

StegaStamp: Invisible Hyperlinks in Physical Photographs

Matthew Tancik*

Ben Mildenhall*

Ren Ng

University of California, Berkeley
 {tancik, bmild, ren}@berkeley.edu

Abstract

Printed and digitally displayed photos have the ability to hide imperceptible digital data that can be accessed through internet-connected imaging systems. Another way to think about this is physical photographs that have unique QR codes invisibly embedded within them. This paper presents an architecture, algorithms, and a prototype implementation addressing this vision. Our key technical contribution is StegaStamp, a learned steganographic algorithm to enable robust encoding and decoding of arbitrary hyperlink bitstrings into photos in a manner that approaches perceptual invisibility. StegaStamp comprises a deep neural network that learns an encoding/decoding algorithm robust to image perturbations approximating the space of distortions resulting from real printing and photography. We demonstrate real-time decoding of hyperlinks in photos from in-the-wild videos that contain variation in lighting, shadows, perspective, occlusion and viewing distance. Our prototype system robustly retrieves 56 bit hyperlinks after error correction – sufficient to embed a unique code within every photo on the internet. Code is available at <https://github.com/tancik/StegaStamp>.

1. Introduction

Our vision is a future in which each photo in the real world invisibly encodes a unique hyperlink to arbitrary information. This information is accessed by pointing a camera at the photo and using the system described in this paper to decode and follow the hyperlink. In the future, augmented-reality (AR) systems may perform this task continuously, visually overlaying retrieved information alongside each photo in the user’s view.

Our approach is related to the ubiquitous QR code and similar technologies, which are now commonplace for a wide variety of data-transfer tasks, such as sharing web addresses, purchasing goods, and tracking inventory. Our approach can be thought of as a complementary solution that avoids visible, ugly barcodes, and enables digital information to be invisibly

and ambiently embedded into the ubiquitous imagery of the modern visual world.

It is worth taking a moment to consider three potential use cases of our system. First, at the farmer’s market, a stand owner may add photos of each type of produce alongside the price, encoded with extra information for customers about the source farm, nutrition information, recipes, and seasonable availability. Second, in the lobby of a university department, a photo directory of faculty may be augmented by encoding a unique URL for each person’s photo that contains the professor’s webpage, office hours, location, and directions. Third, New York City’s Times Square is plastered with digital billboards. Each image frame displayed may be encoded with a URL containing further information about the products, company, and promotional deals.

Figure 1 presents an overview of our system, which we call StegaStamp, in the context of a typical usage flow. The inputs are an image and a desired hyperlink. First, we assign the hyperlink a unique bit string (analogous to the process used by URL-shortening services such as tinyurl.com). Second, we use our StegaStamp encoder to embed the bit string into the target image. This produces an encoded image that is ideally perceptually identical to the input image. As described in detail in Section 4, our encoder is implemented as a deep neural network jointly trained with a second network that implements decoding. Third, the encoded image is physically printed (or shown on an electronic display) and presented in the real world. Fourth, a user takes a photo that contains the physical print. Fifth, the system uses an image detector to identify and crop out all images. Sixth, each image is processed with the StegaStamp decoder to retrieve the unique bitstring, which is used to follow the hyperlink and retrieve the information associated with the image.

This method of data transmission has a long history in both the steganography and watermarking literatures. We present the first end-to-end trained deep pipeline for this problem that can achieve robust decoding even under “physical transmission,” delivering excellent performance sufficient to encode and retrieve arbitrary hyperlinks for an essentially limitless number of images. We extend the traditional learned steganography framework by adding a set of differ-

*Authors contributed equally to this work.

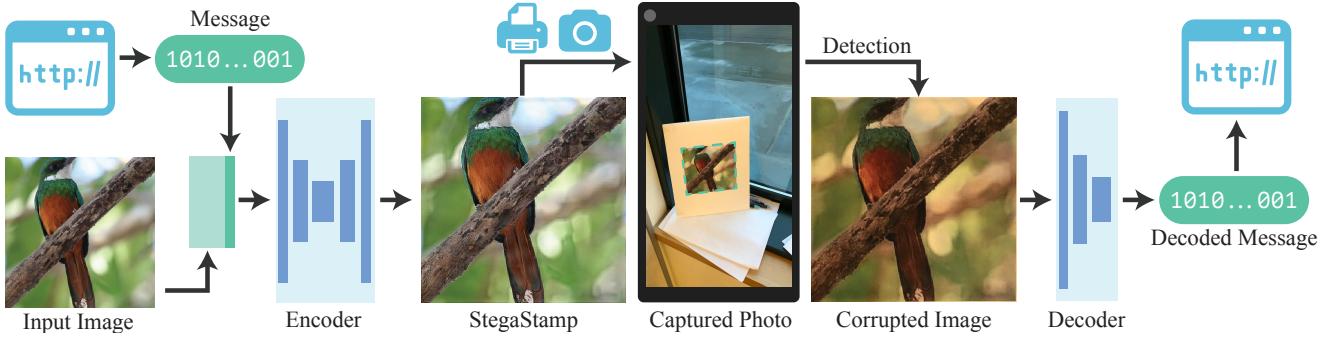


Figure 1: Our deep learning system is trained to hide hyperlinks in images. First, an encoder network processes the input image and hyperlink bitstring into a StegaStamp (encoded image). The StegaStamp is then printed and captured by a camera. A detection network localizes and rectifies the StegaStamp before passing it to the decoder network. After the bits are recovered and error corrected, the user can follow the hyperlink. To train the encoder and decoder networks, we simulate the corruptions caused by printing, reimaging, and detecting the StegaStamp with a set of differentiable image augmentations.

entiable pixelwise and spatial image corruptions between the encoder and decoder that successfully approximate the space of distortions resulting from “physical transmission” (i.e., real printing or display and subsequent image capture). The result is robust retrieval of 95% of 100 encoded bits in real-world conditions while preserving excellent perceptual image quality. This allows our prototype to uniquely encode hidden hyperlinks for orders of magnitude more images than exist on the internet today (upper bounded by 100 trillion).

2. Related Work

2.1. Steganography

Steganography is the act of hiding data within other data and has a long history that can be traced back to ancient Greece. Our proposed task is a type of steganography where we hide a code within an image. Various methods have been developed for digital image steganography. Data can be hidden in the least significant bits of the image, subtle color variations, and subtle luminosity variations. Often methods are designed to evade steganalysis, the detection of hidden messages [18, 34]. We refer the interested reader to surveys [9, 11] that review a wide set of techniques.

