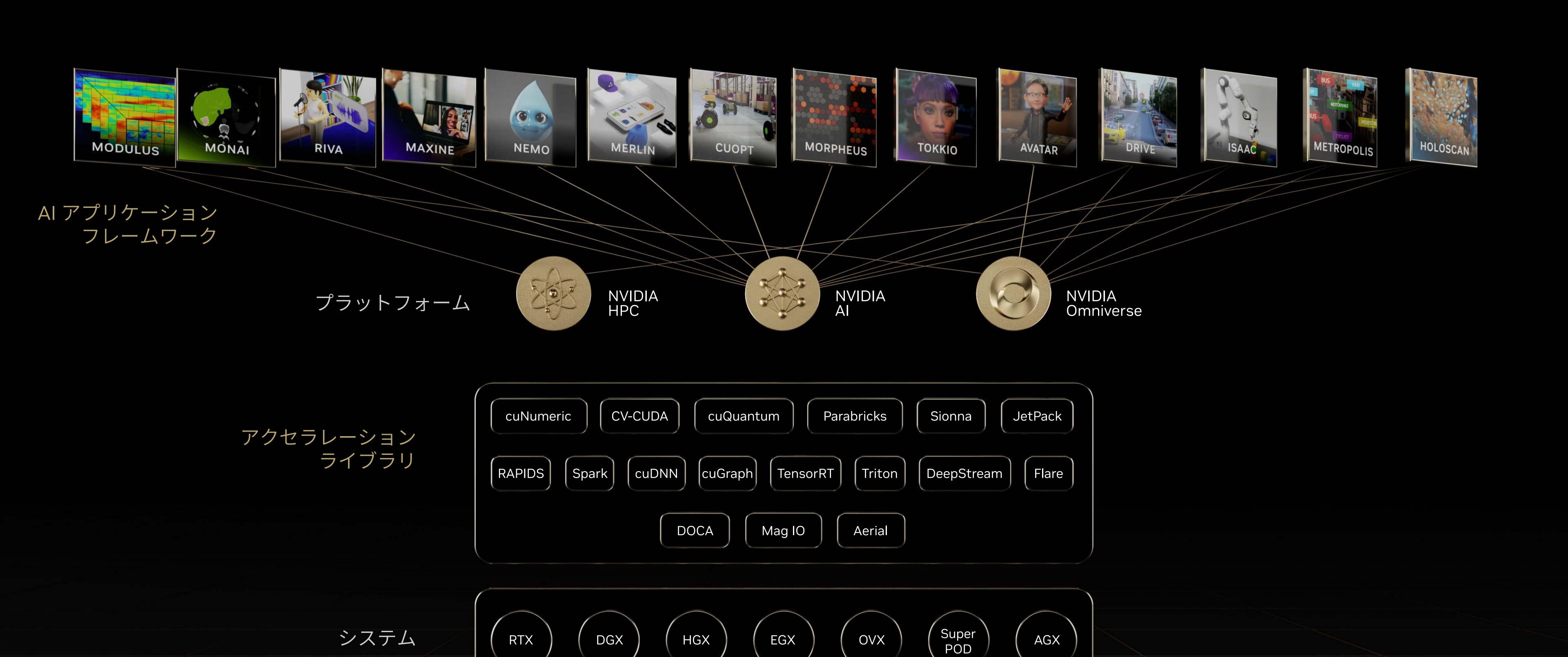


生成AI技術やビジネスの最新動向

エンタープライズ事業本部 事業本部長

井﨑 武士





GPU

CPU

DPU

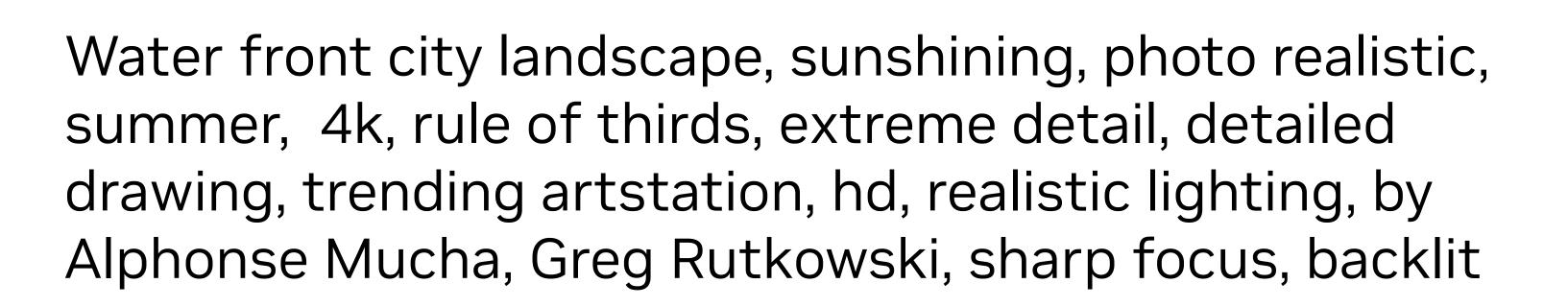
3 つのチップ

生成AIの衝撃

iPhone moment

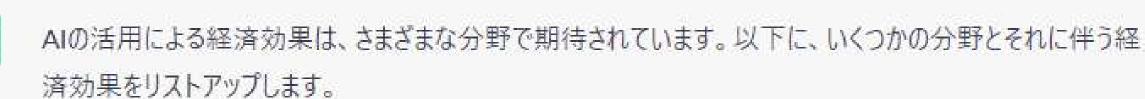
Stable Diffusion





ChatGPT

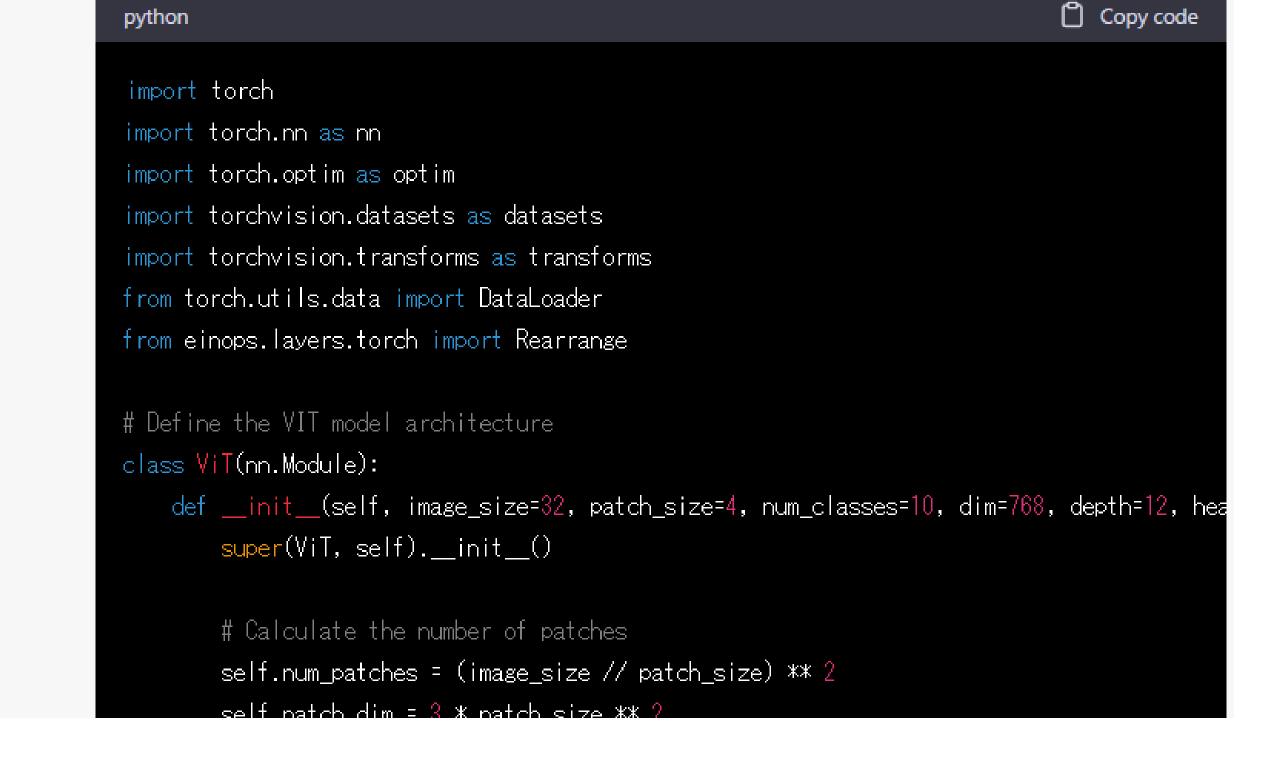
プロの経済アナリストの立場で答えてください。AIの活用による経済効果を分野ごとにリストアップしてください



- 1. 製造業:
 - 自動化とロボット化による生産性向上
 - プロセスの最適化による品質向上
 - 機械学習による予測メンテナンスによる機械の故障
- ・ データ分析による在庫最適化と生産計画の最適(
- 2. 金融業:
 - 高速なデータ処理とパターン認識によるリスク管理
 - モデルベースの予測による株式市場や為替市場の
 - ロボアドバイザーによる個人投資家の助言やポート。
 - 不正検知によるセキュリティの向上
- 3. 医療業界:
 - パターン認識と画像解析による病気の早期検出
 - データマイニングによる臨床試験の最適化と新薬開
 - ヘルスケアデータの分析による予防医療の向上とコ
 - ロボット手術による手術の精度向上と手術時間の

