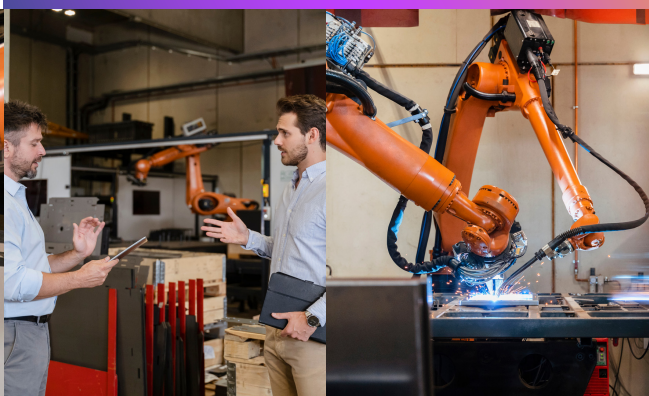


Feel confidently competitive in an AI world



AI model training can take weeks or months; here's how to avoid this common bottleneck.

Artificial intelligence (AI) is increasingly touted as a solution to any computing problem, and the wide availability of consumer-grade generative AI tools has only furthered that reputation. But enterprises are learning that under the glossy surface, getting an AI model into production isn't quite so easy, and a large part of that problem relates to data cleaning, model training, and tuning. How do you improve the training process and get models into production more quickly and accurately, so you can confidently innovate with AI across the organization?

You can start by looking at the usual suspects.

"Cost is one of the biggest complaints," says Rangan Sukumar, a distinguished technologist at Hewlett Packard Enterprise. "Training a large language model from scratch can easily run millions of dollars in compute costs alone.¹ And then once you have trained the model, there are more costs associated with every additional query. When you have millions of users prompting the model, you're dealing with a massive cost issue, not to mention sustainability: An AI model such as GPT-3 is estimated to have consumed 1300 megawatt-hours of electricity² — the equivalent of 1450 U.S. households for a month. So, when it comes to model training, the way we've been doing things so far is just not sustainable."

¹ ["What Large Models Cost You — There Is No Free AI Lunch,"](#) Forbes, Sep 8, 2023

² ["New tools are available to help reduce the energy that AI models devour,"](#) MIT, Oct 5, 2023

The trouble with training

Why is model training so cost-ineffective? Sukumar says the biggest problem facing organizations today is the experimentation required to validate the business impact of the vast array of AI algorithms, models, and infrastructure. “You’ve got hundreds of AI vendors trying to sell you software, services, tools, and models without enough clarity on when and what is relevant to the business,” he says.

To help developers find the shortest path to leveraging AI models, AI vendors need to offer different product experiences for users at various levels of expertise and various stages of the AI journey — but that has been elusive to date, Sukumar says. The training experience for a first-time AI user can and should look completely different from that for a seasoned AI veteran pioneering a new model, he adds, but many users are being treated like pros out of the gate.

Complicating matters is the issue of scale and the uncertainty that surrounds model training. “We have simple models small enough that they can be fit and trained inside the memory of a single GPU, and we have models that go all the way up to requiring thousands of GPUs,” Sukumar notes. This introduces the need for data, model, and pipeline parallelism. Training jobs in these cases must split training data into pieces and partition models across multiple GPUs while coordinating the simultaneous exchange of those pieces among multiple devices, all of which must constantly update one another so they stay on the same page. This is a complex discipline that requires intense levels of skill and, at times, luck — not to mention epic levels of compute and storage power.

This problem is compounded when organizations don’t initially understand the AI model they are trying to train. “Almost everybody thinks they can find an open-source model that’s out in the public, download the model, train it on their data, and it’s going to work like magic,” says Sukumar. “That’s seldom the case. Sometimes the repurposing of a trained model to a specific use case is seamless, but often, if you are trying to make an open-source model work with your sensitive business data, you realize that you must change the model architecture entirely.”

