

Parallel Scalability of Video Decoders

Technical Report

June 30, 2008

Cor Meenderinck, Arnaldo Azevedo and Ben Juurlink
Computer Engineering Department
Delft University of Technology
{Cor, Azevedo, Benj}@ce.et.tudelft.nl

Mauricio Alvarez
Technical University of Catalonia (UPC)
Barcelona, Spain
alvarez@ac.upc.edu

Alex Ramirez
Barcelona Supercomputing Center (BSC)
Barcelona, Spain
alex.ramirez@bsc.es



Abstract

An important question is whether emerging and future applications exhibit sufficient parallelism, in particular thread-level parallelism (TLP), to exploit the large numbers of cores future CMPs are expected to contain. As a case study we investigate the parallelism available in video decoders, an important application domain now and in the future. Specifically, we analyze the parallel scalability of the H.264 decoding process. First we discuss the data structures and dependencies of H.264 and show what types of parallelism it allows to be exploited. We also show that previously proposed parallelization strategies such as slice-level, frame-level, and intra-frame macroblock (MB) level parallelism, are not sufficiently scalable. Based on the observation that inter-frame dependencies have a limited spatial range we propose a new parallelization strategy, called Dynamic 3D-Wave. It allows certain MBs of consecutive frames to be decoded in parallel. Using this new strategy we analyze the limits to the available MB-level parallelism in H.264. Using real movie sequences we find a maximum MB parallelism ranging from 4000 to 7000. We also perform a case study to assess the practical value and possibilities of a highly parallelized H.264 application. The results show that H.264 exhibits sufficient parallelism to efficiently exploit the capabilities of future manycore CMPs.

Contents

1	Introduction	3
2	Overview of the H.264 Standard	5
3	Benchmark	10
4	Parallelizing H.264	11
4.1	Task-level Decomposition	11
4.2	Data-level Decomposition	12
4.2.1	GOP-level Parallelism	13
4.2.2	Frame-level Parallelism for Independent Frames	13
4.2.3	Slice-level Parallelism	14
4.2.4	Macroblock-level Parallelism	15
4.2.5	Macroblock-level Parallelism in the Spatial Domain (2D-Wave)	15
4.2.6	Macroblock-level Parallelism in the Temporal Domain	17
4.2.7	Combining Macroblock-level Parallelism in the Spatial and Temporal Domains (3D-Wave)	19
4.2.8	Block-level Parallelism	21
5	Parallel Scalability of the Static 3D-Wave	22
6	Scalable MB-level Parallelism: The Dynamic 3D-Wave	27
6.1	Static vs Dynamic Scheduling	27
6.2	Implementation Issues of the Dynamic Scheduling	27
6.3	Support in the Programming Model	28
6.4	Managing Entropy Decoding	29
6.5	Open Issues	29
7	Parallel Scalability of the Dynamic 3D-Wave	30
8	Case Study: Mobile Video	38
9	Related Work	40
10	Conclusions	43

1 Introduction

We are witnessing a paradigm shift in computer architecture towards chip multi-processors (CMPs). In the past performance has improved mainly due to higher clock frequencies and architectural approaches to exploit instruction-level parallelism (ILP), such as pipelining, multiple issue, out-of-order execution, and branch prediction. It seems, however, that these sources of performance gains are exhausted. New techniques to exploit more ILP are showing diminishing results while being costly in terms of area, power, and design time. Also clock frequency increase is flattening out, mainly because pipelining has reached its limit. As a result, industry, including IBM, Sun, Intel, and AMD, turned to CMPs.

At the same time we see that we have hit the power wall and thus performance increase can only be achieved without increasing overall power. In other words, power efficiency (performance per watt or performance per transistor) has become an important metric for evaluation. CMPs allow power efficient computation. Even if real-time performance can be achieved with few cores, using manycores improves power efficiency as, for example, voltage/frequency scaling can be applied. That is also the reason why low power architectures are the first to apply the CMP paradigm on large scale [1,2].

It is expected that the number of cores on a CMP will double every three year [3], resulting in an estimate of 150 high performance cores on a die in 2017. For power efficiency reasons CMPs might as well consist of many simple and small cores which might count up to thousand and more [4]. A central question is whether applications scale to such large number of cores. If applications are not extensively parallelizable, cores will be left unused and performance might suffer. Also, achieving power efficiency on manycores relies on TLP offered by applications.

In this paper we investigate this issue for video decoding workloads by analyzing their parallel scalability. Multimedia applications remain important workloads in the future and video codecs are expected to be important benchmarks for all kind of systems, ranging from low power mobile devices to high performance systems. Specifically, we analyze the parallel scalability of an H.264 decoder by exploring the limits to the amount of TLP. This research is similar to the video applications limit study of ILP presented in [5] in 1996. Investigating the limits to TLP, however requires to set a bound to the granularity of threads. Without this bound a single instruction can be a thread, which is clearly not desirable. Thus, using emerging techniques such as light weight micro-threading [6], possibly even more TLP can be exploited. However, we show that considering down to macroblock level granularity, H.264 contains sufficient parallelism to sustain a manycore CMP.

This paper is organized as follows. In Section 2 a brief overview of the H.264 standard is provided. Next, in Section 3 we describe the benchmark we use throughout this paper. In Section 4 possible parallelization strategies are discussed. In Section 5 we analyze the parallel scalability of the best previously proposed parallelization technique. In Section 6 we propose the Dynamic 3D-Wave parallelization strategy. Using this new parallelization strategy in Section 7 we investigate the parallel scalability of H.264 and show it exhibits huge amounts of parallelism. Section 8 describes a case study to assess the practical value of a highly parallelized decoder. In Section 9 we provide an overview of related work. Section 10 concludes the paper and discusses future work.

Currently, the best video coding standard, in terms of compression and quality is H.264 [7]. It is used in HD-DVD and blu-ray Disc, and many countries are using/will use it for terrestrial television broadcast, satellite broadcast, and mobile television services. It has a compression improvement of over two times compared to previous standards such as MPEG-4 ASP, H.262/MPEG-2, etc. The H.264 standard [8] was designed to serve a broad range of application domains ranging from low to high bitrates, from low to high resolutions, and a variety of networks and systems, e.g., Internet streams, mobile streams, disc storage, and broadcast. The H.264 standard was jointly developed by ITU-T Video Coding Experts Group (VCEG) and ISO/IEC Moving Picture Experts Group (MPEG). It is also called MPEG-4 part 10 or AVC (advanced video coding).

```

graph LR
    Input[Compressed video] --> Entropy[Entropy decoding]
    Entropy --> Inverse[Inverse quantization]
    Inverse --> IDCT[IDCT]
    IDCT --> Sum((+))
    Sum --> Deblocking[Deblocking filter]
    Deblocking --> Output[Uncompressed video]
    Output --> FrameBuffer[Frame buffer]
    FrameBuffer --> MCPred[MC prediction]
    IntraPred[Intra prediction] --> Sum
    MCPred --> Sum
    Deblocking -.-> Entropy
  
```

The diagram illustrates the architecture of a video codec, divided into an encoder (top) and a decoder (bottom).

Encoder Path:

- Uncompressed video** enters the system.
- It is split into two paths:
 - One path goes directly to a summation node (+).
 - The other path goes through **ME prediction** (Motion Estimation prediction).
- The output of **ME prediction** is fed into **MC prediction** (Motion Compensation prediction) and **Intra prediction**.
- MC prediction** and **Intra prediction** outputs are combined at a summation node (+) to produce the **residual**.
- The **residual** is added to the original video at the first summation node (+) to produce the **predicted video**.
- The **predicted video** is then processed by **DCT** (Discrete Cosine Transform) and **Quantization**.
- The output of **Quantization** is sent to **Entropy coding**.
- The final output is **Compressed video**.

Decoder Path:

- Compressed video** enters the decoder.
- It is processed by **Entropy coding** and then **Inverse quantization**.
- The output of **Inverse quantization** is processed by **IDCT** (Inverse Discrete Cosine Transform).
- The output of **IDCT** is added to the **residual** at a summation node (+) to produce the **reconstructed video**.
- The **reconstructed video** is then processed by a **Deblocking filter**.
- The output of the **Deblocking filter** is fed back into the **Frame buffer**.
- The **Frame buffer** provides input to **ME prediction** and **MC prediction** in the next frame.

Fig. 2: Block diagram of the encoding process.

A movie picture is called a frame and can consist of several slices or slice groups. A slice is a partition of a frame such that all comprised MBs are in scan order (from left to right, from top to bottom). The Flexible Macroblock Ordering (FMO) feature allows slice groups to be defined, which consist of an arbitrary set of MBs. Each slice group can be partitioned into slices in the same way a frame can. Slices are self contained and can be decoded independently.

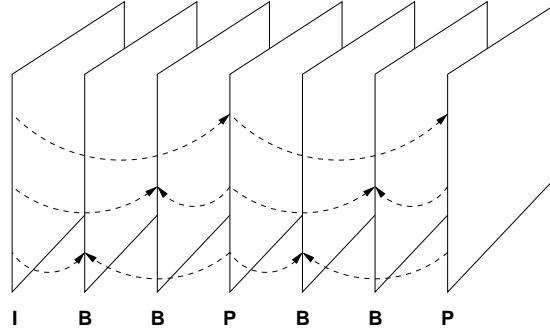


Fig. 3: A typical slice/frame sequence and its dependencies.

H.264 defines three main types of slices/frames: I, P, and B-slices. An I-slice uses intra prediction and is independent of other slices. In intra prediction a MB is predicted based on adjacent blocks. A P-slice uses motion estimation and intra prediction and depends on one or more previous slices, either I, P or B. Motion estimation is used to exploit temporal correlation between slices. Finally, B-slices use bidirectional motion estimation and depend on slices from past and future [9]. Figure 3 shows a typical slice order and the dependencies, assuming each frame consist of one slice only. The standard also defines SI and SP slices that are slightly different from the ones mentioned before and which are targeted at mobile and Internet streaming applications.

The H.264 standard has many options. We briefly mention the key features and compare them to previous standards. Table 1 provides an overview of the discussed features for MPEG-2, MPEG-4 ASP, and H.264. The different columns for H.264 represent profiles which are explained later. For more details the interested reader is referred to [10, 11].

Motion estimation

Advances in motion estimation is one of the major contributors to the compression improvement of H.264. The standard allows variable block sizes ranging from 16×16 down to 4×4 , and each block has its own motion vector(s). The motion vector is quarter sample accurate. Multiple reference frames can be used in a weighted fashion. This significantly improves coding occlusion areas where an accurate prediction can only be made from a frame further in the past.

Intra prediction

Three types of intra coding are supported, which are denoted as Intra_4x4, Intra_8x8 and Intra_16x16. The first type uses spatial prediction on each 4×4

Tab. 1: Comparison of video coding standards and profiles.

	MPEG-2	MPEG-4 ASP	H.264 BP	H.264 MP	H.264 XP	H.264 HiP
Picture types	I, P, B	I, P, B	I, P	I, P, B	I, P, B, SL, SP	I, P, B
Flexible macroblock ordering	No	No	Yes	No	Yes	No
Motion block size	16×16	16×16, 16×8, 8×8	16×16, 16×8, 8×16, 8×8, 8×4, 4×8, 4×4	16×16, 16×8, 8×16, 8×8, 8×4, 4×8, 4×4	16×16, 16×8, 8×16, 8×8, 8×4, 4×8, 4×4	16×16, 16×8, 8×16, 8×8, 8×4, 4×8, 4×4
Multiple reference frames	No	No	Yes	Yes	Yes	Yes
Motion pel accu- racy	1, 1/2	1, 1/2, 1/4	1, 1/2, 1/4	1, 1/2, 1/4	1, 1/2, 1/4	1, 1/2, 1/4
Weighted predic- tion	No	No	No	Yes	Yes	Yes
Transform	16×16	8×8 DCT	4×4 integer DCT	4×4 integer DCT	4×4 integer DCT	4×4, 8×8 inte- ger DCT
In-loop deblocking filter	No	No	Yes	Yes	Yes	Yes
Entropy coding	VLC	VLC	CAVLC, UVLC	CAVLC, UVLC, CABAC	CAVLC, UVLC	CAVLC, UVLC, CABAC

luminance block. Eight modes of directional prediction are available, among them horizontal, vertical, and diagonal. This mode is well suited for MBs with small details. For smooth image areas the Intra_16x16 type is more suitable, for which four prediction modes are available. The high profile of H.264 also supports intra coding on 8×8 luma blocks. Chroma components are estimated for whole MBs using one specialized prediction mode.

Discrete cosine transform

MPEG-2 and MPEG-4 part 2 employed an 8×8 floating point transform. However, due to the decreased granularity of the motion estimation, there is less spatial correlation in the residual signal. Thus, standard a 4×4 (that means 2×2 for chrominance) transform is used, which is as efficient as a larger transform [12]. Moreover, a smaller block size reduces artifacts known as ringing and blocking. An optional feature of H.264 is Adaptive Block size Transform (ABT), which adapts the block size used for DCT to the size used in the motion estimation [13]. Furthermore, to prevent rounding errors that occur in floating point implementations, an integer transform was chosen.

Deblocking filter

Processing a frame in MBs can produce blocking artifacts, generally considered the most visible artifact in prior standards. This effect can be resolved by applying a deblocking filter around the edges of a block. The strength of the filter is adaptable through several syntax elements [14]. While in H.263+ this feature was optional, in H.264 it is standard and it is placed within the motion compensated prediction loop (see Figure 1) to improve the motion estimation.

Entropy coding

There are two classes of entropy coding available in H.264: Variable Length Coding (VLC) and Context Adaptive Binary Arithmetic Coding (CABAC). The latter achieves up to 10% better compression but at the cost of large computational complexity [15]. The VLC class consists of Context Adaptive VLC (CAVLC) for the transform coefficients, and Universal VLC (UVLC) for the small remaining part. CAVLC achieves large improvements over simple VLC, used in prior standards, without the full computational cost of CABAC.

