

Cloud or Colo? The point when Al Startups need to make the switch

by Timothy Prickett Morgan

STHENEXTPLATFORM

Executive Summary

It is much more economical for AI startups running at scale to switch most or all workloads from cloud to a colocation provider. As a rule of thumb, 50 per cent continuous workload utilisation is the cut-over point to make the switch. Before renting one GPU instance on a cloud, AI startups must consider the evolution of computational needs as AI applications move through stages from prototype through production. By getting this right, the startup will avoid cloud lock-in and will keep control of costs when running at production scale.

Introduction

An AI startup, like any other startup, builds a business plan that evaluates a market and delivers what it needs. There are simple and relatively inexpensive ways to get started on the public cloud, with IT infrastructure loaded up with GPU or FPGA accelerators, or hefty CPUs that can do a respectable amount of floating point and integer math that machine learning training and inference workloads require.

Today, any startup can go to one of the big public clouds, which sell AI as-a-Service (AIaaS) for training and inference, in some cases using their homegrown ASICs or AI engines fashioned from FPGA logic. The cloud giants have plenty of raw CPU compute on hand in every imaginable configuration of processor and memory, flash, and block storage, and with fast Ethernet and sometimes InfiniBand network interconnects.

With this capacity on demand, all you need is a credit card, a dream, and an algorithm to test. That's the secret of the general compute and AI-specific platforms that Amazon Web Services (AWS), Microsoft, Google, Baidu, Alibaba, and Tencent have created. It is easy to get started, and even easier to grow capacity to what appears to be – to an AI startup, at least – always-available infinite capacity. However, as we all know, this is an illusion that the cloud giants work hard to maintain. There is no such thing as perfectly fungible, infinite capacity.

In many cases, AI startups have no plans to ever invest in their own IT infrastructure. In the 21st century, the idea seems absurd. Yet in a lot of cases, it makes more sense to build out infrastructure and utilise it more heavily. Budgetary and technical constraints are endemic in public cloud infrastructure, and clients may suffer – at their own expense.

There is some hesitancy in even considering on-premises infrastructure, and understandably so. Noone at an AI startup will be denied funding for running their applications on a major public cloud (just as no-one at a web startup ever got fired for buying an Oracle relational database and big backend Sun Microsystems servers to run it on). To be sure, some employees at AI startups may have worked at research, government, or academic laboratories that bought or built their own distributed computing systems to run HPC and AI workloads. But that can often be a negative experience because of the budgetary woes of building a physical data centre – even a relatively small one – in urban areas where most global research is undertaken.

At some point when the prototype code is done, AI startups will realise they can't run their workloads on a powerful laptop or even a GPU-laden workstation and move to the cloud. It doesn't take long before the bills start piling up for renting the compute and storage from the public cloud. Once that happens, the fear of data gravity (the sheer size of datasets that are used to support machine learning training) and data egress (the astronomical cost and the very slow speed of moving data off a public cloud) finally sinks in.

Maybe for fast growing AI startups, it makes sense to never go "all-in" on the cloud in the first place. There is, of course, a third option that an AI startup can consider – and that is to own or lease their own infrastructure. Better still, they could rent it using a cloud-like pricing model such as HPE GreenLake or Dell's APEX and locate it in a nearby colocation facility, for reduced latency and proximity. This gives the AI startup some of the operational benefits of a cloud data centre and also the locality and control benefits of the on-premises facility.

The Kao Data Campus

This is certainly what Kao Data believes. Founded in 2014, the company runs high performance colocation data centres that cater for AI startups and research organisations based in the UK's "Innovation Corridor" stretching between London and Cambridge. Building a data centre is a serious business, and to that end Kao Data got its initial £33 million funding in February 2017, building a state-of-the-art facility in Harlow, just north of London. (Interestingly, the campus is located on the same site where Nobel Prize winner Sir Charles Kao invented the fibre optic cable while working in the labs at Standard Telephones and Cables, a UK telecoms pioneer.)

About a year after the initial funding, the company opened its first 8.8 megawatt data centre, named Kao Data London One (KDL1), on the 15-acre site. When fully operational, KDL1 will have four data centres running at a combined 40 megawatts across 150,000 square feet – all supported by renewable energy and driving a power usage effectiveness (PUE) of 1.2, even at partial loads.

