



Session Report

AI へのニーズが増大する中で
継続的な監視と評価で精度を維持する

生成 AI 時代の MLOps 実現方法とは？

Google Cloud

カスタマー エンジニアリング データ アナリティクス スペシャリスト
高村 哲貴

Google Cloud

セッションレポート概要

生成 AI モデルのカスタマイズは、AI/ML モデルをゼロから構築するより簡単かもしれませんが、その運用においては共通の課題が多く見られます。ここでは MLOps パイプラインの最適化や AI モデルの管理、評価、セキュリティ フィルターの活用などを通じて、生成 AI 時代に MLOps をどのように適応させるかを解説します。

プレゼンター紹介



Google Cloud

カスタマー エンジニアリング

データ アナリティクス スペシャリスト

高村 哲貴

日系 Sier 企業にて、データベース エンジニア、クラウド アーキテクトを経て、2021 年より現職。現在は、Google Cloud のデータ アナリティクス領域のスペシャリストとして、お客様のデータ活用を技術観点から支援。

目次

- 生成 AI によって変わる MLOps 3
- 従来の MLOps と生成 AI 時代の MLOps の違い 4
- AI モデルのチューニングの重要性 7
- 新しいアーティファクトの管理が必要に 9
- 生成されたアウトプットの評価 & モニタリング 10
- エンタープライズ データとの接続 11
- Google Cloud の生成 AI 評価サービス 13

生成 AI によって変わる MLOps

Google Cloud では、MLOps のことを「ML（機械学習）システムを迅速かつ確実に構築、導入運用するための標準化されたプロセスと機能のセット」と表現しています。MLOps には、一般的なシステム運用とは異なる特有の難しさがあります。それが AI の精度の維持です。

例えば、何かを予測する AI モデルを活用している場合、事前に準備したデータで学習モデルを作り、精度を評価したうえで問題がないものをリリースする手順が一般的です。しかし、データは日々刻々と特徴を変えるため、常に高い精度が維持できるとは限りません。MLOps では、こうした変化に対応しながら一定の精度を保ちながらビジネス価値につなげる必要があります。

では、これらの課題に対して生成 AI はどのような変化をもたらすでしょうか。重要な 2 つのポイントに、「既存の MLOps への投資を捨てる必要はない」「生成 AI 特有のニーズを理解する」が挙げられます。後者の生成 AI 特有のニーズについて、具体的には以下の通りです。

1. マルチタスク モデルやプロンプト、増大する AI インフラストラクチャ ニーズへの対応

生成 AI のモデルは、テキスト生成や画像生成など用途に応じて非常に多くのモデルが存在します。このようなマルチタスクなモデルに対してどれを選ぶか、そしてどのようなプロンプトを与えていくかが大切です。

2. チューニングと厳選されたデータによるカスタマイズ

既存の大規模言語モデル（LLM）の活用だけではなくて自社のデータや、特定のドメインのデータをつなぎ合わせてどのようにチューニングやカスタマイズするかが問われます。

3. 新しいアーティファクトの管理

プロンプト、エンベディング、アダプター層など、生成 AI によって誕生した新しいデータや概念が加わります（ここではアーティファクトと呼びます）。これらを適切に管理しなければなりません。

