Learning to Inpaint for Image Compression

Mohammad Haris Baig*

Department of Computer Science Dartmouth College Hanover, NH

Vladlen Koltun Intel Labs Santa Clara, CA

Lorenzo Torresani Dartmouth College Hanover, NH

Abstract

We study the design of deep architectures for lossy image compression. We present two architectural recipes in the context of multi-stage progressive encoders and empirically demonstrate their importance on compression performance. Specifically, we show that: (a) predicting the original image data from residuals in a multi-stage progressive architecture facilitates learning and leads to improved performance at approximating the original content and (b) learning to inpaint (from neighboring image pixels) before performing compression reduces the amount of information that must be stored to achieve a high-quality approximation. Incorporating these design choices in a baseline progressive encoder yields an average reduction of over 60% in file size with similar quality compared to the original residual encoder.

1 Introduction

Visual data constitutes most of the total information created and shared on the Web every day and it forms a bulk of the demand for storage and network bandwidth [14]. It is customary to compress image data as much as possible as long as there is no perceptible loss in content. In recent years deep learning has made it possible to design deep models for learning compact representations for image data [2, 17, 19, 20, 21]. Deep learning based approaches, such as the work of Rippel and Bourdev [17], significantly outperform traditional methods of lossy image compression. In this paper, we show how to improve the performance of deep models trained for lossy image compression.

We focus on the design of models that produce progressive codes. Progressive codes are a sequence of representations that can be transmitted to improve the quality of an existing estimate (from a previously sent code) by adding missing detail. This is in contrast to non-progressive codes whereby the entire data for a certain quality approximation must be transmitted before the image can be viewed. Progressive codes improve the user's browsing experience by reducing loading time of pages that are rich in images. Our main contributions in this paper are two-fold.

- 1. While traditional progressive encoders are optimized to compress residual errors in each stage of their architecture (residual-in, residual-out), instead we propose a model that is trained to predict at each stage the original image data from the residual of the previous stage (residual-in, image-out). We demonstrate that this leads to an easier optimization resulting in better image compression. The resulting architecture reduces the amount of information that must be stored for reproducing images at similar quality by 18% compared to a traditional residual encoder.
- 2. Existing deep architectures do not exploit the high degree of spatial coherence exhibited by neighboring patches. We show how to design and train a model that can exploit dependences between adjacent regions by learning to inpaint from the available content. We introduce multi-scale convolutions that sample content at multiple scales to assist with inpainting.

^{*}www.cs.dartmouth.edu/~haris/compression.html

We jointly train our proposed inpainting and compression models and show that inpainting reduces the amount of information that must be stored by an additional 42%.

2 Approach

We begin by reviewing the architecture and the learning objective of a progressive multi-stage encoder-decoder with S stages. We adopt the convolutional-deconvolutional residual encoder proposed by Toderici et al. [20] as our reference model. The model extracts a compact binary representation B from an image patch P. This binary representation, used to reconstruct an approximation of the original patch, consists of the sequence of representations extracted by the S stages of the model, $B = [B_1, B_2, \dots B_S]$.

The first stage of the model extracts a binary code B_1 from the input patch P. Each of the subsequent stages learns to extract representations B_s , to model the compressions residuals R_{s-1} from the previous stage. The compression residuals R_s are defined as $R_s = R_{s-1} - \mathcal{M}_s(R_{s-1}|\Theta_s)$, where $\mathcal{M}_s(R_{s-1}|\Theta_s)$ represents the reconstruction obtained by the stage s when modelling the residuals R_{s-1} . The model at each stage is split into an encoder $B_s = \mathcal{E}_s(R_{s-1}|\Theta_s^E)$ and a decoder $\mathcal{D}_s(B_s|\Theta_s^D)$ such that $\mathcal{M}_s(R_{s-1}|\Theta_s) = \mathcal{D}_s(\mathcal{E}_s(R_{s-1}|\Theta_s^E)|\Theta_s^D)$ and $\Theta_s = \{\Theta_s^E, \Theta_s^D\}$. The parameters for the s^{th} stage of the model are denoted by Θ_s . The residual encoder-decoder is trained on a dataset \mathcal{P} , consisting of N image patches, according to the following objective:

$$\hat{\mathcal{L}}(\mathcal{P};\Theta_{1:S}) = \sum_{s=1}^{S} \sum_{i=1}^{N} \|R_{s-1}^{(i)} - \mathcal{M}_s(R_{s-1}^{(i)}|\Theta_s)\|_{2.}^2$$
(1)

 $R_s^{(i)}$ represents the compression residual for the i^{th} patch $P^{(i)}$ after stage s and $R_0^{(i)} = P^{(i)}$.

Residual encoders are difficult to optimize as gradients have to traverse long paths from later stages to affect change in the previous stages. When moving along longer paths, gradients tend to decrease in magnitude as they get to earlier stages. We address this shortcoming of residual encoders by studying a class of architectures we refer to as "Residual-to-Image" (R2I).

2.1 Residual-to-Image (R2I)

To address the issue of vanishing gradients we add connections between subsequent stages and restate the loss to predict the original data at the end of each stage, thus performing *residual-to-image* prediction. This leads to the new objective shown below:

$$\mathcal{L}(\mathcal{P};\Theta_{1:S}) = \sum_{s=1}^{S} \sum_{i=1}^{N} \|P^{(i)} - \mathcal{M}_s(R_{s-1}^{(i)}|\Theta_s)\|_{2.}^2$$
 (2)

Stage s of this model takes as input the compression residuals R_{s-1} computed with respect to the original data, $R_{s-1} = P - \mathcal{M}_{s-1}(R_{s-2}|\Theta_{s-1})$, and $\mathcal{M}_{s-1}(R_{s-2}|\Theta_{s-1})$ now approximates the reconstruction of the original data P at stage s-1. To allow complete image reconstructions to be produced at each stage while only feeding in residuals, we introduce connections between the layers of adjacent stages. These connections allow for later stages to incorporate information that has been recovered by earlier stages into their estimate of the original image data. Consequently, these connections (between subsequent stages) allow for better optimization of the model.

In addition to assisting with modeling the original image, these connections play two key roles. Firstly, these connections create residual blocks [10] which encourage explicit learning of how to reproduce information which could not be generated by the previous stage. Secondly, these connections reduce the length of the path along which information has to travel from later stages to impact the earlier stages, leading to a better joint optimization.

