# IMAGE INPAINTING FOR HIGH-RESOLUTION TEXTURES USING CNN TEXTURE SYNTHESIS

*Pascal Laube, Michael Grunwald, Matthias O. Franz & Georg Umlauf*

Institute for Optical Systems, University of Applied Sciences Constance, Germany

plaube@htwg-konstanz.de

## ABSTRACT

Deep neural networks have been successfully applied to problems such as image segmentation, image super-resolution, coloration and image inpainting. In this work we propose the use of convolutional neural networks (CNN) for image inpainting of large regions in high-resolution textures. Due to limited computational resources processing high-resolution images with neural networks is still an open problem. Existing methods separate inpainting of global structure and the transfer of details, which leads to blurry results and loss of global coherence in the detail transfer step. Based on advances in texture synthesis using CNNs we propose patch-based image inpainting by a CNN that is able to optimize for global as well as detail texture statistics. Our method is capable of filling large inpainting regions, oftentimes exceeding the quality of comparable methods for high-resolution images. For reference patch look-up we propose to use the same summary statistics that are used in the inpainting process.

*Index Terms*— Image inpainting, CNN, texture synthesis, high-resolution

## 1. INTRODUCTION

Image inpainting is the process of filling missing or corrupted regions in images based on surrounding image information so that the result looks visually plausible. Most image inpainting approaches are based on sampling existing information surrounding the inpainting region, wich is called exemplar-based [1–5] inpainting. Recently machine learning techniques have been applied successfully to the problem of texture synthesis and inpainting [6–9]. First introduced by Gatys et al. in [10] texture synthesis CNNs have been shown to surpass well-known methods like the one by Portilla et al. [11] for many textures. Wallis et al. [12] recently showed that artificial images produced from a parametric texture model closely match texture appearance for humans. Especially, the CNN texture model of [10] and the extension by Liu [13] are able to capture important aspects of material perception in humans. For many textures the synthesis results are indiscriminable under foveal inspection. Other methods like the ones by Phatak et al. [14] and Yang et al. [15]

train auto-encoder-like networks, called context-encoders, for inpainting. Inpainting methods using neural networks still suffer two main drawbacks: Due to limited computational resources they are restricted to small inpainting regions and results often lack details and are blurry. For high-resolution textures the inpainting result not only needs to reproduce texture details but also global structure. Applying details after a first coarse inpainting step distorts global statistics. Fig. 1 shows some examples where well-known inpainting methods fail to reproduce global and local structure.
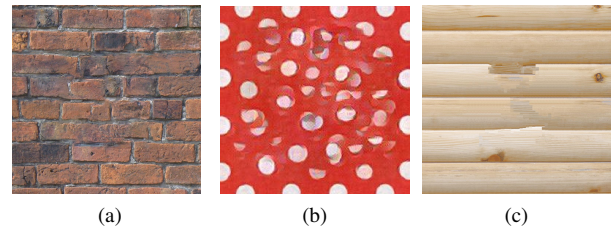


**Fig. 1**: Inpainting results for some of the example textures from Fig. 5 using the methods of (a) Photoshop CS7 which is a combination of methods [5] and [2], (b) the method by Yang et al. [15], (c) and by Criminisi et al. [1].

To resolve the outlined issues we propose an inpainting approach that produces results that reproduce global statistics and contain blur-free details. We fill the inpainting region by synthesis of new texture patch by patch, which enables us to process high-resolution textures. Our inpainting approach creates a smooth transition between the sampling and the inpainting region as well as between patches. Out setup is able to shift focus from optimizing detail to global statistics on different levels of resolutions.

Sections of this paper are arranged as follows. The process of texture synthesis by CNNs is then explained in Sec. 2. In Sec. 3 we present our inpainting approach, followed by an experimental evaluation in Sec. 4. We conclude in Sec. 5.

## 2. TEXTURE SYNTHESIS

First introduced by Gatys et al. [10] CNN texture synthesis uses summary statistics derived from filter responses of convolutional layers, the feature maps, to synthesize new texture.

