

Top-down Flow Transformer Networks

Zhiwei Jia Haoshen Hong Siyang Wang Zhuowen Tu
University of California, San Diego
{zjia, h7hong, siw030, ztu}@ucsd.edu

Abstract

We study the deformation fields of feature maps across convolutional network layers under explicit top-down spatial transformations. We propose top-down flow transformer (TFT) by focusing on three transformations: translation, rotation, and scaling. We learn flow transformation generators that are able to account for the hidden layer deformations while maintaining the overall consistency across layers. The learned generators are shown to capture the underlying feature transformation processes that are independent of the particular training images. We observe favorable experimental results compared to the existing methods that tie transformations to fixed datasets. A comprehensive study on various datasets including MNIST, shapes, and natural images with both inner and inter datasets (trained on MNIST and validated in a number of datasets) evaluation demonstrates the advantages of our proposed TFT framework, which can be adopted in a variety of computer vision applications.

1. Introduction

Recently, deep neural networks [25, 17] have led to tremendous performance improvement on large-scale image classification [36] and other computer vision applications [12, 13, 30, 43, 5]. While Convolutional Neural Networks (CNNs) have shown great promise in solving many challenging vision problems, there remain fundamental questions about the transparency of representations in current CNN architectures. While the explicit role of the top-down process is a critical issue in perception and cognition, it has received less attention within the current CNN literature.

Currently, both training and testing of CNNs is performed in a data-driven manner by passing convolved features from lower layers to the top layers. However, visual perception systems are shown to engage both bottom-up and top-down processes [35, 18]. A top-down process would allow explicit generation and inference of transformations and (high level) configuration changes of images that is otherwise not convenient in a bottom-up process. For exam-

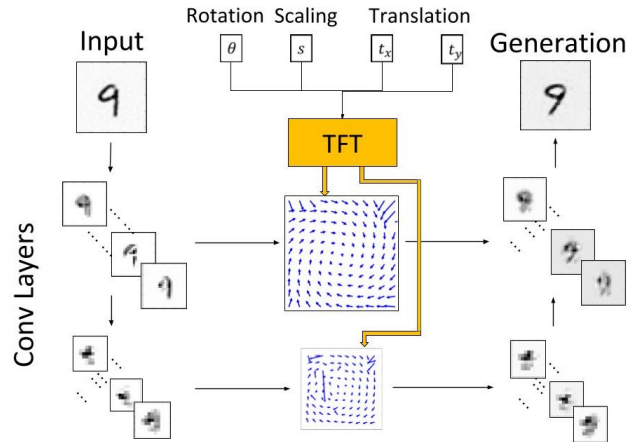


Figure 1. A schematic illustration of our top-down flow transformer framework (TFT).

ple, suppose we wish to train a CNN classifier to detect the translation of an object in an image. A data-driven way to train this CNN would require generating thousands of samples by moving the object around in the image. However, a top-down model, if available, can directly detect translation using two parameters of the translation. Computational models realizing the bottom-up and top-down visual inference have been previously proposed [32, 23, 38]. They, however, are not readily combined into end-to-end deep learning frameworks. Recurrent neural networks (RNN) [7, 20] have feedbacks recursively propagated between the output and input layers but miss explicit top-down generation.

Motivated by the recent development in CNNs [25] that learn effective hierarchical representations, the general pattern theory [16, 44, 46] that provides rigorous mathematical formulation for top-down generations, as well as findings from cognitive perception [15, 4, 11], we seek to build a top-down generator that operates directly on the feature maps of the internal CNN layers to model and account for spatial transformations.

In this paper, we pay particular attention to CNNs under top-down spatial transformations. There often exists

clear flow fields computed between the convolutional layers of the original image and the convolutional layers of the transformed images (after rotation, scaling, and translation) [8]. There is a clear pattern of consistent but nontrivial feature map deformations throughout the convolutional layers which is the key topic to be studied and leveraged here. Our goal is to discover and model operations in CNNs that lead to non-linear activity of the resulting flow fields. Given a source image and a transformed image under translation, rotation, and scaling, the internal CNN feature maps across multiple-layers can be directly computed; flow transformers modeled using an aggregated convolution strategy are themselves learned to perform mappings that transfer the feature maps of source image to the target image across all the intermediate CNN layers.

The training process is supervised since we generate transformed images using different parameters for translation, rotation, and scaling. Given the fact that no supervision is needed to have the explicit correspondences at the feature maps level and that transforming the source images can be readily accomplished by using the explicit transformation parameters, obtaining the training data can be done automatically. The learned top-down flow transformer (TFT) however demonstrates great generalization capable of transforming images that are not seen in the training set — a benefit of having a top-down generator that is not tied to specific images. TFT is therefore distinct from existing work [6, 34, 9] where transformations are learned with strong coupling to the training images that are hard to generalize to novel ones. It can perform all three kinds of flow transformations (rotation, translation and scaling) with arbitrary transformation parameters outside of those used in the training process. Also it generalizes well to various datasets, ranging from small patterns to natural images. Hence, our proposed frameworks are generic to the flow transformations of the feature maps obtained from input images. Moreover, TFT is designed to not be tied to fixed input dimensions of CNN features, and consequently it works for input images of a range of sizes.

