# DSeq-JEPA: Discriminative Sequential Joint-Embedding Predictive Architecture

Xiangteng He[1,2*]    Shunsuke Sakai[3*]    Kun Yuan[4]
Nicolas Padoy[4]    Tatsuhito Hasegawa[3]    Leonid Sigal[1,2]
[1]University of British Columbia    [2]Vector Institute for AI
[3]University of Fukui    [4]University of Strasbourg
xiangteng.he@ubc.ca,mf240599@g.u-fukui.ac.jp

## Abstract

*Image-based Joint-Embedding Predictive Architecture (I-JEPA) learns visual representations by predicting latent embeddings of masked regions from visible context. However, it treats all regions uniformly and independently, lacking an explicit notion of where or in what order predictions should be made. Inspired by human visual perception, which deploys attention selectively and sequentially from the most informative to secondary regions, we propose **DSeq-JEPA**, a **D**iscriminative **Seq**uential **J**oint-**E**mbedding **P**redictive **A**rchitecture that bridges predictive and autoregressive self-supervised learning, integrating JEPA-style latent prediction with GPT-style sequential reasoning. Specifically, DSeq-JEPA (i) first identifies primary discriminative regions based on a transformer-derived saliency map, emphasizing the distribution of visual importance, and then (ii) predicts subsequent regions in this discriminative order, progressively forming a curriculum-like semantic progression from primary to secondary cues – a form of GPT-style pre-training. Extensive experiments across diverse tasks, including image classification (ImageNet), fine-grained visual categorization (iNaturalist21, CUB-200-2011, Stanford-Cars), detection and segmentation (MS-COCO, ADE20K), and low-level reasoning tasks (Clevr/Count, Clevr/Dist), demonstrate that DSeq-JEPA consistently focuses on more discriminative and generalizable representations than I-JEPA variants. Project page: https://github.com/SkyShunsuke/DSeq-JEPA.*
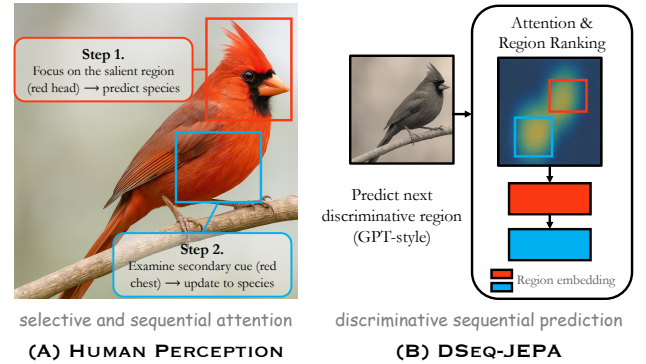
Figure 1. **(A) Humans** perceive visual scenes selectively and sequentially, focusing on discriminative regions such as the red head, and red chest of a Northern Cardinal. **(B) DSeq-JEPA** emulates this process by ranking regions by attention-derived importance and predicting each next discriminative region's embedding in a sequential, GPT-style manner.

## 1. Introduction

Self-supervised learning (SSL) has emerged as a powerful paradigm for learning robust and generalizable visual representations. SSL methods [1–4] generate supervisory signals directly from unlabeled data, drastically reducing dependence on costly annotation for training. Early approaches, such as contrastive learning [2, 3, 5], instance discrimination [6], or masked image modeling [1, 7, 8], have shown that representations learned in this way can rival or even surpass those trained with human supervision. Nevertheless, most existing SSL approaches either contrast different augmented views of the entire image or reconstruct masked pixels at the raw level. While the former emphasizes global invariance and the latter focuses on local appearance recovery, both objectives primarily operate at either the image or pixel level, offering little incentive to model how information is organized within an image. Consequently, they tend to capture either holistic or low-level statistics rather than structured, semantic relationships across regions.

Recently, Joint-Embedding Predictive Architectures (JEPAs) [9] have revisited the foundations of SSL. Instead of contrasting instances or reconstructing pixels, JEPAs directly predict latent embeddings of masked regions from visible context, offering a scalable and information-
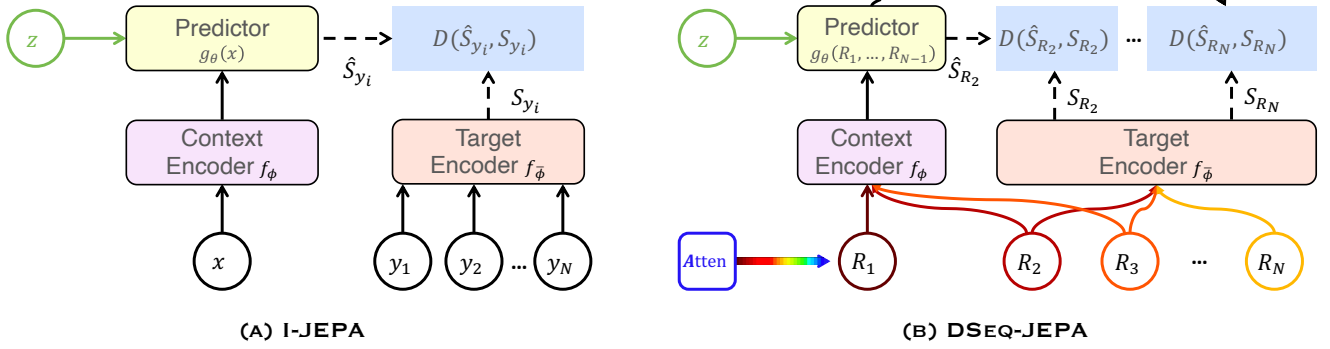
---
*Equal contribution.

Figure 2. **(A) I-JEPA** learns to predict the embeddings of the target regions $(y_1, \ldots, y_N)$ from a single context region $x$, using a predictor network conditioned on latent variables $z$. While **(B) DSeq-JEPA** learns to predict the embeddings of the next discriminative regions $\{R_2, \ldots, R_N\}$ based on its sequence of pre-identified regions in a sequential manner, also using a predictor. The order of the discriminative regions is determined by the attention map of the image.

preserving formulation. The image-based variant, I-JEPA [10], learns *what* to represent by modeling context–target relationships within the same image and has demonstrated impressive performance across a range of visual tasks. Its idea is simple: from a single context region, predict the embeddings of various target regions in the same image, as shown in Figure 2 (A). However, I-JEPA treats all regions equally and lacks an explicit notion of *where* or *in what order* predictions should be made, overlooking that visual information is inherently structured.

Humans, by contrast, perceive scenes *selectively* and *sequentially*, as shown in Figure 1. Visual attention acts as a dynamic filtering mechanism that enables humans to efficiently process the vast amount of information in complex scenes. Models of human vision [11] suggest that human gaze follows a series of information-seeking fixations, with each next fixation location attempting to minimize uncertainty (maximize information) about the stimulus. In other words, our fixations are sequential and focus on the most discriminative regions, shifting attention from primary to secondary cues, naturally guiding fine-grained representation learning. Motivated by this, we ask:

*Can predictive self-supervised models learn not only what to predict, but also where and in what order to predict?*

To answer this question, we revisit predictive self-supervised learning through the lens of discriminative sequencing, moving beyond static prediction toward an ordered reasoning process. Accordingly, we propose **DSeq-JEPA** (Discriminative Sequential Joint-Embedding Predictive Architecture), a new paradigm that learns not only what to predict, but also where and in what order, introducing a discriminative sequential inductive bias (DSeq Bias) that is (loosely) motivated by human perception. This paradigm bridges two previously distinct lines of self-supervised learning: *predictive modeling* (JEPA-style) and *autoregres-*

*sive reasoning* (GPT-style). In doing so, DSeq-JEPA reformulates representation learning as a progressive discovery process of discriminative regions in semantic space through two synergistic components, as illustrated in Figure 2.

