

言語処理学会第31回年次大会 併設ワークショップ JLR2025

日本語言語資源の構築と利用性の向上

RAGEvalに基づく日本語の 実践的RAG評価データセットの構築と検証

白邦裕千, 石井愛 (BIPROGY株式会社)
2025/03/14



BIPROGY

Foresight in sight

目的

顧客の個別ドメイン・想定される利用シーンに即したRAGシステムを
事前検証できるようにするための**RAG評価データセット**の作成

本発表の内容

RAG評価データセット**SynRAG**の構築・検証

- 構築：RAGEvalフレームワークを修正・拡張
- 検証：
 - ドメイン知識を反映した文書および質問をLLMとの協働で作成できることがわかった
 - RAG評価指標による検証により通常のRAGでは難易度が高い問題を示した

→ **従来のと比較し、実用面でより**有用性**の高いデータセットである可能性**

1 背景・目的

2 先行研究

3 手法

4 実験結果・考察

5 まとめ

1 背景・目的

2 先行研究

3 手法

4 実験結果・考察

5 まとめ

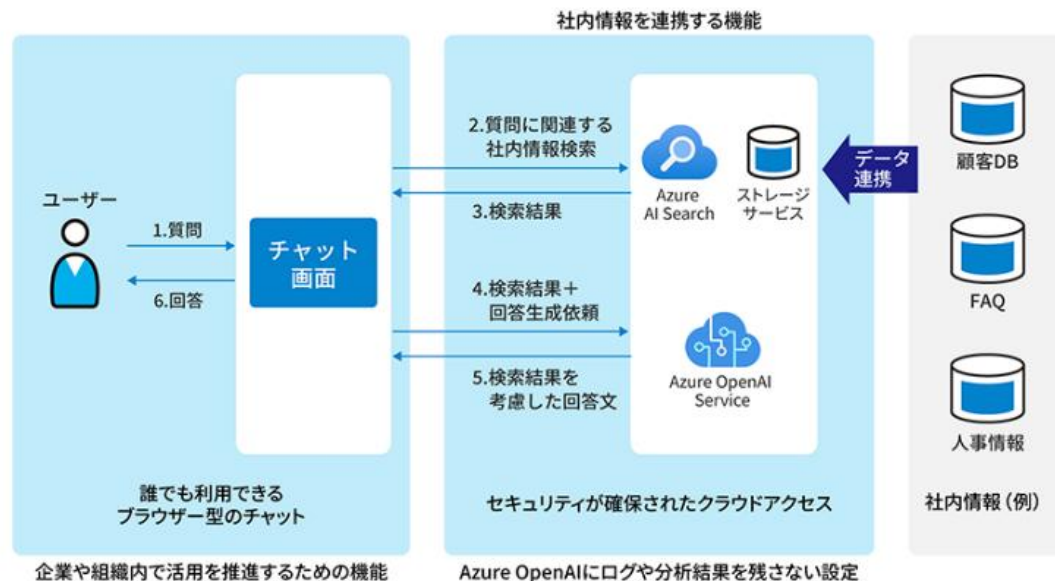
弊社では**Azure Open AI ServiceスターターキットPlus**を提供しており、RAGによる様々なドメインでの社内ナレッジの有効活用をサポートしております。

提供サービス例

- ChatGPT利用環境
- チャットインターフェース
- 活用シナリオ検討

AI適用事例

- ChatGPTによる会話型アプリ
- CoTアプリ
- **RAGアプリ**



顧客の社内へのRAGシステム適用における課題

1. 顧客の社内情報を用いたRAG構築は精度検証のハードルが高い
2. 顧客のドメインに合ったデータセットや、日本語のデータセットがない
3. 既存のデータセットはWikipediaベースのものが多く、LLMが事前学習に用いていない文書群が検索対象の場合に情報を適切に活用できるかを検証できない

本研究の目的

**顧客ごとの個別ドメイン・想定される利用シーンに即したRAGシステムの
事前検証を可能にする**

- **ドメイン**の情報を踏まえた**架空の業務文書を作成**
- それに基づく**利用シーンに即したRAG用データセットを作成**

1 背景・目的

2 先行研究

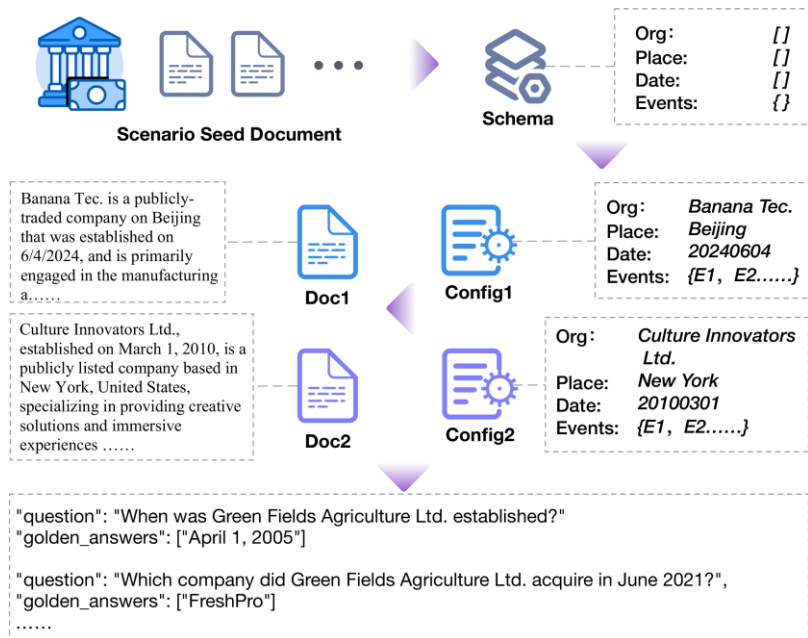
3 手法

4 実験結果・考察

5 まとめ

既存研究 : ドメイン特化・架空なデータセットの作成フレームワークを提案

- 少数のシード文書をもとに、RAGの検証用データを自動生成するプロセスが体系化



RAGEvalのデータセット生成プロセス

1. シード文書をもとにLLMでスキーマを抽出
2. スキーマの値を埋めたConfigを作成
 - 値埋めにはルールベース手法およびLLMを使用
3. Configに沿った文書生成
4. Configと文書を基に質問-参照-回答 (QRA) 生成