In the experiments, we train the proposed top-down flow transformer (TFT) on the MNIST dataset and demonstrate the generated images by “inverting” transformed CNN features of images from several non-MNIST datasets. We perform ablation studies to justify the design choice of our top-down generator. Comparing with the competing transformer [34] shows the immediate benefit our approach.

2. Significance and Related Work

We first discuss the significance of our proposed top-down flow transformer (TFT).

Why a top-down generator? Top-down can play a fundamental role in unraveling, understanding, and enriching the great representation power of deep convolutional neural net-

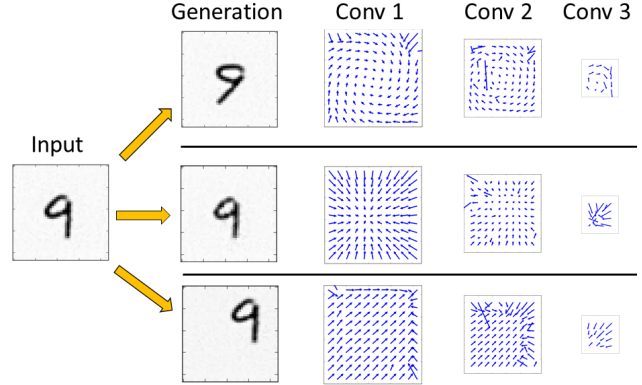


Figure 2. Feature flows of convolutional layers of a CNN after rotation (-30°), scaling ($\times 0.9$), and translation (move to top-right) using TFT: top-down information is reflected in the flow. This phenomenon becomes less clear in deeper layers.

works to bring more transparency and interoperability. The flow maps learned by TFT on novel images show promise for the direction of top-down learning.

Why transform features across the layers? The purpose of our method is to have a generator that can create images subject to spatial transformations from the given parameters. In addition, we are intrigued by the idea to understand how the CNN features change internally with respect to the spatial transformation. This will aid in understanding the representation learned by a CNN in order to improve its robustness and enrich its modeling and computing capabilities. This is not immediately available in the existing frameworks [6, 34] where the models are heavily coupled with the specific training data.

Next, we discuss related work. We first compare TFT with three lines of work that are of high relevance.

1. *Deep image analogy.* The deep visual analogy-making work [34] shows impressive results to learn to transform and compose images. However, method like [34, 9] builds heavily on an encoder-decoder strategy that is strongly tied to the training images; it learns to output results in the image space without modeling the internal feature maps. Applying learned transformer [34] to novel images therefore leads to unsatisfactory results, as shown in the experiment section.

2. *Learning image transformations.* Learning to transform images has been a quite active research area recently. Existing methods that target on building transformation generators [22, 28, 14, 26, 40] trains CNN/RNN to perform transformation on the the output feature or the image space, which is different from our goal of studying the intrinsic flow transformations within the CNN networks. For example, in [22] a spatial transformer was developed to explicitly account for the spatial manipulation of the data. This differentiable module can be inserted into existing CNNs, giving neural networks the ability to actively transform the feature-

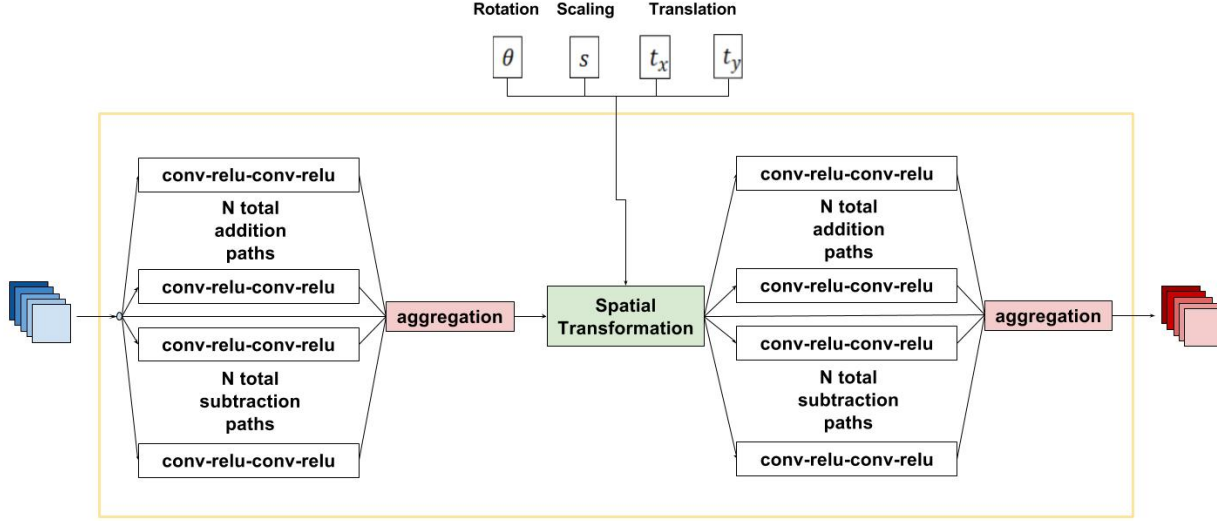


Figure 3. Architecture of our top-down flow transformer (TFT).

