# Chitrarth: Bridging Vision and Language for a Billion People

**Shaharukh Khan, Ayush Tarun, Abhinav Ravi \*, Ali Faraz, Akshat Patidar**
**Praveen Pokala \*, Anagha Bhangare, Raja Kolla, Chandra Khatri \*, Shubham Agarwal \***

Krutrim AI, Bangalore, India
\* Senior Contributors
Contact: {shaharukh.khan, abhinav.ravi, shubham.agarwal1}@olakrutrim.com

## Abstract

Recent multimodal foundation models are primarily trained on English or high resource European language data, which hinders their applicability to other medium and low-resource languages. To address this limitation, we introduce *Chitrarth* (Chitra: Image; Artha: Meaning), an inclusive Vision-Language Model (VLM), specifically targeting the rich linguistic diversity and visual reasoning across 10 prominent Indian languages. Our model effectively integrates a state-of-the-art (SOTA) multilingual Large Language Model (LLM) with a vision module, primarily trained on multilingual image-text data. Furthermore, we also introduce BharatBench, a comprehensive framework for evaluating VLMs across various Indian languages, ultimately contributing to more diverse and effective AI systems. Our model achieves SOTA results for benchmarks across low resource languages while retaining its efficiency in English. Through our research, we aim to set new benchmarks in multilingual-multimodal capabilities, offering substantial improvements over existing models and establishing a foundation to facilitate future advancements in this arena.

## 1 Introduction

With the success and demonstrated effectiveness of *Visual instruction tuning* [Liu et al., 2023, 2024], recent years witnessed a surge of interest in developing general purpose multimodal conversational agents. These unified foundation models excel at algorithmic reasoning and generic perception tasks like image captioning, visual question answering, text-based image retrieval, etc. [Lu et al., 2024a, Laurençon et al., 2024b, Tong et al., 2024, Xue et al., 2024], and more specialized frameworks, for instance, converting Scalable Vector Graphics (SVGs) to code [Rodriguez et al., 2023]. Often, these models rely on pre-trained Large Language Models (LLMs) [Brown et al., 2020, Touvron et al., 2023a, Achiam et al., 2023, Chiang et al., 2023, Touvron et al., 2023b, Gemini et al., 2023, Jiang et al., 2024, Gemma et al., 2024, Dubey et al., 2024] as the transformer [Vaswani et al., 2017] backbones, primarily trained on English or high resource European languages.

This work is driven by two main motivations: 1. *Language diversity gap:* Most Vision Language Models (VLMs) are predominantly trained on English datasets, overlooking the linguistic needs of non-English languages, particularly from the Indian subcontinent. 2. *Lack of low resource language benchmarks:* Absence of corresponding VLM benchmarks hinders the progress for these low resource Indic languages. We aim to address these issues through our research and serve a broader audience, encompassing billions of people.

Few LLMs have been developed specifically for Indic languages [Gala et al., 2024, Labs, 2023, Balachandran, 2023, Kohli et al., 2023], most of which extend and fine-tune text-only English-centric

Multimodal Algorithmic Reasoning, Neural Information Processing Systems (NeurIPS 2024).

LLMs. Naturally, they fail to fully capture the nuances of the language, with the exception of models like [Kallappa et al., 2024, Bendale et al., 2024], trained from scratch. Our model builds upon the recent success of Krutrim LLM [Kallappa et al., 2024], which supports English and 10 other languages including Hindi, Bengali, Telugu, Tamil, Marathi, Gujarati, Kannada, Malayalam, Odia, and Assamese, representing a significant portion of the cultural and linguistic diversity in India.

Another key challenge is the limited availability of low resource data with Indic languages significantly under-represented in Common Crawl despite India (or Bharat) making up 18% of the global population. For instance, Hindi, in spite of being the third most spoken, does not appear among the top 20 languages [Buck et al., 2014, Penedo et al., 2023]. To enhance our model's cross-lingual generalization abilities, we translate the open-source multimodal training datasets into the 10 Indic languages natively supported by the backbone LLM. Developing this multilingual dataset is a substantial endeavor aimed at addressing the disparity between high-resource and relatively low-resource Indian languages in the context of vision-language models.

In this paper, we present our multimodal LLM, which employs the Krutrim multilingual LLM backbone [Kallappa et al., 2024] in conjunction with a pre-trained visual image encoder [Dosovitskiy et al., 2020]. Figure1 demonstrates the multi-lingual capability of our model across major Indian languages. A brief summary of our contribution is provided below:

- We introduce Chitrarth (Chitra: Image; Artha: Meaning), a Multimodal LLM model which leverages images and language modalities for a range of visual tasks such as image captioning, visual question answering in the multilinugal context. We further present optimal training recipe including data composition and architecture configuration.

- We also present *BharatBench*, a comprehensive evaluation benchmark suite designed for 10 under-resourced Indic languages across 3 tasks, which we will make available to the research community upon acceptance.

- Finally, we evaluate Chitrarth and prior baselines on both existing English academic datasets as well as the proposed evaluation framework and demonstrate the effectiveness of our model, using different training strategies and ablations, achieving SOTA results on 3 out of 5 English datasets and propose benchmark results on the derived multi-lingual datasets.

The remainder of the paper is structured as: Section 2 reviews recent related research on VLMs. Section 3 provides a detailed description of our Chitrarth model with information about training data mix in Section 4. Section 5 introduces the BharatBench evaluation framework that we propose, while Section 6 presents the experimental results. Finally, Section 7 offers concluding remarks.

