

New topological technique helps scientists 'see' and search large data sets

August 1 2013



Geographically correct map of the London Underground.

New computational techniques developed at Lawrence Berkeley National Laboratory (Berkeley Lab) may help save scientists from drowning in their own data. Computational scientists at the Lab have figured out how to streamline the analysis of enormous scientific datasets. The analysis uses the same techniques that make complex subway systems understandable at a glance.

They describe their work in a paper published in *PPoPP'13: Proceedings*

of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming.

What's the problem here?

Sophisticated sensors and supercomputers are generating bigger and more complex scientific datasets than ever before. In disciplines like genomics, combustion and [climate science](#), these datasets can range anywhere from tens of terabytes to several [petabytes](#) in size. A petabyte of data is equivalent to the storage consumed by 13.3 years of high-definition television.

To tease out the significant features for analysis, many scientists turn to a branch of mathematics called topology, which characterizes shapes of objects without considering aspects like length or angles—simplifying them the same way a subway [map](#) turns a complex maze of tunnels, trains and stations into colored lines and dots.

But scientific data are becoming so massive and complex that even simplified topological representations are becoming difficult to analyze efficiently. So more and more researchers are turning to massively parallel supercomputers to study their data.

The problem is that existing algorithms for analyzing data topologically do not take the best advantage of a supercomputer's thousands of processors. So Berkeley Lab [computational scientists](#) created a new approach—called distributed merge trees—for analyzing these datasets. This approach will let scientists make better use of next-generation supercomputers, as well as quickly sift out significant features from "noisy" data. In science, noise is irrelevant data generated by obstructing features, such as atmospheric dust in astronomical data, and measurement errors.

"The growth of serial computational power has stalled, so data analysis is becoming increasingly dependent on massively parallel machines," says Gunther Weber, a computational researcher in Berkeley Lab's Visualization Group. "To satisfy the computational demand created by complex datasets, algorithms need to effectively use these parallel computer architectures."

Both Weber and Berkeley Lab postdoctoral researcher Dmitriy Morozov pioneered the distributed merge tree approach to topological data analysis.

Topology: What is it good for?

Anybody who has looked at a pocket map of the London Underground or New York City subway system has seen topology in action. These seemingly simple representations ignore details like distance and physical locations of stations, but still preserve important information like what line a station is on, how different lines are connected and the overall structure of this complicated network.

In topology, what matters most is how things are connected. By disregarding distance, the size of an object no longer matters. The object can be stretched or squeezed and still remain topologically unchanged. So a large complicated structure like London's tube network can be condensed into an easy-to-read pocket-sized map by omitting geography and placing subway stations evenly on a line. Likewise, topology can be used to map the distribution of galaxies in the Universe, or burning regions in a combustion simulation.

Sifting through the noise

Once a massive dataset has been generated, scientists can use the

distributed merge tree algorithm to translate it into a topological map. The algorithm scans the entire scientific dataset and tags values that are of interest to the scientists, as well as merge points or connections in the data.



