

# 生成 AI アプリケーションでデータ活用

Ren Kurosawa

SOLUTIONS ARCHITECT

2024/05/16





# Ren Kurosawa

黒澤 蓮

AWS Japan Solutions Architect

Analytics/ML が Speciality で、普段はWebネイティブ業界のお客様を担当しております！

## 好きなAWSサービス

- AWS Glue
- Amazon Chime SDK
- Amazon Bedrock

## 将来の夢

宇宙でポエムを詠むこと

# 本日のアジェンダ

- 1 自社データを使う意義とユースケース
- 2 プロンプトエンジニアリングとRAGについて
- 3 自前データを使った、モデルの学習
- 4 データの可視化

# 本日のセッションとゴールの確認



汎用的な  
生成 AI



自社のビジネスと顧客を  
理解している生成 AI

# 本日のセッションとゴールの確認



汎用的な  
生成 AI



自社のビジネスと顧客を  
理解している生成 AI

生成 AI  
を活用  
したい



アプリケーション  
にどう  
落とし込むか  
(ユースケース)  
(事例)

# 本日のセッションとゴールの確認



汎用的な  
生成 AI



自社のビジネスと顧客を  
理解している生成 AI

生成 AI  
を活用  
したい



アプリケーション  
にどう  
落とし込むか  
(ユースケース)  
(事例)



差別化/最適化の  
ために1st-Party  
データを使いたい



データ  
貯める



AI に  
フィード  
バック  
/ 分析

# 本日のセッションとゴールの確認

今回のイベント  
の対象



汎用的な  
生成 AI



自社のビジネスと顧客を  
理解している生成 AI

生成 AI  
を活用  
したい



アプリケーション  
にどう  
落とし込むか  
(ユースケース)  
(事例)



差別化/最適化の  
ために1st-Party  
データを使いたい



データ  
貯める  
セッション  
1,2



AI に  
フィード  
バック  
/ 分析

# 本日のセッションとゴールの確認

今回のイベント  
の対象



汎用的な  
生成 AI



自社のビジネスと顧客を  
理解している生成 AI

生成 AI  
を活用  
したい



アプリケーション  
にどう  
落とし込むか  
(ユースケース)  
(事例)

本セッション：軽く



差別化/最適化の  
ために1st-Party  
データを使いたい

本セッション



データ  
貯める

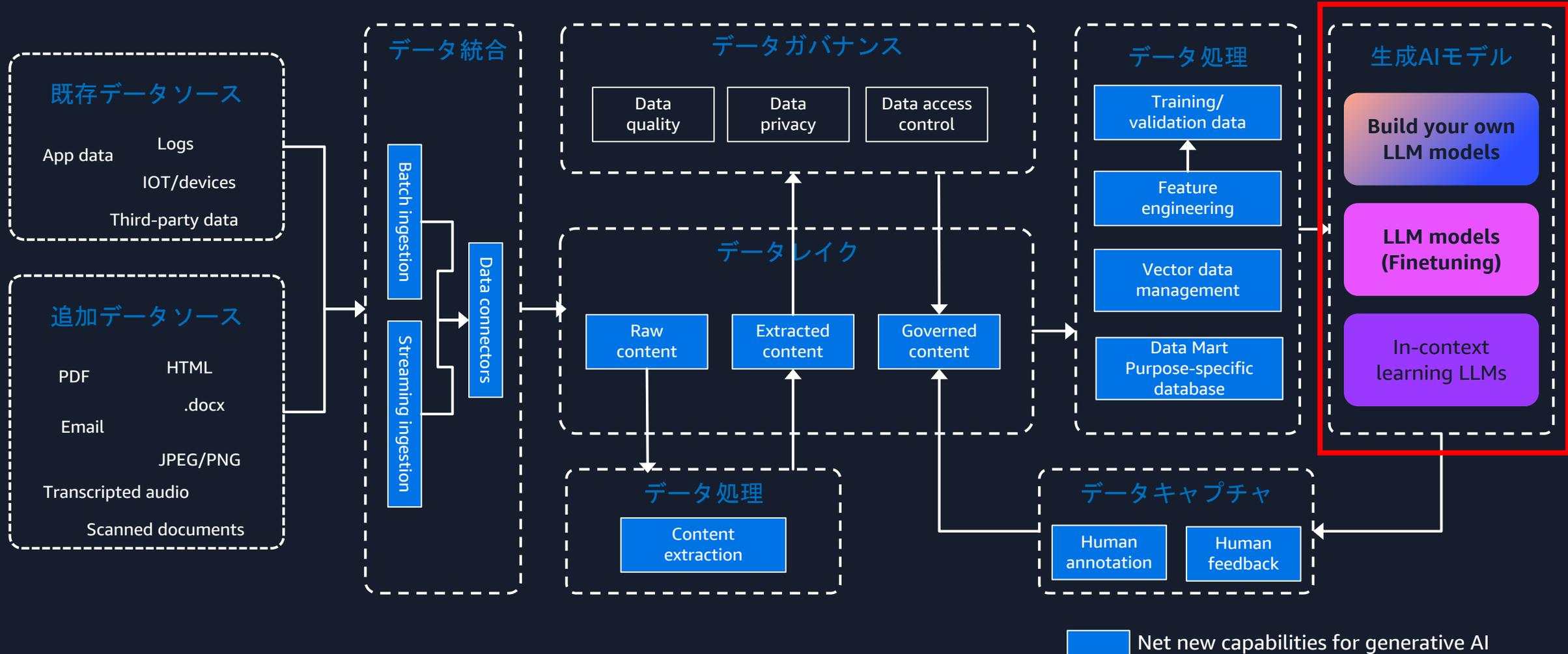
セッション  
1,2



AI に  
フィード  
バック  
/分析

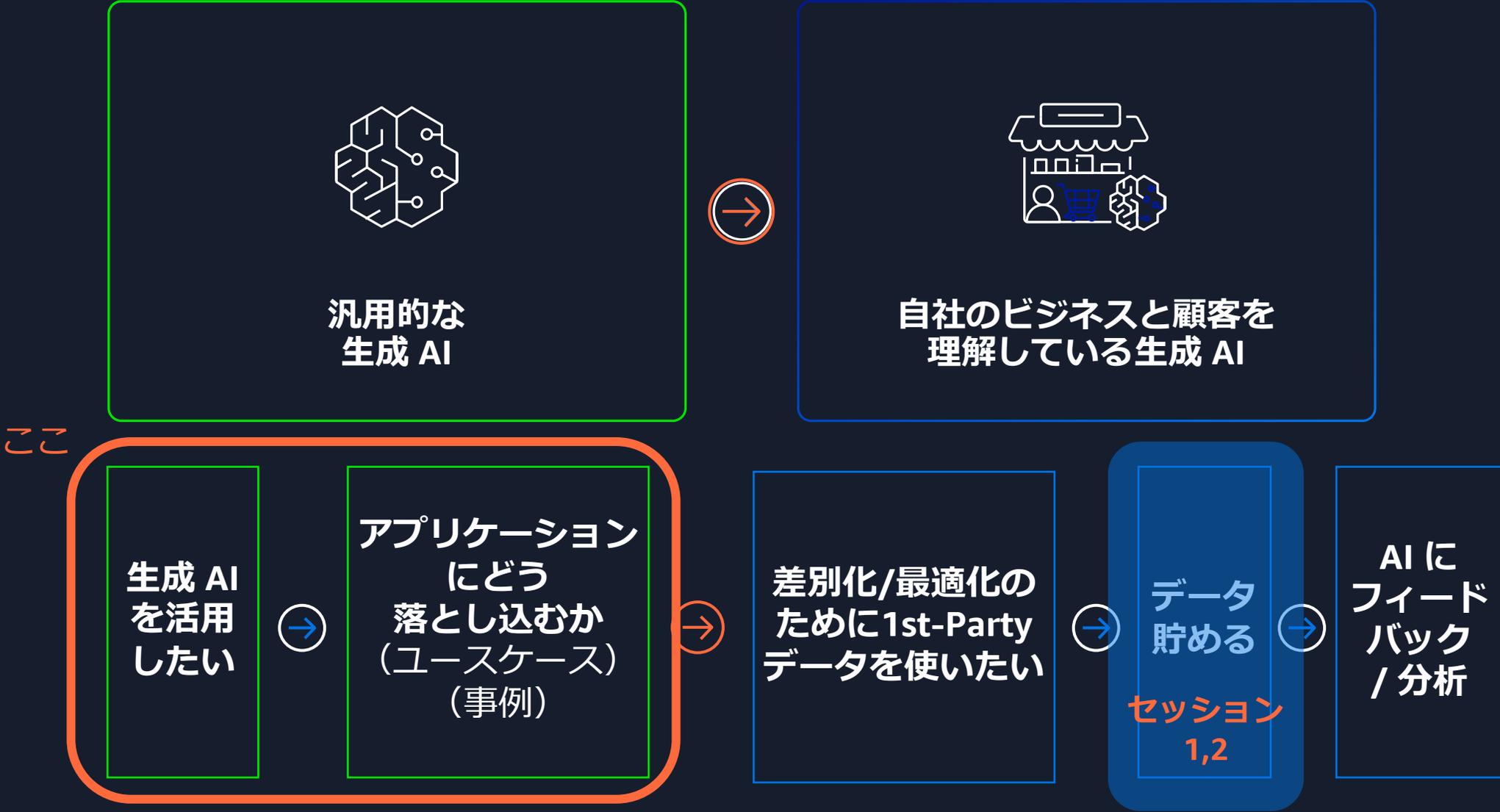
本セッション

# 本セッションの範囲



# 自社データを利用した生成 AI の活用先

# 本日のセッションとゴールの確認



# データを集めることでできるようになること

データを集めることは準備、データ活用は手段、その先の（ビジネス）ゴールは何か



## 業務効率化

在庫管理

異常検知

社内ナレッジの共有



## 顧客分析

施策が成功したか

潜在的なニーズの発掘

売上の予測



## 資産

機械学習データセット

データの売買

社会貢献

# データを集めることでできるようになること

データを集めることは準備、データ活用は手段、その先の（ビジネス）ゴールは何か



## 業務効率化

在庫管理

異常検知

社内ナレッジの共有



## 顧客分析

施策が成功したか

潜在的なニーズの発掘

売上の予測



## 資産

機械学習データセット

データの売買

社会貢献

生成 AI を/へ活用

# 生成 AI (Generative AI) とは

会話、ストーリー、画像、動画、音楽など、  
新しいコンテンツやアイデアを創造できる AI の手法

- 画像生成  
cute cat



- 会話生成

Human
どうすれば人に好かれますか
Assistant
人に好かれるためには、以下のような点を心がけると良いでしょう。
<ul style="list-style-type: none"><li>• 相手の話をじっくり聞く。質問したり相づちを打ったりして関心を示す。</li><li>• 相手の話題に合わせて会話を作る。自分の話に偏らない。</li><li>• 褒める機会があれば積極的に褒める。人は褒められると嬉しくなる。</li><li>• 相手の話を否定せず、共感を示す態度を心がける。</li><li>• 相手の話のポイントを汲み取り、要点をまとめて伝える。</li></ul>

- コード生成

入力:

1から100までの数字を順番に  
出力してください。

ただし、2で割り切れるならFizz,  
3で割り切れるならBuzz,  
2でも3でも割り切れるならFizzBuzz  
と出力してください。

```
1 for i in range(1, 101):
2     output = ""
3     if i % 2 == 0:
4         output += "Fizz"
5     if i % 3 == 0:
6         output += "Buzz"
7     print(output if output != "" else i)
```

