





Technologies d'assistance pour les personnes malvoyantes basées sur la vision : avancées, limites et perspectives

Aela Le Sommer, Panagiotis Papadakis, Christophe Lohr

Défis : Les personnes malvoyantes peinent à se déplacer et à comprendre des environnements inconnus.





Défis : Les personnes malvoyantes peinent à se déplacer et à comprendre des environnements inconnus.



Intérêt : Solutions basées sur des modèles IA



Défis : Les personnes malvoyantes peinent à se déplacer et à comprendre des environnements inconnus.



Intérêt : Solutions basées sur des modèles IA

Langage



Défis : Les personnes malvoyantes peinent à se déplacer et à comprendre des environnements inconnus.



Intérêt : Solutions basées sur des modèles IA Langage

Image captioning, video summarization et Visual-Question Answering (VQA) → solutions prometteuses pour aider les personnes malvoyantes dans leur vie quotidienne.

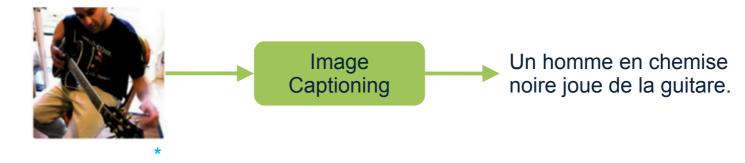
Vision

1. État de l'art

État de l'art

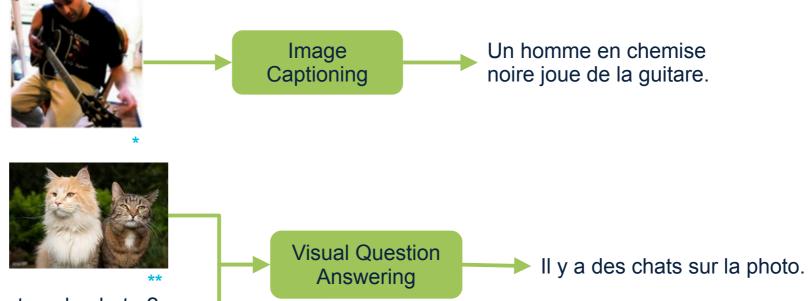
Présentation des tâches

Image Captioning:



Présentation des tâches

Image Captioning:



VQA:

Quels animaux sont sur la photo?

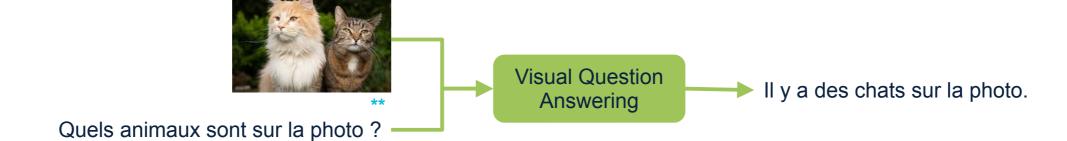
État de l'art

Présentation des tâches

Image Captioning:



VQA:



Video Summarization:





Source de l'image * : Jeux de données COCO

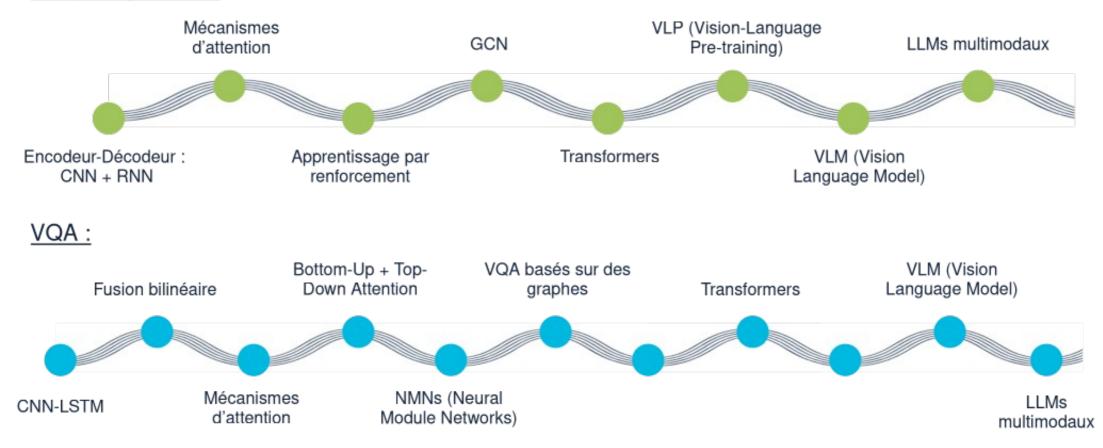
Source de l'image ** : https://cdn.britannica.com/34/235834-050-C5843610/two-different-breeds-of-cats-side-by-side-outdoors-in-the-garden.jpg

