

Autoregressive Image Interpolation via Context Modeling and Multiplanar Constraint

Shihong Deng, Jiaying Liu*, Mading Li, Wenhang Yang and Zongming Guo

Institute of Computer Science and Technology, Peking University, Beijing, P.R. China, 100871

Abstract—In this paper, we propose a novel image interpolation algorithm by context-aware autoregressive (AR) model and multiplanar constraint. Different from existing AR based methods which employ predetermined reference configuration to predict pixel values, the proposed method considers the anisotropic pixel dependencies in natural images and adaptively chooses the optimal prediction context by utilizing the nonlocal redundancy to interpolate pixels. Furthermore, the multiplanar constraint is applied to enhance the correlations within the estimation window by exploiting the self-similarity property of natural images. Similar patches are collected by the combination of patch-wise pixel values and the gradient information. And the inter-patch dependencies are adopted to improve the interpolation. The experimental results show that our method is effective in image interpolation and successfully decreases the artifacts nearby the sharp edges. The comparison experiments demonstrate that the proposed method can obtain better performance than other related ones in terms of both objective and subjective results.

Index Terms—Autoregressive(AR), context modeling, image interpolation, multiplanar, patch similarity

I. INTRODUCTION

Image interpolation aims to generate a high-resolution (HR) image by utilizing the information of a low-resolution (LR) counterpart. Recently it has become a hot research topic in the area of image processing for the wide applications such as video communication, digital photography enhancement, medical analysis and consumer electronics.

Basically, there are two key steps in image interpolation: one is the proper modeling that can easily characterize the image texture and edge information, while the other is the accurate estimation method free from overfitting and visual degradation, *e.g.*, jaggling and ringing. According to different modelings and estimation methods, image interpolation methods can be classified into different categories. Polynomial-based methods, such as Bilinear, Bicubic [1], and Cubic Spline [2], generate results by convolving neighboring pixels with a predetermined kernel for each pixels individually. Although this kind of method is easy to implement and of low computational complexity, it treats every pixel identically and ignores various local structure of natural images, and thus leads to undesirable results. The studies in [3], [4] use the structural information by explicitly detecting edge direction, gradient and isophotes of image as guidance for image interpolation. Meanwhile, AR-based methods are developed for image interpolation. The new edge-directed interpolation (NEDI) [5] is an early work of AR-based methods, which exploits the geometric duality between LR covariance and HR covariance to obtain HR pixels. Based on NEDI, Soft-decision adaptive interpolation (SAI) [6] introduces an additional AR model from cross direction and performs so-called block estimation in which all HR pixels in a local

window are estimated simultaneously. In [7], [8] and [9], similarity modulated model is proposed and different similarity metrics, such as patch-geodesic distance are used to improve the performance of the block estimation for AR model. Furthermore, [8] extends the AR model to general scale interpolation by solving the problem iteratively. For these existing AR-based image interpolation methods, they employ the fixed model reference such as diagonal and cross direction pixels and ignore the complicated and anisotropic dependencies within natural images. In [10] and [11], sparse representation dictionary learning methods have been proposed in order to better utilize the nonlocal dependency in natural images. Although these methods achieve desirable interpolation performance, the space and computational complexity of dictionary-based interpolation methods are much higher compared to traditional ones. Also, dictionary-based methods results are significantly affected by sufficient nonlocal similar patches within external database or the image itself. Performances will degrade greatly when this condition is not satisfied.

In this paper, to overcome the drawbacks of existing AR-based and dictionary-based image interpolation methods, we propose a novel AR-based framework for image interpolation by considering the spatial configuration of the model reference. The proposed image interpolation method adaptively selects optimal reference pixels according to the context of the given position. The reference pixel selecting process incorporates global image correlation to better characterize local structure of the image. A multiplanar constraint is introduced to fully exploit the correlations of all scales within the local estimation window. Finally, the similarity modulated block estimation is deployed based on the patch-geodesic distance, certifying the accuracy and stability of prediction. Comprehensive experimental results demonstrate that the proposed method via context modeling and multiplanar constraint is able to well model the piecewise stationary characteristic of natural images and achieve desirable performance from both objective and subjective perspectives.