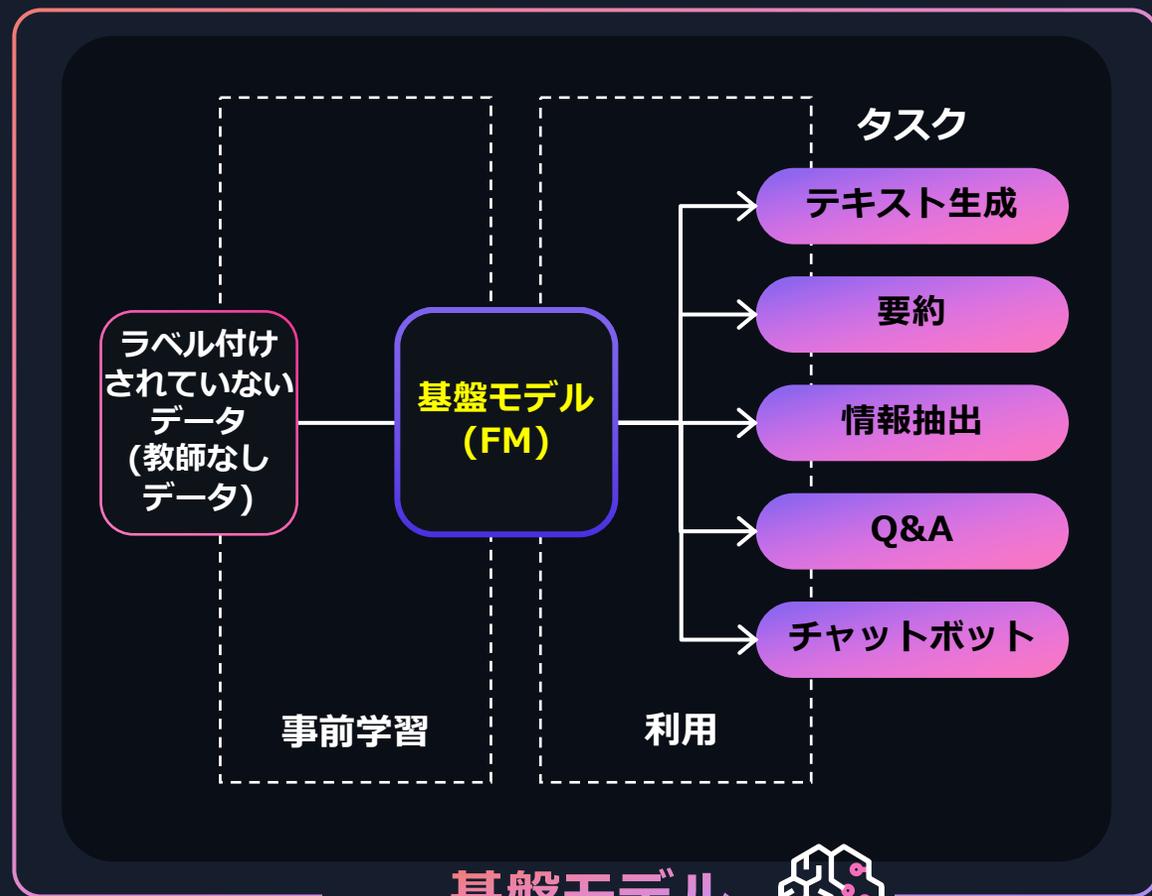
# 生成系 AI を実現する仕組み - 基盤モデル

生成系 AI は「**基盤モデル**」と呼ばれる

「膨大なデータに基づき事前にトレーニングされた大規模な機械学習モデル」によって実現



従来の機械学習モデル



基盤モデル

# 基盤モデルで何ができるか？

# Amazonにおける 生成AI 活用例



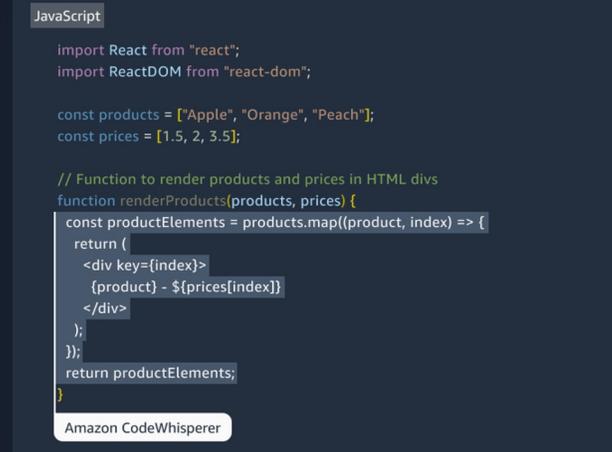
## Amazon の検索機能

基盤モデルを導入し  
検索結果を改善



## Alexa 教師モデル

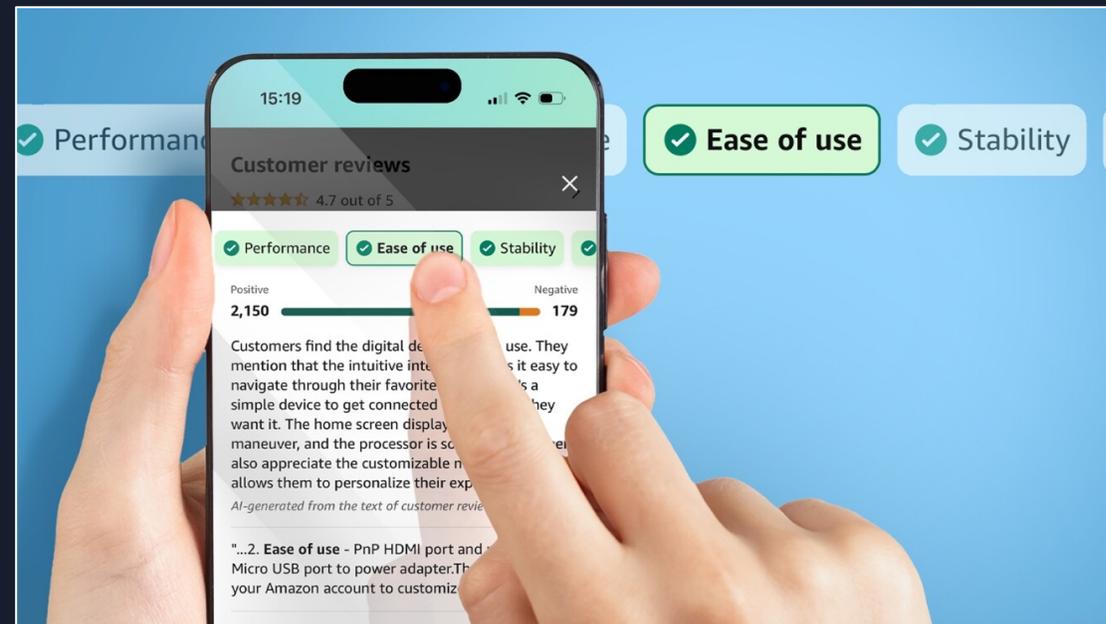
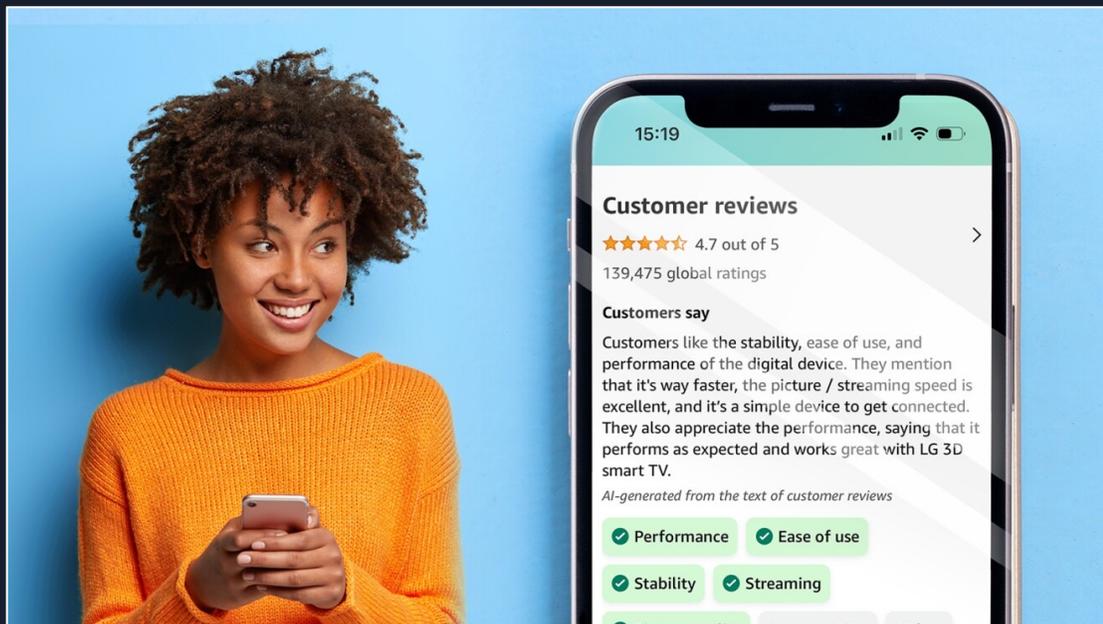
最小限の入力情報から  
多言語による応答を実現



## Amazon CodeWhisperer

コメントとコードから  
推奨コードを自動生成

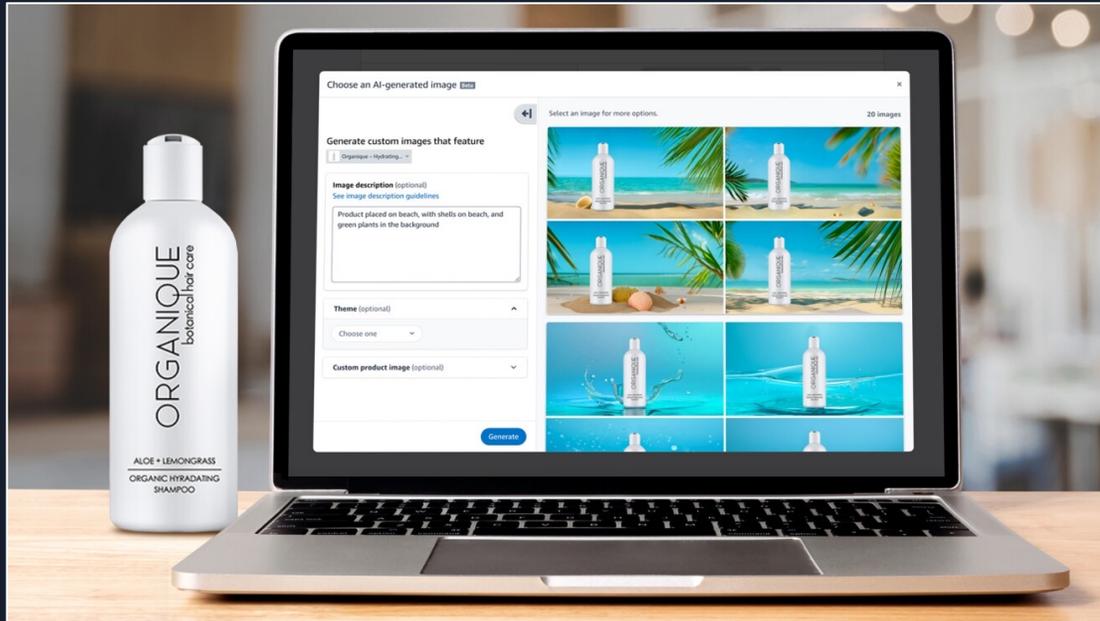
# Amazonにおける 生成AI 活用例



- 信頼できるユーザーのレビューに基づき AI が製品レビューのハイライトを生成
- 製品の特徴を示すアイコンを AI が生成 (性能、使い易さ、安定性など)

- アイコンを選択すると、関連するレビューテキストの詳細をまとめて読める
- 出品者の作業負担を軽減するAI のサービスも開始し、商品の効率的な販売を支援

# Amazonにおける 生成AI 活用例



- 商品の宣伝に失敗した出品者の 75% がクリエイティブの選択を課題に挙げている
- 出品者の労力なしに魅力的な画像を生成しつつ、お客様体験の向上を助ける策を検討

- 出品者がワンクリックで商品に適した背景を生成できるツールのベータ版をリリース
- プロンプトによる画像調整にも対応
- リアルな背景のクリックスルー率は 40% 増

Source: Amazon.com / "Amazon rolls out AI-powered image generation to help advertisers deliver a better ad experience for customers"  
<https://www.aboutamazon.com/news/innovation-at-amazon/amazon-ads-ai-powered-image-generator>

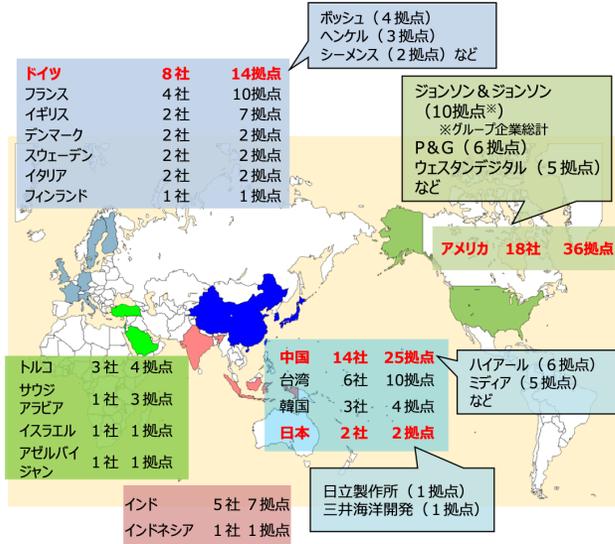
# 資料理解

## 製造業をめぐる新しい国際的な潮流

第5章 環境の変化

- 世界経済フォーラムは、2020年より世界最先端工場「Global Lighthouse」を選出しており、**132のうち日本企業は2拠点**。
- 選出に当たり、経済的効率性だけでなく、デジタル技術を活用した**サプライチェーン全体での最適化を通じ**、生産性の向上・市場ニーズに応じた柔軟な生産の実現や、**環境負荷の低減**等が図られているかどうかが重要視されている。
- **DXやGXによる全体最適化の達成**が、製造事業者の先進性の評価軸となる、国際的な潮流が生まれている。

図 Global Lighthouse選出状況※（本社所在国別）



※(資料) 世界経済フォーラム公表資料より経済産業省作成

### 選出に当たって重視される審査ポイント

インパクト	製造プロセスや製品等の高度化に大きなインパクトを達成しているか
ユースケース	複数の4IRのテクノロジーを活用・統合し、ユースケースを創出しているか
イネーブラー	4IRの活用戦略、IoTのアーキテクチャ、労働力の巻き込み・活用や能力開発などの複数要素をうまく統合・活用しているか
テクノロジープラットフォーム	革新的でスケラブルなテクノロジープラットフォームであるか

※4IRのテクノロジーは、人間の生活・仕事に根本的な変化をもたらす技術革新を可能とする先端的な技術を差し、例えば、ビッグデータ、IoT、AI、ロボットの活用等である。(4IR=第四次産業革命)  
 ※Sustainability Lighthouseの選出に当たっては、「目的：明確なサステナビリティ目標があるか」、「インパクト：複数の環境カテゴリーで改善ができていないか」、「スケール：複数の4IRのテクノロジーを活用し、サステナビリティ面での効果・インパクトを残しているか」の3つの審査ポイントが追加される。

### 選出工場の取組と効果

ハイアール 【中国・家電メーカー】	<ul style="list-style-type: none"> <li>● 冷蔵庫の生産工場において、注文内容をリアルタイムで製造工程に反映させる高度なデジタルプラットフォーム「COSMOplat」を開発。</li> <li>● 柔軟に製造工程を組み換えるマス・カスタマイゼーションを実現した。</li> </ul>
ポッシュ 【ドイツ・自動車部品メーカー】	<ul style="list-style-type: none"> <li>● ビッグデータの活用により、市場調査に要する時間を85%削減</li> <li>● AIによる工程最適化により、エネルギー消費量を37%削減</li> <li>● ロボットの連結と自動化により、組立効率を52%向上</li> </ul>
	<ul style="list-style-type: none"> <li>● 全世界に点在する工場の機械やデバイスをネットワークで連結することで、生産や物流計画の管理、機械の状態の把握等を一本化。</li> <li>● 全社的な全体最適化に成功した。</li> </ul>
	<ul style="list-style-type: none"> <li>● 物流システムの統合により、生産リードタイムを36%削減</li> <li>● 設備に生じる異常の事前検知により、メンテナンスコストを25%削減</li> <li>● エネルギー管理の実現により、消費電力を18%削減</li> </ul>



あなたは製造業のコンサルタントです、この画像を丁寧に確認した上で、内容を要約してください。また、そのようやく結果から得られる示唆と日本の製造業が取り組むべき課題について教えて下さい



この画像は、デジタル技術の活用によりグローバルな製造業の変革が進展していることを示しています。主な内容は以下の通りです。

- 世界的なイニシアチブ「Global Lighthouse」では、132社の先進企業が製造現場のデジタル化を推進している。
- AIやIoTなどの先端技術を活用することで、生産性の向上や製品・サービスの高度化が図られている。
- DX(デジタルトランスフォーメーション)の進展により、製造業の在り方や付加価値が変化している。

この動きから得られる示唆は以下の通りです。

1. デジタル技術の積極的な活用が製造業の競争力強化に不可欠である。
2. 製品・サービスの高付加価値化や新ビジネスモデルの創出が重要になる。
3. データ駆動型のスマート製造への転換が求められる。

日本の製造業が取り組むべき課題としては、以下が考えられます。

- AIやIoTなどの先端技術の一層の導入・活用
- DXを見据えた業務プロセスの見直しと人材育成
- サプライチェーン全体でのデータ連携と共有
- 製品ライフサイクル全体を視野に入れた高付加価値化

グローバルな潮流に遅れを取らないよう、デジタル化を積極的に推進していくことが欠かせません。

<https://www.meti.go.jp/report/whitepaper/mono/2023/index.html>

# 観光画像へのタグ付けと魅力的な説明文生成



建物, 夜景, 和風 [1]

この画像は、日本の伝統的な温泉旅館の夜景を捉えています。木造の建物と赤い灯籠が並ぶ情緒あふれる風景は、日本の古き良き時代の雰囲気を感じさせます。夜の静けさの中で、温泉の湯けむりに包まれながら過ごす時間は、現代の喧騒から離れてゆっくりと心を休める贅沢な体験となるでしょう。日本を象徴する景色の中で、伝統的な温泉文化に触れることができる魅力的な場所です。



通り, 紅葉, 車 [1]

この美しい街路樹が黄金色に染まった風景は、秋の訪れを華やかに告げています。落ち着いた住宅街の一角にあり、ゆったりとした雰囲気が漂う通りです。太陽の光が木々の葉を透かし、幻想的な光景を創り出しています。歩いて散策するのがおすすめです。秋の穏やかな空気に包まれながら、この美しい景色を堪能できるでしょう。地元の人々の生活の一コマを垣間見ることにもできる、魅力的な場所です。

# SNSマーケティング支援

## Instagram Marketing Assistant

Choose an image...



