



**Miguel Oliveira
Inocêncio**

Co-processador da Transformada para AV1

AV1 Transform Co-Processor

DOCUMENTO PROVISÓRIO



**Miguel Oliveira
Inocêncio**

Co-processador da Transformada para AV1

AV1 Transform Co-Processor

Dissertação de Mestrado apresentada à Universidade de Aveiro, para obtenção do grau de Mestre em Engenharia Eletrónica e de Telecomunicações, sob orientação do Professor Doutor António Navarro ...

DOCUMENTO PROVISÓRIO

o júri / the jury

presidente / president

ABC

Professor Catedrático da Universidade de Aveiro (por delegação da Reitora da Universidade de Aveiro)

vogais / examiners committee

DEF

Professor Catedrático da Universidade de Aveiro (orientador)

GHI

Professor associado da Universidade J (co-orientador)

KLM

Professor Catedrático da Universidade N

agradecimientos / acknowledgements

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum...

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris...

Palavras-Chave

Resumo

HEVC, ...

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec non-ummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Keywords

Abstract

HEVC, ...

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec non-ummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Contents

List of Figures	iv
List of Tables	v
Acronyms	viii
Glossary	ix
Nomenclature	xi
1 Introduction	1
1.1 Background and Motivation	1
1.2 Scope	3
1.3 Outline	3
References	4
2 Video Compression Systems	5
2.1 Basic Principles	5
2.1.1 Human Visual System	5
2.1.2 Redundancy Exploitation	8
2.1.3 Basic Video Compression/Decompression System	10
2.2 AV1	15
2.2.1 History and Development	15
2.2.2 Encoding Tools	17
2.2.3 Performance Analysis	19
References	24
3 Video Coding Transforms	25
3.1 Introduction	25
3.2 Background	25
3.2.1 Basis vector/image interpretation	25
3.3 Transformation Kernels	27
3.3.1 Discrete Fourier Transform (DFT)	27
3.3.2 Discrete Walsh-Hadamard Transform (WHT)	28
3.3.3 Discrete Cosine Transform (DCT)	29
3.3.4 Discrete Sine Transform (DST)	31
3.3.5 Asymmetric Discrete Sine Transform (ADST)	32
3.4 Libaom's Integer Transformations	32
3.4.1 Functioning and Implementation	33
3.4.2 Performance and Statistics Analysis	38

References	47
4 Developed Architecture	49
4.1 TO BE DONE	49
Appendices	50
Appendix A Configuration Options	51

List of Figures

1.1	<i>AV1</i> logo	3
2.1	Representation of an human eye	6
2.2	Example of the effect of added noise on figure	7
2.3	Autocorrelation of image 2.2a, with horizontal shifts	8
2.4	Cross-correlation between the first and following nine frames of the <i>Stefan</i> sequence	9
2.5	Representation of chroma subsampling	10
2.6	Simplified Basic Encoder Model	11
2.7	Directional Intra-prediction example	12
2.8	Inter-prediction example	12
2.9	Demonstration of Zig-Zag Scan	13
2.10	Simplified Basic Decoder Model	14
2.11	Processing of 4×4 residue block from <i>transformation</i> to restoring	15
2.12	<i>Alliance for Open Media</i> current members	16
2.13	Description of the recursive partitioning scheme of <i>AV1</i>	17
2.14	<i>AV1</i> bitrate savings	20
3.1	Sequences generated in the first step of table 3.1for the DFT and different DCTs. Filled dots correspond to the original sequence.	30
3.2	Flowchart of the Transform Stage on <i>libaom</i>	34
3.3	Graphical aid for figures 3.4 and 3.5	35
3.4	Block diagram of <i>libaom</i> 's Integer DCT	36
3.5	Block diagram of <i>libaom</i> 's Integer ADST	37
3.6	Description of the Identity transforms in <i>libaom</i>	38
3.7	Average encoding and transform time per resolution, on different quality objectives	40
3.8	Average distribution of used kernels, for all resolutions, according to the quality threshold	41
3.9	Average distribution of vector sizes, for all resolutions, according to the quality threshold	41
3.10	Use of square blocks, same kernel for rows and columns, and symmetric kernels, according to the quality threshold	42
3.11	Different number of bits used on the cosine approximations, throughout different quality sets	42
3.12	Description of the test for comparing impact of number of bits in cosine approximations	44
3.13	Obtained quality for each of the quality objectives, and comparison with different cosine bits approximation	44

3.14	Detail of <i>Parkjoy</i> encodes, through different quality objectives	45
3.15	<i>Quantizer</i> distribution on different quality objectives	46

List of Tables

2.1	BD-rate of <i>VP9</i> and <i>H.264</i> codecs, when compared to <i>AV1</i> (negative corresponds to bitrate savings)	20
2.2	Encoding times of different video encoders, and improvements on <i>AV1</i>	21
3.1	Similarity between the processes of the <i>DFT</i> and the <i>DCT</i>	29
3.2	Sequences used for testing	39

Acronyms

E_{ms} Mean Square Error.

ADST Asymmetric Discrete Sine Transform.

AOM Alliance for Open Media.

AV1 AOM Video 1.

CABAC Context Adaptative Binary Arithmetic Coding.

CMOS Complementary metal–oxide–semiconductor.

Codec Encoder-Decoder.

CRT Cathode Ray Television.

DCT Discrete Cosine Transform.

DFT Discrete Fourier Transform.

FFT Fast Fourier Transform.

fps frames per second.

GPU Graphical Processing Unit.

HEVC High Efficiency Video Coding.

IC Integrated Circuit.

JVT Joint Video Team.

MM Matrix Multiplication.

MPEG Motion Picture Experts Group.

PSNR Peak Signal to Noise Ratio.

QP quantization parameter.

TV television.

UHD Ultra-High-Definition.

VLC variable length codes.

WGN White Gaussian Noise.

WHT Walsh-Hadamard Transform.

Glossary

Bjontegaard-Delta rate (BD-rate) Objective quality metric, that evaluates the bitrate savings according to the obtained PSNR.

Codec Encoder-Decoder. Also referred to the method of compressing and decompressing a video sequence.

H.264/AVC Previous state of the art video codec from Joint Video Team (JVT), released in 2007. As to the writing of this work, it is the most used video compression algorithm..

Interlaced scanning Technique used by televisions for broadcasting and displaying, where only odd or even numbered lines of a frame are transmitted/displayed at a time, alternately.

JPEG Still image compression format, developed by the Joint Photographic Experts Group (JPEG).

libaom Reference software for AV1, released by Google in June 2018.

Pixel Picture Element.

Progressive scanning Technique used by more recent screens, where each frame is displayed as a whole, from top to bottom, and left to right.

RGB Color space based on the addition of Red, Green and Blue components for complex color representation.

VP8/VP9 Open-format video codecs developed by Google, released in 2008 and 2013, respectively.

Nomenclature

$\lceil \mathbf{x} \rceil$ Round to the nearest integer, greater than x .

$\lfloor \mathbf{x} \rfloor$ Round to the nearest integer, lower than x .

\mathbf{x} Round to the nearest integer.

CHAPTER 1

Introduction

1.1 | Background and Motivation

Since the spark of television research in 1887, a tremendous investment has been put into increasing the quality of images, cameras and screens that display them [1].

In the early years of mechanical television (TV), this desire was pursued by making changes to the *Nipkow* disks ¹, up to the decline of the mechanical TV, around the 1930's. The consequential rise of all-electronic TVs started with the capture of images with the same cathode tubes put into Cathode Ray Televisions (CRTs), with broadcasts of the live analog recordings, since there were no available methods of storing images, up to 1955, with the development of the open-reel magnetic tape [2].

The evolution of Complementary metal–oxide–semiconductor (CMOS) technologies however, led to the downfall of cathode ray tubes, and to the rise of image capture to a digital sensor, that allowed better image captures and lower demands in terms of physical storage space. However, with the desire for higher fidelity video, the quantity of information captured also increased. Whether by increasing the sensor resolution, color bit depth or frame rate, the captured video sequences have increased its size throughout the years. For instance, for a video of 640×360 (considered as a low resolution), at 30 frames per second (fps), considering each captured color (RGB) is represented with 8 bits, there is approximately 166 Million bits per second (Mbps) of captured information. This means that a short 5 minute video would occupy more than 6 Giga Bytes (GB) of memory. This aspect gets more severe once higher resolutions are considered. For newer standards such as 4K Ultra-High-Definition (UHD) (3840×2160) or 8K UHD (7680×4320), under the same conditions, a ten minute video would occupy 448 GB and 1792 GB of raw data, respectively.

To further aggravate the situation, video consumption got massively adopted on the average consumer level, and continues to grow, both in the average number of watched hours by users and in the resolutions of the video, making the bandwidth used on the visualization of video footage the highest between all other application. With the development of higher video sizes, increase of the average number of connected devices per user and overall market expansion through the number of consumers, this margin will continue to grow. In fact, according to *Cisco*, by 2022, up to 82% of global IP traffic will be dedicated to video [3, Trends 1 & 4].

This problem has led to the introduction of a new concept: *Video Compression* ², which is the process of reducing the size of a video sequence, while still maintaining its playback capabilities. The Codec takes advantage of redundant information present on the raw data

¹Scanning disks used in mechanical televisions

²Also called *Video Coding*.

to reduce the size of the video, without heavily modifying the original picture or its quality.

The first form of video compression, Interlaced scanning, dates from 1940, and was purely analog. This solution was introduced with the intent of reducing the necessary broadcasting bandwidth for old CRTs, without decreasing the displayed fps. And even though this technique has been implemented over more than seventy years, it has proven to be so efficient that most TV channels today still use interlaced broadcasting.

However, analog television is now obsolete, as well as CRTs. The massive developments in Integrated Circuit (IC) fabrication led to the rise of the digital era we now live in. Therefore, most screens (be it televisions, monitors or cellphones) use digital, Progressive scanning. As such, the use of analog compression techniques wasn't applicable. Accordingly, the evolution of digital video led to the development of digital compression techniques, such as the one presented in this work.

Being purely digital, these methodologies rely on computers and other processors to analyze data and apply the compression algorithms, making them very demanding processes from a computational standpoint. As expected, a high compression ratio is only obtainable by a high complexity algorithm, which also increases with the size of the video (more data leads to more analysis). Since in the early days of digital video, the used resolutions were lower as to the ones used in the present days, the compression algorithms used were not very demanding. However as the pursuit for higher quality video continued, so did the necessity for better compression ratios, and therefore the computational needs also increased. Such complex softwares lead to a high power consumption from the processor executing it, making such implementations unsuitable for portable, battery limited applications, such as cellphones or laptops. Besides this huge factor, such softwares tend to be very slow, specially when a real time compression or decompression is desired.

To amend for these factors, and to increase the reachability of high quality video to as many users as possible, these applications needed to have a viable solution that didn't compromise its usability. Accordingly, a new approach has been implemented on the most recent codec's. Besides the optimization of pure software compression/decompression solutions, there has been a great focus on the development of specialized hardware for such codecs. This solution could redress many of the problems presented previously, making them viable on a mobile implementation, as well as other specialized appliances, since such co-processors usually present a better performance than generic ones. This tendency has already been verified on the implementation choices on recent smartphones [4, p. 14], as well as recent *Nvidia* Graphical Processing Unit (GPU) lineups [5].

Since each compression algorithm tend to be very different from its predecessors, either by making changes to its bitstream or functioning principles, each time a new codec is released, there is a need to backup its development with a new set of hardware implementations. This makes the improvement of video compression techniques a continuous effort, in many engineering branches, as the technology needs to keep up with the demands of consumers, in a variety of applications.

Due to the broad access to video, and its influence in a variety of markets (besides video consumption itself), big companies have made investments on the improvement of video quality, and respective compression algorithms. These investments have provoked somewhat of a "*Codec War*". Since 2010, several video codecs have been deployed, and quickly replaced by a newer version, which presents better compression gains, at a lower quality degradation, such as the replace of *VP8* (released in 2008) with *VP9* (2013).

1.2 | Scope

AOM Video 1 (AV1) is the most recently released³ video codec. It was developed as a Joint Development Foundation [6] project, under the name of Alliance for Open Media (AOM)⁴. This codec took the same objective as its main predecessor, *VP9*, which was to be an open source, royalty free alternative to Motion Picture Experts Group (MPEG)'s state of the art video codec, *High Efficiency Video Coding (HEVC)*.



Figure 1.1: AV1 logo [7]

Upon release, *VP9* rivaled *HEVC*'s performance. However, soon after, the market demanded higher compression performance, giving origin to consortium of enterprises that now represent AOM, and to the development of *AV1*, in 2015. The first release of this coding format was made in March 2018, with the first release of its reference software, *libaom*, being made three months later, in June 2018.

Besides its main objectives, *AV1* was also developed with the intent of being implementable in hardware. Therefore, various design choices were made to make the algorithm low memory consuming, and highly parallelizable.

The desired compression performance was obtained at the cost of a highly complex algorithm (and reference software), that severely outperforms *VP9*, at the cost of much higher compression times [8].

Taken these factors, there is a high demand for dedicated hardware architectures, that can speed up the compression/decompression times and reach real-time usability on live-streaming applications, such as video-conferencing, live-content visualization, etc.

*Thesis objectives

1.3 | Outline

*General outline of the different chapters **Last section to do**

³Currently, there are other codecs being developed, without official bitstream release

⁴Further explained in Chapter 2

References

- [1] Mark Schubin. “What Sparked Video Research in 1877? The Overlooked Role of the Siemens Artificial Eye [Scanning Our Past]”. In: *Proceedings of the IEEE* 105.3 (Mar. 2017), pp. 568–576. ISSN: 0018-9219, 1558-2256.
- [2] Marco Jacobs and Jonah Probell. “A Brief History of Video Coding”. en. In: (), p. 6.
- [3] *Cisco Visual Networking Index: Forecast and Trends, 2017–2022 White Paper*. en. <https://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white-paper-c11-741490.html>.
- [4] Scientiamobile. *Mobile Overview Report April – June 2018*.
- [5] *Video Encode and Decode GPU Support Matrix*. en. <https://developer.nvidia.com/video-encode-decode-gpu-support-matrix>. Nov. 2016.
- [6] *Joint Development Foundation*. en. <http://www.jointdevelopment.org/>.
- [7] aomedia. *Home*. en-US. <https://aomedia.org/>.
- [8] Dan Grois, Tung Nguyen, and Detlev Marpe. “Performance Comparison of AV1, JEM, VP9, and HEVC Encoders”. In: *Applications of Digital Image Processing XL*. Ed. by Andrew G. Tescher. San Diego, United States: SPIE, Feb. 2018, p. 120. ISBN: 978-1-5106-1249-5 978-1-5106-1250-1.

CHAPTER 2

Video Compression Systems

2.1 | Basic Principles

Video Compression Systems have been in development for approximately forty years, with the first video codec, *H.120*, being released in 1984. It was composed of basic operations, which didn't correlate to good compression performances. This has led to a quick downfall of its usage, being aggravated by the release of the *H.261* standard by 1984.

However, the building blocks on which later standards were based are the same as in the first generations, i.e., the strategies implemented on newer standards exploit the same *redundancies* as previous, less efficient, codecs.

By redundancies, it is meant disposable information to the playback of an image sequence. This concept is the key of video compression. Throughout the years, the enhancement of video codecs was based on the improving the algorithms which can reliably represent a video, while maintaining the least of the original information. In other words, the video sequence is analyzed for predictable/identifiable characteristics, e.g. the movement of a subject or the edge of an object, calculates strategies of predicting nearby pixel values through that information and disposes the unused information. This process is mentioned as *redundancy removal*.

This way, to have a better understanding of the functioning behind video codecs, the mentioned redundancies, and respective origins, are presented. Most of such have origin on the way humans perceive vision, being this the first topic of this chapter.

2.1.1 | Human Visual System

Most of the compressed/decompressed video nowadays is directed to content visualization by consumers, with the exception of some network-driven image processing applications, such as automatic video surveillance. Therefore, the compression of video sequences has the intent of making changes to the original data, without serious impact to the users' perception. This process is mentioned as the removal of the *Psychovisual redundancy* [2]. Therefore, a basic understanding of the visual system can clarify many of the design choices made in video compression applications, and why their use doesn't present much impact on the quality of the image, while greatly reducing its size.

The image perception starts in the human eye, represented in figure 2.1. Its different constituents accomplish different tasks, from focusing, to aperture control. Although their importance to the overall functioning of the eye, the part that matters most to the focus of this work is the innermost membrane, the retina.

Once the desired image is properly focused by the lens, an inverse version of it is shined on



Figure 2.1: Representation of an human eye [3]

the aforementioned membrane, which is covered by two types of light sensitive cells, the *cones* and *rods*, which transform the observable image into a series of pulses, that get subsequently processed.

The cones are highly sensitive to color, being responsible for the *photonic* or *bright-light* vision. There are three different types of cones, corresponding to the wavelength they are susceptible to. These are the *S*, *M* and *L* cones, being sensitive to, approximately, the blue, green and red light, respectively, making a somewhat similar capture to the RGB color system.

On the other side, rods aren't stimulated by bright light, being more active on low illumination levels. This aspect makes them responsible for giving a rough overview of the field of view. This is called as *scotopic* or *dim-light* vision. These cells are spread more broadly across the retina, while to the cones, which is also observable in the number of cells (approximately 6 million cones, to 100 million rods).