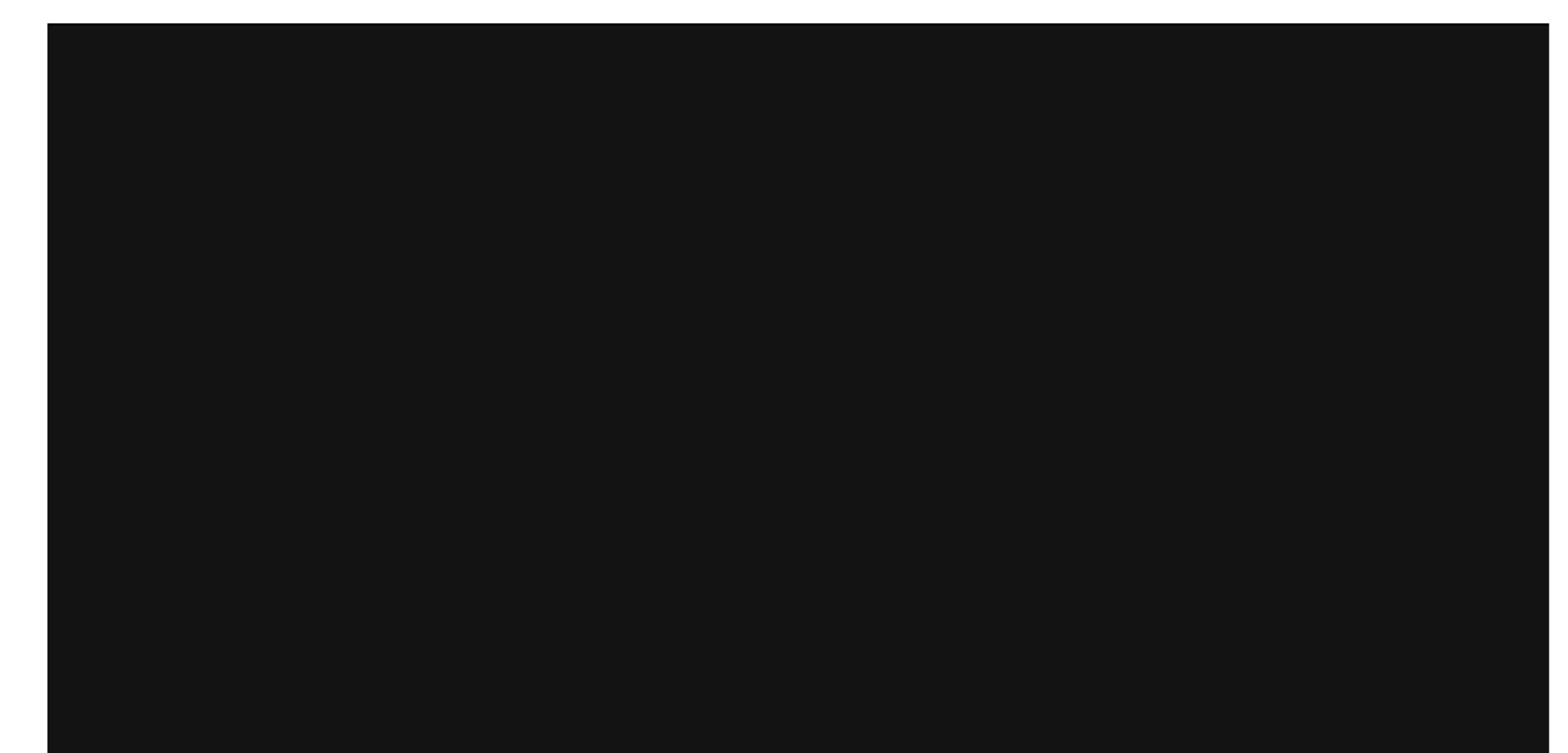
クルマの車種を分類するVITを使ったPyTorchのスクリプト書いて

以下は、PyTorchを使用してVITを使ったクルマの車種分類スクリプトの例です。このスクリプトは、CIFAR-10データセットで訓練されたもので、VITのパラメータは事前に学習されたものを使用しています。





Insight Face 独自の顔画像認識技術とMidjurneyを組み合わせて 顔のスワッピングを行う



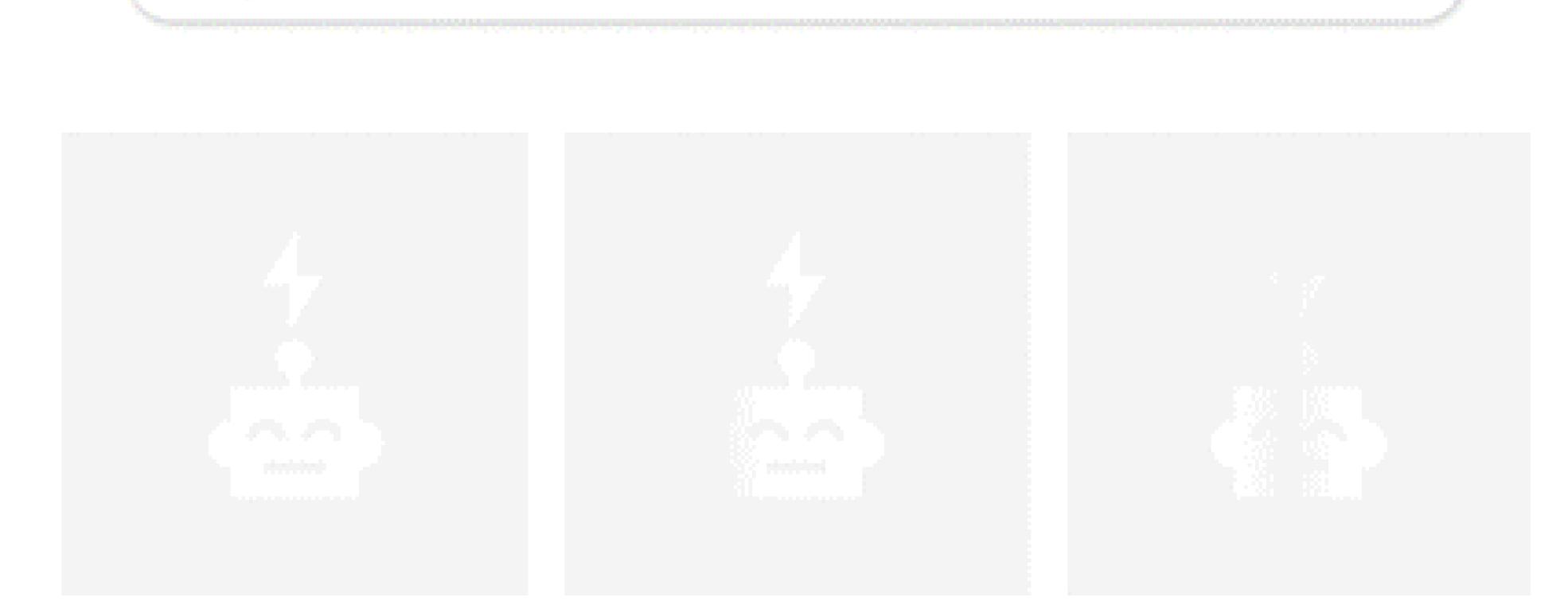
Briefly 電話会議の文字起こしを元に洗練された文書や会議後の 成果物を生成する





chooch

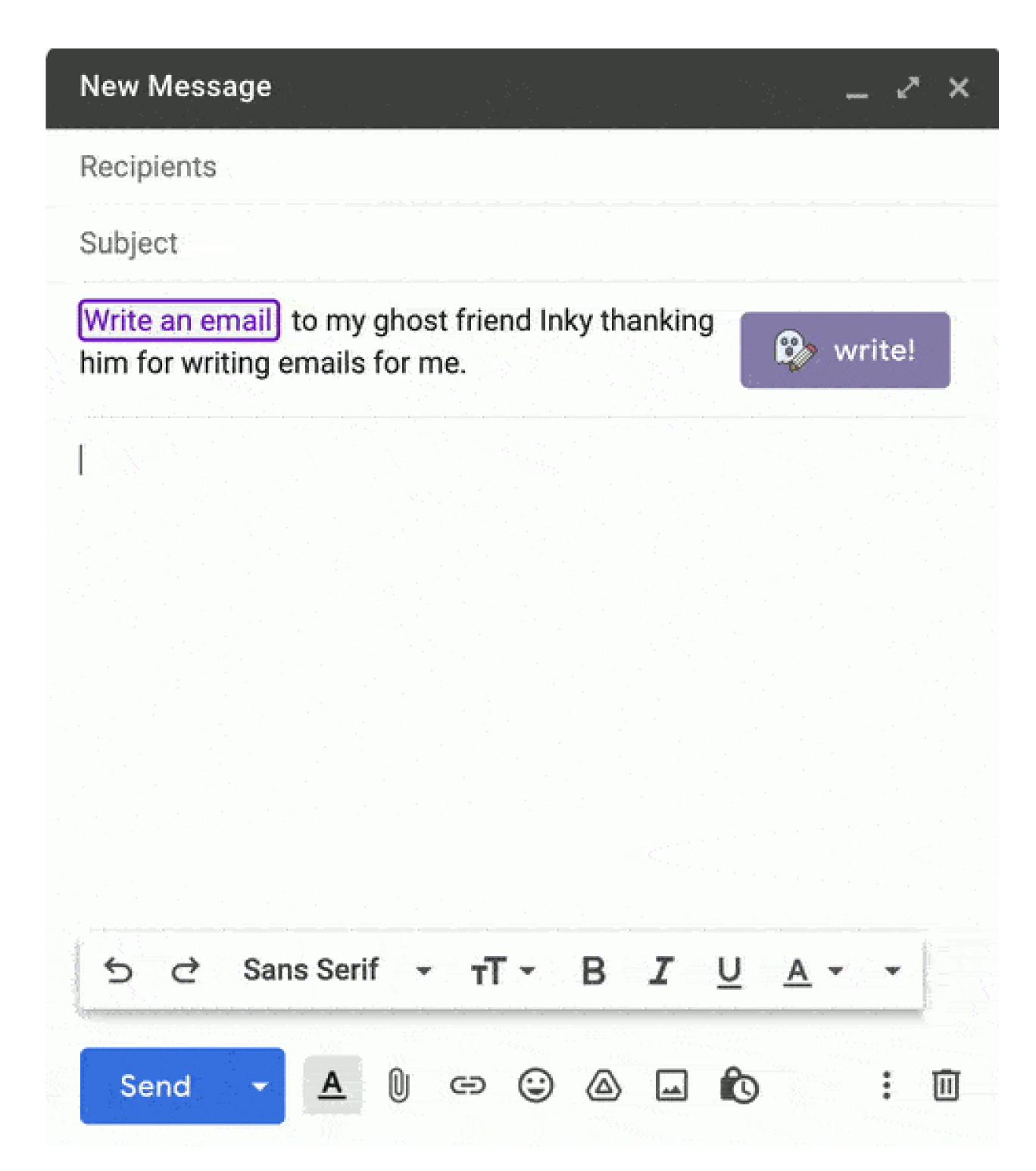
画像データに関するQ&AをChat形式で 行う



BeautifulAl

入力されたプロンプトに応じて、プレゼンテーション スライドを作成する





Ghostwrite

ChatGPTを活用した、e-mail 作成アシスタント

入力文章 Have you tried ChatGPT or Bard? Lots of fun! Try it!