The rest of the paper is organized as follows: Section II gives a brief review of the AR model and shortly introduces the concept of context modeling. Section III describes the details of the proposed interpolation algorithm. Experimental results and analysis are presented in Section IV. Finally, Section V concludes and remarks the whole paper.

II. CONTEXT AUTOREGRESSIVE (AR) MODEL

A. Autoregressive (AR) Model

Due to its ability to describe stochastic structure of sequential data, the autoregressive(AR) model is widely applied in statistic signal processing to model and predict various types of natural signals. Typically, a 2-D image signal $I(x, y)$ can be modeled as an AR process as follows,

$$I(x, y) = \sum_{(i,j) \in \Omega} \psi(i, j)I(x+i, y+j) + \epsilon(x, y), \quad (1)$$

where Ω and $\psi(i, j)$ represent the adjacent neighbors of pixel $I(x, y)$ and their corresponding model parameters (or weights), respectively; $\epsilon(x, y)$ denotes white noise. The formula shows that pixels in images

*Corresponding author

This work was supported by National High-tech Technology R&D Program (863 program) of China under Grant 2014AA015205 and National Natural Science Foundation of China under contract No. 61472001.

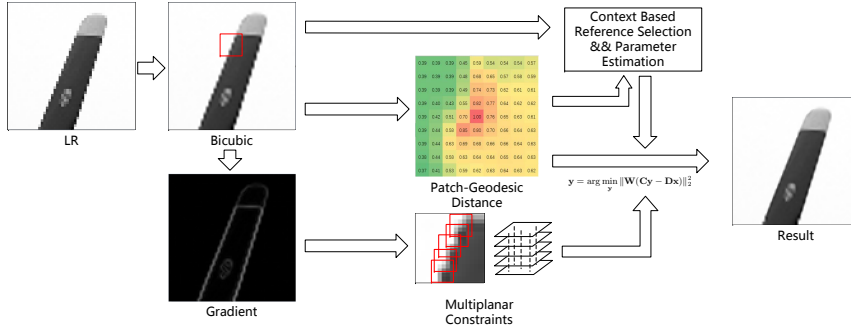


Fig. 1. Flow diagram of the proposed interpolation algorithm.

can be estimated by their adjacent neighbors with the corresponding weights.

B. Context Modeling Prediction

The piecewise based AR predictor introduced above ignores the importance of selecting the proper spatial reference configuration Ω for the model. In fact, most traditional AR models fix their model order and spatial reference pixels empirically, such as four neighbors from cross directions or four neighbors from diagonal directions. The spatial reference configuration Ω is usually a regular rectangular or circular region centered at the pixel to be predicted and takes exactly the literal meaning of adjacent neighbors. This kind of human intuitive configuration indeed follows the piecewise statistical stationary assumption of natural images: the closer a reference pixel is to the center pixel spatially, the more significant relation is between them. Although being valid for many locally consistent images, this configuration and assumption fail to be tenable in regions with various scales textures and anisotropic edges, which commonly and widely exist in natural images. Under this kind of circumstance, a manually predetermined reference configuration of the AR predictor is obviously suboptimal for that it may take adjacent but irrelevant pixels into consideration and miss further but relevant ones in the model estimation. Thus, it would lead to estimation biases.

To overcome the drawback, some existing studies try to model context awareness for AR model. In [12], Wu *et al.* proposed a method to utilize correlation instead of spatial distance between pixels to choose both predictor reference and training set of AR model and apply it in lossless image coding. In this paper, we adopt the concept of *context* to extend piecewise AR model and modify the predictor reference choosing mechanism and then propose a context modeling image interpolation algorithm.

