

# SELF-SUPERVISED SPARSE VISION CONCEPTS FOR IMAGE UNDERSTANDING AND RECONSTRUCTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Self-supervised vision encoders have become critical components of modern machine learning systems. Despite remarkable advances in image understanding, generation, and multimodal alignment, the underlying representation of visual features has remained largely unchanged, constrained by historical architectures and benchmarks. This reliance on dense feature grids introduces redundancy and limits the integration of understanding and generation. We propose a novel framework that represents images with a small number of sparse tokens **in the form of** low-rank matrix factorization. While mathematically simple, this formulation effectively disentangles semantic and spatial information. We demonstrate that vision-only self-supervised learning under this framework yields sparse token representations that simultaneously support high-quality image understanding, detailed pixel-level reconstruction, and fine-grained semantic understanding. Together, these results highlight sparse tokens as a promising alternative to dense grids for efficient and versatile visual representation learning.

## 1 INTRODUCTION

Learning visual representations has been a central pursuit in computer vision since the advent of deep learning (Bengio et al., 2013). Modern vision models transform raw pixels into latent features that power nearly all downstream applications. Architectures have evolved dramatically—from early convolutional networks (Lecun et al., 1998) to ResNets (He et al., 2016) and, more recently, vision transformers (ViTs) (Dosovitskiy et al., 2020). Despite these advances, the geometric format of visual representation has remained largely unchanged: a dense 2D grid of high-dimensional features, each tied to a local patch of the image. This design is natural, since pixels are arranged on a grid.

However, dense representations are highly redundant, as the number of meaningful objects or regions in an image is far smaller than the number of pixels or patches. Prior work has shown that images can be reconstructed from only a small subset of patches (He et al., 2022). Motivated by this redundancy, several methods have already adopted sparse representations in downstream tasks. For instance, Sparse R-CNN (Sun et al., 2021), DETR (Carion et al., 2020), and MaskFormer (Cheng et al., 2021) learn a set of queries from detection or segmentation labels. BLIP-2 (Li et al., 2023) introduced a Q-Former to extract sparse tokens under paired image–text supervision. TiTok (Yu et al., 2024) represents images as sparse 1D tokens for efficient reconstruction and generation, but the semantic quality of these features still lags behind state-of-the-art self-supervised methods.

We take a different path: learning sparse representations directly from images in a fully self-supervised manner (Caron et al., 2021; Oquab et al., 2023; Siméoni et al., 2025; Chen et al., 2020a; Chen & He, 2021; Caron et al., 2020; Grill et al., 2020). Our goal is to obtain a compact set of tokens that support high-quality image understanding and reconstruction at the same time, without human labels or paired data.

To this end, we revisit the fundamentals of visual representation. At its core, perception requires two complementary pieces of information:

1. What objects or visual concepts are present.
2. Where they are located.

The two are unified in forming a holistic representation, but separable in how they behave: the “what” should remain invariant across views, while the “where” changes with viewpoint. Motivated by this principle, we propose to represent an image in the form of low-rank matrix factorization that disentangles these factors. This formulation enables efficient reconstruction from as few as 8 tokens and allows learning the sparse tokens to encode fine-grained semantics in a self-supervised way.

Our contributions can be summarized as follows:

- **Sparse vision representation framework:** We propose STELLAR, an efficient form of latent vision representation modeling an image with only a handful of sparse tokens, by disentangling *what* concepts are present and *where* are they located. The latent representation with as few as 8 tokens can achieve both detailed pixel reconstruction and high-level semantic understanding at the same time.
- **Self-supervised learning method:** We introduce a self-supervised training scheme that learns the sparse latent vision representation without annotation. By discovering and aligning multiple visual concepts across views using optimal transport, we enforces invariance in the “what” while adapting the “where,” inducing rich semantic representation.
- **Observations:** (i) STELLAR surpasses prior sparse representation approaches by jointly achieving strong semantic understanding (IN-1k lin. acc. 79.10%) and high-quality reconstruction (FID 2.60). (ii) The sparse image modeling framework induces fine-grained, region-aware semantics even without explicit supervision on the dense feature map, which transfers effectively to downstream tasks with simple linear probing.

## 2 PRELIMINARIES

Representation learning involves encoding an image  $X$  with a neural network  $\mathcal{E}$  to extract latent features  $\mathbf{Z} = \mathcal{E}(X)$  for downstream tasks. Traditionally, vision representation takes the dense form

$$\mathbf{Z} \in \mathbb{R}^{(h \cdot w) \times d},$$

where  $h$  and  $w$  denote the height and width of the 2D grid partitioning the image. Each grid location is represented by a feature vector  $\mathbf{z}_i := \mathbf{Z}_{i,:} \in \mathbb{R}^d$  for  $1 \leq i \leq h \cdot w$ . Many vision models also include a global representation  $\mathbf{z}_0 \in \mathbb{R}^d$ , obtained either by pooling over the feature map or by introducing a [CLS] token in transformers. Even in variational models such as VAEs or VQ-VAEs, where the latent variables  $\mathbf{z}_i$  are modeled as probability distributions rather than deterministic embeddings, the underlying 2D grid structure remains unchanged.