- **Selective Region Gazing** *(Where to attend)*: Unlike I-JEPA's uniform region sampling, DSeq-JEPA dynamically identifies and prioritizes *discriminative* regions (*e.g.*, bird beak, wing patterns) using a similarity-based saliency map. These regions are selectively attended and iteratively refined during pre-training, mimicking human visual attention that focuses on information gain.
- **Sequential Region Prediction** *(In what order)*: Departing from I-JEPA's flat, independent prediction of regions (BERT-style pre-training [12]), DSeq-JEPA employs a successive prediction process (GPT-style pre-training [13]) where the model sequentially predicts embeddings for regions from the most discriminative to the least. This encourages semantically structured representations, where discriminative cues are captured early, while contextual details are integrated progressively.

We evaluate DSeq-JEPA across diverse visual tasks. On image classification (ImageNet [14]) and fine-grained visual categorization (iNaturalist21 [15], CUB-200-2011 [16], Stanford Cars [17]), DSeq-JEPA consistently outperforms JEPA-based methods. It also generalizes well to dense prediction tasks such as detection and segmentation (MS-COCO [18], ADE20K [19]), and object counting and depth prediction (Clevr/Count, Clevr/Dist [20]). Qualitative analyses further reveal that DSeq-JEPA naturally attends to and predicts next discriminative regions, producing interpretable, fine-grained representations. Overall, our results indicate that predictive self-supervised learning benefits from breaking permutation symmetry over regions and imposing a semantically meaningful prediction order, making DSeq-JEPA a first step towards structured, region-wise

2

autoregressive reasoning within the JEPA framework.

## 2. Related Work

Self-supervised learning (SSL) has gained significant attention in recent years, particularly in the context of visual representation learning. Existing SSL methods can be cast into the framework of Energy-based Models (EBMs) [21], and partitioned into three categories: Joint-Embedding Architectures (JEAs), Generative Architectures (GAs), and Joint-Embedding Predictive Architectures (JEPAs).

**JEAs** are designed to bring the embeddings closer together for compatible inputs, and push the embeddings further apart for incompatible ones. JEA methods include: (i) contrastive learning approaches, *e.g.*, SimCLR [3], MoCo [2], MoCo v2 / v3 [5, 22], SwAV [23], form the foundation of this paradigm, optimizing instance-level discrimination via large batches or memory banks; and (ii) non-contrastive approaches, *e.g.*, BYOL [24], SimSiam [4], further show that meaningful representations can emerge even without negative samples, relying instead on asymmetric prediction and momentum encoders to avoid collapse.

**GAs** aim to reconstruct a target signal $y$ from a related input signal $x$ by employing a decoder network. This process often incorporates additional variables $z$, which may be latent, to enhance the reconstruction. Masked image modeling serves as a prominent technique in this category, with notable examples that include MAE [1], BEiT [7], SimMIM [8], and iGPT [25], where models learn to fill in masked regions of an image. Such methods encourage learning low- and mid-level features that are useful for reconstruction but may not directly align with semantic discriminability.

**JEPAs** [9], on the other hand, formulates self-supervised learning as a predictive problem in latent embedding space, *i.e.*, predicting the embedding of a target signal $y$ from a compatible input signal $x$ through a predictor network that leverages additional (potentially latent) variables $z$ to improve the prediction process. As an instantiation of this architecture, I-JEPA [10] is proposed to improve the semantic level of self-supervised representations. Given an image, a context encoder processes the visible regions while a target encoder extracts latent embeddings for the masked ones. A lightweight predictor then predicts the target embeddings based on the context features, and training minimizes a latent-space regression loss between the predicted and true embeddings. This design discards pixel-level reconstruction and contrastive pairing, enabling scalable, information-preserving pretraining that focuses on high-level semantics rather than low-level texture statistics. C-JEPA [26] integrates I-JEPA with contrastive objectives to stabilize training and better align invariances. These predictive architectures demonstrate strong abstraction capabilities but treat all regions uniformly, lacking an explicit mechanism to determine which regions are most discriminative or in what or-

der to predict them. Our DSeq-JEPA is based on I-JEPA but aims to improve the capabilities of JEPA-based methods for fine-grained representation learning. Compared with I-JEPA, (i) DSeq-JEPA biases prediction toward discriminative regions by prioritizing where to attend, whereas I-JEPA samples target regions uniformly without regard to their discriminativeness, and (ii) DSeq-JEPA performs prediction in a sequential and discrimination-guided manner, explicitly modeling in what order to predict, whereas I-JEPA follows a flat, unordered prediction process.

## 3. Methodology

Our goal is to enhance self-supervised learning by introducing a preferential strategy for predictive region selection and by imposing an explicit order on these predictions, inspired by human visual perception. To this end, we introduce DSeq-JEPA, which operationalizes the DSeq Bias within the I-JEPA family, extending it toward selective and sequential discriminative reasoning, thereby bridging predictive and autoregressive paradigms in self-supervised learning. Figure 3 illustrates that DSeq-JEPA consists of three modules: *(i) Attention Generation:* Produces a spatial saliency map that highlights informative regions. *(ii) Discriminative Region Selection:* Identifies significant regions adaptively based on the saliency map, determining where to look. *(iii) Next Discriminative Region Prediction:* Guides the model to learn fine-grained representations by focusing on pre-identified regions, modeling in what order to predict.

### 3.1. Model Architecture

Following I-JEPA [10], our model also has three modules: a context encoder $f_\phi$, a target encoder $f_{\bar\phi}$, and a predictor $g_\theta$, all built upon the Vision Transformer (ViT) [27]. The context and target encoders follow the standard ViT design, while the predictor is implemented as a lightweight ViT variant. In our experiments, we employ ViT-B/16 and ViT-L/16 as backbones to verify the effectiveness and scalability of DSeq-JEPA. While the overall architecture closely follows I-JEPA, we append an auxiliary class token to both the context and target encoders following MAE [1] to maintain architectural consistency with ViT-based downstream tasks. This token acts as a stable semantic reference that slightly improves the reliability of attention-based region selection without altering the predictive formulation or training objective of I-JEPA. Conceptually, it serves only as a lightweight anchor for attention aggregation rather than a structural modification. We verified that this modification alone brings no gain in Section 4.5.

