

Introduction to Darshan: What to do when you aren't sure what to do

ATPESC 2019

Phil Carns Mathematics and Computer Science Division Argonne National Laboratory

Q Center, St. Charles, IL (USA) July 28 – August 9, 2019





exascaleproject.org

Understanding I/O problems in your application

Example questions:

- How much of your run time is spent reading and writing files?
- Does it get better, worse, or the same as you scale up?
- Does it get better, worse, or the same across platforms?
- How should you prioritize I/O tuning for the most impact?

We recommend using a tool called **Darshan** as a starting point.

This presentation is a introduction; we'll see more detailed examples of Darshan in action later in the day.







What is Darshan?

Darshan is a scalable HPC I/O characterization tool. It captures a concise picture of application I/O behavior with minimal overhead.

- Widely available
 - Deployed at most large supercomputing sites
 - Including all of the DOE Office of Science facilities that we use for ATPESC training
- Easy to use
 - No changes to code or development process
 - Negligible performance impact: just "leave it on"
- Produces a *summary* of I/O activity for every job
 - This is a great starting point for understanding your application's data usage
 - Includes histograms, timers, counters, etc.



How does Darshan work?

Darshan is primarily intended for MPI applications.* It inserts lightweight instrumentation when your program is compiled (for static executables) or when your program is executed (for dynamic executables).

- Records statistics about file accesses
 - These statistics are stored independently in bounded, compact memory at each rank
- Aggregates statistics when the application exits
 - Collect, filter, compress and write one summary file for the job
- The results are left in a log file that you can inspect with command line tools
 - Usually start by generating a summary PDF that plots metrics of interest



Using Darshan on Theta: make sure the software is loaded

These steps are similar on other platforms; check your site documentation!





Using Darshan on Theta: instrument your code

Compile and run your application!

That's all there is to it; Darshan does the rest.*

* Well, Ok, one caveat: your application *must* call MPI_Finalize() (or your programming language's equivalent) to produce a Darshan log.



Using Darshan on Theta: find your log file



All Darshan logs are placed in a central location. This is the path on Theta. **Check your site documentation!**

Go to subdirectory for the year / month / day your job executed.

Be aware of time zone (or just check adjacent days)! Theta, for example, uses the GMT time zone and will roll over to the next day at 7pm local time.

Find the one you want, and copy it somewhere to analyze.



Using Darshan on Theta: generate summary



The main result is a PDF file that summarizes your I/O.





Job analysis example

- The summary PDF contains a few pages of graphs and charts.
- The first page looks like this.
- We'll highlight some key sections in the next slides.



Job analysis example

ior (6/29/2017)				
jobid: 5598836	uid: 52352	nprocs: 256	runtime: 4 seconds	

The header shows basic information about your job:

- Executable name and date
- Job ID, User ID, number of MPI processes, total execution time

I/O performance estimate (at the MPI-IO layer): transferred 79456 MiB at 8083.73 MiB/s I/O performance estimate (at the STDIO layer): transferred 0.1 MiB at 3.86 MiB/s NOTE: STDIO performance appears low, but that's because it didn't transfer enough data to sustain throughput. In this case this I/O was access to a configuration file.

The next lines show an estimate of your I/O performance. It might display one or more of:

- MPI-IO performance (we'll learn more about this later)
- POSIX performance (open/close/read/write)
- STDIO performance (fopen/fclose/fread/fwrite)



Job analysis example



The first graph shows the percentage of execution time that was spent performing I/O.

If the percentage is low, then maybe I/O shouldn't be your top priority for optimization?



Job analysis example



A histogram indicates the distribution of access sizes.

Recall from introduction: if you see many small reads or writes (big spikes on the left hand side), then you are probably not taking advantages of the file system's strongest assets.

File Count Summary						
(estimated by POSIX I/O access offsets)						
type	number of	files	avg. size	max size		
total opened		259	16M	16M		
read-only files		33	16M	16M		
write-only files		226	16M	16M		
read/write files		0	0	0		
created files		226	16M	16M		

A table indicates file counts and sizes for a few different categories of files opened by your application.



Job analysis example



A "timespan" graph on the 2nd page shows the periods of time when the application was reading or writing.

Remember to contact your site's support team for help! The Darshan summary can be a good discussion starter if you aren't sure what's happening in your job or how to improve it..

There are additional graphs in the PDF file with increasingly detailed information not shown here. You can also dump all data from the log in human-readable text format using "darshan-parser".



Darshan: tips and tricks



What if you are doing shared-file IO?



