# Learning and Leveraging World Models in Visual Representation Learning

**Quentin Garrido**[1,2], **Mahmoud Assran**[1], **Nicolas Ballas**[1], **Adrien Bardes**[1,3], **Laurent Najman**[2], **Yann LeCun**[1,4,5]

[1]FAIR at Meta, [2]Univ Gustave Eiffel, CNRS, LIGM, F-77454 Marne-la-Vallée, France, [3]INRIA, [4]Courant Institute, New York University, [5]Center for Data Science, New York University

Joint-Embedding Predictive Architecture (JEPA) has emerged as a promising self-supervised approach that learns by leveraging a *world model*. While previously limited to predicting missing parts of an input, we explore how to generalize the JEPA prediction task to a broader set of corruptions. We introduce Image World Models, an approach that goes beyond masked image modeling and learns to predict the effect of global photometric transformations in latent space. We study the recipe of learning performant IWMs and show that it relies on three key aspects: conditioning, prediction difficulty, and capacity. Additionally, we show that the predictive world model learned by IWM can be adapted through finetuning to solve diverse tasks; a fine-tuned IWM world model matches or surpasses the performance of previous self-supervised methods. Finally, we show that learning with an IWM allows one to control the abstraction level of the learned representations, learning invariant representations such as contrastive methods, or equivariant representations such as masked image modelling.

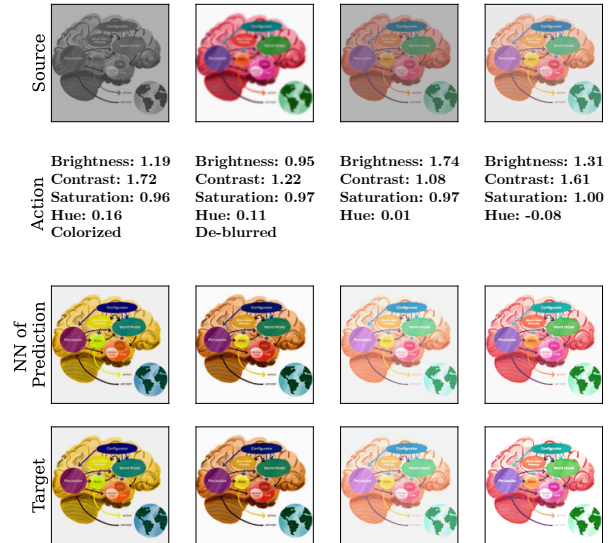**Correspondence:** Quentin Garrido at garridoq@meta.com

## 1 Introduction

Learning and leveraging world models is common practice in reinforcement learning (RL), with demonstrable success in the last few years in particular Ha and Schmidhuber (2018); Hafner et al. (2019, 2023). World models are commonly learned by training a network to predict the consequence of an action, either in input space (Yang et al., 2023), or in latent space (Hu et al., 2023; Hafner et al., 2023). Given such a broad view of world modelling, we seek to explore whether learning and leveraging world models can also be benificial in visual representation learning.
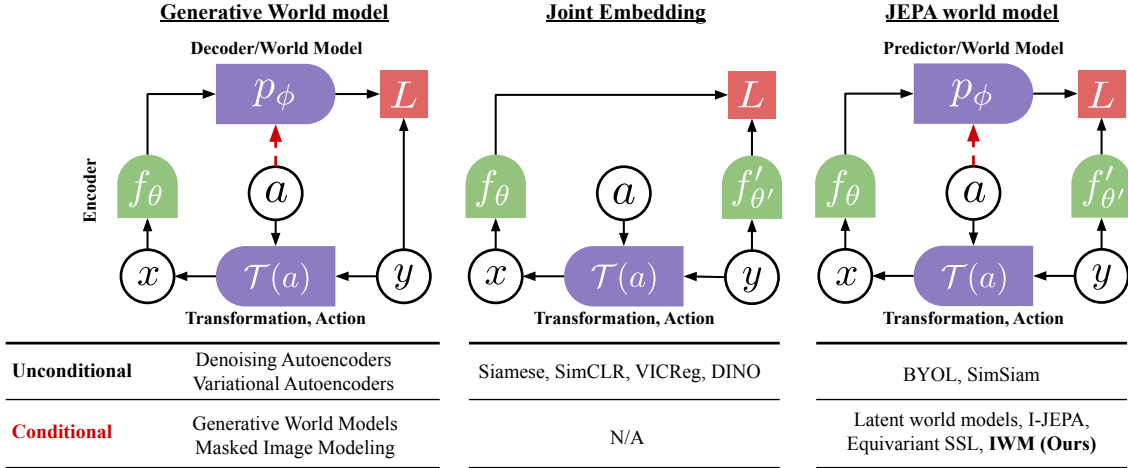
A wide family of self-supervised learning approaches are based on encoder-predictor architectures, wherein the encoder-predictor networks are trained to predict transformations of the data; e.g., masked image modelling (Bao et al., 2021; He et al., 2021), joint-embedding architectures (Grill et al., 2020; Xie et al., 2022; Assran et al., 2023; Baevski et al., 2022), or equivariant prediction objectives (Gupta et al., 2023; Garrido et al., 2023b). If we regard transformations of the data as "actions," then we can easily relate self-supervised learning approaches to world-modelling in reinforcement learning; see figure 2.

For instance, the decoder network in masked au-



**Figure 1 Visualisation of predictions in latent space with a learned Image World Model.** We apply an action on a source image in latent space and retrieve the nearest neighbour of the predicted representation in a bank of 256 images. We see that IWM is capable of modeling transformations and undo corruptions, showing an understanding of the underlying image transformations. Image from: ai.meta.com/blog/yann-lecun-advances-in-ai-research/

toencoders (He et al., 2021) can be thought of as a generative image world model, which learns to infer the effect of the "masking action" $\mathcal{T}(a)$ on an image

**Figure 2** Multiple families of methods with related architectures can be distinguished, in which the conditioning or not of their world model is a key distinction. **Generative World Models** are trained to invert a transformation in input space, leveraging an autoencoder framework. Methods for world modeling and representation learning can be instantiated in this way. **Joint Embedding** methods get rid of the world model but operate in latent space by encoding what is common between transformed inputs. It is the main class of SSL methods. **JEPA World Models** can be seen as a more general framework where a world model is trained in latent space. This family has been very successful both in reinforcement learning and in representation learning, and is where Image World Models (IWM) falls.

$y$; in this case, the transformation parameters $a$ (locations of masked image patches), are also fed to the decoder network. Methods based on joint-embedding predictive architectures (JEPAs), such as I-JEPA (Assran et al., 2023) or data2vec (Baevski et al., 2022), operate similarly, but can be seen as learning a latent image world model, which learns to infer the effect of the masking action on the representation of an image. If one does not condition the predictor on the transformation parameters, then the best we can hope for is learning representations that are invariant to the data transformations, as in BYOL (Grill et al., 2020) and SimSiam (Chen and He, 2020), wherein the image transformations correspond to various photometric and geometric data augmentations.

However, despite some of the apparent similarities between world modelling in reinforcement learning and self-supervised learning from images, the learned world model in reinforcement learning is typically leveraged in downstream tasks, e.g., for planning (Hansen et al., 2022). In contrast, the learned world model in self-supervised learning is typically discarded after pretraining, as the main focus is often on the representation quality of the learned encoder network. This stems from the fact that most downstream tasks in computer vision are unrelated to the world modeling task. Common tasks of interest focus on discriminative aspects and as such, even when the predictor learns useful information, it is simply discarded. We postulate that discarding the world model in representation learning is wasteful, and that just like in RL, we can reuse this world model for downstream tasks. This motivates us to study, in

more depth, learning world models as a paradigm for representation learning. We thus introduce Image World Models (IWM, illustrated to the right of figure 2) as a way to learn both good representations and strong reusable world models. IWM is based on JEPA and extends the usual latent inpainting to also include photometric transformations, allowing us to demonstrate the key aspects in learning a capable world model, which include the choice of predictor conditioning, the strength of the transformations, and the capacity of the world model.

We then focus on leveraging the learned world model for downstream tasks, and find that it can be leveraged through finetuning. Specifically, we find that finetuning the world model on top of the frozen encoder for downstream tasks provides improved performance over encoder finetuning; this is also achieved at a fraction of the cost and number of finetuned parameters. Moreover, only the world model learned by IWM exhibits this behavior; finetuning a randomly initialized network of the same architecture as the predictor does not provide such a performance improvement. This suggests that the world model should be a key part of the inference process, instead of being discarded. Inspired by instruction tuning (Wei et al., 2022; Zhang et al., 2023), we further show that the world model can be finetuned to solve multiple tasks at once, further improving efficiency.

Our study reveals another key aspect of representation learning with world models: the capacity given to the world model has a direct influence on the

level of abstraction of the learned representations. Intuitively, if the predictor is the identity (i.e., no predictor, middle of figure 2), the network will capture high level semantic information, as it will only learn to encode what is shared between the input $y$ and its transformation $x$. This is the driving force behind the representation quality of contrastive learning, where transformations are selected to only preserve the semantics of the image. On the other hand, as the predictor has more capacity and can effectively invert the effect of the transformations, the output of the encoder can retain more information about its input. These two ideas are at the core of equivariant representation learning; a predictor that can apply transformations effectively is equivariant, whereas a predictor that cannot is invariant. We find that a world model that is invariant to transformations performs better in linear evaluation, whereas one that is equivariant correlates with better world model finetuning. This gives a tradeoff between ease of adaption and raw performance. As such, learning representations by learning a world model gives us flexibility in the properties of the representations, making this an attractive representation learning framework.

Our contributions can be summarized as follows:

- We show how to leverage JEPAs to learn an Image World Model (IWM). The key aspects are: complexity of transformations, conditioning on transformations, and capacity of the predictor.
- We show that equivariant world models can be leveraged for discriminative tasks. Finetuning the predictor leads to better performance compared to encoder finetuning, at a fraction of the cost. Inspired by instruction tuning, we also demonstrate that it can be finetuned on several tasks at once.
- We show that controlling the capabilities of the world model gives us representations with different properties. An invariant world model gives us more abstract representations and performs better in linear evaluation, akin to contrastive learning. An equivariant world model preserves more information about the input, giving better peak performance with predictor finetuning.

## 2 Related works

### 2.1 Augmentation invariant Self-Supervised Learning

At the core of contrastive methods lies augmentation invariance. Multiple augmented views of an image should lead to the same representation in latent space. The core of these methods is thus in how to avoid

these representations collapsing. Sample-contrastive methods (Chen et al., 2020a; He et al., 2020; Chen et al., 2020b; Caron et al., 2021; Chen et al., 2021; Yeh et al., 2021; HaoChen et al., 2021; Oquab et al., 2023) avoid this phenomenon by pushing away representations coming from other data points. Dimension-contrastive methods (Bardes et al., 2021; Zbontar et al., 2021; Ermolov et al., 2021; Li et al., 2022; Bardes et al., 2022) avoid collapse by considering the representations as a whole and encouraging maximization of information content. Both dimension- and sample-contrastive methods have been shown to lead to very similar representations (Garrido et al., 2023a). Prediction based methods (Grill et al., 2020; Chen and He, 2020) learn by predicting the augmented representations, but they also lead to invariant representations due to a lack of conditioning on the transformations.

