

# Combined Neural Network-based Intra Prediction and Transform Selection

Thierry Dumas, Franck Galpin, Philippe Bordes  
*Interdigital, Rennes, France*

thierry.dumas@interdigital.com, franck.galpin@interdigital.com, philippe.bordes@interdigital.com

**Abstract**— The interactions between different tools added successively to a block-based video codec are critical to its rate-distortion efficiency. In particular, when deep neural network-based intra prediction modes are inserted into a block-based video codec, as the neural network-based prediction function cannot be easily characterized, the adaptation of the transform selection process to the new modes can hardly be performed manually. That is why this paper presents a combined neural network-based intra prediction and transform selection for a block-based video codec. When putting a single neural network-based intra prediction mode and the learned prediction of the selected LFNST pair index into VTM-8.0,  $-3.71\%$ ,  $-3.17\%$ , and  $-3.37\%$  of mean BD-rate reduction in all-intra is obtained.

**Index Terms**—Transform signaling, intra prediction, neural networks, Versatile Video Coding.

## I. INTRODUCTION

In a block-based video codec featuring multiple transforms, the signaling of the selected inverse transform to be applied to a block of reconstructed transform coefficients at the decoder can take two forms. In the first form, a.k.a explicit signaling, the encoder writes to the bitstream the selected transform index. This way, the decoder identifies the selected inverse transform by reading its index from the bitstream. In the second form, a.k.a implicit signaling, the decoder derives the selected transform index from available information.

Given that the implicit transform signaling does not spend any bit, it makes particular sense in terms of rate-distortion when the used available information correlates with the efficiency of the transforms at compacting the residual block energy into few transform coefficients. This is well illustrated by the Low Frequency Non-Separable Transform (LFNST) [5] in Versatile Video Coding (VVC). During the LFNST training, the 67 VVC intra prediction modes are first divided into groups. Then, for each group, a pair of LFNST matrices is trained on blocks of primary transform coefficients, each of them arising from the application of the DCT2-DCT2 to the residue resulting from the intra prediction of an image block via a mode of this group. Consequently, in the VVC decoder, if a block of reconstructed transform coefficients uses LFNST, its selected intra prediction mode index (the used available information) maps to the index of the pair of LFNST matrices trained on data generated via this mode<sup>1</sup>, i.e. the one probably having the highest energy compaction efficiency in the current case. Note that the selected LFNST matrix index among the chosen pair is explicitly signaled.

Unfortunately, an implicit transform signaling deriving from the selected intra prediction mode index can hardly benefit to an intra prediction mode added to the block-based video codec afterwards as no straightforward correlation exists between the new mode index and the energy compaction efficiencies of the transforms of interest. This happened to the implicit LFNST signaling when introducing the Matrix-based Intra Prediction (MIP) [8] modes in VTM-5.0. To correct this, a fixed mapping from each MIP mode index to one of the 67 VVC modes indices considered by the implicit LFNST signaling appeared in VTM-5.0. Note that, since VTM-6.0, the mapping from each MIP mode index to 0 has replaced the fixed mapping in VTM-5.0, thus showing its low rate-distortion impact.

More critically, if the new intra prediction mode is made of deep neural networks, a solution similar to the above-mentioned mapping from the new mode index to an intra prediction mode index considered by the implicit transform signaling becomes unfeasible for two reasons. Firstly, similarities between the prediction of a block via the new non-linear mode and the predictions of this block via the existing linear modes cannot sometimes be found, preventing the association of modes indices. Secondly, as the characteristics of the deep neural network intra prediction of a block, e.g. horizontal propagation of the pixel intensities from decoded samples around this block into the predicted block, depend on the decoded neighboring samples [4], [6], [10], [11], this kind of mapping cannot be fixed at the encoder and the decoder.

Alternatively, this paper proposes to learn the selection of the transform index from an intermediate representation of the neural network prediction of the block. This way, for a block predicted via the neural network-based mode, the implicit transform signaling adapts to the decoded samples surrounding this block and the non-linear prediction function. Note that, in [9], a neural network predicts the selected transform index from the block of quantized transform coefficients. Differently, our approach aims at integrating into a block-based video codec both a new coding tool, i.e. a neural network-based intra prediction mode, and a complement to the implicit transform signaling now depending on this new coding tool.