The most relevant work to our proposal are methods that utilize deep learning to both encode and decode a message hidden inside an image [5, 21, 43, 47, 51, 54, 44]. Our method assumes that the image will be corrupted by a display-imaging pipeline between the encoding and decoding steps. With the exception of HiDDeN [54] and Light Field Messaging (LFM) [45], small image manipulations or corruptions would render existing techniques useless, as their goal is encoding a large number of bits-per-pixel in the context of perfect digital transmission. HiDDeN introduces various types of noise between encoding and decoding to

increase robustness but focuses only on the set of corruptions that would occur through digital image manipulations (e.g., JPEG compression and cropping). For use as a physical barcode, the decoder cannot assume perfect alignment, given the perspective shifts and pixel resampling guaranteed to occur when taking a casual photo. LFM [45] obtain robustness using a network trained on a large dataset of manually photographed monitors to undo the camera-display corruptions. Our method does not require this time-intensive dataset capture step and generalizes to printed images, a medium for which collecting training data would be even more difficult.

2.2. Watermarking

Watermarking, a form of steganography, has long been considered as a potential way to link a physical image to an Internet resource [2]. Early work in the area defined a set of desirable goals for robust watermarking, including invisibility and robustness to image manipulations [7]. Later research demonstrated the significant robustness benefits of encoding the watermark in the log-polar frequency domain [27, 33, 35, 53]. Similar methods have been optimized for use as interactive mobile phone applications [13, 31, 36]. Additional work focuses on carefully modeling the printer-camera transform [37, 42] or display-camera transform [17, 46, 50] for better information transfer. Some approaches to display-camera communication take advantage of the unique properties of this hardware combination such as polarization [49], rolling shutter artifacts [26], or high frame rate [12]. A related line of work in image forensics explores whether it is possible to use a CNN to detect when an image has been re-imaged [16]. In contrast to the hand-designed pipelines used in previous work on watermarking, our method automatically *learns* how to hide and transmit data in a way that is robust to many different

combinations of printers/displays, cameras, lighting, and viewpoints. We provide a framework for training this system and a rigorous evaluation of its capabilities, demonstrating that it works in many real world scenarios and using ablations to show the relative importance of our training perturbations.

2.3. Barcodes

Barcodes are one of the most popular solutions for transmitting a short string of data to a computing device, requiring only simple hardware (a laser reader or camera) and an area for printing or displaying the code. Traditional barcodes are a one dimensional pattern where bars of alternating thickness encode different values. The ubiquity of high quality cellphone cameras has led to the frequent use of two dimensional QR codes to transmit data to and from phones. For example, users can share contact information, pay for goods, track inventory, or retrieve a coupon from an advertisement.

Past research has addressed the issue of robustly decoding existing or new barcode designs using cameras [29, 32]. Some designs particularly take advantage of the increased capabilities of cameras beyond simple laser scanners in various ways, such as incorporating color into the barcode [8]. Other work has proposed a method that determines where a barcode should be placed on an image and what color should be used to improve machine readability [30].

Another special type of barcode is specially designed to transmit both a small identifier and a precise six degree-of-freedom orientation for camera localization or calibration, e.g., ArUco markers [19, 38]. Hu *et al.* [22] train a deep network to localize and identify ArUco markers in challenging real world conditions using data augmentation similarly to our method. However, their focus is robust detection of highly visible preexisting markers, as opposed to robust decoding of messages hidden in arbitrary natural images.

2.4. Robust Adversarial Image Attacks

Adversarial image attacks on object classification CNNs are designed to minimally perturb an image in order to produce an incorrect classification. Most relevant to our work are the demonstrations of adversarial examples in the physical world [4, 10, 15, 25, 28, 40, 41], where systems are made robust for imaging applications by modeling physically realistic perturbations (i.e., affine image warping, additive noise, and JPEG compression). Jan *et al.* [25] take a different approach, explicitly training a neural network to replicate the distortions added by an imaging system and showing that applying the attack to the distorted image increases the success rate.

These results demonstrate that networks can still be affected by small perturbations after the image has gone through an imaging pipeline. Our proposed task shares some similarities; however, classification targets 1 of $n \approx 2^{10}$ labels, while we aim to uniquely decode 1 of 2^m messages,

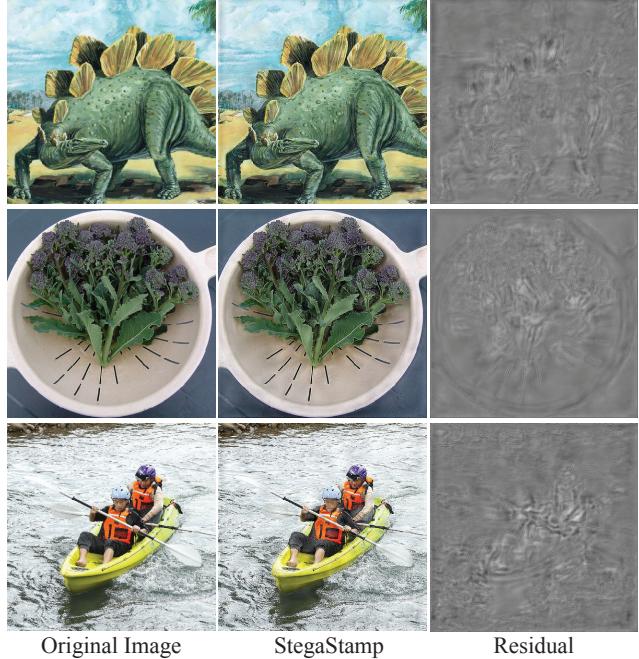


Figure 2: Examples of encoded images. The residual is calculated by the encoder network and added back to the original image to produce the encoded StegaStamp. These examples have 100 bit encoded messages and are robust to the image perturbations that occur through the printing and imaging pipelines.

where $m \approx 100$ is the number of encoded bits. Additionally, adversarial attacks typically do not modify the decoder network, whereas we explicitly train our decoder to cooperate with our encoder for maximum information transferal.

3. Training for Real World Robustness

During training, we apply a set of differentiable image perturbations outlined in Figure 3 between the encoder and decoder to approximate the distortions caused by physically displaying and imaging the StegaStamps. Previous work on synthesizing robust adversarial examples used a similar method to attack classification networks in the wild (termed “Expectation over Transformation”), though they used a more limited set of transformations [4]. HiDDeN [54] used non-spatial perturbations to augment their steganography pipeline against digital perturbations only. Deep ChArUco [22] used both spatial and nonspatial perturbations to train a robust detector specifically for ChArUco fiducial marker boards. We combine ideas from all of these works, training an encoder and decoder that cooperate to robustly transmit hidden messages through a physical display-imaging pipeline.

