

---

# Perceptual Group Tokenizer: Building Perception with Iterative Grouping

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Human visual recognition system shows astonishing capability of compressing  
2 visual information into a set of tokens containing rich representations without label  
3 supervision. One critical driving principle behind it is perceptual grouping [1, 2, 3].  
4 Despite being widely used in computer vision in the early 2010s, it remains a  
5 mystery whether perceptual grouping can be leveraged to derive a neural visual  
6 recognition backbone that generates as powerful representations. In this paper, we  
7 propose *the Perceptual Group Tokenizer*, a model that entirely relies on grouping  
8 operations to extract visual features and perform self-supervised representation  
9 learning, where a series of grouping operations are used to iteratively hypothe-  
10 size the context for pixels or superpixels to refine feature representations. We  
11 show that the proposed model can achieve competitive performance compared  
12 to state-of-the-art vision architectures, and inherits desirable properties including  
13 *adaptive computation without re-training*, and interpretability. Specifically, Percep-  
14 tual Group Tokenizer achieves 79.7% on ImageNet-1K *self-supervised learning*  
15 benchmark with linear probe, marking a new progress under this paradigm.

## 16 1 Introduction

17 Visual recognition mechanisms matter. The pursuit of advanced vision algorithms that encode  
18 an image to meaningful representations dates back to late 80s, with two paradigms marking the  
19 progress over the past 40 years: feature detection [4, 5, 6, 7] and perceptual grouping [8, 9, 10],  
20 where feature detection focuses on specific distinctive patterns, while perceptual grouping considers  
21 similarities among all pixels to produce a compact set of tokens as proxies for image representation.  
22 Ever since the surge of deep learning, feature detection has predominated the vision field and  
23 become the main rationale in representation learning backbone designs and made impressive strides  
24 [11, 12, 6, 13, 14, 15, 7]. The success of the former paradigm is, although striking, raising the  
25 question of whether perceptual grouping can also be used as the driving principle to construct a visual  
26 recognition model.

27 Different from detecting and selecting distinctive features, perceptual grouping emphasizes on  
28 learning feature space where similarity of all pixels can be effectively measured [9, 10]. With such a  
29 feature space, semantically meaningful objects and regions can be easily discovered with a simple  
30 grouping algorithm and used as a compact set to represent an image [9, 10, 16]. This indicates that  
31 image understanding is essentially “pixel space tokenization”, and being able to produce generalizable  
32 feature representations is connected to whether correct contextual pixels are binded together [17, 18].

33 The intriguing properties of perceptual grouping, including natural object discovery, deep connections  
34 with information theory and compression [19], and association with biological vision system [3]  
35 or cognitive science explanations [1], have led to a strong revive recently under deep learning  
36 frameworks [16, 20, 21, 22, 23]. However, these methods are either still focusing on small or toy

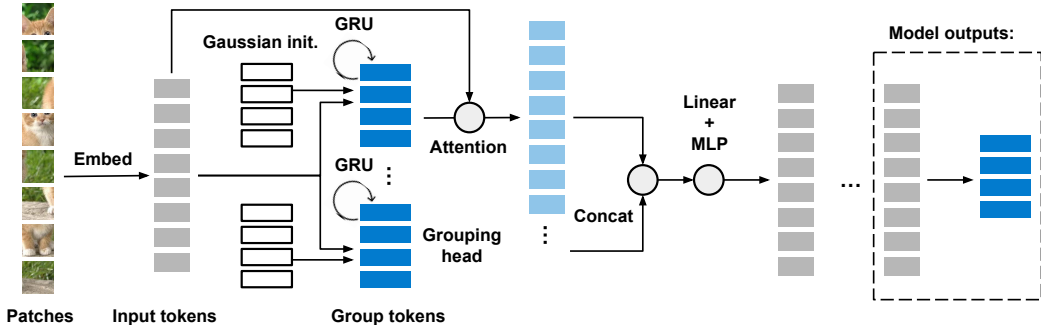


Figure 1: Perceptual Group Tokenizer takes in a sequence of patches (or pixels), generates high-dimensional embedding vectors for all patches, then they pass through a series of grouping layers to refine the embedding vectors as feature representations. Each grouping layer performs  $K$  rounds of binding from input tokens to group tokens. To consider various grouping possibilities, multiple grouping heads are adopted. Each group token provides a useful context for input tokens for feature refinement. The final output of the model contains refined input token, group tokens, and assignments between input tokens and groups tokens.

37 datasets [16, 24, 23], or used as a side add-on component [21] to strengthen the existing vision  
 38 architectures for increased interpretability. Whether perceptual grouping can be used to build models  
 39 and learn representations that are as informative and expressive as those learned by state-of-the-art  
 40 vision architectures remains an open question.

41 In this paper, we propose *Perceptual Group Tokenizer*, a model trained under *self-supervised learning*  
 42 framework and building visual representation *entirely relying on perceptual grouping operations*.  
 43 Given an image, the core of our model is to understand each pixel or patch through hypothesizing  
 44 its contexts with grouping operations. Starting from given input patches, the grouping operation  
 45 performs iterative binding process onto a set of randomly sampled group tokens to determine the  
 46 affinity groups based on similarities. The group tokens are then used as hypothesized contexts to refine  
 47 the feature representation for the image. We show that applying this simple principle can already  
 48 produce expressive representations and works well on self-supervised large dataset pretraining.