From this, it's already observable that the human visual system is more sensitive to differences on the luminosity, than to the color of an object [4], which is a starting point for compressing video, as will be shown later in this chapter. However, many other opportunities come from the processing of the nerve signals, and the *psychovisual* perception that follows.

Although more sensitive to *luminance*, there is a threshold to which the difference between two objects — ΔI — can't be discerned. This relation is mentioned as *contrast sensitivity function*, which is roughly approximated with the *Weber's Law*

$$\frac{\Delta I}{I} \approx \text{constant} \quad (2.1)$$

Analyzing this equation, it's possible to conclude that the darker an object is, the lower the difference in luminance needs to be to distinguish another object. Also, darker images tend to be more susceptible to compression artifacts.

Besides the luminance values, the spatial and temporal frequencies also represent an important role in the perception of such errors.

The image 2.2 gives an example of the dependency with spatial frequency. The first image 2.2a represents the original image, which got added with White Gaussian Noise (WGN),

represented in figure 2.2b. As it is observable, these artifacts are less noticeable on the highly detailed areas (branches and leaves of the tree) than in the smooth ones (sky in the top right corner). The effect of *Weber's law* is also observed if we analyze the effect that the white noise as in the bright sun area, when compared to the darker areas.



(a) Original Image [5]



(b) Image with added WGN

Figure 2.2: Example of the effect of added noise on figure

Temporal frequency dependency, although more challenging to exemplify, is easily understandable. On a sequence of frames with fast camera, or subject movements, the human eye doesn't have the ability to track details or other artifacts, while in slow moving scenes, it can easily identify errors.

These are some of the "*flaws*" of the human visual system, that get exploited during the compression of video. However, other *redundancies*, inherent from the captured images themselves contribute to the reduction of the video size, as will be described in the following

sections.

2.1.2 | Redundancy Exploitation

Even though there are countless observable subjects and sceneries, it's unfair to think of a frame as a random sequence of pixels. Objects tend to represent clusters of pixels with roughly the same values, moving objects follow predictable directions, etc. Such characteristics represent redundancies that can be removed.

2.1.2.1 | Spatial Redundancy

Spatial redundancy comes from the similarity between neighboring pixels, on one frame. This aspect is easily verified through the autocorrelation of an image, as will be shown in the following example.

Taking image 2.2a and calculating its autocorrelation with various horizontal shifts, gives origin to the graph in image 2.3.

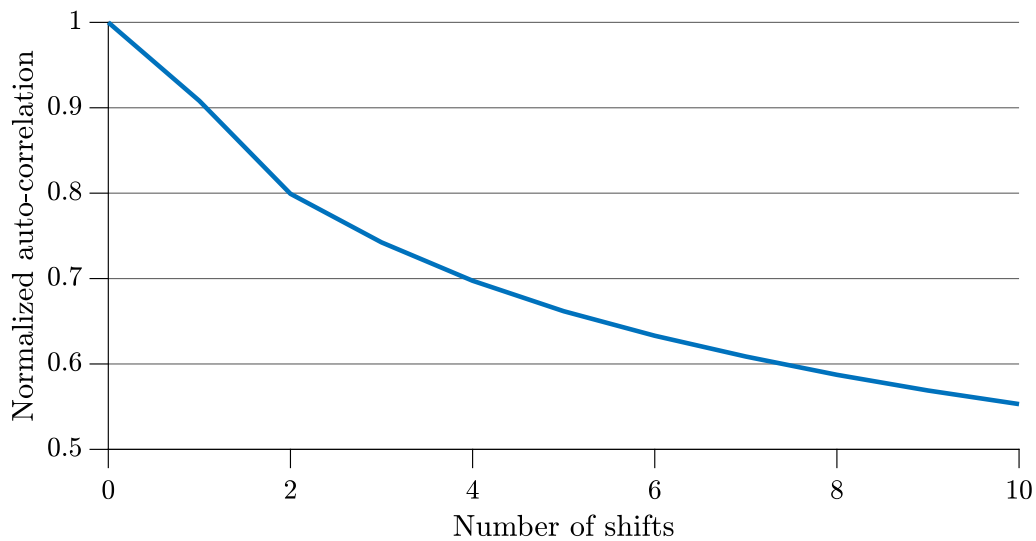


Figure 2.3: Autocorrelation of image 2.2a, with horizontal shifts

As it is observable, for shorter shifts, the normalized autocorrelation is very close to 1, since most of the de-correlation comes for the mismatching edges. Although this relation varies depending on the image, it's safe to assume that it is very similar for the majority of the cases.

Such study gives a promising opportunity for compression, since it means that most pixels can be predicted from its neighbors. This aspect as lead to what now is known as *differential* or *predictive* coding.

On a video compression system, the spatial redundancy is considered in the *intra-prediction* block, which calculates pixels, or pixel blocks, through its surrounds.

2.1.2.2 | Temporal Redundancy

As expected, a series of consecutive frames, on the same subject, tend to be very similar between each other, especially if considered the 30 or 60*fps* desired nowadays.

Making a similar analysis to what was made in section 2.1.2.1, a series of frames of the *Stefan* sequence [6] was considered, and the cross correlation between the first and the following nine was calculated, giving origin to the graph in figure 2.4.

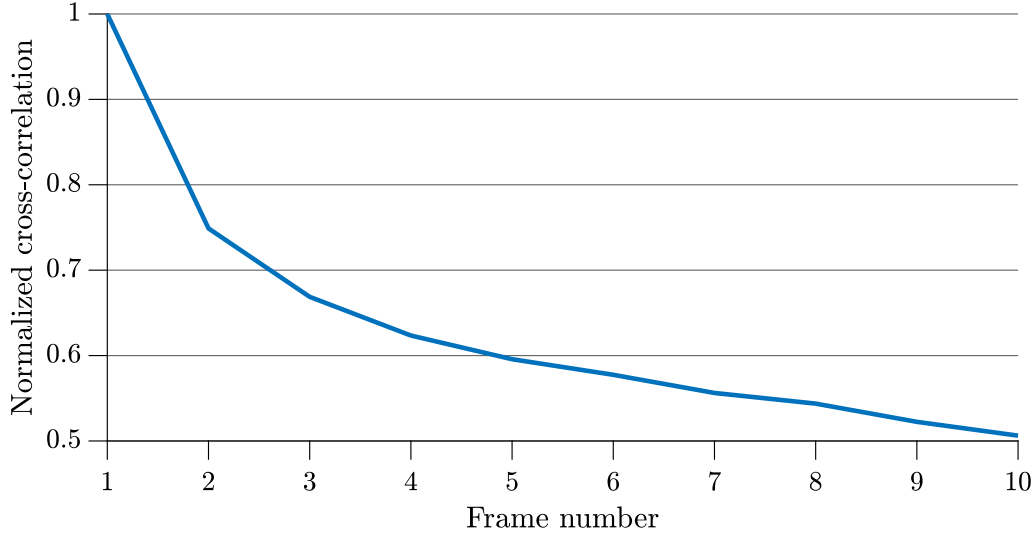


Figure 2.4: Cross-correlation between the first and following nine frames of the *Stefan* sequence

Similarly to what happened in the previous example, the cross correlation between consecutive frames is very high. Even though for faster moving scenes this relation might not be as pronounced, its application on video coding greatly contributes to the compression verified in the latest codecs.

The codec takes advantage of this redundancy in the *inter-prediction* stage, which is composed by the *Motion Estimation* (MC) and *Motion Compensation* (ME) blocks. On this stage, blocks of pixels in nearby frames are analyzed for movement, predicting its position for following frames.

2.1.2.3 | Psychovisual Redundancy

From the redundancies presented in section 2.1.1, video codecs implement compression measures in various stages.

The first measure is the *chroma subsampling*, which takes advantage of the lower perception to color, discarding some of the *chroma* samples, depending on the subsampling chosen.

Typically, a pixel value is represented in one luminance and two chrominance values, on the *YCbCr* color space. The subsampling is defined in through the relation of luminance to chroma samples, being the most common the 4:4:4, 4:2:2 and 4:2:0 standards, represented in figure 2.5. In the first one, no chroma samples are discarded, which means that for each 4 luminance (*Y*) samples, there are 4 *Cb* and 4 *Cr* samples. Correspondingly, in the second

standard, for each 4 Y samples, only 2 of each color components are maintained. The last example, although its misleading term, means that only 1 of each 4 chroma samples are kept.

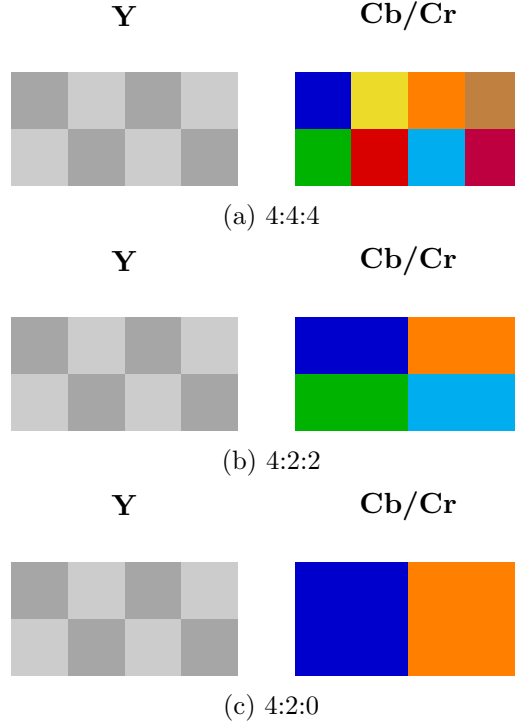


Figure 2.5: Representation of chroma subsampling

From the reduced sensitivity to details (or areas with high spatial frequency), the compression is explored in the *Transform* (T) and *Quantization* (Q) blocks. In the first stage, blocks of pixels are evaluated in their frequency components. These are then evaluated in the second stage, where the least significant ones get discarded. In the decoder, the image is reconstructed with the maintained coefficients, without much impact to the image quality. This process is further explained throughout the work.

On the Quantization block, some work was also developed to account for *Weber's law*, where the quantization depends on the average luminance value of the block. This concept was first introduced in [7], and since then experimented in various codecs, such as HEVC [8].

2.1.2.4 | Coding Redundancy

Coding redundancy is directed to the method of representing information in the digital domain, i.e., the bits themselves, and how they are organized.

It is known that symbol probability plays a major role in information compression, across a wide variety of branches, and video is no exception. Taking this into account, codecs take advantage of coding redundancy in the *Entropy Encoder* stage.

2.1.3 | Basic Video Compression/Decompression System

From the basic principles of the previously mentioned blocks, it is possible to integrate them into two complete compression — *Encoder* — and decompression — *Decoder* — mod-

ules.

2.1.3.1 | Encoder Model

The encoder's objective is to compress a video sequence, turning it into a readable *encoded bitstream*. To do this, the previously presented strategies get implemented on a system based on the schematic of figure 2.6.

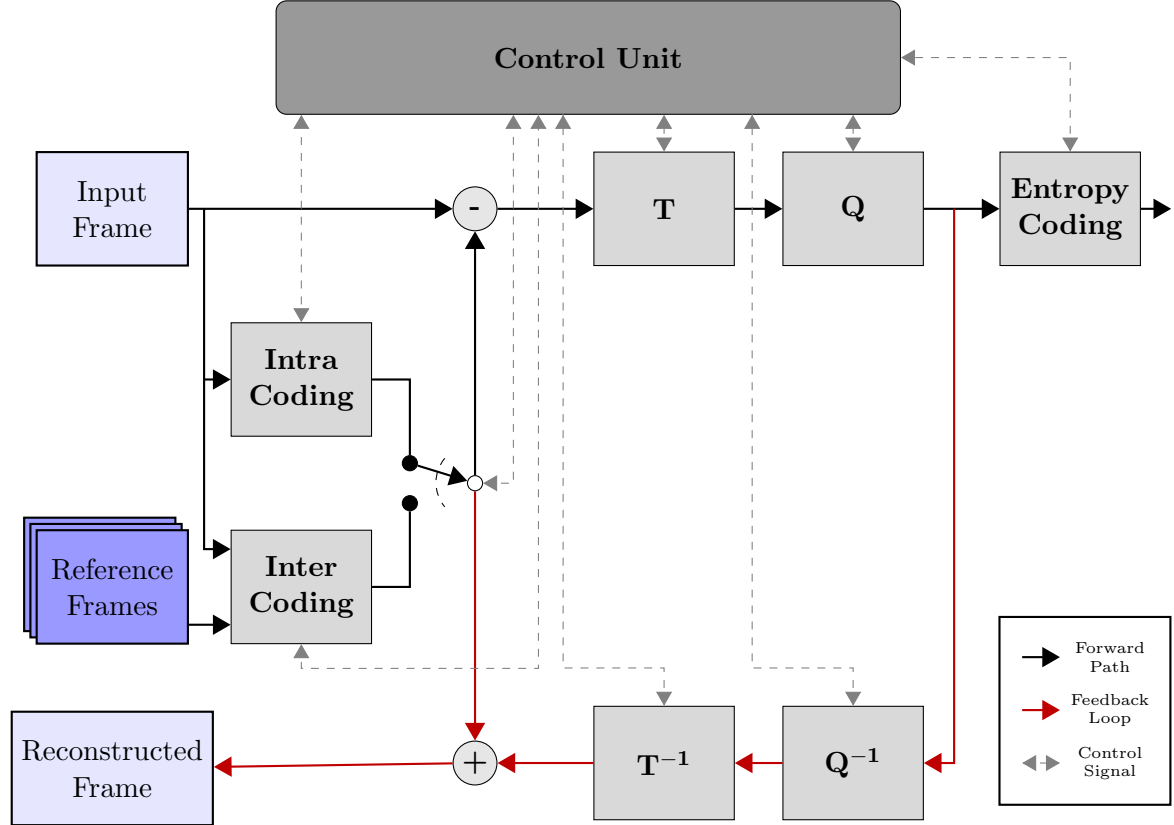


Figure 2.6: Simplified Basic Encoder Model

The encoding process starts on the *Input Frame*, which can be of two types. *I Frames* are encoded using only the information present in themselves, i.e., using only *Intra Prediction/Coding*, while *P Frames* may use predictive coding from previously encoded frames ¹.

The input gets split into blocks, which get fed into the two main blocks of a video encoder: the *Intra* and *Inter* Prediction blocks.

The *Intra Coding* block, as mentioned previously, deals with the spatial redundancy, by predicting the current block from the pixels above and to the left of its upper and left edges. The prediction may be done with various algorithms, ranging from calculating the average from the reference pixels, to replicating these according to a certain direction. One such example is presented in figure 2.7, where pixels B through H get spread across a 4×4 block, diagonally.

¹Most video codecs allow that the encoding sequence isn't the same as the temporal sequence. This allows to reference frames displayed next to the one being encoded.

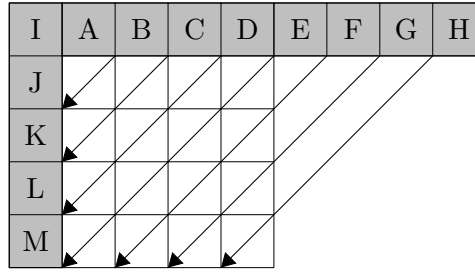


Figure 2.7: Directional Intra-prediction example

Into the *Inter Coding* block, go two inputs. The currently encoding block, as well as a bank of previously encoded frames, named *Reference Frames*. Firstly, the frames inside the buffer get searched for blocks resembling the former input. Once found, this process generates a *motion vector*, corresponding to the difference between the position of the block found in the reference frame, and the position of the currently encoding block, as shown in figure 2.8.

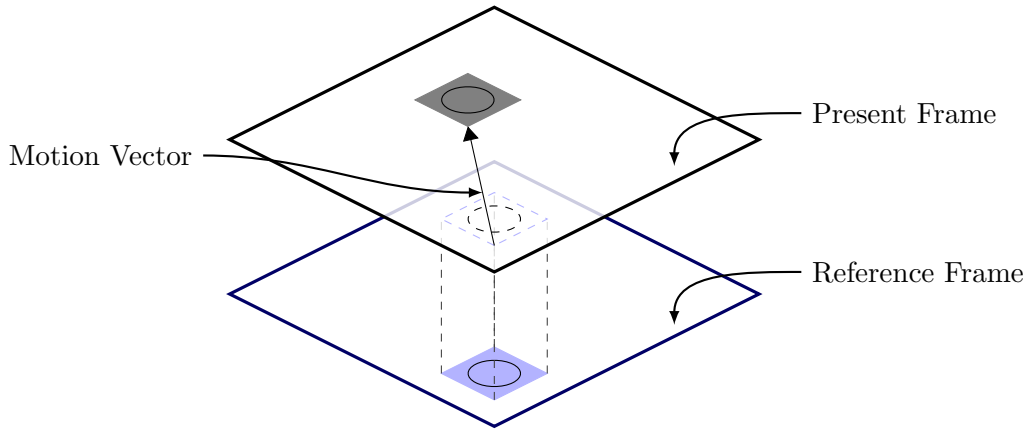


Figure 2.8: Inter-prediction example

In most codecs, the motion vector has a precision below one pixel. This means that the matching block, from the reference frame, may be interpolated from existing pixels. This process is named as *sub-pixel interpolation*, which calculates virtual values between existing pixels.

After the prediction stage, the chosen output between the two processes, i.e., the predicted block, gets subtracted by the current one, giving origin to the *residue*. This corresponds to the pixel value differences between the original and predicted blocks. Lower *residues* indicate more efficient prediction stages.

