

# Top-Down Visual Attention from Analysis by Synthesis

Baifeng Shi  
UC Berkeley

Trevor Darrell  
UC Berkeley

Xin Wang  
Microsoft Research

## Abstract

Current attention algorithms (e.g., self-attention) are stimulus-driven and highlight all the salient objects in an image. However, intelligent agents like humans often guide their attention based on the high-level task at hand, focusing only on task-related objects. This ability of task-guided top-down attention provides task-adaptive representation and helps the model generalize to various tasks. In this paper, we consider top-down attention from a classic Analysis-by-Synthesis (AbS) perspective of vision. Prior work indicates a functional equivalence between visual attention and sparse reconstruction; we show that an AbS visual system that optimizes a similar sparse reconstruction objective modulated by a goal-directed top-down signal naturally simulates top-down attention. We further propose Analysis-by-Synthesis Vision Transformer (AbSViT), which is a top-down modulated ViT model that variationally approximates AbS, and achieves controllable top-down attention. For real-world applications, AbSViT consistently improves over baselines on Vision-Language tasks such as VQA and zero-shot retrieval where language guides the top-down attention. AbSViT can also serve as a general backbone, improving performance on classification, semantic segmentation, and model robustness. Project page: <https://sites.google.com/view/absvit>.

## 1. Introduction

Human visual attention is often *task-guided*, i.e., we tend to focus on different objects when processing different tasks [7, 75]. For example, when we answer different questions about one image, we only attend to the objects that are relevant to the question (Fig. 1 (b-c)). This stands in contrast with the widely-used self-attention [18], which is completely *stimulus-driven*, i.e., it highlights all the salient objects in the image without task-guided selection (Fig. 1 (a)). While the stimulus-driven bottom-up attention has shown promising results in visual representation learning [6], current vision transformers still lack the ability of task-guided top-down attention, which provides task-adaptive representation and

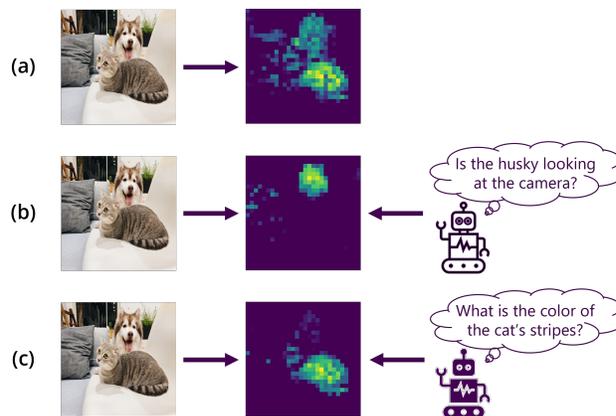


Figure 1. **Top-down vs. bottom-up attention.** (a) Bottom-up attention is stimulus-driven, i.e., any salient objects (dog and cat) in the image may attract attention. (b-c) Top-down attention is task-guided. For example, when the task is to answer a question about a specific object, the attention will only center on that object and ignore the others. In this way, a more focused representation can be extracted for the current goal.

potentially improves task-specific performances [1, 69, 71]. Although some algorithms of top-down attention are proposed in the literature [1, 9, 49, 69, 71], they are incompatible with self-attention-based transformers and principled and unified designs are still missing.

Previous work [5, 10, 36, 37, 53] has studied the mechanism of top-down attention in human vision systems, hypothesizing top-down attention is a result of the human visual system performing Analysis by Synthesis (AbS). AbS [34, 74] is a classic idea that suggests the human visual perception depends on both the input image and a high-level prior about the latent cause of the image, and different priors can lead to different ways to perceive the same image (e.g., visual illusion [35] and bistable perception [58]). This is formulated as Bayesian inference  $\max_{\mathbf{z}} p(\mathbf{h}|\mathbf{z})p(\mathbf{z})$ , where  $\mathbf{h}$  is the input image, and  $\mathbf{z}$  is the latent representation. It is hypothesized that the high-level goal can be formulated as a prior to direct the low-level recognition of different objects through AbS, achieving top-down attention. Still, existing works [10, 46, 73] are conceptual and hardly guide model designs in practice.

In this work, we present a novel perspective on how AbS entails top-down attention, followed by a new Analysis-by-Synthesis Vision Transformer (AbSViT) based on the findings. We start from previous work [59], which shows that visual attention (*e.g.*, self-attention) is functionally equivalent to sparse reconstruction which reconstructs the input using a dictionary containing templates of separate objects in the input. We show that AbS optimizes a similar *sparse reconstruction* objective modulated by a top-down signal. The top-down signal depends on the prior and acts as a preference on which object templates to choose to reconstruct the input. Therefore, only the objects consistent with the high-level prior are selected, equivalent to top-down attention.

Inspired by the connection, we propose AbSViT, a ViT [18] model with prior-conditioned top-down modulation trained to approximate AbS in a variational way. AbSViT contains a feedforward (encoding) and a feedback (decoding) pathway. The feedforward path is a regular ViT, and the feedback path contains linear decoders for each layer. Each inference starts with an initial feedforward run. The output tokens are manipulated by the prior and fed back through the decoders to each self-attention module as top-down input for the final feedforward pass (Fig. 3).

When only pretrained on ImageNet [16], which contains mostly single-object images, AbSViT can attend to different objects in multi-object scenes controllably. For real-world applications, we observe consistent improvements from AbSViT on Vision-Language tasks such as VQA [3] and zero-shot image retrieval, where language is used as a prior to guide attention. For tasks without a strong prior, such as ImageNet classification and semantic segmentation, AbSViT can also serve as a general backbone and achieve substantial improvements. Additionally, the object-centric representation resulting from the top-down attention design enables better generalization to corrupted, adversarial, and out-of-distribution images. We hope this work can encourage future exploration of task-guided attention designs and visual representation learning.

## 2. Related Work

**Top-down visual attention** endows us with the crucial ability to selectively collect information related to the behavioral goal. Several attempts have been made towards understanding the mechanism of top-down attention from experimental observations such as multiplicative tuning [45] and contrast responses [44, 55] in V4, and extra-classical receptive fields in V1 [2, 8, 57]. Other work tries to build a principled computational model for top-down attention [10, 46, 73].

Top-down attention has also found numerous applications in computer vision tasks where additional guidance (*e.g.*, language) is available aside from the image. Previous work employs top-down attention for object detection [48], image

captioning [71], and visual question answering [1, 69]. However, these algorithms are either incompatible with current self-attention-based models or show inferior performance, as indicated by our experiments. Other work [18, 41, 42, 72] uses a feedforward model that takes both image and the high-level guidance (*e.g.*, text tokens or [cls] token) as input, which we show is suboptimal compared to our top-down model design. Dou *et al.* [19] propose to extract image and text features with separate encoders and combine them with a multi-modal fusion module during vision-language pre-training, which works better than using a single multi-modal feedforward model on vision language tasks. However, in this way, the visual encoder is still bottom-up. We show that augmenting it with the proposed top-down attention further improves model performance on standard benchmarks.

### **Top-down attention explained as Analysis by Synthesis.**

Analysis by Synthesis (AbS) is hypothesized as a potential computational model behind top-down attention. Lee [36] starts from a Bayesian inference perspective and explains the top-down modulation in examples such as illusionary contours and shapes from shading. Yu and Dayan [73] focus on the top-down attention in Posner’s task [51] and build a hierarchical model where each layer corresponds to a computational step of Bayesian inference. Subsequent work [10, 53] assumes each object is generated by an appearance variable and a location variable and uses Bayesian inference to perform spatial attention and feature attention. Borji *et al.* [5] adopt a Dynamic Bayesian Network to simulate eye fixation in top-down attention. However, these models do not apply to practical designs in modern deep learning.

**Generative model for discriminative learning.** It has been widely explored in using generative models to assist discriminative learning. Specifically, the belief that representation with strong generative capability can better capture the structure of visual signals has inspired numerous unsupervised learning algorithms, from the early Restricted Boltzmann Machine [29, 30] and Helmholtz Machine [15], to the following auto-encoder models such as DAE [62] and VAE [33]. Recent work [23, 60] has shown impressive results on generative unsupervised learning. Generative models can also help with supervised learning, *e.g.*, by refining object detection [38] or detecting errors in semantic segmentation [66]. Feedforward models with generative feedback are also more robust to input corruptions [31]. In our work, AbSViT also contains a generative feedback path that is able to refine the intermediate representation and attention and thus improves the performance.