Source de l'image *** : Tiwari, V., and Bhatnagar, C. "A survey of recent work on video summarization: approaches and techniques." Multimedia Tools and Applications (2021)

État de l'art

Évolution des architectures des modèles

Image Captioning:





État de l'art

Solutions basées sur la vision disponible sur le marché pour les personnes malvoyantes

- Lunettes intelligentes pour les personnes malvoyantes : Envision Glasses, OrCam MyEye, etc.
- Exemples d'applications mobiles d'assistance basées sur la vision pour les personnes malvoyantes :

App. Mobile	Reconnaissance	Lecture	Description	VQA	Reconnaissance	Nb
	d'objets	de textes	de scène	33.11.03	faciale	Téléchargements
Be My Eyes (Be My AI)			✓	✓		1M+*
Lookout	✓	✓	✓	✓		500k+*
TapTapSee	✓	✓				500k+*
Envision	✓	✓	✓	✓	✓	100k+*
Seeing AI	100	✓	✓	✓	✓	100k+*
Supersense	✓	✓				50k+*
Aipoly	✓	✓				-
Oorion	✓	✓				-
VizWiz				✓		-
VoiceVision			✓			-

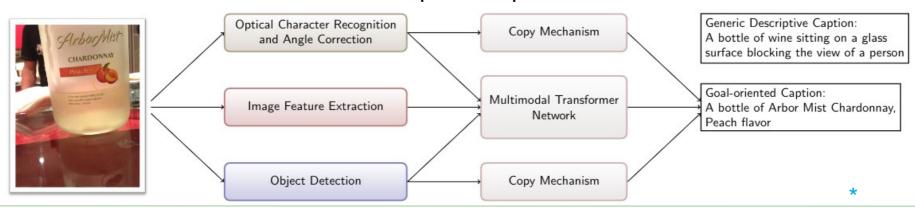
^{*} sur Google Play en février 2025



État de l'art

Approches principales

- (Image Captioning, Video Summarization, VQA) + Text-to-Speech
 - ▶ ⚠ Limites liées à leur structure en plusieurs étapes :
 - Pertes d'informations,
 - Délai de traitement plus important.
 - > Cas d'application pour les malvoyants :
 - Ajouts de modules complémentaires (OCR, détection d'objets, correction d'angle, etc.)
 - · BUT : Améliorer la fiabilité des descriptions/réponses.





État de l'art

Approches principales

- (Image Captioning, Video Summarization, VQA) + Text-to-Speech
 - - Pertes d'informations,
 - Délai de traitement plus important.
 - > Cas d'application pour les malvoyants :
 - Ajouts de modules complémentaires (OCR, détection d'objets, correction d'angle, etc.)
 - · BUT : Améliorer la fiabilité des descriptions/réponses.
- ✓ Image-to-Speech end-to-end → convertit directement l'image en parole
 - **▶** <u>∧</u> Limites :
 - Complexité d'évaluation,
 - → Entraînement nécessitant de grandes quantités de paires image/audio → difficile et coûteux à obtenir.



2. Limites

Limites

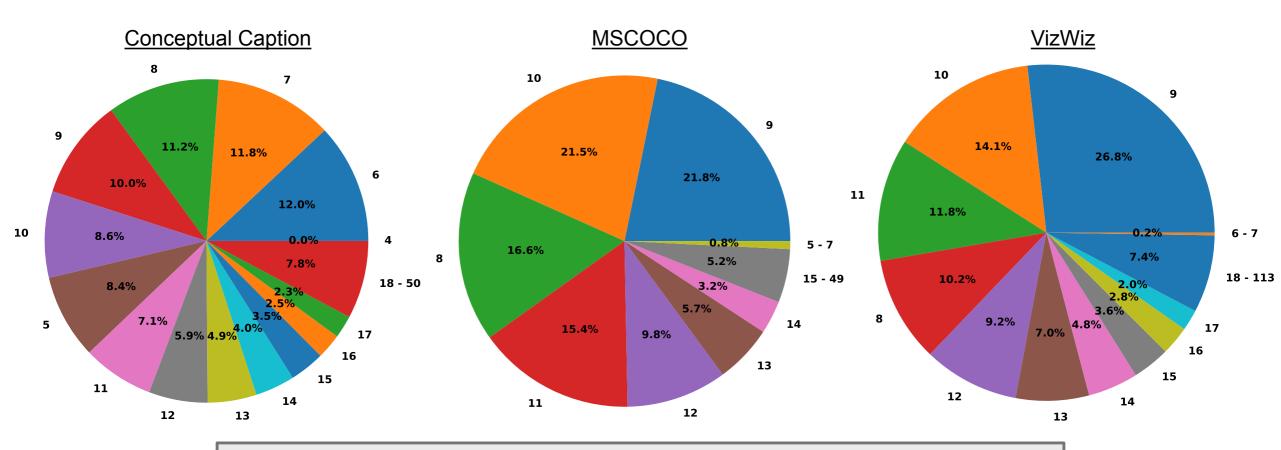
Principaux jeux de données d'image Captioning et de Visual-Question Answering

