



**ISO/IEC JTC 1/SC 29/WG 1  
(& ITU-T SG16)**

**Coding of Still Pictures**

**JBIG**

Joint Bi-level Image  
Experts Group

**JPEG**

Joint Photographic  
Experts Group

**TITLE:** Text of ISO/IEC WD2 XXXXX-1

**SOURCE:** JPEG HDR Ad Hoc Group  
Thomas Richter, Germany  
Walt Husak, USA  
Chaker Larabi, France  
Touradj Ebrahimi, Switzerland  
Alessandro Artusi, UK  
Massimiliano Agostinelli, UK

**PROJECT:** ISO/IEC XXXXX-1  
(JPEG HD Part 1 – JPEG High Dynamic Range Image coding – Specification)

**STATUS:** WD1

**REQUESTED**

**ACTION:** For WD1 ballot

**DISTRIBUTION:** For WG1 review



**INTERNATIONAL ORGANISATION FOR STANDARDISATION**  
**INTERNATIONAL ELECTROTECHNICAL COMMISSION**

---

**Information technology –**

**JPEG digital imaging systems integration –**

**JPEG High Dynamic Range Image Coding  
Systems – Core Coding System Specification**

Draft ITU-T Recommendation | International Standard

---

## CONTENTS

*Page*

## Foreword

This Recommendation | International Standard covers compression and representation of high dynamic range images backwards compatible to ITU.T Rec.81 | ISO/IEC 10918-1. That is, legacy applications conforming to ITU.T Rec. 81 | ISO/IEC 10918-1 will be able to reconstruct streams generated by an encoder conforming to this Recommendation | International Standard, though will possibly not be able to reconstruct such streams in full dynamic range or full quality.

The aim of this standard is to provide a migration path for legacy applications to support, potentially in a limited way, high dynamic range images. Existing tools depending on the existing standards will continue to work, but will only work on a low dynamic range version of the image. This Recommendation | International Standard specifies a coded file format, referred to as JPEG HD, which is designed primarily for storage and interchange of continuous-tone photographic content.

## Introduction

This Recommendation | International standard specifies a coded codestream format for storage of continuous-tone high and low dynamic range photographic content. JPEG HD is a scalable image coding system supporting multiple component images in bitranges from eight up to XXX bits per sample; the codestream is composed in such a way that legacy applications conforming to ITU.T 81 | ISO/IEC 10918-1 are able to reconstruct a lower quality, low dynamic range, eight bits per sample version of the image.

JPEG HD is primarily designed to represent images of high dynamic range content, even though images of low dynamic range using only eight bits per pixel are supported. The goal is to provide a backwards compatible coding specification that allows legacy applications and existing toolchains to continue to operate on codestreams conforming to this Recommendation|International Standard.

Today, the most widely used digital photography format, a minimal implementation of JPEG (specified in ITU-T Recommendation T.81 | ISO/IEC 10918-1), uses a bit depth of 8; each of the three channels that together compose an image pixel is represented by 8 bits, providing 256 representable values per channel or 16 777 216 representable color values. For more demanding applications, it is not uncommon to use a bit depth of 16, providing 65 536 representable values to describe each channel within a pixel, resulting on over  $2.8 \times 10^{14}$  representable color values. In some less common scenarios, even greater bit depths are used. In scenarios when memory or processing power is at a premium, as few as five bits may be used, providing only 32 representable values per channel.

Most common photo and image formats use an 8-bit or 16-bit unsigned integer value to represent some function of the intensity of each color channel. While it might be theoretically possible to agree on one method for assigning specific numerical values to real world colors, doing so is not practical. Since any specific device has its own limited range for color reproduction, the device's range may be a small portion of the agreed-upon universal color range. As a result, such an approach is an extremely inefficient use of the available numerical values, especially when using only 8 bits (or 256 unique values) per channel.

To represent pixel values as efficiently as possible, devices use a numeric encoding optimized for their own range of possible colors or gamut.

JPEG HD has been designed to be backwards compatible to legacy applications while at the same time have a small coding complexity; JPEG HD uses, whenever possible, functional blocks of ITU.T 81|ISO/IEC 10918-1 to represent high dynamic range images. It is optimized for good image quality and compression efficiency while also enabling low-complexity encoding and decoding implementations. The design objectives include:

- Backwards compatible representation of images of various bitdepths, embedded system friendly compression
- Small memory footprint
- High compression quality

The algorithm uses a reversible lapped biorthogonal transform. The transform requires at most 3 non-trivial (multiply plus addition) and 7 trivial (addition or shift) operations per sample (with no divisions) at the highest quality level. In the lowest complexity mode, 1 non-trivial and 4 trivial operations per sample are required for transform. The image is processed in 16x16 macroblocks, allowing a minimal memory footprint for embedded implementations.

The algorithm provides native support for both RGB and CMYK by converting these color formats to an internal luma-dominant format through the use of a reversible color transform. In addition, YUV, monochrome and arbitrary n-component color formats are supported.

INTERNATIONAL STANDARD

ITU-T RECOMMENDATION

**INFORMATION TECHNOLOGY –  
INFORMATION TECHNOLOGY – JPEG HDR IMAGE CODING SYSTEM:  
CODING OF HIGH DYNAMIC RANGE IMAGES**

## **1 Scope**

This Recommendation | International standard specifies a coding format, referred to as JPEG HD, which is designed primarily for continuous-tone photographic content using a sample precision of above eight bits per sample.

## **2 Normative references**

The following Recommendations and International Standards contain provisions which, through reference in this text, constitute provisions of this Recommendation | International Standard. At the time of publication, the editions indicated were valid. All Recommendations and Standards are subject to revision, and parties to agreements based on this Recommendation | International Standard are encouraged to investigate the possibility of applying the most recent edition of the Recommendations and Standards listed below. Members of IEC and ISO maintain registers of currently valid International Standards. The Telecommunication Standardization Bureau of the ITU maintains a list of currently valid ITU-T Recommendations.

### **2.1 Identical Recommendations | International Standards**

ITU-T XXX

### **2.2 Paired Recommendations | International Standards equivalent in technical content**

To be determined.

### **2.3 Additional references**

ITU-T Rec. T.81 | ISO/IEC 10918-1

## **3 Definitions, Abbreviations and Symbols**

### **3.1 Definitions**

For the purposes of this Recommendation | International Standard, the following definitions apply.

**AC coefficient:** Any DCT coefficient for which the frequency is not zero in at least one dimension.

**binary decision:** Choice between two alternatives.

**bit stream:** Partially encoded or decoded sequence of bits comprising an entropy-coded segment.

**block:** An  $8 \times 8$  array of samples or an  $8 \times 8$  array of DCT coefficient values of one component.

**block-row:** A sequence of eight contiguous component lines which are partitioned into  $8 \times 8$  blocks.

**byte:** A group of 8 bits.

**coder:** An embodiment of a coding process.

**coding:** Encoding or decoding.

**coding model:** A procedure used to convert input data into symbols to be coded.

**(coding) process:** A general term for referring to an encoding process, a decoding process, or both.

**compression:** Reduction in the number of bits used to represent source image data.

**continuous-tone image:** An image whose components have more than one bit per sample.

**data unit:** An  $8 \times 8$  block of samples of one component in DCT-based processes; a sample in lossless processes.

**DC coefficient:** The DCT coefficient for which the frequency is zero in both dimensions.

**(DCT) coefficient:** The amplitude of a specific cosine basis function – may refer to an original DCT coefficient, to a quantized DCT coefficient, or to a dequantized DCT coefficient.

**decoder:** An embodiment of a decoding process.

**decoding process:** A process which takes as its input compressed image data and outputs a continuous-tone image.

**dequantization:** The inverse procedure to quantization by which the decoder recovers a representation of the DCT coefficients.

**(digital) reconstructed image (data):** A continuous-tone image which is the output of any decoder defined in this Specification.

**(digital) source image (data):** A continuous-tone image used as input to any encoder defined in this Specification.

**(digital) (still) image:** A set of two-dimensional arrays of samples comprised of one or many components.

**discrete cosine transform; DCT:** Either the forward discrete cosine transform or the inverse discrete cosine transform.

**downsampling (filter):** A procedure by which the spatial resolution of an image is reduced (in hierarchical mode coding).

**encoder:** An embodiment of an encoding process.

**encoding process:** A process which takes as its input a continuous-tone image and outputs compressed image data.

**entropy-coded (data) segment:** An independently decodable sequence of entropy encoded bytes of compressed image data.

**(entropy-coded segment) pointer:** The variable which points to the most recently placed (or fetched) byte in the entropy encoded segment.

**entropy decoder:** An embodiment of an entropy decoding procedure.

**entropy decoding:** A lossless procedure which recovers the sequence of symbols from the sequence of bits produced by the entropy encoder.

**entropy encoder:** An embodiment of an entropy encoding procedure.

**entropy encoding:** A lossless procedure which converts a sequence of input symbols into a sequence of bits such that the average number of bits per symbol approaches the entropy of the input symbols.

**extended (DCT-based) process:** A descriptive term for DCT-based encoding and decoding processes in which additional capabilities are added to the baseline sequential process.

**forward discrete cosine transform;**

**FDCT:** A mathematical transformation using cosine basis functions which converts a block of samples into a corresponding block of original DCT coefficients.

**frequency:** A two-dimensional index into the two-dimensional array of DCT coefficients.

**full progression:** A process which uses both spectral selection and successive approximation (in progressive mode coding).

**grayscale image:** A continuous-tone image that has only one component.

**horizontal sampling factor:** The relative number of horizontal data units of a particular component with respect to the number of horizontal data units in the other components.

**Huffman decoder:** An embodiment of a Huffman decoding procedure.

**Huffman decoding:** An entropy decoding procedure which recovers the symbol from each variable length code produced by the Huffman encoder.

**Huffman encoder:** An embodiment of a Huffman encoding procedure.

**Huffman encoding:** An entropy encoding procedure which assigns a variable length code to each input symbol.

**Joint Photographic Experts Group; JPEG:** The informal name of the committee which created this Specification. The “joint” comes from the ITU-T and ISO/IEC collaboration.

