# JMMMU
# A Japanese Massive Multi-discipline Multimodal Understanding Benchmark

Shota Onohara*[1], Atsuyuki Miyai*[1], Yuki Imajuku*[1], Kazuki Egashira*[1], Jeonghun Baek*[1], Xiang Yue[2], Graham Neubig[2], Kiyoharu Aizawa[1]
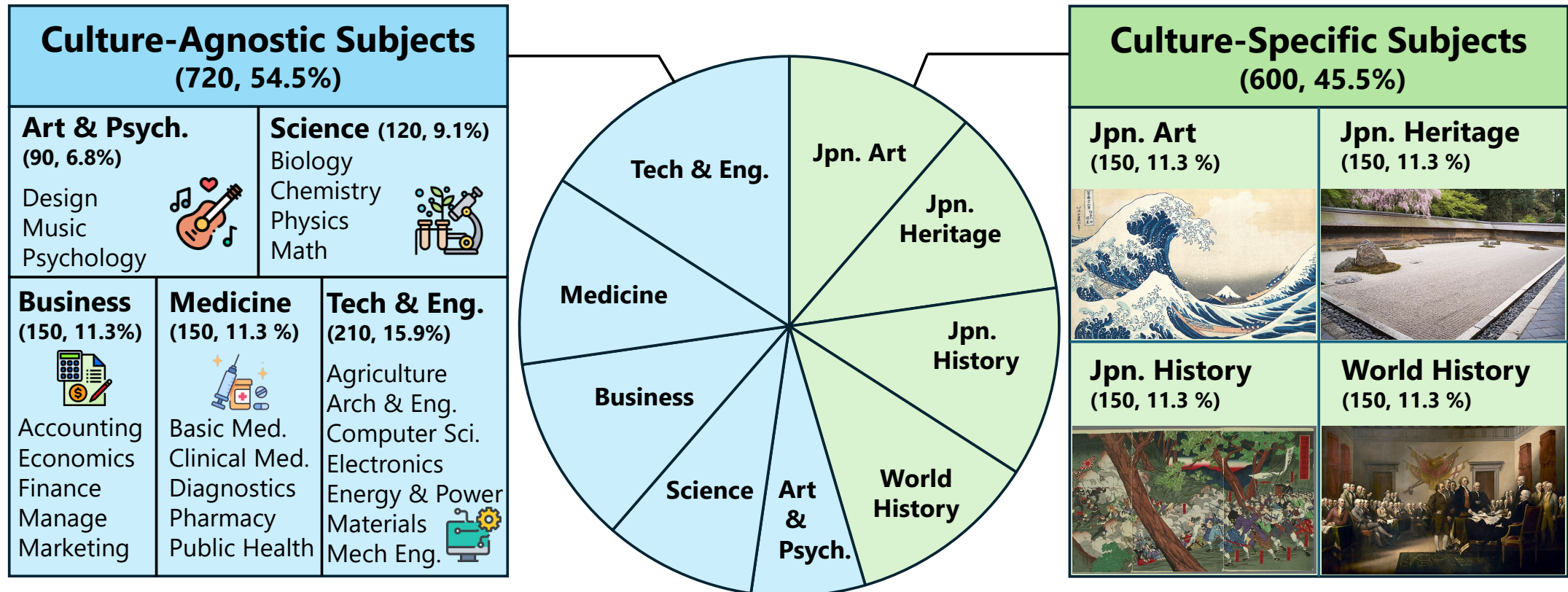
[1]The University of Tokyo, [2]Carnegie Mellon University

*Equal Contributions

# Overview of JMMMU

- A Japanese LMM Benchmark for **college-level** knowledge, **reasoning skills**, and **cultural** understanding

- 28 subjects and 1,320 questions

# LMMs Benchmark

| | Common Knowledge | Expert-level |
|---|---|---|
| **English** | • MMBench<br>• TextVQA<br>• ChartQA | • MMMU<br>• MathVista |
| **Japanese** | • Heron Bench<br>• JA-VG-VQA-500 | • **JMMMU** |

LMMs are mainly developing in English.

⬇

The need for multilingual usage

**Expert-level** benchmarks are necessary for expert AGI.

We create a Japanese expert-level benchmark for expert AGI.