49 Compared to self attention, why can grouping operation work? Analyzing the rationale behind  
 50 it, we build connections from grouping operation to self attention, showing that, if group tokens  
 51 are treated as communication channels, self attention can potentially automatically emerge during  
 52 learning processes as a special case, while the grouping operation can produce even richer interactions  
 53 among tokens. Under this viewpoint, ViT [25] can be considered as a grouping backbone, with a  
 54 fixed number of grouping slots depending on number of input tokens, and the binding is achieved  
 55 through stacking more than one layer with non-shared weights. This provides one explanation on  
 56 why grouping mechanism can be effective on visual representation learning and has the potential to  
 57 be a promising competitive paradigm for vision architecture designs.

58 The primary contribution of this work is proposing a new architecture derived purely by perceptual  
 59 grouping that achieves competitive performance compared to other state-of-the-art architectures  
 60 on *self-supervised learning* benchmarks, contributing to a new paradigm of developing vision  
 61 architectures. We thoroughly analyze the design space of perceptual grouping backbones, show the  
 62 capability of *adaptive computation without re-training*, and visualize the grouping process which  
 63 produces semantically meaningful bindings among patch tokens.

## 64 2 Models

65 In this section, we introduce Perceptual Group Tokenizer (PGT), a visual recognition architecture  
 66 entirely driven by perceptual grouping principles. We discuss the core operations for grouping  
 67 in section 2.1 and the building blocks network architectures in section 2.2 in the main paper, and  
 68 self-supervision loss and more discussion in section 4.1.2 and 4.1.3 in the supplementary material.

69 **2.1 Perceptual grouping**

70 We start with introducing notations for our method. Given an image  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$ , we first reshape  
 71 it as a sequence of small patches<sup>1</sup>. Each patch  $\mathbf{x}_p \in \mathbb{R}^{h \times w \times c}$  has spatial shape  $h \times w$ , where  $h \ll H$   
 72 and  $w \ll W$ , leading to  $N = \frac{HW}{hw}$  number of patches per image. To represent a patch, we embed  
 73 it into a high-dimensional vector  $\mathbf{h} \in \mathbb{R}^d$ . The set of embedded tokens  $\{\mathbf{h}_i\}^N$  is referred as *input*  
 74 *tokens* in later parts, and used as inputs for the following grouping blocks.

75 **Feature refinement through hypothesizing contexts.** One pixel does not have meanings without  
 76 putting it into contexts. At a high level, image understanding or feature learning is equivalent to  
 77 binding the correct contextual pixels at all locations. The core idea of our model is to generate  
 78 many (e.g. over-complete w.r.t number of objects in the image) hypothesized contexts and use the  
 79 hypothesized contexts as cues to refine the feature representation of each patch. This process is  
 80 achieved through a grouping module. Given input tokens  $\{\mathbf{h}_i\}^N$ , the grouping module starts from a  
 81 set of random samples (referred as *group tokens*) from a Gaussian distribution, then performs binding  
 82 process to aggregate information from input tokens to the group tokens, and ends up with a set of  
 83 group tokens  $\mathbf{c}^* = \{\mathbf{c}_j^*\}_{j=1}^M$  representing hypothesized contexts among input tokens. The relation  
 84 between  $\mathbf{h}_i$  and  $\mathbf{c}_j$  is soft assignment, indicating how likely an input token belongs to that context.  
 85 Note that there are often *various ways of generating groupings for an image*, e.g. different semantics,  
 86 colors, textures, etc., we propose the “multi-grouping operation” to hypothesize rich contexts for  
 87 tokens. The overall model is shown in figure 1.

88 **Multi-grouping operation.** The building block of our model is the multi-grouping operation  $\mathcal{G}$ ,  
 89 which contains multiple heads to perform the binding process in parallel. This design encourages the  
 90 model to consider multiple ways of generating groups under different projection spaces. Each head  
 91 owns a separate Gaussian distribution with learnable means and variance, similar to [26, 16]. Starting  
 92 from a set of randomly sampled initial group tokens  $\mathbf{c}_{\text{HEAD}}^{(0)} \sim p(\boldsymbol{\mu}_{\text{HEAD}}, \boldsymbol{\sigma}_{\text{HEAD}})$ , the grouping operation  
 93 uses doubly normalized attention weights to aggregate information from  $\mathbf{h}$ , and the produced group  
 94 tokens  $\mathbf{c}_{\text{HEAD}}^{(1)}$  are used for the next round binding. The attention normalization and feature projection  
 95 are performed in all heads separately.

$$\mathbf{c}_{\text{HEAD}}^{(1)} = \mathcal{G}(\mathbf{c}_{\text{HEAD}}^{(0)}, \mathbf{h}; \theta) \tag{1}$$

...

$$\mathbf{c}_{\text{HEAD}}^* = \mathbf{c}_{\text{HEAD}}^{(K)} = \mathcal{G}(\mathbf{c}_{\text{HEAD}}^{(K-1)}, \mathbf{h}; \theta) \tag{2}$$

96 where after  $K$  steps the final group tokens  $\mathbf{c}^* = \mathbf{c}^{(K)}$  is obtained, and  $\theta$  is learnable parameters in  $\mathcal{G}$ .  
 97 The grouping operator is summarized in algorithm 1.

