

Network Working Group  
Internet-Draft  
Intended status: Informational  
Expires: September 10, 2015

N. Egge  
T. Terriberry  
Mozilla Corporation  
March 9, 2015

Time Domain Lapped Transforms for Video Coding  
draft-egge-videocodec-tdlt-01

Abstract

This proposes the use of Time Domain Lapped Transforms (TDLT) as the transform step for video coding.

Status of This Memo

This Internet-Draft is submitted in full conformance with the provisions of BCP 78 and BCP 79.

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

This Internet-Draft will expire on September 10, 2015.

Copyright Notice

Copyright (c) 2015 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to BCP 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

This document may not be modified, and derivative works of it may not be created, and it may not be published except as an Internet-Draft.

## Table of Contents

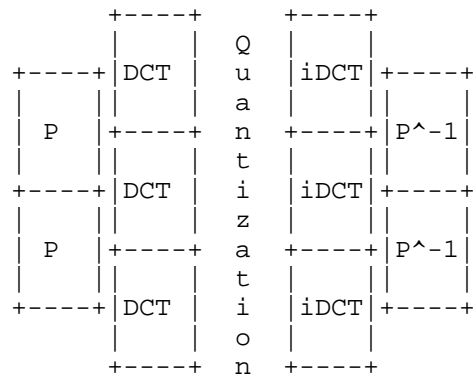
1. Introduction . . . . .	2
2. TDLT Defined . . . . .	2
3. Lapped-Transform Selection . . . . .	4
3.1. Coding Gain . . . . .	5
3.2. Transform Width . . . . .	6
4. Optimal Transform Coefficients . . . . .	6
4.1. Exhaustive Search . . . . .	6
4.2. Stochastic Search . . . . .	7
4.3. Ramp Constraint . . . . .	8
5. Intra Prediction . . . . .	10
6. Motion Compensation . . . . .	11
7. Multiple Block Sizes . . . . .	11
7.1. Variable Sized Lapping . . . . .	12
7.2. Fixed Sized Lapping . . . . .	15
8. IANA Considerations . . . . .	16
9. Security Considerations . . . . .	16
10. Acknowledgments . . . . .	16
11. Informative References . . . . .	16
Authors' Addresses . . . . .	17

## 1. Introduction

This draft outlines a proposal to adapt the Time-Domain Lapped Transforms (TDLT) for use in video coding. Lapped transforms were proposed for video coding at least as far back as 1989 [Malv89]. Like the loop filters more commonly found in recent video coding standards, TDLTs use a post-processing filter that runs between block edges to reduce or eliminate blocking artifacts. Unlike a loop filter, the TDLT filter is invertible, allowing the encoder to run the inverse filter on the input video. This decorrelates blocks before they are passed through a normal block transform and quantization step, improving coding gain (which helps in both smooth and highly textured areas), in addition to reducing blocking artifacts.

## 2. TDLT Defined

The Time-Domain Lapped Transform can be viewed as a set of pre and post filters to an existing block-based DCT transform. The idea is to place an invertible filter along the block boundaries outside an existing block-based DCT encoder.



$$V = \begin{bmatrix} 1 & q_0 \\ & \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ & \\ p_0 & 1 \end{bmatrix} \begin{bmatrix} s_0 & 0 \\ & \\ 0 & s_1 \end{bmatrix}$$

Thus for the 4x8 case, the pre-filter and post-filter are completely described by the four parameters  $q_0$ ,  $p_0$ ,  $s_0$ , and  $s_1$ . In general, any invertible  $V$  matrix may be used. However, factoring  $V$  into a series of lifting steps ensures that it can be implemented efficiently, and can reduce the number of parameters required by the optimization process, since the full flexibility of an arbitrary invertible matrix is not required to achieve good coding gain. [Tran01] proposes two reduced-parameter factorizations, dubbed Type III and Type IV. These are identical in the 4x8 case, but for larger transforms the differ in the order that the  $p_i$  and  $q_i$  steps are applied: interleaved for Type III and ascending and then descending order for Type IV. While Type III appears to give slightly higher coding gain when unconstrained, when coupled with the ramp constraint discussed below and the constraint that all coefficients be dyadic rationals, the number of feasible solutions is much smaller than with Type IV. The increased number of feasible solutions allows Type IV transforms to achieve higher coding gains than Type III when these constraints are imposed. This definition easily extends to the 8x16 and 16x32 TDLT case with similar parameterizations. In general, we use the Type IV factorization from [Tran01]. For a  $V$  matrix of size  $M$ , this has  $(M-1)$   $p_i$  and  $(M-1)$   $q_i$  parameters, and  $M$   $s_i$  parameters. For a transform of size  $N \times 2N$ , this gives a total of  $1.5N-2$  parameters. This is also the number of lifting steps that must be performed to implement the  $V$  portion of the pre- and post-filters.