The standard was designed to suite a broad range of video application domains. However, each domain is expected to use only a subset of the available options. For this reason profiles and levels were specified to mark conformance points. Encoders and decoders that conform to the same profile are guaranteed to interoperate correctly. Profiles define sets of coding tools and algorithms that can be used while levels place constraints on the parameters of the bitstream.

The standard initially defined two profiles, but has since then been extended to a total of 11 profiles, including three main profiles, four high profiles, and four all-intra profiles. The three main profiles and the most important high profile are:

- **Baseline Profile (BP):** the simplest profile mainly used for video conferencing and mobile video.

- **Main Profile (MP):** intended to be used for consumer broadcast and storage applications, but overtaken by the high profile.
- **Extended Profile (XP):** intended for streaming video and includes special capabilities to improve robustness.
- **High Profile (HiP)** intended for high definition broadcast and disc storage, and is used in HD DVD and Blu-ray.

Besides HiP there are three other high profiles that support up to 14 bits per sample, 4:2:2 and 4:4:4 sampling, and other features [16]. The all-intra profiles are similar to the high profiles and are mainly used in professional camera and editing systems.

In addition 16 levels are currently defined which are used for all profiles. A level specifies, for example, the upper limit for the picture size, the decoder processing rate, the size of the multi-picture buffers, and the video bitrate. Levels have profile independent parameters as well as profile specific ones.

3 Benchmark

Throughout this paper we use the HD-VideoBench [17], which provides movie test sequences, an encoder (X264 [18]), and a decoder (FFmpeg [19]). The benchmark contains the following test sequences:

- **rush_hour:** rush-hour in Munich city; static background, slowly moving objects.
- **riverbed:** riverbed seen through waving water; abrupt and stochastic changes.
- **pedestrian:** shot of a pedestrian area in city center; static background, fast moving objects.
- **blue_sky:** top of two trees against blue sky; static objects, sliding camera.

All movies are available in three formats: 720×576 (SD), 1280×720 (HD), 1920×1088 (FHD). Each movie has a frame rate of 25 frames per second and has a length of 100 frames. For some experiments longer sequences were required, which we created by replicating the sequences. Unless specified otherwise, for SD and HD we used sequences of 300 frames while for FHD we used sequences of 400 frames.

The benchmark provides the test sequences in raw format. Encoding is done with the X264 encoder using the following options: 2 B-frames between I and P frames, 16 reference frames, weighted prediction, hexagonal motion estimation algorithm (hex) with maximum search range 24, one slice per frame, and adaptive block size transform. Movies encoded with this set of options represent the typical case. To test the worst case scenario we also created movies with large motion vectors by encoding the movies with the exhaustive search motion estimation algorithm (esa) and a maximum search range of 512 pixels. The first movies are marked with the suffix 'hex' while for the latter we use 'esa'.

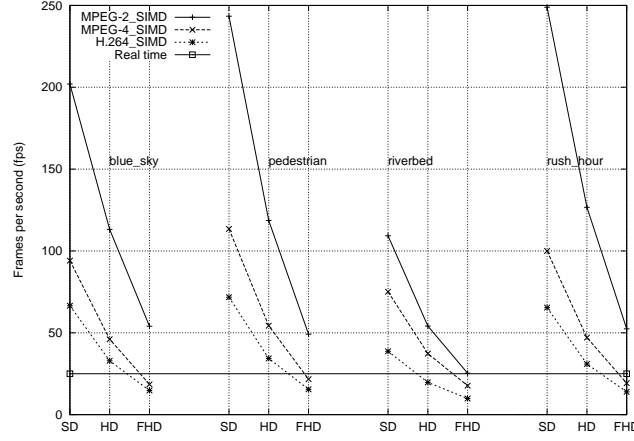


Fig. 4: H.264 decoding performance.

4 Parallelizing H.264

The coding efficiency gains of advanced video codecs like H.264, come at the price of increased computational requirements. The demands for computing power increases also with the shift towards high definition resolutions. As a result, current high performance uniprocessor architectures are not capable of providing the required performance for real-time processing [20–23]. Figure 4 depicts the performance of the FFmpeg H.264 decoder optimized with SIMD instructions on a superscalar processor (Intel IA32 Xeon processor at 2.4 GHz with 512KB of L2 cache) and compare it with other video codecs. For FHD resolution this high performance processor is not able to achieve real-time operation. And, the trend in future video applications is toward higher resolutions and higher quality video systems that in turn require more performance. As, with the shift towards the CMP paradigm, most hardware performance increase will come from the capability of running many threads in parallel, it is necessary to parallelize H.264 decoding.

Moreover, the power wall forces us to compute power efficient. This requires the utilization of all the cores CMPs offer, and thus applications have to be parallelizable to a great extent. In this section we analyze the possibilities for parallelization of H.264 and compare them in terms of communication, synchronization, load balancing, scalability and software optimization.

The H.264 codec can be parallelized either by a task-level or data-level decomposition. In Figure 5 the two approaches are sketched. In task-level decomposition individual tasks of the H.264 Codec are assigned to processors while in data-level decomposition different portions of data are assigned to processors running the same program.

4.1 Task-level Decomposition

In a task-level decomposition the functional partitions of the algorithm are assigned to different processors. As shown in Figure 1 the process of decoding H.264 consists of performing a series of operations on the coded input bitstream.

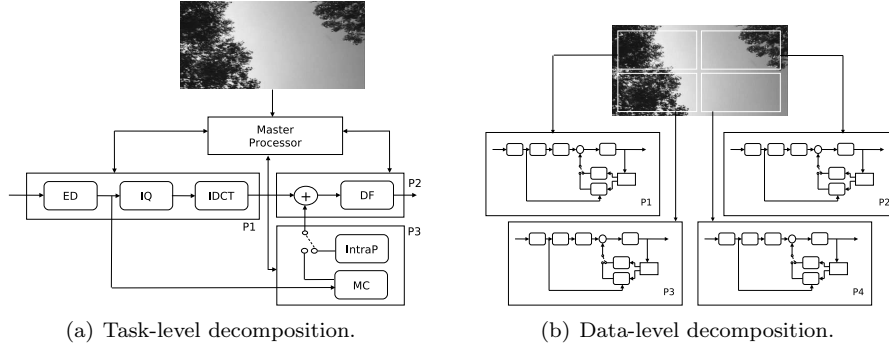


Fig. 5: H.264 parallelization techniques.

Some of these tasks can be done in parallel. For example, Inverse Quantization (IQ) and the Inverse Transform (IDCT) can be done in parallel with the Motion Compensation (MC) stage. In Figure 5(a) the tasks are mapped to a 4-processor system. A control processor is in charge of synchronization and parsing the bitstream. One processor is in charge of Entropy Decoding (ED), IQ and IDCT, another one of the prediction stage (MC or IntraP), and a third one is responsible for the deblocking filter.

Task-level decomposition requires significant communication between the different tasks in order to move the data from one processing stage to the other, and this may become the bottleneck. This overhead can be reduced using double buffering and blocking to maintain the piece of data that is currently being processed in cache or local memory. Additionally, synchronization is required for activating the different modules at the right time. This should be performed by a control processor and adds significant overhead.

The main drawbacks, however, of task-level decomposition are load balancing and scalability. Balancing the load is difficult because the time to execute each task is not known a priori and depends on the data being processed. In a task-level pipeline the execution time for each stage is not constant and some stage can block the processing of the others. Scalability is also difficult to achieve. If the application requires higher performance, for example by going from standard to high definition resolution, it is necessary to re-implement the task partitioning which is a complex task and at some point it could not provide the required performance for high throughput demands. Finally from the software optimization perspective the task-level decomposition requires that each task/processor implements a specific software optimization strategy, i.e., the code for each processor is different and requires different optimizations.

4.2 Data-level Decomposition

In a data-level decomposition the work (data) is divided into smaller parts and each assigned to a different processor, as depicted in Figure 5(b). Each processor runs the same program but on different (multiple) data elements (SPMD). In H.264 data decomposition can be applied at different levels of the data structure (see Figure 6). At the top of the data structure there is the complete video sequence. This video sequence is composed out of Group of Pictures

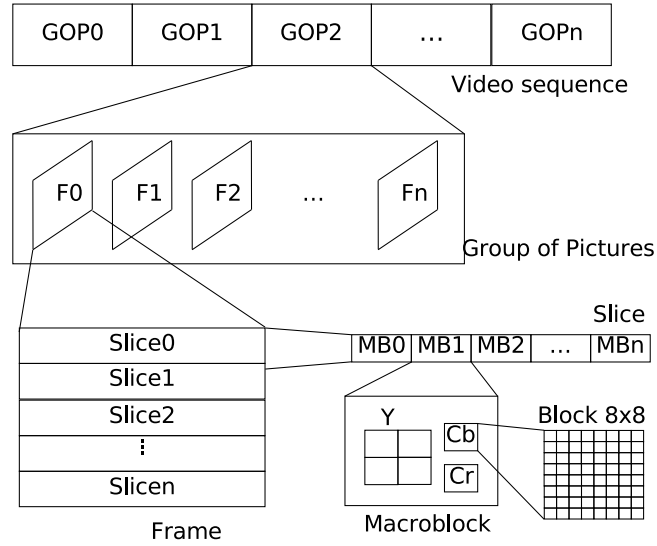


Fig. 6: H.264 data structure.

(GOPs) which are independent sections of the video sequence. GOPs are used for synchronization purposes because there are no temporal dependencies between them. Each GOP is composed of a set of frames, which can have temporal dependencies when motion prediction is used. Each frame can be composed of one or more slices. The slice is the basic unit for encoding and decoding. Each slice is a set of MBs and there are no temporal or spatial dependencies between slices. Further, there are MBs, which are the basic units of prediction. MBs are composed of luma and chroma blocks of variable size. Finally each block is composed of picture samples. Data-level parallelism can be exploited at each level of the data structure, each one having different constraints and requiring different parallelization methodologies.

4.2.1 GOP-level Parallelism

The coarsest grained parallelism is at the GOP level. H.264 can be parallelized at the GOP-level by defining a GOP size of N frames and assigning each GOP to a processor. GOP-level parallelism requires a lot of memory for storing all the frames, and therefore this technique maps well to multicomputers in which each processing node has a lot of computational and memory resources. However, parallelization at the GOP-level results in a very high latency that cannot be tolerated in some applications. This scheme is therefore not well suited for multicore architectures, in which the memory is shared by all the processors, because of cache pollution.

4.2.2 Frame-level Parallelism for Independent Frames

After GOP-level there is frame-level parallelism. As shown in Figure 3 in a sequence of I-B-B-P frames inside a GOP, some frames are used as reference for other frames (like I and P frames) but some frames (the B frames in this case) might not. Thus in this case the B frames can be processed in parallel. To do

so, a control processor can assign independent frames to different processors. Frame-level parallelism has scalability problems due to the fact that usually there are no more than two or three B frames between P frames. This limits the amount of TLP to a few threads. However, the main disadvantage of frame-level parallelism is that, unlike previous video standards, in H.264 B frames can be used as reference [24]. In such a case, if the decoder wants to exploit frame-level parallelism, the encoder cannot use B frames as reference. This might increase the bitrate, but more importantly, encoding and decoding are usually completely separated and there is no way for a decoder to enforce its preferences to the encoder.

4.2.3 Slice-level Parallelism

In H.264 and in most current hybrid video coding standards each picture is partitioned into one or more slices. Slices have been included in order to add robustness to the encoded bitstream in the presence of network transmission errors and losses. In order to accomplish this, slices in a frame should be completely independent from each other. That means that no content of a slice is used to predict elements of other slices in the same frame, and that the search area of a dependent frame can not cross the slice boundary [10, 16]. Although support for slices have been designed for error resilience, it can be used for exploiting TLP because slices in a frame can be encoded or decoded in parallel. The main advantage of slices is that they can be processed in parallel without dependency or ordering constraints. This allows exploitation of slice-level parallelism without making significant changes to the code.

However, there are some disadvantages associated with exploiting TLP at the slice level. The first one is that the number of slices per frame is determined by the encoder. That poses a scalability problem for parallelization at the decoder level. If there is no control of what the encoder does then it is possible to receive sequences with few (or one) slices per frame and in such cases there would be reduced parallelization opportunities. On the other hand, although slices are completely independent of each other, H.264 includes a deblocking filter that can be applied across slice boundaries. This is also an option that is selectable by the encoder, but means that even with an input sequence with multiple slices, if the deblocking filter crosses slice boundaries the filtering process should be performed after the frame processing in a sequential order inside the frame. This reduces the speed-up that can be achieved from slice-level parallelization. Another problem is load balancing. Usually slices are created with the same number of MBs, and thus can result in an imbalance at the decoder because some slices are decoded faster than others depending on the content of the slice.

Finally, the main disadvantage of slices is that an increase in the number of slices per frame increases the bitrate for the same quality level (or, equivalently, it reduces quality for the same bitrate level). Figure 7 shows the increase in bitrate due to the increase of the number of slices for four different input videos at three different resolutions. The quality is maintained constant (40 PSNR). When the number of slices increases from one to eight, the increase in bitrate is less than 10%. When going to 32 slices the increase ranges from 3% to 24%, and when going to 64 slices the increase ranges from 4% up to 34%. For some applications this increase in the bitrate is unacceptable and thus a large number of slices is not possible. As shown in the figure the increase in bitrate depends

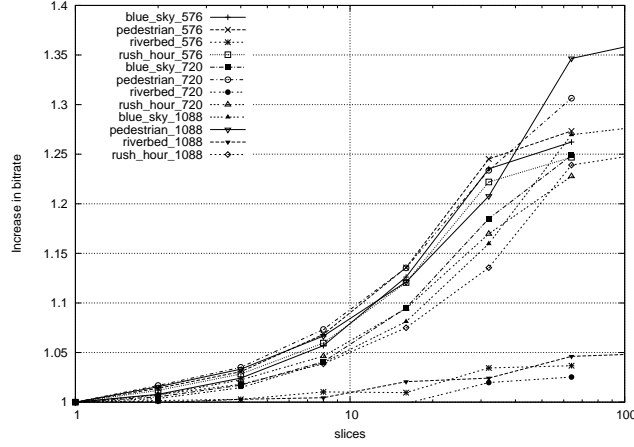


Fig. 7: Bitrate increase due to slices.

heavily on the input content. The riverbed sequence is encoded with very little motion estimation, and thus has a large absolute bitrate compared to the other three sequences. Thus, the relative increase in bitrate is much lower than the others.