98 **Implicit differentiation.** The iterative grouping process unrolls  $K$  steps per operation and leads to  
 99 heavy burden in the training computation graph. Instead of explicitly backpropagating through the  
 100 unrolled graph, we follow [24] and treat the multi-grouping process as a fixed point iteration per head.  
 101 The gradient in the backpropagation is approximated using first-order Neumann series.

102 **2.2 Network architecture**

103 Similar to standard ViT, our model refines the hidden representation  $\mathbf{h}$  using  $L$  model layers. We use  
 104  $\mathbf{h}^l$  to denote the representation after each layer, and explain the design in this section.

105 **Grouping layer.** Each grouping layer takes in  $\mathbf{h}^{l-1}$  as input, and uses the grouping operation in  
 106 equation 1 to generate group tokens  $\mathbf{c}_{\text{HEAD}}^* = \{\mathbf{c}_{j,\text{HEAD}}^*\}_{j=1}^M$ . To use the group tokens to provide  
 107 context for each  $\mathbf{h}_i^{l-1}$ , we perform another attention operation to obtain the attention matrix (only  
 108 normalized over group token axis)  $\mathbf{A} \in \mathbb{R}^{N \times M}$  representing the assignment from input tokens to  
 109 group tokens, and aggregate the feature back to the input token space:

$$\mathbf{h}_{\text{HEAD}}^l = \mathbf{A}[\mathbf{c}_{1,\text{HEAD}}^*; \mathbf{c}_{2,\text{HEAD}}^*; \dots; \mathbf{c}_{M,\text{HEAD}}^*] \tag{3}$$

$$\mathbf{h}^l = \text{Linear}([\mathbf{h}_{\text{HEAD}_1}^l; \dots; \mathbf{h}_{\text{HEAD}_H}^l]) \tag{4}$$

$$\mathbf{h}^l = \mathbf{h}^{l-1} + \text{MLP}(\mathbf{h}^l) \tag{5}$$

<sup>1</sup>We use  $4 \times 4$  patches as inputs in this work. Note that our method is generalizable to either pure pixels or other forms of superpixels given a proper patch-to-vector embedding layer.

Method	Arch	Param.	Linear probe (top-1 acc)
SCLR [31]	RN50W4	375	76.8
SwAV [32]	RN50W2	93	77.3
BYOL [32]	RN50W2	93	77.4
DINO [29]	ViT-B/16	85	78.2
SwAV [32]	RN50W5	586	78.5
BYOL [32]	RN50W4	375	78.6
iBOT [33]	ViT-B/16	85	79.5
BYOL [32]	RN200W2	250	79.6
SCLRv2 [34]	RN152w3+SK	794	79.8
DINO [29]	ViT-B/8	85	<b>80.1</b>
BEiTv2 [35]	ViT-B/16	85	<b>80.1</b>
Ours (PGT-B-256)	PGT-B	70	79.3
Ours (PGT-B-384)	PGT-B	70	79.4
Ours (PGT-B-512)	PGT-B	70	79.6
Ours (PGT-B-768)	PGT-B	70	<b>79.7</b>

Table 1: Comparison with strong baselines on ImageNet-1K under linear probe evaluation protocol. PGT- $X$  represents  $X$  number of group tokens per grouping layer in inference (same trained model with 256 tokens is used). Our model achieves 79.7%, competitive with state-of-the-art vision backbones, and outperforms ResNet architectures.

110 This layer definition follows the standard ViT layer as close as possible, where features from each  
111 head are aggregated through concatenation and a linear layer transformation. Each token  $h$  is further  
112 refined using a follow up multi-layer perceptron.

113 **Grouping blocks.** Similar to previous architecture designs [6, 27], we define blocks for the model.  
114 One block contains multiple grouping layers that share the same hyperparameters setups, i.e. number  
115 of group tokens, group token dimensions. The full model contains three grouping blocks. This  
116 increases the flexibility when exploring model design spaces.

117 See more details in sections 4.1.2 and 4.1.3 in the supplementary material.

### 118 3 Experiments

119 We evaluate the representation learned by our model on standard benchmarks, specifically ImageNet-  
120 1K dataset. Summarized in the main paper in section 3.1. In the supplementary material, we also  
121 thoroughly explore and analyze the design space of perceptual group tokenizer in section 4.2.2, show  
122 the adaptive computation ability in section 4.2.3, demonstrate the generalization ability on semantic  
123 segmentation in section 4.2.4, and visualize the learned attention in section 4.2.5.

#### 124 3.1 Main results

125 **Setup.** The widely-adopted standard benchmark for evaluating self-supervised learning methods  
126 is ImageNet ILSVRC-2012 (ImageNet-1K) [28]. Performance of models are measured by top-1  
127 classification accuracy. The pre-trained backbones are frozen, with a linear classifier trained on top.  
128 For fair comparison, we follow the standard data augmentation used in [29], with the same number of  
129 global views and local views. The model is optimized using AdamW [30] with learning rate 0.0005  
130 and 1024 batch size for 600 epochs, trained with TPUv5 for 21k core hrs (512 cores for 41 hrs).  
131 We use  $4 \times 4$  patches as image tokens, which keeps as much details as possible while maintaining  
132 reasonable computation costs. For machines, we use TPUv5 to run experiments.