## 3. Preliminaries: Attention as Sparse Reconstruction

Shi *et al.* [59] show that a sparse reconstruction (SR) module functionally resembles visual attention. An SR mod-

ule takes an input  $\mathbf{x} \in \mathbb{R}^d$  and outputs  $\mathbf{z} = \mathbf{P}\tilde{\mathbf{u}}^*$  where  $\mathbf{P} \in \mathbb{R}^{d \times d'}$  is the dictionary and  $\tilde{\mathbf{u}}^*$  is the sparse code, *i.e.*,

$$\tilde{\mathbf{u}}^* = \arg \min_{\tilde{\mathbf{u}} \in \mathbb{R}^{d'}} \frac{1}{2} \|\mathbf{P}\tilde{\mathbf{u}} - \mathbf{x}\|_2^2 + \lambda \|\tilde{\mathbf{u}}\|_1. \quad (1)$$

Each atom (column) of  $\mathbf{P}$  contains a template pattern and each element in  $\tilde{\mathbf{u}}$  is the activation of the corresponding template. The objective is to reconstruct the input using as few templates as possible. To solve Eq. (1), one may adopt a first-order optimization [56, 59] with dynamics at time  $t$  of

$$\frac{d\mathbf{u}}{dt} \propto -\mathbf{u} - (\mathbf{P}^T \mathbf{P} - \mathbf{I})\tilde{\mathbf{u}} + \mathbf{P}^T \mathbf{x}, \quad (2)$$

where the optimization is over an auxiliary variable  $\mathbf{u}$  and  $\tilde{\mathbf{u}} = g_\lambda(\mathbf{u}) = \text{sgn}(\mathbf{u})(|\mathbf{u}| - \lambda)_+$  with  $\text{sgn}(\cdot)$  as the sign function and  $(\cdot)_+$  as ReLU. Here  $\mathbf{u}$  is activated by the template matching  $\mathbf{P}^T \mathbf{x}$  between the dictionary and the input, and different elements in  $\mathbf{u}$  inhibit each other through  $-(\mathbf{P}^T \mathbf{P} - \mathbf{I})\tilde{\mathbf{u}}$  to promote sparsity.

To see the connection between visual attention and sparse reconstruction, recall that attention in the human visual system is achieved via two steps [17]: (i) *grouping* features into separate objects or regions, and (ii) *selecting* the most salient objects or regions while repressing the distracting ones. A similar process is also happening in SR, *i.e.*, if each atom in  $\mathbf{P}$  is a template of every single object, then each element in  $\mathbf{u}$  groups the input features belonging to that object through  $\mathbf{P}^T \mathbf{x}$ , while the sparsity constraint promoted by the lateral inhibition  $-(\mathbf{P}^T \mathbf{P} - \mathbf{I})\tilde{\mathbf{u}}$  selects the object that is most activated. As shown in [59], SR modules achieve similar attention effects as self-attention (SA) [61] while being more robust against image corruptions.

Interestingly, it is also pointed out in [59] that under certain constraints (*e.g.*, the key and query transform is the same), SA can be viewed as solving a similar SR problem but without sparsity. After adding the sparsity back, SA is an approximation of

$$\begin{cases} \tilde{\mathbf{U}}^* = \arg \min_{\tilde{\mathbf{U}}} \frac{1}{2} \|\Phi(\mathbf{K})\tilde{\mathbf{U}} - \mathbf{V}\|_2^2 + \lambda \|\tilde{\mathbf{U}}\|_1, \\ \mathbf{Z} = \Phi(\mathbf{Q})\tilde{\mathbf{U}}^*, \end{cases} \quad (3)$$

where  $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{(hw) \times c}$  are the query, key, and value matrices,  $\Phi(\mathbf{Q}), \Phi(\mathbf{K}) \in \mathbb{R}^{(hw) \times d'}$  are the random features [11] that approximate the softmax kernel  $\Phi(\mathbf{Q})_i \Phi(\mathbf{K})_j^T \approx e^{\mathbf{Q}_i \mathbf{K}_j^T}$ ,  $\tilde{\mathbf{U}}^* \in \mathbb{R}^{d' \times c}$  is the sparse code and  $\mathbf{Z}$  is the output. This provides a novel perspective on the mechanism of SA, *i.e.*, it is solving a channel-wise sparse reconstruction of the value matrix  $\mathbf{V}$  using an input-dependent dictionary  $\Phi(\mathbf{K})$ . Visualization of  $\Phi(\mathbf{K})$  shows each atom contains a mask for one single object or region, which means that SA is trying to reconstruct the input with as few masks as possible, thus only the salient objects are selected and highlighted (Fig. 2 (a)).

## 4. Top-Down Attention from AbS

We consider top-down visual attention from an Analysis by Synthesis (AbS) view of vision. We start from the hierarchical AbS formulation of visual perception (Sec. 4.1) and show that it is equivalently optimizing a sparse reconstruction objective that is modulated by a top-down signal, thus entailing top-down attention (Sec. 4.2).

### 4.1. Hierarchical AbS

AbS formulates visual perception as a Bayesian inference process. Given the image generation process  $p(\mathbf{h}|\mathbf{z})$  and a prior  $p(\mathbf{z})$ , where  $\mathbf{h}$  is the image and  $\mathbf{z}$  is the latent code, AbS finds  $\mathbf{z}^* = \arg \max_{\mathbf{z}} p(\mathbf{h}|\mathbf{z})p(\mathbf{z})$ .

In this work, we assume the generation is hierarchical, *i.e.*,  $\mathbf{z}_L \rightarrow \mathbf{z}_{L-1} \rightarrow \dots \rightarrow \mathbf{z}_1 \rightarrow \mathbf{h}$ , where  $\mathbf{z}_\ell$  is the latent at  $\ell$ -th layer. The MAP estimation is

$$\mathbf{z}_L^*, \dots, \mathbf{z}_1^* = \arg \max_{\mathbf{z}_L, \dots, \mathbf{z}_1} p(\mathbf{h}|\mathbf{z}_1) \dots p(\mathbf{z}_{L-1}|\mathbf{z}_L)p(\mathbf{z}_L). \quad (5)$$

For each generation process  $\mathbf{z}_{\ell+1} \rightarrow \mathbf{z}_\ell$  between layer  $\ell$  and  $\ell + 1$ , we further assume that  $\mathbf{z}_\ell$  is constructed by a sparse code  $\tilde{\mathbf{u}}_\ell$  which is generated from  $\mathbf{z}_{\ell+1}$  via a non-linear function  $g_\ell(\cdot)$ , *i.e.*,

$$\begin{aligned} \tilde{\mathbf{u}}_\ell &\sim p(\tilde{\mathbf{u}}_\ell|\mathbf{z}_{\ell+1}) \propto \exp\left\{-\frac{1}{2}\|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{z}_{\ell+1})\|_2^2 - \lambda \|\tilde{\mathbf{u}}_\ell\|_1\right\} \\ \mathbf{z}_\ell &= \mathbf{P}_\ell \tilde{\mathbf{u}}_\ell, \end{aligned} \quad (7)$$

where  $\mathbf{P}_\ell$  is the dictionary. Intuitively, it first generates  $g_\ell(\mathbf{z}_{\ell+1})$  as a blurry and noisy version of  $\mathbf{z}_\ell$ , then find the sparse code  $\tilde{\mathbf{u}}_\ell$  to construct a sharper and cleaner version.

Since  $\mathbf{z}_\ell$  is decided by  $\tilde{\mathbf{u}}_\ell$ , it suffices to optimize the MAP estimation over  $\{\tilde{\mathbf{u}}_\ell\}_{\ell=1}^L$ , *i.e.*,

$$\tilde{\mathbf{u}}_L^*, \dots, \tilde{\mathbf{u}}_1^* = \arg \max_{\tilde{\mathbf{u}}_L, \dots, \tilde{\mathbf{u}}_1} p(\mathbf{h}|\tilde{\mathbf{u}}_1) \dots p(\tilde{\mathbf{u}}_{L-1}|\tilde{\mathbf{u}}_L)p(\tilde{\mathbf{u}}_L). \quad (8)$$

Solving Eq. (8) by simple gradient ascent (of the logarithm) gives the dynamics

$$\frac{d\tilde{\mathbf{u}}_\ell}{dt} \propto \nabla_{\tilde{\mathbf{u}}_\ell} \log p(\tilde{\mathbf{u}}_{\ell-1}|\tilde{\mathbf{u}}_\ell) + \nabla_{\tilde{\mathbf{u}}_\ell} \log p(\tilde{\mathbf{u}}_\ell|\tilde{\mathbf{u}}_{\ell+1}) \quad (9)$$

where  $\tilde{\mathbf{u}}_\ell$  is affected by both  $\tilde{\mathbf{u}}_{\ell-1}$  and  $\tilde{\mathbf{u}}_{\ell+1}$ .

### 4.2. Top-Down Attention from AbS

From AbS (Eq. (6-9)) we can derive the dynamics of  $\tilde{\mathbf{u}}_\ell$  as

$$\frac{d\tilde{\mathbf{u}}_\ell}{dt} \propto \nabla_{\tilde{\mathbf{u}}_\ell} \left( -\frac{1}{2} \|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - (\mathbf{x}_\ell^{bu} + \mathbf{x}_\ell^{td})\|_2^2 - \lambda \|\tilde{\mathbf{u}}_\ell\|_1 - r_\ell(\tilde{\mathbf{u}}_\ell) \right) \quad (10)$$

where  $\mathbf{x}_\ell^{td} = g_\ell(\mathbf{z}_{\ell+1})$  is the top-down signal and  $\mathbf{x}_\ell^{bu} = f_\ell(\mathbf{z}_{\ell-1}) = \mathbf{J}_{g_{\ell-1}}^T \mathbf{z}_{\ell-1}$  is the bottom-up signal where  $\mathbf{J}_{g_{\ell-1}}$  is the jacobian of  $g_{\ell-1}(\mathbf{P}_{\ell-1} \tilde{\mathbf{u}}_{\ell-1})$ , and  $r_\ell(\tilde{\mathbf{u}}_\ell) = \|g_{\ell-1}(\mathbf{P}_{\ell-1} \tilde{\mathbf{u}}_{\ell-1})\|_2^2$  is an additional regularization. Details of the derivation are pushed back to Appendix. One may notice from Eq. (10) that, in AbS each layer is solving a similar