### 2.2 World modeling in visual representation learning

While world modeling is a successful paradigm in reinforcement learning Hafner et al. (2019, 2023) or video prediction Yang et al. (2023); Hu et al. (2023), it has yet to show clear benefits in representation learning. However, multiple families of approaches can be reframed in light of this. Equivariant self-supervised learning methods (Devillers and Lefort, 2022; Park et al., 2022; Garrido et al., 2023b; Gupta et al., 2023; Dangovski et al., 2021) aim to predict transformations of data when such transformations form a group. Masked Image Modeling He et al. (2021); Bao et al. (2021); El-Nouby et al. (2024); Xie et al. (2022) learns representations by predicting masked parts of the image. While these approaches predict in pixel space, their decoders can be seen as instantiations of world models. Similarly, JEPAs (Assran et al., 2023; Baevski et al., 2022) predict masked parts of the image, but in the latent space. Recently, generative approaches have been applied to representation learning Hudson et al. (2023); Clark and Jaini (2023); Chen et al. (2024), and while these approaches seem promising, their performance still remains below contrastive or MIM approaches. Recent work has also shown negative correlations between generation quality and representation quality (Chen et al., 2024). One shared aspect among these works is that the world model (predictor or decoder) is either discarded for evaluations, or only used to augment data (Hudson et al., 2023). We propose to go beyond these practices and show that we can learn a world model that is reusable for downstream tasks while still learning high-quality representations.

## 3 Method

We now describe Image World Models (IWM). It follows a Joint Embedding Predictive Architecture framework (LeCun, 2022) akin to I-JEPA (Assran et al., 2023). In this framework, the predictor is the instantiation of the world model. We consider a world model to be capable if it can apply transformations in latent space, and thus learns equivariant representations. As such, we call a capable world model equivariant [1] and a poor world model invariant.

An appealing aspect of using JEPAs is that approaches which learn equivariant representations using contrastive methods often have to rely on an invariance loss to increase representation quality, whether explicitly (Gupta et al., 2023; Garrido et al., 2023b), or implicitly (Chavhan et al., 2023a). On the other hand, a JEPA style approach does not have this drawback, as the semantic aspect of the representation is learned through latent inpainting. Working in latent space further allows the network to remove unnecessary information, or that which is too hard to predict. This makes the JEPA formulation attractive since, for reconstructive methods, the quality of the reconstruction is not necessarily correlated with representation quality Chen et al. (2024).

To train IWM, the first step is to generate source and target views — $x$ and $y$ respectively in figure 2 — from an image $I$.

**Target $y$.** The target view is generated by applying a random horizontal flip, a crop, and color jitter (brightness, contrast, saturation, hue) to the original image $I$. No destructive augmentations such as grayscale are applied on the target to ensure that the target has as much information as possible. We further elaborate on this choice in appendix C.

**Source $x$.** For the source view, we start from the target $y$ which we further transform. We first apply another color jitter, as well as destructive augmentations: grayscale, blur and solarization. This set of augmentations is the same as the one used in contrastive SSL. Finally, we also mask parts of the image following I-JEPA. We define our mask $M_x$ (a set of indices) as the union of 4 rectangular masks. Confer appendix A for exact implementation details.

**Action $a$.** We denote by $a_{x \to y}$ the transformation parameters associated with the transformation of $x$ to $y$, i.e., the invert of the initial transformation process.

---

$a_{x \to y}$ contains information about the color jitter difference between $x$ and $y$ as well as information on whether or not each destructive augmentation was applied.

**World modeling with $p_\phi$.** The source and target are then fed respectively through an encoder $f_\theta$ and its exponential moving average $f_\theta^{\text{EMA}}$. This gives us representations $z_x = f_\theta(x)$ and $z_y = f_\theta^{\text{EMA}}(y)$. The use of the EMA network is crucial to avoid collapsed solutions. To condition the predictor, acting as our world model, it is fed with geometric information about the target in the form of mask tokens as well as $a_{x \to y}$. We denote these mask tokens as $m_a$, which correspond to the positions in $M_x^C$. The predictor $p_\phi$ then takes as input the embedded source patches $x_c$, transformation parameters $a_{x \to y}$ and mask tokens $m_a$. Its objective is then to match $p_\phi(z_x, a_{x \to y}, m_a) = \hat{z_y}$ to $z_y$.

**Loss.** The loss function used is a squared $L2$ distance between the predictions $\hat{z_y}$ and their targets $z_y$:

$$L(x, y) = \sum_{i \in M_x^C} \| p_\phi (f_\theta(x), a_{x \to y}, m_a)_i - f_\theta^{\text{EMA}}(y)_i \|_2^2.$$

### 3.1 Architecture and nomenclature

Our encoder is a Vision Transformer (Dosovitskiy et al., 2021), in particular we use the ViT-B/16 architecture. Our predictor is based on the same architecture with different depth and embedding dimension. We denote instances of IWM as $\text{IWM}_{X,Y}^Z$ where $X$ is the depth of the predictor, $Y$ its embedding dimension, and $Z$ is either Inv or Equi depending on the capabilities of the world model. For example $\text{IWM}_{18,384}^{\text{Equi}}$ means that the predictor is 18 layers deep, with 384 dimensional embeddings and exhibits equivariant behavior, i.e., has learned a versatile world model.

## 4 Learning an Image World Model for representation learning

### 4.1 Evaluating the quality of the world model

As discussed previously, learning equivariant representations and learning a world model are closely related problems. As such, we can borrow metrics from the equivariance literature to evaluate the quality of a trained world model. We rely on Mean Reciprocal Rank (MRR) (Kipf et al., 2019) as our main metric. To compute it, we generate a bank of augmented target images (256 in practice). We feed the representation of the clean image through the predictor

**Table 1  Influence of predictor conditioning on the quality of the world model.** Both Sequence and Feature conditioning lead to good world models .Gray is our default setting.

| Conditioning: | None | Sequence | Feature |
|---|---|---|---|
| MRR | 0.00 | 0.82 | 0.79 |

**Table 2  Impact of predictor architecture and transformations on MRR.** Learning an effective world model requires complex transformations and adequate predictor capacity. Gray is our default setting. Red and Green respectively indicate invariant and equivariant behavior.

| Predictor: | I-JEPA | IWM | |
|---|---|---|---|
| (depth, dim.): | (12,384) | (12,384) | (18,384) |
| Jitter | 0.00 | 0.11 | 0.25 |
| + Destructive | 0.00 | 0.09 | 0.79 |
| + Strong Jitter | 0.00 | 0.81 | 0.85 |

with the goal of predicting the target image. We then compute the distance between the prediction and the augmented representation bank from which we get the rank of the target in this NN-graph. Averaging the reciprocal ranks over multiple images and transformations gives us MRR which tells us about the quality of the world model. A MRR close to 1 means that the world model is able to apply the transformation, on the contrary a MRR close to 0 means that it cannot.

## 4.2    Learning a strong Image World Model

In order to build a performant IWM, we isolate three key aspects: conditioning the predictor on transformations (or actions), controlling the complexity of the transformations, and controlling the capacity of the predictor. We show that not caring properly for either of those leads to invariant representations.

**World model conditioning.** We study two approaches to condition the predictor on the transformation information.
*Sequence conditioning.* One approach is simply to add tokens representing the transformation to the input of the predictor. Although this seems straightforward, it needs to be implemented in a way that breaks the permutation equivariance of the transformer predictor. To do so, every token is fed through a unique linear layer that allows the network to transform the information in a way that can be disambiguated by the predictor.
*Feature conditioning.* Another option is to mix the information between the transformation and mask tokens by adding the conditioning as extra dimensions, then feeding the mask tokens through a 1x1 convolutional neural network to mix the information in the mask tokens and map back to the right dimensionality.
As we can see in Table 1, no conditioning leads to a world model that cannot apply transformations whereas both conditioning using the sequence or feature axes leads to good world models. We use the feature conditioning in practice as it leads to higher downstream performance.

**Transformation complexity.** We rely on data augmentation as used in contrastive approaches, consisting of color jitter (brightness, hue, contrast, saturation),

grayscale, blur, and solarization. We refer to the last three as destructive since they remove information. Beyond the set of transformations modeled, their strength must also be adequate to learn a useful world model. If the prediction task is too easy, then the predictor will not learn anything useful. As presented in Table 2, the stronger the augmentations, the easier it is to learn a strong world model. We provide more detailed ablations on the augmentations in Appendix C, where we see the trend continuing on a wider range of augmentation scenarios.

**World model capacity.** If the transformation is complex, the predictor needs more capacity to be able to apply it, motivating capacity as a crucial factor in learning Image World Models. As we can see in Table 2, a deeper predictor enables us to learn a strong world model on a wider range of augmentations, and is key to the success of IWM. We study in more detail the influence of depth on achieving a good world model in appendix C. For 12 layers, jitter equivariance is achieved 1 out of 5 times whereas for the 18 layers, it is achieved 4 out of 5 times. As such, predictor capacity is a key component of a strong world model.

## 4.3    Visualizing predictions.

In the same way that we computed MRR, we can compare the predicted representations to a bank of transformed images and look at the image associated to the prediction's nearest neighbor. As we see in Figure 1 the world model learned by IWM is able to properly apply transformations in latent space. We can however see some inaccuracies when inverting grayscale as it is not properly invertible. These visualisations help reinforce the fact that IWM is able to learn strong world models for image transformations. Confer appendix I for more visualizations.

**Table 3 How to predict for predictor finetuning.** Using the teacher improves performance, and the exact prediction task is not crucial. Null latents are more flexible and perform better. For better efficiency, a full prediction is not needed but leads to a small drop in performance. Gray is our default setting.

| Setting | ImageNet Top-1 (%) | Gap |
|---|---|---|
| Default | 82.9 | - |
| + Teacher | 83.2 | + 0.3 |
| + Null latents | 83.3 | + 0.1 |
| + Pred only one token | 82.8 | -0.5 |



**Figure 3 Finetuning efficiency.** When taking into account the number of finetuned parameters, predictor finetuning is significantly more efficient than finetuning the encoder.
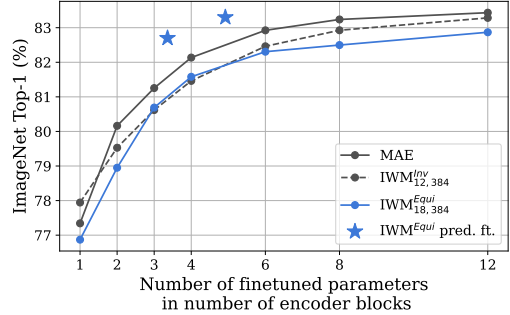
## 5 Leveraging world models for downstream tasks

A limitation of world models learned on images is that the task they solve is not aligned with most downstream tasks. We showed that IWM can apply color jitter or colorize images, but these are not the tasks that drive applications of computer vision. This is in contrast with LLMs where predicting the next token is one of the main applications of such models. We thus study how to leverage a world model in vision, for tasks that go beyond applying transformations. We focus on discriminative tasks such as image classification and image segmentation.

