# Google DeepMind





• For the first time, we show that *multi-modal instruction tuning* can make imagen models to generalize to a large variety of generation tasks.

### **Summary of Contributions**

- We introduce *multi-modal instruction* for image generation, which unifies tasks that inputs different conditions (e.g. Control2Img, Subject2Img, Style2Img, etc.)
- We proposed a *two-stage fine-tuning framework* to adapt pre-trained text2image models to accept multimodal instruction.
- Instruct-Imagen performs in-domain tasks on par with SoTA, and shows strong compositional generalization to novel generation tasks.

**Unifying Generation Tasks with** *Multi-Modal Instruction* We cast existing image generation tasks in the format of multi-modal instruction



## **Two-Stage Multi-Modal Instruction-Tuning**

S1. Retrieval-augmented continue training, S2. Multi-modal instruction-tuning (stage 1 adapts the pre-trained model, and stage 2 fine-tunes the model with target tasks)



Pre-trained Imagen

: Retrieval-augmented Training phase

# Instruct-Imagen: Imagen Generation with Multi-modal Instruction Hexiang Hu\*, Kelvin C.K. Chan\*, Yu-Chuan Su\*, Wenhu Chen\* Yandong Li, Kihyuk Sohn, Yang Zhao, Xue Ben, Boqing Gong, William W. Cohen, Ming-Wei Chang, Xuhui Jia

: Multi-modal Instruction-tuning

### Subject-Driven Generation

![](_page_0_Picture_24.jpeg)

### **Style-Driven Generation**

![](_page_0_Picture_26.jpeg)

![](_page_0_Picture_27.jpeg)

### Mask-Based Editing

![](_page_0_Picture_30.jpeg)

### Mask-Free Editing

![](_page_0_Picture_32.jpeg)

![](_page_0_Picture_33.jpeg)

![](_page_0_Picture_34.jpeg)

# Multi-Turn Editing

![](_page_0_Picture_36.jpeg)

![](_page_0_Picture_39.jpeg)

![](_page_0_Picture_40.jpeg)

![](_page_0_Picture_41.jpeg)

(a) In-domain Evaluation

(b) Zero-shot Evaluation