The next stage, the *Forward Transform*, is the focus of this work. It takes the residue blocks, which may not be the same size of the prediction blocks, and evaluates them according to its spatial frequencies. Its output corresponds to a series of *coefficients*, that are related to the similarity — or *correlation* — between the input block and a series of *basis images*. This process is further explained in chapter 3.

On the *Quantization* stage, the coefficients calculated in \mathbf{T} get scaled according to a *Quantization Matrix*. This stage takes advantage of the eye's lower perception to high frequency details, and scales the higher frequency coefficients by a higher value, than the lower, more

significant ones. In most of the transformed blocks, this leads that only a few low frequency components are maintained, while the others get nullified, since they are not relevant to the reconstruction of the image. Therefore, this stage is the the one that presents the higher loss, although the previously presented also introduce errors. In most of the encoding processes, this stage has the most direct impact on the obtained quality.

The wipe out of the least significant coefficients is particularly efficient when paired with the last stage before the output, the *Entropy Encoder*. On this block, \mathbf{Q} 's output block gets run sequentially via a *zig-zag scan*, which first passes through the lower frequency coefficients, followed by the higher frequency ones.

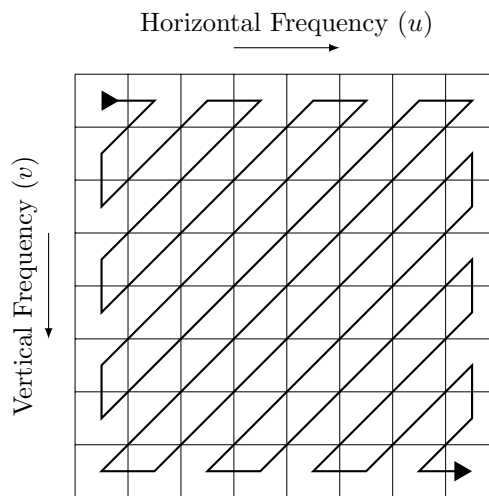


Figure 2.9: Demonstration of Zig-Zag Scan

In most of the cases, this causes that the non-zero coefficients get read first, followed by a sequence of zeros. Such sequence benefits heavily of being encoded with variable length codes (VLC), such as *Huffman Tree Codes* or Context Adaptative Binary Arithmetic Coding (CABAC). Off all the processes, this is the one that doesn't introduce further distortion into the encoded sequence, which is the reason it doesn't get included in the *feedback loop*.

The intent of this loop is to get an exact same copy of the frame reconstructed in the decoder. This reconstructed frame gets used as the reference for intra-prediction, or gets put into the reference frame buffer to be used in a later inter-prediction process.

The output of the encoder is the quantized coefficients, as well as the necessary information to recreate the encoded blocks, such as the type of prediction used, the transformation *kernel* [see p.26], quantization matrix, et al. These encoding parameters are the choices made by the *Control Unit*, which although represented by a block in figure 2.6, may not be a local process, independent from all others.

Since *H.264*, most video codecs standardize the decoding process, specifically, the allowed tools for compressing the video, and how to use them. This means that the encoding process is widely adaptable to the compression objectives, as long as the final product is a bitstream following the norms set on the codec's standard [9]. Therefore, the definition of a *Control Unit* is ambiguous in this context, since such unit can simply represent a set of parameters to be

used throughout the encoding process², or an algorithm that can change between the different capabilities of the codec, in order to achieve an objective, such as a specific distortion rate, or not surpass a maximum bit rate. As expected, different objectives may lead to majorly different results, both in the output video, as well as in the used tools.

2.1.3.2 | Decoder Model

As expected, the decoder (figure 2.10) does the backwards operation of the encoder on figure 2.6. It starts by analyzing the bitstream, separating the control information from the encoded and quantized coefficients.

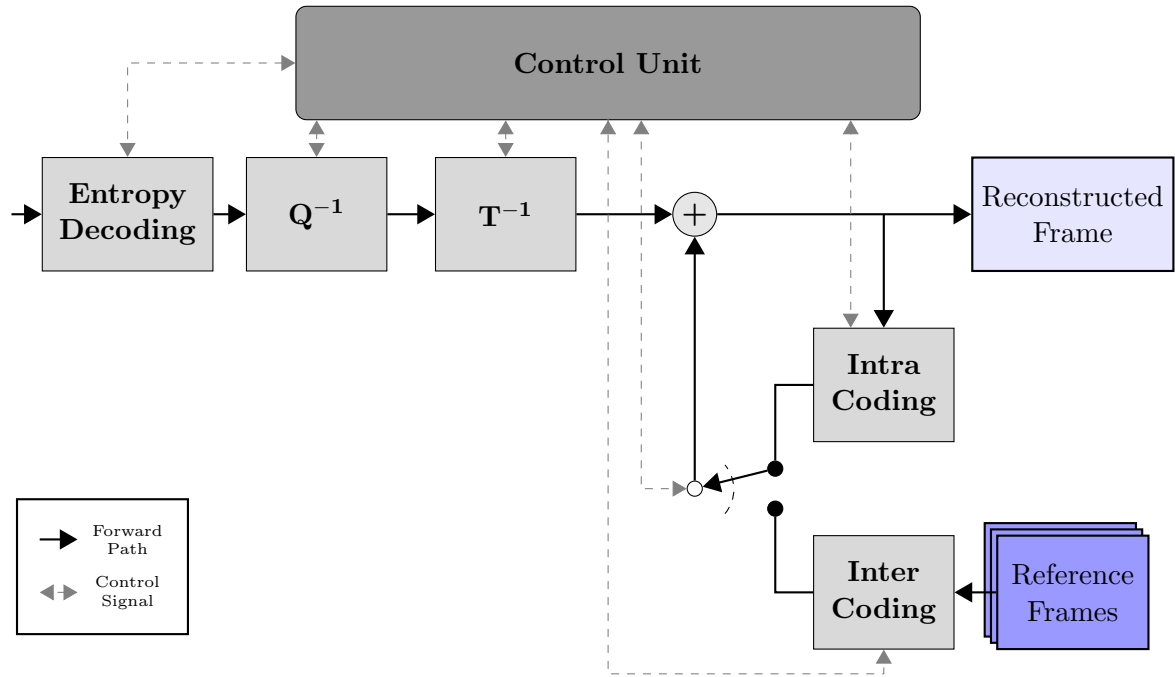


Figure 2.10: Simplified Basic Decoder Model

Having the encoding choices performed by the encoder, the decoder returns the *coding redundancy* to the quantized coefficients, on the *Entropy Decoding* stage. This corresponds to a translation from the varying length code used in codification, back into the raw coefficients.

The *Inverse Quantization* rescales the maintained coefficients, resulting from the previous *Quantization* stage. With this, it is meant that the same quantization matrix used when dividing the transformed coefficients, in the encoder, is now multiplied by the quantized parameters. It must be kept that this operation does not output an exact copy of the transformation coefficients, as a lot of information is permanently lost in \mathbf{Q} . This process can be seen in figure 2.11.

As can also be seen in this figure, the *Inverse Transform* converts the coefficients back into spatial coordinates, therefore getting the restored residue. To obtain the final approximation of the block being decoded, this residue must be added to the same predicted block from the

²One such example would be *lossless* compression modes, which use very a concise conditions on each stage, in order to get the least distortion.

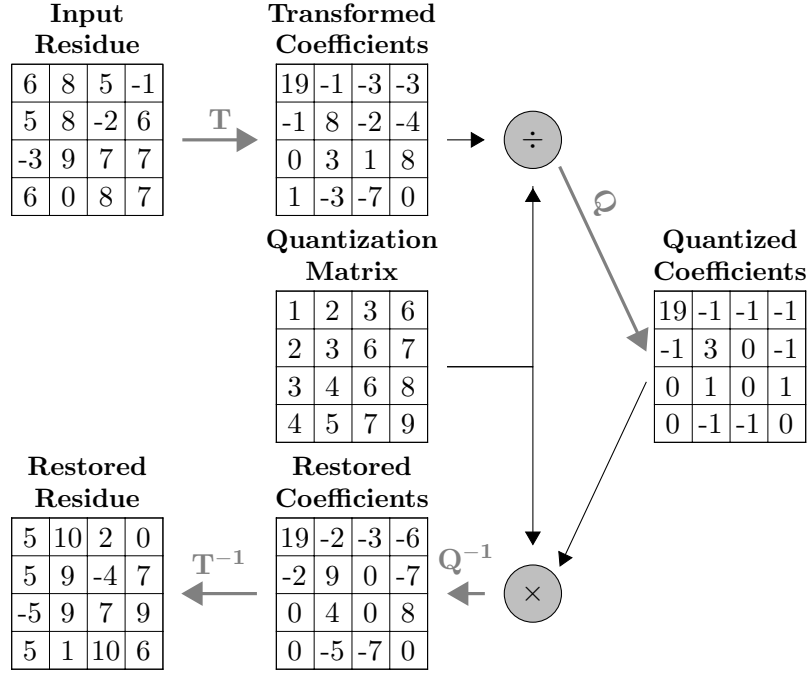


Figure 2.11: Processing of 4×4 residue block from *transformation* to restoring

encoder. To do so, the *Intra* or *Inter Prediction* stages act according to the choices made in the encoding process, as to regenerate this block.

In the decoder, the *Control Unit* represents the process that organizes the different stages, according to the choices done in the encoding stage.

2.2 | AV1

Being the focus of this work, in the following sections, *AV1* is presented on its most relevant aspects, starting with its development process.

2.2.1 | History and Development

The development of this codec started as a need to improve the bandwidth reduction of *VP9*. Therefore, the presentation of *AV1* starts by explaining the guidelines of its predecessor.

VP9 started with project *Webm*, created by *On2*, which got acquired by *Google* in 2010. This project had the objective of developing the first³ open-source, royalty-free video codec. This got support from major video content producers, such as *YouTube*, *Netflix* and *Twitch*, since it represented large savings in licensing payments, from the use of MPEG's standards, which got aggravated from the difficult patenting terms of *HEVC* [12]. After release in 2013, *VP9* got adopted as *YouTube*'s default video codec for video's above *420p*, as well as other web-video consumption services, including *Facebook*.

In 2014, *Google* started working on the next generation of open-source video codecs, *VP10*. However, due to the large interest from other companies which already used the previous

³VP8 got openly released after the acquisition of the company, after closing the development process.

standard, in 2015, the *Alliance for Open Media* was created, and the the development made for this standard got inserted into *AV1*. Alongside *Google*, twelve other companies started *AOM*, including two which also had open video encoder projects, which also majorly contributed to the fast development of *AV1*: *Cisco's Thor* and *Mozilla's Daala*. As the time of writing, 42 companies are official members of *AOM*, englobing a wide range of markets, from video streaming services, to hardware producers.



Figure 2.12: *Alliance for Open Media* current members [13]

By 2016, *AV1* started, with the objective of reaching 30% bitrate decrease, in comparison to *VP9*. After the bitstream freeze in March 2018 and deployment of *libaom* soon later, this first objective was fulfilled. However, the compression performance did not atone for the very high compression and decompression times of the reference software. This left a large margin for improvement, which quickly got explored with the development of other compression and decompression algorithms by the *AOM* members, such as *dav1d*, *rav1e*, *SVT-AV1*, among others. This parallel development gave origin to a competition among the corresponding teams, that benefited the adoption of the standard, since it brought a wide range of possibilities.

With the improvements verified on both encoders and decoders, *AV1* got progressively

more adoption from the industry, getting support from most web browsers, as well as uploads of *AV1* encoded videos to streaming platforms [14].

Besides the advances in software solutions, shortly after the bitstream freeze, IC development companies started to develop hardware solutions. The focus started by hardware decoders for implementation in mobile devices, but some encoder solutions also have been announced. Although some claims of throughput up to *8k 60fps* have been made, third party performance tests still remain to be published [15, 16, 17, 18, 19].

2.2.2 | Encoding Tools

Although the focus of this work revolves around the Transform stage, in this section, *AV1* is presented on its most relevant aspects. Some analogies are also made with *VP9*'s tools, as to justify the performance increases obtained with the most recent generation, as well as the complexity.

2.2.2.1 | Partitioning

At the start of the encoding process, an input frame is divided into *superblocks*. These constitute the starting point of the compression of an image.

These blocks may be of 128×128 pixels, or 64×64 . However, doing operations with such sizes would add complexity, as well as it wouldn't prove to be efficient. Therefore, the *superblock* can be partitioned into various *prediction blocks*. These can range between 128×128 to 4×4 , including rectangular blocks, with 2 : 1 or 4 : 1 ratios. The division of these blocks can be done recursive, where a square block divided into 4 square blocks can originate progressively smaller blocks, according to the schematic in figure 2.13.

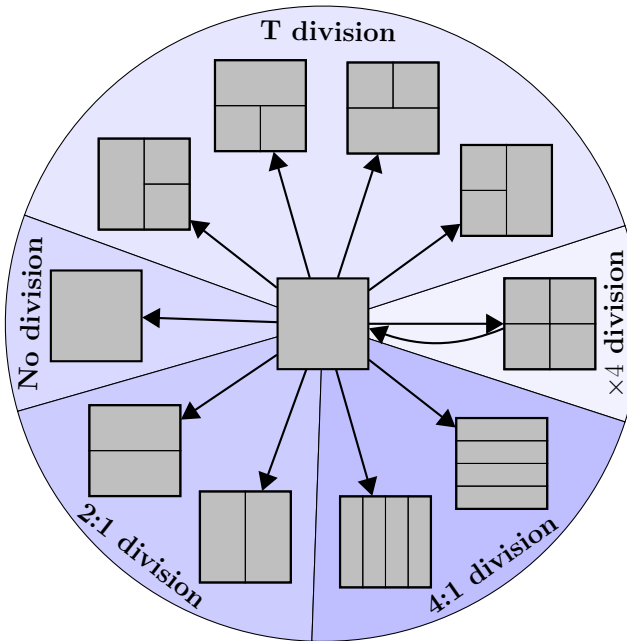


Figure 2.13: Description of the recursive partitioning scheme of *AV1*

VP9 also included a recursive partitioning scheme, but the maximum block size is 64×64 , and each block could only be divided with the $\times 4$ or $2 : 1$ ratios.

2.2.2.2 | Intra-prediction

In figure 2.7, it is presented one of the possible angles from the *directional prediction* mode of *Intra coding*. However, in *AV1*, this stage includes other prediction options, some being revised from previous generations, while others have never been implemented before.

On the *directional mode*, *AV1* improves massively from *VP9*, going from 8 directions to 56. This allows for better maintenance of details, especially on bigger blocks.

As to the *non-directional predictors*, *VP9* two. In *DC mode*, the pixels within a block would get replicated as the average of its references. *True Motion mode* would calculate each pixel as the sum of the one above by the one to the left, and subtract the upper-left diagonal, i.e., $y_{(i,j)} = y_{(i,j-1)} + y_{(i-1,j)} - y_{(i-1,j-1)}$. *AV1*'s *Smooth modes* are similar to the previous *DC*, but it has the possibility of calculating the weighted average of the reference pixels, as well as using just one set of references, horizontal or vertical. *TM mode* gave place to *Paeth*, which makes various calculations similar to *TM*, then considering the most fitted prediction. An hardware architecture for this intra-predictor has already been implemented in [21].

Pallet mode also got revised and included in *AV1*. This mode is paired with other prediction techniques, limiting the pixel values to a set of possible colors. *Pallet* as well as *Intra-block copy* are especially designed for artificial video, such as video game footage, since these kind of videos contained a limited set of colors textures. *Intra-block copy* allows for the replication of a intra-predicted block, similarly to the process in inter prediction.

Finally, *AV1* introduces two new intra prediction modes that haven't been implemented in previous generations. These are *Chroma from Luma* and *Recursive-filter Intra Prediction*. The first is easily understandable through its name. The chroma component of a block is calculated through the corresponding luminance values (see [22]). As to the later, it subdivides a prediction block, and calculates each set of pixels using different filters.

2.2.2.3 | Inter-prediction

This block got major innovations, as well as improvements to previous generations. Regarding the standard techniques, *AV1* improves in the number of motion vector estimation filters, going from two to four, as well as in the number of sub-pixel filters. While *VP9* allowed for three reference frames, the newer inter predictor allows to choose up to seven per frame, in a set of eight reference frames. This highly increases the necessary memory for encoding and decoding, but allows for finer motion estimation.

As to innovations, *AV1* introduces *Warped motion*, which allows to shape the reference block on a trapezoidal manner, *Global motion*, to easily shift an entire frame, as to deal with camera movements, *Wedge mode* which allows to use different prediction schemes in the same block, among others.

Some works have already been published with advances to this stage, such as [23], as well as hardware implementations [24].

2.2.2.4 | Transform

AV1 follows the innovations made in *VP9*, adding more transformation kernels. Besides the regularly implemented Discrete Cosine Transform (DCT), the transformation blocks may

now be transformed using Identity kernels or Asymmetric Discrete Sine Transform (ADST) kernels, which can be implemented in two directions. These different options can be used independently in the columns and rows, giving origin to 16 different options of block transformations. This aspect is further explained in chapter 3.

As to transform sizes, *AV1* allows for extra flexibility, not fixing any of the block's dimensions to a certain value. This way, the block size can vary between 4×4 and 64×64 , including rectangular blocks of 2 : 1 and 4 : 1 ratios.

