# GRAPHCORE'S AI SOFTWARE STACK IS NOW CUSTOMER-DRIVEN

## INTRODUCTION

In a paper published in May 2020, we examined the hardware and software offerings of Graphcore, detailing the development toolchain and noting the commitment to the open-source development community. After 18 months, Graphcore has released the latest version of the Poplar environment (SDK2.4). The company has created critical mass with an open-source development community across various applications to help mature and expand the Graphcore software stack. This updated research paper will examine advancements the company has introduced and the level of community involvement in ongoing software and model engineering.

Software for new processor designs is critical to enabling application deployment and optimizing performance. UK-based startup Graphcore, a systems provider for application acceleration, emphasizes software, dedicating over half its engineering staff to the challenge. Graphcore's second-generation Intelligence Processing Unit (IPU) utilizes the expression of an algorithm as a directed graph, and the company's Poplar software stack translates models and algorithms into those graphs for execution. The software simplifies adapting the IPU-based system for AI and parallel computing, making it vital to success.

All AI hardware startups struggle to break through the adoption impasse: no software ➔ , no customers, ➔ no software. (By "software," we mean optimized models and kernels, in addition to frameworks and tools.) Graphcore, however, seems to be approaching a critical mass of company-engineered and open-source contributions, enabling a customer-centric mindset to drive enhancements. This community has a long way to go before remotely resembling NVIDIA's massive ecosystem to enjoy and harvest. But every journey begins with the first steps, and we are impressed with Graphcore's considerable progress.

## THE MATURING OF THE GRAPHCORE SOFTWARE PLATFORM

The company has made significant strides since we published the original version of our research in May 2020. Beyond the Graphcore-engineered enhancements, the software ecosystem around the IPU is evolving from a push to a pull model: the AI community is now actively contributing to the open-source software Graphcore has released. Graphcore provides an extensive array of learning modules, reference designs, and tools; we would assert that only NVIDIA exceeds Graphcore in developer enablement.

The Graphcore stack is organized into four pillars: tools and apps for the developer ecosystem, the Graphcore engineered Poplar SDK, PopVision performance analyzer tools, and System software for deployment and management.
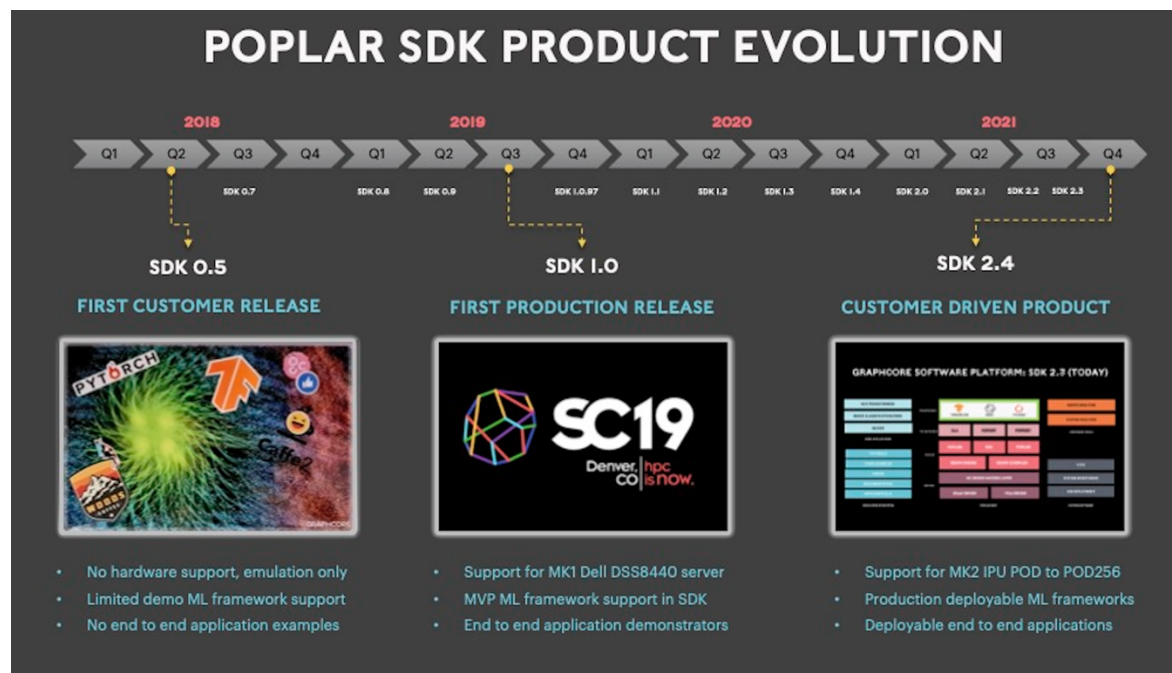


Figure 1: The Poplar SDK development timeline.