# How to design multilingual evaluation?

We evaluate multilingual skills from two aspects.

**Culture-agnostic** subjects
for language skills

| | |
|---|---|
| Cash sales | $3,250 |
| Payments for inventory | 1,760 |
| Investments by owners | 3,000 |
| Supplies used | 175 |
| Cash withdrawals | 260 |
| Inventory received | 2,500 |
| Wages paid | 2,390 |
| Cash balance Dec. 1 | 4,250 |

| | |
|---|---|
| 現金売上 | ¥325,000 |
| 棚卸資産の支払い | 176,000 |
| 所有者による投資 | 300,000 |
| 使用した備品 | 17,500 |
| 現金引き出し | 26,000 |
| 受領した棚卸資産 | 250,000 |
| 支払った賃金 | 239,000 |
| 12月1日の現金残高 | 425,000 |

Translated questions of MMMU

**Culture-specific** subjects
for cultural knowledge

Brand-new questions with cultural images

# Culture-agnostic Subjects

Analyze **the cultural dependency** of MMMU [Yue+, CVPR'24] and selected 24 subjects

| **Art & Design** | **Business** | **Humanities & Social Sci** | **Science** | **Medicine** | **Tech & Eng.** |
|---|---|---|---|---|---|
| ~~Art~~ ~~Art Theory~~ Design Music | Accounting Economics Finance Manage Marketing | ~~History~~ ~~Literature~~ Psychology ~~Sociology~~ | Biology Chemistry Physics Math ~~Geography~~ | Basic Med. Clinical Med. Diagnostics Pharmacy Public Health | Agriculture Arch & Eng. Computer Sci. Electronics Energy & Power Materials Mech Eng. |

● Texts and **images** were translated by experts.

| Year | Inflation, % | Stock Market Return, % | T-Bill Return, % |
|---|---|---|---|
| 1929 | −0.2 | −14.5 | 4.8 |

What was the real return on the stock market in 1932?

| 年度 | インフレ率, % | 株式市場の収益率, % | T-Bill 収益, % |
|---|---|---|---|
| 1929 | −0.2 | −14.5 | 4.8 |

1932年の株式市場の実質リターンは何でしたか？

# Culture-specific Subjects

- 4 subjects (Japanese Art, Japanese History, Japanese Heritage, World History)

- The questions were created by native speakers.

| Japanese Heritage |
| --- |
| **Question:** \<image 1\>の城の名前は何でしょう？<br>(What is the name of the castle in \<image 1\>?)<br><br>**Options:**<br>A. 名古屋城 (Nagoya Castle)<br>B. 弘前城 (Hirosaki Castle)<br>C. 彦根城 (Hikone Castle)<br>**D. 松本城 (Matsumoto Castle)** |

# Results: Language

MMMU(CA)⚛️   JMMMU(CA)🌊

Accuracy [%]

Scores of GPT-4o stay almost **the same** before and after translation.

# Results: Language



MMMU(CA) ⚛️ JMMMU(CA) 🌊

Accuracy [%]

Scores of GPT-4o stay almost the same before and after translation.

Open-source models **drop scores** significantly after translation.

**Weakness** in open-source models

# Results: Culture

GPT-4o    Claude 3.5 Sonnet

Accuracy [%]



1. Similar performance in **MMMU**

2. Similar performance in **CA** subjects

# Results: Culture



GPT-4o     Claude 3.5 Sonnet

Accuracy [%]

3

1. similar performance in MMMU

2. similar performance in CA subjects

3. A Large gap In **CS** subjects

Translation-based evaluation is not enough.

# Summary

- JMMMU is designed to evaluate college-level skills in Japanese.

- This approach can be applied to other languages as well.



