tools and the librarian. Figure 4 shows the multiple components that compose the glossary of technical terms built for the dBlue system.

Based on the glossary, the important keywords in each document are extracted and ranked, and their related glossary terms (variants, synonyms, etc.) are used to enrich the content of the document. The content enrichment is used to create keyword metatags for biased indexing, expand the query terms to include related terms, and enable search for related documents. To improve the user's search experience, keywords are displayed in the search results and navigating through keywords is possible.

### Globalization

As part of the effort to allow different languages to be supported from a single Web application consistent with the vision of "one Web" for all regions, dBlue has a globalization process that consists of two main processes: internationalization and localization.

#### INTERNATIONALIZATION

Internationalization (sometimes abbreviated as i18n) is the process of designing an application so that it can be adapted to various languages and regions without engineering changes. After the internalization of dBlue software components, they can run worldwide with the addition of localized data. Hence, support for new languages doesn't require recompilation. Textual elements, such as status messages and GUI component labels, are stored outside of the source code and retrieved dynamically. Culturally dependent data, such as dates and currencies, appears in formats that conform to the end user's region and language.

The Unicode format, which handles most characters known to mankind, was instrumental in allowing the use of a unique globalized repository where multilingual searchable text and documents are encoded in one unique

format. Unicode was also adopted as a standard format for encoding internal textual data in dBlue.

## LOCALIZATION

Localization (sometimes abbreviated as i18n) is the process of adapting software for a specific region or language by adding locale-specific components and translating text. Usually the most time-consuming part of the localization process is the translation of text. Other types of data, such as sounds and images, may require localization if they are culturally sensitive. Localizers also verify that the formatting of dates, numbers, and currencies conforms to local requirements.

Two innovative approaches in the globalization process are worth mentioning. The first allows documents to be searched, regardless of their language, against a query formulated in user-specific language. This is accomplished in dBlue without extra overhead or the need for a translation at runtime through a specific extension of the inverted index, a core component of most search engines. The second allows the achievement of similar results through dynamic mapping of the user's search query at run time, and use of multi-threading to submit multilingual queries to the search engine. Figure 5 illustrates some aspects of this innovation.

## Remote Site Customization

Another dBlue feature that addresses corporate needs is Remote Site Customization (RSC). IBM, like any other large corporation, has multiple departments that may want to present search results and technical documents to their customers in different ways, adding their own ads, promotions, and so on. The dBlue system enables this by providing the RSC feature, which allows different departments to define their own layouts for search results and technical documents. The idea of RSC is rather simple: each remote site that wants to present the shared system content in a special format is allowed to store and register its own forms. When the system gets a request that specifies this remote site, it will use the appropriate form to build the customized view of the content. Figure 6 shows the six areas that are available for customization in a results page. To assist departments in customizing the layout of Web pages, dBlue provides a Web-based RSC administrative application, which allows the uploading and testing of customized forms.

## Conclusion

dBlue has many advantages. In the near future, customers will be able to ask questions in natural language and the system won't require an exact match of words. In the near future, dBlue will also personalize searching so that once a user fills out a profile, responses will be filtered and ranked based on that profile. Multilanguage searches for documents written in Japanese, Chinese, and French will be supported by late 2002. By 2Q03, it's expected that all languages will be supported from a single Web application consistent with the vision of "one Web" for all regions.

## References

- *IBM WebSphere Software Platform Overview*: [www7b.boulder.ibm.com/wsdd/products/platform/overview.html](http://www7b.boulder.ibm.com/wsdd/products/platform/overview.html)
- *WebSphere Everyplace Suite*: [www.ibm.com/pvc/prod](http://www.ibm.com/pvc/prod)