This leads us to the question of where should such connections be introduced and how should information be propagated? We consider two types of connections to propagate information between successive stages. 1) *Prediction connections* are analogous to the identity shortcuts introduced by He et al. [10] for residual learning. They act as parameter-free additive connections: the output of

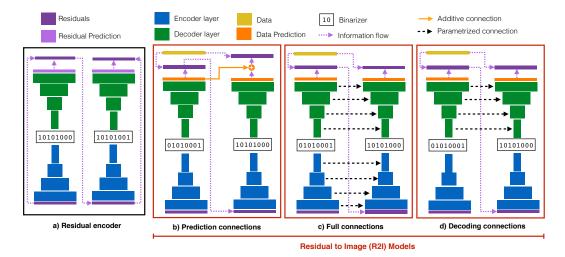


Figure 1: Multiple approaches for introducing connections between successive stages. These designs for progressive architectures allow for varying degrees of information to be shared. Architecture (b-d) do not reconstruct residuals, but the original data at every stage. We call these architectures "residual-to-image" (R2I).

each stage is produced by simply adding together the residual predictions of the current stage and all preceding stages (see Figure 1(b)) before applying a final non-linearity.2) Parametric connections are referred to as projection shortcuts by He et al. [10]. Here we use them to connect corresponding layers in two consecutive stages of the compression model. The features of each layer from the previous stage are convolved with learned filters before being added to the features of the same layer in the current stage. A non-linearity is then applied on top. The prediction connections only yield the benefit of creating residual blocks, albeit very large and thus difficult to optimize. In contrast, parametric connections allow for the intermediate representations from previous stages to be passed to the subsequent stages. They also create a denser connectivity pattern with gradients now moving along corresponding layers in adjacent stages. We consider two variants of parametric connections: "full" which use parametric connections between all the layers in two successive stages (see Figure 1(c)), and "decoding" connections which link only corresponding decoding layers (i.e., there are no connections between encoding layers of adjacent stages). We note that the LSTM-based model of Toderici et al. [21] represents a particular instance of R2I network with full connections. In Section 3 we demonstrate that R2I models with decoding connections outperform those with full connections and provide an intuitive explanation for this result.

2.2 Inpainting Network

Image compression architectures learn to encode and decode an image patch-by-patch. Encoding all patches independently assumes that the regions contain truly independent content. This assumption generally does not hold true when the patches being encoded are contiguous. We observe that the content of adjacent image patches is not independent. We propose a new module for the compression model designed to exploit the spatial coherence between neighboring patches. We achieve this goal by training a model with the objective of predicting the content of each patch from information available in the neighboring regions.

Deep models for inpainting, such as the one proposed by Pathak et al. [15], are trained to predict the values of pixels in the region \hat{W} from a context region \hat{C} (as shown in Figure 2). As there is data present all around the region to be inpainted this imposes strong constraints on what the inpainted region should look like. We consider the scenario where images are encoded and decoded block-by-block moving from left to right and going from top to bottom (similar to how traditional codecs process images [1, 22]). Now, at decoding time only content above and to the left of each patch will have been reconstructed (see Figure 2(a)). This gives rise to the problem of "partial-context inpainting". We propose a model that, given an input region C, attempts to predict the content of the current patch P. We denote by $\hat{\mathcal{P}}$ the dataset which contains all the patches from the dataset \mathcal{P}

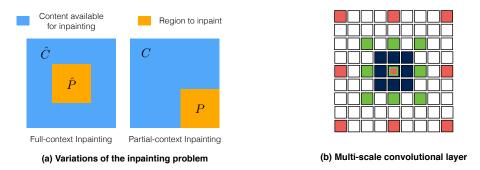


Figure 2: (a) The two kinds of inpainting problems. (b) A multi-scale convolutional layer with 3 dilation factors. The colored boxes represent pixels from which the content is sampled.

and the respective context regions C for each patch. The loss function used to train our inpainting network is:

$$\mathcal{L}_{inp}(\hat{\mathcal{P}};\Theta_I) = \sum_{i=1}^N \|P^{(i)} - \mathcal{M}_I(C^{(i)}|\Theta_I)\|_{2.}^2$$
(3)

The output of the inpainting network is denoted by $\mathcal{M}_I(C^{(i)}|\Theta_I)$, where Θ_I refers to the parameters of the inpainting network.

2.2.1 Architecture of the Partial-Context Inpainting Network

Our inpainting network has a feed-forward architecture which propagates information from the context region C to the region being inpainted, P. To improve the ability of our model at predicting content, we use a multi-scale convolutional layer as the basic building block of our inpainting network. We make use of the dilated convolutions described by Yu and Koltun [24] to allow for sampling at various scales. Each multi-scale convolutional layer is composed of k filters for each dilation factor being considered. Varying the dilation factor of filters gives us the ability to analyze content at various scales. This structure of filters provides two benefits. First, it allows for a substantially denser and more diverse sampling of data from context and second it allows for better propagation of content at different spatial scales. A similarly designed layer was also used by Chen et al. [5] for sampling content at multiple scales for semantic segmentation. Figure 2(b) shows the structure of a multi-scale convolutional layer.

The multi-scale convolutional layer also gives us the freedom to propagate content at full resolution (no striding or pooling) as only a few multi-scale layers suffice to cover the entire region. This allows us to train a relatively shallow yet highly expressive architecture which can propagate fine-grained information that might otherwise be lost due to sub-sampling. This light-weight and efficient design is needed to allow for joint training with a multi-stage compression model.

