

MOTION-ADAPTED THREE-DIMENSIONAL FREQUENCY SELECTIVE EXTRAPOLATION

Andreas Spruck, Markus Jonscher, Jürgen Seiler, and André Kaup

Friedrich-Alexander Universität Erlangen-Nürnberg (FAU)
Multimedia Communications and Signal Processing
Cauerstr. 7, 91058 Erlangen, Germany

ABSTRACT

It has been shown, that high resolution images can be acquired using a low resolution sensor with non-regular sampling. Therefore, post-processing is necessary. In terms of video data, not only the spatial neighborhood can be used to assist the reconstruction, but also the temporal neighborhood. A popular and well performing algorithm for this kind of problem is the three-dimensional frequency selective extrapolation (3D-FSE) for which a motion adapted version is introduced in this paper. This proposed extension solves the problem of changing content within the area considered by the 3D-FSE, which is caused by motion within the sequence. Because of this motion, it may happen that regions are emphasized during the reconstruction that are not present in the original signal within the considered area. By that, false content is introduced into the extrapolated sequence, which affects the resulting image quality negatively. The novel extension, presented in the following, incorporates motion data of the sequence in order to adapt the algorithm accordingly, and compensates changing content, resulting in gains of up to 1.75 dB compared to the existing 3D-FSE.

Index Terms— Frequency Selective Extrapolation, Motion Compensation, Resolution Enhancement

1. INTRODUCTION

The demand for algorithms that are capable of generating high resolution video data from low resolution data grows steadily. The fields of application for such algorithms are manifold and reach from medical imaging over security surveillance to consumer applications like digital remastering of video data. Due to this, there are many different approaches towards this problem. A short overview of the most common methods can be seen in [1]. Another possible method to obtain high resolution image or video data, on which we focus here, is to use a low resolution non-regular sampling sensor for acquisition and applying a reconstruction algorithm on the recorded data afterwards, to recover the image areas where no samples were captured. This method was first introduced in [2] where the two-dimensional frequency selective extrapolation (2D-FSE) [3] was used to extrapolate

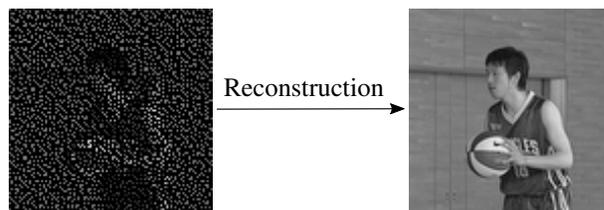


Fig. 1. Schematic representation of the image acquisition using a non-regular sampling sensor.

the non-available areas. As we consider video sequences here, the three-dimensional frequency selective extrapolation (3D-FSE) [4] is used to reconstruct the not directly acquired areas, as it also incorporates the temporal component of the sequence. Moreover, this kind of sampling bears the benefit that aliasing artifacts can be reduced in a visible noticeable amount as shown in [5, 6].

The paper is structured as follows, in Section 2 the basic principle of non-regular sampling will be shown. Section 3 introduces the three-dimensional frequency selective extrapolation. The novel motion adaptive extension will be explained in further detail in Section 4. The simulation results will be presented and discussed in Section 5 before we summarize and conclude our paper in Section 6.

2. RECONSTRUCTION OF NON-REGULAR SAMPLED DATA

The recording scenario that is applied during this work was originally presented in [2]. The sensor introduced there is a low resolution image sensor, where only one fourth of the area of every pixel is sensitive to light. The light sensitive areas are non-regularly distributed among the four quadrants of each pixel. As only a small area of the large low resolution pixel is sensitive to light, this active area can also be regarded as a single pixel of a sensor with a four times higher resolution whose pixels are fully light sensitive. By describing the sensor in that way, it can be regarded as a high resolution sensor covered with a non-regular mask. This mask follows the aforementioned description and leaves only one pixel out of

a 2×2 block unmasked. By this a high resolution image can be acquired while only reading out and storing one fourth of its pixels. As we are processing video data and use the same sensor throughout this work, every frame of the sequence is masked with the same mask in order to obtain a non-regular sub-sampled sequence as presented in [2]. The non-regular sub-sampled sequence can be written as

$$s_{nr}[x, y, t] = s[x, y, t] \cdot b[x, y, t]. \quad (1)$$

Thereby, $s[x, y, t]$ denotes the full high resolution video data and $b[x, y, t]$ denotes the sub-sampling mask as described before. After masking the frames, it becomes obvious that the masked areas need to be reconstructed in order to obtain a satisfying result, as schematically shown in Figure 1. During this work, 3D-FSE is used to reconstruct the signal $s[x, y, t]$ from $s_{nr}[x, y, t]$ following the idea of [2]. This algorithm will be presented in more detail in the next section.