2.2.2.5 | Quantization

Although the simplest stage from the encoding/decoding process, *AV1* developed this stage by allowing a wider set of quantization matrixes to be used within the same frame, as well as updating the choosing criteria. While in *VP9* the *Quantization Parameter (QP)* would be calculated considering the chroma components as one, now both channels (Cb and Cr) are considered independently.

Since *AV1* was targeted at web applications, one other innovation was added to this stage, which is an offset to the quantization matrixes. This is particularly effective on applications where a specific target bitrate is to be achieved.

2.2.2.6 | In-loop Filtering

Although not represented in figures 2.6 and 2.10, recent codecs include some kind of filtering to reduce compression artifacts. In *VP9* there was included a *Deblocking Filter*, which filtered the entire image, as to reduce the edging artifacts from prediction. *AV1* maintains this filter, reducing the necessary memory to implement it.

Besides the revision of the old filter, many others are added, such as the *Constrained Directional Enhancement Filter*, that filters the image directly on the prediction blocks' edges, with the same objective of the *Deblocking filter*. Some further explanation of these filters may be found in [25] and [26].

2.2.3 | Performance Analysis

AV1's decoding specification hasn't change since the release of the standard and freeze of the bitstream. However, this isn't verified on the implementations of the standard. Even *libaom*, which is intended to serve as a guideline for future implementations, has been severely improved since its release in June 2018.

Being so, the comparison of *AV1* throughout these developing months has been divided in two major categories: *Quality* and *Timing*. The first depends on the standard itself, and on how the encoding tools are able to compress the video, while maintaining its playback capabilities. Therefore, if the encoding objectives are maintained throughout the development of the encoders/decoders, this parameter should not vary. However, the same cannot be said of the *Timing Performance*, since as more efficient tools get released, it is expected that the time to encode/decode a video gets reduced, as to reach real-time usability.

According to Moscow State University [27], *AV1* achieved its objective of highly reducing the necessary bitrate. On this test, five *1080p* sequences have been encoded using different implementations of *H.264/AVC*, *H.265/HEVC*, *VP9* and *AV1*. The different softwares have been configured on a similar manner, as to encode the sequences with similar quality, and

the average bitrate per codec was compared relatively to *H.264*. The results are presented in figure 2.14.

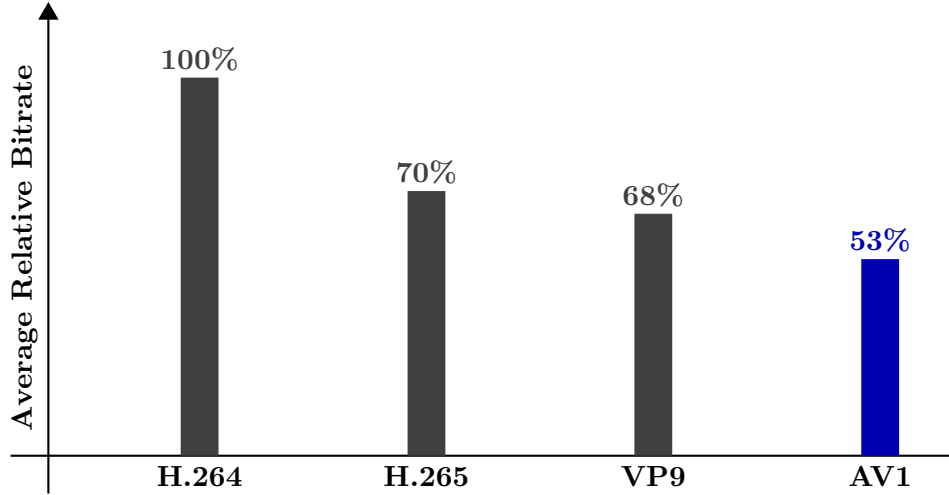


Figure 2.14: *AV1* bitrate savings [27]

These results may vary greatly with the performed tests, as the encoding tools may prove to be more adequate to certain types of videos. On a different test, performed by *Facebook* [28], *AV1* presents a higher performance than the one presented previously, as seen in table 2.1. Here, videos of various resolutions were encoded with *VP9* and *H.264* with equivalent parameters, and the obtained bitstreams are compared to *AV1*'s, according to Bjontegaard-Delta rate (BD-rate).

Codec	Resolutions				Average
	360p	480p	720p	1080p	
VP9	-29.5%	-32.5%	-32.3%	-35.9%	-32.5%
H.264	-43.4%	-49.3%	-51.2%	-57.9%	-50.3%

Table 2.1: BD-rate of *VP9* and *H.264* codecs, when compared to *AV1* (negative corresponds to bitrate savings) [28]

As it can be seen, as the resolution increases, so does the bitrate savings. This leads to believe that if the same test were to be performed with *4K* and *8K* sequences, higher performances would be verified.

As to the encoding times, in two articles from *Streaming Media* [29, 30], it is possible to see the improvements made on *libaom*. In table 2.2 there are presented the encoding times of a 5 second clip, shortly after the reference software was released, August 2018, and in March 2019, under the same conditions. Besides the performance of *libaom*, software for *H.264/AVC*, *H.265/HEVC* and *VP9* is also evaluated.

From these results, it is possible to conclude that *AV1* is a promising codec. When quality and compression gains are considered, it is already verifiable that the codec presents better performances than its predecessors, in some cases even beating its objective of 30% increase over *VP9*. However, when considering the timing issues, the results don't prove as optimistic.

Codec	Encoding Time (s)	
	2018	2019
AV1	226 080	736
H.265		289
VP9		226
H.264		18

Table 2.2: Encoding times of different video encoders, and improvements on *AV1* [29, 30]

As the time of writing, the encoding solutions are still far away from a real-time usability. Although, as better software and hardware solutions get developed, this objective may be achieved in the near future.

References

- [1] Rafael C. Gonzalez and Richard E. Woods. *Digital Image Processing*. eng. Fourth edition. OCLC: 987436552. New York, NY: Pearson, 2018. ISBN: 978-0-13-335672-4.
- [2] Yun Qing Shi and Huifang Sun. *Image and Video Compression for Multimedia Engineering: Fundamentals, Algorithms, and Standards*. eng. 2. ed. Image Processing Series. OCLC: 254564955. Boca Raton, Fla.: CRC Press, 2008. ISBN: 978-0-8493-7364-0.
- [3] *Anatomy of the Eye: Human Eye Anatomy*. en. <https://owlcation.com/stem/Anatomy-of-the-Eye-Human-Eye-Anatomy>.
- [4] Kathy Mullen. “The Contrast Sensitivity of Human Color Vision to Red-Green and Blue-Yellow Chromatic Gratings”. In: *The Journal of physiology* 359 (Mar. 1985), pp. 381–400.
- [5] *Freepik*. en. <https://www.freepik.com>.
- [6] *YUV Sequences*. <http://trace.eas.asu.edu/yuv/>.
- [7] Andrew Watson. “Efficiency of a Model Human Image Code”. In: *Journal of the Optical Society of America. A, Optics and image science* 4 (Jan. 1988), pp. 2401–17.
- [8] Kais Rouis and Chaker Larabi. “Perceptual Video Content Analysis and Application to HEVC Quantization Refinement”. In: Nov. 2018, pp. 1–6.
- [9] “AV1 Bitstream & Decoding Process Specification”. en. In: (), p. 681.
- [10] Luciano Agostini. *Desenvolvimento de Arquiteturas de Alto Desempenho Dedicadas à Compressão de Vídeo Segundo o Padrão H.264/AVC*. 2007.
- [11] Debargha Mukherjee. *AllThingsRTC 2019 - Opening Keynote - Past, Present and Future of AV1*. June 2019.
- [12] Streaming Media. *HEVC Advance Patent Pool Creates Confusion, Lacks Transparency*. en-US. July 2015.
- [13] aomedia. *Home*. en-US. <https://aomedia.org/>.
- [14] Nathan Egge. “Latest Technical and Business Progress with AV1”. en. In: *NAB Streaming Summit* (Apr. 2019), p. 59.
- [15] *Allegro DVT Introduces the Industry First Real-Time AV1 Video Encoder Hardware IP for 4K/UHD Video Encoding Applications*. en-US. <http://www.allegrodvt.com/allegro-dvt-introduces-the-industry-first-real-time-av1-video-encoder-hardware-ip-for-4kuhd-video-encoding-applications/>. Apr. 2019.

- [16] *NGCodec Announces AV1 Support and a 2X Performance Improvement in Broadcast Quality Live Video Encoding.* en-US. <https://ngcodec.com/press-releases/2019-1-7-ngcodec-announces-av1-support-and-a-2x-performance-improvement-in-broadcast-quality-live-video-encoding>.
- [17] Anton Shilov. *Realtek Demonstrates RTD2893: A Platform for 8K Ultra HD TVs.* <https://www.anandtech.com/show/14560/realtek-demonstrates-rtd2893-a-platform-for-8k-ultrahd-tvs>.
- [18] *Realtek Launches Worldwide First 4K UHD Set-Top Box SoC (RTD1311), Integrating AV1 Video Decoder and Multiple CAS Functions - REALTEK.* <https://www.realtek.com/en/press-room/news-releases/item/realtek-launches-worldwide-first-4k-uhd-set-top-box-soc-rtd1311-integrating-av1-video-decoder-and-multiple-cas-functions>.
- [19] *Socionext Implements AV1 Encoder on FPGA over Cloud Service.* en-US. <http://socionextus.com/pressreleases/socionext-implements-av1-encoder-over-cloud-service/>. June 2018.
- [20] Yue Chen et al. “An Overview of Core Coding Tools in the AV1 Video Codec”. In: *2018 Picture Coding Symposium (PCS)*. San Francisco, CA: IEEE, June 2018, pp. 41–45. ISBN: 978-1-5386-4160-6.
- [21] M. Corrêa et al. “A High Throughput Hardware Architecture Targeting the AV1 Paeth Intra Predictor”. In: *2019 IEEE 10th Latin American Symposium on Circuits Systems (LASCAS)*. Feb. 2019, pp. 93–96.
- [22] L. Trudeau, N. Egge, and D. Barr. “Predicting Chroma from Luma in AV1”. In: *2018 Data Compression Conference*. Mar. 2018, pp. 374–382.
- [23] Z. Deng and I. Moccagatta. “Hardware-Friendly Inter Prediction Techniques for AV1 Video Coding”. In: *2017 IEEE International Conference on Image Processing (ICIP)*. Sept. 2017, pp. 948–952.
- [24] R. Domanski et al. “High-Throughput Multifilter Interpolation Architecture for AV1 Motion Compensation”. In: *IEEE Transactions on Circuits and Systems II: Express Briefs* 66.5 (May 2019), pp. 883–887. ISSN: 1549-7747.
- [25] A. Norkin and N. Birkbeck. “Film Grain Synthesis for AV1 Video Codec”. In: *2018 Data Compression Conference*. Mar. 2018, pp. 3–12.
- [26] S. Midtskogen and J. Valin. “The Av1 Constrained Directional Enhancement Filter (Cdef)”. In: *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Apr. 2018, pp. 1193–1197.
- [27] Dr Dmitriy Vatolin et al. “MSU Codec Comparison 2018 Part IV: FullHD Content, High Quality Use Case”. en. In: (2019), p. 42.

- [28] *AV1 Beats X264 and Libvpx-Vp9 in Practical Use Case.* en-US. Apr. 2018.
- [29] *AV1: A First Look.* en-US. Aug. 2018.
- [30] *Good News: AV1 Encoding Times Drop to Near-Reasonable Levels.* en-US. <https://www.streamingmedia.com/Articles/ReadArticle.aspx?ArticleID=130284>. Mar. 2019.

CHAPTER 3

Video Coding Transforms

3.1 | Introduction

As mentioned previously, the basic principle behind the compression of video, is the reduction of inter-pixel/inter-symbol correlation. The various integral blocks of a video compression system try to accomplish this objective through different strategies. The *Intra-frame* and *Inter-frame Prediction* exploit spatial and temporal correlation, respectively. Through the subtraction of the input by the output of one of these blocks, and the attainment of the *residue*, the next compression stage is made in the *Transform* block, which is the focus of this work.

The technique implemented by this process relies on the energy compaction in the frequency domain to reduce the correlation within a frame block, i.e. the input of the Transform block is evaluated on its main frequencies — the *transform coefficients* — on a spatial domain, similarly to the process executed on a *Fourier Transform*. Once each block is quantized on these coefficients, the compression is made with the removal of the least significant ones, on the *Quantization* stage. The intent of the *transform* is to split the image into a set of predefined coefficients, that get transmitted in place of the *residue*.

The objective of this chapter is to give the reader a basic understanding of the theoretical basis behind said Transformations, as well as to introduce the most commonly used ones.

3.2 | Background

3.2.1 | Basis vector/image interpretation

A useful interpretation, and a good starting point to the study of this process, is to see it as the decomposition of the input as a set of basis vectors (1D transforms) or images/matrices (2D transforms). The transformation outputs, T_i , can be seen as the weights of each basis vector/image, \vec{e}_i , that summed return the input, \vec{g} , i.e.

$$\vec{g} = \sum_{i=1}^N T_i \vec{e}_i \quad (3.1)$$

which means that the coefficients are related to the amount of correlation between the input and each basis component, and can be obtained with the *inner product* of the input and each basis vector.

$$T_i = \vec{e}_i^T \vec{g} \quad (3.2)$$

Since each input vector will have different correlation values between the various basis vectors, this operation accomplishes two main objectives:

- De-correlation of the input values
- Signaling of the most important basis vectors.

Considering a 2D image, $g(x, y)$, and its corresponding transformed coefficients, $T(u, v)$, where (x, y) are the pixel coordinates, and (u, v) are the corresponding coordinates in the transform domain, we can obtain an analogous version of equation 3.2 as

$$T(u, v) = \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} g(x, y) f(x, y, u, v) \quad (3.3)$$

Similarly, we can re-obtain the original picture

$$g(x, y) = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} T(u, v) i(x, y, u, v) \quad (3.4)$$

where $f(x, y, u, v)$ and $i(x, y, u, v)$ are the *forward* and *inverse transformation kernels*. To better explain the concept of these, first it's needed to introduce the two following concepts.

3.2.1.1 | Separability

A useful characteristic of 2D Video Coding Transforms is its ability to be independently calculated between rows and columns. This means that given a 2D block as input, the transform coefficients can be calculated first with the *horizontal transform*, and then with the *vertical transform*, or vice-versa.

This aspect is applicable if the following conditions are applied

$$f(x, y, u, v) = f_1(x, u) f_2(y, v) \quad (3.5)$$

$$i(x, y, u, v) = i_1(x, u) i_2(y, v) \quad (3.6)$$

This means that the equation 3.3 is reconstructed as 2 independent and sequential operations

$$T_{temp}(x, v) = \sum_{y=0}^{N-1} g(x, y) f_2(y, v) \quad (3.7)$$

$$T(u, v) = \sum_{x=0}^{M-1} T_{temp}(x, v) f_1(x, u) \quad (3.8)$$

On AV1, due to the various implemented transformation kernels, this aspect is severely explored, since the only way of implementing the combination of different 1D kernels, is to calculate them independently. This aspect is further explained with the following concept.

3.2.1.2 | Symmetry

Taking equation 3.5, a transformation kernel is said to be symmetric if

$$f_1(y, v) = f_2(x, u) \quad (3.9)$$

This characteristic is particularly useful because it makes the forward and inverse transformations expressible as matrix multiplications. Therefore, the equations 3.3 and 3.4 are represented, respectively, as

$$T = F^T G F \quad (3.10)$$

$$G = I^T T I \quad (3.11)$$

where F and I are the forward and inverse transform matrices. This aspect is only possible for square matrix, i.e., input blocks with the same height and width.

This concept isn't exploited in AV1, since the use of different 1D transformation kernels, and rectangular block sizes ($M \neq N$) make the 2D transform asymmetric, and therefore, not executable as matrix multiplication. Consequently, the block transformation is made as 2 separate 1D operations, as shown previously.

Looking now at equation 3.4, we can interpret the inverse transformation kernel as a set of basis images, dependent of the (u, v) pair. By this, it is meant

$$g(x, y) = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} T(u, v) I_{u,v} \quad (3.12)$$

where

$$I_{u,v} = \begin{bmatrix} i(0, 0, u, v) & i(0, 1, u, v) & \dots & i(0, M-1, u, v) \\ i(1, 0, u, v) & i(1, 1, u, v) & \dots & i(1, M-1, u, v) \\ \vdots & \vdots & \dots & \vdots \\ i(N-1, 0, u, v) & i(N-1, 1, u, v) & \dots & i(N-1, M-1, u, v) \end{bmatrix} \quad (3.13)$$

Therefore, the forward and inverse transformation process can be seen as the deconstruction of an input block, into a set of $M \cdot N$ basis images, dependent of the used transformation kernel. As expressed in equations 3.5 and 3.6, this analogy can be made on a 1D space .

Given a general comprehension of the theoretical principles behind the *Transform* block, now the most common transformation kernels are introduced, with focus on the AV1 video codec.

3.3 | Transformation Kernels

3.3.1 | Discrete Fourier Transform (DFT)

Although it isn't implemented in video coding, it's widely used in digital signal processing, and many of the used transformation kernels are approximations of this function.

