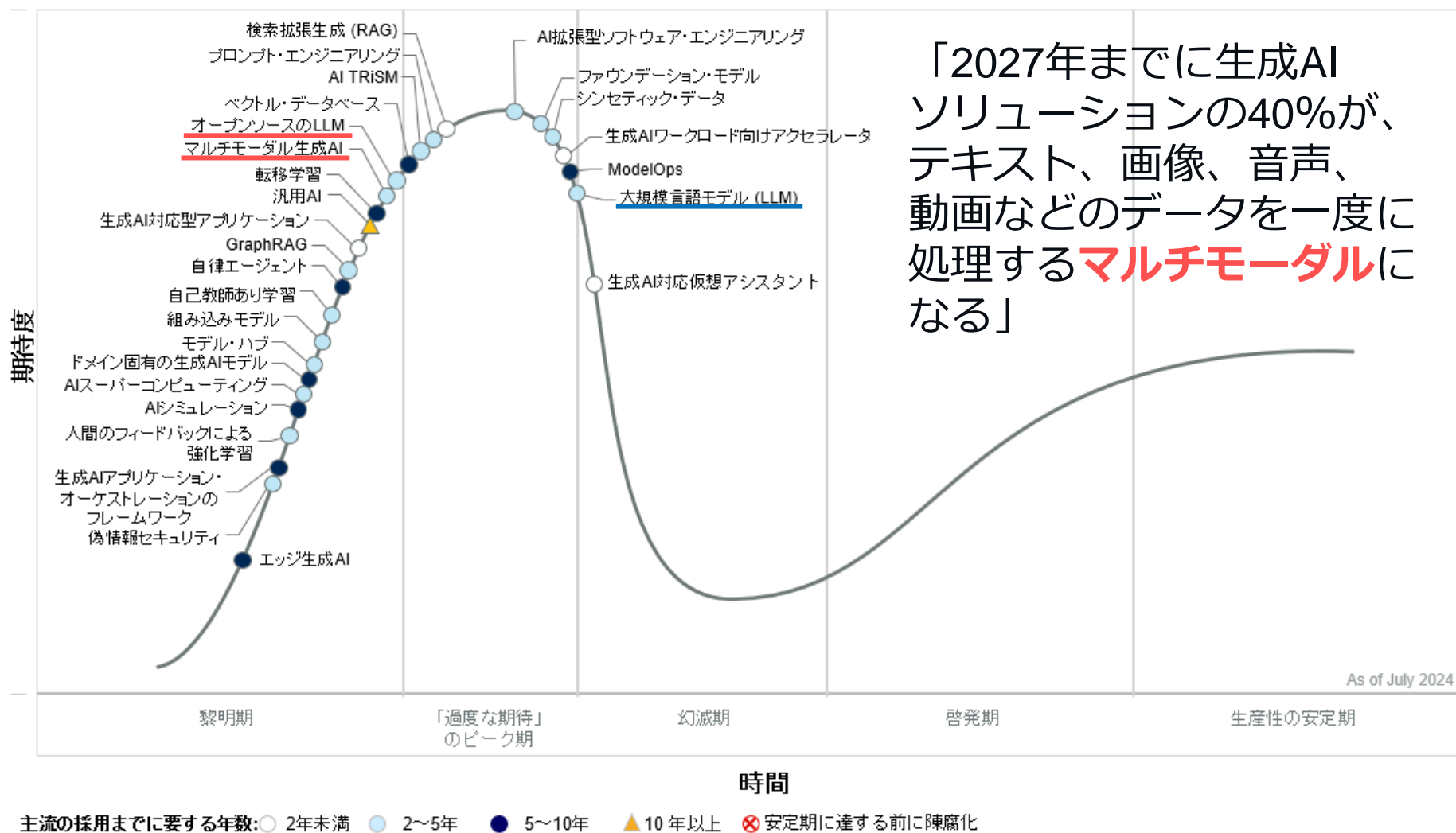


# マルチモーダル基盤モデル

岡崎 直観

東京工業大学 情報理工学院情報工学系  
国立情報学研究所 大規模言語モデル研究開発センター

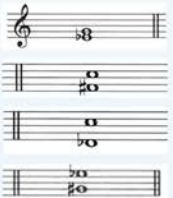
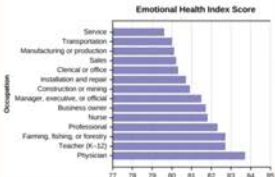
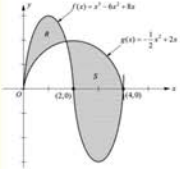
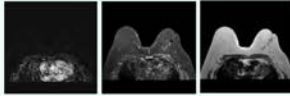