level representation. However, the main goal of [22] is to learn a differentiable transformation field through back-propagation to account for the spatial transformations. *STN learns to perform spatial transformation (without user specified transformation parameters) to match the output features whereas TFT studies the generator (with user specified transformation parameters) modeling the underlying changes to the feature maps in the result of spatial transformations to the images.*

3. *Feature transformation with SIFT-flow* [8]. An earlier attempt [8] (arxiv) has been made to study the feature layer deformation under explicit transformations based on flows computed from the SIFT-flow method [29]. Although the work [8] has a similar big picture idea to our work here but it is very preliminary and builds transformation solely based on the SIFT-flow estimation [29]. It therefore limited in several aspects: (1) only able to morph the features but cannot change the values; (2) may fail if SIFT-flow does not provide reliable estimation; (3) hard to generalize to arbitrary translation, rotation, and scaling. It relies heavily on the specially chosen parameters that do not work well under general situations. Our approach has a much greater learning capability than that in [8].

Next, we also discuss the existing literature in generative modeling. A family of mathematically well defined generators are defined in [16] as the general pattern theory. Its algorithmic implementation however still needs a great deal of further development. Methods [2, 3, 41, 46] developed prior to the deep learning era are inspiring but they have limited modeling capabilities. Deep belief net (DBN) [19] and generative adversarial networks (GAN) [13] do not study the explicit top-down generator for the image transformation. Other generators that perform feed-forward mapping

[6, 39, 45] have transformations as input parameters but the process in [6] maps directly from the input parameter space to the output image space; it cannot be applied to generate novel categories and does not study the intrinsic transformation.

Existing works in flow estimation [29, 5] are used to perform flow estimation, not as a generator for image synthesis. Our frameworks can be applied to generate novel images from images of various datasets.

3. Top-down Flow Transformer

3.1. Architecture

The network comprises three layers: an aggregated feature transformation layer, a linear layer performing spatial transformation, and another aggregated feature transformation layer. Consider the feature maps $f_x \in \mathbb{R}^{n \times n \times m}$ of a convolutional layer with m channels from a CNN that is pretrained for a discriminative task (e.g., image classification) by feeding an image x . The aggregated flow transformation layer here is given as:

$$\mathcal{F}(f_x) = f_x + \sum_{i=1}^N \sigma(w_i) \mathcal{T}_i(f_x) - \sum_{i=N+1}^{2N} \sigma(w_i) \mathcal{T}_i(f_x)$$

where $\sigma(\cdot)$ is the sigmoid function, w_i are some learnable scalars and N is an integer referred as the number of transformation functions. We enforce N addition paths and N subtraction paths for the aggregation process. Each transformation function $\mathcal{T}_i(\cdot)$ is defined as a two-layer convolutional network, each layer has a convolution operation followed by the rectified linear (ReLU) activation. The transformed feature maps are then fed into a linear layer that

applies spatial transformations, including translation, rotation and scaling, to each channel individually (with bilinear interpolations). The specific spatial transformations are modeled by the product of three transformation matrices:

$$M_{rot}(\theta) = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad M_{scale}(s) = \begin{bmatrix} s & 0 & 0 \\ 0 & s & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad \text{and} \quad M_{tran}(t_x, t_y) = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \quad \text{where } \theta,$$

s and (t_x, t_y) are top-down controls of rotation, scaling and translation, respectively. The results are further transformed by another aggregated residual layer to generate the outputs. Figure 2 illustrates the architecture. We build the TFT for each convolutional layer in the pre-trained CNN. For clarity, we denote the collection of the overall transformation parameters as $\Theta = (\theta, s, t_x, t_y)$.

3.2. Model Training

We train our proposed networks in a supervised manner by minimizing the average Euclidean distance $\mathcal{E} = \sum_{x, \Theta} E_{\Theta}(x)$ between the generated feature maps and the ground truth feature maps for any image x and any transformation parametrized by Θ in the training set. The ground truth features maps can be collected automatically from CNN features of input images that are under the corresponding spatial transformations. In specific, for an image x and its transformed version \hat{x} under the transformation parametrized by Θ , $E_{\Theta}(x)$ is given as:

$$E_{\Theta}(x) = \|F_{\Theta}(f_x) - f_{\hat{x}}\|_2^2$$

where $F_{\Theta}(f_x)$ is the generated feature maps by our TFT and $f_{\hat{x}}$ is the ground truth feature maps from the input image \hat{x} .

3.3. Aggregated Flow Transformations

Before delving into the intuition of our networks, we first talk about the flow patterns of CNN features under the three studied transformations.

