Towards Parallel Access of Multi-resolution and Multi-dimensional Scientific Data



Sidharth Kumar, Venkatram Vishwanath, Philip Carns,

Valerio Pascucci ,Robert Latham, Tom Peterka, Michael Papka, Robert Ross



Motivation

IDX enables fast and efficient access to large scale scientific data. Problem with current implementation

- Existing tools for writing IDX data only provides only a serial interface.
- HPC applications fails to utilize available parallel I/O resources. Solution
- Investigate methods for writing IDX data in parallel
- Enable HPC applications to write IDX data with scalable performance

Visus (IDX Format)

Multiple levels of resolution : Easy to query data of desired resolution and dimension



(530 Mib Data Set of Rat Hierarchical Z (HZ) ordering - The crux behind IDX data

format, computation requires spatial coordinates of data.

^	1		Value	Order	z	Lever	
0	0	0	٥	٥	0	0	Data is then reorganized into levels corresponding
0	0	1	1	1	1	1	Data to their reerganized into to toto corresponding
٥	1	0	2	2	2	2	to the formulation:
0	1	1	3	3	3	2	
1	٥	0	4	4	4	3	Level = floor ((log2 (HZ index))) + 1
1	٥	1	5	5	6	3	Level corresponds to different resolution level the dat
1	1	0	6	6	5	3	Level concopondo to amercini resolution level the da
1	1	1	7	7	7	3	is rearranged into.

Visus Serial IDX Writer

A micro benchmark that divides entire data volume into smaller 3D chunks. Each process independently writes to the IDX data set. No concurrent writes due to conflicts in updating metadata and layout. MPI barriers and tokens used to maintain order amongst processes



•The maximum performance is only 9.5 MiB/s. Using IOR, we obtain a peak performance of 218MiB/s for 64 processes writing a total of 8GiB. Thus, we are able to achieve only 4% of the maximum throughput.

·Jumpshot profile of 64 processes writing an IDX file using ViSUS. Each horizontal line represents a process, the yellow regions correspond to waits and black regions correspond to time spent writing data. As expected, we notice that a large portion of the runtime for each process is spent waiting for the I/O token.

PIDX : Prototype API For Parallel IDX Writes

This API is called Parallel IDX (PIDX) and includes functions patterned after ViSUS for creating, opening, reading, and writing IDX data sets. Each of the PIDX functions is a collective operation that accepts an MPI communicator as an argument.



The data in this case is provided in the form of a contiguous, row-major ordered data buffer. Each process calculates an HZ ordering for its sub-volume, reorder the data points accordingly, and write the data points to interleaved portions of the IDX data set.

Optimization Strategies

The PIDX prototype described in the previous section greatly improved I/O performance over serial use of the ViSUS

MPI File Caching

Following figure highlights the time spent by rank 0 when writing data into the four initial HZ levels. The Pink regions represent File Open Time.



Performance of the API is improved by performing MPI file caching while writing data, which prevents unnecessary file opens corresponding to each HZ level.

HZ computation Optimization

Indentified bottlenecks associated with redundant computations as well as inadequate use of the floating point double hummers.

From the graph above: there is almost 220% increase in performance of the API after the two different optimizations.



Single process achieves approximately 6.85 MiB/s, comparable to the speed of a serial writer for an equal volume of data.

The peak aggregate performance of 406 MiB/s is reached with 512 processes. This is approximately 60% of the peak IOR throughput achievable (667MiB/s) on 512 cores of surveyor.

corresponding to file opens and purple corresponds to file close, all the processes writes to the IDX file in parallel.

Future Work

Both with weak and strong scaling, a limit is hit on scalability with our current implementation, falling short of the peak surveyor write performance achieved by IOR.

	revei	Time	Level	Time	
1.596	11	0.129	22	0.304	
1.444	12	0.151	23	0.363	
1.544	13	0.151	24	1.979	
1.530	14	0.148	25	2.140	
1.653	15	0.102	26	2.706	
1.640	16	0.085	27	4.605	
1.628	17	0.077	28	6.760	
0.148	18	0.086	29	10.533	
0.132	19	0.087	30	18.251	
0.134	20	0.140			
0.132	21	0.211			
	1.596 1.444 1.544 1.530 1.653 1.640 1.628 0.148 0.132 0.134 0.132	1.596 11 1.444 12 1.544 13 1.530 14 1.653 15 1.640 16 1.628 17 0.148 18 0.132 19 0.132 20 0.132 21	1.596 11 0.129 1.444 12 0.151 1.544 13 0.151 1.530 14 0.142 1.633 15 0.102 1.640 16 0.085 1.628 17 0.077 0.148 18 0.086 0.132 19 0.087 0.132 20 0.140	1.596 11 0.129 22 1.444 12 0.151 23 1.544 13 0.151 24 1.530 14 0.148 25 1.633 15 0.102 26 1.640 16 0.085 27 1.628 17 0.077 28 0.132 19 0.087 30 0.132 20 0.140	

We found that contention and metadata overhead caused levels 0 through 6 to take a disproportionate amount of time relative to the amount of data that they were writing. We plan to leverage aggregation strategies to better coordinate I/O in the first few cases where processes contribute many data to the same level of the IDX data set.

Time Taken To Write The Various Idx Levels For A 8 Gib Data Volume On 64 Processes

Publication

S. Kumar, V. Vishwanath, P. Carns, V. Pascucci, R. Latham, T. Peterka, M. Papka, R. Ross. "Towards Parallel Access of Multi-dimensional, Multi-resolution Scientific Data," 5th Petascale Data Storage Workshop, Supercomputing 2010