

The internet

– now with a geographic dimension

BY CHRISTIAN S. JENSEN

Christian S. Jensen (born 1963), PhD 1991, D. Tech. 2000, Professor of Computer Science, University of Aarhus, Denmark, 2010. Appointment at Aalborg University, Denmark, and research stays at the University of Maryland (USA), University of Arizona (USA) and Google Inc. (USA). Member of the Danish Academy of Technical Sciences and the Royal Danish Academy of Science and Letters, and recipient of several national and international awards.

The amount of data in electronic form is growing exponentially. At the same time, the IT infrastructure, including the internet, which we use every day, is developing at great speed. For example, smartphones are proliferating rapidly while mobile bandwidth is increasing all the time. At the other end of the infrastructure, we see data centres springing up everywhere. These are buildings with large numbers of processors and hard drives which facilitate the handling of huge volumes of data as cheaply as possible. This trend is continually creating new challenges and opportunities. Christian S. Jensen received the Völum Kann Rasmussen Annual Award for Technical and Scientific Research for his work that includes contributing to the efficient storage of, and searching in, spatio-temporal data, which is data referenced by time and place. Part of this work aims at giving the Internet a geographic dimension. According to Christian S. Jensen, the annual award of DKK 2,500,000, will be used to enable further research into foundations of the internet of the future.



Data centres have large numbers of processors and hard drives that enable the handling of huge volumes of data.
Photo: Robert Scoble

Vast volumes of data

The digital universe, or the total volume of electronic data, is currently doubling every 18-24 months. It has been estimated that it encompassed 1.2 zettabytes in 2010. A zettabyte is 1024 exabytes, which is 1024 petabytes, which is 1024 terabytes. A terabyte is equivalent to what can now be stored on the single platter of a hard disk. In other words, it took 1.2 billion hard disks to store the digital universe as it existed in 2010. In 2020, the digital universe is expected to swell to 35 zettabytes.

It is estimated that there are about 250 million web servers and even more websites on the internet.

The number of documents on the internet exceeds 25 billion. Moreover, it is estimated that the Google search engine alone receives about 3 billion queries a day. These enormous volumes of data make for exciting challenges and new opportunities.

“Google” with a geographic dimension

Before long, the internet will be used more from mobile devices than from stationary computers. At the same time, it is increasingly possible to position mobile devices. It is also possible to attach a geographic location to many websites (such as a restaurant’s website). Studies show that about 20 per cent of all web queries are for results that are



Smartphones are proliferating rapidly. Photo: Cheon Fong Liew

geographically close to the user and thus have “local intent.” This makes it relevant to add a geographical dimension to “Google queries.” A normal query consists of keywords entered by the user. In response to the query, a list of links to web pages matching the search words is returned. Google’s goal is to respond within 200 milliseconds. An important question then is how to also simultaneously take the positions of the users and web pages into account.