3. THREE-DIMENSIONAL FREQUENCY SELECTIVE EXTRAPOLATION (3D-FSE)

The 3D-FSE used here, is a method to reconstruct lost or defective areas of a video based on the preserved data. It is based on 2D-FSE, as described in [3], and was extended to three dimensions in [4]. The 3D-FSE reconstructs the sequence in a block-wise manner. Meaning, that every block $f[m, n, p]$ within the sequence is reconstructed separately. Every block $f[m, n, p]$ has thereby a dimension of $M \times N$ pixels in the spatial direction, and a length of P in temporal direction. The basic principle of the algorithm is to iteratively generate a model of superimposed weighted Fourier basis functions, that is used to fill in the missing part of the signal, which is called loss area \mathcal{B} , and is depicted in Figure 2. The picture content that is held in the available pixels from the block $f[m, n, p]$ within the non-regular sampled block $f_{nr}[m, n, p]$, is called support area \mathcal{A} and serves as a basis for the model generation. In contrast to the 2D-FSE the area that is to be reconstructed extends over several frames, as denoted by the red boxes in Figure 2. The extrapolation area \mathcal{L} is defined by the variables m and n along the spatial axes and p along the time axis. The extrapolation area \mathcal{L} is composed of the support area \mathcal{A} , the loss area \mathcal{B} and the area \mathcal{R} , which holds the previously reconstructed values, as $\mathcal{L} = \mathcal{A} \cup \mathcal{B} \cup \mathcal{R}$ as stated in [4].

The aim of 3D-FSE is to generate a model which replicates the original signal $f[m, n, p]$ in the volume \mathcal{L} . The non-sampled areas of the signal $f_{nr}[m, n, p]$ are then replaced by this model.

The exact description of the algorithm can be found in [4]. During the scope of this paper, the algorithm from [4] was extended by using the optimized processing order from [7]. By incorporating the weighting function $w[m, n, p]$, it is possible to assign different weights to the different areas of the sequence. By this, it is ensured that missing regions are not taken into account for model generation and pixels far

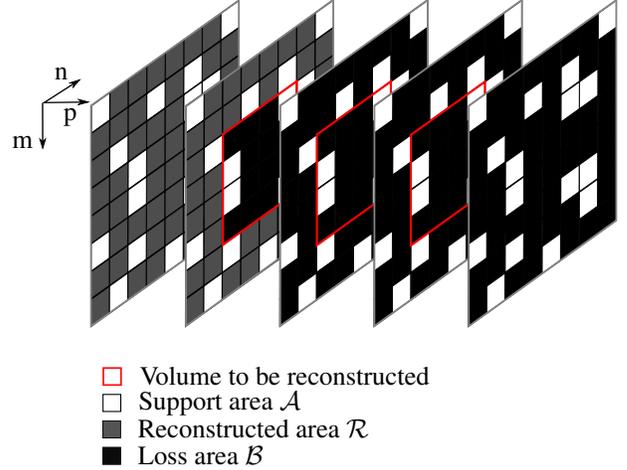


Fig. 2. Schematic representation of the reconstruction area of the 3D-FSE.

away from the currently considered area are weighted with a smaller weight than areas nearby. The spatial weighting function

$$w[m, n, p] = \begin{cases} \rho[m, n, p] & , (m, n, p) \in \mathcal{A} \\ \delta \rho[m, n, p] & , (m, n, p) \in \mathcal{R} \\ 0 & , (m, n, p) \in \mathcal{B} \end{cases} \quad (2)$$

used during the selection of the basis function is defined as in [4]. By varying $\hat{\rho}$ in

$$\rho[m, n, p] = \hat{\rho} \sqrt{\left(m - \frac{M-1}{2}\right)^2 + \left(n - \frac{N-1}{2}\right)^2 + \left(p - \frac{P-1}{2}\right)^2} \quad (3)$$

the decay of the weighting function can be adjusted. Parameter δ controls how much influence previously reconstructed values have on the model generation.

The advantage of the FSE is that the computationally expensive calculations can be performed in the frequency domain as it is shown in [3] and [4]. By doing this, the generation of the model can be sped up drastically. As all operations can be executed in the frequency domain it is sufficient to perform only two transformations, one in the beginning and one in the end [4].