### 3.2. Attention Generation

To determine where to attend, we generate a saliency map that highlights informative regions in the input image. Specifically, we pass the image through the target encoder
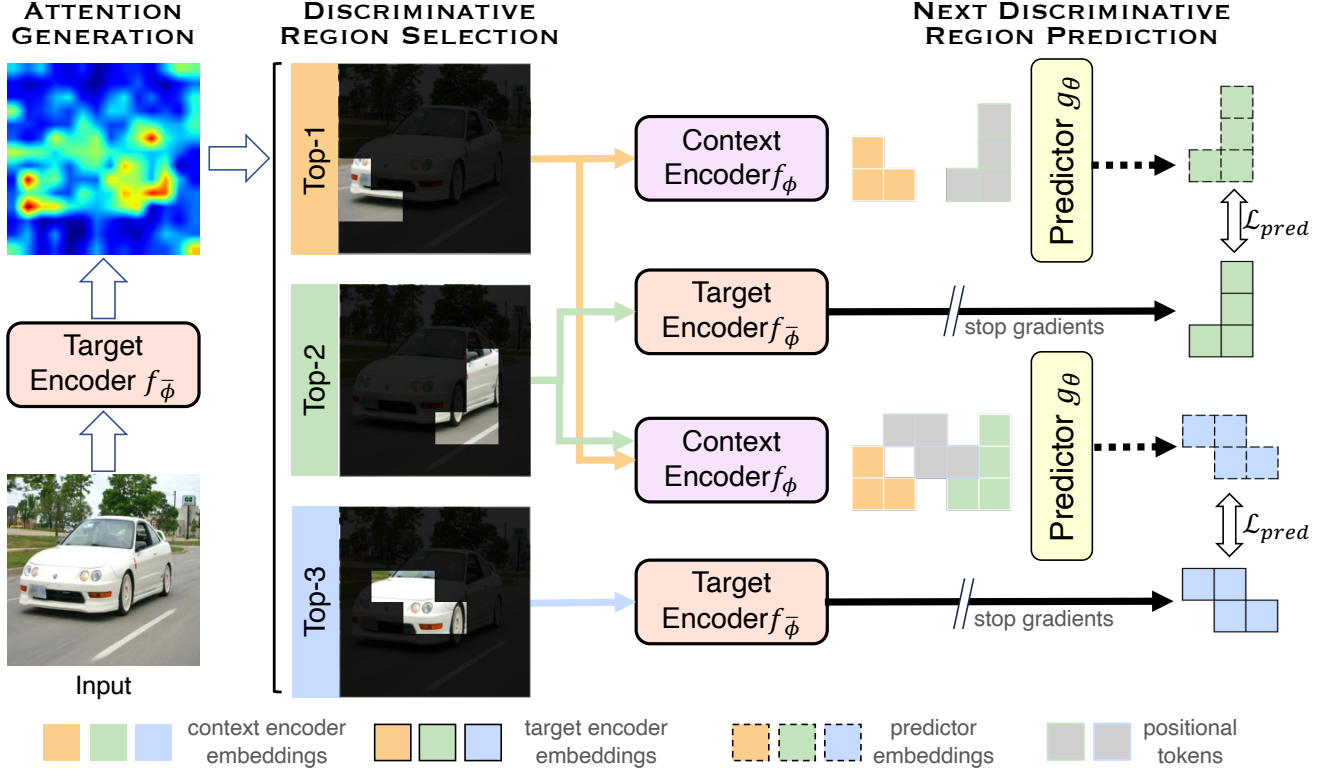
**ATTENTION GENERATION**

**DISCRIMINATIVE REGION SELECTION**

**NEXT DISCRIMINATIVE REGION PREDICTION**

Target Encoder $f_{\bar{\phi}}$

Input

Top-1

Top-2

Top-3

Context Encoder $f_{\phi}$

Target Encoder $f_{\bar{\phi}}$

Context Encoder $f_{\phi}$

Target Encoder $f_{\bar{\phi}}$

Predictor $g_{\theta}$

Predictor $g_{\theta}$

stop gradients

stop gradients

$\mathcal{L}_{pred}$

$\mathcal{L}_{pred}$

context encoder embeddings

target encoder embeddings

predictor embeddings

positional tokens

Figure 3. **Overview of DSeq-JEPA.** We compute a saliency map for each input image, identify the Top-$N$ high-response regions (with $N$=3 for illustration), and feed them to the predictor sequentially. During prediction, each region's embedding is predicted from its preceding regions and positional tokens, and aligned with its target encoder embedding. All encoders and predictors adopt ViT [27] architecture.

$f_{\bar{\phi}}$ and compute a similarity map between the auxiliary class token and all patch embeddings at a specific transformer block. The resulting similarity values indicate each patch's contribution to the global semantics, where regions with higher value correspond to more informative content relevant for understanding the image.

Formally, given an image $\mathbf{x} \in \mathbb{R}^{C \times H \times W}$, we feed it into the target encoder $f_{\bar{\phi}}$, obtaining feature embeddings $\mathbf{s} = f_{\bar{\phi}}(\mathbf{x}) \in \mathbb{R}^{D \times h \times w}$, where $D$ is the embedding dimension, and $h, w$ denote the number of vertical and horizontal patches, respectively. We compute a similarity map $\mathbf{A} \in \mathbb{R}^{h \times w}$ between the embeddings of the class token and all patches at the $l$-th transformer block. While we use this simple similarity-based attention for region selection, the formulation is general and can accommodate other attention mechanisms that provide spatial guidance.

### 3.3. Discriminative Region Selection

Given the saliency map $\mathbf{A}$ from the target encoder, we extract a set of high-saliency regions that indicate discriminative visual evidence. Local maxima of $\mathbf{A}$ are detected and ranked, yielding a series of discriminative regions. This procedure treats internal attention as a learnable notion of

visual importance, enabling the model to localize informative areas without any external saliency supervision. Recent work [28, 29] further supports this view, showing that learnable tokens in ViTs can implicitly function as information banks, an emergent, self-organizing property of transformer representations. We first normalize the saliency map:

$$\tilde{\mathbf{A}} = \frac{\mathbf{A} - \min(\mathbf{A})}{\max(\mathbf{A}) - \min(\mathbf{A})}. \quad (1)$$

An adaptive threshold $\tau$ is then determined by Otsu's method [30], yielding a binary mask:

$$\mathbf{M}(i, j) = \mathbb{1}\left[\tilde{\mathbf{A}}(i, j) \geq \tau\right]. \quad (2)$$

We apply connected-component labeling to $\mathbf{M}$ under a relaxed connectivity criterion (8-neighborhood) to obtain the candidate set $\{R_k\}$. Under this relaxed condition, each $R_k$ may consist of multiple spatially separated (mutually disjoint) sub-regions. Small components ($|R_k| < \alpha hw$) are discarded, where $\alpha$ is a small constant. Each region is assigned a discriminative score:

$$\rho_k = \frac{1}{|R_k|} \sum_{(i,j) \in R_k} \tilde{\mathbf{A}}(i, j). \quad (3)$$

4

Unlike prior works [9, 26], we employ irregular, non-overlapping masks for the extracted regions. We select the top-$(N-1)$ regions with the highest $\rho_k$ values and form $\mathcal{R} = \{R_k\}_{k=1}^N$, where $N$=5 following I-JEPA. The last region $R_N$ is defined as the remaining area: $R_N = \Omega \setminus \bigcup_{k=1}^{N-1} R_k$, where $\Omega$ denotes the full image region. This ensures complete coverage of the image without gaps, yielding more stable and discriminative region-based representation learning.

Since saliency maps are often noisy during the early stages of pre-training, we adopt a probabilistic curriculum to stabilize optimization. The probability $\lambda$ of applying discriminative selection is linearly increased from 0 to 1 over pre-training epochs, so that for each sample we apply I-JEPA's random region sampling with probability $1 - \lambda$ and our discriminative selection with probability $\lambda$. This progressive transition helps the model evolve from diffuse attention to structured saliency, effectively allowing it to learn where to look.