**lossless:** A descriptive term for encoding and decoding processes and procedures in which the output of the decoding procedure(s) is identical to the input to the encoding procedure(s).

**lossless coding:** The mode of operation which refers to any one of the coding processes defined in this Specification in which all of the procedures are lossless (see Annex H).

**lossy:** A descriptive term for encoding and decoding processes which are not lossless.

**marker:** A two-byte code in which the first byte is hexadecimal FF and the second byte is a value between 1 and hexadecimal FE.

**marker segment:** A marker and associated set of parameters.

**MCU-row:** The smallest sequence of MCU which contains at least one line of samples or one block-row from every component in the scan.

**minimum coded unit; MCU:** The smallest group of data units that is coded.

**modes (of operation):** The four main categories of image coding processes defined in this Specification.

**non-interleaved:** The descriptive term applied to the data unit processing sequence when the scan has only one component.

**parameters:** Fixed length integers 4, 8 or 16 bits in length, used in the compressed data formats.

**point transform:** Scaling of a sample or DCT coefficient.

**precision:** Number of bits allocated to a particular sample or DCT coefficient.

**predictor:** A linear combination of previously reconstructed values (in lossless mode coding).

**procedure:** A set of steps which accomplishes one of the tasks which comprise an encoding or decoding process.

**process:** See coding process.

**progressive (coding):** One of the DCT-based processes defined in this Specification in which each scan typically improves the quality of the reconstructed image.

**progressive DCT-based:** The mode of operation which refers to any one of the processes defined in Annex G.

**quantization table:** The set of 64 quantization values used to quantize the DCT coefficients.

**quantization value:** An integer value used in the quantization procedure.

**quantize:** The act of performing the quantization procedure for a DCT coefficient.

**restart interval:** The integer number of MCUs processed as an independent sequence within a scan.

**restart marker:** The marker that separates two restart intervals in a scan.

**sample:** One element in the two-dimensional array which comprises a component.

**sample-interleaved:** The descriptive term applied to the repetitive multiplexing of small groups of samples from each component in a scan in a specific order.

**scan:** A single pass through the data for one or more of the components in an image.

**scan header:** A marker segment that contains a start-of-scan marker and associated scan parameters that are coded at the beginning of a scan.

**sequential (coding):** One of the lossless or DCT-based coding processes defined in this Specification in which each component of the image is encoded within a single scan.

**sequential DCT-based:** The mode of operation which refers to any one of the processes defined in Annex F of ITU-T Rec.81 | ISO/IEC 10918-1.

**table specification data:** The coded representation from which the tables used in the encoder and decoder are generated and their destinations specified.

**transcoder:** A procedure for converting compressed image data of one encoder process to compressed image data of another encoder process.

**(uniform) quantization:** The procedure by which DCT coefficients are linearly scaled in order to achieve compression.

**upsampling (filter):** A procedure by which the spatial resolution of an image is increased (in hierarchical mode coding).

**vertical sampling factor:** The relative number of vertical data units of a particular component with respect to the number of vertical data units in the other components in the frame.

**zero byte:** The X'00' byte.

**zig-zag sequence:** A specific sequential ordering of the DCT coefficients from (approximately) lowest spatial frequency to highest.

## 3.2 Symbols

X Width of the sample grid in positions

Y	Height of the sample grid in positions
Nf	Number of components in an image
$s_{i,x}$	Subsampling factor of component i in horizontal direction
$s_{i,y}$	Subsampling factor of component i in vertical direction
$H_i$	Subsampling indicator of component i in the frame header
$V_i$	Subsampling indicator of component i in the frame header
$v_{x,y}$	Sample value at the sample grid position x,y

### 3.3 Abbreviations

For the purposes of this Recommendation | International Standard, the following abbreviations apply.

ASCII	American Standard Code for Information Interchange
DC	Lowpass
AC	Highpass
LSB	Least Significant Bit
MSB	Most Significant Bit
HDR	High Dynamic Range
LDR	Low Dynamic Range
TMO	Tone Mapping Operator
DCT	Discrete Cosine Transformation

## 4 Conventions

### 4.1 Conformance language

This Recommendation | International Standard consists of normative and informative text.

Normative text is that text which expresses mandatory requirements. The word "shall" is used to express mandatory requirements strictly to be followed in order to conform to this Specification and from which no deviation is permitted. A conforming implementation is one that fulfils all mandatory requirements.

Informative text is text that is potentially helpful to the user, but not indispensable and can be removed, changed or added editorially without affecting interoperability. All text in this Recommendation | International Standard is normative, with the following exceptions: the Introduction, any parts of the text that are explicitly labeled as "informative", and statements appearing with the preamble "NOTE" and behavior described using the word "should". The word "should" is used to describe behavior that is encouraged but is not required for conformance to this Specification.

The keywords "may" and "need not" indicate a course of action that is permissible in a conforming implementation.

The keyword "reserved" indicates a provision that is not specified at this time, shall not be used, and may be specified in the future. The keyword "forbidden" indicates "reserved" and in addition indicates that the provision will never be specified in the future.

### 4.2 Operators

NOTE – Many of the operators used in this Recommendation | International Standard are similar to those used in the C programming language.



#### 4.2.1 Arithmetic operators

+	Addition
−	Subtraction (as a binary operator) or negation (as a unary prefix operator)
++	Increment by one as a unary postfix operator
--	Decrement by one as a unary postfix operator
*	Multiplication
/	Integer division, where the result is truncated towards zero.
÷	Division in mathematical equation where no truncation or rounding is intended
$\frac{x}{y}$	Division in mathematical equation where no truncation or rounding is intended
%	$x\%a$ is defined as the modulus operator for $x \geq 0, a > 0$ $x\%a = -((-x\%a))$ for $x < 0, a > 0$

#### 4.2.2 Logical operators

	Logical OR
&&	Logical AND
!	Logical NOT
$x ? a : b$	If $x$ is TRUE or not equal to 0, evaluates to the value of $a$ ; otherwise, evaluates to the value of $b$
$\in$	$x \in \{A, B\}$ is defined as $(x == A \parallel x == B)$
$\notin$	$x \notin \{A, B\}$ is defined as $(x != A \ \&\& \ x != B)$

TRUE/FALSE convention: In logical expressions, a variable or expression having a non-zero value is evaluated as TRUE and a variable or expression having a zero value is evaluated as FALSE.

#### 4.2.3 Relational operators

>	Greater than
>=	Greater than or equal to
<	Less than
<=	Less than or equal to
=	Equal to
!=	Not equal to

#### 4.2.4 Bit-wise operators

When operating on integer arguments, bit-wise operators operate on a two's complement representation of the integer value using a number of bits sufficient to represent the integer value (with a bit equal to 0 in the MSB of non-negative integer value representations and otherwise with a bit equal to 1 in the MSB). When interpreting the result of a bit-wise operator as an integer, the result is interpreted as a two's complement representation of the integer value. The following bit-wise operators are defined:

&	AND. When operating on integer arguments, operates on a two's complement representation of the integer value. When operating on a binary argument that contains fewer bits than the other argument, the shorter argument is extended by adding more significant bits equal to the MSB of the shorter argument.
	OR. When operating on integer arguments, operates on a two's complement representation of the integer values. When operating on a binary argument that contains fewer bits than the other argument, the shorter argument is extended by adding more significant bits equal to the MSB of the shorter argument.
^	XOR. When operating on integer arguments, operates on a two's complement representation of the integer value. When operating on a binary argument that contains fewer bits than the other argument, the shorter argument is extended by adding more significant bits equal to the MSB of the shorter argument.
x >> b	Arithmetic right shift of a two's complement integer representation of x by b binary digits. Bits shifted into the MSBs as a result of the right shift shall have a value equal to the MSB of x prior to the shift operation.
x << b	Arithmetic left shift of a two's complement integer representation of x by b binary digits. Bits shifted into the LSBs as a result of the left shift shall have a value equal to 0.

#### 4.2.5 Assignment operators

=	Assignment operator
+=	x += a is defined as x = x + a
-=	x -= a is defined as x = x - a
^=	x ^= a is defined as x = x ^ a
*=	x *= a is defined as x = x * a
<<=	x <<= a is defined as x = (x << a)
>>=	x >>= a is defined as x = (x >> a)

#### 4.2.6 Precedence order of operators

Operators are listed below in descending order of precedence. If several operators appear in the same line, they have equal precedence. When several operators of equal precedence appear at the same level in an expression, evaluation proceeds according to the associativity of the operator either from right to left or from left to right.

Operators	Type of operation	Associativity
() , [ ] , .	Expression	Left to Right
–	Unary negation	
*, /	Multiplication	Left to Right
mod	Modulo (remainder)	Left to Right
+, –	Addition and Subtraction	Left to Right
<, >, <=, >=	Relational	Left to Right
&,  , ^	Bitwise operator	Left to Right

#### 4.2.7 Mathematical functions

$\lceil x \rceil$	Ceil of $x$ . Returns the smallest integer that is greater than or equal to $x$ .
$\lfloor x \rfloor$	Floor of $x$ . Returns the largest integer that is lesser than or equal to $x$ .
$ x $	Absolute value, is $-x$ for $x < 0$ , otherwise $x$ .
$\text{sign}(x)$	Sign of $x$ , zero if $x$ is zero, +1 if $x$ is positive, -1 if $x$ is negative.
$\text{clamp}(x, \text{min}, \text{max})$	Clamps $x$ to the range $[\text{min}, \text{max}]$ : returns $\text{min}$ if $x < \text{min}$ , $\text{max}$ if $x > \text{max}$ or otherwise $x$ .

## 5 General

The purpose of this clause is to give an informative overview of the elements specified in this Specification. Another purpose is to introduce many of the terms which are defined in clause 3. These terms are printed in *italics* upon first usage in this clause.

