

combination of both the state-of-the-art experimental capability at SSRL and novel scientific big data mining methods. The nanoscale spectro-microscopy method is capable of generating large scale data at a rate of about 30,000 spectra per second, posing a major challenge for scientists to analyze and interpret the data. Traditionally, scientists rely heavily on the assumption of complete prior knowledge about the principle chemical species in the system before more quantitative further analysis can be done. In the traditional approach, the data are simply fit to the known spectroscopic signatures of the chemical species that are known to be present in the system, which possibly could have left out some components that were unexpected. In this new study, the researchers utilized an unsupervised data mining method known as the DBSCAN (density-based spatial clustering of applications with noise) [4] to automatically perform big data classification, which efficiently (within a few minutes) extracts the scientifically relevant information with very little human interaction, and to finds out what is actually in the sample.

The presented data analytics approach can also be directly applied to the studies in many other fields, including for catalysts, batteries, fuels cells, and optical devices, in which the rational design of hierarchically complex functional materials plays an important role. Investments in new hardware (from new sources, to new detectors) has resulted in commissioning of exciting new facilities and has driven development of novel experimental techniques, and promised an unprecedented deep understanding of complex materials and devices as they function under real world conditions. However, because of hierarchy of scales in a functioning device, and complexity of the materials, these techniques produce data at an exponentially rising rate; a rate that has far surpassed any human's ability to curate and extract scientific knowledge from it in a comparable time frame. Without development of new data analytics approaches, inspired by advances in unsupervised feature extraction of "big data" over the last decades, minimal-loss data compression, and machining learning that can accurately extract hidden features, trends and high-level scientific information with minimal reliance on humans, from complex high-dimension data sets, the promise of an unprecedented new understanding of the materials world will most likely remain only a dream.

References

- [1] X. Jiang, G. Shen, Y. Lai and J. Tian, "Development of an Open 0.3 T NdFeB MRI Magnet", *IEEE Trans. Applied Supercond.* **14**, 1621 (2004).
- [2] M. T. Thompson, "Practical Issues in the Use of NdFeB Permanent Magnets in Maglev, Motors, Bearings, and Eddy Current Brakes", *Proc. IEEE* **97**, 1758 (2009).
- [3] A. Wang, H. Li and C. T. Liu, "On the Material and Temperature Impacts of Interior Permanent Magnet Machine for Electric Vehicle Applications", *IEEE Trans. Magn.* **44**, 4329 (2008).
- [4] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, "A Density-based Algorithm for Discovering Clusters in Large Spatial Databases with Noise", *Proc. 2nd Int. Conf. Knowl. Discov. Data Min.* (1996).

Primary Citation

X. Duan, F. Yang, E. Antonio, W. Yang, P. Pianetta, S. Ermon, A. Mehta and Y. Liu, "Unsupervised Data Mining in Nanoscale X-ray Spectro-microscopic Study of NdFeB Magnet", *Sci. Rep.* **6**, 34406 (2016), DOI: 10.1038/srep34406.

Contacts

Yijin Liu and Apurva Mehta, Stanford Synchrotron Radiation Lightsource

SSRL is primarily supported by the DOE Offices of Basic Energy Sciences and Biological and Environmental Research, with additional support from the National Institutes of Health, National Institute of General Medical Sciences.