To exploit and to better understand the rich feature extraction from discriminative CNNs, we study the flow transformations of feature maps, which implicitly encode both spatial and semantic information of the input images. Under three studied transformations of the input images, namely translation, rotation and scaling, we observe a clear but non-trivial deformation pattern of the corresponding CNN features. For instance, given a CNN (with 3 convolutional layers and 2 fully-connected layers) pre-trained on the MNIST dataset for image classification and one selected image from MNIST, we compare feature maps of the original image and that of the transformed image in all three convolutional layers under those transformations. Illustrated by Table 1, we observe that, from bottom to top, each convolutional layer

contains less and less spatial information about the input image. Furthermore, since convolution operations are shift-invariant, the resulting deformation on CNN features by applying translations to the input image resembles effect of simple translations. The feature transformations are much more non-linear locally, however, for rotation and scaling, while still having global flow patterns.

Table 1. Feature map deformations resulted from rotation, translation, and scaling of the input image. conv1_1 and conv1_2 are two channels of the first conv layer; conv2_1 and conv2_2 similarly.

	image	conv1_1	conv1_2	conv2_1	conv2_2	conv3_1
Original						
Rotation (-40°)						
Scaling (x0.8)						
Translation (up 8, left 8)						

To incorporate the top-down information into our networks, including translation, rotation and scaling parameters, we explicitly apply corresponding spatial transformation to the features maps. Before and after these operations, we perform aggregated residual transformations, similar to that of [42], in order to model the non-linearity of the resulting flow fields. In specific, we arrange both addition and subtraction paths for the transformation functions $\{\mathcal{T}_i(\cdot)\}$ to model a variety of deformation patterns. The convolution-based transformation functions $\mathcal{T}_i(\cdot)$ and the identity connection enforce the aggregated flow transformation function $T(\cdot)$ to only impose local feature operations of the input feature maps. The spatial transformation in between two aggregated feature transformation layer performs the global feature matching. The overall architecture ensures that its generated flow transformations are both clear and non-trivial. Our experiments demonstrates that the proposed TFT generalizes well to several datasets and is generic in terms of performing flow transformations on CNN features of the input images.

3.4. Generating New Images

Upon obtaining the transformed feature maps $\hat{f}_x = F_{\Theta}(f_x)$ for some transformation parametrized by Θ , we can generate images by “inverting” them in CNNs, similar to the process in [10], i.e., by minimizing the representation loss $E(x_{new}) = \|f_{x_{new}} - \hat{f}_x\|_2^2$ with respect to the generated image x_{new} , where $f_{x_{new}}$ is the CNN features of a given image x_{new} . Instead of picking up one convolutional layer, we use a set of layers altogether to generate images, with each layer transformed by a trained TFT. As our TFT can

perform flow transformations to feature maps of all the convolutional layers while remaining consistent across them, we can minimize the representation loss across all layers simultaneously. One resulting benefit is better quality of generated images, because different CNN layers extract and contain different kinds of information of the input images.

In practice, given a set of CNN layers $\{f_x(i)\}$, where i represents the i^{th} convolutional layer, we train the TFT for each one and generate images using a combination of the transformed feature maps $\{\hat{f}_x(i)\}$. We define the combined representation loss as:

$$E(x_{new}) = \sum_i \alpha_i \left\| f_{x_{new}}(i) - \hat{f}_x(i) \right\|_2^2 + \beta \cdot R_{TV}$$

where α_i and β are some constants and $f_x(i)$ represents the feature maps from the i^{th} CNN layer of the input image x . We also add a regularization term R_{TV} defined as:

$$R_{TV} = \sum_{i,j} (x_{i+1,j} - x_{i,j})^2 + (x_{i,j+1} - x_{i,j})^2$$

similar to that in [31].

4. Ablation Studies in Top-down Generator Design

Before the design of our proposed TFT, we come up with two other approaches, namely fully-connected model and gated affine model.

In ablation studies we gain some insights in the feature transformation modeling. In the experiment of the fully-connected model, we find that patterns in deformed features are learnable, yet a completely non-linear model does not explicitly capture the spatial information and thus fails to generalize to arbitrary top-down controls for those three kinds of transformations. In the experiment of the gated affine model, the networks explicitly utilize the top-down spatial information and we realize that, in a big picture, highly linear model like this can acceptably capture the underlying flow transformations of the feature maps, yet it fails to model the local feature transformations which are highly non-linear. The gated affine model, therefore, lacks model complexity. Our proposed TFT absorbs benefits from both models, enabling us to perform top-down feature transformation with both clear flow fields and adequate model complexity.

4.1. Fully-connected Model

The fully-connected model consists of a two layer fully connected net, namely a parameter net that transforms top-down control parameters, and a 4 layer fully-connected net called a generator network which takes the concatenation of the transformed control parameters and CNN features from

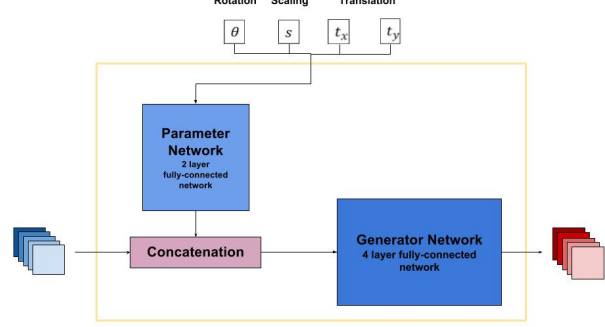


Figure 4. Architecture of the fully-connected model.

the source images and outputs the resulting feature maps. The architecture is illustrated in Figure 3.