In a first step some vectorized texture $\mathbf{x}$ of size $P$ is presented to the analysis CNN. Based on the resulting feature maps one can compute the Gramians which are spatial summary statistics. The Gramian of some network layer $l$ is defined as

$$\mathbf{G}^l_{ij} = \sum_k \mathbf{F}^l_{ik} \mathbf{F}^l_{jk}, \qquad (1)$$

where $\mathbf{F}^l_{ik}$ is feature map $i$ and $\mathbf{F}^l_{jk}$ feature map $j$ at location $k$ given input $\mathbf{x}$. These inner products of filter activations of different layers are then used to define a synthesis loss

$$\mathcal{L}_s(\mathbf{x}, \hat{\mathbf{x}}) = \sum_{l=0}^{L} \frac{1}{2N_l^2 M_l^2} \sum_{i,j} (\mathbf{G}^l_{ij} - \hat{\mathbf{G}}^l_{ij})^2,$$

with $N_l$ feature maps of size $M_l$ at layer $l$. Here $\hat{\mathbf{G}}^l_{ij}$ are the Gramians of a synthesis CNN. Based on this loss some randomly initialized input vector $\hat{\mathbf{x}}$ of the synthesis CNN is optimized to satisfy statistics derived from the analysis CNN. Since Gramians average over feature map positions this leads to a loss of global texture coherence. Berger and Memisevic [16] introduce a second cross-correlation loss by computing Gramians between feature maps $\mathbf{F}^l$ and a spatial translation $T$ of the feature maps $T(\mathbf{F}^l)$. By discarding either rows or columns of feature maps one can now compute correlations of features at some location $k = (x, y)$ and a shifted location $T_{x,+\delta}(k) = (x + \delta, y)$ or $T_{y,+\delta}(k) = (x, y + \delta)$. The horizontally translated Gramian becomes

$$G^l_{x,\delta,ij} = \sum_k T_{x,+\delta}(F^l_{ik}) T_{x,-\delta}(F^l_{jk}), \qquad (2)$$

and $G^l_{y,\delta,ij}$ analogous. The cross-correlation loss $\mathcal{L}_{cc}$ for an arbitrary shift $\delta$ is defined as

$$\mathcal{L}_{cc}(\mathbf{x}, \hat{\mathbf{x}}) = \sum_{l,i,j} \frac{(\mathbf{G}^l_{x,\delta,ij} - \hat{\mathbf{G}}^l_{x,\delta,ij})^2 + (\mathbf{G}^l_{y,\delta,ij} - \hat{\mathbf{G}}^l_{y,\delta,ij})^2}{4N_l^2 M_l^2}.$$

The combined loss is then defined as

$$\mathcal{L}_{s,cc}(\mathbf{x}, \hat{\mathbf{x}}) = w_s \mathcal{L}_s + w_{cc} \mathcal{L}_{cc},$$

with weight factors $w_s$ and $w_{cc}$.

## 3. PATCH-BASED TEXTURE SYNTHESIS FOR IMAGE INPAINTING

### 3.1. Patch-based texture synthesis

Given some image with high-resolution, uncorrupted texture $\Phi$ we propose the application of the synthesis method introduced in Sec. 2 on different scales of resolution to fill the inpainting region $\Omega$ (Fig. 3a). A schematic overview of our setup is given in Fig. 2. We propose to inpaint region $\Omega$ patch by patch with each patch satisfying global as well as detail

statistics. For this purpose, we define a texture loss function that simultaneously evaluates the quality of the synthesized patch $\hat{\mathbf{x}}_d$ in native resolution as well as the quality of an embedding of $\hat{\mathbf{x}}_d$ into a pooled window of its surroundings $\hat{\mathbf{x}}_g$ capturing global information. $\hat{\mathbf{x}}_g$ is initialized with a $Q$-times average-pooled window of the image so that this window fully contains $\Omega$ and the boundary $\Psi$. $Q$ average-pooling-layers are introduced in-between $\hat{\mathbf{x}}_d$ and $\hat{\mathbf{x}}_g$ so that $\hat{\mathbf{x}}_d$ can become a subtensor of $\hat{\mathbf{x}}_g$ at the correct (pooled) position. Depending on the size of $\Omega$, $Q$ needs to be adjusted as a parameter before inpainting. Before generating the next patch $\hat{\mathbf{x}}_d$ at a new location we update $\Omega$ with the synthesis result in $\hat{\mathbf{x}}_d$ and reinitialize $\hat{\mathbf{x}}_g$. Only $\hat{\mathbf{x}}_d$ is optimized in the synthesis process.