133 The main results are summarized in table 3.1. We mainly compare with ResNet and ViT backbones,  
134 the two main stream vision architectures to show that perceptual grouping architecture can also  
135 achieve competitive results on the challenging ImageNet-1K benchmark. Although our model is  
136 trained with 256 group tokens, the model can use different numbers of group tokens in inference (more  
137 experiments in 4.2.2). We evaluate PGT with 256, 384, 512, and 768 number of group tokens and  
138 observe that with PGT-768 the model can achieve 79.7% top-1 accuracy, showing the self-supervised  
139 learned feature of PGT is as good as the ViT architecture.

## References

- [1] Stephen E Palmer. Perceptual grouping: It’s later than you think. *Current Directions in Psychological Science*, 11(3):101–106, 2002.
- [2] Johan Wagemans, James H Elder, Michael Kubovy, Stephen E Palmer, Mary A Peterson, Manish Singh, and Rüdiger Von der Heydt. A century of gestalt psychology in visual perception: I. perceptual grouping and figure–ground organization. *Psychological bulletin*, 138(6):1172, 2012.
- [3] Michael H Herzog. Perceptual grouping. *Current Biology*, 28(12):R687–R688, 2018.
- [4] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [5] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022.
- [8] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000.
- [9] Jasper RR Uijlings, Koen EA Van De Sande, Theo Gevers, and Arnold WM Smeulders. Selective search for object recognition. *International journal of computer vision*, 104:154–171, 2013.
- [10] Pablo Arbeláez, Jordi Pont-Tuset, Jonathan T Barron, Ferran Marques, and Jitendra Malik. Multiscale combinatorial grouping. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 328–335, 2014.
- [11] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [12] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [13] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- [14] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [15] Haozhi Qi, Chong You, Xiaolong Wang, Yi Ma, and Jitendra Malik. Deep isometric learning for visual recognition. In *International conference on machine learning*, pages 7824–7835. PMLR, 2020.
- [16] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *Advances in Neural Information Processing Systems*, 33:11525–11538, 2020.
- [17] Geoffrey Hinton. How to represent part-whole hierarchies in a neural network. *Neural Computation*, pages 1–40, 2022.
- [18] Laura Culp, Sara Sabour, and Geoffrey E Hinton. Testing glom’s ability to infer wholes from ambiguous parts. *arXiv preprint arXiv:2211.16564*, 2022.

- 187 [19] Yi Ma, Harm Derksen, Wei Hong, and John Wright. Segmentation of multivariate mixed data  
188 via lossy data coding and compression. *IEEE transactions on pattern analysis and machine*  
189 *intelligence*, 29(9):1546–1562, 2007.
- 190 [20] Gamaleldin Elsayed, Aravindh Mahendran, Sjoerd van Steenkiste, Klaus Greff, Michael C  
191 Mozer, and Thomas Kipf. Savi++: Towards end-to-end object-centric learning from real-world  
192 videos. *Advances in Neural Information Processing Systems*, 35:28940–28954, 2022.
- 193 [21] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong  
194 Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of*  
195 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18134–18144,  
196 2022.
- 197 [22] Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. Slotformer: Unsuper-  
198 vised visual dynamics simulation with object-centric models. In *The Eleventh International*  
199 *Conference on Learning Representations*, 2022.
- 200 [23] Ondrej Biza, Sjoerd van Steenkiste, Mehdi SM Sajjadi, Gamaleldin F Elsayed, Aravindh  
201 Mahendran, and Thomas Kipf. Invariant slot attention: Object discovery with slot-centric  
202 reference frames. *arXiv preprint arXiv:2302.04973*, 2023.
- 203 [24] Michael Chang, Tom Griffiths, and Sergey Levine. Object representations as fixed points:  
204 Training iterative refinement algorithms with implicit differentiation. *Advances in Neural*  
205 *Information Processing Systems*, 35:32694–32708, 2022.
- 206 [25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
207 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.  
208 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
209 *arXiv:2010.11929*, 2020.
- 210 [26] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
211 *arXiv:1312.6114*, 2013.
- 212 [27] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining  
213 Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings*  
214 *of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021.
- 215 [28] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
216 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual  
217 recognition challenge. *International journal of computer vision*, 115:211–252, 2015.
- 218 [29] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,  
219 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings*  
220 *of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.
- 221 [30] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. 2018.
- 222 [31] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
223 for contrastive learning of visual representations. In *International conference on machine*  
224 *learning*, pages 1597–1607. PMLR, 2020.
- 225 [32] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.  
226 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural*  
227 *information processing systems*, 33:9912–9924, 2020.
- 228 [33] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot:  
229 Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.
- 230 [34] Ting Chen, Simon Kornblith, Kevin Swersky, Mohammad Norouzi, and Geoffrey E Hinton. Big  
231 self-supervised models are strong semi-supervised learners. *Advances in neural information*  
232 *processing systems*, 33:22243–22255, 2020.
- 233 [35] Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and Furu Wei. Beit v2: Masked image  
234 modeling with vector-quantized visual tokenizers. *arXiv preprint arXiv:2208.06366*, 2022.