The original Graphcore software stack was a solid start, but as its engineers supported client projects, they evolved the stack to become more comprehensive and the tools easier to use. Notably, the company's commitment to open-source the entire stack (except the compiler) enabled clients to help build and expand the software stack and available algorithms/models. The end-user engagement and contributions provide a multiplier effect to the internal software engineering efforts in Graphcore.

Let's look at what has changed in the Poplar SDK since our last report. SDK 2.4 has added significant features. The SDK includes ten customer-driven enhancements to frameworks, including enhanced Pytorch integration, TensorFlow Lamb optimizer support, and enhanced scalability for the POD256. Additional support for deployment includes enhanced support for inference processing. We find this telling, so we expect more developments from Graphcore on the inferencing front this year.

One of the many additions that will help customers understand how models are optimized and deployed on IPU clusters are three core applications: open-source, optimized stacks for NLP/Transformers, Image Classification, and the MLPerf models, which Graphcore submitted to the MLPerf V1.1 training benchmarks. Large Language model optimizations were also engineered, covered later in this report.

Figure 2: The 2.4 release includes some ten customer-driven enhancements.

Graphcore is addressing two challenges in software development for its IPU platforms: 1) make it easy to optimize and run existing ML software such as deep neural networks, or DNNs, expressed in high-level frameworks, and 2) enable research and development of entirely new fine-grained parallel workloads to run on an IPU infrastructure. The latter ability is central to the company's strategy and could enable Graphcore to address a much larger set of market segments, including finance and HPC applications.
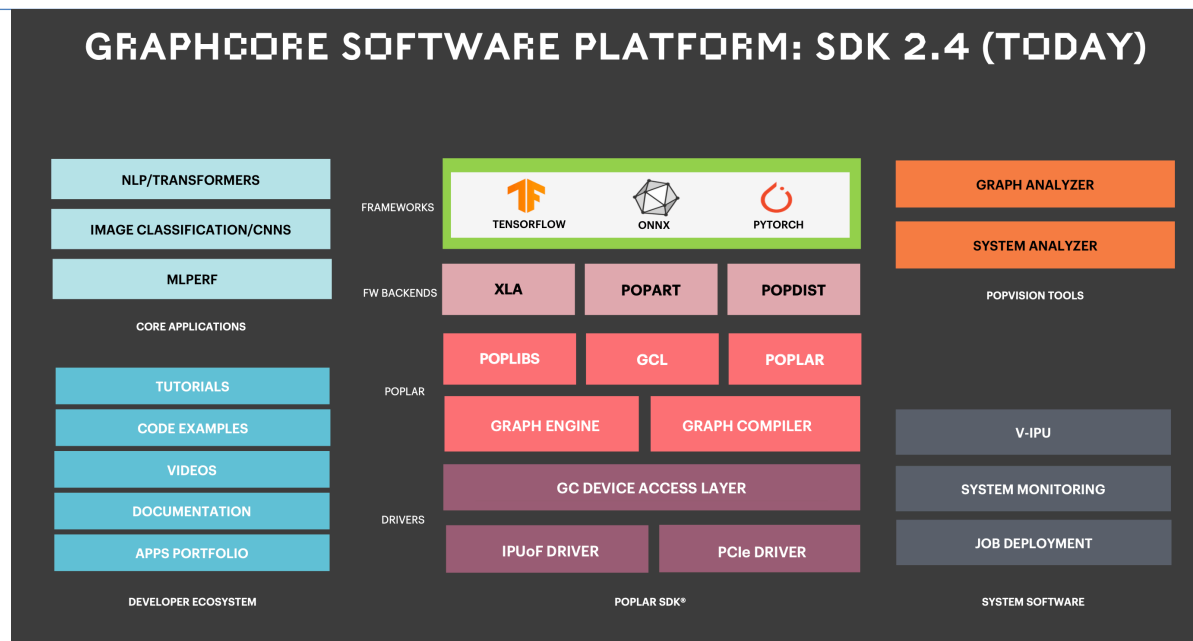
**GRAPHCORE SOFTWARE PLATFORM: SDK 2.4 (TODAY)**

Figure 3: The latest release of the Poplar SDK (V2.4) is complemented by system software, tools, applications, and developer enablement initiatives.

## *PARTNER COLLABORATION*

While Graphcore's hardware and software are innovative and impressive, they alone are insufficient to build Graphcore into a material force on the acceleration landscape. To turn technology into valuable solutions, the company continually provides developer and partner programs to close the gap between potential and realized success. The company's near-constant release of new webinars, tools, tips, examples, and education are an industry best practice; we have not seen such quality and quantity from any AI startup. And the open-source software strategy provides the needed foundation to make this happen.