### 3.4. Next Discriminative Region Prediction

In this stage, we enhance self-supervised representations through a sequential predictive mechanism that models discriminative regions in an autoregressive manner. Specifically, given a set of the most discriminative regions in an image $\{R_k\}_{k=1}^N$, where the index reflects their ordering by discriminativeness, our goal is to enable our predictor to possess such an ability that predicts the representation of the $(k+1)$-th most discriminative region, $R_{k+1}$, conditioned on the previously encoded top-$k$ discriminative regions. The predicted representation is expected to align closely with its true representation, *i.e.*, the output embedding obtained by feeding $R_{k+1}$ to the target encoder. In this process, we predict the representations of regions 2 through $N$, each conditioned on all preceding regions in the sequence. We formulate this prediction process as follows. For each step $k \in \{1, \ldots, N-1\}$, the predictor $g_\theta$ estimates the $(k+1)$-th region embedding as

$$\hat{\mathbf{s}}_{R_{k+1}} = g_\theta(\mathbf{p}_{R_{k+1}}, \mathbf{s}_{R_1}, \ldots, \mathbf{s}_{R_k}), \qquad (4)$$

where $\mathbf{s}_{R_1}, \ldots, \mathbf{s}_{R_k}$ denotes the representations of the $1, \ldots, k$-th most discriminative regions from the context encoder, respectively, $\mathbf{p}_{R_{k+1}}$ is the mask positional token for the predicted discriminative region, and $N$ indicates the number of most discriminative regions. $g_\theta$ is a learned predictor, and $\hat{\mathbf{s}}_{R_{k+1}}$ is the predicted representation of the next discriminative region, which is expected to align closely with its true representation $\mathbf{s}_{R_{k+1}}$ from the target encoder.

Unlike I-JEPA's flat predictions, this formulation introduces causal dependencies among regions, mirroring the progressive reasoning by which humans integrate salient cues into global understanding. It thus unifies JEPA-style latent prediction with GPT-style autoregression, enabling a progressive reasoning process over discriminative regions. From an information-theoretic perspective, sequential prediction can be viewed as an active information-seeking process: each step attends to the region expected to contribute the most new information given the previously encoded regions. This progressive reduction of representational uncertainty aligns with human visual perception, where attention shifts sequentially to maximize information gain.

### 3.5. Optimization Objective

To enforce the alignment between the predicted representation $\hat{\mathbf{s}}_{R_{k+1}}$ and its ground-truth counterpart $\mathbf{s}_{R_{k+1}}$, we adopt the smoothed $\ell_1$ (Huber) loss:

$$\mathcal{L}_{\text{pred}} = \frac{1}{N} \sum_{k=1}^N \sum_{j=1}^D \psi\big((\hat{s}_{R_{k+1}} - s_{R_{k+1}})_j\big), \qquad (5)$$

$$\psi(x) = \begin{cases} \frac{1}{2} x^2, & |x| < \delta, \\ \delta\left(|x| - \frac{1}{2}\delta\right), & |x| \geq \delta, \end{cases} \quad \delta = 1. \qquad (6)$$

This objective encourages stable yet discriminative alignment between predicted and target embeddings, making the sequential prediction both smooth and robust to outliers.

## 4. Experiments

### 4.1. Evaluation Tasks

Following I-JEPA and C-JEPA, we pre-train all models on ImageNet-1K [14], and conduct linear probing and fine-tuning across diverse benchmarks to assess DSeq-JEPA's generalization and discriminative ability: **(1) Image classification:** linear probing evaluation on *ImageNet-1K* [14]. **(2) Fine-grained visual categorization:** *iNaturalist21 (iNat21)* [15], *CUB-200-2011 (CUB)* [16], and *Stanford-Cars (Cars)* [17] to test fine-grained recognition that is sensitive to local regions. **(3) Object detection and instance segmentation**: transfer learning on *MS-COCO* [18] using Mask R-CNN [31] as the detector. **(4) Semantic segmentation**: evaluation on *ADE20K* [19] with the Semantic FPN and UPerNet head [32] as in [1]. **(5) Low-level reasoning**: *Clevr/Count* (object counting) and *Clevr/Dist* (depth prediction) [20] to assess structured scene understanding.

### 4.2. Implementation Details

Our pre-training and downstream evaluation strictly follow the configurations of I-JEPA [10] and C-JEPA [26]. Further details are presented in the supplementary material.

### 4.3. Main Results

Across diverse tasks ranging from coarse-grained recognition to fine-grained categorization, dense prediction, and structured reasoning, DSeq-JEPA consistently achieves the best performance among I-JEPA variants, demonstrating its strong generalization and representational capability.

5

Table 1. **Linear probing evaluation compared with SOTA methods.** We report the Top-1 accuracy of linear probing on image classification (ImageNet) and fine-grained visual categorization (iNat21, CUB, Cars) using the pretrained models, *i.e.,* ViT-B/16 and ViT-L/16. *We use authors' official implementations when available, otherwise we re-implement methods from the papers. †: official code; ‡: our re-implementation; §: numbers reported by the original papers. All models follow the same training/eval protocol.*

| Method | Arch. | Epochs | Image Classification ImageNet [14] | Fine-grained Visual Categorization iNat21 [15] | CUB [16] | Cars [17] | Avg. |
|---|---|---|---|---|---|---|---|
| I-JEPA§ [10] | ViT-B/16 | 600 | 72.9 | - | - | - | - |
| I-JEPA† | ViT-B/16 | 600 | 72.4 | 35.9 | 65.3 | 65.9 | 59.9 |
| **DSeq-JEPA** | ViT-B/16 | 600 | **73.5** $_{+1.1}$ | **36.4** $_{+0.5}$ | **66.2** $_{+0.9}$ | **67.3** $_{+1.4}$ | **60.9** $_{+1.0}$ |
| C-JEPA§ [26] | ViT-B/16 | 600 | 73.7 | - | - | - | - |
| C-JEPA‡ | ViT-B/16 | 600 | 73.5 | 36.2 | 64.9 | 66.1 | 60.2 |
| **DSeq-JEPA+Contrastive** | ViT-B/16 | 600 | **73.8** $_{+0.3}$ | **36.6** $_{+0.4}$ | **66.5** $_{+1.6}$ | **67.4** $_{+1.3}$ | **61.1** $_{+0.9}$ |
| I-JEPA§ [10] | ViT-L/16 | 600 | 77.5 | - | - | - | - |
| I-JEPA† | ViT-L/16 | 600 | 77.1 | 38.8 | 66.9 | 68.1 | 62.7 |
| **DSeq-JEPA** | ViT-L/16 | 600 | **77.9** $_{+0.8}$ | **39.5** $_{+0.7}$ | **68.1** $_{+1.2}$ | **68.9** $_{+0.8}$ | **63.6** $_{+0.9}$ |
| C-JEPA§ [26] | ViT-L/16 | 600 | 78.1 | - | - | - | - |
| C-JEPA‡ | ViT-L/16 | 600 | 78.0 | 39.0 | 65.8 | 67.7 | 62.6 |
| **DSeq-JEPA+Contrastive** | ViT-L/16 | 600 | **78.4** $_{+0.4}$ | **39.7** $_{+0.7}$ | **68.3** $_{+2.5}$ | **68.8** $_{+1.1}$ | **63.8** $_{+1.2}$ |

### 4.3.1. Image Classification

Following I-JEPA and C-JEPA, to validate the effectiveness of DSeq-JEPA, we first conduct experiments on ImageNet-1K under linear probing. Linear probing is a commonly adopted evaluation protocol for self-supervised learning, where a linear classifier is trained on top of the frozen features extracted by the pre-trained model. The classification accuracy reflects the quality of the representation learned through the SSL approach. As shown in Table 1, using a ViT-B/16 backbone, DSeq-JEPA achieves 73.5% Top-1 accuracy, outperforming I-JEPA by +1.1%. The gain is consistent when scaling to the ViT-L/16 model, where DSeq-JEPA reaches 77.9%, exceeding I-JEPA by +0.8%. These improvements confirm that modeling where to attend and in what order to predict leads to more discriminative latent representations, effectively balancing global semantic abstraction and local discriminative focus during pre-training.