2.2.2 Connecting the Inpainting Network with the R2I Compression model

Next, we describe how to use the prediction of the inpainting network for assisting with compression. Whereas the inpainting network learns to predict the data as accurately as possible, we note that this is not sufficient to achieve good performance on compression, where it is also necessary that the "inpainting residuals" be easy to compress. We describe the inpainting residuals as $R_0 = P - \mathcal{M}_I(C|\Theta_I)$, where $\mathcal{M}_I(C|\Theta_I)$ denotes the inpainting estimate. As we wanted to train our model to always predict the data, we add the inpainting estimate to the final prediction of each stage of our compression model. This allows us to (a) produce the original content at each stage and (b) to

discover an inpainting that is beneficial for all stages of the model because of joint training. We now train our complete model as

$$\mathcal{L}_{C}(\hat{\mathcal{P}};\Theta_{I},\Theta_{1:S}) = \mathcal{L}_{inp}(\hat{\mathcal{P}};\Theta_{I}) + \sum_{i=1}^{N} \sum_{s=1}^{S} \|P^{(i)} - [\mathcal{M}_{s}(R_{s-1}^{(i)}|\Theta_{s}) + \mathcal{M}_{I}(C^{(i)}|\Theta_{I})]\|_{2}^{2}$$
(4)

In this new objective \mathcal{L}_C , the first term \mathcal{L}_{inp} corresponds to the original inpainting loss, $R_0^{(i)}$ corresponds to the inpainting residual for example i. We note that each stage of this inpainting-based progressive coder directly affects what is learned by the inpainting network. We refer to the model trained with this joint objective as "Inpainting for Residual-to-Image Compression" (IR2I).

Whereas we train our model to perform inpainting from the original image content, we use a lossy approximation of the context region C when encoding images with IR2I. This is done because at decoding time our model does not have access to the original image data. We use the approximation from stage 2 of our model for performing inpainting at encoding and decoding time, and transmit the binary codes for the first two stages as a larger first code. This strategy allows us to leverage inpainting while performing progressive image compression.

2.3 Implementation Details

Our models were trained on 6,507 images from the ImageNet dataset [7], as proposed by Ballé et al. [2] to train their single-stage encoder-decoder architectures. A full description of the R2I models and the inpainting network is provided in the supplementary material. We use the Caffe library [12] to train our models. The residual encoder and R2I models were trained for 60,000 iterations whereas the joint inpainting network was trained for 110,000 iterations. We used the Adam optimizer [13] for training our models and the MSRA initialization [9] for initializing all stages. We used initial learning rates of 0.001 and the learning rate was dropped after 30K and 45K for the R2I models. For the IR2I model, the learning rate was dropped after 30K, 65K, and 90K iterations by a factor of 10 each time. All of our models were trained to reproduce the content of 32×32 image patches. Each of our models has 8 stages, with each stage contributing 0.125 bits-per-pixel (bpp) to the total representation of a patch. Our models handle binary optimization by employing the biased estimators approach proposed by Raiko et al. [16] as was done by Toderici et al. [20, 21].

Our inpainting network has 8 multi-scale convolutional layers for content propagation and one standard convolutional layer for performing the final prediction. Each multi-scale convolutional layer consists of 24 filters each for dilation factors 1, 2, 4, 8. Our inpainting network takes as input a context region C of size 64×64 , where the bottom right 32×32 region is zeroed out and represents the region to be inpainted.

3 Results

We investigate the improvement brought about by the presented techniques. We are interested in studying the reduction in bit-rate, for varying quality of reconstruction, achieved after adaptation from the residual encoder proposed by Toderici et al. [20]. To evaluate performance, we perform compression with our models on images from the Kodak dataset [8]. The dataset consists of 24 uncompressed color images of size 512×768 . The quality is measured according to the MS-SSIM [23] metric (higher values indicate better quality). We use the Bjontegaard-Delta metric [4] to compute the average reduction in bit-rate across all quality settings.

3.1 R2I - Design and Performance

The table in Figure 3(a) shows the percentage reduction in bit-rate achieved by the three variations of the Residual-to-Image models. As can be seen, adding side-connections and training for the more desirable objective (i.e., approximating the original data) at each stage helps each of our models. That said, having connections in the decoder only helps more compared to using a "full" connection approach or only sharing the final prediction.

	Rate Savings (%)				
Approach	SSIM	MS-SSIM			
R2I Prediction	4.483	5.177			
R2I Full	10.015	7.652			
R2I Decoding	20.002	17.951			
R2I Decoding	20.002	17.9			

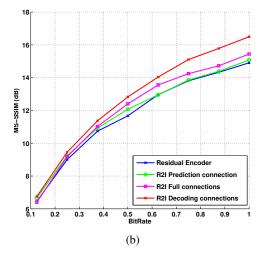


Figure 3: (a) Average rate savings for each of the three R2I variants compared to the residual encoder proposed by Toderici et al. [20]. (b) Figure shows the quality of images produced by each of the three R2I variants across a range of bit-rates.

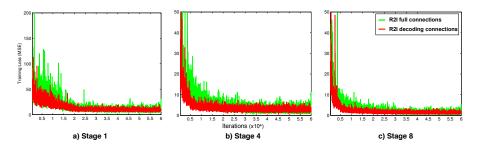


Figure 4: The R2I training loss from 3 different stages (start, middle, end) viewed as a function of iterations for the "full" and the "decoding" connections models. We note that the decoding connections model converges faster, to a lower value, and shows less variance.

The model which shares only the prediction between stages performs poorly in comparison to the other two designs as it does not allow for features from earlier stages to be altered as efficiently as done by the full or decoding connections architectures.

The model with decoding connections does better than the architecture with full connections because for the model with connections at decoding only the binarization layer in each stage extracts a representation from the relevant information only (the residuals with respect to the data). In contrast, when connections are established in both the encoder and the decoder, the binary representation may include information that has been captured by a previous stage, thereby adding burden on each stage in identifying information pertinent to improving reconstruction, leading to a tougher optimization. Figure 4 shows that the model with full connections struggles to minimize the training error compared to the model with decoding connections. This difference in training error points to the fact that connections in the encoder make it harder for the model to do well at training time. This difficulty of optimization amplifies with the increase in stages as can be seen by the difference between the full and decoding architecture performance (shown in Figure 3(b)) because the residuals become harder to compress.

We note that R2I models significantly improve the quality of reconstruction at higher bit rates but do not improve the estimates at lower bit-rates as much (see Figure 3(b)). This tells us that the overall performance can be improved by focusing on approaches that yield a significant improvement at lower bit-rates, such as inpainting, which is analyzed next.