Graphcore is the only AI company that has released the source code of their ML framework backend implementations. Now partners such as Baidu (Paddle Paddle) and Alibaba (ODLA) have built their solutions on top of this. Graphcore enabled these projects by publicly releasing the source code to the PopART framework and Pytorch integration.
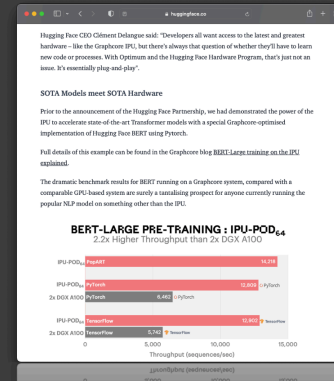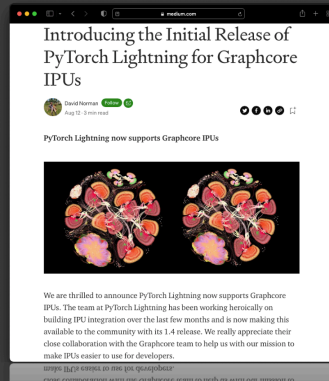
**ECOSYSTEM ENGAGEMENT**

VMWare and Pure Storage recently became Graphcore partners as well. VMWare and Graphcore have collaborated on virtualizing IPU's to enable pooling and sharing of IPU resources, providing familiar management and deployment constructs to data center managers.  VMWare offers a valuable extension of the IPU-POD architecture for Enterprises that want to share these parallel-processing resources across the organization, using tools they are already using for virtualized multi-tenant hosting. This technology alliance is a meaningful indicator of traction and demand for Graphcore solutions instead of substantial deployment announcements, which customers tend to avoid.

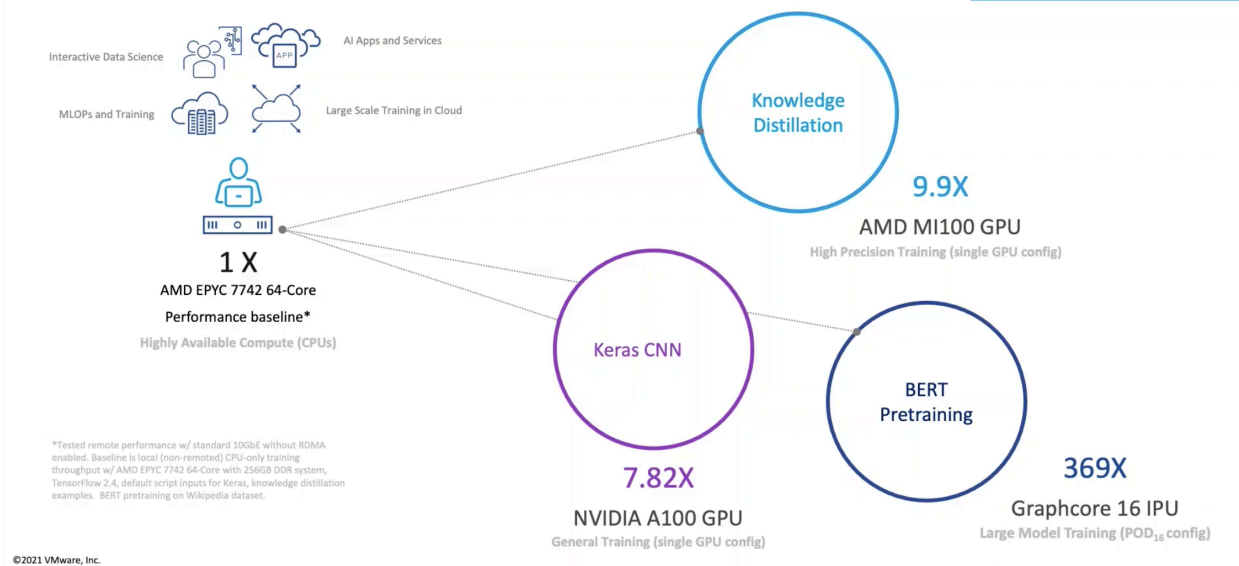**Figure 4: VMWare and Graphcore have collaborated to bring virtualization to IPU users focused on BERT NLP.**

Another example of partner collaboration has brought state-of-the-art flash storage into the growing ecosystem. PureStorage has worked with Graphcore to define and deploy a reference architecture for the PureStorage Flashblade infrastructure. The PureStorage eXternal Fabric Modules scale linearly from one to ten chassis and 7 to 150 blades, providing up to 40x17 TB of storage for Graphcore IPOD customers.

*DEVELOPER COLLABORATION & SUPPORT*

As a result of the open approach and ecosystem enablement, end-users are increasingly embracing the Graphcore platform. The Graphcore Model Garden is an excellent example of ecosystem investments gaining traction with a repository of tested and optimized ML models ready for deployment. Much of this work is a natural outcome of the work done to submit MLPerf benchmark runs.