Drag and drop file here

Limit 200MB per file • JPG, JPEG, PNG

Browse files



2024-03-21\_18-36-36.jpg 87.7KB



Focus item

パンツ



画像をアップロードして、  
紹介したいアイテムを指定

Analyze Image

Generate Instagram posts and hashtags

## 画像の分析と解説

画像全体の説明: この画像は、都会的な街路を背景に、全身黒の服装でポーズをとる男性を捉えています。彼の表情は真剣で、キャップとバックパックを身に付けています。前景には車道と仕切りの壁があり、外国の街並みが窺えます。全体的にストリートスタイルのファッションが際立っています。

パンツの説明: 男性が着用しているパンツは、黒色のワイドシルエットで裾広がりデザインです。生地は厚手で落ち着いた雰囲気を出しています。ヒップからウエストにかけてゆとりがあり、ルーズなスタイリングがされています。足元のスニーカーと相まって、カジュアルでありながらも洗練された印象を与えるアイテムです。

## 提案

**Sonnet** で読み込んだ画像を元に、インスタへの  
投稿用原稿とハッシュタグを出力

インスタグラムの投稿の原稿提案:

1. 都会的なストリートファッションに映える、ワイドシルエットのブラックパンツ。ルーズなシルエットとドレープ感のある生地で、カジュアルさとこなれ感を両立。旬のスタイリングに取り入れて、あなたもストリートの逸材に。
2. 街角のスナップショットから生まれた、最新ブラックパンツ。落ち着いた色味と上質な生地で、あなたのスタイルに洗練された雰囲気をプラス。きちんと感とカジュアル感のバランスが絶妙なこのパンツは、毎日のコーデに欠かせない。
3. ストリートカルチャーから着想を得た、ブラックワイドパンツ。ゆとりのあるシルエットとドレープ感が、さり気ないこなれ感を演出。黒の重厚感とスポーティーなディテールの絶妙なブレンドは、センスの良さを物語る。

適切なハッシュタグ提案:

#ストリートファッション #ワイドパンツ #ブラックパンツ #ルーズシルエット #シックスタイル #カジュアルスタイル #ストリートスナップ #スタイリングアイデア  
#モノトーンコーデ #メンズファッション

# 文字起こし結果の要約

架空の打ち合わせの文字起こし結果から、要約とネクストアクションを抽出しJSONで出力

```
{  
  "summary": "会議では、春らしい具材を使ったおにぎりの提案がなされました。桜の葉、新玉ねぎ、  
  筍、山菜、桜の花、菜の花などの具材が挙げられ、それぞれの味付けや見た目の工夫について議論  
  されました。",  
  "next_action": "出されたアイデアを元に、それぞれの具材を使ったおにぎりの試作品を作り、味や  
  見た目を確認して商品化の可否を検討することになりました。"  
}
```

# 論文の要約 (自然なキャラクターロールプレイも)

太郎「博士、今日はどんな論文を教えてください?」

博士「ほっほっほ、今日はAWSのサーバレスコンピューティングサービス、Lambdaにおける高速コンテナ起動の仕組みについての論文を紹介するぞい」

太郎「Lambdaってクラウドの関数実行サービスですよね。コンテナの起動が速いって、どういうことなんですか?」

博士「おお、鋭い質問じゃ。Lambdaではユーザーの関数をコンテナ内で実行するのじゃが、そのコンテナイメージを起動する際の待ち時間、いわゆるコールドスタート時間を最小限に抑える必要があるのう」

太郎「なるほど。ユーザーにとってコールドスタート時間が長いと使いにくいですからね」

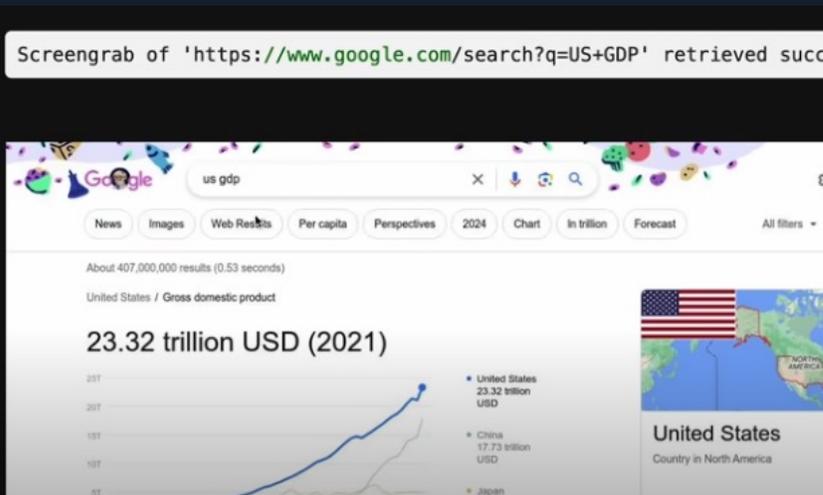
博士「そうじゃ。従来はコンテナイメージ全体を転送・展開する必要があり、イメージサイズが大きくなるとコールドスタート時間が長くなってしまふのじゃ。そこで本論文では、コンテナイメージを小さなチャンクに分割し、必要なチャンクのみを転送・ロードする方式を提案しているのう」

太郎「へえ、スパーローディングって手法ですね。でも大きなイメージを分割するとチャンク数が膨大になりそうですし、セキュリティの問題もありそうです ...」

# 世界のGDPトレンドの分析

## Amazon Bedrockで利用可能なClaude 3 Opusの公式デモ

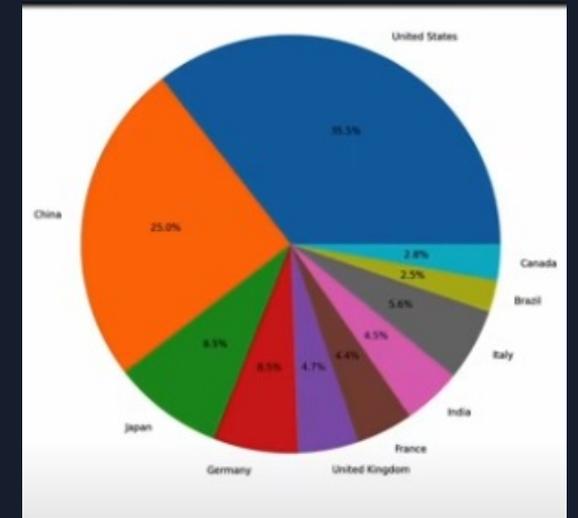
世界のGDPトレンド分析という複雑なタスクに対して、Claude 3が「Webサイト閲覧」「Python実行」「サブエージェントの実行」などの複数のツールを駆使して自律的に分析を実行している



Webサイトを閲覧しグラフを読み取る



国ごとに並列で分析を進める  
サブエージェントを実行



Pythonで結果を可視化

# 動画のハイライト生成 with Bedrock

README MIT license

## Soccer Highlight Generator

Automate the creation of soccer match highlights with the power of Generative AI and AWS. This solution leverages AWS Bedrock (Anthropic's Claude 3 Sonnet model), AWS MediaConvert, Lambda, Step Functions and other AWS services to identify and compile exciting game moments without manual editing.

### Author

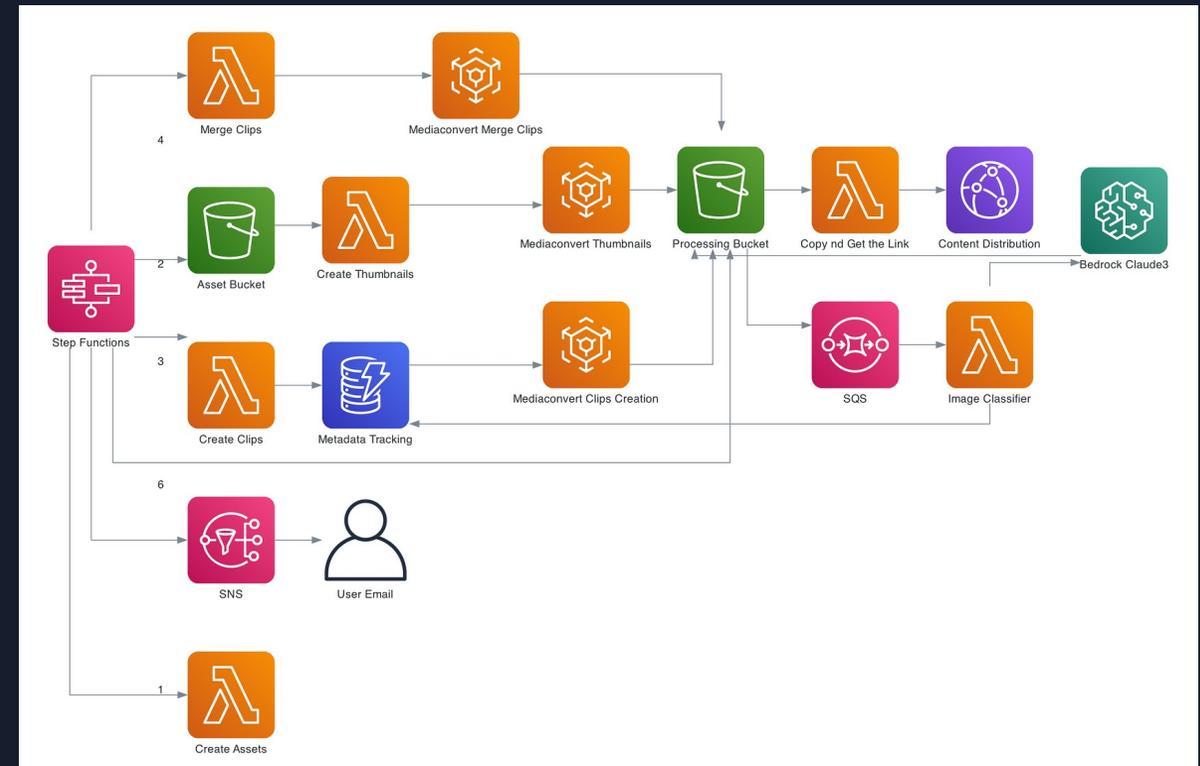
[Pedram Jahangiri](#)

### Getting Started

For a detailed explanation of what this solution does and the benefits it offers, please refer to [my blog](#)

### Prerequisites

- AWS CLI installed and configured with the necessary permissions
- Node.js and npm
- Python 3.11 and pip
- An AWS account with the required services enabled
- Access to Amazon Bedrock foundation models (Before you can use a foundation model in Amazon Bedrock, you must request access to it. Use this Link for detail <https://docs.aws.amazon.com/bedrock/latest/userguide/model-access.html>)



<https://github.com/iut62elec/Soccer-Highlight-Generator-with-GenAI>

# 生成 AI 活用をはじめる上での難しさ

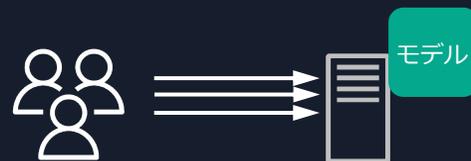
## 生成 AI の進化への追従

複数のモデルを簡単に利用できる環境の構築



## インフラストラクチャ

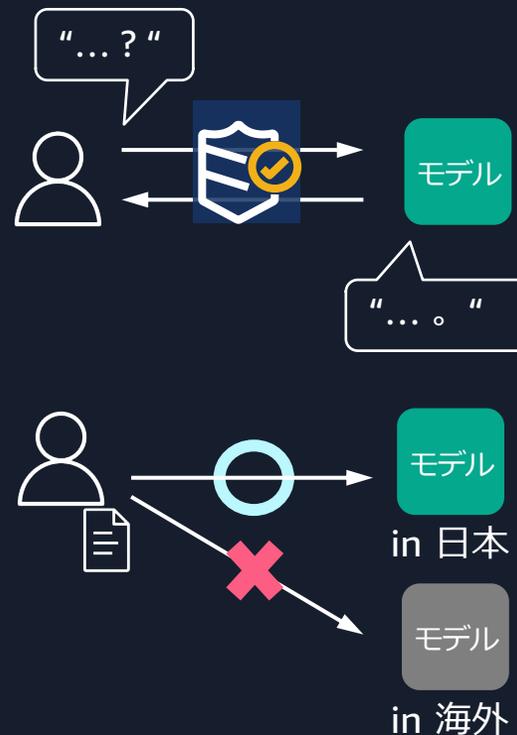
モデルの実行を支えるインフラの管理



- 必要なこと
- 各モデルごとの環境構築と機械学習の知識
  - モデルのデプロイやカスタマイズのための大規模な計算リソース
  - インフラの維持管理
  - 柔軟なスケーリング
  - ...

