# Large-Scale Adversarial Training for Vision-and-Language Representation Learning

Zhe Gan[1], Yen-Chun Chen[1], Linjie Li[1], Chen Zhu[2], Yu Cheng[1], Jingjing Liu[1]

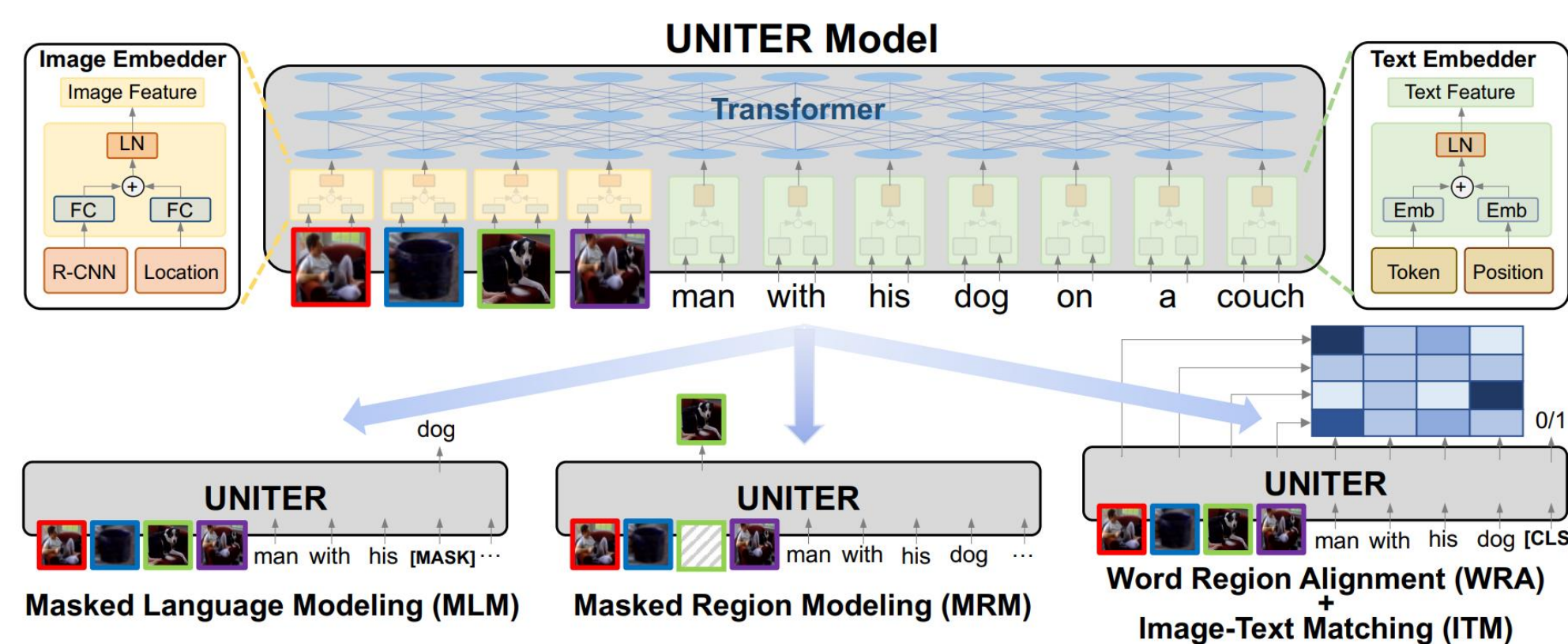[1]Microsoft Dynamics 365 AI Research,    [2]University of Maryland, College Park