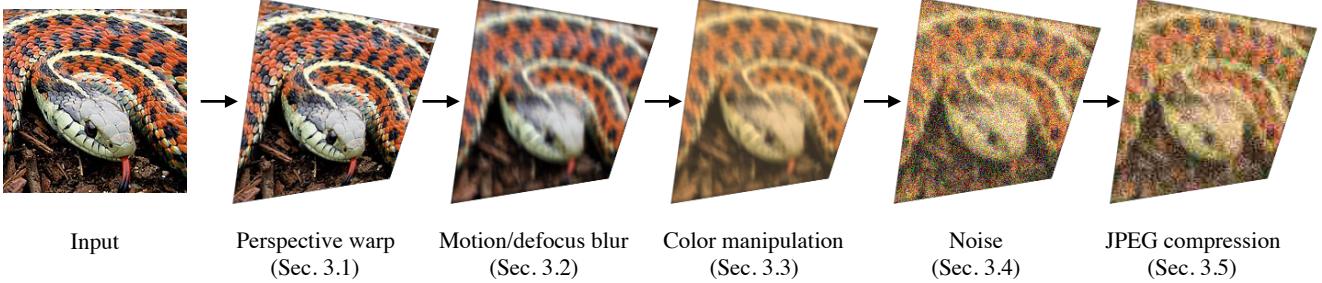


Figure 3: Image perturbation pipeline. During training, we approximate the effects of a physical display-imaging pipeline in order to make our model robust for use in the real world. We take the output of the encoding network and apply the random transformations shown here before passing the image through the decoding network (see Section 3 for details).

3.1. Perspective Warp

Assuming a pinhole camera model, any two images of the same planar surface can be related by a homography. We generate a random homography to simulate the effect of a camera that is not precisely aligned with the encoded image marker. To sample a homography, we randomly perturb the four corner locations of the marker uniformly within a fixed range (up to ± 40 pixels, i.e. $\pm 10\%$) then solve for the homography that maps the original corners to their new locations. We bilinearly resample the original image to create the perspective warped image.

3.2. Motion and Defocus Blur

Blur can result from both camera motion and inaccurate autofocus. To simulate motion blur, we sample a random angle and generate a straight line blur kernel with a width between 3 and 7 pixels. To simulate misfocus, we use a Gaussian blur kernel with its standard deviation randomly sampled between 1 and 3 pixels.

3.3. Color Manipulation

Printers and displays have a limited gamut compared to the full RGB color space. Cameras modify their output using exposure settings, white balance, and a color correction matrix. We approximate these perturbations with a series of random affine color transformations (constant across the whole image) as follows:

1. Hue shift: adding a random color offset to each of the RGB channels sampled uniformly from $[-0.1, 0.1]$.
2. Desaturation: randomly linearly interpolating between the full RGB image and its grayscale equivalent.
3. Brightness and contrast: affine histogram rescaling $mx + b$ with $m \sim U[0.5, 1.5]$ and $b \sim U[-0.3, 0.3]$.

After these transforms, we clip the color channels to $[0, 1]$.

3.4. Noise

Noise introduced by camera systems is well studied and includes photon noise, dark noise, and shot noise [20]. We assume standard non-photon-starved imaging conditions, employing a Gaussian noise model (sampling the standard deviation $\sigma \sim U[0, 0.2]$) to account for imaging noise.

3.5. JPEG Compression

Camera images are usually stored in a lossy format such as JPEG. JPEG compresses images by computing the discrete cosine transform of each 8×8 block in the image and quantizing the resulting coefficients by rounding to the nearest integer (at varying strengths for different frequencies). This rounding step is not differentiable, so we use the trick from Shin and Song [40] for approximating the quantization step near zero with the piecewise function

$$q(x) = \begin{cases} x^3 & : |x| < 0.5 \\ x & : |x| \geq 0.5 \end{cases} \quad (1)$$

which has nonzero derivative almost everywhere. We sample the JPEG quality uniformly within $[50, 100]$.

4. Implementation Details

4.1. Encoder

The encoder is trained to embed a message into an image while minimizing perceptual differences between the input and encoded images. We use a U-Net [39] style architecture that receives a four channel 400×400 pixel input (input image RGB channels plus one for the message) and outputs a three channel RGB residual image. The input message is represented as a 100 bit binary string, processed through a fully connected layer to form a $50 \times 50 \times 3$ tensor, then upsampled to produce a $400 \times 400 \times 3$ tensor. We find that applying this preprocessing to the message aids convergence. We present examples of encoded images in Figure 2.

4.2. Decoder

The decoder is a network trained to recover the hidden message from the encoded image. A spatial transformer network [24] is used to develop robustness against small perspective changes that are introduced while capturing and rectifying the encoded image. The transformed image is fed through a series of convolutional and dense layers and a sigmoid to produce a final output with the same length as the message. The decoder network is supervised using cross entropy loss.

4.3. Detector

For real world use, we must detect and rectify StegaStamps within a wide field of view image before decoding them, since the decoder network alone is not designed to handle full detection within a much larger image. We fine-tune an off-the-shelf semantic segmentation network BiSeNet [48] to segment areas of the image that are believed to contain StegaStamps. The network is trained using a dataset of randomly transformed StegaStamps embedded into high resolution images sampled from DIV2K [1]. At test time, we fit a quadrilateral to the convex hull of each of the network’s proposed regions, then compute a homography to warp each quadrilateral back to a 400×400 pixel image for parsing by the decoder.

4.4. Encoder/Decoder Training Procedure

Training Data During training, we use images from the MIRFLICKR dataset [23] (resampled to 400×400 resolution) combined with randomly sampled binary messages.

Critic As part of our total loss, we use a critic network that predicts whether a message is encoded in a image and is used as a perceptual loss for the encoder/decoder pipeline. The network is composed of a series of convolutional layers followed by max pooling. To train the critic, an input image and an encoded image are classified and the Wasserstein loss [3] is used as a supervisory signal. Training of the critic is interleaved with the training of the encoder/decoder.