We train these networks in the same setting as for TFT. Although this model achieves good numerical and visual results of feature transformations on the MNIST dataset, it has two major drawbacks. First, it lacks the capability of explaining how feature maps are transformed in the sense of deformation fields. Second, it does not generalize well to arbitrary transformation parameters Θ outside of those used in the training process. The TFT, on the other hand, solves these two problems elegantly.

4.2. Gated Affine Model

Our other model within the ablation study is designed based on the assumption that the transformation of feature maps resembles that of the input images, with the only difference that the convolutional layers are subject to some noise in this process. For instance, if an image is rotated by 30 degree, we assume the feature is also rotated, but not necessarily 30 degree. Accordingly, in our gated affine model we use gated transformation signals and a distance matrix to increase the model complexity and the local feature manipulation. This model performs well in translation and scaling transformations on the feature maps but suffers from under-fitting in that of rotation.

The gated affine model takes in the feature maps of an image x , namely f_x and the transformation parameters Θ for the image. It aims to perform a modified version of that transformation with parameters $\hat{\Theta}$ to the feature maps, channel by channel. In specific, firstly, for

$f_x \in \mathbb{R}^{n \times n \times m}$, we define $F \begin{pmatrix} \text{row} \\ \text{col} \\ 1 \end{pmatrix} = \text{row} \times n + \text{col}$

and $F^{-1}(i) = \begin{pmatrix} \lfloor \frac{i}{n} \rfloor \\ i \bmod n \\ 1 \end{pmatrix}$ as flatten and inverse flatten

function for pixel index. Secondly, we compute the estimated transformation parameters $\hat{\Theta}$ for the feature maps, which are gated with the actual transformation parameters Θ , formulated as $\hat{\Theta} = a \cdot \Theta + b$.

Next, we generate a distance matrix to measure the dis-

tance between each point in the original feature maps and the corresponding point in the estimated transformed feature maps. The matrix is defined as:

$$D_A(\hat{\Theta})_{i,j} = \text{dist}(F^{-1}(i), M_A(\hat{\Theta}) \cdot F^{-1}(j))$$

where $A \in \{tran, rot, scale\}$ is a affine operator among translation, rotation and scaling. In our model we use Euclidean metric as dist function.

After that, we use a scaled gated sigmoid mask which shifts and rescales a negative sigmoid function to control every element of D_A between 0 and 1:

$$\text{mask}(M)_{i,j} = \frac{1 + e^{-\mu}}{1 + e^{|\sigma M_{i,j}| - \mu}}$$

where σ, μ are gated with feature f_x .

Dot products with a properly masked distance matrix and flattened feature maps give an estimation of the transformed feature maps by performing local feature manipulations of the feature maps while applying affine transformation to them. σ and μ in the mask together define how a pixel in transformed feature will influence its neighbors.

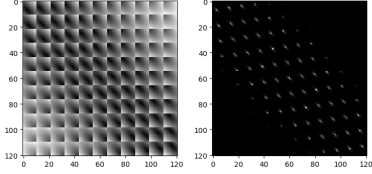


Figure 5. Left: distance matrix with $r = -40^\circ$. Right: masked distance matrix with $\sigma = 2, \mu = -100$ (median of all gated parameters)

The overall gated affine model can be written as

$$A(f_{x_k}, \Theta) = W \cdot \text{mask}(D_A(\hat{\Theta})) \cdot \text{flat}(f_{x_k})$$

for k ranging from 1 to m . The f_{x_k} is the k^{th} feature map of f_x , $\text{flat}(f_{x_k}) \in \mathbb{R}^{n^2}$ is the corresponding flattened feature map, $W \in \mathbb{R}^{n^2 \times n^2}$ is a linear regressor to learn and $A \in \{T, R, S\}$ is a type of affine transformation. Since this model can only learn each kind of affine transformation (rotation, translation or scaling) at a time, we combine three abovementioned generators as

$$GAM(f, \Theta) = T(R(S(f_x, \Theta), \Theta), \Theta)$$

We train this model in a similar setting to that of training TFT. This gated affine model shows the promise of clear flow fields for feature transformations, yet we find that it lacks model complexity for handling flow transformations that are more non-linear, such as that of rotation and, generally speaking, any spatial transformations in higher convolutional layers. The proposed TFT extends the idea from our gated affine model and solves this problem by increasing the model complexity without compromising the promise of clear feature flow fields.

5. Experiments

Table 2. Comparison of transformations on MNIST (within dataset) and notMNIST dataset. Rotations of 60° and 90° is beyond training settings (from -30° to 30°). Images generated by DVAM on notMNIST is vague and even lost its pattern when signals are out-of-bound. Images generated by our model show clear pattern of transformation.