## Motivation & Contribution

- *Multimodal pre-training*, such as ViLBERT, LXMERT and UNITER, has made tremendous progress in Vision-and-Language (V+L) research
- However, aggressive finetuning of pre-trained models often falls into the *overfitting trap*
- *Adversarial training* has shown great potential in improving the generalization ability of BERT for language understanding tasks
- *Our Contribution*: the first known effort to study large-scale adversarial training for V+L tasks
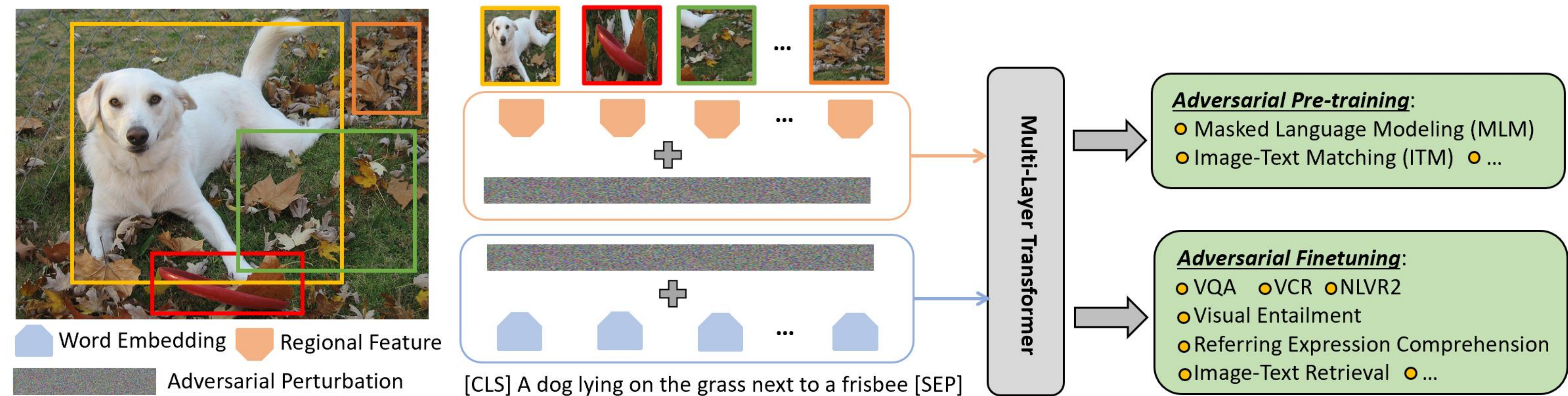
## Algorithm and Backbone (UNITER)



**Algorithm 1** "Free" Multi-modal Adversarial Training used in VILLA.

**Require:** Training samples $\mathcal{D} = \{(x_{img}, x_{txt}, y)\}$, perturbation bound $\epsilon$, learning rate $\tau$, ascent steps $K$, ascent step size $\alpha$

1: Initialize $\theta$
2: **for** epoch $= 1 \ldots N_{ep}$ **do**
3:   **for** minibatch $B \subset X$ **do**
4:     $\delta_0 \leftarrow \frac{1}{\sqrt{N_\delta}} U(-\epsilon, \epsilon),\ g_0 \leftarrow 0$
5:     **for** $t = 1 \ldots K$ **do**
6:       Accumulate gradient of parameters $\theta$ given $\delta_{img,t-1}$ and $\delta_{txt,t-1}$
7:         $g_t \leftarrow g_{t-1} + \frac{1}{K} \mathbb{E}_{(x_{img}, x_{txt}, y) \in B} [\nabla_\theta (\mathcal{L}_{std}(\theta) + \mathcal{R}_{at}(\theta) + \mathcal{R}_{kl}(\theta))]$
8:       Update the perturbation $\delta_{img}$ and $\delta_{txt}$ via gradient ascend
9:         $\tilde{y} = f_\theta(x_{img}, x_{txt})$
10:        $g_{img} \leftarrow \nabla_{\delta_{img}} [L(f_\theta(x_{img} + \delta_{img}, x_{txt}), y) + L_{kl}(f_\theta(x_{img} + \delta_{img}, x_{txt}), \tilde{y})]$
11:        $\delta_{img,t} \leftarrow \Pi_{\|\delta_{img}\|_F \leq \epsilon} (\delta_{img,t-1} + \alpha \cdot g_{img}/\|g_{img}\|_F)$
12:        $g_{txt} \leftarrow \nabla_{\delta_{txt}} [L(f_\theta(x_{img}, x_{txt} + \delta_{txt}), y) + L_{kl}(f_\theta(x_{img}, x_{txt} + \delta_{txt}), \tilde{y})]$
13:        $\delta_{txt,t} \leftarrow \Pi_{\|\delta_{txt}\|_F \leq \epsilon} (\delta_{txt,t-1} + \alpha \cdot g_{txt}/\|g_{txt}\|_F)$
14:     **end for**
15:     $\theta \leftarrow \theta - \tau g_K$
16:   **end for**
17: **end for**