In contrast, *sparse visual representation* aims to represent the image with

$$\mathbf{Z} \in \mathbb{R}^{r \times d}, \quad r \ll h \cdot w.$$

Ideally, the number of sparse tokens  $r$  should be less than an order of magnitude than the total number of dense tokens  $n = h \cdot w$ . In addition, we want the sparse tokens  $\mathbf{Z}$  to serve as a holistic representation of the image  $X$ , i.e.  $\mathbf{Z}$  contains sufficient information to reconstruction the original image, while at the same time possessing rich semantics for downstream tasks. Mathematically, we define such holistic representation as follows:

- **Reconstruction:** There exists a decoder  $\mathcal{D}$  such that the sparse features  $\mathbf{Z} = \mathcal{E}(X)$  can faithfully reconstruct the original image:  $\mathcal{D}(\mathbf{Z}) \approx X$ .
- **Understanding.** For a downstream task with joint distribution  $(X, Y) \sim \mathcal{X} \times \mathcal{Y}$  and loss function  $\mathcal{L}$ , there exists a simple predictor  $f \in \mathcal{F}$  such that, using frozen sparse features  $\mathbf{Z} = \mathcal{E}(X)$ , the expected task loss  $\mathbb{E}_{(X, Y)} [\mathcal{L}(f(\mathbf{Z}), Y)]$  is low.

While  $Y$  can in principle be arbitrary, in downstream tasks it typically reflects human interpretations of the image, such as classification labels, segmentation masks, or textual descriptions. The predictor  $f$  is usually drawn from a simple function class  $\mathcal{F}$ , e.g., a linear layer.

Prior works tend to emphasize only one aspect of “representation.” For example, TiTok (Yu et al., 2024) uses 32 tokens to reconstruct an image with 256 patches at high fidelity, but its sparse features

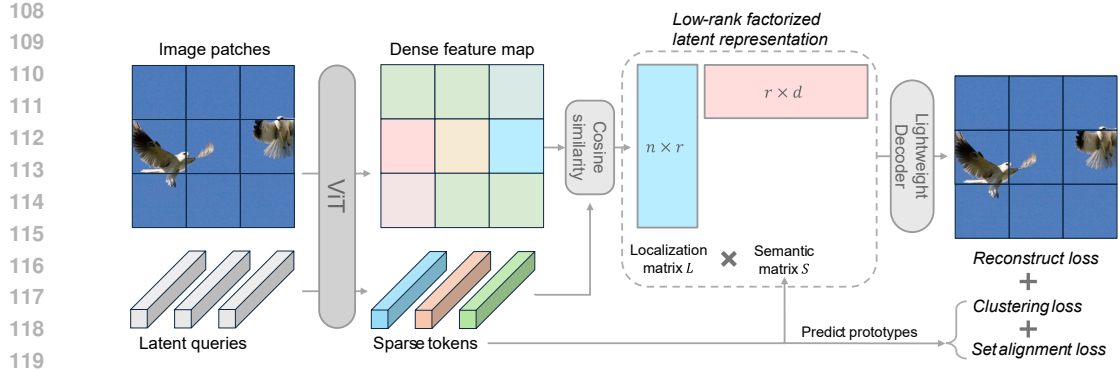


Figure 1: The STELLAR framework. We use a vanilla ViT to extract sparse tokens from an image, and model the latent representation as a low-rank matrix factorization, ensuring reconstruction of the original image. Clustering loss and set alignment loss are applied on the disentangled sparse tokens.

lag far behind other self-supervised vision models in semantic understanding. Conversely, MAE (He et al., 2022) is designed to capture semantics by reconstructing randomly masked patches, yet the resulting reconstructions are often blurry.

This reveals an empirical dilemma in current vision frameworks: models that excel at pixel-level reconstruction often produce weaker semantic representations (Zhang et al., 2022; Chen et al., 2024). Conversely, state-of-the-art SSL methods that achieve strong semantics typically abandon pixel reconstruction to avoid low-level shortcuts (Zhou et al., 2021; Assran et al., 2023; Darcet et al., 2025). In contrast, we demonstrate that by disentangling semantic and localization information, it is possible to learn sparse representations that simultaneously achieve strong image understanding and high-quality reconstruction.

### 3 THE STELLAR FRAMEWORK

#### 3.1 SPARSE IMAGE MODELING

Images depict the physical world, which can be understood as a collection of objects located in space. From this perspective, visual information naturally decomposes into two complementary components: (1) *what* objects or concepts are present, and (2) *where* they are located. Unlike dense grid-based representations, which describe what appears at each individual location, we model an image using a compact set of semantic concepts together with their spatial distributions. Concretely, we represent an image with  $r$  concept embeddings

$$\mathbf{s}_1, \dots, \mathbf{s}_r \in \mathbb{R}^d,$$

where each  $\mathbf{s}_j$  captures a distinct semantic concept. The spatial distribution of these concepts is expressed through weights  $\mathbf{l}_1, \dots, \mathbf{l}_n \in \mathbb{R}^r$ , where  $n$  is the total number of patches. By constraining  $0 \leq l_i \leq 1$  and  $\mathbf{1}^\top \mathbf{l}_i = 1$ , each patch is represented as a convex combination of the concept embeddings:  $\mathbf{v}_i = \sum_{j=1}^r l_{i,j} \mathbf{s}_j$ . Thus, the set  $\mathbf{s}_{j_{j=1}^r}$  acts as a basis for constructing patch-level features. In matrix form, the latent representation of image  $X$  is encoded as:

$$\mathbf{S}, \mathbf{L} = \mathcal{E}(X), \quad \mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_r]^\top \in \mathbb{R}^{r \times d}, \quad \mathbf{L} = [\mathbf{l}_1, \dots, \mathbf{l}_n]^\top \in \mathbb{R}^{n \times r}, \quad (1)$$

which can be combined to form a patch-wise dense representation:

$$\mathbf{V} = \mathbf{L}\mathbf{S}, \quad 0 \leq \mathbf{L} \leq 1, \mathbf{L}\mathbf{1}_r = \mathbf{1}_n. \quad (2)$$

