# Fast Mode Decision for H.264/AVC Using Mode Prediction

Song-Hak Ri and Joern Ostermann

Institut fuer Informationsverarbeitung, Appelstr 9A, D-30167 Hannover, Germany
ri@tnt.uni-hannover.de
ostermann@tnt.uni-hannover.de

**Abstract.** In this paper, we present a new method to speed up the mode decision process using mode prediction. In general, video coding exploits spatial and temporal redundancies between video blocks, in particular temporal redundancy is a crucial key to compress video sequence with little loss of image quality. The proposed method determines the best coding mode of a given macroblock by predicting the mode and its rate-distortion (RD) cost from neighboring MBs in time and space. Compared to the H.264/AVC reference software, the simulation results show that the proposed method can save up to 53% total encoding time with up to 2.4% bit rate increase at the same PSNR.

## 1 Introduction

Video coding plays an important role in multimedia communications and consumer electronics applications. The H.264/AVC is the latest international video coding standard jointly developed by the ITU-T Video Coding Experts Group and the ISO/IEC Moving Picture Experts Group. It can achieve higher coding efficiency than that of previous standards, such as MPEG-4 and H.263 [1].

However, it requires a huge amount of computational loads due to use of the variable block-size motion estimation, intra prediction in P slice coding, quarter-pixel motion compensation, multiple reference frames, etc. The complexity analysis described in [1] shows that examining all possible modes takes the most time out of the total encoding time. Hence, fast mode decision making becomes more and more important.

H.264/AVC Baseline profile employs seven different block sizes for inter frames. The size of a block can be $16 \times 16$, $16 \times 8$, $8 \times 16$, or $8 \times 8$, and each $8 \times 8$ can be further broken down to sub-macroblocks of size $8 \times 8$, $8 \times 4$, $4 \times 8$, or $4 \times 4$, as shown in Fig.1. To encode a given macroblock, H.264/AVC encoder tries all possible prediction modes in the following order; SKIP, Inter$16 \times 16$, Inter$16 \times 8$, Inter$8 \times 16$, Inter$8 \times 8$, Inter$8 \times 4$, Inter$4 \times 8$, Inter$4 \times 4$, Intra$4 \times 4$, Intra$8 \times 8$, Intra$16 \times 16$. The SKIP mode represents the case in which the block size is $16 \times 16$ but no motion and no residual information are coded. Except for SKIP and intra modes, each inter mode decision requires a motion estimation process.

In order to achieve the highest coding efficiency, H.264/AVC uses rate distortion optimization techniques to get the best coding results in terms of maximizing

**Fig. 1.** Variable block sizes in H.264/AVC

coding quality and minimizing coded data bits. The mode decision is made by comparing the rate distortion cost of each possible mode and by selecting the mode with the lowest rate distortion cost as the best one.

The existing fast mode decision algorithms can be classified into two categories: The first class is to find the optimal mode by using some features, such as texture and edge information, which are computed from the raw video data. D. S. Turaga et al [5] and J. Chen et al [2] introduce the so-called mean removed mean absolute difference (mrMAD) and use the feature to make fast intra and inter mode decision. In [6], the $3 \times 3$ Sobel operator is used to get the edge map of a whole frame. The edge map and the gradient are both employed to find the best interpolation direction as the best intra mode. They also use the edge map to determine whether a macroblock is homogeneous in order to find the best inter mode. However, the algorithm has to evaluate all the pixels in the whole frame and it leads to high computational complexity.

The second class is trying to make full use of the relationship among the modes and predicts the best mode by using the already checked modes and their statistical data. A representative method [4] of such class divides all modes into 3 groups. Using one mode from each group, the best group is determined. All modes of the best group are evaluated to determine the best mode selection. Thus the number of candidate modes is greatly reduced. In [3], the most probable mode is predicted based on the observation that most modes are spatially correlated in a given frame. If the predicted mode satisfies some conditions which estimate if the predicted mode is the best mode, the encoder codes the macroblock with the predicted mode. Thus it can skip all of the calculations on other modes.

Based on the analysis above, we propose a novel algorithm to determine the best mode based on RD optimization by using the combination of spatial and temporal mode prediction. We investigate whether it is possible to temporally or spatially predict the best mode.

This paper is organized as follows: Section 2 shows a consideration about the possibility of mode prediction for fast mode decision. In Section 3, we propose a new mode decision scheme by combined mode prediction, and finally, experimental results and conclusions are presented in Section 4 and Section 5, respectively.