C-JEPA combines I-JEPA with contrastive learning. DSeq-JEPA achieves the same or nearly identical results (equal on ViT-B/16 and only 0.1% lower on ViT-L/16). When augmented with the same contrastive regularization, DSeq-JEPA+Contrastive outperforms C-JEPA, highlighting that discriminative region selection and sequential prediction offer complementary benefits and stronger robustness.

We highlight that while average improvements of +0.9% to +1.2% may appear modest, they are consistent and significant (both in statistical and application sense). For comparison, the improvements of C-JEPA [26] over I-JEPA [10] on ImageNet are smaller (+0.8% with ViT-B/16 and +0.6% with ViT-L/16); the average improvement on C-JEPA over I-JEPA is even smaller (+0.3% and -0.1%), so about 1/3 of

what we achieve with DSeq-JEPA.

### 4.3.2. Fine-grained Visual Categorization (FGVC)

The fundamental challenge in the FGVC task lies in distinguishing visually similar sub-categories (*e.g.*, within birds, Indigo Bunting vs. Lazuli Bunting), which requires models to capture subtle discriminative part features such as feather coloration, beak morphology, and tail length [33]. To test the discriminative precision of learned representations, we further evaluate on FGVC benchmarks, including iNaturalist21 (10,000 distinct species) [15], CUB-200-2011 (200 bird species) [16], and Stanford-Cars (196 car models) [17]. DSeq-JEPA consistently outperforms I-JEPA and C-JEPA across all FGVC datasets. For ViT-B/16, it achieves 36.4% on iNat21, 66.2% on CUB, and 67.3% on Cars, surpassing I-JEPA by +0.5%, +0.9%, and +1.4%, and C-JEPA by +0.2%, +1.3%, and +1.2%. When equipped with the same contrastive regularization as C-JEPA, DSeq-JEPA+Contrastive further improves to 36.6% / 66.5% / 67.4%, outperforming C-JEPA by +0.4%, +1.6%, and +1.3%, respectively. Similar trends are observed for ViT-L/16, where DSeq-JEPA and DSeq-JEPA+Contrastive achieve better performance than I-JEPA and C-JEPA. Overall, the consistent performance across both ImageNet and FGVC benchmarks demonstrates the generality of DSeq-JEPA. By modeling where and in what order to attend, DSeq-JEPA learns representations that are not only effective for coarse-grained categorization but also highly transferable to tasks requiring fine-grained discrimination.

Table 2. **Performance on detection and segmentation tasks.** We perform MS-COCO detection/segmentation, and ADE20K semantic segmentation on pre-trained ViT-B/16 models, and report $AP^{box}$, $AP^{mask}$ on MS-COCO, and mIoU on ADE20K. §: numbers reported by C-JEPA [26].

| Method | $AP^{box}$ | $AP^{mask}$ | mIoU |
|---|---|---|---|
| I-JEPA§ [10] | 49.9 | 44.5 | 47.6 |
| **DSeq-JEPA** | **50.5** $_{+0.6}$ | **45.0** $_{+0.5}$ | **48.1** $_{+0.5}$ |
| C-JEPA§ [26] | 50.7 | 45.3 | 48.7 |
| **DSeq-JEPA+Contrastive** | **50.9** $_{+0.2}$ | **45.7** $_{+0.4}$ | **48.9** $_{+0.2}$ |

Table 3. **Performance on object counting (Clevr/Count) and depth prediction (Clevr/Dist)** on pre-trained ViT-L/16 models. §: numbers reported by C-JEPA [26].

| Method | Clevr/Count | Clevr/Dist |
|---|---|---|
| I-JEPA§ [10] | 85.6 | 71.2 |
| **DSeq-JEPA** | **86.4** $_{+0.8}$ | **71.5** $_{+0.3}$ |
| C-JEPA§ [26] | 86.8 | 71.6 |
| **DSeq-JEPA+Contrastive** | **87.1** $_{+0.3}$ | **71.9** $_{+0.3}$ |

### 4.3.3. Detection and Segmentation

To further evaluate the transferability of the learned representations to dense prediction tasks, we perform object detection and instance segmentation on MS-COCO and semantic segmentation on ADE20K using the pre-trained ViT-B/16 models. As shown in Table 2, DSeq-JEPA consistently outperforms I-JEPA across all metrics, achieving +0.6 $AP^{box}$, +0.5 $AP^{mask}$, and +0.5 mIoU improvements. When equipped with the same contrastive regularization, DSeq-JEPA+Contrastive attains the best overall results. While the absolute gains are moderate, this is expected given the strong C-JEPA baseline and the limited sensitivity of dense prediction heads to pre-training variations. Importantly, the consistent improvements across both detection and segmentation demonstrate that discriminative sequential prediction promotes more spatially structured and semantically coherent representations, which generalize well beyond classification to pixel-level prediction.

### 4.3.4. Low-level Reasoning

As shown in Table 3, on low-level reasoning tasks, DSeq-JEPA with ViT-L/16 yields 86.4 on Clevr/Count and 71.5 on Clevr/Dist, improving over I-JEPA by +0.8 and +0.3, respectively. With contrastive regularization, DSeq-JEPA+Contrastive reaches 87.1 / 71.9, surpassing C-JEPA. The fact that improvements persist on these low-level metrics indicates that DSeq-JEPA does not merely enhance semantic categorization, but also strengthens geometric/structural priors in the representation.

### 4.4. Ablation Study

We conduct a series of ablation experiments to analyze the contribution of each component in DSeq-JEPA. Specifically, we study how selective region generation and sequential prediction contribute to overall performance, and whether their combination produces synergistic effects.

Table 4 summarizes the impact of each component using linear probing on ImageNet and iNat21 with pre-trained ViT-B/16 model. The baseline configuration (uniform sampling + flat prediction) corresponds to I-JEPA, where all regions are treated equally and predicted in parallel. Enabling only one of the two components degrades performance. Sequential prediction alone (uniform + sequential) slightly degrades accuracy. This indicates that, without an informative order, autoregression introduces noise: the model commits to an arbitrary prediction path that is not aligned with discriminative content. Conversely, discriminative selection alone (selective + flat) also underperforms, while the model focuses on discriminative regions, predicting them in parallel collapses the structure of inter-region dependencies. From a learning dynamics perspective, our sequential prediction can be viewed as an easy-to-hard curriculum. The most discriminative regions contain the most informative cues, making them easier to predict and more stable. By conditioning subsequent predictions on these strong representations, the model gradually learns to infer less discriminative and more context-dependent regions. In contrast, random or unordered prediction (selective + random sequential) disrupts this natural progression, forcing the model to learn from mixed difficulty signals and thus degrading performance. Importantly, combining both (selective + sequential), DSeq-JEPA achieves the best results. These findings demonstrate that the two modules are *not independently effective, but rather synergistic*: discriminative region selection provides a semantically meaningful trajectory (where to attend), and sequential prediction exploits this trajectory to capture region-to-region causal dependencies (in what order to predict). Together they yield more discriminative and fine-grained representations.

