

# Azure Alを活用した RAG の精度改善

日本マイクロソフト株式会社 Asia Data&Al GTM Manager 小田健太郎



## **Microsoft**

日本マイクロソフト

Asia Azure GTMチーム GTMマネージャー

小田 健太郎

Kentaro Oda

2018年より日本マイクロソフト入社、パートナーマーケティング、業界別の製品戦略リードを経て、2021年よりデータ分析・AI・機械学習製品のプロダクトマーケティングマネージャーとして従事。現在はグローバルアジアチームに所属し、コアプロダクト「Azure AI」の日本、韓国のGo-to-market戦略をリード。



# アジェンダ

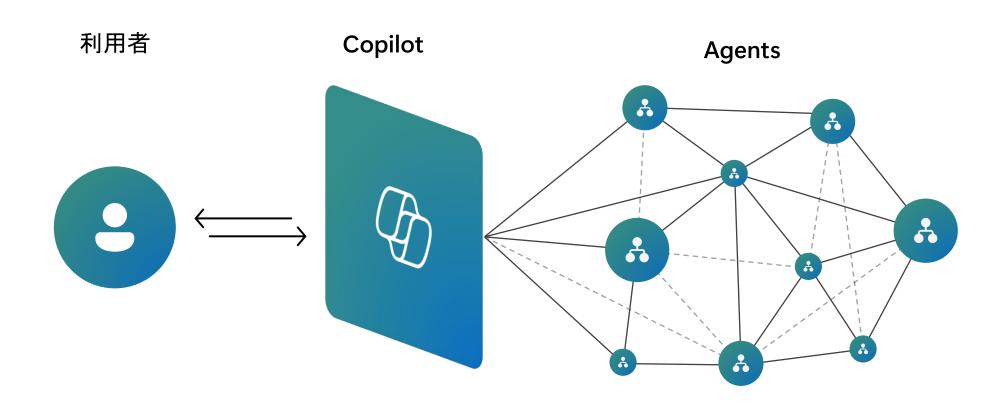
- 1. マイクロソフトが取り組むAIについて
- 2. 精度改善に向けたアプローチ
- 3. 製造業におけるAI利活用の現在地
- 4. 今後どうなる?

# アジェンダ

- 1. マイクロソフトが取り組むAIについて
- 2. 精度改善に向けたアプローチ
- 3. 製造業におけるAI利活用の現在地
- 4. 今後どうなる?



# CopilotとはAIのUI



## **Front Door**

利用者はCopilotをフロントサービス として利用、技術的な細部を意識 せずAIへのアクセスが可能に

## **User Interface (UI)**

利用者からの入力/出力を担い、 AIモデルやエージェントとの橋渡しを 任意/定期的に行うハブとして機能

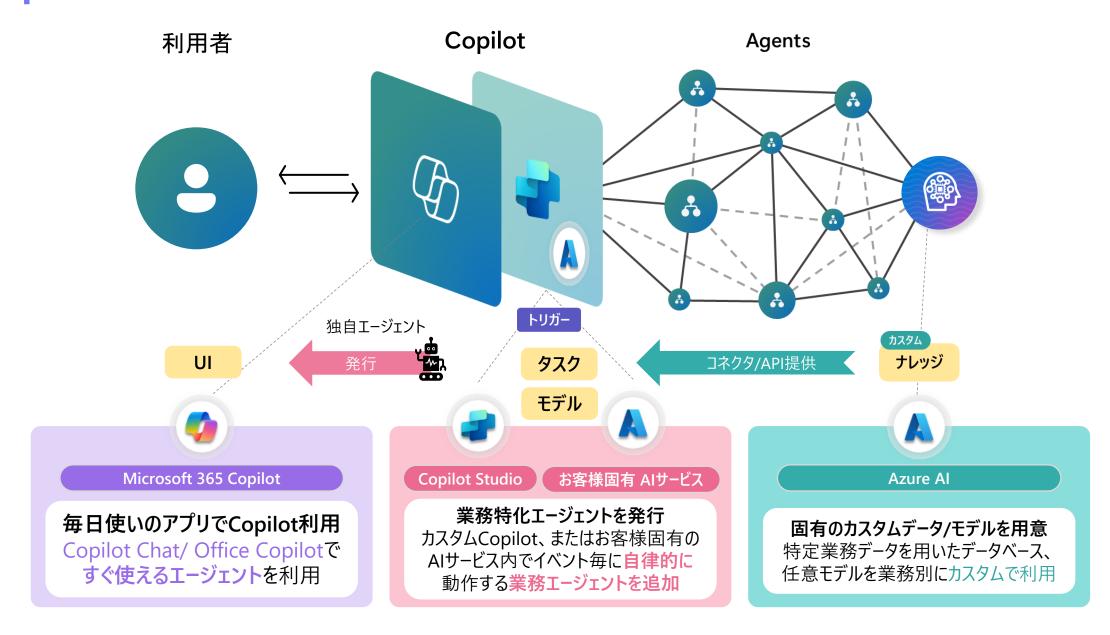
## Core models/actions

ユーザーの質問に対して最適な回答を行うための、必要なデータや知識を検索・ 設計し、最適な回答や結果を導く

# スキルとニーズで使い分け



# CopilotとはAIのUI - 各挙動のイメージ



# スキルとニーズで使い分け



オフィスワーカー向け

すぐ作ってすぐ使える ベーシックAI

Microsoft 365 Copilot

市民開発者向け

業務特化型 ローコードカスタムAI

Copilot Studio with Power Platform



開発者向け

業種特化型 フルカスタムAI

Azure Al

M365 Apps からユーザーが簡単に作成

業務向けエージェントをローコードで作成

高度にフルカスタム・部分/全体最適化

UL

M365 Application (社内)

