# Integrated Feasibility Experiment for Bio-Security: IFE-Bio
# A TIDES Demonstration

Lynette Hirschman, Kris Concepcion, Laurie Damianos, David Day, John Delmore, Lisa Ferro, John Griffith, John Henderson, Jeff Kurtz, Inderjeet Mani, Scott Mardis, Tom McEntee, Keith Miller, Beverly Nunan, Jay Ponte, Florence Reeder, Ben Wellner, George Wilson, Alex Yeh

The MITRE Corporation
Bedford, Massachusetts, USA and
McLean, Virginia, USA
781-271-7789

lynette@mitre.org

## ABSTRACT

As part of MITRE's work under the DARPA TIDES (Translingual Information Detection, Extraction and Summarization) program, we are preparing a series of demonstrations to showcase the TIDES Integrated Feasibility Experiment on Bio-Security (IFE-Bio). The current demonstration illustrates some of the resources that can be made available to analysts tasked with monitoring infectious disease outbreaks and other biological threats.

## Keywords

Translation, information extraction, summarization, topic detection and tracking, system integration.

## 1. INTRODUCTION

The long-term goal of TIDES is to provide delivery of information <u>on demand</u> in real-time from live on-line sources. For IFE-Bio, the resources made available to the analyst include e-mail, news groups, digital library resources, and eventually (in later versions), topic-specific segments from broadcast news. Because of the emphasis on global monitoring, there is a need to process incoming information in multiple languages. The system must deliver the appropriate information <u>content</u> in the appropriate <u>form</u> and in the appropriate <u>language</u> (taken for now to be English). This means that the IFE-Bio system will have to deliver news stories, clusters of relevant documents, threaded discussions, alerts on new events, tables, summaries (particularly over document collections), answers to questions, graphs and geo-spatial temporal displays of information.

The demonstration system for the Human Language Technology Conference in March 2001 represents an early stage of the full IFE-Bio system, with an emphasis on end-to-end processing. Future demonstrations will make use of MITRE's Catalyst architecture, providing an efficient, scalable architecture to facilitate integration of multiple stages of linguistic processing. By June 2001, the IFE-Bio system will provide richer linguistic processing through the integration of modules contributed by other TIDES participants. By June 2002, the IFE-Bio system will include additional functionality, such as real-time broadcast news feeds, new machine translation components, support for question-answering, cross-language information retrieval, multi-document summarization, automatic extraction and normalization of temporal and spatial information, and automated geospatial and temporal displays.

## 2. The IFE-Bio System

The current demonstration (March 2001) highlights the basic functionality required by an analyst, including:

- **Capture** of sources, including e-mail, digital library material, news groups, and web-based resources;

- **Categorizing** of the sources into multiple orthogonal hierarchies useful to the analyst, e.g., disease, region, news source, language;

- **Processing** of the information through various stages, including "zoning" of the text to select the relevant portions for processing; named entity detection, event detection, extraction of temporal information, summarization, and translation from Spanish, Portuguese, and Chinese into English;

- **Access** to the information through use of any mail and news group reader, which allows the analyst to organize, save, and share the information in a familiar, readily accessible environment;

- **Display** of the information in alternate forms, including color-tagged documents, tables, summaries, graphs, and geospatial, map-based displays.

Figure 1 below shows the overall functionality envisioned for the IFE-Bio system, including capture, categorizing, processing, access and display.