Sukumar recalls a published example from 2020 that involved an AI computer vision model that had been trained



on medical imaging data drawn from one hospital. The model worked fine for the hospital that trained it, but when other hospitals tried to use the system, the model would produce an improper diagnosis. The problem stemmed from the calibration differences of CT instruments across hospitals.³

“With the slight change in contrast, the entire algorithm got confused by the new data that it had not seen before,” says Sukumar. “These kinds of issues are everywhere. People start off with assumptions about the data and hope that the AI will just work. It seldom does, so you’re frequently having to go back and fix things.”

Another foundational issue organizations have been wrestling with around AI training is the age-old problem of talent. “The kind of expertise that can build a model from scratch is very rare and very expensive,” says Sukumar. “It probably costs more to hire that kind of expertise than it costs you to train the model using GPUs. Building a model from scratch is extremely hard. To get around this, many organizations are pursuing a strategy of transfer learning and fine-tuning, where an existing, known successful model is adopted and then retrained using new data.”

³ “Common pitfalls and recommendations for using machine learning to detect and prognosticate for COVID-19 using chest radiographs and CT scans,” Nature, Mar 15, 2021

“Hopefully it does slightly better or only slightly worse,” says Sukumar. “If the original model was trained to be 95% accurate and you can live with 90% accuracy when retrained on your own data, that seems to work for many people. That’s where much of the market is today.”

Speaking of data, its quality is yet another perennial issue that impacts every aspect of training. Sukumar says, “It’s such a big issue that when we talk about other problems, we assume the data is already clean.” That’s rarely the case, he says, estimating that up to 70% of the AI training effort can involve getting data ready for use.

“In my initial conversations about training, many customers say their biggest pain point is data preparation,” echoes Michael Woodacre, chief technology officer for high performance computing (HPC) at HPE. “It’s just so much more labor intensive in terms of the human interaction needed.”

Woodacre also emphasizes that you can’t discount the issue of cost. Training an AI model is expensive because it requires vast amounts of computing power, not just to run the GPUs doing the work but to keep them cool as well. “We’re looking at energy efficiency and carbon efficiency, but ultimately, it’s about how much computation you can get done per joule,” he says. “You need an optimized hardware platform, software stack, and the skill set to run parallel code as efficiently as possible.”

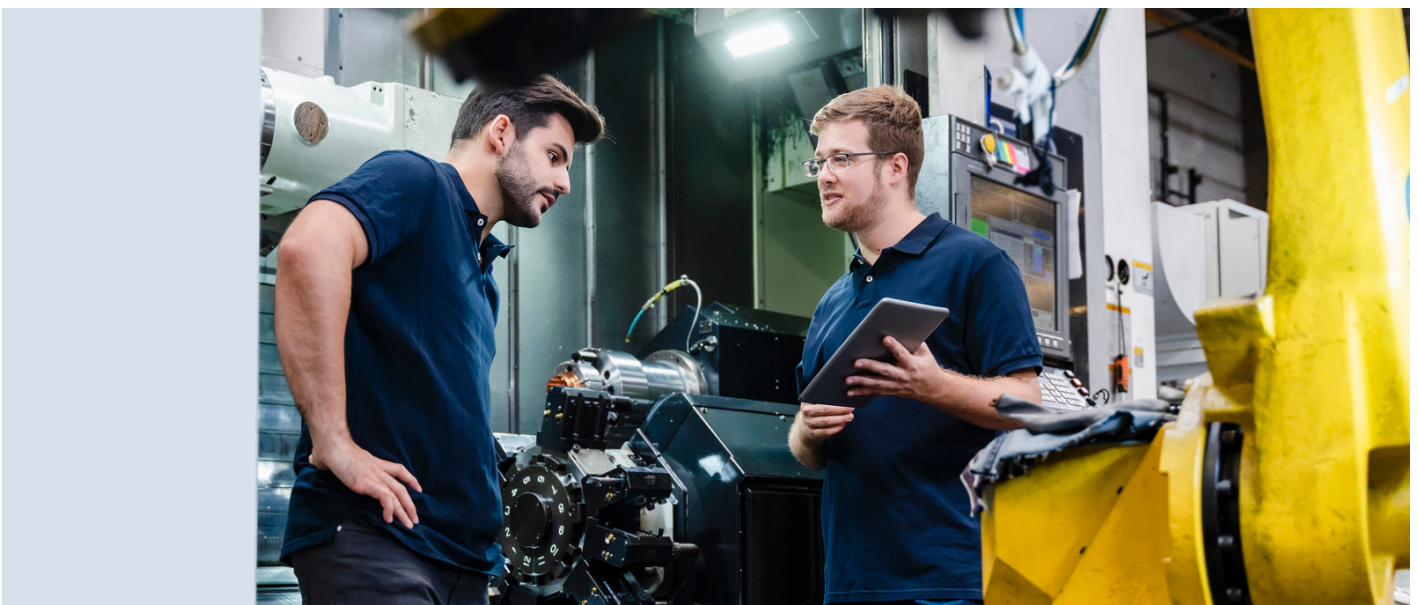
Streamline training by reducing complexity