## 2 Related Work

### 2.1 English-centric VLMs

Recent studies [Laurençon et al., 2024b,a, Tong et al., 2024] have investigated design strategies for multi-stage training pipelines in contemporary VLMs. Typically, these models rely on pre-trained LLMs; however, there are some exceptions where models are trained from scratch [Chameleon, 2024, Lu et al., 2024b]. Prior works like Flamingo[Alayrac et al., 2022] leverage a Perceiver Resampler [Jaegle et al., 2021] to inject visual features into the language model through cross-attention, promoting quick adaptation to various tasks with few labeled examples. The LLaVA family models [Liu et al., 2023, 2024], including LLaVA-1.5 and LLaVA-1.6, demonstrated intriguing multimodal comprehension capabilities by integrating advanced language models with vision encoders through visual instruction tuning. PaliGemma [Beyer et al., 2024], optimized for tasks that require deep integration of visual and textual data, is designed to excel in scenarios where English is the primary language. Florence-2 [Xiao et al., 2024] focuses on handling diverse tasks from simple text prompts addressing the complexity of spatial hierarchy and semantic granularity. The Idefics family [Laurençon et al., 2024b,a] is focused on substantially enhancing capabilities around OCR, document interpretation and visual reasoning functionalities. CogVLM [Wang et al., 2023] drives an intricate fusion of language and vision features unlike other VLMs, which rely on the shallow alignment method. PALI models [Chen et al., 2022] on the other hand explored contrastive pretraining and higher resolution training for the VLM tasks.

આ છબી માટે એક લીટી કેપ્શન લખો.
[Gujarati-English Trans: Write a one line caption for this image.]

એક વૃદ્ધ માણસ, સંભવતઃ એક સંત, પીંછાની પીંછી સાથે પુસ્તક પર લખે છે.
[Gujarati-English Trans: An old man, probably a saint, writes on a book with a feather brush.]

ವ್ಯಕ್ತಿ ಯಾವ ವಸ್ತುವಿನ ಮೇಲೆ ಕುಳಿತಿದ್ದಾನೆ?
[Kannada-English Trans: What object is the person sitting on?]

ಆ ವ್ಯಕ್ತಿಯು ಕೆಂಪು ಹೂವಿನ ಮೇಲೆ ಕುಳಿತಿದ್ದಾನೆ.
[Kannada-English Trans: That person is sitting on a red flower.]

उनके हाथ में क्या है?
[Hindi-English Trans: What's in his hand?]

उनके हाथ में एक पंख है।
[Hindi-English Trans: He has a feather in his hand.]

त्याच्या हातात कोणत्या पक्ष्याचे पंख आहे?
[Marathi-English Trans: Which bird's feather is in his hand?]

त्याच्या हातात मोराचे पंख आहेत.
[Marathi-English Trans: He has peacock feathers in his hands..]

అతను ఏ రంగు లోయర్ గుడ్డ ధరించాడు?
[Telugu-English Trans: What color lower cloth is he wearing?]

చిత్రంలో ఉన్న వ్యక్తి పసుపు రంగు లోయర్ గుడ్డ ధరించి ఉన్నాడు.
[Telugu-English Trans: The person in the picture is wearing a yellow lower cloth.]

What is this image about?

The image features a man with a long white beard sitting on a flower, which is placed on a table. He is wearing an orange robe and appears to be writing on a piece of paper. The man is surrounded by a serene and peaceful atmosphere, with a few birds flying in the background. The scene is a beautiful and artistic representation of the man engaging in his writing activity.

படத்தின் உணர்ச்சித் தொனி என்ன, காட்சி கூறுகள் மூலம் அது எவ்வாறு வெளிப்படுத்தப்படுகிறது?
[Tamil-English Trans: What is the emotional tone of the film, and how is it expressed through visual elements?]

அமைதியாகவும் உள்ளது. இந்த தொனியை உருவாக்கும் காட்சி கூறுகளில் ஒரு தாடி கொண்ட ஒரு மனிதன் ஒரு புத்தகத்தில் எழுதுவது, ஒரு பெரிய புத்தகத்தின் முன் அமர்ந்திருப்பது மற்றும் ஒரு பெரிய புத்தகத்தின் முன் அமர்ந்திருப்பது ஆகியவை அடங்கும். கூடுதலாக, ஒரு புத்தகத்தின் முன் அமர்ந்திருக்கும் மனிதனின் உருவம் அறிவு மற்றும் சிந்தனையின் உணர்வைத் தூண்டுகிறது, இது அமைதியான மற்றும் அமைதியான சூழ்நிலையை மேலும் வலியுறுத்துகிறது.
[Tamil-English Trans: It is also quiet. Visual elements that create this tone include a bearded man writing in a book, sitting in front of a large book, and sitting in front of a large book. Additionally, the image of a man sitting in front of a book evokes a sense of knowledge and contemplation, which further emphasizes the calm and peaceful atmosphere.]

ഈ ചിത്രത്തിലെ ഏതെങ്കിലും പ്രത്യേക ഐക്കണോഗ്രഫി അല്ലെങ്കിൽ ചിഹ്നങ്ങൾ നിങ്ങൾക്ക് തിരിച്ചറിയാനാകുമോ? അവർ എന്തിനെ പ്രതിനിധീകരിക്കുന്നു?
[Malayalam-English Trans: Can you recognize any particular iconography or symbols in this image? What do they represent?]