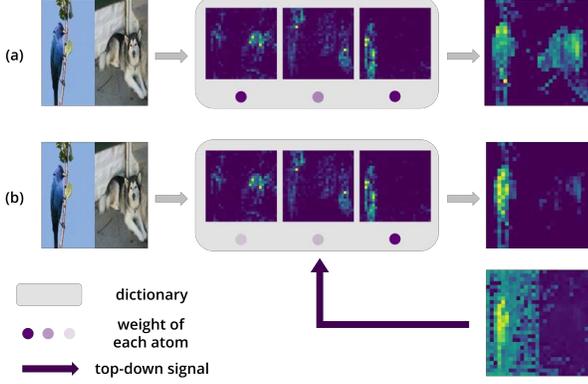


Figure 2. (a) Each atom in the dictionary contains masks for separate objects or regions. The sparse reconstruction tries to use as few masks as possible to reconstruct the input feature map, thus only the salient objects are highlighted. (b) The top-down signal  $\mathbf{x}_\ell^{td}$  puts a bias on the weights of the atoms so that only the objects that agree with  $\mathbf{x}_\ell^{td}$  are selected.

sparse reconstruction problem as in Eq. (1) but with the input of  $\mathbf{x}_\ell^{bu} + \mathbf{x}_\ell^{td}$ , thus simulating attention that is modulated by both bottom-up and top-down signals. This can also be observed by turning Eq. (10) into

$$\frac{d\tilde{\mathbf{u}}_\ell}{dt} \propto -\mathbf{u}_\ell - (\mathbf{P}_\ell^T \mathbf{P}_\ell - \mathbf{I})\tilde{\mathbf{u}}_\ell + \mathbf{P}_\ell^T \mathbf{x}_\ell^{bu} + \mathbf{P}_\ell^T \mathbf{x}_\ell^{td} - \nabla r_\ell(\tilde{\mathbf{u}}_\ell). \quad (11)$$

Comparing with Eq. (2), here  $\tilde{\mathbf{u}}_\ell$  is steered by an additional term  $\mathbf{P}_\ell^T \mathbf{x}_\ell^{td}$  that acts as a bias on which atom in  $\mathbf{P}_\ell$  to choose. For example, if atoms in  $\mathbf{P}_\ell$  are templates of separate objects (like in self-attention), then  $\mathbf{P}_\ell^T \mathbf{x}_\ell^{td}$  highlights the objects that are consistent with the top-down signal (Fig. 2 (b)).

This implies an AbS system naturally entails top-down attention. Intuitively, the prior reflects which objects the output  $\mathbf{z}_L$  should highlight. Then the affected  $\mathbf{z}_L$  is fed back to layer  $L - 1$  through  $g_{L-1}$ , as a top-down signal to direct which objects to select in layer  $L - 1$ . The same process repeats until the first layer. Different priors will direct the intermediate layers to select different objects, achieving top-down attention.

Interestingly, if we consider the analogy between self-attention and sparse reconstruction, Eq. (10) leads to a smooth way of building a top-down version of self-attention, *i.e.*, we only need to add a top-down signal to the value  $\mathbf{V}$ , while keeping other parts such as  $\mathbf{Q}$  and  $\mathbf{K}$  (which decides the dictionary) untouched. We will make it clearer in Sec. 5.

## 5. Analysis-by-Synthesis Vision Transformer

Inspired by the connection between top-down attention and AbS, we propose to achieve top-down attention by building a vision transformer that performs AbS (Eq. (5)), *i.e.*, if the network has input  $\mathbf{h}$  and latent representation  $\mathbf{z}_\ell$  after each layer  $\ell$  (which means  $\mathbf{z}_L$  is the output), the final

latent representation should approximate  $\mathbf{z}_1^*, \dots, \mathbf{z}_L^*$ . Since directly solving Eq. (5) requires an iterative optimization which would be extremely costly, in this work, we adopt a variational approximation to Eq. (5). Specifically, we optimize a variational loss

$$\begin{aligned} \mathcal{L}_{var} &= -\sum_{\ell=0}^{L-1} \log p(\mathbf{z}_\ell | \mathbf{z}_{\ell+1}) - \log p(\mathbf{z}_L) \\ &= \sum_{\ell=0}^{L-1} \left( \frac{1}{2} \|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{z}_{\ell+1})\|_2^2 + \lambda \|\tilde{\mathbf{u}}_\ell\|_1 \right) - \log p(\mathbf{z}_L) \end{aligned} \quad (12)$$

where  $\mathbf{z}_0 = \mathbf{h}$ . However, as stated below, there are several caveats we need to work around when training a network with Eq. (12) in real-world tasks.

**The sparsity regularization.** Since the practical model we build in this work is based on self-attention (Sec. 5.1), which neither has a sparsity constraint nor solves the SR explicitly [59], we remove the sparsity regularization by setting  $\lambda = 0$ , which makes  $-\log p(\mathbf{z}_\ell | \mathbf{z}_{\ell+1}) = \frac{1}{2} \|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{z}_{\ell+1})\|_2^2 = \frac{1}{2} \|\mathbf{z}_\ell - g_\ell(\mathbf{z}_{\ell+1})\|_2^2$ .

**Jointly training the decoder  $g_\ell$ .** Normally, optimizing Eq. (12) requires knowing the generation process  $g_\ell$  beforehand, which in our case is unknown. This can be addressed by training  $g_\ell$  jointly with the whole network, similar to VAE [33]. It is natural to use  $g_\ell$  also as the feedback path of the network, as shown in Sec. 5.1.

**Trade-off between the generative and discriminative power.** The variational loss forces each  $\mathbf{z}_{\ell+1}$  to be capable of generating  $\mathbf{z}_\ell$ . However, we find empirically that enforcing a strong generative power on the feature will harm its discriminative power in the setting of supervised learning. To address this, for each term  $-\log p(\mathbf{z}_\ell | \mathbf{z}_{\ell+1})$  we stop the gradient on  $\mathbf{z}_\ell$  and  $\mathbf{z}_{\ell+1}$ , *i.e.*,  $-\log p(\mathbf{z}_\ell | \mathbf{z}_{\ell+1}) = \frac{1}{2} \|sg(\mathbf{z}_\ell) - g_\ell(sg(\mathbf{z}_{\ell+1}))\|_2^2$ , where  $sg(\cdot)$  is stop-gradient. In this way, only the decoder  $g_\ell$  receives the gradient.

**The variable prior.** Rigorously speaking, variational methods only approximate AbS with a fixed prior  $p(\mathbf{z}_L)$ . However, top-down attention should be able to flexibly attend to different objects by changing different priors. The question is, how can we learn a variational model that generalizes to different priors? In this work, we adopt a simple trick called Meta-amortized VI [65]. Concretely, we assume the prior  $p_\xi(\mathbf{z}_L)$  depends on some parameter  $\xi$ , which can be a sentence or a class prototype cueing what objects to look at in the image. Then we make the model adaptable to  $\xi$  during inference to approximate AbS with prior  $p_\xi(\mathbf{z}_L)$  given any  $\xi$ . See the design details in Sec. 5.1.

After applying these tricks, our variational loss becomes

$$\mathcal{L}_{var} = \frac{1}{2} \sum_{\ell=0}^{L-1} \|sg(\mathbf{z}_\ell) - g_\ell(sg(\mathbf{z}_{\ell+1}))\|_2^2 - \log p_\xi(\mathbf{z}_L), \quad (13)$$

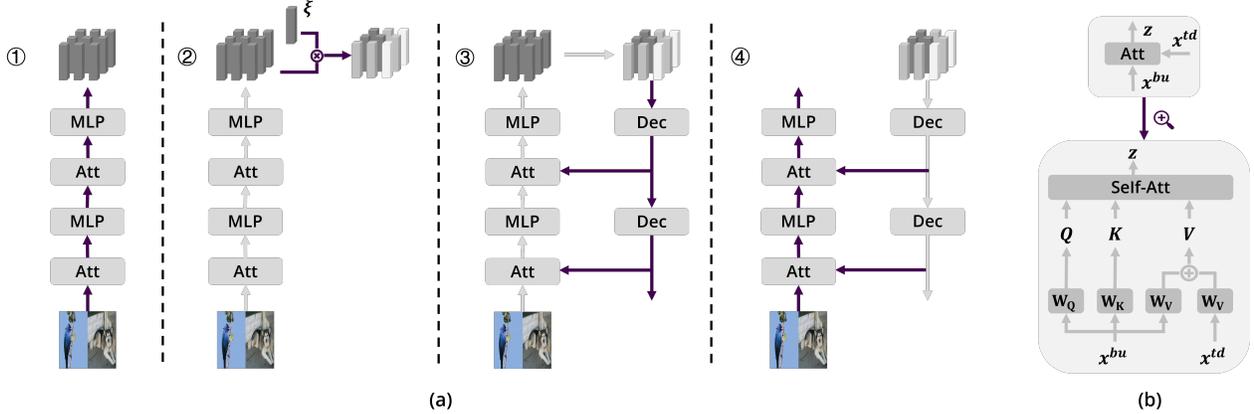


Figure 3. **Design of AbSViT.** (a) Four steps to every single inference. The operations in each step are colored as purple and others as gray. AbSViT first passes the image through the feedforward path. The output tokens are then reweighted by their similarity with the prior vector  $\xi$  and fed back through the decoders to each self-attention module as the top-down input for the final feedforward run. (b) The top-down input to self-attention is added to the value matrix while other parts stay the same.

which contains layer-wise reconstruction loss and a prior loss. We also try cosine similarity instead of  $\ell_2$  distance for reconstruction and get similar results. In Sec. 5.1, we will show how to build a ViT with prior-conditioned top-down modulation and train it with Eq. (13).