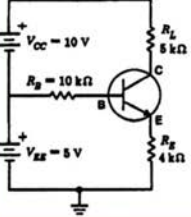
# 大規模言語モデルは幻滅期に突入？



2024年9月10日, Gartner, 「生成AIのハイブ・サイクル: 2024年」を公表. <https://www.gartner.co.jp/ja/newsroom/press-releases/pr-20240910-genai-hc>

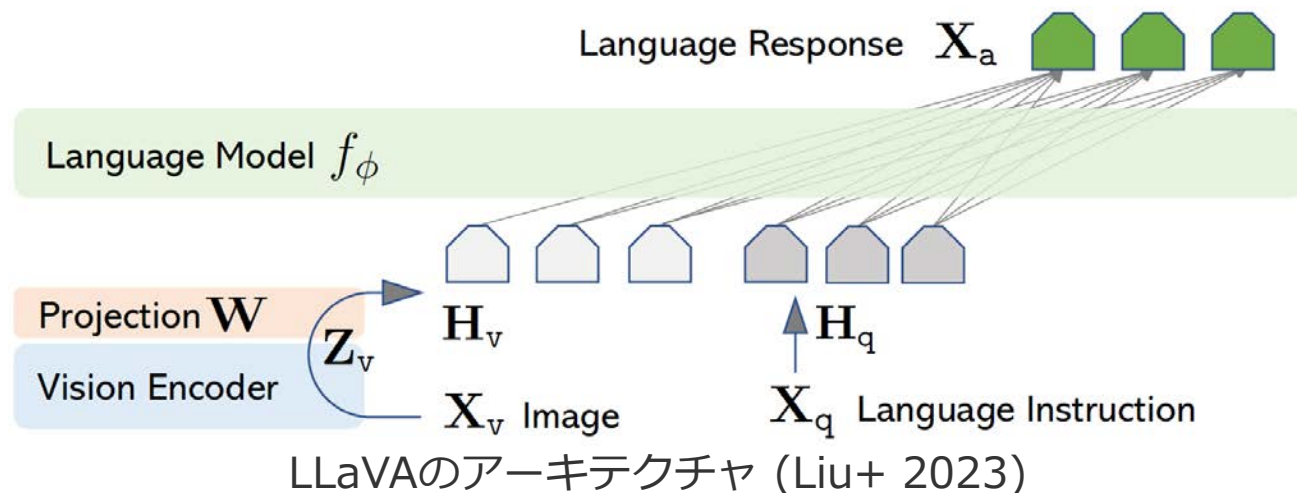
# マルチモーダル基盤モデルとは

言語や画像、音声など、様々な伝達手段の情報を統合的に処理する基盤モデル

Art & Design	Business	Science
<p><b>Question:</b> Among the following harmonic intervals, which one is constructed incorrectly?</p> <p><b>Options:</b></p> <p>(A) Major third <i>&lt;image 1&gt;</i></p> <p>(B) Diminished fifth <i>&lt;image 2&gt;</i></p> <p><b>(C) Minor seventh <i>&lt;image 3&gt;</i></b></p> <p>(D) Diminished sixth <i>&lt;image 4&gt;</i></p> 	<p><b>Question:</b> ...The graph shown is compiled from data collected by Gallup <i>&lt;image 1&gt;</i>. Find the probability that the selected Emotional Health Index Score is between 80.5 and 82?</p> <p><b>Options:</b></p> <p>(A) 0 (B) 0.2142</p> <p><b>(C) 0.3571</b> (D) 0.5</p> 	<p><b>Question:</b> <i>&lt;image 1&gt;</i> The region bounded by the graph as shown above. Choose an integral expression that can be used to find the area of R.</p> <p><b>Options:</b></p> <p><b>(A) <math>\int_0^{1.5} [f(x) - g(x)] dx</math></b></p> <p>(B) <math>\int_0^{1.5} [g(x) - f(x)] dx</math></p> <p>(C) <math>\int_0^2 [f(x) - g(x)] dx</math></p> <p>(D) <math>\int_0^2 [g(x) - x(x)] dx</math></p> 
<p><b>Subject:</b> Music; <b>Subfield:</b> Music;</p> <p><b>Image Type:</b> Sheet Music;</p> <p><b>Difficulty:</b> Medium</p>	<p><b>Subject:</b> Marketing; <b>Subfield:</b> Market Research;</p> <p><b>Image Type:</b> Plots and Charts;</p> <p><b>Difficulty:</b> Medium</p>	<p><b>Subject:</b> Math; <b>Subfield:</b> Calculus;</p> <p><b>Image Type:</b> Mathematical Notations;</p> <p><b>Difficulty:</b> Easy</p>
Health & Medicine	Humanities & Social Science	Tech & Engineering