Jeux de	Tâches	Source des	Nb	Nb des-	Apports potentiels pour des modèles destinés aux malvoyants
données		images	images	cript°/quest°	
VizWiz [11]	IC,	Photo prises par	+ 39k /	+ 195k / +32k	Questions et images réelles issues d'utilisateurs malvoyants
	VQA	des malvoyants	+32k		
MSCOCO	IC	Flickr	+123k	+ 616k	Grandes diversités d'environnement (intérieur, extérieur, objets courant),
[11]					Diversité des annotations
Flickr30k [11]	IC	Flickr	+ 31k	+ 158k	Images d'activités, d'événements et de scènes de la vie quotidienne
Conceptual	IC	Pages Web	+	+ 3.3M	Très grand jeux de données, Grande variété de types d'images (ex :
Captions [11]			3.3M		images naturelles, de produits, dessins animés, dessins, etc.)
NoCaps [11]	IC	OpenImages	+ 15k	+ 166k	Conçu pour tester la généralisation sur de nouveaux objers, ce qui peut
		(validation + test)			être utile pour décrire des scènes inconnues
Visual	IC,	YFCC100M ∩	+ 108k	5.4M	Très dense au niveau des annotations : descriptions de régions, objets,
Genome	VQA	MS-COCO		descriptions de	attributs, relations, graphes de régions, graphes de scènes et VQA
[11, 1]				régions / 1.7M	
GQA [1]	VQA	Visual Genome	+ 113k	+ 22M	Exploite les représentations sémantiques des questions et les graphes de
					scènes pour répondre à la question, permettant un raisonnement structuré
VQAv2 [1]	VQA	COCO	+ 204k	+ 1.1M	Biais réduit : Chaque question est associée à une paire d'images
					similaires qui donnent lieu à deux réponses différente
OK-VQA [1]	VQA	COCO	+ 14k	+ 14k	Images du monde réel nécessitant des connaissances externes pour
					répondre aux questions
TDIUC [17]	VQA	MS-COCO +	+ 167k	+ 1.6M	Grande diversité de type de questions (12), comprenant des questions
		Visual Genome			absurdes pour forcer le système à raisonner sur le contenu de l'image.



Limites

Nombre de mots dans les descriptions



+90 % des descriptions < 20 mots → descriptions générales, peu détaillées

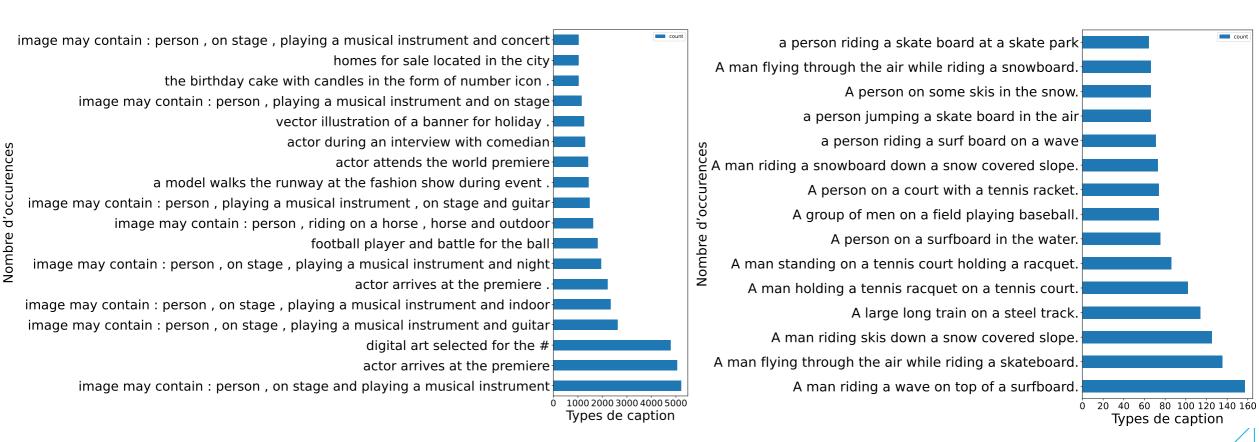


Limites

Descriptions les plus courantes dans les jeux de données (1/2)