Project Page

### Leaderboard on JMMMU

We show a partial leaderboard here. Please find more information in 🏆 HF Leaderboard.

| Model | Overall (1,320) | CS (600) | CA (720) | CA (EN) (720) | Jpn. Art (150) | Jpn. Heritage (150) | Jpn. History (150) | World History (150) | Art & Psych. (90) | Business (150) | Science (120) | Medicine (150) | Tech & Eng. (210) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPT-4o (2024-05-13) | 58.6 | 66.7 | 51.8 | 52.1 | 60.7 | 70.7 | 58.7 | 76.7 | 53.3 | 55.3 | 45.8 | 61.3 | 45.2 |
| Gemini 1.5 Pro | 51.5 | 60.3 | 44.2 | 51.1 | 54.7 | 55.3 | 55.3 | 76.0 | 51.1 | 44.0 | 44.2 | 48.0 | 38.6 |
| Claude 3.5 Sonnet (2024-06-20) | 50.8 | 51.0 | 50.6 | 52.1 | 39.3 | 46.7 | 54.7 | 63.3 | 53.3 | 56.7 | 51.7 | 55.3 | 41.0 |
| LLaVA-OneVision 7B | 40.5 | 43.0 | 38.5 | 45.1 | 36.0 | 30.7 | 37.3 | 68.0 | 41.1 | 36.7 | 31.7 | 38.7 | 42.4 |
| LLaVA-NeXT 34B | 39.8 | 43.2 | 37.1 | 45.7 | 42.0 | 36.0 | 40.7 | 54.0 | 42.2 | 41.3 | 25.0 | 36.7 | 39.0 |
| InternVL2 8B | 38.3 | 42.5 | 34.7 | 43.3 | 41.3 | 38.0 | 35.3 | 55.3 | 40.0 | 36.0 | 34.2 | 34.0 | 32.4 |
| Idefics3 8B | 37.3 | 42.8 | 32.8 | 36.9 | 43.3 | 24.7 | 42.0 | 61.3 | 34.4 | 28.0 | 26.7 | 38.0 | 35.2 |
| CogVLM2 19B | 36.1 | 39.7 | 33.1 | 36.8 | 39.3 | 24.0 | 36.0 | 59.3 | 28.9 | 32.7 | 30.8 | 30.0 | 38.6 |
| Mantis-8B-siglip-llama3 | 35.5 | 39.5 | 32.2 | 36.0 | 42.0 | 30.0 | 35.3 | 50.7 | 37.8 | 28.0 | 31.7 | 37.3 | 29.5 |
| EvoVLM-JP v2 | 38.1 | 45.2 | 32.2 | 33.9 | 44.0 | 40.0 | 42.0 | 54.7 | 32.2 | 28.7 | 28.3 | 38.7 | 32.4 |