### 5.1 Predictor finetuning

For any task, the evaluation head needs to understand the learned latent space and leverage it to solve the problem at hand. This is something our learned predictor can do, suggesting that it has learned useful information that is not necessarily present in the encoder. However, since the predictor is trained to predict another valid representation, its output has no reason to lead to better downstream performance if used as is. This is why the predictor needs to be finetuned to solve discriminative tasks. We thus focus on comparisons with finetuning protocols, following He et al. (2021). All methods studied are pretrained and evaluated on ImageNet Deng et al. (2009) and use ViT-B/16 as encoders.

**Prediction task.** When finetuning the predictor, we still need to use it for a prediction task. In Table 3, we study various ways to define the prediction task and how it impacts performance. The first aspect we notice is that using the teacher network improves performance over the student. Using a random transformation or not is not an important factor, and the most important one is to predict another full image. This makes the evaluation more flexible as we do not have to reuse the pretraining objective for our evaluation. Using a CLS token to aggregate infor-

mation instead of a full image prediction is also a valid strategy, though it lowers the performance by half a point. This techniques has the advantage of being cheaper ($N + 1$ tokens vs $2N$) so it can be a good alternative depending on the use case. Overall, the simplest approach is the best: predicting an untransformed version of the full image. This makes the finetuning protocol easily reusable as it is not dependent on the pretraining task. We provide more detailed ablations in appendix D.

**General Results.** In Table 4, we compare predictor finetuning to encoder finetuning and end-to-end finetuning of both the predictor and encoder, using ViT-B/16 for the encoder. We see that IWM maintains or improves performance over I-JEPA and that an invariant behavior is better in encoder finetuning. Interestingly, predictor finetuning of the equivariant IWM is able to match the performance of finetuning of the invariant model's encoder. This shows that the protocol can be competitive as it trades parameters at inference time for a more computationally friendly adaptation. While this evaluation increases the number of parameters used at inference time, it still amortizes the forward pass through the backbone, something that full finetuning does not do. As such, as soon as multiple tasks are considered, using the finetuned predictor provides a higher throughput than regular finetuning.

When comparing the use of a randomly initialized predictor (i.e., a large evaluation head) versus a pretrained predictor, we see negligible gains for MAE. This suggests that the world model learned by MAE is not better than a randomly initialized network for classification. For I-JEPA and IWM with an invariant world model, we see gains in performance lower than 1 point, suggesting that the world model is not powerful enough to be leveraged. However, when looking at IWM with an equivariant world model, we see a gain of 1.8 points over a random predictor. This shows that the predictor has learned useful information and properties that bring additional benefit to

**Table 4  Finetuning evaluations on ImageNet-1k.** We evaluate prediction based methods by finetuning their encoder, by keeping the encoder frozen and finetuning their predictive world model or by finetuning both. Finetuning the world model is highly effective with IWM when it exhibits an equivariant behavior. This behavior is absent or less clear with other methods, showing the importance of a strong world model.

| Method | Epochs | No predictor | Frozen encoder, tuned predictor | | End to end |
|--------|--------|--------------|--------------|-----------|------------|
| | | Encoder | Random Init. | Pretrained | |
| MAE | 300 | 82.7 | 82.4 | 82.7 (+0.3) | 82.3 |
| | 1600 | **83.6** | **83.0** | 83.1 (+0.1) | 83.3 |
| I-JEPA | 300 | 83.0 | 79.1 | 80.0 (+0.9) | 82.0 |
| $IWM_{12,384}^{Inv}$ | 300 | 83.3 | 80.5 | 81.3 (+0.8) | 82.7 |
| $IWM_{18,384}^{Equi}$ | 300 | 82.9 | 81.5 | **83.3** (+1.8) | **84.4** |

**Table 5  Peak performance achieved from a single pretraining instance.** We compare ImageNet Top-1 accuracy with a frozen encoder or when allowing any evaluation head with any protocol, finetuning or not, with a predictor on top of the encoder or not.

| Method | Epochs | Frozen Encoder | Any protocol |
|--------|--------|----------------|--------------|
| DINO | 1600 | 82.0 | 82.8 |
| MOCOv3 | 300 | 76.4 | 83.2 |
| iBOT | 1600 | 83.0 | 84.0 |
| MAE | 1600 | 83.1 | 83.6 |
| I-JEPA | 300 | 80.0 | 82.0 |
| $IWM_{12,384}^{Inv}$ | 300 | 81.3 | 83.3 |
| $IWM_{18,384}^{Equi}$ | 300 | **83.3** | **84.4** |

**Table 6  Finetuning for segmentation on ADE20k.** Similar to image classification, we observe that predictor finetuning improves performance and outperforms encoder finetuning.

| Method | Encoder | Predictor | End to end |
|--------|---------|-----------|------------|
| I-JEPA | 44.2 | 45.4 | 45.1 |
| $IWM_{12,384}^{Inv}$ | **45.6** | 45.7 | 46.5 |
| $IWM_{18,384}^{Equi}$ | 44.2 | **46.8** | **47.0** |

what the encoder has learned.

The performance can be pushed further by finetuning end-to-end both the encoder and predictor, and IWM is able to outperform every other finetuning protocols. This allows us to get more performance out of a single pretraining since the world model is always trained. We hypothesize that the lack of performance for most approaches on end-to-end finetuning comes from the optimization complexity of finetuning a part of the network (encoder) while training from scratch another part (the predictor). We see in Table 5 that when aggregating the performance over all protocols, leveraging our IWM leads to the best performance with a frozen encoder, that is when allowed to leverage every part of the pretraining. Confer Appendix A for detailed performances.

**Image Segmentation.** We study in Table 6 the performance of I-JEPA and IWM on an image segmentation task on ADE20k. We observe similar trends as in image classification where the invariant model leads to the best encoder. However, finetuning the predictor with an equivariant model leads to significant gain over it, outperforming encoder finetuning by a large margin. Again, we observe gains in end-to-end finetuning. This further validates the potential of our IWM to be leveraged for a wide range of tasks. We

provide additional details in Appendix A.2.

**Efficiency.** In Figure 3, we study the efficiency of predictor finetuning compared to encoder finetuning. We see that when the number of parameters is comparable, and at multiple predictor sizes, predictor finetuning with IWM outperforms encoder finetuning by around 1 point compared to MAE, and by 1.5 points over IWM. This means that predictor finetuning is not only is a competitive protocol performance wise, but also with respect to efficiency of adaptation. We further study the behavior of IWM with a ViT-L/16 in section E. When comparing the end-to-end finetuning of a ViT-B with encoder finetuning of a ViT-L, we observe a gain in performance (84.4% vs 84.3%) with a fraction of the parameters (121M vs 307 M). This further shows how efficient leveraging the world model learned by IWM is, and that reusing all parts of your pretraining can prove as effective as scaling the encoder.

## 5.2  Multitask predictor tuning

We previously discussed efficiency gains when compared to encoder finetuning, but can improve efficiency even further. One of the main goal of representation learning is to obtain representations that can be used for a variety of tasks. And just like the predictor was trained to solve a variety of task (colorization, inpainting, changing color) we show that it can be finetuned on multiple tasks, inspired by prefix tuning (Li and Liang, 2021) and instruction

**Table 7 Multi-task finetuning.** Finetuning the predictor on multiple tasks at once performs similarly as finetuning it on each task separately. This enables the use of a single prediction head for multiple task, amortizing its cost.

| Dataset | Single-task | Multi-task | Difference |
|---------|-------------|------------|------------|
| ImageNet | 80.8 | 79.6 | -1.2 |
| iNat18 | 72.4 | 72.0 | -0.4 |
| SUN397 | 75.6 | 78.2 | +2.6 |
| Places205 | 64.8 | 64.1 | -0.7 |
| Average | 73.4 | 73.5 | +0.1 |

**Table 8 Linear and attentive probing performance on ImageNet-1k.** $\text{IWM}^{\text{Inv}}$ performs similarly to contrastive methods and $\text{IWM}^{\text{Equi}}$ to mask modeling ones.

| Method | Effective Epochs | Linear | Attentive |
|--------|------------------|--------|-----------|
| MoCoV3 | 300 | **76.3** | 76.4 |
| MAE | 300 | 60.2 | 73.5 |
| MAE | 1600 | 68.0 | 76.0 |
| I-JEPA | 300 | 70.0 | 75.0 |
| $\text{IWM}^{\text{Inv}}_{12,384}$ | 300 | 74.5 | **77.0** |
| $\text{IWM}^{\text{Equi}}_{18,384}$ | 300 | 67.5 | 75.1 |

tuning Wei et al. (2022); Zhang et al. (2023) in LLMs. The general idea, that we illustrate graphically in supplementary Figure S2, is to give new learned tokens to the predictor to indicate which task it is trying to solve. This is reminiscent of DyTox Douillard et al. (2022) which uses task tokens for continual learning. For each task, we thus have a task token, as well as a task specific head and/or loss function. All of the task losses are then combined, and the predictor, as well as task specific heads, are updated. We study a simple scenario where the batch is evenly split between tasks, noting that other sampling strategies may lead to further improved performance.

We evaluate in Table 7 $\text{IWM}^{\text{Equi}}_{18,384}$ (pretrained on ImageNet) on ImageNet, iNaturalist18 (Horn et al., 2018), SUN397 (Xiao et al., 2010), and Places205 (Zhou et al., 2014). For each task we train a single-task baseline where the total number of iterations is identical to the multi-task training. As such, training all four single-task baselines has exactly the same cost as the multi-task, although it leads to four different models instead of one. The multi-task predictor is able to achieve similar performance as the single-task predictors, with a moderate drop on most tasks but a significant increase in performance on SUN397. On average it achieves the same performance as the single-task predictors. This further demonstrates the efficiency gains of leveraging good world models, where the parameters are now shared across all tasks, making predictor finetuning lightweight at inference time for every task.

Overall, when a good world model is learned, it can be reused for downstream tasks by finetuning it. This leads to performance rivaling with encoder-finetuning at a fraction of the cost. It can be made even more efficient by doing a multi-task finetuning, highlighting the versatility of this approach.