One of the most challenging tasks developers face when exploring new hardware is optimizing their code and models for an entirely new hardware architecture than what they might be accustomed to. In addition to the Model Zoo, webinars, documentation, etc., Graphcore offers an in-depth memory and performance optimization guide to help developers on this journey. This includes model pipelining and overlapping I/O with computation, among other topics. By providing this level of support, much of which undoubtedly emerged from client engagements, Graphcore can efficiently support the development community and enlarge its ecosystem.
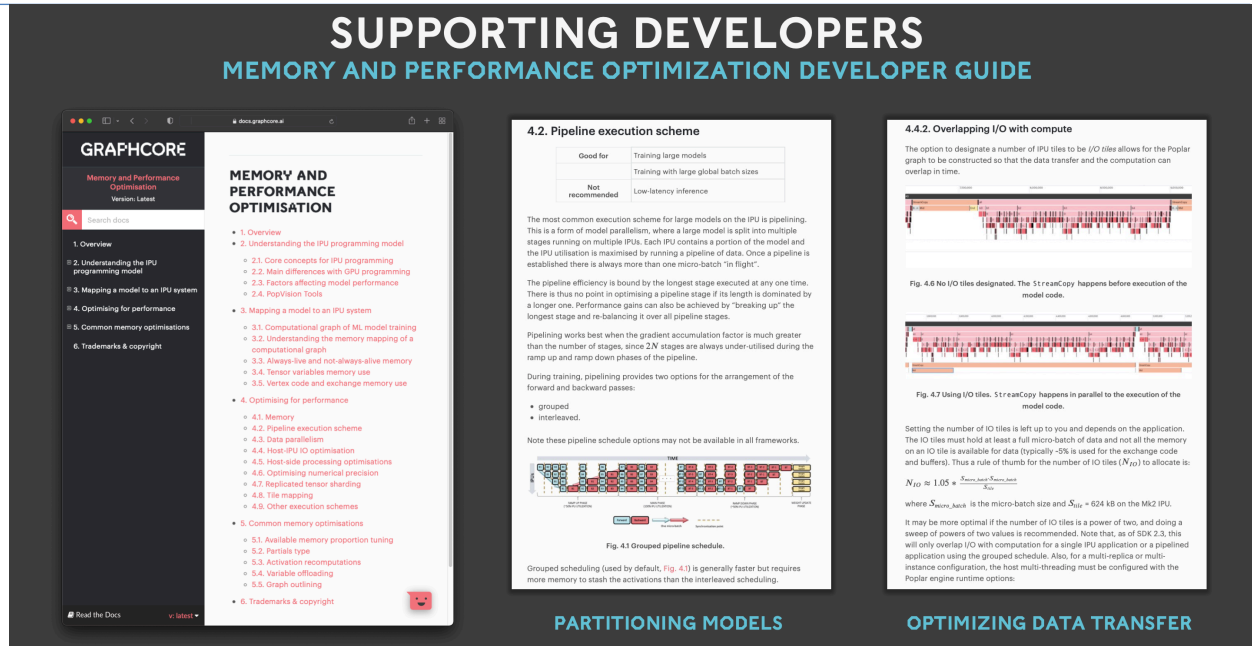
**Figure 5: Graphcore has produced a guide to help developers learn how to tune and optimize their parallel processing applications running on the IPU.**

Graphcore is also adept at leveraging the tremendous body of work from MLPerf submissions and turning the experience and code into learning tools, such as the webinar on BERT-large tuning and Jupyter notebook.

## SOFTWARE PERFORMANCE ENHANCEMENTS

Like other players in the AI Acceleration market, Graphcore has continually delivered significant performance improvements as their software matures, with up to 2.4X improvement in image processing over the last year and a 60% decrease in BERT-Large training time. We anxiously await performance data on inference processing; running models on the same platform used for training has potential synergies and cost savings.