III. THE PROPOSED INTERPOLATION ALGORITHM

In this section, we present the proposed image interpolation algorithm in detail. Firstly, a novel AR model based on context-aware modification is described. Secondly, we introduce the patch-geodesic distance to define the similarity between two pixels. Thirdly, a multiplanar constraint is proposed by exploiting the local structural information and incorporated into the AR model to improve the data fidelity term. Finally, with the above model and similarity, the missing HR pixels are jointly estimated by similarity modulated estimation. Fig. 1 shows the framework of the proposed algorithm. Specifically, a two-pass interpolation is performed as stated in [6].

A. Context-Aware Adaptive Prediction

Interpolation algorithms in [6]-[9] use two sets of AR model parameters $\mathbf{a} = \{a_t\}$ and $\mathbf{b} = \{b_t\}$ ($t = 1, 2, 3, 4$) to predict the missing HR pixels. Parameter \mathbf{a} indicates the weights of pixels from diagonal direction while \mathbf{b} indicates the ones from cross direction. These two sets of parameters are with fixed order and have fixed spatial reference configuration. The AR equations can be represented as:

$$z_i = \sum_{t=1}^4 a_t z_{i \otimes t} + \epsilon_i^{\otimes}, z_i = \sum_{t=1}^4 b_t z_{i \oplus t} + \epsilon_i^{\oplus}, \quad (2)$$

where z_i refers to either LR pixel x_i or HR pixel y_i in local window W ; $z_{i \otimes t}$ and $z_{i \oplus t}$ are reference pixels of z_i from diagonal and cross directions; ϵ_i^{\otimes} and ϵ_i^{\oplus} are the corresponding fitting error. The spatial reference configuration and corresponding mapping relation between LR and HR pixels are illustrated in Fig. 2

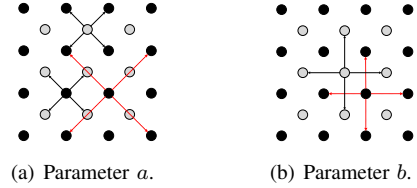


Fig. 2. The spatial reference configuration. Black and red arrows indicate HR-HR/LR and the corresponding LR-LR correlations, respectively.

To better model local image statistics, we extend the fixed spatial reference configurations of AR model to the changeable ones. By choosing the proper model order and spatial reference configuration, we can make the AR prediction more accurate and more stable. The problem is how to determine the model order and the spatial distribution of the reference pixels. To illustrate the problem, we use $\Phi(z_i) = \{\phi_1(z_i), \phi_2(z_i), \dots, \phi_M(z_i)\}$ and W_r to denote the set of possible reference pixels for pixel z_i and the corresponding patch window containing all candidate pixels. Instead of determining the reference merely on pixel-by-pixel geometric distance as previous methods, we firstly obtain z_i 's nearest neighbor set T in terms of patch W_r globally:

$$T = \{x \mid \|W_r(x) - W_r(z_i)\| \leq \tau_T\}, \quad (3)$$

where τ_T is a threshold. The distance type can be either Manhattan distance or Euclidean distance. The missing HR pixels in W_r are generated by bicubic method as initialization. Here we benefit from rich amount of local and nonlocal redundancy of natural images and the size of T should be no less than certain scale. Otherwise, the formula below is not so faithfully confident. Then, using these similar patches which are supposed to have homogeneous texture, we calculate the correlation coefficients between z_i and candidate reference pixels

$$\rho_{\phi_m} = \frac{\sum_{x \in T} \phi_m(x)x - \frac{1}{T} \sum_{x \in T} x \sum_{x \in T} \phi_m(x)}{\sqrt{\sum_{x \in T} x^2 - \frac{1}{T} (\sum_{x \in T} x)^2} \sqrt{\sum_{x \in T} \phi_m(x)^2 - \frac{1}{T} (\sum_{x \in T} \phi_m(x))^2}}. \quad (4)$$