Losses To enforce minimal perceptual distortion on the encoded StegaStamp, we use an L_2 residual regularization L_R , the LPIPS perceptual loss [52] L_P , and a critic loss L_C calculated between the encoded image and the original image. We use cross entropy loss L_M for the message. The training loss is the weighted sum of these loss components.

$$L = \lambda_R L_R + \lambda_P L_P + \lambda_C L_C + \lambda_M L_M \quad (2)$$

We find three loss function adjustments to particularly aid in convergence when training the networks:

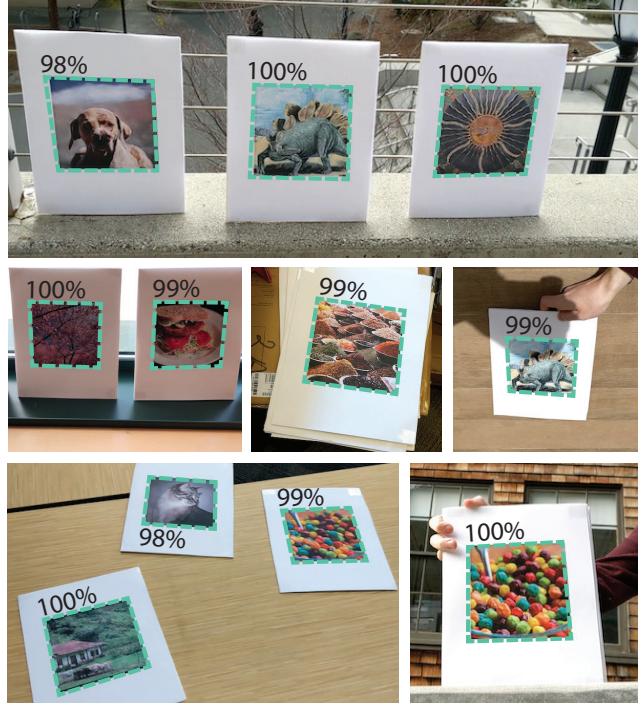


Figure 4: Examples of our system deployed in-the-wild. We outline the StegaStamps detected and decoded by our system and the show message recovery accuracies. Our method works in the real world, exhibiting robustness to changing camera orientation, lighting, shadows, etc. You can find these examples and more in our supplemental video.

1. These image loss weights $\lambda_{R,P,C}$ must initially be set to zero while the decoder trains to high accuracy, after which $\lambda_{R,P,C}$ are increased linearly.
2. The image perturbation strengths must also start at zero. The perspective warping is the most sensitive perturbation and is increased at the slowest rate.
3. The model learns to add distracting patterns at the edge of the image (perhaps to assist in localization). We mitigate this effect by increasing the weight of the L_2 loss at the edges with a cosine dropoff.

5. Real-World & Simulation-Based Evaluation

We test our system in both real-world conditions and synthetic approximations of display-imaging pipelines. We show that our system works in-the-wild, recovering messages in uncontrolled indoor and outdoor environments. We evaluate our system in a controlled real world setting with 18 combinations of 6 different displays/printers and 3 different cameras. Across all settings combined (1890 captured



Figure 5: Despite not explicitly training the method to be robust to occlusion, we find that our decoder can handle partial erasures gracefully, maintaining high accuracy.

images), we achieve a mean bit-accuracy of 98.7%. We conduct real and synthetic ablation studies with four different trained models to verify that our system is robust to each of the perturbations we apply during training and that omitting these augmentations significantly decreases performance.

5.1. In-the-Wild Robustness

Our method is tested on handheld cellphone camera videos captured in a variety of real-world environments. The StegaStamps are printed on a consumer printer. Examples of the captured frames with detected quadrilaterals and decoding accuracy are shown in Figure 4. We also demonstrate a surprising level of robustness when portions of the StegaStamp are covered by other objects (Figure 5). Please see our supplemental video for extensive examples of real world StegaStamp decoding, including examples of perfectly recovering 56 bit messages using BCH error correcting codes [6]. We generally find that if the bounding rectangle is accurately located, decoding accuracy is high. However, it is possible for the detector to miss the StegaStamp on a subset of video frames. In practice this is not an issue, because the code only needs to be recovered once. We expect future extensions that incorporate temporal information and custom detection networks can further improve the detection consistency.

5.2. Controlled Real World Experiments

In order to demonstrate that our model generalizes from synthetic perturbations to real physical display-imaging pipelines, we conduct a series of test where encoded images are printed or displayed, recaptured by a camera, then decoded. We randomly select 100 unique images from the ImageNet dataset [14] (disjoint from our training set) and embed random 100 bit messages within each image. We generate 5 additional StegaStamps with the same source image but different messages for a total of 105 test images. We conduct the experiments in a darkroom with fixed lighting. The printed images are fixed in a rig for consistency and captured by a tripod-mounted camera. The resulting photographs are cropped by hand, rectified, and passed through the decoder.

The images are printed using a consumer printer (HP

			5th	25th	50th	Mean
Webcam	Screen	Enterprise	88%	94%	98%	95.9%
	Printer	Consumer	90%	98%	99%	98.1%
	Screen	Pro	97%	99%	100%	99.2%
	Printer	Monitor	94%	98%	99%	98.5%
	Screen	Laptop	97%	99%	100%	99.1%
	Printer	Cellphone	91%	98%	99%	97.7%
Cellphone	Screen	Enterprise	88%	96%	98%	96.8%
	Printer	Consumer	95%	99%	100%	99.0%
	Screen	Pro	97%	99%	100%	99.3%
	Printer	Monitor	98%	99%	100%	99.4%
	Screen	Laptop	98%	99%	100%	99.7%
	Printer	Cellphone	96%	99%	100%	99.2%
DSLR	Screen	Enterprise	86%	96%	99%	97.0%
	Printer	Consumer	97%	99%	100%	99.3%
	Screen	Pro	98%	99%	100%	99.5%
	Printer	Monitor	99%	100%	100%	99.8%
	Screen	Laptop	99%	100%	100%	99.8%
	Printer	Cellphone	99%	100%	100%	99.8%

Table 1: Real world decoding accuracy (percentage of bits correctly recovered) tested using a combination of six display methods (three printers and three screens) and three cameras. We show the 5th, 25th, and 50th percentiles and mean taken over 105 images chosen randomly from ImageNet [14] with randomly sampled 100 bit messages.

LaserJet Pro M281fdw), an enterprise printer (HP LaserJet Enterprise CP4025), and a commercial printer (Xerox 700i Digital Color Press). The images are also digitally displayed on a matte 1080p monitor (Dell ST2410), a glossy high DPI laptop screen (Macbook Pro 15 inch), and an OLED cellphone screen (iPhone X). To image the StegaStamps, we use an HD webcam (Logitech C920), a cellphone camera (Google Pixel 3), and a DSLR camera (Canon 5D Mark II). All devices use their factory calibration settings. Each of the 105 images were captured across all 18 combinations of the 6 media and 3 cameras. The results are reported in Table 1. Our method is highly robust across a variety of different combinations of display/printer and camera; two-thirds of these scenarios yield a median accuracy of 100% and a 5th percentile accuracy of at least 95% perfect decoding. Our mean accuracy over all 1890 captured images is 98.7%.