When inserting a single additional neural network-based intra prediction mode and the learned prediction of the selected LFNST pair index into VTM-8.0,  $-3.71\%$ ,  $-3.17\%$  and  $-3.37\%$  of mean BD-rate reduction in all-intra is reported.

<sup>1</sup>See "g\_lfnstLut" at [https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware\\_VTM-8.0/blob/master/Tools/Transform/TransformLut/TransformLut.cpp](https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware_VTM-8.0/blob/master/Tools/Transform/TransformLut/TransformLut.cpp)

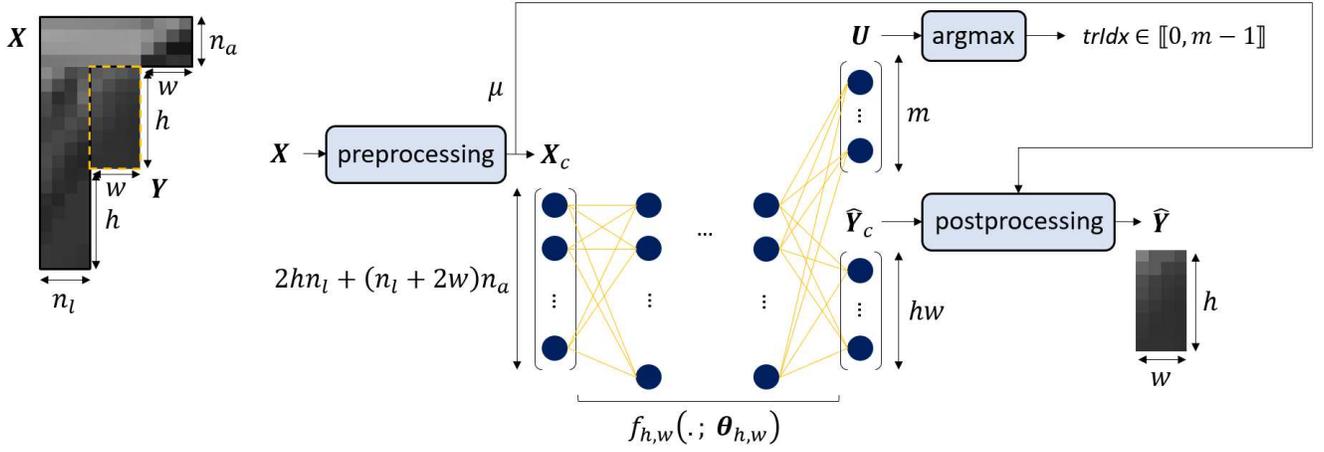


Fig. 1. Prediction  $\hat{Y}$  of a  $w \times h$  block  $Y$  and inference of the selected transform index  $trIdx$  from the block context  $X$  via the neural network  $f_{h,w}(\cdot; \theta_{h,w})$ . Here, the implicit transform signaling features the selection of one transform among  $m$  possible transforms.  $\mu$  gathers pre-processing variables needed by the post-processing step, see an example in Section IV-A.

## II. NEURAL NETWORK-BASED TRANSFORM SELECTION

The method developed in this paper targets the integration of neural network-based intra prediction modes and the learned transform index selection into a block-based video codec. Two factors orient our study towards considering a single neural network-based mode. Firstly, the addition of a single neural network-based intra prediction mode to VVC, this mode having a relatively small signaling cost with respect to those of the 67 VVC intra prediction modes and the MIP modes, has already shown significant rate-distortion gains [3]. Secondly, when the index returned by a learned mapping is used for predictive coding, see Section III-C, the extra encoder running time caused by the additional tests of transforms remains small only in the case of a single neural network-based mode. Despite this orientation, our approach can be easily generalized to multiple new neural network-based intra prediction modes.

In this single neural network-based intra prediction mode, blocks of size  $w \times h$  in the codec are predicted by the neural network  $f_{h,w}(\cdot; \theta_{h,w})$ , parametrized by  $\theta_{h,w}$ . Indeed, as a neural network for intra prediction takes the L-shape context around a block to provide the predicted block, see Figure 1, its architecture must include full-connections, making the number of neural network parameters dependent on the block size. Therefore, the single neural network-based mode contains card ( $Q$ ) neural networks,  $Q$  denoting the set of possible pairs of block height and width in the codec.