### 5.1. AbSViT Design

Fig. 3 (a) shows the proposed AbSViT which is built upon ViT [18]. Every single inference consists of 4 steps: (i) pass the image through the feedforward encoder, (ii) modulate the output tokens with a prior vector  $\xi$ , (iii) send the tokens back through the feedback decoder to intermediate layers, and (iv) run the feedforward path again but with each self-attention layer also receiving the top-down tokens as input.

Within the whole pipeline, the feedforward encoder has the same architecture as regular ViT. For the feedback path, we use a single token-wise linear transform for each layer-wise decoder  $g_\ell$ . The design of token modulation with prior  $\xi$  and the self-attention with top-down input are introduced below:

**Design of token modulation with  $\xi$ .** The purpose is to modify the tokens to carry the information about the prior  $p_\xi$  when fed back to the network. The prior is parameterized by  $\xi$ , which may be a language embedding or a class prototype telling the network which objects to look at. Therefore, we instantiate the modulation as a simple spatial reweighting, *i.e.*,  $\mathbf{z}_L^i \rightarrow \alpha \cdot \text{sim}(\xi, \mathbf{z}_L^i) \cdot \mathbf{z}_L^i$ , where  $\mathbf{z}_L^i$  is the  $i$ -th output token,  $\text{sim}$  is the cosine similarity clamped to  $[0, 1]$ , and  $\alpha$  is a scaling factor controlling the scale of the top-down signal, which is set to 1 by default. In this way, only the tokens with high similarity to  $\xi$  are sent back, and others are (softly) masked out. Note that the design here is for simplicity and may not be suitable for general usage. For example, when

dealing with transparent images where two objects overlap, spatial reweighting cannot separate two objects away.

**Design of self-attention with top-down input.** From the analogy between self-attention and sparse reconstruction (Eq. (3)), the value matrix in SA corresponds to the reconstructed input signal, and the query and key serve as the dictionary. Since the top-down attention in AbS (Eq. (10)) adds a top-down signal to the input while keeping the dictionary untouched, it is natural to design the top-down version of self-attention by simply adding the top-down signal to the value and keep query and key as the same, as illustrated in Fig. 3 (b). We will show in Sec. 6.5 that this is better than an arbitrary design where we add the top-down signal to the query, key, and value.

In this paper, we focus on supervised learning and train the model on two types of tasks. One is Vision-Language (V&L) tasks such as VQA and zero-shot image retrieval, where the language acts as a prior to cue the model where to look at. The other one is image understanding, such as ImageNet classification and semantic segmentation, which do not have a specific prior. When training the network, we optimize the supervised loss as well as the variational loss (Eq. (13)), *i.e.*,

$$\mathcal{L} = \frac{1}{2} \sum_{\ell=1}^L \|sg(\mathbf{z}_\ell) - g_\ell(sg(\mathbf{z}_{\ell+1}))\|_2^2 - \log p_\xi(\mathbf{z}_L) + \mathcal{L}_{sup}, \quad (14)$$

where  $\mathbf{z}_\ell$  is the  $\ell$ -th layer’s output after the whole inference cycle,  $sg$  is stop-gradient, and  $g_\ell$  is the  $\ell$ -th layer’s decoder. The form of prior  $p_\xi$  depends on the task. For V&L tasks,  $\xi$  is the text embedding and we use a CLIP-style prior [52]:

$$p_\xi(\mathbf{z}_L) = \frac{\exp\{\xi^T \mathbf{z}_L\}}{\exp\{\xi^T \mathbf{z}_L\} + \sum_k \exp\{\xi^T \mathbf{z}_k\}}, \quad (15)$$

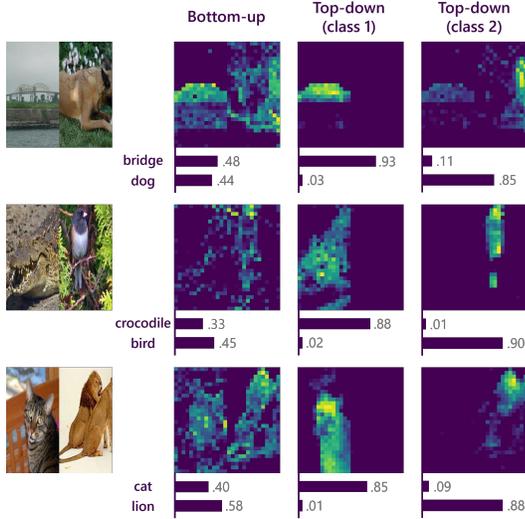


Figure 4. Controllable top-down attention in multi-object images. For each image, bottom-up attention will highlight both objects. In contrast, we can use different class prototypes as the prior to control the top-down attention to focus on different objects, and the classification result also changes accordingly.

where the negative samples  $z_-^k$  are the output from other images. For image classification and segmentation where no specific prior is available, we set  $\xi$  as a trainable query vector that is independent of the input image, and we choose an uninformative prior that does not contribute to the gradient, *i.e.*,  $\nabla \log p_\xi(\mathbf{z}_L) = 0$ .

## 6. Experiments

In this section, we first show that AbSViT achieves controllable top-down attention in multi-object scenes (Sec. 6.1). Then we test AbSViT on Vision-Language tasks such as VQA and zero-shot image retrieval (Sec. 6.2), and also on ImageNet classification and model robustness (Sec. 6.3), as well as semantic segmentation (Sec. 6.4). Finally, we analyze specific designs of AbSViT in Sec. 6.5.

**Datasets.** For VQA, we use VQAv2 [22] for training and testing and compare the attention map with human attention collected by VQA-HAT [14]. For zero-shot image retrieval, we use Flickr30K [50]. For image classification, we train and test on ImageNet-1K (IN) [16], and also test on corrupted images from IN-C [25], adversarial images from IN-A [27], and out-of-distribution images from IN-R [26] and IN-SK [63]. For semantic segmentation, we test on PASCAL VOC [21], Cityscapes [13], and ADE20K [76].

**Experimental setup.** We compare several baselines for goal-directed attention: (i) **PerceiverIO** [32] uses  $e_\xi(\cdot)$  to reweight the tokens from feedforward output just like in AbSViT, but directly outputs the reweighted tokens without any feedback, (ii) **MaskAtt** uses the same soft mask for

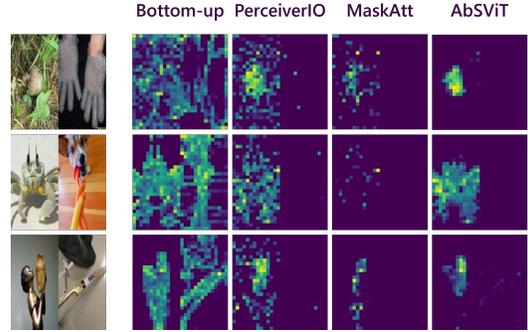


Figure 5. Comparison between different top-down attention algorithms. Prior corresponds to the left image. AbSViT has cleaner attention map than other baselines.

reweighting the output tokens to reweight the value tokens in intermediate self-attention modules, instead of adding the top-down tokens on them, (iii) **Feedback** directly feeds back the output tokens without reweighting. For V&L tasks, we use the METER [19] framework, which contains a vision backbone, a language backbone, and a multimodal fusion module. We use ViT [18] as the vision backbone and replace it with AbSViT or the baseline models. For image classification, we try the backbones of ViT, RVT [43], and FAN [77], which is state of the art on ImageNet and robustness benchmarks. The scaling factor  $\alpha$  is set as 1 during ImageNet pretraining and evaluation and set as 10 for fine-tuning on V&L tasks because we find AbSViT pretrained on supervised single-object classification only learns weak top-down attention in multi-object scenes (Sec. 7.1). See the Appendix for additional implementation details.

### 6.1. Controllable Top-Down Attention of AbSViT

To test the top-down attention in multi-object images, we take a AbSViT pretrained on ImageNet (Sec. 6.3) and create multi-object images by randomly sampling two images from ImageNet and concatenating them side by side. To control the top-down attention, we use the class prototype (from the last linear layer) of the two classes as  $\xi$ . Since in regular ViT, the class prototypes only align with the [cls] token but not with other output tokens, here we use a ViT with global average pooling. We set  $\alpha = 10$ .

To compare the bottom-up and top-down attention, we visualize the norm of output tokens from ViT and AbSViT for each class. As shown in Fig. 4, bottom-up attention highlights both objects while only the target object is selected by top-down attention. Consequently, the classification result, which has a tie between two classes when no prior is available, is biased towards the target class when we turn on the prior. This indicates AbSViT has the ability to control its attention on different objects given different priors. We also compare the top-down attention of AbSViT with several baselines (Fig. 5). We can see that the attention of

Model	VQAv2		Flickr-Zero-Shot		
	test-dev	test-std	IR@1	IR@5	IR@10
BEiT-B-16 [4]	68.45	-	32.24	-	-
CLIP-B-32 [52]	69.69	-	49.86	-	-
ViT-B	67.89	67.92	42.40	77.18	86.82
- PerceiverIO	67.87	67.93	42.52	76.92	86.73
- Feedback	67.99	68.13	42.04	77.38	86.90
- MaskAtt	67.53	67.51	41.89	76.53	86.78
- AbSViT	<b>68.72</b>	<b>68.78</b>	<b>45.28</b>	<b>77.98</b>	<b>87.52</b>

Table 1. Comparison of different top-down attention algorithms on VQA and zero-shot image retrieval. AbSViT achieves consistent improvements on both tasks.