Collection capability for the current IFE-Bio system includes email, news groups, journals, and Web resources. We have a complete copy of the ProMED mailings (a moderated source
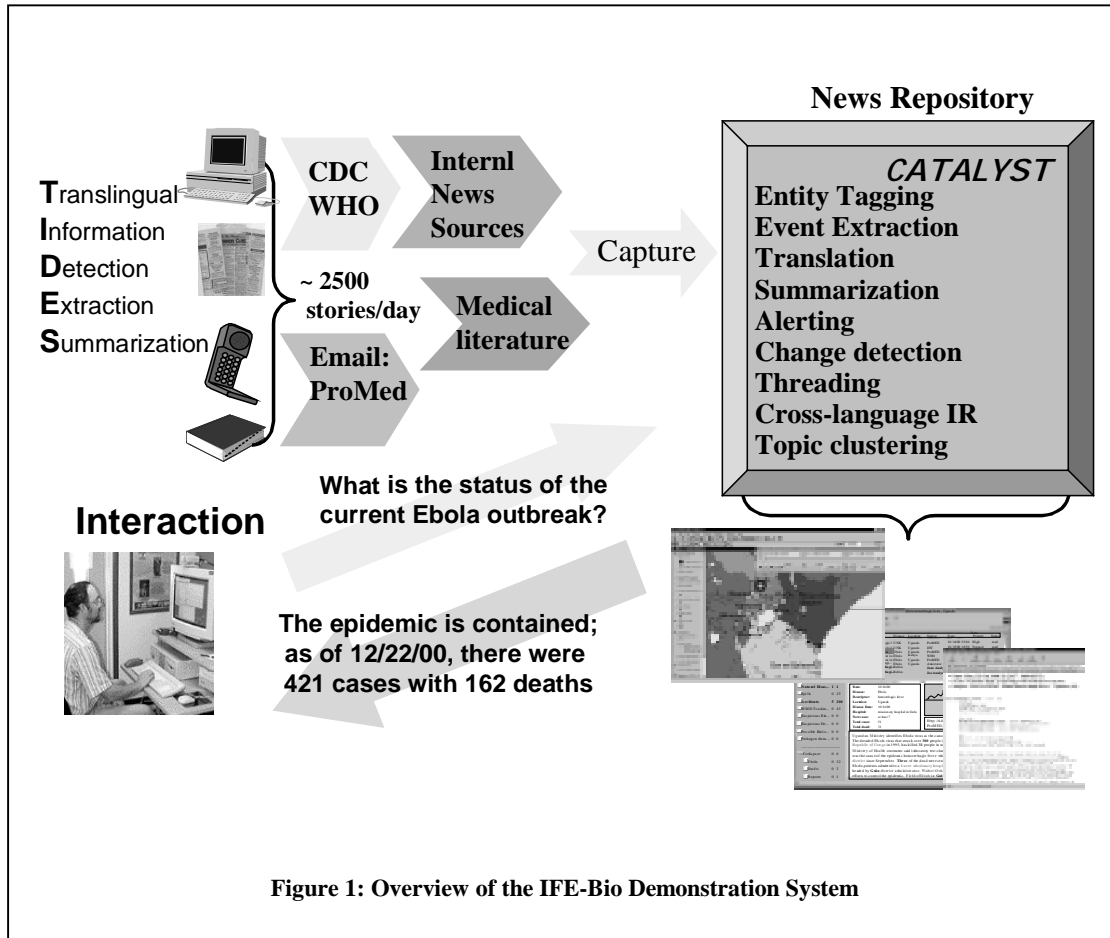
**News Repository**

*CATALYST*
**Entity Tagging**
**Event Extraction**
**Translation**
**Summarization**
**Alerting**
**Change detection**
**Threading**
**Cross-language IR**
**Topic clustering**

**T**ranslingual
**I**nformation
**D**etection
**E**xtraction
**S**ummarization

CDC
WHO

Internl
News
Sources

~ 2500
stories/day

Medical
literature

Email:
ProMed

Capture

**Interaction**

**What is the status of the current Ebola outbreak?**

**The epidemic is contained; as of 12/22/00, there were 421 cases with 162 deaths**

**Figure 1: Overview of the IFE-Bio Demonstration System**

tracking global infectious disease outbreaks), and are routinely collecting other information sources from the World Health Organization and CDC. In addition, we are collecting several general global news feeds. Current volume is around 2000 messages per day; we estimate capacity for the current system at around 4500 messages/day. Once we have integrated a filtering capability, we expect the volume of messages saved in IFE-Bio should drop significantly, since many of the global news services report on a wide range of events and not all need to be passed on to IFE-Bio analysts. The categorizing of sources is done based on the message header. The header is synthesized by extracting key information about disease name, the country, and other relevant information such as type of victim and source of information, as well as date of message receipt.

The processing for the current demonstration system uses a limited subset of the Catalyst architecture capabilities and a number of in-house linguistic modules. The linguistic modules in the current demonstration system include tokenization, sentence segmentation, part-of-speech tagging, named entity detection, temporal extraction (Mani and Wilson 2000) and source-specific event detection. In addition, we have incorporated the CyberTrans embedded machine translation system which "wraps" available machine translation engines to make them available via an e-mail or Web interface (Reeder 2000). Single document summarization is performed by the MITRE WebSumm system (Mani and Bloedorn 1999).

We carefully chose a light-weight interface mechanism for delivery of the information to the analyst. By treating the incoming streams of data as feeds to a news server, the analyst can inspect and organize the information using a familiar news and e-mail browser. The analyst can subscribe to areas of interest, flag important messages, watch specific threads, and create tailored filters for monitoring outbreaks. The stories are crossed-posted to multiple relevant news groups, based on the information in the header, e.g., a story on Ebola in Africa would be cross posted to the Africa regional newsgroup and to the Ebola disease newsgroup. Search by subject and date allow the analyst to select subsets of the messages for further processing, annotation or sharing. The news client provides notification of incoming messages. In later versions, we plan to integrate topic detection and tracking capabilities, to provide improved filtering and routing of messages, as well as detection of new topics. The use of this simple delivery mechanism provides a familiar environment with almost no learning curve, and it avoids issues of platform and operating system dependence.