### 4.5. Effect of Class Token

To ensure architectural consistency with I-JEPA, we introduce an auxiliary class token solely as an anchor for generating saliency maps, without adding extra layers or parameters. As shown in Table 5, incorporating the [CLS] token does not improve I-JEPA's performance, demonstrating that the improvements of DSeq-JEPA originate from the proposed discriminative and sequential learning mechanisms, rather than from any architectural modification.

Table 4. **Ablation study.** We perform linear probing on ImageNet and iNat21 using the ViT-B/16 as backbone. The first row corresponds to the *I-JEPA baseline*, and the last row is our *DSeq-JEPA*.

| Region Generation | | Prediction Strategy | | ImageNet [14] | iNat21 [15] |
| --- | --- | --- | --- | --- | --- |
| *Uniform sampling* | *Discriminative selection* | *Flat* | *Sequential* | | |
| ☑ | | ☑ | | 72.4 | 35.9 |
| ☑ | | | ☑ | 72.3 | 34.9 |
| | ☑ | ☑ | | 72.0 | 35.7 |
| | ☑ | | ☑ (random) | 71.7 | 35.0 |
| | ☑ | | ☑ | **73.5** | **36.4** |

Table 5. **Effect of class token.**

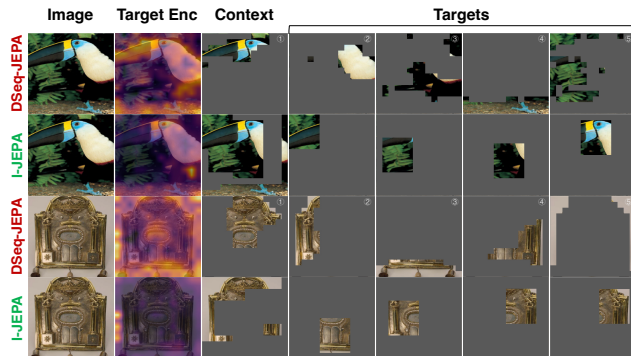| Method | I-JEPA | I-JEPA w/ `[CLS]` |
| --- | --- | --- |
| **ImageNet** | 72.4 | 72.4 |



Figure 4. **Qualitative visualization of learned attention and selected context/target regions using ViT-B/16 model.** For DSeq-JEPA, the numbered regions (①–⑤) correspond to discriminative regions ordered by their estimated importance.

## 4.6. Visualization Analysis

**Comparison with I-JEPA.** As shown in Figure 4, I-JEPA produces relatively diffuse attention patterns and rectangular target masks, often covering multiple unrelated and overlapped areas. In contrast, DSeq-JEPA generates compact, semantically meaningful regions with order, progressively focusing on the most discriminative parts of an object (*e.g.*, the bird's beak and forehead), leading to a clearer region correspondence and improved interpretability.

**Evolution during Pre-training.** Figure 5 visualizes how patch-level clustering in the learned representation space evolves as pre-training progresses. This reflects how DSeq-JEPA's representations self-organize during pre-training, where patches with similar semantics gradually become closer, even under a fully unsupervised setting. As pre-training progresses, from early stages to later ones, the initially fragmented and noisy groupings gradually evolve into



Figure 5. **Evolution of patch-level clustering during pre-training.** From left to right: input image and 4-cluster results from ViT-B/16 checkpoints at 150, 300, 450, and 600 epochs.

coherent, object-aligned clusters. This emergent behavior indicates that DSeq-JEPA explicitly captures fine-grained semantic structure through its attention-guided learning dynamics, without requiring any supervision or external constraints.

## 5. Conclusion

We proposed DSeq-JEPA, extending JEPAs with selective and sequential prediction mechanisms. By explicitly modeling where to attend and in what order to predict, DSeq-JEPA progressively activates discriminative cues and cap-

tures structured, fine-grained representations. Comprehensive evaluations across multiple benchmarks demonstrate consistent gains over JEPAs, confirming DSeq-JEPA's robustness and generality. Our findings highlight that selective attention and sequential reasoning are synergistic inductive biases for predictive representation learning, aligning more closely with human perception. In future work, we will explore integrating DSeq-JEPA into vision–language pre-training to enable more structured visual–textual alignment and selective cross-modal grounding.

# References

[1] 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. 1, 3, 5

[2] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 1, 3

[3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR, 2020. 1, 3

[4] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15750–15758, 2021. 1, 3

[5] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020. 1, 3

[6] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3733–3742, 2018. 1

[7] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021. 1, 3

[8] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9653–9663, 2022. 1, 3

[9] Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62(1):1–62, 2022. 1, 3, 5

[10] Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15619–15629, 2023. 2, 3, 5, 6, 7, 1

[11] Laura Renninger, James Coughlan, Preeti Verghese, and Jitendra Malik. An information maximization model of eye movements. *Advances in neural information processing systems*, 17, 2004. 2

[12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019. 2

[13] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. 2

[14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Proceedings of IEEE/CVF conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2, 5, 6, 8, 1

[15] Grant Van Horn and macaodha. inat challenge 2021 - fgvc8. https://kaggle.com/competitions/inaturalist-2021, 2021. Kaggle. 2, 5, 6, 8, 1

[16] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 2, 5, 6, 1

[17] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE/CVF international conference on computer vision workshops*, pages 554–561, 2013. 2, 5, 6, 1

[18] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European conference on computer vision*, pages 740–755. Springer, 2014. 2, 5, 1

[19] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of IEEE/CVF conference on computer vision and pattern recognition*, pages 633–641, 2017. 2, 5

[20] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of IEEE/CVF conference on computer vision and pattern recognition*, pages 2901–2910, 2017. 2, 5, 1

[21] Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, Fujie Huang, et al. A tutorial on energy-based learning. *Predicting structured data*, 1(0), 2006. 3

[22] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9640–9649, 2021. 3

[23] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. 3

[24] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020. 3

[25] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *International conference on machine learning*, pages 1691–1703. PMLR, 2020. 3

[26] Shentong Mo and Shengbang Tong. Connecting joint-embedding predictive architecture with contrastive self-supervised learning. *Advances in neural information processing systems*, 37:2348–2377, 2024. 3, 5, 6, 7, 1

[27] 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. 3, 4

[28] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In B. Kim, Y. Yue, S. Chaudhuri, K. Fragkiadaki, M. Khan, and Y. Sun, editors, *International Conference on Representation Learning*, volume 2024, pages 2632–2652, 2024. 4

[29] Mark Sandler, Andrey Zhmoginov, Max Vladymyrov, and Andrew Jackson. Fine-tuning image transformers using learnable memory. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12155–12164, 2022. 4

[30] Nobuyuki Otsu. A threshold selection method from gray-level histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-9(1):62–66, 1979. 4

[31] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2961–2969, 2017. 5, 1

[32] 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*, pages 418–434, 2018. 5

[33] Xiangteng He, Yuxin Peng, and Junjie Zhao. Which and how many regions to gaze: Focus discriminative regions for fine-grained visual categorization. *International Journal of Computer Vision*, 127(9):1235–1255, 2019. 6

[34] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2017. 1

# DSeq-JEPA: Discriminative Sequential Joint-Embedding Predictive Architecture

## Supplementary Material

## 6. Implementation Details

### 6.1. Pre-training

**Optimization Strategy.** For optimization, we follow the same setting as I-JEPA [10] and C-JEPA [26], using the AdamW [34] optimizer with cosine learning rate scheduling and warmup strategy, which starts from $1e$-4 and linearly warmup until $1e$-3 during first 15 epochs and ends at $1e$-6. Also, we tuned the weight decay value with cosine annealing, which starts from $4e$-2 and ends at $4e$-1. We initialized the target-encoder weights with the context-encoder weights, and updated the parameters of the target encoder by an exponential moving average (EMA) with decay factors increasing from 0.996 to 1.0.

