

Azure OpenAI Service
リファレンスアーキテクチャ
社内文書検索
2023/9/15

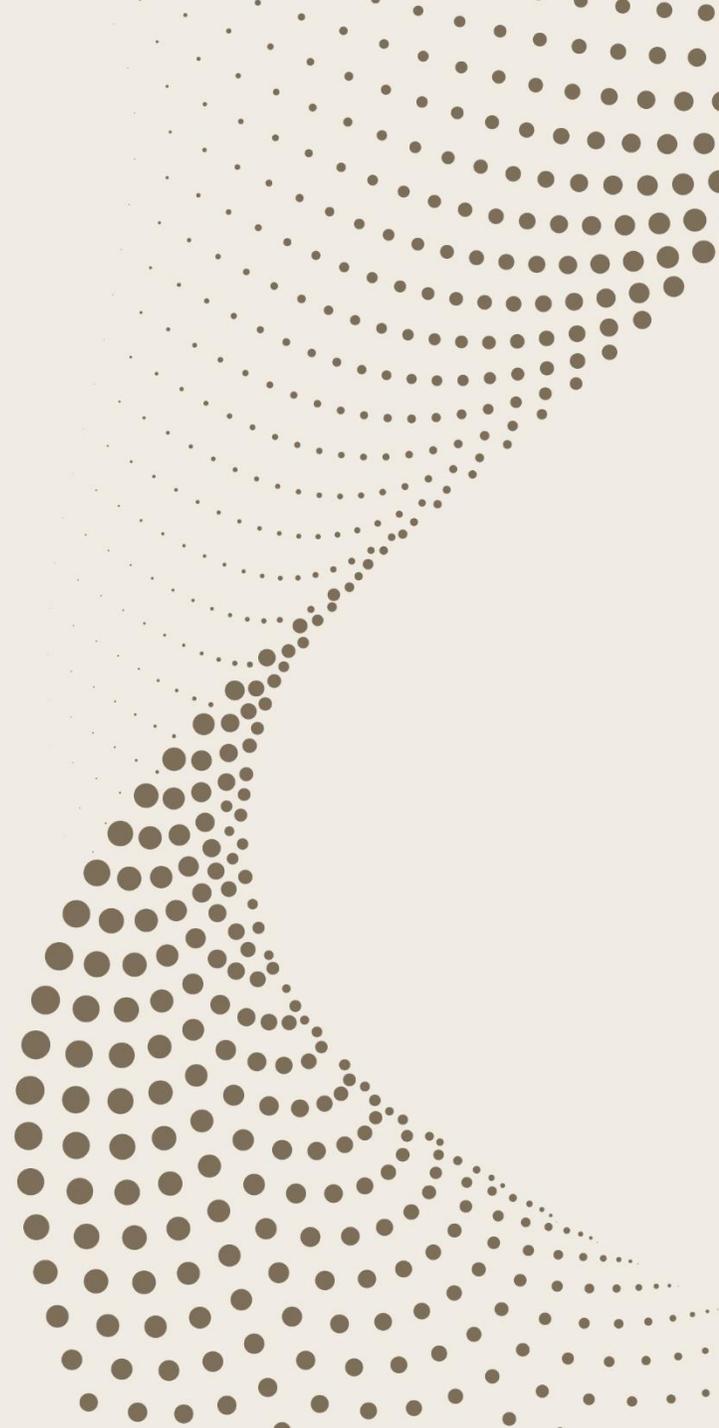


Build Beyond As One.

コンテンツ

1. シナリオ概要
2. デモUI
 - 完成イメージ
 - 現在の実装 メインページ
 - 現在の実装 チャットページ
3. アーキテクチャ
4. 考慮事項

1.シナリオ概要



1.シナリオ概要

- ChatGPTは様々な業種業界での活用が想定され、個人利用でも業務効率化や作業品質向上に繋がる可能性が大いにあります。
- 但し、個人の業務利用が進むと機密情報流出の恐れもあるため**セキュア**な社内公式ChatGPTを実装します。
- “社内ナレッジ検索”を先行着手し、順次、“議事録作成”、“リサーチ業務”、“開発業務”等の機能を実装予定です。