For the synthesis as described in Sec. 2, suitable reference textures $\mathbf{x}_d$ and $\mathbf{x}_g$ are needed. We will describe the reference patch look-up in Sec. 3.2. While $\mathbf{x}_g$ needs to be initialized only once at the beginning of the inpainting process, $\mathbf{x}_d$ is reinitialized with a new reference for every new position of the inpainting patch $\hat{\mathbf{x}}_d$. We further define a boundary loss, that limits the optimization of region $\Psi$ inside $\hat{\mathbf{x}}_d$ in the input domain. We define the boundary loss as

$$\mathcal{L}_b(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{P} \sum (\mathbf{mx} - \mathbf{m}\hat{\mathbf{x}})^2, \qquad (3)$$

where the binary mask $\mathbf{m}$ equals 0, if $\mathbf{m}_i \in \Omega$, and 1 otherwise. The combined loss over both branches together with boundary loss becomes

$$\mathcal{L} = w_d \mathcal{L}_{s,cc}(\mathbf{x}_d, \hat{\mathbf{x}}_d) + w_g \mathcal{L}_{s,cc}(\mathbf{x}_g, \hat{\mathbf{x}}_g) + w_b \mathcal{L}_b(\mathbf{x}_b, \hat{\mathbf{x}}_d),$$

where $w_d$, $w_g$, and $w_b$ are weight terms. $\mathbf{x}_b$ is initialized with $\hat{\mathbf{x}}_d$ before optimization and does change for each new position of $\hat{\mathbf{x}}_d$.

### 3.2. Patch distance by Gramians

For the synthesis of patch $\hat{\mathbf{x}}_d$, suitable reference patches $\mathbf{x}_d$, and $\mathbf{x}_g$ are needed. The initial $\hat{\mathbf{x}}_d$ is a window of the image containing parts of $\Psi$ as well as parts of $\Omega$ while $\hat{\mathbf{x}}_g$ completely contains $\Psi$ and $\Omega$. One now has to find closest patches from $\Phi$ matching $\Psi$ inside $\hat{\mathbf{x}}_d$, and $\hat{\mathbf{x}}_g$ as candidates for $\mathbf{x}_d$, and $\mathbf{x}_g$. Instead of the MSE, we propose to use the distance of texture Gramians as a similarity measure. Since values inside $\Omega$ are unknown we propose masking $\Omega$ for each individual feature map to remove $\Omega$-related correlations from any of the resulting Gramians. Because the network input up to some layer $l$ has passed through both pooling and convolutional layers we need to adapt the feature map masks to compensate for these operations. In a first step, the initial binary mask $\mathbf{m}$ from Eq. (3) needs to be adapted in size to account for the pooling layers. This is done by applying each pooling step of the CNN that has been applied up to layer $l$ to the mask $\mathbf{m}^l$ which is responsible for masking feature maps $\mathbf{F}^l$. In a second step, one needs to account for the propagation of $\Omega$ into $\Psi$ by convolutions. Masks $\mathbf{m}^l$ also need to account for
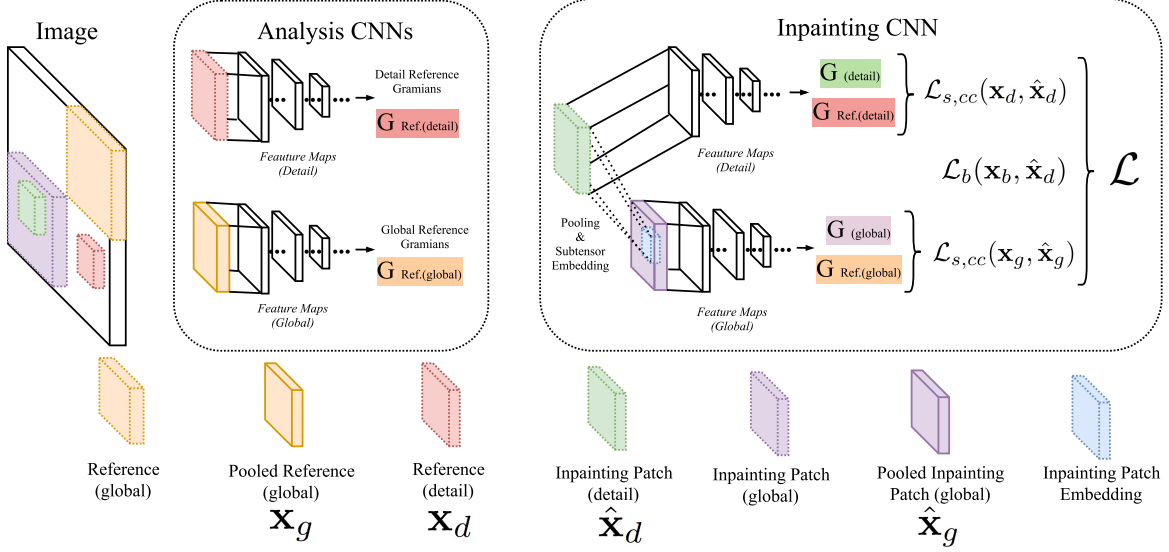
**Fig. 2**: Scheme of our proposed inpainting setup. On the top left the inpainting image together with important image regions is shown. Under "Analysis CNN" the generation of detail as well as reference Gramians is shown. On the top right our "Inpainting CNN" together with the resulting loss terms is shown. The inpainting patch $\hat{\mathbf{x}}_d$ is input to the *detail branch* (top) as well as, after embedding, the *global branch* (bottom). A legend of the involved image elements is given at the bottom.