It has its roots on the *Fourier Transform*, whose forward and inverse transformations are expressed in equations 3.14 and 3.15, respectively.

$$T(u, v) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) e^{-j2\pi(ux+vy)} dx dy \quad (3.14)$$

$$g(x, y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} T(u, v) e^{j2\pi(ux+vy)} du dv \quad (3.15)$$

Once considered a finite number of points, the previous equations become

$$T(u, v) = \frac{1}{MN} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} g(x, y) e^{-j2\pi(\frac{ux}{M} + \frac{vy}{N})} \quad (3.16)$$

$$g(x, y) = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} T(u, v) e^{j2\pi(\frac{ux}{M} + \frac{vy}{N})} \quad (3.17)$$

which corresponds to replacing the kernels in equations 3.3 and 3.4 with

$$f(x, y, u, v) = \frac{1}{MN} e^{-j2\pi(\frac{ux}{M} + \frac{vy}{N})} \quad (3.18)$$

$$i(x, y, u, v) = e^{j2\pi(\frac{ux}{M} + \frac{vy}{N})} \quad (3.19)$$

The position of the multiplication factor, $\frac{1}{MN}$, is irrelevant, and in some works is divided into two terms in the forward and inverse kernels, $\frac{1}{M}$ and $\frac{1}{N}$, or even $\frac{1}{\sqrt{MN}}$.

Because of the use of complex numbers, this operation tends require a high computational effort, whence its disuse in video coding.

3.3.2 | Discrete Walsh-Hadamard Transform (WHT)

This transformation replaces the sum of sines and cosines of the DFT, alternating of positive and negative 1's, depending on the binary representation of the inputs.

Considering the inputs of the transform to be represented with m bits, where $m-1$ is the most significant bit (b_{m-1}), the forward and inverse kernels are represented as

$$f(x, y, u, v) = i(x, y, u, v) = \frac{1}{\sqrt{MN}} (-1)^{\sum_{i=0}^{m-1} [b_i(x)p_i(u) + b_i(y)p_i(v)]} \quad (3.20)$$

where

$$\begin{aligned} p_0(u) &= b_{m-1}(u) \\ p_1(u) &= b_{m-1}(u) + b_{m-2}(u) \\ &\vdots \\ p_{m-1}(u) &= b_1(u) + b_0(u) \end{aligned} \quad (3.21)$$

3.3.3 | Discrete Cosine Transform (DCT)

The most commonly used transform, the *DCT*, was published by Ahmed et al. in 1974 [3]. Since then, it has been adopted in a wide range of applications, being the only transform used in the first generations of video codecs, as well as in *still image compression*, being the basis of the *JPEG* standard.

It is frequently compared to the *DFT*, due to the similarity of their operation. However, as the name implies, the *DCT* relies on the cosine function to create its basis images, which is a *periodic* and *symmetrically even* function. Therefore, as mentioned by [4, A. V. Oppenheim], "*Just as the DFT involves an implicit assumption of periodicity, the DCT involves implicit assumptions of both periodicity and even symmetry*". This is easily observable once considered the equivalent process of both algorithms. Taking an N -point sequence, $g(n)$, the calculation of the *DFT* and *DCT* of such sequence is equivalent to the processes presented at table 3.1.

Step	<i>DFT</i>	<i>DCT</i>
1	Repeat $g(n)$ every N points, giving origin to $\tilde{g}_N(n)$	Concatenate $g(n)$ with a flipped version of itself, creating a $2N$ sequence, $g_{2N}(n)$, and repeat it, giving origin to $\tilde{g}_{2N}(n)$
2	Calculate the <i>Fourier</i> expansion of \tilde{g}_N	Calculate the <i>Fourier</i> expansion of \tilde{g}_{2N}
3	Keep the first N coefficients and set all others to 0	Keep the first N coefficients, and set all others to 0

Table 3.1: Similarity between the processes of the *DFT* and the *DCT*

The main reason behind the heavy adoption of the *DCT* is its great energy compaction on the lower frequencies, where most of the energy in a picture is packed. If the output of the first step of table 3.1 is observed, this aspect is more easily understood. In figure 3.1, a 4 point sequence, corresponding to the filled points, gets replicated throughout the discrete time axis, according to the corresponding transform.

Due to the back-to-head repetition seen in figure 3.1a, there is a disruption every N points, which gives origin to high frequency components in the *Fourier* transform. Therefore, the more continuous behavior obtained with the back-to-back repetition of the *DCT* gives origin to more significant low frequency coefficients. However, there are many ways of creating a periodic and symmetric sequence from a finite number of points. This factor has led to the implementation of different versions of the *DCT*, which differ in minor details between themselves. These differences are consequence of the way the symmetry is obtained, which can be observed in figures 3.1b to 3.1e. The represented implementations are referred to as *DCT-I* to *DCT-IV*, but other possibilities exist. Their definition depends on the overlapping of points when repeating each sequence.

Since the *DCT* in *AV1* is implemented in one dimension, the description of the following kernels is also made in 1D. Therefore, the dimension of the transform, L , is referring either to the blocks' width or height, depending if the operation is made to the rows or columns, respectively (M or N , previously).

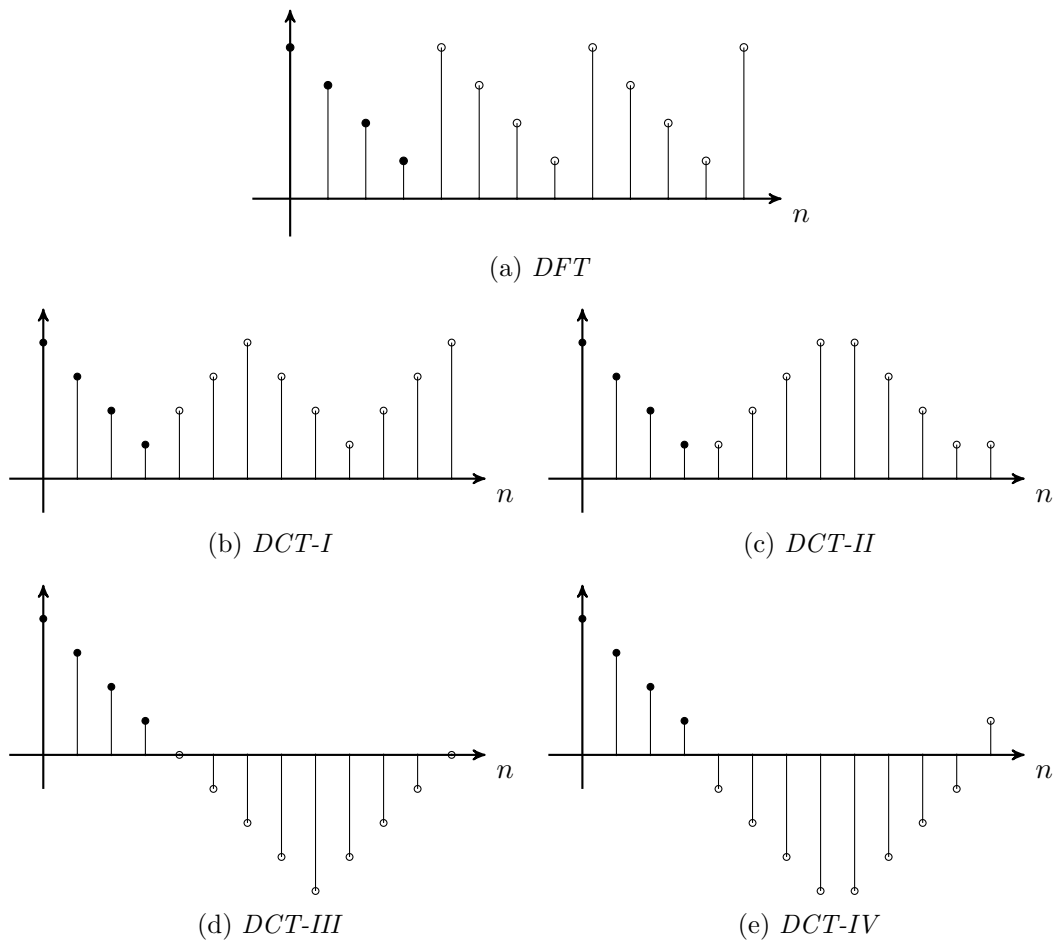


Figure 3.1: Sequences generated in the first step of table 3.1 for the DFT and different DCTs. Filled dots correspond to the original sequence.

DCT-I The sequence created with first version of the DCT has overlapping points at $n = k(L - 1)$, $k = 0, 1, 2, \dots$, making the overall period of the final sequence $2L - 2$.

$$f(x, u) = \frac{2}{L - 1} \alpha(x) \cos\left(\frac{\pi x u}{L - 1}\right) \quad (3.22)$$

where

$$\alpha(x) = \begin{cases} \frac{1}{2}, & x = 0 \vee x = N - 1 \\ 1, & 1 \leq x \leq N - 2 \end{cases} \quad (3.23)$$

The inverse transform becomes

$$i(x, u) = \alpha(u) \cos\left(\frac{x u \pi}{L - 1}\right) \quad (3.24)$$

DCT-II Usually referred to as "the *DCT*", it is by far the most implemented version, being the only one mentioned in many of the studied works.

As seen in figure 3.1c, this version has no overlap on the created sequence, making the period $2L$, and the points of symmetry $kL - \frac{1}{2}$.

$$f(x, u) = i(x, u) = \beta(u) \cos\left(\frac{(2x + 1)u\pi}{2L}\right) \quad (3.25)$$

$$\beta(u) = \begin{cases} \sqrt{\frac{1}{L}}, & u = 0 \\ \sqrt{\frac{2}{L}}, & 1 \leq u \leq N - 1 \end{cases} \quad (3.26)$$

DCT-III Named the *inverse* of DCT-II, due to the switch of the transform and pixel coordinates.

$$f(x, u) = i(x, u) = \beta(u) \cos\left(\frac{(2u + 1)x\pi}{2L}\right) \quad (3.27)$$

$$\beta(u) = \begin{cases} \sqrt{\frac{1}{L}}, & u = 0 \\ \sqrt{\frac{2}{L}}, & 1 \leq u \leq N - 1 \end{cases} \quad (3.28)$$

DCT-IV Is the basis of the *Modified Discrete Cosine Function (MDCT)*, where the input blocks overlap.

$$f(x, u) = i(x, u) = \sqrt{\frac{2}{L}} \cos\left(\frac{(2u + 1)(2x + 1)\pi}{4L}\right) \quad (3.29)$$

3.3.4 | Discrete Sine Transform (DST)

Similarly to the DCT, there is also the possibility to represent a finite sequence as a sum of discrete *sine* functions, giving origin to the *DST*. Contrarily to the former presented transform, this variant uses sinusoidal functions to generate its basis images, which gives origin to *odd symmetric* sequences.

In the same way as its *even* counterpart, there are various different ways off accomplishing such symmetry, which also gives origin to eight different variations of this Transform. However, due to its misuse over the DCT, only the *DST-II* is presented.

$$f(x, u) = i(x, u) = \sqrt{\frac{2}{L+1}} \sin\left(\frac{(j+1)(u+1)\pi}{L+1}\right) \quad (3.30)$$

Equivalently to what happens with the DFT, the odd symmetry of this function gives origin to discontinuities, which are undesirable when coding video blocks, since they lead to less significant low frequency coefficients, and therefore higher quantization errors.

3.3.5 | Asymmetric Discrete Sine Transform (ADST)

The symmetric behavior of previous transforms lead to better performance on evenly spread residue blocks, i.e. when the pixel values post-subtraction (and before transformation) have roughly the same value across the whole block.

However, due to the directional spatial prediction, the residue on one boundary of the block may be different from the other boundaries, since the chosen direction for prediction may diverge from the original block, across the prediction block. This leads to worse energy compression.

In order to combat this aspect, VP9 introduced a new transform called *Asymmetric Discrete Sine Transform (ADST)*, which corresponds to an alternative implementation of the DST with the addition of frequency and phase shifts.

This enhancement provides the developer with a high degree of liberty, since the basis images can be adapted with the variation of the shifts. On AV1, there is only one ADST implementation per block size. However this transformation can be done in two directions, i.e., the input vector can be transformed front-to-back and vice-versa. *AOMedia* named these transforms *ADST* and *Flip-ADST*, according to the direction of the input vector.

3.4 | Libaom's Integer Transformations

In battery driven applications, computing power plays an important role. Consequently, any approach that leads to lower computational costs, without compromising the image quality, tends to get incorporated into a video codec.

When considering the transformation stage, a widely used approach is the use of *integer transforms*. The objective of such operations is to maintain the features of floating point transforms, but severely reducing the complexity, decreasing the used operations to arithmetic additions and integer multiplications. In many implementations, the latter are implemented with bitwise shifts and additions.

From the transforms presented throughout section 3.3, there have been several methods of developing integer counterparts. Most of the fast implementations are based in either *Fast Fourier Transform* algorithms or in the *Walsh-Hadamard Transform* [10, 11]. Since the objective of this work was to develop a *Transform Co-processor* for *libaom*, the focus of this section resolves around these kernels.

3.4.1 | Functioning and Implementation

The first analysis of this section was made through the study of the transformation stage of the reference software. Its main workflow is represented in figure 3.2.

This stage is controlled by a configuration set, which is chosen according to the desired encoding objectives. These parameters control the transformation block width and height (`size_col` and `size_row`¹, respectively), the transformation kernels to use in the rows and columns, the precision to use in the sine and/or cosine coefficient approximations, as well as other parameters for overflow control. Associated to the transformation kernel chosen, the variables `ud_flip` and `lr_flip` are also set. The first one is set to `1` if the block's columns are to be transformed with the *Flip-ADST* kernel. If such choice is applied to the rows, the second variable is, likewise, set to `1`. These variables control if the input rows are flipped vertically, and/or if the coefficients resulting from the column transformation are flipped horizontally².

The choosing of these parameters will not be addressed in this work, since *AV1* allows for a great deal of maneuverability to the designer, as to adjust each encoder/decoder pair to the desired application. In this regard, *libaom* allows for a high number of configuration options, that dramatically change the parameters chosen in the transformation stage, as well as in the rest of the system.

Throughout the represented process, many of the operations are done with sequential, iterative processes, e.g., the input vector selection or the flipping operations. Such operations would greatly benefit of a hardware implementation, since they are easily parallelizable, as the objectives of *AV1* suggested. However, on this work, the focus relies of the *T* block, i.e., the transformation itself.

Independently of the transformation kernel, the operation is done sequently, in various stages. In each of these, the corresponding intermediary coefficients get calculated as function of two of the previous calculated coefficients. These, in most of the stages, are multiplied by a specific integer approximation of cosine/sine value. Such approximations, as mentioned previously, depend on the number of bits on which they are represented.

The arrays on which the calculated cosine and sine values are stored, `cospi` and `sinpi`, respectively, are bi-dimensional. The first dimension, `N`, has 7 positions, where the first corresponds to 10 bit approximations, and the last to 16 bit. The second dimension, `k`, has 64 positions for `cospi` and 5 for `sinpi`, representing the first quadrant of the trigonometric circle. Each position is calculated according to equations 3.31 and 3.32, where k represents the position in the array, and N corresponds to the number of bits. Therefore, `cospi[N][0]` corresponds to $\cos(0)$, `cospi[N][63]` is $\cos(63\pi/128)$, and the following positions can also be associated to a certain angle.

$$\text{cospi}[N][k] = \left\lfloor 2^N \cos\left(\frac{k\pi}{128}\right) \right\rfloor \quad (3.31)$$

$$\text{sinpi}[N][k] = 2^N \left\lfloor \frac{2}{3} \sqrt{2} \sin\left(\frac{k\pi}{9}\right) \right\rfloor \quad (3.32)$$

The `sinpi` array is only used in the shortest length of the *ADST*, which is the reason it only has five positions. All other versions of this kernel use `cospi` to get the desired value.

¹Each number does not correspond to the number of elements in columns and rows, but rather to the number of rows and columns.

²Here, the notation of *horizontally* or *vertically* is set considering a matrix input block. In the 1D transform implemented in *libaom*, this just means that what would be the last coefficient is now the first, and so on.