Indexing

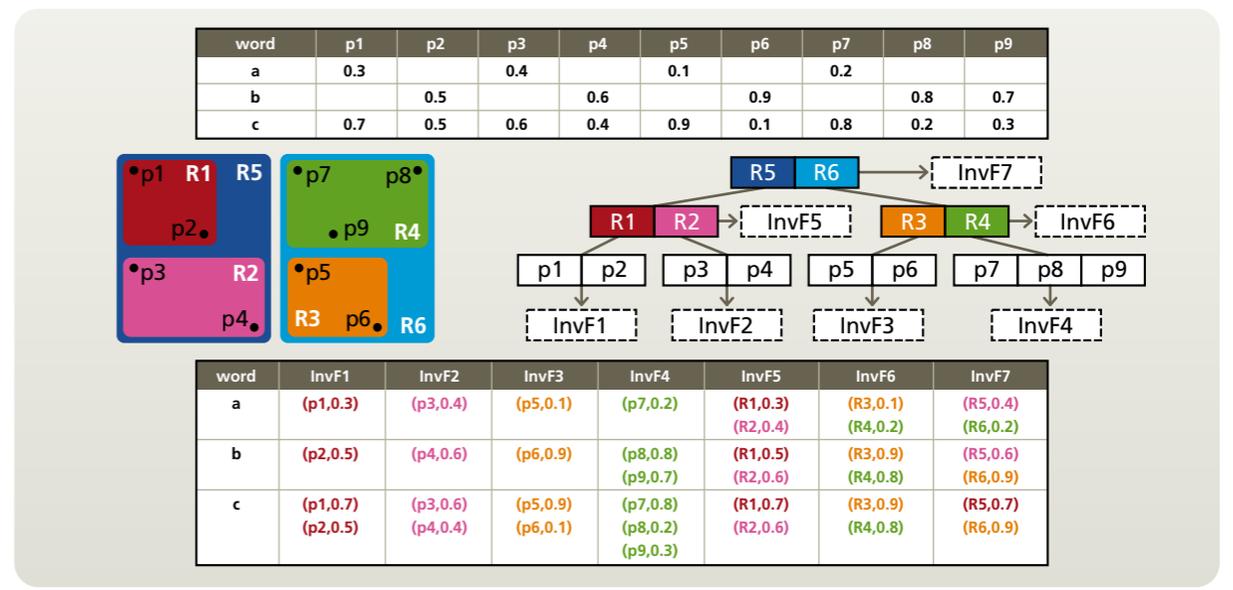
One type of geographic query finds links to web pages that both match the keywords and are close to the user’s location. It is a bad strategy to check every web page for every query. Instead, it is desirable to build indexes that make it possible to quickly disregard web pages that either do not match the keywords or are located far away from the user.

It has long been known how to create indexes that render the finding of web pages based on keywords efficient. These are the kinds of indexes used in search engines. As something new, we have developed indexes that simultaneously take into account the positions of the users and web pages. These new indexes make it possible to find relevant web pages by only looking through very little data, which enables short response times.

Mobile objects

Another challenge arises because mobile users are continuously on the move. Imagine that each of Facebook’s 500 million active users wants to see a list of their 10 closest Facebook friends. Or that tourists want to look up the 10 closest points of interest (cafés, pharmacies, etc.) that best meet their current needs.

The challenge here is to keep all lists updated as the users move, and to do this as efficiently and cheaply as possible. One strategy to solve the problem for the tourists works by first finding the 10 currently best points of interest. This can be done using the indexes described above. Then a safe zone around the tourist is calculated within which the current result does not change. When the tourist moves outside the zone, the tourist’s smartphone sends a message to the data centre where a new result and a new safe zone



Use of the IR-tree for the indexing of nine web pages (p.1 to p.9) with locations. The tree structure captures the spatial containment hierarchy shown on the left. The simplified text associated with the nine pages is shown at the top, and the bottom table shows the contents of the so-called inverted files associated with the nodes in the tree

are calculated and sent to the tourist. It turns out that safe zones can be described by multiplicatively weighted Voronoi cells.

Privacy

The technological advances that enable the services described here also have a downside in terms of access to increased surveillance and disclosure of private information. According to law professor Eva Smith, 82,000 pieces of information were registered about each Dane in 2008, corresponding to 225 pieces of information per day. The concept of location privacy includes aspects such as not wishing to disclose your exact location to a third party, but also not wishing to reveal that you are geographically close to another person or that you are not home.

Research on privacy shows that it is often possible to achieve support for privacy. For example, one

can find out where the “closest pizzeria” is without revealing one’s exact position – one can just ask for all pizzeria locations in the whole of Denmark and sort the results on the phone. The challenge is to return results at the lowest possible cost to the system. One promising strategy would, in this example, be to send pizzeria location queries at an increasing distance from a nearby false location until one is sure that one has received enough information to be able to provide the correct result for the correct location that only the phone knows.

Applications

In addition to making a wide range of location-based internet services possible, research in spatio-temporal data management also has applications in a number of other areas, including intelligent transportation systems, logistics, physical planning, marketing and epidemiology.