propagation of $\Omega$ into $\Psi$ due to convolution. Simply discarding the affected values by setting them to zero in $\mathbf{m}^l$ for each convolutional layer is too restrictive and would lead to masks with all values zero in later layers. We propose to expand $\Omega$ by a smaller individual number of pixels $\mathbf{e}^l$ for each convolutional layer (see Sec. 4). In our experiments this expansion has proven to be sufficient for compensation.

Taking these considerations into account we define our patch distance as

$$\Delta_{\mathbf{G}}(\mathbf{x}, \hat{\mathbf{x}}) = \sum_{l,i,j} (\sum_k \mathbf{m}^l \mathbf{F}_{ik}^l \mathbf{F}_{jk}^l - \sum_k \mathbf{m}^l \hat{\mathbf{F}}_{ik}^l \hat{\mathbf{F}}_{jk}^l)^2.$$

### 3.3. Inpainting

For inpainting we propose a coarse to fine inpainting process with two steps. At each stage $\hat{\mathbf{x}}_d$ is optimized by applying L-BFGS-B [17]. We initialize each color channel in region $\Omega$ with the corresponding color channel mean from $\Phi$. In the coarse inpainting step we focus on optimizing global statistics by setting $w_d = 0$, $w_g = 1$. This leads to $\hat{\mathbf{x}}_d$ satisfying global statistics but at low resolution. Pooling larger input regions introduces color artifacts since loss is shared among pooled pixels as can be seen in Fig. 3b. We eliminate these color artifacts by converting $\Omega$ to greyscale (see Fig. 3c) with RGB weights $r = 0.212$, $g = 0.7154$, $b = 0.0721$. Only this structure is used for initialization of the second stage. In the fine inpainting step we set $w_d = 1$ and $w_g$ to a value in the range of $[0.01, 0.1]$. This ensures focus on the optimization for detail statistics through the *detail branch* while constraining the optimization to also maintain global texture statistics.

