# MarioNette: Self-Supervised Sprite Learning

**Dmitriy Smirnov** [1]  **Michaël Gharbi** [2]  **Matthew Fisher** [2]  **Vitor Guizilini** [3]  **Alexei A. Efros** [2]  **Justin Solomon** [1]

## Abstract

Visual content often contains recurring elements. Text is made up of glyphs from the same font, animations, such as cartoons or video games, are composed of sprites moving around the screen, and natural videos frequently have repeated views of objects. In this paper, we propose a deep learning approach for obtaining a graphically disentangled representation of recurring elements in a completely self-supervised manner. By jointly learning a dictionary of texture patches and training a network that places them onto a canvas, we effectively deconstruct sprite-based content into a sparse, consistent, and interpretable representation that can be easily used in downstream tasks. Our framework offers a promising approach for discovering recurring patterns in image collections without supervision.

## 1. Introduction

Since the early days of machine learning for graphics and vision, the accepted unit of image synthesis has been *the pixel*. But while the pixel grid structure is a natural representation for display hardware and convolutional generators, pixels comprise an extremely low-level representation that does not easily permit high-level reasoning and editing.

In this paper, we propose an alternative approach to deep learning on visual content that establishes an atomic unit that is more interpretable, editable, and meaningful than the pixel: the *sprite*. Our goal is to train a model that simultaneously learns a dictionary of sprites capturing recurring visual elements in a collection of images, and learns to explain the input images as combinations of these sprites.

We are inspired by sprite-based animation, a popular early technique for drawing cartoons that became a mainstay of rendering in early computer games. In sprite-based animation, an artist draws a collection of small patches—a *sprite sheet*—consisting of characters in various poses, texture swatches, static objects, and so on. Then, each frame is created by compositing some subset of these patches onto a canvas. The process of creating additional content requires minimal effort and can be automated procedurally or in real time based on user input.

Given a dataset of images with repeated structure, e.g., portraying the same character or from the same virtual world, we learn to decompose each image into a collection of 2D sprites selected from a learned dictionary, common across the dataset. Our algorithm learns this dictionary and the decomposition simultaneously. Whereas standard CNN-based generators hide their feature representation in their intermediate layers, our model wears its representation "on its sleeve," explicitly placing learned sprite patches on a background canvas. By drawing sprites from a dictionary rather than synthesizing pixels from hidden neural features, it provides readily-interpretable visual representations.

Our contributions include the following:

- We describe a grid-based anchor system along with a learned dictionary of textured patches (with transparency) to extract a sprite-based image representation.
- We propose a method to learn the patch dictionary and the grid-based representation jointly in a differentiable, end-to-end fashion.
- We compare to past work on learned disentangled graphics representations.
- We show how our method generalizes to finding visual patterns in natural images and video.

We will release all code and models upon publication.

## 2. Related Work

Decomposing visual content into semantically meaningful parts for analysis, synthesis, and editing is a long-standing problem. We review the most closely related work.

**Layered decompositions.** In an early work, Wang & Adelson (1994) decompose videos into image layers undergoing temporally-varying spatial warps, effectively grouping the video content by similarity of motion. They show applications to compression. Similarly, Flexible Sprites (Jojic & Frey, 2001) represents videos with full-canvas semi-

---

[1]Massachussetts Institute of Technology [2]Adobe Research [3]Toyota Research Institute. Correspondence to: Dmitriy Smirnov <smirnov@mit.edu>.

transparent layers, where each layer's motion is limited to translations. They solve for the layers and motion parameters using expectation maximization and show applications to video editing. Although Kannan et al. (2005) extend their motion model, we adopt the simpler translation-only motions. Unlike Flexible Sprites, ours are local; we limit motion to a small neighborhood around anchors, which makes inference tractable with many ($\geq 100$) sprites. Our representation factorizes spatially and temporally recurring local patterns, simplifying downstream editing.

Lu et al. (2020) decompose videos of humans into independent video layers, enabling retiming of individual actions, but unlike in sprite-based techniques, motion and appearance are not disentangled. Sbai et al. (2020) use a layered representation as inductive bias in a GAN with layers limited to solid colors. Automatic decompositions of images into "soft layers" according to texture, color or semantic features have been used in image editing (Aksoy et al., 2017; 2018). Gandelsman et al. (2019) use deep image priors (Ulyanov et al., 2018) to separate images into layer pairs. Reddy et al. (2020) discover patterns in images by making compositing differentiable and using it as a strong inductive bias.

**Interpretable generators for neural synthesis.** Neural networks significantly improve the fidelity and realism of generative models (Goodfellow et al., 2014; Karras et al., 2020), but control and interpretability are still wanting (Chen et al., 2016; Härkönen et al., 2020; Bau et al., 2018; 2019). Several works explore interpretability using domain-specific functions made differentiable for gradient-based optimization. For instance, Hu et al. (2018) constrain their generator to a small set of parametric photographic edits for image editing. Li et al. (2018) make well-known imaging operators differentiable to improve fidelity while preserving their efficiency compared to costly neural networks. Mildenhall et al. (2020) use a strong ray-marching prior and rendering model to train a multilayer perceptron that encodes a radiance field for novel view synthesis. Neural textures (Thies et al., 2019) replace RGB textures on 3D meshes with high-dimensional trainable features. Rendering the mesh under new views and decoding the projected features makes view-consistent editing possible. Lin et al. (2018) use spatial transformers in their generator to obtain interpretable, explicit geometric transformations. Our "generator" synthesizes frames by explicitly compositing 2D sprites undergoing rigid motions, enabling direct interpretation and control over appearance and motion.

**Object-centric representations.** Our learned sprites can reveal, segment, and track object instances. Similar to our approach, Slot Attention (Locatello et al., 2020) extracts an object-centric compositional representation of videos in an unsupervised fashion. However, our sprites are easily interpretable—motion and appearance are direct outputs of

our model. Our model scales to significantly more objects per scene and is much faster to train.