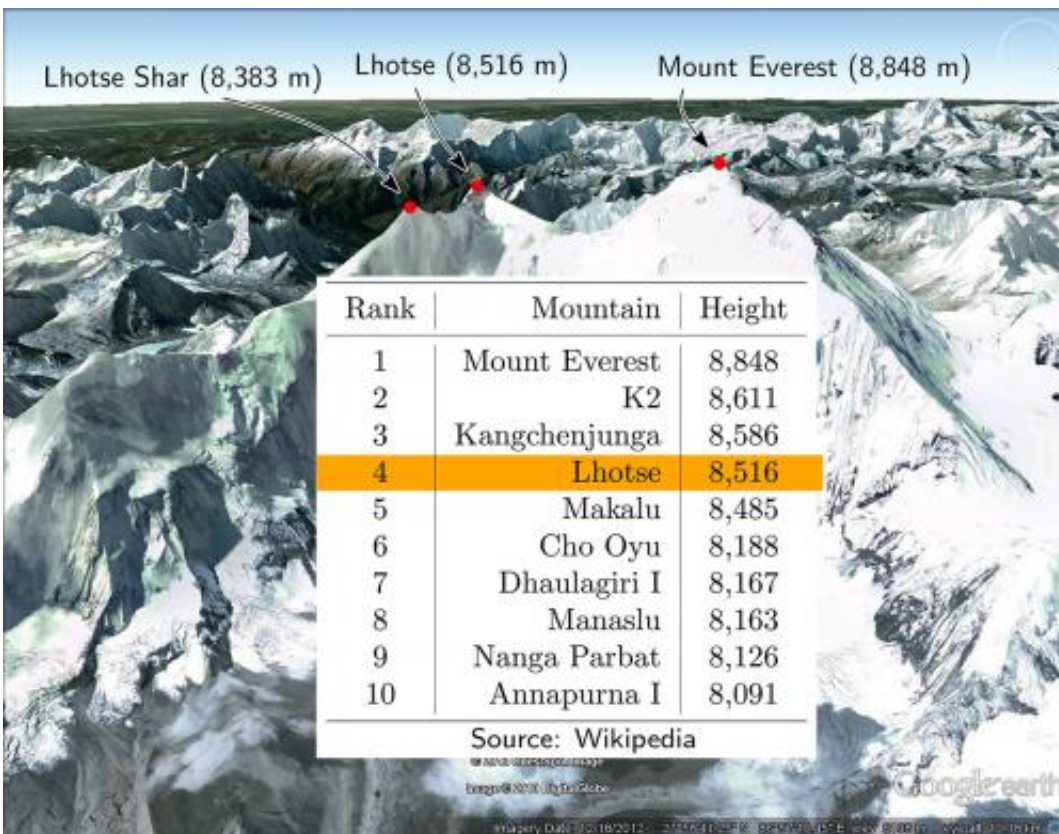
Topological representation of the London Underground.

So if the dataset is a height map of the Himalayas, the distributed merge tree algorithm will initially scan the entire dataset and record all the peaks in the mountain range. Then, it will generate a topological map illustrating how the different mountains in this dataset are connected. Morozov notes that these connections allow researchers to quickly differentiate between "real features" and "noise" in the data.

"A quick Internet search for the six tallest mountains in the Himalayas

will show: Mount Everest, K2, Kangchenjunga, Lhotse, Makalu and Cho Oyu. But, there are many more peaks in this range," says Morozov. "For example, the peak of Lhotse Shar is actually higher than Cho Oyu, but it doesn't count as one of the highest mountains because Lhotse Shar's base—where it merges into Lhotse—is almost as high as its peak."

Thus, a researcher that is only interested in the tallest mountains might consider Lhotse Shar "noise" in the data. A quick search of a topological map generated by the distributed merge tree will show that this mountain merges into Lhotse and disregard it based on the query.

