# DOWN-SAMPLING IN DCT DOMAIN USING LINEAR TRANSFORM WITH DOUBLE-SIDED MULTIPLICATION FOR IMAGE/VIDEO TRANSCODING

Xiang Yu, En-hui Yang, and Haiquan Wang

Dept. of Electrical and Computer Engineering University of Waterloo, Waterloo, Ontario, N2L 3G1, Canada

## ABSTRACT

This paper proposes a designing framework for downsampling compressed images/video frames with arbitrary ratio in the discrete cosine transform (DCT) domain. We first derive a set of DCT-domain down-sampling methods which can be represented by a linear transform with double-sided matrix multiplication (LTDS) in the DCT domain, and show that the set contains a wide range of methods with various complexity and visual quality. Then, based on a pre-selected spatialdomain method, we formulate an optimization problem for finding an LTDS to approximate the given spatial domain method for achieving the best trade-off between the visual quality and the complexity. By selecting a spatial-domain reference method with the popular Butterworth lowpass filtering and bicubic interpolation, the proposed framework discovers LTDSs with better visual quality and lower computational complexity as saving 20%~70% execution time when compared with state-of-the-art methods in the literature.

Index Terms DCT, down-sampling, transcoding.

## 1. INTRODUCTION

As the number of network users with diverse devices is steadily increasing, there has been a great momentum in the multimedia industry for supporting content display in diverse devices all over the network. A big challenge, however, is that these devices have various display resolutions ranging from large resolutions such as those of high-definition TV to small resolutions such as those of smart phones. This leads to research on transcoding, which involves in automatic reformatting. In this paper, we are concerned with transcoding with spatial resolution down-conversion due to its applications to wireless communications. In particular, we will focus on transcoding compressed image/video in the DCT domain because most image/video data to be shared over the network are originally captured with high resolution and coded using DCT-based techniques such as JPEG, DV, etc. A straightforward method for down-sampling DCT images is to concatenate inverse DCT, spatial-domain downsampling and DCT[9]. However, it is not suitable for some real-world transcoding applications because of its high complexity. To tackle the complexity issue, it is desired to perform down-sampling in the DCT domain directly. One category of methods for doing this is referred to as DCT coefficient manipulation [5], [6]. The manipulation methods generally encounter a problem of ringing effect[5], particularly for images with sharp contrast, because truncating DCT coefficients is equivalent to filtering by an ideal filter with a very narrow transition band[9], which leads to ringing effect.

Another category of DCT-domain methods are based on a transform model corresponding to a reference method selected in the spatial domain[1][2][8]. E.g., in [1], an efficient DCT-domain transform was developed based on a specific spatial domain method with the nearest neighbor approach. The omitting of lowpass filtering leads to an efficient DCT-domain transform, but it also limits the quality. On the other hand, a DCT-domain method based on a spatial-domain method with a freely-selected lowpass filter still has a high complexity for software implementation [2].

In this paper, we propose a designing framework for an efficient DCT-domain down-sampling transform corresponding to any given spatial-domain down-sampling method. First, a linear transform with double-sided matrix multiplication for DCT-domain down-sampling is derived corresponding to an arbitrarily selected spatial-domain down-sampling method. An optimization problem is then formulated to minimize a joint cost of the visual quality and the complexity based on the linear transform. By modeling the linear transform as a multiple layer network, a so-called structural learning with forgetting algorithm [7] is used for training, resulting in a desired tradeoff between the quality and the complexity. In our simulations, the popular and well-recognized spatial-domain method with Butterworth lowpass filter and bicubic interpolation is used as the reference. The proposed framework discovers LTDSs with good subjective visual quality and low computational complexity with significant saving of execution time when compared with other methods.

The paper is organized as follows. Section 2 derives LTDSs. An optimization problem is formulated in Section 3 to find

This work was supported in part by the Natural Sciences and Engineering Research Council of Canada under Grants RGPIN203035-02 and RGPIN203035-06, by the Premier's Research Excellence Award, by the Ontario Distinguished Researcher Award, and by the Canada Research Chairs Program.

an LTDS approximation for a given spatial-domain method, while solution is provided in Section 4. Experiments are discussed in Section 5, and conclusion is drawn in Section 6.

