*Image Credit: David MacKenzie*

# Windows Server on ARM is for Azure only

What makes the work of porting code to the ARM instruction set worthwhile is both the scale of cloud and the kind of services Microsoft runs in Azure. The economics works for cloud, not your data centers.

**Mary Branscombe | MAR 21, 2017**

A couple of months after announcing that Windows 10 will be available on ARM laptops based on Qualcomm's Snapdragon 835 when the Creators Update launches this spring, Microsoft also revealed that it has been trying out ARM servers in its Azure data centers and is planning to use them for some very specific production workloads.

That doesn't mean you should expect Microsoft to be shipping Windows Server for ARM any time soon — if at all. "We're optimizing [ARM servers] for data center use, not just for the generation coming out later this year but also the ones beyond and beyond that," Microsoft Distinguished Engineer Leendert van Doorn tells CIO.com.

Yes, Azure runs on Windows Server and the Windows kernel has long supported multiple processor architectures. And, yes, Microsoft has an ARM version of Windows 10. And, yes, Windows 10 and Windows Server are both built from the same source code. But to do anything useful, Windows Server also needs language runtimes, middleware and applications.

*[ Inside a hyperscale data center (how different is it?) ]*

Microsoft has ported only enough of Windows Server and those associated software components to run the workloads for which ARM is a good fit. While he showed that code at the Open Compute Platform Summit earlier this month and confirmed that several Microsoft cloud services will use ARM servers, van Doorn also made it clear that this ARM version of Windows Server is for internal use only.

He didn't rule out, say, a future version of Azure Stack that could take advantage of ARM servers, but the issue is less about technology and more about economics. "The enterprise market isn't growing, and disrupting a stabilized market is difficult. Adopting new things in a growing market like the cloud … that's where you can introduce new technologies, that's where you can disrupt. I could see at some point that things would flow over into enterprise; there's no technical reason why they can't, but it's really about market dynamics."

## ARM is ready for the cloud

The advantage of ARM is that you can have a larger number of cores with a higher number of execution threads at a lower price and with much lower power usage than an Intel Xeon. That's certainly appealing for cloud services like Azure, which are limited more by the amount of power they can get into their data centers than by the numbers of servers they can pack into the building.

But reducing power usage isn't the main motivation here, says van Doorn. "Power is becoming more or less a wash from my perspective. In the end, I'm interested in how can I get the most performance, how can I get the most value out of my systems."

> *"The enterprise market isn't growing, and disrupting a stabilized market is difficult. Adopting new things in a growing market like the cloud … that's where you can introduce new technologies, that's where you can disrupt."*
> *— Leendert van Doorn*

Until now, ARM wasn't powerful enough. "You had a lot of processors, you had a lot of cores, but the performance per thread was really low. That is changing with these products coming out

this year. We're seeing really very high-end per-thread performance and we're still seeing lots of threads," says van Doorn.

Not only does this new generation of ARM servers have the necessary performance, but thanks to the fast pace of smartphone development they're part of a nimble, competitive ecosystem. "We're seeing interesting connectivity options [with ARM], especially around new bus standards. We're seeing interesting integration options," van Doorn says. "We see multiple partners with high-end processors competing with each other. We like that because it drives innovation."

## The ARM advantage

Moore's Law, which has driven performance improvements for the x86 architecture for so long, has always been about both increasing the density and lowering the cost of transistors. That trend is slowing down, both technically and economically. ARM will face that slowdown itself, but that's further down the line.

ARM also has a potential advantage over the familiar x86 architecture because it already uses multiple different instruction sets, developed as the chips have evolved. By contrast, when AMD developed 64-bit x86 chips it opted for compatibility, which made it easier for developers to move from 32-bit to 64-bit, but means that creating another instruction set would be disruptive and likely inefficient. Even when you use a 64-bit x86 chip to run 32-bit software, it's using the same processor instructions, just with a prefix telling the chip whether to run them in 32-bit or 64-bit address mode.

"If you look at ARM's 16-bit and 32-bit and 64-bit instruction set architectures," van Doorn explains, "they're completely different sets of instructions. There is no overlap. So for ARM, it's really easy to create a new instruction set without breaking all the other instruction sets. With x86 that's much harder to do, because you've got one big instruction set, so you can't do fundamental shifts in how you program processors."

"We feel the ARM instruction set architecture lends itself more to an evolution than any of the other instruction sets we're working with," van Doorn says. That, coupled with innovations in the ARM ecosystem, makes it a good choice for what Microsoft can do with servers now, and in the future — and it wasn't a difficult job to get Windows Server running.