So, what can you do to improve the training process? According to Woodacre and Sukumar, quite a bit. “These issues are similar to the challenges we’ve been working on for a few decades in the HPC community,” says Woodacre. “Now, the AI world is working to develop the solutions to similar challenges within a five-year time span, rapidly learning from the HPC experience.”

Adds Sukumar, “The cycle of AI starting off as an experiment and then scaling up to larger data sets and bigger models and then actually deploying those models in production isn’t trivial. The complexity of doing DataOps, DevOps, and MLOps all together is bigger than anything IT has dealt with before.”

To ease that complexity, Sukumar says services are being developed, including many by HPE, that can abstract many of those challenges for enterprises. A range of delivery models can, for example, orchestrate prepackaged AI solutions on either cloud or on-premises environments on behalf of customers — or provide access to GPUs on demand. AI jobs can be launched with just a few lines of code or batched for large-scale training on supercomputers.

“We can reduce complexity by minimizing the toolset that users have to interact with,” says Sukumar, which is an increasing necessity in a daunting world that now features hundreds of open-source AI tools. “We can pick the best winners for you and then give you the infrastructure and software to manage it all.” A strong partner can also help advise





on the proper model to use for an AI project, balancing cost and performance while taking steps to future-proof the system.

Employee training sessions and workshops can help overcome skills gaps among internal staff, but fine-tuning training operations can require more in-depth consultations. “We work directly with people who are building the state of the art and pushing the limits of technology,” says Sukumar, “helping them get that last bit of performance or energy efficiency out of their AI operations.”

Lastly, several tools are available to help enterprises manage the data side of things, says Woodacre. “Your model quality is directly impacted by your data quality. You must clean the data used for training and keep track of what happens every time you make an update, taking into consideration issues like data drift. And our MLDE/MLDM software suites — machine learning development environment and machine learning data management — are really focused on model development and

training, giving people the tools not only to train models but also to track and annotate data so they can better understand what data is being used by their AI models. For example, if you have regulatory requirements governing your model, you’ve got to know what’s going into it so you can reproduce results from the AI model when needed.”

Racing to stay competitive in your AI journey

Remember that you’re not alone if you find the world of AI moving faster than you’d like. This technological capability is in constant flux, and no single person can master every aspect of the field.

“It’s incredibly hard for individuals to keep up with it all,” says Woodacre. “People should feel comfortable reaching out to get the expertise they need to support them on their AI journey. Otherwise, you can become overwhelmed by the complexity.”

Learn more at

[HPE.com/AI](https://hpe.com/AI)

Visit **HPE GreenLake**



Chat now