	Rate Savings (%)		
Approach	SSIM	MS-SSIM	
R2I Decoding	20.002	17.951	
R2I Decoding Sep-Inp	27.379	27.794	
R2I Decoding Joint-Inp	63.353	60.446	

⁽a) Impact of inpainting on the performance at compression. All bit-rate savings are reported with respect to the residual encoder by Toderici et al. [20]

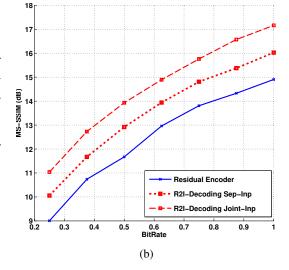


Figure 5: (a) Average rate savings with varying forms of inpainting. (b) The quality of images with each of our proposed approaches at varying bit-rates.

3.2 Impact of Inpainting

We begin analyzing the performance of the inpainting network and other approaches on partial-context inpainting. We compare the performance of the inpainting network with both traditional approaches as well as a learning-based baseline. Table 1 shows the average SSIM achieved by each approach for inpainting all non-overlapping patches in the Kodak dataset.

Approach	PDE-based	Exemplar-based	Learn	ing-based
	[3]	[6]	Vanilla network	Inpainting network
SSIM	0.4574	0.4611	0.4545	0.5165

Table 1: Average SSIM for partial-context inpainting on the Kodak dataset [8]. The vanilla model is a feed-forward CNN with no multi-scale convolutions.

The vanilla network corresponds to a 32-layer (4 times as deep as the inpainting network) model that does not use multi-scale convolutions (all filters have a dilation factor of 1), has the same number of parameters, and also operates at full resolution (as our inpainting network). This points to the fact that the improvement in performance of the inpainting network over the vanilla model is a consequence of using multi-scale convolutions. The inpainting network improves over traditional approaches because our model learns the best strategy for propagating content as opposed to using hand-engineered principles of content propagation. The low performance of the vanilla network shows that learning by itself is not superior to traditional approaches and multi-scale convolutions play a key role in achieving better performance.

Whereas inpainting provides an initial estimate of the content within the region it by no means generates a perfect reconstruction. This leads us to the question of whether this initial estimate is better than not having an estimate? The table in Figure 5(a) shows the performance on the compression task with and without inpainting. These results show that the greatest reduction in file size is achieved when the inpainting network is jointly trained with the R2I model. We note (from Figure 5(b)) that inpainting greatly improves the quality of results obtained at lower and at higher bit rates.

The baseline where the inpainting network is trained separately from the compression network is presented here to emphasize the role of joint training. Traditional codecs [1] use simple non learning-based inpainting approaches and their predefined methods of representing data are unable to compactly encode the inpainting residuals. Learning to inpaint separately improves the performance

as the inpainted estimate is better than not having any estimate. But given that the compression model has not been trained to optimize the compression residuals the reduction in bit-rate for achieving high quality levels is low. We show that with joint training, we can not only train a model that does better inpainting but also ensure that the inpainting residuals can be represented compactly.

3.3 Comparison with Existing Approaches

Table 2 shows a comparison of the performance of various approaches compared to JPEG [22] in the 0.125 to 1 bits-per-pixel (bpp) range. We select this range as images from our models towards the end of this range show no perceptible artifacts of compression.

The first part of the table evaluates the performance of learning-based progressive approaches. We note that our proposed model outperforms the multi-stage residual encoder proposed by Toderici et al. [20] (trained on the same $6.5 \, \mathrm{K}$ dataset) by 17.9 % and IR2I outperforms the residual encoder by reducing file-sizes by 60.4 %. The residual-GRU, while similar in architecture to our "full" connections model, does not do better even when trained on a dataset that is $1000 \, \mathrm{times}$ bigger and for $10 \, \mathrm{times}$ more training time. The results shown here do not make use of entropy coding as the goal of this work is to study how to improve the performance of deep networks for progressive image compression and entropy coding makes it harder to understand where the performance improvements are coming from. As various approaches use different entropy coding methods, this further obfuscates the source of the improvements.

The second part of the table shows the performance of existing codecs. Existing codecs use entropy coding and rate-distortion optimization. We note that even without using either of these powerful post processing techniques, our final "IR2I" model is competitive with traditional methods for compression, which use both of these techniques. A comparison with recent non-progressive approaches [2, 19], which also use these post-processing techniques for image compression, is provided in the supplementary material.

Approach	Number of Training Images	Progressive	Rate Savings (%)
Residual Encoder [20]	6.5K	Yes	2.56
Residual-GRU [21]	6 M	Yes	33.26
R2I (Decoding connections)	6.5K	Yes	18.53
IR2I	6.5K	Yes	51.25
JPEG-2000 [18]	N/A	No	63.01
WebP [1]	N/A	No	64.98

Table 2: Average rate savings compared to JPEG [22]. The savings are computed on the Kodak [8] dataset with rate-distortion profiles measuring MS-SSIM in the 0-1 bpp range.

We observe that a naive implementation of IR2I creates a linear dependence in content (as all regions used as context have to be decoded before being used for inpainting) and thus may be substantially slower. In practice, this slowdown would be negligible as one can use a diagonal scan pattern (similar to traditional codecs) for ensuring high parallelism thereby reducing run times. Furthermore, we perform inpainting using predictions from the first step only. Therefore, the dependence only exists when generating the first progressive code. For all subsequent stages, there is no dependence in content, and our approach is comparable in run time to similar approaches.

4 Conclusion and Future Work

We study a class of "Residual to Image" models and show that within this class, architectures which have decoding connections perform better at approximating image data compared to designs with other forms of connectivity. We observe that our R2I decoding connections model struggles at low bit-rates and we show how to exploit spatial coherence between content of adjacent patches via inpainting to improve performance at approximating image content at low bit-rates. We design a new

model for partial-context inpainting using multi-scale convolutions and show that the best way to leverage inpainting is by jointly training the inpainting network with our R2I Decoding model.

One interesting extension of this work would be to incorporate entropy coding within our progressive compression framework to train models that produce binary codes which have low-entropy and can be represented even more compactly. Another possible direction would be to extend our proposed framework to video data, where the gains from our discovery of recipes for improving compression may be even greater.

5 Acknowledgements

This work was funded in part by Intel Labs and NSF award CNS-120552. We gratefully acknowledge NVIDIA and Facebook for the donation of GPUs used for portions of this work. We would like to thank George Toderici, Nick Johnston, Johannes Balle for providing us with information needed for accurate assessment. We are grateful to members of the Visual Computing Lab at Intel Labs, and members of the Visual Learning Group at Dartmouth College for their feedback.