## プライバシーとセキュリティ

データ漏洩の防止、国内にデータを留めたい



# Generative AI tools and services



## APPLICATIONS THAT LEVERAGE LLMs AND FMs

 Amazon Q Business

 Amazon Q Developer

 Amazon Q in QuickSight

 Amazon Q in Connect

## TOOLS TO BUILD WITH LLMs AND OTHER FMs

 **Amazon Bedrock** | Guardrails | Agents | Customization Capabilities

## INFRASTRUCTURE FOR MODEL TRAINING & INFERENCE

 GPUs  Trainium  Inferentia  SageMaker

 UltraClusters  EFA  EC2 Capacity Blocks  Nitro  Neuron

# Generative AI tools and services



## APPLICATIONS THAT LEVERAGE LLMs AND FMs

-  Amazon Q Business
-  Amazon Q Developer
-  Amazon Q in QuickSight
-  Amazon Q in Connect

## TOOLS TO BUILD WITH LLMs AND OTHER FMs

-  **Amazon Bedrock** | Guardrails | Agents | Customization Capabilities

## INFRASTRUCTURE FOR MODEL TRAINING & INFERENCE

-  GPUs
-  Trainium
-  Inferentia
-  SageMaker
-  UltraClusters
-  EFA
-  EC2 Capacity Blocks
-  Nitro
-  Neuron

# Amazon Bedrock



基盤モデルを活用した  
生成 AI アプリケーションを  
簡単に構築、拡張できる方法



API を介してさまざまな基盤モデルにアクセス、  
インフラ管理は不要



お客様の業務用途に適した基盤モデルを選択  
Amazon、AI21 Labs、  
Anthropic、Cohere、Meta、Stability AI



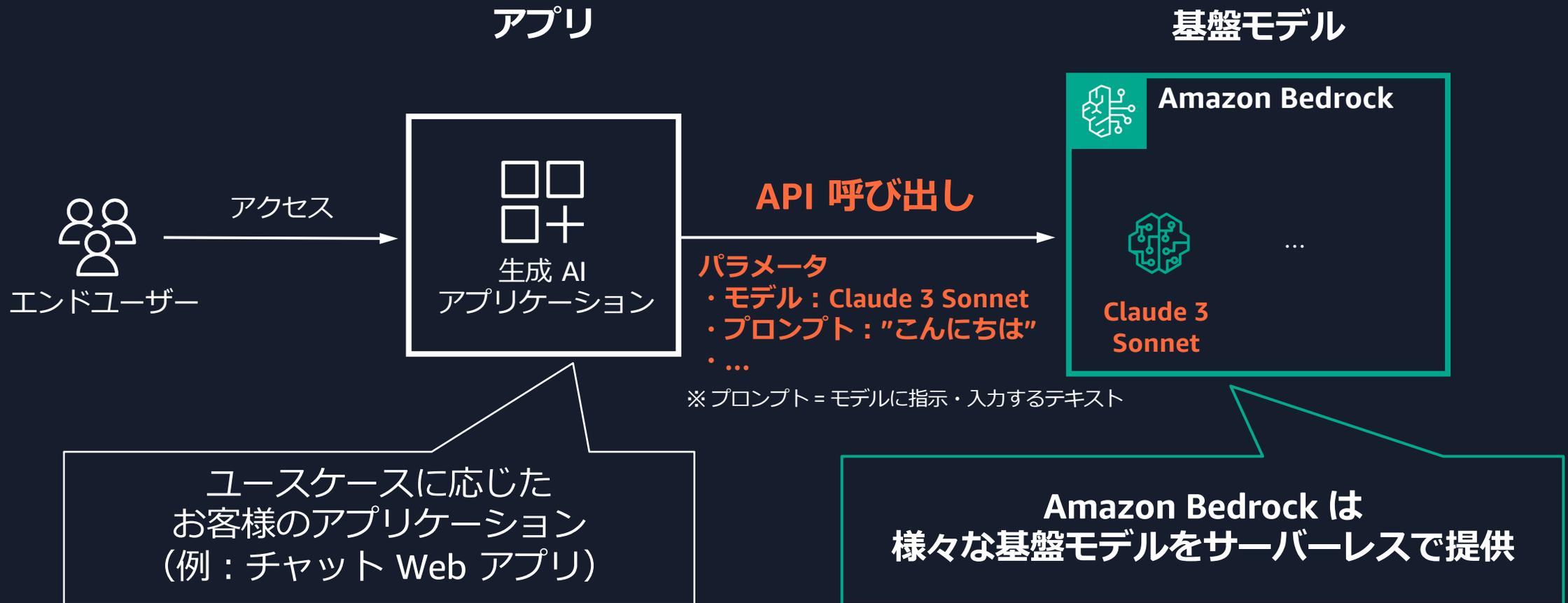
データセキュリティやコンプライアンスを実現



エージェント機能、RAG 機能、非公開でのモデル  
のカスタマイズなど基盤モデルの効果を高める  
さまざまな機能を提供

# Amazon Bedrock の利用イメージ

基盤モデルをサーバーレスで提供しておりアプリケーションから API を通じて利用



# Amazon Bedrock

幅広い基盤モデルの選択肢をご提供

AI21 labs

ANTHROPIC

Mistral AI

co:here

∞ Meta AI

stability.ai

amazon

JURASSIC

CLAUDE

MISTRAL  
& MIXTRAL

COMMAND  
& EMBED

LLAMA 2

SDXL

AMAZON TITAN

テキスト

テキスト  
& ビジョン

テキスト

テキスト

テキスト

画像

テキスト

Jurassic-2 Ultra  
Jurassic-2 Mid

**Claude 3 Opus**  
Claude 3 Sonnet  
Claude 3 Haiku

**Mistral Large**  
Mistral 7B  
Mixtral 8X7B

**Command R+**  
**Command R**  
Command  
Command Light

**Llama 3 70B**  
**Llama 3 8B**  
Llama 2 70B  
Llama 2 13B

SDXL 1.0

**Titan Text Premier**  
Titan Text Express  
Titan Text Lite

テキスト

Claude 2.1  
Claude 2.0  
Claude Instant

埋め込み

Embed - Multilingual  
Embed - English

画像

Titan Image Generator

埋め込み

Titan Multimodal Embeddings  
**Titan Text Embeddings V2**  
Titan Text Embeddings



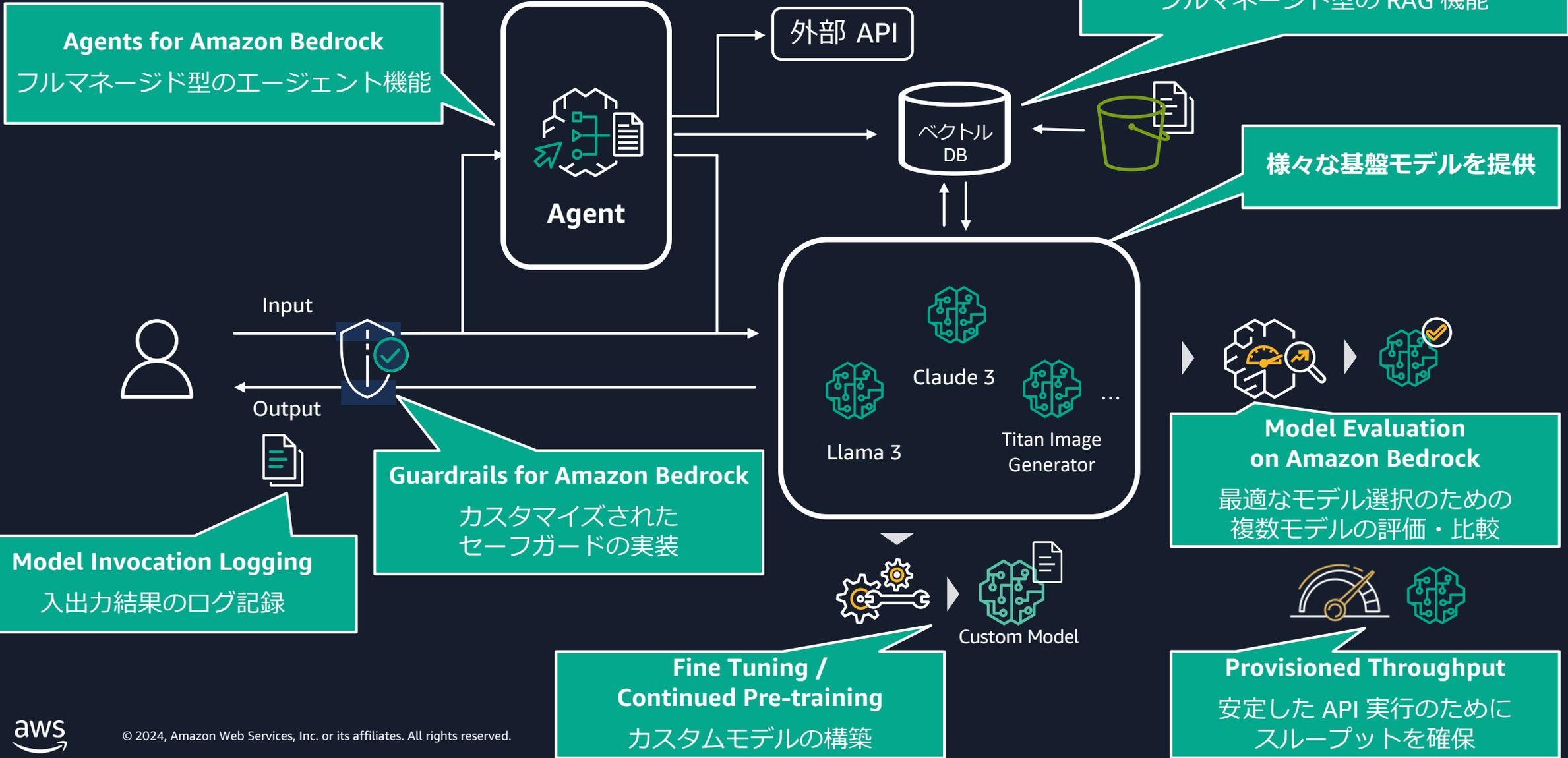


Amazon Bedrock は  
様々な基盤モデルを  
API 経由で利用できる サービス



基盤モデルだけでなく  
生成 AI アプリの構築・運用をサポートする  
幅広い機能を提供するサービス

# Amazon Bedrock の機能全体像



# 本日のセッションとゴールの確認



汎用的な  
生成 AI



自社のビジネスと顧客を  
理解している生成 AI

ここ

生成 AI を活用したい

アプリケーションにどう落とし込むか (ユースケース) (事例)