**EVOLVING PERFORMANCE**
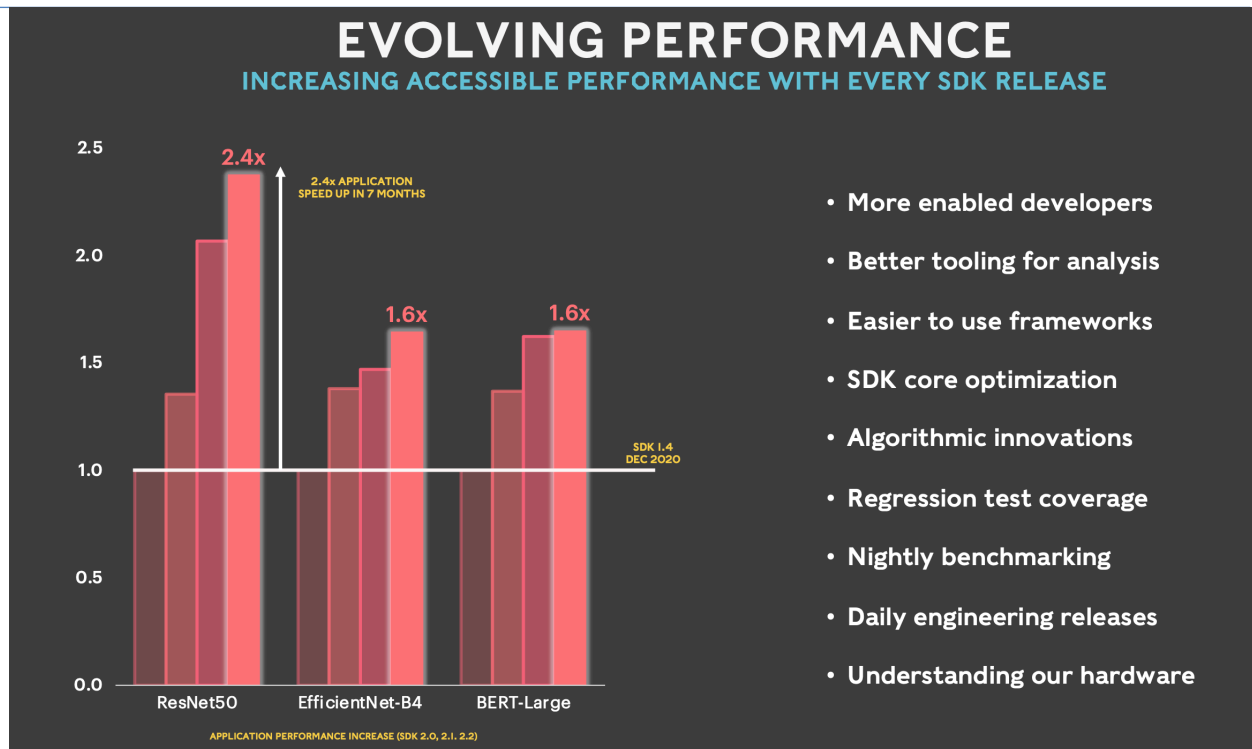INCREASING ACCESSIBLE PERFORMANCE WITH EVERY SDK RELEASE

*Figure 6: Graphcore has continually improved the performance of MLPerf and other benchmarks, and we expect them to continue to make progress. The company has also improved the scalability of these and other models with the POD256.*

### LARGE LANGUAGE MODEL TRAINING

As we have noted, the updated SDK 2.4 comes with a rich set of tools, training, webinars, videos, podcasts, and documentation. Graphcore also provides a model garden to enable experimentation and adoption. These all assist developers and deliver content for Graphcore marketing to attract developers to the platform.

One of the most recent Graphcore communication campaigns promoted Large Language Model (LLMs) development on IPU-PODs. LLMs have become one of the latest innovations for AI, spurred by Open.AI GPT-3, and the subject of a media frenzy. While many of these multi-billion-parameter models have been trained on thousands of GPUs, the Graphcore architecture provides a cost-effective solution thanks to the exchange memory on each IPU-Machine, complementing the fast on-die memory of each IPU.

A recent Graphcore blog discussed how the 128- and 256-node clusters support models with tens of billions of parameters in the IPU-POD$_{256}$ utilizing the 16 TB of memory in these configurations, which are now fully supported and characterized by Graphcore. Graphcore shared the scalability of training the ViT vision transformer and GPT-2 (which, unlike GPT-3, has been open sourced), demonstrating roughly a 12X performance increase with 16X more IPU nodes. Beyond that, Graphcore believes the