4. 生成された出力の評価と監視

生成 AI によって生成されたアウトプットに対して、どのように評価・監視していくかを考えなければなりません。

5. エンタープライズデータとの接続

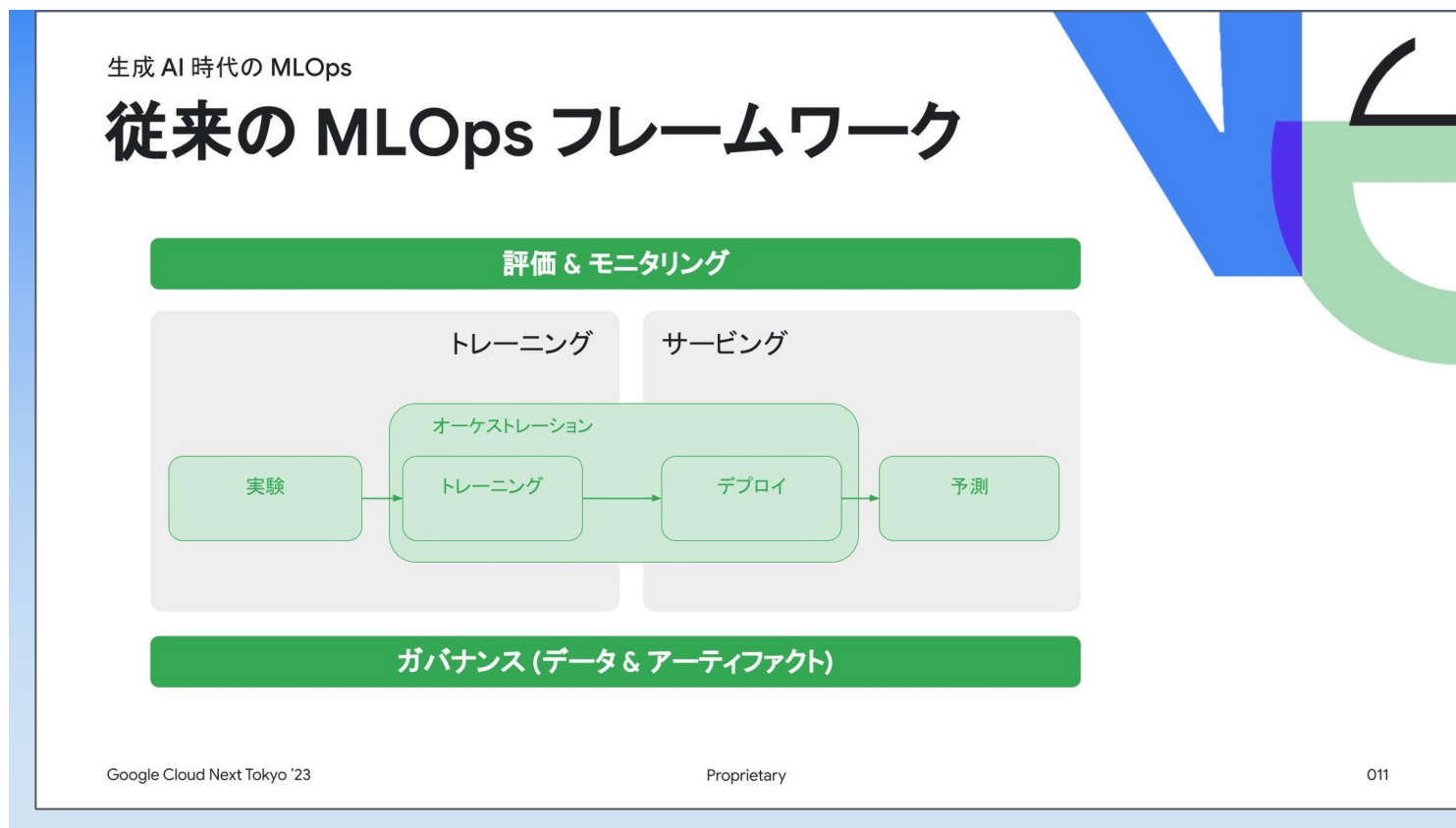
ハルシネーション（事実に基づかない誤った情報を AI が生成してしまうこと）を回避する手段として、自社のエンタープライズ データに接続したいというニーズがあります。自社データと LLM をつなぐことで、より特定のドメイン知識を持ったり、最新の情報を踏まえた回答やアクションを AI が行えます。こうした特徴が、今までの MLOps と異なる点です。

このようなニーズに対して、Google Cloud の AI ソリューション「Vertex AI」は予測 AI (Predictive AI) だけではなく、生成 AI も含めた統合プラットフォームとして、生成 AI 向けの固有のニーズに対応するためにアップグレードを続けています。その具体的な機能は次項から解説していきます。

従来の MLOps と生成 AI 時代の MLOps の違い

生成 AI 時代の MLOps は、従来の MLOps と具体的にどのように変わってくるのでしょうか。

下図は従来の MLOps のフレームワークです。簡易的ですが、実験・トレーニングを行い、次にデプロイして予測します。それらを支えるガバナンスのレイヤーと、それを評価・モニタリングするレイヤーが存在します。

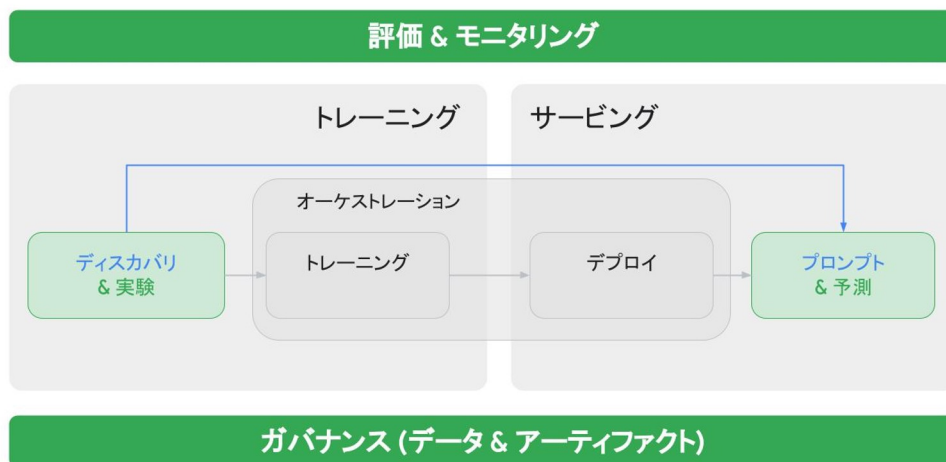


従来の MLOps フレームワーク

この既存の MLOps のフレームワークに対して、生成 AI によって変化したものを示したのが下図です。生成 AI における MLOps を考える上で必要な概念として、まずディスカバリ、プロンプトの2つがあることが分かります。