<p><b>Question:</b> You are shown subtraction <i>&lt;image 1&gt;</i>, T2 weighted <i>&lt;image 2&gt;</i> and T1 weighted axial <i>&lt;image 3&gt;</i> from a screening breast MRI. What is the etiology of the finding in the left breast?</p> <p><b>Options:</b></p> <p>(A) Susceptibility artifact</p> <p>(B) Hematoma</p> <p><b>(C) Fat necrosis</b> (D) Silicone granuloma</p> 	<p><b>Question:</b> In the political cartoon, the United States is seen as fulfilling which of the following roles? <i>&lt;image 1&gt;</i></p> <p><b>Option:</b></p> <p>(A) Oppressor</p> <p>(B) Imperialist</p> <p><b>(C) Savior</b> (D) Isolationist</p> 	<p><b>Question:</b> Find the VCE for the circuit shown in <i>&lt;image 1&gt;</i>. Neglect VBE</p> <p><b>Answer:</b> 3.75</p> <p><b>Explanation:</b> ...<math>I_E = [(V_{EE}) / (R_E)] = [(5 \text{ V}) / (4 \text{ k-ohm})] = 1.25 \text{ mA}</math>; <math>V_{CE} = V_{CC} - I_{E} R_L = 10 \text{ V} - (1.25 \text{ mA}) 5 \text{ k-ohm}</math>; <math>V_{CE} = 10 \text{ V} - 6.25 \text{ V} = 3.75 \text{ V}</math></p> 
<p><b>Subject:</b> Clinical Medicine; <b>Subfield:</b> Clinical Radiology;</p> <p><b>Image Type:</b> Body Scans: MRI, CT.;</p> <p><b>Difficulty:</b> Hard</p>	<p><b>Subject:</b> History; <b>Subfield:</b> Modern History;</p> <p><b>Image Type:</b> Comics and Cartoons;</p> <p><b>Difficulty:</b> Easy</p>	<p><b>Subject:</b> Electronics; <b>Subfield:</b> Analog electronics;</p> <p><b>Image Type:</b> Diagrams;</p> <p><b>Difficulty:</b> Hard</p>

## MMMUの問題の例 (Liu+ 2024)

# 画像と言語のマルチモーダル基盤モデル: LLaVA (Liu+ 2023)



入力画像の例<sup>[1]</sup>

## 学習手順

1. 言語モデル $f_\phi$ と画像エンコーダを事前に学習しておく
2. 画像の特徴量 $z_v$ を言語の空間に変換する行列 $W$ を学習する

$$H_v = Wz_v, \quad z_v = g(X_v)$$

※画像 $X_v$ 、質問 $X_q$ 、答え $X_a$ が組になった学習データを用いる

3. 変換行列 $W$ と言語モデル $f_\phi$ を同時に学習する (画像エンコーダのパラメータは固定しておく)

User: What is unusual about this image?

LLaVA: The unusual aspect of this image is a man ironing clothes on the back of a minivan or van. This is not a typical place to perform this activity, as one would usually iron clothes in a more stationary and safe location, such as a home, using a regular ironing board. The scene depicted in the image is peculiar as it involves a makeshift ironing setup on a vehicle, which can be both unsafe and unconventional. Additionally, it is not clear how the man is able to maintain balance and stability while ironing clothes in such an unstable environment.

