

**Title:** Connecting large experimental facility and computing facility with streaming analytics

**NAME:** John Wu, Alex Sim, and Dula Parkinson

**AFFILIATION:** Lawrence Berkeley National Laboratory

**EMAIL:** John.Wu@nslcr.gov

DOE is the leading sponsor of research activities in the physical sciences. Its research portfolio includes large experimental facilities and large computing facilities, with a high-speed network serving as the glue that binds all others into a grand tool of discovery. However, so far, the connections among experimental facilities and computing facilities are quite weak, where many experimental facilities are sending users back to their home institutions with thumb-drives holding the data records collected at their facilities. As data rate increases at the experimental facilities, many runs are producing more data than could fit on a thumb-drive. In addition, to fully utilize the experimental facilities, many experiment designs are calling for near-real-time feedback to provide control signals. Some of these experiments require a significant amount of computing power than could be easily provided onsite at the experimental facility. For these reasons, it is urgent to develop the capability to provide online stream analysis capability.

In the past two decades, the volume of data has grown exponentially and this growth in data volume was the primary driver in the network usage. In the next decade, we anticipate a dramatic rise of real-time applications. In commercial application, this is appearing as the Internet of Things. In DOE science community, these real-time use cases will be initially dominated by distributed collaboration on large-scale experiments, such as ITER and LSST, where many participants may contribute to the analysis of an ongoing experiment and provide feedback to the control of the next stage of the experiments. In time, the experiments that currently produce modest amounts of data are expected to produce a much large volume of data, because of the increased resolution of sensors used. A large number of experimental facilities are also developing support for more dynamic experimental conditions, for example, some stations at ALS are supporting heating and cooling of experiment samples to a wide range of temperatures to observe the dynamic properties in real-time. These real-time observations frequently demand immediate analysis to provide feedback about experimental conditions. These analysis operations could benefit from the computing power available in the various computer centers. Furthermore, some of the data necessary for the analysis might be located far away. In these cases, there is a need to transfer a potentially large amount of data in a limited time window. As the work on Internet of Things proceed, many of the large science experience will adopt the automated measurement devices that can make its own decisions in real-time and therefore increasing the need for supporting real-time applications.

Due to the limitation of speed of light, the distributed computing facilities could only support near-real-time use cases, which we call online use cases<sup>1</sup>. Currently, the DOE network technology is primarily designed for moving bulk data among the various large computing facilities, to support low-latency workloads would require fundamentally different approach. Emerging technologies such as software defined networking and named data networking promise to enable dynamic routing of data and opportunistic harvesting of computing capabilities. Such features may allow a larger variety of computing resource to be available for analyzing the experimental and observational data, and effectively reducing the response time of the online data analysis tasks.

---

<sup>1</sup> Use cases with strict real-time needs will need to be collocated computing resources. Such use cases are not considered in this white paper.