## 2   Mode Prediction

Video coding is achieved by reducing spatial and temporal redundancies between video frames. This implies indirectly that a mode of a given macroblock (MB hereafter) also might be correlated to that of MBs neighboring in space and time. It was noted that there was a spatial mode correlation between a given MB and its neighboring MBs and, therefore, it is possible to spatially predict a mode of the MB [3].

Since a video sequence contains, in general, more redundancies in time domain than in space domain, we stipulate that temporal mode correlation is higher than spatial mode correlation. Thus we consider spatial, temporal and spatial-temporal prediction of the best mode for a given MB.

In order to do that, we must answer these two questions:

1. How *high* is the correlation of spatial and temporal modes?
2. Is it necessary to consider *all* modes for the mode prediction?

Let's mark the current MB as $X$, collocated MB of $X$ in the previous frame as $X_{-1}$ and neighboring MBs as $A$, $B$, $C$ and $D$ (see Fig.2).

To compare correlation of both mode predictions, let's define the following 3 events:

$E_S$ : Modes of 2, 3 or 4 MBs out of $A$, $B$, $C$ and $D$ are the same as
      the RD-optimal mode of $X$.
$E_T$ : The mode of $X_{-1}$ is the same as the RD-optimal mode of $X$.
$E_C$ : $E_S \cup E_T$

Here, $E_S$, $E_T$ and $E_C$ mean spatial, temporal and combined mode events. Table 1 shows the probabilities ($P_S$, $P_T$ and $P_C$) of each event for the video sequences *container*, *mother&daughter*, *stefan*, *mobile*, *foreman* and *coastguard*.

In Table 1, it is found that the probability of spatial mode prediction is lower than the probability of temporal mode prediction and, of course, combined mode prediction is also greater than spatial or temporal mode correlation. In the case of sequences such as *container* and *mother&daughter*, which are characterized by slow and smooth motion, the probability of a spatial mode event is similar to the temporal mode event. In the case of some sequences, such as *foreman* and *coastguard*, which are characterized by fast motion, the probability of a spatial mode event is far lower than that of the temporal mode event. The table tells us that by using combined mode correlation, the encoder can predict the best mode of a given MB more frequently than by using spatial mode correlation. From now on, the combined mode prediction will be called mode prediction.

To answer to the second question, let's calculate the probability of an event where the predicted mode of a given MB is SKIP, Inter16×16, Inter16×8, Inter8×16, Sub8×8, Intra4×4 and Intra16×16, under the condition that $X$ has the same prediction mode with $X_{-1}$ (see Table 2). Let's mark the

**Fig. 2.** The current MB $X$, its collocated MB $X_{-1}$ of the previous frame, and neighboring MBs, $A$, $B$, $C$ and $D$

**Table 1.** Comparison of an occurrence probability of spatial, temporal and combined mode events, QP (Quantization Parameter)=28, QCIF, 100 frames

| Sequences | $P_S(\%)$ | $P_T(\%)$ | $P_C(\%)$ |
|---|---|---|---|
| container | 46.9 | 53.6 | 68.7 |
| mother-daughter | 33.8 | 40.3 | 56.5 |
| stefan | 13.8 | 31.2 | 41.7 |
| mobile | 10.2 | 27.4 | 35.5 |
| foreman | 8.8 | 29.8 | 34.2 |
| coastguard | 8.5 | 24.0 | 32.1 |

**Table 2.** Statistics of modewise-temporal mode correlation in case $X$ and $X_{-1}$ have the same RD-optimal mode (unit=%), QP=28, QCIF, 100 frames

| Sequences | $P_{SKIP}$ | $P_{Inter16\times16}$ | $P_{Inter16\times8}$ | $P_{Inter8\times16}$ | $P_{Sub8\times8}$ | $P_{Intra4\times4}$ | $P_{Intra16\times16}$ |
|---|---|---|---|---|---|---|---|
| container | 79.5 | 9.4 | 3.6 | 3.7 | 3.7 | 0.0 | 0.1 |
| mother-da. | 64.9 | 16.6 | 6.1 | 6.9 | 5.2 | 0.1 | 0.0 |
| stefan | 21.1 | 32.2 | 10.3 | 15.1 | 20.2 | 0.6 | 0.5 |
| mobile | 17.4 | 27.2 | 13.1 | 10.2 | 31.7 | 0.1 | 0.4 |
| foreman | 12.3 | 39.8 | 13.5 | 11.7 | 22.4 | 0.2 | 0.1 |
| coastguard | 3.4 | 19.7 | 16.0 | 14.2 | 46.6 | 0.0 | 0.0 |

probability $P(\text{SKIP}|X=X_{-1})$ as $P_{SKIP}$, $P(\text{Inter}16\times16|X=X_{-1})$ as $P_{Inter16\times16}$, $\dots, P(\text{Intra}4\times4|X=X_{-1})$ as $P_{Intra4\times4}$ and $P(\text{Intra}16\times16|X=X_{-1})$ as $P_{Intra16\times16}$.