SPACE (Lin et al., 2020) also seeks meaningful decompositions into object layers using unsupervised learning. Our method achieves a higher recall of recurring sprite patterns, while SPACE tends to embed them in the background layer, providing no control. MONet (Burgess et al., 2019) decomposes images into multiple object regions using an attention mechanism. Recent works use pre-determined parametric primitives as building blocks in their image representations (Smirnov et al., 2020; Li et al., 2020).

Using our sprite decompositions on video games, we can learn about their dynamics and gameplay mechanics, which can benefit downstream learning agents (Justesen et al., 2019). GameGAN (Kim et al., 2020) trains a Generative Adversarial Network to synthesize new frames from a controller's input signal. They split their rendering into a static and a dynamic component, but otherwise render full frames, with no factorization into individual parts. Unlike our sprites, their generator is difficult to interpret: appearance and dynamics are deeply entangled within its parameters.

**Compression.** Appearance consistency and motion compensation are central to video compression techniques (Bachu & Chari, 2015). Recent compression methods apply deep learning (Lu et al., 2018; Lombardo et al., 2019). Our algorithm models videos as compositions of finitely many moving sprites, factoring out redundancy in the input signal. We learn a static full-canvas background sprite and a collection of recurring foreground sprites. This draws inspiration from works like DjVu (Haffner et al., 1999), which compresses scanned documents by separating them into two layers: a background layer compressed using a wavelet approach and the foreground text. DjVu exploits the recurrence of glyphs to compress the foreground. Digipaper (Huttenlocher et al., 1999) uses a similar approach, with three layers but allows for color text. Image epitomes (Jojic et al., 2003) summarize and compress an image's shape and appearance content into a miniature texture (or epitome). The original is reconstructed by sampling from the epitome using a smooth mapping. Our sprite dictionary fills a similar role but provides superior editing control.

## 3. Method

We start with an input sequence of $n$ color frames $\{I_1, \ldots, I_n\}$ with resolution $w \times h$. Our goal is to decompose each frame $I_i \in \mathbb{R}^{3 \times w \times h}$ into a set of possibly overlapping sprites, organized into $\ell$ depth layers, selected from a finite-size dictionary. The dictionary is a collection of trainable latent codes $\{z_1, \ldots, z_m\}$ that are decoded into RGBA sprites using a neural network generator (§3.1).

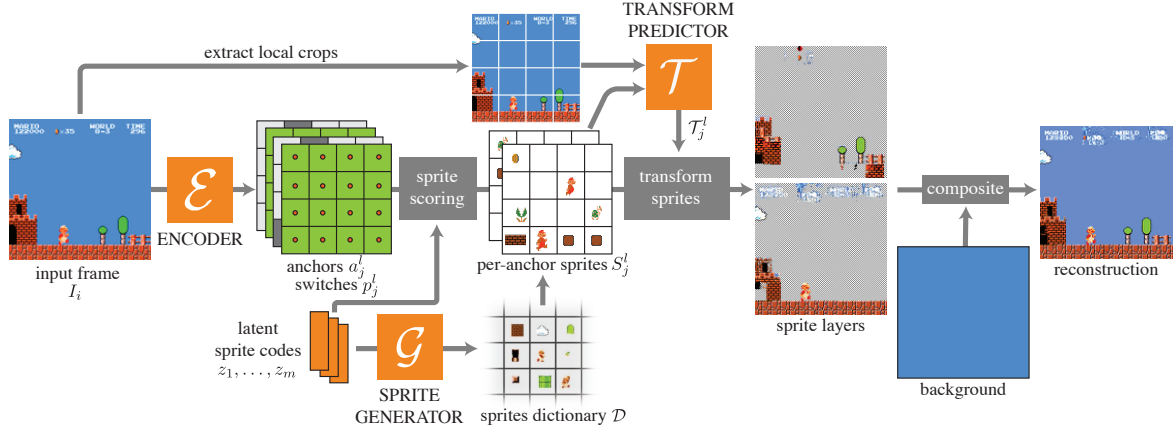Our training pipeline is illustrated in Figure 1. We start by

*Figure 1.* Overview of our method. We jointly learn a sprite dictionary, represented as a set of trainable latent codes decoded by a generator, as well as an encoder network that embeds an input frame into a grid of latent codes (which we call anchors). Comparing the anchor embeddings to the dictionary codes lets us assign a sprite to each grid cell. Our encoder also outputs a binary switch per anchor to turn sprites on and off. After compositing, we obtain a reconstruction of the input. We train our networks in a self-supervised fashion by optimizing a reconstruction loss.

processing each input frame with a convolutional encoder to produce $\ell$ grids of feature vectors, one grid per depth layer (§3.2). The grids are typically lower resolution than the input frame, with a downsampling factor proportional to the sprite size. We call the center of each grid cell an *anchor*. We compare each anchor's feature vector against the dictionary's latent codes, using a softmax scoring function, to select the best matching sprite for each anchor (§3.3). Using our sprite generator, we decode each anchor's matching sprite. This gives us a grid of sprites for each of the $\ell$ layers. To factorize image patterns that may not align perfectly with our anchor grid, we allow the sprites to move in a small neighborhood around their anchors (§3.4). We composite the layers from back to front onto the output canvas to obtain our final reconstruction (§3.5). The background is modeled as a special learnable sprite that covers the entire canvas.

We train the dictionary latent codes, frame encoder, and sprite generator jointly on all the frames, comparing our reconstruction to the input (§3.6). This self-supervised training procedure yields a representation that is sparse, compact, interpretable, and well-suited for downstream editing and learning applications.

### 3.1. Dictionary and sprite generator

The central component of our representation is a global *dictionary* of $m$ textured patches or sprites $\mathcal{D} = \{P_1, \ldots, P_m\}$, where each $P_i \in \mathbb{R}^{4 \times k \times k}$ is an RGBA patch. Our sprites have an alpha channel, which allows them to be partially transparent, with possibly irregular (i.e., non-square) boundaries. This is useful for representing animations with multiple depth layers. The dictionary is shared among all frames; we reconstruct frames using only sprites from the dictionary.