ചിത്രത്തിൽ, ഒരു വെളുത്ത താടിയുള്ള ഒരു മനുഷ്യൻ ഒരു പീഠത്തിൽ ഇരിക്കുന്നു, ഒരു പേപ്പർ കഷണം എഴുതുന്നു. ഈ രംഗത്തിൽ ഒരു പ്രത്യേക ചിഹ്നമോ ഐക്കണോഗ്രാഫിയോ ഉണ്ട്: ഒരു മയിൽ തൂവൽ. മയിൽ തൂവലിന്റെ സാന്നിധ്യം സൂചിപ്പിക്കുന്നത് ഈ രംഗം ഒരു പ്രത്യേക സാംസ്കാരികമോ മതപരമോ ആയ പ്രാധാന്യമുള്ള ഒരു വ്യക്തിയെ ചിത്രീകരിക്കുന്നു എന്നാണ്. മയിൽ തൂവലുകൾ പലപ്പോഴും സൗന്ദര്യം, ജ്ഞാനം, ആത്മീയത എന്നിവയുമായി ബന്ധപ്പെട്ടിരിക്കുന്നു. ഈ ചിഹ്നത്തിന്റെ സാന്നിധ്യം സൂചിപ്പിക്കുന്നത് ചിത്രം ഒരു പ്രത്യേക സംസ്കാരത്തിൽ നിന്നോ മതത്തിൽ നിന്നോ ഉള്ള ഒരു വ്യക്തിയെ ചിത്രീകരിക്കുന്നു എന്നാണ്, ഒരുപക്ഷേ ഒരു ഹിന്ദു അല്ലെങ്കിൽ ബുദ്ധ വ്യക്തിയായിരിക്കാം, കാരണം ഈ രണ്ട് മതങ്ങളിലും മയിൽ തൂവലുകൾ പ്രതീകാത്മക പ്രാധാന്യം വഹിക്കുന്നു.

[Malayalam-English Trans: In the picture, a man with a white beard is sitting on a pedestal, writing on a piece of paper. There is a special symbol or iconography in this scene: a peacock feather. The presence of peacock feathers indicates that the scene depicts a person of particular cultural or religious significance. Peacock feathers are often associated with beauty, wisdom, and spirituality. The presence of this symbol indicates that the image depicts a person from a particular culture or religion, perhaps a Hindu or Buddhist person, as peacock feathers carry symbolic significance in both of these religions.]

এই ব্যক্তি একটি ময়ূর পালক সঙ্গে কি করছেন?
[Bengali-English Trans: What is this person doing with a peacock feather?]

ছবির ব্যক্তিটি একটি ময়ূরের পালক ধরে আছেন, যা ইঙ্গিত করে যে তিনি এটি ব্যবহার করে কিছু লেখার বা আঁকার প্রক্রিয়ায় থাকতে পারেন। ময়ূরের পালকগুলি প্রায়শই শৈল্পিক অভিব্যক্তির জন্য ব্যবহৃত হয়, বিশেষত ভারতীয় শিল্পের ক্ষেত্রে, যেখানে তারা প্রতীকী এবং সাংস্কৃতিক তাৎপর্য ধারণ করে। এই প্রসঙ্গে, ব্যক্তিটি হয়তো একটি শিল্পকর্ম তৈরি করছেন, একটি বার্তা লিখছেন বা ময়ূরের পালক ব্যবহার করে একটি নির্দিষ্ট ধারণা প্রকাশ করছেন।
[Bengali-English Trans: The person in the picture is holding a peacock feather, indicating that he may be in the process of writing or drawing something using it. Peacock feathers are often used for artistic expression, especially in Indian art, where they hold symbolic and cultural significance. In this context, the person may be creating a work of art, writing a message or expressing a specific idea using peacock feathers.]

Figure 1: **Multi lingual capability of Chitrarth model across major Indian languages.** For the same underlying image, we present question-answer pairs in English and several Indian languages - Gujarati, Kannada, Hindi, Marathi, Telugu, Tamil, Malayalam, and Bengali (in order). Questions are highlighted in purple, and responses are shown in orange (provided with English translations). The model accurately understands and identifies the 'image of a saint writing a book with a feather' and correctly addresses related questions in different languages.

## 2.2 Multi-lingual VLMs

Qwen-VL [Bai et al., 2023] is a multilingual VLM, trained on English and Chinese data, supporting diverse instructions and multi-image context analysis. InternVL 1.5 [Chen et al., 2024] proposed an enhanced vision encoder and a superior bilingual dataset, i.e., English and Chinese. Phi-3 family [Abdin et al., 2024] offer multilingual, multimodal, and long-context support in 11 languages, including English, across the world but do not cover Indian languages. PALO [Maaz et al., 2024] is the closest VLM to our research, however supporting only 3 Indian languages Hindi, Urdu, and Bengali apart from the other high-to-medium resource language offerings. To our knowledge, no other open-source multimodal LLMs include low-resource Indic languages in the training mix. In contrast, our work introduces a multilingual VLM system that supports ten Indian languages.
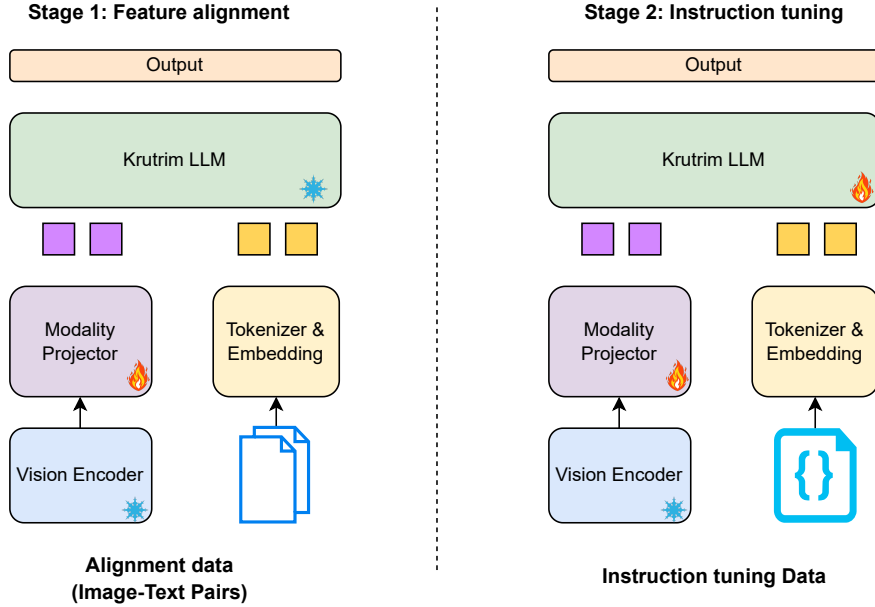