References

- [1] WebP a new image format for the web. https://developers.google.com/speed/webp/. Accessed: 2017-04-29.
- [2] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. In ICLR, 2017.
- [3] Marcelo Bertalmio, Guillermo Sapiro, Vincent Caselles, and Coloma Ballester. Image inpainting. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 417–424. ACM Press/Addison-Wesley Publishing Co., 2000.
- [4] Gisle Bjontegaard. Improvements of the bd-psnr model. In ITU-T SC16/Q6, 35th VCEG Meeting, Berlin, Germany, July 2008, 2008.
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv* preprint arXiv:1606.00915, 2016.
- [6] Antonio Criminisi, Patrick Pérez, and Kentaro Toyama. Region filling and object removal by exemplar-based image inpainting. *IEEE Transactions on image processing*, 13(9):1200–1212, 2004.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition*, 2009. CVPR 2009. IEEE Conference on, pages 248–255. IEEE, 2009.
- [8] Eastman Kodak Company. Kodak lossless true color image suite, 1999. http://r0k.us/graphics/kodak/.
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [11] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 448–456, 2015.
- [12] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014.
- [13] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In ICLR, 2014.
- [14] Mary Meeker. Internet Trends Report 2016. KPCB, 2016.

- [15] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016.
- [16] Tapani Raiko, Mathias Berglund, Guillaume Alain, and Laurent Dinh. Techniques for learning binary stochastic feedforward neural networks. In *ICLR*, 2015.
- [17] Oren Rippel and Lubomir Bourdev. Real-time adaptive image compression. In *International Conference on Machine Learning*, 2017.
- [18] Athanassios Skodras, Charilaos Christopoulos, and Touradj Ebrahimi. The jpeg 2000 still image compression standard. *IEEE Signal processing magazine*, 18(5):36–58, 2001.
- [19] L. Theis, W. Shi, A. Cunningham, and F. Huszár. Lossy image compression with compressive autoencoders. In *ICLR*, 2017.
- [20] George Toderici, Sean M. O'Malley, Sung Jin Hwang, Damien Vincent, David Minnen, Shumeet Baluja, Michele Covell, and Rahul Sukthankar. Variable rate image compression with recurrent neural networks. In *ICLR*, 2016.
- [21] George Toderici, Damien Vincent, Nick Johnston, Sung Jin Hwang, David Minnen, Joel Shor, and Michele Covell. Full resolution image compression with recurrent neural networks. *arXiv* preprint arXiv:1608.05148, 2016.
- [22] Gregory K. Wallace. The JPEG still picture compression standard. *Communications of the ACM*, 34(4), 1991.
- [23] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *Signals, Systems and Computers*, 2004. Conference Record of the Thirty-Seventh Asilomar Conference on, volume 2, pages 1398–1402. IEEE, 2003.
- [24] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *ICLR*, 2016.

A Progressive vs Non-Progressive Coding Approaches

Table 3 shows a comparison of the performance of various approaches compared to JPEG [22] in the 0.125 to 1 bits-per-pixel (bpp) range. We compare our progressive model with non-progressive approaches.

Balle et al. [2] proposed an end-to-end optimization framework which required them to train 6 separate models, one for each target quality of images they desire. Theis et al. [19] trained three models for generating codes for images at low, high, and medium quality. These approaches perform entropy coding on the hidden representation extracted by their model. Both approaches explicitly optimize for the discovery of hidden-representations that can be compactly represented in binary form using entropy coding.

In contrast, we train a single model that can generate images of increasing quality using progressive transmission. In the future we aim to investigate the introduction of entropy estimation within the training objective of progressive models as done by Balle et al [2] and Theis et al [19] for non-progressive approaches.

Approach	Number of Training Images	Progressive	Rate Savings (%)	Entropy Coding
Decoder-only PIC IR2I	6.5K 6.5K	Yes Yes	$18.53 \\ 51.25$	No No
Balle [2]	6.5K	No	142.87 39.35	Yes
Theis [19]	0.7K	No		Yes
JPEG-2000 [18]	N/A	No	63.01	Yes
WebP [1]	N/A	No	64.98	Yes

Table 3: Average % savings of bit-rates for various approaches compared to JPEG [22]. The savings are computed from rate-distortion profiles with the MS-SSIM metric for images in the 0-1 bpp range on the Kodak [8] dataset.

B Qualitative Analysis

In the following sections we provide a qualitative analysis of different variants of our R2I architecture and also compare compression results with and without inpainting. We then show a quantitative comparison of our R2I model with non-progressive approaches for compression. Finally, we describe in detail the architectures of the different models proposed.

Table 4 shows a comparison of the images generated by the three variants of our R2I model and the original residual encoder. We invite the reader to observe that at higher bit-rates (closer to 1 bpp) compression artifacts are barely noticeable. When reaching this high bit-rates the model R2I with decoding connections almost always outperforms all other models considered in this evaluation.

Table 5 provides a visual comparison of the images generated by the R2I model using decoding connections both with and without inpainting. We consider both approaches to inpainting discussed in the paper i.e. separately trained and jointly trained. The table also includes a comparison with the original residual encoder proposed by Toderici et al. [20].

C Architecture Specification

First we describe some notation common to all models. All convolutional layers use filters of size 3×3 except for Conv* which uses 1×1 size filters. All of our deconvolutional layers use filters of size 2×2 with group size set to be equal to the number of channels in the input. Elt-Sub stands for an element-wise subtraction operation and Elt-Add stands for an element-wise addition operation. We train all architectures with patches of input size $32\times 32\times 3$. BN denotes batch-normalization as described by Ioffe and Szegedy [11]

C.1 Residual Encoder

Table 6 describes our adaptation of the residual encoder proposed by Toderici et al. [20]. The loss is applied at the end of each stage to the output of the layer conv_pred_s and is computed with respect to residual_(s-1). The model has 24.2 million parameters.

C.2 R2I Models

We now describe the architectures for our R2I models. The highlighted part of each table shows the salient changes from the residual encoder adapted from Toderici et al. [20] described in subsection C.1. For each of the models, for stage s=1 we add zeros for the connections coming in as there is no earlier stage.