RAGEvalデータセットの課題

スキーマ・Configで文書同士の関係性がない
→ 実践的な複雑な質問作成が難しい

用意されたスキーマを日本語化して作成した文書とQRAデータセットの例

文書

■ 文書同士の関係性や整合性が無い

イノベーション研究社	財務報告書_2017	""【財務報告】\nイノベーション研究社は、東京都千代田区に拠点を置き、新素材の研究開発およびライセンス供与を主な事業とする、東京証券取引所第一部上場の企業です。 \n2017年、イノベーション研究社は、事業構造的
イクスプローラーズ株式会社	財務報告書_2018	""【財務報告】\nイクスプローラーズ株式会社は、東京都千代田区に拠点を置き、旅行パッケージの企画・販売や観光ガイド、宿泊施設の運営を行う未上場の観光企業です。 \nE..
エンターテインメントジャパン株式会社	財務報告書_2019	""【財務報告】\nエンターテインメントジャパン株式会社は、映画製作、配信サービス、イベント運営を手掛ける、東京都渋谷区に拠点を置く上場エンターテインメント企業です。 \nエンターテインメントジャパン株式会社は、2019年において敦..
オリエント工業株式会社	財務報告書_2019	""【財務報告】\nオリエント工業株式会社は、東京都品川区に拠点を置き、精密機械および工業用製品の製造を主な事業とする、東京証券取引所第一部上場の製造業企業です。 \n2019..
クリーンワークス株式会社	財務報告書_2019	""【財務報告】\nクリーンワークス株式会社は、東京都港区に拠点を置き、住宅および商業施設向けの清掃・管理業務を提供する未上場の清掃・メンテナンス企業です。 \nク

- 業種間の比較が出来ない
- 同一企業の時系列による比較が出来ない

QRAデータセット

■ 不自然な質問

西区中央病院の入院記録によると、河野学の出生地はどこですか？
逗子市立病院の入院記録によると、林美羽の母の健康状態は何ですか？
平塚市総合病院と豊島総合病院の入院記録によると、吉田詩織と石井大樹の手術内容は何ですか？
北区総合病院と千葉県立病院の入院記録によると、菊地守と中野心春のどちらが先に入院しましたか？
2020年にメディアインク株式会社が改定したポリシーは何ですか？
2021年に未来学園株式会社が締結したライセンス契約の相手はどのような機関ですか？

- 表面的な質問
- 実用上起こりえないであろう質問

■ 不自然なデータ構造

成田市立総合病院の入院記録により、遠藤聡の治療計画を要約してください。
成田市立総合病院の入院記録によると、遠藤聡の2型糖尿病の診断根拠はいつありますか？
:
成田市立総合病院の入院記録により、遠藤聡の入院後の経過

- 1病院1患者（1業種1企業）

1 背景・目的

2 先行研究

3 手法

4 実験結果・考察

5 まとめ

実際に起こり得る質問を生成するため、既存の生成プロセスを改善

①想定質問起点の文書設計, ②構造化データ生成ロジックの精緻化 でデータセットを高度化

RAGEvalを改善したデータセット生成プロセス

1. 少数のシード文書をもとにLLMでスキーマを作成
 - 1-1: LLM・専門家との反復的対話プロセスに基づき、想定質問リストを作成
 - 1-2: 想定質問から、文書の依存関係を定義・Configを更新
2. ルールベース手法およびLLMでスキーマの値を埋めたConfigを作成
 - 2-1: ドメイン固有の制約と文書同士の依存関係を考慮したルールベースのプログラム構築
3. Configに沿った文書生成
 - 3-1: LLMによる文書レビュー、修正点の反映（必要に応じて2-1～3-1を反復）
4. Configと文書を基に質問-参照-回答（QRA）生成

新たに追加したプロセス

①

②

① 想定質問起点の文書設計

1-1: LLM・専門家との反復的対話プロセスに基づき、想定質問リストを作成

架空の入院記録を参照しながら答える質問応答データセットを作成します。
入院記録を業務で用いる人が実際に調べたいような問題を作成します。

その入院記録を参照しなければ答えられない問題、もしくは、その入院記録と一般的な知識を掛け合わせて答える問題の例を挙げてください。

入院記録は、添付のような形式で同じ患者の単一の記録が存在する場合と、複数の記録が存在する場合があります。

...

問題をカテゴリごとに分類して質問と回答生成の例を示してください。

入院記録を用いた質問応答データセットのカテゴリと例を示します。

1. 単一記録からの事実確認

質問: 患者P502596の入院時の主訴は何でしたか？また、発症からどのくらいの期間で入院に至りましたか？

回答: 入院時の主訴は頭痛、下痢、筋肉痛でした。4月28日に症状が発現し、5月4日に入院したため、発症から6日後に入院に至りました。

2. 複数記録の比較分析

質問: 患者の2回の入院における肝機能検査（AST/ALT）の値を比較し、どのように変化したか説明してください。

回答: 1回目の入院時（2024年1月）はAST 19U/L、ALT 27U/L と正常範囲内でしたが、2回目の入院時（2024年5月）...

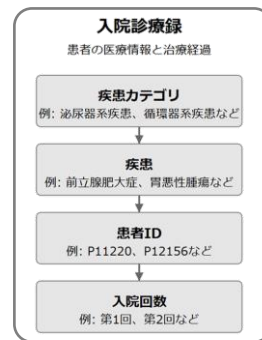
② 構造化データ生成ロジックの精緻化

2-1: ドメイン固有の制約と文書同士の依存関係を考慮したルールベースのプログラム構築

- ドメイン固有の制約（疾患ごとの制約条件）

カテゴリ	性別	年齢範囲	疾患数
女性特有疾患 (female_specific)	女性	20~50歳	13種 <small>乳房悪性腫瘍、子宮頸部悪性腫瘍、妊娠管理など</small>
男性特有疾患 (male_specific)	男性	40~90歳	1種: 前立腺肥大症または炎症
高齢者疾患 (elderly)	制約なし	65~95歳	3種 <small>老年期精神病、パーキンソン病、白内障</small>
新生児疾患 (neonatal)	制約なし	0歳	4種 <small>早産児および未熟児、分娩損傷など</small>