- 235 [36] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
236 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning  
237 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 238 [37] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei  
239 Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. Rethinking semantic segmentation  
240 from a sequence-to-sequence perspective with transformers. In *Proceedings of the IEEE/CVF*  
241 *conference on computer vision and pattern recognition*, pages 6881–6890, 2021.
- 242 [38] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image trans-  
243 formers. *arXiv preprint arXiv:2106.08254*, 2021.

## 244 4 Supplementary Material

### 245 4.1 More about models

#### 246 4.1.1 Algorithm

247 We provide the pseudo code for the perceptual grouping algorithm as below.

---

**Algorithm 1** Multi-grouping operation.

---

```
def multi_grouping(h_key, h_value, steps, mu, sigma, num_tokens, num_heads):  
    """ Input tensors:  
        h_key and h_value are projected multi-head tensors with shape [num_heads x N x d].  
    """  
    # Initial M group tokens.  
    group_tokens = Normal(mean=mu, std=sigma, nsamples=num_tokens)  
    group_tokens = group_tokens.reshape(num_heads, num_tokens, d) #[num_heads x M x d]  
  
    # Binding process.  
    for step in range(steps):  
        # Implicit differentiation.  
        if step == steps - 1:  
            group_tokens = stop_gradient(group_tokens)  
        # Attention operation for group assignment.  
        attn_matrix = attention(group_tokens, h_key) #[num_heads x N x M]  
        attn_matrix /= attn_matrix.sum(-2, keep_dim=True)  
        h_updates = einsum("hij,hid->hd", attn_matrix, h_value) #[num_heads x M x d]  
        group_tokens = gru_cell(h_updates, group_tokens)  
        # Grouped mlp/layernorm performs independent mlp/layernorm for each head.  
        group_tokens = grouped_mlp(grouped_layer_norm(group_tokens)) + group_tokens  
  
    return group_tokens
```

---

#### 248 4.1.2 Self-supervision loss

249 Following the student-teacher self-supervision loss [29, 36], we use a moving average of online  
250 network (student model) as the teacher model to perform representation learning. To summarize  
251 group tokens outputted from the final layer, we use one multi-head attention layer with a learnable  
252 token to attend to all group tokens. The produced single vector is treated as the feature representation  
253 for the image and is input to the loss function.

#### 254 4.1.3 Discussion

255 Our proposed model, perceptual group tokenizer, is free  
256 of self attention operation and relies purely on grouping  
257 operations. In this section, we link the grouping process  
258 to several techniques and discuss the rationale on why the  
259 model can be effective on representation learning.

260 **Group tokens as “communication channels”.** The core  
261 of feature representation learning is how information is ex-  
262 changed among pixels. In perceptual grouping backbones,  
263 we can consider the set of group tokens as communication channels, where information from different  
264 input tokens are aggregated in various ways. Each group token represents a high-order channel that  
265 links input tokens with high affinity under certain projected space to exchange information among  
266 them. As a thought experiment, if each input token is solely assigned to a different group token  
267 (given enough group tokens), then the perceptual grouping layer is equivalent to one self attention  
268 layer (up to some engineering design difference). While self attention layers mainly rely on pairwise  
269 communications, grouping operation, hypothetically, can automatically learn and emerge both pair-  
270 wise and higher-order information exchange through the group token communication channels. This  
271 can also be linked to traditional *factor graphs* in probabilistic graphical models. Through the lens of  
272 that, grouping is forming factor nodes automatically through the learning processes. Under properly  
273 designed loss and grouping operation, it has the potential to be more effective if adopting a per-layer  
274 comparison between self attention and grouping operation.

275 **Efficiency.** Due to the flexibility in customizing number of group tokens (controlled by initial  
276 number of samples), grouping operation does not require a strict  $O(N^2)$  operation and is  $O(NM)$

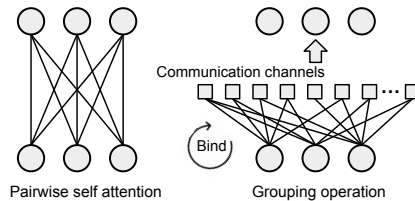


Figure 2: Operation comparison.



	Descend	Flat	Ascend
Token size	[576, 384, 192]	[384, 384, 384]	[192, 384, 576]
Accuracy	62.0	63.1	<b>63.4</b>
Token shape	[192, 128, 64]	[128, 128, 128]	[64, 128, 192]
Accuracy	63.6	<b>63.7</b>	63.1

Table 2: Exploring the design choices for PGT. Token size: dimensions for group tokens in three grouping blocks. Token shape: number of tokens for group tokens in three grouping blocks. Accuracy measured on ImageNet-1K under linear probe protocol. Results indicate progressively large group token dimensions with flat or descend number of tokens arrangements work the best.

277 on complexity. Furthermore, we show that *even in inference time*, number of group tokens can be  
278 adaptively customized, given an already trained model.