Table 7 describes our R2I model with prediction connections. The loss is applied at the end of each stage to the output of the layer output_s and is computed with respect to data. The model has 24.2 million parameters.

Table 8 describes our R2I model with "full" connections. The loss is applied at the end of each stage to the output of the layer output_s and is computed with respect to data. The model has 49.8 million parameters.

Table 9 describes our R2I model with "decoding" connections. The loss is applied at the end of each stage to the output of the layer conv_pred_s and is computed with respect to data. The model has 29.2 million parameters.

C.3 Inpainting Network

Table 10 describes our Inpainting Network. The model takes data of size 64×64 as input with the bottom right 32×32 zeroed out. The inpainting loss is applied to pred and conv_RGB_crop(before application of the non-linearity) is shared with successive stages of the model as each stage applies a non-linearity to the final output before applying it's loss.

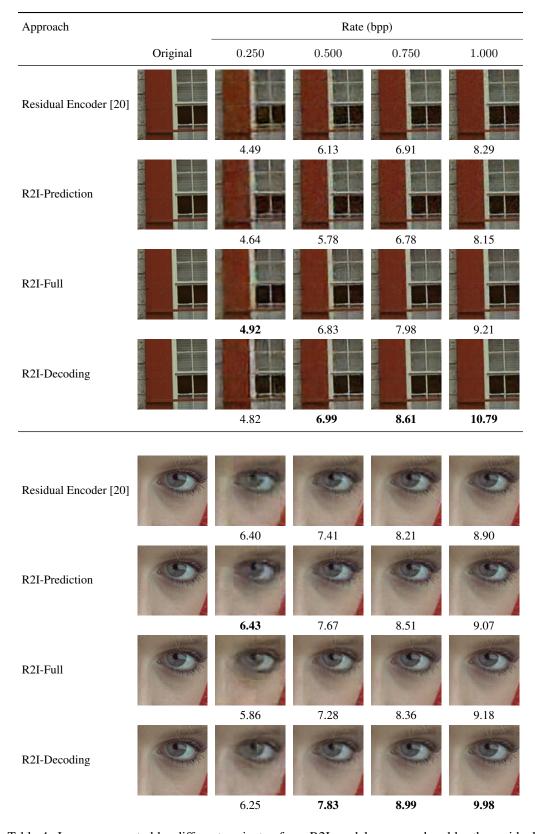


Table 4: Images generated by different variants of our R2I models proposed and by the residual encoder [20] at various bit-rates. The numbers beneath each row of images represent the MS-SSIM [23] (in decibals(dB)) for each of the images shown. For the first example, note how even at the highest bit-rate only R2I decoding connections is able to capture and render with good detail the blinds in the window.

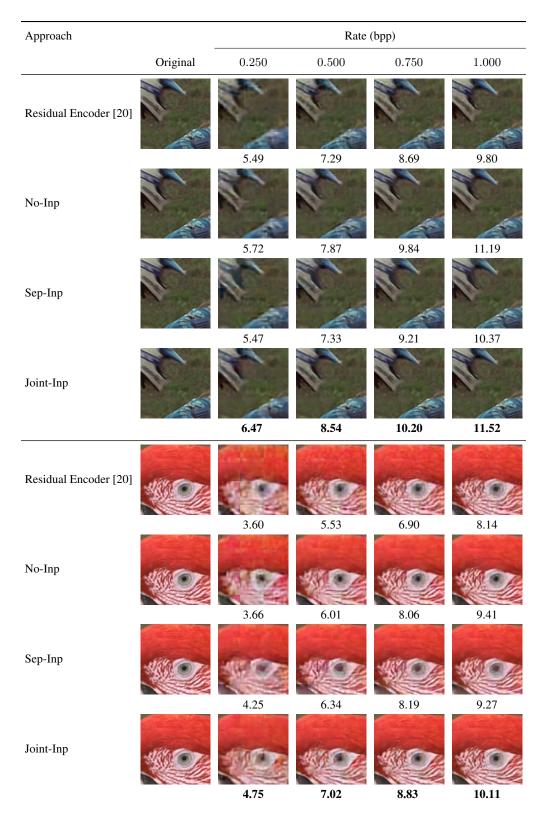


Table 5: Images generated from the R2I decoding connections model with and without inpainting and by the residual encoder [20] at various bit-rates. The numbers beneath each row of images represent the MS-SSIM [23] (in decibals(dB)) for each of the images shown. Note, how for the first example the joint inpainting model captures the blue stripe on the riders helmet much better than other models. We also observe how the detail surrounding the eye of the parrot is much better displayed in the joint-inp model.

Layer	Type	Stride	# Filters	BN	Non-Linearity	Input	Output
1	Conv	1	64	Y	ReLU	residual_(s-1)	conv_1_s
2	Conv	2	128	Y	ReLU	conv_1_s	conv_2_s
3	Conv	1	128	Y	ReLU	conv_2_s	conv_3_s
4	Conv	2	256	Y	ReLU	conv_3_s	conv_4_s
5	Conv	1	256	Y	ReLU	conv_4_s	conv_5_s
6	Conv	2	256	Y	ReLU	conv_5_s	conv_6_s
7	Conv*	1	8	N	TanH	conv_6_s	conv_7_s
8	Bin	-	-	-	-	conv_7_s	bin_pred_s
9	Conv	1	256	N	-	bin_pred_s	conv_8_s
10	Conv	1	256	Y	ReLU	conv_8_s	conv_9_s
11	DeConv	2	256	N	-	conv_9_s	deconv_10_s
12	Conv	1	128	Y	ReLU	deconv_10_s	conv_11_s
13	DeConv	2	128	N	-	conv_11_s	deconv_12_s
14	Conv	1	64	Y	ReLU	deconv_12_s	conv_13_s
15	DeConv	2	64	N	-	conv_13_s	deconv_14_s
16	Conv	1	3	N	TanH	deconv_13_s	conv_pred_s
17	Elt-Sub	-	-	N	-	{ conv_pred_s,	residual_s
						residual_(s-1)}	

Table 6: The architecture for stage s of a residual encoder adapted from Toderici et al. [20].

