



Intel IPU's Fundamental Role In Your Cloud Strategy

CORPORATE PARTICIPANTS

Xiaojun (Shawn) Li

Intel – Sales Director, Next Wave OEM & eODM

Brian Niepoky

Intel – Director, Connectivity Group Marketing

Sabrina Gomez

Intel – Director, Programmable Solutions Group Marketing

PRESENTATION

Xiaojun Li

Welcome everyone to the Intel Network Builders Insights Series. I am Shawn Li, Sales Director, Next Wave OEM and eODM Network and Communications Sales Organization at Intel Corporation. I am your host for today's webinar. Thank you for taking the time to join us today for our webinar titled Intel IPU's Fundamental Role In Your Cloud Strategy.

Before we get started, I want to point out some of the features of the BrightTALK tool that may improve your experience. There is a Questions tab below your viewer, I encourage our live audience to please ask questions at any time. Our presenters will hold answering them until the end of the presentation.

Below your viewer screen, you will also find Attachments tab with additional documents and reference materials, which pertain to this presentation.

Finally, at the end of the presentation, please take the time to provide feedback using the Ratings tab. We value your thoughts and will use the information to improve our future webinars.

Intel Network Builders Insights Series takes place live every month, so please check the channel to see what is coming and access our growing library of recorded content. In addition to the resources you see here, we also offer a comprehensive NFV and 5G training program through the Intel Network Builders University. We can find the link to this program in the Attachments tab as well as a link to the Intel Network Builders Newsletter.

Today, we are pleased to welcome Brian Niepoky and Sabrina Gomez from Intel. Brian leads the Intel Connectivity Group Marketing and focuses on ethernet controller and adapters including Intel's new Infrastructure Processor Unit product line that we will discuss today. He is excited about the future of data-centric networking and believes that Intel offers the breadth of connectivity products to support a wide variety of the data center networks.

Sabrina Gomez is a Director of Marketing, Platforms, and Product Sales at Intel, where she leads Intel Product Marketing and the ecosystem for programmable acceleration platforms, as well as worldwide, regional FPGA product marketing. She and her global team are responsible for FPGA-based platforms and the solutions go-to marketing targeting cloud, networking, edge, ecosystem partner management, and regional marketing for all PSG FPGA products. Sabrina joined Intel with the acquisition of Altera Corporation.

Welcome, Brian, and Sabrina, and thank you for taking the time to join us today. Brian, I will hand it over to you to start off. Thank you.

Brian Niepoky

Thanks, Shawn. We are in the middle of an infrastructure revolution from the edge to the cloud that affects how we deliver compute. It's a huge transformational change that impacts how we enable IT, but also how we build data centers. When we started this revolution, the systems on the servers made up the cloud-- excuse me. There we go. The systems or the servers that made up the cloud data centers looked pretty much like the systems in the classic data center. But now we're seeing this architecture diverging.

Intel IPU's Fundamental Role In Your Cloud Strategy

In the classic data center, everything is owned by one party. In the cloud, the workloads and the systems are owned by different ones, tenants, and cloud service providers. Here is an example of a typical server in a classic enterprise data center. The physical structure, the hypervisor, and the applications are all owned by one entity. Think of a bank. All software runs on the CPU.

But for the servers that are built for the cloud infrastructure, a different architecture has emerged. They have a dedicated processor that runs the infrastructure functions in the cloud. We call this new category of processor an IPU or Infrastructure Processing Unit. The cloud service provider's software runs on this IPU, and the revenue generating guest software runs on the CPU. For example, a bank's financial app running on the CPU would be now cleanly separated from the cloud service provider's infrastructure software running on the IPU.

It's a bit like a hotel and a single family home. In my home, I want to be able to easily move from the living room to the kitchen or the dining room. In a hotel, the guest rooms, the dining hall, and the kitchen are separate. The area where the hotel staff works is different than the area where the hotel guests are. You may need a badge to go from one area to another. And we're seeing the same type of trend in the cloud infrastructure.

The advantages of IPU's are threefold. First, there's a strong separation of infrastructure function and the tenant workloads that allow tenants to take full control of the CPU. That's what I show on this foil. But there are two more. The second, the cloud operator can offload the infrastructure tasks to the IPU maximizing utilization for public clouds, it enables them to fully rent out the CPU cores and maximize revenue. And third, the IPU allows for a fully diskless server architecture in the cloud data center. And I'll explain a little bit more in some detail.