## 279 4.2 More about experiments

### 280 4.2.1 Main results

281 **Architecture details.** In the experiments, we mainly evaluate two variants of PGT: the main model  
282 and a tiny version for exploring design choices. On the ImageNet-1K benchmark, we report the  
283 numbers of our main model. Three grouping blocks are used, with 10 grouping layers in each block.  
284 The dimension for input token is 384, with 256 group tokens per layer. The dimensions for group  
285 tokens are 98, 192, and 288 for the three blocks, respectively. There are 6 grouping heads used. For  
286 number of grouping iterations, we observe three rounds are enough. The MLP hidden size for each  
287 layer is 384 as well, i.e. MLP multiplication factor is 1. The final multihead attention layer uses a  
288 learnable token with 2048 dimensions to summarize all group tokens outputs from the model.

### 289 4.2.2 Ablations

290 To explore design choices of PGT, we adopt a tiny version with 3 blocks, 2 layer in each block (6  
291 layers in total), 256 hidden size for input tokens, and 3 number of grouping iterations. The learnable  
292 token in MAP head has 512 dimensions. There are around 10M parameters in a PGT-tiny model.

293 **Group token layouts.** Given a fixed number of budget on group tokens, we explore three choices on  
294 how they should be arranged across grouping blocks and layers: descend, flat and ascend. Intuitively,  
295 more group tokens will have higher capacity of capturing smaller parts and detailed visual features,  
296 while less group tokens are more prone to carry global information. As shown in table 4.2.2 bottom  
297 row, flat or descend number of group tokens performs the best. In practice, we find that using flat  
298 (same number of group tokens in three grouping blocks) version has better stability in training.

299 **Group token dimension shapes.** Similar to token number arrangements, we explore how group  
300 token dimensions should be set. Under three choices, progressively increasing the dimension size in  
301 the later layers performs the best, shown in first row of table 4.2.2. This also aligns with the intuition  
302 that later layers contain more information and requires higher capacity to represent groups.

303 **Multi-grouping vs single grouping.** We further tests whether multi-head grouping helps improve  
304 performance. As a fair comparison, we use 6 heads and 128 group tokens per head for multi-grouping  
305 model, and 1 head with  $6 \times 128$  group tokens for the single grouping model. We find that adopting  
306 multi-head design can improve the performance from 62.2% to 66.3%, a 4.1% accuracy boosts,  
307 showing that having multiple heads indeed helps with representation learning.

308 **Grouping distribution entropy.** Will grouping process collapse to some specific group token  
309 during training? We visualize the entropy of marginal distribution over tokens  $p(c)$  and conditional  
310 distribution  $p(c|x)$  in figure 3 and 4. Interestingly, we observe that conditional probability, i.e.  
311 the assignment to group tokens, tends to become more certain during training, while the marginal  
312 distribution remains having descend entropy. This indicates that collapse does not happen in the  
313 self-supervised training process.

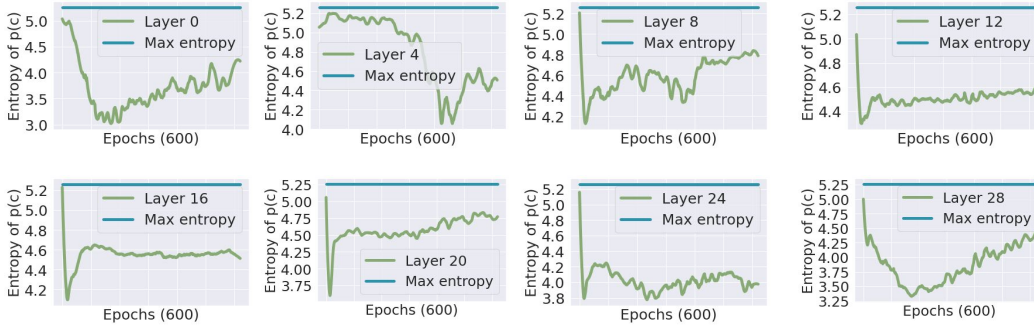


Figure 3: The entropy curves of marginal distribution  $p(c)$  grouping across different layers.

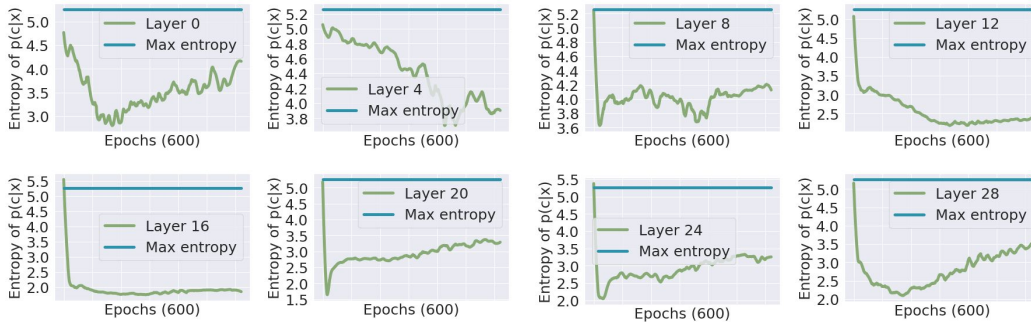


Figure 4: The entropy curves of conditional distribution  $p(c|x)$  grouping across different layers.