**Training Protocol.** We pre-train DSeq-JEPA for 600 epochs. The input images were resized to $224 \times 224$ pixels, with a global batch size of 2048. Data augmentation techniques included random resized cropping with a random scale within $[0.3, 1.0]$. Following I-JEPA, we do not use hand-crafted data augmentations such as horizontal flipping, Gaussian blur, and color distortion. In Section 3.2, the saliency map is obtained from the 10-th layer, and in Section 3.3, $\alpha = 0.15$ is used to discard small connected regions. Masks for all discriminative regions were sampled at a scale between 0.15 and 0.2, ensuring a minimum of 10 patches per mask. We set the number of regions to 5, the same with I-JEPA. Overlapping between masks is not allowed. We set the window size to 3 in the max pooling operation.

**Computational Resources.** We used 8 H200 GPUs with 140 GB of memory. The use of 8× H200 was opportunistic, not required.

### 6.2. Downstream Evaluation

**Image Classification.** We use the ImageNet-1K [14] validation subset for evaluating learned representations in large-scale benchmark. As in [10, 26], we train a linear head on the full training set, where we use concatenated feature maps from the last four transformer blocks. The linear head was configured with batch normalization and optimized using SGD with a momentum of 0.9, Nesterov acceleration, and a weight decay of $5 \times 10^{-4}$. Training was conducted for 28 epochs with a base learning rate of 0.01, scaled automatically by batch size, and scheduled via a multi-step decay at epochs 8, 16, and 24 (learning rate values: $0.01 \rightarrow 0.001 \rightarrow 0.0001 \rightarrow 0.00001$). A batch size of

32 per replica was used, with random resized crop augmentation (crop size 224 from resized 256-pixel images).

**Fine-grained Visual Categorization.** For fine-grained visual categorization, we use iNaturalist21 [15], CUB-200-2011 [16], Stanford-Cars [17], and as ImageNet linear-probing, we train a linear classifier in same setting.

**Detection and Segmentation.** For object detection, we fine-tune a Mask R-CNN detector [31] with an FPN on MS-COCO 2017 [18], reporting $AP^{box}$ and $AP^{mask}$ on val2017 following the standard COCO protocol. The convolutional backbone is replaced by a ViT encoder pretrained with DSeq-JEPA , as in [1, 10]. We fine-tune the entire model end-to-end for 25 epochs with AdamW, a base learning rate of $1.0 \times 10^{-4}$, weight decay 0.1, using linear warmup in the first epoch followed by cosine decay, a global batch size of 16. For the ViT backbone we use stochastic depth with maximum drop-path rate 0.1, and otherwise adopt default Mask R-CNN hyperparameters from the underlying detection framework.

For semantic segmentation on ADE20K, we adapt a UPerNet-style decoder with a ViT-Base backbone pretrained using DSeq-JEPA. We extract feature maps from pre-trained model at 3, 5, 7, and 11 layers, and pass them into decoder. We train and evaluate with an input crop size of $512 \times 512$, and use a sliding-window inference strategy with a $512 \times 512$ crop and a stride of $341 \times 341$. The model is optimized with AdamW with a learning rate of $1 \times 10^{-4}$, a weight decay of 0.05.

**Low-level Tasks.** For compositional reasoning on Clevr [20], we evaluate DSeq-JEPA representations on the Clevr-Count and Clevr-Distance tasks from VTAB-1k. We use a ViT-based backbone pretrained with DSeq-JEPA and freeze the target encoder during downstream training. Each image is resized to an input resolution of $224 \times 224$ and fed into the backbone; we take the final-layer [CLS] token as a global image representation and pass it through a task-specific linear classification head. Clevr/Count and Clevr/Distance are treated as 8-way and 6-way classification problems, respectively. We follow the standard VTAB-1k data split, using 800 images for training, 200 images for validation, and 15,000 images for testing for each task. The linear heads are optimized with AdamW with a learning rate of $1 \times 10^{-3}$, a weight decay of 0.05, a batch size of 256, and we train for 100 epochs using a cosine learning rate schedule.

**Algorithm 1:** Discriminative Region Selection

---

**Input:** Saliency map $\mathbf{A}$, current epoch $t$, total pre-training epochs $T$

**Output:** Discriminative regions $\{R_k\}_{k=1}^{N}$

1  Compute $\lambda \leftarrow \min\big(1, \max(0, t/T)\big)$;
2  Normalize $\mathbf{A}$ to $\tilde{\mathbf{A}}$ via Equation (1);
3  Obtain binary mask $\mathbf{M}$, connected regions $\{R_k\}$;
4  Compute discriminative score $\rho_k$ for each region via Equation (3);
5  Select top-$N$ regions $\mathcal{R} = \{R_k\}_{k=1}^{N}$;
6  **for** $k = 1$ **to** $N$ **do**
7  $\quad$ Sample $b_k \sim \text{Bernoulli}(p)$;
8  $\quad$ **if** $b_k = 1$ **then**
9  $\quad\quad$ $R_k := R_k$;
10 $\quad$ **else**
   $\quad\quad$ /* Randomly sample region with center $(u,v)$ and size $(w,h)$ */
11 $\quad\quad$ $R_k := \text{R(u,v,w,h)}$;
12 **return** $\{R_k\}_{k=1}^{N}$

---

# 7. Implementation of Discriminative Region Selection

For completeness, Algorithm 1 summarizes the implementation of our discriminative region selection with a probabilistic curriculum. Given the saliency map $\mathbf{A}$ from the target encoder, we first normalize it, extract connected components, and rank the resulting regions $\{R_k\}$ by their discriminative scores (Equation (1)-Equation (3) in the main paper). We then select the top-$N$ regions as discriminative candidates. During pre-training epoch $t$, we compute a mixing probability $\lambda$ that linearly increases from 0 to 1, and for each region index $k$ we either use the discriminative region $R_k$ with probability $\lambda$ or a randomly sampled block region $R(u, v, w, h)$ with probability $1 - \lambda$. This per-sample Bernoulli mixing implements the curriculum described in Section 3.3 of the main paper.

Table 6. **Effect of different prediction orders (ViT-B/16).**

| Order scheme | ImageNet |
|---|---|
| I-JEPA (flat) | 72.0 |
| Random order | 71.7 |
| Spatial order | 72.7 |
| **Sequential order (ours)** | **73.5** |

# 8. Additional Experimental Analyses

## 8.1. Effect of Prediction Order

Table 6 revisits the ablation on prediction order using ViT-B/16 pre-trained on ImageNet. The results for I-JEPA (flat), random order, and our sequential order are identical to those reported in the main paper (Table 4); here we additionally evaluate a spatial order variant, where the same set of discriminative regions is sorted by the coordinates of their centers and predicted from left to right and top to bottom (row-major scan), *i.e.*, according to image geometry structure rather than discriminativeness. We observe that random order performs slightly worse than the flat I-JEPA baseline (71.7% vs. 72.0%), spatial order brings only a small improvement (72.7%), while our discriminative sequential order achieves the best accuracy (73.5%). This shows that autoregression alone is not sufficient: when the prediction order carries no useful information (random), it can even hurt performance; a purely geometric order provides some signal but limited benefit, whereas a semantically meaningful prediction trajectory over discriminative regions is crucial for DSeq-JEPA to fully improve the learned representations.