Conceptual Caption

MSCOCO





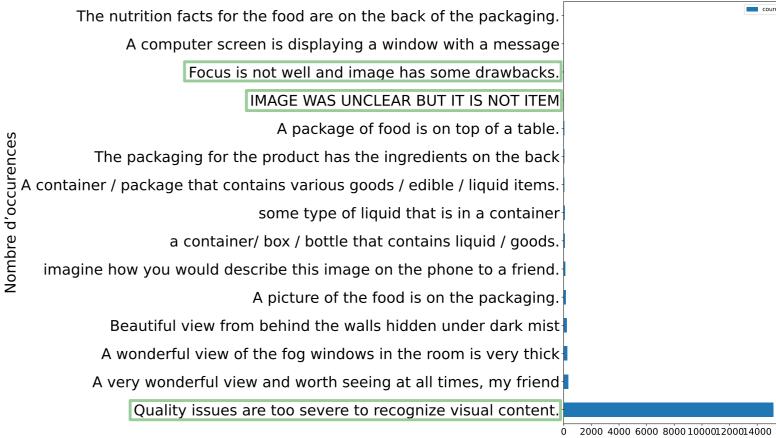
Limites

Descriptions les plus courantes dans les jeux de données (2/2)

VizWiz

Prend en compte les problèmes de qualité des photos ...

... MAIS descriptions qui restent vagues.