生成 AI 時代の MLOps

マルチタスクモデル & プロンプト



Google Cloud Next Tokyo '23

Proprietary

012

「ディスカバリ」と「プロンプト」に従来との違いが見られる

「ディスカバリ」とは、利用可能な LLM の中から、それぞれのユースケースに合うモデルを探し出すことを示します。「プロンプト」については、生成 AI が返すデータの精度を向上させるためにインプットをどのように扱い管理すればよいかということなのです。これに対して Google Cloud では、それぞれに対応する以下の2つのソリューションを提供しています。

1. Vertex Model Garden

Vertex Model Garden の中には、ファーストパーティのモデルだけではなく、サードパーティも含めた非常に多くのモデルが登録されています。例えば、Llama2 や Stable Diffusion などのモデルも含めてそれぞれのユースケースに合うモデルを探することができます。

2. Vertex Generative AI Studio

自身が選んだモデルに対して、どのようなプロンプトを与えたことでどのようなアウトプットがされるのかを確認することも、スムーズなモデル選定に欠かせません。Vertex Generative AI Studio は、非常にシンプルな操作と GUI の画面によって、選定したモデルに対して、好きなプロンプトを入力し、そのアウトプットを確認するという行為を簡単に行えます。

生成 AI 時代の MLOps

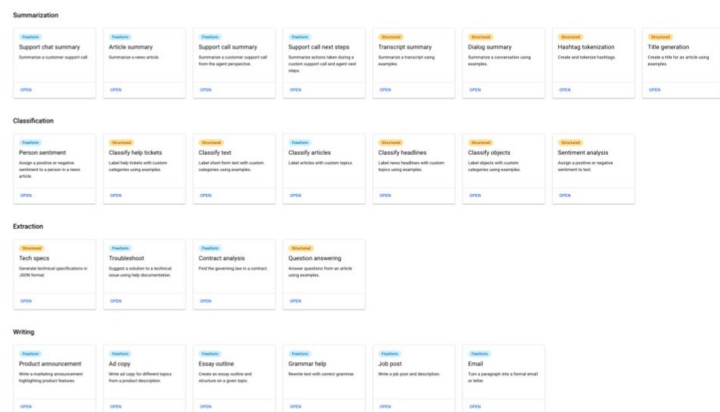
事前学習モデルの選定とプロンプト

Vertex Model Garden

利用可能な事前トレーニング済みのオープン基盤 AI モデルを見つけて試す

Vertex Generative AI Studio

明確かつ簡潔で有益な指示を含む入力プロンプトを慎重に作成する



生成 AI ニーズに対応した AI インフラストラクチャ
NVIDIA GPUs A100, L4 and H100 | TPU v3, v4 and v5e

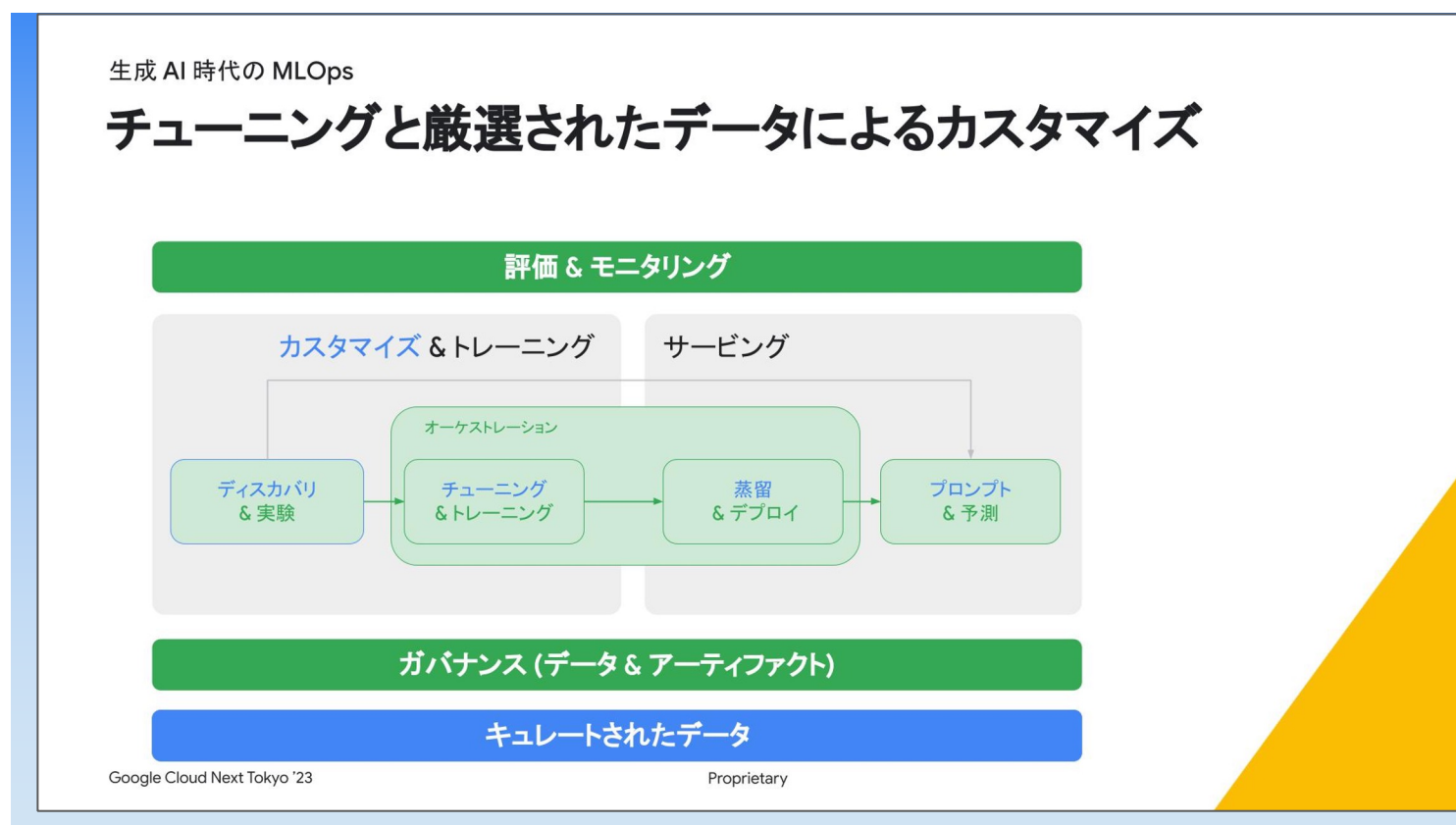
Proprietary

013

LLM の選定とプロンプトの管理を簡便化する2つの機能

AI モデルのチューニングの重要性

もちろん、生成 AI によって MLOps に生じる変化は上記だけではありません。このほか注目したいのが、下図に示す「チューニング」「蒸留」「キュレートされたデータ」というキーワードです。蒸留とは、大規模なモデルを教師モデルとしてその知識を小さいモデルに学習させることを意味し、「モデル圧縮」とも呼ばれます。



生成 AI ならではの MLOps で押さえておきたいチューニング、蒸留、キュレートされたデータ

この中のチューニングに関して、Google Cloud では、先述の Vertex Generative AI Studio を用いて 2 種類の方法で実現することができます。

1つ目は教師あり学習です。ファイン チューニングと呼ぶ場合もあります。各々が持っている特定のタスクに特化させるために必要なデータを学習させることで、より特定の知識を持った、出力の振る舞いを変えたりすることができます。