Compared to a canonical dense representation of shape  $n \times d$ ,  $\mathbf{V} = \mathbf{L}\mathbf{S}$  can be considered as a form of low-rank matrix approximation from the sparse representation. Critically, we enforce that this low-rank approximation can reconstruct the original image through some decoder  $\mathcal{D}$ :

$$\mathcal{D}(\mathbf{L}\mathbf{S}) \approx X. \quad (3)$$

This *low-rank approximated reconstruction* ensures that the sparse tokens  $\{s_j\}_{j=1}^r$  capture sufficient information to recover the image when combined with their spatial distributions. While the form in equation 2 resembles the low-rank structure used in convex semi-nonnegative matrix factorization (Ding et al., 2008), we do not perform NMF or any matrix factorization algorithm on the feature map (i.e.  $LS \approx \mathcal{E}(X)$ ). Instead, both  $S$  and  $L$  are learnable latent variables output directly from the forward pass of the encoder, and their product is decoded back to the original image ( $\mathcal{D}(LS) \approx X$ ), allowing an autoencoder-style training. Finally, a compact sparse representation is then obtained by concatenating the concept embeddings with the transposed localization matrix:

$$Z = [S, L^T] \in \mathbb{R}^{r \times d^*}, \quad d^* = d + n. \quad (4)$$

We refer to our framework defined by the form of  $S$ ,  $L = \mathcal{E}(X)$  and  $\mathcal{D}(LS) \approx X$  as **S**parse **T**oken **E**xtraction and **L**ocalization with **L**ow-rank **A**pproximated **R**econstruction (**STELLAR**). We note that the framework is flexible and does not prescribe any specific encoder or decoder architecture. In this work, we adopt a simple design with common modules to obtain  $S$  and  $L$  as described below.

As illustrated in Fig. 1, the encoder includes a ViT and  $r$  learnable latent query vectors, which are passed to the ViT alongside patchified image tokens. Processed by the ViT jointly, the latent queries produce sparse tokens  $S \in \mathbb{R}^{r \times d}$ , and the image patches output a dense feature map  $U \in \mathbb{R}^{n \times d}$ .

To obtain the localization matrix  $L \in \mathbb{R}^{n \times r}$  associated with the sparse tokens, we project both  $S$  and  $U$  into a shared embedding space and compute their pairwise cosine similarities, followed by a softmax normalization with temperature  $t$ :

$$L = \text{softmax} \left( \frac{\text{cosim}(UW_1, SW_2)}{t} \right), \quad (5)$$

where  $W_1$  and  $W_2$  are learnable linear projections, and  $t$  controls the sharpness of the spatial distribution. This mapping is structurally similar to the attention weights obtained in a single-head cross-attention layer, up to the use of L2 normalization and an explicit temperature parameter. We adopt this simple formulation to compute the  $L$  matrix, and found it to be stable and effective for learning sparse concept localization across all experiments.

All together, the encoder  $\mathcal{E}$  includes a ViT,  $r$  learnable latent query vectors, and projection layers  $W_1, W_2$ . The decoder  $\mathcal{D}$  is a lightweight ViT reconstructing the image patches from  $LS$ .

Although the low-rank form  $LS$  is simple, it effectively disentangles high-level semantic concepts from low-level spatial localization. This yields two key benefits for both reconstruction and representation learning:

- The concept matrix  $S$  no longer needs to encode spatial information, and can instead focus purely on learning what objects or visual concepts are present. Through the linear combination  $LS$ , these concepts can be flexibly allocated across spatial locations to form a dense semantic map, enabling efficient reconstruction.
- Because the semantic embeddings in  $S$  are independent of location, we can freely apply image transformations while enforcing consistency in the learned concepts. This invariance induces robust high-level semantic features that transfer well to image understanding tasks.

Next we introduce how to learn semantic-rich latent representation as shown in Fig. 2.

### 3.2 LEARNING VISION CONCEPT VOCABULARY

To encourage sparse tokens to represent transferable vision concepts, we structure them into  $K$  learnable prototypes  $c_1, \dots, c_K \in \mathbb{R}^p$ . A backbone encoder  $\mathcal{E}$  maps a mini-batch of  $m$  images into sparse features  $S^1, \dots, S^m$ . Each token is projected onto the unit sphere  $\mathbb{S}^{p-1}$  via a normalized projector  $h: \mathbb{R}^d \rightarrow \mathbb{S}^{p-1}$ , and its similarity to prototypes  $C = [c_1, \dots, c_K]$  gives logits

$$\lambda_j^i = [c_1 \cdot h(s_j^i), \dots, c_K \cdot h(s_j^i)], \quad i = 1, \dots, m, \quad j = 1, \dots, r. \quad (6)$$