# 2. DCT-DOMAIN LINEAR TRANSFORM WITH DOUBLE-SIDED MATRIX MULTIPLICATION

This section derives a result that for a wide range of spatialdomain down-sampling methods with certain properties, a concatenation of inverse DCT, spatial-domain down-sampling and DCT can be implemented as an LTDS in the DCT domain.

Denote  $t_{ss}$  as a DCT matrix. Consider to down-sample an M×N DCT image  $C_{MN}$ . First, apply inverse DCT to recover the original image  $X_{MN}$  as

$$X_{\rm MN} = t_{\rm SS}' \boxdot C_{\rm MN} \boxdot t_{\rm SS}, \qquad (1)$$

where  $\Box$  denotes block-wise multiplications. Then, a spatialdomain method is used to down-sample  $X_{MN}$  to an I×J image  $x_{U}$ . Finally, a DCT-domain down-sampling result is

$$\boldsymbol{V}_{\rm IJ} = \boldsymbol{t}_{\rm SS} \boxdot \boldsymbol{x}_{\rm IJ} \boxdot \boldsymbol{t}_{\rm SS}^{\prime}. \tag{2}$$

Consider a general spatial-domain down-sampling method. First apply a 2D low-pass filter based on the 2D discrete Fourier transform (DFT), i.e.,

$$\tilde{\boldsymbol{X}}_{\rm MN} = \boldsymbol{A}_{\rm MM}^* ((\boldsymbol{A}_{\rm MM} \cdot \boldsymbol{X}_{\rm MN} \cdot \boldsymbol{B}_{\rm NN}) \otimes \boldsymbol{F}_{\rm MN}) \cdot \boldsymbol{B}_{\rm NN}^*, \qquad (3)$$

where  $\hat{X}_{MN}$  is the low-pass filtering output,  $A_{MM}$  is a standard M×M DFT transform matrix and the \* represents the conjugate.  $B_{NN}$  is a standard N×N DFT transform matrix.  $F_{MN}$  is a low-pass filtering matrix in the DFT domain. The symbol  $\otimes$  denotes element-wise multiplications. Then, consider to implement a 2D interpolation with two 1D linear interpolations along the width and the height separately. Denote  $E_{IM}$  and  $G_{NJ}$  as the interpolation matrixes. The down-sampling output is

$$\boldsymbol{x}_{\mathrm{IJ}} = \boldsymbol{E}_{\mathrm{IM}} \cdot \boldsymbol{X}_{\mathrm{MN}} \cdot \boldsymbol{G}_{\mathrm{NJ}}.$$
 (4)

Combining equations (1) to (4), we have

$$V_{\rm IJ} = t \odot [E_{\rm IM} A^*_{\rm MM} [(A_{\rm MM} (t' \odot C_{\rm MN} \odot t) B_{\rm NN}) \otimes F_{\rm MN}] B^*_{\rm NN} G_{\rm NJ}] \odot t', \quad (5)$$

where  $t = t_{ss}$ . The  $\Box$  in (5) can be replaced by applying a result of  $t \Box C_{MN} = T_{MM,t} \cdot C_{MN}$ , where  $T_{MM,t} = \text{diag}(t, \cdots, t)$ . The  $\otimes$  in (5) may be removed under one condition that the lowpass filter takes a form of

$$\boldsymbol{F}_{\rm MN} = \boldsymbol{L}_{\rm M1} \cdot \boldsymbol{R}_{\rm 1N}. \tag{6}$$