セキュリティ観点

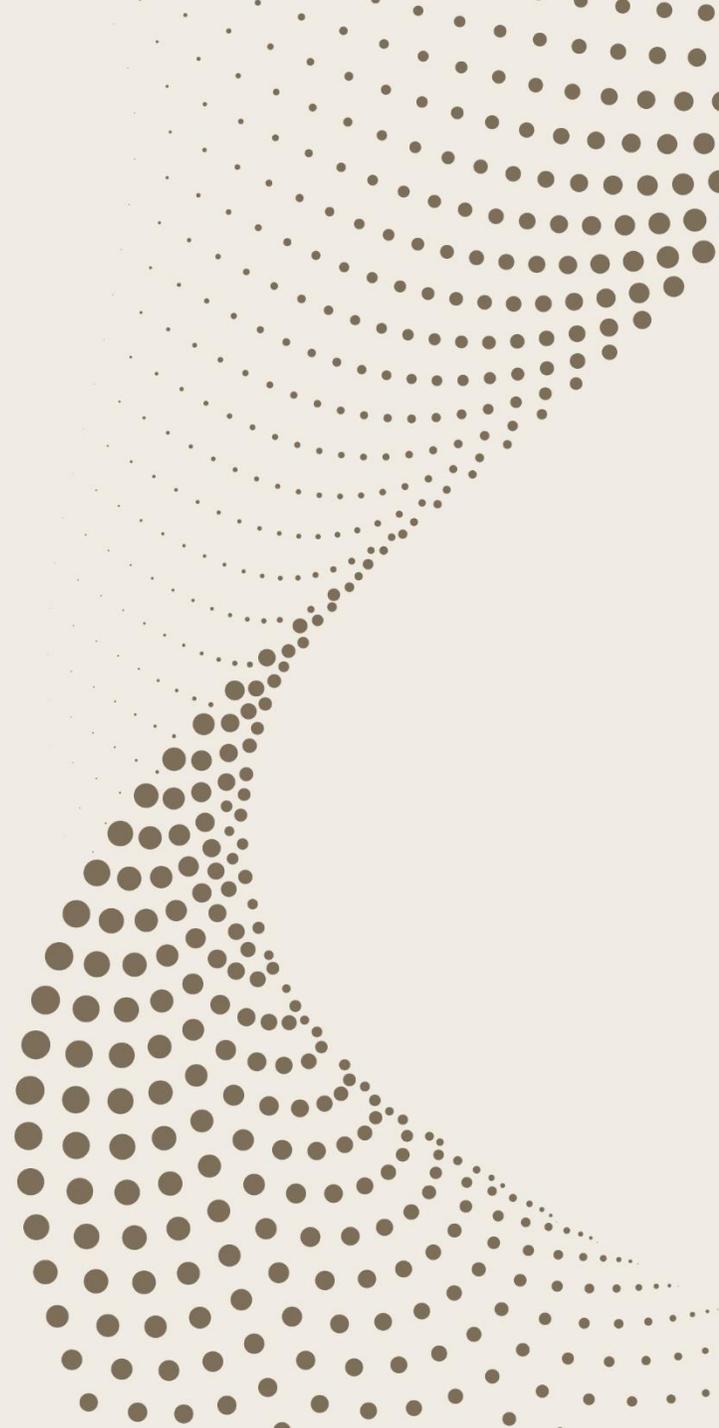
個人がOpenAI社のChatGPTを利用すると顧客情報漏洩等の事故が起こる危険性があるため、