**CLEAR**



差別化/最適化の  
ために1st-Party  
データを使いたい



データ  
貯める

セッション  
1,2



AI に  
フィード  
バック  
/ 分析

# 1<sup>st</sup> party データを利用した生成 AI

# 新たな生成 AI のデータパターン

1



## 検索拡張生成 (RAG)

事前にトレーニングされたモデルに、自社のドメイン固有のコンテキストデータを加えた結果を生成

例: ドメイン固有の、要件が限定された仮想エージェント

2



## 既存モデルの ファインチューニング

既存のモデルをベースに、ドメイン固有のデータによる追加学習を実施

例: ドメイン特化型の知識エージェント

3



## 独自モデルの作成

ドメイン固有データによるモデルの関するモデルのトレーニングを実施

例: 特定のデータで学習を行った、高度にドメインに特化したアプリケーション

# 新たな生成 AI のデータパターン

0

プロンプトエンジニアリング

1



## 検索拡張生成 (RAG)

事前にトレーニングされたモデルに、自社のドメイン固有のコンテキストデータを加えた結果を生成

例: ドメイン固有の、要件が限定された仮想エージェント

2



## 既存モデルの ファインチューニング

既存のモデルをベースに、ドメイン固有のデータによる追加学習を実施

例: ドメイン特化型の知識エージェント

3



## 独自モデルの作成

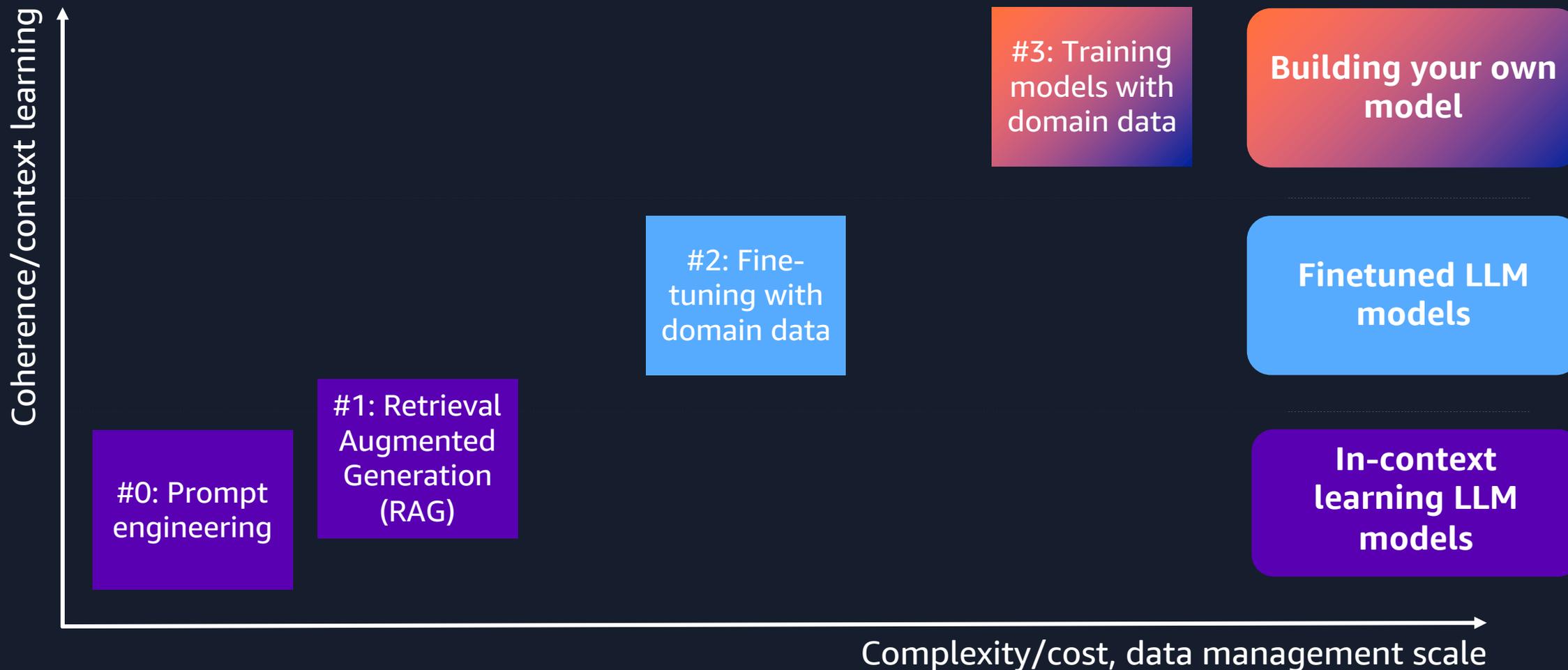
ドメイン固有データによるモデルのトレーニングを実施

例: 特定のデータで学習を行った、高度にドメインに特化したアプリケーション

一般のモデルを利用し、  
入力するデータをカスタマイズ

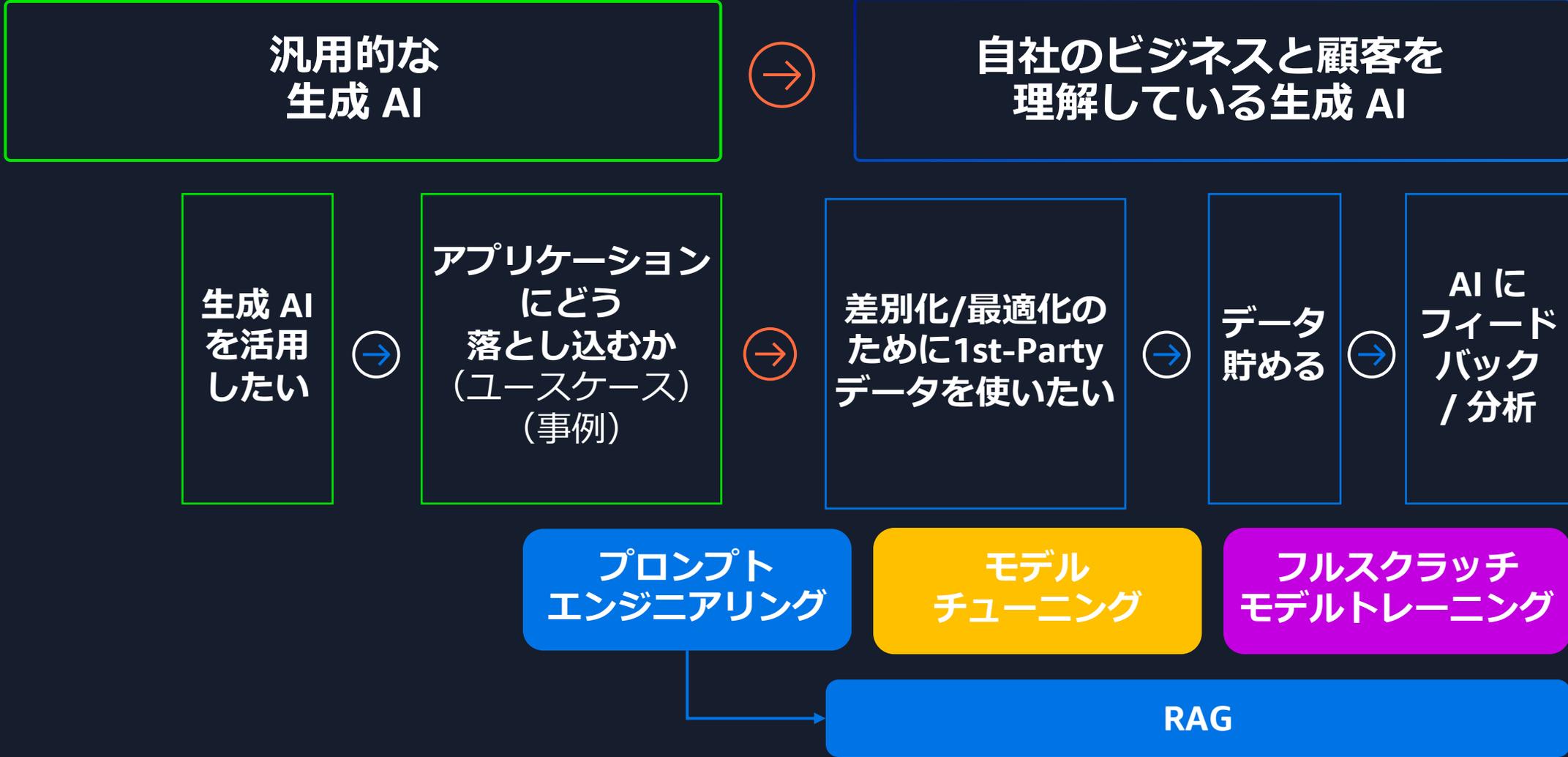
モデルそのものをカスタマイズする

# 新たな生成 AI のデータパターンと難易度



\* Qualitative comparison and not drawn to scale

# 本日のセッションとゴールの確認



# プロンプトエンジニアリング

# プロンプトエンジニアリングとは

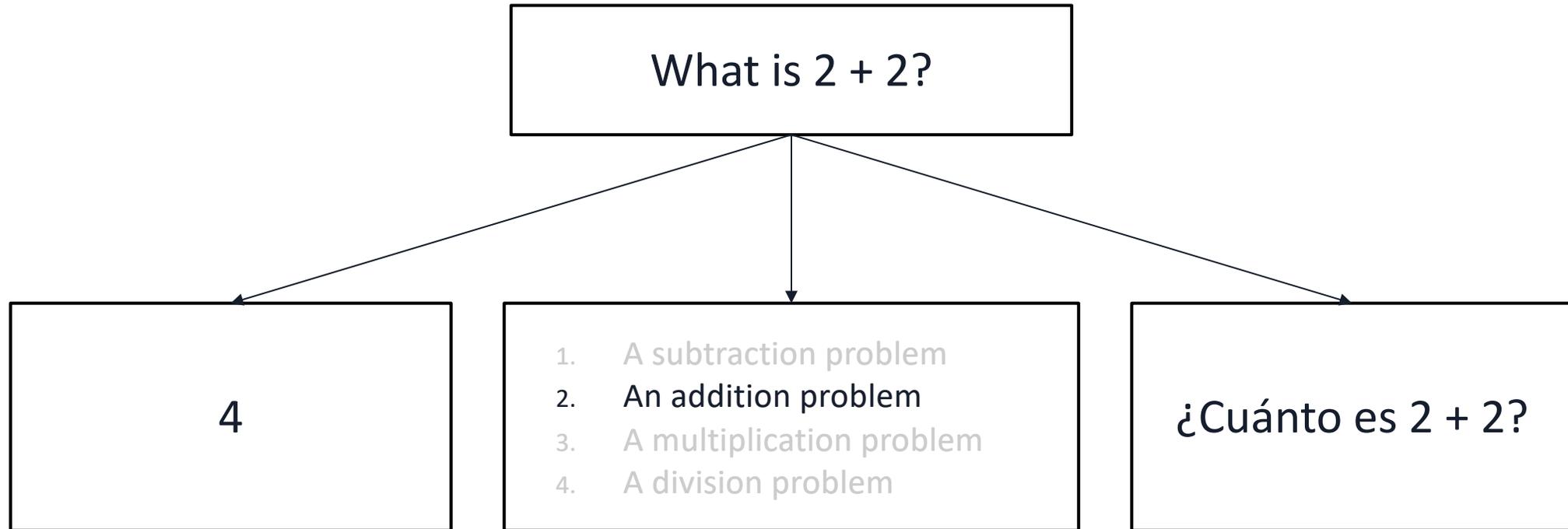
生成系 AI で適切な出力を得るためには、入力の指示文章(プロンプト)をうまく設計する必要があります

以下のようなプロンプトが有効とされています

- タスクを具体的に指示する
- 否定系を使わない
- 望ましい出力の例を含める

など

# What is prompt engineering?



**Prompt engineering** is the process of **controlling model behavior** by **optimizing your prompt to elicit high performing LLM responses** (as assessed by rigorous evaluations tailored to your use case).

# Example:

## Parts of a prompt

1. Task context
2. Tone context
3. Background data, documents, and images
4. Detailed task description & rules
5. Examples
6. Conversation history
7. Immediate task description or request
8. Thinking step by step / take a deep breath
9. Output formatting
10. Prefilled response (if any)

User	<p>You will be acting as an AI career coach named Joe created by the company AdAstra Careers. Your goal is to give career advice to users. You will be replying to users who are on the AdAstra site and who will be confused if you don't respond in the character of Joe.</p> <p>You should maintain a friendly customer service tone.</p> <p>Here is the career guidance document you should reference when answering the user: &lt;guide&gt;{{DOCUMENT}}&lt;/guide&gt;</p> <p>Here are some important rules for the interaction:</p> <ul style="list-style-type: none"><li>- Always stay in character, as Joe, an AI from AdAstra careers</li><li>- If you are unsure how to respond, say "Sorry, I didn't understand that. Could you repeat the question?"</li><li>- If someone asks something irrelevant, say, "Sorry, I am Joe and I give career advice. Do you have a career question today I can help you with?"</li></ul> <p>Here is an example of how to respond in a standard interaction:</p> <pre>&lt;example&gt; User: Hi, how were you created and what do you do? Joe: Hello! My name is Joe, and I was created by AdAstra Careers to give career advice. What can I help you with today? &lt;/example&gt;</pre> <p>Here is the conversation history (between the user and you) prior to the question. It could be empty if there is no history: &lt;history&gt; {{HISTORY}} &lt;/history&gt;</p> <p>Here is the user's question: &lt;question&gt; {{QUESTION}} &lt;/question&gt;</p> <p>How do you respond to the user's question?</p> <p>Think about your answer first before you respond. Put your response in &lt;response&gt;&lt;/response&gt; tags.</p>
Assistant (prefill)	<response>