Consequently, (5) can be rewritten as a linear transform,

$$\boldsymbol{V}_{\rm IJ} = \boldsymbol{D}_{\rm IM}^{\rm o} \cdot \boldsymbol{C}_{\rm MN} \cdot \boldsymbol{W}_{\rm NJ}^{\rm o}, \qquad (7)$$

where

$$W_{\rm NJ}^{\circ} = T_{\rm NN,l}B_{\rm NN}R_{\rm NN}^{*}B_{\rm NN}G_{\rm NJ}T_{\rm JJ,t'},$$

and  $L_{\text{MM}}$  and  $R_{\text{NN}}$  are diagonal matrixes with diagonal elements being  $L_{\text{M1}}$  and  $R_{1\text{N}}$ , respectively. Motivated by (7), we call any DCT domain transform in the form of  $D_{IM} \cdot C_{MN} \cdot W_{NJ}$  a DCT-domain LTDS, where  $D_{IM}$ and  $W_{NJ}$  are arbitrary matrixes with respective dimensions. We are interested in the set *S* of all DCT-domain LTDSs. Clearly, as shown in (7), the set contains all spatial domain methods satisfying (4) and (6). Some LTDSs in the set have high complexity while others have low complexity. Given any spatial-domain method, in the following we want to find its LTDS approximation in the set, which gives the best trade-off between the quality and complexity.

#### **3. PROBLEM FORMULATION**

Fix a spatial-domain method A. Applying A to any DCT image  $C_{\text{MN}}$ , one gets a corresponding down-sampled DCT image  $V_{\text{II}}$ . Our purpose is to find an LTDS in the set S that, when applied to  $C_{\text{MN}}$ , outputs an image with similar quality to  $V_{\text{II}}$  with a low complexity.

#### 3.1. Visual Quality Measurement

Quality measure is the very basis to formulate our optimization problem for finding optimal LTDSs corresponding to a pre-selected spatial-domain down-sampling method. There have been two major objective quality measures used in the literature for image down-sampling. The first one is to measure the quality with a reference image obtained using a standard spatial-domain down-sampling method [1]. The second one is to up-sample the down-sampled image to the original resolution. Then, the quality is measured by the MSE between the up-sampled image and the original one [5].

In this paper, we apply the first measure, which will naturally allow us to approach the visual performance of a preselected spatial-domain method by setting the down-sampling output of the pre-selected method as the reference. Then, we may be able to approach the best visual quality by preselecting a spatial-domain method with the best visual quality if there is one. Since we are interested in viewing a downsampled image in its own resolution, there is no up-sampling process involved, and hence the first measure is more appropriate than the second measure for our purpose.

Specifically, the quality measure in the paper is, for any  $g(\cdot) \in S$ , the quality of  $g(C_{MN})$  is measured by  $|V_{II} - g(C_{MN})|^2$ , where  $V_{II}$  is obtained using a given spatial-domain method.

#### 3.2. A complexity model

In general, the complexity for computing  $D_{IM} \cdot C_{MN} \cdot W_{NJ}$  is related to two facts. The first is the number of non-zero elements in  $D_{IM}$  and  $W_{NJ}$ . The second is how multiplications may be implemented by additions and shifts. Initially, we consider a complexity model of  $r_f = |D_{IM}| + |W_{NJ}|$ , where  $|\cdot|$  defines the  $l_1$  norm of a matrix. By a learning with forgetting stage of SLF[7], the minimization of  $|D_{IM}| + |W_{NJ}|$  will force as many as possible elements to be zero.

The complexity model  $r_{\rm f}$  is further adjusted by considering to implement a multiplication with a series of additions and shifts. First, we introduce a quantization procedure as follows,

$$\mathbf{Q}(x) = \sum_{i=-2}^{i=15} a_i \cdot 2^{-i}, \ a_i \in \{1, -1, 0\}$$
(8)

where 
$$\{a_i\} = \arg \min_{|x-\sum (a_i 2^{-i})| \le |x|\eta} \sum |a_i|$$
,

where  $\eta$  is a small constant, and  $x \in (-8, 8)^{-1}$ . Essentially, this quantization procedure leads to an approximation within a given neighboring region with the minimal number of ones in the binary representation. Thus, the corresponding multiplication may be implemented with the minimal number of shifts and additions. Overall, the complexity model becomes,