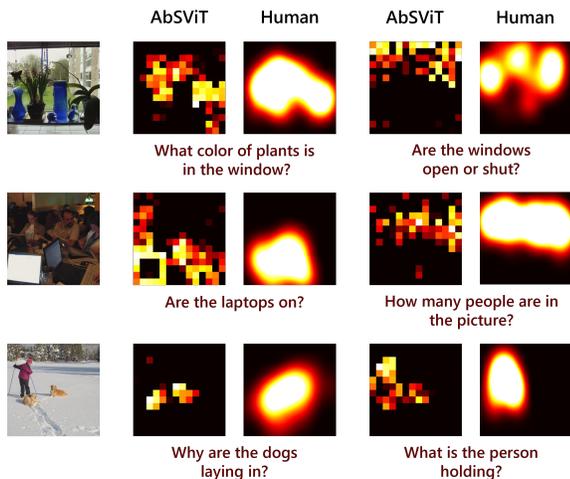


Figure 6. Comparison of attention map from AbSViT and human attention on VQA. AbSViT’s attention is adjustable to different questions and is consistent with human attention.

PerceiverIO focuses coarsely on the target object but is noisy, possibly because it lacks a feedback mechanism. MaskAtt, on the other hand, tends to miss parts of the object, implying that masking attention is less suitable for ViTs.

## 6.2. AbSViT for Vision-Language Tasks

We test AbSViT on two V&L tasks, VQA, and zero-shot image retrieval. We use the METER framework and replace the vision backbone with ViT-B, AbSViT-B, and other baselines. All the vision backbones are pretrained on ImageNet (Sec. 6.3). Results are shown in Tab. 1.

On VQAv2, AbSViT surpasses the baselines on both test splits and reaches the same performance as the unsupervised model (BEiT-B). At the same time, PerceiverIO has no improvement over ViT, probably because the multimodal fusion in METER can already perform token reweighting. The pure feedback network helps a little, mainly due to the feature refinement during the feedback loop. It is worth noticing that MaskAtt, a strategy frequently used in previous work, actually hurts performance when added to the vision

Model	P/F	Clean	IN-C (↓)	IN-A	IN-SK	IN-R
PiT-Ti [28]	5/0.7	72.9	69.1	6.2	34.6	21.6
ConViT-Ti [20]	6/1.4	73.3	68.4	8.9	35.2	22.4
PVT-Ti [64]	13/1.9	75.0	79.6	7.9	33.9	21.5
GFNet-Ti [54]	8/1.3	74.6	65.9	6.3	40.4	27.0
ViT-Ti [18]	6/1.3	72.5	71.1	7.5	33.0	20.1
- AbS	7/2.6	<b>74.1</b>	<b>66.7</b>	<b>10.1</b>	<b>34.9</b>	<b>22.6</b>
RVT-Ti [43]	9/1.3	78.1	58.8	13.9	42.5	29.1
- AbS	11/2.7	<b>78.6</b>	<b>55.9</b>	<b>17.3</b>	<b>43.2</b>	<b>29.9</b>
FAN-Ti [77]	7/1.3	77.5	59.8	13.1	42.6	29.9
- AbS	9/2.9	<b>78.3</b>	<b>57.4</b>	<b>16.5</b>	<b>42.8</b>	<b>31.2</b>
PiT-S [28]	24/2.9	80.9	52.5	21.7	43.6	30.8
PVT-S [64]	25/3.8	79.9	66.9	18.0	40.1	27.2
Swin-T [39]	28/4.5	81.2	62.0	21.6	41.3	29.1
ConvNext-T [40]	29/4.5	82.1	53.2	24.2	47.2	33.8
ViT-S [18]	22/4.2	80.1	54.6	19.2	41.9	28.9
- AbS	26/9.8	<b>80.7</b>	<b>51.6</b>	<b>24.3</b>	<b>43.1</b>	<b>30.2</b>
RVT-S [43]	22/4.3	81.9	50.5	26.0	47.0	34.5
- AbS	26/10.4	81.9	<b>48.7</b>	<b>31.1</b>	<b>48.5</b>	<b>35.6</b>
FAN-S [77]	28/5.3	82.8	49.1	29.3	47.4	35.6
- AbS	32/11.4	<b>83.0</b>	<b>47.4</b>	<b>34.0</b>	<b>48.3</b>	<b>36.4</b>
PiT-B [28]	74/12.5	82.4	48.2	33.9	43.7	32.3
PVT-L [64]	61/9.8	81.7	59.8	26.6	42.7	30.2
Swin-B [39]	88/15.4	83.4	54.4	35.8	46.6	32.4
ConvNext-B [40]	89/15.4	83.8	46.8	36.7	51.3	38.2
ViT-B [18]	87/17.2	80.8	49.3	25.2	<b>43.3</b>	31.6
- AbS	99/38.9	<b>81.0</b>	<b>48.3</b>	<b>28.2</b>	42.9	31.7
RVT-B [43]	86/17.7	80.9	52.1	26.6	<b>39.6</b>	26.1
- AbS	100/39.5	80.9	<b>51.7</b>	<b>28.5</b>	39.3	26.0
FAN-B [77]	54/10.4	83.5	45.0	33.2	51.4	39.3
- AbS	62/21.8	<b>83.7</b>	<b>44.1</b>	<b>38.4</b>	<b>52.0</b>	<b>39.8</b>

Table 2. Results on ImageNet classification and robustness benchmarks. AbSViT improves performance across different benchmarks and backbones. P/F: # of parameters and FLOPs. ↓: lower is better.

transformer. On zero-shot image retrieval, AbSViT also has higher performance than all other baselines. Especially, it has an improvement of  $\sim 3\%$  over bottom-up ViT on IR@1.

We also visualize the attention map of AbSViT on VQA and compare it to human attention. As shown in Fig. 6, AbSViT can adjust its attention to the objects related to the question. The attention map is also consistent with human attention. Nevertheless, the attention map of AbSViT is still not precise enough. For example, in the last example, when the question is “What is the person holding?”, the top-down attention highlights both the person and the dogs. Since the model is only pretrained on ImageNet, it may be further improved by CLIP [52] pretraining.

## 6.3. Image Classification and Robustness

We test AbSViT on ImageNet classification and robustness benchmarks (Tab. 2). We report mCE (lower the better) [25] for IN-C and accuracy for other datasets. On clean images, AbSViT consistently improves over baselines, with a similar number of parameters although higher FLOPs. The clean accuracy on FAN-B is improved to 83.7%, reaching the same level as ConvNext-B with fewer parameters. On

Model	Clean	IN-C (↓)	IN-A	IN-R	IN-SK
ViT-Ti	72.5	71.1	7.5	33.0	20.1
- PerceiverIO	72.8	70.4	8.0	32.8	20.5
- Feedback	73.4	67.8	9.7	34.6	22.4
- MaskAtt	72.5	70.6	8.3	33.4	20.5
- AbS	<b>74.1</b>	<b>66.7</b>	<b>10.1</b>	<b>34.9</b>	<b>22.6</b>
RVT-Ti	78.1	58.8	13.9	42.5	29.1
- PerceiverIO	78.3	57.8	13.7	42.8	29.8
- Feedback	79.1	55.7	18.2	44.1	31.3
- MaskAtt	77.9	59.0	13.5	43.0	29.7
- AbS	<b>79.5</b>	<b>54.8</b>	<b>18.7</b>	<b>44.5</b>	<b>32.5</b>

Table 3. Comparison of different top-down attention algorithms on ImageNet classification and robustness.

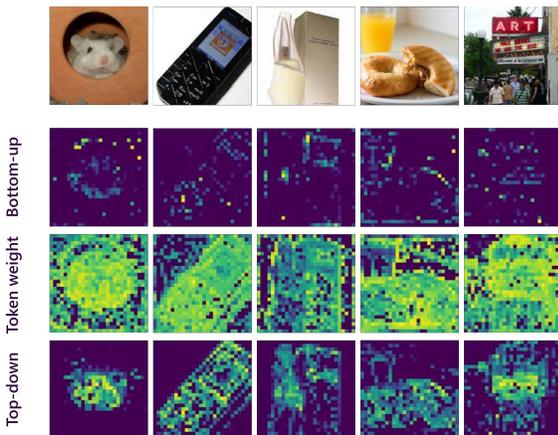


Figure 7. Visualization of the bottom-up attention, token weights, and the top-down attention in AbSViT. The bottom-up attention is noisy and fails to detect the complete foreground object. In AbSViT, the query mask can coarsely detect the foreground object and reweight tokens fed back to direct the top-down attention to better extract the foreground object.

corrupted (IN-C) and adversarial (IN-A) images, AbSViT boosts the performance by about 1-5% across all the scales. Especially, the performance on FAN-B is raised by 1% and 5% for IN-C and IN-A, reaching a new state-of-the-art result. On out-of-distribution images, AbSViT also improves by 3% on Tiny and Small models and 0.5% on FAN-B.