もう1つが人間のフィードバックからの強化学習（Reinforcement Learning Human Feedback：RLHF）です。生成 AI 系のアプリを使う際によく目にすると思いますが、回答に対してサムズアップ・サムズダウンを選択したり、または 2 種類の回答が出てきた際に、どちらが好みかを選択したりする行為です。そのモデルをより人間の好みや価値観に近づけるチューニング手法です。

生成 AI 時代の MLOps

チューニングと厳選されたデータによるカスタマイズ

Vertex Generative AI Studio によるチューニング



教師あり学習 Supervised Tuning

特定のタスクのパフォーマンスを向上させたり、指示が不十分な場合に特定の出力要件を遵守したりするのに役立ちます。

人間のフィードバックからの強化学習 Reinforcement Learning Human Feedback (RLHF)

好み（サムズアップ/ダウンなど）と一貫した報酬を予測し、モデルの動作を好みに合わせます

データキュレーション

一般 & ドメイン コーパス | タスクの組み合わせ | 専門的なタスク | 人間のフィードバック

Proprietary

015

Vertex Generative AI Studio を用いた 2 つのチューニング方法

データ キュレーションの側面では、例えば、特定のドメインに特化した学習をさせたければ、そのドメインのコーパスとなるデータや専門的なタスクのデータ、RLHF のための人間のフィードバックのデータなどを準備して学習させることでチューニング可能です。

新しいアーティファクトの管理が必要に

生成 AI における MLOps では、新しいアーティファクトを適切に管理していかなければなりません。その具体的な対象は、チューニングのジョブや、エンベディング、チューニング済みモデルなどです。

Google Cloud では、これらを管理・実行するためのソリューションを提供しています。例えばチューニング ジョブに関しては、「Vertex Pipelines」を活用します。これは従来の MLOps から使用されているソリューションですが、生成 AI のチューニングに対しても有効です。

次に、チューニングされたモデルのラインアップやバージョン管理には「Vertex Model Registry」が有効です。モデルの管理のほか評価までを手軽に実行できます。

最後に、エンベディングに関しては、「Vertex Feature Store」を活用できます。エンベディングされたデータの保存・管理・展開を容易に行えます。

生成 AI 時代の MLOps

新しいアーティファクトを管理する

ジョブ、エンベディング、アダプター層のチューニング

生成 AI チューニング ジョブの調整と管理を **Vertex Pipelines** で行い、データセットからモデルまでのリネージを追跡する



生成 AI のチューニングモデルを **Vertex Model Registry** で管理、評価メトリクスへのアクセスとワンクリックのデプロイ

Type	Source
Imported	Model Garden
Imported	Custom training
Tabular	AutoML training
Large model	Generative AI Studio

