HPC User Forum Update

# Interviews with HPC Community Leaders: Gary Grider, Los Alamos National Laboratory

Thomas Gerard and Steve Conway
September 2020

## IN THIS UPDATE

After the global pandemic forced Hyperion Research to cancel its 2020 HPC User Forums, we decided to reach out to the HPC community by publishing a series of interviews with select HPC thought leaders. Our hope is that these seasoned leaders' perspectives on HPC's past, present and future will be interesting and beneficial to others. To conduct the interviews, Hyperion Research engaged insideHPC Media. We welcome comments and questions addressed to Steve Conway, sconway@hyperionres.com or Earl Joseph, ejoseph@hyperionres.com.

This interview is with Gary Grider, leader of the High Performance Computing (HPC) Division at Los Alamos National Laboratory. As Division Leader, Gary is responsible for all aspects of high performance computing technologies and deployment at Los Alamos. Additionally, Gary is responsible for managing the R&D portfolio keeping the new technology pipeline full to provide solutions to problems in the lab's HPC environment, through funding of university and industry partners. Gary also helps manage the U.S. government investments in data management, mass storage, and IO. Gary has 26 granted patents, with 17 pending in the data storage area and has been working in HPC and HPC-related storage since 1984.

The HPC User Forum was established in 1999 to promote the health of the global HPC industry and address issues of common concern to users. More than 75 HPC User Forum meetings have been held in the Americas, Europe and the Asia-Pacific region since the organization's founding in 2000.

# GARY GRIDER INTERVIEWED BY DOUG BLACK, insideHPC

**Black**: I'm Doug Black, editor-in-chief of InsideHPC and today is part of our Hyperion Research series of interviews. We're speaking with Gary Grider, HPC Division Leader at Los Alamos National Laboratory. Gary, thanks so much for joining us today.

**Grider**: You're welcome.

**Black**: Tell us, if you would, a little bit about your area of expertise and also about some of the more important work you've done of late and what its potential implications and impacts might be.

**Grider**: Sure. Right now, I'm a manager, and I've been a manager, essentially the lead overall for HPC, since 2012 at Los Alamos, so my technical chops are from before that, unfortunately. I've been at Los Alamos since '89, so I've been here a while, and inside the DOE complex since '84. I guess my expertise and claim-to-fame is in the data storage area. Lee Ward and I, Lee Ward is at Sandia [Sandia National Laboratory], were the ones that went to DOE and got the money for Lustre, and I was with Peter Braam every Thursday morning for calls with his team to get Lustre built back in the late 2000's and, of course, that's used everywhere, so I feel pretty proud about that. Of course, when people complain about Lustre I feel a little bit responsible, but it's used everywhere and arguably it's been a great success. Certainly, those people at Cluster felt these were some really smart guys and they did a great job working on that for us and Peter is an awesome resource.

Another thing that's interesting is that Lee Ward and I, again, went to DOE and got money for what was called the NNSA ASC University Alliance Program for I/O and we started a number of projects. One of the projects we started was a project at University of California, Santa Cruz to build a test file system for testing out metadata at scale and security in metadata at scale in file systems. They asked if they could spend some time building a user-space object-oriented file system to be able to play around with and test ideas on what we asked them to do. There was a graduate student there named Sage Weil and, essentially, that's where Ceph came from. So, while I didn't give birth to Ceph, in some way Lee Ward and I did give birth to Ceph in that we got the money and sort of directed University of California, Santa Cruz to build what they were going to build. Not necessarily to build a product that was eventually going to become part of Red Hat and used all over the world, but to do some interesting R&D for us and that was the result. So, that is another example of an interesting technology that I was there at the beginning of.

We also were responsible for going to the High End Interagency Working Group. I don't know if you recall what that was during the Earth Simulator days, but that was a way for agencies to collaborate on what they funded in technologies and we did a bunch of projects with about 20 different universities. One of the ones that turned out really well was one we did with Rutgers and Stony Brook [University] and others. It was called Cache-Oblivious B-Trees but it eventually became Tokutek, which was a fast-insert database that was stuck underneath the Berkeley DB in Linux and eventually became $B^\varepsilon$-tree F S [file system], which is really a slick file system that's out there today that does inserts faster than anything else.