When you try
ChatGPT for the
first time and
realize it's amazing



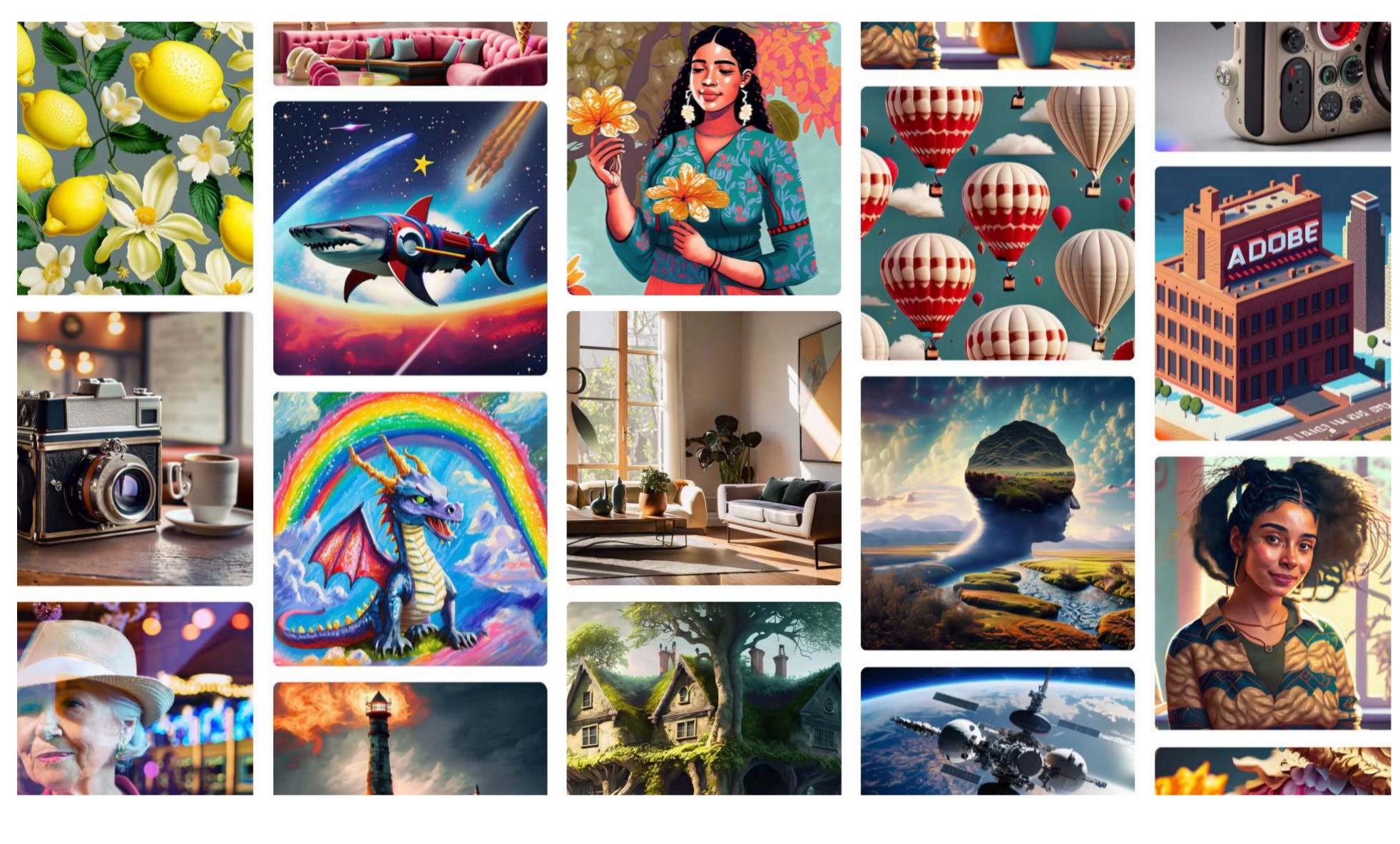
When someone asks you why they should try Bard, but you just can't put into words how great it is