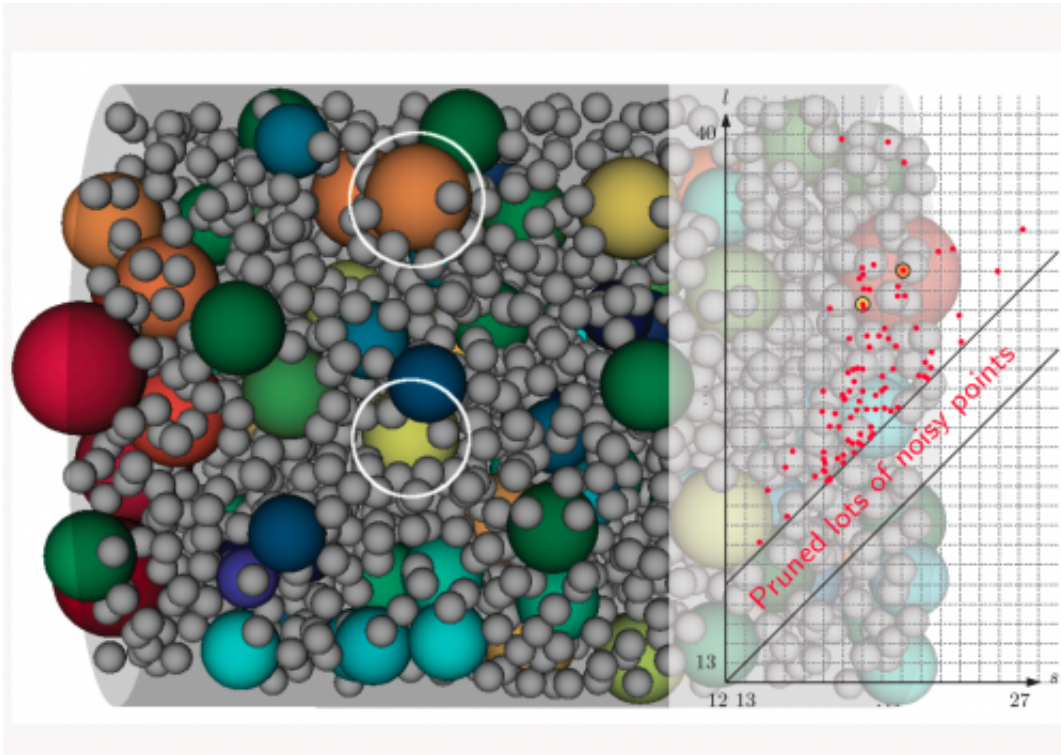


"This analogy applies to many areas of science as well," says Morozov. "In a combustion simulation, researchers can use our algorithm to create a topological map of the different fuel consumption values within a flame. This will allow scientists to quickly pick out the burning and non-burning regions from an ocean of 'noisy' data."

Parallelizing the search

According to Weber, distributed merge trees take advantage of massively parallel computers by dividing big topological datasets into blocks, and then distributing the workload across thousands of nodes. A supercomputer like the National Energy Research Scientific Computing Center's (NERSC) Hopper contains about 6,384 nodes, and each of which contains 24 processor cores.

In this approach, each block of data is assigned to a single node. Additionally, each node also stores an aggressively simplified global map, much like subway maps. The nodes only communicate to exchange information pertaining to a shared boundary.



Generated with the distributed merge tree algorithm, this visual characterizes a porous material at a glance. The colored spheres represent the pockets in the material, while the grey spheres represent solid material. The graph (right) shows how "prominent" individual pockets are by plotting the radius of each pore on the vertical axis, and radius of the largest sphere that can leave/escape the pocket on the horizontal axis. Prominence is the difference in sizes of the largest sphere that fits in and the largest sphere that can escape. The algorithm eliminates "noisy" data by ignoring pockets close to the diagonal, i.e., pockets that are not very prominent.

"Although the individual nodes simplify non-local parts of the map, each portion of the map is still available in full resolution on some node, so that the combined analysis of the map is the same as the unsimplified map," says Morozov. "By identifying thresholds, or tags, scientists can ensure that the desired level of detail required for analysis is not disregarded as the global map is simplified."

"Today, most researchers will only have one node keep track of the 'big picture', but the memory on a single compute node is often insufficient to store this information at the desired detail for analysis," says Weber. "These problems are further exacerbated by the trend in supercomputing to add more processor cores to a chip without adding correspondingly more memory. As a result, each new generation of parallel computers is operating with less memory per core than the previous one."

According to Weber, once a global topological representation is resident on a single node, it is difficult to parallelize queries for features and derived quantities. This in-turn leads to processor underutilization and slower queries.

"By reducing the tree size per node while maintaining a full accurate representation of the merge tree, we speed up the topological analysis and make it applicable to larger datasets, " says Weber. "This is also an important step in making topological analysis available on massively parallel, distributed memory architectures."

More information: dl.acm.org/citation.cfm?id=2442526

Provided by US Department of Energy

Citation: New topological technique helps scientists 'see' and search large data sets (2013, August 1) retrieved 13 March 2024 from <https://phys.org/news/2013-08-topological-technique-scientists-large.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.