業務アプリケーション(社内+社外)

UI / プロンプト / イベントベース自律動作

Any組込み・新規アプリ (社内+社外)

トリガー

UI/プロンプト

Any & カスタム可

モデル

ナレッジ

標準(最新モデルに準ずる)

標準/カスタム可 LLM+生成オーケストレート●

Any & カスタム可 LLM

SharePoint / OneDrive

Webサイト / Graph コネクタ

Power Platform コネクタ / Dataverse

カスタムプロンプト / カスタム会話 / アクション

Sales/Service ロールベース アドオン可

SharePoint / OneDrive

Webサイト / Graph コネクタ

Power Platform コネクタ / Dataverse

カスタムプロンプト / カスタム会話 / アクション

カスタム可

Any & カスタム可 データソース

Power Platform マネージド環境 / 管理センター

**Any & カスタム可 環境** 

運用

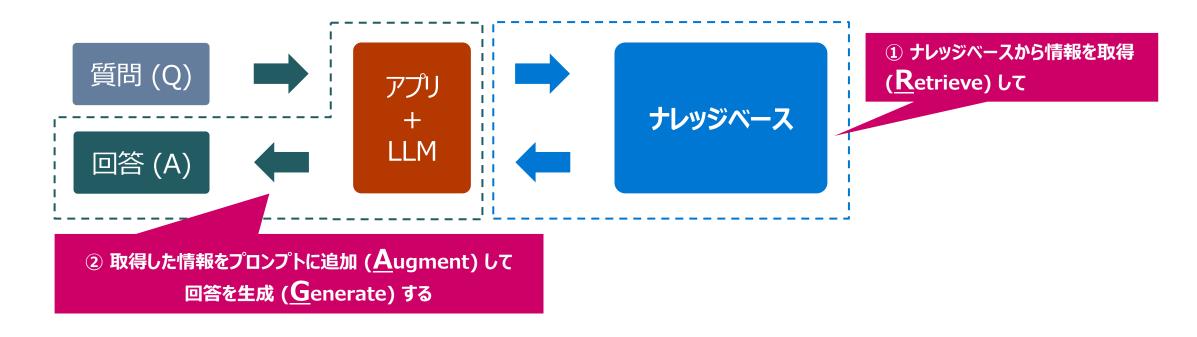
M365管理センター

## アジェンダ

- 1. マイクロソフトが取り組むAIについて
- 2. 精度改善に向けたアプローチ
- 3. 製造業におけるAI利活用の現在地
- 4. 今後どうなる?

# RAG (Retrieval Augmented Generation)とは?

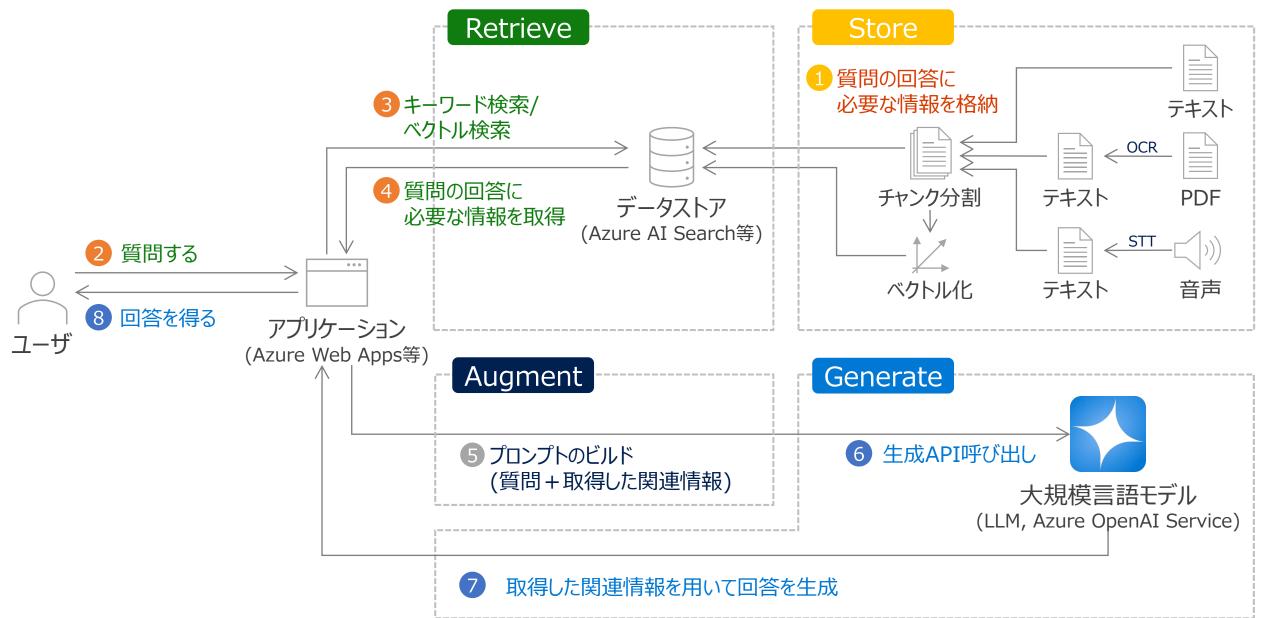
Meta の研究者によって考案された大規模言語モデル (LLM) によるハルシネーションを低減するための手法。 一言で表現するとナレッジベースの外部化、RAG ではユーザーからの問いに対してバックエンドに置いたナレッジベース (Azure Al Searchなど) で検索を行い、その結果を引用してユーザーのからの問いかけに対する回答を生成する。