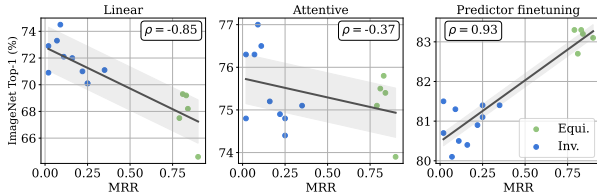## 6 Image World Models enable flexible representations

To complete our analysis of IWM for representation learning, we study how it performs on lightweight evaluation protocols that are commonly used in self-supervised learning. We focus on linear Chen et al. (2021) and attentive probing Chen et al. (2023).

As we see in Table 8, when IWM learns an invariant world model, it achieves a behavior akin to contrastive approaches such as MoCov3 with significant performance gains in linear evaluation compared to MIM or other JEPA based approaches. Similarly, when IWM learns an equivariant world model, its behavior is akin to MIM methods such as MAE with lower performance in linear evaluation but more competitive performance in attentive probing.
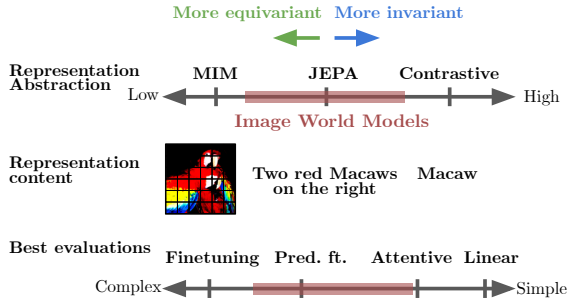
This suggests that a big difference between methods is not necessarily in the quality of the representation but in their abstraction level, i.e., how easy it is to extract information from them. Linear probing being one of the simplest evaluations, attentive being slightly more elaborate and finetuning being a more complex protocol.

In Figure 4, we see clear links between most suited evaluation protocols and equivariance of the world model. More invariant world models excel in linear evaluation and equivariant world models shine with larger evaluation heads such as in predictor finetuning. We also note that the richer representations stemming from equivariant world models lead to better performance on OOD datasets(see appendix F). This allows us to place families of approaches on a spectrum of representation abstraction in Figure 5. Contrastive methods occupy the high abstraction end of the spectrum, with information that is easily extractible with a simple protocol. However they suffer from lower peak performance when ignoring the adaptation cost, as seen in Table 5. On the opposite end lies Masked Image Modeling, which offers stronger performance with complex evaluations such as fine-

**Figure 4** While the level of equivariance influences performance in Linear and Predictor finetuning setting, it is hardly correlated to Attentive probing. This suggests that there is a trade-off in terms of the level of abstraction of the representation, and that different evaluation protocols evaluate different properties.



**Figure 5  Image World Models allow representation modularity.** Different families of methods offer representations with different properties, but IWM allows exploring the whole spectrum.

tuning but suffers in linear probing as information is not as easily accessible. By varying the equivariance of the world model, IWM is able to occupy the spectrum in between contrastive approaches and MIM, as we can see in Figure 4 and Table 8 with $\text{IWM}_{12,384}^{\text{Inv}}$ and $\text{IWM}_{18,384}^{\text{Equi}}$ being the two extremes of the IWM spectrum.

This spectrum can be summarized by the SSL ethos of "Learning what is predictible". Learning with a weak world model means that it cannot model the world properly and the encoder removes the information that cannot be predicted. On the other hand, if the world model is very powerful, the representation does not need to be as abstract or semantic as it can find a way to predict representations in any situation. This means that learning a world model offers a measurable way to control the level of abstraction of the representations.

## 7    Conclusion and future perspectives

We introduced IWM, an approach to learn self-supervised visual representations with world models. With an in-depth study, we provided guidelines and key components for learning a good image world model. Conditioning the world model with the image transformation is crucial to avoid collapsing to classical SSL behavior. Using strong transformations is also key to ensure that the world model learns to

model more complex behavior and be useful. Finally, enough capacity is needed for modeling complex behaviors. We showed that only a capable world model can be reused for discriminative task. This led to our predictor finetuning protocol that matches encoder finetuning at a fraction of the cost, showing that world models are versatile evaluation heads. We further adapted it to solve multiple tasks at once without losing performance. Finally, we studied how learning a world model impacts representation quality. A capable world model learns rich representations that improve performance on downstream tasks such as image classification and semantic segmentation. Additionally, learning an invariant world model led to better representations for linear evaluation. While MIM and Contrastive approaches are two ends of a spectrum in terms of representation abstraction, Image World Models allow us to interpolate between them. As such, we believe that learning image world models is a very promising framework for visual representation learning.

## 8    Broader impact statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none of which we feel must be specifically highlighted here.

## References

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.

Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. Data2vec: A general framework for self-supervised learning in speech, vision and language. In *International Conference on Machine Learning*, pages 1298–1312. PMLR, 2022.

Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.

Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.

Adrien Bardes, Jean Ponce, and Yann LeCun. VICRegl: Self-supervised learning of local visual features. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and

Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. https://openreview.net/forum?id=ePZsWeGJXyp.

Mathilde Caron, Hugo Touvron, Ishan Misra, Herve Jegou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021.

Ruchika Chavhan, Henry Gouk, Da Li, and Timothy Hospedales. Quality diversity for visual pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5384–5394, October 2023a.

Ruchika Chavhan, Jan Stuehmer, Calum Heggan, Mehrdad Yaghoobi, and Timothy Hospedales. Amortised invariance learning for contrastive self-supervision. In *The Eleventh International Conference on Learning Representations*, 2023b. https://openreview.net/forum?id=nXOhmfFu5n.

Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607. PMLR, 2020a.

Xiaokang Chen, Mingyu Ding, Xiaodi Wang, Ying Xin, Shentong Mo, Yunhao Wang, Shumin Han, Ping Luo, Gang Zeng, and Jingdong Wang. Context autoencoder for self-supervised representation learning. *International Journal of Computer Vision*, pages 1–16, 2023.

Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, 2020.

Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.

Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *ICCV*, 2021.

Xinlei Chen, Zhuang Liu, Saining Xie, and Kaiming He. Deconstructing denoising diffusion models for self-supervised learning, 2024.

Kevin Clark and Priyank Jaini. Text-to-image diffusion models are zero-shot classifiers. *arXiv preprint arXiv:2303.15233*, 2023.

MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. https://github.com/open-mmlab/mmsegmentation, 2020.

Ekin Dogus Cubuk, Barret Zoph, Jon Shlens, and Quoc Le. Randaugment: Practical automated data augmentation with a reduced search space. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18613–18624. Curran Associates, Inc., 2020. https://proceedings.neurips.cc/paper_files/paper/2020/file/d85b63ef0ccb114d0a3bb7b7d808028f-Paper.pdf.

Rumen Dangovski, Li Jing, Charlotte Loh, Seungwook Han, Akash Srivastava, Brian Cheung, Pulkit Agrawal, and Marin Soljačić. Equivariant contrastive learning. *arXiv preprint arXiv:2111.00899*, 2021.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.

Alexandre Devillers and Mathieu Lefort. Equimod: An equivariance module to improve self-supervised learning. *arXiv preprint arXiv:2211.01244*, 2022.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.

Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.

Alaaeldin El-Nouby, Michal Klein, Shuangfei Zhai, Miguel Angel Bautista, Alexander Toshev, Vaishaal Shankar, Joshua M Susskind, and Armand Joulin. Scalable pre-training of large autoregressive image models. *arXiv preprint arXiv:2401.08541*, 2024.

Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-supervised representation learning, 2021.

Quentin Garrido, Yubei Chen, Adrien Bardes, Laurent Najman, and Yann LeCun. On the duality between contrastive and non-contrastive self-supervised learning. In *The Eleventh International Conference on Learning Representations*, 2023a. https://openreview.net/forum?id=kDEL91Dufpa.

Quentin Garrido, Laurent Najman, and Yann Lecun. Self-supervised learning of split invariant equivariant representations. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 10975–10996. PMLR, 23–29 Jul 2023b. https://proceedings.mlr.press/v202/garrido23b.html.

Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning. In *NeurIPS*, 2020.

Sharut Gupta, Joshua Robinson, Derek Lim, Soledad Villar, and Stefanie Jegelka. Structuring representa-

tion geometry with rotationally equivariant contrastive learning, 2023.

David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems 31*, pages 2451–2463. 2018.

Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.

Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models. *arXiv preprint arXiv:2301.04104*, 2023.

Nicklas Hansen, Yixin Lin, Hao Su, Xiaolong Wang, Vikash Kumar, and Aravind Rajeswaran. Modem: Accelerating visual model-based reinforcement learning with demonstrations. *arXiv preprint arXiv:2212.05698*, 2022.

Jeff Z HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable guarantees for self-supervised deep learning with spectral contrastive loss. *NeurIPS*, 34, 2021.

Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.

Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. *arXiv preprint arXiv:2111.06377*, 2021.

Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2018.

Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving, 2023.

Drew A. Hudson, Daniel Zoran, Mateusz Malinowski, Andrew K. Lampinen, Andrew Jaegle, James L. McClelland, Loic Matthey, Felix Hill, and Alexander Lerchner. Soda: Bottleneck diffusion models for representation learning, 2023.

Thomas Kipf, Elise Van der Pol, and Max Welling. Contrastive learning of structured world models. *arXiv preprint arXiv:1911.12247*, 2019.

Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62(1), 2022.

Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation, 2021.

Zengyi Li, Yubei Chen, Yann LeCun, and Friedrich T Sommer. Neural manifold clustering and embedding. *arXiv preprint arXiv:2201.10000*, 2022.

Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. https://openreview.net/forum?id=Bkg6RiCqY7.

Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.

Jung Yeon Park, Ondrej Biza, Linfeng Zhao, Jan Willem van de Meent, and Robin Walters. Learning Symmetric Embeddings for Equivariant World Models, June 2022. http://arxiv.org/abs/2204.11371. arXiv:2204.11371 [cs].

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022.

Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010.

Tete Xiao, Liu Yingcheng, Bolei Zhou, Jiang Yuning, and Sun Jian. Unified perceptual parsing for scene understanding. In *ECCV*, 2018.

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, 2022.

Mengjiao Yang, Yilun Du, Kamyar Ghasemipour, Jonathan Tompson, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. *arXiv preprint arXiv:2310.06114*, 2023.

Chun-Hsiao Yeh, Cheng-Yao Hong, Yen-Chi Hsu, Tyng-Luh Liu, Yubei Chen, and Yann LeCun. Decoupled contrastive learning. *arXiv preprint arXiv:2110.06848*, 2021.

Yang You, Igor Gitman, and Boris Ginsburg. Large batch training of convolutional networks. *arXiv preprint arXiv:1708.03888*, 2017.

Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.

Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *ICML*, pages 12310–12320. PMLR, 2021.