- 文書同士の依存関係



Configの例：財務ドメイン（統合報告書）

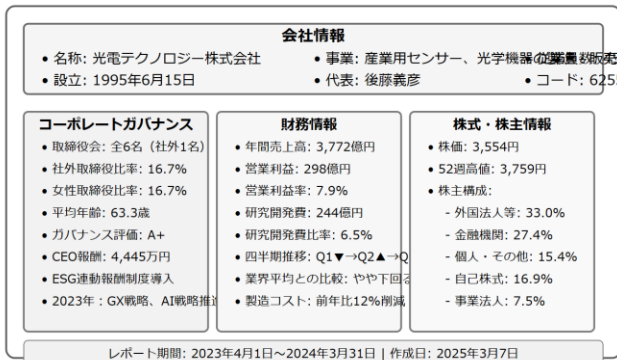


データセット作成

Claude 3.7 SonnetでConfigを可視化

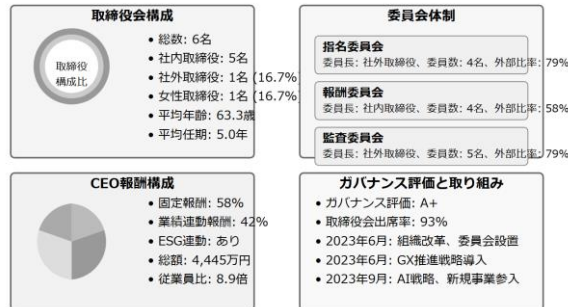
■ 全体

光電テクノロジー株式会社 概要

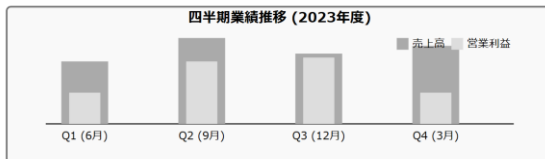


■ 内容詳細

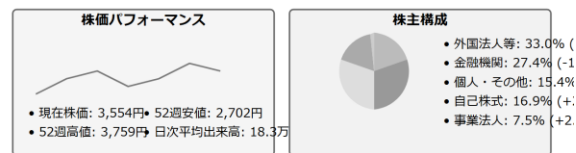
光電テクノロジー株式会社 ガバナンス構造



光電テクノロジー株式会社 財務データ



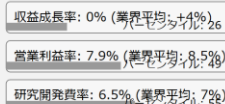
光電テクノロジー株式会社 株主構成



業績ハイライト

- 売上高: 3,772億円 (前年同期比 横ばい)
- 営業利益: 298億円
- 営業利益率: 7.9%
- 研究開発費: 244億円 (6.5%)
- Q2における需要回復 (前期比14.3%増)
- Q4に新製品向け研究開発投資増強

業界比較



大株主情報

株主名 持株比率 前年比

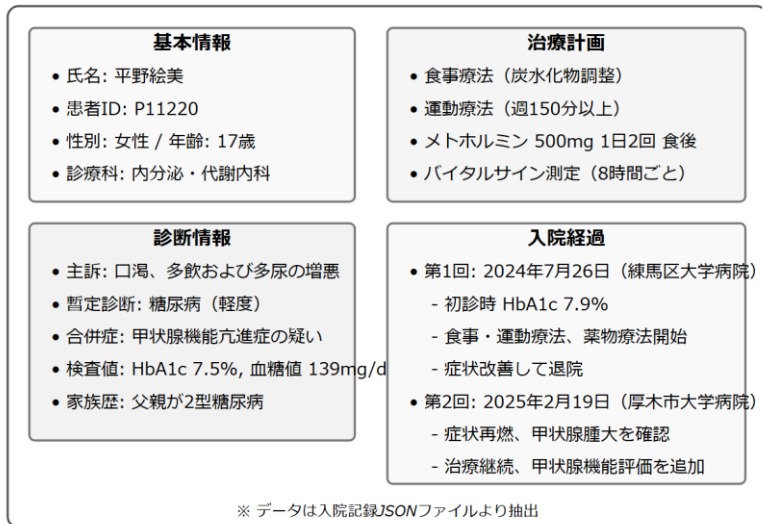
株主4 (外国法人等)	20.4%	+1.2%
株主5 (金融機関)	18.6%	+1.7%
株主3 (金融機関)	11.3%	-0.7%
株主1 (外国法人等)	8.1%	+0.5%
株主2 (外国法人等)	7.6%	+1.3%

株主エンゲージメント

- 投資家ミーティング:
 - 機関投資家: 69回
 - 個人投資家: 4回
- 証券アナリスト: 27

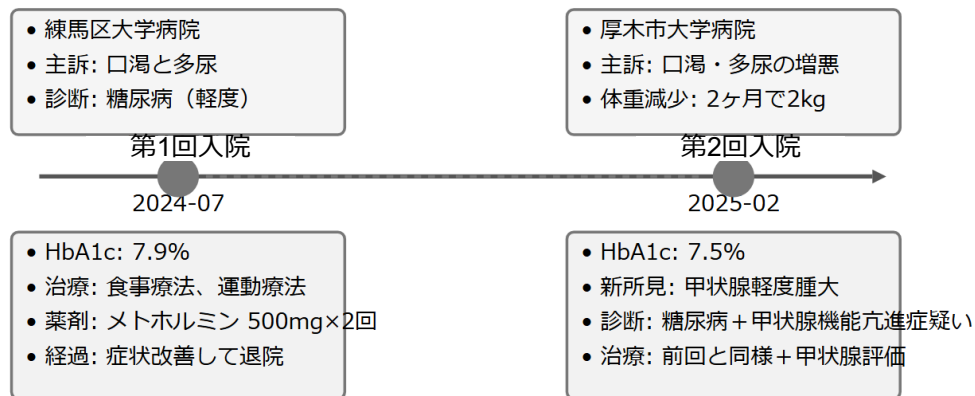
■ 全体

入院記録データ構造



■ 時系列複数データ（複数回入院）

平野絵美さん(P11220)の入院記録時系列



Claude 3.7 SonnetでConfigを可視化

データセットの質・検索/回答生成性能の評価・検証を実施

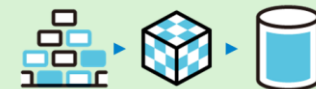
■ データセット生成

- ドメイン/業務文書：
 1. 財務ドメイン/統合報告書
 2. 医療ドメイン/入院診療録
- LLM：

モデル	用途
Claude 3.5, 3.7	<ul style="list-style-type: none">• 想定質問リスト作成• スキーマ更新• 構造データ生成プログラム生成
ChatGPT Deep Research / o3-mini-high	生成された文書のレビュー
Azure OpenAI Service API GPT-4o	<ul style="list-style-type: none">• Config生成• 文書生成• QRA生成