### 3. Lapped-Transform Selection

We would like to find good candidate transform coefficients that perform well within a video coding framework. There are several metrics we can use for evaluating pre-filter parameters. Including

1. Coding Gain - how well energy is compacted into only a few coefficients
2. Side Band Attenuation - how much energy from frequencies outside the passband leaks into each basis function
3. Transform Width - how wide are the basis functions and how much ringing they will cause
4. Orthogonality - how linearly independent the basis functions are

Of these, the most important by far is coding gain as it allows us to directly measure the improvement in bits between different candidate transforms. At high bit rates using an efficient quantizer, every 6.02 dB improvement in coding gain saves a bit of entropy per coefficient.

### 3.1. Coding Gain

Coding gain is a useful metric for comparing different candidate transforms. Roughly speaking, it is the measure of how well energy is compacted into only a few coefficients. The formula for coding gain of the lapped transform can be found in [Terr12]. Using an AR(1) model with  $r=0.95$ , we have

$$C_g = 10 \cdot \text{Log}_{10} \left| \frac{1}{\prod_i ((G \cdot \text{AR}(1) \cdot G^T)[i,i] \cdot (H^T \cdot H)[i,i])} \right|$$

where  $G$  is the analysis filter of the lapped transform:

$$G = \begin{bmatrix} & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \end{bmatrix} \begin{bmatrix} P & 0 \\ & \\ & \\ & \\ & \\ & \\ & \\ & \end{bmatrix}$$

and  $H$  is the synthesis filter of the lapped transform:

$$H = \begin{bmatrix} & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & & \end{bmatrix} \begin{bmatrix} 0 & \\ i\text{DCT} & \\ 0 & \end{bmatrix}$$

In [Terr12] the coding gain of the non-lapped DCT is compared with the optimal non-lapped Karhunen-Loeve transform for the same AR(1) model with  $r=0.95$ .

	4 point	8 point	16 point
DCT	7.5701 dB	8.8259 dB	9.4555 dB
KLT	7.5825 dB	8.8462 dB	9.4781 dB

Similarly, in [Tran01] the coding gain of the TDLT using fast factorizations with real coefficients produced by unconstrained optimization are

	4x8	8x16	16x32
Type III TDLT	8.6349 dB	9.6115 dB	9.9496 dB
Type IV TDLT	8.6349 dB	9.6005 dB	9.9057 dB

### 3.2. Transform Width

In general, the wider the transform, the higher the coding gain: a 16-point DCT will always have a higher coding gain than a 4-point DCT. In the case of lapped transform, the width of the transform is more than just counting the number of points, it involves the shape of the basis functions. At equal coding gain, a narrower transform is better because it causes a smaller amount of ringing around edges. We define the width of the transform as

$$w = C * \left| \frac{\sum_{ij} (H[i,j]^2 * (j-N+1/2)^4)}{\sum_{ij} (H[i,j]^2)} \right|^{1/4},$$

where  $C=2.991$  is a constant calibrated such that the width of the 1024-point non-overlapped DCT is equal to 1024.

## 4. Optimal Transform Coefficients

Of the four metrics described in Section 3 we chose to optimize our transform parameters for the highest coding gain.

To avoid the use of floating point operations, we use dyadic rationals to represent the parameters of our TDLT. These are the  $p$ 's,  $q$ 's and  $s$ 's that describe the  $V$  matrix in the pre-filter. We chose a base of  $2^6$  because it offered enough resolution to find good approximations of the optimal values for the  $p$ 's,  $q$ 's, and  $s$ 's and still allowed us to fit the results of multiplications in a 16 bit word. Increasing the base to  $2^8$  improves the achievable coding gain of the 4x8 transform by less than 0.002 dB. On the other hand, dropping it even one bit to  $2^5$  lowers the coding gain by 0.037 dB.