Soft assignments follow as

$$q_{j,k}^i = \frac{\exp(\lambda_{j,k}^i/\tau)}{\sum_{k'=1}^K \exp(\lambda_{j,k'}^i/\tau)}, \quad (7)$$

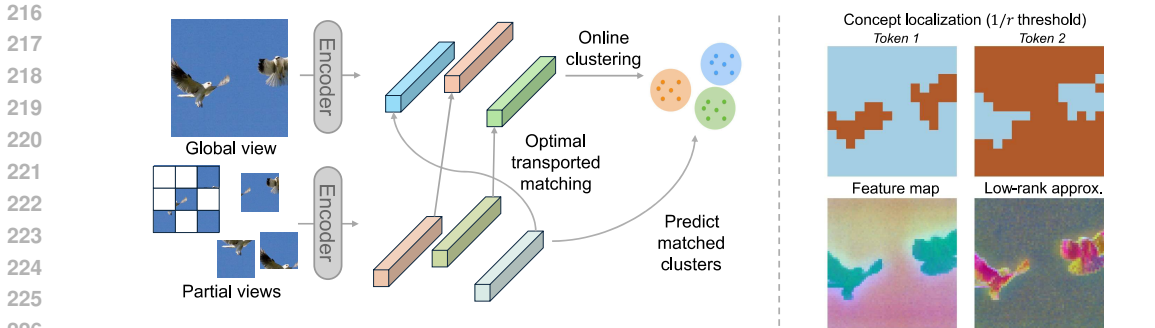


Figure 2: Left: Concept clustering and alignment workflow. Right: visualization of learned representation.

where  $\tau$  controls sharpness. Direct entropy minimization of  $q_j^i$  is unstable due to non-convexity and empty clusters. Following SwAV (Caron et al., 2020) and CAPI Darcet et al. (2025), we compute balanced assignments  $\tilde{q}_j^i$  with Sinkhorn-Knopp (temperature  $\tilde{\tau} > \tau$ ), and optimize

$$\mathcal{L}_{\text{cluster}} = \frac{1}{mr} \sum_{i=1}^m \sum_{j=1}^r \sum_{k=1}^K \tilde{q}_{j,k}^i \log q_{j,k}^i. \tag{8}$$

Unlike DINOv2 and SwAV, which use Sinkhorn only for balancing teacher targets, or CAPI, which optimizes prototypes separately, we minimize  $\mathcal{L}_{\text{cluster}}$  end-to-end along with other objectives.

### 3.3 SET CONCEPTS ALIGNMENT

STELLAR produces a fixed-size set of sparse tokens invariant to cropping, masking, or resolution. To align features across augmented views without inherent ordering, we use optimal transport as in Fig. 2. Given global view features  $s_1, \dots, s_r$  and partial-view features  $s'_1, \dots, s'_r$ , the cost matrix

$$\Theta_{j'j} = \|s'_{j'} - s_j\|_2. \tag{9}$$

We solve for an assignment matrix  $P$  via entropy-regularized optimal transport:

$$\min_{P \geq 0} \sum_{j',j} P_{j'j} \Theta_{j'j} - \epsilon H(P), \quad P \mathbf{1}_r = P^T \mathbf{1}_r = \frac{1}{r} \mathbf{1}_r, \tag{10}$$

with  $H(P) = -\sum_{j',j} P_{j'j} \log P_{j'j}$ . We solve for  $P$  by Sinkhorn algorithm, and define the matching  $\sigma(j') := \operatorname{argmax}_j P_{j'j}$ . We then compute prototype assignments for the partial-view tokens  $q'_{j'} = \operatorname{softmax}(C^T h(s'_{j'}) / \tau)$ , and minimize the set concept alignment loss

$$\mathcal{L}_{\text{align}} = \frac{1}{r} \sum_{j'=1}^r \sum_{k=1}^K \tilde{q}_{\sigma(j'),k} \log q'_{j',k}. \tag{11}$$

Optionally, the CLS token is treated as another sparse feature with its own projector and prototypes, but not used for reconstruction. In addition, we optionally use an exponential moving average (EMA) updated momentum encoder to encode the target assignments  $\tilde{q}$  in equation 8 and equation 11. We observed that using a momentum encoder is essential in the warm-up stage when training from scratch, but suboptimal in subsequent training. We provide detailed results in ablation study. All together, we jointly optimize the following to learn latent representation  $S$  and  $L$ :

- Reconstruction:  $\mathcal{L}_{\text{recon}} = \ell(\mathcal{D}(\mathbf{L}\mathbf{S}), X)$  via a lightweight decoder  $\mathcal{D}$ .
- Sparse concept clustering:  $\mathcal{L}_{\text{cluster}}$  on prototype assignments.
- Set concept alignment:  $\mathcal{L}_{\text{align}}$  between global and partial views.
- KoLeo regularization (Sablayrolles et al., 2018) on the sparse tokens from the same image:

$$\mathcal{L}_{\text{KoLeo}} = -\frac{1}{r} \sum_{j=1}^r \log \left( \frac{1}{2} \min_{j' \neq j} \|\bar{s}_j - \bar{s}_{j'}\|_2 \right), \quad \bar{s}_j := s_j / \|s_j\|.$$



proaches include PCA, low-rank matrix recovery (Candes & Plan, 2009), and dictionary learning or sparse coding (Olshausen & Field, 1997; Mairal et al., 2008; Tošić & Frossard, 2011), where signals are expressed through a small set of basis vectors. In deep learning, low-rank constraints have been widely applied for efficiency: for example, low-rank factorization of neural network weights (Sainath et al., 2013), low-rank approximations of attention maps (Katharopoulos et al., 2020), and parameter-efficient fine-tuning methods such as LoRA (Hu et al., 2022). In contrast to these works, STELLAR applies low-rank factorization directly to the feature map from a single image, disentangling spatial and semantic information.