"To bootstrap an ecosystem, it helps if you have another ecosystem that's developing all the tools and software for it," van Doorn says. "x86 server really benefited from x86 on PCs. In ARM That same thing exists on the mobile side. At Microsoft, we've greatly benefited from that ecosystem in porting Windows Server, because we've got a unified source code base where both client and server versions are built out of the same source tree, so doing a Windows Server version was relatively straightforward."

## ARM is good for PaaS not IaaS

Before you get enthusiastic about the idea of running Windows Server workloads on your own cheaper, lower power servers, remember that will mean either porting all of your applications

and code to a new instruction architecture or running them in emulation and losing the performance benefits of those cores.

It isn't that ARM couldn't do virtualization, but that it's not the best tool for the job. "We already have a very good vehicle to run x86 workloads, namely x86 processors. So ARM doesn't provide any benefit, and if we do want to run x86 workloads on ARM we have to go and binary translate it or emulate it, and all the performance is gone," says van Doorn. "Where ARM does shine is platform-as-a-service (PaaS) workloads."

*[ To the cloud! Real-world container migrations ]*

What makes the work of porting code to the ARM instruction set worthwhile for Azure is both the scale of cloud and the kind of cloud services Microsoft runs in Azure. Those ARM servers aren't going to power infrastructure services that let you run virtual machines (VMs); they're going to run platform services like search and indexing, storage, databases, big data and machine learning.

Replacing an expensive Xeon chip that's overkill on a storage node reduces costs, which is why you often find ARM chips powering network-attached storage (NAS) hardware. But the real win is in the heavily parallel processing needed by services like SQL Azure, HDInsight and Azure Data Lake, Azure IoT Suite and Cortana Analytics, Azure Stream Analytics, DocumentDB, Azure ML and Microsoft Cognitive Services, Azure Search and Azure CDN. All those cores and threads make ARM servers ideal for efficiently crunching through multitudes of relatively small pieces of training data for machine learning.

Running these PaaS workloads on ARM means when you 'use' ARM servers on Azure, you're never going to know you're on ARM. And if Azure Stack does ever bring ARM servers to your data center, or the evolution of instruction sets makes ARM servers the hardware choice of the future, you're still not going to be using it for infrastructure-as-a-service (IaaS) because IaaS is for legacy applications. If you ever buy them, you're going to use ARM servers for PaaS.

## Data centers get diverse

It's not new for businesses to have a mix of hardware in the data center, even though it hasn't been a recent trend. Instead of fading away like RISC servers, mainframes have kept their niche — especially in finance. In fact, market intelligence firm IDC says mainframes are becoming increasingly connected to other enterprise systems, using web APIs and Java and Linux.

AMD's return to the data center with its new 32-core Naples server chip may blur the line between CPU and system on a chip (SoC), but even if you directly attach GPUs for running virtual desktops or machine learning workloads, you're still dealing with the familiar x86 architecture.

But there are other options, and not just ARM.

The OpenPOWER Foundation — a technology alliance IBM created to try and break into the Intel-dominated data center market — is hoping that businesses will accept more diverse

systems. Google is a founding member, and the POWER architecture will be in Google data centers this year. But POWER will sit alongside Intel's new Skylake Xeon processors, NVidia and AMD Radeon graphics processing units (GPU), ARM servers and Google's own custom tensor processing unit (TPU) ASICs.

The operational complexity of that kind of heterogeneous environment is different from the usual hyperscale cloud approach, which uses as much identical infrastructure as possible for economies of scale in both purchasing and management. But it's driven by the same customer (and internal) demand that has Microsoft introducing GPU-based virtual machines and field-programmable gate array (FPGA) boards and now ARM servers to speed up search and AI workloads. Similarly, Facebook is putting Tesla GPUs into its own Open Compute Project (OCP) designs for machine learning servers.

Some of the reason for the heterogeneous hardware approach is picking the best architecture for specific workloads, like machine learning. "When I look at the spectrum of all these servers, it's turning into 'how can I match the right hardware to my workloads'," says van Doorn. "That's an interesting new thing. It always used to be that you had one kind of processor and you needed to take your workload and match it to that. The scale of the cloud is the other side of this; it makes economic sense to optimize a particular workload and you can now optimize hardware for that, and that's very exciting."