Given the composition of the single additional mode, each of its neural network must own a different mapping to the selected transform index. More precisely, the neural network  $f_{h,w}(\cdot; \theta_{h,w})$  takes a pre-processed version  $X_c$  of the context  $X$  of decoded samples around a  $w \times h$  block  $Y$  to return both a prediction  $\hat{Y}_c$  of  $Y$  before post-processing and a vector  $U$  of unscaled log-probability of each transform to be selected, see Figure 1. Then, the post-processing turns  $\hat{Y}_c$  into the final prediction  $\hat{Y}$  of  $Y$  and the selected transform index  $trIdx$  according to the neural network corresponds to the position of the maximum in  $U$ . Note that a fully-connected

architecture is displayed in Figure 1. Yet, Figure 1 can be adapted to a convolutional architecture like the one in Section IV-A. Note also that the penultimate neural representation in  $f_{h,w}(\cdot; \theta_{h,w})$  is chosen as input to the fully-connected layer returning  $U$  because, in the conditions of the experiments in Section IV-A, it has been observed that this input yields the best accuracy of the classification of the selected transform.

## III. NEURAL NETWORK-BASED LFNST SELECTION

This section applies the generic combined neural network-based intra prediction and transform index selection in Section II to the implicit LFNST signaling in VVC. First, the indexing of the secondary transforms of LFNST is changed to suit the neural network-based transform index selection framework, see Section III-A. Then, Sections III-B and III-C describe respectively the training of the neural network-based LFNST selection and the signaling of the proposed method in VVC.

### A. Modified indexing of the secondary transforms of LFNST

The LFNST signaling mixes an implicit signaling and an explicit one. Regarding the implicit signaling, the index of the intra prediction mode selected to predict the current Coding Block (CB) determines the index of the pair of LFNST matrices among four pairs and whether the primary transform coefficients resulting from the application of the DCT2 horizontally and the DCT2 vertically to the residue of prediction are transposed, see Table I. Note that this implicit signaling stems from the LFNST training summarized in Section I. Regarding the explicit signaling,  $lfnstIdx = 0$  means that the encoding/decoding of the current CB does not use LFNST.  $lfnstIdx \in \{1, 2\}$  indicates the selected LFNST matrix index among the pair given by the implicit LFNST signaling.

As seven different pairs of transformations of the current primary transform coefficients before their potential transposition into secondary transform coefficients are actually possible, see Table I, the index  $trPairIdx \in [0, 6]$  is introduced to connect  $U$  to the implicit LFNST signaling. Moreover, the constraint of following the same implicit signaling for the different

TABLE I

DEFINITION OF THE INDEX  $trPairIdx$  LINKING  $\mathbf{U}$  TO THE IMPLICIT LFNST SIGNALING. IN THE ROW “WIDE ANGLE INTRA MODE INDEX”, EACH RANGE BETWEEN BRACKETS REFERS A GROUP OF INDICES OF VVC INTRA PREDICTION MODES, EXCLUDING THE MIP MODES. (\*) DENOTES  $\{0, 1\} \cup$  THE MIP MODES INDICES. (\*\*) DENOTES  $[-14, -1] \cup [2, 12]$ . THE “TRANSFORM SET INDEX” INDEXES THE PAIRS OF LFNST MATRICES.

wide angle intra mode index	(*)	(**)	[[13, 23]]	[[24, 34]]	[[35, 44]]	[[45, 55]]	[[56, 83]]
transform set index [5]	0	1	2	3	3	2	1
transposition of the primary transform coeffs	false	false	false	false	true	true	true
$trPairIdx$	0	1	2	3	4	5	6

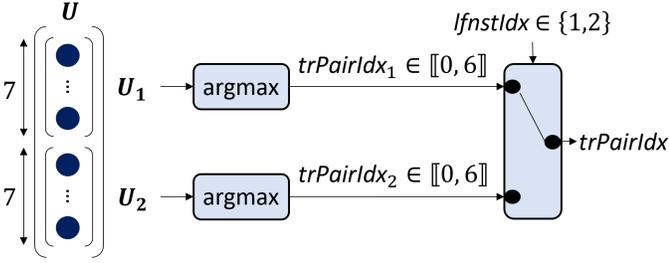


Fig. 2. Computation of the index  $trPairIdx$  of the pair of secondary transforms selected according to the neural network model from  $\mathbf{U}$ .

values  $lfnstIdx \in \{1, 2\}$  can be removed to gain flexibility. Then,  $trPairIdx$  is duplicated into  $\{trPairIdx_{lfnstIdx}\}_{lfnstIdx \in \{1, 2\}}$ .