Instead of optimizing for the RGBA pixel values directly, we represent the dictionary as a set of trainable latent codes $\{z_1, \ldots, z_m\}$, with $z_i \in \mathbb{R}^d$. We decode these latent codes into RGBA sprites using a fully-connected sprite generator $P_i = \mathcal{G}(z_i)$. This latent representation allows us to define a similarity metric over the latent space, which we use to pair anchors with dictionary sprites to best reconstruct the input frame (§3.3). At test time, we can forego the sprite generator and edit the RGBA sprites directly. Unless otherwise specified, we set the latent dimension to $d = 128$ and the patch size to $k = 32$.

We randomly initialize the latent codes from the standard normal distribution. Our sprite generator first applies zero-mean unit-variance normalization—Layer Normalization (Ba et al., 2016), without an affine transformation—to each latent code $z_i$ individually, followed by one fully-connected hidden layer with $8d$ features, Group Normalization (Wu & He, 2018), and a ReLU activation. We obtain the final sprite using a final fully-connected layer with a sigmoid activation to keep the RGBA values in $[0, 1]$. Latent code normalization is crucial to stabilize training and keep the latent space in a compact subset as the optimization progresses. See §4.2 for an ablation study of this and other components.

### 3.2. Layered frame decomposition using sprite anchors

Given an input frame, we seek a decomposition that best explains the image, using sprites from the dictionary. We exploit the translation invariance and locality in our representation; our sprites are "attached" to a regular grid of reference points that we call *anchors*, inspired by (Redmon et al., 2016; Girshick, 2015). Each anchor has at most one sprite; we call it *inactive* if it has none. We give the sprites freedom of motion around their anchors to better factorize
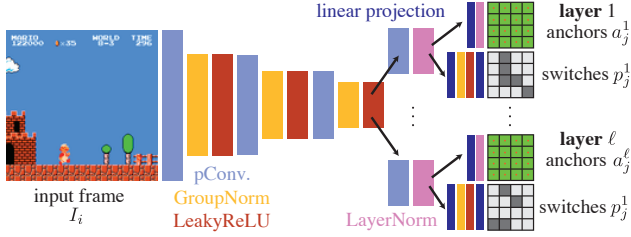
*Figure 2.* Encoder architecture.



*Figure 3.* Layered sprite decomposition with local anchors. We compare each anchor's embedding to all the dictionary latent codes (not shown) to assign at most one sprite per anchor. Our model also predicts a local transformation of each placed sprite in a small neighborhood around its anchor. To allow for occlusions between sprites, we allow for multiple sprite layers, which we compose back to front to obtain the final image.

recurring image structures that may not be aligned with the anchor grid. This local—or, Eulerian (Wu et al., 2012)—viewpoint makes inference tractable and avoids the pitfalls of tracking the possibly complex global motion of *all* the sprites across the image canvas (a Lagrangian viewpoint). To enable multiple layers with sprites occluding each other, we estimate $\ell > 1$ anchor grids for each frame (we use $\ell = 2$ in our experiments). Figure 3 illustrates our layered anchor grids and local sprite transformations.

We use a convolutional encoder $\mathcal{E}$ to map the $w \times h$ RGB frame $I_i$ to $\ell$ grids of anchors, with resolution $\frac{2w}{k} \times \frac{2h}{k}$. Each anchor $j$ in layer $l$ is represented by a feature vector $a_j^l \in \mathbb{R}^d$—which characterizes the local image appearance around the anchor, and an active/inactive switch probability $p_j^l \in [0, 1]$. Our frame encoder uses $\log_2(k) - 1$ downsampling blocks, which use partial convolutions (Liu et al., 2018) with kernel size 3 and stride 2 (for downsampling), Group Normalization, and Leaky ReLU activation. It produces a tensor of intermediate features for each layer, which are normalized using Layer Normalization. From these, we obtain the anchor switches using an MLP with one hidden layer of size $d$ followed by Group Normalization and Leaky ReLU. We obtain anchor features using a linear projection followed by a second Layer Normalization. The encoder architecture is illustrated in Figure 2.

### 3.3. Per-anchor sprite selection

Once we have the layered anchor grids for the input frame, we need to assign sprites to the active anchors. We do this by scoring every dictionary element $i$ against each anchor $j$ at layer $l$, using a softmax over dot products between dictionary codes and anchor features:

$$s_{ij}^l = \frac{\exp\left(a_j^l \cdot z_i / \sqrt{d}\right)}{\sum_{k=1}^m \exp\left(a_j^l \cdot z_k / \sqrt{d}\right)}. \tag{1}$$

Recall that both the anchor features and dictionary latent codes are individually normalized using a Layer Normalization operator. Restricting both latent spaces to a compact subset helps stabilize the optimization and avoid getting stuck in local optima. During training, each anchor's sprite is a weighted combination of the dictionary elements,
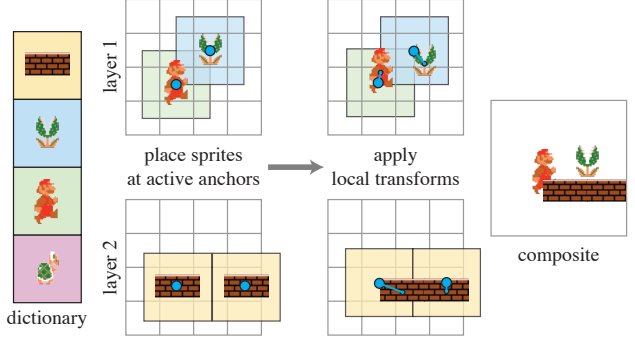
masked by the anchor's active probability:

$$S_j^l = p_j^l \sum_{i=1}^m s_{ij}^l P_i. \tag{2}$$

This soft selection of dictionary patches allows gradients to propagate to both the dictionary and the anchor features during training. Except for with natural image and video datasets, at test time, we use hard selections, i.e., for each anchor, we pick the sprite ($S_j^l := P_i$) with highest score $s_{ij}^l$. We binarize the switches at test time $p_j^l \in \{0, 1\}$.

