

Letter from the Special Issue Editor

One thing that I have always found fascinating about computer science is how a relatively technical discipline, with arcana not understood by the majority of the world's population, can relatively quickly produce inventions that impact the lives of large portions of that population. Big data processing is an example - it was big data processing, among other things, that helped the internet grow from a technical curiosity to a daily essential part of billions of people's lives. Big data processing enables the planning and tracking of the movement of products around the globe, manages the financial records of the world's economy, and impacts the world in myriad other ways.

However, we have begun to see that big data processing has the potential to impact the world in another, less desirable way. Its insatiable desire for more CPU cycles, more disk accesses, and more network bits means that ever-increasing energy resources are needed to power and cool the massive arrays of servers that actually do the work. This appetite for energy has driven the creation of new types of datacenters, server hardware and cooling technologies. But more work is needed to control energy usage that continues to grow.

Luckily for computer science researchers, there is an opportunity to do interesting research on this problem while helping the world as a whole. In this issue, we highlight four efforts that illustrate the challenges, and the opportunities, for energy-aware big data processing.

- The real electric grid is not just a smooth, unending stream of electrons, and renewable energy sources in particular provide varying amounts of power based on the current wind or solar radiation. Krioukov et al explore how batched analytic workloads are well suited to become "supply-following" - doing more work when clean renewable energy is available, and drawing less power when the only supplies are from dirty, non-renewable sources.
- Often, computing query results "as fast as possible" is not strictly required. Lang, Kandhan and Patel suggest that query processors can pick the most energy efficient plan that is still "fast enough" to meet SLA requirements.
- Recent batched data architectures, such as Dryad or MapReduce, are explicitly designed to use whatever resources are available. In reality, however, different jobs have different needs for CPU, disk and memory. Xiong and Kansal describe how batch systems can schedule work at a fine grain to use only necessary server resources, allowing unneeded resources to be used for other jobs or not to be used at all.
- Of course, most efforts to reduce energy usage are moot if we cannot effectively measure the usage! While TPC is popular for measuring performance, energy aware versions of TPC can be used to measure consumption. Poess and Nambiar describe techniques to analyze TPC results and estimate the associated power usage.

The first three articles share a common motivation of adapting query scheduling and planning to minimize usage of dirty, wasteful energy. The fourth article demonstrates how mature, well-used benchmarking techniques can be extended to evaluate the effectiveness of these query scheduling and planning techniques. It is my hope that these works, in addition to standing in their own right, will spark discussion and further research in how to manage energy consumption even as we continue to process ever larger amounts of data.

I would like to thank the issue authors for contributing and revising their work. I would also like to thank Dave Lomet for his advice and assistance throughout this process.

Brian Cooper
Google, Inc.
Mountain View, California, USA