

ふじさん

fujijoho group monthly magazine

12

~ 2025年指針 ~

△富士情報

開 点

[今月のひとこと]

TPU GPU

- ・国字
- ・冬至



水かけ菜の収穫

写真提供：都留市 産業課



今月のひとこと

TPU GPU

社長 渡辺直企

11月19日にGoogleが新しいAI、Gemini3をリリースしました。このAIはほとんどすべてのテストにおいてChatGPT5.1、Grok4.1などの競合モデルを上回るスコアを記録しています。これを受けたOpenAIのサム・アルトマンは同日「Congrats to Google on Gemini 3! Looks like a great model.」とお祝いをツイートしていました。12月2日にウォール・ストリート・ジャーナルはサム・アルトマンが社内向けに「コード・レッド（緊急事態）」宣言したと報じています。その後12月11日にOpenAIは性能を知的労働で専門家を超えるレベルまで大幅強化したGPT5.2を発表しています。

OpenAIをはじめMetaのLlamaやxAIのGrokなどほとんどのAIはNVIDIAのGPU(Graphics Processing Unit)を使用しています。その名の通り元々は3Dグラフィックスの描画を高速に処理するために作られ、CPUに比べ膨大な並列処理を得意としています。2009年スタンフォード大学のアンドリュー・ン教授は、機械学習においてCPUの代わりにGPUを使用することで1億パラメータのニューラルネットワークの学習に必要な時間を数週間から約一日に短縮できると発表しました。このとき使用しているのは白黒32x32ピクセルの画像でした。

機械学習では用途別に画像認識向けのCNN(Convolutional Neural Network)、音声認識、株価予測、機械翻訳など時系列データ向けのRNN(Recurrent neural network)、SNSの人間関係、交通網などグラフ(網)学習向けのGNN(Graph neural networks)など多くのニューラルネットワークを開発していました。2017年GoogleはTransformerアーキテクチャを発表しました。TransformerはRNN同様言語など時系列データを扱いますが、並列処理を強化し、データ全体を俯瞰することで、RNNが苦手だった長文理解や長期記録が可能となり、汎用性を実現しました。また、言語のみならず画像なども処理が可能となりました。TransformerはLLM(Large Language Model)の基礎となっています。

OpenAIは2018年にTransformerを予測機能(生成)に特化した生成AI、GPT-1(Generative Pre-trained Transformer)を発表しました。OpenAIをはじめ多くの生成AIベンダーはNVIDIAのGPUで生成AIを開発するためにCUDA(Compute Unified Device Architecture)という開発環境を使用しています。一方、Googleは2013年から独自で開発しているTPU(Tensor Processing Unit)を使用しています。Tensorは行列を多次元に一般化した概念で、TPUはニューラルネットワークを行列とみなした計算に特化したプロセッサです。GPUは画像処理向けですので4x4などの小規模な行列演算を得意としています。GPUで生成AIの大規模な行列演算する場合は、小さい行列演算に分解して計算する必要があります。計算のたびにメモリを使用する必要があります。一方、TPUは256x256という大規模な行列演算をメモリの使用を最小限に押さえて処理可能なので高効率になります。生成AIの処理に関してはTPUはGPUよりも1.2~2倍程度、消費電力に優れています。GPUはTPUに比べ汎用的で、処理できる機能も多い一方、生成AI向けの半導体という観点では効率性に課題があります。GPT-4やGemini2.5のパラメータは1兆を超え、博士号レベルの知識を得るに至っています。GPT-5やGemini3のエージェント機能の拡充で自律的に仕事を任せられるようになってきました。今後、一層便利になるので、積極的に活用していくようにしていきたいと考えています。