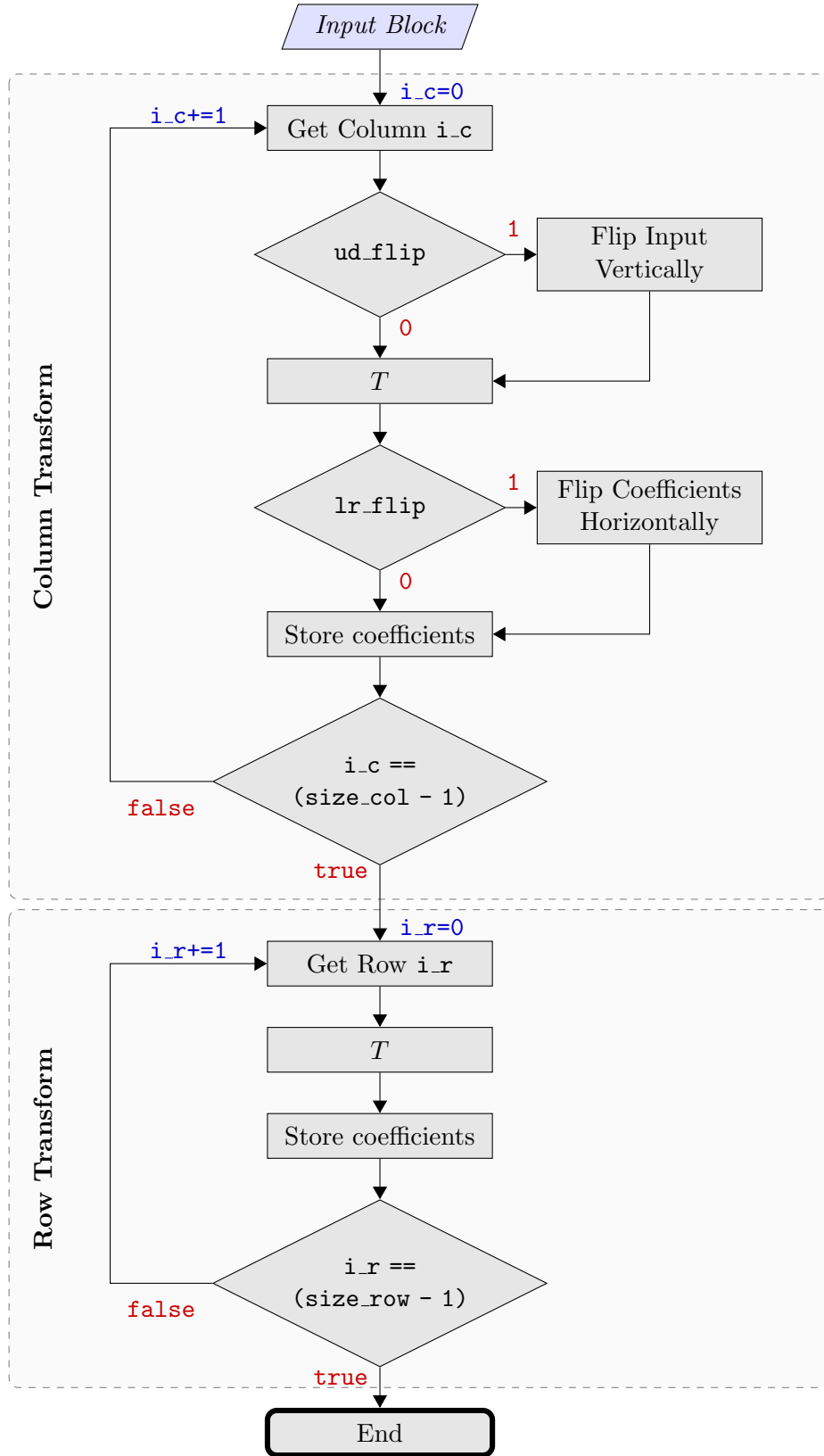


Figure 3.2: Flowchart of the Transform Stage on *libaom*

Most of the temporary coefficients inside each stage are calculated with the function `half_btf`, which performs the operation represented in equation 3.33. This function takes the two previously calculated coefficients, two values from the previously introduced arrays, as well as the number of bits used to represent these, and maps the result from the multiplications and sum of the first inputs to the desired number of bits.

$$\text{half_btf}(w_0, in_0, w_1, in_1, bits) \cong \left\lfloor \frac{w_0 in_0 + w_1 in_1 + 2^{bits-1}}{2^{bits}} \right\rfloor \quad (3.33)$$

Although the code implementation is sequential, the *8 length* transformation kernels are represented in figures 3.4 and 3.5 as parallel block diagrams, with the diverse stages in series. While *AV1* accepts *transform block* sizes varying between 4 and 64, the method of transformation is similar between the different sizes.

Both pictures start with the input vector components, i.e., `x0` to `x7`. The following sum's represent the addition of the two pointing values, in case the that the arrow guiding these doesn't present any further coefficient. If such is verified, the operation to be realized is the one presented in equation 3.33. The value near each arrow is referred to the equivalent `cospi` position, that multiplies by the result coming from the arrow's origin. Figure 3.3 presents a visual aid for the following figures.

The identity transforms, *IDTX*, are the simplest of the ones implemented in *libaom*, since they consist of a scale factor, which varies throughout the transform sizes.

On the 4 and 16 length transforms, the scaling factor includes a 12-bit integer approximation of the square root of 2, which is calculated through

$$N_{\sqrt{2}} = \left\lfloor 2^{12} \sqrt{2} \right\rfloor = 5793 \quad (3.34)$$

Being so, the input also suffers an additional mapping, similar to the operation in 3.33. These operations are demonstrated in figure 3.6.

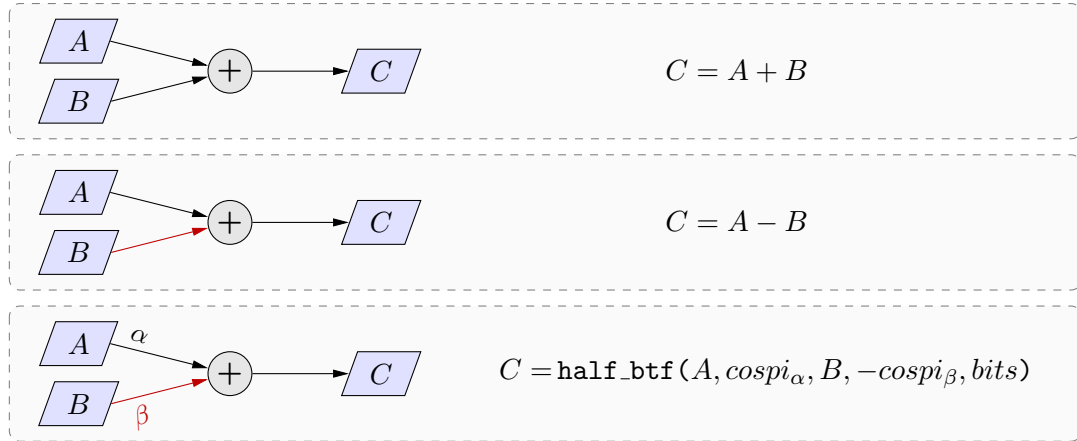


Figure 3.3: Graphical aid for figures 3.4 and 3.5

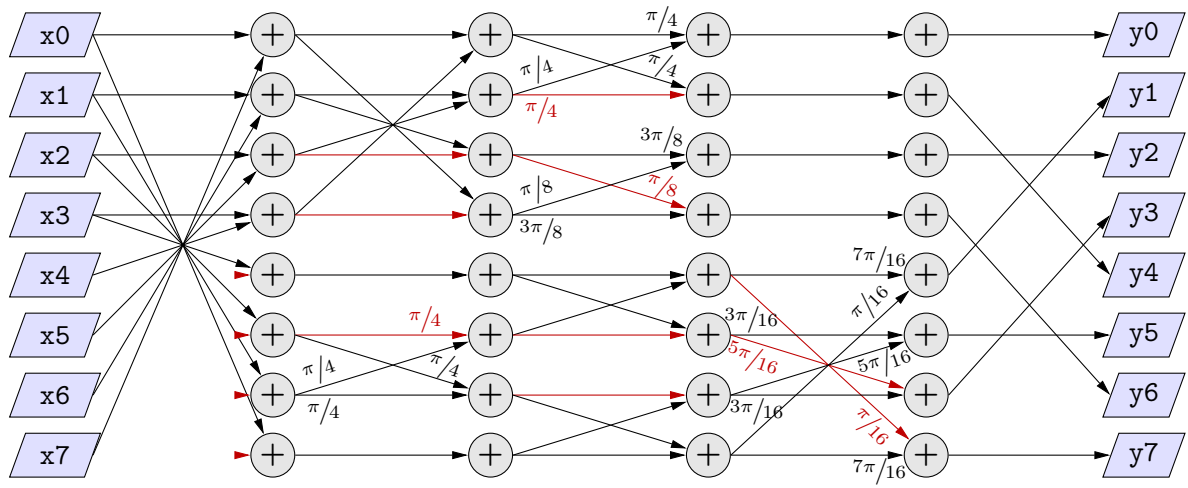


Figure 3.4: Block diagram of *libaom*'s Integer DCT

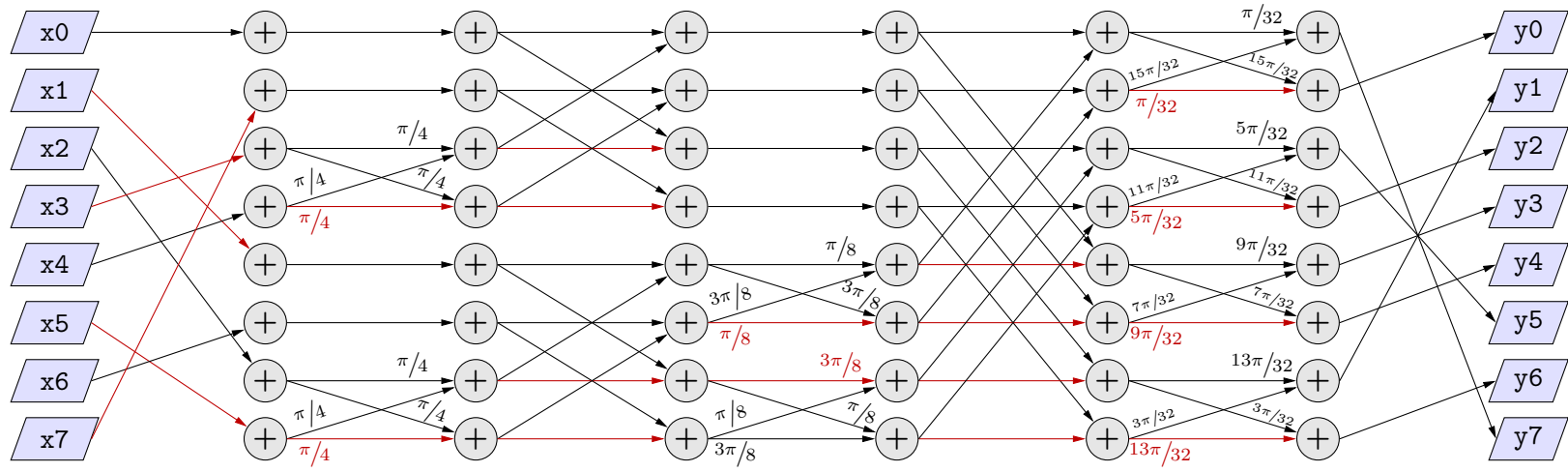


Figure 3.5: Block diagram of *libaom*'s Integer ADST

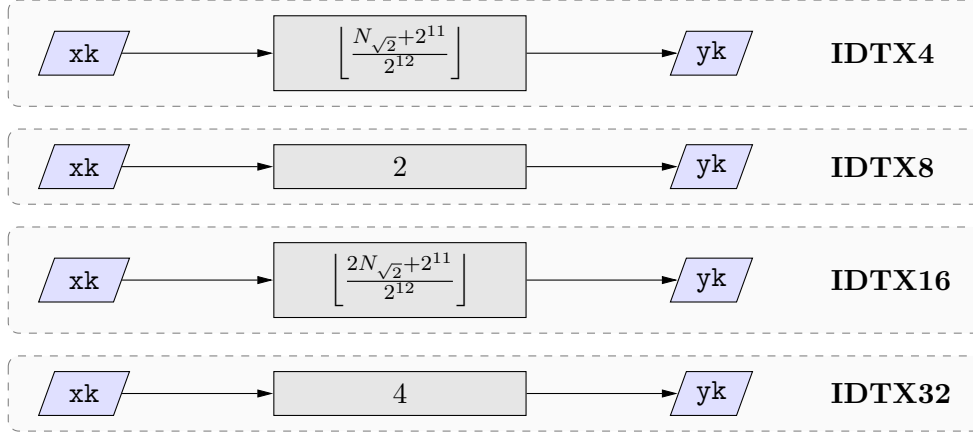


Figure 3.6: Description of the Identity transforms in *libaom*

With the forward transformations explained and represented graphically, it's easily understandable that the corresponding inverses correspond to the backwards operation in figures 3.4, 3.6 and 3.5. With this, it is meant that only the direction of the arrows change, and the corresponding procedure is done right-to-left, i.e., the output's position, y , is now the input.

3.4.2 | Performance and Statistics Analysis

From the obtained understanding of the tools and characteristics of the *Transform* stage, various tests were performed in order to gather information on what were the most commonly used options, and the corresponding impact of such on encoding performance.

With these tests, it was intended to get to know:

- the time spent per encode in the *Transform* stage;
- used kernels and sizes;
- bits used on cosine approximation.

For this, *libaom* was modified in order to retrieve these statistics, without impacting the encoding performance. Namely, the timing test was performed independently from the others, in order to get an accurate result of the time spent, since the writing of the transformation options would cause a major impact during the transform stage.

The tests were performed on different sequences with various resolutions and encoding difficulties. These were obtained in *Xiph's* test media [12], and are presented in table 3.2.

As mentioned previously in section 2.1.3.1, the definition of the encoding objectives greatly impacts the performance of the encoder, both in terms of compression gain, and obtained video quality. These objectives are typically defined according to a certain performance or quality metric. *AV1* reference encoder, *aomenc*, provides four different metrics, corresponding to constant or constrained quality (q and cq) and constant or varying bitrate (cbr and vbr). These options allow for the easy adaptation of each encode to the corresponding use case. For instance, for a network application, a certain bit rate objective, or range, is more adequate than the definition of a certain quality, since in this case the obtained bit rate may depend on the scene.

Label	Resolution		Sequence Name
	Height	Width	
CIF	288	352	<i>Waterfall</i>
			<i>Flower</i>
			<i>Bridge Close</i>
HD	720	1280	<i>Ducks take off</i>
			<i>Parkrun</i>
			<i>Shields</i>
FHD	1080	1920	<i>Parkjoy</i>
			<i>Dinner</i>
			<i>Factory</i>
UHD	2160	3840	<i>Into tree</i>
			<i>Old Town Cross</i>
			<i>Crowd Run</i>

Table 3.2: Sequences used for testing

For the performed tests, the variation of used tools and performance was evaluated by varying the desired objective on a constant quality mode, which is controlled through a subjective parameter, `cq-level`. This may vary between 0 and 63, corresponding the latter to the lowest quality encode. Three different quality parameters were tested, *60*, *25*, and *5*, corresponding to a *Low*, *Medium* and *High* quality sets, respectively.

Besides the quality objectives, the encoder was also configured to use the highest computing power, `cpu-used=8`, as well as a single pass encode, `passes=1`. These options were set in order to get the lowest encoding times, since their impact on the obtained results didn't justify the added complexity.

The resulting command for configuring *aomenc* for the developed tests, encoding the first 10 frames of each video, becomes

```
./aomenc <INPUT-FILE> -h <HEIGHT> -w <WIDTH> -o <OUTPUT-FILE> --limit=10 -p 1 --
cpu-used=8 --i420 --q-hist=64 --end-usage=q --cq-level=<CQ-LEVEL>
```

The results presented in the following sections are derived from the test encodes, which were performed on a *Ryzen 7 2700*, clocked at 3.9GHz.

3.4.2.1 | Timing Analysis

On the performed *AV1* timing analysis, the main aspect to evaluate was the impact of the time spent on *Transform* stage, for different quality thresholds.

For a certain resolution, it is expected that the time spent on this stage remains approximately constant, regardless of the desired quality. Therefore, the higher the total encoding time, the least impact the *Transform* stage would have.

This result was verified in the performed test, whose results are represented in figure 3.7.

In this figure, each tall bar corresponds to the average time spent per encode, in each of the resolutions, varying the desired quality objective. The corresponding smaller bar represents

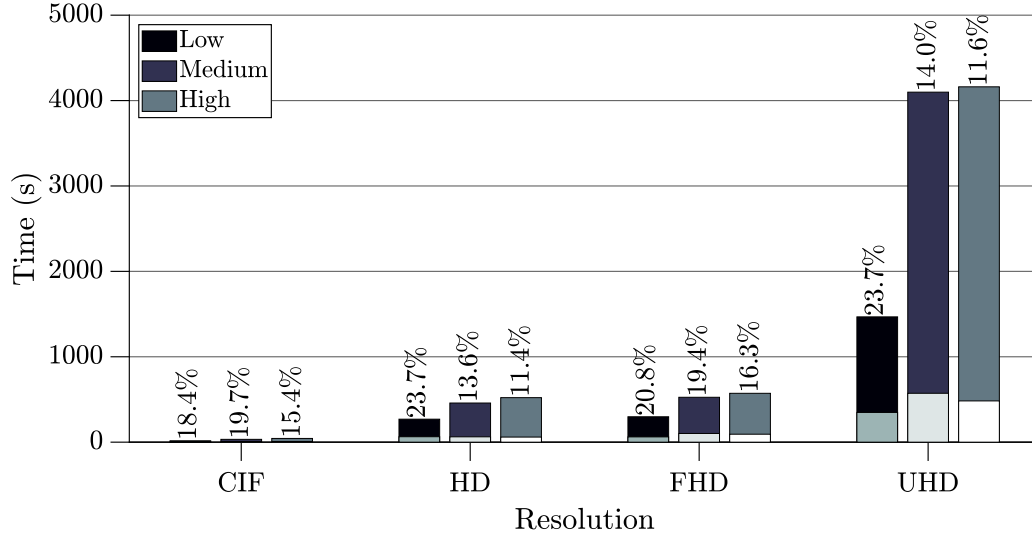


Figure 3.7: Average encoding and transform time per resolution, on different quality objectives (dark colours represent total encoding time, while lighter are the respective time spent on transform stage)

the percentage of time spent during transformation, also expressed numerically by the number on top of each bar.

As anticipated, the desired quality of each encode plays an important role in the necessary total time, majorly considering the *Low* to *Medium* quality objectives, since the latter, in average, spends *double* the time of the former. However, once considered the *High* quality set, the encoder takes 14% more time to encode the same sequence on a *Medium* quality objective.