There are three elements specified in this Specification:

- An *encoder* is an embodiment of an *encoding process*. An encoder takes as input *digital source image data* and *encoder specifications*, and by means of a specified set of *procedures* generates as output a *codestream*.
- A *decoder* is an embodiment of a *decoding process*. A decoder takes as input a *codestream*, and by means of a specified set of *procedures* generates as output *digital reconstructed image data*.
- The *codestream* is a compressed image data representation which includes all necessary data to allow a (full or approximate) reconstruction of the sample values of a digital image. Additional data might be required that define the interpretation of the sample data, such as color space or the spatial dimensions of the samples.

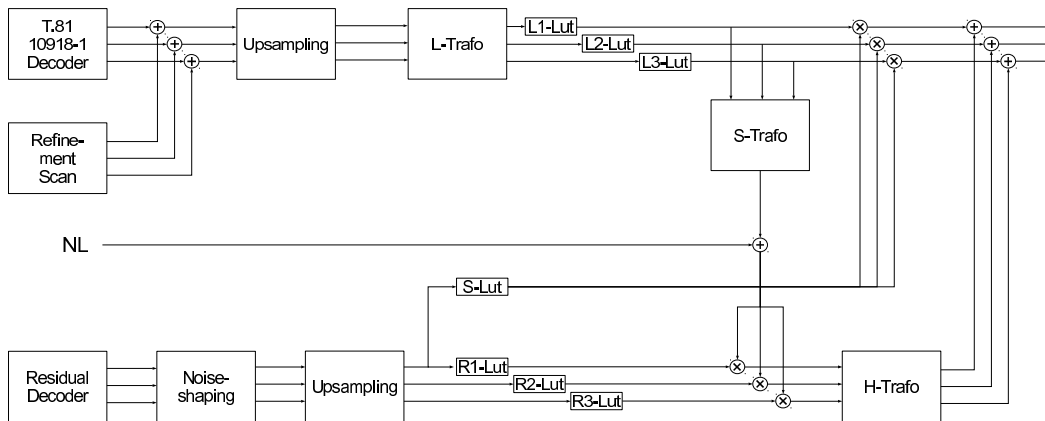


Figure 1 – Overview on the decoding process of high-dynamic range images

Figure 1 illustrate the general case for which the *continuous-tone* source and reconstructed image data consist of multiple *components*. (A *colour* image consists of multiple components; a *grayscale* image consists only of a single component.) A significant portion of this Specification is concerned with how to handle multiple-component images in a flexible, application-independent way.

## 5.1 Low and High Dynamic Range Images

This Specification distinguishes between low and high dynamic range images. The first type of images consists of unsigned integer samples of a sample precision of exactly eight bits per sample. Such images are represented in a codestream that is compatible to ITU-T Rec. T.81|ISO/IEC 10918-1, even though this Recommendation | International Standard consists only of a subset of the tools available in ITU-T Rec. T.81|ISO/IEC 10918-1. High dynamic range images consist of samples that are either floating point, or use bit precisions between nine and 16 bits per sample. Such images require additional tools defined in this Recommendation | International Standard. These additional tools are implemented in such a way that decoders conforming to ITU-T Rec. T.81|ISO/IEC 10918-1 are still able to decode a low-dynamic range version of the original image, but will not have access to additional data that comprises the original image.

Specifically, this Recommendation | International Standard depends only the DCT coding process defined in Section 4.3 of ITU-T Rec. T.81|ISO/IEC 10918-1, and the Baseline, Sequential and Progressive *Scan Types* defined in Annexes F and G thereof. Additional restrictions apply such as images in this Recommendation | International Standard may only use up to four components and chroma subsampling factors of one and two. The process of chroma subsampling is defined in detail in Annex A.

The additional tools defined here first generate a tone mapped version of the high dynamic range image and encode this image by means already defined in the earlier standard. The encoder also computes a *residual image* between the low dynamic and high dynamic range image and transmits this data, along with all necessary metadata in side channels. The side channel is build on top of *Application Markers* defined in Annex B of the earlier standard, and thus would remain invisible to legacy decoders.

## 5.2 Functional overview on the decoding process

A high-level overview on the decoding process is given by Figure 1: A low dynamic range representation of the image is represented by a ITU-T Rec. T.81|ISO/IEC 10918-1 compliant codestream is decoded into a low dynamic range representation (top left). An additional scan type, defined in Annex XXX may refine this data. The resulting samples are then upsampled on a common sample grid. The upsampling process is defined in Annex A. Following that, the output data is processed by an inverse decorrelation transformation, the L-Transformation, which is either the ICT, the FCT or the identity transformation. The inverse decorrelation transformation is defined in Annex C. The output of this transformation is then processed by an inverse tone mapping operation denoted by L1, L2 and L3, resulting in an approximate version of the HDR image. The B-Transformation computes the luminance of the approximate HDR image, resulting in a single value for each RGB-triplet. The B transformation can be either the forwards ICT transformation, or a null-transformation that outputs zero for any input. A noise-floor value, encoded in the JPEG Extensions markers, is added to the luminance value.

The residual image to the high-dynamic range image is represented either by an ITU-T-Rec.T.81|ISO/IEC 10918-1 codestream, or by a additional scan type denoted "Residual Scan" defined in Annex G. The reconstructed residual sample values are optionally processed by an inverse noise shaping algorithm removing quantization noise before performing an upsampling step that is identical to the low-dynamic range upsampling defined in Annex A. An additional inverse tone mapping step is applied to the residual signals as well, denoted in the figure by R1 through R3. The output luminance and chroma residuals are multiplied by the LDR luminance computed from the S-Transformation, and then undergo a second inverse decorrelation transformation called the H-Transformation. This transformation may either be the inverse ICT, the RCT or the identity transformation, all of which are defined in Annex C.

The residual luma signal is in addition mapped by a third tone mapping operation into a scale value by which the approximate high-dynamic range sample values are multiplied. Finally, the output of the H-transformation is added to the approximate high-dynamic sample values resulting in the decoded image.

## 5.3 Example Applications

A **purely additive** HDR decoding process can be implemented by setting the B transformation to the null-tansform, the noise floor value to one and the R1 through R3 transformations to identity. The S-tonemapper is in this application scenario a constant table that generates on the output one. The L1 through L3 tables are arbitrary and define here an inverse tone mapping process from the LDR to the HDR regime. If the H-transformation is identical the L-transformation, the residual stream represents a simple additive error residual that is added to the decoded, inversely tonemapped LDR signal. It is suggested to use in this application the residual scan defined in Annex XXX to encode the additive residual signal.

A **lossless or near lossless** image coder can be realized by additionally setting the H-transformation to the inverse RCT transformation. Furthermore, the FCT must be selected as L-transformation and the DCT in the 10918-1 decoder decoding the LDR signal must be the exact fixpoint DCT defined in Annex XXX. Since the RCT is exactly invertible, the R-tables are

identities and the residual scan does not create loss either, the error residual can be selected exactly to match the error signals of the fully specified forwards FCT and forwards fixpoint DCT transformation, resulting in a lossless image codec.

A **multiplicative** HDR image decoding process in which the residual image represents an exponent signal that scales the LDR image can be realized as follows: The L-transformation should here be selected as either the FCT or the ICT, the L1 through L3 curves implement here a mapping from sRGB to linear gamma and may be selected to be identical to the sRGB nonlinearities. The S-Transformation extracts the luminance of the approximate HDR signal and scales by that the chroma signals of the residual signal. The R1 through R3 curves should here be selected such that the luminance signal is completely suppressed, i.e. the luma R-table is completely zero, and the chroma curves may implement a simple linear function. The H-transformation is here identical to the inverse ICT transformation, creating additive chroma residuals. The S-curve implements in this setting an inverse logarithm, i.e. an exponential function reconstructing the exact scale of the approximate HDR signal. The advantage of this approach is that the residual image behaves close to a natural image and can be compressed by legacy 10918-1 encoders, the drawback is that lossless image compression cannot be implemented.

In addition to these three application profiles many other profiles can be configured. The parameters specifying the L,S and R transformations and curves are recorded in JPEG Extension markers invisible to legacy decoders and so is the residual signal and the refinement signal. A legacy decoder would only reconstruct the low dynamic range portion of the image recorded in a codestream that is compliant to the interchange format specified in ITU.T Rec.T 80|ISO/IEC 10918-1.

## 5.4 Encoder requirements

An encoding process converts source image data to compressed image data. This includes first generating a low dynamic range image, and representing it by a coding process specified in Annex F or Annex G of ITU.T Rec.T 81 | ISO/IEC 10918-1, and then generating a residual image which is encoded by one of the processes defined in this Recommendation|International Standard.

In order to comply with this Specification, an encoder shall satisfy at least one of the following two requirements. An encoder shall with appropriate accuracy, convert source image data to compressed image data which comply with the codestream format syntax specified in Annex B for the encoding process(es) embodied by the encoder. Compliance tests for the above requirements are specified in this Recommendation | International Standard.

NOTE – There is **no requirement** in this Specification that any encoder which embodies one of the encoding processes specified here shall be able to operate for all ranges of the parameters which are allowed for that process. An encoder is only required to meet the compliance tests and to generate the compressed data format according to Annex B for those parameter values which it does use.

## 5.5 Decoder requirements

A decoding process converts compressed image data to reconstructed image data. For that, it has to follow the decoding operation specified in ITU.T Rec.T 81 | ISO/IEC 10918-1, using either the Baseline, Sequential or Progressive Scan process of the earlier Recommendation | International Standard defined there in Annexes F and G, and merge the resulting sample values with decoded residual values reconstructed from data included in the side channel represented by one or multiple application markers. The actual merging process is driven by parameters also recorded in Application Markers, the format of which is specified in this Recommendation | International Standard.

In order to comply with this Specification, a decoder shall satisfy all three of the following requirements. A decoder shall

- a) with appropriate accuracy, convert a codestream conforming to this Recommendation | International Standard **without considering the side channel** into a low dynamic range digital image.
- b) **may additionally**, reconstruct with appropriate accuracy, convert a conforming codestream including the side information included in Application Markers, into a high dynamic range digital image.