### 314 4.2.3 Out-of-distribution adaptive computation

315 One surprising and powerful ability of PGT is adaptive computation. Because the initial group tokens  
 316 are sampled from a Gaussian distribution, the number of group tokens can be flexibly customized  
 317 in inference time given a trained model. This leads to an out-of-distribution adaptive computation  
 318 ability customizable according to needs, e.g. computation or tasks. We mainly test PGT-Tiny with a  
 319 grid evaluation of different number of group tokens used in training and inference. PGT-B model  
 320 with 256 group tokens is also tested under different inference token budgets. Results are summarized  
 321 in table 4.2.3. Our model shows strong out-of-distribution generalizability, indicated by the results.  
 322 Surprisingly, *with more number of tokens, the performance can be increased*. When using the larger  
 323 main model PGT-B to perform adaptive inference, with only 12.5% of the number of group tokens  
 324 compared to training, the performance can still be maintained at 71.18% with only a 8% drop on  
 325 top-1 accuracy.

### 326 4.2.4 Semantic segmentation on ADE20k

327 To evaluate the generalizability of pretrained feature produced by PGT, we test the transfer perfor-  
 328 mance of semantic segmentation with ADE20k. Following the standard setup, we finetune our model  
 329 with the same data augmentation for 128 epoch. The baseline method uses DINO + ViT-B/16, and  
 330 is fine-tuned with a SETR-PUP segmentation head [37]. For our model, we only add one linear  
 331 classification layer after the pre-trained PGT for fine-tuning. To adapt to more objects and complex  
 332 scenes in the segmentation datasets, we use 1024 group tokens for inference, benefiting from the  
 333 adaptive computation ability of our model. We find that our model (PGT-B + Linear) can obtain  
 334 44.5% on mean IoU while the baseline (ViT-B/16 + SETR-PUP) achieves 44.1% [38], leading to a  
 335 0.4% improvements.

tr/inf	16	32	64	128	256	384
PGT-Ti-16	57.44 ( $\times 1$ )	58.28 ( $\times 2$ )	<b>58.49</b> ( $\times 4$ )	58.47 ( $\times 8$ )	58.48 ( $\times 16$ )	58.39 ( $\times 24$ )
PGT-Ti-32	57.33 ( $\times \frac{1}{2}$ )	<u>59.86</u> ( $\times 1$ )	60.82 ( $\times 2$ )	<b>61.01</b> ( $\times 4$ )	60.99 ( $\times 8$ )	60.89 ( $\times 12$ )
PGT-Ti-64	53.02 ( $\times \frac{1}{4}$ )	59.20 ( $\times \frac{1}{2}$ )	<u>61.68</u> ( $\times 1$ )	62.55 ( $\times 2$ )	62.91 ( $\times 4$ )	<b>62.92</b> ( $\times 6$ )
PGT-Ti-128	44.86 ( $\times \frac{1}{8}$ )	56.63 ( $\times \frac{1}{4}$ )	61.80 ( $\times \frac{1}{2}$ )	<u>63.88</u> ( $\times 1$ )	64.66 ( $\times 2$ )	<b>64.82</b> ( $\times 3$ )
PGT-Ti-256	27.20 ( $\times \frac{1}{16}$ )	47.35 ( $\times \frac{1}{8}$ )	58.80 ( $\times \frac{1}{4}$ )	<u>63.32</u> ( $\times \frac{1}{2}$ )	65.13 ( $\times 1$ )	<b>65.52</b> ( $\times \frac{2}{3}$ )
PGT-Ti-384	26.09 ( $\times \frac{1}{24}$ )	43.02 ( $\times \frac{1}{12}$ )	55.35 ( $\times \frac{1}{6}$ )	61.67 ( $\times \frac{1}{3}$ )	64.57 ( $\times \frac{2}{3}$ )	<b>65.50</b> ( $\times 1$ )
PGT-B-256	58.18 ( $\times \frac{1}{16}$ )	71.18 ( $\times \frac{1}{8}$ )	76.51 ( $\times \frac{1}{4}$ )	78.39 ( $\times \frac{1}{2}$ )	<u>79.16</u> ( $\times 1$ )	<b>79.47</b> ( $\times \frac{2}{3}$ )

Table 3: Out-of-distribution adaptive computation by selecting different numbers of initially sampled tokens. Row: number of tokens used for training. Column: number of tokens used for inference. Top-1 accuracy is reported under linear evaluation protocol using ImageNet-1K. The reported performance of first six rows is obtained using a tiny version of PGT, and last row is the main model. Number of group tokens is the same for underlined numbers in training and inference. **Bold numbers** are the best results.

### 336 4.2.5 Grouping visualization

337 We further visualize the generated attention maps during grouping processes to inspect the behaviour  
338 of grouping operations. In figure 5, the patch to group token attention maps across all grouping  
339 iterations are shown. We find that even the first iteration step can sometimes generate meaningful  
340 attention maps. With more iterations, attention maps are more focused on meaningful regions. Figure  
341 6 shows attention maps across different layers of PGT-B. We observe that early layers tend to capture  
342 fine-grained elements, while the last layer focuses on semantic information. Multiple grouping heads  
343 are indeed capturing various ways of grouping image features, for example, in the first image, first  
344 group focuses on color and lights, second head relies on spatial cues, while the last one potentially  
345 captures textures.