4.2.4 Macroblock-level Parallelism

There are two ways of exploiting MB-level parallelism: in the spatial domain and/or in the temporal domain. In the spatial domain MB-level parallelism can be exploited if all the intra-frame dependencies are satisfied. In the temporal domain MB-level parallelism can be exploited if, in addition to the intra-dependencies, inter-frame dependencies are satisfied.

4.2.5 Macroblock-level Parallelism in the Spatial Domain (2D-Wave)

Usually MBs in a slice are processed in scan order, which means starting from the top left corner of the frame and moving to the right, row after row. To exploit parallelism between MBs inside a frame it is necessary to take into account the dependencies between them. In H.264, motion vector prediction, intra prediction, and the deblocking filter use data from neighboring MBs defining a structured set of dependencies. These dependencies are shown in Figure 8. MBs can be processed out of scan order provided these dependencies are satisfied. Processing MBs in a diagonal wavefront manner satisfies all the dependencies and at the same time allows to exploit parallelism between MBs. We refer to this parallelization technique as 2D-Wave.

Figure 9 depicts an example for a 5×5 MBs image (80×80 pixels). At time slot T7 three independent MBs can be processed: MB (4,1), MB (2,2) and MB (0,3). The figure also shows the dependencies that need to be satisfied in order to process each of these MBs. The number of independent MBs in each frame depends on the resolution. In Table 2 the number of independent MBs for different resolutions is stated. For a low resolution like QCIF there are only 6 independent MBs during 4 time slots. For High Definition (1920×1088) there

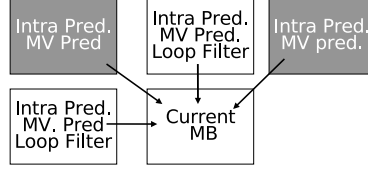


Fig. 8: Dependencies between neighboring MBs in H.264.

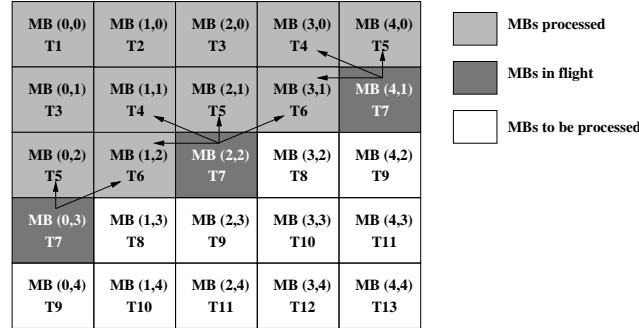


Fig. 9: 2D-Wave approach for exploiting MB parallelism in the spatial domain. The arrows indicate dependencies.

are 60 independent MBs during 9 slots of time. Figure 10 depicts the available MB parallelism over time for an FHD resolution frame, assuming that the time to decode a MB is constant.

MB-level parallelism in the spatial domain has many advantages over other schemes for parallelization of H.264. First, this scheme can have a good scalability. As shown before the number of independent MBs increases with the resolution of the image. Second, it is possible to achieve a good load balancing if a dynamic scheduling system is used. That is due to the fact that the time to decode a MB is not constant and depends on the data being processed. Load balancing could take place if a dynamic scheduler assigns a MB to a processor once all its dependencies have been satisfied. Additionally, because in MB-level parallelization all the processors/threads run the same program the same set of software optimizations (for exploiting ILP and SIMD) can be applied to all processing elements.

However, this kind of MB-level parallelism has some disadvantages. The

Tab. 2: Maximum parallel MBs for several resolutions using the 2D-Wave approach. Also the number of times slots this maximum is available is stated.

Resolution		MBs	Slots
QCIF	176×144	6	4
CIF	352×288	11	8
SD	720×576	23	14
HD	1280×720	40	6
FHD	1920×1088	60	9

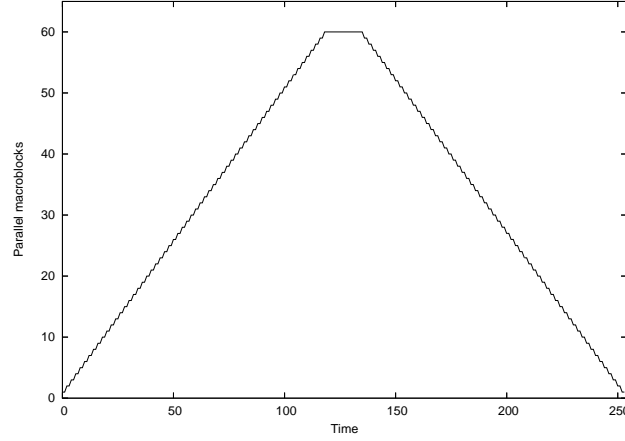


Fig. 10: MB parallelism for a single FHD frame using the 2D-Wave approach.

first one is that the entropy decoding can not be parallelized using data decomposition due to the fact that H.264 entropy coding adapts its probabilities using data from consecutive MBs. As a result the lowest level of data that can be parsed from the bitstream are slices. Individual MBs can not be identified without performing entropy decoding. That means that in order to decode independent MBs, they should be entropy decoded first and in sequential order. Only after entropy decoding has been performed the parallel processing of MBs can start. The effect of this disadvantage can be mitigated by using special purpose instructions or hardware accelerators for the entropy decoding process.

The second disadvantage is that the number of independent MBs does not remain constant during the decoding of a frame (see Figure 10). At the start of the frame there are few independent MBs, after some time a steady state is reached and a maximum number of independent MBs can be processed, but at the end of the frame the number of MBs starts to decline. Therefore, it is not possible to sustain a certain processing rate during the decoding of a frame. Using the parallelization strategy we propose in Section 6, this problem is overcome.

4.2.6 Macroblock-level Parallelism in the Temporal Domain

In the decoding process the dependency between frames is in the Motion Compensation (MC) module only. MC can be regarded as copying an area, called the reference area, from the reference frame, and then to add this predicted area to the residual MB to reconstruct the MB in the current frame. The reference area is pointed to by a Motion Vector (MV). The MV length is limited by two factors: the H.264 level and the motion estimation algorithm. The H.264 level defines, amongst others, the maximum length of the vertical component of the MV [8]. However, not all encoders follow the guidelines of the levels, but use the absolute limit of the standard which is 512 pixels vertical and 2048 pixels horizontal. In practice the motion estimation algorithm defines a maximum search range of dozens of pixels, because a larger search range is very computationally demanding and provides only a small benefit [25].

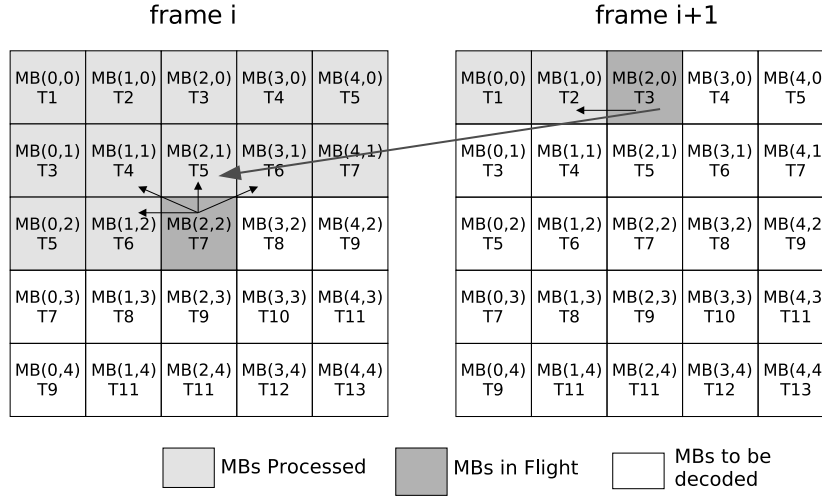


Fig. 11: MB-level Parallelism in the Temporal Domain in H.264.

When the reference area has been decoded it can be used by the referencing frame. Thus it is not necessary to wait until a frame is completely decoded before decoding the next frame. The decoding process of the next frame can start after the reference areas of the reference frames are decoded. Figure 11 shows an example of two frames where the second depends on the first. MBs are decoded in scan order and one at a time. The figure shows that MB (2, 0) of frame $i + 1$ depends on MB (2, 1) of frame i which has been decoded. Thus this MB can be decoded even though frame i is not completely decoded.

The main disadvantage of this scheme is the limited scalability. The number of MBs that can be decoded in parallel is inversely proportional to the length of the vertical motion vector component. Thus for this scheme to be beneficial the encoder should be enforced to heavily restrict the motion search area which in far most cases is not possible. Assuming it would be possible, the minimum search area is around 3 MB rows: 16 pixels for the co-located MB, 3 pixels at the top and at the bottom of the MB for sub-sample interpolations and some pixels for motion vectors (at least 10). As a result the maximum parallelism is 14, 17 and 27 MBs for STD, HD and FHD frame resolutions respectively.

The second limitation of this type of MB-level parallelism is poor load-balancing because the decoding time for each frame is different. It can happen that a fast frame is predicted from a slow frame and can not decode faster than the slow frame and remains idle for some time. Finally, this approach works well for the encoder who has the freedom to restrict the range of the motion search area. In the case of the decoder the motion vectors can have large values (even if the user ask the encoder to restrict them as we will show later) and the number of frames that can be processed in parallel is reduced.

This approach has been implemented in the X264 open source encoder (See Section 9). In Figure 12 it is shown the performance of X264 encoder in a ccNUMA multiprocessor (SGI Altix with Itanium 2 processors) for the blue_sky FHD input sequence. With 4 threads the speed-up is 3.23X (efficiency of 80%), with 16 threads speed-up is 10.54X (efficiency of 66%) and with 32 threads

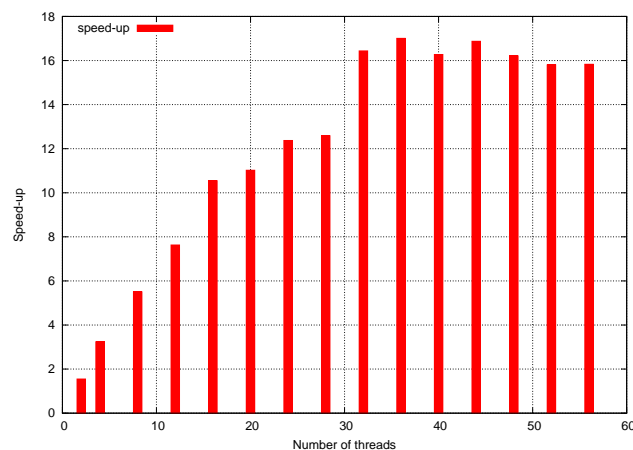


Fig. 12: Parallel scalability of X264 encoder on an Itanium2 ccNUMA multiprocessor

the speed-up is 16.4X (efficiency of 51%) a point from which the performance improvement is saturated.

4.2.7 Combining Macroblock-level Parallelism in the Spatial and Temporal Domains (3D-Wave)

None of the single approaches described in the previous sections scales to future manycore architectures containing 100 cores or more. There is a considerable amount of MB-level parallelism, but in the spatial domain there are phases with a few independent MBs, and in the temporal domain scalability is limited by the height of the frame.

In order to overcome these limitations it is possible to exploit both temporal and spatial MB-level parallelism. Inside a frame, spatial MB-level parallelism can be exploited using the 2D wave scheme mentioned previously. And between frames temporal MB-level parallelism can be exploited simultaneously. Adding the inter-frame parallelism (time) to the 2D-Wave intra-frame parallelism (space) results in a combined 3D-Wave parallelization. Figure 13 illustrates this way of parallel decoding of MBs.

3D-Wave MB decoding requires a scheduler for assigning MBs to processors. This scheduling can be performed in a static or a dynamic way. Static 3D-Wave exploits temporal parallelism by assuming that the motion vectors have a restricted length, based on that it uses a fixed spatial offset between decoding consecutive frames. This static approach has been implemented in previous work (See Section 9) for the H.264 encoder. Note that it is more suitable for the encoder because it has the flexibility of selecting the size of the motion search area.

But, considering the fact that the decoder cannot influence the MV choice of the encoder, this static approach is not capable of exploiting all the MB-level parallelism. For the static approach in a decoder to work correctly, the worst case MV length has to be considered, which is 512 vertical pixels for full HD. As we show later, MVs are typically very short and thus a lot of parallelism would

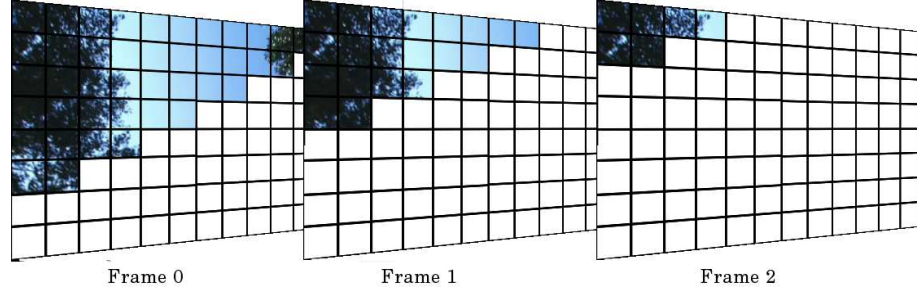


Fig. 13: 3D-Wave strategy: frames can be decoded in parallel because inter frame dependencies have limited spatial range.