エンベディングの保存、管理、展開を **Vertex Feature Store** で行い、他の ML Data と合わせて管理する

Feature Table				CREATE
All	307	Search		
		Name	Feature type	Sour
Feature Table				
User	17	FeatureTable1	--	--
Destination	8	text_snippet	string	1234t
Fare	11	text_snippet_e	Embedding	1234t

Google Cloud Next Tokyo '23

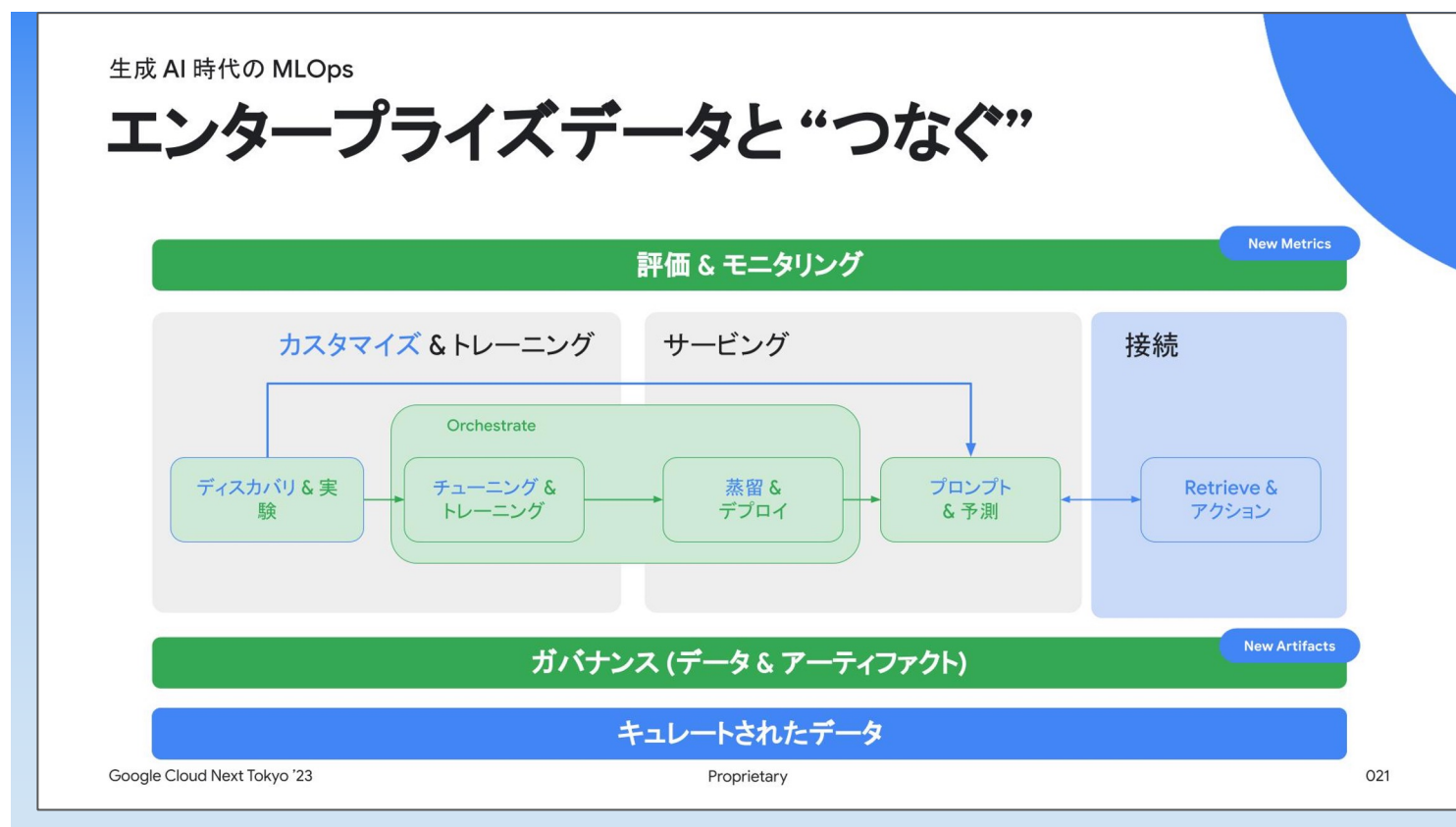
Proprietary

017

新しいアーティファクトと対応する Google Cloud のソリューション

エンタープライズ データとの接続

生成 AI ならではの MLOps のフレームワークとして、最後に言及したいのが、下図の「接続」のプロセスにある「Retrieve & アクション」です。



エンタープライズ データとの接続を考慮する必要がある

先述したように、社内の情報検索などで生成 AI を使えるように、自社固有のデータと LLM を接続したいという強いニーズが多く企業から寄せられています。これを実現するにあたり Google Cloud では 3 つの観点が必要だと考えています。