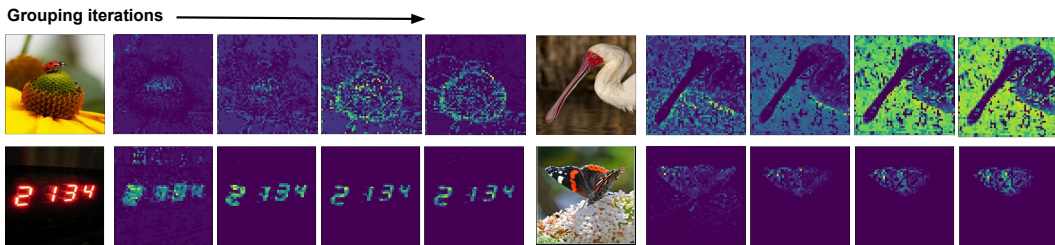


Figure 5: Attention maps produced along grouping processes in the last layer of PGT.

### 346 4.3 Conclusion

347 In this paper, we propose Perceptual Group Tokenizer (PGT), a new visual recognition architecture  
348 entirely built through perceptual grouping principles. The proposed model shows strong performance  
349 on self-supervised learning benchmark ImageNet-1K with linear probe evaluation, and has desirable  
350 properties such as adaptive computation and interpretability in each operation. This work potentially  
351 opens a new paradigm for designing visual recognition backbones and hopes to inspire more research  
352 progress along this direction. One limitation of the proposed model is its relatively expensive  
353 computation cost due to iterative binding processes. This can be potentially addressed by other  
354 grouping operations, for example grouping operations with closed-form solutions. We leave this to  
355 future works.

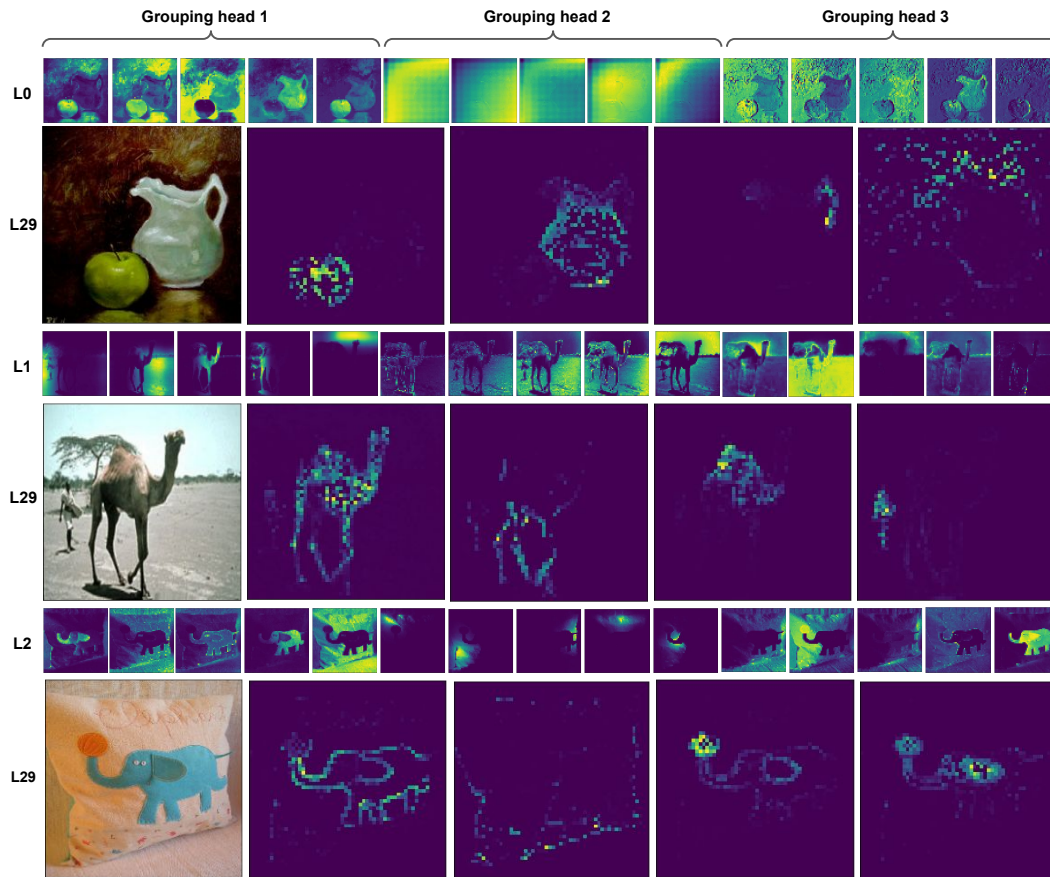


Figure 6: Visualization of attention maps of each group tokens across layers and grouping head.  $L$  indicates layer indices. Five group tokens for each grouping head. Smaller images are for early layers, arranged as five group tokens per grouping head. Zoom in for better viewing. Large images are for the last layer.