### 4.1. Exhaustive Search

For the smaller lapped transforms, it is possible to simply do an exhaustive search and check all possible transform candidates to find the one with the best coding gain. The limitation that the  $p$ 's,  $q$ 's, and  $s$ 's all be dyadic rationals allows us to simply enumerate all reasonable values. Additional constraints allowed us to further reduce the search space. Because the  $p$ 's and  $q$ 's are liftings steps that represent rotations in the plane their values are between -1.0

and 1.0. Likewise the limitation that the pre- and post-filter steps be reversible requires that the scale factors be greater than or equal 1.0, otherwise information would be lost during the transform. Finally, all things equal we prefer smaller scale factors as it makes quantizing and encoding the coefficients cheaper. We thus cap the scale factors at 2.0. Based on some limited experimentation, scale factors larger than this do not appear to produce useful transforms according to our metrics, anyway.

With a dyadic rational base of  $2^6$ , the number of possible candidates to consider is

$$\begin{aligned} |C| &= (2*(2^6)+1)^{(|p|+|q|)} * (2^6+1)^{|s|} \\ &= (2*(2^6)+1)^{(2*(N/2-1))} * (2^6+1)^{(N/2)} \end{aligned}$$

Thus for the transform sizes we are interested in, the number of candidates is tractable only for the 4x8 case:

		N	C
+-----+-----+			
4x8	TDLT	4	68161536
8x16	TDLT	8	$7.731400 * 10^{19}$
16x32	TDLT	16	$9.947082 * 10^{43}$
+-----+-----+			

An exhaustive search for parameters that give the optimal coding gain for the 4x8 TDLT are below:

+-----+-----+		+-----+-----+		+-----+-----+	
p_0	-11/64	q_0	36/64	s_0	91/64
+-----+-----+		+-----+-----+		+-----+-----+	
				s_1	85/64
				+-----+-----+	

#### 4.2. Stochastic Search

For the larger lapped transforms, doing an exhaustive search is not possible. Instead we formulate the optimization problem as an integer programming problem and use a robust industrial solver to find optimal integer values for the p's, q's, and s's.

For the 8x16 TDLT, the parameters are below:

+-----+-----+		+-----+-----+		+-----+-----+	
p_0	-23/64	q_0	48/64	s_0	90/64
p_1	-18/64	q_1	34/64	s_1	73/64
p_2	-6/64	q_2	20/64	s_2	72/64
+-----+-----+		+-----+-----+		+-----+-----+	
				s_3	75/64
				+-----+-----+	

For the 16x32 TDLT, the parameters are below:

p_0	-24/64	q_0	50/64	s_0	90/64
p_1	-23/64	q_1	40/64	s_1	74/64
p_2	-17/64	q_2	31/64	s_2	73/64
p_3	-12/64	q_3	22/64	s_3	71/64
p_4	-14/64	q_4	18/64	s_4	67/64
p_5	-13/64	q_5	16/64	s_5	67/64
p_6	-7/64	q_6	11/64	s_6	67/64
				s_7	72/64

In order to confirm that the integer approximations found are in fact optimal, we can compare them with the optimal real valued coding gains for the three lapped-transforms we are proposing. In [Tran01] a numeric solver was used to find optimal values for a Type IV lapped transform.

	4x8	8x16	16x32
Real Valued	8.6349 dB	9.6005 dB	9.9057 dB
Approximate	8.63473 dB	9.60021 dB	9.89338 dB
Loss	0.00017 dB	0.00029 dB	0.01232 dB

#### 4.3. Ramp Constraint

It is also possible to constrain the lapped transform so that it is (1,2)-regular [DT03], i.e., that it has one vanishing moment in the analysis filter and two vanishing moments in the synthesis filter. This allows the synthesis filter to reconstruct any piecewise linear function solely from the DC coefficients. This causes the shape of the DC basis function to be a symmetric linear ramp. This can be particularly useful when it matches the shape of other windowing functions used in the codec. For example, a linear window is commonly used with Overlapped Block Motion Compensation (OBMC), which is one possible approach for avoiding blocking artifacts in the motion-compensation stage of the codec. More vanishing moments are possible, allowing reconstruction of piecewise quadratic or even higher-order functions, but these require additional overlap stages.