#### 参考:

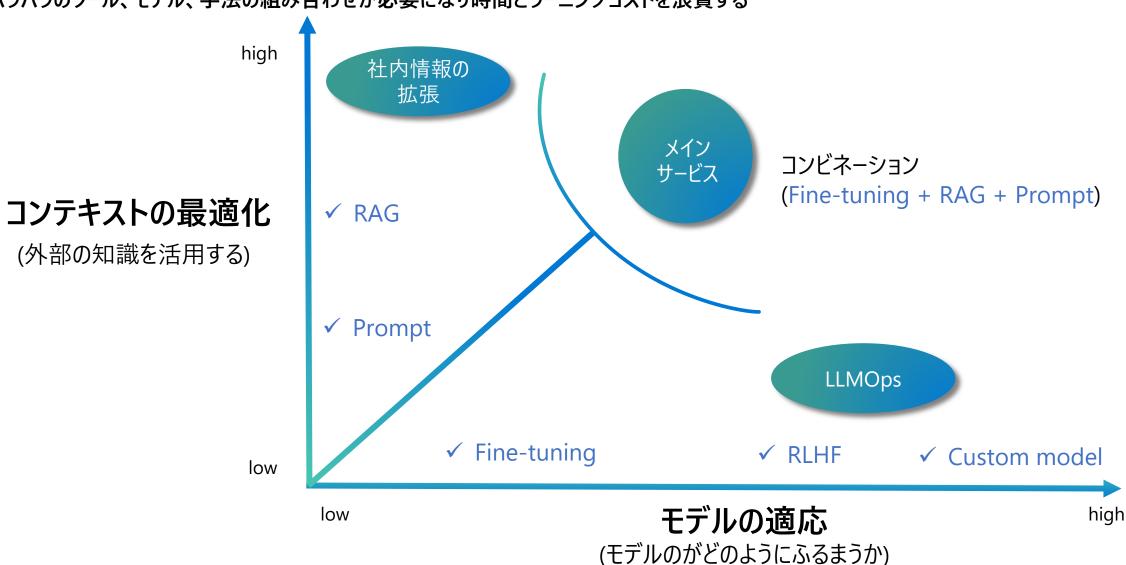
- Retrieval Augmented Generation: Streamlining the creation of intelligent natural language processing models (meta.com)
- 検索により強化された生成 (RAG) | Prompt Engineering Guide (promptingguide.ai)
- <u>幻覚 (人工知能) Wikipedia</u>

# RAG (Retrieval Augmented Generation) の動作



# AI利活用に対してどこから始めるか?

エージェント以外にも検討事項や要素技術の選定を考慮する点が多い。結果としてデータ整備、 バラバラのツール、モデル、手法の組み合わせが必要になり時間とラーニングコストを浪費する



# AI利活用に対してどこから始めるか?

エ−ジェント以外にも検討事項や要素技術の選定を考慮する点が多い。結果としてデ−タ整備、 バラバラのツール、モデル、手法の組み合わせが必要になり時間とラーニングコストを浪費する Optimization Flow

high All of the above RAG RAG コンテントを トレーニングデータ Add HyDE retrieval + に追加 Fact-checking step Fine-tune model Add simple retrieval 別のトライ をする **Prompt Engineering** Fine-tuning Add Few Shot Prompt 評価 トライする モデルの適応 low high

コンテキストの最適化

(外部の知識を活用する)

low

(モデルのがどのようにふるまうか)

# RAGで精度を出すために検討が必要なポイント

以下の各ステップの検討・チューニング・評価を繰り返すことが回答精度向上につながる

## Store

PDFファイルの テキスト化

Officeファイルの テキスト化

音声データの テキスト化

動画データの テキスト化

画像データの テキスト化 チャンク分割 (チャンクチューニング)

オーバラップ

テーブル構造抽出

エンリッチメント

カテゴリ付け

## Retrieve

検索アルゴリズム

アルゴリズムの パラメータ チューニング

スコアリング プロファイル

カスタム アナライザ

類似度 チューニング

## Augment

システム メッセージ定義

ユーザ メッセージ定義

検索クエリ生成

ユーザへの聞き返し

仮説的文書埋め込み

**Function Calling** 

## Generate

モデル選択 (回答生成)

モデル選択 (埋め込みモデル)

マルチモーダル

ファイン チューニング

## Evaluation

# ステップごとのRAGの精度影響因子

精度向上施策を打つ前に、原因を特定することが極めて重要

## 入力情報の加工

ユーザ入力に検索のための 情報が足りない、整理されていない

#### AI技術の強みを教えて。



#### ドキュメント・クエリマッチング

入力された内容と検索対象が 意図した類似になっていない

#### 初心者用バットがほしい



類似はしているが 意図が拾えてない

## 類似度ヒットしたドキュメント

初心者用バット、ほしいですよね!