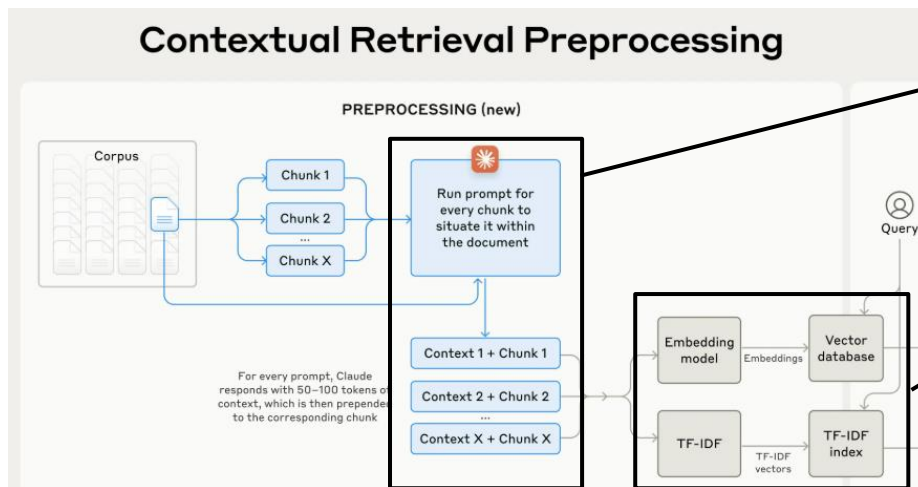
■ 評価用RAGシステム

- 検索処理：
 - ◆ 密ベクトル検索：
 - 検索エンジン：Numpyによるコサイン類似度検索
 - 埋め込みモデル：OpenAI API text-embedding-3-small
 - ◆ 疎ベクトル検索：
 - 検索エンジン：Elastic Search
 - 検索アルゴリズム：BM25
 - 日本語トークナイザ：analysis-kuromoji
- 回答生成
 - ◆ Azure OpenAI Service API GPT-4o



データセットの評価のため、以下の2種類のRAGシステムを実装

1. 一般的な構成のRAG：密ベクトル検索+LLM
2. Contextual Retrieval：



- ドキュメントとチャンクをLLMに渡して、チャンクの文脈（Context）を生成する
- 作成した文脈をチャンクの前頭に付与する

- 密ベクトル検索を実行する
- 疎ベクトル検索を実行する
- 上記検索結果をマージする

参考：[Introducing Contextual Retrieval](#)

3. RAG評価指標



データセットの評価には先行研究と同じ評価指標を使用

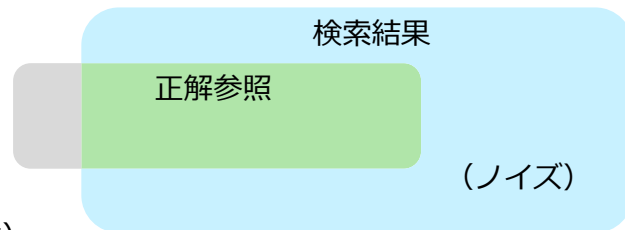
検索

■ Recall

- $(\text{検索結果に占める正解参照の数}) / (\text{正解参照の数})$
→ 検索の正確性の指標

■ EIR (Effective Information Rate)

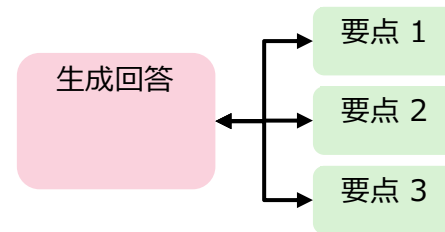
- $(\text{検索結果に占める正解参照と一致する単語数}) / (\text{検索結果の単語数})$
→ 検索の効率性の指標



回答生成

生成回答と正解要点を比較し、生成した回答が以下のいずれかに該当するかを評価する※ (keypoint metrics)

- Completeness (完全性) : 正解要点を含んでいる
- Hallucination (幻覚) : 正解要点と矛盾しているか
- Irrelevance (無関係性) : 完全性、幻覚どちらにも該当しない場合
→ 生成回答の質・信頼性の指標



※ 評価にはLLM（本研究ではGPT-4o）を使用

1 背景・目的

2 先行研究

3 手法

4 実験結果・考察

5 まとめ

■ 財務ドメイン（統合報告書）

◆ 構造的特徴

- 文章構造：叙述形式
- 平均文書長：12,000字程度

◆ 情報的特徴

- 時系列の数値データ比較
- 表データ（Markdown）を含む
- 企業の多様性の表現
 - ◆ 業種、規模、成長段階
- 財務指標間の数学的整合性
- ガバナンス構造と業績の関連性

■ 医療ドメイン（入院診療録）

◆ 構造的特徴

- 文章構造：箇条書き形式
- 平均文書長：2,000字程度

◆ 情報的特徴

- 時系列の数値データ比較
- コーパス内の分布の計算
- 疾患進行の時系列的一貫性
- 医学的整合性
 - ◆ 臨床検査値・バイタルサイン
- 治療効果と症状変化の因果関係

※ RAGEvalの生成プロセス改善前後を明確にするため、改善後のデータセットを“SynRAG”と表記



■ 財務ドメイン:

大和化成株式会社 2023年度 統合報告書

1. 企業概要

企業概要

トップメッセージ

皆さま、こんにちは。大和化成株式会社代表取締役兼社長は1972年の創業以来、高機能樹脂や化学製品を中心に、環境対応型新素材の開発を通じて、社会に新たな価値を提供し、持続可能な社会の実現に貢献するという私たちの使命を、次世代型バイオプラスチックの商業化に成り果てています。

私たちは「イノベーションを通じて、社会に新たな価値を提供し、持続可能な社会の実現に貢献する」という使命を、次世代型バイオプラスチックの商業化に成り果てています。

企業プロフィール

社名: 大和化成株式会社
設立: 1972年4月1日
本社所在地: 東京都中央区日本橋2丁目5番1号
業種: 化学 (製造業)
上場コード: 7237
従業員数: 約4,500人