Although these results represent an interesting analysis from a performance standpoint, for the focus of this work, the most interesting analysis comes from the percentage of time spent during the *Transform* stage. As expected, this time stays roughly the same, independently of the quality objective. However its impact to the total encoding time decreases as the quality increases.

This aspect leads to conclude that there is a necessity to develop fast and efficient architectures for the *Transform* block, since it corresponds to a relevant percentage of the total encode time, regardless of the quality. And although lower quality encodes benefit more of such improvements when comparing to higher qualities, both cases would gain, since such architecture could be used on a high variety of encoders.

3.4.2.2 | Configuration Set Analysis

This analysis is divided in different segments, dedicated to each of the transform options analyzed. Some of which, although not entirely relevant to the aim of this work, may prove useful for the implementation of different architectures.

On the distribution of used kernels verified in figure 3.8, the most relevant aspect is the clear dominance of the *DCT* among the others. However, as quality increases, the distribution starts to spread out.

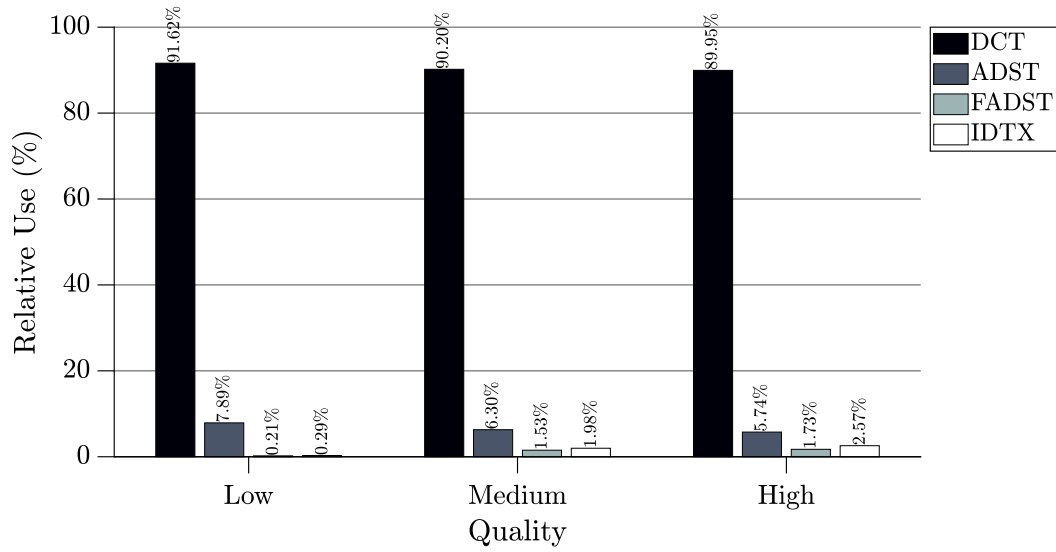


Figure 3.8: Average distribution of used kernels, for all resolutions, according to the quality threshold

A contrary behavior is verified on the transform vector size. As seen in figure 3.9, as quality increases, the smaller vector sizes (namely the 4 length vector), get used more frequently. This is easily understandable, as smaller blocks present lower losses during the *Quantization* stage.

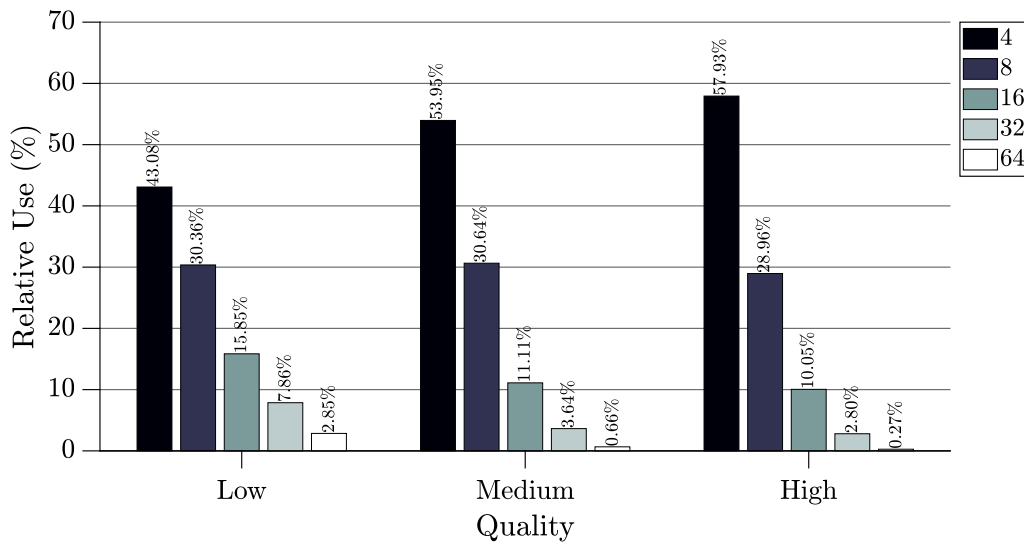


Figure 3.9: Average distribution of vector sizes, for all resolutions, according to the quality threshold

Although *AV1* supports asymmetric transform blocks, in figure 3.10 it's possible to verify that the encoder, in most of the block transformations, doesn't take advantage of such, using square blocks, as well as the same kernel for both directions. This behavior remains similar throughout the different resolutions.

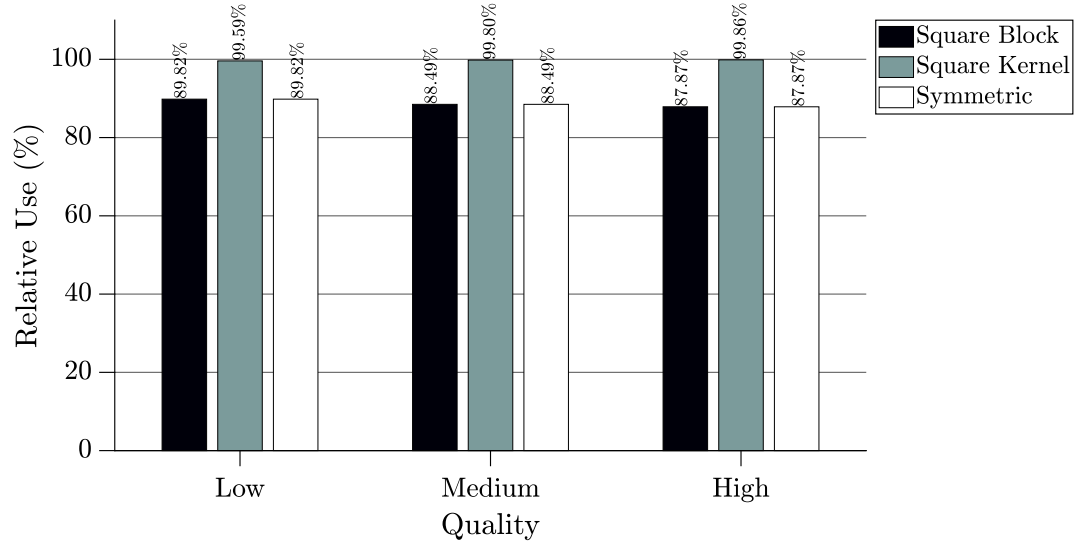


Figure 3.10: Use of square blocks, same kernel for rows and columns, and symmetric kernels, according to the quality threshold

This aspect may prove a starting point for improvement of the transform stage, since, as mentioned previously, symmetric transforms may be implemented with *Matrix Multiplication (MM)*. Therefore, using fast *MM* architectures for symmetric blocks, and complementing with the algorithm implemented in *libaom* for asymmetric blocks, the transform time may decrease.

The final analyzed transform option is the number of bits used in the cosine representation. In figure 3.11, the distribution is represented for the different quality objectives.

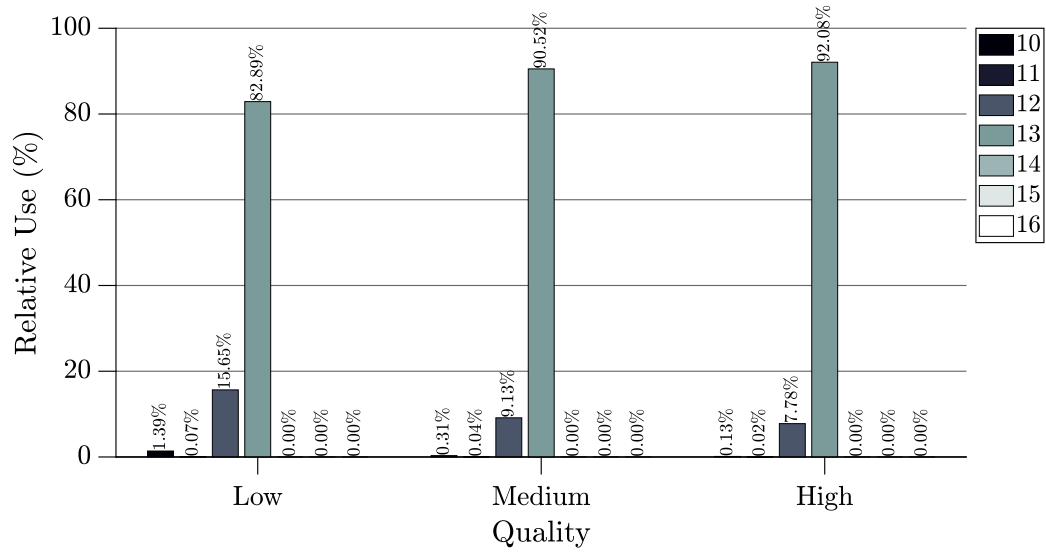


Figure 3.11: Different number of bits used on the cosine approximations, throughout different quality sets

Various conclusions can be derived from this data. Firstly, as expected, as quality increases, so does the number of bits for cosine representation, as seen by the increase of the percentage on the 13 bit representation.

However, most representations hardly get used, since most of the transformations use 13, 12, 10 and, very few times, 11 bits, without using any of the higher representations (on the performed tests).

Analyzing this information, it may be thought that the number of bits used in the cosine contributes for the overall quality of the compressed sequence. In the following section, this aspect is tested.

3.4.2.3 | Quality Analysis

In this test, besides evaluating the obtained quality for each tested `cq-level`, the impact of the cosine bits representation also was evaluated.

To evaluate the impact of the number of bits used in the cosine approximations, *aomenc* was modified to force either 10 or 16 bits throughout the encoding operation, for both forward (T) and inverse transformations (T^{-1}). *Aomdec* (reference decoder), since it acts according to the specified *Bitstream Decoding Format* [13], it may not be modified. Nonetheless, it uses 12 bit representation, regardless of the choices made in the decoder.

From the gathered reconstructed sequences, \hat{g} , the Peak Signal to Noise Ratio (PSNR) of each one was calculated, according to equation 3.35.

$$PSNR = 10 \log_{10} \left(\frac{255^2}{E_{ms}} \right) \quad (3.35)$$

E_{ms} corresponds to the *Mean Squared Error* of the reconstructed video. Considering a single $M \times N$ monochrome frame, this error is given by equation 3.36.

$$E_{ms} = \frac{1}{MN} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} (g_{x,y} - \hat{g}_{x,y})^2 \quad (3.36)$$

However, since the test revolves around a sequence of K reconstructed frames, with three chroma channels per bit, c , E_{ms} becomes

$$E_{ms} = \frac{1}{3KMN} \sum_{k=0}^{K-1} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} \sum_{c=0}^{3-1} (g_{k,x,y,c} - \hat{g}_{k,x,y,c})^2 \quad (3.37)$$

The workflow of the performed test is represented in figure 3.12.

The average results from all resolutions, for each of the quality objectives is represented in figure 3.13.

As expected, as `cq-level` increases, so does the obtained quality. Also, considering the encoding time differences verified in section 3.4.2.1, the smaller PSNR gap between *Medium* and *High* qualities was also expected. However, the difference between these two parameters depends on the encoded video, as shown in figure 3.14, where the difference between *Medium* and *High* encodes is $2dB$.

The differences between obtained PSNRs can be easily explained through analysis of the *Quantization* stage in each quality objective. Looking at the distribution of the *Quantizer/Quantization parameter (QP)* ³ throughout the different `cq-level`'s (figure 3.15), it is

³Parameter that indicates the *quantization matrix* to use (higher values indicate more severe quantization).

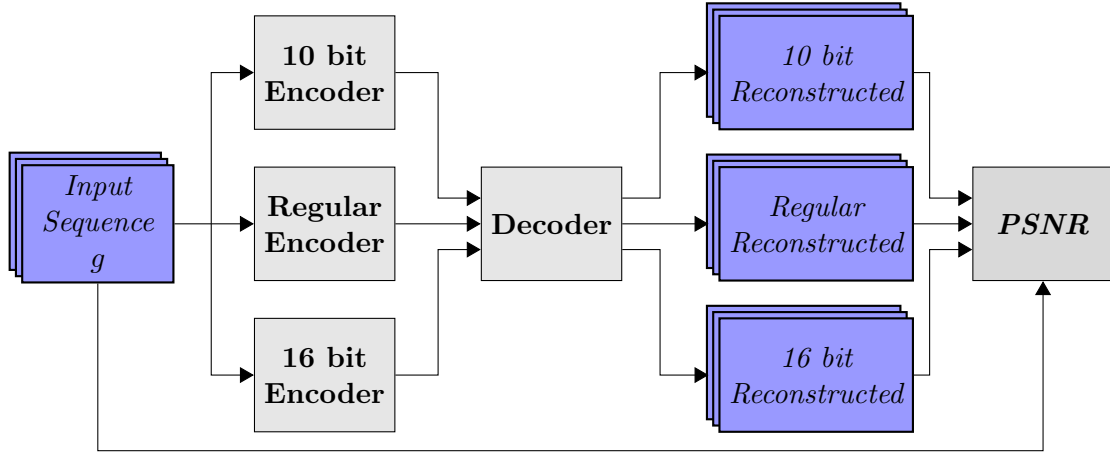


Figure 3.12: Description of the test for comparing impact of number of bits in cosine approximations

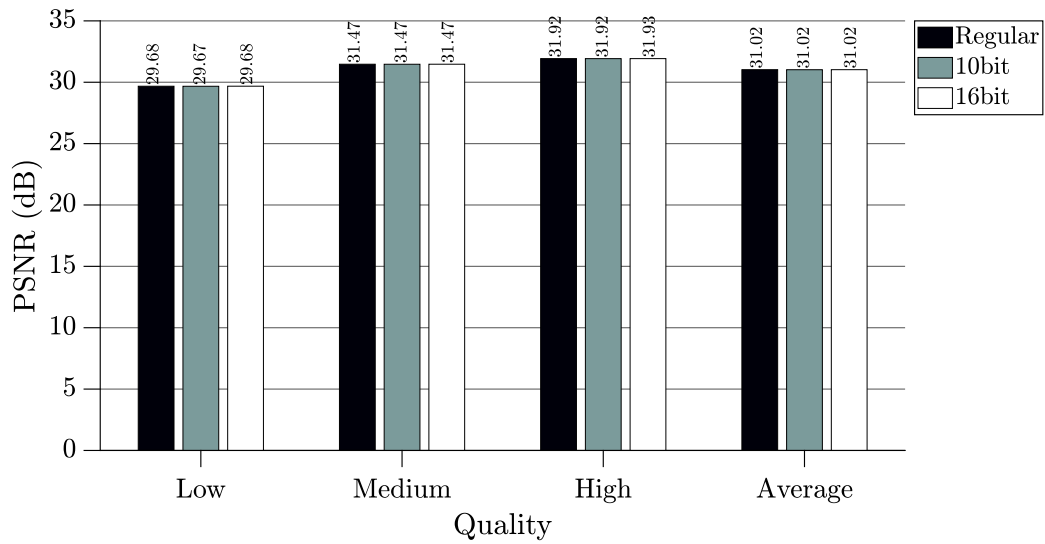
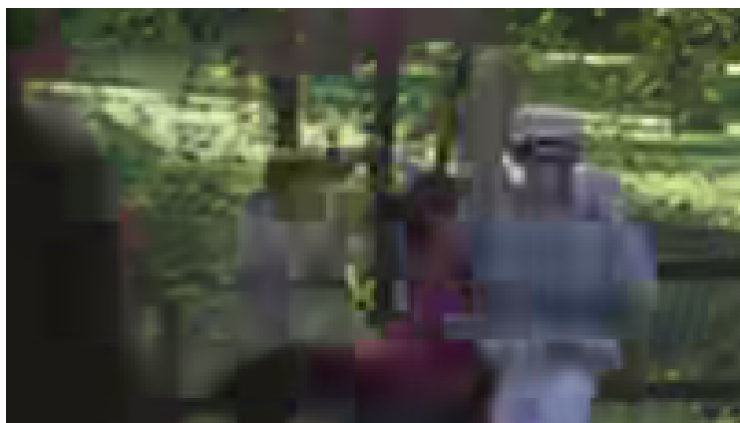
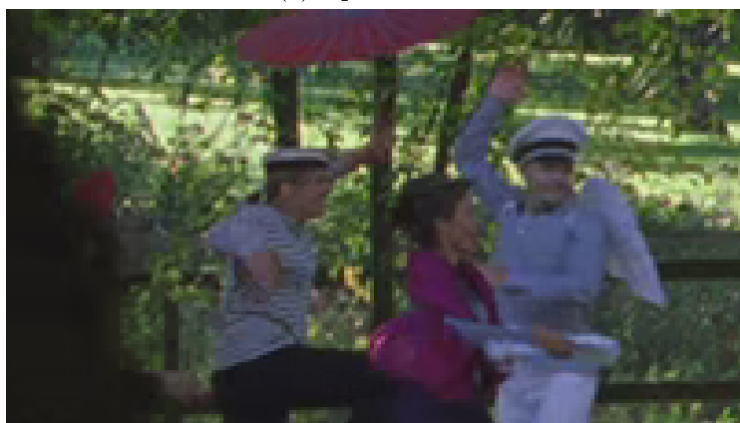


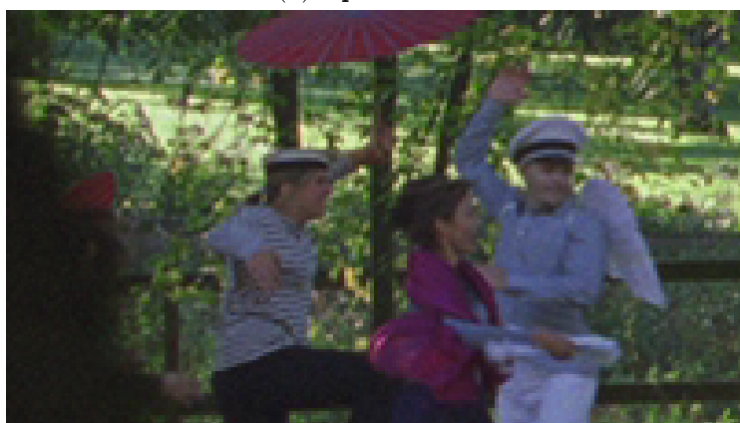
Figure 3.13: Obtained quality for each of the quality objectives, and comparison with different cosine bits approximation



(a) cq-level=60



(b) cq-level=25



(c) cq-level=5

Figure 3.14: Detail of *Parkjoy* encodes, through different quality objectives

possible to verify that this stage deeply adapts to the desired quality objective, increasing QP for lower qualities.