子供が初心者バットを欲しがってる

壊れてしまった初心者用バット

## 検索実行

検索エンジンの精度が悪い。

スコア

検索対象ドキュメント

- **0.702 1 初心者**でも扱いやすいように、特別に軽量化されています。
- 0.**401 2** かなり振りやすいので初めてでも 扱いやすい**バット**といえます。
- 0.780 3 クリケットの初心者は、バットと同じ要領でスイングしてしまいます

#### コンテキストベース回答

検索結果を 正しく解釈できていない

## プロンプト

User Questionに回答せよ。

- # User Question たまごはコレステロールが多く健康に悪いですよね?
- # Context

表1のように卵の摂取は長年健康への悪 影響が懸念されていた。しかし、

> 途中で文章が切れている。 図表が取り込めていない

#### 对 策

原

大

- > 一般情報+聞き返し指示
- ▶ 検索実行条件プロセス定義
- ▶ クエリ拡張

- ▶ ドキュメント加工
- > 類似度チューニング

- > ハイブリッド検索
- > フィルタリング
- > リランク

- > チャンクチューニング
- > データの構造化
- ▶ コンテキストの要約

# 製造業における変革のトレンド

## データがすべての基盤

企業内のすべての現場の運用 プロセスにデータとAIを導入する (デジタル運用プラットフォーム)

## "人"中心に考えたAI活用

AIは、先進的な技術を活用して 人間の作業をサポートし、向上 させることで、従業員の健康を 保護し、人の総合力を強化する

## フィードバックループの確立

トレーサビリティや材料の品質、 生産計画の変更など、 デジタルでの分析から現実への フィードバックループを確立する

# アジェンダ

- 1. マイクロソフトが取り組むAIについて
- 2. 精度改善に向けたアプローチ
- 3. 製造業におけるAI利活用の現在地
- 4. 今後どうなる?

## データのサイロ化を解消し、AI活用を見据えたデータ基盤を構築

## <Before>

管理とセキュリティに リソースを大量消費

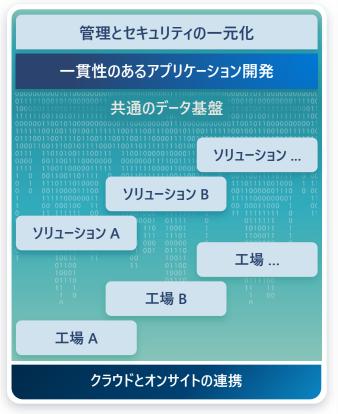
グローバルなオペレー ションをエンドツーエン ドで可視化できない

> サイロ化によるステー クホルダー間のコラボ レーションの欠如

ソリューション間での データとインサイトの 統合がない

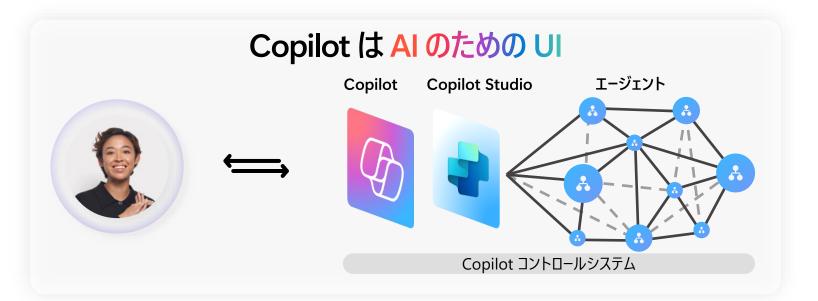


## <After>





# AI エージェント によるビジネスプロセスの変革



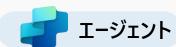




人間の拡張

プライベート、パーソナルアシスタント

すべての従業員に Copilot を



Copilot に接続または自律

タスクやプロセスの自動化

すべてのビジネスプロセスにエージェントを



エージェント

急な要求に対応する製品の在庫と適正な在庫を持つ

業務特化フルカスタマイズ

タスクやプロセスの自動化

すべてのビジネスプロセスにエージェントを

# 自然言語での工場データの操作

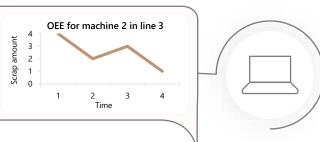


昨日ライン3のOEEが下がった 理由を教えてください。

**ライン3**のOEE のステータス概要を示します。 可用性、パフォーマンス、および出力データを使用 します。このデータは、マシン 2 のパフォーマンスが 低下していることを示唆しています。



以下はその分析です 履歴データに基づくマシン 2 の可用性:



## 曖昧さ回避とデータのリレーション

データモデルは固有の製造設備に基づいて構築。 工場固有の略称を認識し、ライン3 に存在する可能性の あるマシンなどのエンティティ関係を理解します。

## 根本原因の特定

Copilotは、迅速にLLM推論機能を活用。 製造オントロジーに基づいて階層構造で形成された マシンのデータを特定します。

## 履歴データの分析

Copilotは、ユーザーが探している履歴データの種類を理解し、エンティティとリレーションシップを探索して分析情報を提供します。

## 製造業におけるM365 Copilot利活用シナリオの発掘例



主要プロセス

#### AI 導入前

AI 導入後

アフター サービス

技術者が入手できる情報が限られていることで、サービスの遅延、不正確な診断、コールバックが発生し、損失を伴うダウンタイムや顧客満足度の低下を招いている

メンテナンス活動の実施

Copilot が履歴データと診断に基づいてメンテナンスに関する助言を提供し、問題を的確に特定できるようにする。 さらに、技術者を専門家につなぎ、その場で支援を受けられるようにする

予測の収集と要約、メールや会議を通じたチームフィードバックの整理、S&OP会議の事前準備などを手作業で行っている

販売·事業計画

Copilot が利害関係者からのフィードバック収集を促進し、最新の予測や要約を生成することで、S&OP プロセスの 合理化を支援。レビュー中に、Copilot が重要な意思決定を記録し、アクション プランを策定して、タスクを効率的に割り当てる

製诰

プラント マネージャーが資産のパフォーマンスを 手動で確認し、メンテナンスや診断を管理している

資産管理の最適化

Copilot が、パフォーマンスのモニタリング、問題の診断、メンテナンス手順のガイド、 メンテナンス後レビューの簡素化を支援。

従業員、ツール、リソースの計画と管理を手作業で 行っているため、稼働率と処理能力が低下している 工場および倉庫のタスクのサポート

Copilot がデータを分析し、体系的なアイデア創出を促進し、デザイン スケッチを生成して、フィードバック収集を効率化

ガイド付きデモ (英語) ⊃

録画デモ⊃

研究開発

データ分析、非体系的なブレーンストーミング、プロトタイプ 作成、フィードバック収集を手動で行っており、非効率 新製品のアイデア創出

Copilot がデータを分析し、体系的なアイデア創出を促進し、デザイン スケッチを生成して、フィードバック収集を効率化

録画デモ⇒

Azure OpenAIで進化する、三菱重工のデジタル戦略

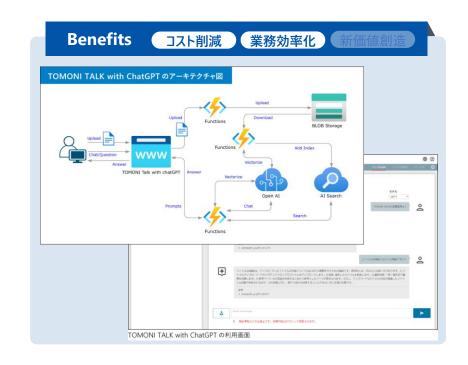
# 三菱重工、Azure OpenAI を活用した独自アーキテクチャでデジタルイノベーションを加速



## 三菱重工業株式会社

## 活用イメージ

- ●三菱重工業株式会社は、生成AIとAzure OpenAIを活用して、業務効率の向上とデジタルイノベーションを推進している。2023年春に立ち上げられた「TOMONI TALK with ChatGPT」プロジェクトは、その生成AIの社内活用を推し進める原動力となった。Azure OpenAI Serviceを導入することで、チャットへの入力内容、プロンプト、社内情報が外部へ流出する懸念を解消した。
- ●Azure OpenAlを活用して、要約や翻訳機能を備えたチャットアプリケーションが試作され、2023年7月に社内向けにリリースさた。リリース直後から社内ユーザーが急増し、トラブルシューティングやメールの整理、会議室予約機能など、生成Alの活用アイデアが次々と寄せられた。
- ●さらに、Azure OpenAl Service on your dataを評価し、RAG (Retrieval Augmented Generation)手法を用いることで、社内保有データを活用した高精度な回答生成が可能となった。



導入事例の詳細はこちら

https://customers.microsoft.com/ja-jp/story/1772636025129504407-mhi-azure-openai-service-discrete-manufacturing-ja-japan

沢井製薬、「SAWAI DX」により業務プロセス改革に取り組む

# Azure OpenAl Service を活用し、セキュアな環境で研究員の相談役となる Al エージェント



## 沢井製薬株式会社

## 活用イメージ

- ●沢井製薬は、生成 AI を活用した業務プロセス改革に取り組んでいる。製剤研究部では、過去の膨大な報告書データを活用しきれず、情報検索に多くの時間を費やしていたところ、Azure OpenAI Service を導入し、ベテラン研究員のように相談相手として機能するAIエージェントを構築した。質問を解釈し自律的に回答を生成する仕組みを実現しましたことで、部門全体での利用頻度は月間 100 回に及び、社内でも高い評価を得ている。
- ●生成 AI の導入により、若手研究者や部署異動者の成長速度が向上し、製剤開発の生産性や品質が向上した。特に、生成 AI がリファレンスを提示することで、研究者が必要な情報に迅速にアクセスできるようになった。
- ●また、Azure OpenAl Service はクラウド上で閉域網を構築し、セキュリティを担保した環境で利用できるため、機密情報の漏えいを防ぎつつ先進的な生成 Al を活用できる点が評価されている。