For our approach inpainting order is not important as long as

the first patch overlaps with $\Psi$ and consecutive patches overlap. Overlapping a patch by $\frac{1}{4}$ of its own size with surrounding texture has proven to be sufficient for smooth boundary transition. We chose to fill $\Omega$ in a top to bottom, left to right fashion. To ensure a smooth transition in-between patches we apply image quilting [18] on overlaps. As a result of our experiments we set $w_s = 1e6$ and $w_{cc} = 1e7$ for inpainting of 8-bit color images. Choosing $w_b$ in the range $[5, 25]$ has shown to be sufficient. The large difference between Gramian-based loss weights and weights related to loss in pixel space results from different value ranges.
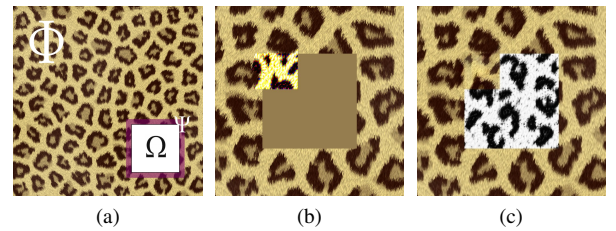


**Fig. 3**: (a) Example image (2048x2048px) with inpainting region $\Omega$, boundary $\Psi$ and texture $\Phi$. (b) First patch of the coarse inpainting step. (c) Fine inpainting of texture detail after coarse inpainting.

## 4. EXPERIMENTAL EVALUATION

We present inpainting results of exemplar high resolution textures. All textures have a resolution of 2048x2048px while the inpainting region $\Omega$ is of size 512x512px. We use ImageNet pre-trained VGG-19 CNNs for analysis as well as syn-

thesis with input size 256x256px. We use layers *conv1_1*, *pool1*, *pool2*, *pool3* and *pool4* for computing global as well as detail statistics. For very stochastic textures we propose to use *pool3*, *pool4* and *pool5* to compute global statistics since this leads to improved texture scale in the coarse inpainting step. For patch distance computation we define pixel expansions $\mathbf{e} = (1, 1, 2, 3, 2)$, and for shift $\delta$ of translated Gramians $\mathbf{G}^l_{x,\delta}$ and $\mathbf{G}^l_{y,\delta}$ we define $\boldsymbol{\delta} = (6, 6, 5, 4, 3)$. We use $Q = 2$ pooling layers. To find suitable reference patches $\mathbf{x}_d$ and $\mathbf{x}_g$ region $\Phi$ is searched at a step size of 64px. Inpainting of the exemplar textures was done using a Nvidia GeForce 1080 Ti and took roughly 8 min strongly depending on the number of iterations of the L-BFGS-B optimization. In Fig. 6 we present results of our inpainting approach for inpainting $\Omega$ of the example textures in Fig. 5. While many methods have difficulties maintaining global as well as local texture characteristics our results look reasonable on both scales. Using the difference of masked Gramians as a metric for patch distance has major benefits for our inpainting approach over using simple MSE. Since we are not dependent on reference textures $\mathbf{x}_d$ or $\mathbf{x}_g$ exactly matching $\Psi$ inside $\hat{\mathbf{x}}_d$ or $\hat{\mathbf{x}}_g$ in terms of MSE, we can reduce the number of samples taken from $\Phi$ in reference patch look-up. Due to the averaging of feature information inside Gramians, global spacial information is lost. This enables the Gramian to represent texture invariant to rotation and translation to some degree (see Fig. 4). Because our loss term $\mathcal{L}$ is based on the difference of Gramians this further ensures that $\Psi$ inside $\hat{\mathbf{x}}$ already satisfies target statistics to some extent. When choosing $w_d$ and $w_g$ one needs
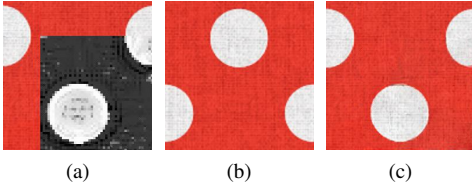


**Fig. 4**: (a) Inpainting patch $\hat{\mathbf{x}}_d$. (b) Closest reference patch from $\Phi$. (c) Inpainting result.

to be aware of the trade-off introduced. While higher $w_g$ ensures persistence of global statistics it also introduces artifacts as a result of pooling $\hat{\mathbf{x}}_d$ before subtensor embedding and vice versa. Higher $w_d$ lays larger emphasis on details while possibly violating global structure. This trade-off is further influenced by the number of poolings $Q$.