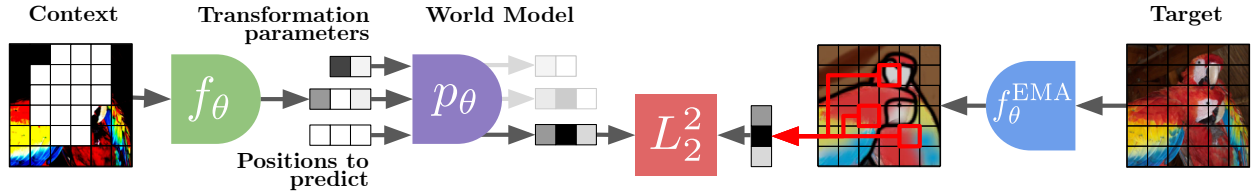Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk

minimization. In *International Conference on Learning Representations*, 2018. https://openreview.net/forum?id=r1Ddp1-Rb.

Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. Instruction tuning for large language models: A survey, 2023.

Bolei Zhou, Agata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. In *NeurIPS*, 2014.

Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *IJCV*, 2019.

**Figure S1  IWM (Image World Model).** Starting from an image, two augmented views are produced: the source and the target. The source view is partially masked to form the context and then encoded to be used as conditioning for the world model, instantiated by the predictor. The target is encoded through an exponential moving average of the encoder, and target positions are sampled as the masked patches of the source image. Conditioned on the transformation parameters between the source and target, the encoded source image, and the positions to predict, the predictor is trained to predict the target representations.

# A   Experimental details

## A.1   Pretraining

We provide a more detailed architecture for IWM in figure S1.

**Architecture and optimization.** All of our models use a ViT-B/16 encoder trained for 300 epochs on ImageNet. We use the AdamW optimizer Loshchilov and Hutter (2019) with $1 \times 10^{-3}$ as our learning. We further use $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate follows a linear warmup for 40 epochs and then a cosine annealing. We use an iteration per epoch scale of 1.25 for the scheduler, which stretches the scheduler and makes the training end before the end of the schedule. Not having a 0 learning rate near the end of training was found beneficial in our experiments. We use a cosine weight decay schedule which goes from 0.04 to 0.4.

**Source and target.** In practice we build the source and target separately by first applying a random crop of scale between 0.3 and 1. We then apply a horizontal flip with probability 0.5. We will call the resulting image $I'$.

**Target transformations.** Starting from $I'$ we then apply a color jitter with probability 0.8, brightness maximum strength 0.4, contrast maximum strength 0.4, hue maximum strength 0.1, and saturation maximum strength 0.2.

**Source transformations.** Starting from $I'$ we apply a color jitter with probability 0.8, brightness maximum strength 0.4, contrast maximum strength 0.4, hue maximum strength 0.1, and saturation maximum strength 0.2. A gaussian blur of radius between 0.1 and 2 is applied with probability 0.2, solarization with probability 0.2 and grayscale with probability 0.2. These augmentations correspond to the ones used in BYOL (Grill et al., 2020). We then generate a mask $M_x$ as the union of 4 masks of area between 0.15 and 0.2 of the image, with aspect ratios between 0.75 and 1.5. All of the patches in $M_x$ are then dropped from the source $x$.

**Predictor conditioning.** We rely on the feature mixing strategy. Consider a mask token $m \in \mathbb{R}^d$ and $a \in \mathbb{R}^k$ a vector of $k$ scalars corresponding to augmentation parameters. We first add position embeddings to $m$ to indicate which patch of the target it needs to predict. We then concatenate $m$ and $a$ and feed them through a three layer fully-connected network with ReLU activation and dimensions $d, d, d$. This gives us a mask token that contains information about all of the transformation. Both the geometric aspect of where to predict and details on the photometric augmentations.

## A.2   Evaluation

For all evaluations on image classification, the augmentations applied to compute the validation accuracy are a resize to 256 followed by a 224 by 224 center crop. All hyperparameters reported are the optimal ones, chosen after careful tuning for every method.

*Linear.*   We take inspiration from the protocol of Chen et al. (2021). We train for 90 epochs on ImageNet. We sample random crops of images with scale between 0.08 and 1, then apply a horizontal flip with probability
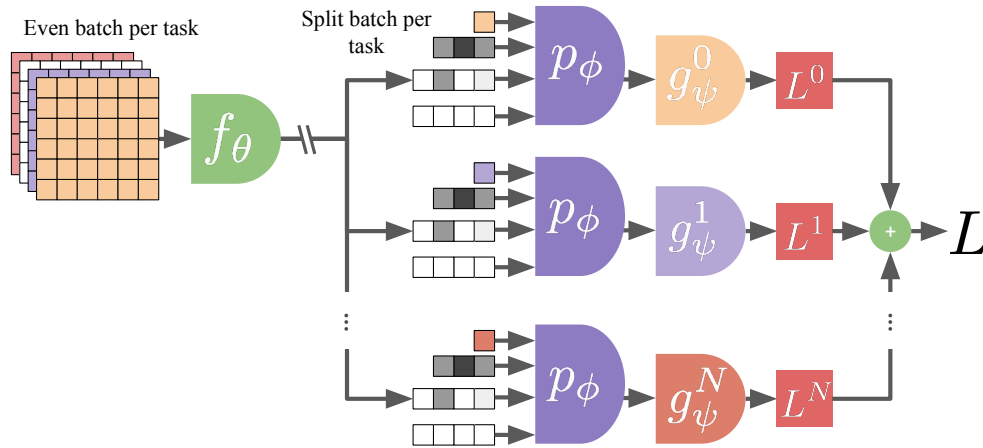
0.5.

The features are average pooled along the sequence axis to obtain a global representation which is then fed to a linear layer. We use a batch size of 16,384, with the LARS (You et al., 2017) optimizer and a learning rate of 6.4 with a warmup of 10 epochs. The learning rate then follows a cosine annealing schedule. Weight decay is set to 0 and momentum to 0.9.

*Attentive.*  The attentive head is taken from Chen et al. (2023). It consists of a cross attention block where the attention is computed between an additional token the unpooled representations. This allows an adaptive pooling strategy. We train for 90 epochs on ImageNet. We sample random crops of images with scale between 0.3 and 1, then apply a horizontal flip with probability 0.5. We also apply the same augmentations as used for the source transformations besides masking. We use a batch size of 1024 and AdamW optimizer with a learning rate of $1 \times 10^{-4}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. It follows a cosine annealing schedule. We use a weight decay of 0.01 kept constant during training.

*Encoder finetuning.*  We append a linear layer to the end of the encoder as for the linear evaluation and train for 100 epochs on ImageNet. We use the same RandAugment (Cubuk et al., 2020) strategy as MAE (He et al., 2021) as well as CutMix (Yun et al., 2019) and MixUp (Zhang et al., 2018). For RandAugment we use the string 'rand-m9-mstd0.5-inc1'. We use random erasing with probability 0.25 in pixel mode. We use a mixup $\alpha$ of 0.8, cutmix $\alpha$ of 1 and label smoothing of 0.1.

For the optimization we use AdamW with a learning rate of $2 \times 10^{-3}$ with 5 epochs of warmup followed by a cosine annealing schedule, weight decay of 0.005 and a batch size of 1024. We also use a drop path rate of 0.2 through the encoder and a layer wise learning rate decay of 0.65.

*Predictor finetuning.*  When finetuning the predictor we use an attentive head on top of the predictor output. We plug the predictor on top of the teacher network and it is tasked with predicting the whole target image, with null transformation parameters. We use the same augmentation protocol as for encoder finetuning. We train for 100 epochs on ImageNet with a batch size of 1024. We use AdamW for the optimizer, a learning rate of $1 \times 10^{-3}$ with a 5 epoch warmup then cosine annealing schedule. We use a weight decay of 0.1, no layer wise lr decay and a drop path rate of 0.2 through the predictor. Importantly, if the predictor is pretrained we divide it's learning rate by 10, and keep it identical to the attentive if head if random.



**Figure S2** Multitask tuning of the predictor. We sample a batch uniformly across task which is then fed through the predictor with an additional task token, indicating which task is being solved. The predictions are then fed through a task specific head and losses are summed.

*Multitask predictor finetuning.*  To give a clearer presentation of the protocol, we provide a graphical version of multitask predictor finetuning in figure S2. For the training in itself, we follow the same protocol as for predictor finetuning but train for the equivalent of 50 ImageNet epochs. The batch size used is 512 for each

task, where the batch is independently split between tasks. When training on a single task, we simply use 512 as the batch size and also train for 50 ImageNet epochs.

*End to end finetuning.* We follow the protocol of predictor finetuning but tweak certain parameters. First, the encoder also gets his learning rate divided by 10 like the predictor. The factors are treated separately and ablated for all methods. We use a 0.9 layer decay across the combination of predictor and encoder. The learning rate used is $2 \times 10^{-3}$ and all other parameters are identical to predictor finetuning.

*Segmentation.* We give here details about our protocol for semantic segmentation evaluations. We use the MMSegmentation library Contributors (2020). We fine-tune our pretrained models (either encoder only, predictor only, or end-to-end) with an UperNet head Xiao et al. (2018) on the ADE20k semantic segmentation dataset Zhou et al. (2019) for 160k iterations and report the validation mIoU. We concatenate the last 4 layers of the predictor, or encoder for encoder only finetuning, and feed the result to the segmentation head. At training time we resize the images at the pretraining resolution. At testing time we do not resize the images and interpolate the positional embeddings to the original resolution. For all setups and methods we pick the best run among several learning rate values: $1e - 5$, $2e - 5$ and $3e - 5$. We use a weight decay of 0.01 and a linear learning rate decay schedule.

## B  Complete finetuning results

**Table S1** Complete results of tables 4 and 5.

| Method | Epochs | No predictor | Frozen encoder, tuned predictor | | End to end |
|---|---|---|---|---|---|
| | | Encoder | Random Init. | Pretrained | |
| DINO | 1600 | 82.8 | 82.0 | N/A | 82.1 |
| MOCOv3 | 300 | 83.2 | 56.4 | N/A | 79.4 |
| iBOT | 1600 | **84.0** | **83.0** | N/A | 82.8 |
| | 300 | 82.7 | 82.4 | 82.7 (+0.3) | 82.3 |
| MAE | 600 | 83.2 | 82.8 | 83.0 (+0.2) | 83.1 |
| | 1600 | 83.6 | **83.0** | 83.1 (+0.1) | 83.3 |
| I-JEPA | 300 | 83.0 | 79.1 | 80.0 (+0.9) | 82.0 |
| $\text{IWM}_{12,384}^{\text{Inv}}$ | 300 | 83.3 | 80.5 | 81.3 (+0.8) | 82.7 |
| $\text{IWM}_{12,384}^{\text{Equi}}$ | 300 | 82.7 | 81.3 | 82.7 (+1.4) | 83.3 |
| $\text{IWM}_{18,384}^{\text{Equi}}$ | 300 | 82.9 | 81.5 | **83.3** (+1.8) | **84.4** |

We provide complete results for table 4 and table 5 in table S1. Some interesting behaviors are $\text{IWM}_{12,384}^{\text{Equi}}$ and MoCov3 in predictor finetuning. For $\text{IWM}_{18,384}^{\text{Equi}}$, we see the same behavior as $\text{IWM}_{12,384}^{\text{Equi}}$ but with slightly lower performance. This is consistent across all evaluations. Yet, even when accounting for scale of the predictor to compare with I-JEPA and $\text{IWM}_{12,384}^{\text{Inv}}$, all of our previous conclusions still hold. For MoCov3, it was the only method which did not perform well when attaching a random predictor to it. While we do not have conclusive evidence, we hypothesize that it is related to the low norm of its output. Adding a normalization between the encoder and predictor did not help.