(1) エンベディングとベクトル検索

テキストや画像などマルチモーダルに対応したエンベディングの仕組みをもつことで、より高度化されたセマンティック検索が可能になります。さらにベクトル検索の技術によって低レイテンシーかつスケーラブルな検索を実現できます。Google Cloud では、具体的にこのためのソリューションとして「Vertex Embeddings APIs」を提供しています。

(2) グラウンディング

グラウンディング（Grounding）とは、生成 AI が外部の情報を文脈として理解し、その情報をもとに正確な結果を返す手法です。Google Cloud でも、ファーストパーティのデータだけでなく、サードパーティのデータと接続してグラウンディングに関する技術のケイパビリティを高めようとしています。

(3) エクステンション

これは、例えば LLM の回答やチャット、外部の予約システムなど複数のシステムを組み合わせることでシームレスなユーザー体験を提供するというように、現実世界のアクションと AI を接続しようとする考え方です。

Google Cloud の生成 AI 評価サービス

最後に、Google Cloud の生成 AI 評価サービスについて紹介します。生成 AI の正しさや正確さ、安全性が評価できなければ、お客様のビジネスに直接大きなインパクトを与えることとなります。それゆえに「評価は MLOps ライフサイクル全体のコア」といっても過言ではありません。正しく評価やモニタリングを行いながら、生成 AI が本番環境で稼働する状態を作ることが、MLOps を考える上で、最も重要なポイントとなります。

Google Cloud では、生成 AI の評価に関して、「Automatic Metrics」「AutoSxS」「Safety Bias」の3つのサービスを発表しています。それぞれを以下にて解説します。

Google Cloud は 3つの生成 AI 評価サービスを発表

1

Automatic Metrics

rougeLSum	0.1345
bleu	0.763

リファレンスデータに基づいてタスク特化型のメトリックスでモデルのパフォーマンスを評価

- 高速 & 効率的
- 学术界で使われる標準的な方法と多くのオープンベンチマーク

Google Cloud Next Tokyo '23

2

AutoSxS

model_a_win_rate	0.35
model_b_win_rate	🏆 0.65

2つのモデルのパフォーマンスをarbiterモデルで比較する

- 間の評価による貴したパフォーマンスを実現しながら、より速く、より安価に、オンデマンドで利用可能

Proprietary

3

Safety Bias

Safety attribute: Toxicity	
GenderID: female	0.89
GenderID: male	0.12
GenderID: transgender	0.02
GenderID: non-binary	0.67