Figure 2: **Chitrarth model features a fully autoregressive architecture with a two-stage training process.** In Stage 1, the model is trained using images and their descriptions, aligning visual and linguistic embeddings through image-caption pairs. In Stage 2, model is fine-tuned on multimodal instruction-following and domain-specific academic datasets.

## 3 Chitrarth: A Multilingual Vision-Language Model

In this section, we outline the architecture of our proposed Chitrarth model. Chitrarth is an autoregressive VLM where the input image is tokenized into visual tokens, combined with textual instruction tokens and fed into the large language model (LLM). Inspired by the versatile and widely followed LLaVA [Liu et al., 2023, 2024] framework, our model incorporates several key components, as illustrated in Figure 2, where we use pre-trained Krutrim LLM [Kallappa et al., 2024] instead, as the autoregressive multi-lingual LLM backbone.

For multimodal training, we start by encoding images using a vision encoder. The modality projection layer (adapter/connector) maps the vision embeddings into the LLM embedding space, producing a sequence of visual tokens. The multilingual LLM then generate responses based on these visual tokens. The Krutrim LLM [Kallappa et al., 2024] supports a context length of 4096 tokens, of which 576 tokens (14X14 patch size results in 729 tokens) are allocated for image representation after the modality projection. We explore different configurations for the projection layer, including a single-layer projection [Liu et al., 2023, 2024] and a two-layer MLP vision-language connector with non-linearity [Liu et al., 2023]. Additionally, we experiment with various vision encoders, including

the pre-trained CLIP ViT-L/14@336px [Radford et al., 2021] and SigLIP-SO400M [Zhai et al., 2023]. Our model is trained in multiple stages:

**Stage 1: Pre-Training (PT) for Feature Alignment.** In this stage, we conduct pre-training using image-text pairs, with the projector layer being trained while keeping the vision encoder and LLM fixed. Each sample is treated as a single-turn conversational instruction for tuning.

**Stage 2: Instruction Tuning (IT).** In this stage, we maintain the vision encoder in a frozen state, following the approach used in LLaVA models [Liu et al., 2023, 2024]. However, unlike the previous stage, we also update the weights of the LLM in addition to tuning the modality projection layer. The objective of this stage is to develop a general-purpose multimodal agent (chatbot) capable of comprehending and executing complex instructions across multiple conversational turns. We describe the datasets used in both the stages in the next section.

## 4 Dataset

Figure 3 illustrates the language distribution of our data mix for both the training stages, which we describe in more detail below:

**Stage 1:** For Stage 1 adapter Pre-Training (PT), we use the 1.2 million-sample ShareGPT4V-PT dataset [Chen et al., 2023], which demonstrated consistent superior performance compared to other PT datasets, such as LLaVA-Pretrain-LCS-558K [Liu et al., 2023], in our preliminary experiments. This dataset was subsequently translated into the ten Indic languages supported by the Krutrim LLM. Specifically, we use the open-source model, IndicTrans2 [Gala et al., 2023] for this text-only translation task. IndicTrans2 outperformed other translation services (Yandex, ChatGPT, Google Translate, and Bard) in small-scale in-house qualitative human evaluation (win rates 93% and 80% for Bengali and Marathi respectively). We ensure the pre-training data remained at 1.2M points, with half of the data in English, and sample translations across different languages in an equal ratio to create a balanced multilingual dataset. This approach was designed to preserve linguistic diversity and computational efficiency, thereby ensuring robust performance in English while developing capabilities in the Indic languages. The balanced dataset mitigates potential biases towards any single language, fostering equitable performance across all supported languages.

**Stage 2:** The Stage 2 Instruction Tuning (IT) dataset is notably more intricate. The core element of this dataset is the complete English version of LLaVA-1.5-665K [Liu et al., 2024]. Additionally, we translate LLaVA-Instruct-150K [Liu et al., 2023] into ten languages using the methodology outlined in Stage 1. Our dataset also incorporates the Cauldron dataset [Laurençon et al., 2024b], which includes 50 academic vision-language tasks along with its corresponding in-house translations. Furthermore, we add a substantial collection of images reflecting Indian cultural diversity comprising prominent personalities, monuments, artwork, culinary dishes, and more; transformed into multilingual pluralistic instruction tuning data, analogous to the open-source English-based LLaVA-IT datasets. Lastly, our dataset features high-quality, text-only English proprietary data. The final composition of the dataset includes approximately 880K English and 90K samples in multiple languages, ensuring a balanced and diverse dataset. This comprehensive range of content supports the development of a model capable of generating and understanding complex descriptions across various domains and visual scenarios, thereby enhancing its reasoning capabilities.