This regularity can be enforced solely by enforcing a series of constraints on the scale factors,  $s_i$ .



$$s_0 = N \cdot (1 - q_0)$$

$$s_i = \frac{N}{2^i + 1} * \left\lfloor \frac{(q_{i-1} - 1) * p_{i-1} - q_i}{N} \right\rfloor, \text{ for } i > 0$$

Since  $2^i + 1$  is odd, but we want  $s_i$  to be a dyadic rational value, the remainder of the expression must be evenly divisible by  $(2^i + 1)$ . A similar set of constraints can be derived for Type III, but they involve more of the  $p$ 's and  $q$ 's per  $s_i$  value, and thus have far fewer admissible solutions when coupled with the dyadic rational constraint.

The additional restrictions described above greatly reduce the number of combinations to consider, both because there are fewer parameters (the  $s_i$ 's can no longer be chosen independently) and because there are fewer combinations of parameter values which produce dyadic rational coefficients. With these constraints, the number of combinations is small enough that an exhaustive search is now tractable for the 8x16 TDLT.

	N	C
4x8 TDLT	4	442
8x16 TDLT	8	331677320

An exhaustive search for parameters that give the optimal coding gain under the ramp and dyadic rational constraints for the 4x8 and 8x16 TDLT are below:

<table><tr><td>p_0</td><td>-16/64</td></tr></table>	p_0	-16/64	<table><tr><td>q_0</td><td>41/64</td></tr></table>	q_0	41/64	<table><tr><td>s_0</td><td>92/64</td></tr><tr><td>s_1</td><td>93/64</td></tr></table>	s_0	92/64	s_1	93/64												
p_0	-16/64																					
q_0	41/64																					
s_0	92/64																					
s_1	93/64																					
<table><tr><td>p_0</td><td>-24/64</td></tr><tr><td>p_1</td><td>-20/64</td></tr><tr><td>p_2</td><td>-4/64</td></tr></table>	p_0	-24/64	p_1	-20/64	p_2	-4/64	<table><tr><td>q_0</td><td>53/64</td></tr><tr><td>q_1</td><td>40/64</td></tr><tr><td>q_2</td><td>24/64</td></tr></table>	q_0	53/64	q_1	40/64	q_2	24/64	<table><tr><td>s_0</td><td>88/64</td></tr><tr><td>s_1</td><td>75/64</td></tr><tr><td>s_2</td><td>76/64</td></tr><tr><td>s_3</td><td>76/64</td></tr></table>	s_0	88/64	s_1	75/64	s_2	76/64	s_3	76/64
p_0	-24/64																					
p_1	-20/64																					
p_2	-4/64																					
q_0	53/64																					
q_1	40/64																					
q_2	24/64																					
s_0	88/64																					
s_1	75/64																					
s_2	76/64																					
s_3	76/64																					

Unfortunately, in the 16x32 TDLT case the number of combinations is still not tractable, even with these additional constraints. Again, we use an integer programming model to solve for the integer parameters that optimize coding gain in this context.

p_0	-32/64	q_0	59/64	s_0	80/64
p_1	-28/64	q_1	53/64	s_1	72/64
p_2	-24/64	q_2	46/64	s_2	73/64
p_3	-32/64	q_3	41/64	s_3	68/64
p_4	-24/64	q_4	35/64	s_4	72/64
p_5	-13/64	q_5	24/64	s_5	74/64
p_6	-2/64	q_6	12/64	s_6	74/64
				s_7	70/64

	4x8	8x16	16x32
Dyadic	8.63473 dB	9.60021 dB	9.89338 dB
Ramp + Dyadic	8.59886 dB	9.56161 dB	9.78294 dB
Loss	0.03587 dB	0.0386 dB	0.11044 dB

## 5. Intra Prediction

Since the final pixel values of a block are not available until after the post-filter runs, they cannot be used to predict neighboring blocks. There are a number of possible solutions to this. For example, one could simply use pixels from outside the overlap region. However, as these pixels are farther away, they are poorer predictors, and the extra distance reduces the range of prediction directions which have enough neighbors available to form an adequate extrapolation [OP11].

An alternate approach is to perform the prediction in the frequency domain. Initial experiments suggest that this is just as effective as prediction in the time domain, and has similar computational requirements [Eggel3]. However, because the frequency domain coefficients of a neighboring block are impacted both by what size DCT was used, and the lapping across all four of its edges, directional predictors can only be so good. At low rates, this meant more bits were spent correcting an incorrect predictor than were saved by coding only a directional mode.

