

Idefics2: A Powerful 8B Vision-Language Model





These are two cats
lying on a pink couch.

Answer



Multimodal LLM



Image

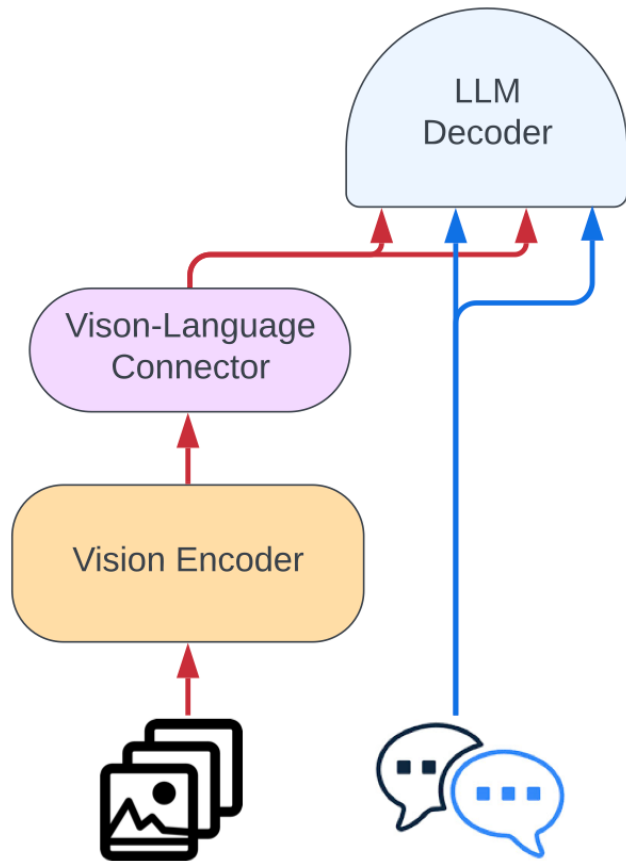


Question



What are these?

Idefics2



The model is built on top of two pre-trained models:

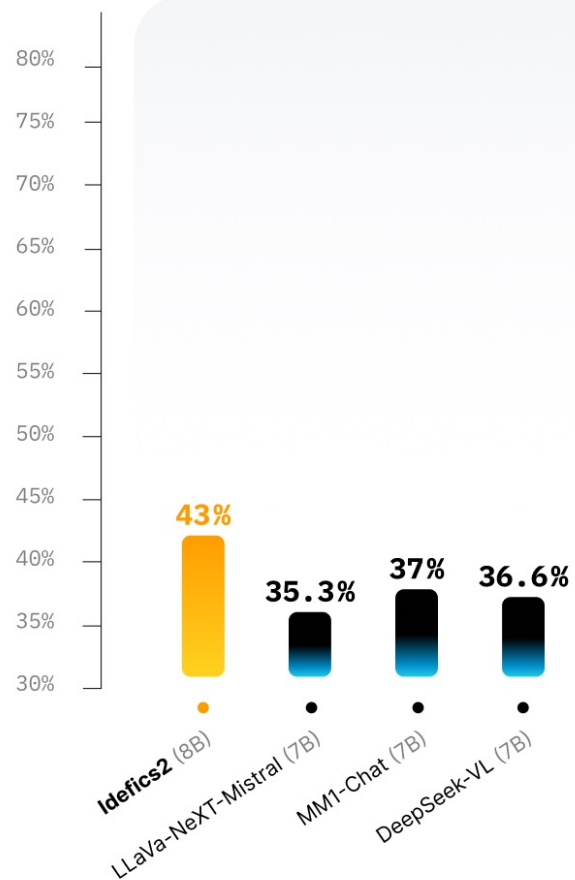
- [Mistral-7B-v0.1](#) and
- [siglip-so400m-patch14-384](#)