So, on this foil, in the servers within IPU, infrastructure and tenant workloads are cleanly separated with the tenant workloads running on the CPU and the infrastructure software running on the IPU. The immediate result is better isolation between the two. Spikes in the infrastructure load will no longer lead to performance issues with the CPU, but more importantly, it now allows a tenant to take full control of the CPU. So, for example, a tenant, like a bank, can bring their own hypervisor and run it on the CPU, but the IPU can still confine it to a virtual network segment or specific storage volumes. It's why we think the IPU is the control point of the data center going forward.

The second advantage of the IPU is about the infrastructure function offload. Modern applications today are often structured as microservices that incur a substantial communication overhead. From this example, from Facebook, you can see a 31 to 83% overhead for microservices. In some cases, the majority of all CPU cycles are spent on this infrastructure overhead and the IPU can help reduce this overhead. With an IPU, the cloud operator can offload these infrastructure tasks to the IPU. Thanks to the IPU's hardened accelerators, it can process these functions very efficiently. The optimized performance and the cloud operator can now rent out 100% of the CPUs to their guests, which also helps to maximize revenue.

And I know this foil is a little busy, but the third advantage is the diskless server architecture that you can enable migration from your historical/traditional data center cloud where you have disks attached to each server. As tenant demand for disk space is hard to predict, you have overprovisioning servers-- you have to overprovision servers and end up with underutilized disk capacity. With an IPU, you can move to a diskless model. All storage is on the central service as shown on the right side of this foil. When a customer starts a workload on a server, the CSP creates a virtual volume on the storage service. And here, we're showing a 250-gigabyte segment. Via the management network, the cloud service provides-- provider tells the IPU to create a new NVMe solid-state disk based on that virtual volume. As this virtual NVMe SSD shows up on the PCI Express bus, it works just like a regular SSD and works with most operating systems and hypervisors. And plus, we can boot from it.

From a performance perspective, the actual storage traffic between the storage server and the workload and the OS on the server happens on the fast path without the involvement of the host CPU cores. So, this results in low latency and high throughput with maximum flexibility, and with a key advantage of you're able to maximize your storage resources and not have to overprovision, and therefore, underutilize local resident storage services.

Intel IPU's Fundamental Role In Your Cloud Strategy

So, to recap, with the strong separation of infrastructure and tenant, and the acceleration that allows for efficient offloads of infrastructure functions, and the ability to move to a diskless architecture, we think the IPU will be a central component for the future data center architectures.

There are two common architectures for infrastructure processing units today. The first is dedicated ASIC IPU's. And second are FPGA-based IPU's. Each type has their own advantages. FPGA-based IPU's give the ability to implement new protocols quickly, react to rapidly changing requirements and implement proprietary protocols. On the other hand, dedicated ASIC IPU's maximize performance and efficiency.

But both are very different from the classic SmartNICs, which lack the capability of executing infrastructure control planes. A SmartNIC cannot act as a control point in the data center, because the infrastructure software remains on the host CPU. Whereas an IPU with cores enables full offload of the infrastructure function from the CPU and provides the CSP with necessary isolation from the tenant workload to be the control point of the system. Because there's not a one-size-fits-all for the infrastructure acceleration, Intel will continue to invest in both types of IPU's as well as SmartNICs.

Sabrina will describe our recently announced FPGA IPU and SmartNICs in more detail now. Sabrina.

Sabrina Gomez

Thank you, Brian. Hi folks, a pleasure to be here.

Before I jump into the FPGA-based offerings, I wanted to talk a little bit about our partnerships. Intel is the volume leader in the IPU market today with our engagements with the majority of tier one CSPs, where they are deploying Xeon, FPGA, and ethernet components from Intel. We have been evolving our products in partnership with these major hyperscalers like Microsoft and Baidu. We've also deployed with JD.com where they've leveraged our FPGA-based C5000X platform that is available today. And in addition, we're partnering with major software vendors like VMware on Project Monterey. So, partnerships is-- and our engagements with the key folks in this market are instrumental to our strategy.

And now, I'll share details on our most recently announced FPGA-based IPU, code-named Oak Springs Canyon.

Oak Springs Canyon expands our offering beyond the C5000X that I mentioned to provide 2x100-gig capability. It supports workloads like OVS, NVMe over Fabric, and RoCE solutions. In addition to that, with the FPGA on the card, it is customizable. This is a key feature and value of FPGAs, if you're new to FPGAs, enabling customers to achieve faster time to market and the ability to customize their solutions amidst changing standards. And we all know we're in that timeframe, where a lot of change is happening with our standards and flexibility is very important. It supports standard software frameworks like DPDK and SPDK, and it also supports Intel OFS. Intel Open FPGA Stack is a source accessible software stack for FPGAs that essentially abstracts away the FPGA, provides TOE side drivers, BMC, and other functionality to really improve the time to market and the customization of the FPGA-based solutions.