•

### Annex A

## Component Subsampling and Expansion of Subsampling

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

Kommentar [tr1]: This should probably be moved into one of the additional parts as subsampling expansion was not really part of the original JPEG, though it is part of any (software!) decoder I'm aware of. Walt would prefer to exclude it in part-1. I'm personally open for all options.

## 5.6 Component dimensions and subsampling factors (normative)

An image is defined to consist of  $N_f$ , each of which is identified by a unique identifier  $C_i$  defined in the frame header of the of the codestream format specified in Annex B, where  $N_f$  shall be either one or three. A component consists of a rectangular array of samples  $x_i$  wide and  $y_i$  samples high. The component dimensions are derived from the image dimensions  $X$  and  $Y$ , also parameters recorded in the frame header. These two parameters define a sample grid of  $X$  grid points wide and  $Y$  grid points high; however, not every grid point need to carry a sample value. For each component, two subsampling factors  $s_{i,x}$  and  $s_{i,y}$  define the spacing between sample points of component  $i$  relative to the sample grid. Namely only those sample grid points whose horizontal coordinate is an integer multiple of  $s_{i,x}$  and whose vertical coordinate  $s_{i,y}$  carry sample values that are represented in the codestream, where the left topmost grid coordinate is  $(0,0)$  is  $(0,0)$  and coordinates increase from left to right and from top to bottom. Sample grid points whose sample value is not recorded in the codestream shall to be interpolated from their surrounding samples by mechanisms specified in subclause subclause A.2.]

The subsampling factors  $s_{i,x}$  and  $s_{i,y}$  are not directly represented in the binary codestream or any of its markers, but shall be derived from the parameters  $H_i$  and  $V_i$  recorded in the frame header. If  $N_f$  equals one, i.e. the image consists of a single component,  $H_i$  and  $V_i$  shall be one, and  $s_{i,x}$  and  $s_{i,y}$  are both one. If  $N_f$  equals three, Table A-1 defines the relation between  $H_i$ ,  $V_i$  and  $s_{i,x}$  and  $s_{i,y}$ . No other combinations of  $H_i$  and  $V_i$  than those listed in this table shall be used.

$H_i$	$V_i$	$H_2$	$V_2$	$H_3$	$V_3$	$s_{1,x}$	$s_{1,y}$	$s_{2,x}$	$s_{2,y}$	$s_{3,x}$	$s_{3,y}$
1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	1	2	1	1	1	1	2	1	2
2	2	1	2	1	2	1	1	2	1	2	1
2	2	1	1	1	1	1	1	2	2	2	2
All other values reserved for ITU/ISO purposes											

Table A-1: Subsampling Values

**NOTE** – ISO/IEC 10918-1 allowed more complicated component arrangements and relations between grid positions and sample positions that are not allowed in this Recommendation | International Standard. However, the definitions given here are special cases of the more general relations provided in 10918-1 and both definitions agree whenever both are defined.

## 5.7 Expansion of subsampled components (normative)

Whenever the subsampling factors  $s_{i,x}$  and  $s_{i,y}$  are not both one, sample values at sample grid positions not populated by actual samples shall be interpolated from surrounding positions. The precise algorithm how this interpolation has to be performed is not specified in this Recommendation | International Standard and various implementation choices exist. This subclause provides an example for a simple bi-linear interpolation that is sufficient for most applications, though more sophisticated choices might provide better subjective image quality.

## 5.8 Bilinear expansion of subsampled components (informative)

In a first step, check for each component  $i$  whether  $s_{i,y}$  is two or one. If  $s_{i,y}$  is one, perform no activity. If  $s_{i,y}$  equals two, interpolate samples  $v$  at odd lines from even lines as follows:

$$v_{x,2y+1} = \lfloor (v_{x,2y} + v_{x,2y+2} + y \bmod 2) / 2 \rfloor$$

If  $2y+2$  is larger or equal  $Y$ , the horizontal sample grid dimension, set  $v_{x,2y+1}$  equal to  $v_{x,2y}$ .

In a second step, check for each component  $i$  whether  $s_{i,x}$  is two or one. If  $s_{i,x}$  is one, no further activity is required and the algorithm terminates. Otherwise, if  $s_{i,x}$  is equal to two, interpolate samples at odd positions from samples at even positions accordingly:

$$v_{2x+1,y} = \lfloor (v_{2x,y} + v_{2x+2,y} + x \bmod 2) / 2 \rfloor$$

Again, if  $2x+2$  is larger or equal than  $X$ , the vertical sample grid dimension, set  $v_{2x+1,y}$  to  $v_{2x,y}$ .

Kommentar [tr2]: This is actually not how JPEG was originally designed. Samples are not co-located as in JPEG 2000. It is unfortunately how one popular implementation works.

Kommentar [tr3]: Do we actually need this?

Kommentar [tr4]: Ditto – Walt would prefer to not include it. No problem with me.

Kommentar [tr5]: Actually, the algorithm shown here is probably not the best given that samples in JPEG are actually not co-located, but centered. It is however the method implemented by IJG. Any better options?

## Annex B

### Codestream Syntax

(This annex forms an integral part of this Recommendation | International Standard)

This annex defines the compressed bitstream syntax which, structurally, consists of an ordered collection of marker segments and entropy coded data segments. Marker segments specify parameters necessary to reconstruct the sample values from the entropy coded data segments. Because all of these constituent parts are represented with byte-aligned codes, each compressed data format consists of an ordered sequence of 8-bit bytes. For each byte, a most significant bit (MSB) and a least significant bit (LSB) are defined.

NOTE – The codestream syntax defined here agrees mostly with the "interchange format" defined in ITU-T Rec. T.81 | ISO/IEC 10918-1, with some additional constraints on the parameters in the marker segments and some additional markers carrying information that is irrelevant for the older standard. Such additional data is encapsulated in a way that will be ignored by decoders complying to the older standard.

#### 5.9 Parameters

Parameters are integers, with values specific to the encoding process, source image characteristics, and other features selectable by the application. Parameters are assigned either 4-bit, 1-byte, or 2-byte codes. Except for certain optional groups of parameters, parameters encode critical information without which the decoding process cannot properly reconstruct the image. The code assignment for a parameter shall be an unsigned integer of the specified length in bits with the particular value

of the parameter.

For parameters which are 2 bytes (16 bits) in length, the most significant byte shall come first in the compressed data's ordered sequence of bytes. Parameters which are 4 bits in length always come in pairs, and the pair shall always be encoded in a single byte. The first 4-bit parameter of the pair shall occupy the most significant 4 bits of the byte. Within any 16-, 8-, or 4-bit parameter, the MSB shall come first and LSB shall come last. This encoding is commonly known as "big endian" representation of unsigned integers.

#### 5.10 Markers

Markers serve to identify the various structural parts of the compressed data formats. Most markers start marker segments

containing a related group of parameters; some markers stand alone. All markers are assigned two-byte codes: an 0xff byte followed by a byte which is not equal to 0x00 or 0xff. Any marker may optionally be preceded by any number of fill bytes, which are bytes assigned code 0xff.

NOTE – Because of this special code-assignment structure, markers make it possible for a decoder to parse the compressed data and locate its various parts without having to decode other segments of image data.

#### 5.11 Marker assignments

All markers shall be assigned two-byte codes: a 0xff byte followed by a second byte which is not equal to 0x00 nor 0xff. The second byte is specified in Table B-1 for each defined marker. An asterisk (\*) indicates a marker which stands alone, that is, which is not the start of a marker segment. Most of the marker segments used by this Recommendation | International Standard are defined in ITU-T Rec.T.81 | ISO/IEC 10918-1, though some marker segments defined there shall not be used in this Recommendation|International standard. For completeness, these markers are also included in Table B-1. Furthermore, care must be taken that the parameters of some markers are more constraint than in the earlier Recommendation | International Standard, and the meaning of several markers changed slightly. Such changes are also indicated in the table, and a full specification of the changes follows.

Code Assignment	Symbol	Description	Defined in
Start of Frame markers valid in this Recommendation   International Standard			
0xFFC0	SOF <sub>0</sub>	Baseline DCT Process	This Recommendation   Internatinonal Standard
0xFFC1	SOF <sub>1</sub>	Extended Sequential DCT Process	

0xFFC2	SOF <sub>2</sub>	Progressive DCT Process	
Start of Frame markers defined in T.81 10918-1 that shall not be used in this Recommendation   International Standard			
0xFFC3	SOF <sub>3</sub>	Additional Start of Frame Markers, shall not be used.	ITU.T Rec.T81   ISO/IEC 10918-1
0xFFC5	SOF <sub>5</sub>		
0xFFC6	SOF <sub>6</sub>		
0xFFC7	SOF <sub>7</sub>		
0xFFC9	SOF <sub>9</sub>		
0xFFCA	SOF <sub>10</sub>		
0xFFCB	SOF <sub>11</sub>		
0xFFCD	SOF <sub>12</sub>		
0xFFCE	SOF <sub>14</sub>		
0xFFCF	SOF <sub>15</sub>		
Table Specifications			
0xFFC4	DHT	Define Huffman Tables	ITU.T Rec.T81   ISO/IEC 10918-1
0xFFDB	DQT	Define Quantization Tables	ITU.T Rec.T81   ISO/IEC 10918-1
Restart Interval Termination			
0xFFD0 through 0xFFD7	RST <sub>m</sub> <sup>*</sup>	Restart modulo 8 count 'm'	ITU.T Rec.T81   ISO/IEC 10918-1
0xFFDD	DRI	Define Restart Interval	ITU.T Rec.T81   ISO/IEC 10918-1
Other Markers			
0xFFD8	SOI <sup>*</sup>	Start of Image	ITU.T Rec.T81   ISO/IEC 10918-1
0xFFD9	EOI <sup>*</sup>	End of Image	
0xFFDA	SOS	Start of Scan	
0xFFFE	COM	Comment	
JPEG Extension Markers			
0xFFE0	APP <sub>0</sub>	JFIF Application Marker	ISO/IEC 10918-5 and ISO/IEC 10918-4
0xFFE1	APP <sub>1</sub>	EXIF Application Marker	JEITA CP-3451 and ISO/IEC 10918-4
0xFFE2	APP <sub>2</sub>	ICC Profile Marker	ISO/IEC 10918-6 and ISO/IEC 10918-4
0xFFE3 through 0xFFE8	APP <sub>3</sub> through APP <sub>8</sub>	Application Markers 3 through 8	ISO/IEC 10918-4
0xFFE9	APP <sub>9</sub>	JPEG Extensions Marker	This Recommendation   International Standard
0xFFEA	APP <sub>10</sub>	Application Marker 10	ISO/IEC 10918-4
0xFFEB	APP <sub>11</sub>	JPEG High Dynamic Range Extensions Marker	This Recommendation   International Standard
0xFFEC	APP <sub>12</sub>	Application Marker 12	ISO/IEC 10918-4
0xFFED	APP <sub>13</sub>	Component Decorrelation Control Marker	This Recommendation   International Standard
0xFFEE through 0xFFEF	APP <sub>14</sub> through APP <sub>15</sub>		



