

S1 EP26 - Transforming Cloud Data Centers with CXL

Thursday, October 20, 2022 · 9:35

This week, join Sarvagya Kochak, Senior Director, Product Marketing, and podcast host Chris Banuelos to discuss the fundamentals of Compute Express Link (CXL) technology. Marvell, along with a consortium of industry leaders, is transforming the next generation of cloud data centers with a new industry standard for connecting processors, memory, and accelerators. With CXL based memory expansion and pooling, Marvell is uniquely positioned to transform the data center with architectures that have significantly improved scale, flexibility, and performance. Hear why CXL is imperative to the data center, its use cases, and future trends.

Speakers

Sarvagya Kochak

Sarvagya Kochak, Senior Director, Product Marketing

Host

Christopher Banuelos

Senior Manager of Global Social Media Marketing

C Christopher Banuelos 00:04

Welcome to the Marvell Essential Technology Podcast. I'm your host, Chris Banuelos. On today's episode I talk with Sarvagya Kochak Senior Director of Product Marketing on Compute Express Link, also known as CXL. We discuss the ins and outs of CXL, why it's important, some of the essential use cases, as well as what are some of the future trends related to CXL. To stay up to date on future episodes, please be sure to subscribe to the Marvell Essential Technology Podcast. Hey Sarvagya it's great to have you on today's episode, I'm looking forward to our discussion around compute Express link, otherwise known as CXL. Let's go ahead and first establish, what is CXL?

S Sarvagya Kochak 00:56

Chris, it's great to be here on your show today. And I'm super excited to talk about CXL otherwise known as Compute Express Link. It's a new high performance interconnect that leverages the PCI physical subsystem. And one of the key benefits it provides is byte addressable cache coherence. Now, that's a lot of words, but it's a fancy way of saying that a processor can now access memory, which is not connected over its DDR bus and treated just as if it was system memory. This capability is what is exciting so many people in the ecosystem. And there are literally over 100 companies that have signed up to be a part of the CXL consortium and drive the specification forward.

C Christopher Banuelos 01:49

And why is CXL so important?

S Sarvagya Kochak 01:51

Yeah, so the reason why CXL is really important to the industry today is because it will enable the next generation of infrastructure in the data center. Up to now one of the key things that's not been possible is attaching memory to devices and to processors and accelerators over a standard common unified interface. With CXL that's now possible.

C Christopher Banuelos 02:19
What are some of the use cases for CXL enabled products?

S Sarvagya Kochak 02:23
Yeah, so for for CXL we're looking at multiple usage models, the cool thing is that once you start taking memory and putting it on a CXL interface, there's a lot of things that it can enable and a lot of different ways in which you can see and use this memory. So we'll start off with the most basic use case, we call this memory expansion. And just as the name suggests, it's conceptually very simple that you already have some amount of memory attached to your host devices. Now, with CSL, you can expand on that memory without adding additional burden to the host memory subsystem. So for instance, with a CXL memory expander device, you can now increase the overall memory bandwidth that goes in to the host, you can also add memory capacity to the host. This is very useful, especially in a world where we are going from one generation to another with an increasing number of core counts and memory not scaling at the same rate. So now CXL can come and help augment all those issues.

C Christopher Banuelos 03:34
And what are some of the other use cases?

S Sarvagya Kochak 03:36
Yeah, so now, beyond memory expansion, the next step and in terms of the evolution of how memory will be used in the subsystem is what we call memory pooling. With Memory pooling, you can now think of a large number of hosts all talking to the same memory resource. This is actually very, very cool because an architecture like this has not really been possible up to date. And CXL fundamentally helps enable that. So you can think of, let's say four, or maybe eight servers within a rack that may be talking to almost something like the form of a memory appliance and be able to share the same memory in that appliance across all of those servers. And across all of those CPUs. In addition to being able to share that you can also dynamically partition the amount of memory capacity as well as memory bandwidth going into each of those hosts. So this is this gives a lot of flexibility, especially in public clouds, where user workloads may vary over time. And utilization also changes at the same time. Your memory is no longer tied down to the CPUs. So this can actually have a very significant net impact on TCOs. Now, if we take the next step along the evolution of Memory pooling. It's not just limited to CPUs, you can very well see use cases emerging where GPUs or accelerators can also share this memory pool along with CPUs. And that's one of the key drivers for CXL three. So you can actually share the same memory address range across all of these different types of devices. And we expect that this is going to drive the next generation of applications and the acceleration of use cases in the cloud, which is emerging with the deployment of CXL based solutions.