$$r_{q} = (|\boldsymbol{D}_{IM}|_{|d_{im}| < d_{0}} + |\boldsymbol{W}_{NJ}|_{|w_{nj}| < w_{0}}) + \rho \cdot (|\boldsymbol{D}_{IM} - Q(\boldsymbol{D}_{IM})| + |\boldsymbol{W}_{NJ} - Q(\boldsymbol{W}_{NJ})|), \quad (9)$$

where  $\rho$  is a constant, and  $Q(\boldsymbol{D}_{\text{IM}})$  and  $Q(\boldsymbol{W}_{\text{NJ}})$  mean to apply  $Q(\cdot)$  to each element of  $\boldsymbol{D}_{\text{IM}}$  and  $\boldsymbol{W}_{\text{NJ}}$ .

# **3.3.** Joint Optimization of Visual Quality and Complexity

A joint optimization problem is then formulated as follows,

$$\min_{\boldsymbol{D}_{\mathrm{IM}},\boldsymbol{W}_{\mathrm{NJ}}} ||\boldsymbol{D}_{\mathrm{IM}}\boldsymbol{C}_{\mathrm{MN}}\boldsymbol{W}_{\mathrm{NJ}} - \boldsymbol{V}_{\mathrm{IJ}}||^2 + \lambda \cdot r.$$
(10)

Clearly, knowledge of down-sampling in the spatial domain has been utilized to set up a benchmark of the quality for designing a down-sampling algorithm in the DCT domain. The objective is to find an LTDS with a desired trade-off between the visual quality and the complexity of r, which takes value of  $r_{\rm f}$  or  $r_{\rm q}$ .

#### 4. PROBLEM SOLUTION

The optimization problem (10) is solved by training the 3layer network shown in Figure 1 using the SLF algorithm [7]. Specifically, the learning procedure includes two stages, i.e., learning with forgetting and learning with selective forgetting, with  $r = r_f$  and  $r = r_q$ , respectively. Given training data ( $C_{MN}$ ,  $V_{U}$ ), the learning with forgetting stage is as follows:

1. Pass  $C_{\rm MN}$  forward to compute the network outputs.

$$\boldsymbol{Y}_{\text{IN}} = \boldsymbol{D}_{\text{IM}} \cdot \boldsymbol{C}_{\text{MN}} \quad \Rightarrow \quad \boldsymbol{Z}_{\text{IJ}} = \boldsymbol{Y}_{\text{IN}} \cdot \boldsymbol{W}_{\text{NJ}}$$

2. Compute the network error, and propagate it backward.

$$\Delta \mathbf{Z}_{\mathrm{II}} = \mathbf{Z}_{\mathrm{II}} - \mathbf{V}_{\mathrm{II}} \quad \Rightarrow \quad (\Delta \mathbf{Y})_{\mathrm{IN}} = (\Delta \mathbf{Z})_{\mathrm{II}} \cdot (\mathbf{W}^{\mathrm{t}})_{\mathrm{IN}}$$

3. Compute the learning amount for D and W.

$$\Delta \boldsymbol{D} = (\Delta \boldsymbol{Y})_{\text{IN}} \cdot (\boldsymbol{C}^{\text{t}})_{\text{NM}} + \lambda_1 \cdot \text{sgn}(\boldsymbol{D}_{\text{IM}}) \Delta \boldsymbol{W} = (\boldsymbol{Y}^{\text{t}})_{\text{NI}} \cdot (\Delta \boldsymbol{Z})_{\text{IJ}} + \lambda_1 \cdot \text{sgn}(\boldsymbol{W}_{\text{NJ}})$$
(11)

where  $sgn(\cdot)$  is the sign function.



Fig. 1. A network structure and connections.