As seen in Table 2, when the predicted mode of $X$ equals the actual best mode, which can be calculated by the exhaustive mode decision of JM reference software, the occurrence probabilities of Intra4×4 and Intra16×16 are very low. This probability is also very low at other QP values, too. Therefore, we don't use the predicted modes, Intra4×4 and Intra16×16, as candidates for the best mode of a given MB, if the predicted mode is Intra4×4 or Intra16×16.

## 3   Fast Mode Decision by Mode Prediction

The most important thing for applying mode prediction to fast mode decision is to make sure that the predicted mode is the best mode for a given MB. So far, there have been several ways to decide whether the predicted mode is the best mode of the MB or not.

The most common method [3] adopts a threshold value derived from the RD cost which is already calculated. The threshold is set to an average of RD costs of neighboring MB with the same mode and it is compared with the RD cost of the given MB $X$ with the predicted mode to estimate if it is the best mode. Another method [7] adopts the square of the quantization parameter (QP) as a threshold to decide whether the predicted mode is to be used.

For the sequence $foreman$, the RD cost difference between the spatially predicted mode and the optimal mode is shown in Fig 3. The size of this difference does not necessarily depend on the actual RD cost. Therefore, a threshold based on neighboring MBs or QP should not be used for evaluating the quality of the predicted mode.



**Fig. 3.** Relationship between RD cost of $X$ and average RD cost of neighboring MBs with the same mode (spatial RD cost prediction)

We use the RD cost of $X_{-1}$ as the threshold. Fig.4 intuitively shows a relationship between the actually optimal RD cost of $X$ and the optimal RD cost of $X_{-1}$ when the optimal mode of $X$ is the same as one of $X_{-1}$. From Fig. 5 and Table 3, it also should be noted that the correlation of both RD costs is great even in the case that the optimal mode of $X_{-1}$ is not the same as one of $X$, which means that optimal RD cost of a MB can be predicted by one of the previous MB.

**Fig. 4.** Relationship between RD cost of $X$ and RD cost of $X_{-1}$ when the best mode of $X$ is the same as one of $X_{-1}$



**Fig. 5.** Relationship between RD cost of $X$ and RD cost of $X_{-1}$

Such a relationship can be seen in the comparison of the following three correlation coefficients; correlation coefficient ($\rho_S$) between the spatially predicted RD cost and the optimal RD cost, correlation coefficient ($\rho_{T'}$) between the actually optimal RD cost of $X$ and the optimal RD cost of $X_{-1}$ when the optimal mode of $X$ is the same as one of $X_{-1}$, and correlation coefficient ($\rho_T$) between the actually optimal RD cost of $X$ and the optimal RD cost of $X_{-1}$. Table 3 shows that a temporal correlation is greater than a spatial one.

**Table 3.** Comparison of three correlation coefficients in QCIF and CIF format

| Sequences | QCIF | | | CIF | | |
|---|---|---|---|---|---|---|
| | $\rho_S$ | $\rho_{T'}$ | $\rho_T$ | $\rho_S$ | $\rho_{T'}$ | $\rho_T$ |
| foreman | 0.722 | 0.949 | 0.922 | 0.683 | 0.952 | 0.939 |
| coastguard | 0.772 | 0.942 | 0.933 | 0.560 | 0.934 | 0.921 |
| stefan | 0.870 | 0.969 | 0.957 | 0.779 | 0.975 | 0.972 |
| mother-daughter | 0.814 | 0.979 | 0.964 | 0.789 | 0.987 | 0.976 |
| mobile | 0.485 | 0.974 | 0.970 | 0.358 | 0.964 | 0.965 |
| container | 0.764 | 0.988 | 0.976 | 0.508 | 0.993 | 0.983 |
| average | 0.738 | 0.967 | 0.954 | 0.613 | 0.968 | 0.959 |



**Fig. 6.** Probabilities at which the chosen, the eliminated candidate or other mode is the same as the RD-optimal mode (average probability in the case of other mode)

Another problem in using mode prediction might be error propagation, due to the misprediction of the best mode. To prevent the propagation of mode prediction errors, it is expected that an exhaustive mode decision will be carried out periodically. In the experiment of the proposed algorithm, a phenomena of error propagation is likely to happen more frequently in video sequences with smooth and slow motion, resulting in some increase of the total bit rate.