	Original Image	Rotation (30°)	Rotation (60°)	Rotation (90°)	Scaling (x0.9)	Scaling (x1.1)	Translation	Combination
DVAM								
Ours								
DVAM								
Ours								

We train our top-down flow transformer (TFT) on the MNIST dataset and evaluate the learned generator on images from several datasets, including MNIST, notMNIST, Kimia-99 [1], MPEG-7 Shape [27], COIL-20 [33]. Under all three studied transformation, namely rotation, translation and scaling, we generate new images from transformed feature maps by applying TFT. We compare our results to those of [34] (DVAM) on MNIST and notMNIST dataset. We achieve better results both graphically and numerically. Our framework is able to perform out-of-bound transformations, i.e., transformations with arbitrary parameters, those that are out of the range of parameters used in the training process. Our model can generalize well to new datasets as the learned flow transformation is generic. Moreover, as TFT is not tied to fixed size of the input CNN features, it can perform flow transformations for arbitrary size of the input images, while existing method [34] cannot.

5.1. Training and Evaluation on MNIST

We resize each image of dimension 28 x 28 in the MNIST dataset to 44 x 44 by zero-padding the original images so that each image has enough space to perform translation and scaling. For data in the training set, we perform a combination of all three studied transformations to the input images and output the CNN feature maps across different convolutional layers. Specifically, for translation, we perform two dimensional shifting of the input images, with each axis ranging from +7 to -7 pixels, i.e., 225 combinations. For rotation, we perform with 13 different angle, namely rotate the input images by $5^\circ, 10^\circ, 15^\circ, 20^\circ, 25^\circ, 30^\circ$ clockwise and counterclockwise as well as 0° . For scaling, we choose three scaling factors: 0.9, 1.0 and 1.1 (1.0 indicates no scaling). These four parameters are the top-down transformation parameters in the training set. We form 3-tuples $(f_{x_{ori}}, f_{x_\Theta}, \Theta)$, where $f_{x_{ori}}$ is the feature

maps of the original images, f_{x_Θ} is the feature maps of the corresponding images after performing spatial transformation parameterized by Θ . Our training set is a collection of these 3-tuples. In practice, we apply the combinations of three transformations with randomly generated parameters from the range specified above.

Regarding choosing pre-trained CNN frameworks for obtaining feature maps used in the training process, we use a simplified version of DenseNet [21] pre-trained on MNIST. Specifically, the network only contains one dense block and each layer is a convolution followed by ReLU. We pick the number of initial filters and the growth rate as both 16 and the filter size as 5×5 with stride size 1. In our experiments, we train three TFT for the first three convolutional layers of the DenseNet, respectively, and use all these three layers to generate new images. We also use a classical CNN (3 convolutional layers, each is conv+relu+max-pooling, and 2 fully-connected layers) pre-trained on MNIST, and train 3 more TFT for the corresponding convolutional layer for better comparison of the resulting flow fields of TFT and that of the method in [34] (as they have the same feature map size).

The learning process of our TFT is supervised and the objective function is simply the Euclidean distance between the generated feature maps and the target ones, as mentioned in section 3.2. We train our networks by 200k steps with a L_2 regularizer of coefficient 0.0001. We use the ADAM optimizer [24] with learning rate of 0.0001. We set the batch size to be 128.

We also train the networks proposed in [34], namely DVAM in the same settings by using the same training set. We compare our graphical results of generated images from MNIST to those of DVAM in Table 2. When evaluated on MNIST, our approach outperforms theirs in terms of both feature flow representation and out-of-bound transformations.

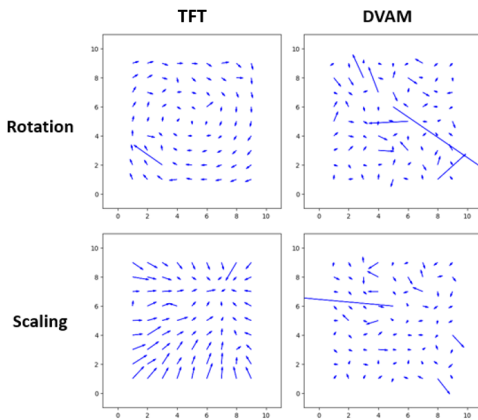


Figure 6. Feature flows of TFT and DVAM in [34]. Flow generated by our model has clear pattern while the flow generated by DVAM does not.

5.1.1 Learned Feature Flow

One critical advantage of our top-down flow transformer is its clear representation of learned feature flow. We ingest an identity matrix (a matrix with everything zero but only a pixel activated as 1) as input features for each channel and compute the center of mass of transformed features in each channel as the end point of the flow vector of this activated pixel. We perform this flow experiment on MNIST with rotation and scaling. As a comparison, we also apply the same method and transformations to draw feature flows in [34]. To better compare the results, we train the TFT based on feature maps obtained from a pre-trained CNN that is of similar width and height to the encoder-decoder used in [34]. The learned flow transformations in our model has clearer rotation and scaling deformation patterns than that of [34], as illustrated by Figure 6.