Finally, the system makes use of several different devices to display the information appropriately. Figure 2 shows the layout of the Netscape news browser interface. It includes the list of newsgroups that have been subscribed to (on the left), the list of messages from the chosen newsgroup (on top), and a particular message with color-coded named entities (including disease terms displayed in red, so that they are easy to spot in the message).
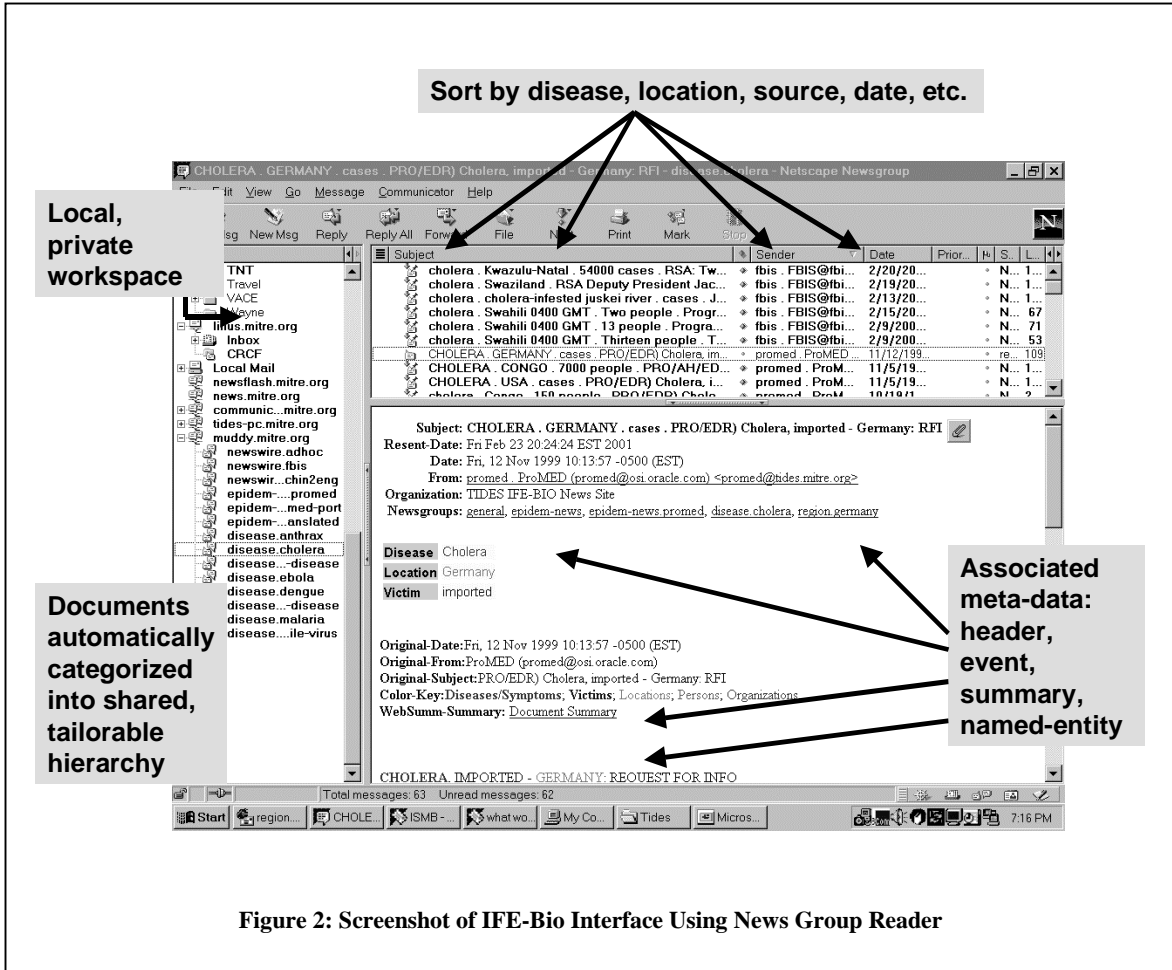
Sort by disease, location, source, date, etc.