## 5 EXPERIMENTS

We train STELLAR on ImageNet-1K (Deng et al., 2009) in a self-supervised setting. The encoder is a vanilla ViT (Dosovitskiy et al., 2020) augmented with 8–24 learnable latent queries that produce sparse tokens (without positional encoding). The [CLS] token is retained but not used for reconstruction. A lightweight 6-layer ViT serves as the decoder, predicting either MaskGIT-VQGAN tokens (Esser et al., 2021; Chang et al., 2022) or raw pixels (ablation). We initialize the ViT backbone with MAE pre-trained weights to accelerate training, enabling the model to focus on the sparse latent queries, projection layers, and decoder. Alternatively, we experimented with a momentum encoder warm-up. MAE initialization is used by default, with other methods analyzed in our ablations.

### 5.1 SPARSE TOKENS FOR UNDERSTANDING AND RECONSTRUCTION

We evaluate reconstruction with FID (Heusel et al., 2017) and representation quality with linear probing on mean-pooled sparse features. Table 1 and Figure 3 shows results across token counts (8, 16, 24), compared with TiTok (Yu et al., 2024) and MAE (He et al., 2022). STELLAR achieves strong reconstruction even with few tokens (rFID = 3.68 with 8 tokens; 2.60 with ViT-H, 16 tokens), approaching MaskGIT-VQGAN (2.28) without decoder finetuning. For linear probing, STELLAR maintains robust accuracy and does not drop with more tokens as in TiTok (Yu et al., 2024). Reconstruction improves with the number of token, but plateaus after 16. We used 16 tokens in all other experiments unless specified. We also see in Figure 3 that STELLAR preserves the location of the objects. Interestingly, it also automatically removes the dark edge in the bottom example, indicating it is reconstructing from high-level semantics rather than memorizing low-level details. Overall, STELLAR balances efficient reconstruction and discriminative understanding at high quality.

Table 1: Reconstruction FID and linear probing accuracy (%) of sparse tokens on IN1K. Model sizes are ViT-B by default, with larger sizes indicated in parentheses.

	VQGAN	TiTok		MAE		STELLAR (ours)			
# tokens	256	32 (L)	64	16	32	8	16	24	16 (H)
rFID ↓	2.28	2.75	1.99	150.73	131.01	3.68	3.14	3.19	2.60
lin. acc.	-	33.42	32.87	44.43	56.52	72.97	73.26	72.17	79.10

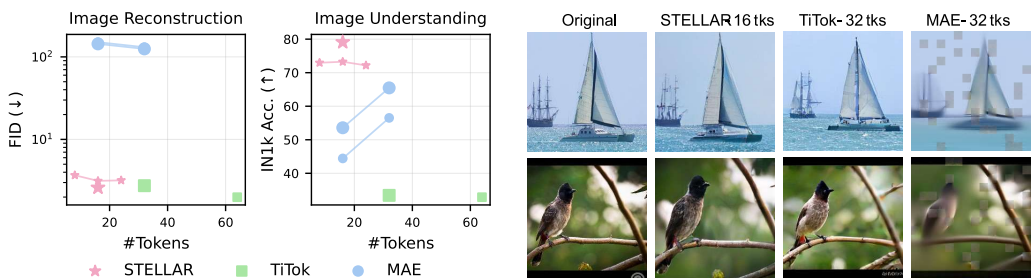


Figure 3: Sparse representation for image understanding and reconstruction. Left: reconstruction quality (FID) and semantic quality (lin. acc.) v.s. number of tokens. Right: reconstruction examples.







## REFERENCES

- 540  
541  
542 Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. Structured pruning of deep convolutional neural  
543 networks. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 13(3):1–18,  
544 2017.
- 545 Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent,  
546 Armand Joulin, Mike Rabbat, and Nicolas Ballas. Masked siamese networks for label-efficient  
547 learning. In *European conference on computer vision*, pp. 456–473. Springer, 2022.
- 548  
549 Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat,  
550 Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding  
551 predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and  
552 Pattern Recognition*, pp. 15619–15629, 2023.
- 553 Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. Data2vec:  
554 A general framework for self-supervised learning in speech, vision and language. In *International  
555 conference on machine learning*, pp. 1298–1312. PMLR, 2022.
- 556  
557 Hangbo Bao, Li Dong, and Furu Wei. Beit: BERT pre-training of image transformers. *arXiv preprint  
558 arXiv:2106.08254*, 2021.
- 559 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new  
560 perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828,  
561 2013.
- 562  
563 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative com-  
564 ponents with random forests. In *European Conference on Computer Vision*, 2014.
- 565  
566 Emmanuel J Candes and Yaniv Plan. Accurate low-rank matrix recovery from a small number of  
567 linear measurements. In *2009 47th Annual Allerton Conference on Communication, Control, and  
568 Computing (Allerton)*, pp. 1223–1230. IEEE, 2009.
- 569 Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and  
570 Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on  
571 computer vision*, pp. 213–229. Springer, 2020.
- 572  
573 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.  
574 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural  
575 information processing systems*, 33:9912–9924, 2020.
- 576  
577 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
578 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of  
579 the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- 580  
581 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative  
582 image transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern  
583 recognition*, pp. 11315–11325, 2022.
- 584  
585 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
586 contrastive learning of visual representations. In *International conference on machine learning*,  
587 pp. 1597–1607. PMLR, 2020a.
- 588  
589 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of  
590 the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- 591  
592 Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum  
593 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- 594  
595 Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision  
596 transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp.  
597 9640–9649, 2021.

