

### HIGHLIGHTS





# **UNI-PERCEIVER V2: A GENERALIST MODEL** FOR LARGE-SCALE VISION AND VISION-LANGUAGE TASKS

HAO LI\*, JINGUO ZHU\*, XIAOHU JIANG\*, XIZHOU ZHU⊠, HONGSHENG LI, CHUN YUAN, XIAOHUA WANG, YU QIAO, XIAOGANG WANG, WENHAI WANG, JIFENG DAI \* Equal Contribution  $\bowtie$  Corresponding Author

$$f_{\text{image}}(x) = \text{Concat}\left(\{q_i^{\text{global}}\}_{i=1}^M, \{q_j^{\text{proposal}}\}_{j=1}^N\right)$$

$$q^{\text{global}} = \text{Concat}\left(\left\{\operatorname{AttnPool}_{i}(\mathcal{F}_{L})\right\}_{i=1}^{M'}, \operatorname{Flatten}(\mathcal{F}_{L})\right)$$

$$q_j^{\text{proposal}} = q_j^{\text{sem}} + \mathcal{B}(q_j^{\text{box}}) + \mathcal{M}(q_j^{\text{mask}})$$

$$\Rightarrow \begin{cases} \mathbf{g}_{t} \leftarrow \omega_{k} \frac{\nabla L_{t,k}(\theta_{t-1})}{\|\nabla L_{t,k}(\theta_{t-1})\|} \\ \mathbf{m}_{t} = (1 - \beta_{1}) \mathbf{m}_{t-1} + \frac{\beta_{1}}{s_{k}} \mathbf{g}_{t} \\ \mathbf{n}_{t} = (1 - \beta_{2}) \mathbf{n}_{t-1} + \frac{\beta_{2}}{s_{k}} \mathbf{g}_{t}^{2} \\ \theta_{t} = \theta_{t-1} - \alpha \frac{\mathbf{m}_{t}}{\sqrt{\mathbf{n}_{t}} + \varepsilon} \end{cases}$$

## EXPERIMENTS

Methods	#params	Image Classification	Object Detection	Instance Segmentation	Image Captioning COCO		Text Retrieval		Image Retrieval	
		ImageNet-1k	COCO	COCO			COCO	Flickr30k	COCO	Flickr30k
		Acc	mAP	mAP	B@4	CIDEr	R@1	R@1	R@1	R@1
Pix2Seq v2	132M	_	<u>46.5</u>	<u>38.2</u>	34.9	_	_	_	_	_
UniTab	185M	_	_	_	_	115.8	_	_	_	_
Unified-IO LARGE	776M	71.8	_	_	_	-	_	_	_	_
Unified-IO <sub>XL</sub>	2.9B	79.1	_	_		122.3	_	_	-	_
Flamingo-3B	3.2B	_	_	_	_	-	65.9	<u>89.3</u>	48.0	<u>79.5</u>
Uni-Perceiver BASE	124M	79.2	_	_	32.0	_	64.9	82.3	50.7	71.1
Uni-Perceiver LARGE	354M	82.7	_	_	35.3	_	67.8	83.7	54.1	74.2
Uni-Perceiver-MoE <sub>BASE</sub>	167M	80.3	_	_	33.2	_	64.6	82.1	51.6	72.4
Uni-Perceiver-MoE LARGE	505M	<u>83.4</u>	_	_	<u>35.5</u>	_	<u>67.9</u>	83.6	<u>55.3</u>	75.9
Uni-Perceiver-v2 <sub>BASE</sub>	308M	86.3	58.6	50.6	35.4	116.9	71.8	88.1	55.6	73.8
Uni-Perceiver-v2 <sub>LARGE</sub>	446M	<b>87.2</b> (+3.8)	<b>61.9</b> (+15.4)	<b>53.6</b> (+15.4)	<b>36.5</b> (+1.6)	<b>122.5</b> (+0.2)	<b>75.0</b> (+7.1)	<b>89.3</b> (+0.0)	<b>58.5</b> (+3.2)	<b>79.6</b> (+0.1)

### Comparison of Uni-Perceiver v2 with generalist models and commonly-recognized strong task-specific models



## **ABLATION STUDIES**

## Ablation of task collaboration and interference

Tasks	COCO Detection	ImageNet-1k Classification	CO Retr	CO ieval	COCO Caption	Pretrained Method	Pretrained Data	COCO Detection	ImageNet-1k Classification	COCO Retrieval	COCO Caption
Single Task	50.1	76.1	50.0	37.6	30.2	Supervised	IN-1k	45.7	76.8	51.2 38.9	27.3
All Tacks	/0.8	76.3	46.0	347	28.9	Supervised	IN-21k	48.3	80.1	55.1 41.2	30.2
	<b>H7.0</b>	70.5	10.0	01./	20.7	Supervised	IN-1k & COCO	49.9	76.9	51.3 38.8	30.6
w/o Detection	-	76.6 (+0.3)	47.0(+1.0)	34.6(-0.1)	30.4(+0.5)	$M_0C_0 v^2$	IN-1k	48.3	75.0	54.8 40.5	29.6
w/o Classification	$50.1_{(+0.3)}$	-	51.6 (+5.6)	38.6(+3.9)	$25.9_{(-3.0)}$		CLIP data	17.2	73.8	55 2 /1 2	22.0
w/o Retrieval	$49.5_{(-0.3)}$	$76.3_{(+0.0)}$	_	-	$27.4_{(-1.5)}$			47.2	75.0	55.5 41.5	52.0
w/o Captioning	$49.7_{(-0.1)}$	$76.3_{(+0.0)}$	51.2 (+5.2)	38.3(+3.6)	-						
All Tasks w/ MoE	$49.9_{(+0.1)}$	76.9 (+0.6)	51.3 (+5.3)	38.8(+4.1)	30.6(+0.7)						

## Ablation of different representation types

Representation	COCO	ImageNet-1k	COCO	COCO	Task	Gather	TRCN	COCO	ImageNet-1k	COCO	COCO
Types	Detection	Classification	Retrieval	Caption	Sampling	Feature	IDGIN	Detection	Classification	Retrieval	Caption
Global	_	76.8	46.3 34.6	28.8	mixed			49.6	76.7	40.1 31.9	27.6
Regional	48.2	75.9	52.3 39.2	31.2	unmixed			49.2	76.6	39.8 30.9	27.5
Global + Regional	49.9	76.9	51.3 38.8	30.6	unmixed	$\checkmark$		49.3	76.8	50.4 37.3	27.6
	<b>'</b>	•	,		unmixed	$\checkmark$	$\checkmark$	49.9	76.9	51.3 38.8	30.6



### **Comparison of Uni-Perceiver v2 with existing generalist models**

### Ablation of different pre-trained image encoders

### Ablation of sampling and optimization strategies