### 3.4. Local sprite transformations

In real animations, sprites will rarely perfectly align with our regular anchor grid so, to avoid learning several copies of the same sprites (e.g., all sub-grid translations of a given image pattern), we allow the sprites to move around their anchors. In our implementation, we only allow 2D translations of up to $1/2$ the sprite size on each side of the anchor, i.e., $\mathcal{T}_j^l = (x_j^l, y_j^l) \in [-k/2, k/2]^2$.

We use a convolutional network to predict the horizontal and vertical translation offsets from the anchor's sprite and a crop of the input frame centered around the anchor, with identical spatial dimensions. This network follows the architecture of $\mathcal{E}$ followed by an MLP with a single hidden layer of size $d$, Group Normalization, and Leaky ReLU. Specifically, we concatenate the image crop and the anchor's sprite $S_j^l$ along the channel dimension, and pass this tensor through this network to obtain the $x_j^l$ and $y_j^l$ offsets.

We represent the spatial offsets as discrete pixel shifts in $\{-k/2, \ldots, k/2\}$ using a softmax classification output over the $k + 1$ possible outcomes (independently for each spatial dimension). We found this discrete encoding better behaved than using a continuous offset with a spatial trans-
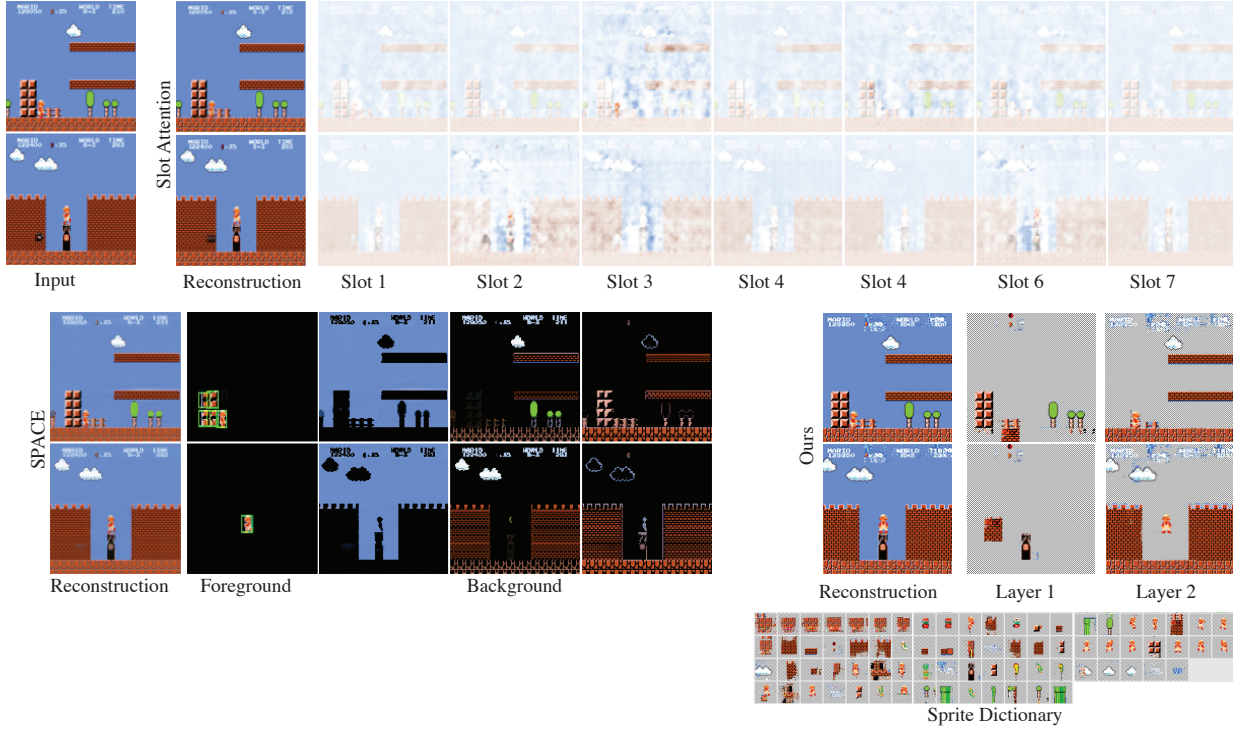
*Figure 4.* Qualitative comparison to SPACE (Lin et al., 2020) and Slot Attention (Locatello et al., 2020). While all three methods obtain a good reconstruction, SPACE only recognizes a small subset of sprites in the frame and Slot Attention does not yield a meaningful decomposition, distributing the input across all slots. We decompose the entire foreground and produce a learned sprite dictionary.

former (Jaderberg et al., 2015). In particular, the discrete encoding allows the gradient signal to propagate to all shifts rather than the weak local gradient from bilinear interpolation in spatial transformers.

### 3.5. Compositing and reconstruction

Each anchor in our layered representation is now equipped with a sprite $S_j^l$ and a transformation $\mathcal{T}_j^l$. We filter out inactive anchors ($p_j^l = 0$). For each layer $l$, we transform the active sprites in their anchor's local coordinate system and render them onto the layer's canvas, initialized as fully transparent. Because of the local transformation, neighboring sprites within a layer may overlap. When this happens, we randomly choose an ordering, as in Figure 3. This random permutation encourages our model to avoid overlapping sprites or make sprite colors agree in the overlap region.

As shown in Figure 1, we optionally learn a background texture to capture elements that cannot be explained using repeatable sprites. This can be thought of as a special patch of resolution equal to or greater than that of a single frame. If the background is larger than a frame, the reconstruction process specifies a (discrete) position offset in the background from which to crop. In some experiments, we use a simpler background model: a single fixed solid color, determined by analyzing the input frame before training. In this variant, we sample 100 random frames, cluster the pixel

values into 5 clusters using $k$-means, and choose the largest cluster center as background color. Finally, we combine the background and sprite layers using standard back to front alpha compositing (Porter & Duff, 1984).