# Prompt Engineering: Few-shot Prompting

指示を出す際に、望ましい回答を例示することで、意図した出力を得ます

アンケートの回答結果を、  
下記の例に沿って分類してください。

<example>

最寄り駅から10分は歩く必要がありました：施設

音響が良くなく、音がひび割れていました：施設

ドリンクの種類が豊富でした：料理

従業員の対応はとても親切でした：人

軽食はどれもおいしかったです：料理

</example>

<question>

館内の清掃はとても行き届いていました：

</question>



質問のContextを別段落で追加することも多い

# RAG



# 新たな生成 AI のデータパターン

0

プロンプトエンジニアリング

1



## 検索拡張生成 (RAG)

事前にトレーニングされたモデルに、自社のドメイン固有のコンテキストデータを加えた結果を生成

例: ドメイン固有の、要件が限定された仮想エージェント

2



## 既存モデルの ファインチューニング

既存のモデルをベースに、ドメイン固有のデータによる追加学習を実施

例: ドメイン特化型の知識エージェント

3



## 独自モデルの作成

ドメイン固有データによるモデルのトレーニングを実施

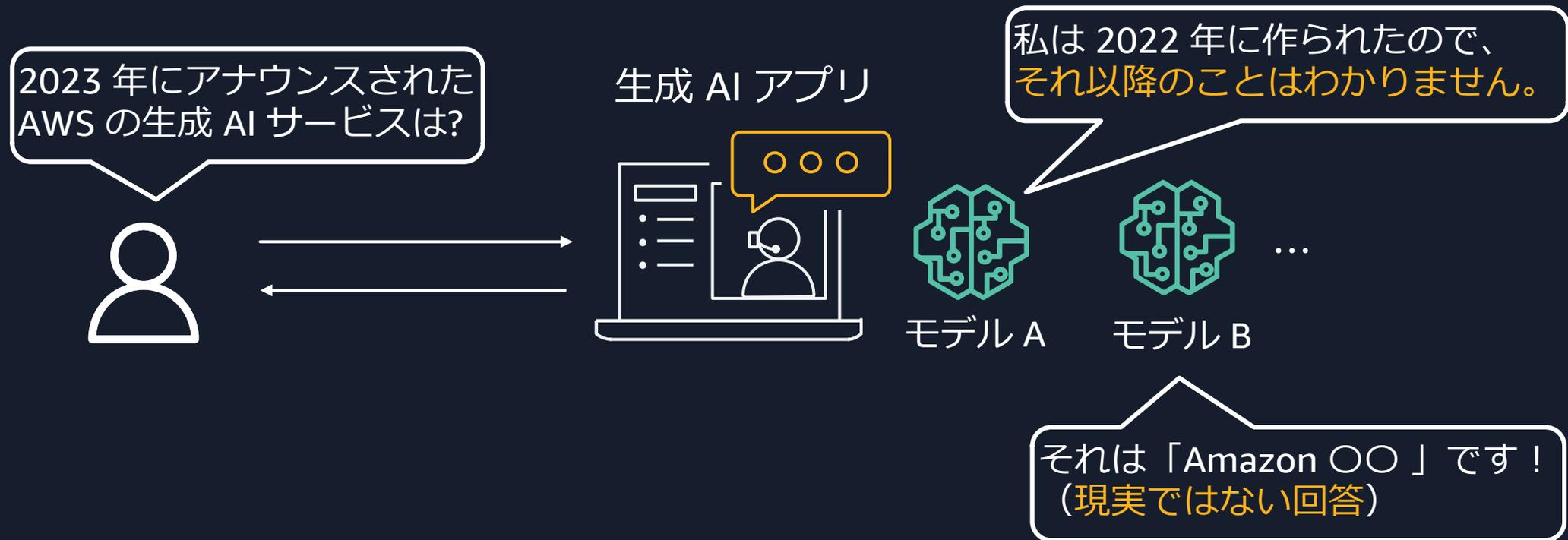
例: 特定のデータで学習を行った、高度にドメインに特化したアプリケーション

一般のモデルを利用し、  
入力するデータをカスタマイズ

モデルそのものをカスタマイズする

# よくある課題：お客様独自のデータの活用

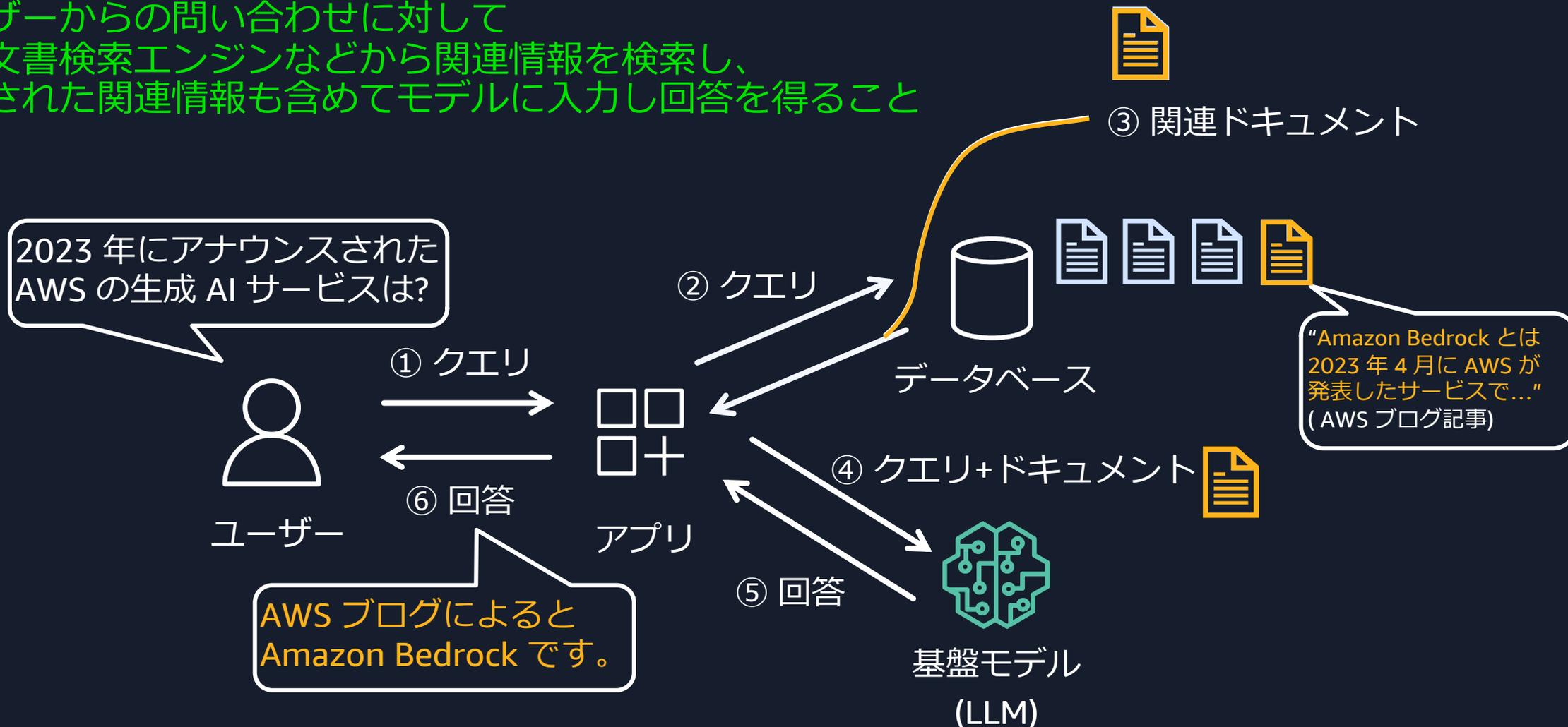
- 基盤モデルは学習データに含まれていないような情報は回答できない。場合によっては不正確な回答をしてしまう。(=ハルシネーション)
- 社内の情報やドキュメントなどを回答に含めたい場合に陥る課題



# RAG とは

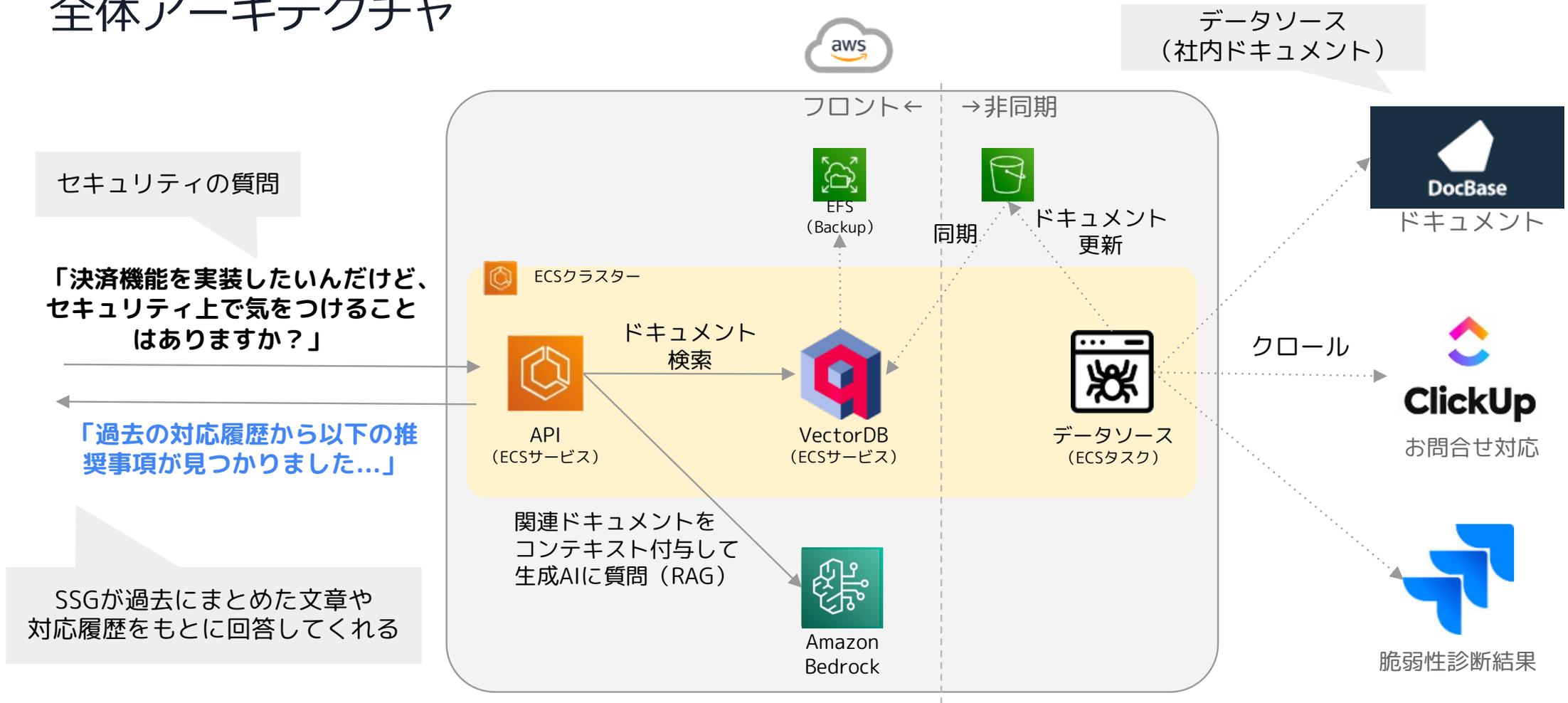
## Retrieval-Augmented Generation (検索拡張生成)

ユーザーからの問い合わせに対して社内文書検索エンジンなどから関連情報を検索し、取得された関連情報も含めてモデルに入力し回答を得ること



# CyberAgent SSG 様のセキュリティデスク RAG 事例

## 全体アーキテクチャ



# 現状のAWSのRAG俯瞰図

完全フル  
マネージド

## Amazon Q Business

[Retriever] ・ Native ・ Kendra + [Data Source]

## Knowledge Base for Amazon Bedrock

[VectorDB] ・ OpenSearch SL ・ PineCorn ・ Redis ・ Aurora (pgVector) ・ MongoDB  
+ Bedrock (FM)

## LangChain

[フルマネージド検索] ・ Kendra

[VectorDB] ・ OpenSearch SL ・ Document DB ・ Memory DB ・ pgVector (Postgres DB)

[Embedding] ・ Titan ・ Cohere (Bedrock)

+ LLM

手動作成

# Knowledge bases for Amazon Bedrock

検索拡張生成 (RAG) の  
フルマネージドネイティブサポート



テキスト文書を埋め込み形式に  
自動的に変換



変換した埋め込み形式データを  
ベクトルデータベースに保存

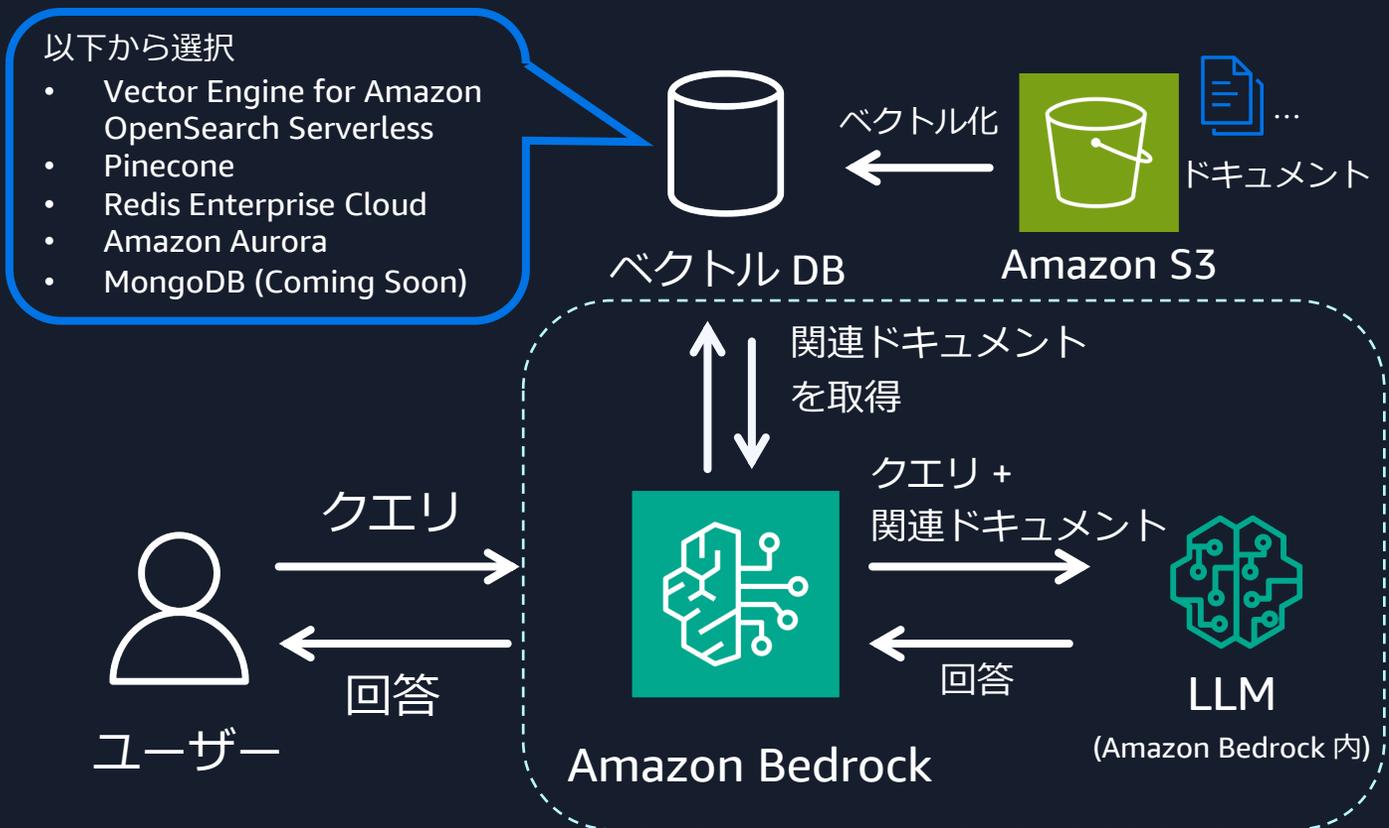


保存した埋め込み形式データを  
取得してプロンプトを拡張

# Knowledge Bases for Amazon Bedrock

フルマネージド型の RAG 機能を提供

基盤モデルと自社データソースを組み合わせた RAG (検索拡張生成) をフルマネージドに実現可能に



2023年にアナウンスされた AWS の生成系 AI サービスは?

2023年にアナウンスされたAWSの生成系AIサービスはAmazon Bedrockです。[1]

[Show result details >](#)

Amazon Bedrockのメリットは何ですか?