H Liu, C Li, Q Wu, Y J Lee. 2023. [Visual Instruction Tuning](#). *NeurIPS*.

[1] <https://www.barnorama.com/wp-content/uploads/2016/12/03-Confusing-Pictures.jpg>

# マルチモーダルWGでの取り組み

## 目的

- 言語や画像、音声など、様々なモダリティの情報を統合的に扱える基盤モデルを開発する
- 大規模言語モデルの言語に関する知能を外界と接続していく（グラウンディング）
- マルチモーダル基盤モデルの有用な応用、およびその評価方法を研究・確立する

## 今年度の取り組み

- VLMの学習データの整備・構築
- VLMの試作（学習ノウハウの蓄積）
  - LLMに画像のモダリティを追加
- VLMの評価基盤の検討・確立
- 医学文献に対するVLMの応用



## 来年度以降の計画

- 日本の文化や風景に対応した画像特徴量抽出器（エンコーダ）の採用・開発
- 透明性の高いマルチモーダル基盤モデルのフルスクラッチでの構築
- 動画などの新たなモダリティへの対応
- マルチモーダル基盤モデルの応用の開拓

# 今年度の取り組み (1/4): VLMの学習データの整備・構築

## 日本語版MMC4

- 画像とテキストのインターリーブ・データ
- Common Crawlの日本語ウェブページから画像をダウンロード・フィルタリングして構築

42) 口コミを見ると、土・日はいつも満席、比較的平日の方が座れるとのことでした。¥n¥n

43) ## ペットOKは外のテラス席のみ¥n¥n私達が選んだ席は、広瀬川に面したオープンスタイルのテラス席です。 [類似度: 0.264]

44) カフェに着いた時は、少し小雨が降ってる様子でしたが、幸いなことに食事の途中で小雨はやんでくれました。¥n¥n [類似度: 0.255]

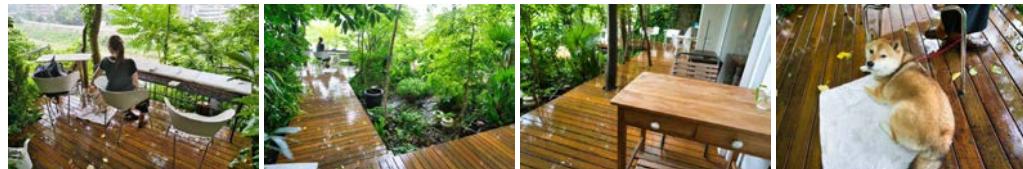
45) デッキに屋根はありませんので小雨の影響でデッキは濡れていたので、持参したこまめ用のタオルを敷いています。¥n¥n [類似度: 0.209]

46) カウンターに注文をしにいくご主人をじっとみつめるこまめ。¥n¥n

47) デッキの席は全部で10席程度ありました。

48) デッキに行くまではペットも店内を通ります。 [類似度: 0.263]

