

# World's first stream aggregation technology to rapidly process both historical and incoming data

November 19 2012

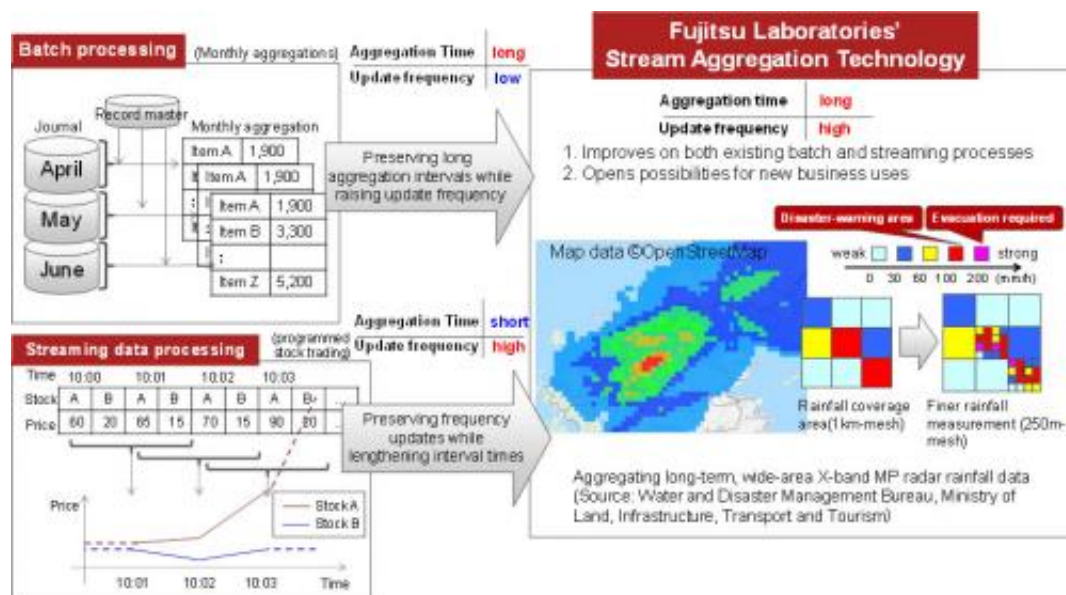


Figure 1: Characteristics of existing technologies and the new technology.

Fujitsu Laboratories announced development of the world's first stream aggregation technology able to rapidly process both stored historical data and incoming streams of new data in a big data context.

The nature of big data requires that enormous volumes of data be processed at a high speed. When data is aggregated, longer aggregation times result in larger data volumes to be processed. This means

computation times lengthen, which causes frequent updating operations to become more difficult. This is why improving the frequency of updates when aggregation times are lengthened has so far been challenging. Fujitsu Laboratories has therefore developed a technology that returns computation results quickly and manages snapshot operations, without re-doing computations or re-reading a variety of data types that change over time. As a result, even with high-frequency updating and long aggregation times, data can be processed 100 times faster than before.

This technology promises to improve both large volumes of batch processing and the processing of streaming data. Furthermore, in [meteorology](#), it is now possible to show concentrated downpours in specific areas. As well as the utility gained for future weather forecasting, it may also have uses in new fields that demand the ability to process [longitudinal data](#) in real time.

Details of this technology will be announced at a special workshop lecture of the Special Interest Group on Software Enterprise Modeling (SWIM) of the Institute of Electronics, Information and Communication Engineers (IEICE) held on Friday, November 30, at the Takanawa campus of Tokai University in Japan.

Many companies are interested in using advanced ICT technology to improve their competitive position by rapidly processing large volumes of data. Some uses are large-scale batch processes performed periodically on [transaction data](#), or processing streaming data in real time based on changing [stock prices](#).

In the data processing of such activities, aggregating computations is essential. In large-volume batch processing, however, there are differences in the aggregation times and update frequency. Typically, large-volume batch processes that emphasize throughput operate on

aggregation times lasting weeks or months. Streaming data processes emphasize response, on the other hand, and are in units of seconds or minutes. Update times roughly correspond with these.

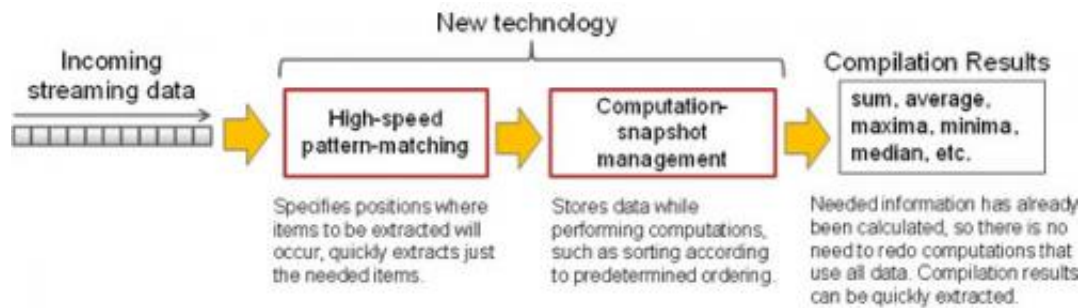


Figure 2: Flow of stream aggregation in this technology.