Using a test set comprised of the cellphone camera + consumer printer combination, we compare variants of our method (described further in Section 5.3) to Baluja [5], HiDDeN [54], and LFM [44] in Figure 6. The variants of our model use the same architecture but are trained with different augmentations; the names *None*, *Pixelwise*, *Spatial*, and *All* indicate which categories of perturbations were applied during training. We see that Baluja [5], trained with a minimal amount of augmented noise (similar to our *None* variant) performs no better than guessing. HiDDeN [54] incorporates

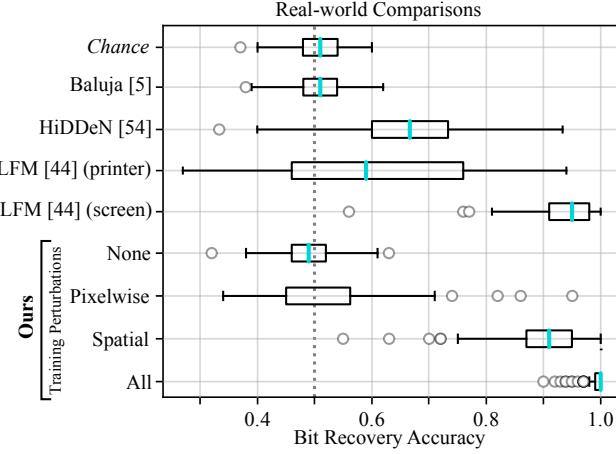


Figure 6: Real world comparisons of variants of our method described in Section 5.3 and competing methods, using the cellphone camera + consumer printer pipeline from Table 1. We show the distribution of random guessing (with its mean of 0.5 indicated by the dotted line) to demonstrate that the no-perturbations ablation and Baluja [5] perform no better than chance. HiDDeN [54] uses pixelwise perturbations along with random masking. Adding spatial perturbations is critical for achieving high real-world performance. LFM [44] works well on screens but fails to generalize to printed media.

augmentations into their training pipeline to increase robustness to perturbations. Their method is trained with a set of pixelwise perturbations along with a “cropping” augmentation that masks out a random image region. However, it lacks augmentations that spatially resample the image, and we find that its accuracy falls between our *Pixelwise* and *Spatial* variants. LFM [44] specifically trains a “distortion” network to mimic the effect of displaying and recapturing an encoded image, trained on a dataset they collect of over 1 million images from 25 display/camera pairs. In this domain (“screen”), we find LFM performs fairly well. However, it does not generalize to printer/camera pipelines (“printer”). Please refer to the supplement for testing details regarding the compared methods. Among our own ablated variants, we see that training with spatial perturbations alone yields significantly higher performance than only using pixelwise perturbations; however, *Spatial* still does not reliably recover enough data for practical use. Our presented method (*All*), combining both pixelwise and spatial perturbations, achieves the most precise and accurate results by a large margin.

5.3. Synthetic Ablation Test

We test how training with different subsets of the image perturbations from Section 3 impacts decoding accuracy in a synthetic experiment (Figure 7). We evaluate both our base model (trained with all perturbations) and three additional

Metric	Message length			
	50	100	150	200
PSNR \uparrow	29.88	28.50	26.47	21.79
SSIM \uparrow	0.930	0.905	0.876	0.793
LPIPS \downarrow	0.100	0.101	0.128	0.184

Table 2: Image quality for models trained with different message lengths, averaged over 500 images. For PSNR and SSIM, higher is better. LPIPS [52] is a learned perceptual similarity metric, lower is better.

models (trained with no perturbations, only pixelwise perturbations, and only spatial perturbations). Most work on learned image steganography focuses on hiding as much information as possible, assuming that no corruption will occur prior to decoding (as in our “no perturbations” model).

We run a more exhaustive synthetic ablation study over 1000 images to separately test the effects of each training-time perturbation on accuracy. The results shown in Figure 7 follow a similar pattern to the real world comparison test. The model trained with no perturbations is surprisingly robust to color warps and noise but immediately fails in the presence of warp, blur, or any level of JPEG compression. Training with only pixelwise perturbations yields high robustness to those augmentations but still leaves the network vulnerable to any amount of pixel resampling from warping or blur. On the other hand, training with only spatial perturbations also confers increased robustness against JPEG compression (perhaps because it has a similar low-pass filtering effect to blurring). Again, training with both spatial and pixelwise augmentations yields the best result.

5.4. Practical Message Length

Our model can be trained to store different numbers of bits. In all previous examples, we use a message length of 100. Figure 8 compares encoded images from four separately trained models with different message lengths. Larger message are more difficult to encode and decode; as a result, there is a trade off between recovery accuracy and perceptual similarity. The associated image metrics are reported in Table 2. When training, the image and message losses are tuned such that the bit accuracy converges to at least 95%.

We settle on a message length of 100 bits as it provides a good compromise between image quality and information transfer. Given an estimate of at least 95% recovery accuracy, we can encode at least 56 error corrected bits using BCH codes [6]. As discussed in the introduction, this gives us the ability to uniquely map every recorded image in history to a corresponding StegaStamp. Accounting for error correcting, using only 50 total message bits would drastically reduce the number of possible encoded hyperlinks to under one billion. The image degradation caused by encoding 150 or 200 bits is much more perceptible.

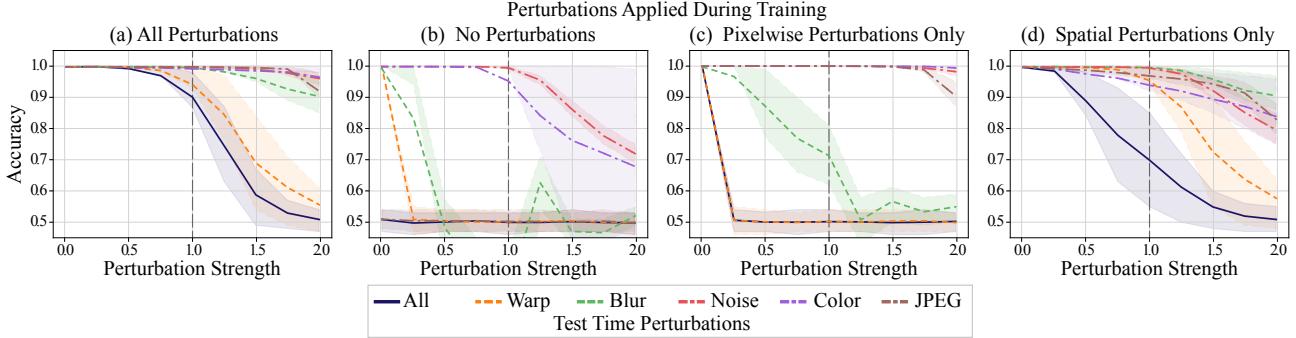


Figure 7: Synthetic ablation tests showing the effect of training with various image perturbation combinations on bit recovery robustness. “Pixelwise” perturbations (c) are noise, color transforms, and JPEG compression, and “spatial” perturbations (d) are perspective warp and blur. To test robustness across a range of possible degradation, we parameterize the strength of each perturbation on a scale from 0 (weakest) to 1 (maximum value seen during training) to 2 (strongest). Models not trained against spatial perturbations (b-c) are highly susceptible to warp and blur, and the model trained only on spatial perturbations (d) is sensitive to color transformations. The lines show the mean accuracies and the shaded regions shows the 25th-75th percentile range over 100 random images and messages. See Section 5.3 for details.