C Christopher Banuelos 05:38
Sarvagya, I wanted to get your perspective on another use case. And I remember a conversation that we had a couple of weeks back here at the office, and it was around memory acceleration. Can you talk a little bit more about that?

S Sarvagya Kochak 05:50
Yeah. So you know, we spoke about the benefits of just taking memory as is and moving it further away from all these host devices. Now, if you look at this subsystem, potentially, and how it will evolve, because of the amount of bandwidth that's involved in moving data back and forth between memory and the hosts, and the amount of power that consumes and the additional latency that it adds, we can actually see that it will make more sense to add some processing capabilities near this disaggregated memory pool itself, so that the host CPU can actually be freed up from some of the processing and the workload that it needs to over there. And the memory itself can become more intelligent and smart so that the CXL controller can be doing some data manipulation and help in terms of workloads that we see evolving in the future.

C Christopher Banuelos 06:51
And Marvell's technology is going to help enable all of these use cases you just described, correct?

S Sarvagya Kochak 06:56

Yes, at Marvell, we are investing in all the fundamental building blocks for developing these complex silicon solutions that will support very high performance, CXL interfaces, memory interfaces, as well as computing capabilities.

C Christopher Banuelos 07:16

My last question is, where do you see the future going? And what's next?

S Sarvagya Kochak 07:20

Yeah, so So what's next is, you know, we, of course, we want to productize, all of these things that we just discussed. But in addition to that, you know, we're looking at a complete end to end solution stack that would also in addition to these controllers include electrooptical solutions, so that we can provide very high performance interconnects and links to build out a complete CXL fabric within the data center. So we can look at, you know, a CXL fabric existing within a rack, you can think of a fabric going across racks. And there is a lot of investigation going on in the market today with regards to the right way in which people want to build these solutions. And based on all the discussions we've been having so far, I can tell you one thing for sure that there is no right way, everybody has a certain problem that they are trying to solve for. And we are going to see a very, very unique ecosystem of solutions and products that are going to be catered towards addressing these emerging user needs.

C Christopher Banuelos 08:24

It sounds like it's a really exciting time for you and your team. And it's great to see you working on this technology.

S Sarvagya Kochak 08:31

Yes, this is a very, very exciting time at Marvell. And it's definitely one of the most exciting spaces, I would say in the industry right now. There is a lot of enthusiasm in the end users for adopting this technology. There are a lot of POCs that are in flight right now, we demonstrated memory pooling at Flash Memory Summit (FMS). And we'll also be showcasing it at the Open Compute Summit (OCP). So, a very exciting time and expect to hear a lot more from us in the future as we make formal announcements around product.

C Christopher Banuelos 09:08

Thanks for joining today's podcasts are Sarvagya, really good to talk to you and look forward to hearing more.

S Sarvagya Kochak 09:13

Thank you, Chris. It's been a real pleasure and looking forward to coming back and talking about this more in the future.

C Christopher Banuelos 09:21

Thank you for listening to the Marvell Essential Technology Podcast. As always, please feel free to visit our website to learn more, and we'll see you on the next episode.



To deliver the data infrastructure technology that connects the world, we're building solutions on the most powerful foundation: our partnerships with our customers. Trusted by the world's leading technology companies for 25 years, we move, store, process and secure the world's data with semiconductor solutions designed for our customers' current needs and future ambitions. Through a process of deep collaboration and transparency, we're ultimately changing the way tomorrow's enterprise, cloud, automotive, and carrier architectures transform—for the better.

Copyright © 2022 Marvell. All rights reserved. Marvell and the Marvell logo are trademarks of Marvell or its affiliates. Please visit www.marvell.com for a complete list of Marvell trademarks. Other names and brands may be claimed as the property of others.