OSC leverages Intel's Agilex FPGA, together with a Xeon D SoC. Agilex is the industry's leading FPGA in the market in power, efficiency, and performance. And this is key for data centers. In concert with the Xeon-based server, this platform provides the performance necessary to offload 2x100-gig workloads. And it offers a rich software ecosystem optimized for x86. Oak Springs Canyon also features a hardened crypto block in the Agilex FPGA, and this allows you to secure all the infrastructure traffic, the storage and networking, all at line rate performance, which is another key value-add to the FPGA-based product with Xeon. Oak Springs Canyon is designed to meet the needs of the next wave CSP deployments at 100-gig.

The next FPGA-based platform I'll discuss is the Arrow Creek Acceleration Development Platform. This platform is also based on Agilex FPGA as well as our E810 100-gig ethernet controller. This is targeted to the comms service provider market. It builds upon the success of Intel's PAC N3000, which you may have heard about, which is deployed today in some of our top comms service providers worldwide, namely Rakuten and Verizon, for example.

Intel IPU's Fundamental Role In Your Cloud Strategy

Arrow Creek will help telco providers to offer the flexibility, the accelerated workloads like Juniper Contrail, OVS, and SRv6. And like I mentioned for OSC, it has that flexibility of the FPGA to enable faster time to market and customization amidst changing standards.

Arrow Creek also supports Intel OFS that enables remote update and root of trust. Its interfaces, form factor, and power envelope are designed with the telcos in mind. And as Brian mentioned, this is important in our portfolio where we have infrastructure acceleration solutions for both cloud and comms.

With these FPGA-based solutions in our portfolio, along with ASIC-based IPU that Brian is going to cover, we really have that breadth of product for the markets and the key areas.

And with that, I'll hand it back over to Brian.

Brian Niepoky

Thanks, Sabrina. I'm really excited to talk about the new product code-named Mount Evans.

Mount Evans is our first 200-gig ASIC infrastructure IPU, Infrastructure Processing Unit. We have architected and developed Mount Evans hand-in-hand with a top cloud provider and is currently under test today. This has provided tremendous insight into our development requirements for network and scale.

Also, as Sabrina mentioned, Intel has been working closely with other cloud providers through our FPGA-based solutions, and our learnings with those products has influenced many of the Mount Evans architecture and design tradeoffs. Mount Evans has been designed for performance at scale under real world workloads.

In order to be hyperscale-ready, we designed a security-- we designed in security and isolation from the ground-up throughout this chip. On the technology front, Mount Evans is loaded with innovation.

To start with, the focal point of the product is what we believe to be a best-in-class packet processing engine that supports a large number of existing use cases like vSwitch offloads and firewalls, as well as providing significant headroom for future use cases. Another technology created by extending Intel's proven high perform-- Intel's proven high performance Optane NVMe Controller enables Mount Evans to emulate NVMe devices.

A third technology innovation that I'm excited about is our Next Gen Reliable Transport protocol. We have co-innovated on this technology with our CSP partner to solve the long tail latency problems on a lossy network.

A fourth enablement technology that can be used across a large variety of use cases is our advanced crypto and compression accelerators, leveraging our high performance QuickAssist Technology.

Finally, at Intel, we really want to make IPU's a compelling technology across market segments beyond the cloud. And this, first and foremost, means enabling software developers to do what they do best. We start with innovative performance hardware designed for flexibility and ease of programmability. We add to this the expertise that came in through our Barefoot acquisition, driving the use of P4 language in the industry as a standard framework for programming network data planes onto IPU's. We will extend well known SDKs like DPDK and SPDK to take advantages of the IPU capabilities for data and storage processing.

On this foil, I'm showing a high level block diagram of Mount Evans. As you can see, Mount Evans is organized as a networking subsystem on the left and a compute subsystem on the right. I'll walk through a few of these key features for the next few minutes.

In the networking subsystem, Mount Evans supports 200-gigabyte per second of throughput, connecting up to four Xeon hosts together. We recognize that cloud performance needs will drive many applications like storage, messaging, and HPC to migrate to RDMA-based protocols. Mount Evans supports this with implementations of both RoCEv2 and the new reliable transport technology that I mentioned earlier. Our Optane derived NVMe engine exposes high performance NVMe devices to the host processor, enabling infrastructure providers to use the IPU to implement their storage protocol of choice whether it is hardware-accelerated NVMe over Fabrics, or a custom software backend on the compute subsystem.