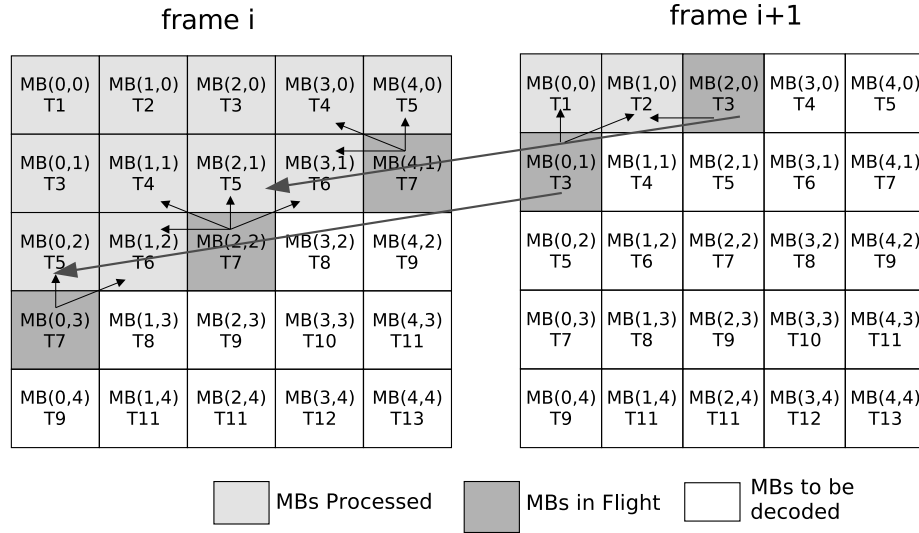


Fig. 14: 3D-Wave strategy: intra- and inter-frame dependencies between MBs

be unexploited.

Our proposal is the Dynamic 3D-Wave approach, which to the best of our knowledge has not been reported before. It is not based on restricted values of the motion vectors, instead it uses a dynamic scheduling system in which MBs are scheduled for decoding when all the dependencies (intra-frame and inter-frame) have been satisfied. As shown in Figure 14 as soon as the search area in the reference frame has been decoded a MB in a different frame can be processed. The dynamic system is best suited for the decoder taken into account that the length of motion vectors is not known in advance and that the decoding time of each MB is input dependent. The Dynamic 3D-Wave system results in a better thread scalability and a better load-balancing. A more detailed analysis of the Dynamic 3D-Wave is presented in Section 6.

4.2.8 Block-level Parallelism

Finally, the finest grained data-level parallelism is at the block-level. Most of the computations of the H.264 kernels are performed at the block level. This applies, for example, to the interpolations that are done in the motion compensation stage, to the IDCT, and to the deblocking filter. This level of data parallelism maps well to SIMD instructions [20, 26–28]. SIMD parallelism is orthogonal to the other levels of parallelism described above and because of that it can be mixed, for example, with MB-level parallelization to increase the performance of each thread.

5 Parallel Scalability of the Static 3D-Wave

In the previous section we suggested that using the 3D-Wave strategy for decoding H.264 would reveal a large amount of parallelism. We also mentioned that two strategies are possible: a static and a dynamic approach. Zhao [29] used a method similar to the Static 3D-Wave for encoding H.264 which so far has been the most scalable approach to H.264 coding. However, a limit study to the scalability of his approach is lacking. In order to compare our dynamic approach to his, in this section we analyze the parallel scalability of the Static 3D-Wave.

The Static 3D-Wave strategy assumes a static maximum MV length and thus a static reference range. Figure 15 illustrates the reference range concept, assuming a MV range of 32 pixels. The hashed MB in Frame 1 is the MB currently considered. As the MV can point to any area in the range of $[-32, +32]$ pixels, its reference range is the hashed area in Frame 0. In the same way, every MB in the wave front of Frame 1 has a reference range similar to the presented one, with its respective displacement. Thus, if a minimum offset, corresponding to the MV range, between the two wavefronts is maintained, the frames can be decoded in parallel.

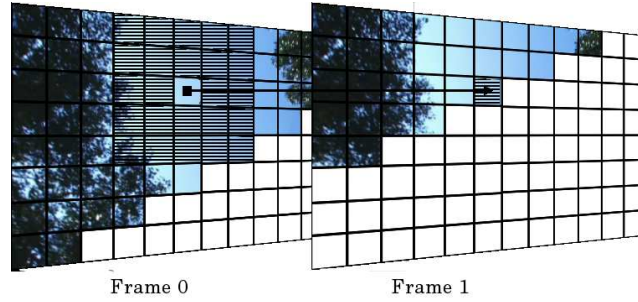


Fig. 15: Reference range example: Hashed area on frame 0 is the reference range of the hashed MB of frame 1.

As in Section 4.2, to calculate the number of parallel MBs, it is assumed that processing a MB requires one time slot. In reality, different MBs require different processing times, but this can be solved using a dynamic scheduling system. Furthermore, the following conservative assumptions are made to calculate the amount of MB parallelism. First, B frames are used as reference frames. Second, the reference frame is always the previous one. Third, only the first frame of the sequence is an I frame. These assumptions represent the worst case scenario for the Static 3D-Wave.

The number of parallel MBs is calculated as follows. First, we calculate the MB parallelism function for one frame using the 2D-Wave approach, as in Figure 10. Next, given a MV range, we determine the required offset between the decoding of two frames. Finally, we added the MB parallelism function of all frames using the offset. Formally, let $h_0(t)$ be the MB parallelism function of the 2D-Wave inside frame 0. This graph of this function is depicted in Figure 10 for a full HD frame. Then the MB parallelism function of frame i is computed as $h_i(t) = h_{i-1}(t - offset)$ for $i > 1$. The MB parallelism function of the total

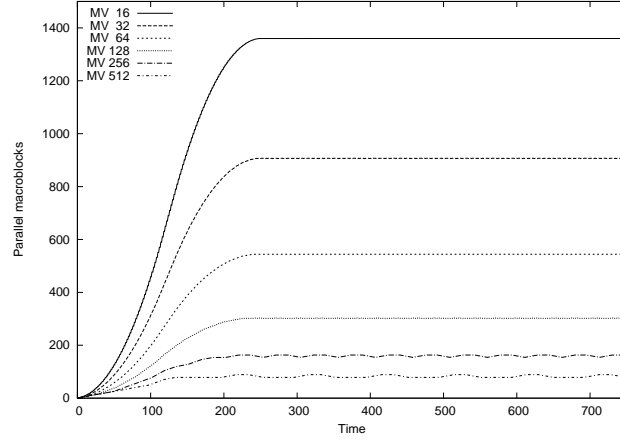


Fig. 16: Static 3D-Wave: number of parallel MBs for FHD resolution using several MV ranges.

Static 3D-Wave is given by $H(t) = \sum_i h_i(t)$.

The offset is calculated as follows. For motion vectors with a maximum length of 16 pixels, it is possible to start the decoding of the second frame when the MBs (1,2) and (2,1) of the first frame have been decoded. Of these, MB (1,2) is the last one decoded, namely at time slot T6 (see Figure 9). Thus, the next frame can be started decoding at time slot T7, resulting in an offset of 4 time slots. Similarly, for maximum MV ranges of 32, 64, 128, 256, and 512 pixels, we find an offset of 9, 15, 27, 51, and 99 time slots, respectively. In general, for a MV range of n pixels, the offset is $3 + 3 \times \lceil n/16 \rceil$ time slots.

First, we investigate the maximum amount of available MB parallelism using the Static 3D-Wave strategy. For each resolution we applied the Static 3D-Wave to MV ranges of 16, 32, 64, 128, 256, and 512 pixels. The range of 512 pixels is the maximum vertical MV length allowed in level 4.0 (HD and FHD) of the H.264 standard.

Figure 16 depicts the amount of MB-level parallelism as a function of time, i.e., the number of MBs that could be processed in parallel in each time slot for a FHD sequence using several MV ranges. The graph depicts the start of the video sequence only. At the end of the movie the MB parallelism drops similar as it increased at startup. For large MV ranges, the curves have a small fluctuation which is less for small MV ranges. Figure 17 shows the corresponding number of frames in flight for each time slot. The shape of the curves for SD and HD resolutions are similar and, therefore, are omitted. Instead, Table 3 presents the maximum MB parallelism and frames in flight for all resolutions.

The results show that the Static 3D-Wave offers significant parallelism if the MV length is restricted to 16 pixels. However, in most cases this restriction cannot be guaranteed and the maximum MV length has to be considered. In such a case the parallelism drops significant towards the level of the 2D-Wave. The difference is that the Static 3D-Wave has a sustained parallelism while the 2D approach has little parallelism at the beginning and at the end of processing a frame.

A potential problem of the 3D-Wave strategy is the large number of frames

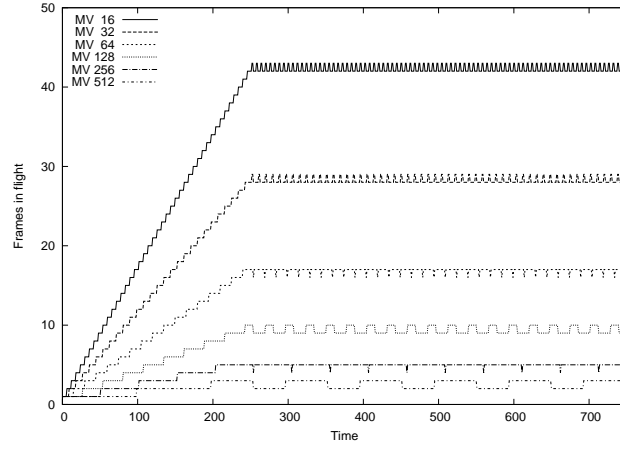


Fig. 17: Static 3D-Wave: number of frames in flight for FHD resolution using several MV ranges.

Tab. 3: Static 3D-Wave: maximum MB parallelism and frames in flight for several MV ranges and all three resolutions.

MV range	Max # parallel MBs			Max # frames in flight		
	SD	HD	FHD	SD	HD	FHD
512	-	36	89	-	2	3
256	33	74	163	3	4	5
128	60	134	303	5	7	10
64	108	240	544	8	12	17
32	180	400	907	13	19	29
16	276	600	1360	20	28	43

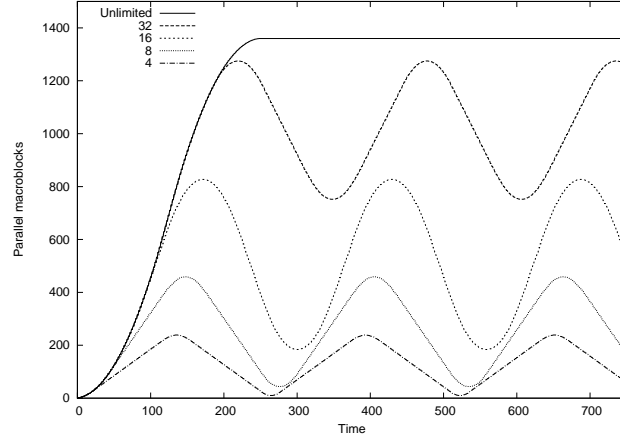


Fig. 18: Static 3D-Wave: number of parallel MBs for FHD resolution, for several limits of frames in flight.

that are decoded simultaneously. As Figure 17 shows when the MV range is 16 pixels, there can be up to 43 frames in flight simultaneously. This can be a problem in systems with limited memory since each active frame has to be stored in memory. To investigate this potential limitation, we analyze the tradeoff between the number of frames in flight and the amount of MB parallelism.

In this experiment the number of frames in flight is limited, i.e., the decoding of the next frame is started only if the number of frames in flight is less than or equal to the specified maximum. If not, it has to wait until a frame is completely decoded.

Figure 18 depicts the number of parallel MBs for an FHD sequence for several limits of frames in flight. The analysis was performed assuming a MV range of 16 pixels. This MV range was chosen because it clearly shows the effect of limiting the number of frames in flight, as it has the maximum amount of MB parallelism.

Limiting the number of frames in flight has two different effects. The first one is a proportional decrease of the MB parallelism. For unlimited frames in flight, after the ramp up there are 1360 MBs available. Limiting the number of active frames to 32, the average number of MBs available to be processed in parallel is reduced to 1013. For 16, 8, and 4 active frames, the results show an average of 505, 251, and 124 parallel MBs, respectively.

The second effect of limiting the number of frames in flight is that the number of available MBs fluctuates over time. While there is just a very small (or no) fluctuation when there is no limit on the number of active frames, the fluctuation is large when the number of active frames is limited. As shown in Table 4 for FHD resolution and a MV range of 16 pixels, when the limit is 32 frames, the minimum number of MBs is 752 while the maximum number of 1275. For a limit of 4 frames the minimum is 10 while the maximum is 239. This fluctuation is the result of the superposition of peaks in the 2D-Wave MB parallelism curves.

The fluctuation on the number of available parallel MBs can result in underutilization of computational resources. A scheduling technique can be used to reduce the fluctuation. To minimize the fluctuation, the decoding of the

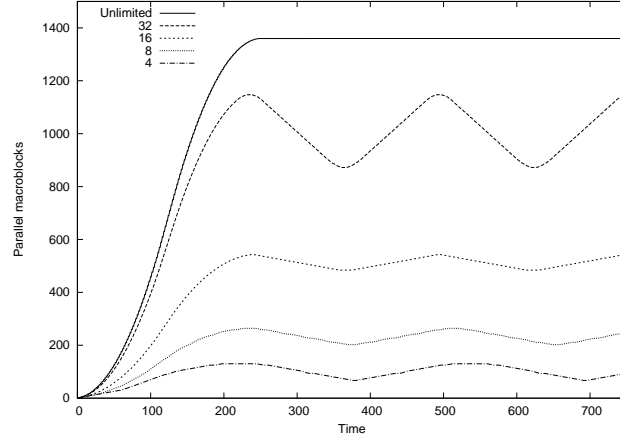


Fig. 19: Static 3D-Wave: available MB parallelism for several limits of frames in flight using scheduling.

