| cases | doc_1 | | doc_2 | | decision | id |
|---|---|---|---|---|---|---|
| | authors | • Xiaoyi Dong<br>• Jianmin Bao<br>• Ting Zhang<br>• Dongdong Chen<br>• Weiming Zhang<br>• Lu Yuan<br>• Dong Chen<br>• Fang Wen<br>• Nenghai Yu | authors | • Xiaoyi Dong<br>• Jianmin Bao<br>• Ting Zhang<br>• Dongdong Chen<br>• Weiming Zhang<br>• Lu Yuan<br>• Dong Chen<br>• Fang Wen<br>• Nenghai Yu<br>• Baining Guo | DUPLICATES | 240 |
| | title | PeCo: Perceptual Codebook for BERT Pre-training of Vision Transformers | title | PeCo: Perceptual Codebook for BERT Pre-training of Vision Transformers | | |
| | publication_date | 2021-11-24 00:00:00 | publication_date | 2021-11-24 18:59:58+00:00 | | |
| | source | SupportedSources.SEMANTIC_SCHOLAR | source | SupportedSources.ARXIV | | |
| | journal | ArXiv | journal | None | | |
| | volume | abs/2111.12710 | volume | | | |
| | doi | | doi | | | |
| | urls | • https://www.semanticscholar.org/paper/3e38f4b4055abecbac2e618df2ecb33554073e08 | urls | • http://arxiv.org/pdf/2111.12710v3<br>• http://arxiv.org/abs/2111.12710v3<br>• http://arxiv.org/pdf/2111.12710v3 | | |
| | id | id-2127133131246859522 | id | id-1611205866241915138 | | |
| | abstract | This paper explores a better prediction target for BERT pre- training of vision transformers. We observe that current prediction targets disagree with human perception judgment. This contradiction motivates us to learn a perceptual prediction target. We argue that perceptually similar images should stay close to each other in the prediction target space. We surprisingly ï¬nd one simple yet effective idea: enforcing percep- tual similarity during the dVAE training. Moreover, we adopt a self-supervised transformer model for deep feature extrac- tion and show that it works well for calculating perceptual similarity. We demonstrate that such learned visual tokens in- deed exhibit better semantic meanings, and help pre-training achieve superior transfer performance in various downstream tasks. For example, we achieve 84.5% Top-1 accuracy on ImageNet-1K with ViT-B backbone, outperforming the com- petitive method BEiT by +1.3% under the same pre-training epochs. Our approach also gets signiï¬cant improvement on object detection and segmentation on COCO and semantic segmentation on ADE20K. Equipped with a larger backbone ViT-H, we achieve the state-of-the-art ImageNet accuracy ( 88.3% ) among methods using only ImageNet-1K data. | abstract | This paper explores a better prediction target for BERT pre-training of vision transformers. We observe that current prediction targets disagree with human perception judgment.This contradiction motivates us to learn a perceptual prediction target. We argue that perceptually similar images should stay close to each other in the prediction target space. We surprisingly find one simple yet effective idea: enforcing perceptual similarity during the dVAE training. Moreover, we adopt a self-supervised transformer model for deep feature extraction and show that it works well for calculating perceptual similarity.We demonstrate that such learned visual tokens indeed exhibit better semantic meanings, and help pre-training achieve superior transfer performance in various downstream tasks. For example, we achieve $\textbf{84.5\%}$ Top-1 accuracy on ImageNet-1K with ViT-B backbone, outperforming the competitive method BEiT by $\textbf{+1.3\%}$ under the same pre-training epochs. Our approach also gets significant improvement on object detection and segmentation on COCO and semantic segmentation on ADE20K. Equipped with a larger backbone ViT-H, we achieve the state-of-the-art ImageNet accuracy (\textbf{88.3\%}) among methods using only ImageNet-1K data. | | |
| | versions | | versions | | | |