memory capacity of the POD$_{256}$, which supports 16 terabytes of IPU Exchange Memory, opens the possibility to support models with trillions of parameters. These brain-scale models can also benefit from the new compute approaches made possible by our IPU processors.
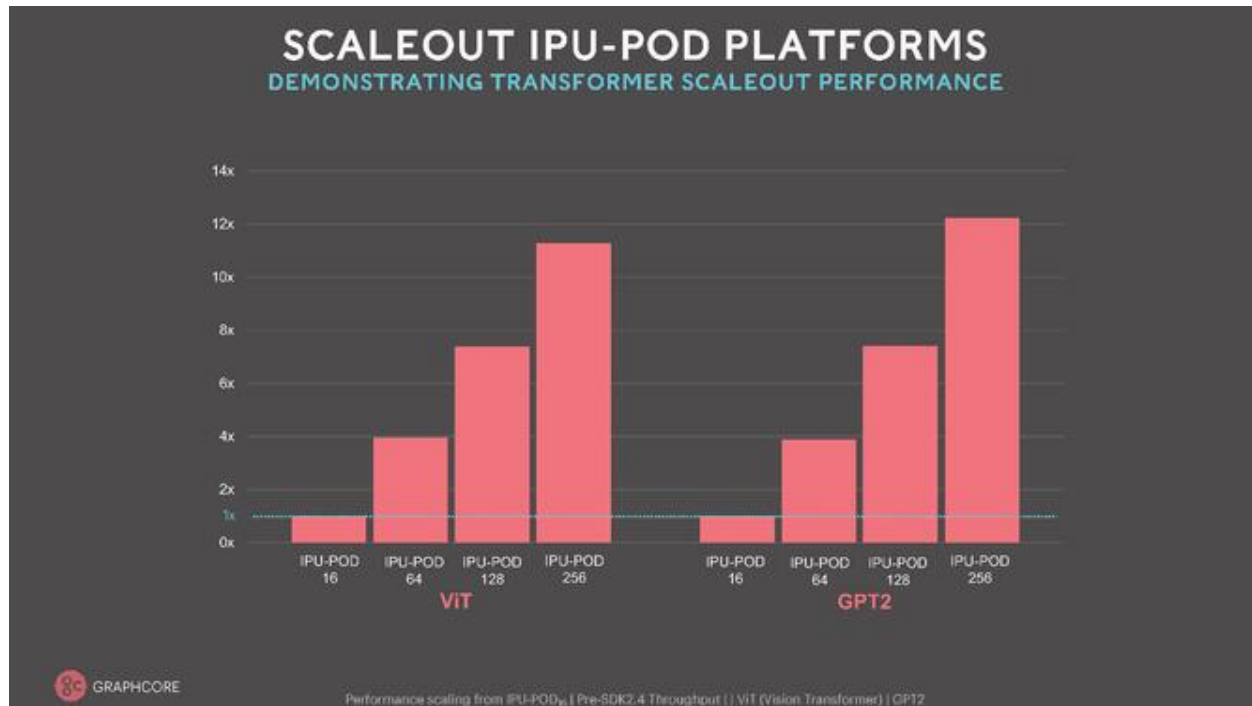


**Figure 7: Graphcore has recently demonstrated excellent performance and scalability when training Large Language Models such as the Vision Transformer and GPT-2.**

The Graphcore Bulk Synchronous Parallel Processing approach lends itself well to large models, with optimized data communications and phased execution enabling efficient distributed processing of hundreds of billions of parameters. By combining pipelining, tensor, and model parallelism, Graphcore can support massive models without burdening the developers with the job of manually managing these parallel queues.
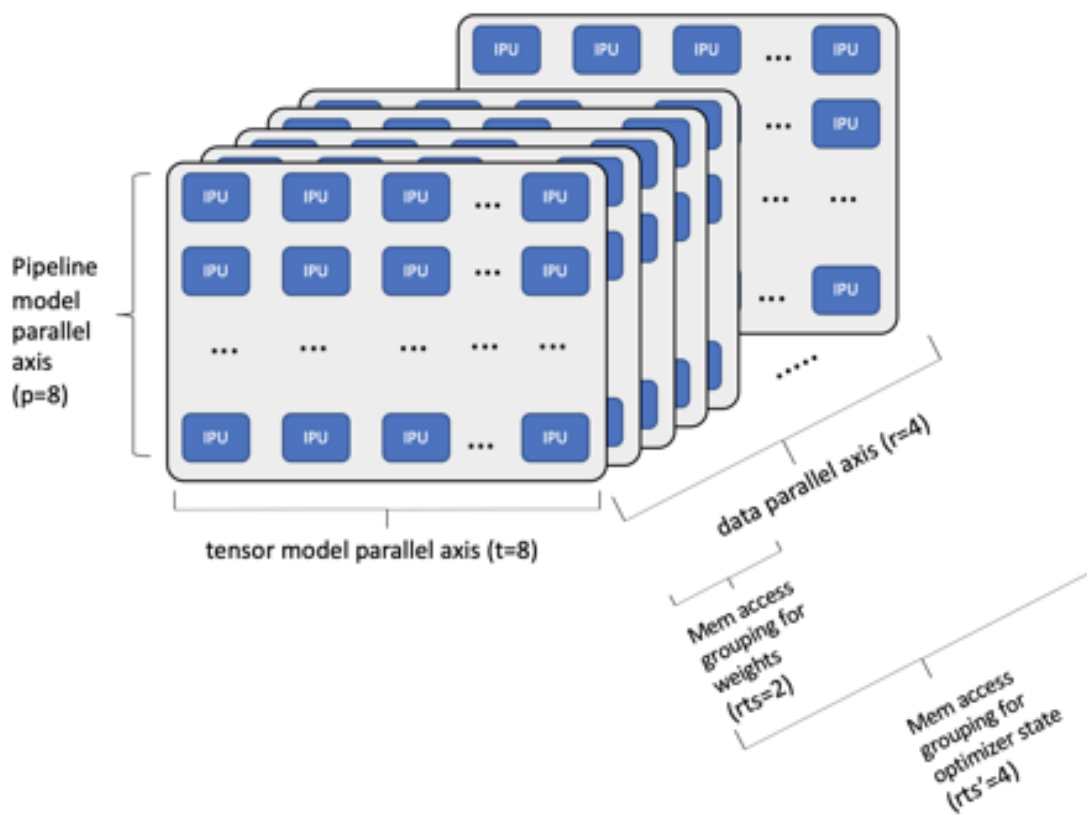
Figure 8: Large Language Model deployment on Graphcore IPUs.