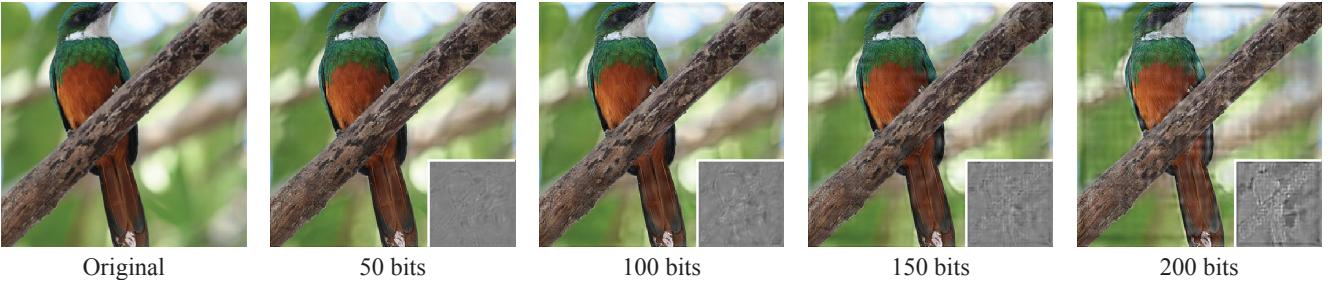


Figure 8: Four models trained to encode messages of different lengths. The inset shows the residual relative to the original image. The perceptual quality decreases as more bits are encoded. We find that a message length of 100 bits provides good image quality and is sufficient to encode a virtually unlimited number of distinct hyperlinks using error correcting codes.

5.5. Limitations

Though our system works with a high rate of success in the real world, it is still many steps from enabling broad deployment. Despite often being very subtle in high frequency textures, the residual added by the encoder network is sometimes perceptible in large low frequency regions of the image. Future work could improve upon our architecture and loss functions to generate more subtle encodings.

Additionally, we find our off-the-shelf detection network to be the bottleneck in our decoding performance during real world testing. A custom detection architecture optimized end to end with the encoder/decoder could increase detection performance. The current framework also assumes that the StegaStamps will be single, square images for the purpose of detection. We imagine that embedding multiple codes seamlessly into a single, larger image (such as a poster or billboard) could provide even more flexibility.

6. Conclusion

We have presented an end-to-end deep learning framework for encoding 56 bit error corrected hyperlinks into arbitrary natural images. Our networks are trained through an image perturbation module that allows them to generalize to real world display-imaging pipelines. We demonstrate robust decoding performance on a variety of printer, screen, and camera combinations in an experimental setting. We also show that our method is stable enough to be deployed in-the-wild as a replacement for existing barcodes that is less intrusive and more aesthetically pleasing.

7. Acknowledgments

We thank Coline Devin and Cecilia Zhang for acting in our supplemental video and Utkarsh Singhal and Pratul Srinivasan for useful feedback. BM is supported by a Hertz Fellowship and MT is supported by NSF GRFP.