Fig. 7 visualizes the attention map of ViT and AbSViT, as well as token weights generated in  $e_{\xi}(\cdot)$ . The bottom-up attention in ViT is often noisy and only partly detects the foreground object. On the other hand, the query  $\xi$  in AbSViT learns to coarsely detect the foreground and reweight the feedforward output tokens, which are fed back and generate top-down attention that better detects the foreground object.

We compare AbSViT with several baseline algorithms for goal-directed attention in Tab. 3. One may see that a pure feedback model already improves the clean accuracy and robustness, and AbSViT further boosts the performance by better extracting the foreground object. Due to a similar

Model	PASCAL VOC	Cityscapes	ADE20K
ResNet-101 [12]	77.1	<b>78.7</b>	42.9
ViT-B	80.1	75.3	45.2
AbSViT-B	<b>81.3 (+1.2)</b>	76.8 (+1.5)	<b>47.2 (+2.0)</b>

Table 4. Semantic segmentation results on three datasets.

reason, PerceiverIO without feedback also slightly improves the performance. On the other hand, MaskAtt is sometimes harmful (on Clean, IN-C, and IN-A for RVT), implying that a mask attention design is unsuitable for vision transformers.

## 6.4. Semantic Segmentation

We evaluate the performance of AbSViT as a backbone for semantic segmentation on three datasets (PASCAL VOC, Cityscapes, and ADE20K). We compare with two baseline backbones, regular ViT and ResNet-101. We use UperNet [67] as the segmentation head for all the backbones. Results are shown in Tab. 4. We can see that when using AbSViT as the backbone, we can achieve 1.2-2.0% improvements over the ViT baseline with approximately the same number of parameters. This indicates that AbSViT can be used as a general backbone for different vision tasks.

## 6.5. Justification of Model Design

The design of AbSViT follows the principle of AbS. For example, AbSViT adds the top-down signal only to the value matrix considering the analogy between self-attention and sparse reconstruction (Sec. 5.1). At the same time, an arbitrary design may also add it to the query and key. We also optimize the variational loss to approximate AbS instead of just building a top-down model and training with the supervised loss. In this section, we show the advantage of these “destined” designs compared with an arbitrary design, which also justifies the proposed guiding principle of AbS.

We first try an arbitrary design of self-attention with top-down input by adding the top-down signal on the query, key, and value instead of only on the value. We name this design as AbSViT-QKV. We compare AbSViT and AbSViT-QKV on image classification and robustness (Tab. 5), and we can see that AbSViT is superior to AbSViT-QKV on every benchmark. This is consistent with our analysis in Sec. 4.2 that the sparse reconstruction AbS is optimizing has an additional top-down input (corresponding to V), while the dictionary (corresponding to Q and K), which contains templates for separate objects, is fixed.

We also test the effect of the variational loss  $\mathcal{L}_{var}$ , which ensures the model is approximating AbS. We compare AbSViT with its counterpart without  $\mathcal{L}_{var}$ , *i.e.*, a top-down model trained with only supervised loss. As shown in Tab. 6, adding  $\mathcal{L}_{var}$  largely improves the clean accuracy and robustness. Note that, as discussed in Sec. 5.1, we do not have a prior loss  $-\log p(\mathbf{z}_L)$  for image classification, which means

Model	Clean	IN-C ( $\downarrow$ )	IN-A	IN-R	IN-SK
AbSViT-QKV	73.3	68.0	9.4	33.8	21.2
AbSViT	<b>74.1</b>	<b>66.7</b>	<b>10.1</b>	<b>34.9</b>	<b>22.6</b>

Table 5. The predicted design of top-down self-attention (AbSViT) is better than an arbitrary design (AbSViT-QKV).

	$\mathcal{L}_{var}$	Clean	IN-C ( $\downarrow$ )	IN-A	IN-R	IN-SK
AbSViT	$\times$	73.1	69.0	9.5	33.5	20.8
AbSViT	$\checkmark$	<b>74.1</b>	<b>66.7</b>	<b>10.1</b>	<b>34.9</b>	<b>22.6</b>

Table 6. Ablation on the variational loss  $\mathcal{L}_{var}$ .

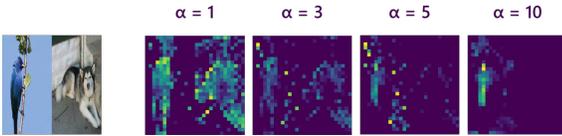


Figure 8. Visualization of top-down attention with different scaling factor  $\alpha$ . Prior corresponds to the bird. The top-down attention gets more and more biased on the bird when increasing  $\alpha$ .

the improvement completely comes from the reconstruction loss  $\frac{1}{2} \sum_{\ell=1}^L \|sg(\mathbf{z}_\ell) - g_\ell(sg(\mathbf{z}_{\ell+1}))\|_2^2$  which forces the decoder to reconstruct  $\mathbf{z}_\ell$  from  $\mathbf{z}_{\ell+1}$ . This implies that a generative model (“synthesis”) is important to high-quality top-down attention in visual recognition (“analysis”).

## 7. Limitations and Future Work

### 7.1. ImageNet Classification Is a Poor Teacher of Top-Down Attention

AbSViT is trained to focus on different objects given different priors in multi-object images. However, ImageNet classification targets single object classification without any prior, making it unsuitable for pretraining top-down attention. We find that the ImageNet-supervised AbSViT only learns weak top-down attention. A simple trick to augment the top-down attention for downstream tasks such as VQA is manually setting a larger scaling factor  $\alpha$  (e.g.,  $\alpha = 10$ ). In Fig. 8, we visualize the top-down attention with different  $\alpha$ . We can see that, with a prior corresponding to the bird, the attention under  $\alpha = 1$  still highlights both the bird and the dog but is more and more biased towards the bird as we increase  $\alpha$ . For future exploration, we may learn stronger top-down attention through object-level unsupervised learning [24, 68] or vision-language pretraining [47, 70].

### 7.2. How Many Syntheses Do We Need for Analysis?

In Sec. 5, we mention that enforcing strong generative capability on the features  $\mathbf{z}_\ell$  will downgrade the discriminative power regarding classification accuracy. There is a similar observation in recent self-supervised learning work [23], where reconstruction-based algorithms have worse linear



Figure 9. Examples of images decoded from the bottom-up, top-down, or the combination of bottom-up and top-down signals. The decoder can reconstruct the whole image from the bottom-up signal while failing to generate anything recognizable from the top-down signal alone. When decoding from the combination of bottom-up and top-down signals, only the foreground object is reconstructed.

probing performance [6]. However, the empirical results in Tab. 6 indicate that at least some degree of generative power is still helpful. This echoes the classical debate of how much generative capability (“synthesis”) we need for visual discrimination (“analysis”). As a starting point, we measure the generative power of the ImageNet-pretrained AbSViT (Fig. 9). Specifically, we train a linear decoder that projects the bottom-up input  $\mathbf{x}_0^{bu}$  of the first layer to the original image and then visualize the image decoded from the bottom-up signal  $\mathbf{x}_0^{bu}$ , the top-down signal  $\mathbf{x}_0^{td}$ , or their combination  $\mathbf{x}_0^{bu} + \mathbf{x}_0^{td}$ . We can see that the bottom-up signal contains full information about the original image and gives a perfect reconstruction. On the other hand, the top-down signal has lost most of the information, which is reasonable considering that  $\mathbf{x}_0^{td}$  itself is decoded from the last layer’s feature. Intriguingly, when we combine the bottom-up and the top-down signals, it can reconstruct only the foreground object, implying AbSViT can selectively preserve partial information in the image, and the selection process is adaptive to different priors. This leaves the question of whether a *selective* generation process is the best companion of the discriminative model and how to control the selective process under different priors adaptively.

## 8. Conclusion

We consider top-down attention by explaining from an Analysis-by-Synthesis (AbS) view of vision. Starting from previous work on the functional equivalence between visual attention and sparse reconstruction, we show that AbS optimizes a similar sparse reconstruction objective but modulates it with a goal-directed top-down modulation, thus simulating top-down attention. We propose AbSViT, a top-down modulated ViT model that variationally approximates AbS. We show that AbSViT achieves controllable top-down attention and improves over baselines on V&L tasks as well as image classification and robustness.

**Acknowledgement.** The authors would like to thank Tianyuan Zhang and Amir Bar for their valuable suggestions. Baifeng Shi and Trevor Darrell are supported by DARPA and/or the BAIR Commons program.