### 3.6. Training procedure

Our pipeline is fully differentiable. We train the latent codes dictionary, sprite generator, frame encoder, transformation predictor and background layer jointly by minimizing a simple $L_2$ loss between our reconstructions and the ground truth frames. Additionally, we impose a Beta$(0.5, 0.5)$ prior on the switches, dictionary element scores, and spatial offset probabilities to favor values close to 0 or 1, by adding a log-likelihood term to the loss, weighed by a factor of 0.0001.

We optimize the parameters using the AdamW (Loshchilov & Hutter, 2019) optimizer, for 42 hours on a single GeForce GTX 1080 GPU, using batch size 8 and learning rate 0.0001, except for the background module, which we train with learning rate 0.001 when used.

## 4. Experimental Results

We evaluate our self-supervised decomposition on several datasets of video games and natural images, compare to related work, and conduct an ablation study to motivate our design choices.
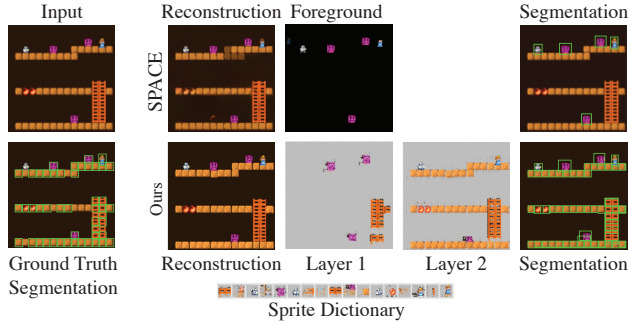
*Figure 5.* Qualitative comparison to SPACE (Lin et al., 2020) on the synthetic game dataset. We show ground truth sprite segmentations as well as the segmentations obtained from SPACE foreground and from our learned sprites. While SPACE only learns several of the sprites, we reconstruct the entire foreground using our learned dictionary.

In all figures, we render a high-frequency checkerboard in areas where a sprite is transparent. Unless otherwise specified, the sprite order within a dictionary is determined by sorting along a 1-dimensional $t$-SNE embedding of the sprite latent codes. We find this sorting tends to group together semantically similar sprites, making the dictionary easier to interpret and manipulate.

## 4.1. Comparisons

While, to the best of our knowledge, no prior works specifically target unsupervised sprite-based reconstruction, we provide comparisons to two recent state-of-the-art methods that obtain similarly disentangled representations.

In Figure 4, we compare to SPACE (Lin et al., 2020) and Slot Attention (Locatello et al., 2020). The former decomposes a scene into a foreground layer consisting of several objects as well as a background, segmented into three layers. The latter deconstructs a scene into several discrete "slots." We train both methods to convergence using their default parameters. While both achieve good reconstructions of the input frames, SPACE only recognizes a few sprites in its foreground layer, and Slot Attention does not achieve a semantically meaningful decomposition. In contrast, our method models the entire scene using learned sprites.

Additionally, we evaluate on a synthetically-generated sprite-based game from (Dubey et al., 2018), which is made of sprites on a solid background. We compare quantitatively to SPACE in Table 1 and show qualitative results in Figure 5. Since we have ground truth segmentation of each scene into sprites, we are able to compute accuracy using SPACE foreground objects and our learned sprites. While SPACE consistently discovers several sprites in each frame, yielding a high precision, it defers many sprites to the background, which results in a low recall. Our method is able to learn the sprites as dictionary elements, resulting in a
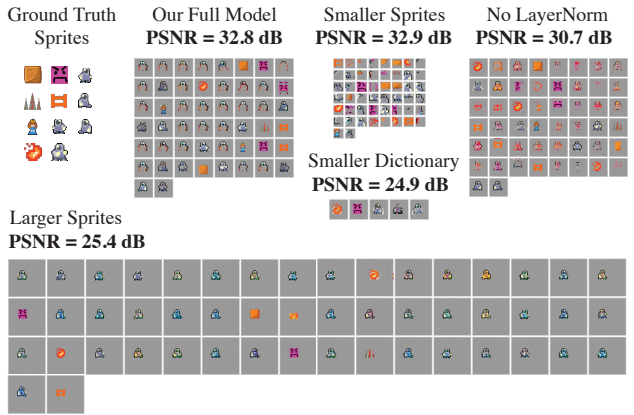


*Figure 6.* Ablation study on floating sprites dataset. With proper sprite size, our full model recovers the ground truth sprites and achieves a high reconstruction quality (measured by PSNR). Reducing sprite size preserves quality, but the sprites we learn correspond to sub-parts of the originals. Reducing the number of sprites increases compression but can also reduce quality if the dictionary is too small. Removing Layer Normalization of the dictionary and anchor latent codes hurts reconstruction and sprite quality.

significantly higher recall. This is additionally highlighted in the IoU for foreground, where our method outperforms SPACE by $20\times$. The main factors that detract from our precision and recall are the fact that we sometimes split a sprite into two or represent two sprites with one (e.g., we capture the staircase using a single sprite consisting of two staircase elements as well as another sprite with a side element). Quantitative measures of accuracy cannot reflect such segmentation ambiguities. Our method creates visually pleasing outputs, but our reconstructions use exact sprites from the dictionary, so small misalignments and appearance discrepancies can incur a high error. SPACE fixes these mismatches in the background layer, and achieves a lower quantitative error.

## 4.2. Ablation

We train each model for 250,000 iterations (about 15 hours) on a synthetic dataset of floating sprites. We take 11 ground truth sprites used in the synthetic dataset, and, for each

| Method | Precision | Recall | IoU | PSNR |
|--------|-----------|--------|------|------|
| SPACE  | **97.5%** | 2.17%  | 0.0421 | **31.9** |
| Ours   | 53.5%     | **35.6%** | **0.843** | 23.5 |