Markers defined in other specifications that shall not be used in this Recommendation   International Standard			
0xFFC8	JPG	Reserved for JPEG Extensions	ITU.T Rec.T81   ISO/IEC 10918-1
0xFFCC	DAC	Define Arithmetic Coding Conditions	
0xFFDC	DNL	Define Number of Lines	
0xFFDE	DHP	Define Hierarchical Progression	
0xFFDF	EXP	Expand Reference Components	
All other values		Reserved for ITU/ISO purposes, shall not be used in this Recommendation   International Standard	

**Table B-1: Markers and Marker Segments**

## 5.12 Marker segments

A marker segment consists of a marker followed by a sequence of related parameters. The first parameter in a marker segment is the two-byte length parameter. This length parameter encodes the number of bytes in the marker segment, including the length parameter and excluding the two-byte marker. The marker segments identified by the SOF and SOS marker codes are referred to as headers: the frame header and the scan header respectively.

## 5.13 Entropy-coded data segments

An entropy-coded data segment contains the output of an entropy-coding procedure. It consists of an integer number of bytes created by the Huffman coding procedure. An entropy coded data segment stands either for itself and is then started by an SOS marker, or it is contained in the APP<sub>9</sub> or APP<sub>11</sub> JPEG Extensions marker. In the latter cases, the size of the marker segment containing the entropy coded data includes the size of the entropy coded data, the size of additional parameters or marker segments included in the APP<sub>9</sub> or APP<sub>11</sub> marker segment and the size of the 16-bit marker size. The marker itself is not included in the count.

Entropy coded segments shall be always a multiple of eight bits long, and shall be padded by one bits at the end of the segment to fill an entire byte. In order to ensure that a marker does not occur within an entropy-coded segment, the encoder shall insert a zero byte into the output bitstream following any any 0xff byte generated by either the Huffman encoder or by the above termination algorithm. Decoders shall remove such "stuffed" zero bytes before feeding the resulting bitstream into the Huffman decoder.

## 5.14 High-Level Syntax

The high-level syntax of the codestream defined in this Recommendation|International Standard shall follow the syntax of the Interchange Format specified in Subclauses B.2 and B.3 of ITU.T Rec.T 81 | ISO/IEC 10918-1 where only the subset of the markers specified in Table B-1 shall be used.

## 5.15 Frame header syntax

Table B-2 specifies the frame header which shall be present at the start of a frame. This header specifies the dimensions of the image, the components in the frame, and the sampling factors for each component, and specifies the destinations from which the quantized tables to be used with each component are retrieved.

Parameter	Size (in Bits)	Value	Meaning
Lf	16	8+3*Nf	Size of the marker segment

			not including the marker
P	8	8	Precision of the low dynamic range representation of the image contained in the legacy codestream
Y	16	1-65535	Height of the sample grid in lines
X	16	1-65535	Width of the sample grid in lines
Nf	8	1 or 3	Number of visible components in the image
C <sub>i</sub>	8	0-255	Component identifiers
H <sub>i</sub>	4	1 or 2	Specification of the component subsampling factors, see Table A-1 for allowable values and the relation between H <sub>i</sub> V <sub>i</sub> and s <sub>i,x</sub> and s <sub>i,y</sub>
V <sub>i</sub>	4	1 or 2	
Tq <sub>i</sub>	8	0-3	Quantization Table Destination selector

**Table B-2: Frame Header Parameters and values**

NOTE – The frame header syntax defined here is identical to the syntax defined in ITU.T Rec.T 81| ISO/IEC 10918-1 except that some parameters are restricted. Specifically, the bit precision of the components must always be eight, the height must be greater than zero indicating that the number of lines of the image is known, the number of components must be either one or three and the horizontal and vertical sampling factors are constraint to the values listed in table A-1.

## 5.16 Scan Header Syntax

The scan header defines the parameters necessary to decode a single scan over the image. For the baseline or sequential frame syntax, this is the only scan that represents the low dynamic range representative of the image. Multiple scans are possible in the progressive mode. Residual data necessary to recover the high dynamic range representation of the image is recorded in one or several JPEG extensions markers; encoding or decoding this image representation may require additional scans over the image that are not represented by the scan header defined in this subclause.

The syntax of the scan header follows the definition for the same marker in ITU.T Rec. T81 | ISO/IEC 10918-1, though some of the parameters are here more constrained than in the older Recommendation | International Standard. The parameters and sizes of the Scan header are defined in Table B-3.

Parameter	Size (in Bits)	Values			Meaning
		Baseline	Extended Sequential	Progressive	
Ls	16	6+2*Ns			Size of the marker segment (not including the marker)
Ns	8	1 or 3		1 or 3 (if the progressive progressive scan is a subsequent scan and includes	Number of components in the scan

			includes AC components, only only 1 is allowed here)	
$Cs_j$	8	0-255 (shall be a subset of the component identifiers defined in the frame header)		Component identifiers of the components defined by the scan
$Td_j$	4	0-1	0-3	Huffman DC table specification
$Ta_j$	4	0-1	0-3	Huffman AC table specification
Ss	8	0	0-63	Start of Spectral Selection
Se	8	63	Ss-63	End of Spectral Selection
Ah	4	0	0-9 If nonzero, must be equal to Al+1	Successive Approximation high bit position
Al	4	0	0-8	Successive Approximation low bit position

**Table B-3: Scan Header Parameters and Values**

NOTE – The scan header syntax defined here is identical to the syntax defined in ITU-T Rec.T 81| ISO/IEC 10918-1, though a scan may contain at most three components.

### 5.17 Component Decorrelation Control marker

The APP<sub>13</sub> controls whether the Multi-Component Decorrelation Transformation described in Annex C is applied to the components before subsampling and entropy coding. This transformation typically improves the compression performance by representing data in an RGB type colorspace in a luma/chroma format that is close to YCbCr. This marker shall only be inserted into the codestream as parts of the Table/misc segment if the number of components in the frame (Nf) is three. The marker parameters and sizes are defined in Table B-4, the organization of the marker is shown in Figure B-1.

NOTE – This marker is identical to the "Adobe Color Management" marker defined in ISO/IEC 10918-6, though the possible values the parameters inside the marker segment are more restricted. The intended use is also slightly different: Whereas the marker in ISO/IEC 10918-6 defines an output color space of the encoded samples, it controls here the decorrelation transformation from the internal representation to the external color format, though not necessarily the color space itself. Legacy software, however, will assume an RGB-type output colorspace in the absence of any additional information and will implement a color transformation close to the decorrelation transformation specified in Annex C, though with the intent to change the colorspace and not to undo a component decorrelation transformation. The presence of this marker, with suitable parameters, will then suppress this change in the colorspace, and thus implement the desired functionality, namely to disable the multi-component decorrelation transformation.

0xFFED	La	ACid	ver	f1	f2	cc
--------	----	------	-----	----	----	----

**Figure B-1 : Organization of the Component Decorrelation Control marker**

Parameter	Size (bits)	Value	Meaning
-----------	-------------	-------	---------

Kommentar [tr6]: This should possibly not go into part-1, even though all software implementations I am aware of support it, and support it exactly in the way described here. Walt would be happy not to, I do not mind too much, though I believe it would make sense here on the rationale of standardizing a living specification. Maybe make it informative here and normative later.

There is one corner point, namely if both this marker and the JFIF marker is present. Then, it seems, IJG ignores this marker. Should be specify this corner case?

La	16	14	Size of the marker segment (not including the marker)
ACid	40	0x41 0x64 0x6f 0x62 0x65 (ASCII encoding of "Adobe")	Identifies the use of the APP <sub>13</sub> marker segment for the purpose of the Component Decorrelation Control marker. Readers shall ignore the APP <sub>13</sub> marker for the purpose of color decorrelation if this value does not match.
ver	16	0x64	A version number. Only version 100 is defined.
f1	16	0	Flags 1. Only the value 0 is defined in this Recommendation   International Standard.
f2	16	0	Flags 2. Only the value 0 is defined in this Recommendation   International Standard.
cc	8	0 or 1	Color Decorrelation Control Byte. If this value is zero, the color decorrelation transformation(s) specified in Annex C shall not be applied. If this value is one, or the marker is not present, the color decorrelation transformation shall be employed.  All other values are reserved and their meaning is not defined by this Recommendation   International Standard.

**Table B-4: Component Decorrelation Control marker parameters and sizes**

### 5.18 JPEG Extensions Marker Segment

The APP<sub>9</sub> marker is reserved by this Recommendation | International Standard to include additional parameters and entropy coded data that remains invisible to legacy decoders. All extensions share a common syntax that is defined in this and the following subclauses.