Tab. 4: Static 3D-Wave: comparison of the available MB parallelism, with and without scheduling technique, for FHD resolution and a MV range of 16 pixels.

Frames in Flight	Regular			With Scheduling		
	Max	Min	Avg	Max	Min	Avg
unlimited	1360	1360	1360	1360	1360	1360
32	1275	752	1013	1148	872	1010
16	827	184	505	544	484	514
8	459	44	251	264	202	233
4	239	10	124	130	67	98

next frame should start T_{frame}/n_{active} time slots after the start of decoding the current frame, where T_{frame} is the number of time slots to process a frame and n_{active} is the maximum number of active frames. Figure 19 shows the MB parallelism time curve resulting from this scheduling strategy. This scheduling technique has a small effect on the average amount of MB parallelism, but reduces the fluctuation significantly. Table 4 compares the maximum, minimum and average MB parallelism with and without scheduling. As can be seen from the table, the amount of MB parallelism is proportional to the limit of the number of frames in flight.

6 Scalable MB-level Parallelism: The Dynamic 3D-Wave

In the previous section we showed that the best parallelization strategy proposed so far has limited value for decoding H.264 because of poor scalability. In Section 4.2.7 we proposed the dynamic 3D-Wave as a scalable solution and explained its fundamentals. In this section we further develop the idea and discuss some of the important implementation issues.

6.1 Static vs Dynamic Scheduling

The main idea of the 3D-wave algorithm is to combine spatial and temporal MB-level parallelism in order to increase the scalability and improve the processor efficiency. But in turn it poses a new problem: the assigning of MBs to processors. In static 3D-wave MBs are assigned to processors based on a fixed description of the data flow of the application. For spatial MB-level parallelism it assigns rows of MBs to processors and start new rows only when a fixed MB offset has been completed. For temporal MB-level parallelism it only starts a new frame when a known offset in the reference frame has been processed. If the processing of the dependent MB (in the same or in a different frame) is taking more time than the predicted by the static schedule the processor has to wait for the dependent MB to finish. This kind of static scheduling algorithms work well for applications with fixed execution times and fixed precedence constraints. In the case of H.264 decoding these conditions are not met. First, the decoding time of a MB depends on the input data of the macroblock and its properties (type, size, etc). Second, the temporal dependencies of a macroblock are input dependent and can not be predicted or restricted by the decoder application. As a result, a static scheduler can not discover all the available parallelism and can not sustain a high efficiency and scalability.

The dynamic 3D-wave tries to solve the above mentioned issues by using a dynamic scheduler for assigning independent MBs to threads. The main idea is to assign a MB to a processor as soon as its spatial and temporal dependencies have been resolved. The differences in processing time and the input dependent precedence constraining are considered by trying to process the next ready MB not the next MB in scan, row or frame order. This implies keeping-track both spatial and temporal MB dependencies at run time and generating a dynamic input-dependent data-flow. The global result is a reduction in the time that a thread spend waiting for ready-to-process MBs.

6.2 Implementation Issues of the Dynamic Scheduling

This dynamic scheduling can be implemented in various ways. Several issues play a role in this matter such as thread management overhead, data communication overhead, memory latency, centralized vs distributed control, etc. Many of these issues are new to the CMP paradigm and are research projects in itself.

A possible implementation that we are investigating is based on a dynamic task model (also known as a work-queuing model) and is depicted in Figure 20. In this model a set of threads is created/activated when a parallel region is encountered. In the case of the Dynamic 3D-Wave a parallel region is the decoding of all MBs in a frame/slice. Each parallel region is controlled by a frame manager, which consist of a thread pool, a task queue, a dependence

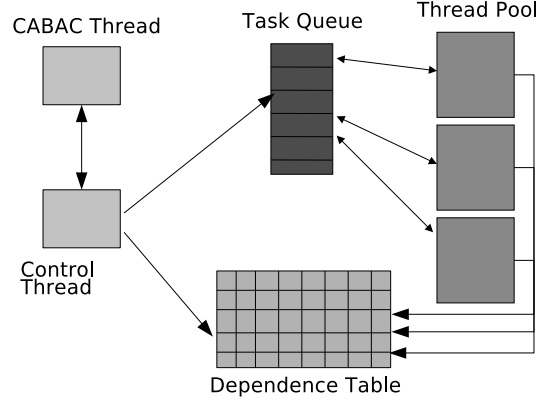


Fig. 20: Workqueue model for Dynamic 3D-Wave scheduling

table and a control thread as showed in Figure 20. The thread pool consists of a group of worker threads that wait for work on the task queue. The generation of work on the task queue is dynamic and asynchronous. The dependencies of each MB are expressed in a dependence table. When all the dependencies (inter- and intra-frame) for a MB are resolved a new task is inserted on the task queue. Any thread from the task queue can take this task and process it. When a thread finishes the processing of a task it updates the table of dependencies and if it finds that at least one MB has resolved its dependencies the current thread can decode it directly. If that thread wakes-up more than one MB it process one directly and submit the others into the task queue. If at the end of the update of the dependence table there are not a new ready MBs the thread goes to the task pool to ask for a new MB to be decoded. This technique has the advantage of distributed control, low thread management overhead, and also exploits spatial locality, reducing the data communication overhead. In this scheme the control thread is responsible only for handling all the initialization and finalization tasks of each that are not parallelizable, avoiding centralized schedulers that impose a big overhead.

Thread synchronization and communication overhead can also be minimized by taking advantage of clustering patterns in the dependencies and by using efficient hardware primitives for accessing shared structures and managing thread queues [6, 30]. Furthermore, to reduce data communication groups of MBs can be assigned to cores, possibly speculating on the MVs. If cores and memory are placed in a hierarchy, memory latency can be reduced by decoding spatially neighboring MBs in the same leaf of the hierarchy.

6.3 Support in the Programming Model

Currently most of the implementation of this dynamic system should be done with low-level threading APIs like POSIX threads. That is due to the fact that in current programming models for multicore architectures there is not a complete support for a dynamic work-queue model for irregular applications like H.264 decoding. For handling irregular loops with dynamic generation of threads there

has been some proposals of extensions to current parallel programming models like OpenMP [31, 32]. A feature that is missing is the possibility of specifying the structured dependencies among the different units of work. Currently this has to be implemented manually by the application programmer.

6.4 Managing Entropy Decoding

The different set of dependencies in the kernels imply that it is necessary to exploit both task- and data-level parallelism efficiently. Task-level decomposition is required because entropy decoding has parallelism at the slice and frame level only. We envision the following mapping to a heterogeneous CMP. Entropy decoding is performed on a high frequency core optimized for bit-serial operations. This can be an 8 or 16 bit RISC like core and therefore can be very power efficient. Several of these cores can operate in parallel each decoding a slice or a frame. The tasks for these cores can be spawned by another core that scans the bitstream, detects the presence or absence of slices, and provides each entropy decoder the start address of its bitstream. Once entropy decoding is performed, the decoding of MBs can be parallelized over as many cores there are available and as long as there are MBs available for decoding. The more cores are used, the lower the frequency and/or voltage can be assuming real-time performance. Thus the power efficiency scales with the number of cores.

One issue that has to be mentioned here is Amdahl's law for multicores, which states that the total speedup is limited by the amount of serial code. Considering a MB-level parallelization there is serial code at the start of each frame and slice. However, using the 3D-Wave strategy, this serial code can be hidden in the background it can be processed during the time that other frames or slices provide the necessary parallelism.

6.5 Open Issues

Despite the fact that there are many issues to be resolved in this area, it is of major importance to highly parallelize applications. It seems to be clear that in the future there will be manycores. However, what is currently being investigated by computer architecture researchers, just like the authors are doing, is what types of cores should be used, how the cores should be connected, what memory structure should be used, how the huge amount of threads should be managed, what programming model should be used, etc. To analyze design choices like these, it is necessary to analyze the characteristics of highly parallel applications. A scalability analysis like this is of great value in this sense.

7 Parallel Scalability of the Dynamic 3D-Wave

Now we have established a parallelization strategy that we expect to be scalable, we move on to our main goal: analyzing the parallel scalability of H.264 decoding. Remember we perform this analysis in the context of emerging manycore CMPs, that offer huge amounts of threads. In the previous section we explained how we envision to decouple entropy decoding by using task-level parallelism. Thus it scales with the number of slices and frames used in the 3D-Wave MB decoding. However, the parallelism is relatively small compared to the parallelism provided by the MB decoding. Therefore omitting the parallelism in entropy decoding does not significantly influence the results. Of course this assumes the entropy decoders are fast enough to support real-time decoding.

To investigate the amount of parallelism we modified the FFmpeg H.264 decoder to analyze real movies, which we took from the HD-VideoBench. We analyzed the dependencies of each MB and assigned it a timestamp as follows. The timestamp of a MB is simply the maximum of the timestamps of all MBs upon which it depends (in the same frame as well as in the reference frames) plus one. Because the frames are processed in decoding order¹, and within a frame the MBs are processed from left to right and from top to bottom, the MB dependencies are observed and it is assured that the MBs on which a MB B depends have been assigned their correct timestamps by the time the timestamp of MB B is calculated. As before, we assume that it takes one time slot to decode a MB.

The way the inter-frame dependencies are handled needs a detailed explanation. The pixels required for motion compensation are not only those of the sub-block the motion vector is pointing to. If the motion vector is fractional, the pixels used for MC are computed by interpolation and for this the surrounding pixels of the sub-block are required. Also the standard defines that deblocking filtered pixels should be used for MC. The bottom and right four pixels of each MB are filtered only when the bottom or left MB is decoded. Thus, although a MV points to a sub-block that lies entirely in a MB A , the actual dependencies might go as far as the right and/or bottom neighbor of MB A . We took this into account by adding to each MV the pixels needed for interpolation and checking if the DF of the neighboring blocks would affect the required pixels.

To start, the maximum available MB parallelism is analyzed. This experiment does not consider any practical or implementation issues, but simply explores the limits to the parallelism available in the application. We use the modified FFmpeg as described before and for each time slot we analyze, first, the number of MBs that can be processed in parallel during that time slot. Second, we keep track of the number of frames in flight. Finally, we keep track of the motion vector lengths. All experiments are performed for both normal encoded and esa encoded movies. The first uses a hexagonal (hex) motion estimation algorithm with a search range of 24 pixels, while the later uses an exhaustive search algorithm and results in worst case motion vectors.

Table 5 summarizes the results for normal encoded movies and shows that a huge amount of MB parallelism is available. Figure 21 depicts the MB parallelism time curve for FHD, while Figure 22 depicts the number of frames in flight. For the other resolutions the time curves have a similar shape. The

¹ The decoding order of frames is not equal to display order. The sequence I-B-B-P as in Figure 3 is decoded as I-P-B-B sequence.

Tab. 5: Maximum available MB parallelism, and frames in flight for normal encoded movies. Also the maximum and the average motion vectors (in square pixels) are stated.

	576				720			
	MBs	frames	max mv	avg mv	MBs	frames	max mv	avg mv
rush_hour	1199	92	380,5	1,3	2911	142	435	1,8
riverbed	1511	119	228,6	2	3738	185	496	2,2
pedestrian	1076	95	509,7	9,2	2208	131	553,5	11
blue_sky	1365	99	116	4,4	2687	135	298	5,1
	1088							
	MBs	frames	max mv	avg mv				
rush_hour	5741	211	441	2,2				
riverbed	7297	245	568,9	2,6				
pedestrian	4196	205	554,2	9,9				
blue_sky	6041	214	498	5,6				

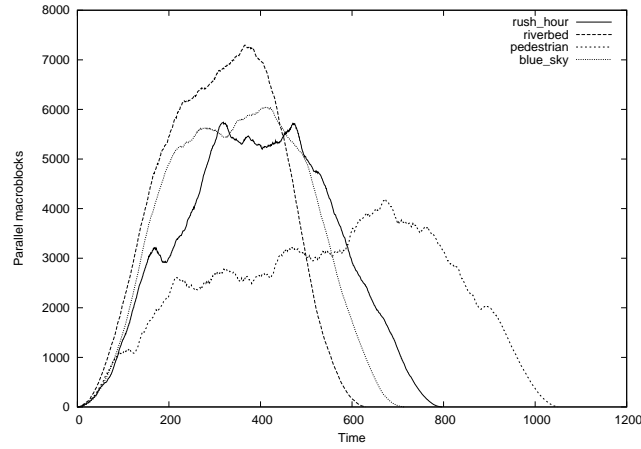


Fig. 21: Number of parallel MBs for FHD resolution using a 400 frames sequence.

MB parallelism time curve shows a ramp up, a plateau, and a ramp down. For example, for blue_sky the plateau starts at time slot 200 and last until time slot 570. Riverbed exhibits so much MB parallelism that it has a small plateau. Due to the stochastic nature of the movie, the encoder mostly uses intra-coding, resulting in very few dependencies between frames. Pedestrian exhibits the least parallelism. The fast moving objects in the movie result in many large motion vectors. Especially objects moving from right to left on the screen causes large offsets between MBs in consecutive frames.

Rather surprising is the maximum MV length found in the movies (see Table 5). The search range of the motion estimation algorithm was limited to 24 pixels, but still lengths of more than 500 square pixels are reported. According to the developers of the X264 encoder this is caused by drifting [33]. For each

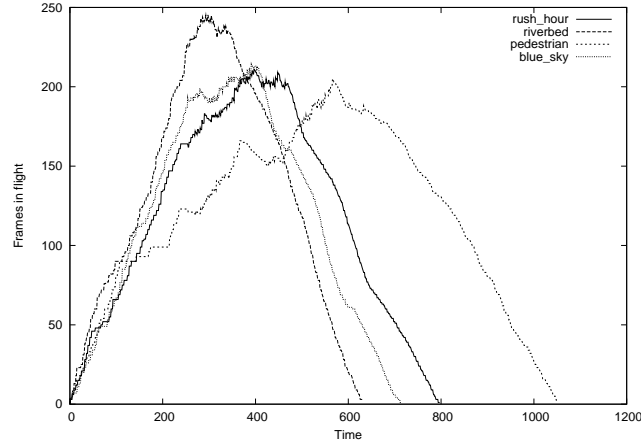


Fig. 22: Number of frames in flight for FHD resolution using a 400 frames sequence.

