

论文讲解







KOSMOS-2: Grounding Multimodal Large Language Models to the World

We introduce KOSMOS-2, a Multimodal Large Language Model (MLLM), enabling new capabilities of perceiving object descriptions (e.g., bounding boxes) and grounding text to the visual world. Specifically, we represent refer expressions as links in Markdown, i.e., "[text span](bounding boxes)", where object descriptions are sequences of location tokens.







Input Prompt

Completion



the left eye of emoji



Question: How many cows are here? Answer:

COME TO CARVARY STR

Question: What does the sign say? Answer:



two cows present in the image.



The sign says "Welcome to Carnaby Street."



KOSMOS-2



Sin







Construction of Web-Scale Grounded Image-Text Pairs



Figure 3: The pipeline of constructing web-scale grounded image-text pairs.





Input Representations

In total, $P \times P$ location tokens are introduced, and these tokens are added to word vocabulary to enable unified modeling with texts.

<s> <image> Image Embedding </image> <grounding> It <box><loc₄₄><loc₈₆₃></box> seats next to a campfire <box><loc₄><loc₁₀₀₇></box> </s>

- "What is it <box><loc1><loc2></box>? It is {*expression*}."
- "What is this <box><loc₁><loc₂></box>? This is {*expression*}."
- "Describe this object <box><loc₁><loc₂></box>. This object is {*expression*}."
- "It <box><loc₁><loc₂></box> is {*expression*}."
- " This <box><loc₁><loc₂></box> is {*expression*}."
- " The object <box><loc₁><loc₂></box> is {*expression*}."

Table 9: Instruction templates used for expression generation.



COSA: Concatenated Sample Pretrained Vision-Language Foundation Model

- Due to the limited scale and quality of video-text training corpus, most vision language foundation models employ image-text datasets for pretraining and primarily focus on modeling visually semantic representations while disregarding temporal semantic representations and correlations.
- To address this issue, we propose COSA, a Concatenated Sample pretrained visionlanguage foundation model. COSA jointly models visual contents and event-level temporal cues using only image-text corpora.
- We achieve this by sequentially concatenating multiple image-text pairs as inputs for pretraining.





Figure 1: Visualizations of the traditional image-text and video-text model pretraining pipeline and the proposed unified COSA, which transforms the image-text corpus into a synthetic long-form video-text corpus online through random sample concatenation.





Single Sample Training (SST)



Concatenate Sample Training (COSA)



Figure 2: Visualizations of the training framework for COSA (bottom). In contrast to the conventional single sample training framework (SST), COSA takes the on-the-fly transformed pseudo long-form video-paragraph corpus as input. Circles and squares in the figure represent global and patch features, respectively.



THANKS