References

- [1] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *CVPR Workshops*, 2017. 5
- [2] Adnan M. Alattar. Smart images using digimarc’s watermarking technology, 2000. 2
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *ICML*, 2017. 5
- [4] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. In *ICML*, 2018. 3
- [5] Shumeet Baluja. Hiding images in plain sight: Deep steganography. In *NeurIPS*, 2017. 2, 6, 7, 11
- [6] Raj Chandra Bose and Dwijendra K Ray-Chaudhuri. On a class of error correcting binary group codes. *Information and Control*, 1960. 6, 7, 11
- [7] G. W. Braudaway. Protecting publicly-available images with an invisible image watermark. In *ICIP*, 1997. 2
- [8] Orhan Bulan, Henryk Blasinski, and Gaurav Sharma. Color qr codes: Increased capacity via per-channel data encoding and interference cancellation. In *Color and Imaging Conference*, 2011. 3
- [9] Abbas Cheddad, Joan Condell, Kevin Curran, and Paul Mc Kevitt. Digital image steganography: Survey and analysis of current methods. *Signal Processing*, 90(3), 2010. 2
- [10] Shang-Tse Chen, Cory Cornelius, Jason Martin, and Duen Horng Polo Chau. Shapeshifter: Robust physical adversarial attack on faster r-cnn object detector. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2018. 3
- [11] Ingemar Cox, Matthew Miller, Jeffrey Bloom, Jessica Fridrich, and Ton Kalker. *Digital watermarking and steganography*. Morgan kaufmann, 2007. 2
- [12] Hao Cui, Huanyu Bian, Weiming Zhang, and Nenghai Yu. Unseencode: Invisible on-screen barcode with image-based extraction. In *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, pages 1315–1323. IEEE, 2019. 2
- [13] L. A. Delgado-Guillen, J. J. Garcia-Hernandez, and C. Torres-Huitzil. Digital watermarking of color images utilizing mobile platforms. *IEEE MWSCAS*, 2013. 2
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009. 6
- [15] Kevin Eykholt, Ivan Evtimov, Earlence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno, and Dawn Song. Robust physical-world attacks on deep learning visual classification. In *CVPR*, 2018. 3
- [16] W. Fan, S. Agarwal, and H. Farid. Rebroadcast attacks: Defenses, reattacks, and redefenses. *EUSIPCO*, 2018. 2
- [17] H. Fang, W. Zhang, H. Zhou, H. Cui, and N. Yu. Screen-shooting resilient watermarking. *IEEE TIFS*, 2019. 2
- [18] Jessica Fridrich, Tomáš Pevný, and Jan Kodovský. Statistically undetectable jpeg steganography: Dead ends, challenges, and opportunities. In *Proceedings of the 9th Workshop on Multimedia & Security*, 2007. 2
- [19] Sergio Garrido-Jurado, Rafael Muoz-Salinas, Francisco Madrid-Cuevas, and Rafael Medina-Carnicer. Generation of fiducial marker dictionaries using mixed integer linear programming. *Pattern Recognition*, 2015. 3
- [20] Samuel W. Hasinoff. Photon, poisson noise. In *Computer Vision: A Reference Guide*. 2014. 4
- [21] Jamie Hayes and George Danezis. Generating steganographic images via adversarial training. In *NeurIPS*, 2017. 2
- [22] Danying Hu, Daniel DeTone, and Tomasz Malisiewicz. Deep charuco: Dark charuco marker pose estimation. In *CVPR*, 2019. 3
- [23] Mark J. Huiskes and Michael S. Lew. The mir flickr retrieval evaluation. In *MIR ’08: Proceedings of the 2008 ACM International Conference on Multimedia Information Retrieval*. ACM, 2008. 5
- [24] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *NeurIPS*, 2015. 5
- [25] Steve T. K. Jan, Joseph Messou, Yen-Chen Lin, Jia-Bin Huang, and Gang Wang. Connecting the digital and physical world: Improving the robustness of adversarial attacks. In *AAAI*, 2019. 3
- [26] Kensei Jo, Mohit Gupta, and Shree K. Nayar. Disco: Display-camera communication using rolling shutter sensors. *ACM Trans. Graph.*, 2016. 2
- [27] X. Kang, J. Huang, and W. Zeng. Efficient general print-scanning resilient data hiding based on uniform log-polar mapping. *IEEE TIFS*, 2010. 2
- [28] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016. 3
- [29] Yue Liu, Ju Yang, and Mingjun Liu. Recognition of qr code with mobile phones. In *Chinese Control and Decision Conference*. IEEE, 2008. 3
- [30] Emi Myodo, Shigeyuki Sakazawa, and Yasuhiro Takishima. Method, apparatus and computer program for embedding barcode in color image, 2013. US Patent 8,550,366. 3
- [31] Takao Nakamura, Atsushi Katayama, Masashi Yamamuro, and Noboru Sonehara. Fast watermark detection scheme from camera-captured images on mobile phones. *IJPRAI*, 2006. 2
- [32] Eisaku Ohbuchi, Hiroshi Hanaizumi, and Lim Ah Hock. Barcode readers using the camera device in mobile phones. In *International Conference on Cyberworlds*. IEEE, 2004. 3
- [33] Shelby Pereira and Thierry Pun. Robust template matching for affine resistant image watermarks. *IEEE Transactions on Image Processing*, 2000. 2
- [34] Tomáš Pevný, Tomáš Filler, and Patrick Bas. Using high-dimensional image models to perform highly undetectable steganography. In *International Workshop on Information Hiding*, 2010. 2
- [35] Anu Pramila, Anja Keskinarkaus, and Tapio Seppänen. Watermark robustness in the print-cam process. In *IASTED SPPRA*, 2008. 2
- [36] Anu Pramila, Anja Keskinarkaus, and Tapio Seppänen. Toward an interactive poster using digital watermarking and a mobile phone camera. *Signal, Image and Video Processing*, 2012. 2

- [37] Anu Pramila, Anja Keskinarkaus, and Tapio Sepp  nen. Increasing the capturing angle in print-cam robust watermarking. *Journal of Systems and Software*, 135:205–215, 2018. [2](#)
- [38] Francisco Romero Ramirez, Rafael Muoz-Salinas, and Rafael Medina-Carnicer. Speeded up detection of squared fiducial markers. *Image and Vision Computing*, 2018. [3](#)
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*. Springer International Publishing, 2015. [4](#)
- [40] Richard Shin and Dawn Song. Jpeg-resistant adversarial images. In *NeurIPS Workshop on Machine Learning and Computer Security*, 2017. [3, 4](#)
- [41] Chawin Sitawarin, Arjun Nitin Bhagoji, Arsalan Moseinia, Mung Chiang, and Prateek Mittal. Darts: Deceiving autonomous cars with toxic signs. *arXiv preprint arXiv:1802.06430*, 2018. [3](#)
- [42] K. Solanki, U. Madhow, B. S. Manjunath, S. Chandrasekaran, and I. El-Khalil. ‘Print and scan’ resilient data hiding in images. *IEEE TIFS*, 2006. [2](#)
- [43] Weixuan Tang, Shunquan Tan, Bin Li, and Jiwu Huang. Automatic steganographic distortion learning using a generative adversarial network. *IEEE Signal Processing Letters*, 2017. [2](#)
- [44] Eric Wengrowski and Kristin Dana. Light field messaging with deep photographic steganography. In *CVPR*, 2019. [2, 6, 7, 11](#)
- [45] Eric Wengrowski and Kristin Dana. Light field messaging with deep photographic steganography. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1515–1524, 2019. [2](#)
- [46] E. Wengrowski, W. Yuan, K. J. Dana, A. Ashok, M. Gruteser, and N. Mandayam. Optimal radiometric calibration for camera-display communication. In *WACV*, 2016. [2](#)
- [47] Pin Wu, Yang Yang, and Xiaoqiang Li. Stegnet: Mega image steganography capacity with deep convolutional network. *Future Internet*, 2018. [2](#)
- [48] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *ECCV*, 2018. [5, 12](#)
- [49] W. Yuan, K. Dana, M. Varga, A. Ashok, M. Gruteser, and N. Mandayam. Computer vision methods for visual mimo optical system. 2011. [2](#)
- [50] W. Yuan, K. J. Dana, A. Ashok, M. Gruteser, and N. Mandayam. Spatially varying radiometric calibration for camera-display messaging. In *GlobalSIP*, 2013. [2](#)
- [51] Kevin Alex Zhang, Alfredo Cuesta-Infante, and Kalyan Veeramachaneni. Steganogan: Pushing the limits of image steganography. *arXiv preprint arXiv:1901.03892*, 2019. [2](#)
- [52] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. [5, 7](#)
- [53] D. Zheng, J. Zhao, and A. El Saddik. Rst-invariant digital image watermarking based on log-polar mapping and phase correlation. *IEEE TCSV*, 2003. [2](#)
- [54] Jiren Zhu, Russell Kaplan, Justin Johnson, and Li Fei-Fei. Hidden: Hiding data with deep networks. In *ECCV*, 2018. [2, 3, 6, 7, 11](#)

A. StegaStamp Examples

See Figure 9 for additional examples of encoded images and their residuals.

B. Supplemental Videos

<https://youtu.be/E8OqgNDBGO0>

This video provides an overview of StegaStamps with example use cases and a condensed demonstration of in-the-wild results.