A signal free technique was developed for doing limited intra prediction in the frequency domain when using lapped transforms. Note that when the spatial prediction mode is exactly horizontal or vertical, applying the filters described in this draft along the orthogonal direction is the identity. Thus it is possible to look at the horizontal coefficients of the neighboring block to the left, and the vertical energy of the neighboring block above and simply use the coefficients where the energy is larger. When this technique is

coupled with a quantization and coefficient coder that makes signaling no predictor cheap [Vali15], this becomes an effective frequency domain intra predictor.

Finally, a technique was developed for intra predicting chroma frequency domain coefficients from decoded coincident luma coefficients [Egge15]. While this technique does not strictly require the use of lapped transforms, because the block size extent (and thus the lapping region) for both the chroma and luma planes is the same, the use of a lapped transform does not change the effectiveness of this technique.

## 6. Motion Compensation

There have been several lapped transform proposals that perform block-by-block motion compensation by simply expanding the size of the prediction region for each block [TT01], [OPT11]. However, in addition to increasing the amount of motion-compensated prediction pixels that must be computed by a factor of four, this also increases the number of applications of the pre- and post-filter by a factor of four, since this must now be done separately for each block, using the motion-compensated frame difference for that block.

An alternate approach is simply perform motion compensation of the frame in a completely separate step, prior to any transform, using any method desired [Terr15]. The lapping can then be applied to this motion-compensated prediction, producing per-block predictors. This still allows the prediction mode (inter, intra, bi-prediction, etc.) to be chosen on a block-by-block basis. It also interacts well with other techniques designed to operate in the frequency domain, such as the Pyramid Vector Quantization (PVQ) proposed elsewhere.

The downside is that motion estimation in the encoder needs to be performed for regions slightly beyond the current block. However, this is already required by blocking-artifact-free motion compensation techniques, such as Overlapped Block Motion Compensation (OBMC). Experience with OBMC has shown that an encoder can mostly ignore look-ahead and still get acceptable results, unlike other techniques, such as control-grid interpolation (CGI).