A JPEG Extensions marker segment consists of the marker, the marker length that includes the marker length, all parameters and entropy coded data but not the marker itself, a two-byte common identifier and a four byte type field, plus all parameters and data encapsulated by the marker. Figure B-2 shows the organization of the marker, Table B-4 defines the fields, the field lengths and the contents

0xFFE9	Le	Common Identifier	Type Identifier	Payload Data
--------	----	-------------------	-----------------	--------------

**Figure B-2: Organization of the JPEG Extensions Marker Segment**

Parameter	Size (bits)	Value	Meaning
Le	16	8..65535	Length of the marker segment, including the size, the common identifier, the type identifier and the payload, but not including the marker itself.
CI	16	0x4A50 (ASCII encoding of "JP")	The special value 0x4A50 (ASCII: 'J' 'P') allows readers to distinguish the JPEG Extensions marker segment from other uses of the APP <sub>0</sub> marker. Readers shall ignore the APP <sub>0</sub> marker for the purpose of encoding JPEG extensions if this value does not match.
TI	32	See subclauses B.XX through B.YY	The Type Identifier describes the type of the payload data that follows. Subclauses B.XX through B.YY define the payload data used in this Recommendation/International Standard.
Data	Varies	Varies	Defined in subclauses B.XX through B.YY.

**Table B-5: JPEG Extensions Marker Parameters and Sizes**

#### 5.18.1 JPEG Extensions Marker Segment: Entropy Coded Data for Refinement Coding

Entropy coded segments from refinement coding are encapsulated in this variant of the JPEG Extensions marker. This data extends the precision of the DCT transformed coefficients by additional least significant bits that improve the precision of the DCT coefficients beyond the sample depths available for legacy decoders. To form the entropy coded bitstream, the contents of the Entropy Coded data from all JPEG Extensions markers of this type shall be concatenated by a decoder in the order the markers appear in the codestream, and then form an input for the refinement entropy decoding process defined in Annex XX.

The structure of this marker is defined in Figure B-3, the parameters and sizes in Table B-6.

0xFFE9	Le	Common Identifier	Type Identifier	Entropy Coded Data
--------	----	-------------------	-----------------	--------------------

**Figure B-3: Organization of the JPEG Extensions Marker Segment for Refinement Coding**

Parameter	Size (bits)	Value	Meaning
Le	16	8..65535	Length of the marker segment, including the size, the common identifier, the type identifier and the payload, but not including the marker itself.
CI	16	0x4A50	Identifies this marker segment

		(ASCII encoding of "JP")	as JPEG Extensions marker
TI	32	0x46494e45 (ASCII encoding of "FINE")	Identifies this marker as carrying entropy coded data for refinement coding.
Data	Varies	Varies	Entropy coded data segment of variable lengths.

**Table B-6: JPEG Extensions Marker for Refinement Coding, Parameters and Sizes**

#### 5.18.2 JPEG Extensions Marker Segment: Entropy Coded Data for Residual Coding

Entropy coded segments from residual coding are encapsulated in this variant of the JPEG Extensions marker. This data provides entropy coded data that, when decoded and merged with the low-dynamic range data described by the data following the SOS marker, describes a high-dynamic range image. The algorithm for decoding residual data is defined in Annex XX, the operation for merging the decoded low-dynamic range data with the high-dynamic range data in Annex YY.

To form the entropy coded bitstream, the contents of the Entropy Coded data from all JPEG Extensions markers of this type shall be concatenated by a decoder in the order the markers appear in the codestream, and then form an input for the residual entropy decoding process defined in Annex XX.

The structure of this marker is defined in Figure B-4, the parameters and sizes in Table B-7.

0xFFE9	Le	Common Identifier	Type Identifier	Entropy Coded Data
--------	----	-------------------	-----------------	--------------------

**Figure B-3: Organization of the JPEG Extensions Marker Segment for Residual Coding**

Parameter	Size (bits)	Value	Meaning
Le	16	8..65535	Length of the marker segment, including the size, the common identifier, the type identifier and the payload, but not including the marker itself.
CI	16	0x4A50 (ASCII encoding of "JP")	Identifies this marker segment as JPEG Extensions marker
TI	32	0x52455349 (ASCII encoding of "RESI")	Identifies this marker as carrying entropy coded data for residual coding.
Data	Varies	Varies	Entropy coded data segment of variable lengths.

**Table B-6: JPEG Extensions Marker for Residual Coding, Parameters and Sizes**

#### 5.18.3 JPEG Extensions Marker Segment: Inverse Tone Mapping Marker

This marker segment defines an inverse tone mapping curve in the form of a look-up table. This look-up table is used to map the low-dynamic range image encoded by the entropy coded data following the SOS markers into the high-dynamic range regime and by that implements an inverse tone mapping. The inverse tone mapping operation is part of

the merging process combining the low-dynamic and residual image information to reconstruct a high-dynamic range image; it is described in more detail in Annex XX.

The marker segment organization is defined in Figure B-4, the parameters and sizes in Table B-7.

0xFFE9	Le	Common Identifier	Type Identifier	M	h	D <sub>i</sub>
--------	----	-------------------	-----------------	---	---	----------------

**Figure B-3: Organization of the JPEG Extensions Marker: Inverse Tone Mapping Curve**

Parameter	Size (in Bits)	Value	Meaning
Le	16	$9+2 \cdot 2^{h+8}$	Length of the marker segment (not including the marker)
CI	16	0x4A50 (ASCII encoding of "JP")	Identifies this marker segment as JPEG Extensions marker
TI	32	0x544f4e45 (ASCII encoding of "TONE")	Identifies this marker as carrying entropy coded data for residual coding.
M	4	0..15	Tone mapping table identifier. Up to 16 tone mapping curves can be defined.
h	4	0..4	Number of hidden DCT bits represented by refinement coding.
D <sub>k</sub>	16	0..65535	Output high dynamic range sample value for low dynamic range sample value $i$ . The marker includes $2^{h+8}$ table entries each 16 bits wide.

**Table B-7: JPEG Extensions Marker containing an Inverse Tone Mapping Curve**

#### 5.18.4 JPEG Extensions Marker Segment: Specifications of the Merging Process

This type of the JPEG extensions marker segment defines the process and parameters for merging low dynamic range image with the residual image to reconstruct the high dynamic range image. The details of this merging process is defined in Annex XX. Figure B-4 describes the organization of this marker segment, Table B-8 the parameters and parameter sizes.

0xFFE9	Le	CI	TI	Te	Rb	rI	Ct	S	Lt	h	m
--------	----	----	----	----	----	----	----	---	----	---	---

**Figure B-4: Organization of the JPEG Extensions Marker: Specification of the Merging Process**

Parameter	Size (in bits)	Value	Meaning
-----------	----------------	-------	---------

Le	16	12..16	Length of the marker segment (not including the marker)
CI	16	0x4A50 (ASCII encoding of "JP")	Identifies this marker segment as JPEG Extensions marker
TI	32	0x53504543 (ASCII encoding of "SPEC")	Identifies this marker as carrying entropy coded data for residual coding.
Te <sub>i</sub>	1	0 or 1	<p>Inverse Tone Mapping Enable. Exactly four Te<sub>i</sub> parameters shall be present. If the i<sup>th</sup> component carries high dynamic range data, Te<sub>i</sub> shall be 1, otherwise 0.</p> <p>Note that i is the absolute component number, not the component identifier C<sub>i</sub>.</p>
Rb	4	0..7	Number of additional bits available for high dynamic range data. The bit precision of the high dynamic range data shall be computed as 8+Rb.
rl	2	0	Reserved for ITU/ISO purposes. Shall be set to zero.
Ct	2	0..3	Identifies the index of the quantization table that shall be used for dequantizing the chroma components of the residual data.
S	2	0 or 1	Enables or disables noise shaping of the residual data. Noise shaping is enabled if this parameter is 1, disabled if the parameter is zero. See Annex XXX for details.
Lt	2	0..3	Identifies the index of the quantization table that shall be used for dequantizing the luma components of the residual data.
h	8	0..4	Number of hidden DCT bits bits represented by refinement coding. The value of h here shall be



			identical to the h parameter in all 'TONE' JPEG Extension markers.
$m_i$	8	0..15	For every nonzero $T_e$ , bit a $m_i$ parameter shall be present in the marker segment. The parameter $m_i$ selects for component i one out of sixteen inverse tone mapping curves defined by the 'TONE' JPEG Extension marker.

**Table B-8: JPEG Extensions Marker containing an Inverse Tone Mapping Curve**

## Annex C

### Multi-Component Decorrelation

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

This Annex describes several multiple component transformations. These multiple component transformations are used to improve compression efficiency. They are not related to multiple component transformations used to map colour values for display purposes. One multiple component transformation may be used for lossy or lossless coding. The other is an irreversible approximation of the first and may only be used for lossy coding, though may be more efficient to implement. The third transformation is an exactly reversible transformation that, depending on the settings defined in the codestream, acts on the residual data only; it is used for lossy and lossless encoding in one of the profiles of this Recommendation | International Standard. The last transformation is an additional lossy transformation that is required in one of the profiles for lossy reconstruction of high dynamic range images.

#### 5.19 Irreversible inverse multi component transformation (ICT) (normative)

This transformation is applied for reconstructing a low-dynamic range version of the image from the data encoded in the entropy coded data segment following SOS markers. It shall be used only if the number of components in the frame,  $N_f$ , is equal to three, and either the component decorrelation control marker, see Annex B, is not present or the control value  $cc$  in this marker is set to one. In case the number of components is one or a component decorrelation control marker with the  $cc$  parameter set to zero is present, this transformation shall be replaced by clamping the reconstructed data to the range  $[0, 2^{8+h}-1]$  without applying any further transformation, where  $h$  is the number of hidden DCT bits defined in the JPEG Extensions Specification marker.

