# Observability for AI Factories –
## What's Missing from This AI Reference Architecture?

By **Iain Kenney**, Senior Director and Head of Product Management

It is great that NVIDIA has partnered up with a range of companies listed in the **blog** by Bob Pette to help enterprises build their own AI Factory. Especially as so many Enterprises are seeking to bring next-generation compute back under their control including the rapidly expanding AI workloads.

AI Factories/Workloads bring a new level of requirements in terms of compute and ultimately networking. With NVIDIA providing these Enterprise RA's as they are called it should help guide the architects to define the final designs based on the unique needs of these enterprise organizations.

Several key aspects are missing for each of the Enterprise RAs posted online by the various partners.

According to a **blog** by Andy Patrizio commenting on the original NIVIA blog, "The one thing that the reference architecture does not cover is storage…"

Well, that is true and certainly, an area I am sure is being addressed by people like me at the various storage vendors listed.

However, truly I think the biggest part that is missed completely by NVIDIA and all their Enterprise RA's that I read through is Observability.

This new class and design of Enterprise AI data centers bring challenges in computing through networking speed and burstiness of those networks, they bring challenges in terms of power required to run all these machines and GPU's and next-generation processors, etc. Even if all these issues (and 100 others) can be solved in theory and initially bring up…. how will the NetOps teams running these systems see what's going on during normal operation?

Now I am sure that you may be thinking… well I am sure the secret sauce management software is going to be log files and application tracing etc. Yes, I am quite sure that is true but as any NetOps team worth their salt knows, the actual Packets are the single source of truth for any network. Especially one where the team builds the train while learning to drive the train and specifically making that train travel at light speed because management said it should! 🙂

cPacket

So maybe you are a network architect who wants to implement the latest generation data center based on one of these Enterprise RAs from the NVIDIA partner vendor that you have worked with for many years. That's great and I am sure you will get a lot of support from that vendor and if you are big enough from NVIDIA itself.  However, don't forget about Observability.

- You need to see what is going on in the microbursts that are critical to the operation of this new cluster approach.

- You need to be able to spot the packet loss at that critical moment in your daily operation.

- You need to be able to capture packets all the time to make sure that you can recreate any issues that you notice from within the cluster.

Etc..etc..etc…  So don't wait until you have an issue to try and extend the budget to add the taps, brokers, and capture devices… you already know how that conversation with your budget approvals team is going to go.

So why not get ahead of these questions and build a practical and effective observability approach into your budget before you start deploying?  It should be able to tackle the issues above and will make you look like the hero when you start deployment and go live.

cPacket can help you refine your approach and augment your Enterprise RA to give your organization the leg up in this race to deploy next-generation AI Factory.

## About cPacket Networks

cPacket Networks de-risks IT I&O through network-aware service and security assurance across hybrid and multi-cloud environments. Our AIOps-ready Intelligent Observability Platform provides single-pane-of-glass analytics and deep network visibility required for complex IT environments enabling Fortune 500 organizations around the world to keep their business running. cPacket solutions are fully reliable, tightly integrated, and consistently simple. Our cutting-edge technology enables network, application, and security teams to proactively identify issues before negatively impacting the business. The result: increased service agility, enhanced experience assurance, and faster transactional velocity. Learn more at cpacket.com.

BL-020725