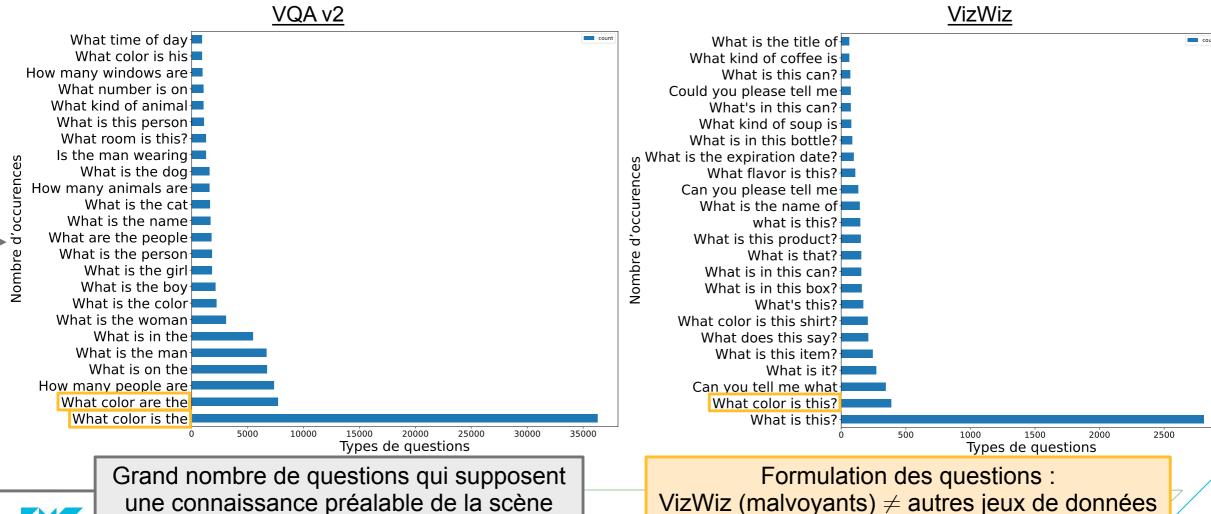


Types de caption



Limites

Types de questions les plus posées dans les jeux de données

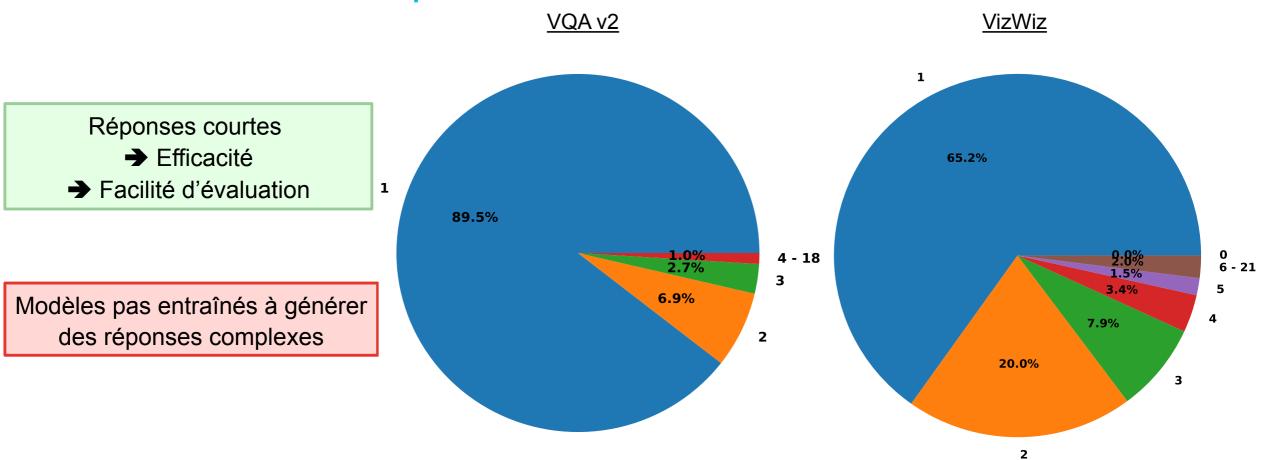


IMT Atlantique Bretagne-Pays de la Loir Ezgla Mines-Télécon

19

Limites

Nombre de mots dans les réponses





Limites

Répartition des réponses dans les jeux de données

VQA v2 <u>VizWiz</u>







Limites

Principales problématiques



Limites

Principales problématiques

Mauvaise qualité des photos prises par les utilisateurs malvoyants,





Photos provenant du jeu de données VizWiz prises par des malvoyants



Limites

Principales problématiques

- Mauvaise qualité des photos prises par les utilisateurs malvoyants,
- Complexité d'évaluation des modèles en raison de leur nature générative,





Photos provenant du jeu de données VizWiz prises par des malvoyants



Limites

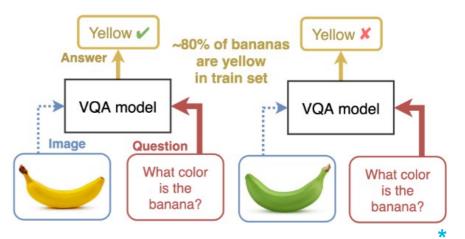
Principales problématiques

- Mauvaise qualité des photos prises par les utilisateurs malvoyants,
- Complexité d'évaluation des modèles en raison de leur nature générative,
- Biais dans les modèles,





Photos provenant du jeu de données VizWiz prises par des malvoyants



Exemple de biais hérités par les modèles

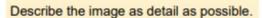


Limites

Principales problématiques

- Mauvaise qualité des photos prises par les utilisateurs malvoyants,
- Complexité d'évaluation des modèles en raison de leur nature générative,
- Biais dans les modèles,
- Hallucinations des modèles,







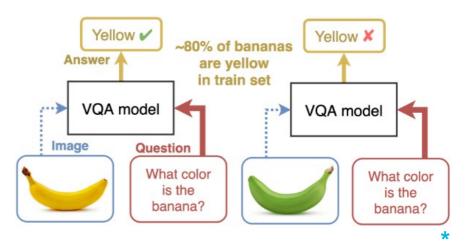
The picture shows a long-haired man in a suit sitting on the steps in the city, checking his watch. Beside him are two green cups, a laptop, and some scattered documents, with a bicycle parked in front of him. A small dog on the step is cuirously observing him.

Exemple d'hallucinations (objets, attributs, relations)





Photos provenant du jeu de données VizWiz prises par des malvoyants



Exemple de biais hérités par les modèles



3. Pistes d'exploration

Pistes d'exploration

Poser le problème comme un problème d'intelligence incarnée (embodied intelligence)

BUT:

- Pallier les limites des benchmarks existants,
- Contrôle total sur la scène,
- Facilité de génération de scénarios variés adaptés au cas d'application,
- · Annotations automatiques, précises et exhaustives (segmentation, coordonnées, matériaux, etc.)
- Contourner les problèmes de confidentialité,



Pistes d'exploration

Poser le problème comme un problème d'intelligence incarnée (embodied intelligence)

BUT:

- Pallier les limites des benchmarks existants,
- Contrôle total sur la scène,
- Facilité de génération de scénarios variés adaptés au cas d'application,
- · Annotations automatiques, précises et exhaustives (segmentation, coordonnées, matériaux, etc.)
- Contourner les problèmes de confidentialité,
- Distinction des niveaux de gravité des hallucinations



Pistes d'exploration

Poser le problème comme un problème d'intelligence incarnée (embodied intelligence)

BUT:

- Pallier les limites des benchmarks existants,
- Contrôle total sur la scène,
- Facilité de génération de scénarios variés adaptés au cas d'application,
- · Annotations automatiques, précises et exhaustives (segmentation, coordonnées, matériaux, etc.)
- Contourner les problèmes de confidentialité,
- Distinction des niveaux de gravité des hallucinations
- Intégrer l'utilité des réponses dans l'évaluation