# C  Impact of data augmentation

**Table S2** Impact of data augmentation strategy on IWM's performance. In all settings, destructive augmentations are never applied to the target.

| Predictor | Strengths | | | | Probabilities | | | | MRR | Linear | Attentive | Pred. ft. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Bright. | Contrast | Sat. | Hue | Jitter | Blur | Gray. | Solarize | | | | |
| | 0.4 | 0.4 | 0.2 | 0.1 | 0.8 | 0.2 | 0.2 | 0.1 | 0.09 | 74.5 | 77.0 | 81.3 |
| | 0.4 | 0.4 | 0.2 | 0.1 | 0.8 | | | | 0.11 | 72.1 | 76.5 | 80.5 |
| $IWM_{12,384}$ | 0.4 | 0.4 | 0.2 | 0.1 | 0.8 | 0.4 | 0.4 | 0.2 | 0.22 | 71.0 | 74.9 | 80.9 |
| | 0.5 | 0.5 | 0.4 | 0.2 | 0.8 | 0.2 | 0.2 | 0.1 | 0.81 | 69.3 | 75.5 | 82.7 |
| | 0.5 | 0.5 | 0.4 | 0.2 | 0.8 | | | | 0.07 | 73.3 | 76.3 | 80.1 |
| | | | | | | 0.2 | 0.2 | 0.1 | 0.02 | 72.9 | 76.3 | 80.7 |
| | 0.4 | 0.4 | 0.2 | 0.1 | 0.8 | 0.2 | 0.2 | 0.1 | 0.79 | 67.5 | 75.1 | 83.3 |
| | 0.4 | 0.4 | 0.2 | 0.1 | 0.8 | | | | 0.25 | 70.1 | 74.8 | 81.4 |
| $IWM_{18,384}$ | 0.4 | 0.4 | 0.2 | 0.1 | 0.8 | 0.4 | 0.4 | 0.2 | 0.85 | 56.1 | 74.5 | 83.1 |
| | 0.5 | 0.5 | 0.4 | 0.2 | 0.8 | 0.2 | 0.2 | 0.1 | 0.85 | 34.3 | 71.0 | 81.7 |
| | 0.5 | 0.5 | 0.4 | 0.2 | 0.8 | | | | 0.83 | 69.2 | 75.8 | 83.3 |
| | | | | | | 0.2 | 0.2 | 0.1 | 0.02 | 70.9 | 74.8 | 81.5 |

We study in table S2 the impact of augmentations used during pretraining, along with the depth of the predictor. We notice that depth is a deciding factor in the quality of the learned world model, where 4 out 5 scenarios with color are able to achieve color equivariance for the 18 layer predictor, compared to only 1 for the 12 layer predictor. The strength of the augmentations also plays a role and too weak augmentations do not lead to an equivariant model.

*On the asymmetry of augmentations.*  The asymmetry of augmentations is both a conceptual choice, to make the augmentations used more similar to contrastive approaches, but also a practical one. When learning an invariant world model with symmetric augmentations we noticed a drop in performance of 2 points on ImageNet in attentive probing and 1.5 points on linear probing. While this drop is not catastrophic, it is sufficient to recommend using asymmetric augmentations. As the depth of the predictor decreases, we expect this gap to widen.
On the other hand, when looking at an equivariant predictor, we did not notice any notable change in performance. This suggests that learning world models can also help improve stability over the choice of augmentations. The predictor does not have to be designed by keeping in mind which information may get removed but only by whether or not it can apply the transformation.

# D  Impact of the prediction task on predictor finetuning performance

In order to use predictor finetuning to solve downstream tasks, we need to apply a prediction task. We aim at giving a more combinatorial view of table 3 in this appendix.

**Table S3** Predictor finetuning performance which different prediction tasks.

| Method | Null latents | On teacher | Pred only one token | Accuracy |
|---|---|---|---|---|
| | ✓ | ✓ | | 83.3 |
| | ✓ | ✓ | ✓ | 82.8 |
| | ✓ | | | 83.1 |
| | ✓ | | ✓ | 82.6 |
| $IWM_{12,384}$ | | ✓ | | 83.2 |
| | | ✓ | ✓ | 82.8 |
| | | | | 82.9 |
| | | | ✓ | 82.9 |

We can see in table S3 that the conclusions drawn from table 3 still hold over a larger setting. Notably, using null latents is more flexible while not changing performance, using the teacher always gives a small boost in performance, and predicting only one token lowers performane by roughly half a point.

# E  Scaling to larger models

In order to scale to larger models, such as a ViT-L/16 encoder, multiple challenges need to be overcome. Notably, both the depth and the width of the predictor must be scaled in order to increase the number of parameters of the predictor to a suitable number. Scaling the width can lead to instabilities and hyperparameters such as the EMA schedule become more important. We noticed that a ratio of predictor weights/encoder weights of around 0.3 is suitable to learn a good world model.

**Table S4**  With a ViT-L/16 encoder, we observe a similar trend as with the base model. Significant gains are observed with a good world model, allowing it to surpass encoder finetuning.

| Method | Epochs | Encoder | Predictor | End to end |
|---|---|---|---|---|
| I-JEPA | 300 | 84.1 | 79.9 | |
| $\text{IWM}^{\text{Inv}}_{18,384}$ | 300 | 84.3 | 81.5 | |
| $\text{IWM}^{\text{Equi}}_{36,512}$ | 300 | 83.7 | 85.0 | 85.4 |

We study in table S4 the performance when scaling to a larger ViT-L/16. We see that the observation we made with the smaller ViT-B/16 still hold. The invariant model is the best on encoder finetuning, and predictor finetuning improves the performance significantly. Here again, end-to-end finetuning leads to performance gains.

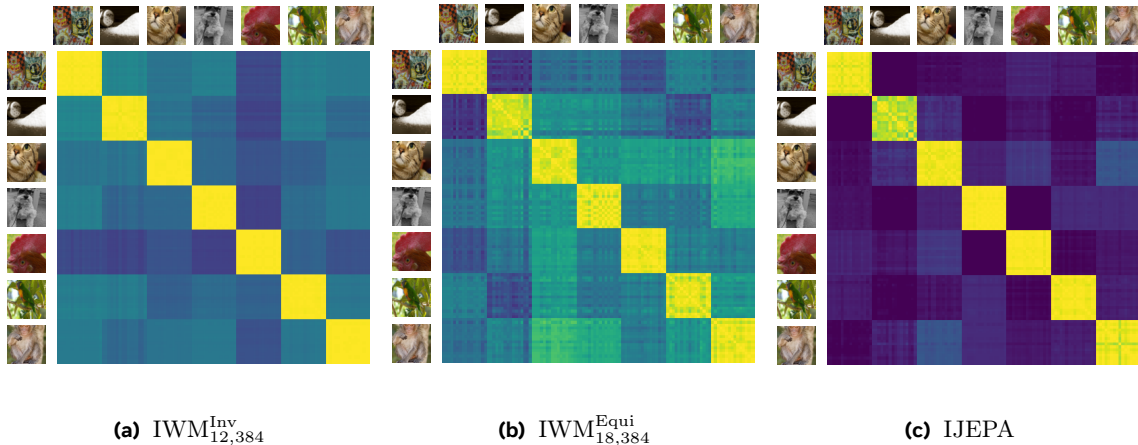# F  Evaluation on downstream datasets beyond ImageNet

We evaluated I-JEPA and IWM on iNaturalist18 (Horn et al., 2018), SUN397 (Xiao et al., 2010) and Places205 (Zhou et al., 2014) using attentive probing. We train our models for 50 epochs on iNaturalist18, 12 for Places205 and 28 for SUN397.

**Table S5**  When evaluating with attentive probing on downstream task, being equivariant improves performance across the board. All methods use ViT-B/16 encoders and were pretrained for 300 epochs on ImageNet

| Method | ImageNet | iNat18 | SUN397 | Places205 |
|---|---|---|---|---|
| MAE | 73.5 | 50.1 | 70.2 | 60.3 |
| I-JEPA | 75.0 | 50.4 | 69.2 | 58.3 |
| $\text{IWM}^{\text{Inv}}_{12,384}$ | **77.0** | 51.6 | 71.0 | 59.4 |
| $\text{IWM}^{\text{Equi}}_{18,384}$ | 75.1 | **54.2** | **71.7** | **60.5** |

As we can see in table S5, IWM consistently improves over I-JEPA and MAE when pretraining all methods for 300 epochs. We notice that while $\text{IWM}^{\text{Equi}}_{18,384}$ is not the top performing model on ImageNet, it significantly outperforms it's invariant counterpart, with gains of 2.6 points on iNaturalist, 0.7 points on SUN397 and 1.1 point on Places205. This suggests that while the richness of the representation of an equivariant model is not optimal for in domain performance, it helps improve generalisation to downstream tasks.

## G  Visualizing representation differences between invariant and equivariant behavior



**(a)** $\text{IWM}^{\text{Inv}}_{12,384}$  **(b)** $\text{IWM}^{\text{Equi}}_{18,384}$  **(c)** IJEPA

**Figure S3** Difference in embedding space between invariant and equivariant behaviours. Each image is augmented 16 times and we compute the similarity matrix between all images. The yellow regions indicate high similarities between samples originating from the same image. We can see more variations in the equivariant model, or in I-JEPA where invariance is not enforced. This suggests that augmentations influence the representation more in these models.

As we see in figure S3, the invariant model collapses augmented views to very similar embeddings, as shown by the high similarity in the diagonal blocks. On the other hand the equivariant model shows more variation, which shows that augmentation information is more present in the representation. Interestingly, I-JEPA has a behaviour in between because it was not trained to be either invariant or equivariant. I-JEPA has no force controlling how information is kept or removed from the representation.