49) ちなみに店内での飲食はペットNGです。¥n¥n



43

44

45

48

<https://kurashiba.com/sendai-cafe-mozart>

## 指示チューニングデータの自動生成

- 日本で撮影された画像11,810枚に対してGPT-4oを用いて質問・応答のペアを生成



<https://huggingface.co/datasets/ThePioneer/japanese-photos>

Q: この画像にはどんな建物が写っていますか？

A: 画像には「LaLaport」と書かれた建物が写っています。

Q: この建物のファサードには何がプロジェクションされていますか？

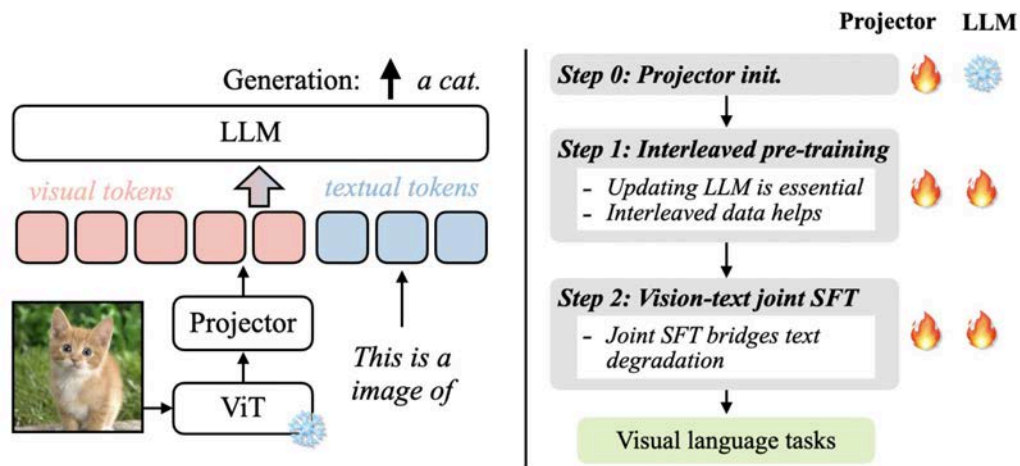
A: 建物のファサードにはアニメのキャラクターがプロジェクションされています。

Q: この画像の中で目立つ大きなフィギュアは何ですか？

A: 目立つ大きなフィギュアはガンダムのロボットです。

.....

# 今年度の取り組み (2/4): 日本語LLMのVLM化



採用したアーキテクチャ (VILA: Lin+ 2024)  
ViTはOpen AI CLIP (336px)、Projectorは2層のFFN、LLMはLlama 2 7Bを用いた



Q: 鳥の羽の色は何色でしょうか?  
A: 鳥の羽の色は青です。  
Q: 映像に写っている少年は何をしていますか?  
A: 画像の少年は、ハサミを使って紙を切っています。



Q: あなたは車を運転しており、交差点で右折しようと思っています。このまま交差点に向かってよいでしょうか?  
A: いいえ

Heron Bench (Inoue+ 2024) での生成例

学習段階	学習に用いたデータ	学習部分	学習時間 (7Bモデル)
Step 0	LLaVAのPretrainデータ (595K)	Projector	8xA100, 3.5h
Step 1	mmc4-core (25M), coyo (25M)	LLMとProjector	16x8xA100, 30h
Step 2	LLaVA1.5にinstruction tuningデータに他のデータを混ぜたもの	LLMとProjector	8xA100, 25h

J Lin et al. 2024. VILA: [On Pre-training for Visual Language Models](#). CVPR.

Y Inoue et al. 2024. Heron-Bench: A Benchmark for Evaluating Vision Language Models in Japanese. [arXiv:2404.07824](#).