Amazon Bedrockのメリットとして、主要な基盤モデルの選択ができること、データを使用してモデルを簡単にカスタマイズできること、APIを動的に呼び出してタスクを実行するフルマネージド型エージェントが構築できること、RAGのネイティブサポートによりFMの機能を拡張できること、データセキュリティとコンプライアンス認証があることが挙げられています。[1]

[Show result details >](#)

# Knowledge Base for Amazon Bedrock

- Vector database は、 OpenSearch Serverless、 Redis Enterprise Cloud、 Pinecone、 Aurora、 MongoDB から選択できる



Vector Engine For Amazon  
OpenSearch Serverless



Redis Enterprise  
Cloud



Pinecone



Amazon  
Aurora

COMING SOON



MongoDB

# 現状のAWSのRAG俯瞰図

完全フル  
マネージド

## Amazon Q Business

[Retriever] ・ Native ・ Kendra + [Data Source]

## Knowledge Base for Amazon Bedrock

[VectorDB] ・ OpenSearch SL ・ PineCorn ・ Redis ・ Aurora (pgVector) ・ MongoDB  
+ Bedrock (FM)

## LangChain

[フルマネージド検索] ・ Kendra

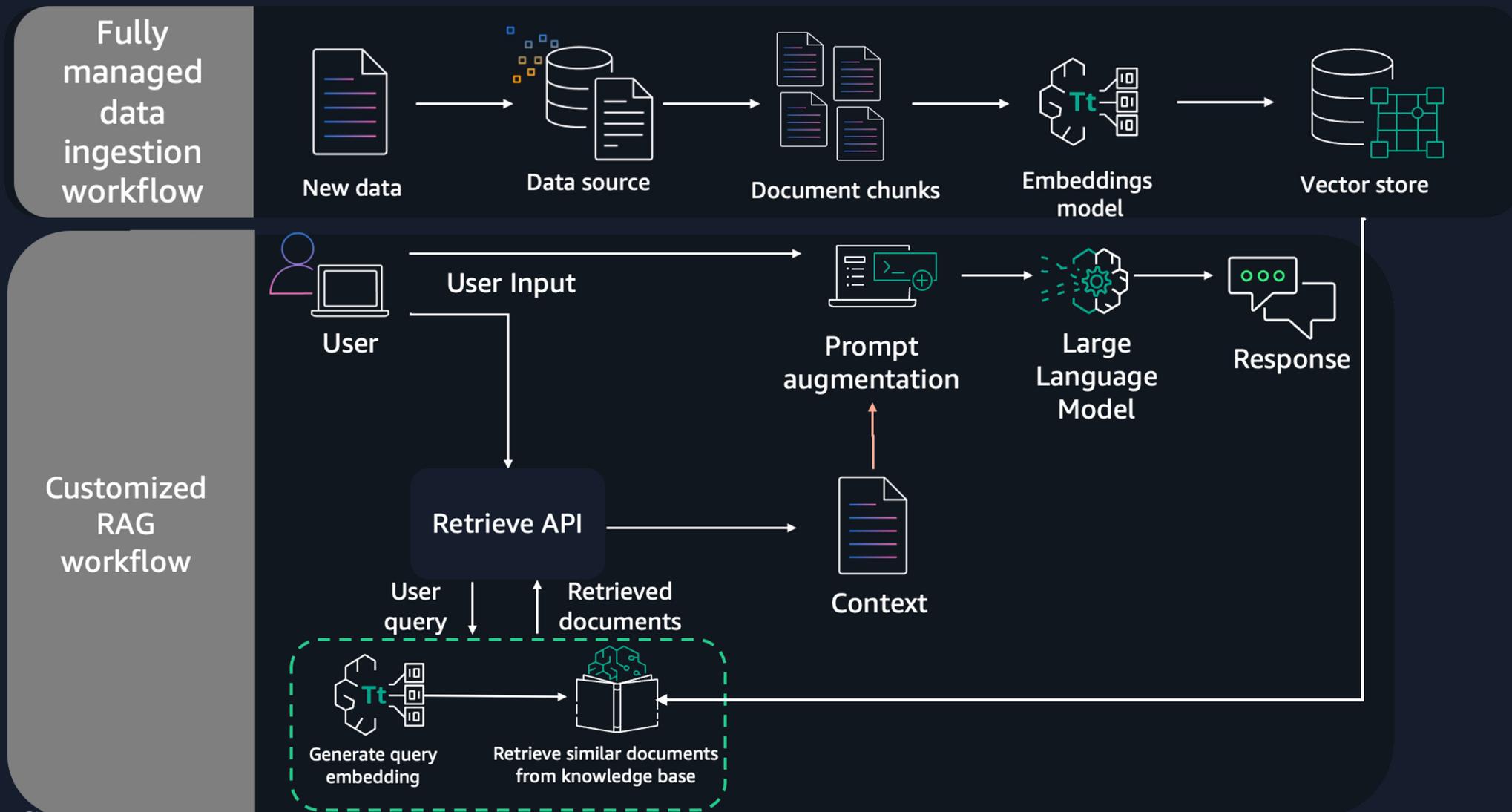
[VectorDB] ・ OpenSearch SL ・ Document DB ・ Memory DB ・ pgVector (Postgres DB)

[Embedding] ・ Titan ・ Cohere (Bedrock)

+ LLM

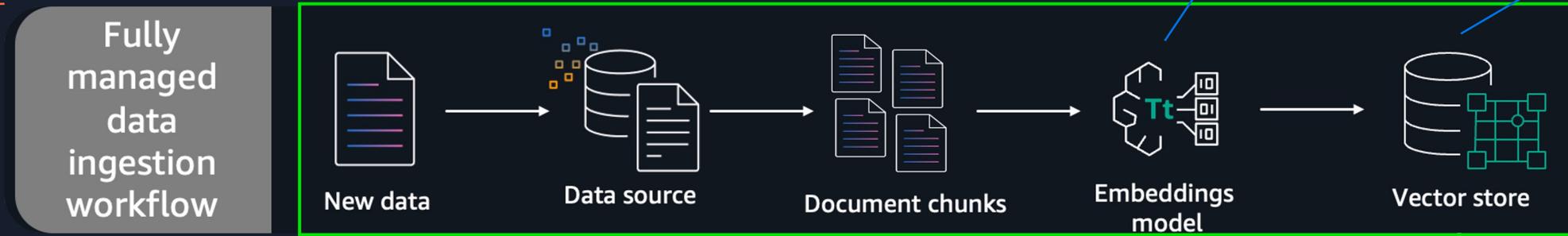
手動作成

# RAG の詳細な仕組み



# RAG の詳細な仕組み

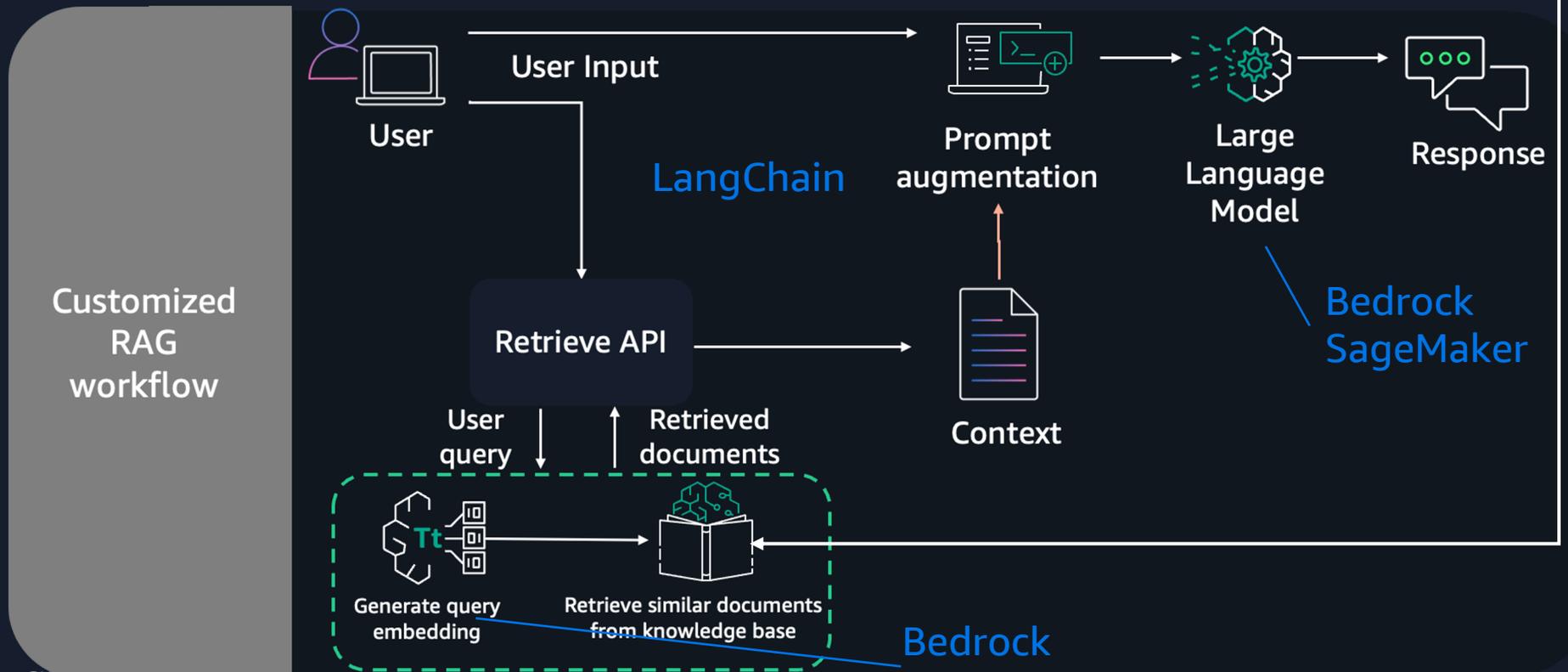
Amazon Q / Knowledge Base for Amazon Bedrock



Bedrock SageMaker

OpenSearch Document DB  
Memory DB  
pgVector

Kendra



LangChain

Bedrock SageMaker

Bedrock SageMaker



# 基盤モデルのカスタマイズ

# 新たな生成 AI のデータパターン

0

1



## 検索拡張生成 (RAG)

事前にトレーニングされたモデルに、自社のドメイン固有のコンテキストデータを加えた結果を生成

例: ドメイン固有の、要件が限定された仮想エージェント

2



## 既存モデルの ファインチューニング

既存のモデルをベースに、ドメイン固有のデータによる追加学習を実施

例: ドメイン特化型の知識エージェント

3



## 独自モデルの作成

ドメイン固有データによるモデルのトレーニングを実施

例: 特定のデータで学習を行った、高度にドメインに特化したアプリケーション

プロンプトエンジニアリング

一般のモデルを利用し、  
入力するデータをカスタマイズ

モデルそのものをカスタマイズする

# ドメイン学習による基盤モデルの最適化

## 経済情報特化の生成AI、日経が開発 40年分の記事学習

生成AI [+ フォローする](#)

2024年4月24日 11:00 (2024年4月26日 13:53更新)

保存



性能の指標となる「パラメーター数」が130億と700億の2種類を開発した

日本経済新聞社は24日、経済情報に特化した生成AI（人工知能）の基盤技術を開発したと発表した。大規模言語モデルと呼ばれるもので、約40年分の日経グループの新聞や雑誌の記事を学習させた。記事の要約機能などで活用を見込む。

パラメーター数が700億のモデルは米メタが2023年7月に公開した大規模言語モデル「Llama（ラマ）2」をベースに、「ファインチューニング（追加学習）」と呼ぶ手法を使って開発した。ラマ2は外部の企業や開発者による改良が広く認められている。

開発したモデルは最新ニュースに関する知識を備え、記事の要約などで性能を発揮することを確認した。現在はメタが4月に公開した最新版の「ラマ3」をベースに、パラメーター数が80億のモデルの学習にも着手している。

日経イノベーション・ラボは19年に言語モデルの研究を始め、著作権上の課題やAIがもっともらしく嘘をつくハルシネーション（幻覚）の問題に対処する方法を探ってきた。

<https://www.nikkei.com/article/DGXZQOUC1941R0Z10C24A4000000/>

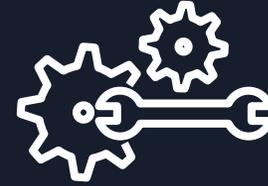
# Fine Tuning / Continued Pre-training

自身のデータセットを学習して基盤モデル自体をカスタマイズする機能



## Fine Tuning

ラベル付きデータを新たに学習して  
特定のタスクの精度を高める手法



## Continued Pre-training

大量のラベル無しデータを用いて  
モデルに新たなドメイン知識を習得させる手法



再学習用のコード不要でマネージドな環境でモデルのカスタマイズを実行

# Fine-tuning と Continued pre-training の違い



{"prompt": "日本の首都は?", "completion": "東京なのだ"}  
{"prompt": "アメリカの首都は?", "completion": "ワシントンなのだ"}  
{"prompt": "日本の次の元号は?", "completion": "うーん、わからないのだ"}

{"input": "<raw text>"  
{"input": "<raw text>"  
{"input": "<raw text>"

# Generative AI tools and services



## APPLICATIONS THAT LEVERAGE LLMs AND FMs

-  Amazon Q Business
-  Amazon Q Developer
-  Amazon Q in QuickSight
-  Amazon Q in Connect

## TOOLS TO BUILD WITH LLMs AND OTHER FMs

-  **Amazon Bedrock** | Guardrails | Agents | Customization Capabilities

## INFRASTRUCTURE FOR MODEL TRAINING & INFERENCE

-  GPUs
-  Trainium
-  Inferentia
-  **SageMaker**
-  UltraClusters
-  EFA
-  EC2 Capacity Blocks
-  Nitro
-  Neuron

# Amazon Bedrock Custom Models

Bedrock に搭載されている一部のモデルを対象として  
Custom Models 機能でモデルのファインチューニング / Continued pre-training  
が可能に

## Amazon Bedrockでモデルのファインチューニングが可能に

- ラベル付きの小規模なデータセットで、モデルを特定のタスクに特化させる手法がファインチューニング
- Amazon Bedrockで選択できる以下のモデルで、ファインチューニングを実行可能に
  - Cohere Command Lite, Meta Llama 2, Amazon Titan Text Lite/Express
- バージニアとオレゴンにて一般利用開始



一般利用開始

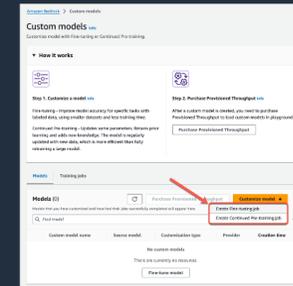


9

© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

## Amazon BedrockのContinued pre-trainingを発表

- Amazon Titan Text Express/Liteについて、安全で管理された環境を利用して独自のラベルなしデータによるカスタマイズが可能に
  - ドメイン知識を必要とするアプリケーションを構築するには、一般的な公開データでトレーニングされた基盤モデルでは不十分な場合がある
  - 用途に応じた知識を基盤モデルに与えるために、ユーザが提供する独自データを加えて継続的に事前トレーニングを行うことが可能
- バージニアとオレゴンでプレビューを開始



プレビュー



10

© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Amazon SageMaker JumpStart

## 数百のモデルが活用可能

AI21 labs

Meta AI

cohere

Hugging Face

stability.ai

LightOn

databricks

alexa

### モデル

Jurassic-2 Ultra, Mid

状況に応じた回答

テキスト要約

言い換え

文法エラー訂正

### タスク

テキスト生成

長文生成

テキスト要約

言い換え

チャット

情報抽出

### モデル

Llama 2 7B, 13B, 70B

### タスク

質問回答

チャット

テキスト要約

言い換え

感情分析

テキスト生成

### モデル

Cohere

Command XL

### タスク

テキスト生成

情報抽出

質問回答

テキスト要約

### モデル

Falcon-7B, 40B

Open LLaMA

RedPajama

MPT-7B

BloomZ 176B

Flan T-5 models (8 種類)

DistilGPT2

GPT NeoXT

Bloom models

(3 種類)