Pistes d'exploration

- Évaluation de la qualité et sélection des images :
 - Actuellement:
 - Évaluation de la qualité de l'image avant Image Captioning/VQA.
 - · ☑ Bonne qualité → génération d'une réponse/description
 - Mauvaise qualité → Notification à l'utilisateur → nouvelle photo demandée



Pistes d'exploration

- Évaluation de la qualité et sélection des images :
 - Actuellement:
 - Évaluation de la qualité de l'image avant Image Captioning/VQA.
 - · ☑ Bonne qualité → génération d'une réponse/description
 - Mauvaise qualité → Notification à l'utilisateur → nouvelle photo demandée
 - Procédé long et contraignant au quotidien.
 - Piste d'amélioration :
 - Capture de l'environnement (série de photos) → filtrage des images de bonne qualité → Image
 Captioning/VQA sur la série
 - · BUT:
 - · Réduire le risque de réponses/descriptions peu fiables liées à la mauvaise qualité de l'image,
 - Limiter le nombre de reprises de photos.
 - · <u>A</u> Limites : Étapes de calculs supplémentaires







MERCI

Aela Le Sommer

Lunettes intelligentes pour les personnes malvoyantes

Envision Glasses:

- ✓ Instant Text,
- ✓ Scan Text,
- ✓ Batch Scan,
- Call a Companion,
- ✓ Call Aira,
- ✓ Describe Scene,
- ✓ Detect Light,
- ✓ Recognise Cash,
- ✓ Detect Colors,
- ✓ Find People,
- Find Objects,
- Teach a Face and Explore.



OrCam MyEye 3 Pro:

- ✓ Al Assistant-Just Ask
- ✓ Summarize text into main topics
- Utilizes intext and outsourced data
- Reads handwriting
- ✓ Zoom in&out
- ✓ Change contrast
- ✓ Extract and copy text
- Converts image text to digital text
- Connect to any screen

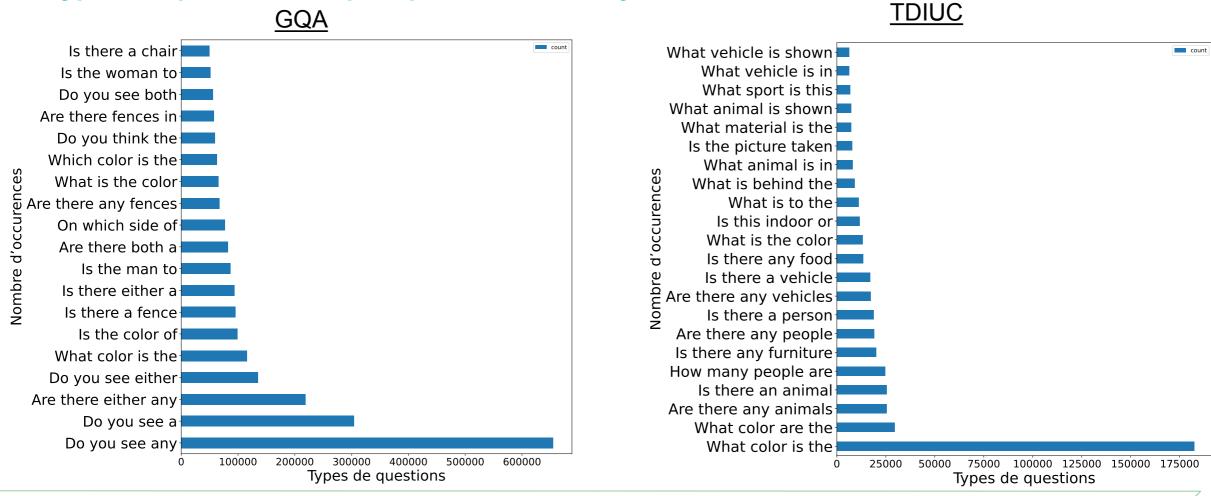
- ✓ Supports over 140 languages
- ✓ Voice commands
- ✓ Text reading
- Smart reading
- ✓ Recognize Faces
- ✓ Identifying Products
- ✓ Money Notes
- ✓ Barcodes
- ✓ Colors