<https://youtu.be/jpbRhOH3D9Y>

This video is a compilation of multiple in-the-wild captures. The first set of clips visualizes the output bounding polygons along with the percentage of bits recovered correctly out of 100. We filter the output to only show detections where the bit accuracy is greater than 70 percent. We note that the messages are regularly recovered with greater than 90% accuracy when they are properly detected. The second set of clips demonstrates the used of BCH error correction [6] to robustly detect and correct recovered codes. The transmitted data consists of 56 message bits and 40 error correcting bits. When the accuracy is greater than 95% (fewer than 5 corrupted bits), the original 56-bit message can be recovered exactly. If too many bits are corrupted, the error correcting fails and we filter out the proposal. The video represents successfully decoded StegaStamps with green polygons. The decoded code is printed above the polygon. Note that for most real world applications, it is only necessary to recover the code in a single video frame to count it as successfully scanned.

C. Comparison Details

We compare our method to Baluja [5], HiDDeN [54], and LFM [44]. Baluja was designed to hide images within images, which differs from our task of hiding a bitstring within an image. To account for this, we convert our 100 bit message into a 10×10 grid of ones and zeros that is upscaled to the resolution of the cover image. During decoding we round the model output to 0 and 1 and take the mode within each upscaled block. As the original model was trained to hide natural images, we retrain the model from scratch to hide our bitstring grids.

HiDDeN was trained to hide 30 bit messages in 128×128 pixel images. We observed a significant drop in accuracy when we trained a model to hide 100 bit messages in 400×400 pixel images, therefore we report accuracy results on the 30 bit in 128^2 image version.

LFM [44] was trained to encode 1024 bit messages as 4×4 pixel blocks in a 256×256 pixel image. To encode our 100 bit message, we allocated 9 blocks for each message bit (we therefore only use a 244×244 pixel subset of the

		Mean Acc. \uparrow	bits/MP \uparrow
	Baluja [5]	0.51	0.5
	HiDDeN [54]	0.65	125
	LFM [44] (printed)	0.61	287
	LFM [44] (screen)	0.93	1109
Ours	None	0.49	0.1
	Pixelwise	0.51	0.2
	Spatial	0.89	318
	All	0.99	571

Table 3: Quantitative comparison of other methods and our ablations. We show numbers in terms of fraction of bits correctly recovered (mean accuracy) as well as bits-per-megapixel (bits/MP). Higher is better for both metrics. The bits/MP metric normalizes the message length and image sizes between different methods. All methods except “LFM [44] (screen)” (cellphone camera/cellphone screen) are reported on the cellphone camera/consumer printer pipeline. We report LFM’s results in this additional case because it was explicitly designed for screen/camera transmission.

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Baluja [5]	24.61	0.926	0.256
HiDDeN [54] (native)	31.07	0.940	0.070
HiDDeN [54]	24.55	0.775	0.202
LFM [44]	20.89	0.910	0.315
Ours	27.25	0.927	0.194

Table 4: Quantitative comparison of encoded image quality, indicating how well hidden the message is. For HiDDeN [54] we show both the metrics for the original lower resolution (native 128×128) and upsampling to our compared resolution of 400×400 with bicubic interpolation. At full resolution, our method produces an encoded image most similar to the original in all metrics.

image). We average and round the 9 block predictions to recover the message bit.

Each compared method encodes a different length message into a different size image. However, if we treat the mean bit recovery accuracy (first column in Table 3) as the crossover probability p in a binary symmetric channel, we can use information theory to calculate the channel capacity (with unit “bits”):

$$C(p) = 1 - (-p \log_2 p - (1-p) \log_2 1-p) \quad (3)$$

If we divide $C(p)$ by the number of pixels N_{pix} in the original image, we get the expected number of bits-per-pixel transmitted by that method. Multiplying $\frac{C(p)}{N_{pix}}$ by 10^6 yields our bits-per-megapixel metric in the second column of Table 3.



Figure 9: Additional examples of encoded images and their residuals.

D. Architecture Details

Network architectures for our encoder (Table 5) and decoder (Table 6). Our detector uses the BiSeNet [48] architecture.

E. Code

The code and pretrained networks can be found at <https://github.com/tancik/StegaStamp>.

Layer	k	s	chns	in	out	input
inputs			6			image + secret
conv1	3	1	6/32	1	1	inputs
conv2	3	2	32/32	1	2	conv1
conv3	3	2	32/64	2	4	conv2
conv4	3	2	64/128	4	8	conv3
conv5	3	2	128/256	8	16	conv4
up6	2	1	256/128	16	8	upsample(conv5)
conv6	3	1	256/128	8	8	conv4 + up6
up7	2	1	128/64	8	4	upsample(conv6)
conv7	3	1	128/64	4	4	conv3 + up7
up8	2	1	64/32	4	2	upsample(conv7)
conv8	3	1	64/32	2	2	conv2 + up8
up9	2	1	32/32	2	1	upsample(conv8)
conv9	3	1	70/32	1	1	conv1 + up9 + inputs
conv10	3	1	32/32	1	1	conv9
residual	1	1	32/3	1	1	conv10

Table 5: Our encoder network architecture. **k** is the kernel size, **s** the stride, **chns** the number of input and output channels for each layer, **in** and **out** are the accumulated stride for the input and output of each layer, **input** denotes the input of each layer with + meaning concatenation and “upsample” performing $2\times$ nearest neighbor upsampling. A ReLU is applied after each layer except the last.

Layer	k	s	chns	in	out	input
conv1	3	2	3/32	1	2	image
conv2	3	2	32/64	2	4	conv1
conv3	3	2	64/128	4	8	conv2
fc0			320000			flatten(conv3)
fc1			320000/128			fc0
fc2			128/6			fc1
image_warped			3/3			transf(image, fc2)
conv1	3	2	3/32	1	2	image_warped
conv2	3	1	32/32	2	2	conv1
conv3	3	2	32/64	2	4	conv2
conv4	3	1	64/64	4	4	conv3
conv5	3	2	64/64	4	8	conv4
conv6	3	2	64/128	8	16	conv5
conv7	3	2	128/128	16	32	conv6
fc0			20000			flatten(conv7)
fc1			20000/512			fc0
secret			512/100			fc1

Table 6: Our decoder network architecture. We indicate convolutional layers with the prefix “conv” and fully connected layers with the prefix “fc.” The first half of the network outputs an affine warp that is applied using a differentiable spatial transformer layer (“transf”). The warped result is fed into the second part of the network. A ReLU is applied after each layer except the last layer before the spatial transformer.