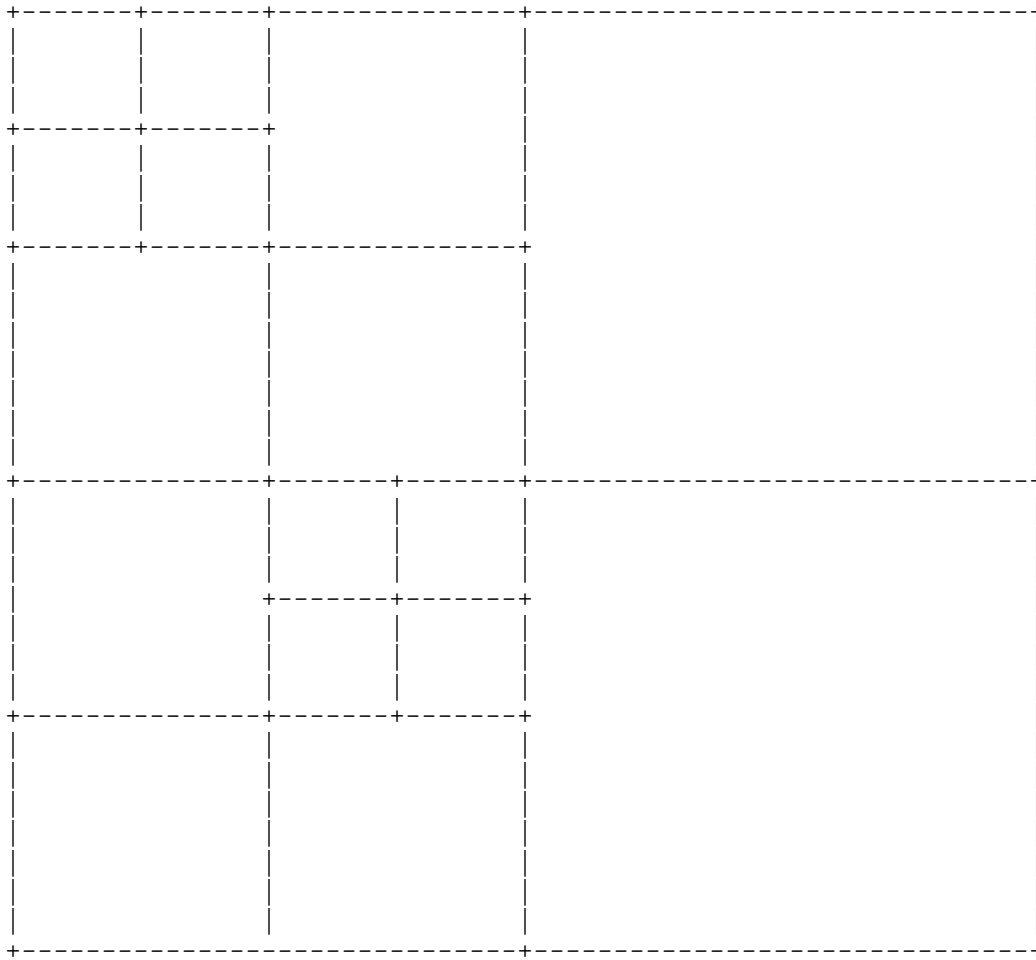
## 7. Multiple Block Sizes

Multiple block size support is important for lapped transforms, since the larger support region increases their susceptibility to ringing artifacts compared to a non-overlapped transform with the same number of coefficients (though it is greatly reduced compared to a non-overlapped transform with a support region of the same size).

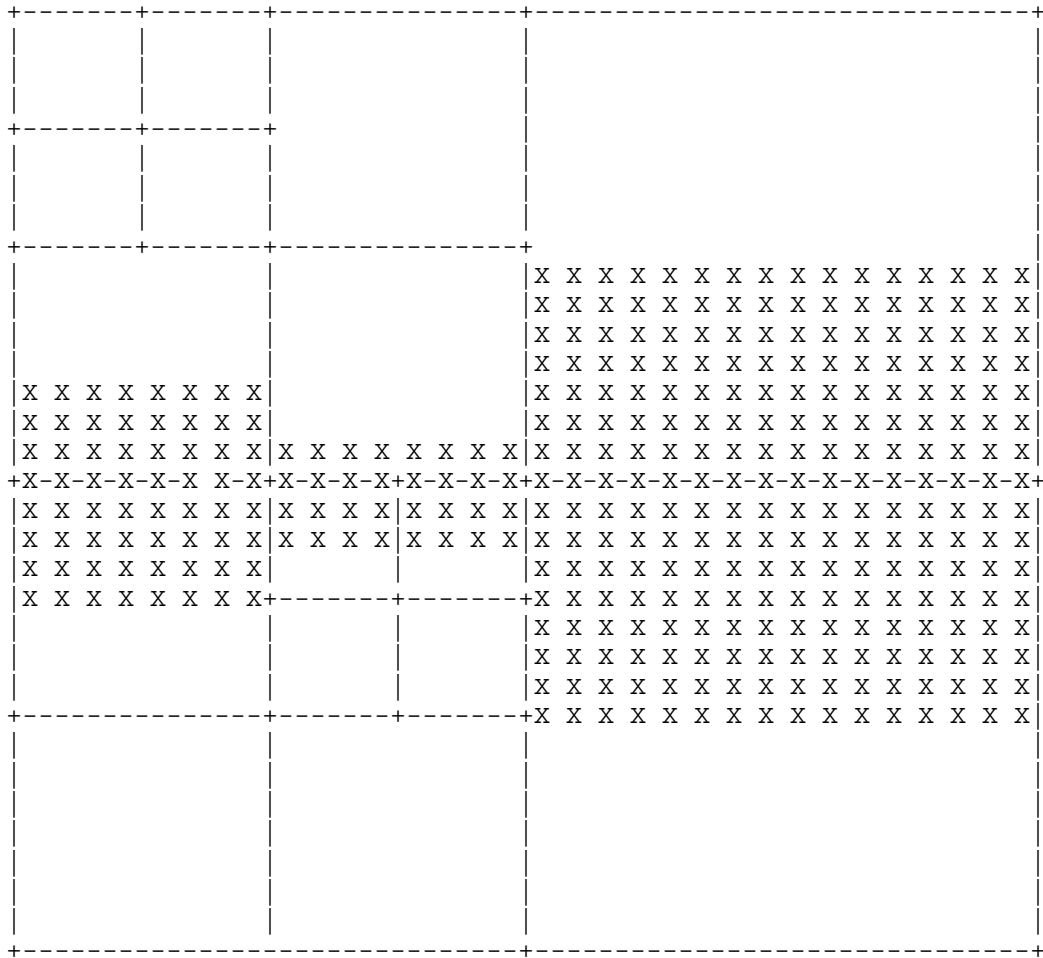
### 7.1. Variable Sized Lapping

The most obvious approach is to require that the size of the overlap filter be constrained by the smallest block adjacent to a given edge. This requires some amount of look-ahead in the encoder, but has the benefit of using the largest lapping possible in regions where all blocks are the same size while not introducing discontinuities where blocks of different sizes meet. Note that this has an effect on the coding syntax, as the block size decision for the block below the one being coded must be made and communicated to the decoder prior to coding. Using this convention no additional information need to be communicated other than the block size decision to completely describe how the variable sized lapping should be applied.