The last problem of mode prediction is that using only one predicted mode to decide upon the best mode could be unstable. It has been observed that sometimes temporal mode prediction shows better result than spatial one, and also vice versa. Therefore we apply two mode candidates, $m_t$ and $m_s$, predicted temporally and spatially to a given MB and choose the mode with the lower RD cost. Fig. 6 shows the three probabilities: the red curve is the probability that the chosen candidate, $m_t$ or $m_s$, is the best mode, the blue curve the probability that the other mode, $m_t$ or $m_s$, is the best one and the green curve shows the probability that a different mode is the best mode. As one can see, the probability that the chosen candidate is the best mode is far higher than the probability of a different mode.

The proposed algorithm is as follows:

Step 1: if the current frame is an exhaustive mode decision frame, check all modes and stop mode decision.

Step 2: get two predicted modes, $m_t$ and $m_s$, from temporal and spatial mode predictions.

Step 3: get the RD cost, $RD_{Pred}$, of the collocated MB in the previous frame with its already known best mode.

Step 4: if both predicted modes are the same, apply it to the current MB, otherwise, compare the two RD costs by applying both and choose the better one.

Step 5: if the chosen RD is lower than the threshold, $TH = \alpha \cdot RD_{Pred}$, set it to the best mode and stop, otherwise check all other modes (here, $\alpha$ is a positive constant derived from experiment).

## 4    Experimental Results

The proposed fast mode decision scheme was implemented in H.264/AVC reference software JM 10.1 baseline profile for performance evaluation. The experimental conditions are as follows:

Software & Profile : H.264/AVC reference software JM 10.1 Base-line

Sequences :          $container$, $coastguard$, $stefan$, $foreman$, $mobile$, $mother\&daughter$

Video Format :       QCIF, CIF

ME Strategy :        Full Motion Estimation

The proposed algorithm was evaluated based on the exhaustive RDO mode decision of H.264/AVC in the following performance measures:

- Degradation of image quality in term of average Y-PSNR: $\triangle$PSNR (dB)
- Increase of bit rate: +Bits (%)
- Prediction rate: PR (%)

$$PR = \frac{N_{Pred}}{N_{Total}} \times 100(\%),$$

where, $N_{Total}$ is total number of MBs and $N_{Pred}$ is the number of the mode prediction successes.

- Encoding time saving: TS (%)

$$TS = \frac{T_{REF} - T_{PROP}}{T_{REF}} \times 100(\%),$$

where, $T_{REF}$ and $T_{PROP}$ are the total encoding times of REFerence and PROPosed method, respectively.

In the experiment, the exhaustive mode decision is implemented at an interval of 10 frames, to prevent error propagation. $\alpha$ is set to 1.1.

We compared the performance of the proposed method with that of an alternative method which is based on spatial mode prediction [3]. For the QCIF video format, Table 4 shows that the proposed algorithm can achieve 44% of

**Table 4.** The comparison in the performance measures, QP=24, QCIF, 100 frames, where the alternative method is spatial mode prediction based method [3]

| sequences | Alternative | | | Proposal | | | |
|---|---|---|---|---|---|---|---|
| | △PSNR(dB) | +Bits(%) | TS(%) | △PSNR(dB) | +Bits(%) | TS(%) | PR(%) |
| mobile | -0.05 | 3.6 | 24.6 | 0.00 | 1.4 | 42.5 | 45.6 |
| stefan | -0.06 | 4.7 | 29.6 | 0.04 | 1.9 | 43.2 | 49.6 |
| foreman | -0.01 | 3.5 | 26.5 | -0.05 | 2.9 | 35.3 | 37.8 |
| mother-da. | -0.03 | 3.9 | 35.3 | -0.02 | 1.2 | 46.4 | 51.5 |
| container | -0.02 | 3.7 | 26.2 | 0.01 | 1.7 | 35.1 | 40.7 |
| coastguard | -0.01 | 2.3 | 23.3 | -0.02 | 2.0 | 35.3 | 38.9 |
| average | -0.03 | 3.6 | 27.6 | -0.01 | 1.9 | 39.6 | 44.0 |