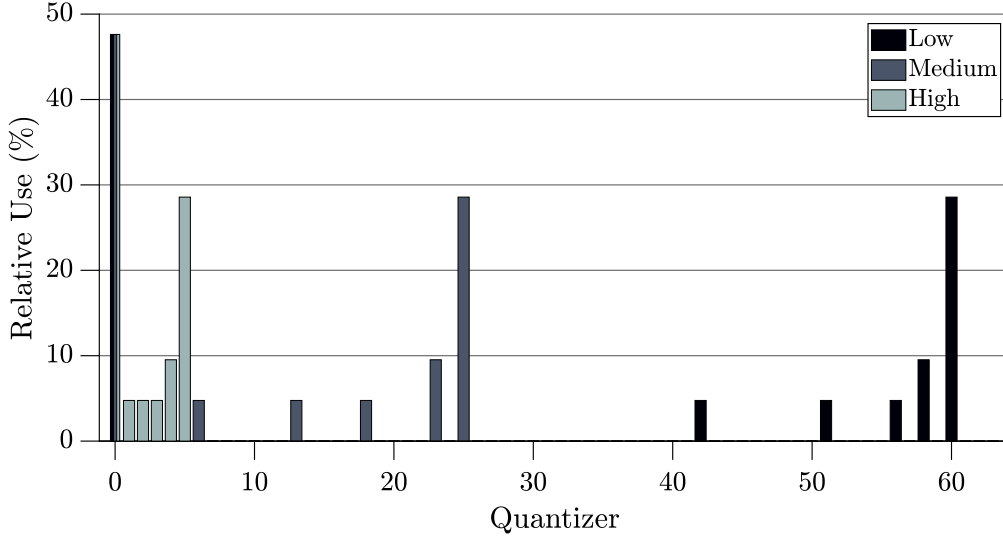


Figure 3.15: *Quantizer* distribution on different quality objectives

However, the differences in the used QP's don't justify the increased encoding time throughout the different quality objectives, since the *Quantization* stage's complexity shouldn't vary depending on the desired quality, similarly to what happens on the *Transform* stage. The time difference is mainly caused in the higher complexity blocks, the *Inter* and *Intra Prediction* stages, as the encoder adapts the processes in these blocks depending on the desired quality. For instance, there is no need for the encoder to make an exceptionally precise prediction, when most of the transform coefficients are discarded.

Considering now the obtained quality for each of the three different encodes (*Regular*, *10 bit* and *16 bit*), it is possible to observe that the number of bits on cosine approximation doesn't contribute to the obtained quality, regardless of the desired objective, contrary to what was verified in figure 3.11. Accordingly, it would be safe to assume that the cosine approximations could be fixed on a certain number of bits, without major impact to the video quality.

This factor presents a major point for exploring faster *Transform* block architectures. The architectures in figures 3.4 and 3.5 are highly dependent on the `half_btf` function (equation 3.33), which is adaptable to the number of bits used for cosine, in each block transformation. However, with the use of a fixed number of bits, this function could be simplified, since the multiplications and divisions performed in it could be implemented with a fixed number of shifts and additions.

References

- [1] Yun Qing Shi and Huifang Sun. *Image and Video Compression for Multimedia Engineering: Fundamentals, Algorithms, and Standards*. eng. 2. ed. Image Processing Series. OCLC: 254564955. Boca Raton, Fla.: CRC Press, 2008. ISBN: 978-0-8493-7364-0.
- [2] Rafael C. Gonzalez and Richard E. Woods. *Digital Image Processing*. eng. Fourth edition. OCLC: 987436552. New York, NY: Pearson, 2018. ISBN: 978-0-13-335672-4.
- [3] N. Ahmed, T. Natarajan, and K. R. Rao. “Discrete Cosine Transform”. In: *IEEE Transactions on Computers* C-23.1 (Jan. 1974), pp. 90–93. ISSN: 0018-9340.
- [4] Alan V. Oppenheim, Ronald W. Schaffer, and John R. Buck. *Discrete-Time Signal Processing*. eng. 2nd ed. Upper Saddle River, NJ: Prentice Hall, 1998. ISBN: 978-0-13-754920-7.
- [5] *Discrete Cosine Transform - MATLAB Dct*.
- [6] William K. Pratt. *Digital Image Processing: PIKS Inside*. eng. 3. ed. A Wiley-Interscience Publication. OCLC: 248389925. New York: Wiley, 2001. ISBN: 978-0-471-37407-7.
- [7] Jingning Han, Yaowu Xu, and Debargha Mukherjee. “A Butterfly Structured Design of the Hybrid Transform Coding Scheme”. en. In: *2013 Picture Coding Symposium (PCS)*. San Jose, CA, USA: IEEE, Dec. 2013, pp. 17–20. ISBN: 978-1-4799-0294-1 978-1-4799-0292-7.
- [8] Benny Bing. *Next-Generation Video Coding and Streaming*. Hoboken: Wiley, 2015. ISBN: 978-1-119-13332-2 978-1-119-13333-9.
- [9] Soo-Chang Pei and Jian-Juin Ding. “The Integer Transforms Analogous to Discrete Trigonometric Transforms”. In: *IEEE Transactions on Signal Processing* 48.12 (Dec. 2000), pp. 3345–3364. ISSN: 1053-587X.
- [10] S. Wolter et al. “Parallel Architectures for 8*8 Discrete Cosine Transforms”. In: *[Proceedings] 1992 IEEE International Symposium on Circuits and Systems*. Vol. 1. San Diego, CA, USA: IEEE, 1992, pp. 149–152. ISBN: 978-0-7803-0593-9.
- [11] Yonghong Zeng et al. “Integer DCTs and Fast Algorithms”. In: *IEEE Trans. Signal Process.* 49.11 (Nov./2001), pp. 2774–2782. ISSN: 1053587X.
- [12] *Xiph.Org :: Test Media*. <https://media.xiph.org/>.
- [13] “AV1 Bitstream & Decoding Process Specification”. en. In: (), p. 681.

CHAPTER 4

Developed Architecture

4.1 | TO BE DONE

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Appendices

Appendix A | *aomenc* Configuration Options

Usage: `./aomenc <options> -o dst_filename src_filename`

Options:

- `--help` Show usage options and exit
- `-c <arg>`, `--cfg=<arg>` Config file to use
- `-D`, `--debug` Debug mode (makes output deterministic)
- `-o <arg>`, `--output=<arg>` Output filename
 - `--codec=<arg>` Codec to use
- `-p <arg>`, `--passes=<arg>` Number of passes (1/2)
 - `--pass=<arg>` Pass to execute (1/2)
 - `--fpf=<arg>` First pass statistics file name
 - `--limit=<arg>` Stop encoding after n input frames
 - `--skip=<arg>` Skip the first n input frames
 - `--good` Use Good Quality Deadline
- `-q`, `--quiet` Do not print encode progress
- `-v`, `--verbose` Show encoder parameters
 - `--psnr` Show PSNR in status line
 - `--webm` Output WebM (default when WebM IO is enabled)
 - `--ivf` Output IVF
 - `--obu` Output OBU
- `-P`, `--output-partitions` Makes encoder output partitions. Requires IVF output!
 - `--q-hist=<arg>` Show quantizer histogram (n-buckets)
 - `--rate-hist=<arg>` Show rate histogram (n-buckets)
 - `--disable-warnings` Disable warnings about potentially incorrect encode settings.
- `-y`, `--disable-warning-prompt` Display warnings, but do not prompt user to continue.
 - `--test-decode=<arg>` Test encode/decode mismatch
 - off, fatal, warn

Encoder Global Options:

- `--yv12` Input file is YV12
- `--i420` Input file is I420 (default)
- `--i422` Input file is I422
- `--i444` Input file is I444
- `-u <arg>`, `--usage=<arg>` Usage profile number to use
- `-t <arg>`, `--threads=<arg>` Max number of threads to use
 - `--profile=<arg>` Bitstream profile number to use
- `-w <arg>`, `--width=<arg>` Frame width
- `-h <arg>`, `--height=<arg>` Frame height
 - `--forced_max_frame_width` Maximum frame width value to force
 - `--forced_max_frame_height` Maximum frame height value to force
 - `--stereo-mode=<arg>` Stereo 3D video format
 - mono, left-right, bottom-top, top-bottom, right-left
 - `--timebase=<arg>` Output timestamp precision (fractional seconds)
 - `--fps=<arg>` Stream frame rate (rate/scale)
 - `--global-error-resilient=<` Enable global error resiliency features
- `-b <arg>`, `--bit-depth=<arg>` Bit depth for codec (8 for version <=1, 10 or 12 for version 2)

```

8, 10, 12
--lag-in-frames=<arg> Max number of frames to lag
--large-scale-tile=<arg> Large scale tile coding (0: off (default), 1: on
)
--monochrome Monochrome video (no chroma planes)
--full-still-picture-hdr Use full header for still picture

Rate Control Options:
--drop-frame=<arg> Temporal resampling threshold (buf %)
--resize-mode=<arg> Frame resize mode
--resize-denominator=<arg> Frame resize denominator
--resize-kf-denominator=<a Frame resize keyframe denominator
--superres-mode=<arg> Frame super-resolution mode
--superres-denominator=<ar Frame super-resolution denominator
--superres-kf-denominator= Frame super-resolution keyframe denominator
--superres-qthresh=<arg> Frame super-resolution qindex threshold
--superres-kf-qthresh=<arg Frame super-resolution keyframe qindex
threshold
--end-usage=<arg> Rate control mode
vbr, cbr, cq, q
--target-bitrate=<arg> Bitrate (kbps)
--min-q=<arg> Minimum (best) quantizer
--max-q=<arg> Maximum (worst) quantizer
--undershoot-pct=<arg> Datarate undershoot (min) target (%)
--overshoot-pct=<arg> Datarate overshoot (max) target (%)
--buf-sz=<arg> Client buffer size (ms)
--buf-initial-sz=<arg> Client initial buffer size (ms)
--buf-optimal-sz=<arg> Client optimal buffer size (ms)

Twopass Rate Control Options:
--bias-pct=<arg> CBR/VBR bias (0=CBR, 100=VBR)
--minsection-pct=<arg> GOP min bitrate (% of target)
--maxsection-pct=<arg> GOP max bitrate (% of target)

Keyframe Placement Options:
--enable-fwd-kf=<arg> Enable forward reference keyframes
--kf-min-dist=<arg> Minimum keyframe interval (frames)
--kf-max-dist=<arg> Maximum keyframe interval (frames)
--disable-kf Disable keyframe placement

AV1 Specific Options:
--cpu-used=<arg> CPU Used (0..8)
--dev-sf=<arg> Dev Speed (0..255)
--auto-alt-ref=<arg> Enable automatic alt reference frames
--sharpness=<arg> Loop filter sharpness (0..7)
--static-thresh=<arg> Motion detection threshold
--single-tile-decoding=<ar Single tile decoding (0: off (default), 1: on)
--tile-columns=<arg> Number of tile columns to use, log2
--tile-rows=<arg> Number of tile rows to use, log2 (set to 0 while
threads > 1)
--arnr-maxframes=<arg> AltRef max frames (0..15)
--arnr-strength=<arg> AltRef filter strength (0..6)
--tune=<arg> Distortion metric tuned with

```

```

                                psnr, ssim, cdef-dist, daala-dist
--cq-level=<arg> Constant/Constrained Quality level
--max-intra-rate=<arg> Max I-frame bitrate (pct)
--max-inter-rate=<arg> Max P-frame bitrate (pct)
--gf-cbr-boost=<arg> Boost for Golden Frame in CBR mode (pct)
--lossless=<arg> Lossless mode (0: false (default), 1: true)
--enable-cdef=<arg> Enable the constrained directional enhancement filter
                    (0: false, 1: true (default))
--enable-restoration=<arg> Enable the loop restoration filter (0: false,
                    1: true (default))
--disable-trellis-quant=<a Disable trellis optimization of quantized
                    coefficients (0: false (default) 1: true)
--enable-qm=<arg> Enable quantisation matrices (0: false (default), 1:
                    true)
--qm-min=<arg> Min quant matrix flatness (0..15), default is 8
--qm-max=<arg> Max quant matrix flatness (0..15), default is 15
--enable-dist-8x8=<arg> Enable dist-8x8 (0: false (default), 1: true)
--frame-parallel=<arg> Enable frame parallel decodability features (0:
                    false (default), 1: true)
--error-resilient=<arg> Enable error resilient features (0: false (
                    default), 1: true)
--aq-mode=<arg> Adaptive quantization mode (0: off (default), 1: variance
                    2: complexity, 3: cyclic refresh)
--deltaq-mode=<arg> Delta qindex mode (0: off (default), 1: deltaq 2:
                    deltaq + deltalf)
--frame-boost=<arg> Enable frame periodic boost (0: off (default), 1: on)
--noise-sensitivity=<arg> Noise sensitivity (frames to blur)
--tune-content=<arg> Tune content type
                    default, screen
--cdf-update-mode=<arg> CDF update mode for entropy coding (0: no CDF
                    update; 1: update CDF on all frames(default); 2: selectively update
                    CDF on some frames
--color-primaries=<arg> Color primaries (CICP) of input content:
                    bt709, unspecified, bt601, bt470m, bt470bg,
                    smpte240, film, bt2020, xyz, smpte431,
                    smpte432, ebu3213
--transfer-characteristics Transfer characteristics (CICP) of input
                    content:
                                unspecified, bt709, bt470m, bt470bg, bt601,
                                smpte240, lin, log100, log100sq10, iec61966,
                                bt1361, srgb, bt2020-10bit, bt2020-12bit,
                                smpte2084, hlg, smpte428
--matrix-coefficients=<arg> Matrix coefficients (CICP) of input content:
                    identity, bt709, unspecified, fcc73, bt470bg,
                    bt601, smpte240, ycgco, bt2020ncl, bt2020cl,
                    smpte2085, chromncl, chromcl, ictcp
--chroma-sample-position=< The chroma sample position when chroma 4:2:0
                    is signaled:
                                unknown, vertical, colocated
--min-gf-interval=<arg> min gf/arf frame interval (default 0, indicating
                    in-built behavior)
--max-gf-interval=<arg> max gf/arf frame interval (default 0, indicating
                    in-built behavior)

```

```

--sb-size=<arg> Superblock size to use
                    dynamic, 64, 128
--num-tile-groups=<arg> Maximum number of tile groups, default is 1
--mtu-size=<arg> MTU size for a tile group, default is 0 (no MTU
                    targeting), overrides maximum number of tile groups
--timing-info=<arg> Signal timing info in the bitstream (model only works
                    for no hidden frames, no super-res yet):
                    unspecified, constant, model
--film-grain-test=<arg> Film grain test vectors (0: none (default), 1:
                    test-1 2: test-2, ... 16: test-16)
--film-grain-table=<arg> Path to file containing film grain parameters
--enable-ref-frame-mvs=<ar Enable temporal mv prediction (default is 1)
-b <arg>, --bit-depth=<arg> Bit depth for codec (8 for version <=1, 10 or 12 for
                    version 2)
                    8, 10, 12
--input-bit-depth=<arg> Bit depth of input
--sframe-dist=<arg> S-Frame interval (frames)
--sframe-mode=<arg> S-Frame insertion mode (1..2)
--annexb=<arg> Save as Annex-B

```

Stream timebase (--timebase):

The desired precision of timestamps in the output, expressed in fractional seconds. Default is 1/1000.

Included encoders:

av1 - AOMedia Project AV1 Encoder v0.1.0 (default)

Use --codec to switch to a non-default encoder.