Intel IPU's Fundamental Role In Your Cloud Strategy

The programmable packet processor delivers leadership performance for use cases like vSwitch offload, firewalls, telemetry functions all the while supporting up to 200 million packets per second performance on real world implementations.

Finally, Mount Evans provides inline IPsec to secure every packet being sent across the network while supporting up to 16 million connections.

On the right side, our compute complex is built on the Arm Neoverse architecture, using N1 Ares Cores. These 16 high frequency cores come with a large system level cache backed by three LP DDR-- DDR4 controllers. The compute complex is tightly coupled with the network subsystem allowing the NSS accelerators to use the system level cache as a last level cache, providing high bandwidth, low latency connections between the two, and enabling a flexible combination of hardware and software packet processing.

Our Lookaside crypto and compression engine is derived from Intel's QuickAssist technology that you can see in our Xeon roadmap, but we've adapted it for IPU use models. This includes the support of the Zstandard compression algorithm.

Finally, our dual-core management processor provides an interface to the platform and orchestration layers, supporting robust system manageability. As I mentioned earlier, we designed Mount Evans from a software-first mind-set. Enabling applications on IPU's requires a robust software foundation. I already shared a few of these details on using P4 language for programming network data lanes and extending well known SDKs, like DPDK and SPDK. We'll share more details about our software plans in the next few months.

A few key points to summarize. There's a major change in data center architectures occurring today to maximize resources by offloading storage, networking, and management functions to an IPU, allowing CPUs to process data, and cloud service providers to run all of those CPU cores.

As Sabrina mentioned, Intel is the volume leader at major CSPs today, and we are building upon our next generation IPU and SmartNIC products to address the challenges found in existing data center deployments.

And finally, we want to leave with you that we have existing SmartNIC and IPU platforms available today for both cloud and comms service providers. You can find out more information about these shipping products in the below link.

Thank you all for your time today.

Xiaojun Li

Great. Thank you, Brian, and Sabrina. We have some questions here.

Question one. "What is the difference between the IPU and the DPU? IPU versus SmartNIC and difference between the IPU and the foundational NIC?"

Brian Niepoky

OK, I think I can get started on-- there's a lot packed in that question, let me get started with it.

So, you know, the difference between a DPU and an IPU really is about a clarification on naming for us to get started with. First and foremost, we believe that, you know, the purpose of this device is to offload infrastructure processing, so we think that Infrastructure Processing Unit is just a better name to be able to describe what an IPU does.

Secondly, we also think that it's important that this device remain nimble. And so, we want to make sure that we have an IPU that is sized for infrastructure functions but doesn't assume any more workload, so that we leverage the CPU or the GPU or other accelerators in the platform to do what they do best. And the IPU is really about all processing infrastructure processes.

Xiaojun Li

Great.

*Intel IPU's Fundamental Role In Your Cloud Strategy***Brian Niepoky**

I think the second part of your question was the difference between an IPU and a SmartNIC—

Xiaojun Li

That's correct.

Brian Niepoky

—if I heard that correctly.

Xiaojun Li

That's correct.

Brian Niepoky

So, the main difference between an IPU and a SmartNIC is that a SmartNIC is-- works as a companion to the CPU, and the infrastructure functions are not offloaded on a SmartNIC as they are on an IPU. The infrastructure functions can be accelerated with hardened accelerators in a SmartNIC. And so, it is a companion device that works closely with the CPU.

On the converse, the IPU actually has cores that has its own operating system and its own stack, so the infrastructure functions are entirely offloaded. This provides a few benefits.

One, it provides that isolation and separation from the CPU where the tenant may rent out that space, and it could be an untrusted environment versus the IPU, where you, as a cloud service provider, entirely controls that device.

And second, it's a control point. It is where you, as the cloud service provider, can allocate and provision storage services for the CPU. You can bring up that CPU in a bare metal environment.

A SmartNIC doesn't have that capability to be a control point to the system, because it is, again, a companion device to the CPU.

Was there another part to that question?

Xiaojun Li

Yes. What's different between IPU and foundational NIC?

Brian Niepoky

Ah, OK. And so, Intel has been in the ethernet business for over 35 years, and throughout that entire time, we have developed ethernet controllers, host-based controllers that sit on a client or server system, as well as ethernet adapters. Those controllers and adapters really only have the capability of working with a driver that plugs into a host system. They really have limited acceleration features, so they're not SmartNICs and they have no capability to offload infrastructure functions because they don't have an operating stack on them.

