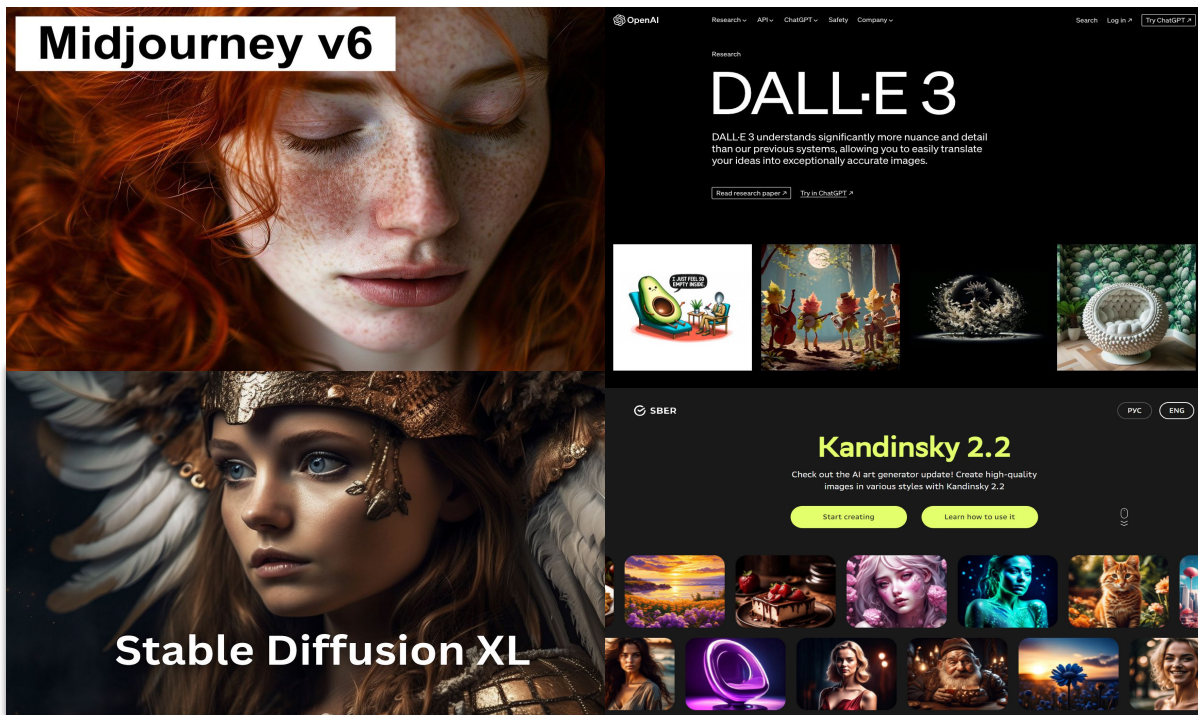# Annotation d'images

## Mounir Bendali-Braham

# L'objectif de l'annotation d'images



c.f. slide "Références" pour les sources de ces images

# Qu'est-ce que l'annotation d'images ?



Captioning Model

A happy dog is standing in the ocean

# Comparaison de quelques modèles d'annotation

Picture taken from FuseCap's article

# Vue d'ensemble des modèles d'annotation: FuseCap

From FuseCap article "how a good a bad captionner describe an image"



**Original:** Two men with eye glasses looking at something

**Ours:** Two bespectacled men, one with black glasses and a black and brown beard, the other with silver glasses and short brown hair, sit together with an open blue laptop on a table in front of them. A gray cat lounges nearby
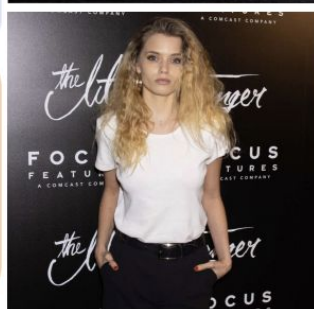
**Original:** Mhmm, some clouds in the sky

**Ours:** A woman wearing dark sunglasses stands next to a red car with a black license plate reading 166882, PRI. The car has off and round headlights, a chrome and silver bumper, a black tire, and a red door. The cloudy and white sky is visible in the background.