Local, private workspace

Documents automatically categorized into shared, tailorable hierarchy

Associated meta-data: header, event, summary, named-entity

**Figure 2: Screenshot of IFE-Bio Interface Using News Group Reader**
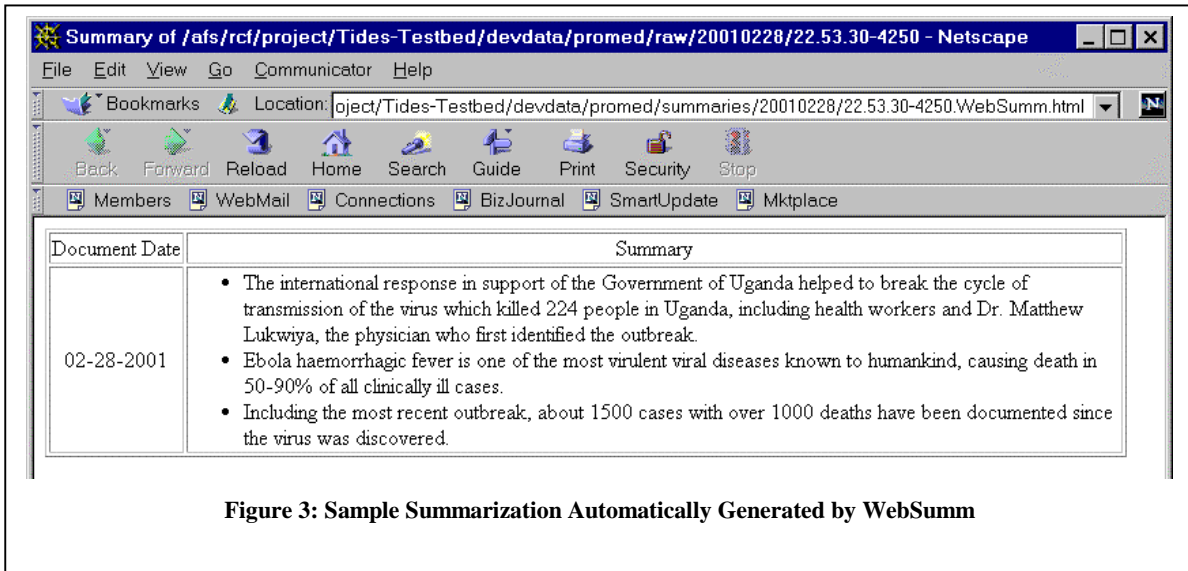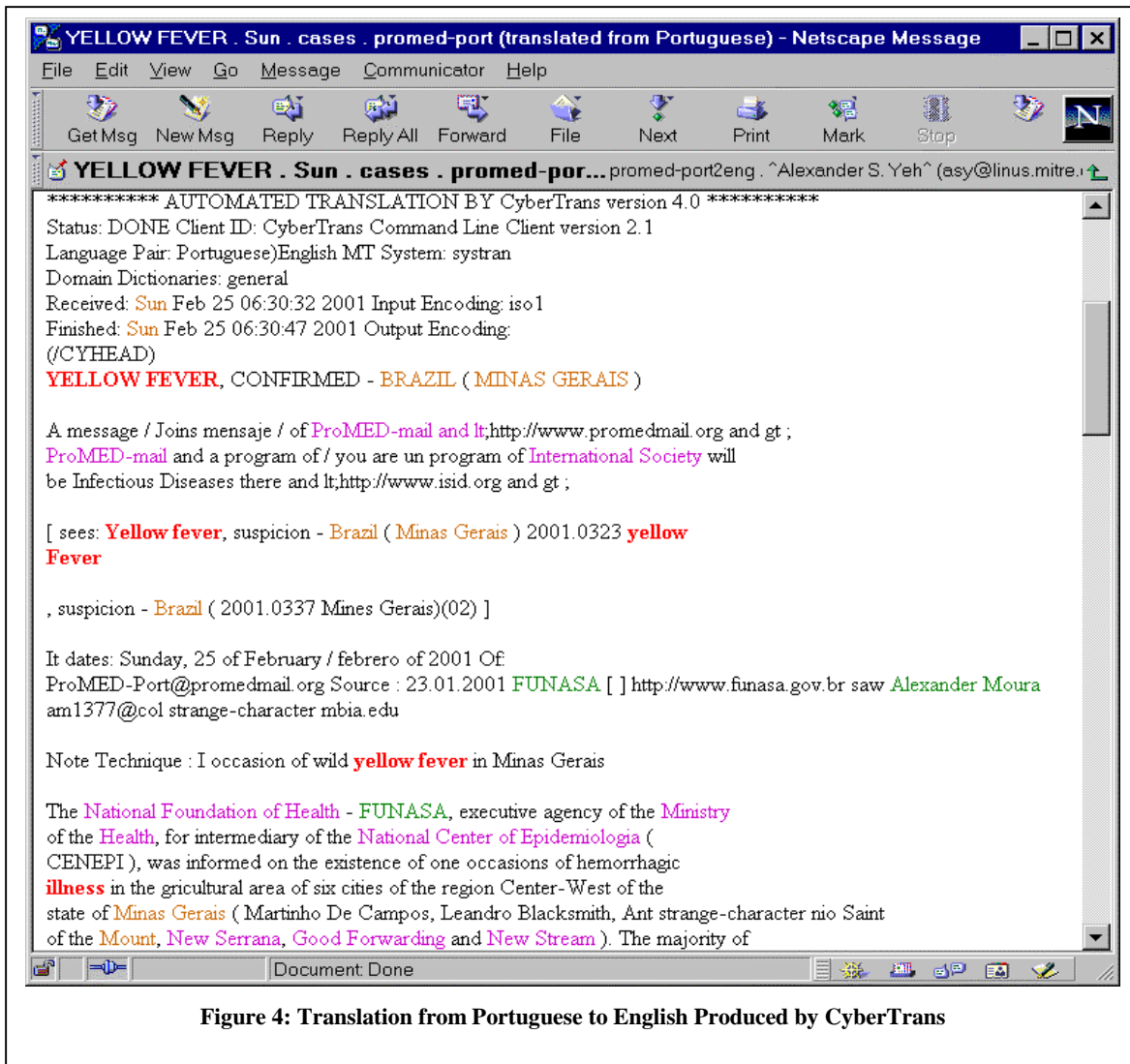
**Figure 3: Sample Summarization Automatically Generated by WebSumm**