4. Learn with error propagation and forgetting.

$$\begin{array}{lll} \boldsymbol{D}_{\scriptscriptstyle \mathrm{NJ}}^{(n+1)} & = & \boldsymbol{D}_{\scriptscriptstyle \mathrm{NJ}}^{(n)} - \alpha \cdot \Delta \boldsymbol{D} \\ \boldsymbol{W}_{\scriptscriptstyle \mathrm{NJ}}^{(n+1)} & = & \boldsymbol{W}_{\scriptscriptstyle \mathrm{NJ}}^{(n)} - \alpha \cdot \Delta \boldsymbol{W} \end{array}$$

where  $\alpha$  is a constant, the superscripts (n) and (n+1) accord to the *n*th and (n+1)th iterations. Note that the superscripts are omitted in steps 1 to 4 for simplicity.

5. Repeat steps 1 to 4 until the decrement of  $J_f$  is smaller than a given threshold.

The above learning stage ends with a skeleton structure but a large distortion of  $||D \cdot C \cdot W - V||^2$ . The selective forgetting stage is then used to tune the structure for a better trade-off between the distortion and the complexity. Compared with the above procedure, the selective forgetting stage is mostly the same, except the computation of (11). See [3] for details.

The design procedure is summarized as follows,

- Generate a training set: Choose some pictures. Given a pixel-domain down-sampling method, apply it to these images. Compute DCT transform for the original images and the down-sampled versions as the training set.
- 2. Learning with forgetting: Construct the 3-layer structure with  $D_{IM}$  and  $W_{NI}$ . Find a skeleton structure using the learning with forgetting algorithm.
- 3. Learning with selective forgetting: Refine  $D_{IM}$  and  $W_{NJ}$  with the learning with selective forgetting algorithm.

#### 5. SIMULATION RESULTS

The proposed design algorithm has been implemented and applied to discover down-sampling methods in the DCT domain with good quality and low complexity. A spatial-domain method with the 10th order Butterworth low-pass filtering and bicubic interpolation is selected to generate reference images for evaluating the visual quality among LTDSs.

The LTDS corresponding to  $d_0 = w_0 = 0.1$  makes a good choice for down-sampling by 2:1 in terms of a better quality and a lower complexity when compared with other algorithms. Table 1 provides a comprehensive comparison for our obtained LTDS with four other algorithms in the literature. The method in [5] is developed for down-sampling by a factor of 2 based on DCT coefficient manipulation while the method

<sup>&</sup>lt;sup>1</sup>This range is set empirically as observation shows that elements in  $D_{IM}$  and  $W_{NJ}$  have magnitudes strictly smaller than 1.

**Table 1**. Performance comparison of five methods for down-sampling JPEG images/DV frames with a ratio of 2:1. Complexity is measured with number of operations per pixel in the original image, while computation time is reported based on our computer with 3.4Ghz P-IV CPU. The visual quality is measured by subjective criterions for a testing set of 20 images. Note that a 'yes' means that the corresponding effect shows up for some, not necessarily all, images in the whole set, while a 'no' means that the corresponding effect has not been observed for all images in the set. Visual illustration is omitted due to the page limitation and the fact that all images need to be illustrated in original sizes in order not to introduce resizing distortion. See [3] for visual illustration.

	Complexity			Visual quality			Computation time
	MUL	ADD	SHL	Ringing	Aliasing	Other artifacts	per image
LTDS with $(d_0 = w_0 = 0.1)$	0	5.16	3.66	no	no	no	1.6ms
Method in [5]	1.25	1.25	0	yes	no	no	2.1ms
L/M [6]	3.31	8.68	2.2	yes	no	no	6.3ms
Bilinear average in [8]	3.75	5.81	0.38	yes	yes	no	2.9ms
Fast algorithm in [8]	0	2.72	0.72	yes	yes	severe	0.9ms

**Table 2.** Performance comparison for down-sampling JPEG images/DV frames with a ratio of 3:2.

	C	omplexit	y	Visual quality	Computation time
	MUL	ADD	SHL	Ringing effect	per image
LTDS	0.26	10.4	14.3	no	3.2ms
L/M [6] (II)	1.94	8.06	0	yes	4.1ms

in [6] shares a similar spirit of DCT coefficient manipulation but it is extended to support arbitrary down-sampling ratios. The work in [8] targeted down-sampling by 2:1 with  $8 \times 8$ DCT. The fast algorithm in [8] is the only one more efficient than the obtained LTDS, yet it has severe artifacts. In term of visual quality, the method in [5] is similar with the obtained LTDS as it only shows a slight ringing effect, while the obtained LTDS shows a 20% faster execution speed in our simulation.