MB the starting MV of the algorithm is predicted using the result of surrounding MBs. If a number of MBs in a row use the motion vector of the previous MB and add to it, the values accumulate and reach large lengths. This drifting happens only occasionally, and does not significantly affect the parallelism using dynamic scheduling, but would force a static approach to use a large offset resulting in little parallelism.

Tab. 6: Maximum MB parallelism, and frames in flight for esa encoded movies. Also the maximum and the average motion vectors (in square pixels) are stated.

	576				720			
	MBs	frames	max mv	avg mv	MBs	frames	max mv	avg mv
rush_hour	505	63	712,0	6,9	708	84	699,3	9
riverbed	150	27	704,9	62,3	251	25	716,2	63,1
pedestrian	363	50	691,5	23,5	702	70	709,9	27,7
blue_sky	138	19	682,0	10,5	257	23	711,3	13
	1088							
	MBs	frames	max mv	avg mv				
rush_hour	1277	89	706,4	9,4				
riverbed	563	37	716,2	83,9				
pedestrian	1298	88	708,5	23,6				
blue_sky	500	24	705,7	15,7				

To evaluate the worst case condition the same experiment was performed for movies encoded using the exhaustive search algorithm (esa). Table 6 shows that the average MV length is substantially larger than for normal encoded movies. The exhaustive search decreases the amount of parallelism significantly. The time curves in Figures 23 and 24 have a very jagged shape with large peaks

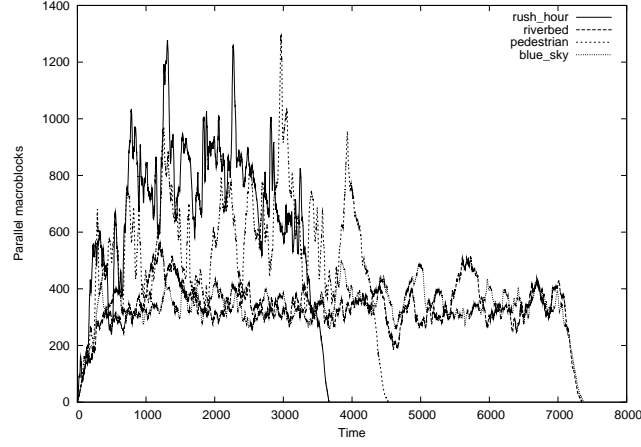


Fig. 23: Number of parallel MBs for FHD using esa encoding.

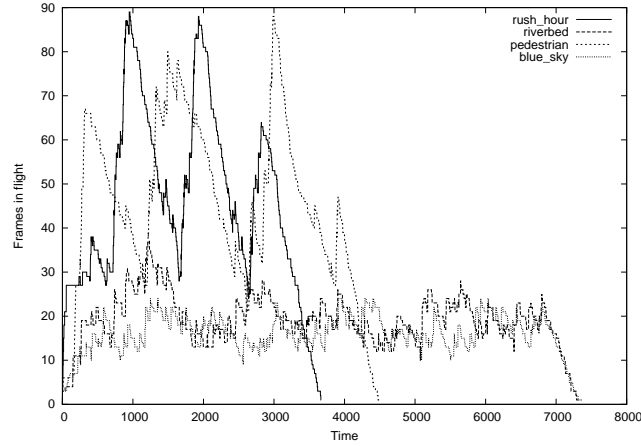


Fig. 24: Number of frames in flight for FHD using esa encoding.

and thus the average on the plateau is much lower than the peak. In Table 7 the maximum and the average parallelism for normal (hex) and esa encoded movies is compared. The average is taken over all 400 frames, including ramp up and ramp down. Although the drop in parallelism is significant, the amount of parallelism is still surprisingly large for this worst case scenario.

As we are comparing two encoding methods, it is interesting to look at some other aspects as well. In Table 8 we compare the compression ratio between hex and esa encoded movies. Only for FHD a considerable improvement is observed. However, the cost of this improvement in terms of encoding time is large. A 300 frame FHD sequence takes less than 10 minutes to encode using hexagonal motion estimation and running on a modern fast machine (Pentium D, 3.4 GHz, 2GB memory). The same sequence takes about two days when using the esa option, which is almost 300 times slower. Thus we can assume that this worst

Tab. 7: Comparison of the maximum and average MB parallelism for hex and esa encoded movies (FHD only).

	FHD					
	max MB			avg MB		
	hex	esa	Δ	hex	esa	Δ
rush_hour	5741	1277	-77.8%	4121.2	892.3	-78.3%
riverbed	7297	563	-92.3%	5181.0	447.7	-91.4%
pedestrian	4196	1298	-69.1%	3108.6	720.8	-76.8%
blue_sky	6041	500	-91.7%	4571.4	442.6	-90.3%

case scenario is unlikely to occur.

Tab. 8: Compression improvement of esa encoding relative to hex.

	SD	HD	FHD	average
rush_hour	-0.7%	0.0%	-0.3%	-0.3%
riverbed	1.4%	1.9%	3.2%	2.1%
pedestrian	0.6%	1.6%	1.2%	1.1%
blue_sky	-1.4%	-1.5%	1.4%	-0.5%
average	0.0%	0.5%	1.4%	0.6%

The analysis above reveals that H.264 exhibits significant amounts of MB parallelism. To exploit this type of parallelism on a CMP the decoding of MBs needs to be assigned to the available cores, i.e., MBs map to cores directly. However, even in future manycores the hardware resources (cores, memory, and NoC bandwidth) will be limited. We now investigate the impact of resource limitations.

We model limited resources as follows. A limited number of cores is modeled by limiting the number of MBs in flight. Memory requirements are mainly related to the number of frames in flight. Thus limited memory is modeled by restricting the number of frames that can be in flight concurrently. Limited NoC bandwidth is captured by both modeled restrictions. Both restrictions decrease the throughput, which is directly related to the inter-core communication.

The experiment was performed for all four movies of the benchmark, for all three resolutions, and both types of motion estimation. The results are similar, thus only the results for the normal encoded blue_sky movie at FHD resolution is presented.

First, the impact of limiting the number of MBs in flight is analyzed. Figure 25 depicts the available MB parallelism for several limits on the number of MBs in flight. As expected, for smaller limits, the height of the plateau is lower and the ramp up and ramp down are shorter. More important is that for smaller limits, the plateau becomes very flat. This translates to a high utilization rate of the available cores. Furthermore, the figure shows that the decoding time is approximately linear in the limit on MBs in flight. This is especially true for smaller limits, where the core utilization is almost 100%. For larger limits, the core utilization decreases and the linear relation ceases to be true. In real

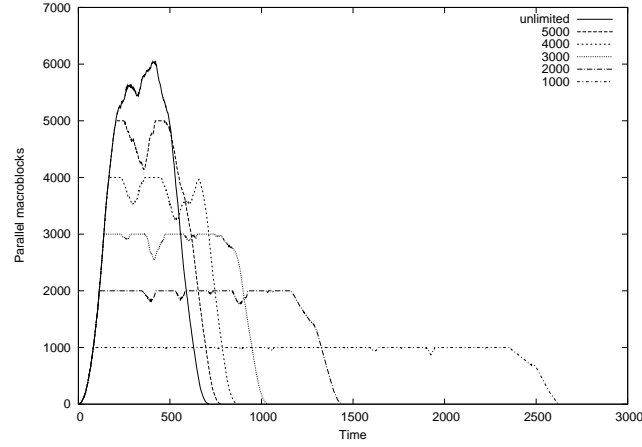


Fig. 25: Available MB parallelism in FHD blue_sky, for several limits of the number of MBs in flight.

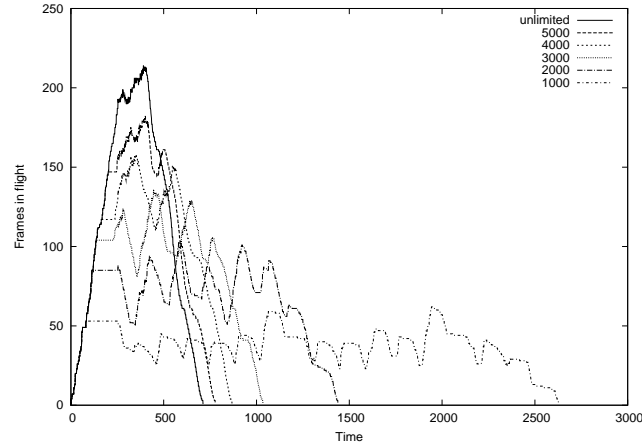


Fig. 26: Number of frames in flight for FHD blue_sky, for several limits of the number of MBs in flight.

systems, however, also communication overhead has to be taken into account. For smaller limits, this overhead can be reduced using a smart scheduler that assigns groups of adjacent MBs to cores. Thus, in real systems the communication overhead will be less for smaller limits than for larger limits.

Figure 26 depicts the number of frames in flight as a function of time when the number of MBs in flight is limited. Not surprisingly, because limiting the number of MBs in flight limits the number of frames that can be in flight, the shape of the curves are similar to those of the available MB parallelism, although with a small periodic fluctuation.

Next, we analyze the impact of restricting the number of frames concurrently in flight. The MB parallelism curves are depicted in Figure 27 and show large

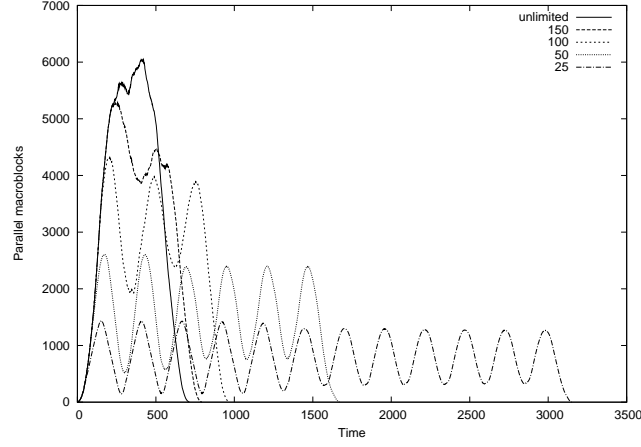


Fig. 27: Available MB parallelism in FHD blue_sky for several limits of the number of frames in flight.

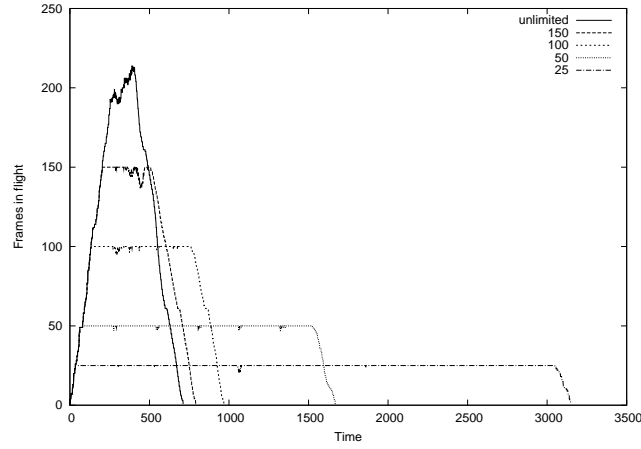


Fig. 28: Number of actual frames in flight for FHD blue_sky, for several limits of the number of frames in flight.

fluctuations, possibly resulting in under utilization of the available cores. These fluctuations are caused by the coarse grain of the limitation. At the end of decoding a frame, a small amount of parallelism is available. The decoding of a new frame, however, has to wait until the frame currently being processed is finished. Figure 28 shows the actual frames in flight for several limits. The plateaus in the curves translate to full memory utilization.

From this experiment we can conclude that for systems with limited number of cores dynamic scheduling is able to achieve near optimal performance by (almost) fully utilizing the available computational power. On the contrary, for systems where memory is the bottleneck, additional performance losses might occur because of temporal underutilization of the available cores. In Section 8

we perform a case study, indicating that memory will likely not be a bottleneck.

We have shown that using the Dynamic 3D-Wave strategy huge amounts of parallelism are available in H.264. Analysis of real movies revealed that the number of parallel MBs ranges from 3000 to 7000. This amount of MB parallelism might be larger than the number of cores available. Thus, we have evaluated the parallelism for limited resources and found that limiting the number of MBs in flight results in a equivalent larger decoding time, but a near optimal utilization of the cores.

8 Case Study: Mobile Video

So far we mainly focused on high resolution and explored the potential available MB parallelism. The 3D-Wave strategy allows an enormous amount of parallelism, possibly more than what a high performance CMP in the near future could provide. Therefore, in this section we perform a case study to assess the practical value and possibilities of a highly parallelized H.264 decoding application.

For this case study we assume a mobile device such as the iPhone, but in the year 2015. We take a resolution of 480×320 just as the screen of the current iPhone. It is reasonable to expect that the size of mobile devices will not significantly grow, and limitations on the accuracy of the human eye prevent the resolution from increasing.

CMPs are power efficient, and thus we expect mobile devices to adopt them soon and project a 100-core CMP in 2015. Two examples of state-of-the-art embedded CMPs are the Tiler Tile64 and the Clearspeed CSX600. The Tile64 [2] processor has 64 identical cores that can run fully autonomously and is well suited to exploit TLP. The CSX600 [1] contains 96 processing elements, which are simple and exploit DLP. Both are good examples of what is already possible today and the eagerness of the embedded market to adopt the CMP paradigm. Some expect a doubling of cores every three years [3]. Even if this growth would be less the assumption of 100 cores in 2015 seems reasonable.