収益構造の変化と要因

2023年度、収益の伸びは成長は、中国市場の需要減速やエネルギーコストの上昇といった外部要因によるものでした。一方で、製品ポートフォリオの改善や輸出拡大戦略が収益を一定水準に保つ要因となっています。また、「エコポリマーX」の市場投入により、環境対応型製品の売上比率が増加しつつあります。

3. 研究開発活動

R&D支出の推移と戦略的意義

2023年度の研究開発費は合計 ¥212,460百万円 (年間売上上の6.4%) に達し、業界平均を上回りました。特に、FY2023-Q3では高性能樹脂の次世代製品に注力した結果、研究開発費が大幅に増加しています (Q2比 55.4%増)。

四半期	R&D支出 (¥)	R&D比率 (%)
FY2023-Q1	¥53,076百万円	6.4%
FY2023-Q2	¥41,962百万円	5.1%
FY2023-Q3	¥65,231百万円	7.5%
FY2023-Q4	¥52,189百万円	6.4%

4. 業界比較と市場ポジション

業種: 31種類

業種ごとに以下の条件を設定

- イベント、リスク要因
- 財務指標 (成長率など) のレンジ

企業の規模感 (4段階) × 成長ステージ (4段階)

■ 医療ドメイン:

【入院記録】

患者基本情報:

- 氏名: 中野誠
- 性別: 男性
- 年齢: 57歳
- 病院名: 市川市大
- 患者ID: P69314
- 生年月日: 1963年
- 診療科: 総合診療
- 主治医: 三浦通 専
- 担当看護師: 鈴木

入院管理情報:

- 患者ID: P69314
- 入院回数: 初回
- 過去の入院歴: なし

現病歴:

患者は2021年5月14日早朝起床時、突然の呼吸困難をきたす。発症当日は強い労作後に症状をきたす程度で、SpO2は95%程度で特記すべき異常は確認されず、徐々に症状が軽くなり、尿・排便は正常で

• 血液検査:

- CRP: 2.4 mg/dL (軽度上昇)
- WBC: 9,326 / μ L (軽度上昇)
- AST: 28 U/L
- ALT: 17 U/L
- 血小板数: $215 \times 10^5 / \mu\text{L}$
- BUN: 14.3 mg/dL
- クレアチニン: 0.9 mg/dL

診断:

気管支喘息 (軽症、急性増悪なし)。

診断根拠:

呼吸困難および喘鳴の症状、聴診での喘鳴所見、SpO₂の軽度低下 (95%)、軽度の炎症反応 (CRP 2.4 mg/dL) の上昇を総合して診断。

治療計画:

- 酸素療法: SpO₂をモニタリングしながら継続。
- 吸入療法: 必要に応じて吸入 β 刺激薬の投与を検討する。

疾患: 55種類

疾患ごとに性別や年代などの制約を定義

再発度合 (3段階) によって再入院データを作成



Case：化学メーカーの統合報告書

■ 情報密度

- 企業概要, 事業概要, 経営戦略
- 財務分析レポート
- コーポレートガバナンス
- 将来の見通し
 - ESG戦略, 投資戦略

アークケミカル株式会社 2022年度 統合報告書

主要財務指標の概要

2022年度のアークケミカル株式会社の財務パフォーマンスは、以下の通りです：

- 売上高：4,294億6,050万円（前年比横ばい）
- 営業利益：515億6,447万円（前年比+6.8%）
- 営業利益率：12.9%（前年比+1.1ポイント）
- 研究開発費：340億6,502万円（売上比率7.7%）

売上高は前年並みとなったものの、営業利益は製造プロセスの効率化と高付加価値製品の販売拡大により堅調に推移しました。また、研究開発への積極的な投資により、環境対応型製品の市場投入が成功し、収益性向上に寄与しました。

四半期ごとの業績推移と変動要因

Q1（2022年4月～6月）

- 売上高：4,294億6,050万円
- 営業利益：489億9,504万円（営業利益率11.8%）
- 国内市場の需要減少と主要取引先の在庫調整の影響で売上は微減しましたが、コスト管理効率化により営業利益率は堅調に推移しました。

Q2（2022年7月～9月）

設備投資・R&D投資計画

- 設備投資：
 - 新工場の生産能力を2024年度中にフル稼働させ、年間生産能力を20%向上。
 - 既存プラントの改修投資を進め、エネルギー効率を20%改善。
- R&D投資：
 - 次世代素材（軽高強度樹脂、バイオプラスチック）の研究開発に集中投資。
 - AIを活用した材料開発プロセスの効率化を推進。

M&A・アライアンス戦略

- ターゲット市場：東南アジアおよび中東地域を中心に、現地の特殊化学品メーカーとのM&Aや戦略的提携を検討。
- 技術提携：環境対応型製品の技術開発において、大学や研究機関との共同プロジェクトを拡大。

新規事業・新市場開拓

- バイオプラスチックを中心とした新規事業を育成し、5年以内に年間売上500億円規模に成長させる。
- エネルギー産業（特に再生可能エネルギー関連）向けの新規素材市場に参入。

1. 企業概要

アークケミカル株式会社 企業概要

トップメッセージ

アークケミカル株式会社の代表を務めます坂口秀です。当社は1991年の創業以来、化学品の製造・販売を中心に、持続可能な未来を築くべく成長を続けてまいりました。材の研究開発をさらに一歩進め、新たな環境対応型化学製品の市場投入に成功しました。これらの成功は、社員一人ひとりの努力と、株主・投資家の皆様からのご謝しております。

現在、化学産業は大きな変革期を迎えています。気候変動問題や環境規制の強化が可能なイノベーション」という信念を掲げ、化学産業の新たな未来を切り開くリーダーとして、投資家の皆様におかれましては、これからも当社の挑戦を温かく見ただけますようお願い申し上げます。

企業プロフィール

社名：アークケミカル株式会社

設立：1991年6月15日

本社所在地：東京都中央区日本橋3丁目8番10号

上場市場：東京証券取引所プライム市場（コード：7523）

アークケミカル株式会社の企業理念は「化学の力で未来を創る」です。この理念のしい製品の開発を通じて社会課題の解決に貢献し、持続可能な社会の実現を目指し、革新と信頼で業界をリードし、次世代の可能性を切り拓くこと」です。そして「化学産業を通じて人々の生活を豊かにすること」を掲げています。