Given the new indexing, at the VVC encoder, for a given  $w \times h$  block, the neural network  $f_{h,w}(\cdot; \theta_{h,w})$  computes from the pre-processed version of the context of this block the unscaled log-probability of each of the seven pairs of secondary transforms for  $lfnstIdx = 1$  and the unscaled log-probability of each of the seven pairs of secondary transforms for  $lfnstIdx = 2$ . Then, for each  $lfnstIdx \in \{1, 2\}$ ,  $trPairIdx_{lfnstIdx}$  is the position of the largest unscaled log-probability. Finally, the selected pair index  $trPairIdx$  in  $\{trPairIdx_1, trPairIdx_2\}$  depends on the value of  $lfnstIdx$  found by the encoder. Note that, at the VVC decoder, the same procedure applies, except that the value of  $lfnstIdx$  is read from the bitstream. Thus, for LFNST, the top-right part of Figure 1 becomes Figure 2.

### B. Training of the neural network-based LFNST selection

During the training of the neural network-based LFNST selection, as two learned LFNST selections for each value  $lfnstIdx \in \{1, 2\}$  are considered separately, for each training example, the objective function should involve the ground truth selected secondary transform indices for  $lfnstIdx = 1$  and  $lfnstIdx = 2$  respectively as classification labels. Therefore, the objective function  $\mathcal{L}(S_{h,w}; \phi_{h,w})$  to be minimized over the parameters  $\phi_{h,w}$  in the branch of  $f_{h,w}(\cdot; \theta_{h,w})$  dedicated

$$diffIdx = trExpldx > trPairIdx ? trExpldx - 1 : trExpldx$$

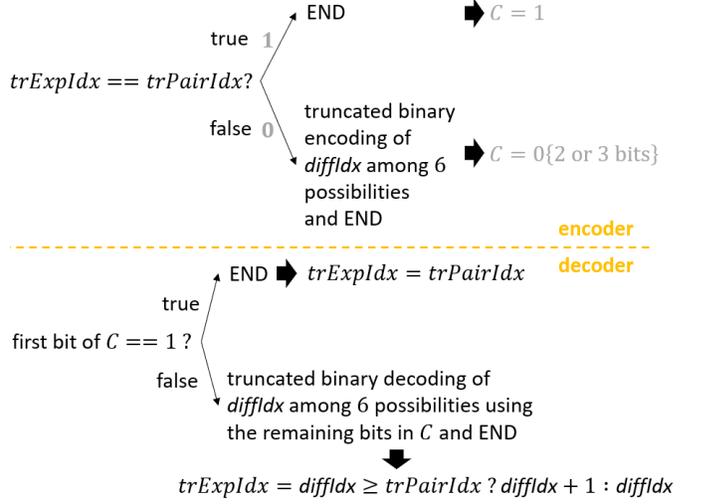


Fig. 3. Predictive encoding and decoding of the pair index  $trExpldx$  with respect to its prediction  $trPairIdx$  from the neural network model.  $C$  is the code of the remainder of the predictive encoding written to the bitstream.

to the LFNST selections is expressed as

$$\mathcal{L}(S_{h,w}; \theta_{h,w}) = \frac{1}{N} \sum_{(\mathbf{X}_c, i_1, i_2) \in S_{h,w}} -\mathcal{H}(\mathbf{X}_c, i_1, i_2; \theta_{h,w})$$

$$\mathcal{H}(\mathbf{X}_c, i_1, i_2; \theta_{h,w}) = \log(\sigma(\mathbf{U}[0:7])_{i_1}) + \log(\sigma(\mathbf{U}[7:])_{i_2})$$

$$\{\hat{\mathbf{Y}}_c, \mathbf{U}\} = f_{h,w}(\mathbf{X}_c; \theta_{h,w}) \text{ and } N = \text{card}(S_{h,w})$$

where  $S_{h,w}$  denotes the training set of triplets of the pre-processed version  $\mathbf{X}_c$  of the context of a  $w \times h$  block, the index  $i_1$  of the selected secondary transform for  $lfnstIdx = 1$  when encoding this block via VVC, and the index  $i_2$  of the selected secondary transform for  $lfnstIdx = 2$ .  $\sigma$  denotes the softmax. Note that the array indexation  $\mathbf{U}[0:7]$  excludes the coefficient of  $\mathbf{U}$  of index 7, as in C++. The training hyperparameters will be detailed right before the experiments in Section IV-A as they depend on the chosen neural network architectures.