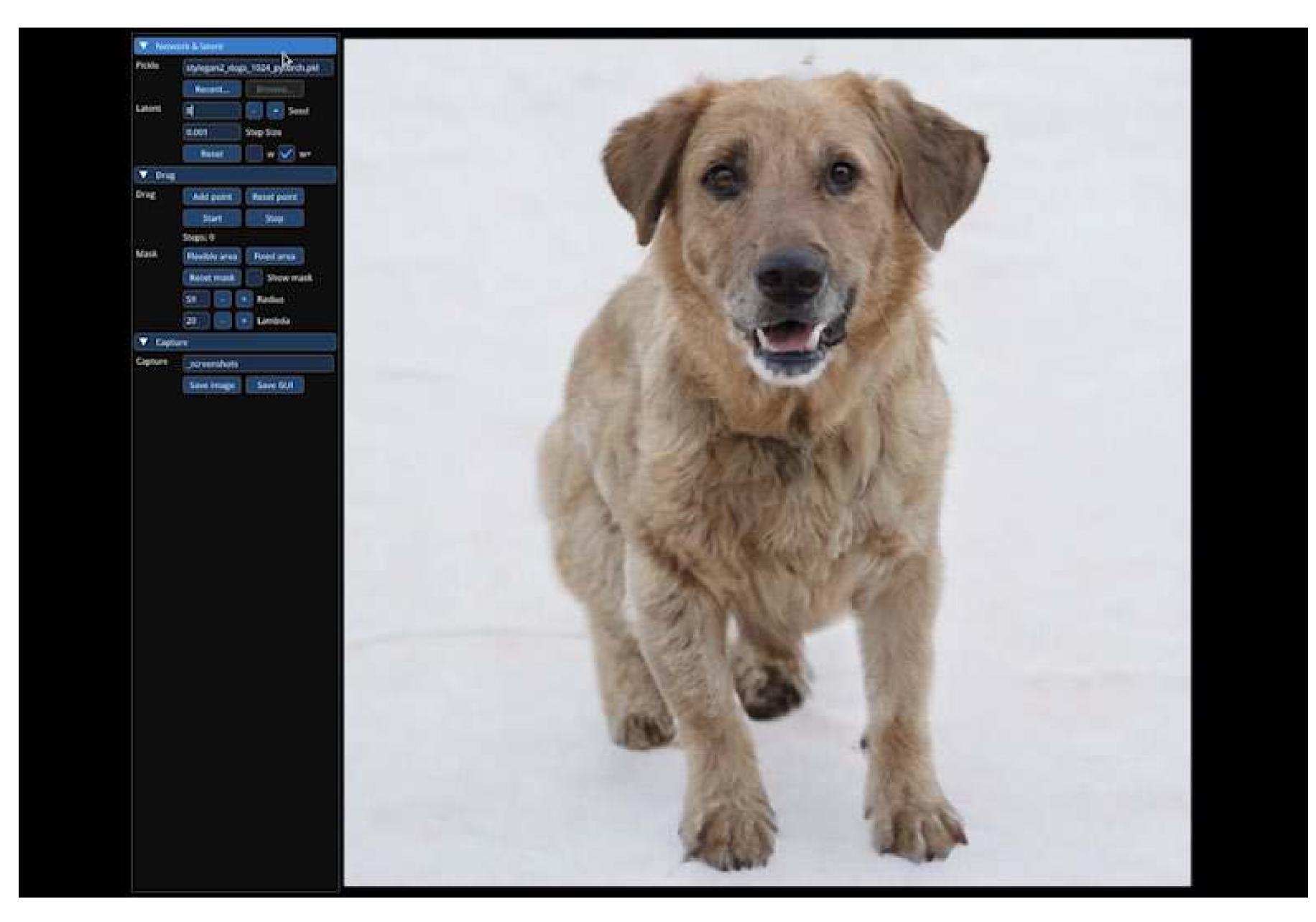


Super Meme

入力文章に応じたミーム(静止画やGIF)を生成 110以上の言語に対応





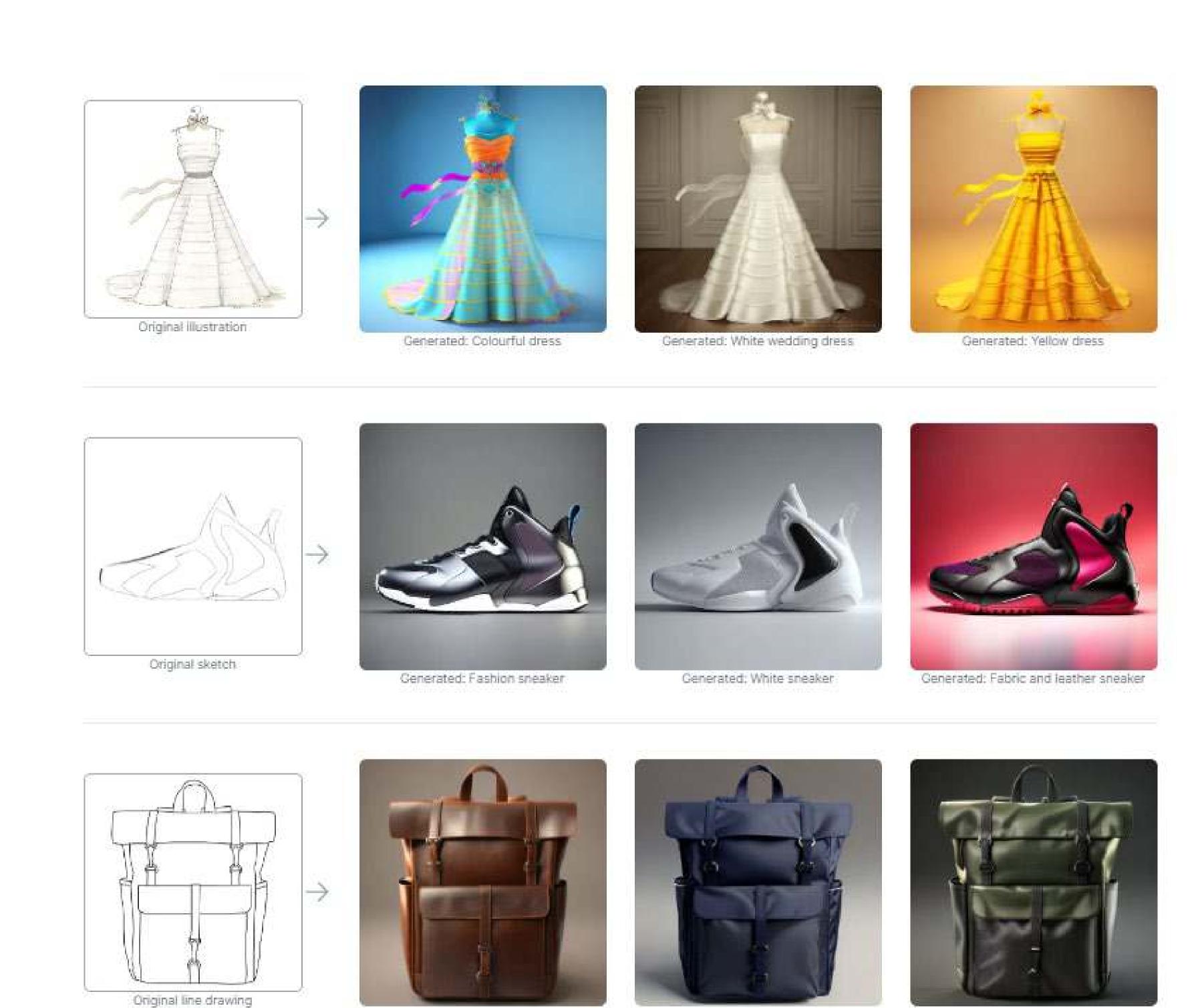


Adobe.comより

Adobe Firefly ベータリリース

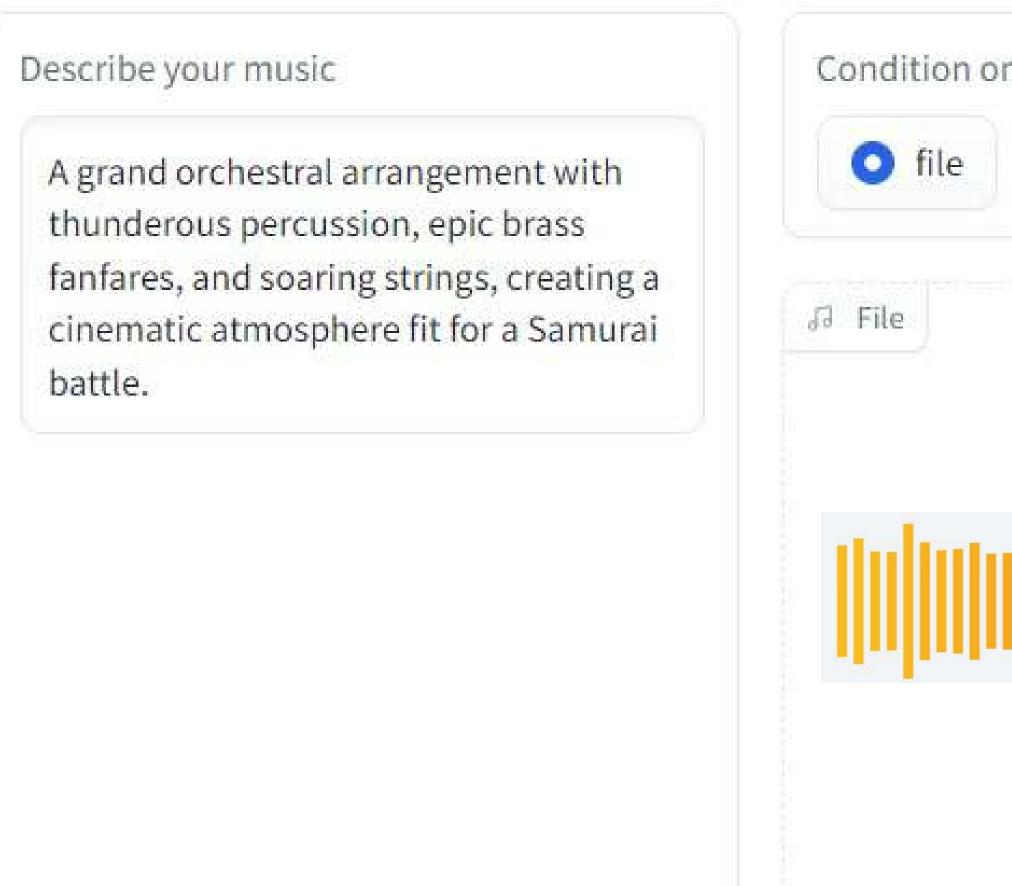
DragGAN

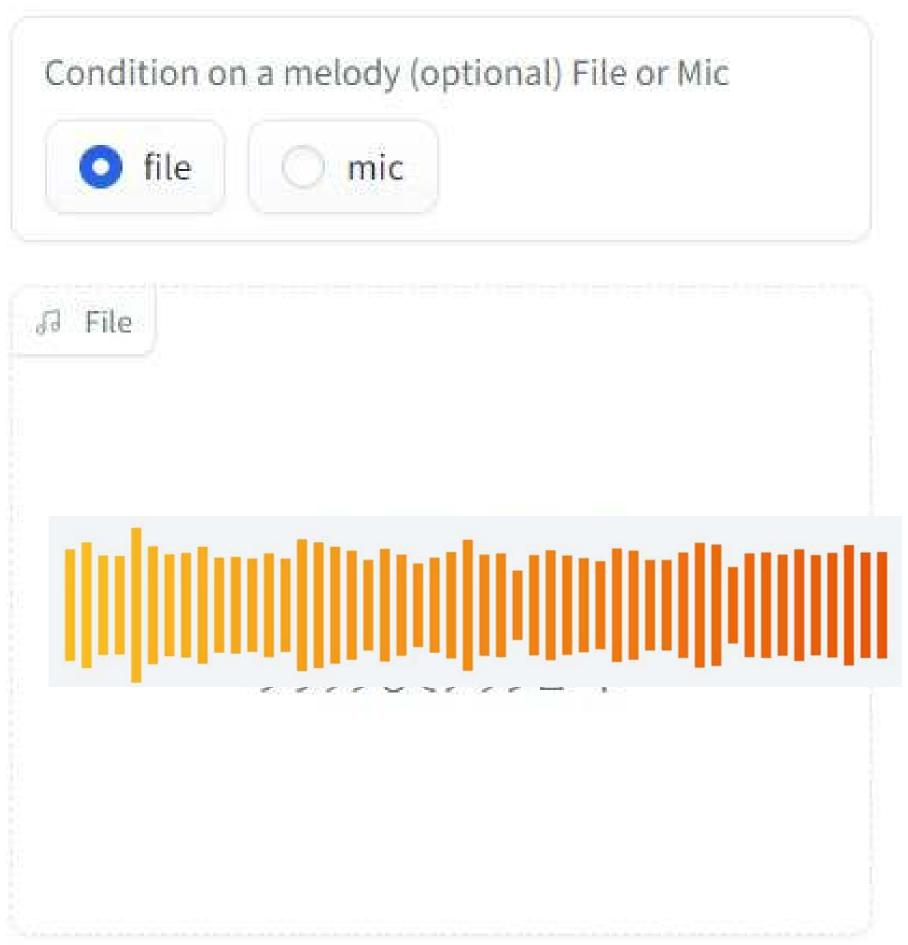
"DragYourGAN:InteractivePoint-basedManipulationontheGenerative ImageManifold", XINGANGPAN, et.al, SIGGRAPH'23