**Original:** save yourself the expense of a professional arrangement .

**Ours:** Floral Arrangement: A colorful assortment of sunflowers, yellow, white, orange, and purple flowers, and green leaves arranged on a black and wood table.

**Original:** <PERSON> 2018 : <PERSON>: The Little Stranger Premiere —01

**Ours:** A woman with blond, long hair wearing a black belt and pants attends the premiere of The Little Stranger in 2018.

# Vue d'ensemble des modèles d'annotation: MiniGPT-4



https://minigpt-4.github.io/

# Vue d'ensemble des modèles d'annotation: Qwen-VL (1)

# Vue d'ensemble des modèles d'annotation: Qwen-VL (2)



https://arxiv.org/pdf/2308.12966.pdf

# Jeux de données d'annotation d'images

A computer screen with a Windows message about Microsoft license terms.

A can of green beans is sitting on a counter in a kitchen.

A photo taken from a residential street in front of some homes with a stormy sky above.

A blue sky with fluffy clouds, taken from a car while driving on the highway.

A hand holds up a can of Coors Light in front of an outdoor scene with a dog on a porch.

A digital thermometer resting on a wooden table, showing 38.5 degrees Celsius.

A Winnie The Pooh character high chair with a can of Yoohoo sitting on it in front of a white wall.
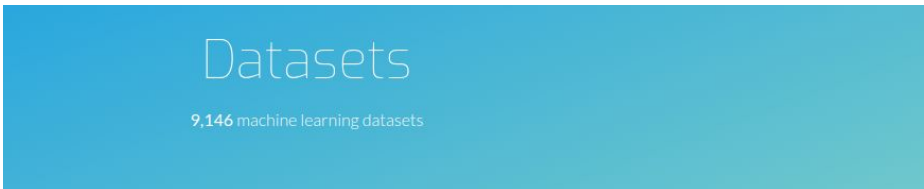
A cup holder in a car holding loose change from Canada.

https://vizwiz.org/tasks-and-datasets/image-captioning/

# D'autres jeux de données en annotation d'images

https://paperswithcode.com/datasets?task=image-captioning



Datasets

9,146 machine learning datasets

⚠️ Share your dataset with the ML community!

63 dataset results for  Image Captioning ✕

**MS COCO** (Microsoft Common Objects in Context)
The MS COCO (Microsoft Common Objects in Context) dataset is a large-scale object detection, segmentation, key-point detection, and captioning dataset. The dataset consists of…
9,475 PAPERS • 88 BENCHMARKS

**Flickr30k**
The Flickr30k dataset contains 31,000 images collected from Flickr, together with 5 reference sentences provided by human annotators.
676 PAPERS • 9 BENCHMARKS

**Conceptual Captions**
Automatic image captioning is the task of producing a natural-language utterance (usually a sentence) that correctly reflects the visual content of an image. Up to this point, the resour…
297 PAPERS • 2 BENCHMARKS

# Captioning with ClipInterrogator

https://colab.research.google.com/github/pharmapsychotic/clip-interrogator/blob/main/clip_interrogator.ipynb

# Captioning with MiniGPT-4

https://github.com/camenduru/MiniGPT-4-colab

# Captioning with Qwen-VL

https://github.com/camenduru/Qwen-VL-Chat-colab/tree/main

# Références

- Image midjourney https://mid-journey.ai/midjourney-v6-release/
- Image Stable Diffusion XL
  https://generativeai.pub/stable-diffusion-xl-is-here-whats-new-4e6ed27df70c
- Image Dall-E 3 (capture du site web https://openai.com/dall-e-3)
- Image Kandinsky v2.2
  https://www.gadgetvoize.com/2023/07/12/sber-presents-its-neural-networks-new-version-kandinsky-2-2/