### C. Signaling in VVC

Up to now, for a CB predicted via the neural network-based intra prediction mode, the index  $trPairIdx$  of the pair of secondary transforms selected according to the neural network model has directly specified the implicit LFNST signaling. This is called the “inference” scheme.

Alternatively, a predictive coding, called “prediction” scheme, can be constructed from the learned LFNST selection. In this case, a new syntax element  $trExpldx$  replaces  $trPairIdx$  in Table I.  $trPairIdx$  in Figure 2 becomes a prediction of  $trExpldx$ . At the VVC encoder, for a CB predicted via the neural network-based mode, for  $lfnstIdx \neq 0$ , the best value of  $trExpldx$  in terms of rate-distortion is found, and the remainder of the predictive coding of  $trExpldx$  with respect to  $trPairIdx$  is written to the bitstream. Figure 3 details the predictive coding used in this work.

Note that, in both the “inference” and “prediction” schemes, for a CB predicted by an intra prediction mode different from the single additional neural network-based intra prediction mode, the transform signaling in VVC remains unchanged.

#### IV. EXPERIMENTS

Now that the proposed learned LFNST selection is specified, its relevance in terms of rate-distortion in VVC can be studied, see Section IV-A. Then, our single neural network-based intra prediction mode with the learned LFNST selection inside VVC is compared to the state-of-the-art, see Section IV-B.

##### A. Rate-distortion analysis of the learned LFNST selection

If, as said in the second paragraph of Section II, in the neural network-based mode, blocks of each possible size in VVC are predicted by a different neural network, the number of neural networks integrated into VVC would be large, and, as a common neural network for intra prediction may contain numerous parameters, the neural network parameters would incur an excessive memory footprint. To circumvent this, only the blocks of each size in  $\{4 \times 4, 8 \times 4, 16 \times 4, 32 \times 4, 8 \times 8, 16 \times 8, 16 \times 16, 32 \times 32\}$  are predicted by a different neural network in the single additional mode. The single additional mode thus comprises 8 neural networks. To maximize the usage rate of the neural network-based mode, the following steps are added. The context  $\mathbf{X}$  of a  $32 \times 16$  block is downsampled horizontally by 2 before the pre-processing step in Figure 1 and the neural network prediction after the post-processing step is interpolated horizontally by 2, making the prediction of this block via  $f_{16,16}(\cdot; \theta_{16,16})$  feasible. The same goes for a  $64 \times 64$  block, but the horizontal and vertical downsampling and interpolation factors are 2, and  $f_{32,32}(\cdot; \theta_{32,32})$  is used for prediction. Besides, for  $(h, w) \in \{(8, 4), (16, 4), (32, 4), (16, 8), (32, 16)\}$ , the context  $\mathbf{X}$  of a  $w \times h$  block is transposed before the pre-processing step and the neural network prediction after the post-processing step is transposed, allowing the prediction of this block via  $f_{w,h}(\cdot; \theta_{w,h})$ .

From now on, the following parametrization of the context of a  $w \times h$  block applies. If  $\min(h, w) \leq 8$ ,  $n_a = n_l = \min(h, w)$ . Otherwise,  $n_a = h/2$  and  $n_l = w/2$ . Moreover, the pre-processing and post-processing steps in Figure 1 correspond to those detailed in [3].

For the training in Section III-B, the RGB images in the ILSVRC 2012 training dataset and those in DIV2K converted into  $YCbCr$  are encoded via VTM-8.0 with Quantization Parameter (QP) drawn from  $\{22, 27, 32, 37\}$  for each image. Each neural network training runs for 800000 iterations with batch size 100, ADAM, and 0.0002 as learning rate.