- 594 Xinlei Chen, Zhuang Liu, Saining Xie, and Kaiming He. Deconstructing denoising diffusion models  
595 for self-supervised learning. *arXiv preprint arXiv:2401.14404*, 2024.  
596
- 597 Yabo Chen, Yuchen Liu, Dongsheng Jiang, Xiaopeng Zhang, Wenrui Dai, Hongkai Xiong, and  
598 Qi Tian. Sdae: Self-distillated masked autoencoder. In *European conference on computer vision*,  
599 pp. 108–124. Springer, 2022.
- 600 Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need  
601 for semantic segmentation. *Advances in neural information processing systems*, 34:17864–17875,  
602 2021.  
603
- 604 Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-  
605 attention mask transformer for universal image segmentation. In *Proceedings of the IEEE/CVF*  
606 *conference on computer vision and pattern recognition*, pp. 1290–1299, 2022.
- 607 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo  
608 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban  
609 scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern*  
610 *recognition*, pp. 3213–3223, 2016.  
611
- 612 Timothée Darcet, Federico Baldassarre, Maxime Oquab, Julien Mairal, and Piotr Bojanowski.  
613 Cluster and predict latent patches for improved masked image modeling. *arXiv preprint*  
614 *arXiv:2502.08769*, 2025.
- 615 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-  
616 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
617 pp. 248–255. Ieee, 2009.
- 618 Chris HQ Ding, Tao Li, and Michael I Jordan. Convex and semi-nonnegative matrix factorizations.  
619 *IEEE transactions on pattern analysis and machine intelligence*, 32(1):45–55, 2008.  
620
- 621 Xiaoyi Dong, Jianmin Bao, Ting Zhang, Dongdong Chen, Weiming Zhang, Lu Yuan, Dong Chen,  
622 Fang Wen, and Nenghai Yu. Bootstrapped masked autoencoders for vision bert pretraining. In  
623 *European Conference on Computer Vision*, pp. 247–264. Springer, 2022.
- 624 David L Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–  
625 1306, 2006.  
626
- 627 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
628 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An  
629 image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*  
630 *arXiv:2010.11929*, 2020.
- 631 Alaaeldin El-Nouby, Michal Klein, Shuangfei Zhai, Miguel Angel Bautista, Alexander Toshev,  
632 Vaishaal Shankar, Joshua M Susskind, and Armand Joulin. Scalable pre-training of large au-  
633 toregressive image models. *arXiv preprint arXiv:2401.08541*, 2024.  
634
- 635 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image  
636 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recogni-*  
637 *tion*, pp. 12873–12883, 2021.
- 638 Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman.  
639 The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):  
640 303–338, 2010.  
641
- 642 Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
643 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
644 et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural*  
645 *information processing systems*, 33:21271–21284, 2020.
- 646 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
647 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
770–778, 2016.



- 702 Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, and Hervé Jégou. Spreading vectors for  
703 similarity search. *arXiv preprint arXiv:1806.03198*, 2018.  
704
- 705 Tara N Sainath, Brian Kingsbury, Vikas Sindhwani, Ebru Arisoy, and Bhuvana Ramabhadran. Low-  
706 rank matrix factorization for deep neural network training with high-dimensional output targets.  
707 In *2013 IEEE international conference on acoustics, speech and signal processing*, pp. 6655–  
708 6659. IEEE, 2013.
- 709 Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose,  
710 Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv*  
711 *preprint arXiv:2508.10104*, 2025.  
712
- 713 Korsuk Sirinukunwattana, Josien P. W. Pluim, Hao Chen, Xiaojuan Qi, Pheng-Ann Heng, Yun Bo  
714 Guo, Li Yang Wang, Bogdan J. Matuszewski, Elia Bruni, Urko Sanchez, Anton Böhm, Olaf  
715 Ronneberger, Bassem Ben Cheikh, Daniel Racoceanu, Philipp Kainz, Michael Pfeiffer, Martin  
716 Urschler, David R. J. Snead, and Nasir M. Rajpoot. Gland segmentation in colon histology im-  
717 ages: The glas challenge contest, 2016.
- 718 Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka,  
719 Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with  
720 learnable proposals. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
721 *recognition*, pp. 14454–14463, 2021.  
722
- 723 Chenxin Tao, Xizhou Zhu, Weijie Su, Gao Huang, Bin Li, Jie Zhou, Yu Qiao, Xiaogang Wang, and  
724 Jifeng Dai. Siamese image modeling for self-supervised vision representation learning. In *Pro-*  
725 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2132–  
726 2141, 2023.
- 727 Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical*  
728 *Society Series B: Statistical Methodology*, 58(1):267–288, 1996.  
729
- 730 Ivana Tošić and Pascal Frossard. Dictionary learning. *IEEE Signal Processing Magazine*, 28(2):  
731 27–38, 2011.
- 732 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.  
733
- 734 Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning  
735 for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer*  
736 *Vision and Pattern Recognition*, pp. 3024–3033, 2021.
- 737 Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu.  
738 Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF*  
739 *conference on computer vision and pattern recognition*, pp. 9653–9663, 2022.  
740
- 741 Xiaoyu Yang, Lijian Xu, Hongsheng Li, and Shaoting Zhang. One leaf reveals the season:  
742 Occlusion-based contrastive learning with semantic-aware views for efficient visual representa-  
743 tion. *arXiv preprint arXiv:2411.09858*, 2024.
- 744 Qihang Yu, Mark Weber, Xueqing Deng, Xiaohui Shen, Daniel Cremers, and Liang-Chieh Chen.  
745 An image is worth 32 tokens for reconstruction and generation. *Advances in Neural Information*  
746 *Processing Systems*, 37:128940–128966, 2024.  
747
- 748 Le Zhang, Qian Yang, and Aishwarya Agrawal. Assessing and learning alignment of unimodal  
749 vision and language models. In *Proceedings of the Computer Vision and Pattern Recognition*  
750 *Conference*, pp. 14604–14614, 2025.
- 751 Mingtian Zhang, Tim Z Xiao, Brooks Paige, and David Barber. Improving vae-based representation  
752 learning. *arXiv preprint arXiv:2205.14539*, 2022.  
753
- 754 Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10  
755 million image database for scene recognition. *IEEE transactions on pattern analysis and machine*  
*intelligence*, 40(6):1452–1464, 2017a.