4. MOTION COMPENSATED WEIGHTING

The problem occurring with the static weighting function as presented before and used in [4] and [7], is, that due to motion within the sequence, the region of interest at which the maximum of the weighting function is placed, moves out of the maximum of the decaying weighting function over time, as can be seen in Figure 3. By the fact, that the content of the maximally weighted area varies over time, content that is actually not present in the part of the sequence, that is to be reconstructed, might be weighted with a high weight, resulting in admission of basis functions to the model that are

not part of the original signal. This leads to unsharp edges and even ghosting artifacts, resulting in a strong degradation of the overall reconstruction quality. Due to this, the fixed spatial weighting function as presented in (2) is replaced by a more sophisticated one in this paper.

Different approaches to solving this problem have been made in the past. In [8], the extrapolation volume \mathcal{L} was selected, by incorporating motion data, such that the selected area of every frame shows the same image region over time. Thereby it is guaranteed that the weighting function assesses the same regions of the image the maximal weight, as the selection of the support area of the 3D-FSE is performed in a motion compensated manner. In [9] another approach is presented. There, a first reconstruction using the two-dimensional frequency selective reconstruction (2D-FSR) from [10] is performed. After this reconstruction, a motion estimation is performed to search for pixels in the loss area \mathcal{B} of the currently considered frame that are contained in the support area \mathcal{A} in previous or succeeding frames. If such pixels are found, they are copied into the loss area of the current cube. By this, the size of the loss area can be reduced, which results in a far better reconstruction quality of the following 2D-FSR [9].

In this paper we pursue a different approach. We shift the weighting function according to the motion in the sequence such that in every frame the maximum of the weighting function is placed over the same content, as it can be seen in the bottom line of Figure 3. For the motion estimation, the sub-sampled sequence is reconstructed with a bilinear interpolation in a first step in order to enable a motion estimation. Afterwards a optical flow implementation following [11] is used to estimate the motion on the interpolated data. For the estimation of the motion, the whole interpolated frame is regarded. This results in a motion vector field holding a motion vector for every pixel in the frame relative to the currently considered frame. Out of this vector field the cube currently considered by the 3D-FSE is extracted. These vectors are averaged to one overall motion vector for each slice of the cube

$$\begin{pmatrix} \bar{v}_x[p] \\ \bar{v}_y[p] \end{pmatrix} = \frac{1}{MN} \cdot \sum_{\forall m,n} \begin{pmatrix} v_x[m,n,p] \\ v_y[m,n,p] \end{pmatrix} \quad (4)$$

representing the averaged motion between two succeeding frames within the considered volume. The maximum of the weighting function is shifted in the direction the averaged motion vector $(\bar{v}_x[p] \ \bar{v}_y[p])^T$ is pointing. By this, it is ensured that the weighting function emphasizes the same image region over time. The overall shape of the weighting function decaying in spatial and temporal direction stays thereby unchanged. Following this description the motion compensated weighting function can be given as

$$\tilde{w}[m,n,p] = \begin{cases} \tilde{\rho}[m,n,p] & , (m,n,p) \in \mathcal{A} \\ \delta \tilde{\rho}[m,n,p] & , (m,n,p) \in \mathcal{R} \\ 0 & , (m,n,p) \in \mathcal{B} \end{cases} \quad (5)$$

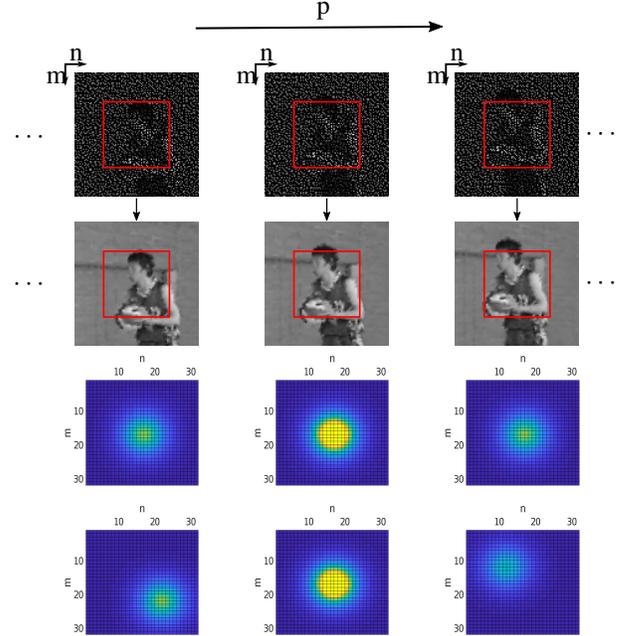


Fig. 3. First row: Non-regular sampled sequence, second row: linear interpolated sequence, third row: static spatial weighting function for corresponding frame, fourth row: motion compensated spatial weighting function for corresponding frame.