- ①利用ルールを策定
- ②アクセス管理や監視機能のある安全な環境を整備

業務効率化観点

過去の提案書、成果物、プロジェクト開発標準などのナレッジを読み込ませることで、企業独自の情報に基づく回答作成や、議事録作成支援など社内業務の効率化をはかる

2. デモUI

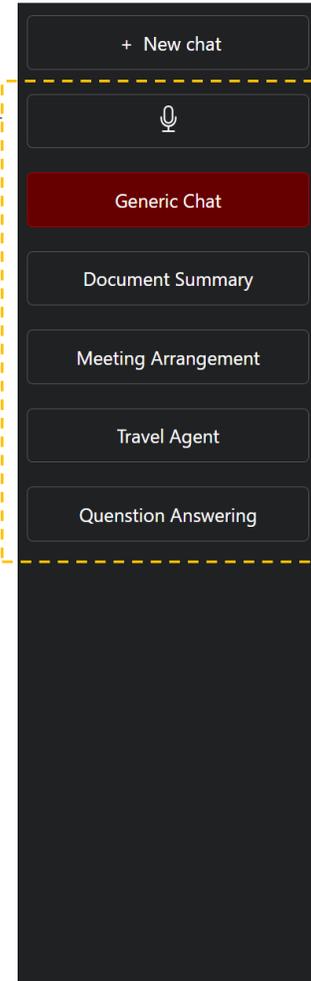


2. デモUI 完成イメージ

- ChatGPTに似た独自のフロントエンドを用意※1し利用頻度が多いユースケースはメニュー化する想定です。

※1.社員がPlaygroundを直接利用することは難しいため

サイドバーに様々な機能(メニュー)を追加していく想定



ABeam Chat LLM

☀ Examples	⚡ Capabilities	⚠ Limitations
Show how to list prime numbers with python	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
Generate poem titled 'a horse walking on the moon'	Sending prompt and response to OpenAI is avoided	May occasionally produce harmful instructions or biased content
What do I learn how to go with AI first?	Similar functionality as public ChatGPT website	Limited knowledge of world and events after 2021

Send a message...

This chatbot may produce inaccurate information about people, places, or facts.

2. デモUI 現在の実装 メインページ

- 現在は社内ナレッジを自然文で検索できるシステムを構築済みです。

現在はナレッジ
検索機能のみ実装

+ New chat

Search Document

ABeam Chat LLM

☀ Examples	⚡ Capabilities	⚠ Limitations
Show how to list prime numbers with python	Remembers what user said earlier in the conversation	May occasionally generate incorrect information
Generate poem titled 'a horse walking on the moon'	Sending prompt and response to OpenAI is avoided	May occasionally produce harmful instructions or biased content
What do I learn how to go with AI first?	Similar functionality as public ChatGPT website	Limited knowledge of world and events after 2021

Send a message...

This chatbot may produce inaccurate information about people, places, or facts.

2. デモUI 現在の実装 チャットページ

- 質問を行うとチャット画面に遷移し、検索結果が表示されます。

The screenshot displays a chat interface with a dark sidebar on the left containing a '+ New chat' button and a red 'Search Document' button. The main chat area shows a user's input: 'AIやIoTに関する提案事例を教えてください' (Please tell me about proposal cases related to AI and IoT). A response from the chatbot follows: '関連する過去問合せとして以下の事例が見つかりました。' (We found the following cases as related past inquiries). Two search results are shown, each with a title, URL, and key details:

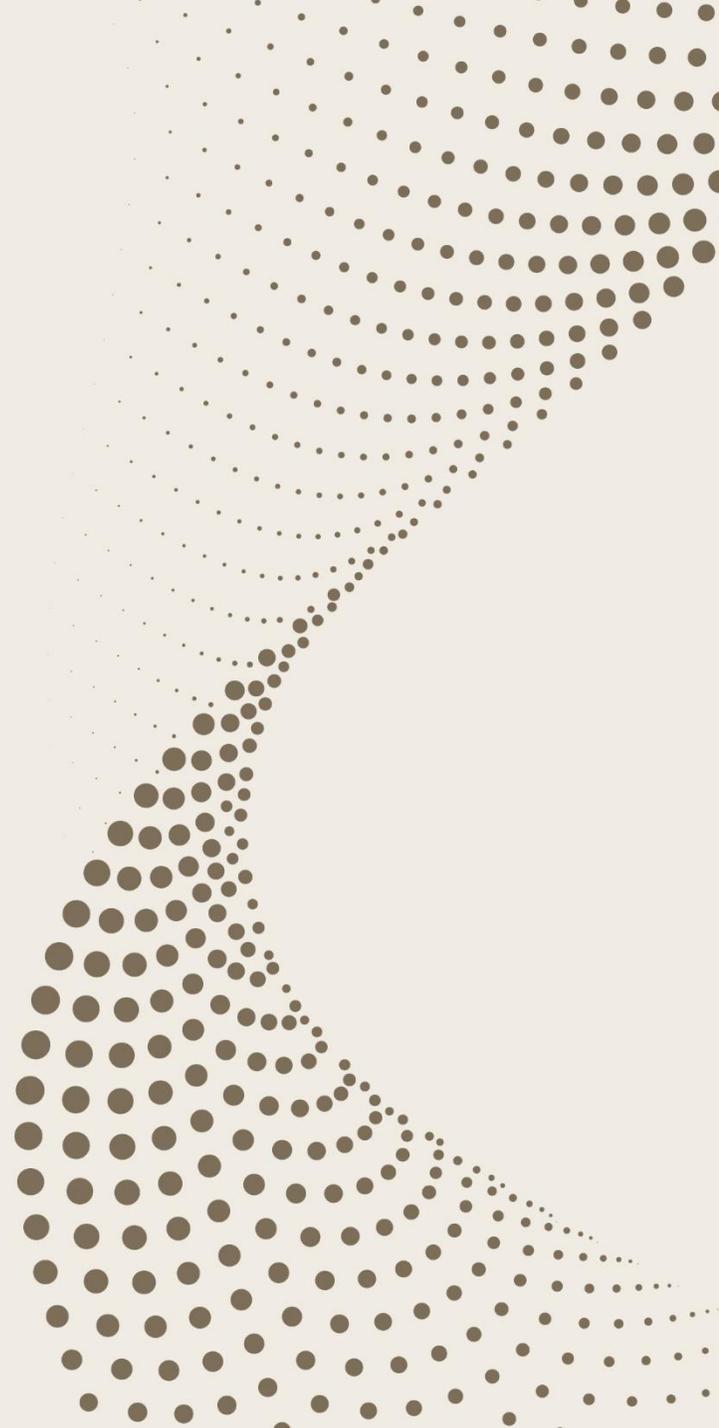
- [1]:000社 AIモデル導入事例走**
https://dummyurl/AI_01.xxxx
 - 提案日時: 2020年8月
 - クライアント: 株式会社000社
 - 業種: 自動車
 - 提案概要: 当社では以前、競争力の向上と効率化を図るため、AIモデルの導入に関する提案を検討しました。このAIモデルは、データ解析から顧客対応まで多岐にわたる業務に適用可能で、労力と時間の節約、意思決定の支援に役立つという期待がありました。また、カスタマイズ可能な設計により、ビジネスに合ったソリューションを提供できる点が魅力でした。当社はROI向上を期待し、競争力強化を目指していました。導入に関しては包括的なサポートも考慮し、スムーズな移行を回りました。このAIテクノロジーは革新的で、ビジネスに大きな影響を与えました。
- [2]:XXX社 IoTプロジェクト**
https://dummyurl/IoT_01.xxxx
 - クライアント名: XXX株式会社
 - 業種: 保険業界
 - 背景: クライアントから、AIやIoTを利用した品質検査の過去事例紹介してほしいという要望を頂きました。
 - 提案概要: 当社はIoTテクノロジーを活用し、店頭の顧客の動きをAIによってリアルタイムで分析する新たな取り組みを検討中です。このシス

At the bottom, there is a 'Send a message...' input field and a disclaimer: 'This chatbot may produce inaccurate information about people, places, or facts.'

Annotations in the image:

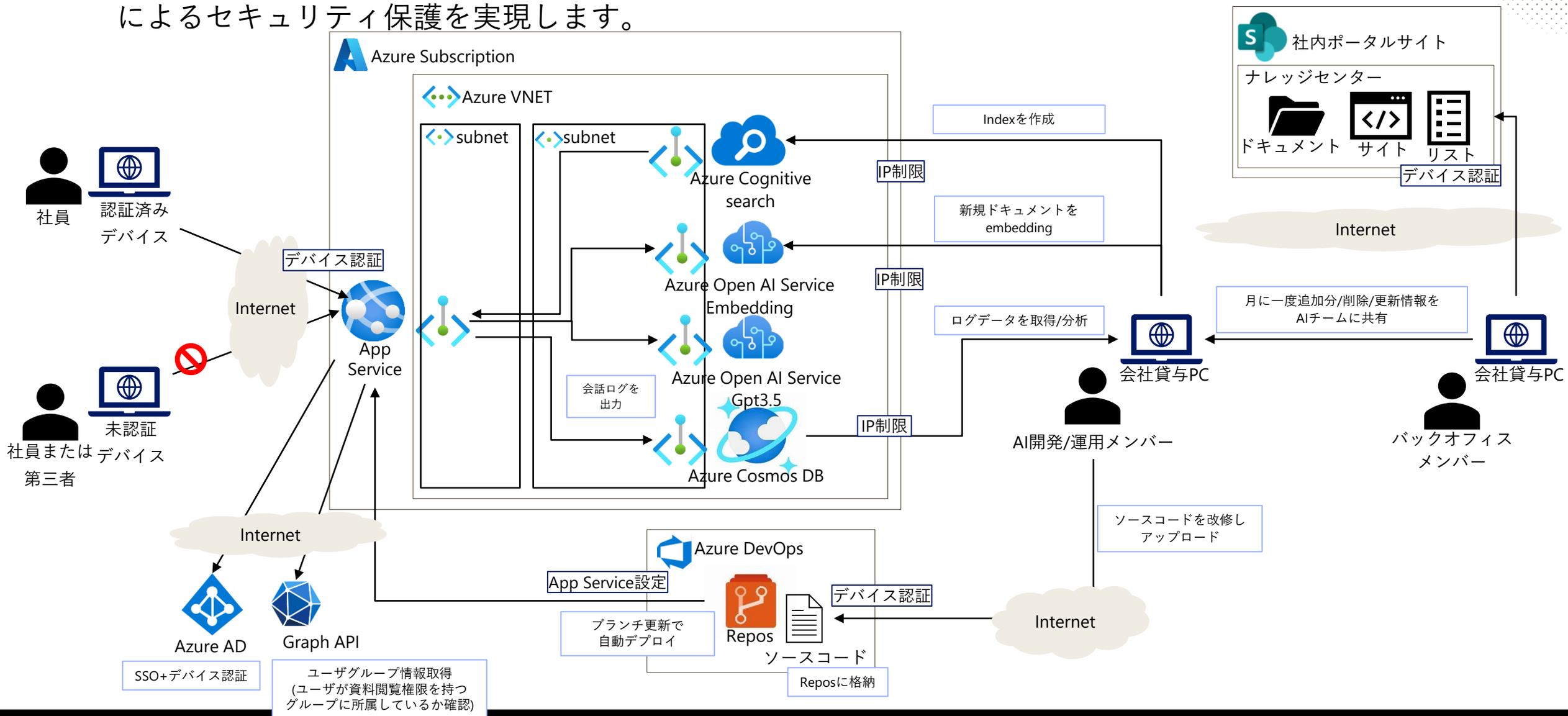
- A dashed yellow box around the user's input is labeled 'ユーザーが自然文で文章を検索' (User searches for text in natural language).
- A dashed yellow box around the first search result is labeled '関連する文章を提示' (Present related text).