So, it really is our standard ethernet controller and adapter business. You may be familiar with the X520 or our 700 or 800 Series at Intel that supports that.

And many customers, as far as their connectivity and networking needs, are perfectly suited and can continue to use these foundational ethernet controllers and adapters, as well as looking at both the SmartNIC and the IPU options that are available today and coming out in the market in the future.

Xiaojun Li

Intel IPU's Fundamental Role In Your Cloud Strategy

Good. Thank you. Thank you for answering those long questions.

And another one, hopefully this one is shorter, "Why does Intel have both FPGA and ASIC-based IPU's when the competition is focused on ASIC-only based IPU's?"

Sabrina Gomez

I can take a stab at that one.

Xiaojun Li

Great, thank you, Sabrina.

Sabrina Gomez

So, as we talked about a little bit, one of the key takeaways as we talk about IPU and Intel infrastructure acceleration is our breadth of portfolio. And we purposefully have that breadth of portfolio for a couple of main reasons.

One is our ability to achieve faster time to market amidst changing standards. We know based on our partnerships and early engagements with the tier one CSPs and several large telcos that that flexibility that is offered by an FPGA-based solution is key in order to enable them to be in the leading edge in the market. That coupled with our ASIC-based pieces of the portfolio really rounds it out for us.

So, in its-- it goes back to the essential cycle or value of FPGAs with ASICs. And as you know, if you've been in the FPGA space, it's historically been a cycle where new capabilities and functionalities are implemented in FPGA, and where you continue to need that customization flexibilities, the FPGA continues to be deployed, somewhat future-proofing solutions where that's needed.

Where-- when and where you need to optimize and move to a lower power, lower cost fixed solution, that's where we move into the ASIC space. So, we're really realizing those benefits and that breadth in our portfolio, and that's why we offer both FPGA-based as well as ASIC-based.

Xiaojun Li

Good. Thank you, Sabrina.

One more question. "Why does Mount Evans use Arm core? Why did not Intel choose to use Atom core?" This is an interesting question. Who will want to take it?

Brian Niepoky

Yes, so I can take that. So, Mount Evans used Arm cores because Intel looked at that design point. Intel uses a variety of architectures and IPU's in our many products, and we pick a core-- sometimes it comes down to the IPU availability, the performance, and the time to market. In this case, that Neoverse N1 core met our selection criteria. But for our FPGA-based IPU's, we use IA SoCs like our Xeon Ds, as Sabrina described earlier with both Oak Springs Canyon and Arrow Creek.

So, we'll continue to make choices based on delivering the best overall product for the customer.

Xiaojun Li

Got it. Thank you.

And another question. "Is Mount Evans P4 compatible package processing based on the Tofino chip from Intel or Barefoot?"

Brian Niepoky

Intel IPU's Fundamental Role In Your Cloud Strategy

So, we leverage the same open source P4 programmable abstraction layer with Mount Evans as we do with Tofino, but the actual design of the hardware underneath them are different.

Xiaojun Li

OK, thank you.

So, the last question currently is, "Do you have anything available today?"

Sabrina Gomez

I can start off with that one, Brian, and you can add in. So, I'll tie it back to the first question when we were talking about the portfolio and the different types. We have IPU, SmartNIC, foundational NICs.

So, the foundational NICs, we have a rich portfolio offering of standard NICs, as Brian mentioned, Intel NICs. We leverage those for various platforms. As I mentioned, they're actually leveraged on the Arrow Creek platform that we announced, the E810, so many, many options in terms of the standard Intel ethernet NICs.

We also have infrastructure acceleration SmartNICs. I talked about one, I mentioned it, the N3000. We also have another N5010 and information on those products is at intel.com/IPU. Those feature the FPGA-- Stratix FPGA or prior generations, Intel ethernet components and offer infrastructure acceleration solutions. Several of them are also qualified at OEMs and available today and deployed, it's in the work with telcos.

In terms of IPU, we have FPGA-based IPU platforms available, namely the C5000X. This is a platform that features the Stratix 10 FPGA as well as a Xeon SoC processor on-board. It is being deployed in a couple of-- next wave CSPs, and we leverage our ecosystem to deploy that. We also offer ecosystem partners that will help to customize that platform based on what the specific end CSP or CoSP needs. And again, information on that platform is available on intel.com/IPU.

Xiaojun Li

Great, great, thank you. OK, that's it for the question and answering.

And thanks, Brian, and Sabrina. Thank you all for joining us today. Please do not forget to give our team a rating for the live recording, so we may consistently improve the quality of our webinars.

Thank you, again, for joining us today. This concludes our webcast. Thank you.