JMMMU is combined to LMMs-Eval.
You can easily evaluate your model🎉

# Appendix

# Detailed Results

| Models | Overall | CS | CA | CA (EN) | Jpn. Art | Jpn. Heritage | Jpn. History | World History | Art & Psych. | Business | Science | Health & Medicine | Tech & Eng. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | (1,320) | (600) | (720) | (720) | (150) | (150) | (150) | (150) | (90) | (150) | (120) | (150) | (210) |
| Random | 24.8 | 25.0 | 24.6 | 24.6 | 25.0 | 25.0 | 25.0 | 25.0 | 25.4 | 25.0 | 22.8 | 25.6 | 24.3 |
| **Open Source** | | | | | | | | | | | | | |
| LLaVA-OV-0.5B | 26.0 | 23.3 | 28.2 | 29.4 | 22.7 | 22.7 | 24.0 | 24.0 | 26.7 | 27.3 | 24.2 | 30.7 | 30.0 |
| InternVL2-2B | 28.3 | 29.2 | 27.6 | 31.9 | 31.3 | 22.7 | 30.7 | 32.0 | 30.0 | 30.0 | 30.8 | 25.3 | 24.8 |
| xGen-MM | 28.6 | 28.2 | 28.9 | 35.7 | 30.0 | 20.7 | 22.7 | 39.3 | 32.2 | 21.3 | 22.5 | 36.7 | 31.0 |
| Phi-3v | 29.5 | 26.5 | 31.9 | 37.6 | 31.3 | 18.7 | 29.3 | 26.7 | 26.7 | 28.7 | 25.8 | 37.3 | 36.2 |
| LLaVA-1.6-13B | 31.1 | 33.7 | 29.0 | 29.9 | 32.0 | 24.0 | 32.0 | 46.7 | 25.6 | 28.7 | 30.0 | 34.0 | 26.7 |
| Idefics2-8B | 31.9 | 37.0 | 27.6 | 35.1 | 40.7 | 24.0 | 30.0 | 53.3 | 32.2 | 22.7 | 22.5 | 32.0 | 29.0 |
| Phi-3.5v | 32.4 | 34.3 | 30.8 | 39.2 | 37.3 | 27.3 | 35.3 | 37.3 | 27.8 | 31.3 | 30.0 | 36.7 | 28.1 |
| †LLaVA CALM2 | 34.9 | 41.5 | 29.4 | 29.9 | 42.7 | 36.7 | 40.0 | 46.7 | 27.8 | 26.0 | 26.7 | 34.0 | 31.0 |
| Mantis 8B | 35.5 | 39.5 | 32.2 | 36.0 | 42.0 | 30.0 | 35.3 | 50.7 | 37.8 | 28.0 | 31.7 | 37.3 | 29.5 |
| CogVLM2-19B | 36.1 | 39.7 | 33.1 | 36.8 | 39.3 | 24.0 | 36.0 | 59.3 | 28.9 | 32.7 | 30.8 | 30.0 | 38.6 |
| Idefics3-8B | 37.3 | 42.8 | 32.8 | 36.9 | 43.3 | 24.7 | **42.0** | 61.3 | 34.4 | 28.0 | 26.7 | 38.0 | 35.2 |
| †EvoVLM JP v2 | 38.1 | **45.2** | 32.2 | 33.9 | **44.0** | **40.0** | **42.0** | 54.7 | 32.2 | 28.7 | 28.3 | **38.7** | 32.4 |
| InternVL2-8B | 38.3 | 42.5 | 34.7 | 43.3 | 41.3 | 38.0 | 35.3 | 55.3 | 40.0 | 36.0 | **34.2** | 34.0 | 32.4 |
| LLaVA-1.6-34B | 39.8 | 43.2 | 37.1 | 45.7 | 42.0 | 36.0 | 40.7 | 54.0 | **42.2** | **41.3** | 25.0 | 36.7 | 39.0 |
| LLaVA-OV-7B | **40.5** | 43.0 | **38.5** | 45.1 | 36.0 | 30.7 | 37.3 | **68.0** | 41.1 | 36.7 | 31.7 | **38.7** | **42.4** |
| **Proprietary** | | | | | | | | | | | | | |
| Claude 3.5 Sonnet | 50.8 | 51.0 | 50.6 | 52.1 | 39.3 | 46.7 | 54.7 | 63.3 | **53.3** | **56.7** | **51.7** | 55.3 | 41.0 |
| Gemini 1.5 Pro | 51.5 | 60.3 | 44.2 | 51.1 | 54.7 | 55.3 | 55.3 | 76.0 | 51.1 | 44.0 | 44.2 | 48.0 | 38.6 |
| GPT-4o | **58.6** | **66.7** | **51.8** | 52.1 | **60.7** | **70.7** | **58.7** | **76.7** | **53.3** | 55.3 | 45.8 | **61.3** | **45.2** |
| **Text Only** | | | | | | | | | | | | | |
| GPT-4o text | 38.1 | 35.5 | 40.3 | 44.9 | 32.7 | 32.0 | 35.3 | 42.0 | 38.9 | 36.0 | 41.7 | 45.3 | 39.5 |

# An Error in GPT-4o

**Question:**
で見られることで有名な動物は何でしょう？
(What animal is known to be seen in <image 1>.)
**Option:**
A. ヤンバルクイナ (Okinawa rail)
B. イリオモテヤマネコ (Iriomote wildcat)
C. イヌワシ (golden eagle)
D. トキ (crested ibis)

**GPT-4o's Answer:**
この画像は日本の石垣島を示しています。
石垣島で見られることで有名な動物は「イリオモテヤマネコ」です。したがって、正解はB. イリオモテヤマネコです。
This image shows Ishigaki Island, Japan. The Iriomote wildcat is a well-known animal found on Ishigaki Island. Therefore, the correct answer is B. Iriomote wildcat.

**Ground Truth:** D.トキ (crested ibis)