3.アーキテクチャ

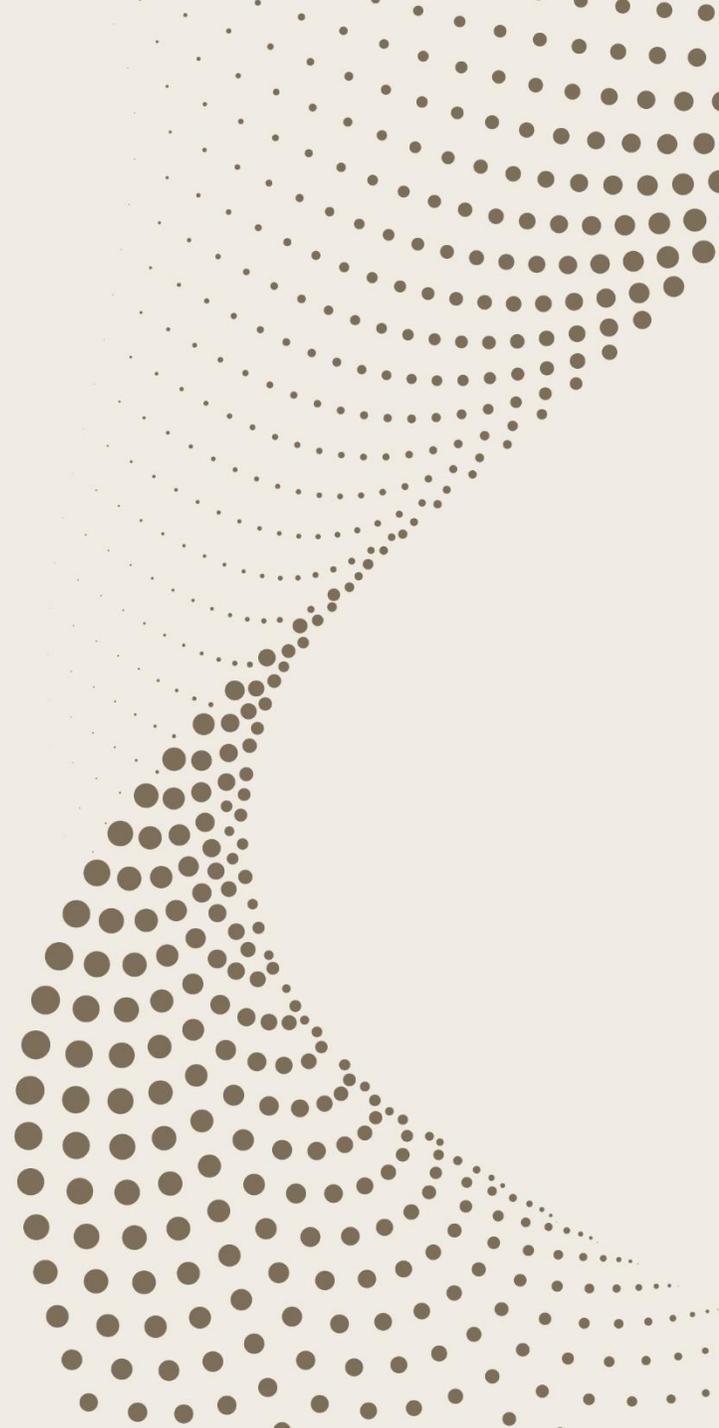


アーキテクチャ図

- デバイス認証でユーザーアクセスを制御し、ナレッジをSharepointサイトに格納することで、権限管理によるセキュリティ保護を実現します。



4. 考慮事項



4. Azure Well-Architected Framework 観点での考慮事項 (1/2)

[Microsoft Azure Well-Architected Framework](#)の観点に則った以下の点を考慮した実装・運用は、非機能要件の充足に有効です。

信頼性	可用性	<ul style="list-style-type: none">• Azureのサービスでは「可用性ゾーン」や「リージョン」といった単位で可用性を設計しており、これらを適切に組み合わせることで、ビジネスクリティカルなワークロードの信頼性を実現するように設計することが可能です。詳細は「Azure リージョンと可用性ゾーンとは」をご参照ください。• このシナリオで用いられているAzure Open AI、Azure Cosmos DB for PostgreSQL 等のコンポーネントはそれぞれゾーン冗長、リージョン冗長、geoレプリケーションなど高可用性のオプションや構成を利用可能です。必要となる可用性に応じて導入を検討してください。複数リージョン間/複数ゾーン間でAct-Act構成を取る場合にはAzure Front Door、Azure Application Gatewayなどの利用をご推奨します。
	回復性	<ul style="list-style-type: none">• Azure Cosmos DB は、定期的にデータのバックアップを自動的に取ります。自動バックアップは、データベース操作のパフォーマンスや可用性に影響を与えずに取得されます。すべてのバックアップは別々のストレージサービス内に個別に保存されます。これらのバックアップは、リージョンの障害からの回復性を確保するためにグローバルにレプリケートされます。詳細については、「Azure Cosmos DB と信頼性」を参照してください。
セキュリティ		<ul style="list-style-type: none">• セキュリティは、重要なデータやシステムの意図的な攻撃や悪用に対する保証を提供します。詳細については、「セキュリティの重要な要素の概要」を参照してください。• このシナリオでは、Azure ADを使用してユーザーを認証します。セキュリティで保護されたソリューションの設計に関する一般的なガイダンスについては、「Azure のセキュリティのドキュメント」を参照してください。• このシナリオにおけるWebAppsとAzure Open AI Serviceの間の通信ではPrivate Linkを使用し、内部からアクセスしています。