Kao Data's industrial-scale facilities have a column-less design, allowing for the tuning of racks and rows, which can accommodate specific power densities and network topologies. Further, it uses environmentally friendly indirect evaporative cooling to support its customers sustainably. Its data centres are architected explicitly to house OCP-Ready[™] servers, storage, and switches based on Open Compute Project (OCP) designs. However, companies can house any racks they want on the concrete slab floors.

KDL1 houses one of the most important and famous new intensive computing systems in the world, the Cambridge-1 supercomputer. This system is at the heart of the quantum computing, digital biology and AI healthcare research that NVIDIA, Arm and Cambridge-1's founding partners will be undertaking. The data centre is also DGX-Ready certified to support NVIDIA SuperPOD and smaller DGX system configurations for other clients, and Megaport enabled, providing the perfect platform for rapid, hybrid computing between heavy HPC iron in the data centre and public cloud infrastructure.



NVIDIA Cambridge-1, located at Kao Data's Harlow campus

KDL1 has a Technology Cell as a base unit of consumption, which has capacity for 34 racks of compute, storage, and networking, hot aisle containment, and 350 kilowatts of power. The Technology Suites come in two sizes, offering either 1,100 kilowatts and 4,500 square feet of space or 2,200 kilowatts and 9,000 square feet of space, which top out at 221 racks or 442 racks, respectively. This is about the same rack capacity of large, industrial-scale deployments, whether they are GPU-heavy AI or HPC systems. The OCP racks, which stand at 58U high, obviously pack more of a punch than a standard 42U rack. From a security perspective, it offers cage services, closed circuit TV monitoring, and dedicated biometric verification locks can be added to the cells and suites for security above and beyond the data centre perimeter.

What is good for NVIDIA and Arm is also good for all manner of AI startups who might be breaking the budget using public cloud. For example, InstaDeep, the first new customer for Kao Data in 2021, runs its hybrid CPU-GPU system in KDL1. The company has created a protein design platform called DeepChain that uses machine learning and reinforcement learning to automatically validate the efficacy of simulated proteins without human interaction. Other startups are lining up after becoming disenchanted with the cost of public cloud infrastructure, and some will start out from the get-go with their own production systems in KDL1, after first building their AI and business models on cloudy infrastructure. Others, however, just want to move faster than the public clouds can or ever will.

Contrary to popular opinion, the major public clouds are not on the cutting edge of all technologies. In reality, they often get access to the latest CPUs a quarter or so ahead of the broader market, but they are not always on the front of the line with GPU and FPGA technologies or with switching technology. They tend to hang back a bit, even if they are ahead of the mainstreaming across large enterprises where a lot of the volume in the market resides.

Given the rapid change in the AI space these days, having access to the latest and greatest technologies is expensive, but it is often quite easy to justify the investment over the N-1 or N-2 technology that is commonly available on the large public clouds. If you can get 10X the performance for 2X or 3X the cost, you try to do that if you can as a startup because that investment pays for itself. Not every startup does this math correctly but understanding the AI application lifecycle and how capacity demands are likely to grow is vital for success.

Doing the maths on Colo versus Cloud for Al

While Moore's Law is still working for GPU and FPGA accelerators, the leaps in transistor density are stretching out in time and that means the pace of the lowering of the cost of transistors has slowed. This is beginning to translate into a rising cost of creating increasingly denser semiconductor devices, at a time when AI model complexity is exploding on an exponential scale, as NVIDIA demonstrated recently:



The Evolution of AI Model Complexity

As machine learning models become more complex, the number and size of datasets also continues to grow, and this compounds the amount of computing that companies must budget for. As AI becomes integrated into some applications, it is reasonable to expect that all applications will have some element of AI associated with them, too, which compounds the growth even further.

The public clouds are, therefore, great for running back office and Web infrastructure applications, and the economics often work out. However, with high performance computing systems with GPU acceleration, the costs of renting such capacity can be prohibitively high. Also, Web-scale infrastructure in essence comprises virtual compute, virtual storage, and virtual networking scattered around a 100,000-server data centre or multiple data centres in a cloud region. This is a sub-optimal platform for running AI applications that simply cannot deliver the same deterministic genuine performance that HPC customers expect and need.