Furthermore, within the same CNN, the feature flows of different convolution layers showcases a coherent representation of the top-down information ingested into our model. However, it is also evident that the top-down information plays a smaller role in the feature transformation of the higher layers as feature flows become less coherent in higher layers, namely convolution layer 3 in the CNN experimented.

5.1.2 Out-of-bound Transformations

Another advantage of our model is its robustness to out-of-bound parameters. We train our model on MNIST with rotation signal ranging from -30° to 30° , yet we test the rotation transformation with 60° and 90° . In comparison, we also perform this experiment on [34]. We can see that our model succeeds on the out-of-bound transformations while the model in [34] fails as illustrated in 3rd and 4th columns of Table 2.

5.2. Evaluation on non-MNIST Datasets

5.2.1 Evaluation on notMNIST Dataset

We further test our model and the analogy network (DVAM) in [34] trained on the same MNIST dataset to investigate their inter dataset performance. In this experiment, as well as all the other inter dataset experiments, we use our trained TFT based on the feature maps from the simplified version of DenseNet, as discussed in section 5.1. To achieve similar visual and numerical effects as MNIST dataset we normalize the notMNIST dataset using max norm. The results demonstrate that our networks have learned flow transformations of the CNN features that explicitly utilize top-down information and generalize well to new types of data. Numerical data is in Table 4 and visual reconstructions of transformed features are illustrated in Table 2

Table 3. Comparison of transformations on images from a natural image.


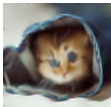
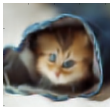
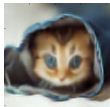
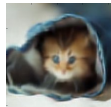
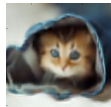





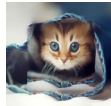
	Original	Rotation (30°)	Rotation (-30°)	Scaling (x1.3)	Scaling (x0.75)	Translation (up 30)
FlowPCA [8]						
Ours						

Table 4. Mean squared pixel prediction error according to different affine transformations of DVAM and our model on notMNIST

Model	translation	rotation	scaling	combination
DVAM	0.091315	0.077126	0.056829	0.062878
Ours	0.001384	0.004627	0.005122	0.006134

5.2.2 Evaluation on Kimia-99, MPEG-7, and COIL-20 Dataset

Our model can be easily extended to feature maps of images with different sizes while [34] fails to do so. Our model also has good performance in various other datasets, ranging from small patterns to real world images. We apply the TFT trained on MNIST to images from Kimia-99 [37, 1], MPEG-7 Shape [27] and COIL-20 [33] datasets. The results are illustrated in Table 5.

5.2.3 Evaluation on Natural Images






















































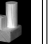


So far all our inter dataset experiments perform on grayscale images without background. We also apply our learned TFT for natural images. We use the same DenseNet-based network pre-trained on MNIST and the same TFT as in the previous experiments. Since the pre-trained network based on MNIST only accepts single channel images, we perform flow transformation to the colored images channel by channel. We compare our results with that in [8], illustrated by Table 3. We use a mask on the input CNN features when applying our method.

6. Conclusions and Discussion

We have developed top-down feature transformer (TFT) that learns a top-down generator by studying the internal transformations across CNN layers. The learned transformer is illustrated on both within and across datasets which demonstrates its clear advantage over those heavily learned through data-driven techniques. TFT points to a promising direction within the study of a CNN’s internal representation and top-down processes.

Potential application to ImageNet classification. In [8], online data augmentation was implemented and adopted

Table 5. Generated images by TFT from the Kimia [37], MPEG7 [27], and COIL datasets [33]. Images from different datasets are transformed using top-down transformers learned from the MNIST dataset. This is an inter dataset evaluation which demonstrates the effectiveness of our top-down model being generic and not tied to the training data.

Input	Rotation			Scaling		Translation	Compositional
	30°	60°	90°	0.9	1.1	varied	varied
							
							
							
							
							
							
							

to training AlexNet [25] on ImageNet [36] where features of the ImageNet images under small perturbations (translation, rotation, and scaling) were generated as augmented data. This online data augmentation procedure is different from the standard data augmentation process [25] in which augmented images are generated before hand. Online data augmentation can be efficient in speed and space. A small but visible improvement over the baseline AlexNet was observed in [8]. Similarly, we expect our TFT model to be able to help improve ImageNet classification as well.

Acknowledgments This work is supported by NSF IIS-1618477 and NSF IIS-1717431.