In the following, set  $I_0, I_1$  and  $I_2$  be the sample values reconstructed from the upsampled components 0, 1 and 2, and let  $L_0, L_1$  and  $L_2$  be the sample values of the reconstructed low dynamic range version of the image. Let  $h$  be the number of hidden DCT bits as defined by the  $h$  parameter of the JPEG Extensions Specification ('SPEC') marker, or zero if this marker is not present. Then:

$$\begin{aligned} L_0 &= \text{clamp}(I_0 + 1.40200 * (I_2 - 2^{7+h}), 0, 2^{8+h} - 1) \\ L_1 &= \text{clamp}(I_0 - 0.7141362859 * (I_2 - 2^{7+h}) - 0.3441362861 * (I_1 - 2^{7+h}), 0, 2^{8+h} - 1) \\ L_2 &= \text{clamp}(I_0 + 1.772 * (I_1 - 2^{7+h}), 0, 2^{8+h} - 1) \end{aligned}$$

NOTE – This transformation is intentionally identical to the YUV to RGB transformation specified in ISO/IEC 10918-6, though its purpose is here slightly different.

#### 5.20 Inverse fixpoint multi component transformation (FCT) (normative)

This transformation shall be applied in the lossless process for reconstructing a low-dynamic range version of the image from the data encoded in the entropy coded data segment following SOS markers. It shall only be used if the number of components in the frame,  $N_f$ , is equal to three, and either the component decorrelation control marker, see

**Kommentar [tr7]:** All this comes from the Stuttgart implementation. Not everything here is "nice", I'm not so happy about SPEC, but it works. TONE and RESI should possibly remain as is, but should definitely go into a separate part.

FINE and the "h" field are debatable – there are maybe some applications where this makes sense, but it typically buys little.

**Kommentar [tr8]:** This should go, at least in parts, into a separate part. What might remain in part-1 is probably the ICT.

Annex B, is not present or the control value cc in this marker is set to one. In case the number of components is one or a component decorrelation control marker with the cc parameter set to zero is present, this transformation shall be replaced by rounding the fixpoint data to the nearest integer and clamping the result to the range  $[0, 2^{8+h}-1]$ .

This transformation assumes that the reconstructed untransformed sample values are the output of the fixpoint DCT approximation defined in Annex G and thus are integer values that are premultiplied with the factor 16 – or equivalently, the input data to this transformation is given in fixpoint with four fractional bits.

Implementations that do not desire to implement the lossless process may either use the irreversible inverse multiple component transformation, or this transformation which is a fixpoint approximation of the former transformation. The fixpoint inverse multi component transformation and the irreversible inverse multi component transformation are identical within the desired implementation precision.

In the following, set  $I_0, I_1$  and  $I_2$  be the sample values reconstructed from the upsampled components 0, 1 and 2, and let  $L_0, L_1$  and  $L_2$  be the sample values of the reconstructed low dynamic range version of the image that are preshifted by 4 bits, i.e. premultiplied by 16. Let  $h$  be the number of hidden DCT bits as defined by the  $h$  parameter of the JPEG Extensions Specification ("SPEC") marker, or zero if this marker is not present. Then:

$$\begin{aligned} L_0 &= \text{clamp}(\lfloor (2^{13} * I_0 + 11495 * (I_2 - 2^{7+4+h}) + 2^{16}) / 2^{17} \rfloor, 0, 2^{8+h} - 1) \\ L_1 &= \text{clamp}(\lfloor (2^{13} * I_0 - 5850 * (I_2 - 2^{7+4+h}) - 2819 * (I_1 - 2^{7+4+h}) + 2^{16}) / 2^{17} \rfloor, 0, 2^{8+h} - 1) \\ L_2 &= \text{clamp}(\lfloor (2^{13} * I_0 + 14516 * (I_1 - 2^{7+4+h}) + 2^{16}) / 2^{17} \rfloor, 0, 2^{8+h} - 1) \end{aligned}$$

NOTE – The fixpoint inverse multi component transformation is an implementation of the irreversible inverse multi component transformation when the input data is represented in fixpoint with four fractional bits and the transformation coefficients are represented with 13 fractional bits. The number of fractional bits has been selected such that 32 bit integer arithmetic is sufficient to implement the transformation.

## 5.21 Irreversible inverse high-dynamic range multiplicative residual multi component transformation (HCT) (normative)

This transformation shall be used for the reconstruction of the high-dynamic range chroma components in the multiplicative merging process. It is almost identical to the irreversible multi component transformation, though requires only two chroma components of the inversely tonemapped chroma signals, here denoted by  $C_1$  and  $C_2$ . The output signals  $M_0$ ,  $M_1$  and  $M_2$  are added to the gamma-corrected low-dynamic range image and scaled by the inversely tone-mapped luma signal. For details of the multiplicative process, consult Annex XXX.

$$\begin{aligned} M_0 &= 1.3984 * C_2 \\ M_1 &= -0.3054 * C_1 - 0.7038 * C_2 \\ M_2 &= 1.7969 C_1 \end{aligned}$$

## 5.22 Reversible inverse multi component transformation (RCT) (normative)

This transformation shall be applied to the additive residual data if the number of components,  $N_f$ , is equal to three. The output of this transformation is then added to the output of the inverse tone mapping process to reconstruct the high-dynamic range data. This transformation is applied regardless of whether the component decorrelation control marker is present, or the cc parameter within this marker. If the number of components is equal to one, this transformation shall be replaced by the identity.

In the following, let  $I_0, I_1$  and  $I_2$  be the reconstructed residual data generated by the residual scan as defined in Annex XXX, and  $R_0, R_1$  and  $R_2$  the reconstructed residual color data to be added to the inversely tonemapped low dynamic range data. Then set:

$$\begin{aligned} R_1 &= I_0 - \lfloor (I_1 + I_2) / 4 \rfloor \\ R_0 &= R_1 + I_2 \\ R_2 &= R_1 + I_1 \end{aligned}$$

NOTE – This transformation is a rough integer approximation of the irreversible inverse multi component transformation that has an exact inverse. It is identical to the RCT defined in ISO/IEC 15444-1.

## 5.23 Irreversible forward multi component transformation (ICT) (informative)

This transformation is used to decorrelate the low-dynamic range versions of the image before encoding the components by a legacy 10918-1 conforming encoder. This transformation is only applied when the number of

Kommentar [tr9]: This is what is found in the Dolby code currently. I have no idea why that should not be identical to the ICT, but it isn't. Can't we fix this?

Isn't there an offset shift missing? Is this part of the tone mapping?

Dolby, please add and clarify.

components,  $N_f$  indicated in the frame header is 3, and either the component decorrelation marker is absent or its  $cc$  parameter is set to one. Otherwise, the irreversible forward multi component transformation will be replaced by the identity transformation.

The input signals  $L_0$ ,  $L_1$  and  $L_2$  are already expected to be tone-mapped and hence in the low-dynamic range regime, all three in the range  $[0, 2^{8+h}-1]$  where  $h$  is the number of hidden DCT bits recorded in the JPEG Extensions Specification marker ('SPEC'). The output components  $I_0$ ,  $I_1$  and  $I_2$  are the input of the legacy 10918-1 process.

$$\begin{aligned} I_0 &= 2.9900 * L_0 & + & 0.58700 * L_1 & + & 0.11400 * L_2 \\ I_1 &= -0.1687358916 * L_0 & - & 0.3312641084 * L_1 & + & 0.5 * L_2 + 2^{7+h} \\ I_2 &= 0.5 * L_2 & - & 0.4186875892 * L_1 & - & 0.08131241085 * L_2 + 2^{7+h} \end{aligned}$$

NOTE – This transformation is, up to rounding errors, identical to the inverse of the irreversible inverse multi component transformation and thus identical to the transformation specified in 10918-6.

## 5.24 Fixpoint forward multi component transformation (FCT) (informative)

The transformation given in this subclause provides an implementation choice for the ICT, and by that – within the precision of a fixpoint representation – the forward transformation of the FCT defined in subclause C-XX. This transformation will only be used if the number of Components in the image, as indicated by the  $N_f$  parameter of the frame header, is three, and if either the component decorrelation marker is absent or its  $cc$  parameter is set to one. If  $N_f$  is either one or the component decorrelation marker is present and has its  $cc$  parameter set to zero, the fixpoint forward multi component transformation will be replaced by a simple conversion of the input data to the intermediate fixpoint format required for the fixpoint DCT transformation specified in Annex XXX. This process will be used in the lossless coding process, and may also be used for lossy compression.

Denote by  $L_0$ ,  $L_1$  and  $L_2$  the low-dynamic range sample values of the tone mapped high-dynamic range source image, represented by integer values in the range  $[0, 2^{8+h}-1]$  where  $h$  is the number of hidden DCT bits recorded in the JPEG Extensions Specification marker ('SPEC'). Denote by  $I_0$ ,  $I_1$  and  $I_2$  the output components, represented in a fixpoint format that is preshifted by four bits, i.e. four fractional bits are present. This format is required as input for the fixpoint DCT transformation that is used for the lossless encoding process. Then:

$$\begin{aligned} I_0 &= \lfloor (2449 * L_0 + 4809 * L_1 + 934 * L_2 + 2^8) / 2^9 \rfloor \\ I_1 &= \lfloor (-1382 * L_0 - 2714 * L_1 + 4096 * L_2 + 2^{13+7+h} + 2^8) / 2^9 \rfloor \\ I_2 &= \lfloor (4096 * L_0 - 3430 * L_1 - 666 * L_2 + 2^{13+7+h} + 2^8) / 2^9 \rfloor \end{aligned}$$

NOTE – This transformation is not exactly the inverse of the fixpoint inverse multi component transformation, though the additive residual coding process not only includes the residuals due to the finite implementation precision of the DCT, but also due to the precision loss in the FCT.

## 5.25 Irreversible forward high-dynamic range multiplicative residual multi component transformation (HCT) (normative)

TO BE SUPPLIED BY DOLBY.

Kommentar [tr10]: Needs to be clarified by Dolby.

## 5.26 Reversible forward multi component transformation (RCT) (informative)

The transformation supplied by this subclause decorrelates the additive error residuals of the lossless and lossy additive coding process and is the exact inverse of the RCT defined in subclause XXX. This transformation is applied whenever  $N_f$ , the number of components indicated in the frame header, equals 3, regardless of whether the component decorrelation marker is present, or what the value of the  $cc$  field of this marker is. If  $N_f$  is equal to one, the error residuals are coded directly and no further transformation is performed before the residual entropy coding process is applied to them.