To reinforce this point, NVIDIA houses its Selene supercomputer in a data centre in Santa Clara, next to its headquarters, while the Cambridge-1 supercomputer is in Kao Data's facility in Harlow. NVIDIA could have just run its applications in the public cloud but chooses not to. Similarly, Google doesn't actually run its machine learning and analytics workloads on the Google Cloud on top of Kubernetes containers, but rather on internal gear that uses its predecessor Borg/Omega container and cluster controller. Amazon and Microsoft are moving more and more of their infrastructure to the public cloud, but it is debatable if they will do so on self-contained, isolated instances of their infrastructure that are closer to AWS Outposts or Azure Stack, rather than the raw AWS or Azure infrastructure.

Even the public clouds don't want to deal...

There is something to be said for dedicated infrastructure that has cloud-style management, flexibility, and pricing, but literally sharing that infrastructure with unpredictable noisy neighbours will present problems – in a multi-tenant cloud computing environment, demand often exceeds supply, sometimes because one tenant has monopolised resources. The potential for performance degradation increases as a startup begins to scale up their operations on the hockey stick curve that everyone hopes to ride when they start a new company.

To be sure, getting free instances from the public clouds is economically attractive at the front end, and AI startups are smart to take advantage of them. Spot pricing for capacity is significantly less expensive than reserved capacity, which is again less expensive than precious on-demand capacity and is great for early proof of concept and prototyping work. But at some point, picking a specific cloud provider to run AI applications using its tooling, its compute, data, and networking services represents as much of a vendor lock-in as buying a proprietary minicomputer or mainframe platform did back in the 1980s and 1990s. And with data egress charges – the cost of moving data off the clouds – being astronomical, moving large datasets can be slow in time and high in cost. Fast moving and fluid AI startups can really get stuck; and pitting two public clouds against each other to drive down some of the costs does not change the fact that they are stuck.

Admittedly, it can be tricky to figure out the inflection point where colocation makes more sense than cloud capacity, as the calculations to figure out when to do what with AI infrastructure is different for every organisation at every point in the evolution of their workloads. That said, some general principles prevail.

Before renting one GPU instance on a cloud, AI startups must consider the evolution of computational needs as AI applications move through stages from prototype through production. They must recognise the scale issues before they get buried by them, and in general, it looks like this:



In the beginning, whether they buy a workstation, place cheap GPU graphics cards in existing servers, or rent GPU-accelerated instances on the public cloud, AI startups pay for the capacity they use, which in many cases, isn't very much. It's easy and not particularly economically painful. They take years and years of data, spend months or weeks cleaning it up so it can be used to train AI models, and then spend weeks or days training their models.

You can measure the capacity needs on any number of metrics and the graph above stays the same, more or less. Megaflops of floating point capacity, megawatts of power consumption and heat dissipation, or the number of times the complete dataset is processed through training, which is called an epoch. Capacity needs grow over time, as the AI application moves from prototype to product development to production. Dealing with this scale ramp is tough enough, without having to deal with finding the lowest common denominator for public cloud services and trying not to use APIs and services (particularly storage) that will render the organisation effectively locked in.

There is, however, a more reasonable way to do this, and that is to use a mix or hybrid of on-premises infrastructure (perhaps best deployed in a colocation facility since startups neither have nor want their own data centres) and cloud services. In general, AI startups will get better performance from a cluster of machines that are in the same physical place and that do not have to share resources with noisy neighbours that most definitely impact performance. When an AI workload is in actual production and is meant to drive revenue in some fashion, predictable performance, as well as better performance, both become crucial. The fact is that it is very hard to get either on the cloud at a reasonable price and all the time.

Both the performance and the pricing are unpredictable since AI startups may have to do instance substitutions within a cloud or across them, depending on the availability of particular instance capacity with the public clouds. In short, once you get entangled in the web of a public cloud, it is hard to get out. It is better, perhaps, to not get tangled up in the first place.

As a rule of thumb, anytime you are using cloud capacity at scale more than 50 per cent of any appreciable long term – say a year or two or three – even with three-year reserved instances and other methods that the public clouds use to cut prices, it is far more economical to have a base production load running on iron you own (or rent or lease from an OEM), and then use multiple clouds at a highly abstracted level to deal with production peaks. Like this:



In this scenario, the AI startup can buy the time to increase on-premises production capacity, simply by bursting into the public cloud. But more importantly, this environment does not need to change as OEM or cloud vendors come and go. The trick is to avoid tying into anyone or any plan. And if conditions change, the AI startups will be able to adapt by increasing spending on-premises (in this case, meaning the need to run infrastructure in a colo facility) or on the cloud as they see fit at any given moment.

