

GRAPHCORE

BOW POD₂₅₆

Explore | Build | **Grow**

Bow Pod systems deliver high performance and efficiency for machine intelligence deployment at scale. They are designed to accelerate the large and complex models of today while also providing a platform for innovators to explore and invent the solutions of tomorrow.

The Bow Pod₂₅₆ system is the solution for innovators ready to grow their capacity to supercomputing scale. It delivers massive efficiency and productivity gains by enabling large model training runs to be completed in hours or minutes instead of months or weeks. Bow Pod₂₅₆ delivers AI at scale for production deployment in enterprise data centres, as well as private and public clouds.

Latest generation IPU

The Bow Pod₂₅₆ system features 64 Bow-2000 machines, each containing 4 of our pioneering Bow IPU processors. This innovative IPU is the world's first processor to be manufactured using Wafer-on-Wafer (WoW) technology, taking the benefits of the proven IPU technology to the next level.

System Specifications

Processors	256 Bow IPUs
1U blade units	64 Bow-2000 machines
Separate cores	376,832
Threads	> 2 million
Performance	89.6 petaFLOPS FP16.16 22.4 petaFLOPS FP32
Memory	230.4 GB In-Processor-Memory™ Up to 16,384 GB Streaming Memory™
Software	Poplar® SDK

Performance and efficiency

Bow Pod₂₅₆ delivers up to 89.6 petaFLOPS of AI compute as well as industry leading efficiency, all thanks to the use of innovative silicon technologies, a compute and memory architecture focused on efficiency and scale-out, and a software- and application-first approach in solution deployment.

Smooth deployment and short time to market

The whole system, hardware and software, has been architected together. Bow Pod₂₅₆ supports all standard frameworks and protocols to enable straightforward integration into existing data centre environments, as well as private and public clouds.

A wide selection of market leading server platforms and high-performance storage appliances designed for AI have been tested and validated to offer choice, in addition to short configuration and deployment times for system aggregators.

Innovators can focus on deploying their AI workloads at scale, using familiar tools and frameworks while unlocking cutting-edge performance and efficiency.

Host-Link	100 GE RoCEv2
System Weight	1800 kg + Host servers and switches
System Dimensions	64U + Host servers and switches
Host server	Selection of approved host servers from Graphcore® partners.
Storage	Selection of approved solutions from Graphcore partners.
Thermal	Air-Cooled

Disaggregation for customised compute

Machine intelligence workloads have very diverse compute demands. For production deployment, optimising the ratio of AI to host compute can help maximise performance, while improving total cost of ownership. Bow Pod systems allow flexible mapping of the number of servers and switches to the requisite number of Bow-2000 machines, so deployment is better tailored to production AI workloads. Bow Pod₂₅₆ supports multiple server configurations.

Communication architecture built for scaling

Efficient data access and transfer can unlock greater AI performance. IPU-Fabric is an innovative communication architecture for system-wide data transfer, extending high-speed interconnect within individual Bow IPUs, across Bow-2000s, between Bow Pods and throughout the data centre. IPU-Fabric delivers high-performance low-latency communication to maximise AI application efficiency and is built to work with standard data centre communication technologies.

Platform for AI developers

TensorFlow, PyTorch, PaddlePaddle, and many other popular ML frameworks are supported and available as open source, along with the comprehensive PopLibs™ library, for community driven

collaboration and innovation. For developers who want full control to exploit maximum performance, the Graphcore Poplar SDK enables direct IPU programming in C++.

Designed for deployment at scale

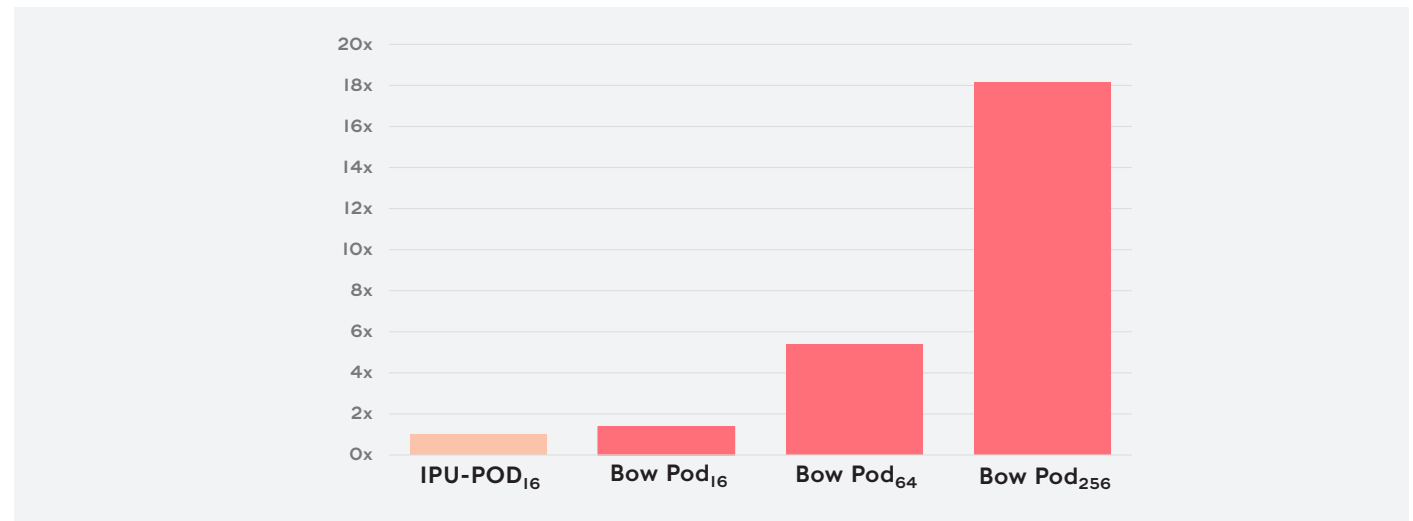
Pre-built Docker containers with Poplar SDK tools and frameworks images let innovators get up and running fast. Various common frameworks for container orchestration, platform visualisation and provisioning are also supported, including Slurm, Kubernetes and OpenStack.

Software First

Fully integrated and IPU-optimised, Poplar software leverages the unique characteristics of the IPU architecture to build AI applications of unrivalled performance and flexibility. Poplar allows effortless scaling of models from one to thousands of IPUs without adding development complexity, allowing innovators to focus on the accuracy and performance of the application.

Access to AI expertise

A wealth of experience and support for installation, production and application development is available globally from Graphcore AI experts and from our elite partner network.



Ready to experience the next level in Machine Intelligence?

Connect with our partners below to assess your AI infrastructure requirements and solution fit. Still have questions? Contact Graphcore directly at info@graphcore.ai