With the calculated correlation coefficients ρ_{ϕ_m} , $1 \leq m \leq M$, we reorder the candidate reference set

$$\Phi(z_i) = \{\phi_1(z_i), \phi_2(z_i), \dots, \phi_M(z_i), \rho_{\phi_1} \geq \rho_{\phi_2} \geq \dots \rho_{\phi_M}\}. \quad (5)$$

With this ranking, we are able to easily pick the reference from candidate set sequentially with the increase of model order. Please note that there is a tradeoff between the estimation bias and the prevention of overfitting. A larger candidate set may bring in more highly correlative reference pixels. However, relative larger patch window naturally decreases the number of similar patch and reduces the confidence of the model. The larger patch window will also arise the problem of ignoring fine-grained textures and edges at the same

time. Thus, we just extend the candidate set from four neighbor pixels to eight pixels, and choose first four of them in accordance with the originally order of the AR model. Also, considering the scale transformation between HR-HR/LR and LR-LR correlations, we only apply the context modeling to parameter \mathbf{b} . The different reference scale between HR-HR/LR and LR-LR correlations of parameter \mathbf{a} may enlarge the fitting error. The context modeling AR equations are revised as follows:

$$z_i = \sum_{t=1}^4 a_t z_{i \otimes t} + \epsilon_i^{\otimes}, z_i = \sum_{t=1}^4 b_t z_{i \odot t} + \epsilon_i^{\odot}, \quad (6)$$

where \odot denotes context modeling spatial reference configuration. The process are clearly illustrated in Fig. 3

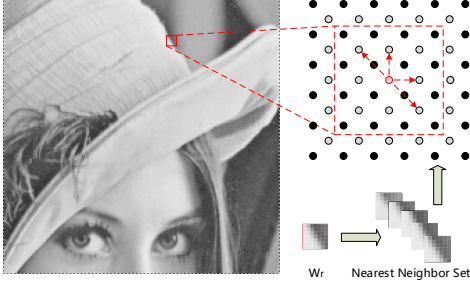


Fig. 3. Context-Based Parameter b .

B. Geodesic Distance Based Similarity

In order to characterize the similarity between two pixels, we use the patch based metric in [9] named patch-geodesic distance to determine whether two pixels are in the same stable region. The basic idea of geodesic distance comes from the connected component consisting of pixels with similar intensities. This modified patch-based geodesic distance takes both the center and its neighboring pixels into consideration, and performs well in deciding whether two pixels share the same AR parameters.

The patch-geodesic distance $D(x, c)$ is defined as minimum value of the accumulated patch difference along all connecting paths:

$$D(x, c) = \min_{P=\{p_1, \dots, p_n\} \in \mathcal{P}_{x,c}} d(P), \quad (7)$$

$$d(P) = \sum_{i=2}^n \|M(p_i) - M(p_{i-1})\|_1,$$

where $\mathcal{P}_{x,c}$ is the set of paths connecting x and c , $M(\cdot)$ is the operator to extract the patch values centered at a pixel and the distance type adopt Manhattan distance.

Then, the distance is converted to a pixel-level similarity:

$$w(x, c) = \exp \left\{ \frac{-D(x, c)}{\beta} \right\}, \quad (8)$$

where β is the parameter controlling the shape of the exponential function. Based on whether x is HR or LR pixels, w is expressed as w^H or w^L . They are latterly incorporated into similarity modulated estimation in Section III-D.

C. Multiplanar Constraint

When the interpolation operation is performed in a local window and HR pixels are simultaneously estimated by the AR model, pixel values are constrained by their neighboring pixels. As stated in Section III-A, the neighbor size is not large and thus there is a magnitude difference between estimation window and neighbor window. The AR parameters merely present a small part of the correlations between pixels. Context modeling AR still lacks enough correlations within estimation window to get a promising interpolation HR image.