4. Azure Well-Architected Framework 観点での考慮事項 (2/2)

コスト最適化

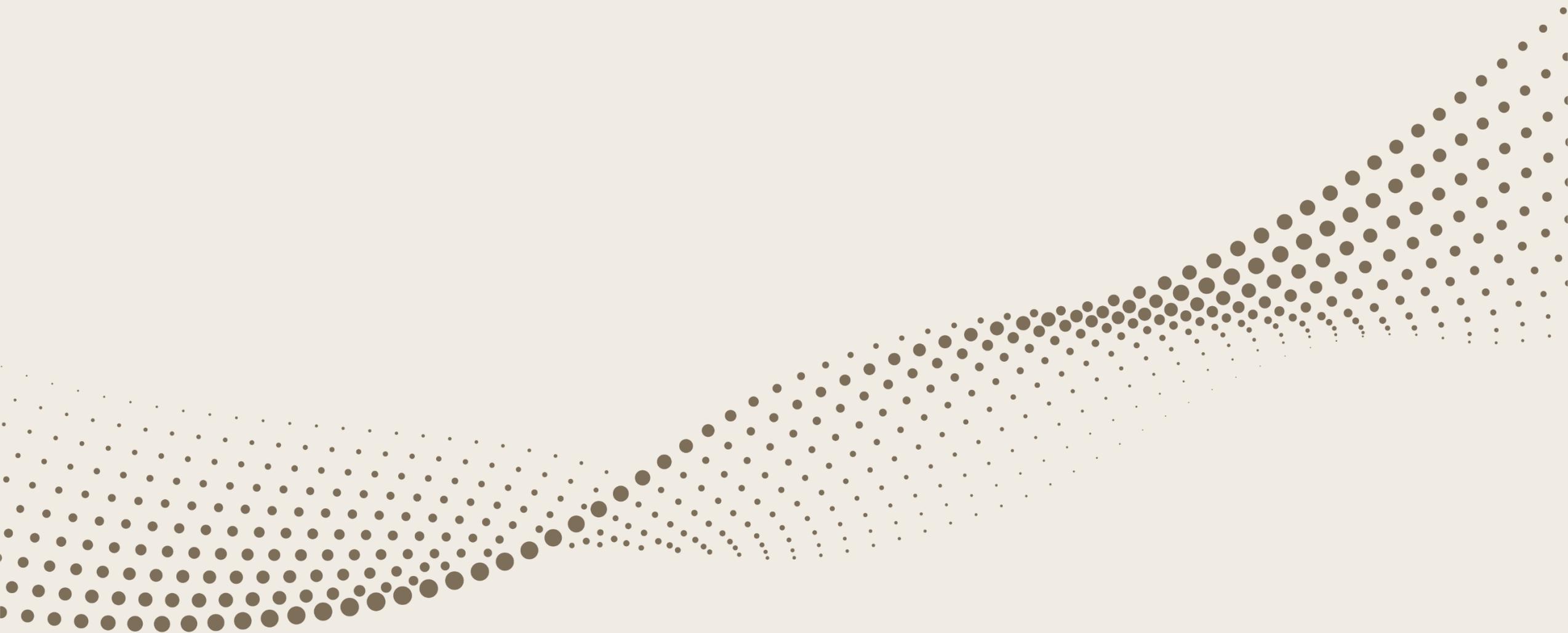
- 不要な費用を削減し、運用効率を向上させる方法を検討することです。詳しくは、[コスト最適化の柱の概要](#)に関する記事をご覧ください。

オペレーション エクセレンス

- システムの健全性の担保、トラブルの解決、利用動向の監視を行うためには適切な監視とログ収集が必要となります。詳細は「[ワークロードの監視](#)」をご参照ください。
- ソフトウェアのアップデートや脆弱性への対応など、ソフトウェア/インフラ設計の改修を円滑に進められるよう、DevOpsプロセスを確立してください。詳細は「[リリース エンジニアリングの継続的インテグレーション](#)」をご参照ください。

パフォーマンス 効率

- アプリケーションの負荷が高まることを見越し、スケーラビリティの確保をあらかじめ検討することは重要です。詳細は「[スケーリング用のアプリケーションを設計する](#)」をご参照ください
- App Serviceは負荷に応じて水平にスケールさせることが可能です。詳細については「[自動スケーリングを有効にする方法](#)」をご参照ください。
- 特定のユーザーにAzure Open AIの利用が集中することを避けたい場合にはAPI Managementによるクォータ導入などをご検討ください。設定の詳細は「[Azure API Management を使用した高度な要求スロットル](#)」をご参照ください。



アビーム、ABeam及びそのロゴは、アビームコンサルティング株式会社の日本その他の国における登録商標です。
本文に記載されている会社名及び製品名は各社の商号、商標又は登録商標です。©2023 ABeam Consulting Ltd.



Build Beyond As One.