Code is available at
https://github.com/zhegan27/VILLA

## The Proposed VILLA Framework



[CLS] A dog lying on the grass next to a frisbee [SEP]

Word Embedding    Regional Feature    Adversarial Perturbation

**Adversarial Pre-training**:
- Masked Language Modeling (MLM)
- Image-Text Matching (ITM)    ○ ...

**Adversarial Finetuning**:
- VQA  ○ VCR  ○ NLVR2
- Visual Entailment
- Referring Expression Comprehension
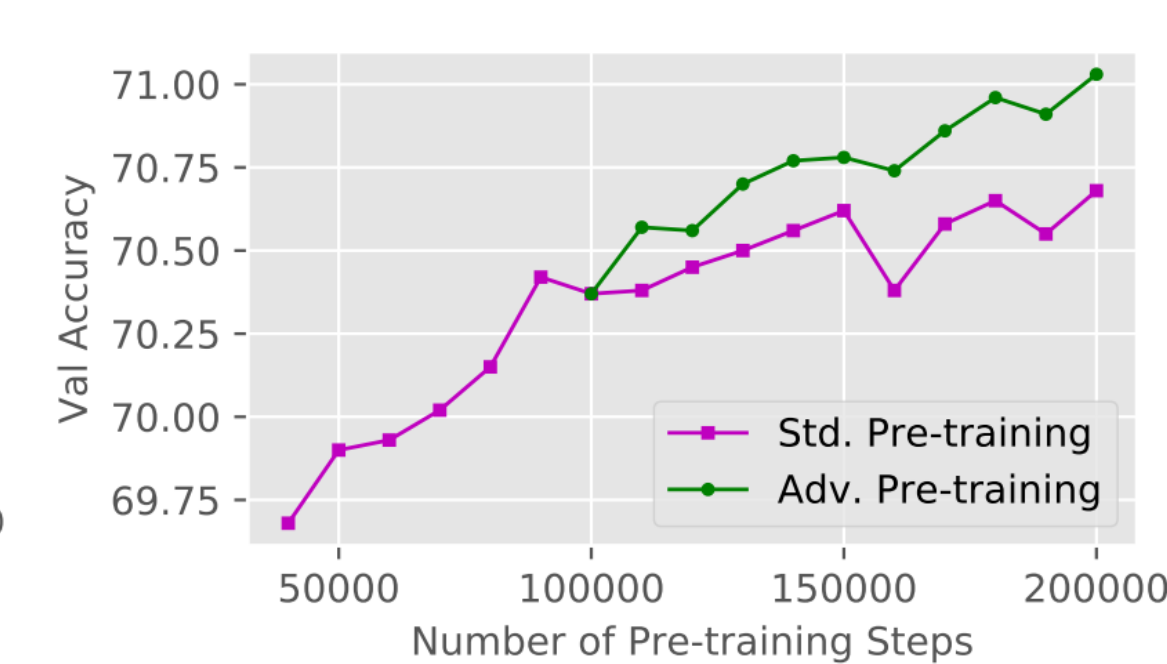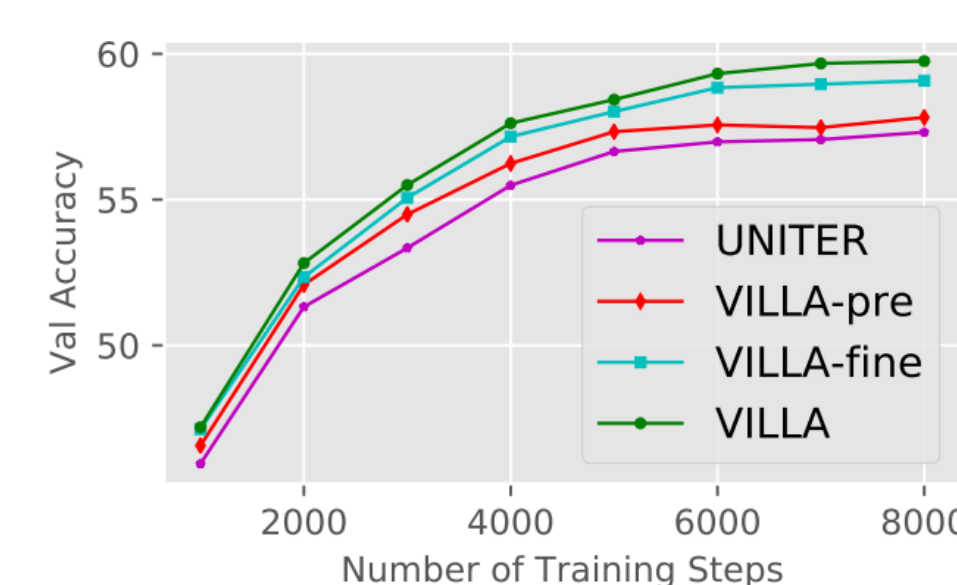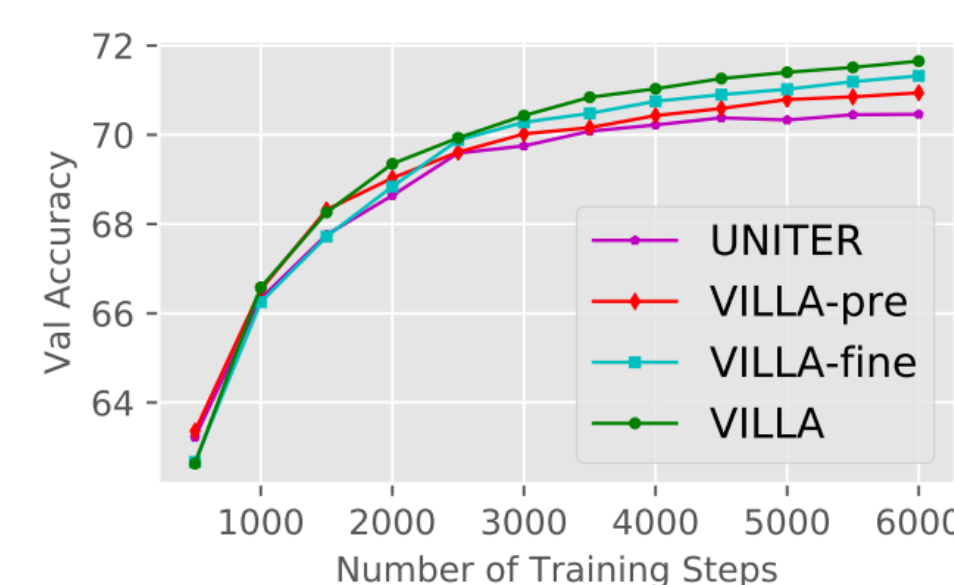- Image-Text Retrieval  ○ ...