## 5 BharatBench Evaluation Suite

Although recent efforts have advanced text-only multilingual evaluation [Ahuja et al., 2023, Singh et al., 2024], there is still a lack of evaluation framework for multimodal multilingual scenarios. We introduce BharatBench, a benchmark designed to assess the image understanding capabilities of multilingual Vision-Language Models (VLMs). Expanding upon LLaVa-Bench (In-the-Wild) [Liu et al., 2023], initially adapted for Hindi and Bengali by [Maaz et al., 2024], we further broadened the benchmark to cover eight additional low resource languages. This extension now forms part of our comprehensive benchmark suite. Furthermore, we include translated versions of prominent VLM evaluation datasets, such as MMVet [Yu et al., 2023] and POPE [Li et al., 2023] covering all ten languages in our study, in addition to English.
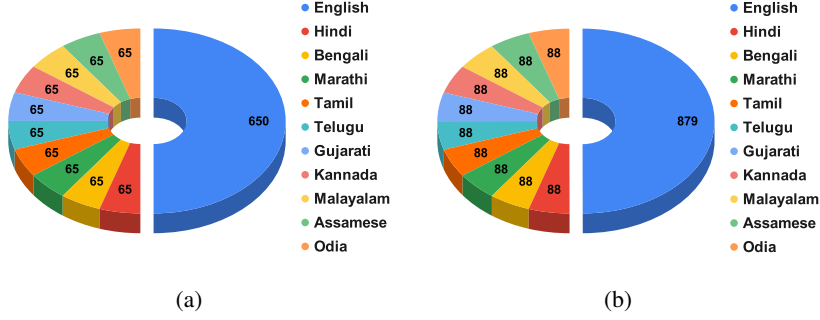
Figure 3: **Language distribution in data mix.** (a) Stage 1 data consists of 1.2M ShareGPT4V in the original English version (650K) and remaining Indian language translations (65K each) (b) Stage 2 data involves 879K samples in English and 88K for each respective language, discussed in Section 4.

In essence, we intentionally chose to extend existing benchmarks through translation, which not only facilitates the creation of valuable multi-way parallel data but also addresses data scarcity issues and leverages the inherent quality of established evaluation frameworks [Singh et al., 2024]. This methodology enhances our ability to evaluate and advance multimodal models in a multilingual context. We followed similar guiding principles while creating the training datasets described earlier.

# 6 Experiments

## 6.1 Implementation

We use PyTorch [Paszke et al., 2019] based HuggingFace Transformers [Wolf et al., 2019] for our experiments. Our Stage 1 and 2 tuning use hyperparameters consistent with those of the LLaVA model [Liu et al., 2023], unless otherwise specified. Particularly, we train the model for 1 epoch in both the stages with an overall batch size of 256 in Stage 1 and 128 in Stage 2. We used cosine LR scheduler with Adam optimizer and a learning rate of 2e-3 and 2e-5 in both the stages respectively. We consider IDEFICS 2 [Laurençon et al., 2024b] and PALO [Maaz et al., 2024] as respective English and multi-lingual baselines and report results from their published work. All our models are trained on 8 × H100 GPUs which takes around 8 hours for Stage 1 and 18 hours for Stage 2 tuning.

## 6.2 English academic benchmarks

We also evaluate our model using a range of English academic benchmarks, including VQA-v2 [Goyal et al., 2017] and GQA [Hudson and Manning, 2019] for visual perception, VizWiz [Gurari et al., 2018] for zero-shot generalization on questions posed by visually impaired users, and TextVQA [Singh et al., 2019] for text-rich visual question answering. We also use POPE [Li et al., 2023] to assess hallucination tendencies, MME [Fu et al., 2023] for yes/no question responses, and LLaVA-Bench (In-the-Wild) [Liu et al., 2023] and MM-Vet [Yu et al., 2023] for visual conversation capabilities. Evaluation scores are reported following prior works.

| Bench | Telugu | Hindi | Bengali | Malayalam | Kannada | Assamese | Tamil | Marathi | Gujarati | Odia | English |
|---|---|---|---|---|---|---|---|---|---|---|---|
| POPE | 79.9 | 78.68 | 83.24 | 85.29 | 85.52 | 55.59 | 83.28 | 79.17 | 84.75 | 82.03 | 87.63 |
| LLaVA-Bench | 54.8 | 51.5 | 53.7 | 55.5 | 58.1 | 59.1 | 58.3 | 52.80 | 55.90 | 62.80 | 67.90 |
| MMVet | 43.76 | 38.85 | 33.24 | 25.36 | 46.19 | 37.29 | 34.31 | 40.96 | 39.03 | 19.67 | 30.49 |

Table 1: **Performance of Chitrarth on BharatBench Evaluation framework.** Our model is unique in its ability to handle all included languages, setting a baseline for future research.

## 6.3 Results

On the English academic datasets, our model depicts State-of-the-art (SOTA) results for POPE, VQAv2 and GQA compared to the baseline models, while remaining competitive on TextVQA and

**(A)**

এই মহিলার সৌন্দর্য নিয়ে একটি কবিতা লিখুন

সুন্দরী মহিলাটি দাঁড়িয়ে আছে, তার সাদা পোশাকের সাথে মিলে যাওয়া একটি সবুজ স্কার্ট পরে, তার মাথায় একটি ওড়না দিয়ে তার কমনীয়তা প্রদর্শন করেছে। তার চোখ বন্ধ, তার হাত তার বুকের কাছে ধরে, তার সৌন্দর্যের প্রতি মনোযোগ দেয়। তার সৌন্দর্য তার চারপাশের প্রতিটি জিনিসকে আলোকিত করে, তার চারপাশের বাতাসে একটি অলৌকিক উজ্জ্বলতা যোগ করে। তার সৌন্দর্য প্রকৃতির সৌন্দর্যের মতো, যা প্রশংসা এবং প্রশংসা করার জন্য।
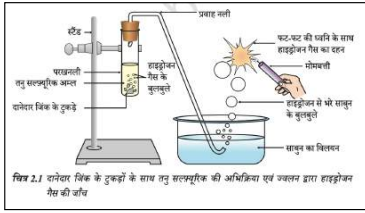
**(B)**

Extract following attributes of the given cloth: 1. Neck type 2. Occasion to wear 3. Sleeve length 4. Hemline 5. Dominant colour

1. Neck type: The neck of the dress is scooped.
2. Occasion to wear: The dress is suitable for a casual occasion or a party.
3. Sleeve length: The dress has long sleeves.
4. Hemline: The hemline of the dress is rounded.
5. Dominant colour: The dominant colour of the dress is pink.