The emphasis on batch processes and streaming processes is different, and therefore the process needs to be adapted according to application.

Large-volume batch processing handles large volumes of historical data, so each round of processing re-reads all data, which creates long delays before results are ready.

The constant flow of data is held in a buffer—known as a window—and therefore each round of processing does not need to re-read any earlier data. Depending on the type of computation, however, the process does need access to all the data in that window in order to obtain computation results. For this reason, the duration of one round of computations will be proportionate to the window length, which diminishes responsiveness.

When using both historical (stored) and current (realtime streaming) data, with conventional processing methods, it has been difficult to simultaneously lengthen the aggregation intervals and raise the frequency

of updates for the reasons outlined above.

Fujitsu Laboratories has developed a fast stream aggregation technology for long aggregation intervals and frequent updates, based on a combination of the two technologies described below.

## **1. Rapid pattern matching technology:**

This is a technology that efficiently and directly picks out relevant items from an incoming stream of data. The conventional technique begins by analyzing the structure of input data and temporarily accumulating all input data in the memory. Next, it performs an extraction process of the items needed for aggregation to extract data. Structural analysis and item extraction is necessarily a two-step process. This technology is different in that it specifies the positions where items to be extracted will appear based on pattern matching, skipping over unneeded items thereby speeding up the process. Also, because pattern-matching is flexible, as well as using it with fixed-format data (such as CSV data) that conventional techniques use, it can work with other forms of data having recursive or hierarchical structures (such as XML data).

## **2. Snapshot operation management technology:**

This is a technology that quickly returns computation results to deal with a variety of data types that change over time, without re-reading or re-computing data. The conventional technique is to store in memory an incoming stream of data following its time sequence. This technology stores the data even as it performs required computations, such as sorting according to a predefined order. It is always managed based on its computed state (snapshot operation), and therefore never needs to redo computations that involve all the data, including not only sums and averages but also minima, maxima, and medians. This lets it quickly pick

out computation results.

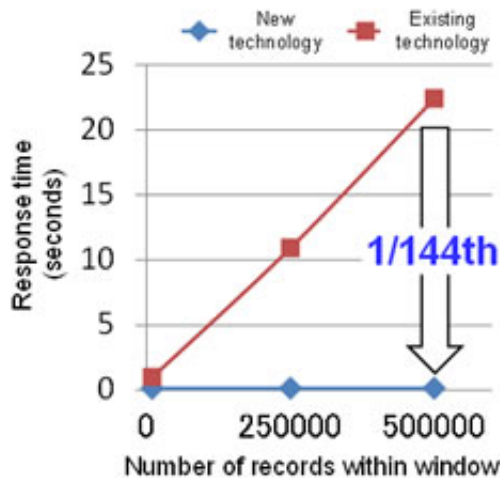


Figure 3: Comparison of response times between new and existing technologies.

The response time for aggregation results when using a window length of 500,000 records was shown to be roughly 100 times faster than the commonly used open-source Complex Event Processing engine. It was also demonstrated that response time does not depend on window length (Figure 3).

This technology is expected to have applications with regard to the utilization of high-precision sensor data. Fujitsu Laboratories conducted verification of the technology using rainfall data generated by XRAIN(1), a project conducted by the Water and Disaster Management Bureau of the Ministry of Land, Infrastructure, Transport and Tourism. In the case of aggregating rainfall volume data collected over several hours from 500,000 locations in the Kansai region of western Japan, every several minutes a window of approximately 100 million records needs to be processed. The test conducted by Fujitsu Laboratories

confirmed the technology's ability to execute data aggregation within intervals and no variation in aggregation times, and that the smooth movement of the rainfall area could be replicated, even for such a wide range of data (Figure 1). More than a sudden downpour, the actual volume of rainfall is what is strongly associated with disasters, and now, areas that require vigilance due to concentrated downpours can be readily verified.

Moreover, applications are anticipated for existing batch processing and stream processing. By enhancing the real-time [aggregation](#) of sales data, for example, it becomes possible to further strengthen production and inventory management.

Fujitsu plans to incorporate the new technology into its Big Data Platform and Big Data Middleware in fiscal 2013.

Source: Fujitsu

Citation: World's first stream aggregation technology to rapidly process both historical and incoming data (2012, November 19) retrieved 26 April 2024 from <https://phys.org/news/2012-11-world-stream-aggregation-technology-rapidly.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.