Autres jeux de données : VQA

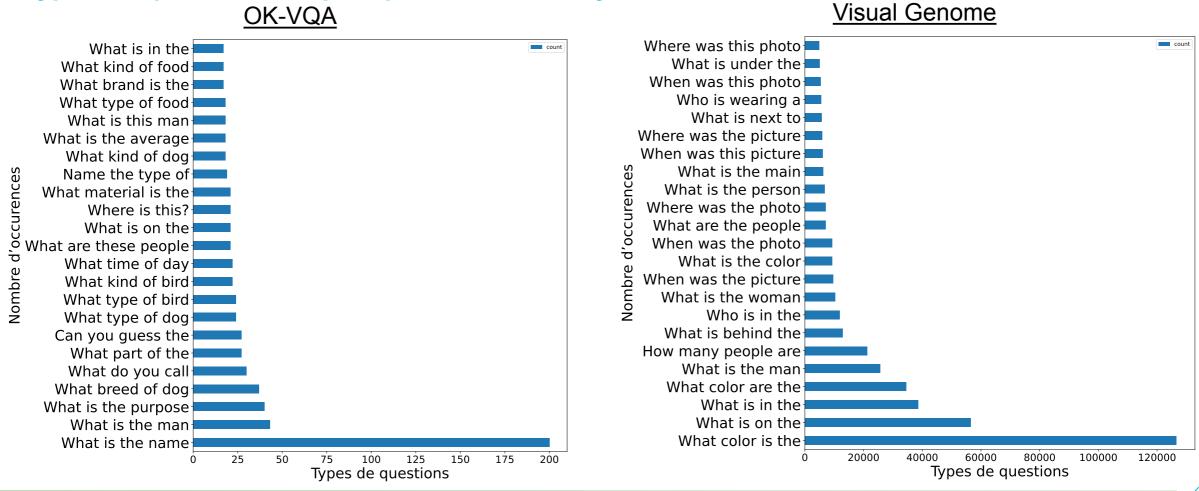
Types de questions les plus posées dans les jeux de données