**(C)**

यहाँ कौन सी प्रक्रिया दिखाई गई है?

छवि में एक प्रयोगशाला सेटिंग में एक रासायनिक प्रतिक्रिया को दर्शाया गया है। प्रतिक्रिया में एक बीकर, एक फ्लास्क और एक कांच के कटोरे का उपयोग शामिल है। प्रतिक्रिया के दौरान, एक बुलबुला बनता है, जो दर्शाता है कि एक रासायनिक प्रतिक्रिया हो रही है। प्रयोगशाला में एक व्यक्ति भी मौजूद होता है, जो संभवतः प्रयोग की निगरानी करता है या प्रतिक्रिया के परिणाम की प्रतीक्षा करता है। प्रयोगशाला में विभिन्न उपकरणों और उपकरणों की उपस्थिति से पता चलता है कि यह एक वैज्ञानिक प्रयोग है जिसमें रासायनिक प्रतिक्रियाओं का अध्ययन या प्रदर्शन शामिल है।

**(D)**

What is this UI about?

The UI is a screen displaying a person's meditation statistics. It shows the number of minutes spent meditating, the number of meditation sessions, and the number of points earned. The screen also includes a bar chart to visually represent the data. Additionally, there is a picture of a person's head, possibly the user, and a small image of a man in the top right corner of the screen.

Explain individual details.

1. Number of minutes spent meditating: This indicates the total amount of time the user has spent meditating.
2. Number of meditation sessions: This shows the number of times the user has meditated.
3. Points earned: This represents the user's progress and achievements in their meditation practice.
4. Bar chart: This visual representation of the data helps the user understand the data better by showing the distribution of the meditation minutes, sessions, and points.

**(E)**

ఈ పరిశ్రమ నేపథ్యంలో పని చేస్తున్న కార్మికులందరికీ, భద్రతా ప్రమాదాలు, ఐపా మొదలైనవాటిని సూచించండి.

ఈ పరిశ్రమ నేపథ్యంలో, కార్మికులు భద్రతా ప్రమాదాలు మరియు ప్రమాదాలను ఎదుర్కొంటారు, ముఖ్యంగా నిర్మాణం మరియు మరమ్మతు పనుల సమయంలో. చిత్రంలో, నలుగురు కార్మికులు నిర్మాణ ప్రాంతంలో పని చేస్తున్నారు, ఒక వ్యక్తి పడిపోయిన హార్డ్ టోపీని తీసుకుంటాడు. ఈ ప్రమాదాలలో పడిపోవడం, వస్తువులు లేదా పరికరాలు పడిపోవడం, భారీ వస్తువులను ఎత్తడం లేదా మోసుకెళ్ళడం మరియు ప్రమాదకరమైన పదార్థాలు లేదా యంత్రాలతో పనిచేయడం వంటివి ఉండవచ్చు. ఈ ప్రమాదాలను తగ్గించడానికి, కార్మికులు హెల్మెట్లు, చేతి తొడుగులు, భద్రతా బూట్లు మరియు భద్రతా అద్దాలు వంటి తగిన భద్రతా పరికరాలను ధరించాలి. అదనంగా, వారు తమ పరిసరాల గురించి తెలుసుకోవాలి, భద్రతా మార్గదర్శకాలను అనుసరించాలి మరియు ప్రమాదాలు లేదా గాయాలను నివారించడానికి వారి సహోద్యోగులతో సమర్థవంతంగా కమ్యూని

**(F)**

This is a CC camera feed. Pls give following values.

1. Accident /non accident
2. Severity
3. Vehicles involved

1. Accident
2. Severity: High
3. Vehicles involved: Cars, motorcycles, and a person on a scooter

Figure 4: **Multilingual VLM Capabilities.** Our model demonstrates robust performance across various languages in: a) Creative writing, b) Fine-grained attribute extraction, c) Explaining scientific diagrams, d) Screen reading/OCR, e) Anomaly and hazard detection, and f) Real-time accident and incident monitoring.
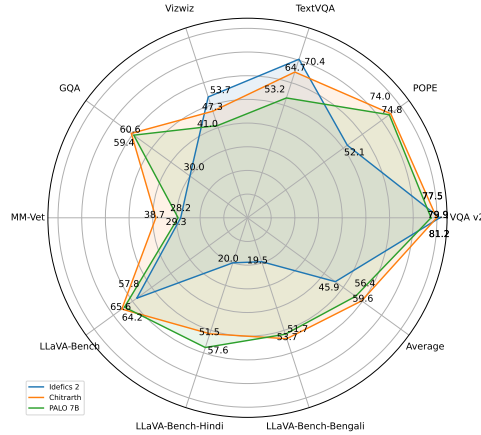
Figure 5: **Performance against SOTA VLMs on different academic multimodal tasks.** Our model consistenly outperforms IDEFICS 2 (7B) and PALO 7B on different benchmarks while remaining competitive on TextVQA and Vizwiz.
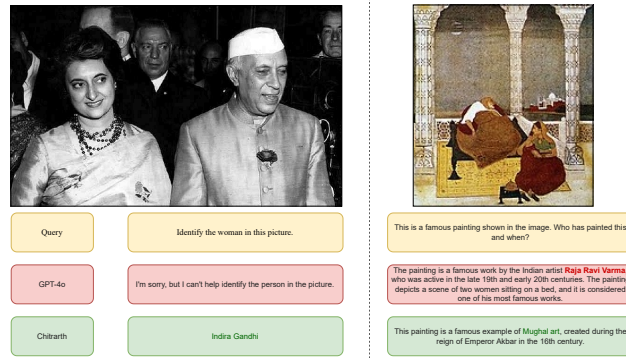