アークケミカルは、設立当初より高機能樹脂分野を中心に事業を展開し、2000年代成功。2010年代には環境対応型製品を戦略的に拡充し、現在では化学業界全体を牽引した。30年以上の歴史を持つ当社は、革新と安定を両立させた経営基盤を築いてい

事業概要

アークケミカル株式会社は、次の3つの柱を中心に事業を展開しています。

1. 高機能樹脂および特殊化学品の製造・販売：自動車・電子機器・建設分野などに高品質な素材を提供。耐久性、軽量化、環境負荷低減を実現する製品群が
2. 次世代素材の研究開発：社内外の研究機関と連携し、再生可能エネルギー分野を推進。特に、カーボンニュートラルを実現する材料技術が注目されています。
3. 環境対応型化学製品の提供：グリーンケミストリーを基盤に、従来の化学製



Case：新生児の破傷風発症

■ 時系列一貫性

- 再感染による再発
- 前回入院記録を保持

■ 医学的整合性

- 疾患と症候との対応

破傷風（Tetanus）とは



■ 症状

感染して3日から3週間からの症状のない期間があった後、口を開けにくい、首筋が張る、体が痛いなどの症状があらわれます。その後、体のしびれや痛みが体全体に広がり、全身を弓なりに反らせる姿勢や呼吸困難が現れたのちに死亡します。

参考：厚生労働省検査所 FORTH：破傷風

基本情報

- 患者名: 藤田里奈
- 性別: 女性
- 年齢: 生後12日
- 記録者: 菊地彩香 研修医
- 入院施設: 八千代市大学病院
- 患者ID: P79541
- 生年月日: 2023年10月22日
- 診療科: 新生児科
- 担当看護師: 伊藤悠人 看護師
- 主治医資格: 周産期専門医

入院管理情報

- 患者ID: P79541
- 入院回数: 初回入院
- 前回入院日: 記録なし
- 入院日時: 2023年11月3日 15:00

主訴

哺乳力低下、開口障害（軽度）、軽度の筋強直

診断

- 暫定診断: 新生児破傷風

診断根拠:

顔部発赤所見（発赤、浸出液）、痙攣、筋強直、髄液検査結果に基づく。

鑑別診断:

- 新生児敗血症: 血液検査で白血球数が上昇し、外可能。
- 髄膜炎: 髄膜刺激徴候を現時点で認めず除外。

基本情報

- 患者氏名: 藤田里奈
- 性別: 女性
- 年齢: 8か月
- 記録者: 松本五郎 医局長
- 病院名: 足立区大学病院
- 患者ID: P79541
- 生年月日: 2023年10月22日
- 診療科: 新生児科
- 主治医: 松本五郎 医局長 (周産期専門医)
- 担当看護師: 池田竜也 看護師
- 入院回数: 2回目

入院管理情報

- 前回の入院日: 2023年11月3日 15:00
- 今回の入院日: 2024年6月19日 11:45

既往歴

- 発育・発達: 年齢相応
- 疾患歴: 前回、新生児破傷風の診断で2023年11月3日に入院（軽症で改善）
- 感染症歴: 特記事項なし
- 予防接種歴: 破傷風トキソイド未接種
- 手術・外傷歴: 出生時、臍帯切断後の特別な処置なし
- 輸血歴: 特記事項なし
- アレルギー歴: 薬物アレルギーなし

診断

- 新生児破傷風（再発）

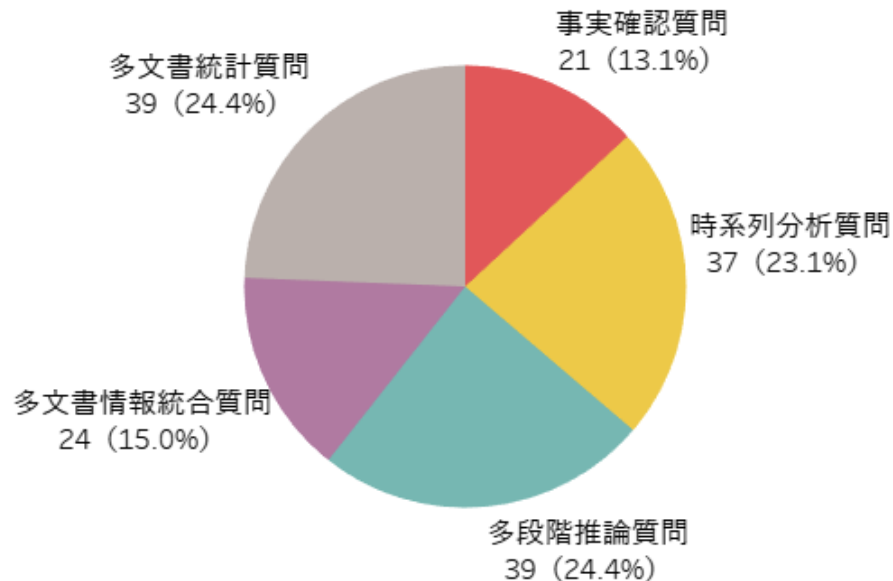
診断根拠

- 顔部の発赤と浸出液、Clostridium tetani陽性の培養結果
- 筋強直、哺乳力低下の臨床症状
- 前回の疾患歴と一致した所見

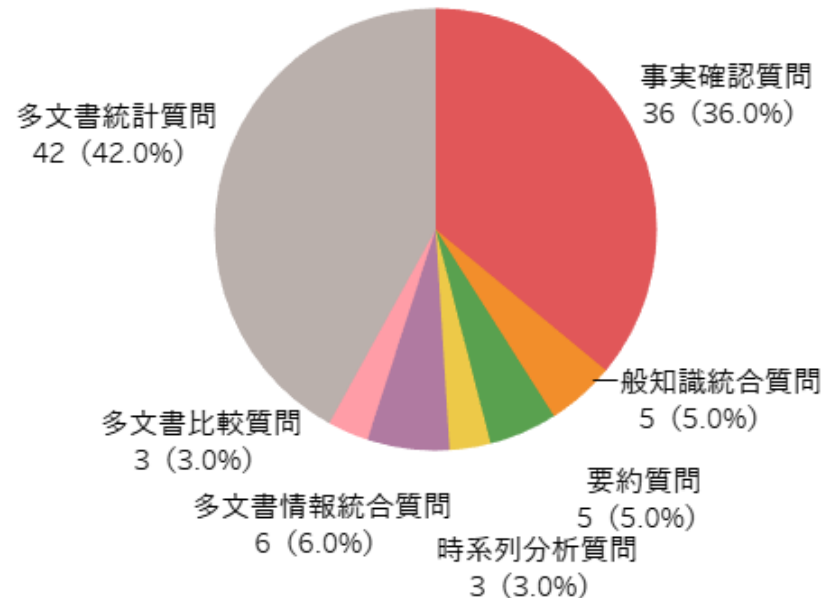
SynRAGデータセット（QRAデータセット）：質問タイプの分布



■ 財務ドメイン:



■ 医療ドメイン:



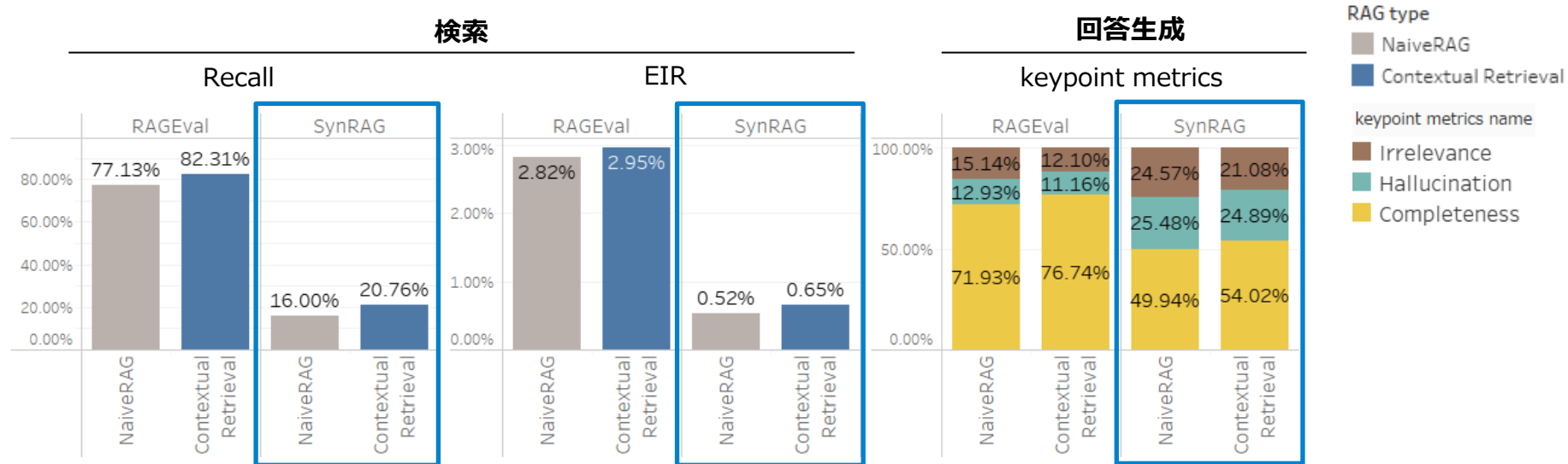


予め想定質問を作成しておくことで、実践的かつ網羅的な質問を作成

事実確認質問	日本陸運ネットワーク株式会社の設立年は何年ですか？
一般知識統合質問	患者ID P73022の痙攣発作及び開口障害の症状は、 破傷風の診断基準 と比較してどのように評価されますか？
時系列分析質問	<ul style="list-style-type: none">クローバー・デジタル株式会社の2020年から2023年における研究開発費（R&D費用）の変化はどのような傾向を示し、それが営業利益率や新製品開発にどのような影響を与えたのかを分析してください。平野健二の重要な検査値の推移パターンを教えてください。
多段階推論質問	エコトランスポート株式会社の取締役会の構成とガバナンス体制の改善が、 業績および株主還元 にどのように 影響した と考えられますか？
多文書比較質問	柴田佳子の各入院時の治療反応を 比較 するとどのような違いがあるか？
多文書統計質問	住血吸虫症患者群 における肝脾腫の 出現頻度 はどのくらいですか？



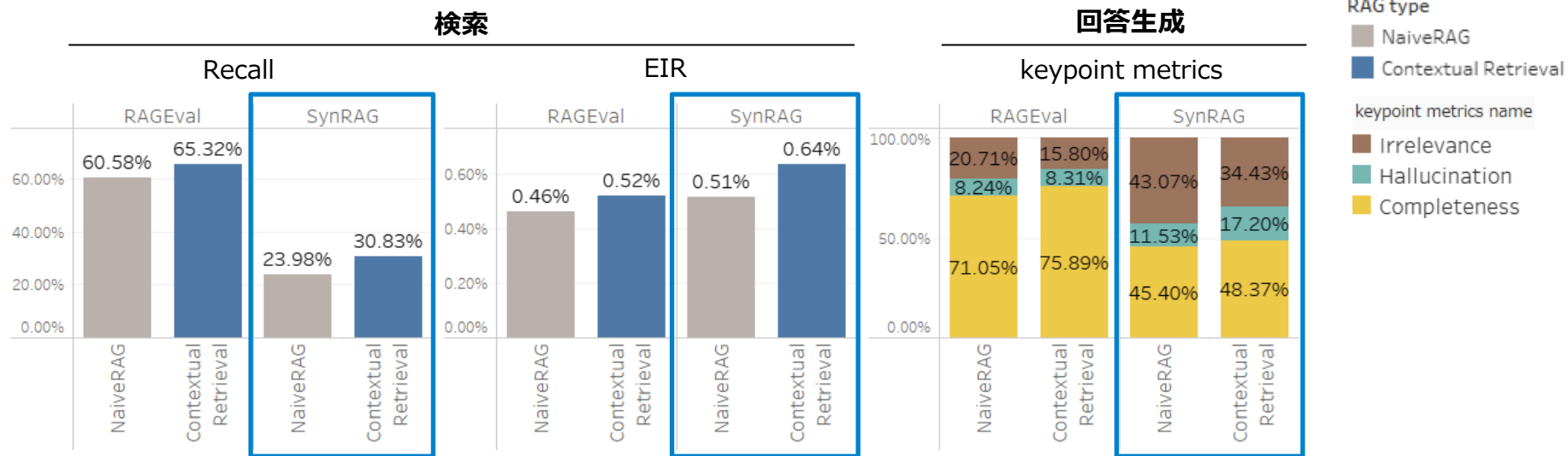
SynRAGでは文書長が大きく、検索難易度が高いことが結果に表れたと考えられる



※ チャンクサイズは1,000、チャンクオーバーラップは0で固定。
 チャンク数はそれぞれRAGEvalデータセット: 289, SynRAGデータセット: 1,664
 検索上位5件を取得