## 5. CONCLUSION

In this work, we presented a new CNN-based method for inpainting that can be applied to large-scale, high-resolution textures. Texture analysis and inpainting are done on two scales, one for global structure and one for details. This avoids the problems of blurry or missing details from which previous CNN approaches suffered while plausibly continu-



**Fig. 5**: Examples for evaluation with inpainting region $\Omega$.



**Fig. 6**: Closeup on results using our method to inpaint region $\Omega$ from Fig. 5.

ing global image structure. In principle, our network architecture can be extended to include a hierarchy of more than two interacting scales. The design of such a multi-resolution architecture could be an interesting line of research that we plan to pursue in the future.

# 6. REFERENCES

[1] Antonio Criminisi, Patrick Pérez, and Kentaro Toyama, "Region filling and object removal by exemplar-based image inpainting," *IEEE Transactions on image processing*, vol. 13, no. 9, pp. 1200–1212, 2004.

[2] Yonatan Wexler, Eli Shechtman, and Michal Irani, "Space-time completion of video," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 29, no. 3, 2007.

[3] Vivek Kwatra, Irfan Essa, Aaron Bobick, and Nipun Kwatra, "Texture optimization for example-based synthesis," *ACM Transactions on Graphics (ToG)*, vol. 24, no. 3, pp. 795–802, 2005.

[4] Alexei A Efros and Thomas K Leung, "Texture synthesis by non-parametric sampling," in *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*. IEEE, 1999, vol. 2, pp. 1033–1038.

[5] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph.*, vol. 28, no. 3, pp. 24–1, 2009.

[6] Chuan Li and Michael Wand, "Combining markov random fields and convolutional neural networks for image synthesis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2479–2486.

[7] Leon A Gatys, Alexander S Ecker, and Matthias Bethge, "Image style transfer using convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2414–2423.

[8] Justin Johnson, Alexandre Alahi, and Li Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision*. Springer, 2016, pp. 694–711.

[9] Alexey Dosovitskiy and Thomas Brox, "Generating images with perceptual similarity metrics based on deep networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 658–666.

[10] Leon Gatys, Alexander S Ecker, and Matthias Bethge, "Texture synthesis using convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2015, pp. 262–270.

[11] Javier Portilla and Eero P Simoncelli, "A parametric texture model based on joint statistics of complex wavelet coefficients," *International journal of computer vision*, vol. 40, no. 1, pp. 49–70, 2000.

[12] Thomas SA Wallis, Christina M Funke, Alexander S Ecker, Leon A Gatys, Felix A Wichmann, and Matthias Bethge, "A parametric texture model based on deep convolutional features closely matches texture appearance for humans," *Journal of Vision*, vol. 17, no. 12, pp. 5–5, 2017.

[13] Gang Liu, Yann Gousseau, and Gui-Song Xia, "Texture synthesis through convolutional neural networks and spectrum constraints," in *Pattern Recognition (ICPR), 2016 23rd International Conference on*. IEEE, 2016, pp. 3234–3239.

[14] 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*, 2016, pp. 2536–2544.

[15] Chao Yang, Xin Lu, Zhe Lin, Eli Shechtman, Oliver Wang, and Hao Li, "High-resolution image inpainting using multi-scale neural patch synthesis," *arXiv preprint arXiv:1611.09969*, 2016.

[16] Guillaume Berger and Roland Memisevic, "Incorporating long-range consistency in cnn-based texture generation," *arXiv preprint arXiv:1606.01286*, 2016.

[17] Ciyou Zhu, Richard H Byrd, Peihuang Lu, and Jorge Nocedal, "Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization," *ACM Transactions on Mathematical Software (TOMS)*, vol. 23, no. 4, pp. 550–560, 1997.

[18] Alexei A Efros and William T Freeman, "Image quilting for texture synthesis and transfer," in *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*. ACM, 2001, pp. 341–346.