

Paper Review: Fast Error-bounded Lossy HPC Data Compression with SZ

Summary

The purpose of this study was to create a novel compression scheme (SZ) that can be used to manage the growing amount of data required to run HPC systems. Many current compression techniques can guarantee no compression errors yet suffer from low compression ratio, or vice versa. SZ attempts to fit/predict data based on bestfit curve models. When data cannot be predicted, it uses an optimized lossy compression. This combined technique has proved effective compared to tools like Gzip, ZFP, etc.

Strengths

I like that the researchers evaluated applications from several different scientific domains in the study. This shows that their compression technique can be useful for a variety of purposes. I also enjoyed seeing some pseudocode that showed how the compression algorithm worked. As a computer scientist this helps me a lot more than just reading long paragraphs and cryptic equations.

Shortcomings

Improvements

Question(s) for Presenter

Why is it so important to convert the n-dimensional array into a 1d array for this technique? Why couldn't they just check the data in row by row from left to right?

Additional Questions

- What is the motivation behind SZ?
 - SZ is a novel HPC data compression scheme with strictly bounded errors and low overheads. The motivation is to improve upon current/previous compressions techniques that sacrifice compression ration for compression error rate or vice versa.
- What are the different curve-fitting methods SZ uses to predict data? Describe.
 - Preceding Neighbor Fitting (PNF)
 - Simplest prediction model, just uses the preceding value to fit the current value.
 - Linear-Curve Fitting (LCF)
 - Assumes that the current value can be estimated by the linear line constructed using its previous two consecutive values.
 - Quadratic-Curve Fitting (QCF)
 - Assumes that the current value can be predicted precisely by a quadratic curve that is constructed by the previous three consecutive values.
- What error bound types does SZ support? Do you think they are sufficient, or should other error bounding metrics be added?
 - Absolute and relative error bound.
 - I don't feel like I know enough to advocate for or against the chosen bounding metrics. However, since absolute is set to constant and relative is based on a linear function of

the global data, perhaps another metric could be added which would be based off of a quadratic function based on the same data somehow.