# 今年度の取り組み (3/4): VLM評価基盤の構築

## 日本語指示追従性評価

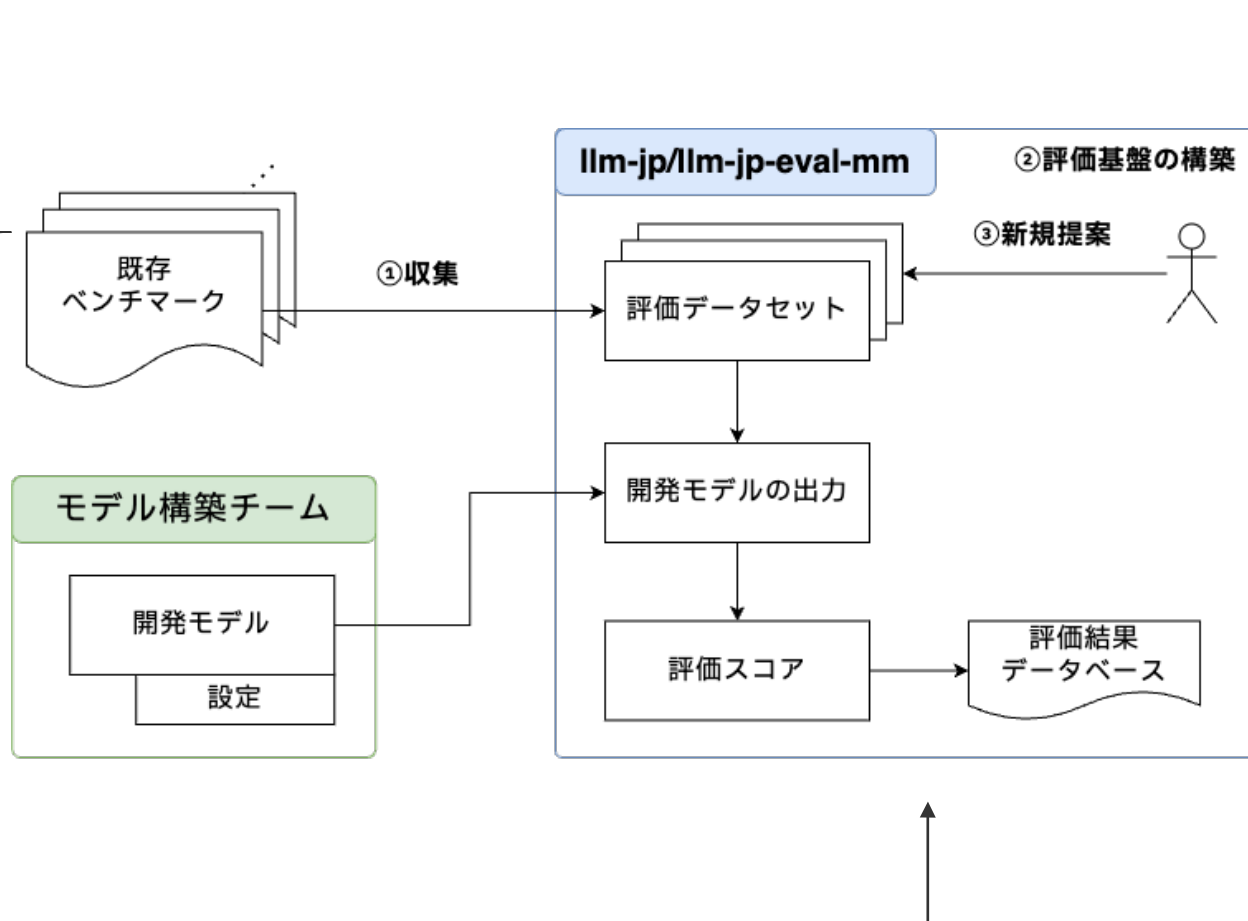
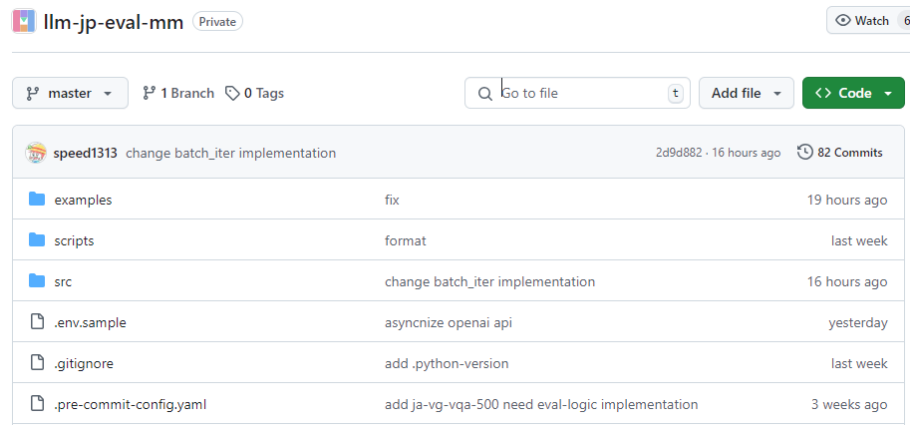
- turing-motors/Japanese-Heron-Bench
- SakanaAI/JA-VLM-Bench-In-the-Wild

## 日本語VQA

※ 淡色は検討中

- JDocQA (Onami+ 2024)
- GazeVQA (Inadumi+ 2024)
- SakanaAI/JA-VG-VQA-500

## VLMベンチマーク基盤の開発



E Onami, S Kurita, T Miyanishi, T Watanabe. 2024. [JDocQA: Japanese Document Question Answering Dataset for Generative Language Models](#). LREC-COLING.

S Inadumi, S Kawano, A Yuguchi, Y Kawanishi, K Yoshino. 2024. [A Gaze-grounded Visual Question Answering Dataset for Clarifying Ambiguous Japanese Questions](#). LREC-COLING.