**Table 5.** The comparison in the performance measures, QP=24, CIF, 100 frames, where the alternative method is spatial mode prediction based method [3]

| sequences | Alternative | | | Proposal | | | |
|---|---|---|---|---|---|---|---|
| | △PSNR(dB) | +Bits(%) | TS(%) | △PSNR(dB) | +Bits(%) | TS(%) | PR(%) |
| mobile | -0.10 | 3.1 | 29.7 | -0.01 | 2.9 | 52.9 | 58.3 |
| stefan | -0.08 | 3.7 | 28.6 | -0.02 | 2.5 | 47.4 | 55.5 |
| foreman | -0.09 | 3.5 | 38.2 | -0.05 | 2.1 | 41.6 | 49.8 |
| mother-da. | -0.03 | 3.9 | 46.3 | -0.03 | 1.6 | 46.0 | 52.4 |
| container | -0.10 | 3.3 | 37.5 | -0.05 | 2.6 | 48.2 | 55.6 |
| coastguard | -0.10 | 2.3 | 26.1 | -0.10 | 1.8 | 51.5 | 63.6 |
| average | -0.08 | 3.3 | 34.4 | -0.04 | 2.4 | 47.9 | 55.9 |



**Fig. 7.** Comparison of several RD (PSNR vs. bitrate) plots, QP=24, 28 and 32

average time savings in total encoding time with 0.01dB of PSNR degradation and 1.9% of extra bits. For the CIF video format, Table 5 shows better performance compared to that of the QCIF case resulting in about 48% of time saving with 0.04dB PSNR degradation and 2.4% of extra bits. The proposed algorithm shows better performance than that of the alternative algorithm [3], which is achieving about 27% of average time savings, 0.03dB of PSNR degradation and 3.6% of extra bits for QCIF sequences, and 35% of average time savings, 0.08dB of PSNR degradation and 3.3% of extra bits for CIF sequence.

Table 4 and Table 5 also show that prediction rate of the best mode depends on the contents and resolutions of video sequence, that is, how slow or fast motion is, and how fine the spatial resolution is. Fig 7 shows the rate-distortion performance of the three algorithms, the exhaustive method, alternative method (based on spatial mode prediction) and the proposed method. The curve shows that the proposed algorithm has better RD efficiency than the alternative method, achieving similar efficiency to the exhaustive method.

## 5 Conclusions

In this paper, we proposed a new method to speed up mode decision process using mode prediction. The proposed method determines the best coding mode of a given macroblock by predicting the mode and its RD cost from neighboring MBs in time and space. Compared to the H.264/AVC reference software, the simulation result shows that the proposed method can save up to 53% of the total encoding time with up to 2.4% bit rate increase at the same PSNR.

## References

1. Ostermann, J., Bormans, J., List, P., Marpe, D., Narroschke, M., Pereira, F., Stockhammer, T., Wedi, T.: Video Coding with H.264/AVC: Tools, Performance, and Complexity, IEEE Circuits and Systems Magazine, Vol. 4, 1(2004) 7-28
2. Chen, J., Qu, Y., He, Y.: A Fast Mode Decision Algorithm in H.264, 2004 Picture Coding Symposium (PCS2004), San Francisco, CA, USA (2004)
3. Chang, C.-Y., Pan, C.-H., Chen, H.: Fast Mode Decision for P-Frames in H.264, 2004 Picture Coding Symposium (PCS2004), San Francisco, CA, USA (2004)
4. Yin, P., Tourapis, H.-Y.C., Tourapis, A. M., Boyce, J.: Fast mode decision and motion estimation for JVT/H.264, in Proceedings of International Conference on Image Processing (2003) 853-856
5. Turaga, D. S., Chen, T.: Estimation and Mode Decision for Spatially correlated Motion Sequences, IEEE Transactions on Circuits and Systems for Video Technology, Vol. 11, 10(2001) 1098-1107
6. Lim, K. P., Wu, S., Wu, D. J., Rahardja, S., Lin, X., Pan, F., Li, Z. G.: Fast Inter Mode Selection, JVT-I020, JVT 9th Meeting, san Diego, USA (2003)
7. Arsura, E., Caccia, G., Vecchio, L. D., Lancini, R.: JVT/H.264 Rate-Distortion Optimization based on Skipping Mechanism and Clustering Mode Selection using MPEG7 Transcoding Hints, 2004 Picture Coding Symposium (PCS2004), San Fransisco, CA, USA (2004)