with

$$\tilde{\rho}[m,n,p] = \hat{\rho} \sqrt{(m - \frac{M-1}{2} - \bar{v}_x[p])^2 + (n - \frac{N-1}{2} - \bar{v}_y[p])^2 + (p - \frac{P-1}{2})^2} \quad (6)$$

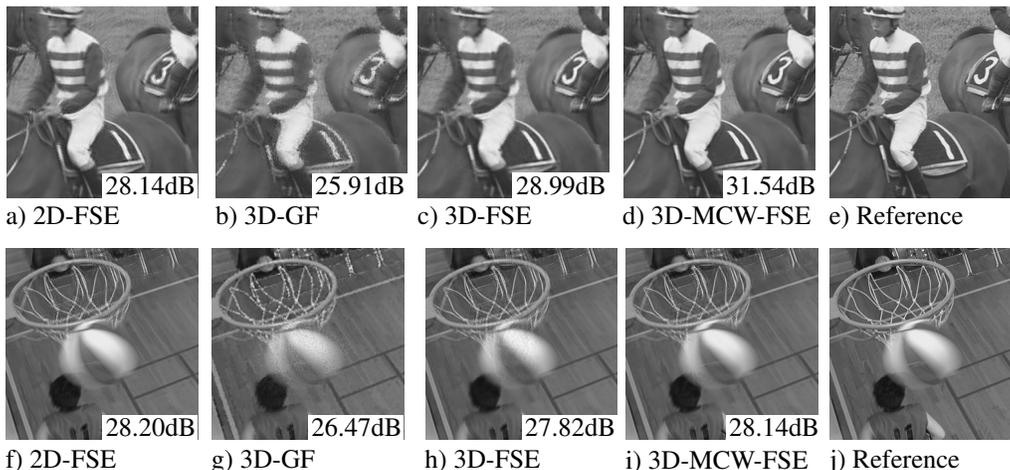
where $\bar{v}_x[p]$ and $\bar{v}_y[p]$ are the horizontal and the vertical component of the averaged motion vector of the according frame p in relation to the current frame, respectively.

5. SIMULATIONS AND ANALYSIS

For testing the proposed method, the classes C and D of the HEVC test-set [14] are used. These classes consist of four sequences each with a resolution of 832×480 and 416×240 pixels, respectively. The first fifty frames of each sequence are considered. The 3D-FSE incorporating the novel motion compensated weighting function (3D-MCW-FSE) is compared to the common 3D-FSE as presented in [7]. Furthermore, the presented extension of the 3D-FSE algorithm is compared to the 3D-Gap Filling algorithm (3D-GF) by Garcia et al. [12, 13]. Moreover these three dimensional algorithms are compared to the 2D-FSE [3]. For the 2D-FSE, 3D-FSE and the 3D-MCW-FSE the same parameters are used: A block-size of 4×4 respectively $4 \times 4 \times 1$, a border width of 14 and FFT of size $32 \times 32 \times 32$ is used. The weighting functions decay with $\hat{\rho} = 0.7$, γ and δ are set to 0.5.

Table 1. COMPARISON OF THE ACHIEVED PSNR-VALUES OF THE SIMULATION RESULTS.

Sequence	2D-FSE [3]	3D-GF [12, 13]	3D-FSE [7]	3D-MCW-FSE
Basketball Pass	30.18 dB	29.70 dB	31.49 dB	31.64 dB
Blowing Bubbles	28.04 dB	28.04 dB	29.93 dB	30.19 dB
BQ Square	21.02 dB	22.17 dB	23.77 dB	23.84 dB
Race Horses	28.40 dB	27.30 dB	28.94 dB	30.69 dB
Basketball Drill	31.55 dB	29.81 dB	31.27 dB	31.80 dB
BQ Mall	28.72 dB	27.12 dB	29.24 dB	29.72 dB
Party Scene	23.74 dB	24.13 dB	26.26 dB	26.38 dB
Race Horses	28.58 dB	27.09 dB	28.97 dB	30.62 dB
Average	27.53 dB	26.92 dB	28.73 dB	29.36 dB

**Fig. 4.** Comparison of results of the 2D-FSE, 3D-GF, 3D-FSE and 3D-MCW-FSE for the sequences Race Horses (top) and Basketball Drill (bottom).