## References

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018. 1, 2
- [2] Alessandra Angelucci, Jonathan B Levitt, Emma JS Walton, Jean-Michel Hupe, Jean Bullier, and Jennifer S Lund. Circuits for local and global signal integration in primary visual cortex. *Journal of Neuroscience*, 22(19):8633–8646, 2002. 2
- [3] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433, 2015. 2
- [4] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 7
- [5] Ali Borji, Dicky Sihite, and Laurent Itti. An object-based bayesian framework for top-down visual attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pages 1529–1535, 2012. 1, 2
- [6] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 1, 9
- [7] Marisa Carrasco. Visual attention: The past 25 years. *Vision research*, 51(13):1484–1525, 2011. 1
- [8] James R Cavanaugh, Wyeth Bair, and J Anthony Movshon. Nature and interaction of signals from the receptive field center and surround in macaque v1 neurons. *Journal of neurophysiology*, 88(5):2530–2546, 2002. 2
- [9] Gang Chen. Where to look: A unified attention model for visual recognition with reinforcement learning. *arXiv preprint arXiv:2111.07169*, 2021. 1
- [10] Sharat Chikkerur, Thomas Serre, Cheston Tan, and Tomaso Poggio. What and where: A bayesian inference theory of attention. *Vision research*, 50(22):2233–2247, 2010. 1, 2
- [11] Krzysztof Choromanski, Valerii Likhoshesterov, Davidohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020. 3
- [12] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mms Segmentation>, 2020. 8, 15
- [13] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 6
- [14] Abhishek Das, Harsh Agrawal, Larry Zitnick, Devi Parikh, and Dhruv Batra. Human attention in visual question answering: Do humans and deep networks look at the same regions? *Computer Vision and Image Understanding*, 163: 90–100, 2017. 6
- [15] Peter Dayan, Geoffrey E Hinton, Radford M Neal, and Richard S Zemel. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995. 2
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2, 6
- [17] Robert Desimone, John Duncan, et al. Neural mechanisms of selective visual attention. *Annual review of neuroscience*, 18(1):193–222, 1995. 3
- [18] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 1, 2, 5, 6, 7
- [19] Zi-Yi Dou, Yichong Xu, Zhe Gan, Jianfeng Wang, Shuohang Wang, Lijuan Wang, Chenguang Zhu, Pengchuan Zhang, Lu Yuan, Nanyun Peng, et al. An empirical study of training end-to-end vision-and-language transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18166–18176, 2022. 2, 6, 15
- [20] Stéphane d’Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *International Conference on Machine Learning*, pages 2286–2296. PMLR, 2021. 7
- [21] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>. 6
- [22] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017. 6

- [23] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. [2](#), [9](#)
- [24] Olivier J Hénaff, Skanda Koppula, Evan Shelhamer, Daniel Zoran, Andrew Jaegle, Andrew Zisserman, João Carreira, and Relja Arandjelović. Object discovery and representation networks. *arXiv preprint arXiv:2203.08777*, 2022. [9](#)
- [25] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019. [6](#), [7](#)
- [26] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021. [6](#)
- [27] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15262–15271, 2021. [6](#)
- [28] Byeongho Heo, Sangdoon Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh. Rethinking spatial dimensions of vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11936–11945, 2021. [7](#)
- [29] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with neural networks. *science*, 313(5786):504–507, 2006. [2](#)
- [30] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006. [2](#)
- [31] Yujia Huang, James Gornet, Sihui Dai, Zhiding Yu, Tan Nguyen, Doris Tsao, and Anima Anandkumar. Neural networks with recurrent generative feedback. *Advances in Neural Information Processing Systems*, 33:535–545, 2020. [2](#)
- [32] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021. [6](#)
- [33] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. [2](#), [4](#)
- [34] David C Knill and Whitman Richards. *Perception as Bayesian inference*. Cambridge University Press, 1996. [1](#)
- [35] T Sing Lee. Analysis and synthesis of visual images in the brain: evidence for pattern theory. *IMA VOLUMES IN MATHEMATICS AND ITS APPLICATIONS*, 133:87–106, 2003. [1](#)
- [36] Tai Sing Lee. Top-down influence in early visual processing: a bayesian perspective. *Physiology & behavior*, 77(4-5):645–650, 2002. [1](#), [2](#)
- [37] Tai Sing Lee and David Mumford. Hierarchical bayesian inference in the visual cortex. *JOSA A*, 20(7):1434–1448, 2003. [1](#)
- [38] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. [2](#)
- [39] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. [7](#)
- [40] 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. [7](#)
- [41] Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Hierarchical question-image co-attention for visual question answering. *Advances in neural information processing systems*, 29, 2016. [2](#)
- [42] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 375–383, 2017. [2](#)
- [43] Xiaofeng Mao, Gege Qi, Yuefeng Chen, Xiaodan Li, Ranjie Duan, Shaokai Ye, Yuan He, and Hui Xue. Towards robust vision transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12042–12051, 2022. [6](#), [7](#), [15](#)
- [44] Julio C Martinez-Trujillo and Stefan Treue. Attentional modulation strength in cortical area mt depends on stimulus contrast. *Neuron*, 35(2):365–370, 2002. [2](#)
- [45] Carrie J McAdams and John HR Maunsell. Effects of attention on orientation-tuning functions of single neurons in macaque cortical area v4. *Journal of Neuroscience*, 19(1):431–441, 1999. [2](#)
- [46] M Berk Mirza, Rick A Adams, Karl Friston, and Thomas Parr. Introducing a bayesian model of selective attention based on active inference. *Scientific reports*, 9(1):1–22, 2019. [1](#), [2](#)
- [47] Jishnu Mukhoti, Tsung-Yu Lin, Omid Poursaeed, Rui Wang, Ashish Shah, Philip HS Torr, and Ser-Nam Lim. Open vocabulary semantic segmentation with patch aligned contrastive learning. *arXiv preprint arXiv:2212.04994*, 2022. [9](#)

- [48] Aude Oliva, Antonio Torralba, Monica S Castelhana, and John M Henderson. Top-down control of visual attention in object detection. In *Proceedings 2003 International Conference on Image Processing (Cat. No. 03CH37429)*, volume 1, pages I–253. IEEE, 2003. 2
- [49] Bo Pang, Yizhuo Li, Jiefeng Li, Muchen Li, Hanwen Cao, and Cewu Lu. Tdaf: Top-down attention framework for vision tasks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 2384–2392, 2021. 1
- [50] Bryan A Plummer, Liwei Wang, Chris M Cervantes, Juan C Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In *Proceedings of the IEEE international conference on computer vision*, pages 2641–2649, 2015. 6
- [51] Michael I Posner. Orienting of attention. *Quarterly journal of experimental psychology*, 32(1):3–25, 1980. 2
- [52] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021. 5, 7
- [53] Rajesh PN Rao. Bayesian inference and attentional modulation in the visual cortex. *Neuroreport*, 16(16):1843–1848, 2005. 1, 2
- [54] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. *Advances in Neural Information Processing Systems*, 34:980–993, 2021. 7
- [55] John H Reynolds and David J Heeger. The normalization model of attention. *Neuron*, 61(2):168–185, 2009. 2
- [56] Christopher J Rozell, Don H Johnson, Richard G Baraniuk, and Bruno A Olshausen. Sparse coding via thresholding and local competition in neural circuits. *Neural computation*, 20(10):2526–2563, 2008. 3
- [57] Michael P Sceniak, Michael J Hawken, and Robert Shapley. Visual spatial characterization of macaque v1 neurons. *Journal of neurophysiology*, 85(5):1873–1887, 2001. 2
- [58] Paul R Schrater and Rashmi Sundaeswara. Theory and dynamics of perceptual bistability. *Advances in neural information processing systems*, 19, 2006. 1
- [59] Baifeng Shi, Yale Song, Neel Joshi, Trevor Darrell, and Xin Wang. Visual attention emerges from recurrent sparse reconstruction. *arXiv preprint arXiv:2204.10962*, 2022. 2, 3, 4
- [60] Shengbang Tong, Xili Dai, Yubei Chen, Mingyang Li, Zengyi Li, Brent Yi, Yann LeCun, and Yi Ma. Unsupervised learning of structured representations via closed-loop transcription. *arXiv preprint arXiv:2210.16782*, 2022. 2
- [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [62] Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010. 2
- [63] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *Advances in Neural Information Processing Systems*, 32, 2019. 6
- [64] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 568–578, 2021. 7
- [65] Mike Wu, Kristy Choi, Noah Goodman, and Stefano Ermon. Meta-amortized variational inference and learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 6404–6412, 2020. 4
- [66] Yingda Xia, Yi Zhang, Fengze Liu, Wei Shen, and Alan L Yuille. Synthesize then compare: Detecting failures and anomalies for semantic segmentation. In *European Conference on Computer Vision*, pages 145–161. Springer, 2020. 2
- [67] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European conference on computer vision (ECCV)*, pages 418–434, 2018. 8, 15
- [68] Tete Xiao, Colorado J Reed, Xiaolong Wang, Kurt Keutzer, and Trevor Darrell. Region similarity representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10539–10548, 2021. 9
- [69] Huijuan Xu and Kate Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *European conference on computer vision*, pages 451–466. Springer, 2016. 1, 2
- [70] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18134–18144, 2022. 9
- [71] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057. PMLR, 2015. 1, 2

- [72] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. Stacked attention networks for image question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 21–29, 2016. [2](#)
- [73] Angela J Yu and Peter Dayan. Inference, attention, and decision in a bayesian neural architecture. *Advances in neural information processing systems*, 17, 2004. [1](#), [2](#)
- [74] Alan Yuille and Daniel Kersten. Vision as bayesian inference: analysis by synthesis? *Trends in cognitive sciences*, 10(7): 301–308, 2006. [1](#)
- [75] Li Zhaoping. *Understanding vision: theory, models, and data*. OUP Oxford, 2014. [1](#)
- [76] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. [6](#)
- [77] Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animesh Anandkumar, Jiashi Feng, and Jose M Alvarez. Understanding the robustness in vision transformers. In *International Conference on Machine Learning*, pages 27378–27394. PMLR, 2022. [6](#), [7](#), [15](#)