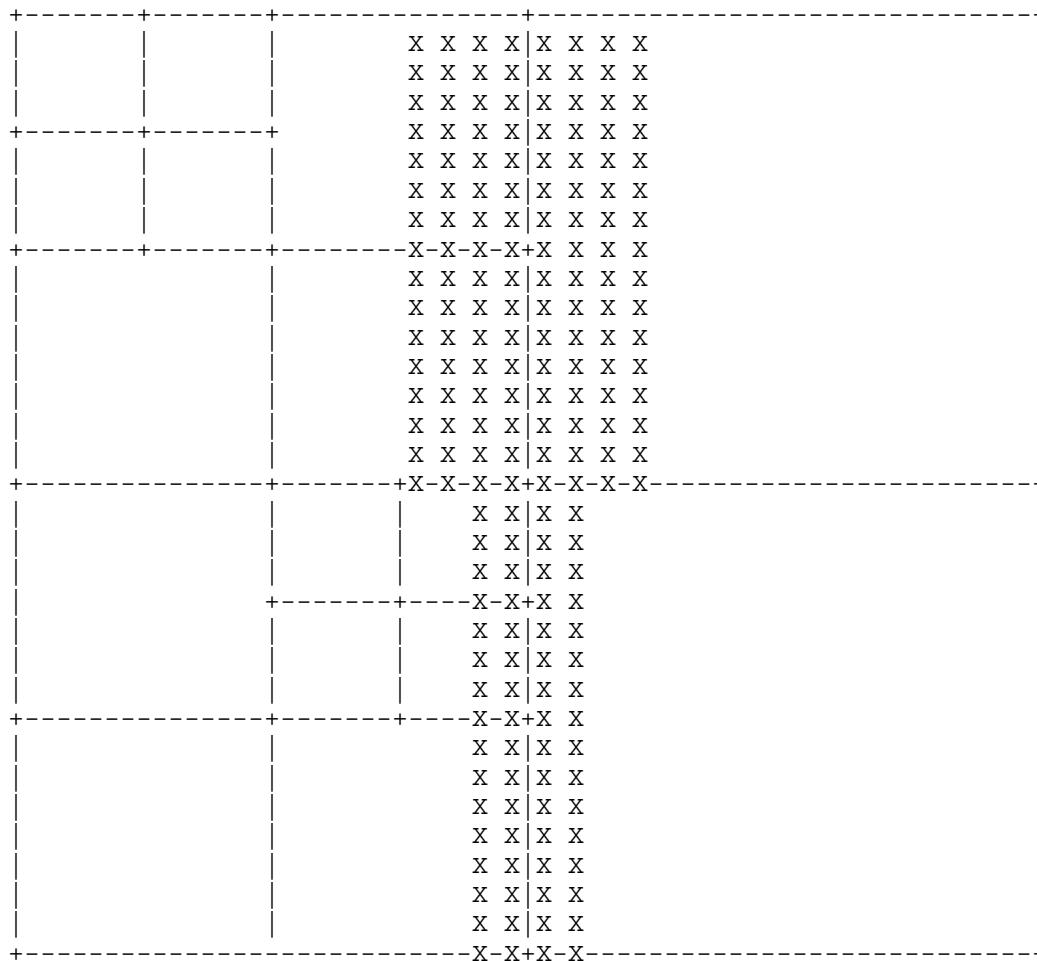
Consider an example image that is 32x32 with the following block size decisions. We apply the lapping recursively to blocks of 32x32 at a time until we reach a block that is not subdivided into smaller blocks. At each step in the recursions, we apply a filter vertically across the block edges that run left to right splitting the block in half. We then apply a horizontal filter across the block edges that split the block in half top to bottom.



Block size decision for a 32x32 frame.