Autres jeux de données : VQA

Types de questions les plus posées dans les jeux de données

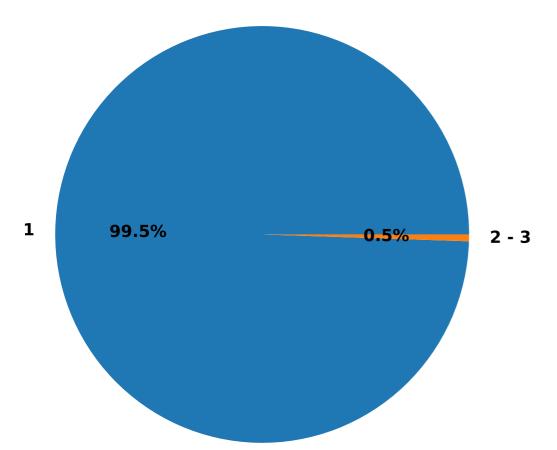




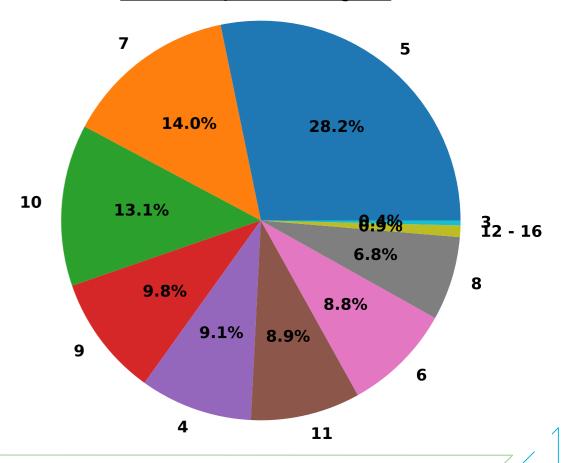
Autres jeux de données : VQA

Nombre de mots dans les réponses

GQA: Réponses courtes



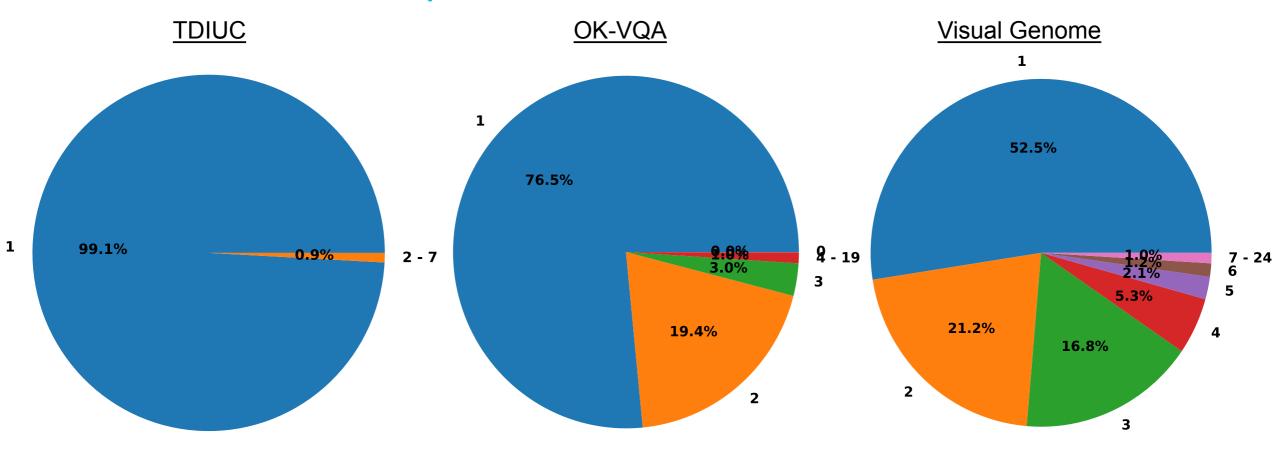
GQA: Réponses longues





Autres jeux de données : VQA

Nombre de mots dans les réponses

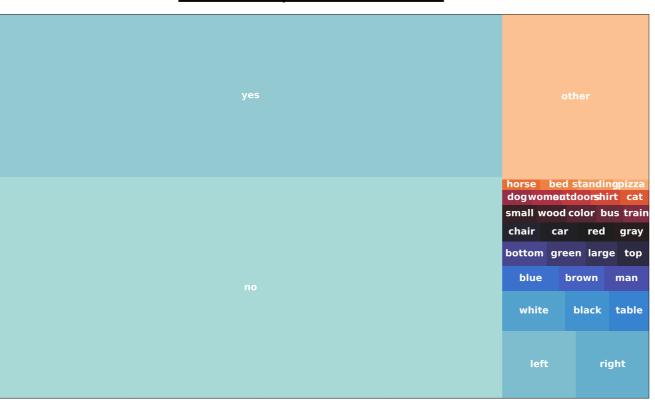




Autres jeux de données : VQA

Répartition des réponses dans les jeux de données

GQA: Réponses courtes



GQA: Réponses longues

Réponses	Nb d'apparitions	Pourcentage d'apparitions
No, there are no fences.	219099	1.53
No, there are no cars.	78880	0.55
Yes, there is a window.	60311	0.42
No, there are no glasses.	47528	0.33
No, there are no helmets.	43327	0.30
No, there are no fences or cars.	36227	0.25
The vehicle is a car.	35888	0.25
Yes, there is grass.	35585	0.24
No, there are no cars or fences.	34614	0.24
No, there are no chairs	33826	0.23
The piece of furniture is a chair.	30121	0.21

⚠ Réponses qui semblent très structurées



Autres jeux de données : VQA

Répartition des réponses dans les jeux de données

TDIUC

Visual Genome



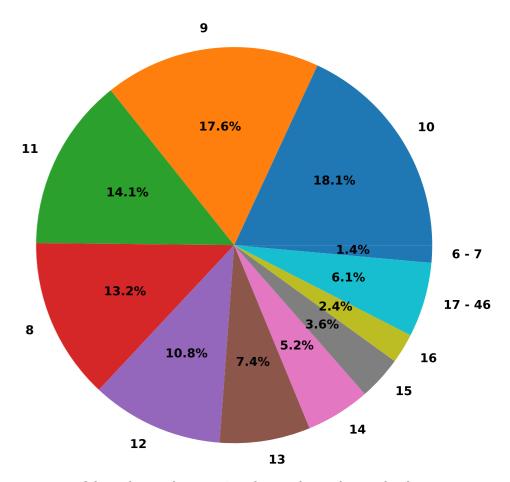


Autres jeux de données : image captioning

Jeux de données NoCaps

A batter, a catcher, and an umpire are at home plate. A woman in a miniskirt with a floral tattoo on her thigh. A man wearing a white shirt and a black tie. A group of people are rollerblading together. A bottle of wine sitting next to a wine glass. A man is wearing a black shirt and black pants. A man sitting on the ground with bikes behind him. Nombre d'occurences The people are playing baseball on the large field. A man in a shirt and tie playing the trombone. A group of people are playing instruments together. A man wearing sunglasses fires a shotgun into the air. Two people are playing with a ball in a pool. A white car is parked on the side of the road. A ladybug is crawling on a blade of grass. Three tea light candles are lit on a table. A man in a suit and tie at a podium. A man wearing a helmet is riding a bicycle. A woman is wearing a tiara and a necklace. A blue vintage vehicle with people near it. A man in a white helmet is riding a bicycle. A dragonfly sitting on a blade of grass. A man in a wheel chair playing table tennis. 0.0 0.5 1.0 1.5 2.0 2.5 3.0 Types de caption

Descriptions les plus courantes dans le jeu de données



Nombre de mots dans les descriptions