The iPhone is available with 8GB of memory. Using Moore's law for memory (doubling every two years), in 2015 this mobile device would contain 128 GB. It is unclear how much of the flash drive memory is used as main memory, but let us assume it is only 1%. In 2015 that would mean 1.28GB and if only half of it is available for video decoding, then still almost 1400 frames of video fit in main memory. It seems that memory is not going to be a bottleneck.

We place a limit of 1 second on the latency of the decoding process. That is a reasonable time to wait for the movie to appear on the screen, but much longer would not be acceptable. The iPhone uses a frame rate of 30 frames/second, thus limiting the frames in flight to 30, causes a latency of 1s. This can be explained as follows. Let us assume that the clock frequency is scaled down as much as possible to save power. When frame x is started decoding, frame $x - 30$ has just finished and is currently displayed. But, since the framerate is 30, frame x will be displayed 1 second later. Of course this is a rough method and does not consider variations in frame decoding time, but it is a good estimate.

Putting all assumptions together, the resulting Dynamic 3D-Wave strategy has a limit of 100 parallel MBs and up to 30 frames can be in flight concurrently. For this case study we use the four normal encoded movies with a length of 200 frames and a resolution of 480×320 .

Figure 29 presents the MB parallelism under these assumptions as well as for unlimited resources. The picture shows that even for this small resolution, all 100 cores are utilized nearly all time. The curves for unrestricted resources show that there is much more parallelism available than the hardware of this case study offers. This provides opportunities for scheduling algorithms to reduce communication overhead.

Figure 30 depicts the number of frames in flight. The average is approximately 12 frames with peaks of up to 19. From the graphs we can also conclude that the decoding latency is approximately 0.5 second, as around 15 frames are

in flight.

From this case study we can draw the following conclusions:

- Even a low resolution movie exhibits sufficient parallelism to fully utilize 100 cores efficiently.
- For mobile devices, memory will likely not be a bottleneck.
- For mobile devices with a 100-core CMP, start-up latency will be short (0.5s).

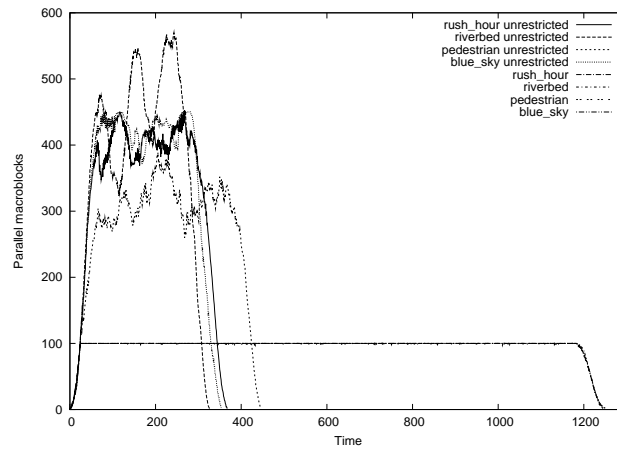


Fig. 29: Number of parallel MBs for the mobile video case study. Also depicted is the available MB parallelism for the same resolution but with unrestricted resources.

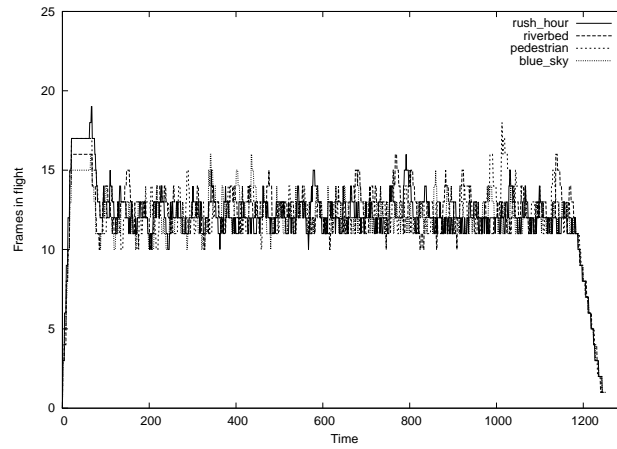


Fig. 30: Number of frames in flight for the mobile video case study.

9 Related Work

A lot of work has been done on the parallelization of video Codecs but most of it has been focused on coarse grain parallelization for large scale multiprocessors systems.

Shen et al. [34] have presented a parallelization of the MPEG-1 encoder for the Intel Paragon MIMD multiprocessor at the GOP level. Their approach scale to a large number of processors but with a large latency. Bilas et al. [35] has described a parallel implementation of the MPEG-2 decoder evaluating GOP and slice-level parallelism on a shared memory multiprocessor. They compare the two levels in terms of speedup, load balancing, memory usage and synchronization overhead.

Some other works have reported MB level parallelism at the encoder side for the ME stage by performing ME in each processing element and replicating the search window. Akramullah et al. [36] have reported a parallelization of the MPEG-2 encoder for cluster of workstations and Taylor et al. [37] have implemented an MPEG-1 encoder in a multiprocessor made of 2048 custom DSPs.

Other works have analyzed function-level parallelism. Lin et al. [38] and Cantineau et al. [39] have reported functional parallelization of the H.263 and MPEG-2 encoders respectively. They have used a multiprocessor with one control processor and four DSPs. Oehring et al. [40] have reported an analysis of the parallelization of the MPEG-2 decoder exploiting function level parallelism, and using a simulator for SMT processors with 4 threads.

There are some works on slice level parallelism for previous video Codecs. Lehtoranta et al [41] have described a parallelization of the H.263 encoder at the slice level for a multiprocessor system made of 4 DSP cores with one master processor and 3 computing processors. Lee et al [42] have also reported an implementation of the MPEG-2 decoder exploiting slice-level parallelism for HDTV input videos.

Yadav et al [43] have reported a study on the parallelization of the MPEG-2 decoder for a multiprocessor SoC architecture. They studied slice-level, function-level (error computation and prediction) and stream-level parallelism (VLD, IQ, IDCT). Jacobs et al [43] have analyzed the thread parallelism of MPEG-2, MPEG-4 and H.264 decoders. For MPEG-2 and MPEG-4 they have taken into account MB-level parallelism. For these cases the parallel version scales with the height of the frame. For the H.264 encoder they discarded the idea of implementing MB-level parallelism and went for slice-level.

In the specific case of H.264 there has been some previous works on multi-threaded versions for GOP, frame, slice, function and MB level parallelism.

Gulati et al. [44] describe a system for encoding and decoding H.264 on a multiprocessor architecture using a task-level decomposition approach. The multiprocessor includes eight DSPs and four control processors. The decoder is implemented using one control processor and 3 DSPs. The control processor performs the master control of the application and an initial parsing of the bitstream. Entropy decoding, inverse quantization, and IDCT are mapped to the first DSP, motion compensation and intra-prediction to the second, and finally the deblocking filter is mapped to the third DSP. Using this mapping, a pipeline for processing MBs is implemented. This system achieves real-time operation for low resolution video inputs, using the baseline profile. For more

demanding applications (HD video, High profiles of H.264) this scheme can not scale in performance because it is very complicated to further divide the current mapping of tasks to use more processors. In a similar way Schoffmann et al. [45] propose a macroblock pipeline model for H.264 decoding. Using an Intel Xeon 4-way SMP 2-SMT architecture their scheme achieves a 2X speed-up, but as mentioned previously macroblock pipelining does not scale for manycore architectures.

Frame-level parallelism has been described for the H.264 encoder in the work of Chen et al. [46]. In this case a combination of frame-level and slice-level parallelism is proposed. First, they exploit frame level parallelism. (to gain something from it they do not allow to use B-frames as references and use a static P-B-B-P-B-B sequence of frames). When the limit of frame-level parallelism has been reached they explore slice-level parallelism. By using this approach 4.5X speed-up is achieved in a machine with 8 cores and 9 slices. Their approach, however, does not scale to more processors because the limits in frame-level parallelism and the bit-rate increase due to slices.

Frame-level parallelism of dependent frames has been implemented in the X264 open-source encoder [18] and in the H.264 Codec that is part of the Intel Integrated Performance Primitives [47]. In these cases the encoder limits the vertical motion search range in order to allow a next frame to start the encoding before the current one completes. This approach does not include macroblock-level parallelism and thus only one macroblock per frame at a time is decoded. Results of scalability of the X264 version are presented in Figure 12. As mentioned previously this approach can not be exploited efficiently in the decoder because there is not a guarantee that the motion vectors have limited values. The scalability is limited by the height of the frame, and may be more reduced if some advanced coding techniques like adaptive placement of B frames, scene cut and rate control are used, because they require a motion estimation prepass that has to be implemented in a sequential way, becoming the limiting factor according to Amhdal law.

In [48] a scheme is proposed for exploiting slice-level parallelism in the H.264 decoder by modifying the encoder. The main idea is to overcome the load balancing disadvantage by making an encoder that produces slices that are not balanced in the number of MBs, but in their decoding time. The main disadvantages of this approach is that it requires modifications to the encoder in order to exploit parallelism at the decoder, and the inherent loss of coding efficiency due to having a large number of slices.

In [49] the authors proposed an encoder that combines GOP-level and slice-level parallelism for encoding real-time H.264 video using clusters of workstations. Frame-level parallelism is exploited by assigning GOPs to nodes in a cluster, and after that, slices are assigned to processors in each cluster node. This methodology is only applicable to the encoder, which has the freedom of selecting the GOP size. It scales to a small number of processors in each node, because of the loss in coding efficiency due to slices.

MB-level parallelism has been proposed in several works. Van der Tol et al. [50] have proposed the exploitation of MB-level parallelism for optimizing the H.264 decoding. The analysis was made for a multiprocessor system consisting of 8 Trimedia processors with private L1 caches and a shared L2 cache. Also the paper proposed the grouping of MBs into data partitions and a special memory allocation technique for allocating those partitions in the cache in order

to optimize the access to the MBs. The paper also suggests the combination of MB-level with frame-level parallelism to increase the availability of independent MBs. The use of frame-level parallelism is determined statically by the length of the motion vectors. Chen et al. [51] have evaluated a similar approach: a combination of MB parallelism and frame-level parallelism for the H.264 encoder on Pentium machines with SMT and CMP capabilities. In the above mentioned works the exploitation of frame-level parallelism is limited to two consecutive frames and the identification of independent MBs is done statically by taking into account the limits of the motion vectors. Although this combination of MB and frame-level parallelism increases the amount of independent MBs, it requires that the encoder puts some limits on the length of the motion vectors in order to define statically which MBs of the next frame can start execution while decoding the current frame.

A combination of frame-level parallelism and macroblock-level parallelism for H.264 encoding has been discussed by Zhao et al. [29]. In this approach multiple frames are processed in parallel similar to X264: a new frame is started when the search area in the reference frame has been fully encoded. But in this case macroblock-level parallelism is also exploited by assigning threads to different rows inside a frame. This scheme is a variation of what we call in this paper "Static 3D-Wave". It is static because it depends on a fixed value of the motion vectors in order to exploit frame-level parallelism. This works for the encoder who has the flexibility of selecting the search area for motion estimation, but in the case of the decoder the motion vectors can have big values and the static approach should take that into account. Another limitation of this approach is that all the MBs in the same row are processed by the same processor/thread. This can result in poor load balancing because the decoding time of MBs is not constant. The Dynamic 3D-Wave approach allows a thread that has finished a fast MB to start a new row when the next MB in the current row can not be decoded because it is waiting for a slow MB in the previous row. Additionally this work is focused on the performance improvement on the encoder for slow resolutions with a small number of processors, discarding an analysis for a high number of cores as an impractical case and not taking into account the peculiarities of the decoder. Finally, they do not address the scalability issues of their approach in the context of manycore architectures and high definition applications.

Chong et al [52] has proposed another technique for exploiting MB level parallelism in the H.264 decoder by adding a prepass stage. In this stage the time to decode a MB is estimated heuristically using some parts of the compressed information of that MB. Using the information from the preparsing pass a dynamic schedule of the MBs of the frame is calculated. MBs are assigned to processors dynamically according to this schedule. By using this scheme a speedup of 3.5X has been reported on a 6 processors simulated system for small resolution input videos. Although they present a dynamic scheduling algorithm it seems to be not able of discovering all the MB level parallelism that is available in a frame. The preparsing scheme presented can be beneficial for the Dynamic 3D-Wave algorithm.

10 Conclusions

In this paper we have investigated if future applications exhibit sufficient parallelism to exploit the large number of cores expected in future CMPs. This is important for performance but maybe even more important for power efficiency. As a case study, we have analyzed the parallel scalability of the H.264 video decoding process, currently the best coding standard.

First, we have discussed the parallelization possibilities and showed that slice-level parallelism has two main limitations. First, using many slices increases the bitrate and, second, not all sequences contain many slices since the encoder determines the number of slices per frame. It was also shown that frame-level parallelism, which exploits the fact that some frames (B frames) are not used as reference frames and can therefore be processed in parallel, is also not very scalable because usually there are no more than three B frames between consecutive P frames and, furthermore, in H.264 B frames can be used as reference frames.

More promising is MB-level parallelism, which can be exploited inside a frame in a diagonal wavefront manner. It was observed that MBs in different frames are only dependent through motion vectors which have a limited range. We have therefore proposed a novel parallelization strategy, called Dynamic 3D-Wave, which combines MB parallelism with frame-level parallelism.

Using this new parallelization strategy we performed a limit study to the amount of MB-level parallelism available in H.264. We modified Ffmpeg and analyzed real movies. The results revealed that motion vectors are small on average. As a result a large amount of parallelism is available. For full HD resolution the total number of parallel MBs ranges from 4000 to 7000.

