

About Us





What to expect?

#multimodal-ml - Channel
@multimodal - Role

July Edition is Live



NEWSLETTER

Learning Resources



Paper Reading Session



Guest Speaker Session



Coding Session

First Monday

Monthly Newsletter

Industry

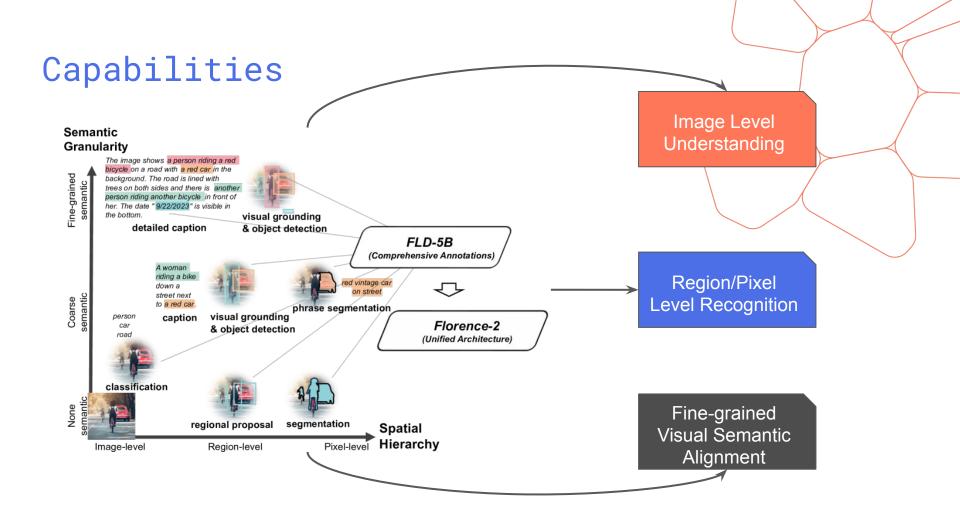
Research

Second and Fourth Friday 10 AM ET And More...

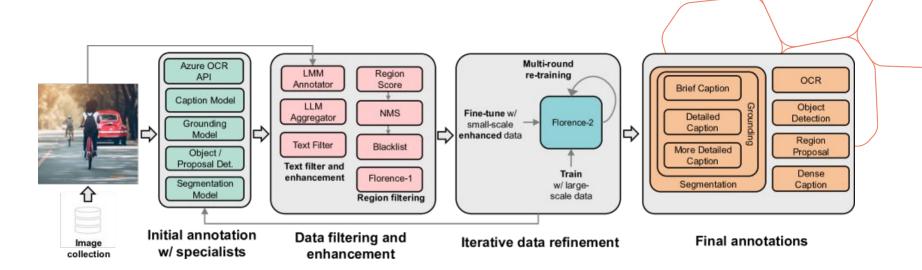
Fine-Tuning Florence-2 for DocVQA



Florence-2 101 A Quick Overview

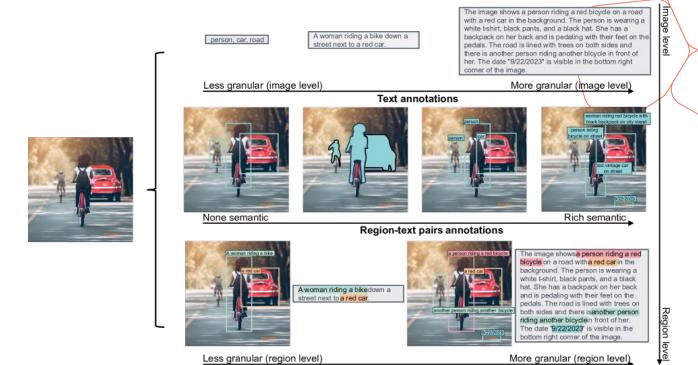


Data Engine



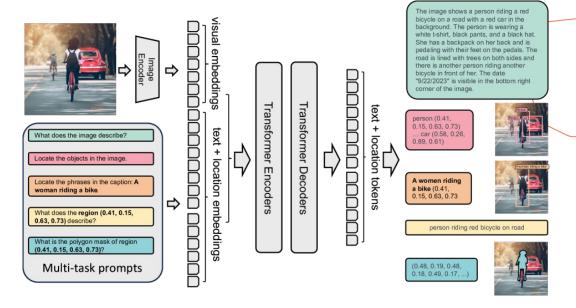
126M Images | **500M** Text Annotations | **1.3B** Region-Text Annotations | **3.6B** Text-Phrase-Region Annotations

Dataset Sneak-Peek



Text-phrase-region annotations

Model Architecture



Model	Image Encoder (DaViT)				Encoder-Decoder (Transformer)			
	dimensions	blocks	heads/groups	#params	encoder layers	decoder layers	dimensions	#params
Florence-2-B	[128, 256, 512, 1024]	[1, 1, 9, 1]	[4, 8, 16, 32]	90M	6	6	768	140M
	[256, 512, 1024, 2048]			360M	12	12	1024	410M

Supported Tasks

Task	Annotation Type	Prompt Input	Output	
Caption	Text	Image, text	Text	
Detailed caption	Text	Image, text	Text	
More detailed caption	Text	Image, text	Text	
Region proposal	Region	Image, text	Region	
Object detection	Region-Text	Image, text	Text, region	
Dense region caption	Region-Text	Image, text	Text, region	
Phrase grounding	Text-Phrase-Region	Image, text	Text, region	
Referring expression comprehension	Region-Text	Image, text	Text, region	
Open vocabulary detection	Region-Text	Image, text	Text, region	
Referring segmentation	Region-Text	Image, text	Text, region	
Region to text	Region-Text	Image, text, region	Text	
Text detection and recognition	Region-Text	Image, text	Text, region	

Available Models

Model	Model size	Model Description	
Florence-2-base[HF]	0.23B	Pretrained model with FLD-5B	_
Florence-2-large[<u>HF</u>]	0.77B	Pretrained model with FLD-5B	
Florence-2-base-ft[<u>HF</u>]	0.23B	Finetuned model on a colletion of downstream tasks	
Florence-2-large-ft[<u>HF</u>]	0.77B	Finetuned model on a colletion of downstream tasks	

Florence-2 In Action

Inference : Notebook

Fine Tuning: Notebook

Resources

- Github Repository from Microsoft
- HuggingFace Blog
- Fine-tuned Model on DocVQA
- Demo