モデルの安全性プロファイルが特定のアイデンティティグループに対して偏っているかどうかを理解

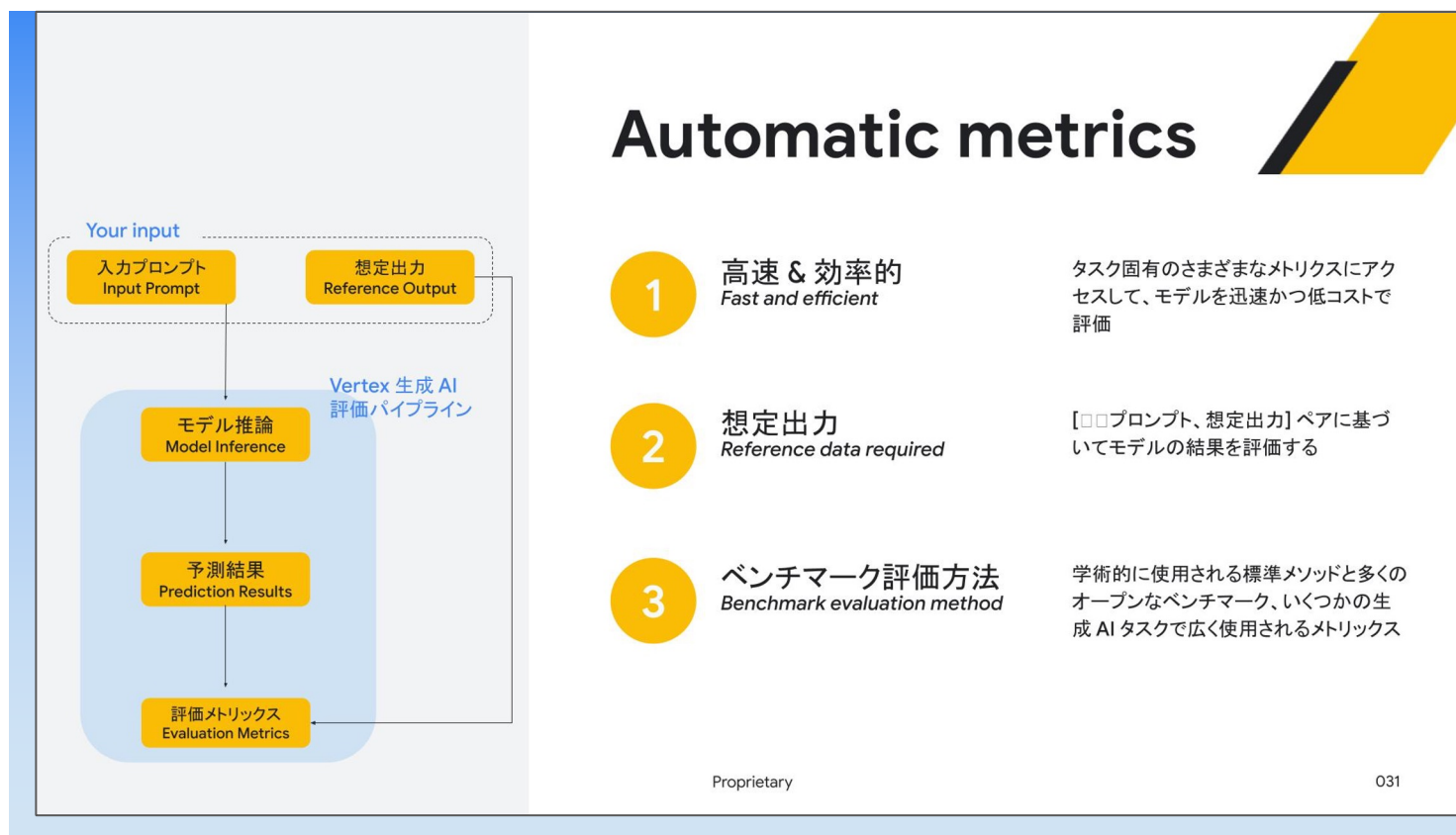
030

生成 AI を評価する 3 つのサービス

(1) Automatic Metrics

さまざまな生成 AI のタスクやその固有のタスクに最適な評価のメトリックスを、Google Cloud が自動で選定します。これにより AI モデルを定量的に評価することが可能です。例えば、要約なら「ROUGE-LSum」、またはテキスト生成なら「BLEU」など、適切なメトリックスの選定とそれに対する定量評価を自動で実行します。

具体的には、入力にはプロンプトとグラウンド トゥールズ（正解データ、想定している出力のデータ）を与えます。そうすると、Google Cloud の中でモデルの推論と予測結果、またそれに対する評価メトリックスのスコアリングが提示されます。