So, just how big is the pricing gap between renting AI hardware on the public cloud and buying machinery to host within a colocation facility? Let's take a DGX-1 from NVIDIA and an eight GPU instance running on AWS as a baseline.

TCO of Cloud versus Colo

For the colocation facility, the cost of a DGX-1 machine and its storage plus switching is on the order of \$238,372. If you round that up and depreciate it using a straight line method over two years, that's \$10,000 a month. Then, add in 10 kilowatts of power and colocation rent, and that is another \$2,000 a month or so.

On AWS, a DGX-1 equivalent instance, the p3dn.24xlarge, costs \$273,470 per year on-demand and \$160,308 on a one-year reserved instance contract. Comparably, Microsoft Azure charges about 30 per cent less for an equivalent instance, but AWS is the touchstone in the public cloud. Add in AWS storage services to drive the AI workloads, and it is around a cool \$1 million to rent this capacity for two years.

Assuming 80 per cent utilisation – meaning a lot of work is going through these systems – here is how on-demand, reserved, and colo costs stack up for two-year and three-year scenarios:

AWS Cloud vs Colo approximate comparison framework, 2-year costs					
Utilisation (%)	On-demand Cloud	1-yr Reserve Cloud	Colo/In-house		
80	\$800,000	\$586,200	\$288,000		

AWS Cloud vs Colo approximate comparison framework, 3-year costs					
Utilisation (%)	On-demand Cloud	3-yr Reserve Cloud	Colo/In-house		
80	\$1,200,000	\$463,194	\$312,000		

Utilisation rates matter in scenarios comparing the cloud to colocation, of course. You pay for the colo whether you are using it or not. So, let's scale that cloud utilisation up, from a low of 10 per cent to a high of 100 per cent, to show how on-demand prices look small at low utilisation but can quickly skyrocket as utilisation climbs. At full utilisation on the cloud – meaning it is running all the time, non-stop – the prices are very high compared to colocation over the two-year and three-year terms shown below:

Cloud vs Colo approximate comparison framework, 2-year costs					
Utilisation (%)	On-demand Cloud	1-yr Reserve Cloud	Colo/In-house		
10	\$100,000 = \$1,000,000 / 10	\$586,200 = \$1,000,000 *18.30/31.218	\$288,000 = 24* \$12,000		
25	\$250,000	\$586,200	\$288,000		
50	\$500,000	\$586,200	\$288,000		
80	\$800,000	\$586,200	\$288,000		
100	\$1,000,000	\$586,200	\$288,000		

AWS Cloud vs Colo approximate comparison framework, 3-year costs				
Utilisation (%)	On-demand Cloud	1-yr Reserve Cloud	Colo/In-house	
10	\$150,000 = \$1,500,000 / 10	\$463,194 = \$1,500,000 *9.64/31.218	\$312,000 = 36* \$8,666	
25	\$375,000	\$463,194	\$312,000	
50	\$750,000	\$463,194	\$312,000	
80	\$1,200,000	\$463,194	\$312,000	
100	\$1,500,000	\$463,194	\$312,000	

The cells in these tables in green show the best choices, and the cells in yellow show the choice to make when capital outlay has to be avoided and takes precedence over having to spend more on the operational budget.

While these numbers are illustrative of how much colo is, versus the cloud, rather than a hard IT budget, the conclusions are inescapable. Many AI startups will unfortunately learn this the hard way – as did many startups during the Internet 2.0 era with more generic public clouds. What's clear is that you don't need artificial intelligence to see what the right strategy is; human intelligence will do.



SPONSORED BY:

About Kao Data

Founded in 2014, <u>Kao Data</u> develops and operates advanced data centres for high performance colocation. From our hyperscale inspired campus in the heart of the UK Innovation Corridor between London and Cambridge - we provide cloud, HPC, AI and enterprise customers with a world-class home for their compute.

Our Harlow campus - built on the site of Sir Charles Kao's pioneering discovery of fibre optic cable in 1966 - is a development of four state-of-the-art, OCP-Ready[™], carrier neutral data centres. When fully completed the 15 acre, £230m-plus campus will support an ITE load of over 40MW, across 150,000sq ft of technical space – all powered by 100 per cent renewable energy. Backed by Legal & General and Goldacre - Noé Group, Kao Data is one of the largest campus developments in the UK and represents the future in sustainable, efficient and scalable computing - providing an industry blueprint to develop further best-in-class data centres.

kaodata.com