Generated: Leather backpack

newarc.ai





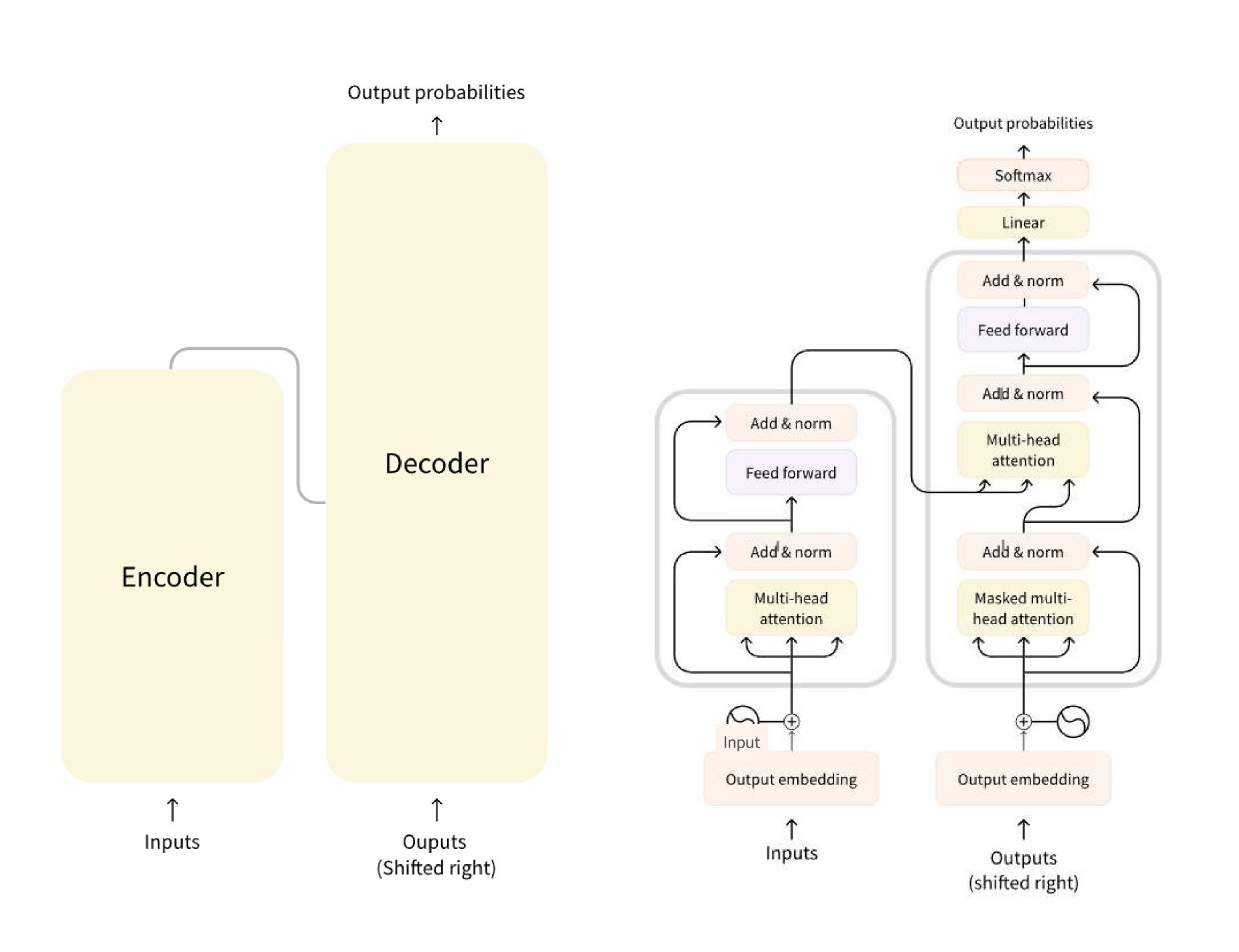
Generate

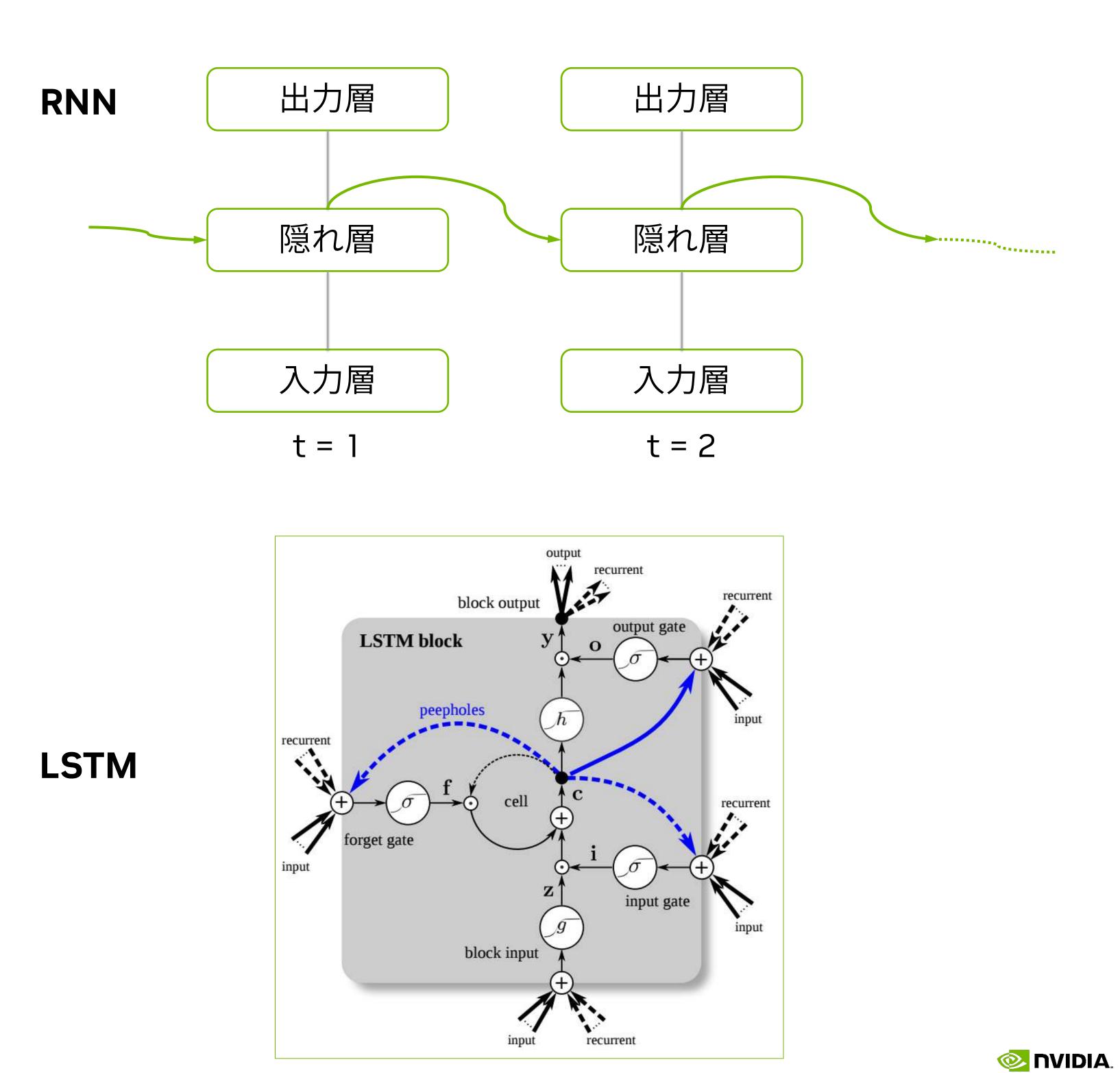
A grand orchestral arrangement with thunderous percussion, epic brass fanfares, and soaring strings, creating a cinematic atmosphere fit for a Samurai battle.