As this work does not relate to the enhancement of the neural network prediction of a block, the objective function on block prediction is picked from [3] and the neural network architectures in [3] are simply adapted to the learned LFNST selection. Moreover, the architectures are reduced to decrease the modified VVC encoder and decoder running times, see Tables II to V. Note that, when  $\min(h, w) > 8$ ,  $\mathbf{X}_c$  is split into

TABLE II  
ARCHITECTURE OF  $f_{h,w}(\cdot; \theta_{h,w})$  WHERE  $\min(h, w) \leq 8$ . IN THE COLUMN “INPUT”, A NUMBER REFERS TO THE INDEX OF THE LAYER WHOSE OUTPUT IS THE INPUT TO THE CURRENT LAYER. THE LAYERS OF INDICES 3 AND 4 RETURN  $\hat{\mathbf{Y}}_c$  AND  $\mathbf{U}$  RESPECTIVELY.

layer index	input	layer type	number of neurons	non-linearity
1	$\mathbf{X}_c$	fully-connected	1200	LeakyReLU
2	1	fully-connected	1200	LeakyReLU
3	2	fully-connected	$hw$	-
4	2	fully-connected	14	-

TABLE III  
ARCHITECTURE OF THE BRANCH OF  $f_{16,16}(\cdot; \theta_{16,16})$  TAKING THE ABOVE PORTION  $\mathbf{X}_0$  OF  $\mathbf{X}_c$ . FOR  $f_{32,32}(\cdot; \theta_{32,32})$ , THE SAME ARCHITECTURE APPLIES BUT THE STRIDES IN BOLD BECOME (2, 2).

layer index	input	layer type	filter size	nb of filters	stride	non-linearity
1	$\mathbf{X}_0$	convolutional	$3 \times 3 \times 1$	32	(2, 2)	LeakyReLU
2	1	convolutional	$3 \times 3 \times 32$	64	(2, 2)	LeakyReLU
3	2	convolutional	$3 \times 3 \times 64$	128	(1, 2)	LeakyReLU
4	3	convolutional	$3 \times 3 \times 128$	128	(1, 2)	LeakyReLU
5	4	flattening	-	-	-	-

two portions, see [3], each portion being fed into a different convolutional branch of  $f_{h,w}(\cdot; \theta_{h,w})$ .

Finally, for a given CB, the intra signaling of the neural network-based mode explained in [3] is re-used here.

To assess the relevance of the proposed learned LFNST selection, the “inference” and “prediction” schemes must be compared in terms of rate-distortion against a baseline in which, for a given CB predicted via the neural network-based mode, if  $lfnstIdx \in \{1, 2\}$ , the pair of LFNST matrices of transform set index 0 is always chosen, without transposing the primary transform coefficients. This is called the “default” scheme. Another interesting baseline, called “fully explicit LFNST” corresponds to the “prediction” scheme without the neural network-based prediction of  $trExplDx$ . This means that, at the VVC encoder, for a given CB predicted by the neural

TABLE IV  
ARCHITECTURE OF THE BRANCH OF  $f_{16,16}(\cdot; \theta_{16,16})$  TAKING THE LEFT PORTION  $\mathbf{X}_1$  OF  $\mathbf{X}_c$ . FOR  $f_{32,32}(\cdot; \theta_{32,32})$ , THE SAME ARCHITECTURE APPLIES BUT THE STRIDES IN BOLD BECOME (2, 2).

layer index	input	layer type	filter size	nb of filters	stride	non-linearity
6	$\mathbf{X}_1$	convolutional	$3 \times 3 \times 1$	32	(2, 2)	LeakyReLU
7	6	convolutional	$3 \times 3 \times 32$	64	(2, 2)	LeakyReLU
8	7	convolutional	$3 \times 3 \times 64$	128	(2, 1)	LeakyReLU
9	8	convolutional	$3 \times 3 \times 128$	128	(2, 1)	LeakyReLU
10	9	flattening	-	-	-	-

TABLE V  
ARCHITECTURE OF THE PART OF  $f_{16,16}(\cdot; \theta_{16,16})$  MERGING THE OUTPUTS OF THE TWO BRANCHES IN TABLES III AND IV. THE LAYERS OF INDICES 13 AND 14 RETURN  $\hat{\mathbf{Y}}_c$  AND  $\mathbf{U}$  RESPECTIVELY. FOR  $f_{32,32}(\cdot; \theta_{32,32})$ , THE SAME ARCHITECTURE APPLIES BUT THE LAYER OF INDEX 12 CONTAINS  $hw$  NEURONS INSTEAD OF 500.