756 Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene  
757 parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and*  
758 *pattern recognition*, pp. 633–641, 2017b.  
759  
760 Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot:  
761 Image bert pre-training with online tokenizer. *arXiv preprint arXiv:2111.07832*, 2021.  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

## 810 A STATEMENTS

### 811 A.1 ETHICS STATEMENT

812 This work adheres to the ICLR Code of Ethics. Our research is based on publicly available datasets  
813 (e.g., ImageNet-1K, ADE20K, Cityscapes, Pascal VOC) and does not involve human subjects, pri-  
814 vate data, or personally identifiable information. We follow standard licensing terms for all datasets  
815 used. The proposed framework, STELLAR, is a general-purpose method for self-supervised repre-  
816 sentation learning and is not designed for harmful or sensitive applications.  
817

### 818 A.2 REPRODUCIBILITY STATEMENT

819 We have made significant efforts to ensure the reproducibility of our results. Detailed descriptions of  
820 the STELLAR framework, training objectives, and experimental setups are provided in the main text  
821 and appendix. We report all datasets used, including preprocessing steps and evaluation protocols,  
822 and we describe ablation studies to clarify the contribution of each component. We will release the  
823 source code and pretrained model checkpoints after the internal review process to further facilitate  
824 reproducibility.  
825

### 826 A.3 USE OF LARGE LANGUAGE MODELS

827 Large language models (ChatGPT) were used exclusively for language refinement, including pol-  
828 ishing grammar, phrasing, and clarity of the manuscript. They were not used for research ideation,  
829 methodological design, experimental implementation, data analysis, or drawing conclusions. All  
830 scientific contributions of this work are entirely by the authors.  
831

## 832 B IMPLEMENTATION DETAILS

### 833 B.1 STELLAR TRAINING

834 We trained STELLAR with ViT models at size base, large, and huge, [along with the latent queries,](#)  
835 [projection layers, clustering head, and a 6-layer ViT decoder.](#) In the default setting, we initialized the  
836 [ViT part in the encoder](#) from public MAE checkpoint, and trained for 150 epochs for STELLAR-B,  
837 100 epochs for STELLAR-L, and 50 epochs for STELLAR-H. We used 16 NVIDIA A100-80GB  
838 with batch size 128 each, totaling 2048. We used AdamW(Loshchilov, 2017) with base learning  
839 rate  $1.5 \times 10^{-4}$  for STELLAR-B, and  $5 \times 10^{-5}$  for STELLAR-L and STELLAR-H.  
840

841 For concept clustering, we used 16384 prototypes for sparse and CLS tokens each. The projector  
842 is a 2-layer MLP before the prototype layer. We used 3 steps of Sinkhorn-Knopp algorithm. The  
843 temperature in sparse-dense cosine similarity softmax is 0.06. We used 6-8 random masked views  
844 to align the sparse tokens, and additional 6-8 local crops to align the CLS token.  
845

846 In the ablation study of model initialization and training strategy (Table 5), we trained the model  
847 from scratch and used exponential moving average (EMA) updated momentum encoder to encode  
848 the target prototype assignments in the warm-up stage. We EMA updated the full encoder (ViT,  
849 latent queries, projection, clustering head with momentum 0.996. The momentum encoder was used  
850 to encode a global view of the image into target prototype assignments, for both clustering loss and  
851 alignment loss. The masking ratio was 0.6 in the warm-up stage, and 0.8 during standard training.  
852 We trained the model with 150 epochs of EMA warm-up and 75 epochs of standard training.  
853

### 854 B.2 EVALUATION PROTOCOL

855 For STELLAR and all baseline models, we evaluated the frozen feature from the pretrained model  
856 with linear probing. We used layer norm in classification tasks, and batch norm in segmentation  
857 tasks, followed by a single linear layer predicting the class of the image or patch. For all benchmarks,  
858 we split 10% from the training set for validation. We tuned hyper-parameter with learning rate  
859  $1 \times 10^{-5}, 2 \times 10^{-5}, 5 \times 10^{-5}, 1 \times 10^{-4}, 2 \times 10^{-4}, 5 \times 10^{-4}, 1 \times 10^{-3}, 2 \times 10^{-3}, 5 \times 10^{-3}, 1 \times 10^{-2},$   
860 and batch size 64, 128, 256, 512, 1024, 2048, 4096, 8192.  
861



## C.2 SEMANTICS FROM DIFFERENT FEATURES

We conducted linear probing of different mean-pooled features of different types, and compared in Table 9. Sparse feature showed strongest global understanding quality.