Here, we introduce the multiplanar constraint to utilize larger scale correlations within estimation window during image interpolation. Similar patches, whose scales are approximately the same as AR

neighbor window, are collected. The distance function between two patches is defined as follows:

$$dis(\mathbf{x}_s, \mathbf{x}_t) = \|\mathbf{x}_s - \mathbf{x}_t\|_2^2 + \eta \|\nabla \mathbf{x}_s - \nabla \mathbf{x}_t\|_2^2, \quad (9)$$

where \mathbf{x} denotes pixel values of a patch in vector form centered at pixel x , ∇ denotes the gradient operator and η is a parameter used to balance the contribution of two terms. Using Eq. (9), we can obtain the similar patch set compared with the center patch of the estimation window:

$$S = \left\{ \mathbf{x} \mid \exp \left\{ \frac{-dis(\mathbf{x}, \mathbf{x}_c)}{\alpha} \right\} \geq \tau_S \right\}, \quad (10)$$

where τ_S is a threshold and α is the parameter to control the exponential function. Based on whether center pixel x is HR or LR pixel, S can be divided into two subsets, S_H and S_L .

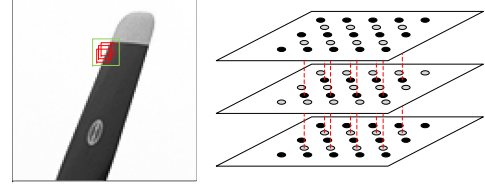


Fig. 4. Multiplanar Constraint.

As shown in Fig. 4, by overlapping and interleaving S_H and S_L we can discover that two patches, one from S_H and the other from S_L , are staggered in means of HR and LR pixels. That is to say, the same position in patches from two different set S_H and S_L lies different type of pixels. Intuitively, HR pixels should be similar with LR pixels (in average value) at the same position in patch. Then, we construct the multiplanar constraint term by utilizing the feature as follows:

$$E_m = \sum_{m \in S_L} \sum_{i \in P_m^L} \left[y_i^{P_m^L} - \frac{1}{N} \sum_{n=1}^N x_i^{P_n^H} \right]^2 + \sum_{n \in S_H} \sum_{i \in P_n^H} \left[y_i^{P_n^H} - \frac{1}{M} \sum_{m=1}^M x_i^{P_m^L} \right]^2, \quad (11)$$

where M and N represent the sizes of S_L and S_H , respectively; P_m^L and P_n^H denote the patches from S_L and S_H ; y_i and x_i denote the HR and LR pixels at position i in corresponding patches.

D. Similarity Modulated Estimation

In order to better characterize the piecewise stationarity in a local window, we incorporate the similarity introduced in Section III-B into AR model as weighting terms.

1) *Parameter Estimation*: The model parameters \mathbf{a} and \mathbf{b} have the similar form, and can be estimated by weighted linear least squares:

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \sum_{x_i \in \mathcal{W}} \left[w_i^L (x_i - \sum_{t=1}^4 a_t x_{i \otimes t}) \right]^2. \quad (12)$$

By introducing l_2 norm into the objective function, Eq. (12) can be written in vector form:

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{a}} \|\mathbf{W}_l(\mathbf{A}\mathbf{a} - \mathbf{x})\|_2^2 + \lambda \|\mathbf{a}\|_2^2, \quad (13)$$

And the analytical solution for Eq.(13) is:

$$\mathbf{a} = (\mathbf{A}^T \mathbf{W}_l^2 \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{W}_l^2 \mathbf{x}, \quad (14)$$

where \mathbf{W}_l is the diagonal matrix of similarity w_i^L , \mathbf{A} is a matrix whose k -th row is the value of corresponding reference.