関連詳細/公式情報はこちら

https://www.microsoft.com/ja-jp/customers/story/21162-sawai-pharmaceutical-azure

# アジェンダ

- 1. マイクロソフトが取り組むAIについて
- 2. 精度改善に向けたアプローチ
- 3. 製造業におけるAI利活用の現在地
- 4. 今後どうなる?





パーソナル エージェント



組織のためのエージェント



ビジネスプロセスの ためのエージェント



組織間を超えた エージェント

# 精度と能力の向上

AIの精度及び能力を向上させる手法を理解し、AIを構成することが重要。自社のデータを適切にと組み合わせることで、業務プロセスにおけるAI活用の効果が最大化します。

"道具を与える" (エージェント)

- 外部Tool(予測AI、Code InterpreterのPython処理、外部APIなど)を道具としてAIエージェント に与え、予測結果を考慮したプロアクティブな行動や動的な計算処理が可能
- エージェントに予測AIを組み合わせることで、リアクティブ(事後対応)的な動作からプロアクティブ (事前対応)的な動作へと高度化。 例:需要予測に基づく生産計画の立案など

精<u>度/</u> 能力向上 "外部知識を与える" (RAG)

- 企業が持つ社内データ(業務ナレッジなど)をAIエージェントが参照(検索)できるように整備することで、社内データを活用したAIエージェントのタスク実行が可能となる
- ベーシックな検索エンジン/ベクトルDBを利用したNative RAGだけでなく、ナレッジグラフを利用したGraphRAGを掛け合わせたHybrid RAGにより、AIエージェントの精度向上が見込める

"役割を与え 振る舞いを整える" (プロンプト エンジニアリング)

- システムプロンプトやAIエージェントの役割ごとの指示を設計し、LLMが望ましい回答を出すように誘導。 ユースケースやモデルの特性に応じたプロンプト設計が必要
- AIエージェントにおいては、プロンプトをテンプレート化、パイプライン化し、業務プロセスを自律的にAIが 実行できるように調整。ロングコンテキストを活用した知識背景の埋め込みも検討

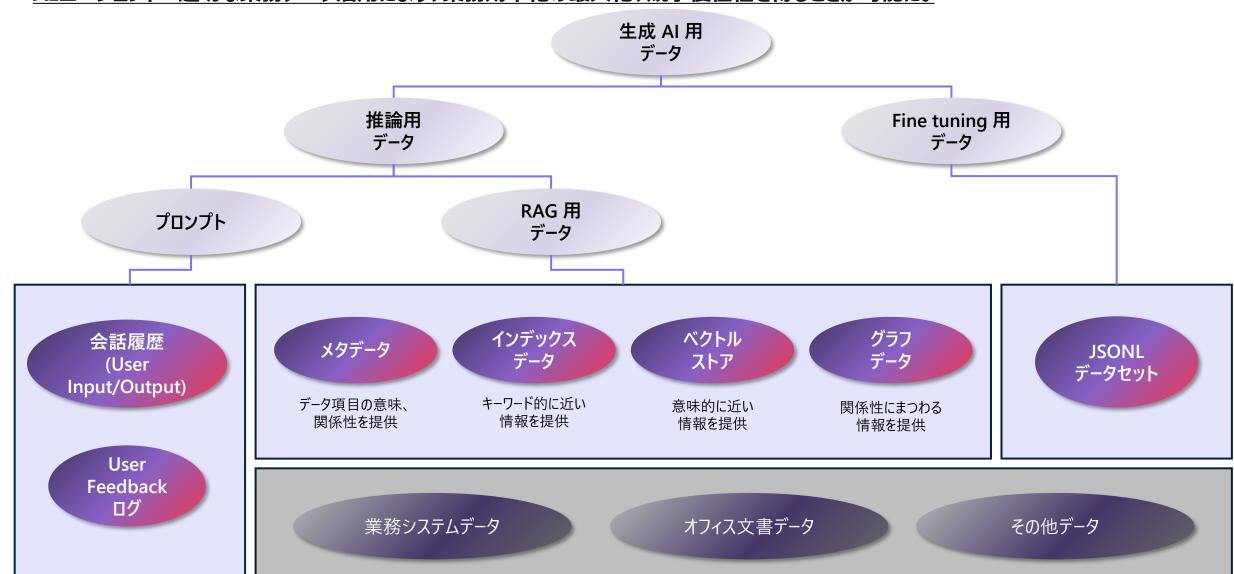


"Alを内部から鍛える" (Fine-Tuning)

- LLM/SLMに対する微調整/追加学習。特定のデータセットやドメインに適合するように予め生成AIに 学習。ラベル付きの学習データを手動で作成するだけでなく、教師用モデルを利用した合成データ(蒸留) の活用も検討。
- 出力形式・トーンの調整、タスク精度の強化、トークンの節約などが見込める。