## THE BENEFITS OF HOST DISAGGREGATION

Depending on the scalar processing and management required, AI models and other parallel applications may require vastly different accelerator/server ratios. Consequently, instead of designing an IPU server on a PCIe card inserted into a 1- or 2-socket server, as is the case for almost all AI accelerators, the IPU-Machines are disaggregated from the host server(s), communicating over 100Gb Ethernet. Graphcore recently shared some data to illustrate the potential saving of this approach, and it is compelling. Why spend more money on CPUs if your application doesn't need them. And the infrastructure is dynamic, allowing different CPU to IPU ratios concurrently and re-configurations on the fly.
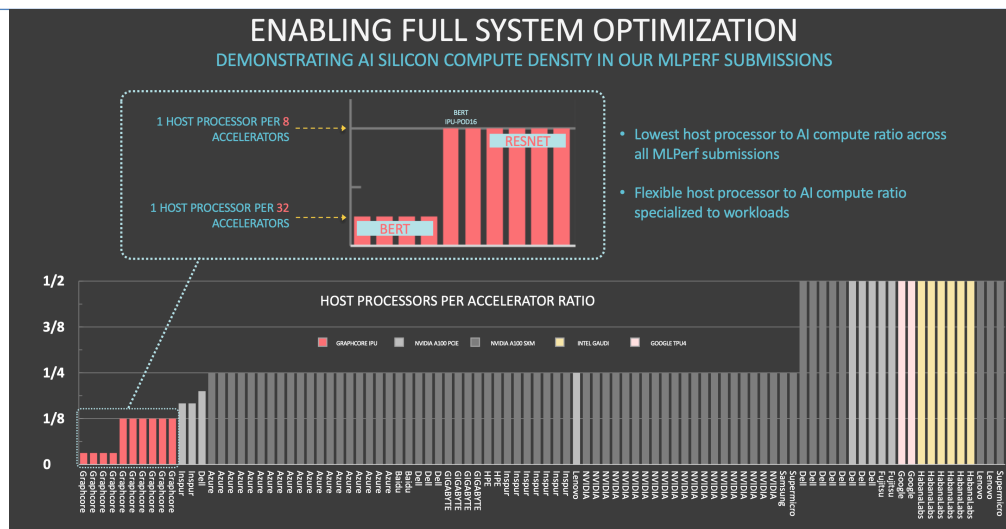
**Figure 9: Graphcore illustrates how server disaggregation enables a more cost-effective training infrastructure, with a variable number of CPUs per accelerator.**

This plug-and-play approach of CPUs and IPUs enables:

1. Optimized performance for models that require more servers than, say, two CPUs for 4 or 8 accelerators,

2. Lower costs for models that require fewer CPUs by avoiding over-provisioning,

3. A flexible data center infrastructure that can handle both, and

4. Servers can reside in shared utility racks for optimal rack power utilization and serviceability.

## THE ROAD AHEAD

Innovation in AI hardware and software presents one of the fastest-paced technologies today, and Graphcore has ambitious plans for both. In the software area, the company plans to add more applications such as conditional sparsity and GNNs. On the framework front, Graphcore has invested significantly in supporting the two main machine learning frameworks of Pytorch and Tensorflow alongside providing a lightweight and flexible runtime that supports ONNX and experimentation in PopART. These projects are open-source and can be leveraged as an underlying platform to support further ecosystem engagement and growth for programming frameworks in 2022.

New and enhanced support for Keras, Huggingface, TIMM, JAX, Pytorch Lightning, ODLA, and Paddle Paddle are all planned to increase IPU platforms' adoption further and make Graphcore products available to more developers. Alongside these, there will be product enhancements to support inference deployments with an inference-focused