The results of the aforementioned simulations are shown in Table 1. The given values are the average result of three reconstructions of differently sampled versions of the same sequences. As seen there, the 3D-MCW-FSE is in all cases the best performing algorithm. The highest gain can be achieved for sequences that contain much motion, as for example the Race Horses or Basketball Drill sequence. For sequences, that contain only little motion the 3D-FSE approaches the 3D-MCW-FSE as the weighting functions are nearly identical in this case. This effect can be observed with the BQ Square sequence for example. Figure 4 shows some results of the simulations for the sequences Race Horses and Basketball Drill. For better visibility enlarged excerpts are displayed here. As can be seen in Figure 4 d) and i), the 3D-MCW-FSE achieves much sharper edges and is capable of reconstructing finer details than the other algorithms. Comparing Figure 4 h) and i) one can observe, that the ball is distorted with structured artifacts in h). These artifacts are caused by the motion within the sequence, as the net, which is visible in previous and succeeding frames, contributes to the model up to a cer-

tain degree. This leads to a visual noticeable degradation of the image quality. Using the 3D-MCW-FSE with the motion compensated spatial weighting function, proposed here, results in a noticeable better reconstruction quality.

6. CONCLUSION AND OUTLOOK

In this paper, the novel three-dimensional frequency selective extrapolation with motion compensated spatial weighting (3D-MCW-FSE) was presented. For video data that was recorded using a non-regular sampling sensor the proposed 3D-MCW-FSE was able to achieve a visually noticeable gain of up to 1.75 dB over the existing 3D-FSE. This gain is achieved by incorporating motion data that is obtained using optical flow.

Focus of further research will be to incorporate a scene change detection algorithm in addition, to prevent negative effects due to image content from neighboring scenes within a sequence.

7. REFERENCES

- [1] S. C. Park, M. K. Park, and M. G. Kang, "Super-resolution image reconstruction: a technical overview," in *IEEE Signal Processing Magazine*, May 2003, vol. 20, pp. 21–36.
- [2] M. Schöberl, J. Seiler, S. Foessel, and A. Kaup, "Increasing imaging resolution by covering your sensor," in *18th IEEE International Conference on Image Processing*, Sept 2011, pp. 1897–1900.
- [3] J. Seiler and A. Kaup, "Complex-valued frequency selective extrapolation for fast image and video signal extrapolation," in *IEEE Signal Processing Letters*, Nov 2010, vol. 17, pp. 949–952.
- [4] Katrin Meisinger and André Kaup, "Spatiotemporal selective extrapolation for 3D signals and its applications in video communications," *IEEE Transactions on Image Processing*, vol. 16, no. 9, pp. 2348–2360, 2007.
- [5] Gilles Hennenfent and Felix J Herrmann, "Irregular sampling—from aliasing to noise," in *69th EAGE Conference and Exhibition incorporating SPE EUROPEC 2007*, 2007.
- [6] Yui Maeda and Junichi Akita, "A CMOS image sensor with pseudorandom pixel placement for clear imaging," in *Intelligent Signal Processing and Communication Systems, 2009. ISPACS 2009. International Symposium on*. IEEE, 2009, pp. 367–370.
- [7] J. Seiler, S. Schöll, W. Schnurrer, and A. Kaup, "Optimized processing order for 3D hole filling in video sequences using frequency selective extrapolation," in *Picture Coding Symposium*, Dec 2016, pp. 1–5.
- [8] J. Seiler and A. Kaup, "Motion compensated three-dimensional frequency selective extrapolation for improved error concealment in video communication," *Journal of Visual Communication and Image Representation*, vol. 22, no. 3, pp. 213–225, 2011.
- [9] M. Jonscher, K. Jaskolka, J. Seiler, and A. Kaup, "Recursive frequency selective reconstruction of non-regularly sampled video data," in *Picture Coding Symposium (PCS)*, Dec 2016, pp. 1–5.
- [10] J. Seiler, M. Jonscher, M. Schöberl, and A. Kaup, "Resampling images to a regular grid from a non-regular subset of pixel positions using frequency selective reconstruction," in *IEEE Transactions on Image Processing*, Nov 2015, vol. 24, pp. 4540–4555.
- [11] Gunnar Farneback, "Two-frame motion estimation based on polynomial expansion," *Image analysis*, pp. 363–370, 2003.
- [12] Damien Garcia, "Robust smoothing of gridded data in one and higher dimensions with missing values," *Computational statistics & data analysis*, vol. 54, no. 4, pp. 1167–1178, 2010.
- [13] Guojie Wang, Damien Garcia, Yi Liu, Richard De Jeu, and A Johannes Dolman, "A three-dimensional gap filling method for large geophysical datasets: Application to global satellite soil moisture observations," *Environmental Modelling & Software*, vol. 30, pp. 139–142, 2012.
- [14] F. Bossen et al., "Common test conditions and software reference configurations," in *11th Meeting: Joint Collaborative Team on Video Coding of ITU-T SG*, 2011, vol. 16.