Table 9: Semantics in different features

Feature	sparse	cls	dense
IN-1K lin. acc (%)	73.26	72.23	72.21

## C.3 CONCEPT ALIGNMENT WITH LANGUAGE

Inspired by Zhang et al. (2025), we used frozen feature from STELLAR and aligned with the text tower of CLIP (Radford et al., 2021) with a single attention pooled probing layer. The evaluation on vision language tasks with comparison to baseline models are shown in Table 10.

Table 10: Language alignment evaluation.

	IN-1K 0-shot		MS COCO		Winoground		MMVP
	@1	@5	T2I	I2T	Text	Image	Avg.
MAE	23.18	50.43	11.28	13.46	20.75	9.00	19.26
iBOT	50.01	80.43	20.79	29.38	24.75	12.00	18.52
STELLAR	51.53	80.04	17.94	22.34	26.25	8.25	19.26
CLIP	72.7	-	43.0	59.7	30.5	11.5	20.0

## C.4 FINETUNING

We performed finetuning for STELLAR on ImageNet-1K classification and ADE20K segmentation, and compared with baseline models. We used the same evaluation protocol as in Sec. B.2, with the backbone unfrozen and finetuned for 75 epochs. We used ViT-B for all models. The finetuning results are shown in Table 12. STELLAR showed consistent performance gain across different tasks, and close to the top model iBOT with slight difference.

Table 11: Finetuning performance in ImageNet-1K classification accuracy and ADE20K segmentation mIOU (%). We show in parentheses the gain over the respective linear probing results.

Model	ImageNet-1K Acc.	ADE20K mIOU
DINO	79.58 (+3.12)	39.22 (+12.35)
MAE	77.75 (+11.43)	40.33 (+9.42)
iBOT	80.72 (+9.14)	42.76 (+10.97)
STELLAR	80.05 (+6.78)	41.98 (+10.65)

## C.5 EFFICIENCY ANALYSIS

To analyze the efficiency of the STELLAR framework, we printed the processing time of the main components in the STELLAR framework with one A100 GPU at different batch sizes. Encoding the main global view of the image takes up most of the processing time, followed by encoding the masked views (8 views at 80% masking ratio) and decoding to the original image. The Sinkhorn-Knopp algorithm used for clustering and the Sinkhorn algorithm used in optimal transport matching take up much less amount of time, and their total processing time stay at similar level when increasing the batch size.

In comparison to the Sinkhorn matching algorithm we used in our experiments, we show the processing time using an alternative Hungarian matching algorithm commonly used in previous literature such as Sparse R-CNN (Sun et al., 2021), DETR (Carion et al., 2020) and MaskFormer (Cheng et al., 2021). As the implementation of the exact matching is not scalable with GPU parallelization, it’s computational time increases linearly with the batch size. At batch size 64, it is already 6 times of the encoder processing, while the Sinkhorn algorithm is over 100 times faster. For this reason, we added a small entropy regularization term in the bipartite matching objective, allowing us to use the Sinkhorn algorithm for efficient matching with GPU parallelization.

Table 12: Processing time (s) of the main components in the STELLAR framework with one A100 GPU at different batch sizes. In comparison to the Sinkhorn matching algorithm we used in our experiments, we show the processing time using an alternative Hungarian matching algorithm commonly used in previous literature (shown in gray).

Batch size	4	8	16	32	64
Encoder	$8.2 \times 10^{-3}$	$9.1 \times 10^{-3}$	$1.4 \times 10^{-2}$	$2.0 \times 10^{-2}$	$3.2 \times 10^{-2}$
Decoder	$4.6 \times 10^{-3}$	$6.8 \times 10^{-3}$	$8.8 \times 10^{-3}$	$1.2 \times 10^{-2}$	$1.5 \times 10^{-2}$
Mask encoding	$7.9 \times 10^{-3}$	$8.9 \times 10^{-3}$	$1.1 \times 10^{-2}$	$1.8 \times 10^{-2}$	$1.7 \times 10^{-2}$
SK clustering	$3.4 \times 10^{-4}$	$3.4 \times 10^{-4}$	$3.4 \times 10^{-4}$	$3.7 \times 10^{-4}$	$3.9 \times 10^{-4}$
<b>Sinkhorn matching</b>	$1.4 \times 10^{-3}$	$1.4 \times 10^{-3}$	$1.4 \times 10^{-3}$	$1.4 \times 10^{-3}$	$1.2 \times 10^{-3}$
Hungarian matching	$5.7 \times 10^{-3}$	$1.7 \times 10^{-2}$	$4.0 \times 10^{-2}$	$9.0 \times 10^{-2}$	$1.8 \times 10^{-1}$