Music gen

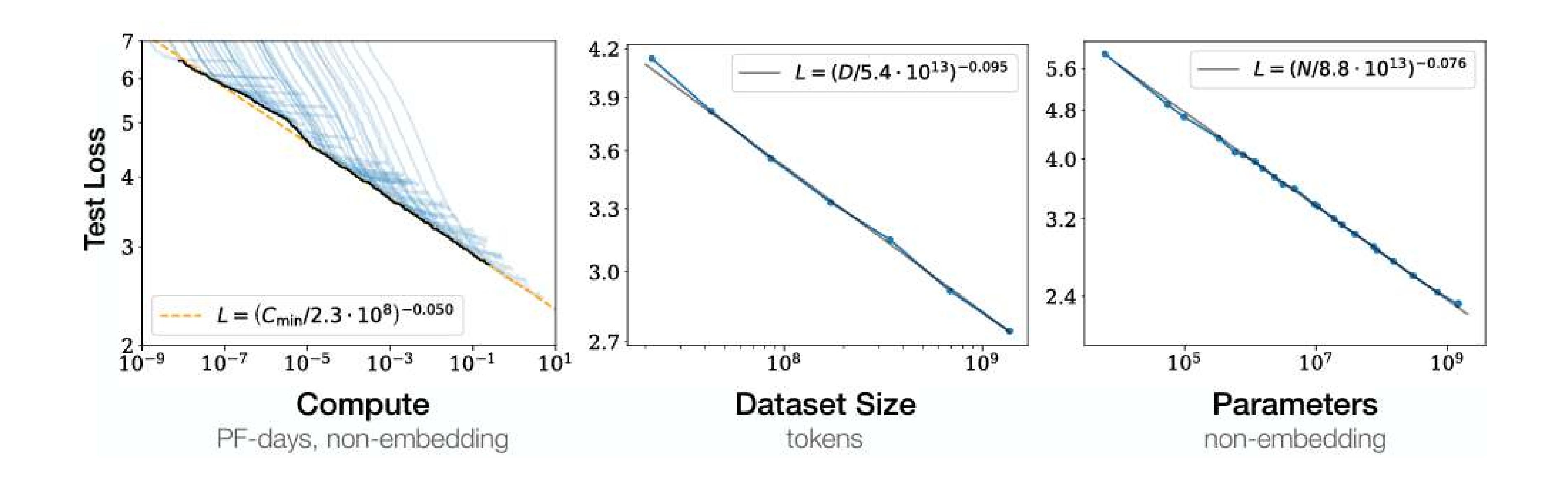
トランスフォーマー

アーキテクチャ



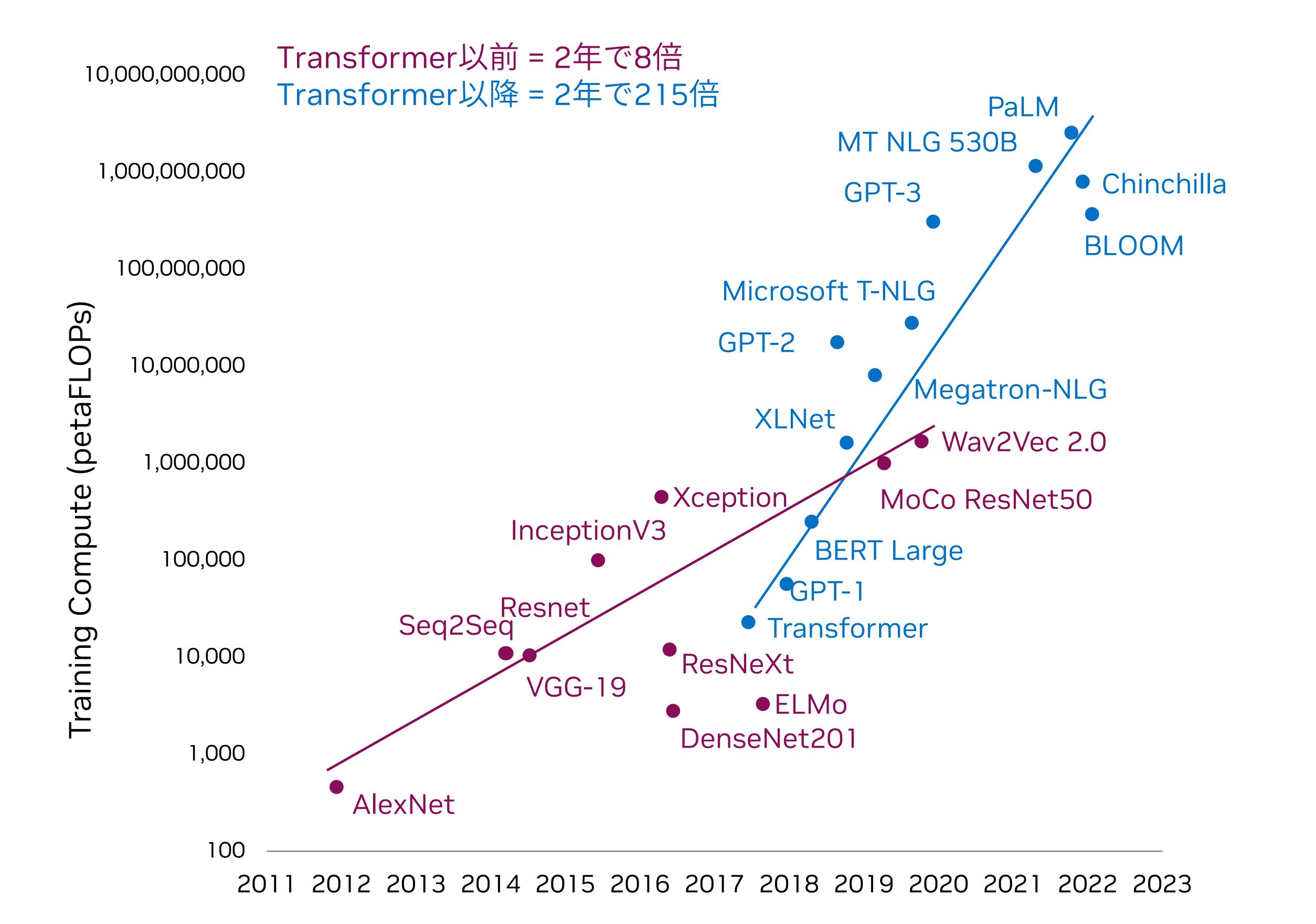


自然言語モデルのスケール則



指数関数的に増大するLLMモデルサイズ

LLMが必要な計算能力は指数関数的に伸びている



代表的なLLMと学習

モデル	発表時期	モデルサイズ (B)	事前学習データサイズ	ハードウェア	学習時間
OPT (Meta)	May-2022	175	180B Token	992 80G A100	
GLM (清華大学)	Oct-2022	130	400B Token	768 40G A100	60日
BLOOM (BigScience)	Nov-2022	176	366B Token	384 80G A100	105日
LLaMA (Meta)	Feb-2023	65	1.4T Token	2048 80G A100	21日
MT-NLG (MS/NVIDIA)	Jan-2022	530	270B Token	4480 80G A100	_



NVIDIA H100

世界のAIインフラを支える新たなエンジン

最高の AI/HPC 性能

4PF FP8 (6X) | 2PF FP16 (3X) | 1PF TF32 (3X) | 60TF FP64 (3X) 3TB/s (1.5X), 80GB HBM3 memory

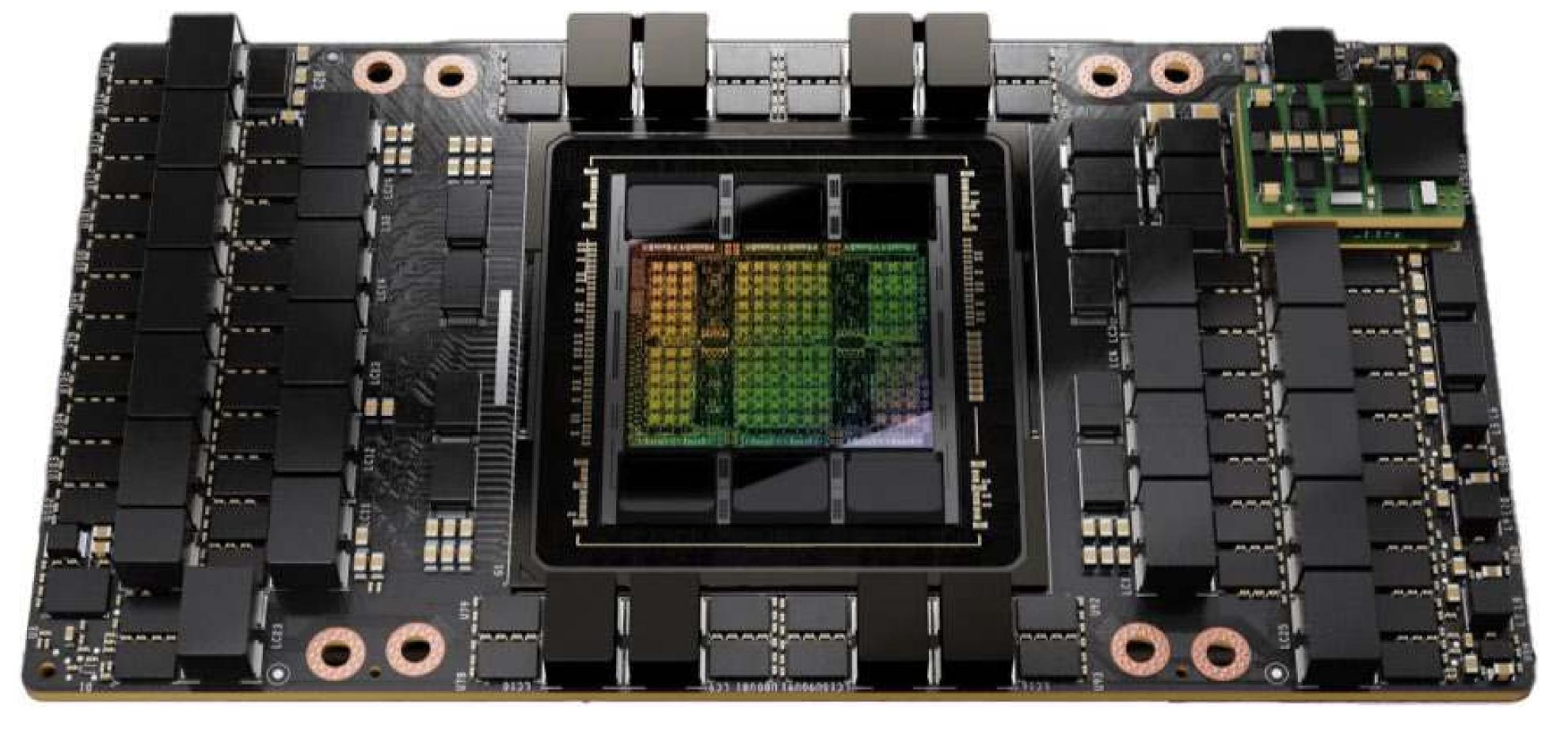
TRANSFORMER モデルへの最適化

6X faster on largest transformer models

高い稼働率とセキュリティ

7 Fully isolated & secured instances, guaranteed QoS 2nd Gen MIG | Confidential Computing

史上最速でスケーラブルなインターコネクト 900 GB/s GPU-2-GPU connectivity (1.5X) | 128GB/s PCI Gen5



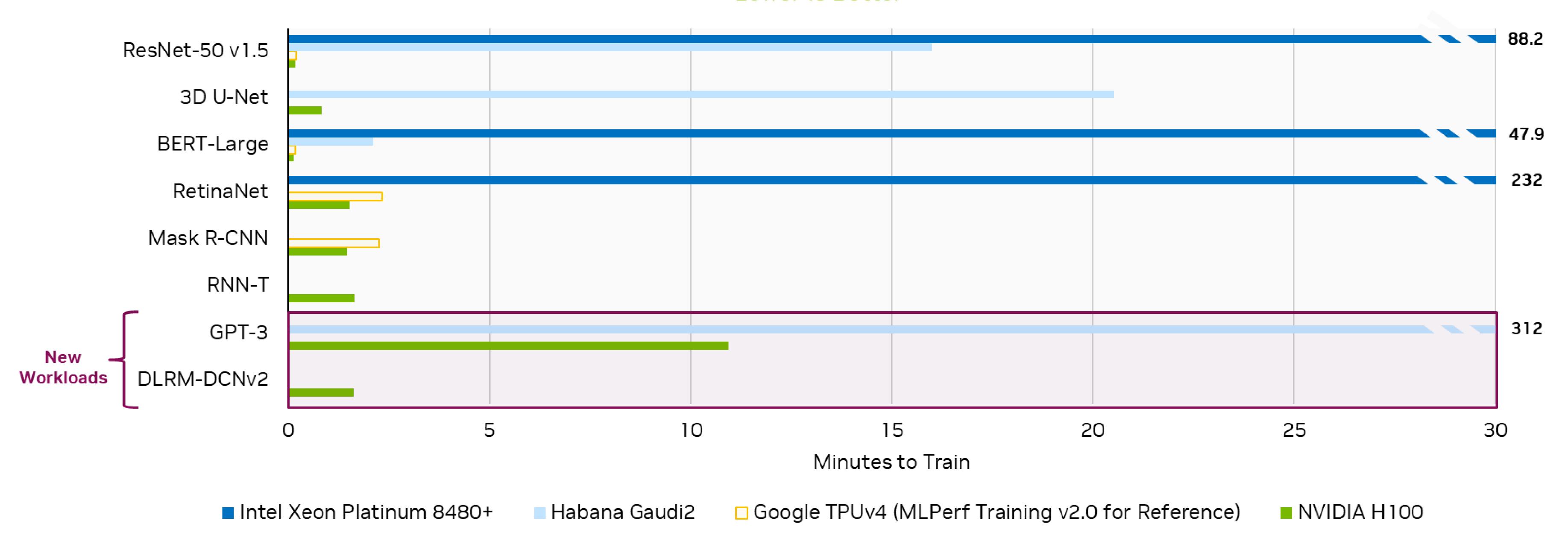
Custom TSMC 4N Process | 4.9 TB/s Total External B/W