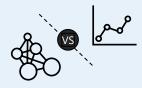
# AI 開発や利用を促進するデータの種類

業務特化型のAIエージェントを構築するには、エージェントが利用するデータを適切に管理し、活用する必要。 AIエージェント×適切な業務データ活用により、業務効率化の最大化、競争優位性を得ることが可能に。



# AI利活用企業が抱える課題 ※MSのお客様傾向

人/プロセス

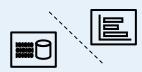


- データガバナンス体制・運用整備
- 社内のデータ活用に推進における協力体制の構築
- PoCから本番運用への脱却

データ

- 非構造データをどのように管理していくべきか?
- 生成AIによる回答精度をあげるために必要なデータ加工
- 機密データの扱いとアクセス管理

テクノロジー



- 従来のデータ基盤をどのように拡充していくべきか?
- ビジネス部門からの要求に迅速に対応するための仕組み作り

# エージェントに飛びつく前に考えてほしいことD

- ・AI アプリ (AI モデル) のバックボーンとして、"データ" が回答精度を決定
  - 回答精度を支える中心要素は AI モデルであり、そのモデルを支える"土台"がデータ (RAG, Fine tuning etc)
  - データの品質・探索のしやすさ・セキュリティの確保が早期に判断できるかどうかが成否を左右
- •非構造・半構造データに対するニーズの高まり
  - 文章、画像、音声、ログデータなど、多様かつ大量の半構造・非構造データを扱う必要性がある
  - •構造化データと異なり、非構造・半構造データに対するデータガバナンスは未成熟。 ナレッジ不足によって意思決定の判断や対応が鈍化し、失敗リスクが高まりやすい。
  - •また、そのデータを主に扱うデータサイエンティスト / アナリストに対する理解・共感が不十分という企業課題が多い。 つまり、どういうデータを提供すれば喜ばれるか (AI 開発が加速するか?) の理解に乏しく、これも鈍化リスクを助長する。
- AI 向けのデータガバナンスとデータマネジメントの知見を持つかどうか?が "ミッシングピース"
  - •データガバナンス = (例) 料理店舗のコンセプト・ポリシーを策定すること
  - •データマネジメント = (例) 実際の調理を実行すること

## データガバナンスとデータマネジメントの目的の違い

- データガバナンスは、データ利用者が効率的にデータを利用するための、一連の原則、標準、及び管理監督です。 組織体制、ルールとプロセス、手段である製品機能によって実現します。
- 一方で、データマネジメントは、データを収集、保存、管理、分析、活用するための日々の運用や技術的な活動を指します。

## データガバナンスの主な目的:ポリシー策定

- **データ利用・管理の標準ガイドライン策定** データの利用や管理に関するガイドラインを統一化する。
- データ在りかの可視性向上 データ資産の位置、所有権、ライフサイクルを明確化する。
- データ品質の確保データの正確性、一貫性、完全性を確保する。
- データのセキュリティ・ポリシー策定 個人情報や機密データを適切に保護すべく、 準拠すべき法規制 (例:GDPR、CCPA) の策定
- **意思決定者の透明性を担保** データ利用に関する責任やルールを明確にし、 利害関係者間の合意を形成する。

## データマネジメントの主な目的:処理自体

- 多様なデータ・クレンジング処理の実装 生成 AI に必要なデータの保存、処理、 配信を最適化する。適切な処理エンジン選択。
- AI 活用促進に向けた
   構造・非構造データの配置最適化
   AI 開発や分析のために、データを利用しやすくする。
   具体的にはチューニングや最適化されたフォーマット選定。
- データセキュリティ
   データの機密性をマスキングしたり、アクセス制御をポリシーにアラインする形で具体的に実装する。

配置先のデータストア選択

• データの保全とバックアップ、可用性 データの損失や破損を防ぎ、復旧可能な体制を整える。

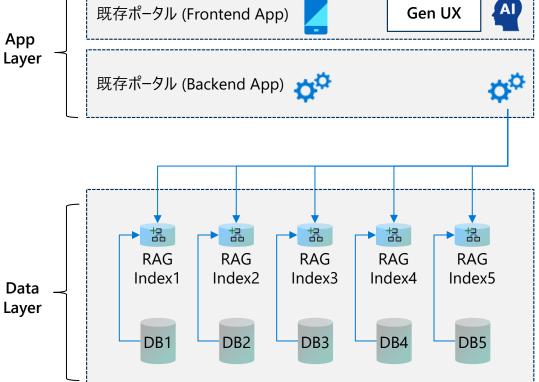
## AI 活用に向けたデータアーキテクチャのアプローチアイディア

AI 活用に適したデータ領域新設によって、リスク低減された AI Agent アプリ開発を推進 (パターン① < パターン②)