So those are three. There's a bunch more: NFS Version 4 and PNFS came from DOE funding that Lee and I helped manage at the University of Michigan. I could go on and on. I've been working in this area for a long, long time and a lot of products have come from DOE funding in this space and I was sort of the lead in that space for Los Alamos and Lee was the lead for Sandia. We kind of guided this whole area for 15 years or more and produced a lot of things that people use today.

**Black**: Wow. That's great stuff. It must be very rewarding to see technology you had a direct hand in, being so widely used.

**Grider**: Yes, of course, that was the goal and that was the idea and it's nice to accomplish goals once in a while.

**Black**: Can you tell us about any current project development work you're focused on?

**Grider**: Well, not me directly because, again, I have a couple hundred people that report to me and we run one of the largest computing centers in the world for HPC, so I do a lot of people management. The things that we're working on at Los Alamos that I think are interesting and that I have had a hand in shaping, and one thing we do a little different than most other labs - not all but most - is we tend to buy systems for small numbers of applications to run on half or more of the systems for the entire life of the system, which is kind of unusual for a large science site to do. One would never consider buying a $200,000,000 supercomputer and letting one application run on half of it for two or three years, but that's exactly what we do at Los Alamos. So, we kind of live on the edge of the DOE space in that regard. You've got NERSC and some of the Office of Science labs on one end of the spectrum that have lots of applications and would never give half the machine to anybody for very long at all and we live at exactly the other end of the spectrum where we do that all the time. In fact, that's the way we run our big machines. And one of the things we wanted to do, we ran it on Trinity, but we had been trying to run it on Sequoia at Livermore [Lawrence Livermore National Laboratory] and before that on Cielo at Los Alamos, and it's sort of a 20 year-old weapons issue that we've been trying to resolve for that long. We really needed to run it to get to the bottom of it, so Trinity, which is our current big machine on the floor - 2 petabytes of DRAM and 20,000 nodes which, by flops, isn't much but we don't really care about that. It's still probably the largest memory footprint of anybody out there and it's almost five years old.

We took half of that machine and ran a problem on it for over seven months to get the answer. The machine was designed for problems like that, in fact that one in particular, for sure. We realized that if you're going to run a problem like that with a memory footprint of a petabyte and you're going to do checkpoints every four hours and you're going to do analysis and the like, you're going to produce hundreds of petabytes of data per month probably, or something of that ilk, which is really a lot. So, we had to invent a few things to get that to occur. One of the things we sort of invented was a thing called the burst-buffer, which you may have heard of. It became a Cray product called DataWarp, which is an ephemeral flash-based file system. So, there's four petabytes of flash on this machine and you can checkpoint to it in about a minute, which means it runs at multiple terabytes per second, something like four. So, that's interesting, I think, that we are able to run these really large problems for a really long time because we can do checkpoint and analysis real quickly because we have sort of mastered the art of using flash inside the machine for that purpose.

We also realized that we needed to keep a lot of data around because if this particular application wanted to go backwards two months and start two months into the past and go in a different direction we need to keep a fair amount of those state checks around for a while, months and months, up to a year or something. That turns out to be a lot of data and we tried a lot of the object storage systems that are out in the world and had some success but it never really got to where we needed to be, so we built this thing called Campaign Store, which is kind of an object storage server but on steroids. It uses two tiers of erasure because we had a day where we lost 400 disk drives in the same day, which is not a lot compared to the population of disk drives but is still a lot of disk drives to lose in the same day. We built it so it could protect itself from such things and keep a lot of data around – hundreds of

petabytes around near the machine - to be able to restart and go backwards in time and the like, which are all things that our scientists want to do.

Another big thing that we are super interested in is efficiency. We don't really buy machines for flops at Los Alamos, we buy machines for running a problem in the most efficient way that we can to try to get answers in human times. So you won't find the word flop in any of our RFPs in the last decade. It's really about apps and workflows and how we can get the most number of those workflows through a system in its lifetime when the workflows are these gargantuan things that are petabytes of DRAM in size and the like. We're sort of a little different in that regard. We try really hard to service those applications. We don't have a lot of applications, but they do important work certifying the weapons complex.

Those are some of the things we've done recently, working on these efficient machines and making I/O efficient to be able to run these gargantuan things for long periods at a time and the like. That's some of it.