## A. Derivation of Eq. (10)

From Eq. (6-7) we have

$$p(\tilde{\mathbf{u}}_\ell | \tilde{\mathbf{u}}_{\ell+1}) \propto \exp\left\{-\frac{1}{2}\|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{P}_{\ell+1} \tilde{\mathbf{u}}_{\ell+1})\|_2^2 - \lambda \|\tilde{\mathbf{u}}_\ell\|_1\right\}. \quad (16)$$

Then Eq. (10) is derived by

$$\begin{aligned} \frac{d\tilde{\mathbf{u}}_\ell}{dt} &\propto \nabla_{\tilde{\mathbf{u}}_\ell} \log p(\tilde{\mathbf{u}}_{\ell-1} | \tilde{\mathbf{u}}_\ell) + \nabla_{\tilde{\mathbf{u}}_\ell} \log p(\tilde{\mathbf{u}}_\ell | \tilde{\mathbf{u}}_{\ell-1}) \\ &= -\nabla_{\tilde{\mathbf{u}}_\ell} \frac{1}{2} \|\mathbf{P}_{\ell-1} \tilde{\mathbf{u}}_{\ell-1} - g_{\ell-1}(\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell)\|_2^2 - \nabla_{\tilde{\mathbf{u}}_\ell} \frac{1}{2} \|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{P}_{\ell+1} \tilde{\mathbf{u}}_{\ell+1})\|_2^2 - \nabla_{\tilde{\mathbf{u}}_\ell} \lambda \|\tilde{\mathbf{u}}_\ell\|_1 \\ &= \mathbf{P}_\ell^T \mathbf{J}_{\ell-1}^T (\mathbf{P}_{\ell-1} \tilde{\mathbf{u}}_{\ell-1} - g_{\ell-1}(\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell)) - \mathbf{P}_\ell^T (\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{P}_{\ell+1} \tilde{\mathbf{u}}_{\ell+1})) - \nabla_{\tilde{\mathbf{u}}_\ell} \lambda \|\tilde{\mathbf{u}}_\ell\|_1 \\ &= -\mathbf{P}_\ell^T \left( \mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{P}_{\ell+1} \tilde{\mathbf{u}}_{\ell+1}) - \mathbf{J}_{\ell-1}^T \mathbf{P}_{\ell-1} \tilde{\mathbf{u}}_{\ell-1} \right) - \nabla_{\tilde{\mathbf{u}}_\ell} \lambda \|\tilde{\mathbf{u}}_\ell\|_1 - \mathbf{P}_\ell^T \mathbf{J}_{\ell-1}^T g_{\ell-1}(\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell) \\ &= -\nabla_{\tilde{\mathbf{u}}_\ell} \frac{1}{2} \|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - g_\ell(\mathbf{z}_{\ell+1}) - \mathbf{J}_{\ell-1}^T \mathbf{z}_{\ell-1}\|_2^2 - \nabla_{\tilde{\mathbf{u}}_\ell} \lambda \|\tilde{\mathbf{u}}_\ell\|_1 - \nabla_{\tilde{\mathbf{u}}_\ell} \frac{1}{2} \|g_{\ell-1}(\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell)\|_2^2 \\ &= -\nabla_{\tilde{\mathbf{u}}_\ell} \frac{1}{2} \|\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell - (\mathbf{x}_\ell^{td} + \mathbf{x}_\ell^{bu})\|_2^2 - \nabla_{\tilde{\mathbf{u}}_\ell} \lambda \|\tilde{\mathbf{u}}_\ell\|_1 - \nabla_{\tilde{\mathbf{u}}_\ell} \frac{1}{2} \|g_{\ell-1}(\mathbf{P}_\ell \tilde{\mathbf{u}}_\ell)\|_2^2. \end{aligned} \quad (17)$$

We informally use  $\nabla$  for subgradients as well.

## B. Additional Results on Natural Images

In Fig.4-5, we show examples of top-down attention on artificial images. Here we show more results on natural images containing multiple objects. We borrow the LVIS dataset and collect images that contain object categories that also appear in ImageNet. We demonstrate that given different prior, AbSViT is able to focus on different objects in the same image (Fig. 10). We also compare AbSViT’s top-down attention with several baseline methods (Fig. 11) and observe that AbSViT has cleaner attention maps than other methods.

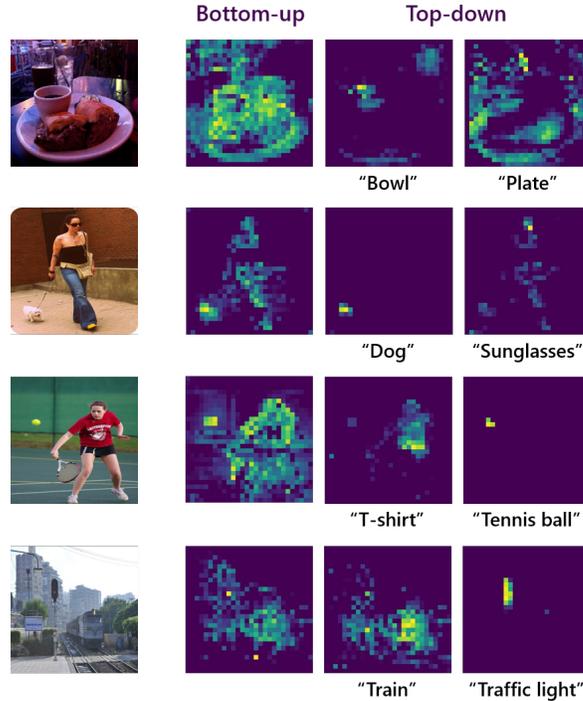


Figure 10. Visualization of top-down attention on natural images. From left to right, we show the original images, the bottom-up attention, as well as the top-down attention regarding to different objects in each image.

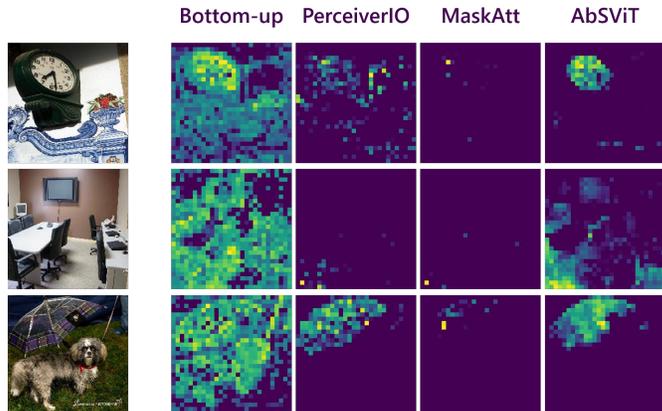


Figure 11. Comparison of top-down attention map between AbSViT and different baselines.

### C. Additional Implementation Details

**ImageNet Pretraining.** The ViT and RVT baselines as well as our AbSViT model are trained using the recipe in [43], and FAN is trained using the recipe in its original paper [77]. Specifically, we use AdamW optimizer to train AbSViT for 300 epochs, with a batch size of 512, a base learning rate of  $5e-4$ , and 5 warm-up epochs. One may use different batch-size and adjust the learning rate by the linear scaling rule. We use a cosine learning rate scheduling and weight decay of 0.05. We use the default setting of data augmentation, which includes Mixup, Cutmix, ColorJittering, AutoAugmentation, and Random Erasing. For AbSViT, the weights of supervised loss and variational loss are set as 1 and 0.1.

**Robustness against Image Corruptions.** We evaluate model robustness against image corruption on ImageNet-C, which contains a total of 19 corruption types. We follow [43] and evaluate 15 types of corruption including Brightness, Contrast, Defocus Blur, Elastic Transform, Fog, Frost, Gaussian Noise, Glass Blur, Impulse Noise, JPEG Compression, Motion Blur, Pixelate, Shot Noise, Snow, and Zoom Blur. Note that other work (e.g. [77]) tests on a different subset of corruption types. To make a fair comparison, all the models are tested under the aforementioned 15 corruption types.

**Semantic Segmentation.** We use MMSegmentation [12] as our test bed. We take the ImageNet pretrained ViT-B and AbSViT-B and finetune them on semantic segmentation on PASCAL VOC, Cityscapes, and ADE20K. For all the experiments, we use UperNet [67] as the decoder head and FCNHead as the auxiliary head. We train on 2 GPUs with a total batch size of 16, using AdamW optimizer, a learning rate of 0.00006, and weight decay of 0.01. We train for 20k, 40k, and 160k iterations for three datasets, respectively. We use image resolution of  $512 \times 512$  for PASCAL VOC and ADE20K, and  $512 \times 1024$  for Cityscapes.

**V&L Finetuning.** Following [19], the whole model contains a pretrained visual encoder, a pretrained text encoder, and a multimodal encoder to merge vision and language. We use the ImageNet pretrained ViT or AbSViT for the visual encoder, a pretrained RoBERTa for the text encoder, and the multimodal encoder is trained from scratch. We use a learning rate of  $1e-5$  for visual and text encoders and  $5e-5$  for the multimodal encoder. For top-down attention, we use the `[cls]` token as the prior  $\xi$ . Since the text and visual tokens are not aligned initially, we train a linear transform to project the text tokens into the same space as the visual tokens. This is trained by the prior loss, which is set as a CLIP-style loss (Eq. (15)) to align the text and visual tokens.