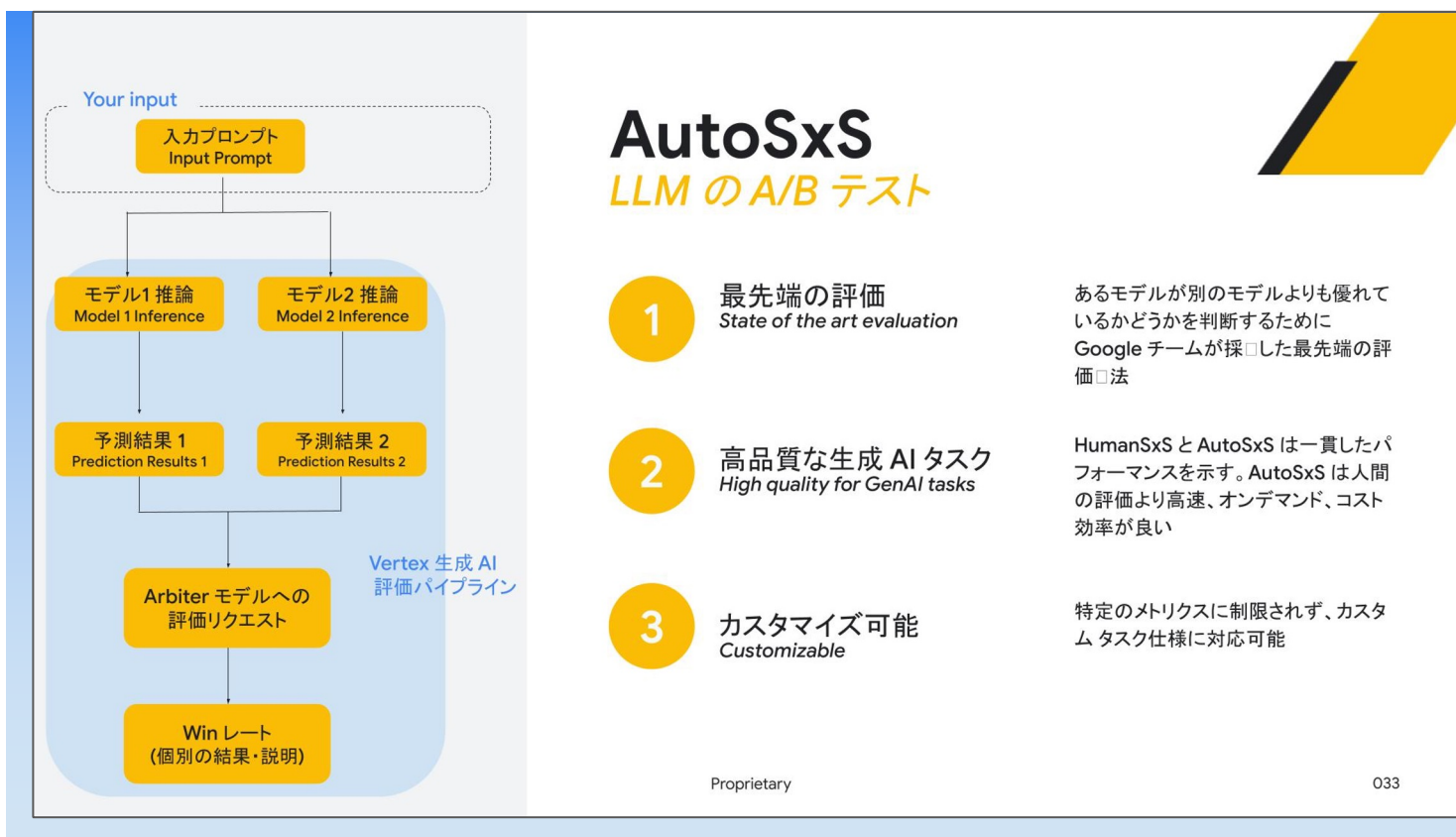


複数存在するメトリックスから適切なものを自動で選定してくれる

(2) AutoSxS (オートサイドバイサイド)

これは、ABテストのような形でモデルA・モデルBのうちどちらの精度が高いかを提示するサービスです。

こちら、まずは入力としてプロンプトを与え、さらに評価したいモデルを2つ選びます。各モデルに対して推論と予測結果が自動的に算出されるので、次に Arbiter モデルという当社が持つ第3の評価モデルを使って評価リクエストを送ることで Win レートが算出されます。Arbiter モデルによる評価は、Google が実際に社内でも使っている評価の手法です。



LLM の評価に特化した AB テスト

(3) Safety Bias

これは今後ローンチされる予定の機能であり、安全性バイアスを評価するものです。さまざまなアイデンティティ (ID) グループに分けながら、グループの中でどのスコアリングが良いか悪いか、つまり AI の安全性を事前に測定します。

下図では横軸、つまり列で ID グループが分かれており、それぞれのグループごとに AI モデルがどのような評価であるか (特定のグループで AI が提示する回答にバイアスや偏りがないかなど) を把握できます。皆様が作成されたチューニング済みモデルが本当に安全かどうか、特定の ID グループで危険なスコアが出ていないかを確認できます。

安全性バイアスの評価 Safety bias evaluation

Vertex AI のバイアス評価ツールを使用すると、お客様は調整されたモデル内の安全属性のパフォーマンスにおける潜在的なバイアスを評価できます

Safety attribute	Overall	Gender: Female	Gender: Male	Gender: Transgender	Gender: Non-binary
Death, Harm & Tragedy	0.062	0.023	0.018	0.014	0.007
Drugs	0.208	0.098	0.038	0.042	0.03
Finance	0.531	0.196	0.175	0.083	0.077
Firearms & Weapons	0.369	0.13	0.071	0.121	0.047
Hate	0.323	0.105	0.111	0.075	0.032
Health	0.191	0.056	0.045	0.02	0.07
Insult	0.569	0.021	0.076	0.093	0.379
Legal	0.176	0.033	0.01	0.034	0.099
Obscene	0.484	0.11	0.001	0.077	0
Politics	0.426	0.006	0.043	0.007	0
Public Safety	0.585	0.191	0.19	0.162	0
Religion & Belief	0.658	0.197	0.071	0.057	0.333
Sexual	0.233	0.074	0.071	0.005	0.083
Toxic	0.591	0.16	0.005	0.144	0.282
Violent	0.256	0.013	0.081	0.002	0.16
War & Conflict	0.081	0.011	0.004	0.012	0.054

33.3% of records show bias against Religion & Belief and Gender: Non-binary

Google Cloud Next Tokyo '23

Proprietary

035

安心して利用できる AI かどうかを事前に確認できる

今回は、生成 AI 時代こそ留意しておくべき MLOps のニーズについて触れながら、それに対応する Google Cloud の Vertex AI のソリューションも紹介してきました。Vertex AI は既存の MLOps のケイパビリティを拡張しながら、生成 AI のニーズに対応するように進化を遂げているので、お使いの既存の MLOps 投資を捨てる必要はありません。ぜひ、皆様の MLOps または生成 AI 活用の取り組みにご活用いただければと思います。

参照リンク

1. [Google Cloud のジェネレーティブ AI の概要](#)
2. [生成 AI 時代の MLOps 実現方法とは？ アーカイブ動画視聴ページ](#)

製品、サービスに関するお問い合わせ



goo.gl/CCZL78

Google Cloud の詳細については、上記 URL もしくは QR コードからアクセスしていただくか、同ページ「お問い合わせ」よりお問い合わせください。

© Copyright 2024 Google

Google は、Google LLC の商標です。その他すべての社名および製品名は、それぞれ該当する企業の商標である可能性があります。