○ *Adversarial pre-training and finetuning*    ○ *Perturbations in the embedding space*    ○ *Enhanced adversarial training algorithm*

## Experimental Results

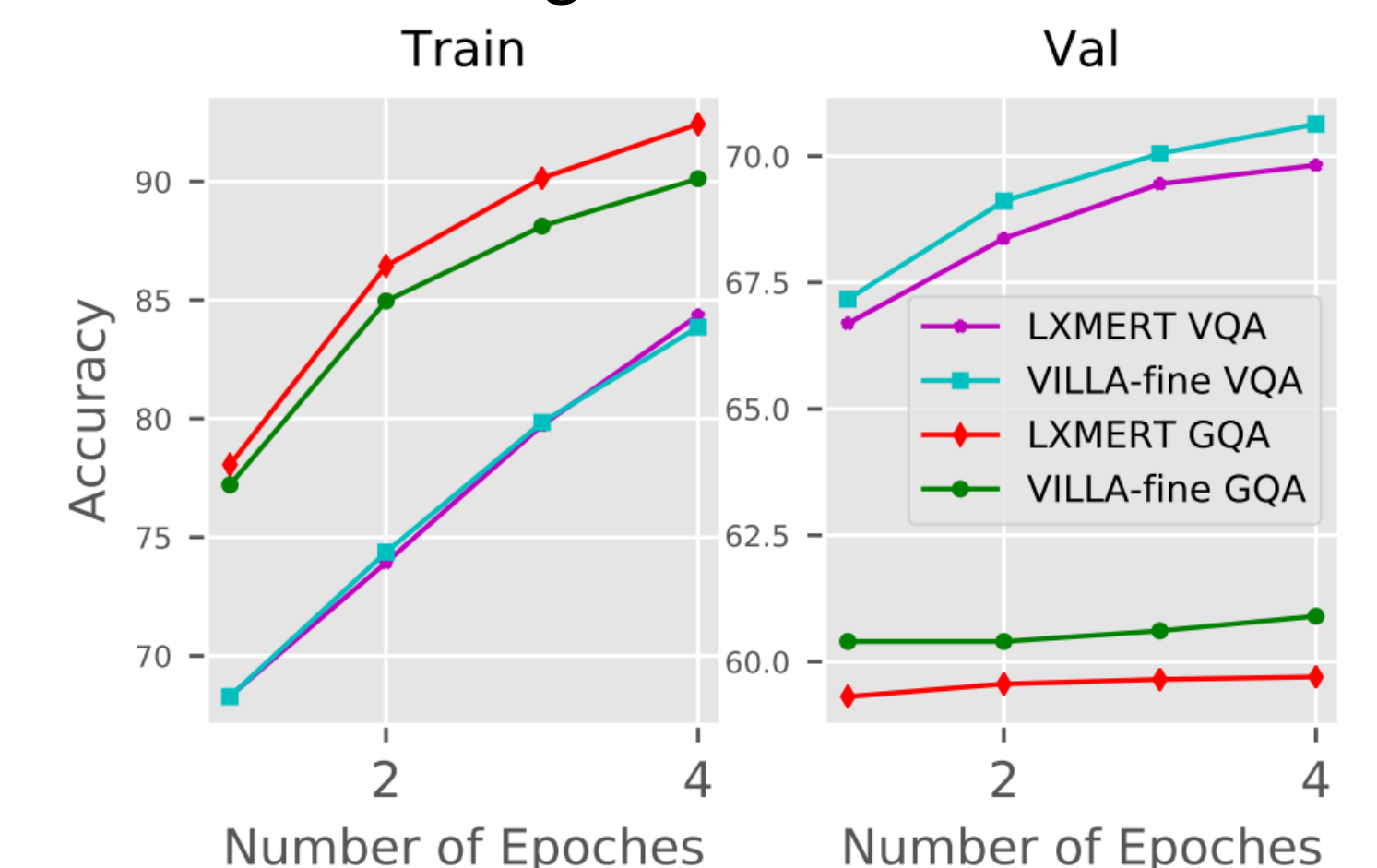○ New state of the art on a wide range of V+L tasks (see paper for details)

| Task | VQA | VCR | NLVR2 | VE | RefCOCOg | RefCOCO+ | Flickr30k IR | Flickr30k TR | VQA-Rep. |
|------|-----|-----|-------|-----|----------|----------|--------------|--------------|----------|
| UNITER | 74.02 | 62.8 | 79.98 | 79.38 | 75.77 | 66.70 | 75.56 | 87.30 | 64.56 |
| VILLA | 74.87 | 65.7 | 81.47 | 80.02 | 76.71 | 66.84 | 76.26 | 87.90 | 65.35 |

○ Both adversarial pre-training (**VILLA-pre**) and finetuning (**VILLA-fine**) contribute to performance boost



(a) VQA

(b) VCR

○ Adversarial training serves as effective regularizer



○ Adversarial training on image or text modality alone is already effective
○ VILLA captures richer visual coreference and visual relation knowledge than UNITER
○ VILLA learns more accurate and sharper attention maps than UNITER
○ VILLA is more robust to paraphrases than UNITER

○ VILLA can be readily extended to other pre-trained V+L models, such as LXMERT