layer index	input	layer type	number of neurons	non-linearity
11	5 and 10	concatenation	-	-
12	11	fully-connected	500	LeakyReLU
13	12	fully-connected	$hw$	-
14	12	fully-connected	14	-

TABLE VI

MEAN BD-RATE REDUCTIONS IN % OF VTM-8.0 WITH THE SINGLE ADDITIONAL NEURAL NETWORK-BASED MODE W.R.T VTM-8.0. ONLY THE FIRST FRAME OF EACH SEQUENCE IN THE JVET CTC [2] IS CONSIDERED. THE LARGEST ABSOLUTE MEAN BD-RATE REDUCTION IN LUMINANCE IS IN BOLD.

Video class	default			fully explicit LFNST			inference			prediction		
	Y	C <sub>b</sub>	C <sub>r</sub>	Y	C <sub>b</sub>	C <sub>r</sub>	Y	C <sub>b</sub>	C <sub>r</sub>	Y	C <sub>b</sub>	C <sub>r</sub>
A1	-4.41	-3.78	-3.45	-5.04	-4.36	-4.35	-5.23	-4.37	-4.59	<b>-5.46</b>	-4.66	-4.74
A2	-2.02	-1.58	-1.91	-2.35	-2.00	-2.29	-2.51	-2.48	-2.13	<b>-2.55</b>	-2.08	-2.62
B	-2.54	-2.35	-2.07	-3.03	-2.27	-2.20	-3.15	-2.40	-2.83	<b>-3.31</b>	-2.78	-2.98
C	-2.53	-1.74	-2.01	-2.88	-2.76	-2.79	-3.04	-2.05	-2.16	<b>-3.06</b>	-2.75	-2.65
D	-2.66	-2.38	-1.22	-3.19	-2.78	-2.72	-3.25	-3.05	-2.57	<b>-3.42</b>	-3.46	-2.51
E	-3.99	-3.32	-3.79	-4.48	-3.23	-3.81	-4.46	-3.87	-4.71	<b>-4.66</b>	-3.97	-4.38
F	-1.46	-1.58	-1.90	-1.68	-1.35	-1.50	-1.71	-2.03	-2.01	<b>-1.97</b>	-2.02	-1.73
Mean	-3.01	-2.49	-2.55	-3.46	-2.84	-2.97	-3.58	-2.91	-3.17	<b>-3.71</b>	-3.17	-3.37

TABLE VII

MEAN ENCODER AND DECODER RUNNING TIMES OF VTM-8.0 INCLUDING THE SINGLE ADDITIONAL NEURAL NETWORK-BASED MODE WITH RESPECT TO VTM-8.0 ON THE JVET CTC IN ALL-INTRA. 100% MEANS THE SAME RUNNING TIME AS VTM-8.0.

	default	fully explicit LFNST	inference	predictive
Encoder	369%	387%	400%	499%
Decoder	3330%	3551%	3854%	5052%

network-based mode, for  $lfnstIdx \in \{1, 2\}$ , the best value of  $trExpIdx$  in terms of rate-distortion is found, and this value is written to the bitstream via a truncated binary encoding. In the experiments, only the first frame of each video sequence of the JVET CTC [2] is considered. The configuration is all-intra.

The most striking remark is that the “default” scheme is much worse than the three other transform selection schemes in terms of rate-distortion, see Table VI. Besides, on the luminance channel, the “prediction” scheme adds  $-0.25\%$  of mean BD-rate reduction with respect to the “fully explicit LFNST”. The “inference” scheme yields  $-0.12\%$  of additional mean BD-rate reduction with respect to the “fully explicit LFNST”. These three observations prove the relevance of the proposed neural network-based transform selections. As the “inference” scheme adds no explicit signaling to VVC, on the luminance channel, the difference of  $-0.57\%$  in mean BD-rate reduction between the “inference” scheme and the “default” scheme can be viewed as the pure rate-distortion gain of the neural network-based transform selections. This gain arises from bit-rate drop at close PSNRs. Table VII reports the mean encoder and decoder running times of each tested scheme.

Note that an approach involving different neural networks for each QP (or range of QPs) was not studied in this work as this type of approach increases significantly the memory footprint of the neural network parameters inside VVC.