*"As we scale the public cloud, you have to have the services and the hardware so you can run any of the workloads that today run in the enterprise."*
*— Kushagra Vaid*

That hardware optimization applies to more than just ARM servers. "It isn't that ARM is going to be the winner or GPU is going to be the winner. They all have their own specific application domains where they shine — and they're all interoperable, so ARM can have GPUs, it will have FPGAs. We build our own smart network interface controller (NIC) with an FPGA on it that we use for acceleration in our data centers, and that fits in ARM servers."

Cloud needs a mix of hardware because it runs so many different workloads, points out Kushagra Vaid, the general manager for Azure hardware infrastructure. "As we scale the public cloud, you have to have the services and the hardware so you can run any of the workloads that today run in the enterprise. Those enterprise workloads are diverse; they can be something simple like a line-of-business app, to something complicated like high-performance computing (HPC) or machine learning. Now we have to do all that in one cloud, so the hardware has to be there to run these different workloads."

"It doesn't necessarily mean we'll have massive amounts of heterogeneity; it means we have to find what is the right hardware design that can do the job most efficiently for a given workload," says Vaid. But he also notes that workloads are getting more diverse rather than less. "Just five years ago, x86 could pretty much handle all kinds of workloads except maybe HPC so people did

CUDA, but between GPU and x86 you could pretty much handle everything. Now there's been an explosion of these next generation workloads. Machine learning is the poster child, and machine learning and AI do not map that well even to GPUs. That's why there are so many machine learning startups trying to solve the problem more efficiently."

## Standards and supply chains

It's also worth remembering that using ARM servers gives the large cloud builders a way of fostering competition among vendors. As the dominant supplier with something of a stranglehold on the data center market, Intel hasn't had much incentive to lower prices.

Unlike Intel, ARM doesn't make and sell chips; it licences the instructions to create and optimize those central processing units (CPU) to multiple companies. There are numerous ARM server vendors, and their rivalry will help drive technical improvements as well as keep prices competitive, although van Doorn notes that driving innovation is more important than pushing down prices because Microsoft can get cost savings "by being good at workloads."

Having multiple suppliers also avoids the worry of depending on a single vendor to supply an entire cloud. Microsoft is working with several ARM licensees, including Qualcomm and Cavium. "There are clearly supply chain motivations here but they transcend cost," he says. "It's really about 'Can we get the right products; can we build out the right volume we need?' If any of our partners have a hiccup, what implications does it have for us?"

Handling those questions demands standardization. Some of that depends on the Server Base System Architecture (SBSA) ARM introduced in 2014. "SBSA is really about trying to reduce the fragmentation you see in the ARM space. We leverage that a lot. As a result of SBSA, we have the same Windows Server version running on Cavium boxes and Qualcomm boxes and there is no difference between the two," van Doorn says.

In fact, designing ARM servers to work with Windows is actually helping to drive standardization further, he says. Vendors can modify Linux to make it work with their hardware in ways they can't modify Windows Server. "This is enforcing standardization and unification in a way I don't see happening in the Linux ecosystem, because they're internally fragmented." Microsoft is also working with ARM to make system building easier (in part, using lessons it learned building Surface and Windows Phone). "Going forward, buses will be enumerable; we will be able to figure out what devices are in a system."

Microsoft is also able to avoid a lot of the usual complexity of having a mix of hardware because the Qualcomm and Cavium motherboards fit in the Project Olympus server design Microsoft announced last year. This is a universal design where you can drop in an Intel, AMD or ARM motherboard, as well as GPU and FPGA accelerators.

"That reduces the friction of taking some of these new technologies into our data centers," van Doorn points out. "It's one thing to get a completely new instruction set architecture like ARM into the data center; it's another thing to figure out what the physical box goes into and to make sure it complies with all our data center standards. Olympus makes it so we don't have to worry about the power supply or the out of band management or the network connectivity or even

what chassis. This whole LEGO block approach really helps us adopt, deploy and monetize innovation earlier."

Instead of waiting to launch a new service until the engineers can integrate new hardware, Azure can put that hardware into the same rack and plug in the same power supply and provision and monitor it with the same tools as all the other servers.

All those advantages are enough to outweigh the work of moving to ARM, at least for Microsoft. Putting that mix of architectures into your own data centers is rather less likely though.

Using PaaS services lets businesses get the advantages of multiple architectures without the headaches of managing and integrating them, and that's likely to be the way most organizations buy into the ARM architecture.

---

*Mary Branscombe is a freelance journalist who has been covering technology for over two decades and has written about everything from programming languages, early versions of Windows and Office and the arrival of the web to consumer gadgets and home entertainment.*

*"Windows Server on ARM is for Azure only" originally appeared on CIO.com*