*Table 1.* Quantitative comparison to SPACE (Lin et al., 2020). We compute object discovery precision and recall per sprite using an IoU threshold. Additionally, we report IoU with respect to foreground segmentation as well as PSNR to evaluate reconstruction quality. Our method is able to recognize significantly more sprites than SPACE, resulting in a higher recall and IoU.
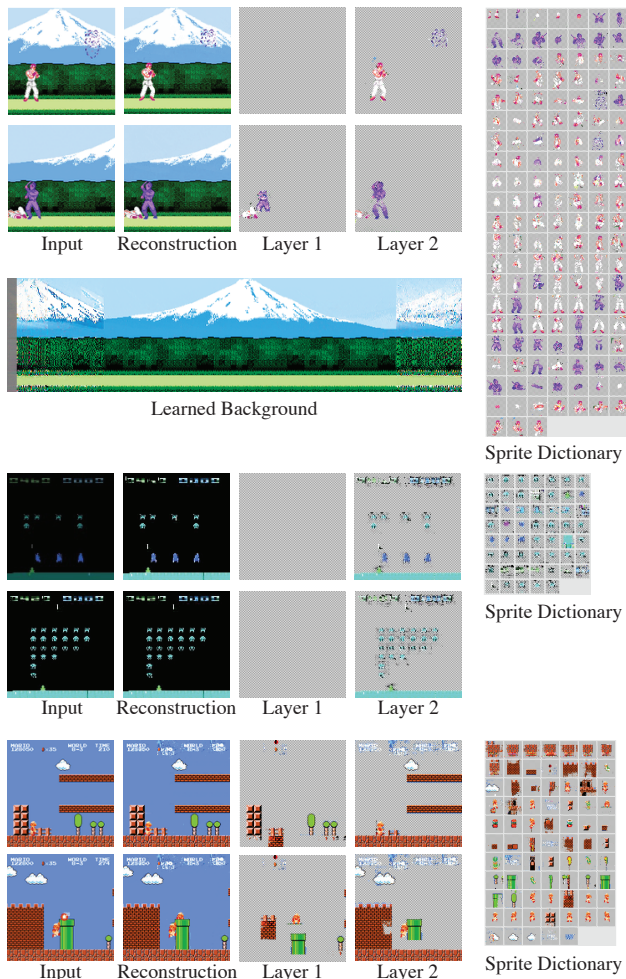
Input    Reconstruction    Layer 1    Layer 2

Learned Background

Sprite Dictionary

Input    Reconstruction    Layer 1    Layer 2

Sprite Dictionary

Input    Reconstruction    Layer 1    Layer 2

Sprite Dictionary

*Figure 7.* Video game decompositions. We show results of our learned sprite decomposition for three 2D games. Our self-supervised technique recovers a compact dictionary of semantically meaningful sprites representing characters (or their body parts) and props. The first example shows our learned background texture; the other two use a solid color as background.

frame, randomly place between 5 and 15 of them on a white canvas. We train a full model, one with smaller $16 \times 16$ sprites, another with larger $64 \times 64$ sprites, a model with an underparameterized dictionary (just five elements), and a model without Layer Normalization. We show the learned dictionary and reconstruction PSNR in Figure 6.

### 4.3. Sprite-based game deconstruction

We train on gameplay of Fighting Hero (one level, 5,330 frames), Super Mario Bros. (one level, 2,220 frames), and Space Invaders (5,000 frames). We use patch size $k = 32$ for Mario and Fighting Hero and $k = 16$ for Space Invaders. For Fighting Hero, we learn a background texture, as described in §3.5.

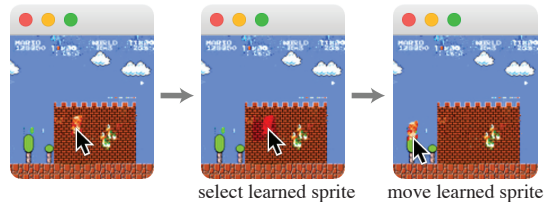The sprites, background, and example frame reconstructions



select learned sprite    move learned sprite

*Figure 8.* GUI to edit a game using our learned decomposition.

are shown in Figure 7. Our model is able to successfully disentangle foreground from background and recover a reasonable sprite sheet for each game. Having reverse-engineered the games, we can use the decomposition for applications such as editing. In Figure 8, we demonstrate a GUI that allows the user to move sprites around the screen.

### 4.4. Natural image decomposition

Although our method is designed with sprite-based animation in mind, it can also generalize to natural images and videos. In Figure 9, we show results of training our model on a dataset of a tennis practice video (4,000 frames). The model learns parts of the tennis player's body (head, limbs, shirt, etc.) as sprites and captures most of the tennis court background in the learned background texture. By simply selecting the sprites in the dictionary that make up the tennis player, we successfully segment each frame of the dataset.

We also can use our model to discover recurring patterns in a single natural image. We train on random $128 \times 128$ crops of a $768 \times 512$ photograph from the 2013 US Presidential Inauguration, which contains many repeating elements such as stairs, columns, and people. With just a dictionary of 50 $32 \times 32$ sprites (51,200 pixels), we are able to recover much of the detail in the original 393,216 pixel image.

### 4.5. Limitations

To highlight some of our approach's limitations, we apply it to automatic font glyph discovery. We train on six scanned pages of *Moby Dick*, each of approximately $500 \times 800$ resolution, taking a random $128 \times 128$ crop of a random page at train time. In Figure 11, we illustrate an excerpt of input text, our reconstruction, and the learned sprite dictionary. We remove patches that are used less than once on average and sort the dictionary by frequency of occurrence, in decreasing order from left to right, top to bottom for visualization.