**Figure 4: Translation from Portuguese to English Produced by CyberTrans**

There are multiple display modalities available. The message in Figure 2 contains a short tabular display in the beginning, identifying disease, region and victim type. Below that is a URL to a document summary, created by MITRE's WebSumm system (see Figure 3 for a sample summary). If an incoming message is in a language other than English, then CyberTrans is called to run code set and language identification modules, and the language is translated into English for further processing. Figure 4 below shows a sample translated message; note that there are a number of untranslated words, but it is still possible to get the gist of the message.

In addition, we are working on a mechanism to provide geographic and eventually, temporal display of outbreak information. Figure 5 shows the stages of processing involved. Stage 1 shows onamed entity and temporal tagging to identify the items of interest. These are combined into disease events by further linguistic processing; the result is shown in the table in Stage 2. This spreadsheet of events serves as input for a map-based display, shown in Stage 3. The graph plots number of new cases and number of cumulative cases over time. In the map, the size of the outer dot represents total number of cases to date, and the inner dot represents new cases. This allows the analyst to visualize spread of the disease, as well as the stage of the outbreak (spreading or subsiding).

## 3. REFERENCES

[1] Mani, I. and Bloedorn, E. (1999). "Summarizing Similarities and Among Related Documents". Information Retrieval 1(1): 35-67.

[2] Mani, I. and Wilson, G. (2000). "Robust Temporal Processing of News," Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL'2000), 69-76. New Brunswick, New Jersey. Association for Computational Linguistics.

[3] Reeder, F. (2000) "At Your Service: Embedded MT as a Service", NAACL Workshop on Embedded MT, March, 2000.

## 1. Annotate entities of interest via XML

| Disease | Source | Country | City_name | Date | Cases | New_cases | Dead |
|---------|--------|---------|-----------|------|-------|-----------|------|
| Ebola | PROMED | Uganda | Gula | 26-Oct-2000 | 182 | 17 | 64 |
| Ebola | PROMED | Uganda | Gula | 5-Nov-2000 | 280 | 14 | 89 |
| Ebola | PROMED | Uganda | Gulu | 13-Oct-2000 | 42 | 9 | 30 |
| Ebola | PROMED | Uganda | Gulu | 15-Oct-2000 | 51 | 7 | 31 |
| Ebola | PROMED | Uganda | Gulu | 16-Oct-2000 | 63 | 12 | 33 |
| Ebola | PROMED | Uganda | Gulu | 17-Oct-2000 | 73 | 2 | 35 |
| Ebola | PROMED | Uganda | Gulu | 18-Oct-2000 | 94 | 21 | 39 |
| Ebola | PROMED | Uganda | Gulu | 19-Oct-2000 | 111 | 17 | 41 |

## 2. Assemble entities into events

## 3. Display events...



○ Total Cases
● New Cases

**Figure 5: Steps in Extraction to Support Temporal and Geospatial Displays of Disease Outbreak**