Similarly, \mathbf{b} can be computed by:

$$\mathbf{b} = (\mathbf{B}^T \mathbf{W}_l^2 \mathbf{B} + \lambda \mathbf{I})^{-1} \mathbf{B}^T \mathbf{W}_l^2 \mathbf{x}. \quad (15)$$

2) *Block Estimation*: We perform the interpolation in a local window W with the assumption that pixels in window share the same parameters. The objective function contains two terms: AR matching error term and multiplanar constraint term

$$E = \sum_{i \in W} \left[w_i^H \left(y_i - \sum_{t=1}^t a_t y_{i \otimes t} \right) \right]^2 + \sum_{i \in W} \left[w_i^L \left(x_i - \sum_{t=1}^t a_t x_{i \otimes t} \right) \right]^2 + \lambda_1 \sum_{i \in W} \left[w_i^H \left(y_i - \sum_{t=1}^t b_t y_{i \odot t} \right) \right]^2 + \lambda_2 E_m, \quad (16)$$

where y and x refer to HR pixel and LR pixel, λ_1 is the Lagrange multiplier for parameter \mathbf{b} , and λ_2 is a user-defined parameter to balance the weight of multiplanar constraint term.

To minimize the objective function, we can reduce Eq.(16) to the vector form:

$$\mathbf{y} = \arg \min_{\mathbf{y}} \|\mathbf{W}(\mathbf{C}\mathbf{y} - \mathbf{D}\mathbf{x})\|_2^2. \quad (17)$$

And the analytical solution can be obtained:

$$\mathbf{y} = (\mathbf{C}^T \mathbf{W}^2 \mathbf{C})^{-1} \mathbf{C}^T \mathbf{W}^2 \mathbf{D} \mathbf{x}. \quad (18)$$

Estimations are performed in high frequency areas and block estimation process only outputs the center pixel. For pixels in low frequency areas, we use the bicubic interpolation method whose results are already good enough.

IV. EXPERIMENTAL RESULTS

The proposed interpolation is implemented on MATLAB 8.6 platform. The proposed algorithm is compared with bicubic interpolation method and four state-of-the-art interpolation algorithms, including NEDI, SAI, IPAR and NARM. We test our method on a large number of images from the Kodak and USC-SIPI image databases.

TABLE I
PSNR(dB) RESULT OF FIVE INTERPOLATION METHODS

Images	Bicubic	NEDI	SAI	IPAR	NARM	Proposed
Child	35.49	34.56	35.63	35.70	35.52	35.72
Lena	34.01	33.72	34.76	34.79	35.09	34.80
Tulip	33.82	33.76	35.71	35.85	36.04	35.91
Cameraman	25.51	25.44	25.99	26.06	26.05	26.11
Monarch	31.93	31.80	33.08	33.34	34.10	33.40
Airplane	29.40	28.00	29.62	30.05	30.05	30.09
Caps	31.25	31.19	31.64	31.67	31.77	31.69
Status	31.36	31.01	31.78	31.94	31.72	31.97
Sailboat	30.12	30.18	30.69	30.85	30.64	30.90
Bike	25.41	25.25	26.28	26.31	26.23	26.33
Ruler	11.98	11.49	11.37	11.81	11.88	13.22
Slope	26.74	26.54	26.63	26.78	26.89	27.10
Average	28.92	28.58	29.43	29.60	29.67	29.77

The input LR images are obtained by downsampling the original HR images directly with the factor of two. Then, different interpolation methods are applied to generate the HR images from the input LR images. Peak Signal-to-Noise Ratio (PSNR) is used to evaluate the experimental results. Table I shows the results of six different interpolation algorithms. We can see that, compared with other AR based method, the proposed method perform better than the second best IPAR method with 0.17dB in terms of PSNR. Even compared with dictionary-based method NARM, the proposed method perform competitively or even better.

Specifically, the proposed method achieves desirable performance in interpolating images with sharp, long edges and large scale texture. *Ruler* and *Slope* achieve the highest gain compared with the secondbest AR based method with 1.41dB and 0.32dB, respectively.