## H  On the meaning and role of invariance in Self-Supervised learning

One of the key component of the success of self-supervised learning is augmentation invariance Chen et al. (2020a). We can say that we have learned invariant representations if $\forall a, \ f_\theta(x) = f_\theta(\mathcal{T}(a,x))$. However there are many scenarios that satisfy this property. The two main ones that we are interested in are:
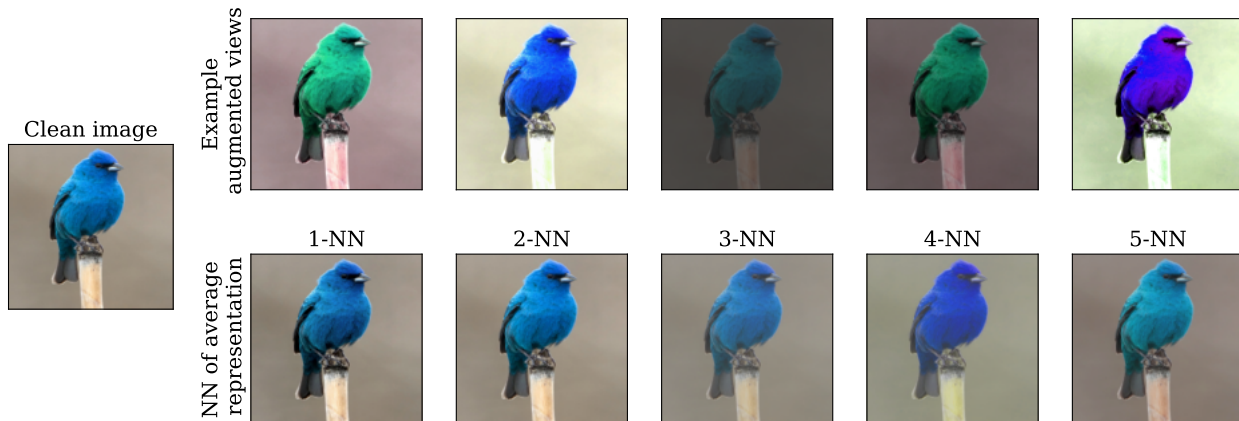
- Any augmented view leads to the same information as the clean image
- The encoder removes the information related to the transformation

In the first case, the representations still contain all of the information about the input, whereas in the second we are removing information that can be deemed superfluous. In the case of contrastive methods, the focus is usually on removing information. Indeed if an image and its grayscale version are made to have the same representation, the encoder must remove color information. This is one of the key drivers of performance of such methods. By removing information until only the semantics of the image remains, the representations will be easy to leverage for a task such as classification.

We can thus wonder if the first invariance scenario also leads to improved performance, and if we can even leverage it. As we have demonstrated how IWM is able to preserve information, and we have a predictor that can apply transformations, we can marginalize over augmentations to create invariant representations in an efficient way. Here, we do not need to apply the encoder on all augmented views, but can directly use the predictor which is more compute efficient. If we consider a set of randomly sampled augmentations $A$ such that $\text{card}(A) = N$ we can compute an invariant representation as

$$z_x^{\text{Inv}} = \frac{1}{N} \sum_{i=1}^{N} p_\phi\left(f_\theta(x), A_i, m_{A_i}\right)$$

We can then visualize which image has representation most similar to $z_x^{\text{Inv}}$ and see if using $z_x^{\text{Inv}}$ improves performance on classification task.



**Figure S4** Retrieval of invariant representations computed using 256 augmentations in latent space. In the top row we visualize some of the corresponding image and on the bottom the nearest neighbours of the invariant representation. We can notice that the nearest neighbour is the original non-augmented image, followed by images with small transformations.

As we can see in figure S4, the images that have representations which are most similar with $z_x^{\text{Inv}}$ are the clean image and images with small transformations. We also know that our encoder preserves augmentation related information and is thus not invaraint to transformations. Combining these two facts tells us that the marginalizaiton process creates a clean representations, akin to the first kind of invariance.
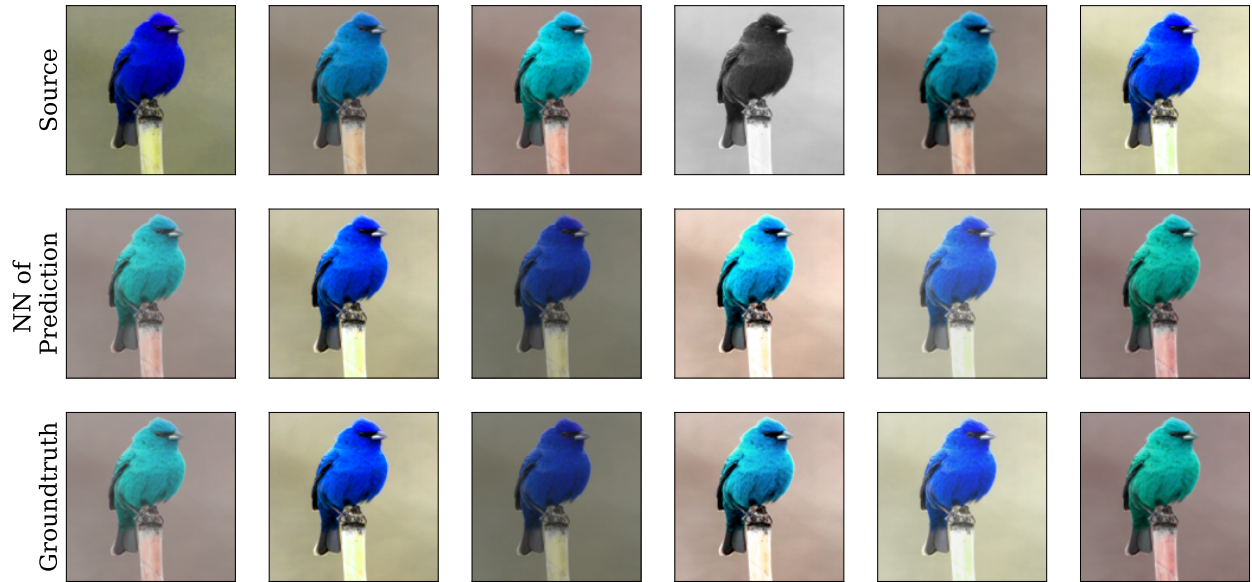
**Table S6** Linear evaluation on marginalized representations. Using more augmented prediction to create an invariant representation does not improve performance.

| Number of predictions ($N$) | 1 (default) | 8 | 16 | 32 | 128 |
|---|---|---|---|---|---|
| ImageNet Top-1 accuracy (%) | 64.5 | 64.3 | 64.6 | 64.6 | 64.4 |

However, when looking at table S6 we can see that no performance gain is present when using invariant representations obtained by marginalizing over predictions. This is true even with 128 augmented views, which already increases the compute budget by a factor of around 64. As such, using invariant representations that preserve the content of the image is not necessarily beneficial for downstream evaluation.

Overall, the key to the success of augmentation invariance in contrastive learning is not just in building invariant representations, but in the way that the representations are invariant. Building invaraince by removal of information has been shown to be very effective (Chen et al., 2020a), whereas we see here that invariance by always predicting the representation of the clean image is not necessarily helpful. This does not mean that equivariant representations cannot build invariances that are useful from downstream tasks, as the contrary was shown in Chavhan et al. (2023b), but that we have to be careful in how we create invariant representations.

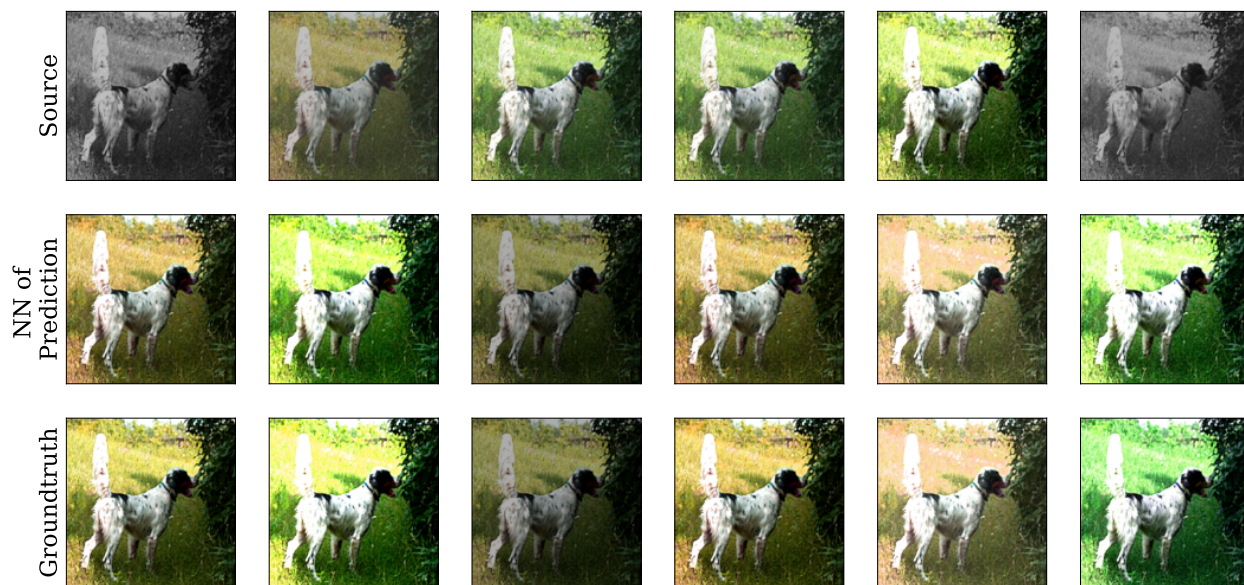# I    Additional qualitative evaluations of the world model