This dataset differs significantly from others that we test on. Each input frame consists of many densely packed sprites (approximately 100 glyphs in each $128 \times 128$ crop). Additionally, many individual glyphs consist of smaller repeating elements. We hypothesize that because of this, combined with the fact that there are no motion cues between frames that our model can leverage, we are unable to achieve a
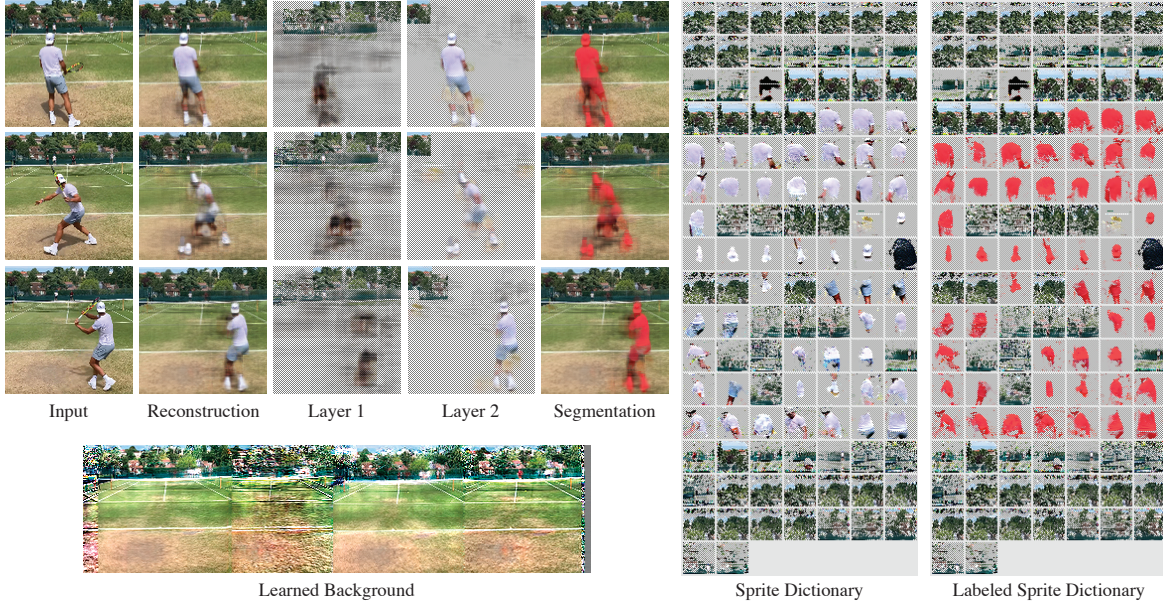
Input        Reconstruction        Layer 1        Layer 2        Segmentation

Learned Background

Sprite Dictionary        Labeled Sprite Dictionary

*Figure 9.* Segmentation of natural video. Despite its simplistic motion and appearance model, our approach can be applied to real-world videos. By selecting the few sprites corresponding to the tennis player, we can quickly obtain a segmentation of the full video sequence.



Input        Reconstruction        Sprite Dictionary

*Figure 10.* Reconstruction of a single natural image. In this example, we factorize the recurring elements of a single high-resolution image into a compact learned sprite dictionary.



Input

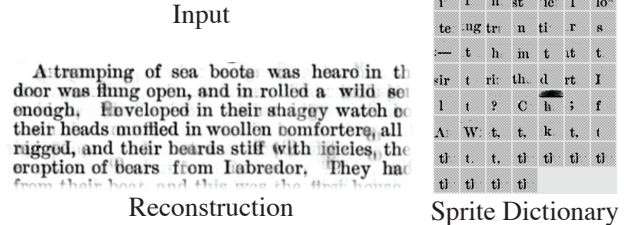Reconstruction        Sprite Dictionary

*Figure 11.* Reconstruction of a scanned excerpt of text, illustrating some limitations of our method. Without strong motion cues from individual letters, our models learns sprites comprising more than a single glyph. Because text is tightly packed, 2 layers of sprites are often insufficient for perfect reconstruction when several glyphs appear close to a single anchor.

perfect reconstruction, learning certain sprites with multiple glyphs and others with just partial glyphs. Incorporating priors tailored to regularly structured and dense data like text is a potential future direction.

## 5. Conclusion

We present a self-supervised method to jointly learn a patch dictionary and a frame-encoder from a video, where the encoder explains frames as compositions of elements from the dictionary, anchored on a regular grid. By generating multiple layers of alpha-masked sprites and predicting per-sprite local transformation, our model recovers fine-scale motion, achieves high-quality reconstructions, and produces semantically meaningful, well-separated sprites. Applied to content with significant recurrence, such as text or 2D video games, our approach recovers structurally significant patches like glyphs or moving characters.

Understanding recurring patterns and their relationships is central to machine learning. Learning to act intelligently,

in video games or in the physical world, requires breaking experiences down into recurrent elements between which knowledge can be transferred effectively. In this work we have focused on a simplified video domain. In the future, we would like to significantly expand the range of deformations applied to the learned dictionary elements, such as appearance or shape changes. We believe our work opens significant avenues for future research to explore recurrences and object relationships in more complex domains.

# References

Aksoy, Y., Aydin, T. O., Smolić, A., and Pollefeys, M. Unmixing-based soft color segmentation for image manipulation. *ACM Transactions on Graphics (TOG)*, 36 (2):1–19, 2017.

Aksoy, Y., Oh, T.-H., Paris, S., Pollefeys, M., and Matusik, W. Semantic soft segmentation. *ACM Transactions on Graphics (TOG)*, 37(4):1–13, 2018.

Ba, J. L., Kiros, J. R., and Hinton, G. E. Layer normalization. *arXiv:1607.06450*, 2016.

Bachu, S. and Chari, K. M. A review on motion estimation in video compression. In *International Conference on Signal Processing and Communication Engineering Systems*, pp. 250–256. IEEE, 2015.

Bau, D., Zhu, J.-Y., Strobelt, H., Zhou, B., Tenenbaum, J. B., Freeman, W. T., and Torralba, A. GAN dissection: Visualizing and understanding generative adversarial networks. *arXiv:1811.10597*, 2018.

Bau, D., Zhu, J.-Y., Wulff, J., Peebles, W., Strobelt, H., Zhou, B., and Torralba, A. Seeing what a GAN cannot generate. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4502–4511, 2019.

Burgess, C. P., Matthey, L., Watters, N., Kabra, R., Higgins, I., Botvinick, M., and Lerchner, A. MONet: Unsupervised scene decomposition and representation. *arXiv:1901.11390*, 2019.

Chen, X., Duan, Y., Houthooft, R., Schulman, J., Sutskever, I., and Abbeel, P. InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. *arXiv:1606.03657*, 2016.

Dubey, R., Agrawal, P., Pathak, D., Griffiths, T. L., and Efros, A. A. Investigating human priors for playing video games. In *ICML*, 2018.