Combining macroblock and frame-level parallelism was done before for H.264 encoding in a static way. To compare this with our new dynamic scheme we analyzed the scalability of the Static 3D-Wave approach for decoding. Our results show that this approach can deliver a parallelism of 1360 MBs, but only if motion vectors are strictly shorter than 16 pixels. In practice this condition cannot be guaranteed and the worst case motion vector length has to be assumed. This results in a parallelism of 89 MBs for full HD. It is clear that the Dynamic 3D-Wave outperforms the static approach.

The amount of MB parallelism offered by the Dynamic 3D-Wave might be larger than the number of cores available. Moreover, the number of frames in flight to achieve this parallelism might require more memory than available. Thus, we have evaluated the parallelism for limited resources and found that limiting the number of MBs in flight results in a larger decoding time, but a near optimal utilization of cores. Limiting the number of frames in flight, on the other hand, results in large fluctuations of the available parallelism, likely causing underutilization of the available cores and thus performance loss.

We have also performed a case study to assess the practical value and possibilities of the highly parallelized H.264 decoding. In general, the results show that our strategy provides sufficient parallelism to efficiently exploit the capabilities of future manycore CMPs.

The large amount of parallelism available in H.264 indicates that future manycore CMPs can effectively be used. Although we have focused on H.264, other video codecs and multimedia applications in general exhibit similar features and we expect that they can also exploit the capabilities of manycore

CMPs. However, the findings of this paper also pose a number of new questions, which we intend to investigate in the future.

First, we will extend the parallelism analysis to include variable MB decoding time. We will investigate its impact on the dynamic scheduling algorithm, the amount of available parallelism, and performance. Second, we will investigate MB clustering and scheduling techniques to reduce communication and synchronization overhead. We are currently working on an implementation of the 3D-Wave strategy for an aggressive multicore architecture and the MB scheduling and clustering algorithm is an important part of that. Finally, we will investigate ways to accelerate entropy coding. More specifically, we plan to design a hardware accelerator tailored for contemporary and future entropy coding schemes.

Acknowledgments

This work was partially supported by the European Commission in the context of the SARC integrated project #27648 (FP6), the Ministry of Science of Spain and European Union (FEDER funds) under contract TIC-2004-07739-C02-01, and HiPEAC, European Network of Excellence on High-Performance Embedded Architecture and Compilation. The authors would like to thank Jan Hoogerbrugge from NXP labs for the valuable discussions on parallelizing H.264.

References

- [1] ClearSpeed, “The CSX600 Processor.” [Online]. Available: <http://www.clearspeed.com>
- [2] Tiler, “TILE64(TM) Processor Family.” [Online]. Available: <http://www.tiler.com>
- [3] P. Stenström, “Chip-multiprocessing and Beyond,” in *Proc. Twelfth Int. Symp. on High-Performance Computer Architecture*, 2006, pp. 109–109.
- [4] K. Asanovic *et al.*, “The Landscape of Parallel Computing Research: A View from Berkeley,” EECS Department University of California, Berkeley, Tech. Rep. UCB/EECS-2006-183, December 2006.
- [5] H. Liao and A. Wolfe, “Available parallelism in video applications,” in *proc. 30th Annual Int. Symp. on Microarchitecture (Micro '97)*, 1997.
- [6] K. Kissell, “MIPS MT: A Multithreaded RISC Architecture for Embedded Real-Time Processing,” in *”proc. High Performance Embedded Architectures and Compilers (HiPEAC) Conference*, 2008.
- [7] T. Oelbaum, V. Baroncini, T. Tan, and C. Fenimore, “Subjective Quality Assessment of the Emerging AVC/H.264 Video Coding Standard,” in *Int. Broadcast Conference (IBC)*, 2004.
- [8] “International Standard of Joint Video Specification (ITU-T Rec. H. 264—ISO/IEC 14496-10 AVC),” 2005.
- [9] M. Flierl and B. Girod, “Generalized B Pictures and the Draft H. 264/AVC Video-Compression Standard,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 587–597, 2003.
- [10] T. Wiegand, G. J. Sullivan, G. Bjontegaard, and A. Luthra, “Overview of the H.264/AVC Video Coding Standard,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 560–576, July 2003.
- [11] A. Tamhankar and K. Rao, “An Overview of H. 264/MPEG-4 Part 10,” in *Proc. 4th EURASIP Conference focused on Video/Image Processing and Multimedia Communications*, 2003, p. 1.
- [12] H. Malvar, A. Hallapuro, M. Karczewicz, and L. Kerofsky, “Low-Complexity Transform and Quantization in H. 264/AVC,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 598–603, 2003.
- [13] M. Wien, “Variable Block-Size Transforms for H. 264/AVC,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 604–613, 2003.
- [14] P. List, A. Joch, J. Lainema, G. Bjntegaard, and M. Karczewicz, “Adaptive Deblocking Filter,” *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 614–619, 2003.

- [15] D. Marpe, H. Schwarz, and T. Wiegand, "Context-Based Adaptive Binary Arithmetic Coding in the H.264/AVC Video Compression Standard," *IEEE Trans. on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 620–636, 2003.
- [16] G. Sullivan, P. Topiwala, and A. Luthra, "The H.264/AVC Advanced Video Coding Standard: Overview and Introduction to the Fidelity Range Extensions," in *Proc. SPIE Conference on Applications of Digital Image Processing XXVII*, 2004, pp. 454–474.
- [17] M. Alvarez, E. Salami, A. Ramirez, and M. Valero, "HD-VideoBench: A Benchmark for Evaluating High Definition Digital Video Applications," in *IEEE Int. Symp. on Workload Characterization*, 2007. [Online]. Available: <http://personals.ac.upc.edu/alvarez/hdvideobench/index.html>
- [18] "X264. A Free H.264/AVC Encoder." [Online]. Available: <http://developers.videolan.org/x264.html>
- [19] "The FFmpeg Libavcoded." [Online]. Available: <http://ffmpeg.mplayerhq.hu/>
- [20] M. Alvarez, E. Salami, A. Ramirez, and M. Valero, "A Performance Characterization of High Definition Digital Video Decoding using H.264/AVC," in *Proc. IEEE Int. Workload Characterization Symposium*, 2005, pp. 24–33.
- [21] J. Ostermann, J. Bormans, P. List, D. Marpe, M. Narroschke, F. Pereira, T. Stockhammer, and T. Wedi, "Video Coding with H.264/AVC: Tools, Performance, and Complexity," *IEEE Circuits and Systems Magazine*, vol. 4, no. 1, pp. 7–28, 2004.
- [22] V. Lappalainen, A. Hallapuro, and T. D. Hamalainen, "Complexity of Optimized H.26L Video Decoder Implementation," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 717–725, 2003.
- [23] M. Horowitz, A. Joch, and F. Kossentini, "H.264/AVC Baseline Profile Decoder Complexity Analysis," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 704–716, 2003.
- [24] M. Flierl and B. Girod, "Generalized B pictures and the draft H.264/AVC video-compression standard," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 13, no. 7, pp. 587–597, July 2003.
- [25] Y. Liu and S. Orintara, "Complexity Comparison of fast Block-Matching Motion Estimation Algorithms," in *Proc. IEEE Int. Conf. on Acoustics, Speech, and Signal Processing*, 2004.
- [26] X. Zhou, E. Q. Li, and Y.-K. Chen, "Implementation of H.264 Decoder on General-Purpose Processors with Media Instructions," in *Proc. SPIE Conf. on Image and Video Communications and Processing*, 2003.
- [27] H. Shojania, S. Sudharsanan, and C. Wai-Yip, "Performance improvement of the h.264/avc deblocking filter using simd instructions," in *Proc. IEEE Int. Symp. on Circuits and Systems ISCAS*, May 2006.

- [28] J. Lee, S. Moon, and W. Sung, "H.264 Decoder Optimization Exploiting SIMD Instructions," in *Asia-Pacific Conf. on Circuits and Systems*, Dec. 2004.
- [29] Z. Zhao and P. Liang, "Data partition for wavefront parallelization of H.264 video encoder," in *IEEE International Symposium on Circuits and Systems. ISCAS 2006*, 21-24 May 2006.
- [30] P. Schaumont, B.-C. C. Lai, W. Qin, and I. Verbauwhede, "Cooperative multithreading on embedded multiprocessor architectures enables energy-scalable design," in *DAC '05: Proceedings of the 42nd annual conference on Design automation*, 2005, pp. 27–30.
- [31] X. Tian, Y.-K. Chen, M. Girkar, S. Ge, R. Lienhart, and S. Shah, "Exploring the use of hyper-threading technology for multimedia applications with intel openmp compiler," *Proceedings of the International Parallel and Distributed Processing Symposium, 2003*, 22-26 April 2003.
- [32] E. Ayguad and N. Coptý and A. Duran and J. Hoeflinger and Y. Lin and F. Massaioli and E. Su and P. Unnikrishnan and G. Zhang, "A proposal for task parallelism in OpenMP," in *Proceedings of the 3rd International Workshop on OpenMP*, June 2007.
- [33] "X264-devel – Mailing list for x264 developers," July-August 2007, subject: out-of-range motion vectors. [Online]. Available: <http://mailman.videolan.org/listinfo/x264-devel>
- [34] K. Shen, L. A. Rowe, and E. J. Delp, "Parallel implementation of an MPEG-1 encoder: faster than real time," in *Proc. SPIE Vol. 2419, p. 407-418, Digital Video Compression: Algorithms and Technologies 1995*, 1995, pp. 407–418.
- [35] A. Bilas, J. Fritts, and J. Singh, "Real-time parallel mpeg-2 decoding in software," *Parallel Processing Symposium, 1997. Proceedings., 11th International*, pp. 197–203, 1-5 Apr 1997.
- [36] S. Akramullah, I. Ahmad, and M. Liou, "Performance of software-based mpeg-2 video encoder on parallel and distributed systems," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 7, no. 4, pp. 687–695, Aug 1997.
- [37] H. Taylor, D. Chin, and A. Jessup, "A mpeg encoder implementation on the princeton engine video supercomputer," *Data Compression Conference, 1993. DCC '93.*, pp. 420–429, 1993.
- [38] W. Lin, K. Goh, B. Tye, G. Powell, T. Ohya, and S. Adachi, "Real time h.263 video codec using parallel dsp," *International Conference on Image Processing, 1997.*, vol. 2, pp. 586–589 vol.2, 26-29 Oct 1997.
- [39] O. Cantineau and J.-D. Legat, "Efficient parallelisation of an mpeg-2 codec on a tms320c80 video processor," *1998 International Conference on Image Processing, 1998. ICIP 98.*, pp. 977–980 vol.3, 4-7 Oct 1998.

- [40] H. Oehring, U. Sigmund, and T. Ungerer, "Mpeg-2 video decompression on simultaneous multithreaded multimedia processors," *International Conference on Parallel Architectures and Compilation Techniques, 1999.*, pp. 11–16, 1999.
- [41] O. Lehtoranta, T. Hamalainen, and J. Saarinen, "Parallel implementation of h.263 encoder for cif-sized images on quad dsp system," *The 2001 IEEE International Symposium on Circuits and Systems, ISCAS 2001*, vol. 2, pp. 209–212 vol. 2, 6-9 May 2001.
- [42] C. Lee, C. S. Ho, S.-F. Tsai, C.-F. Wu, J.-Y. Cheng, L.-W. Wang, and C. Wang, "Implementation of digital hdtv video decoder by multiple multimedia video processors," *International Conference on Consumer Electronics, 1996*, pp. 98–, 5-7 Jun 1996.
- [43] R. S. Ganesh Yadav and V. Chaudhary, "On Implementation of MPEG-2 Like Real-Time Parallel Media Applications on MDSP SoC Cradle Architecture," in *Lecture Notes in Computer Science. Embedded and Ubiquitous Computing*, Jul. 2004.
- [44] A. Gulati and G. Campbell, "Efficient Mapping of the H.264 Encoding Algorithm onto Multiprocessor DSPs," in *Proc. Embedded Processors for Multimedia and Communications II*, vol. 5683, no. 1, March 2005, pp. 94–103.
- [45] O. L. Klaus Schffmann, Markus Fauster and L. Bszrmenyi, "An Evaluation of Parallelization Concepts for Baseline-Profile Compliant H.264/AVC Decoders," in *Lecture Notes in Computer Science. Euro-Par 2007 Parallel Processing*, August 2007.
- [46] Y. Chen, X. Tian, S. Ge, and M. Girkar, "Towards efficient multi-level threading of h.264 encoder on intel hyper-threading architectures," in *Proc. 18th Int. Parallel and Distributed Processing Symposium*, 2004.
- [47] "Intel Integrated Performance Primitives." [Online]. Available: <http://www.intel.com/cd/software/products/asmona/eng/perflib/ipp/302910.htm>
- [48] M. Roitzsch, "Slice-Balancing H.264 Video Encoding for Improved Scalability of Multicore Decoding," in *Work-in-Progress Proc. 27th IEEE Real-Time Systems Symposium*, 2006.
- [49] A. Rodriguez, A. Gonzalez, and M. P. Malumbres, "Hierarchical parallelization of an h.264/avc video encoder," in *Proc. Int. Symp. on Parallel Computing in Electrical Engineering*, 2006, pp. 363–368.
- [50] E. van der Tol, E. Jaspers, and R. Gelderblom, "Mapping of H.264 Decoding on a Multiprocessor Architecture," in *Proc. SPIE Conf. on Image and Video Communications and Processing*, 2003.
- [51] Y. Chen, E. Li, X. Zhou, and S. Ge, "Implementation of H. 264 Encoder and Decoder on Personal Computers," *Journal of Visual Communications and Image Representation*, vol. 17, 2006.

-
- [52] J. Chong, N. R. Satish, B. Catanzaro, K. Ravindran, and K. Keutzer, “Efficient parallelization of h.264 decoding with macro block level scheduling,” in *2007 IEEE International Conference on Multimedia and Expo*, July 2007, pp. 1874–1877.