Apply the filter vertically across the horizontal internal edge.



Apply the filter horizontally across the vertical internal edge.

The filters are then applied recursively in this manner to the four quadrants of the block. By applying the filters recursively this way, we have prevented any discontinuities from appearing where block is split but its neighbor is not.

## 7.2. Fixed Sized Lapping

One of the challenges using variable sized lapping is that changing the block size decision (either splitting a block into four blocks a quarter as big, or merging four blocks into one four times the size) can have an impact on the coding performance outside the block considered. This makes computing the optimal block size decision for

a frame computationally difficult as traditional rate-distortion optimization (RDO) algorithms exploit this locality to iteratively improve an initial decision.

One way to simplify the problem is to assume a fixed sized lapping across the entire image. If only the 4-point filter is used across block boundaries, then it is possible to compare the distortion of an 8x8 block with that of four 4x4 blocks by simply computing the mean squared error (MSE) of the 64 spatial domain coefficients after applying the inverse lapped transform. Because changing the block size decision, and thus the interior lapping has no impact on the lap decision on the border of the 8x8 block, then just looking at the rate and distortion of the interior coefficients is sufficient.

This approach has does not leverage the additional coding gain and deblocking achieved by using larger lapping filters but may make up for this by allowing computationally cheap block size decision heuristics in real-time encoding environments.

#### 8. IANA Considerations

This document has no actions for IANA.

#### 9. Security Considerations

This draft has no security considerations.

#### 10. Acknowledgments

Thanks to Greg Maxwell and Jean-Marc Valin for their assistance in the experimentation and other valuable contributions to this document.

#### 11. Informative References

- [DT03] Dai, W. and T. Tran, "Regularity-Constrained Pre- and Post-Filtering for Block DCT Based Systems", IEEE Transactions on Signal Processing 51(10):2568--2581, October 2003.
- [Egge13] Egge, N., "Intra-Prediction in Daala", October 2013, <<http://people.xiph.org/~unlord/Daala-Intra.pdf>>.
- [Egge15] Egge, N. and J. Valin, "Predicting Chroma from Luma with Frequency Domain Intra Prediction", February 2015, <[http://people.xiph.org/~unlord/spie\\_cfl.pdf](http://people.xiph.org/~unlord/spie_cfl.pdf)>.



- [Malv89] Malvar, H. and D. Staelin, "The LOT: Transform Coding Without Blocking Effects", IEEE Transactions on Acoustics, Speech, and Signal Processing , April 1989, <<http://research.microsoft.com/apps/pubs/default.aspx?id=102073>>.
- [OP11] de Oliveria, R. and B. Pesquet-Popescu, "Intra-Frame Prediction with Lapped Transforms for Image Coding", Proc. of the 36th IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP'11) pp. 805--808, May 2011.
- [OPT11] de Oliveria, R., Pesquet-Popescu, B., and M. Trocan, "Inter Prediction Using Lapped Transforms for Advanced Video Coding", Proc. of the 18th IEEE International Conference on Image Processing (ICIP'11) pp. 3705--3708, September 2011.
- [Tran01] Tran, T., "Lapped Transform via Time-Domain Pre- and Post-Processing", IEEE Transactions on Signal Processing , October 2001.
- [TT01] Tran, T. and C. Tu, "Lapped Transform Based Video Coding", Proc. of the 24th SPIE Conference on Applications of Digital Image Processing vol. 4472, pp. 319--333, July 2001.
- [Terr12] Terriberry, T., "Introduction to Video Coding Part 1: Transform Coding", February 2012.
- [Terr15] Terriberry, T., "Adaptive Motion Compensation Without Blocking Artifacts", February 2015, <<https://people.xiph.org/~tterribe/daala/vbsobmc.pdf>>.
- [Vali15] Valin, J., "Perceptual Vector Quantization for Video Coding", February 2015, <[http://jmvalin.ca/video/spie\\_pvq.pdf](http://jmvalin.ca/video/spie_pvq.pdf)>.

## Authors' Addresses

Nathan E. Egge  
Mozilla Corporation  
650 Castro Street  
Mountain View, CA 94041  
USA

Email: [negge@dgql.org](mailto:negge@dgql.org)

Timothy B. Terriberry  
Mozilla Corporation  
650 Castro Street  
Mountain View, CA 94041  
USA

Phone: +1 650 903-0800  
Email: tterribe@xiph.org