We also provide some visual comparison samples from different image interpolation algorithm in Fig. 5. Compared with other AR based methods, the experimental results of the proposed method are

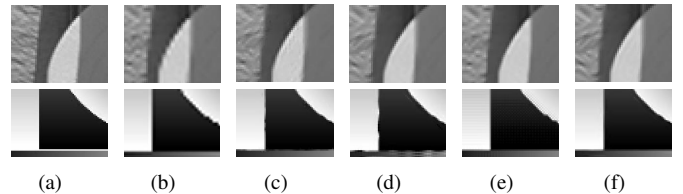


Fig. 5. Subjective image quality. Interpolation results using different method on several images. From top to bottom: *Caps*, *Slope*. From left to right: ground truth, Bicubic, NEDI, IPAR, NARM, proposed.

free from annoying artifacts around the sharp edge and look more similar to the ground truth.

V. CONCLUSION

In this paper, we propose a new AR model for image interpolation by incorporating context-awareness. The proposed context-aware image interpolation methods can obtain better AR parameters by adaptively selecting the reference pixels from a larger candidate set ranked by the correlation coefficient. It includes some reference pixels that are irrelevant with prediction so as to reduce the noisy information. Meanwhile, it includes some closely related reference pixels that are ignored in traditional models due to their long distance from the center pixel, increasing the model precision and stability. Also, we design the multiplanar constraint to enhance the correlations within the estimation window, not only preserving more structural information from LR image, but also greatly reducing the ill-posed condition of normal equation when solving the least square problem of the objective function. The experimental result show that the proposed algorithm achieves better performance than other existing ones.

REFERENCES

- [1] R. G. Keys, "Cubic convolution interpolation for digital image processing," *Acoustics, Speech and Signal Processing, IEEE Transactions on*, vol. 29, no. 6, pp. 1153–1160, 1981.
- [2] H. S. Hou and H. Andrews, "Cubic splines for image interpolation and digital filtering," *Acoustics, Speech and Signal Processing, IEEE Transactions on*, vol. 26, no. 6, pp. 508–517, 1978.
- [3] Q. Wang and R. K. Ward, "A new orientation-adaptive interpolation method," *Image Processing, IEEE Transactions on*, vol. 16, no. 4, pp. 889–900, 2007.
- [4] C. M. Zwart and D. H. Frakes, "Segment adaptive gradient angle interpolation," *Image Processing, IEEE Transactions on*, vol. 22, no. 8, pp. 2960–2969, 2013.
- [5] X. Li and M. T. Orchard, "New edge-directed interpolation," *Image Processing, IEEE Transactions on*, vol. 10, no. 10, pp. 1521–1527, 2001.
- [6] X. Zhang and X. Wu, "Image interpolation by adaptive 2-d autoregressive modeling and soft-decision estimation," *Image Processing, IEEE Transactions on*, vol. 17, no. 6, pp. 887–896, 2008.
- [7] J. Ren, J. Liu, W. Bai, and Z. Guo, "Similarity modulated block estimation for image interpolation," in *Image Processing (ICIP), 2011 18th IEEE International Conference on*. IEEE, 2011, pp. 1177–1180.
- [8] M. Li, J. Liu, J. Ren, and Z. Guo, "Adaptive general scale interpolation based on weighted autoregressive models," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 25, no. 2, pp. 200–211, 2015.
- [9] W. Yang, J. Liu, S. Yang, and Z. Guo, "Novel autoregressive model based on adaptive window-extension and patch-geodesic distance for image interpolation," in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 1211–1215.
- [10] W. Dong, L. Zhang, R. Lukac, and G. Shi, "Sparse representation based image interpolation with nonlocal autoregressive modeling," *Image Processing, IEEE Transactions on*, vol. 22, no. 4, pp. 1382–1394, 2013.
- [11] Y. Romano, M. Protter, and M. Elad, "Single image interpolation via adaptive nonlocal sparsity-based modeling," *Image Processing, IEEE Transactions on*, vol. 23, no. 7, pp. 3085–3098, 2014.
- [12] X. Wu, G. Zhai, X. Yang, and W. Zhang, "Adaptive sequential prediction of multidimensional signals with applications to lossless image coding," *Image Processing, IEEE Transactions on*, vol. 20, no. 1, pp. 36–42, 2011.