### B. Comparison to the state-of-the-art

Most of the previous approaches on the neural network-based intra prediction for block-based video coding integrate their tool into HEVC [4], [6], [7], [10]–[12], the ancestor of VVC. For comparison, our proposed neural network-based transform selection should also be integrated into HEVC. However, this makes little sense in HEVC as, unlike in VVC, the implicit transform signaling in HEVC is extremely limited. Indeed, the DST7-DST7 applying to  $4 \times 4$  luminance blocks predicted in intra instead of the DCT2-DCT2 is the main

implicit transform signaling in HEVC. For a rate-distortion comparison between the neural network-based intra prediction used in this paper without the neural network-based transform selection and several neural network-based intra prediction tools in the literature, all inside HEVC, please see [3].

To our knowledge, as of now, only [1] presents a neural network-based intra prediction tool tested in a recent version of VVC. A neural network-based intra prediction mode for chrominance featuring an attention mechanism is put into VTM-7.0. Using the same test data as in Table VI,  $-0.15\%$ ,  $-0.68\%$ ,  $-0.53\%$  of mean BD-rate reduction is reported in all-intra. The mean encoder and decoder running times of their modified VTM-7.0 w.r.t VTM-7.0 are 120% and 947%.

## V. CONCLUSION

This paper has introduced a combined neural network-based intra prediction and transform selection for a block-based video codec. As the neural network-based transform selection depends on both the decoded samples around the current block and the neural network intra prediction function, it can modelize all the intra-transform correlations. When integrated into VTM-8.0, this approach yields large rate-distortion gains.

## REFERENCES

- [1] M. G. Blanch, S. Blasi, A. Smeaton, N. O’Connor, and M. Mrak, “Chroma intra prediction with attention-based CNN architectures,” *ICIP*, 2020.
- [2] J. Boyce, K. Suehring, X. Li, and V. Seregin, “JVET common test conditions and software reference configurations,” *Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11, 10<sup>th</sup> meeting, San Diego*, April 2018.
- [3] T. Dumas, F. Galpin, and P. Bordes, “Iterative training of neural networks for intra prediction,” *IEEE Transactions on Image Processing*, vol. 30, pp. 697–711, November 2020.
- [4] Y. Hu, W. Yang, M. Li, and J. Liu, “Progressive spatial recurrent neural network for intra prediction,” *IEEE Transactions on Multimedia*, vol. 21, no. 12, pp. 3024–3037, December 2019.
- [5] M. Koo, M. Salehifar, J. Lim, and S.-H. Kim, “Low Frequency Non-Separable Transform (LFNST),” *PCS*, 2019.
- [6] J. Li, B. Li, J. Xu, R. Xiong, and W. Gao, “Fully-connected network-based intra prediction for image coding,” *IEEE Transactions on Image Processing*, vol. 27, no. 7, pp. 3236–3247, July 2018.
- [7] Y. Li, L. Li, Z. Li, J. Yang, N. Xu, D. Liu, and H. Li, “A hybrid neural network for chroma intra prediction,” in *ICIP*, 2018.
- [8] J. Pfaff, B. Stallenberger, M. Schäfer, P. Merkle, P. Helle, T. Hinz, H. Schwarz, D. Marpe, and T. Wiegand, “Affine linear weighted intra prediction,” *Joint Video Experts Team (JVET) of ITU-T SG 16 WP 3 and ISO/IEC JTC 1/SC 29/WG 11, 14<sup>th</sup> meeting, Geneva*, March 2019.
- [9] S. Puri, S. Lasserre, and P. Le Callet, “CNN-based transform index prediction in multiple transforms framework to assist entropy coding,” *EUSIPCO*, 2017.

- [10] H. Sun, L. Yu, and J. Katto, "Fully neural network mode-based intra prediction of variable block size." *VCIP*, 2020.
- [11] Y. Wang, X. Fan, S. Liu, D. Zhao, and W. Gao, "Multi-scale convolutional neural network-based intra prediction for video coding," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 7, pp. 1803–1815, July 2020.
- [12] L. Zhu, S. Kwong, Y. Zhang, S. Wang, and X. Wang, "Generative adversarial network-based intra prediction for video coding," *IEEE Transactions on Multimedia*, vol. 22, no. 1, January 2020.