Figure 6: **Performance on images with Indian context.** Chitrarth is able to better understand the images of Image context such as the prominent lady figure (Late Indian Prime Minister Indira Gandhi) in left as well as historical artwork compared to generic and incorrect responses from GPT-4o.

Vizwiz (see radar graph in Figure 5). On the LLaVA-Bench (Bengali) our model outperforms the multi-lingual baseline PALO and achieves SOTA results of 53.7 points. Table 1 presents results on BharatBench across various languages, demonstrating that ours is the only model capable of handling all included languages, establishing baseline results for future research. Figure 4 showcases selected outputs from our top-performing Multimodal LLM across various languages. The model excels in tasks such as creative writing, fine-grained attribute extraction, explaining scientific diagrams, and screen reading/OCR, while also demonstrates strong capabilities in anomaly and hazard detection, as well as real-time accident and incident monitoring. In our manual qualitative evaluation, we observe that our model is able to better understand the images of Indian context such as the prominent lady figure in Figure 6, compared to generic and incorrect responses from GPT-4o. This could be attributed to the inclusion of high quality culturally rich images in Stage 2. A further quantitative analysis around this would be interesting but out of scope of this work.

We conducted an ablation study evaluating various vision encoders and found that SigLIP-SO400M consistently outperforms CLIP ViT-L/14@336px across all English benchmarks, achieving faster convergence (see Figure 7 ). Notably, SigLIP-SO400M yields improvements of 11 points on TextVQA and 13 points on LLaVA-Bench compared to CLIP ViT-L/14@336px. Figure 8 explores the impact of multilingual training data on the English academic benchmarks. We compare our model's performance when trained with only English, bilingual, and multilingual data across both
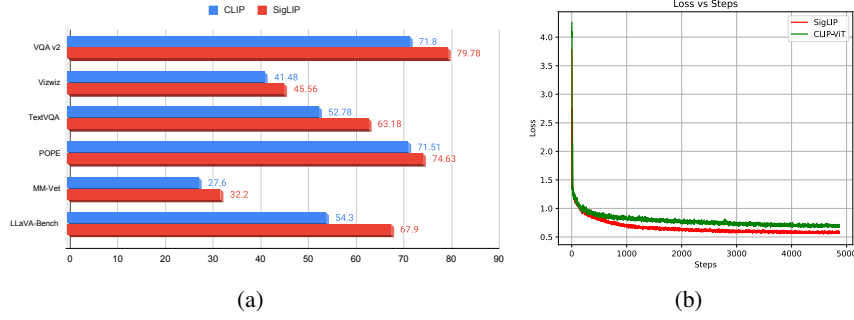
Figure 7: **Ablation on visual encoder choice.** a) SigLIP as the vision encoder consistently performs better than CLIP in the same training regime. b) SigLIP based model also achieve faster convergence as depicted in Stage 1 loss curve. Stage 2 follows a similar pattern.
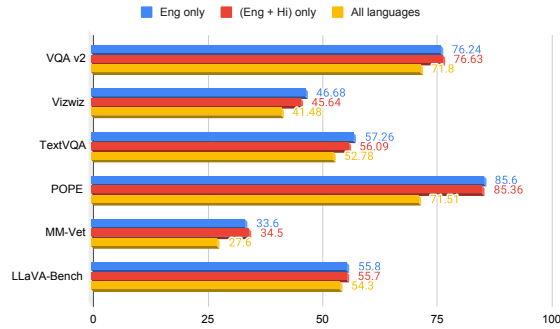


Figure 8: **Impact of Multi-lingual training Data.** Expanding the number of languages in the training data enhances multilingual capabilities but results in decreased scores on academic English datasets.

stages. Consistent with the findings of [Scao et al., 2022], expanding the range of languages in the training data improves multilingual capabilities but leads to decreased performance on academic English datasets. This underscores a key challenge in balancing cross-lingual performance.

## 7    Conclusion

This paper presents Chitrarth, a multilingual multimodal LLM that is able to have image grounded conversations in English as well as across multiple Indian languages. Our model encodes images using a pre-trained vision encoder [Dosovitskiy et al., 2020] and autoregressively generates response using a pre-trained multi-lingual LLM. Empirically, our model outperforms previous baselines for different multimodal tasks. As part of this work, we also introduce BharatBench, a multimodal evaluation framework and provide benchmark results for low resource languages. We anticipate that our research will significantly contribute to the advancement of VLMs for Indian languages, thereby providing substantial benefits to a population exceeding one billion people.