toolkit and broader developer productivity support for Jupyter notebooks that can be used with standard ML ops platforms.
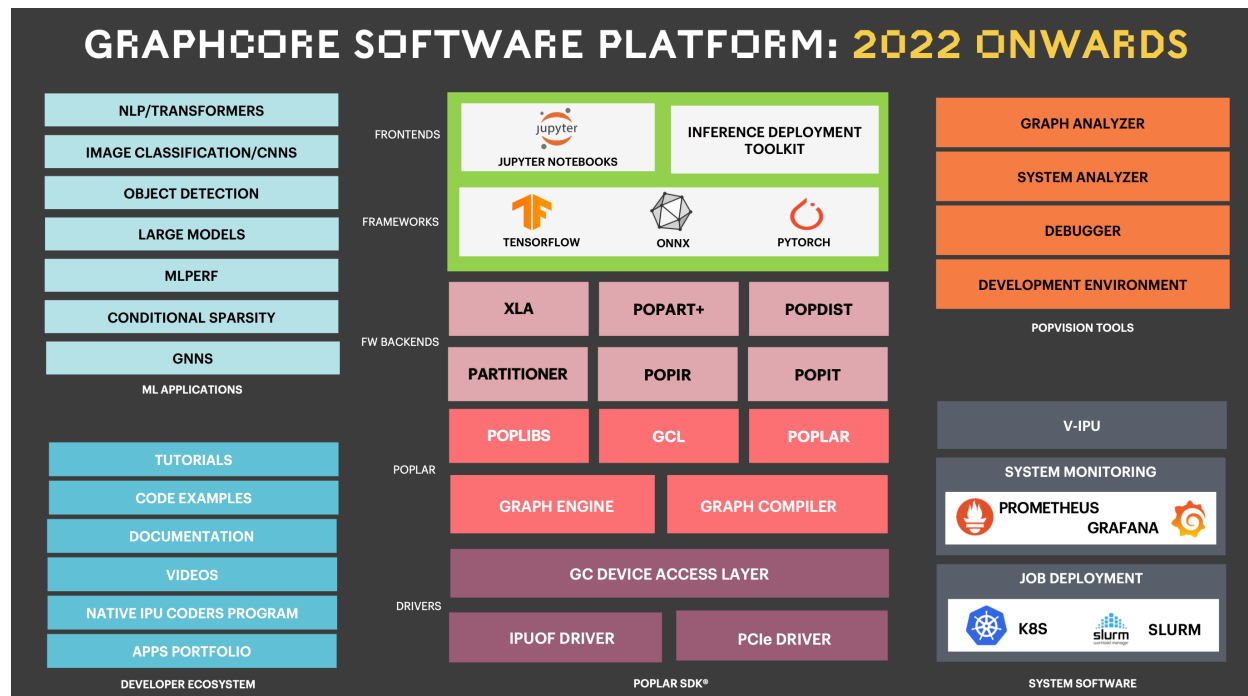


Figure 10: Graphcore has shared a high-level roadmap of software it is currently developing.

## CONCLUSION: GRAPHCORE HAS ENABLED A SELF-EVOLVING ECOSYSTEM FOR PARALLEL PROGRAMMING ON THEIR HARDWARE AND IS CONSEQUENTLY GAINING MOMENTUM.

Developers require fast, scalable accelerators to handle the massive computational loads of larger models such as natural language processing and conversational interfaces to train large neural networks. But AI developers also need a robust development environment that meets them in their AI development journey. The latest version of Graphcore software delivers and is benefiting from client contributions to the open-source base and is being picked up by end-users and partners.

More broadly, computationally intensive applications are emerging which require a high-speed parallel processor that goes beyond the matrix-multiplication operations common to neural networks. Graphcore appears to have developed just such a general-purpose, high-performance hardware and software platform to meet these needs. The emergence of Large Language Models with performance data indicates that Graphcore intends to stay on the leading edge of model development and optimization.

## IMPORTANT INFORMATION ABOUT THIS PAPER

### AUTHOR

Karl Freund, Founder and Principal Analyst at Cambrian-AI Research

### PUBLISHER

Karl Freund, Cambrian-AI Research, LLC.

### INQUIRIES

Contact us if you would like to discuss this report, and Cambrian-AI Research will respond promptly.

### CITATIONS

This paper can be cited by accredited press and analysts but must be mentioned in the context, displaying the author's name, author's title, and "Cambrian-AI Research." Non-press and non-analysts must receive prior permission from Cambrian-AI Research for any citations.

### LICENSING

This document, including any supporting materials, is owned by Cambrian-AI Research. This publication may not be reproduced, distributed, or shared in any form without Cambrian-AI Research's prior written permission.

### DISCLOSURES

Graphcore Inc. Cambrian-AI Research commissioned this paper to provide research, analysis, advising, and consulting to many high-tech companies mentioned in this paper. No employees at the firm hold any equity positions with any companies cited in this document.

### DISCLAIMER

The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions, and typographical errors. Cambrian-AI Research disclaims all warranties as to the accuracy, completeness, or adequacy of such information and shall have no liability for errors, omissions, or inadequacies in such information. This document consists of the opinions of Cambrian-AI Research and should not be construed as statements of fact. The views expressed herein are subject to change without notice.

Cambrian-AI Research provides forecasts and forward-looking statements as directional indicators and not as precise predictions of future events. While our forecasts and forward-looking statements represent our current judgment on what the future holds, they are subject to risks and uncertainties that could cause actual results to differ materially. You are cautioned not to place undue reliance on these forecasts and forward-looking statements, which reflect our opinions only as of the date of publication for this document. Please keep in mind that we are not obligating ourselves to revise or publicly release the results of any revision to these forecasts and forward-looking statements considering new information or future events.