Layer	Туре	Stride	# Filters	BN	Non-Linearity	Input	Output
1	Conv	1	64	Y	ReLU	residual_(s-1)	conv_1_s
2	Conv	2	128	Y	ReLU	conv_1_s	conv_2_s
3	Conv	1	128	Y	ReLU	conv_2_s	conv_3_s
4	Conv	2	256	Y	ReLU	conv_3_s	conv_4_s
5	Conv	1	256	Y	ReLU	conv_4_s	conv_5_s
6	Conv	2	256	Y	ReLU	conv_5_s	conv_6_s
7	Conv*	1	8	N	TanH	conv_6_s	conv_7_s
8	Bin	-	-	-	-	conv_7_s	bin_pred_s
9	Conv	1	256	N	-	bin_pred_s	conv_8_s
10	Conv	1	256	Y	ReLU	conv_8_s	conv_9_s
11	DeConv	2	256	N	-	conv_9_s	deconv_10_s
12	Conv	1	128	Y	ReLU	deconv_10_s	conv_11_s
13	DeConv	2	128	N	-	conv_11_s	deconv_12_s
14	Conv	1	64	Y	ReLU	deconv_12_s	conv_13_s
15	DeConv	2	64	N	-	conv_13_s	deconv_14_s
16	Conv	1	3	N		deconv_14_s	conv_pred_s
17	Elt-Add	-	-	N	-	{ conv_pred_s,	output_s
					TanH	conv_pred_(s-1)}	-
18	Elt-Sub	-	-	N	-	{ output_s,	residual_s
						data}	

Table 7: The architecture for stage s of the R2I prediction connection model.

Layer	Type	Stride	# Filters	BN	Non-Linearity	Input	Output
1	Conv	1	64	Y	ReLU	residual_(s-1)	conv_1_s
2	Conv	1	64	Y	ReLU	conv_1_(s-1)	conv_1_s'
3	Elt-Add	-	-	N	TanH	$\{ conv_1_s,$	conv_1c_s
						conv_1_s'}	
4	Conv	2	128	Y	ReLU	conv_1c_s	conv_2_s
5	Conv	1	128	Y	ReLU	conv_2_(s-1)	conv_2_s'
6	Elt-Add	-	-	N	TanH	$\{ conv_2_s,$	conv_2c_s
						conv_2_s'}	
7	Conv	1	128	Y	ReLU	conv_2c_s	conv_3_s
8	Conv	1	128	Y	ReLU	conv_3_(s-1)	conv_3_s'
9	Elt-Add	_	_	N	TanH	{ conv_3_s,	conv_3c_s
						conv_3_s'}	
10	Conv	2	256	Y	ReLU	conv_3c_s	conv_4_s
11	Conv	1	256	Y	ReLU	conv_4_(s-1)	conv_4_s'
12	Elt-Add	_	_	N	TanH	{ conv_4_s,	conv_4c_s
	210 1100			- '		conv_4_s'}	00111_10_0
13	Conv	1	256	Y	ReLU	conv_4c_s	conv_5_s
14	Conv	1	256	Y	ReLU	conv_5_(s-1)	conv_5_s'
15	Elt-Add	-	200	N	TanH	{ conv_5_s,	conv_5c_s
13	Lit / idd			11	Tunit	conv_5_s'}	CONV_CC_B
16	Conv	2	256	Y	ReLU	conv_5c_s	conv_6_s
17	Conv	1	256	Y	ReLU	conv_6_(s-1)	conv_6_s
18	Elt-Add	-	-	N	TanH	{ conv_6_s,	conv_6c_s
10	Lit-Add	_	_	14	Taiii	conv_6_s'}	CONV_OC_S
19	Conv*	1	8	N	TanH	conv_6c_s	conv_7_s
20	Bin	-	-	-		conv_7_s	
21	Conv				-		bin_pred_s
22		1 1	256	N	- Dalii	bin_pred_s	conv_8_s
	Conv		256	Y	ReLU	conv_8_s	conv_9_s
23	Conv	1	256	Y	ReLU	conv_9_(s-1)	conv_9_s'
24	Elt-Add	-	-	N	TanH	{ conv_9_s,	conv_9c_s
25	D. C	0	050	N		conv_9_s'}	1 0
25	DeConv	2	256	N	- D I I I	conv_9_s	deconv_9_s
26	Conv	1	256	Y	ReLU	deconv_9_(s-1)	deconv_9_s'
27	Elt-Add	-	-	N	TanH	{ deconv_9_s,	deconv_9c_s
•	~	_	100			deconv_9_s'}	
28	Conv	1	128	Y	ReLU	deconv_9c_s	conv_10_s
29	Conv	1	128	Y	ReLU	conv_10_(s-1)	conv_10_s'
30	Elt-Add	-	-	N	TanH	{ conv_10_s,	conv_10c_s
						conv_10_s'}	
31	DeConv	2	128	N	-	conv_10c_s	deconv_11_s
32	Conv	1	128	Y	ReLU	deconv_11_(s-1)	deconv_11_s'
33	Elt-Add	-	-	N	TanH	$\{ \ {\tt deconv_11_s},$	deconv_11c_s
						deconv_11_s'}	
34	Conv	1	64	Y	ReLU	deconv_11c_s	conv_12_s
35	Conv	1	64	Y	ReLU	conv_12_(s-1)	conv_12_s'
36	Elt-Add	-	-	N	TanH	$\{ conv_12_s,$	conv_12c_s
						conv_12_s'}	
37	DeConv	2	64	N	-	conv_12c_s	deconv_13_s
38	Conv	1	128	Y	ReLU	deconv_13_(s-1)	deconv_13_s'
39	Elt-Add	-	-	N	TanH	{ deconv_13_s,	deconv_13c_s
						deconv_13_s'}	
40	Conv	1	3	N	TanH	deconv_13c_s	conv_pred_s
41	Elt-Sub	-	-	N	-	{ conv_pred_s,	data
						data}	

Table 8: The architecture for stage s of the R2I "full" connection model.