MLPerf Training v3.0

全てのワークロードにおいて最速を記録

Time to Train

Lower is Better



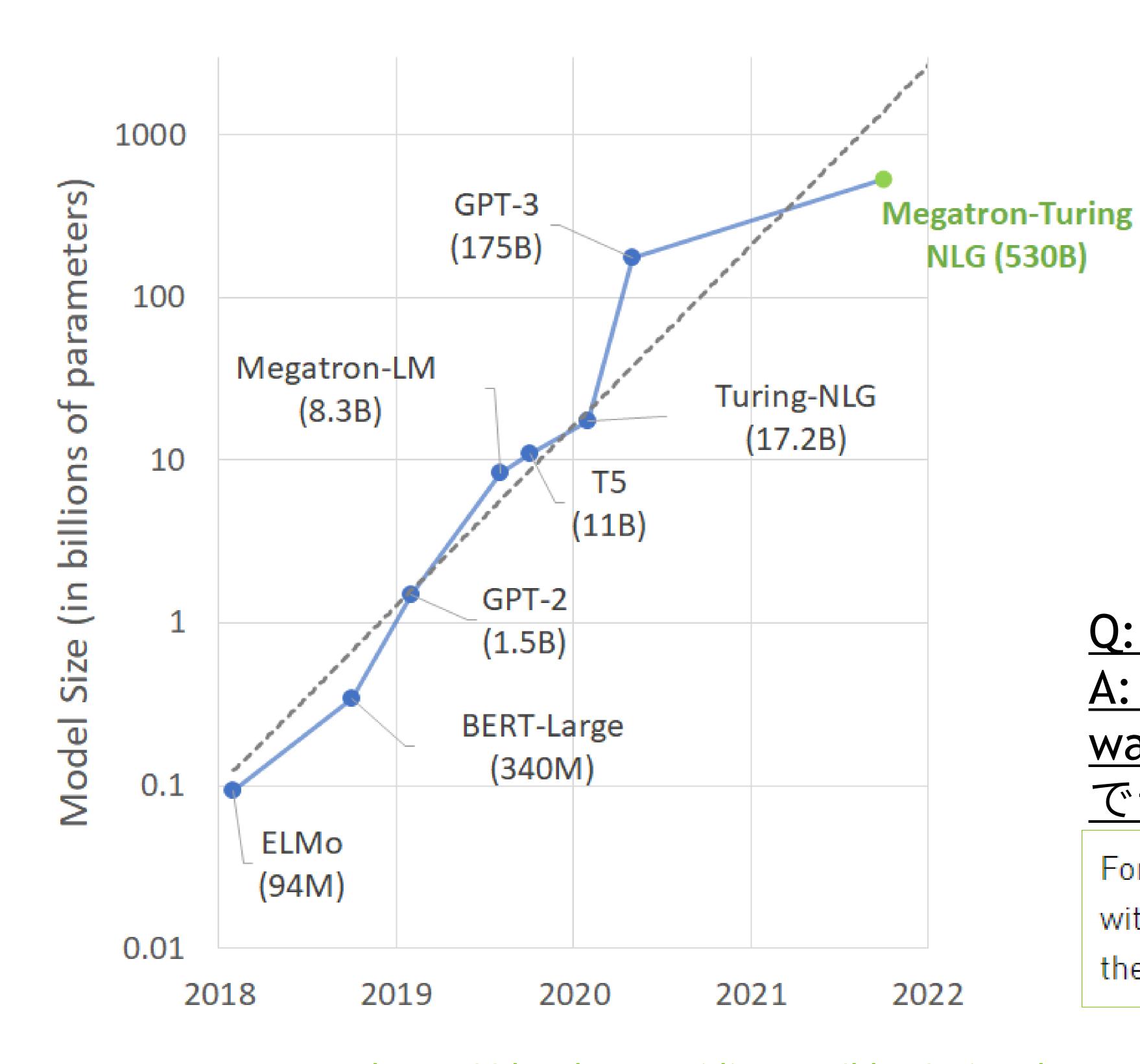
MLPerf Training v3.0. Fastest time to train on each benchmark by each submitter's platform

| Format: Chip count, Submitter, MLPerf-ID | ResNet-50 v 1.5: 3584x NVIDIA+CoreWeave 3.0-2002, 8x Intel-HabanaLabs 3.0-2011 | 3D U-Net: 768x NVIDIA 3.0-2075, 8x Intel-Habana Labs 3.0-2016 | BERT-Large: 3,072x NVIDIA+CoreWeave 3.0-2001, 64x Intel-HabanaLabs 3.0-2015, 32x Intel 3.0-2011 | RetinaNet: 768x NVIDIA 3.0-2011 | 2,048x Google 2.0-2010 | Mask R-CNN: 384x NVIDIA 3.0-2066, 2,048x Google 2.0-2010 | RNN-T: 512x NVIDIA 3.0-2070 | GPT-3: 3,584x NVIDIA+CoreWeave 3.0-2003, 384x Intel-HabanaLabs 3.0-2014 | DLRM-dcnv2: 128x NVIDIA 3.0-2065.



モデルがGPUに載らない?

分散学習の必要性



MT-NLGは530Bパラメータ

- 単純計算で 2,120 GB (in FP32)
- モデルをすべてメモリにロードするだけで、 8xA100 (80GB) サーバが、3 台強必要
- ワーキングメモリも当然必要

Q: 実際どう扱っているのか?

A: 280xA100 (=35 nodes) で 8-way tensor parallel & 35-way pipeline parallel のモデル並列&4480GPUのクラスタでデータ並列

For example, for the 530 billion model, each model replica spans 280 NVIDIA A100 GPUs, with 8-way tensor-slicing within a node and 35-way pipeline parallelism across nodes. We then use data parallelism from DeepSpeed to scale out further to thousands of GPUs.

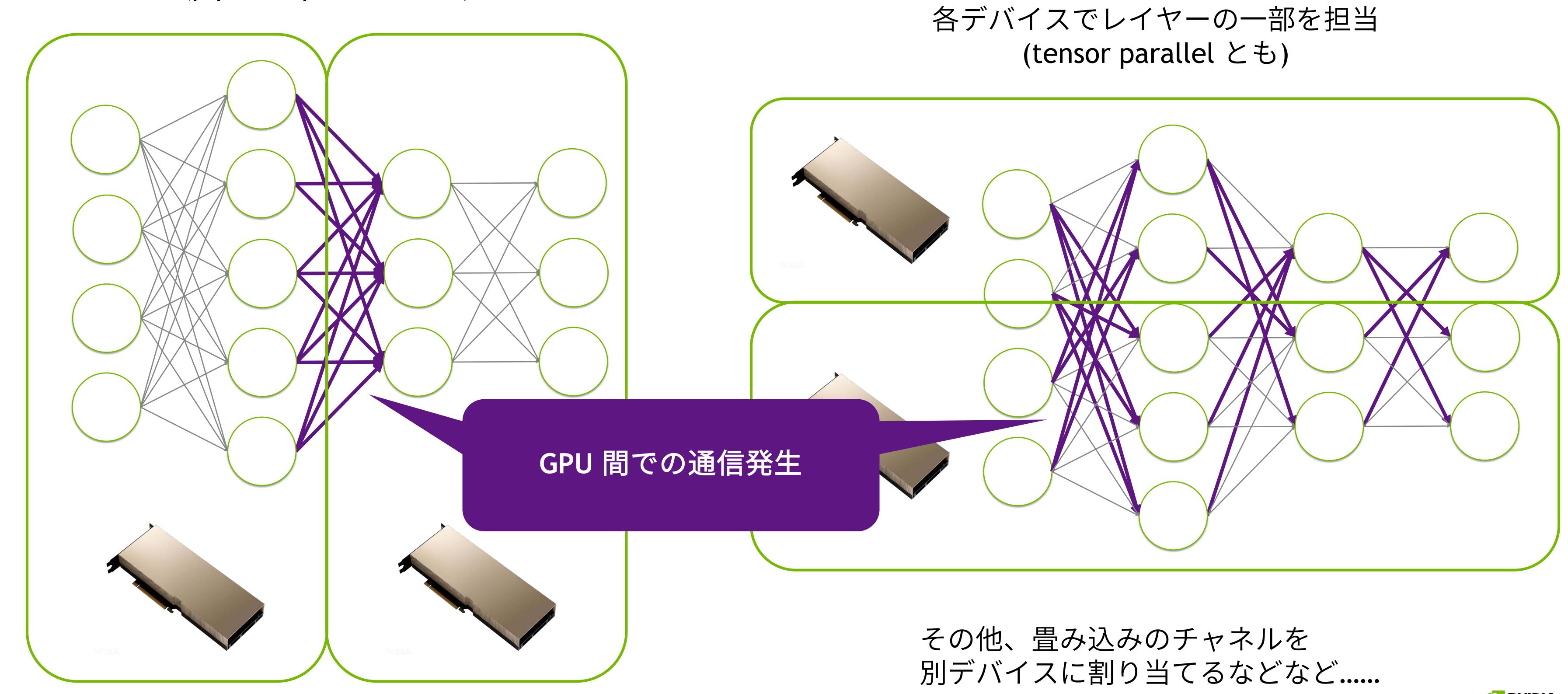
https://developer.nvidia.com/blog/using-deepspeed-and-megatron-to-train-megatron-turing-nlg-530b-the-worlds-largest-and-most-powerful-generative-language-model/



モデル並列の考え方

大規模なモデルを扱えるサイズに分割する

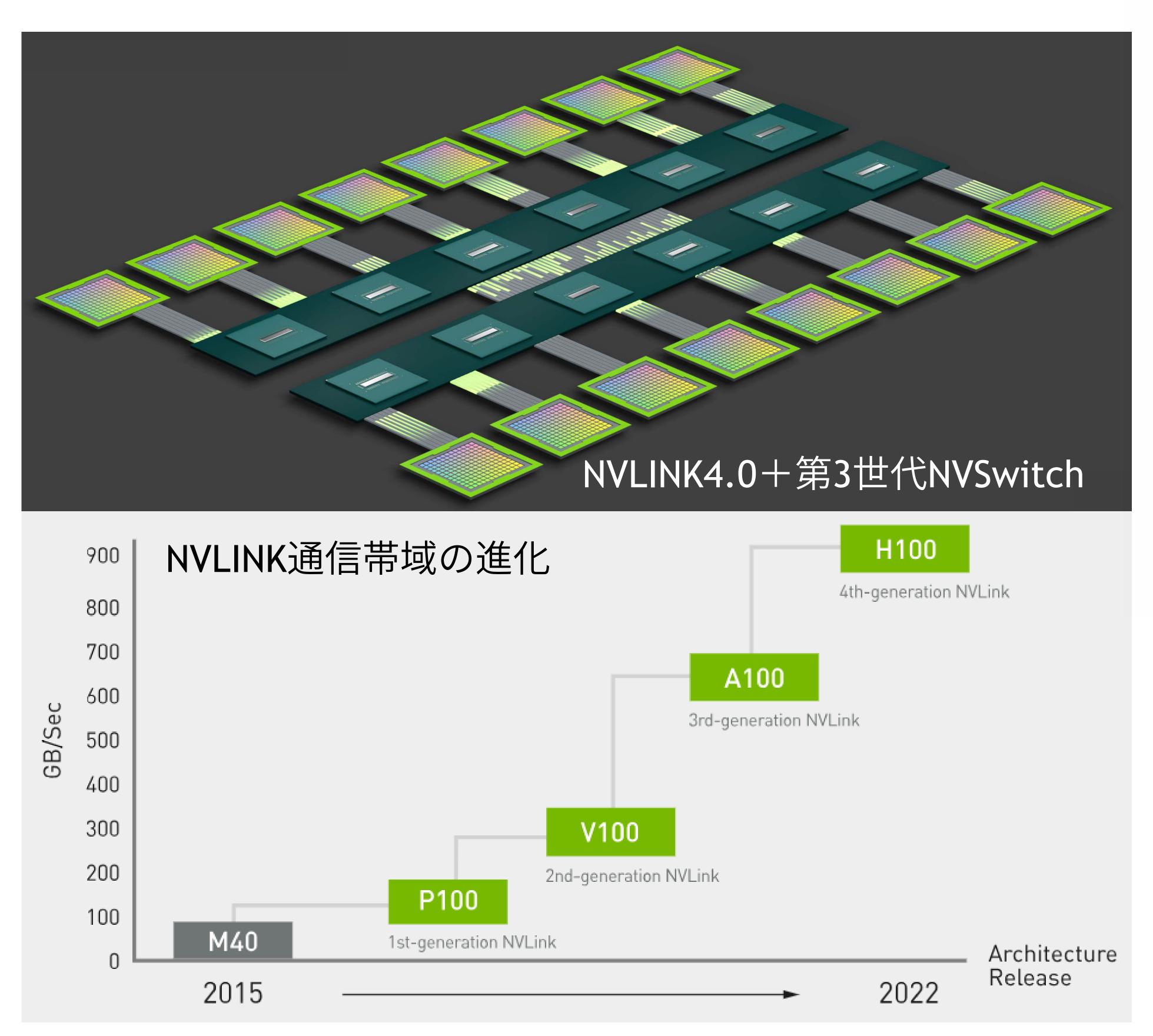
レイヤーごとに別デバイスへ割り当て (pipeline parallel とも)