The proposed framework is applicable for generating downsampling algorithms with arbitrary ratios. Experiments have been conducted to generate down-sampling algorithms with a ratio of 3:2. The obtained LTDS for down-sampling by 3:2 with  $d_0 = w_0 = 0.02$  is compared with the L/M method in [6], since the method in [8] is for 2:1 only and the work in [5] targets for down-sampling by a factor of 2. The result is shown in Table 2, with a focus on the complexity. Essentially, the obtained LTDS has a lower complexity than the method in [6]. Specifically, there are two algorithms proposed for downsampling by 3:2 in [6], referred to as case I and case II, with case II being the enhanced version of case I. Experimental results by the computation time show that the obtained LTDS is more efficient than the case II algorithm.

## 6. CONCLUSIONS

In this paper, we have proposed a design framework for efficient down-sampling in the DCT domain by jointly optimizing the visual quality and the computational complexity. First, a DCT-domain LTDS is derived corresponding to a concatenation of inverse DCT, spatial-domain down-sampling, and DCT. Based on the LTDS model, an optimization problem is formulated to find an optimal DCT domain method for a given spatial-domain method. Then, the problem is solved using an automatic machine learning algorithm called structure learning with forgetting, resulting in a desired trade-off between the visual quality and the computational complexity. Experiments on JPEG image and DV video transcoding show that the down-sampling methods discovered by the proposed design framework outperforms other methods with a significant reduction of computation time by  $20\% \sim 70\%$  when similar or better visual quality is maintained.

#### 7. REFERENCES

- N. Merhav, V. Bhaskaran, "Fast Algorithms for DCT-Domain Image Down-sampling and for Inverse Motion Compensation", *IEEE Transactions on Circuits and Systems for Video Technol*ogy, pp.468-476, Vol. 7, No. 3, June 1997.
- [2] J.B. Lee, A. Eleftheriadis, "2-D Transform-Domain Resolution Translation", *IEEE Transactions on Circuits and Systems for Video Technology*, pp.704-714, Vol. 10, No. 5, August 2000.
- [3] X. Yu, E.-h. Yang, and H.Q. Wang, "Down-sampling Design in DCT Domain with Arbitrary Ratio for Image/Video Transcoding," Submited to IEEE Trans. on Image Processing.
- [4] S.F. Chang, D. G. Messerschmitt, "Manipulation and Compositing of MC-DCT COmpressed Video", *IEEE Journal on Selected Areas in Communications*, ppl-11, Vol. 13, No. 1, January 1995.
- [5] R. Dugad, and N. Ahuja, "A Fast Scheme for Image Size Change in the Compressed Domain", *IEEE Transactions on Circuits and Systems for Video Technology*, pp.461-474, Vol.11, No. 4, April 2001.
- [6] Y.S. Park, and H.W. Park, "Arbitrary-Ratio Image Resizing using Fast DCT of Composite Length for DCT-based Transcoder" IEEE Transactions on Image Processing, pp.494-500, Vol. 15, No.2, Feb. 2006.
- [7] M. Ishikawa, "Structral Learning with Forgetting" Neural Networks, Vol. 9, No. 3, pp.509-521, 1996.
- [8] B. K. Natarajan, V. Bhaskaran, "A Fast Approximate Algorithm for Scaling Down Digital Images in the DCT Domain", in *Proc. IEEE Int. Conf. Image Processing*'1995, pp. 241-243.
- [9] W. K. Pratt, *Digital Image Processing*, John Wiley & Sons, Inc, 1991.