Layer	Type	Stride	# Filters	BN	Non-Linearity	Input	Output
1	Conv	1	64	Y	ReLU	residual_(s-1)	conv_1_s
2	Conv	2	128	Y	ReLU	conv_1_s	conv_2_s
3	Conv	1	128	Y	ReLU	conv_2_s	conv_3_s
4	Conv	2	256	Y	ReLU	conv_3_s	conv_4_s
5	Conv	1	256	Y	ReLU	conv_4_s	conv_5_s
6	Conv	2	256	Y	ReLU	conv_5_s	conv_6_s
7	Conv*	1	8	N	TanH	conv_6_s	conv_7_s
8	Bin	-	-	-	-	conv_7_s	bin_pred_s
9	Conv	1	256	N	-	bin_pred_s	conv_8_s
10	Conv	1	256	Y	ReLU	conv_8_s	conv_9_s
11	DeConv	2	256	N	-	conv_9_s	deconv_9_s
12	Conv	1	256	Y	ReLU	deconv_9_(s-1)	deconv_9_s'
13	Elt-Add	-	-	N	TanH	{ deconv_9_s, deconv_9_s'}	deconv_9c_s
14	Conv	1	128	Y	ReLU	deconv_9c_s	conv_10_s
15	DeConv	2	128	N	-	conv_10c_s	deconv_11_s
16	Conv	1	128	Y	ReLU	deconv_11_(s-1)	deconv_11_s'
17	Elt-Add	-	-	N	TanH	{ deconv_11_s, deconv_11_s'}	deconv_11c_s
18	Conv	1	64	Y	ReLU	deconv_11c_s	conv_12_s
19	DeConv	2	64	N	-	conv_12c_s	deconv_13_s
20	Conv	1	128	Y	ReLU	deconv_13_(s-1)	deconv_13_s'
21	Elt-Add	-	-	N	TanH	{ deconv_13_s, deconv_13_s'}	deconv_13c_s
22	Conv	1	3	N	TanH	deconv_13c_s	conv_pred_s
23	Elt-Sub	-	-	N	-	$\{ ext{ conv_pred_s}, \\ ext{ data} \}$	data

Table 9: The architecture for stage \boldsymbol{s} of the R2I "decoding" connection model.

Layer	Type	Stride	# Filters	Dilation	BN	Non-Linearity	Input	Output
1	Conv	1	24	1	Y	ReLU	data	conv_0_1
2	Conv	1	24	2	Y	ReLU	data	conv_0_2
3	Conv	1	24	4	Y	ReLU	data	conv_0_4
4	Conv	1	24	8	Y	ReLU	data	conv_0_8
5	Concat	-	-	-	N	_	$\{ conv_0_1, conv_0_2, $	conv_0
		-	-			_	conv_0_4,conv_0_8 }	
6	Conv	1	24	1	Y	ReLU	conv_0	conv_1_1
7	Conv	1	24	2	Y	ReLU	conv_0	conv_1_2
8	Conv	1	24	4	Y	ReLU	conv_0	conv_1_4
9	Conv	1	$\overline{24}$	8	Y	ReLU	conv_0	conv_1_8
10	Concat	-	-	-	N	-	{ conv_1_1,conv_1_2,	conv_1
10	Concui	_	_		11	_	conv_1_4,conv_1_8 }	00111_1
11	Conv	1	$\overline{24}$	1	Y	ReLU	conv_1	conv_2_1
12	Conv	1	24	2	Y	ReLU	conv_1	conv_2_1
13	Conv	1	$\frac{24}{24}$	4	Y	ReLU	_	
					Y		conv_1	conv_2_4
14	Conv	1	24	8		ReLU	conv_1	conv_2_8
15	Concat	-	-	-	N	-	{ conv_2_1,conv_2_2,	conv_2
1.6	-	-	-		* 7	-	conv_2_4,conv_2_8 }	
16	Conv	1	24	1	Y	ReLU	conv_2	conv_3_1
17	Conv	1	24	2	Y	ReLU	conv_2	conv_3_2
18	Conv	1	24	4	Y	ReLU	conv_2	conv_3_4
19	Conv	1	24	8	Y	ReLU	conv_2	conv_3_8
20	Concat	-	-	-	N	-	{ conv_3_1,conv_3_2, conv_3_4,conv_3_8 }	conv_3
21	Conv	1	24	1	Y	ReLU	conv_3	conv_4_1
22	Conv	1	24	2	Y	ReLU	conv_3	$conv_4_2$
23	Conv	1	24	4	Y	ReLU	conv_3	conv_4_4
24	Conv	1	24	8	Y	ReLU	conv_3	conv_4_8
25	Concat	_	-	_	N	_	{ conv_4_1,conv_4_2,	conv_3
		_	-			_	conv_4_4,conv_4_8 }	_
26	Conv	1	24	1	Y	ReLU	conv_4	conv_5_1
27	Conv	1	24	2	Y	ReLU	conv_4	conv_5_2
28	Conv	1	24	4	Y	ReLU	conv_4	conv_5_4
29	Conv	1	24	8	Y	ReLU	conv_4	conv_5_8
30	Concat	-	-	-	N	-	{ conv_5_1,conv_5_2,	conv_5
30	Concai	_	-	_	14	_	conv_5_4,conv_5_8 }	COIIV_5
31	Conv	1	24	1	v	Dal II		20mm 6 1
32	Conv Conv			1 2	Y Y	ReLU	conv_5	conv_6_1
		1	24			ReLU	conv_5	conv_6_2
33	Conv	1	24	4	Y	ReLU	conv_5	conv_6_4
34	Conv	1	24	8	Y	ReLU	conv_5	conv_6_8
35	Concat	-	-	-	N	-	{ conv_6_1,conv_6_2, conv_6_4,conv_6_8 }	conv_6
36	Conv	1	24	1	Y	ReLU	conv_6	conv_7_1
37	Conv	1	24	2	Y	ReLU	conv_6	$conv_7_2$
38	Conv	1	24	4	Y	ReLU	conv_6	conv_7_4
39	Conv	1	$\overline{24}$	8	Y	ReLU	conv_6	conv_7_8
40	Concat	-	-	-	N	-	{ conv_7_1,conv_7_2,	conv_7
		-	-	a.		-	conv_7_4,conv_7_8 }	_
41	Conv	1	3	1	N	-	conv_7	conv_RGB
42	Crop	3	3	-	-		$\mathtt{conv}_\mathtt{RGB}$	conv_RGB_cro
43	-	-	-	-	-	TanH	conv_RGB_crop	pred

Table 10: The architecture for the Inpainting Network.