Table 7. **Sensitivity of DSeq-JEPA to the number of regions N for ViT-B/16 pre-trained on ImageNet.** We report linear probing accuracy on ImageNet.

| N | 3 | 5 | 7 |
|---|---|---|---|
| **ImageNet** | 72.9 | **73.5** | 73.4 |

## 8.2. Sensitivity to the Number of Regions $N$

We fix the number of regions to $N = 5$ in the main experiments, following the I-JEPA setting. To assess sensitivity to this choice, we vary $N \in \{3, 5, 7\}$ for DSeq-JEPA with ViT-B/16 and pre-train on ImageNet under the same schedule. Table 7 reports linear probing accuracy on ImageNet. We find that DSeq-JEPA is robust for the number of discriminative regions, and the overall performance fluctuates little. Using too few regions ($N = 3$) slightly under-utilizes contextual information, while increasing to $N = 7$ brings no further benefit and even reduces accuracy by 0.1% compared to $N = 5$, despite the additional complexity. Overall, it can be concluded that DSeq-JEPA is a stable self-supervised learning algorithm.

## 8.3. Per-step Prediction Difficulty

To better understand the effect of discriminative sequencing, we analyze the prediction loss at different steps of the region sequence. Using the DSeq-JEPA ViT-B/16 checkpoint at epoch 450, we sample 10,000 ImageNet images and, for each Top-$k$ region prediction, compute the average
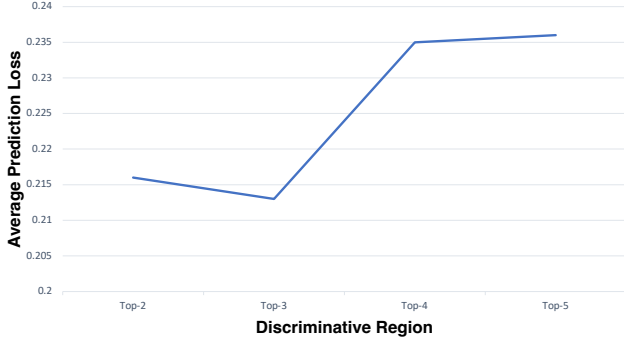
Figure 6. **Per-step prediction difficulty for DSeq-JEPA (ViT-B/16)**. We report the average prediction loss on 10,000 ImageNet images when predicting the Top-2, Top-3, Top-4, and Top-5 discriminative regions, using the checkpoint at epoch 450.

loss. As shown in Figure 6, early steps (Top-2 and Top-3) exhibit lower prediction error, while later steps (Top-4 and Top-5) are clearly harder, with noticeably higher loss. This pattern is consistent with our interpretation that ordering regions by discriminativeness induces an implicit easy-to-hard curriculum: the model first predicts highly informative, stable regions and then gradually moves to weaker or more context-dependent regions. In turn, the sequential predictor leverages this trajectory to structure the learned representations.

Table 8. **Pre-training time comparison on ImageNet-1K with ViT-B/16.**

| Method | I-JEPA | C-JEPA | DSeq-JEPA |
|--------|--------|--------|-----------|
| Time(h) | 24.2 | 24.5 | 26.5 |

## 8.4. Pre-training Time Overhead

Table 8 reports the pre-training time on ImageNet with ViT-B/16 under the same hardware and training configuration. Compared to I-JEPA (24.2h) and C-JEPA (24.5h), DSeq-JEPA requires 26.5h, *i.e.*, roughly a 9% increase in wall-clock time. This overhead is modest and mainly comes from computing the similarity-derived saliency map and selecting a small number of discriminative regions ($N = 5$). Given the consistent gains across ImageNet, FGVC, dense prediction, and low-level reasoning benchmarks, we view this additional cost as practically acceptable.

## 8.5. Additional Visualizations

In Figure 5, we visualize the evolution of patch-level clustering on a few *training* images with 4 clusters and observe that DSeq-JEPA progressively organizes patches into semantically meaningful groups. To verify that this behavior generalizes beyond the training set and to explore
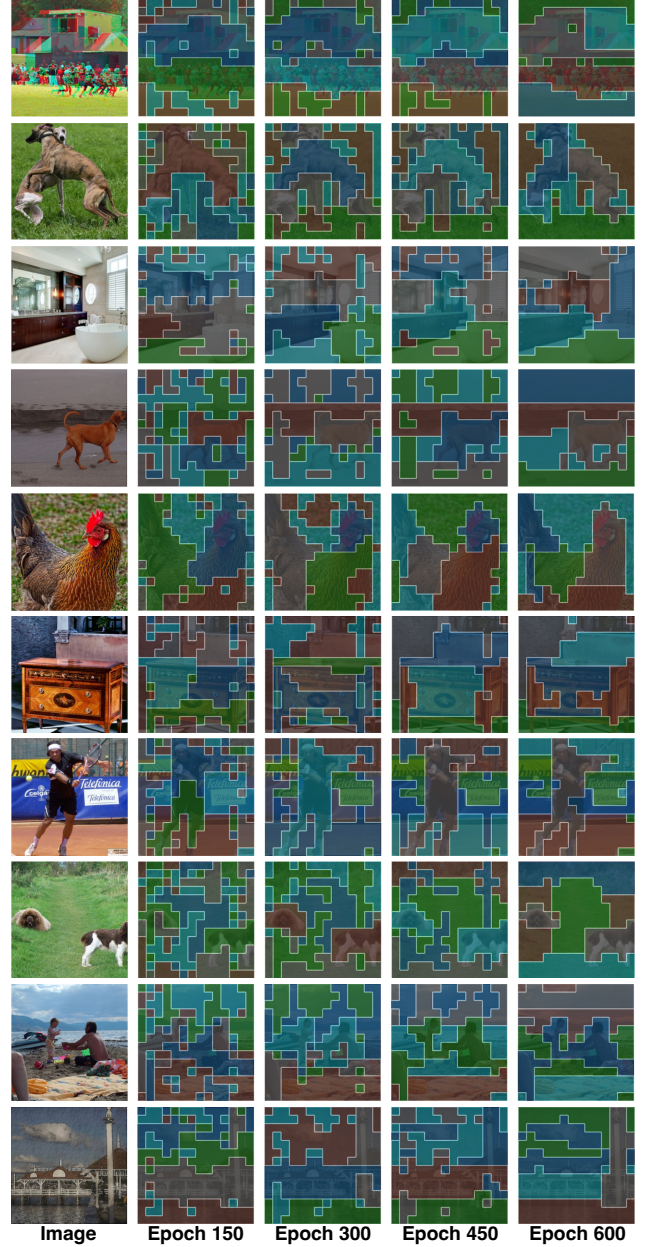


Figure 7. **Evolution of patch-level clustering during pre-training.** From left to right: input image and 5-cluster results from ViT-B/16 checkpoints at 150, 300, 450, and 600 epochs.

a slightly finer partition, Figure 7 shows analogous visualizations on held-out *test* images with 5 clusters using the ViT-B/16 checkpoints at 150, 300, 450, and 600 epochs. Consistent with the training examples, DSeq-JEPA encourages patches with similar semantics to be grouped together into compact, coherent clusters. These results further support our claim that discriminative sequencing leads to more structured and interpretable patch-level representations.

3