- Your timeline might look like this.
- No per-process information available because the data was aggregated by Darshan to save space/overhead.
- It is shown as just one line for "all processes" in bottom graph.
- Does this matter? The level of detail that you need depends on what you want to learn about your application.



What if you are doing shared-file IO?



#!/bin/bash -l		
#SBATCH -p debug		
#SBATCH - A m888		
#SBATCH -N 8		
#SBATCH -t 0:15:00		
#SBATCH -J ior-example		
#SBATCH -o ior-example.o%j		
#SBATCHmail-type=BEGIN,EN	ID	
#SBATCHmail-user=carns@mc	s.anl.gov	
#SBATCH -C haswell		
#SBATCH -L SCRATCH		
export DARSHAN_DISABLE_SHARE	D_REDUCTION=1	
<mark>s</mark> tripe_medium \$SCRATCH/ior srun -n 256 ./ior/src/ior -f ~	ior-shared.conf	

- You can optionally set an environment variable in your job, DARSHAN_DISABLE_SHARED_REDUCTION, that tells Darshan *not* to summarize shared file access.
- Every rank will report it's own timespans.
- This increases overhead and log size, but can be very helpful in some cases.





Detailed trace data



xport DXT_ENABLE_I0_TRACE=4

stripe_medium \$SCRATCH/ior srun -n 256 ./ior/src/ior -f ior-shared.conf

- What if that still isn't enough detail? You can also capture a full trace with timestamp, offset, and size of every I/O operation on every rank.
- Set the DXT_ENABLE_IO_TRACE environment variable in your job to enable this feature.
- This causes additional overhead and larger files, but captures precise access data.

You can dump the trace information with "darshan-dxt-parser."

	pcarns@cor 5599892_6-2	L12:~/W 29-6544	огкіпд/ 3-43686	other/nersc-d 3287276195393	arsnan-seminar 2_1.darshan	201//logs> (darsnan-dxt-pa	arser pcarn	s_lor_ld	
	# ******* # DXT_POSI # *******	******* X modul ******	******* e data ******	***********	**************	**				
	# DXT, fil # DXT, ran # DXT, wri # DXT, mnt # DXT, Lus # DXT, Lus	e_id: 1 k: 0, h te_coun _pt: /g tre str tre OST	1542722 ostname t: 16, lobal/c ipe_siz obdidx	479531699073 : nid00511 read_count: : scratch1, fs_ e: 1048576, I : 49 185 115	, file_name: /g 16 _type: lustre _ustre stripe_c 7 135 3 57 95	lobal/cscrat ount: 24 43 27 191 1	ch1/sd/pcarns 163 51 15 153	/ior/ior.d	at-summa 1 239 79	
	25 137 47 # Module	Rank	Wt/Rd	Segment	Offset	Length	Start(s)	End(s)	[OST]	
	X_POSIX	O	write	0	Θ	1048576	0.7895	0.8267	[49]	
	X_POSIX	0	write	1	1048576	1048576	0.8267	0.9843	[185])e
17	AT X_POSIX	0	write	2	2097152	1048576	0.9843	1.0189	[115]	
	X_POSIX	0	write	3	3145728	1048576	1.0189	1.0250	[7] ^{ONAL}	LADUK

Detailed trace data



- You can also plot trace data using "dxt_analyzer.py".
- Example on the left:
 - Looks similar to the timespan plots that you already get in the normal Darshan summary.
 - But it plots every individual operation precisely, rather than just showing ranges of times that each process was performing I/O.
 - Closer inspection of the log can identify exactly when and where problematic access patterns were issued.



Darshan: a recap

- We've gone through a lot of Darshan usage tips.
 - Probably more detail than you wanted whoops!
 - The important thing is to remember the possibilities; you can refer back to these slides later for details.
- Key takeaways:
 - Tools are available to help you understand how your application accesses data.
 - The simplest starting point is Darshan.
 - It's likely already instrumenting your application, or can quickly be made to do so.
 - Refer to documentation and site support for help interpreting.



Darshan: the future

- Darshan is actively supported and updated.
- What can you expect in the future?
 - Support for new HPC systems as they are deployed
 - Support for more kinds of applications (not just MPI)
 - Ability to collect data from applications that crash before MPI_Finalize()





Darshan hands on exercises

- The hands on exercises include 3 Darshan examples that you can try tonight or as time permits during the day:
 - helloworld: a simple application that you can run to test out the Darshan toolchain.
 - warpdrive and fidgetspinner: applications with A and B versions that you can compare to spot the performance differences (and their cause).

The warpdrive and fidgetspinner examples will be easier to understand after seeing some of the later presentations that include details about MPI-IO and performance tuning.

Check with the instructors to share what you find!







Thank you!







exascaleproject.org