Gandelsman, Y., Shocher, A., and Irani, M. "Double-DIP": Unsupervised image decomposition via coupled deep-image-priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11026–11035, 2019.

Girshick, R. Fast R-CNN. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1440–1448, 2015.

Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks. *arXiv:1406.2661*, 2014.

Haffner, P., Bottou, L., Howard, P. G., and LeCun, Y. DjVu: Analyzing and compressing scanned documents for internet distribution. In *Proceedings of the Fifth International Conference on Document Analysis and Recognition. ICDAR'99 (Cat. No. PR00318)*, pp. 625–628. IEEE, 1999.

Härkönen, E., Hertzmann, A., Lehtinen, J., and Paris, S. Ganspace: Discovering interpretable GAN controls. *arXiv:2004.02546*, 2020.

Hu, Y., He, H., Xu, C., Wang, B., and Lin, S. Exposure: A white-box photo post-processing framework. *ACM Transactions on Graphics (TOG)*, 37(2):1–17, 2018.

Huttenlocher, D., Felzenszwalb, P., and Rucklidge, W. Digipaper: A versatile color document image representation. In *Proceedings 1999 International Conference on Image Processing (Cat. 99CH36348)*, volume 1, pp. 219–223. IEEE, 1999.

Jaderberg, M., Simonyan, K., Zisserman, A., and Kavukcuoglu, K. Spatial transformer networks. *arXiv:1506.02025*, 2015.

Jojic, N. and Frey, B. J. Learning flexible sprites in video layers. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 1, pp. I–I. IEEE, 2001.

Jojic, N., Frey, B. J., and Kannan, A. Epitomic analysis of appearance and shape. In *Proceedings Ninth IEEE International Conference on Computer Vision*, pp. 34–41 vol.1, 2003.

Justesen, N., Bontrager, P., Togelius, J., and Risi, S. Deep learning for video game playing. *IEEE Transactions on Games*, 12(1):1–20, 2019.

Kannan, A., Jojic, N., and Frey, B. J. Generative model for layers of appearance and deformation. In *AISTATS*, volume 2005, pp. 1. Citeseer, 2005.

Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. Analyzing and improving the image quality of StyleGAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8110–8119, 2020.

Kim, S. W., Zhou, Y., Philion, J., Torralba, A., and Fidler, S. Learning to simulate dynamic environments with GameGAN. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2020.

Li, T.-M., Gharbi, M., Adams, A., Durand, F., and Ragan-Kelley, J. Differentiable programming for image processing and deep learning in halide. *ACM Transactions on Graphics (TOG)*, 37(4):1–13, 2018.

Li, T.-M., Lukáč, M., Gharbi, M., and Ragan-Kelley, J. Differentiable vector graphics rasterization for editing and learning. *ACM Transactions on Graphics (TOG)*, 39 (6):1–15, 2020.

Lin, C.-H., Yumer, E., Wang, O., Shechtman, E., and Lucey, S. ST-GAN: Spatial transformer generative adversarial networks for image compositing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9455–9464, 2018.

Lin, Z., Wu, Y.-F., Peri, S. V., Sun, W., Singh, G., Deng, F., Jiang, J., and Ahn, S. SPACE: Unsupervised object-oriented scene representation via spatial attention and decomposition. In *International Conference on Learning Representations*, 2020.

Liu, G., Reda, F. A., Shih, K. J., Wang, T.-C., Tao, A., and Catanzaro, B. Image inpainting for irregular holes using partial convolutions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

Locatello, F., Weissenborn, D., Unterthiner, T., Mahendran, A., Heigold, G., Uszkoreit, J., Dosovitskiy, A., and Kipf, T. Object-centric learning with slot attention. *arXiv:2006.15055*, 2020.

Lombardo, S., Han, J., Schroers, C., and Mandt, S. Deep generative video compression. In Wallach, H., Larochelle, H., Beygelzimer, A., d Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 32*, pp. 9287–9298. 2019.

Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.

Lu, E., Cole, F., Dekel, T., Zisserman, A., Salesin, D., Freeman, W. T., and Rubinstein, M. Layered neural rendering for retiming people in video. *ACM Trans. Graph.*, 2020.

Lu, G., Ouyang, W., Xu, D., Zhang, X., Cai, C., and Gao, Z. Dvc: An end-to-end deep video compression framework. *arXiv:1812.00101*, 2018.

Mildenhall, B., Srinivasan, P. P., Tancik, M., Barron, J. T., Ramamoorthi, R., and Ng, R. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, pp. 405–421. Springer, 2020.

Porter, T. and Duff, T. Compositing digital images. In *Proceedings of the 11th Annual Conference on Computer Graphics and Interactive Techniques*, pp. 253–259, 1984.

Reddy, P., Guerrero, P., Fisher, M., Li, W., and Mitra, N. J. Discovering pattern structure using differentiable compositing. *ACM Trans. Graph. (Proceedings of SIGGRAPH Asia 2020)*, 39(6):262:1–262:15, 2020. doi: https://doi.org/10.1145/3414685.3417830.

Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 779–788, 2016.

Sbai, O., Couprie, C., and Aubry, M. Unsupervised image decomposition in vector layers. In *IEEE International Conference on Image Processing*, Abu Dhabi, UAE, 2020. IEEE.

Smirnov, D., Fisher, M., Kim, V. G., Zhang, R., and Solomon, J. Deep parametric shape predictions using distance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 561–570, 2020.

Thies, J., Zollhöfer, M., and Nießner, M. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019.

Ulyanov, D., Vedaldi, A., and Lempitsky, V. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9446–9454, 2018.

Wang, J. Y. A. and Adelson, E. H. Representing moving images with layers. *IEEE Transactions on Image Processing*, 3(5):625–638, 1994.

Wu, H.-Y., Rubinstein, M., Shih, E., Guttag, J., Durand, F., and Freeman, W. Eulerian video magnification for revealing subtle changes in the world. *ACM Transactions on Graphics (TOG)*, 31(4):1–8, 2012.

Wu, Y. and He, K. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19, 2018.