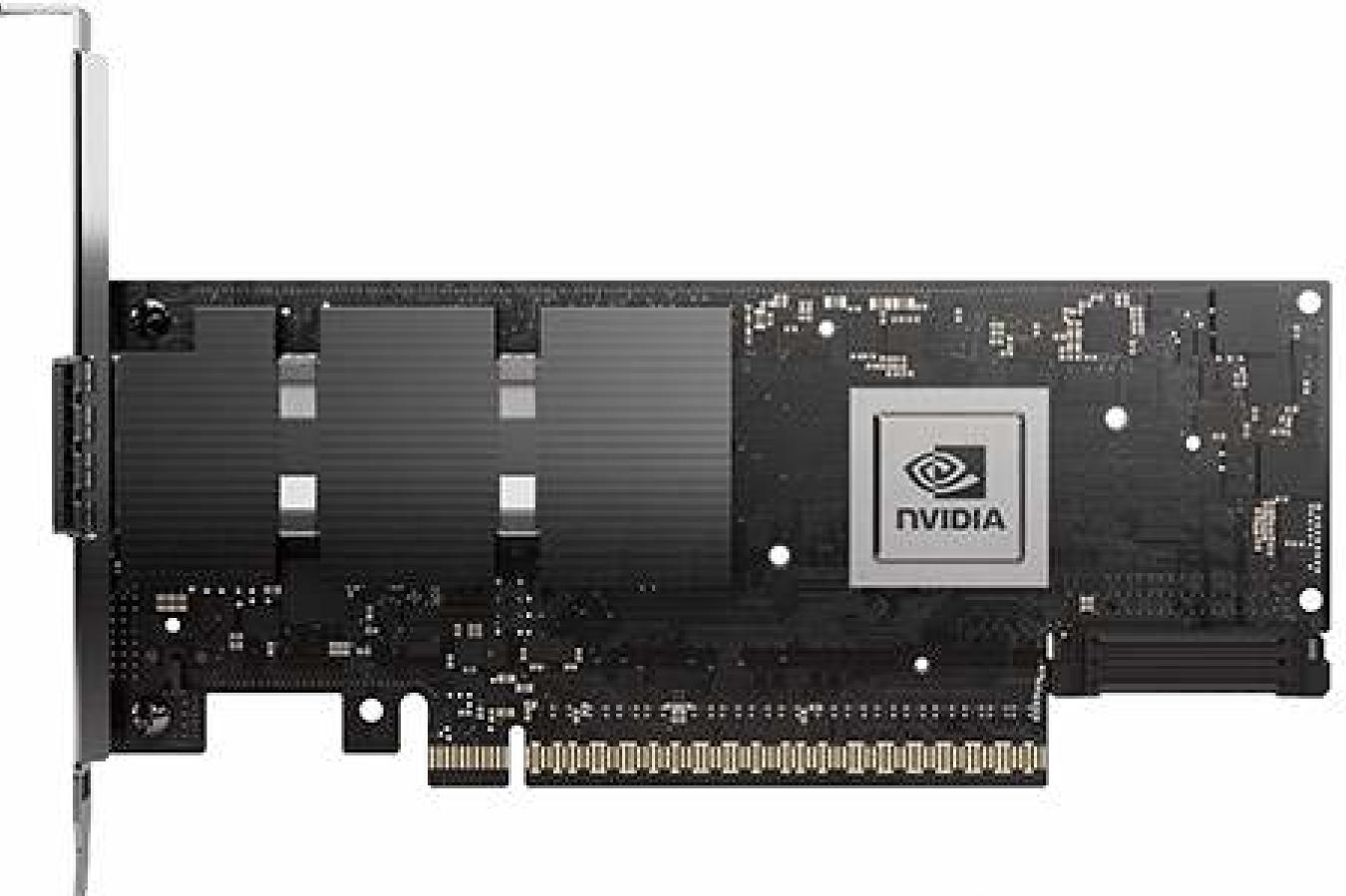


GPU-GPU インターコネクト/ノード間通信

ノード内



ノード間



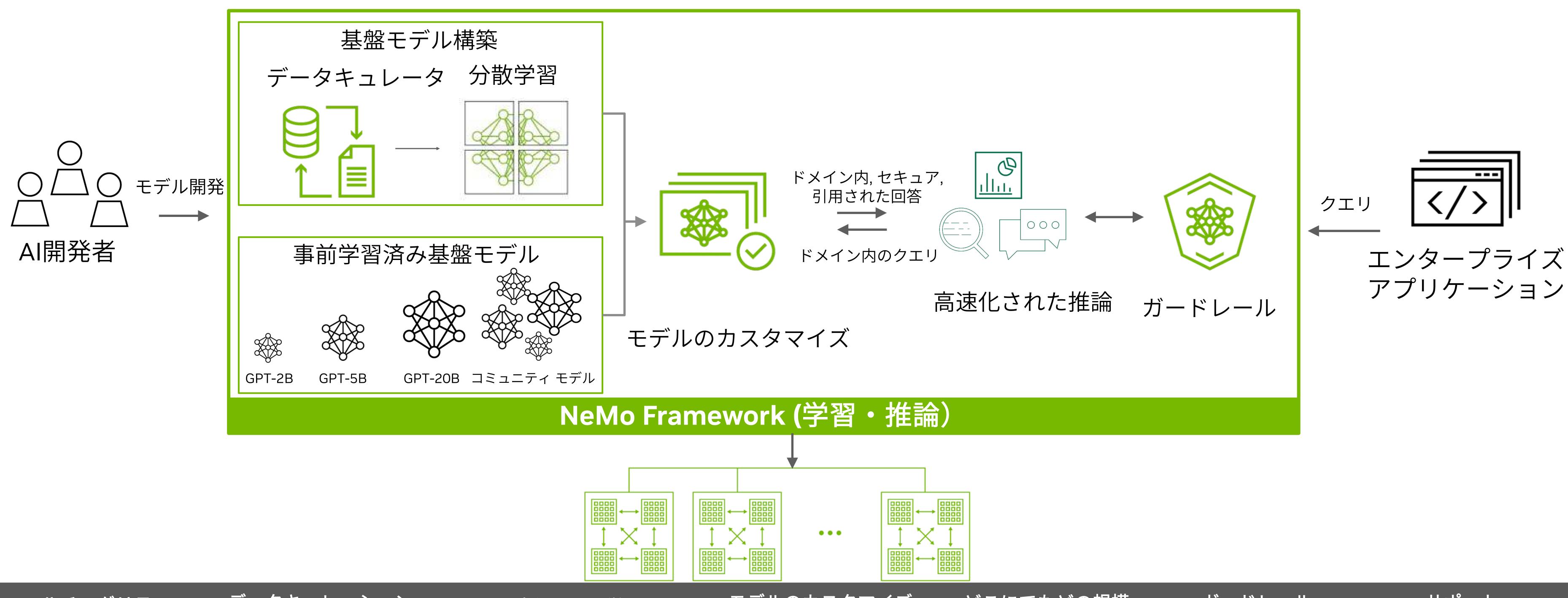
ConnectX-7 Infiniband アダプター

NDR 400Gb/s InfiniBand PCIe Gen5 最大x32 レーン GPUDirect® RDMA GPUDirect Storage In-Network Computing



NeMo Framework

会話、言語、画像の生成AIモデルを構築、カスタマイズ、デプロイ するためのエンド・ツー・エンドのワークフロー



マルチモダリティ 対応

言語、画像の生成AI モデル構築 データキュレーション @ Scale

抽出、重複、大規模 非構造化データの フィルタ情報, @ scale 最適化された学習

数千のノードで並列化 されたモデルと学習 データによる高速化された 学習とスループット モデルのカスタマイズ

P-tuning, SFT, Adapters, RLHF, AliBiに よるカスタマイズ どこにでもどの規模 にも展開可能

あらゆる場所に拡張 可能な最適化された 推論実行 ガードレール

安全、セキュリティ 要求を満たすアプリ ケーションをNeMo ガードレールで実現 サポート

NVIDIA AI Enterprise と顧客専門家によりプ ロジェクトをオント ラックに

NeMo Framework 性能 – 学習

	3000億トークンを学習するための日数 (A100) – BF16					
	800 GPU (100 DGX A100)	480 GPU (60 DGX A100)	160 GPU (20 DGX A100)	64 GPU (8 DGX A100)		
GPT-3: 126M	0.07	0.12	0.37	0.92		
GPT-3: 5B	0.8	1.3	3.9	9.8		
GPT-3: 20B	3.6	6	18.1	45.3		
GPT-3: 40B	6.6	10.9	32.8	82		
GPT-3: 175B	28	46.7	140	349.9		