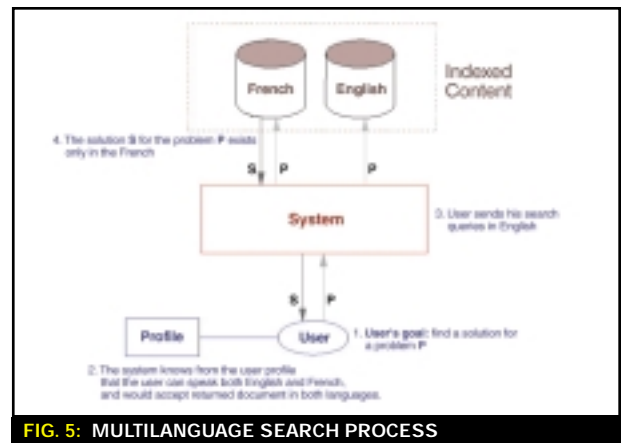


FIG. 5: MULTILANGUAGE SEARCH PROCESS

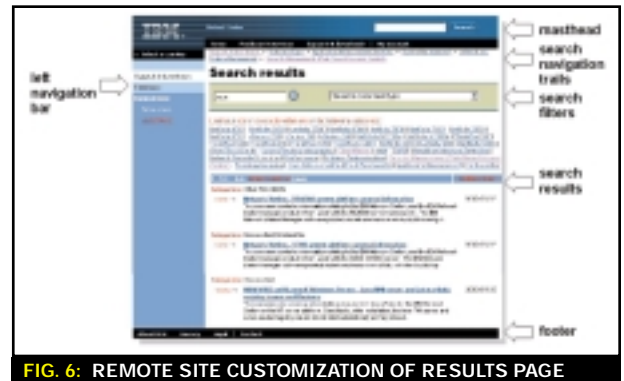


FIG. 6: REMOTE SITE CUSTOMIZATION OF RESULTS PAGE



## ABOUT THE AUTHOR

Dr. Lev Kozakov is a research staff member at IBM T.J. Watson Research Center and is a member of the dBlue project architecture and research teams. He has worked in many areas, including dynamic systems, applied statistics, information management systems, man-machine interface, medical software, computer telephony, and design patterns. Lev holds a number of patents and is the author of several publications.

## E-MAIL

kozakov@us.ibm.com

- [ucts/wes/index.shtml](http://www7b.boulder.ibm.com/wsdd/index.shtml)
- Booth, Alan E. *Extending the Reach of Enterprise Applications with Transcoding and Machine Translation*: [www7b.software.ibm.com/wsdd/library/techarticles/0206\\_booth/booth.html](http://www7b.software.ibm.com/wsdd/library/techarticles/0206_booth/booth.html)
- *WebSphere Technology for Developers (overview & downloads)*: [www7b.software.ibm.com/wsdd/downloads/wstechnology\\_tech\\_preview.html](http://www7b.software.ibm.com/wsdd/downloads/wstechnology_tech_preview.html)
- Snell James. *Implementing Web services with IBM WebSphere Version 4.0*: [www7b.software.ibm.com/wsdd/library/techarticles/0108\\_snell/0108\\_snell.html](http://www7b.software.ibm.com/wsdd/library/techarticles/0108_snell/0108_snell.html)
- *DB2 Product Family Overview*: [www.ibm.com/software/data/db2/](http://www.ibm.com/software/data/db2/)
- *About the Internationalization Activity*: [www.w3.org/International/about.html](http://www.w3.org/International/about.html)
- *Internationalization in Java*: <http://java.sun.com/docs/books/tutorial/i18n/>
- *The Unicode Home Page*: [www.unicode.org](http://www.unicode.org)
- Park, Youngja; Byrd, Roy J.; and Boguraev, Branimir K. *Automatic glossary extraction: Beyond terminology identification*, IBM Research Technical Report RC22421, 2002. IBM T.J. Watson Research. *The Talent (Text Analysis and Language Engineering) project*: [www.research.ibm.com/talent](http://www.research.ibm.com/talent)
- Boguraev, Branimir K. and Neff, Mary S. (2000). *Lexical Cohesion, Discourse Segmentation and Document Summarization*. RIAO-2000. April.
- Park, Youngja and Byrd, Roy J. (2001). *Hybrid text mining for finding terms and their abbreviations*. EMNLP-2001.
- Chu-Carrol, Jennifer; Prager, John; Rabin, Yael; and Cesar, Christian. (2002). *A Hybrid Approach to Natural Language Web Services*. EMNLP-2002. 