# 今年度の取り組み (4/4): 医学文献に対するVLMの応用

## J-Stageの医学文献からのコーパス構築

- PDF論文からテキストを検出・抽出し、レイアウトの分析、および読み順を解析する
- PDFの処理にSurya<sup>[1]</sup>を用いる

Page-header  
医学文献

Figure  
小腫瘍影



Text  
64 mg/dl, 血糖 105 mg/dl.  
尿検査: 比重 1.030, 蛋白 100 mg/dl, プドウ糖 2+, ケトン体 1+, 潜血 2+, ウロビリノゲン normal, ビリルビン-, 亜硝酸塩-, 白血球 1+, 沈査白血球 <1/HPF, 沈査白血球 <5-9/HPF.

Text  
血液培養: 2セットとも陰性。  
胸部CT: 胆嚢底部小腫瘍影以外には異常はなく、胆嚢癌の可能性が否定できなかった(図1)。  
心電図: 特記事項なし。

Text  
血小板低下, 肝障害, 腎障害, 炎症反応高値を認めた。

Section-header  
入院時検査所見

Text  
入院時検査所見では炎症反応増多を認め、急性期DICスコアは4点とDIC基準を満たしていたが、PTやAPTTの異常はなく血小板減少に関してはDICとしては非典型的と考えた。急性腎不全。

Text  
X-1日に大腿部筋肉痛と発熱が出現した。X日前医を受診し、血小板減少, 急性腎不全, 胆汁うっ滞型の肝障害を指摘され、播種性血管内凝固症候群(DIC)疑いとして当院を紹介受診となった。

Text  
住歴: 2型糖尿病。  
家族歴: 特記事項なし。

Caption  
腹部単純CT所見  
胆嚢底部に小腫瘍影を認める

## レイアウト解析の例<sup>[2]</sup>

[1] <https://github.com/VikParuchuri/surya>

[2] 小林 俊諒, 豊嶋 弘一, 山田 英嗣, 南 博仁, 藤枝 敦史, 玉木 茂久. Jarisch-Herxheimer反応が診断の端緒となったWeil病の1例, 日本内科学会雑誌, 112巻, 6号, pp. 1005-1011, 2023.

## マルチモーダル試験コーパスの構築

- 医師国家試験問題から問題を抽出し、別冊として収録されている画像と対応付ける

24 52歳の女性。2週間前に受けた人間ドックの腹部超音波検査で胆嚢の異常を指摘されたため精査目的で来院した。自覚症状はない。既往歴に特記すべきことはない。身長158 cm、体重64 kg。BMI 25.6。体温36.2℃。腹部は平坦、軟で、圧痛を認めない。血液所見: 赤血球458万、Hb 13.7 g/dL、Ht 41%、白血球7,300。血液生化学所見: 総ビリルビン0.9 mg/dL、AST 20 U/L、ALT 18 U/L、LD 148 U/L(基準124~222)、ALP 86 U/L(基準38~113)、 $\gamma$ -GT 28 U/L(基準9~32)、CEA 1.1 ng/mL(基準5以下)、CA19-9 14 U/mL(基準37以下)。CRP 0.1 mg/dL。腹部超音波像(別冊No. 8)を別に示す。

この患者の治療方針で適切なのはどれか。

- 経過観察
- 抗菌薬投与
- 腹腔鏡下胆嚢摘出術
- 薬物による抗癌治療
- 経皮経肝胆嚢ドレナージ



別冊  
No. 8

第118回医師国家試験問題 A問題

# まとめ

- 目的
  - 言語や画像、音声など、様々なモダリティの情報を統合的に扱える基盤モデルを開発する
- 今年度の取り組み
  - VLMの学習データの整備・構築、VLMの試作、VLM評価基盤の検討・確立
  - 医学文献に対するVLMの応用に向けて
- 主なメンバー
  - 岡崎、河原、栗田、小田、相澤、金澤、笹川、前田、Zhishen、杉浦、Yin、村岡
  - 定例ミーティングに参加して頂いている方々
- 活動場所
  - 定例ミーティング: 隔週月曜日 9:00-10:00
  - [#マルチモーダル](#) Slackチャンネル
  - マルチモーダルWGの活動に興味のある方はぜひお声がけください!