**Limitations and Future Work:** We use an automated translation pipeline for creating multi-lingual training data which may introduce biases from large language models (LLMs), potentially leading to misrepresentations of cultural symbols and gestures, impacting content accuracy. Addressing these biases requires additional evaluation and targeted training, which we plan to address in the future work. Building on our promising results across 10 low-resource languages, we plan to broaden the language scope in the future research to enhance linguistic diversity and inclusivity in our Vision-Language Models (VLMs). In our current training pipeline, we keep the vision encoder frozen throughout both training stages. However, recent research [Laurençon et al., 2024b, Tong et al., 2024] suggests that unfreezing the vision encoder could enhance representation learning. We plan to investigate this approach in future work with higher resolution vision encoders, along with expanding our model's ability to interpret multiple images within a conversational context.

## Acknowledgements

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

Kabir Ahuja, Harshita Diddee, Rishav Hada, Millicent Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, et al. Mega: Multilingual evaluation of generative ai. *arXiv preprint arXiv:2303.12528*, 2023.

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.

Abhinand Balachandran. Tamil-llama: A new tamil language model based on llama 2, 2023.

Abhijit Bendale, Michael Sapienza, Steven Ripplinger, Simon Gibbs, Jaewon Lee, and Pranav Mistry. Sutra: Scalable multilingual language model architecture. *arXiv preprint arXiv:2405.06694*, 2024.

Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.

Christian Buck, Kenneth Heafield, and Bas Van Ooyen. N-gram counts and language models from the common crawl. In *LREC*, volume 2, page 4, 2014.

Team Chameleon. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.

Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023.

Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022.

Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *arXiv preprint arXiv:2404.16821*, 2024.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6, 2023.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ArXiv*, abs/2010.11929, 2020. URL `https://api.semanticscholar.org/CorpusID:225039882`.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. MME: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023.

Jay Gala, Pranjal A Chitale, A K Raghavan, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar M, Janki Atul Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, Pratyush Kumar, Mitesh M Khapra, Raj Dabre, and Anoop Kunchukuttan. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL `https://openreview.net/forum?id=vfT4YuzAYA`.

Jay Gala, Thanmay Jayakumar, Jaavid Aktar Husain, Aswanth Kumar M, Mohammed Safi Ur Rahman Khan, Diptesh Kanojia, Ratish Puduppully, Mitesh M. Khapra, Raj Dabre, Rudra Murthy, and Anoop Kunchukuttan. Airavata: Introducing hindi instruction-tuned llm, 2024.

Team Gemini, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

Team Gemma, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017.

Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3608–3617, 2018.

Drew A Hudson and Christopher D Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019.

Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR, 2021.

Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

Aditya Kallappa, Palash Kamble, Abhinav Ravi, Akshat Patidar, Vinayak Dhruv, Deepak Kumar, Raghav Awasthi, Arveti Manjunath, Shubham Agarwal, Kumar Ashish, Gautam Bhargava, Chandra Khatri, and Krutrim Team. Krutrim LLM: Multilingual foundational model for over a billion people. *arXiv preprint arXiv:2502.09642*, 2024.

Guneet Singh Kohli, Shantipriya Parida, Sambit Sekhar, Samirit Saha, Nipun B Nair, Parul Agarwal, Sonal Khosla, Kusumlata Patiyal, and Debasish Dhal. Building a llama2-finetuned llm for odia language utilizing domain knowledge instruction set, 2023.

Telugu Labs. Navarsa: Indic llms based on gemmma, 2023.

Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better understanding vision-language models: insights and future directions. *arXiv preprint arXiv:2408.12637*, 2024a.

Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024b.

Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2023.

Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.

Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Yaofeng Sun, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024a.

Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26439–26455, 2024b.

Muhammad Maaz, Hanoona Rasheed, Abdelrahman Shaker, Salman Khan, Hisham Cholakal, Rao M Anwer, Tim Baldwin, Michael Felsberg, and Fahad S Khan. Palo: A polyglot large multimodal model for 5b people. *arXiv preprint arXiv:2402.14818*, 2024.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*, 2023.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

Juan A Rodriguez, Shubham Agarwal, Issam H Laradji, Pau Rodriguez, David Vazquez, Christopher Pal, and Marco Pedersoli. Starvector: Generating scalable vector graphics code from images. *arXiv preprint arXiv:2312.11556*, 2023.

Teven Le Scao, Thomas Wang, Daniel Hesslow, Lucile Saulnier, Stas Bekman, M Saiful Bari, Stella Biderman, Hady Elsahar, Niklas Muennighoff, Jason Phang, et al. What language model to train if you have one million gpu hours? *arXiv preprint arXiv:2210.15424*, 2022.

Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326, 2019.

Harman Singh, Nitish Gupta, Shikhar Bharadwaj, Dinesh Tewari, and Partha Talukdar. Indicgenbench: A multilingual benchmark to evaluate generation capabilities of llms on indic languages. *arXiv preprint arXiv:2404.16816*, 2024.

Shengbang Tong, Ellis Brown, Penghao Wu, Sanghyun Woo, Manoj Middepogu, Sai Charitha Akula, Jihan Yang, Shusheng Yang, Adithya Iyer, Xichen Pan, et al. Cambrian-1: A fully open, vision-centric exploration of multimodal llms. *arXiv preprint arXiv:2406.16860*, 2024.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.

Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural Information Processing Systems*, 2017. URL `https://api.semanticscholar.org/CorpusID:13756489`.

Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface's transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.

Bin Xiao, Haiping Wu, Weijian Xu, Xiyang Dai, Houdong Hu, Yumao Lu, Michael Zeng, Ce Liu, and Lu Yuan. Florence-2: Advancing a unified representation for a variety of vision tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4818–4829, 2024.

Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, et al. xgen-mm (blip-3): A family of open large multimodal models. *arXiv preprint arXiv:2408.08872*, 2024.

Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. MM-VET: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.

Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. URL `https://arxiv.org/abs/2303.15343`.