**Figure S5**  Randomly selected retrieval samples of our world model. For each image, we generate 256 augmented views and apply transformations in latent space. We then retrieve the nearest neighbor of the prediction and visualize whether it is close to the groundtruth or not. The learned world model performs well in most settings but has some inaccuracies with inverting grayscale.
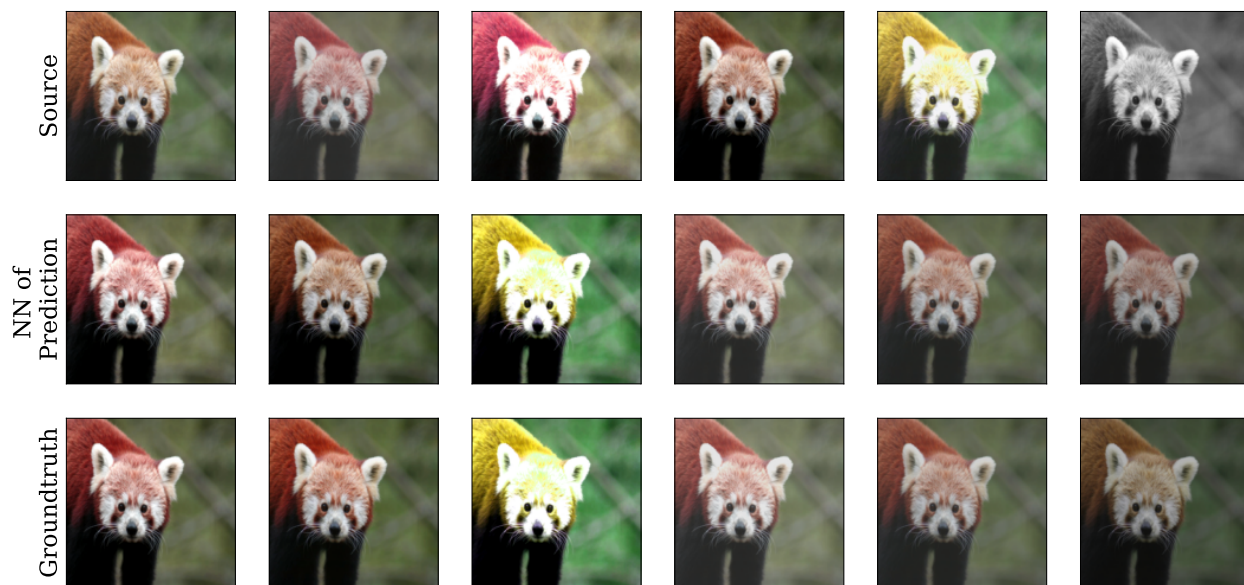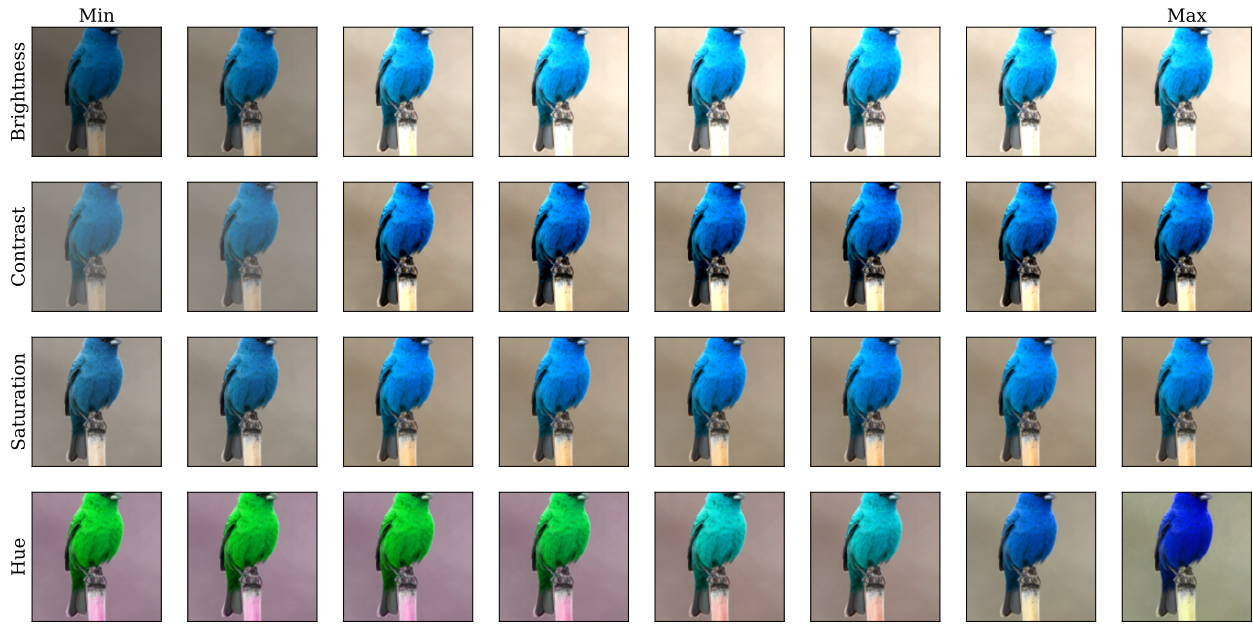


**Figure S6**  Randomly selected retrieval samples of our world model. For each image, we generate 256 augmented views and apply transformations in latent space. We then retrieve the nearest neighbor of the prediction and visualize whether it is close to the groundtruth or not. The learned world model performs well in most settings but has some inaccuracies with inverting grayscale.
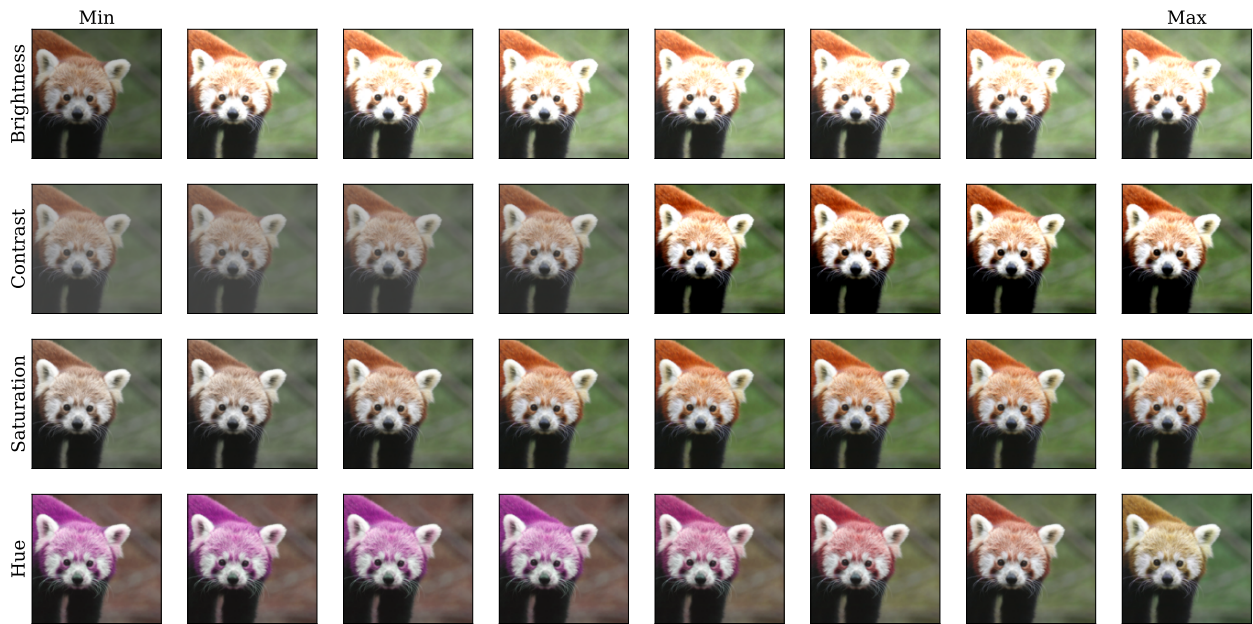
**Figure S7** Randomly selected retrieval samples of our world model. For each image, we generate 256 augmented views and apply transformations in latent space. We then retrieve the nearest neighbor of the prediction and visualize whether it is close to the groundtruth or not. The learned world model performs well in most settings but has some inaccuracies with inverting grayscale.
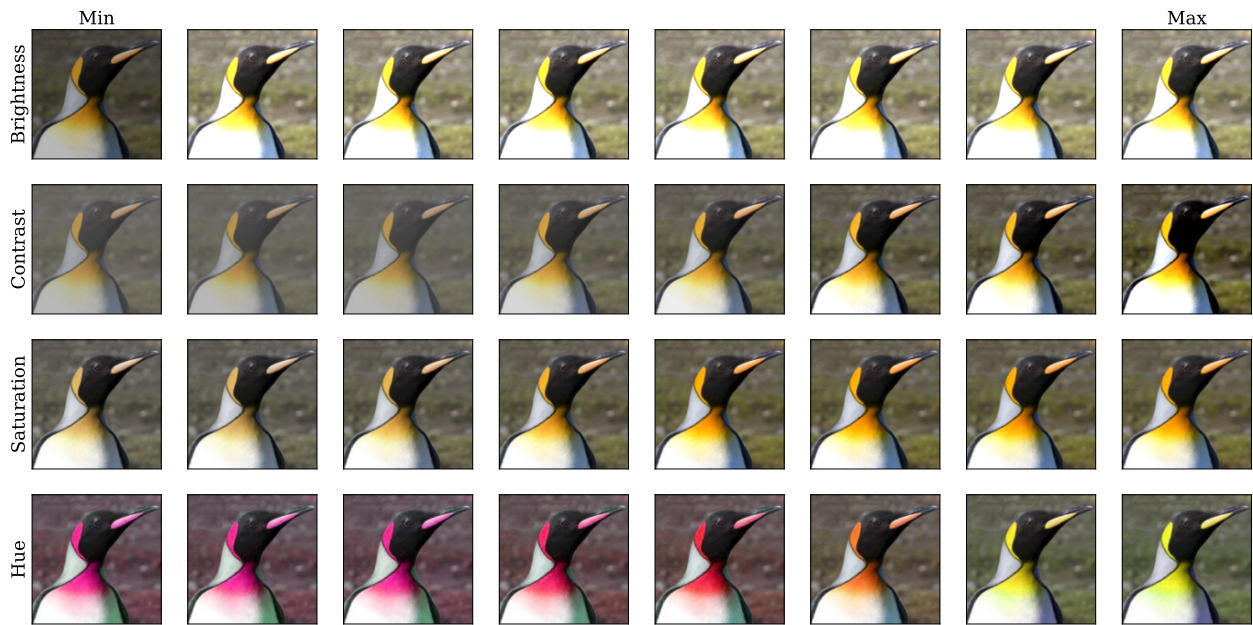


**Figure S8** Randomly selected retrieval samples of our world model. For each image, we generate 256 augmented views and apply transformations in latent space. We then retrieve the nearest neighbor of the prediction and visualize whether it is close to the groundtruth or not. The learned world model performs well in most settings but has some inaccuracies with inverting grayscale.

**Figure S9** Application of the world model on precise transformations. For each parameter, we vary its value on a grid to see whether the model is able to predict small changes. The model is able to show the gradient of transformations, highlighting again the capabilities of the world model. We can still notice some imperfections however, as the model was only trained on combinations of augmentations. To make changes more visible, we used a model trained with a strong color jitter for this figure.



**Figure S10** Application of the world model on precise transformations. For each parameter, we vary its value on a grid to see whether the model is able to predict small changes. The model is able to show the gradient of transformations, highlighting again the capabilities of the world model. We can still notice some imperfections however, as the model was only trained on combinations of augmentations. To make changes more visible, we used a model trained with a strong color jitter for this figure.

**Figure S11** Application of the world model on precise transformations. For each parameter, we vary its value on a grid to see whether the model is able to predict small changes. The model is able to show the gradient of transformations, highlighting again the capabilities of the world model. We can still notice some imperfections however, as the model was only trained on combinations of augmentations. To make changes more visible, we used a model trained with a strong color jitter for this figure.