## パターン①:現行データ基盤拡張

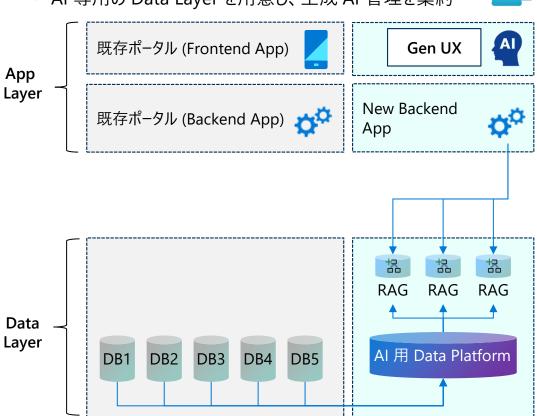
- 既存環境制約に該当し、開発/運用コスト増になる傾向あり
- 既存アプリとの依存性が高く、リリースタイミングの遅延あり
- 生成 AI 用データの管理面で煩雑化するリスクあり





## パターン②: AI の近場に AI 活用向けデータ領域新設

- 既存環境から切り離し、技術的環境制約リスクを最小化 (開発コスト最適)
- 既存アプリとの依存性が低く、柔軟なリリースタイミング確保
- AI 専用の Data Layer を用意し、生成 AI 管理を集約

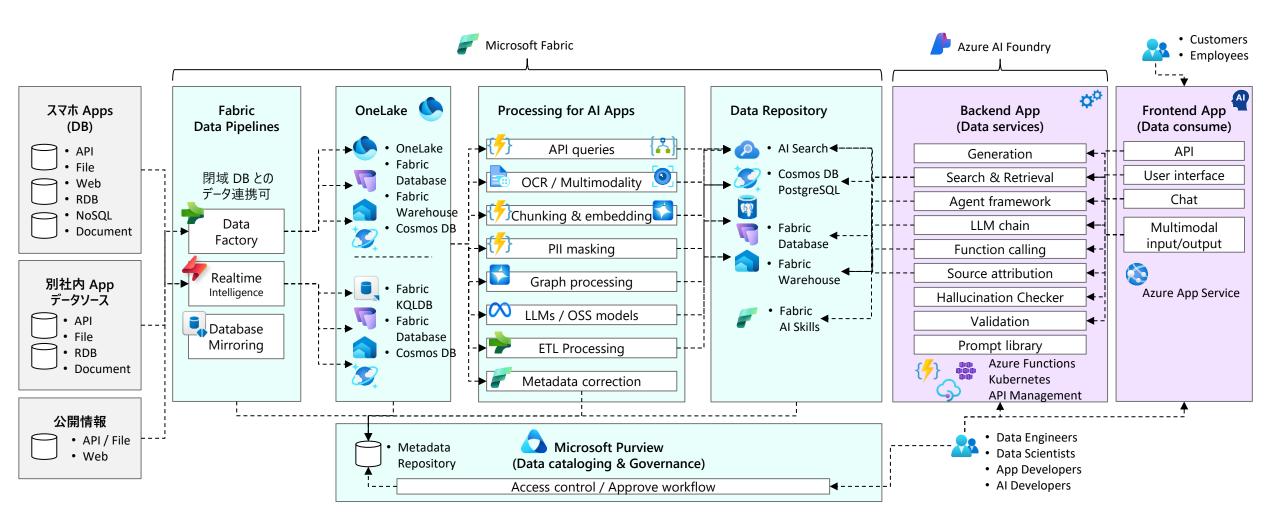


# データ特性/活用用途に合わせたサービス選択

## 適切なDataサービスを選択することで、AIエージェント構築時の効率性を最大限に高めることができます

カテゴリ	対象データ	データ特性/活用用途	サービス例
Fine-Tuning	ラベル付き学習データ /合成(蒸留)データ	システムプロンプト、ユーザプロンプト、回答のセットを学習デー タとしてJSON形式で整備	Azure AI Foundry Microsoft Fabric(OneLake) / Azure Data Lake Storage
プロンプト/メモリ	会話履歴	LLMおよび各Toolで呼びされるAPIへの入力と出力	Azure CosmosDB ※RDB系のJSON型も可
	プロンプトストア	エージェントが共通で利用するプロンプトテンプレート	
	ツール情報/API仕様	エージェントが利用するTools(API)の仕様情報	
RAG	検索インデックス	全文検索(キーワード)、ベクトル検索、セマンティック検索に対 応した高度な検索インデックス	Azure AI Search
	ベクトルDB	RDBおよびNoSQL DBにおけるスケーラブルなベクトル検索。 DBならではの信頼性のあるトランザクショナルな検索	Azure SQL / Azure CosmosDB / Azure DB for PostgreSQL
	ナレッジグラフ	グラフ構造(エッジとノード)を利用したナレッジグラフの構築。 GraphRAGでの利用	
	業務データ (Text to SQL)	業務データ(構造化データ)に対する自然言語での動的な問い合わせ(集計処理など)	Microsoft Fabric (AI Skill)
予測AI	業務データ (予測AI)	Lakehouse / Warehouseによるデータクレンジングと統合。 特徴量エンジニアを行い、予測AIモデルを構築	Azure Machine Learning / Microsoft Fabric/ Azure Databricks
	オープンデータ	オープンデータやSNSデーによる特徴量の拡張	

# Al Ready "データマネジメント"の重要性





Thanks!

