

Science

The New York Times



HOME	SEARCH Go to Advanced Search	GO TO MEMBER CENTER
JOB MARKET	Past 30 Days	Welcome, loreto3
REAL ESTATE	E-Mail This Article	Printer-Friendly Format
AUTOMOBILES	<small>sponsored by</small>	STARBUCKS.COM

NEWS

- International
- National
- Nation Challenged
- Politics
- Business
- Technology
- Science
 - Earth Science
 - Life Science
 - Physical Science
 - Social Science
 - Space
 - Columns
- Health
- Sports
- New York Region
- Education
- Weather
- Obituaries
- NYT Front Page
- Corrections

OPINION

- Editorials/Op-Ed
- Readers' Opinions



FEATURES

- Arts
- Books
- Movies
- Travel
- Dining & Wine
- Home & Garden
- Fashion & Style
- New York Today
- Crossword/Games
- Cartoons
- Magazine
- Week in Review
- Photos
- College
- Learning Network

SERVICES

- Archive
- Classifieds
- Theater Tickets
- NYT Mobile
- NYT Store
- E-Cards & More
- About NYTDigital
- Jobs at NYTDigital
- Online Media Kit
- Our Advertisers

MEMBER CENTER

- Your Profile
- E-Mail Preferences
- News Tracker
- Premium Account
- Site Help

NEWSPAPER

- Home Delivery
- Customer Service
- Electronic Edition
- Media Kit

Most E-Mailed Articles

April 30, 2002

Fun With Your Zip Program: Sort Through Texts, and More

By BRUCE SCHECHTER

One of the basic truths of the digital age is that almost anything - the plays of Shakespeare, the genetic sequence of DNA, or the twitching of a seismograph needle - can be reduced to a sequence of ones and zeroes. More striking is the discovery that these sequences are largely full of hot air - redundancies that add nothing to their meaning. Clever computer programs can "zip" or compress these files, streamlining them for speedier transmission. Zipping programs have long been a boon to computer users with slow modem connections.

But now a group of Italian physicists has shown how these same programs can be used to analyze and categorize text quickly. Using little more than the zipping programs found on most personal computers, they can easily distinguish between texts written in 10 different languages and almost unfailingly tell which of a large group of texts were written by the same author.

Writing in the January issue of Physical Review Letters, the scientists - Dr. Dario Benedetto, Dr. Emanuele Caglioto and Dr. Vittorio Loreto - explain their work with an analogy to Morse code.

To keep the number of dits and dahs to a minimum, Samuel Morse considered how often each letter was used in an average English message. The letter e is the most common in English so Morse encoded it as a single dit. The next most common letter is t, so he assigned that a single dah. A relatively uncommon letter like Q takes four taps to encode: dah dah dit dah.

Compression programs work in a similar fashion, except they

TIMES NEWS TRACKER **NEW**

Topics	Alerts
Language and Languages	<input type="button" value="Create"/>
Physics	<input type="button" value="Create"/>
Science and Technology	<input type="button" value="Create"/>
Computer Software	<input type="button" value="Create"/>

Create Your Own | Manage Alerts
Take a Tour

Sign Up for Newsletters

NYT STORE

Scientists at Work



Buy this book for \$19.95 .

DIRECT INVESTOR JACKIE PEARSON

Be a direct investor

ADVERTISER LINKS

Revlon Run/Walk
Saturday May 4th

Find More Low Fares!

Compression programs work in a similar fashion, except they invent a new code for each message based on patterns unique to that message. The program might, for example, find that a text uses the word "compression" frequently and save space by substituting a two-letter abbreviation.

Experience Orbitz!

Sale to Europe from \$199, offer ends 5/2

The Italian physicists understood that a compression scheme invented to compress a text written in English would do a poor job on one written in Italian.

REPRINTS & PERMISSIONS

[Click here to order Reprints](#)
Permissions of this Article

"Transmitting an Italian text with a Morse code optimized for English will result in the need of transmitting an extra number of bits," they wrote. They conjectured that just how many extra bits it takes would be a measure of the distance between English and Italian.

To demonstrate this, the researchers used a zip program to compress a text written in one language. They then appended to the original text some text written in Italian or another language and compressed that document. As predicted, the compression program did not do as good a job when the languages of the two texts were different. They tried the same trick on a group of texts all written in Italian, but by a variety of different authors. They found they could distinguish between the authors more than 90 percent of the time.

The scientists performed a further test of their technique by analyzing a single text that has been translated into many different languages - in this case the Universal Declaration of Human Rights. The researchers used their method to measure the linguistic "distance" between more than 50 translations of this document. From these distances, they constructed a family tree of languages that is virtually identical to the one constructed by linguists.

The researchers say linguistics is just a "playground" for them to sharpen their techniques. The same methods, they say, might help create order out of the rapidly accumulating libraries of DNA and protein sequences, earthquake catalogs and other geophysical data. They might even lead to a solution to one of the most troublesome problems of modern computer science: filtering the junk e-mail from your in box.



It's easy to follow the top stories with home delivery of The New York Times newspaper. [Click Here](#) for 50% off.



Copyright 2002 The New York Times Company | [Privacy Information](#)