**Black**: Great stuff. Well, you've probably touched on some of the themes here I'm about to ask about, but looking at the big picture, what are your thoughts about where HPC is headed? Are there trends in place that you're excited about? I'm sure there are. Also, are there areas that give you concern?

**Grider**: I view this world as leveraging things. HPC hasn't been the biggest ship on the sea for quite a while. Certainly, it was between the 40's and 80's but, you know, late 90's and beyond it certainly wasn't and it isn't now. Cloud and other kinds of things are massive in scale compared to HPC. So we've always sort of been about leveraging ever since we became the smaller ship in the sea. There have been some exciting things to leverage, like the Campaign Store, we leveraged technology from cloud, erasure, and stuff like that. Leveraging flash, of course, the burst-buffer was a natural thing to do. We're trying to figure out how to leverage NVMe. NVMe over Fabrics look like really cool things to leverage. Smart networks are coming along, there is silicon everywhere that you can leverage to offload part of your problem to networks and the like. We're investigating that as well as offloading to storage and I think that's exciting. Finally, people are starting to talk about mixing some of the memory bandwidth problems and processors and that's kind of exciting. There are things that are coming up that seem quite leverageable and will move us ahead.

The part that I think bothers me a little bit is that classical simulation, which is about 80% of what we do at the lab if not more, seems to be falling out of favor in being the exciting thing to work on. Of course, there's AI and machine learning and quantum and all kinds of things which we all have to be involved in and we all have to pay attention to, but it feels like simulation has taken not just a back seat but is in a trailer hooked up to the back of the car. There are still massive problems to solve in that space and we haven't really come close to solving them in human times. If you think of me saying we ran a problem for seven months on 10,000 nodes and a petabyte of DRAM, that's not exactly human learning time, right? That's seven months to get a result. You don't get to talk about that once a week with your friends because it takes too long to get the answer. Those kinds of problems are going to become routine in the world that we live in. I see very few people interested in that problem. I see way more people interested in going where the money is, which is natural. Go after AI, there's where the money is. Go do interesting quantum things, that's where the money is. We're doing that, and that's all great and that's wonderful and we have to do that, but I just worry that simulation is getting a bit of a raw deal these days and we're building machines that aren't terribly well designed and usable for large-scale, complex simulations. There are not enough people doing it, apparently, or not enough

pressure, or it's just not sexy enough anymore to draw the "A-team" at any company. That's my big concern, that big, massive, extremely complex simulation is falling by the wayside.

**Black**: AI is certainly the shiny, bright new object, although I've also heard of this convergence of AI and simulation which sounds like it has great potential.

**Grider**: Maybe.

**Black**: Tell us a little bit about your background in HPC. How did you get involved in HPC in the first place?

**Grider**: Well, that's a long story. I came out of college as an electrical engineer who was classically trained in field theory, so I worried about transmission line field theory kind of stuff. Digital was just beginning in math departments. Everything was still analog then. I went to work for Sandia in the Wind Energy Program the first time there was a Wind Energy Program back in the 80's, designing magnets and the like. That ended when OPEC decided they'd make oil cheap again for a decade and put all of that stuff out of business. I was looking around trying to figure out what to do and I knew a lot about magnetics, and someone suggested I get involved in data storage. That's kind of how I entered the digital world and I worked at Sandia in that area for a little while and then came to Los Alamos as a contractor, and eventually as an employee, mostly working in large-scale data storage needs and applications. That was in the HPC organization in both labs, so that's kind of how I got into HPC and that's how I got into storage as well. It was by happenstance.

**Black**: Well, a great career and still going. We wish you the best of luck in your future endeavors and thanks so much for joining us today.

**Grider**: Thank you for your time.

## About Hyperion Research, LLC

Hyperion Research provides data-driven research, analysis and recommendations for technologies, applications, and markets in high performance computing and emerging technology areas to help organizations worldwide make effective decisions and seize growth opportunities. Research includes market sizing and forecasting, share tracking, segmentation, technology and related trend analysis, and both user & vendor analysis for multi-user technical server technology used for HPC and HPDA (high performance data analysis). We provide thought leadership and practical guidance for users, vendors and other members of the HPC community by focusing on key market and technology trends across government, industry, commerce, and academia.

## Headquarters

365 Summit Avenue
St. Paul, MN 55102
USA
612.812.5798
www.HyperionResearch.com and www.hpcuserforum.com