### タスク

機械翻訳

質問回答

テキスト要約

### モデル

Stable Diffusion XL 1.0

2.1 base

高精細化

インペインティング

### タスク

テキスト入力から  
フォト・リアルな  
画像生成

生成画像の品質改善

### 特長

Stable Diffusion 2.1  
ファイン・チューン対応

### モデル

Lyra-Fr

10B, Mini

### タスク

テキスト生成

キーワード抽出

情報抽出

質問回答

テキスト要約

感情分析

テキスト分類

### モデル

Dolly

### タスク

質問回答

チャット

言い換え

感情分析

テキスト生成

### モデル

AlexaTM 20B

### タスク

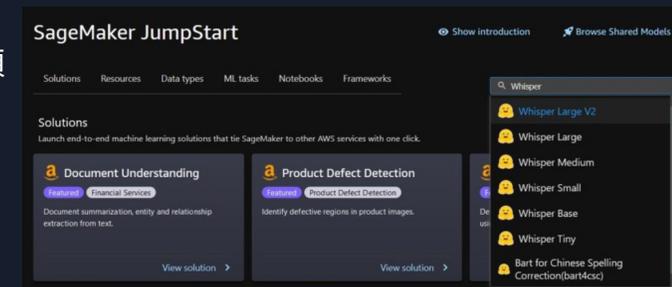
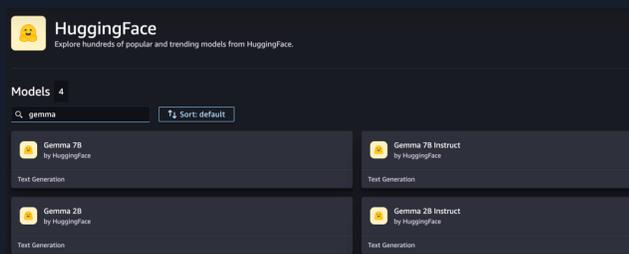
機械翻訳

質問回答

テキスト要約

注釈付与

データ生成



# Amazon SageMaker JumpStart 🦉 Fine-tuning / Continued Pre-training

The screenshot displays the Amazon SageMaker JumpStart interface. On the left, a search for 'calm2' yields 38 results, listing various models such as CyberAgentLM2-7B-Chat, In-context learning with..., Llama 2 7B Neuron, Llama 2 7B Chat Neuron, Llama 2 13B Neuron, Llama 2 13B Chat Neuron, Llama 2 70B Chat, Llama 2 7B, and Llama 2 7B Chat. On the right, the configuration page for 'CyberAgentLM2-7B-Chat (CALM2-7B-Chat)' is shown, with tabs for Deploy, Train, Notebook, and Model details. The 'Train Model' section includes instructions on creating a training job and configuring data sources. The 'Training data set' is set to 's3://jumpstart-cache-prod-us-west-2/training-datasets/genuq/small/' and the 'Validation data set' is 's3://bucketName/path-to-folder/'. A 'Train' button is visible at the bottom of the configuration page.

モデルを選んで、データを用意してクリックするだけ

# 独自モデルのフルスタック開発



## APPLICATIONS THAT LEVERAGE LLMs AND FMs

-  Amazon Q Business
-  Amazon Q Developer
-  Amazon Q in QuickSight
-  Amazon Q in Connect

## TOOLS TO BUILD WITH LLMs AND OTHER FMs

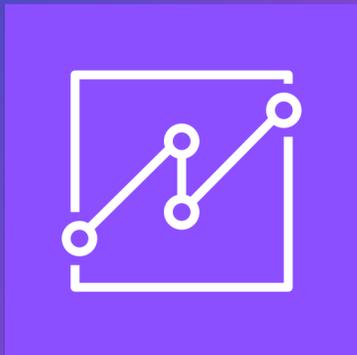
-  **Amazon Bedrock** | Guardrails | Agents | Customization Capabilities

## INFRASTRUCTURE FOR MODEL TRAINING & INFERENCE

-  GPUs
-  Trainium
-  Inferentia
-  SageMaker
-  UltraClusters
-  EFA
-  EC2 Capacity Blocks
-  Nitro
-  Neuron

# データの可視化と生成 AI





# Amazon QuickSight

UNIFIED BI SERVICE  
AT HYPERSCALE



あらゆるデータ活用を実現し  
BI環境を統合



オートスケーリングで高パフォーマンスを実現



生成系AI機能で分析をさらに増強



従量課金制でコスト削減

# Amazon Q in QuickSight

あらゆるユーザのデータ活用の意思決定過程を加速する

## AIにより加速された ダッシュボード作成

ビジネスアナリストがより早くビルドを実現

## AIがオンデマンドでデータの質問に答える

ビジネスユーザによるインサイト抽出を簡単に

## AIアシスタントによる データストーリー

ビジネスユーザによる発見や所見を共有可能に

## カスタムアプリにAIを拡張

開発者にも高度なAI機能をすばやく利用可能

Powered by  
Amazon Bedrock

※英語のみのサポート

※2024年5月現在US East (N. Virginia), US West (Oregon), EU(Frankfurt)のみでの提供,



# AI-powered ダッシュボード作成におけるエクスペリエンス

A NEW DASHBOARD BUILDING EXPERIENCE POWERED BY GENERATIVE BI

## ビジュアルのビルド

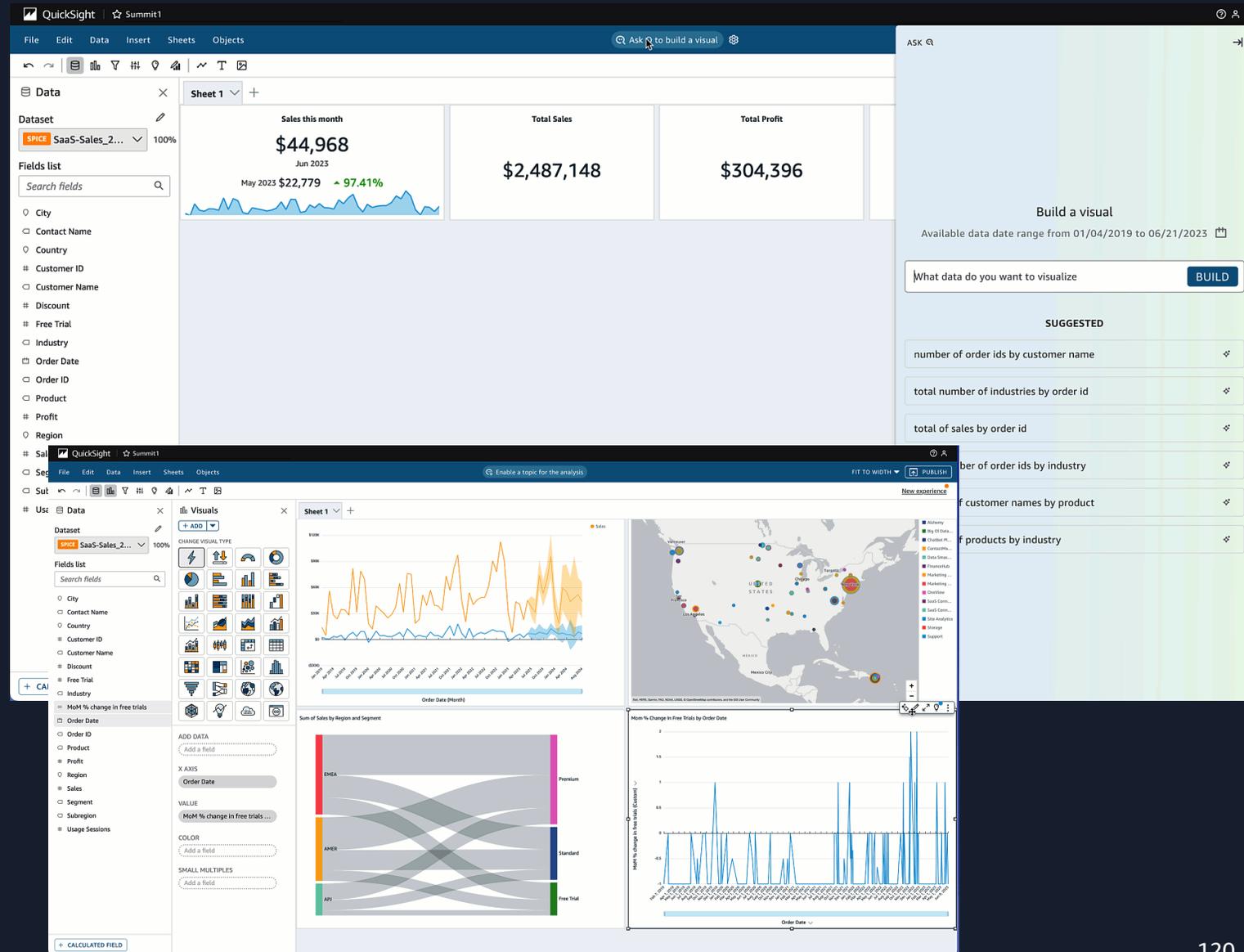
ダッシュボードやレポートのビジュアルのビルドを自然言語ですばやく実現

## 計算のビルド

文法を学んだり調べることなく、自然言語で計算式を作成

## ビジュアルの設定

設定したいフォーマットを自然言語で表現することにより、ビジュアルの更新をすばやく実現



# AIがオンデマンドでデータへの質問に答える

DEEP INSIGHTS AT YOUR FINGERTIPS

## ダッシュボードのエグゼクティブサマリー

自然言語によるダッシュボードの主要なインサイト（上位ムーバー、外れ値など）を速攻で要約

エキスパートでなくても容易に、パワフルなQ&Aが可能

Suggested questionsとwhat's in my dataにて質問例が提示

複数のビジュアルによる返答とインサイトサマリーのナラティブにて答えのコンテキストが説明

曖昧な質問サポートやdid-you-meanによる代替案にて、ファクトを繰り返して発見できる

The screenshot displays a dashboard with an AI chat interface. The chat window shows a user asking "leona eldridge conversions by month" and the AI responding with a detailed analysis of conversion rates, a line chart, and a bar chart. The bar chart shows conversion rates by source: LinkedIn Ads (0.32), Email Marketing (0.23), Facebook Ads (0.19), Influencer Marketing (0.18), Mobile App Ads (0.15), Google Ads (0.14), SMS Marketing (0.06), and Organic Search (0).

Account Manager	Source	Total Impressions	Total Conversions
Leona Eldridge	Email Marketing	6,017	28
Leona Eldridge	Facebook Ads	12,393	43
Leona Eldridge	Google Ads	2,393	7
Leona Eldridge	Influencer Marketing	10,280	30
Leona Eldridge	Instagram Ads	42,915	62

# AIアシスタントによる データストーリー

## データ解釈の共有

データの意味を説明し意思決定へ導くための結論を示す

## AIを使用してストーリーを生成

数個の単語だけで、データを分析し、親密でパワフル、洞察力あるナラティブを生成

## 洗練されたコンテンツの作成

AIによる長い文章の制御、ナラティブのテキストのカスタマイズ、ビジュアルのテーマ適用により、コンテンツが洗練

## 統制され常に最新の状態

すばやく容易に共有でき、いつでもデータを更新可能

QuickSight | A Data-Driven Journey: Optimizing Marketing Perfo...

File Edit Insert PREVIEW SHARE

## A Data-Driven Journey: Optimizing Marketing Performance Through Insights

Prepared by [Name]

### Introduction

As we examine how to maximize success across our marketing efforts, we will analyze campaign results over time to understand what strategies have been most effective. By identifying our goal is to optimize processes and targeting approaches to continue improve

### Revenue and Conversions Leader Board

Revenue and Conversions Leader Board				
Nestor Pallin 206	Leona Eldridge 183	Tasha Doug... 156	Karol Brown... 152	Celka Dui... 139

Build story Beta

Describe your data story in simple language and add the visuals you want to include.

Build a story about marketing campaign performance over time. Describe top campaigns and account managers. How can we improve overall campaign success?

Revenue and Co... Revenue by State  
Impressions and... Conversions per ...

+ ADD VISUALS

BUILD

# カスタムアプリにAIを拡張

QUICKLY ADD GENERATIVE BI CAPABILITIES TO WEB APPS

## 統合を簡単に実現

数行のコードのみで、AIの経験なくても既存のWebアプリに複数ビジュアルを返すQ&A機能を追加

## どのデータでも利用可能

特定のユースケースとデータのコンテキストを理解するため、質問を再定義したり、代替の質問が提案される。

## シームレス・エクスペリエンス

ブランドに合わせ、サイズや場所、インタラクションをカスタマイズ可

Data Q&A

Software Sales | PINBOARD | SHARE | FEEDBACK

Top customers by sales [ASK]

Interpreted as: Top customers by sales

There are 99 unique customers that contributed to a total sales revenue of \$2,752,804. Of these customers, Anthem has the best performance, generating \$60,228 in total sales, while Safeway ranks at the bottom with total sales of \$8,103. Anthem is responsible for the single largest deal on record; an order in the amount of \$23,460 for the product Big OI Database on March 18, 2019, the transaction was completed by Diane Murray in Germany.

Review for accuracy

Did you mean...

Top customers by total usage sessions

Top customers by profit

Sales by Country

\$2,752,804 Sales | 99 Customer

Sales by Customer

Order ID	Order Date	Customer Contact	Customer Name	Region	Country
API-2018-103800	Jan 4, 2019	Brian Dickens	Siemens	APJ	Japan
API-2018-112326	Jan 5, 2019	Liam Churchill	CVS Health	APJ	Australia
EMEA-2018-141817	Jan 6, 2019	Ella Springer	Mitsubishi	EMEA	France
EMEA-2018-106654	Jan 7, 2019	Keith Coleman	Exxon Mobil	EMEA	Finland
EMEA-2018-167199	Jan 7, 2019	Andrea Hughes	Royal Dutch Shell	EMEA	Ireland

Powered by QuickSight Q

# まとめ



# まとめ

**なぜ自社データを生成 AI に活用させる必要があるのか？**

他社との差別化や、自社のアプリを最適化するために自社の持つ1st party データを生成 AI に利用することで、AI が他社にはないドメイン知識を持ち、タスクの精度が上がる

**プロンプトエンジニアリング / RAG**

生成 AI へ適用するデータ活用で最も簡単に試せる方法はプロンプトエンジニアリング。RAGを利用することで質問された内容に類するドキュメントを検索して回答に利用できる。

**基盤モデルのカスタマイズ**

少量データで特定のタスクに特化させる Fine-tuning や、大量のラベルなしデータを用いてドメイン知識をつける Continued pre-training などの方法で既存モデルを自社データでカスタマイズできる。

**可視化と生成 AI**

QuickSight を利用することで自然言語によるダッシュボードの作成や、データレポートを自動で生成し、既存アプリケーションに埋め込むことが可能に。

# Thank you!