回答のIrrelevanceの割合が高い要因として、関連性が無く、内容が似たドキュメントが誤って検索されていることが考えられる

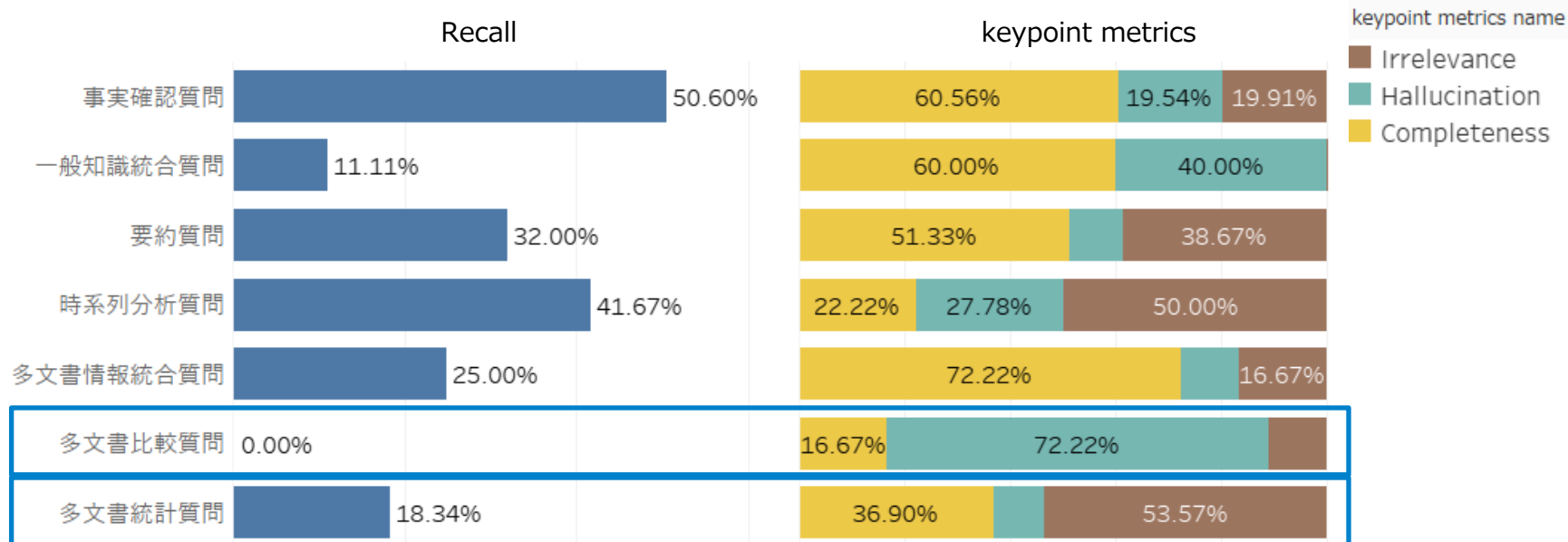


※ チャンクサイズは1,000、チャンクオーバーラップは0で固定。
 チャンク数はそれぞれRAGEvalデータセット: 289, SynRAGデータセット: 1,664
 検索上位5件を取得



Cotextual Retrievalによる各質問タイプごとの精度指標は以下の通り

- 精度が特に低い質問タイプを調査し、質問の難易度を評価





専門知識の理解及び意味的解釈に基づいて、マルチホップな検索が要求される

➤ 高い検索精度が求められる質問

質問	高橋香織さん（患者ID: P69511）の各入院時の治療反応を比較してください。
参照文書数	3
正解参照	1. 適切な抗菌薬治療により7～10日以内の改善を期待。再発リスクは中程度と評価。 2. 14日程度での症状改善が期待される。 3. 患者はICUで管理され、酸素療法および抗菌薬治療を継続予定。
検索結果	対象の計3回の入院記録のチャンクは取得出来ているものの、正解参照となるチャンクを全て取得することは出来ていない。 → 治療反応の記述部分の検索には専門知識（「予後」）の理解が必要



生成モデルの制御能力、統合・集計性能が要求される

➤ RAGシステムの適切な判断能力が問われる質問

質問	尿路結石症の患者群における男性と女性の性別分布はどのようになっていますか？
参照文書数	7
正解参照	1. '性別: 男性', 2. '性別: 男性', … , 7. '性別: 男性'
正解要点	1. 尿路結石症の患者群では男性が6名。 2. 尿路結石症の患者群では女性が1名。
検索結果	尿路結石症の患者に関する検索はできているものの、対象者全てを検索できていない。
回答生成結果	<p>「男性は女性に比べて尿路結石を発症するリスクが高いとされています。 具体的な性別分布は地域や研究によって異なることがありますが、一般的には…」</p> <p>→ ・ 検索結果ではなく、<u>事前学習された一般知識に基づく回答がされる</u> ・ 類似問題においては<u>文脈から回答根拠を得ようとする</u></p>

1 背景・目的

2 先行研究

3 手法

4 実験結果・考察

5 まとめ

ドメインに特化・利用シーンに即したRAG評価データセットを構築・検証

- RAGEvalフレームワークを改善し、新たに**SynRAGデータセット**を作成
 - 医療ドメイン, 財務ドメイン
- データセットの質・検索／回答生成性能の評価・検証を実施
 - 文書：多様性・情報密度・整合性が改善されたことを確認
 - QRAデータセット：RAGシステムの性能を評価するのに十分な質問の種類・難易度

→ **特定のドメイン・利用シーンに即しただけでなく、**

RAGシステム性能評価用としても実用的なデータセット

今後の展望

- データセット・フレームワーク修正版の公開
- 他ドメイン（P&P等の社内文書, 製造, 小売,...）への適用
- データセットの追加評価：より高度なRAGシステム, 一般的なRAG評価指標



BIPROGY

| Foresight in sight