<https://huggingface.co/blog/idefics2>

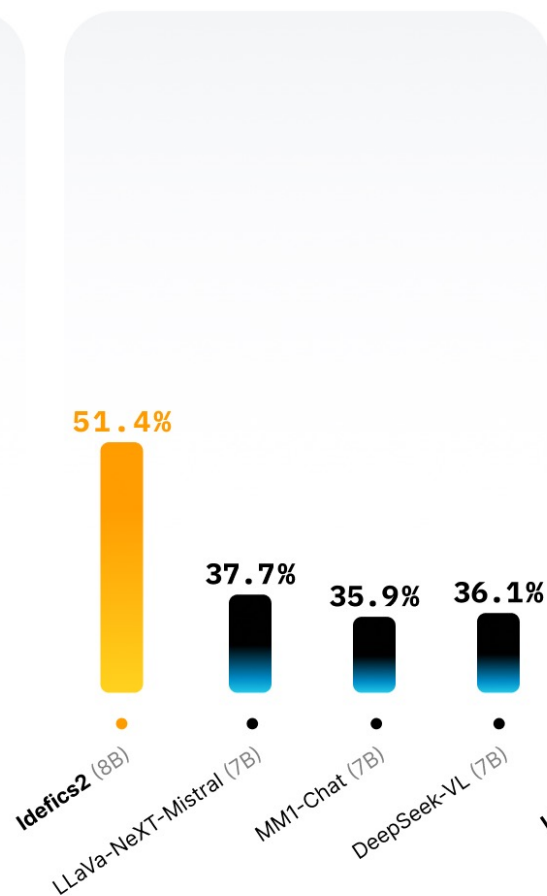
Idefics2

- 8B parameters
- Open license (Apache 2.0)
- Enhanced OCR (Optical Character Recognition) capabilities
- Achieves top-of-class performance on Visual Question Answering benchmarks, competing with much larger models.
 - [LLava-Next-34B](#)
 - [MM1-30B-chat](#)

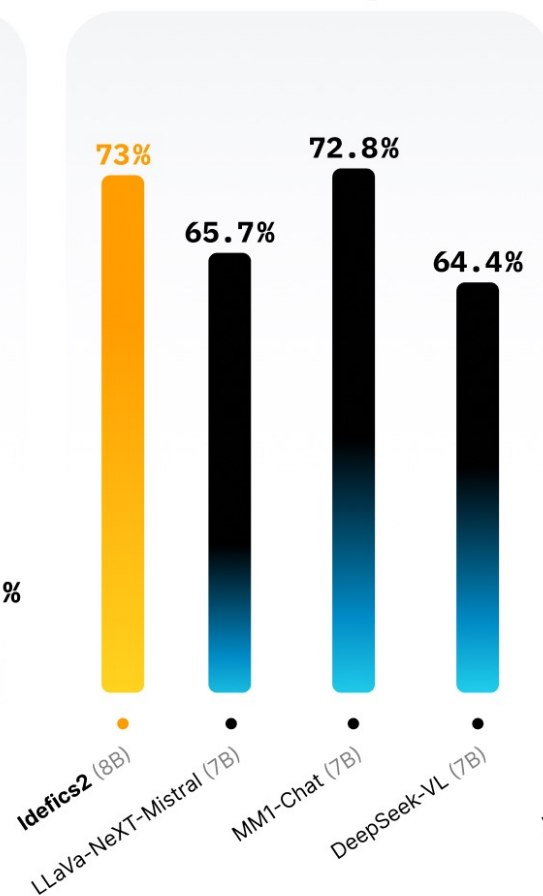
MMMU



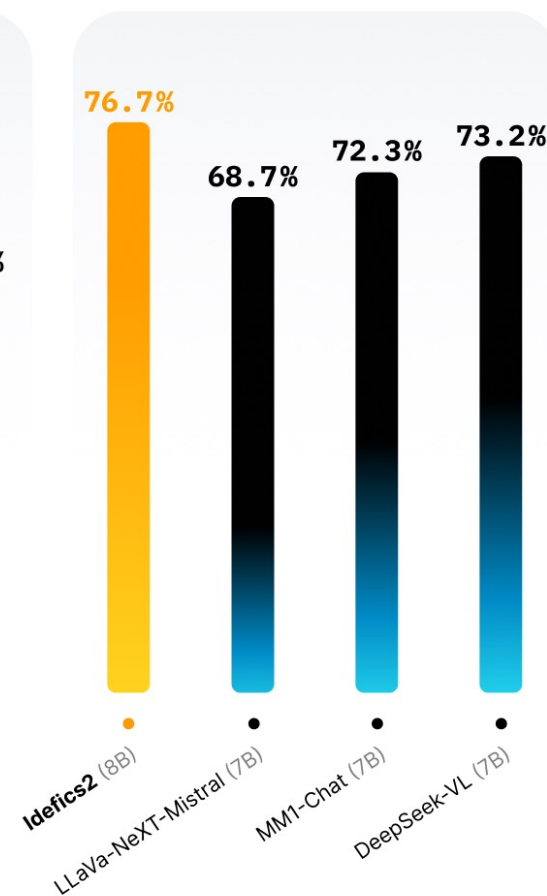
Mathvista



Textvqa



MMbench



Training Data

- Interleaved web documents (Wikipedia, [OBELICS](#))
- Image-caption pairs (Public Multimodal Dataset, LAION-COCO)
- OCR data ([PDFFA \(en\)](#), [IDL](#) and [Rendered-text](#), and image-to-code data ([WebSight](#)))

Colab

- [https://github.com/NielsRogge/Transformers-Tutorials/blob/master/Idefics2/Fine_tune_Idefics2_for_JSON_extraction_use_cases_\(PyTorch_Lightning\).ipynb](https://github.com/NielsRogge/Transformers-Tutorials/blob/master/Idefics2/Fine_tune_Idefics2_for_JSON_extraction_use_cases_(PyTorch_Lightning).ipynb)
- <https://colab.research.google.com/drive/1NtcTgRbSBKN7pYD3Vdx1j9m8pt3fhFDB?usp=sharing>