References

- [1] S. Belongie, J. Malik, and J. Puzicha. Shape matching and object recognition using shape contexts. *TPAMI*, 24(4):509–522, 2002.
- [2] A. Blake and A. Yuille. Active vision. 1993.
- [3] T. F. Cootes, G. J. Edwards, and C. J. Taylor. Active appearance models. *IEEE Transactions on pattern analysis and machine intelligence*, 23(6):681–685, 2001.
- [4] P. C. Dodwell. The lie transformation group model of visual perception. *Perception & Psychophysics*, 34(1):1–16, 1983.
- [5] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. van der Smagt, D. Cremers, and T. Brox. FlowNet: Learning optical flow with convolutional networks. In *ICCV*, 2015.
- [6] A. Dosovitskiy, J. T. Springenberg, and T. Brox. Learning to generate chairs with convolutional neural networks. In *CVPR*, 2015.
- [7] J. L. Elman. Distributed representations, simple recurrent networks, and grammatical. *Machine Learning*, 7:195–225, 1991.
- [8] P. W. Gallagher, S. Tang, and Z. Tu. What happened to my dog in that network: Unraveling top-down generators in convolutional neural networks. *arXiv preprint arXiv:1511.07125*, 2015.
- [9] J. R. Gardner, P. Upchurch, M. J. Kusner, Y. Li, K. Q. Weinberger, K. Bala, and J. E. Hopcroft. Deep manifold traversal: Changing labels with convolutional features. In *ECCV*, 2015.
- [10] L. A. Gatys, A. S. Ecker, and M. Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.
- [11] J. J. Gibson. A theory of direct visual perception. *Vision and Mind: selected readings in the philosophy of perception*, pages 77–90, 2002.
- [12] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *CVPR*, 2014.
- [13] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, 2014.
- [14] K. Gregor, I. Danihelka, A. Graves, D. Rezende, and D. Wierstra. Draw: A recurrent neural network for image generation. In *ICML*, 2015.
- [15] R. L. Gregory. *The intelligent eye*. Weidenfeld and Nicolson, 1980.
- [16] U. Grenander. *General pattern theory-A mathematical study of regular structures*. Clarendon Press, 1993.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [18] H. Hill and A. Johnston. The hollow-face illusion: Object-specific knowledge, general assumptions or properties of the stimulus? 36:199–223, 01 2007.
- [19] G. E. Hinton, S. Osindero, and Y. W. Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18:1527–1554, 2006.
- [20] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 1997.
- [21] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [22] M. Jaderberg, K. Simonyan, A. Zisserman, et al. Spatial transformer networks. In *NIPS*, 2015.
- [23] D. Kersten, P. Mamassian, and A. Yuille. Object perception as bayesian inference. *Annual Review of Psychology*, 55(1):271–304, 2004.
- [24] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [25] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [26] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum. Deep convolutional inverse graphics network. In *Advances in Neural Information Processing Systems*, pages 2539–2547, 2015.
- [27] L. J. Latecki, R. Lakamper, and T. Eckhardt. Shape descriptors for non-rigid shapes with a single closed contour. In *CVPR*, volume 1, pages 424–429, 2000.
- [28] C.-H. Lin and S. Lucey. Inverse compositional spatial transformer networks. *arXiv preprint arXiv:1612.03897*, 2016.
- [29] C. Liu, J. Yuen, and A. Torralba. Sift flow: Dense correspondence across scenes and its applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):978–994, 2011.
- [30] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. *CVPR*, 2015.
- [31] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *CVPR*, 2015.
- [32] D. Marr. *Vision: A Computational Investigation Into the Human Representation and Processing of Visual Information*. MIT Press, 2010.
- [33] S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (coil-20). Technical report, 1996.
- [34] S. E. Reed, Y. Zhang, Y. Zhang, and H. Lee. Deep visual analogy-making. In *NIPS*, 2015.
- [35] J. Ridley Stroop. Studies of interference in serial verbal reactions. 121:15–23, 03 1992.
- [36] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- [37] D. Sharvit, J. Chan, H. Tek, and B. B. Kimia. Symmetry-based indexing of image databases. In *Content-Based Access of Image and Video Libraries, 1998. Proceedings. IEEE Workshop on*, pages 56–62. IEEE, 1998.
- [38] Z. Tu, X. Chen, A. L. Yuille, and S. C. Zhu. Image parsing: unifying segmentation, detection, and recognition. In *ICCV*, 2003.
- [39] J. Wu, T. Xue, J. J. Lim, Y. Tian, J. B. Tenenbaum, A. Torralba, and W. T. Freeman. Single image 3d interpreter network. In *ECCV*, 2016.
- [40] W. Wu, M. Kan, X. Liu, Y. Yang, S. Shan, and X. Chen. Recursive spatial transformer (rest) for alignment-free face recognition. In *ICCV*, 2017.

- [41] Y. N. Wu, Z. Si, H. Gong, and S.-C. Zhu. Learning active basis model for object detection and recognition. *International journal of computer vision*, 90(2):198–235, 2010.
- [42] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5987–5995. IEEE, 2017.
- [43] S. Xie and Z. Tu. Holistically-nested edge detection. In *ICCV*, 2015.
- [44] A. L. Yuille, P. W. Hallinan, and D. S. Cohen. Feature extraction from faces using deformable templates. *International journal of computer vision*, 8(2):99–111, 1992.
- [45] Q. Zhang, R. Cao, Y. N. Wu, and S.-C. Zhu. Growing interpretable part graphs on convnets via multi-shot learning. In *AAAI*, 2017.
- [46] S.-C. Zhu, D. Mumford, et al. A stochastic grammar of images. *Foundations and Trends® in Computer Graphics and Vision*, 2(4):259–362, 2007.