Otherwise, let  $I_0$ ,  $I_1$  and  $I_2$  be the (integer) error residuals of the additive coding process, and  $R_0$ ,  $R_1$  to  $R_2$  the output of the transformation which is then quantized and coded by the residual coding process of Annex XXX. Then set:

$$\begin{aligned} R_0 &= \lfloor (I_0 + 2 * I_1 + I_2) / 4 \rfloor \\ R_1 &= I_2 - I_1 \\ R_2 &= I_0 - I_1 \end{aligned}$$

NOTE – This transformation is the inverse of the RCT defined in subclause XXX and, unlike the FCT, exactly invertible. It is a rough approximation of the RCT and identical to the forward RCT defined in ISO/IEC 15444-1.

## Annex D

### Tone Mapping Operations

This Annex forms a normative and integral part of this Recommendation | International Standard.)

Kommentar [tr11]: All this has to go into a separate part once we have the request for subdivision.

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

#### 5.27 Inverse Tone mapping (ITMO) (normative)

This subclause defines the inverse tone mapping of the additive residual coding process. Inverse tonemapping is applied to the output of the inverse multi component decorrelation transformation. Inverse tonemapping reconstructs from the low dynamic range image encoded by the codestream defined in ITU.T Rec.T 81| ISO/IEC 10918-1 an image in the high-dynamic range regime. In the additive coding process, residual data decoded from data in the JPEG Extensions Residual Data marker is added to the output of the inverse tone mapping process defining the reconstructed sample value.

The following steps shall be performed for inverse tone mapping of component i:

- Test the Tone Mapping Enable Bit  $T_e$  in the JPEG Extensions Specification marker whether inverse tone mapping shall be performed.
- If this bit is zero, multiply the sample value  $L_i$  given as the output of the inverse irreversible color transformation by  $2^{k_{b-h}}$  to reconstruct the high-dynamic range intermediate output  $H_i$ .  $R_b$  is the number of additional high dynamic range bits recorded in the JPEG Extensions Specification marker and  $h$  is the number of hidden DCT bits also specified in the same marker.  $R_b$  shall be larger or equal than  $h$ .
- If this bit is non-zero, find the tone mapping table index  $m_i$  in the JPEG Extensions Specification marker for component i,
  - Then locate the JPEG Extensions Inverse Tone Mapping Marker whose identifier  $M$  is equal to  $m_i$ . If no such marker exists, the codestream is nonconforming.
  - The intermediate high dynamic range output  $H_i$  is then defined as the value of the  $L_i$ -th table entry of the  $D_k$  table included in the JPEG Extensions Inverse Tone Mapping Marker found in the step above, i.e.

$$H_i = D_k[L_i],$$

where the brackets  $[x]$  evaluate the table  $D_k$  at entry  $x$ .

#### 5.28 Forward Tone mapping (FTMO) (informative)

This Recommendation | International Standard does not define a way how to find a forward tone mapping curve, nor how to generate its inverse required for the process described in subclause XXX. In general, the forward tone mapping process generates for each high-dynamic range sample value  $X_i$  a low dynamic sample value  $L_i$ , which is then encoded by a coding process defined in ITU.T Rec. T.80 | ISO/IEC 10918-1. It is, however, recommended for optimum coding efficiency that the algorithm for determining the sample value  $L_i$  and the inverse tone mapping curve  $D_k$  specified in the JPEG Extensions Inverse Tone Mapping Marker are selected such that the error residual  $R_i = X_i - D_k[L_i]$  is as small as possible.

A suitable algorithm would, for example, first construct a forward tone mapping curve  $F_k$  depending on the rendering intent of the image to be encoded. In a second step, an approximate inverse of  $F_k$  would be encoded in the JPEG Extensions Inverse Tone Mapping Marker, i.e.  $D_k$  would be selected such that  $|D_k[F_k[x]] - x|$  is as small as possible. The table  $D_k$  would then be encoded in the codestream, and the sample values  $L_i$  would be computed by  $L_i = F_k[X_i]$ .

TODO: Tone Mapping in the Multiplicative (Dolby) process.

Kommentar [tr12]: To be provided by Dolby.

## Annex E

### Splitting and Merging Low Dynamic Range and Residual Image

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

## 5.29 Merging LDR and residual image to a HDR image (normative)

This subclause specifies the reconstruction process of a high-dynamic range image from a LDR image and a residual image decoders shall follow. This process consists of the following steps, see also Figure XXX:

- Reconstruction of an image from the ITU.T Rec.T.81|ISO/IEC 10918-1 legacy codestream and the refinement codestream if the number of hidden DCT bits  $h$  is larger than zero. The  $h$  parameter is found in the JPEG Extensions Specification marker, the process how to reconstruct a common LDR image from the legacy codestream and the refinement extensions is specified in Annex XXX. If lossless decoding is desired, the fixpoint DCT specified in Annex YYY shall replace the DCT process specified in ITU.T Rec.T.81|ISO/IEC 10918-1. If no refinement data is present, i.e. if  $h=0$ , then the decoding specifications in ITU.T Rec.T.81|ISO/IEC 10918-1 shall be followed, where the fixpoint DCT in Annex YYY shall replace the DCT process in the legacy standard. The output of this step consists of one or three components that are neither upsampled nor inversely multi component decorrelated. The numerical representation of the sample values is implementation defined, unless lossless decoding is desired. In the latter case, the decoding process will, by means of the FDCT of Annex YYY generate sample values in a fixpoint representation with four fractional bits, i.e. output sample values are multiplied by a factor of 16.
- The upsampling process specified in Annex YYY shall be followed to generate samples for all positions on the sample grid. Upsampling is only valid for lossy coding. For lossless coding, no upsampling shall take place and the subsampling factors of all components shall be (1,1).
- The inverse multi-component decorrelation transformation shall be applied if the image consists of more than one component and the component decorrelation control marker is not present or its  $cc$  parameter is zero. This decorrelation transformation is specified in Annex XXX, and the output of this process are three sample values per sample grid location in the range  $[0, 2^{8+h}-1]$ . In case lossless decoding is required, the FCT specified in Annex XXX shall be applied, otherwise either the ICT or the FCT are suitable. If the codestream contains only one component, or the inverse multi component decorrelation transformation is disabled by a component decorrelation control marker whose  $cc$  value is one, then the upsampled sample values shall be converted to integers and clamped to the range  $[0, 2^{8+h}-1]$ . If lossless decoding is desired, this conversion process from fixpoint to integer shall be implemented as

$$I = \text{clamp}(\lfloor (L + 2^3) / 2^4 \rfloor, 0, 2^{8+h}-1)$$

where  $L$  is the output of the FCT and  $I$  is the reconstructed sample value.

- Each of the output components shall be inversely tonemapped according to the  $T_i$  parameter of the JPEG Extensions specification marker, where  $i$  is the component index, counting up from zero, in the order they appear in the frame header. If  $T_i$  is zero, the reconstructed sample value  $I$  shall be multiplied by  $2^{Rb-8-h}$  to reconstruct a high-dynamic range approximation value  $A$ :

$$H_i = 2^{Rb-8-h} \times I_i \quad (\text{for } T_i = 0)$$

If  $T_i$  is one, the table driven inverse tone mapping process specified in Annex XXX shall be used to find  $H_i$  from  $I_i$ . For that, the decoder shall locate the JPEG Extensions tone mapping marker whose  $M$  value is identical to  $T_i$ .

If  $T_i$  is two, a gamma-transformation shall be used to compute  $H_i$  from  $I_i$ . For this, the decoder shall compute the value  $\gamma_i$  from the  $Tg_i$  value of the JPEG Extensions specification marker as follows:

$$\gamma_i = Tg_i / 64$$

The output of the inverse gamma transformation is then found by

$$H_i = \begin{cases} I_i / 12.92 & \text{if } I_i / 2^{8+h} \leq 0.03928 \\ 2^{8+h} \times ((I_i / 2^{8+h} + 0.055) / 1.055)^{\gamma_i} & \text{otherwise} \end{cases}$$

- The decoder shall compute for each pixel the luminance prediction value  $Lu$  from the output sample values  $H_i$ . If the  $Tu$  parameter of the JPEG Extensions marker is zero, the  $Lu$  value shall be zero regardless of the input sample values. If  $Tu$  is one and the number of components in the image is one, the value of  $Lu$  shall be

Kommentar [tr13]: Not yet what the code does. It uses another DCT with zero fractional bits. Fix!

Kommentar [tr14]: This is not yet there.

identical to  $H_0$ , the only sample value of the image. If  $T_u$  is one and the image contains three components, the luminance prediction value  $L_u$  shall be computed by applying the forwards ICT specified in Annex XXX to the triple  $(H_0, H_1, H_2)$  and setting  $L_u$  to the first output:  $L_u = I_0$ . The outputs  $I_1$  and  $I_2$  of the ICT shall be discarded.

- The noise floor shall be added to the luminance prediction value  $L_u$ . The noise floor is derived from the  $T_{nl}$  parameter of the JPEG Extensions Specification marker:

$$L_p = L_u + T_{nl} / 128$$

- The residual image shall be reconstructed. If JPEG Extensions Residual Markers are present in the codestream, the process defined in Annex XXX shall be used. The output of this process is inversely quantized and an adaptive low-pass filter is applied, following the specifications of Annex XXX. The output of this process are one or three integer sample values  $R_i$  per sample grid point. Subsampling factors shall be all one if the JPEG Extensions Residual Marker is used to encode entropy coded data.
- If JPEG Extensions Quotient Markers are present, the decoding process defined in Annex YYY shall be used. The output of this data shall be upsampled following the specifications of Annex XXX, the output of this process are one or three integer sample values  $R_i$  per sample grid position.
- An inverse tonemapping process shall be applied to the first sample  $R_0$  giving the scale factor  $\mu$ .
  - If the  $T_m$  parameter of the JPEG Extensions marker is 1, and inverse tone mapping is applied to find  $\mu$  from  $R_0$ . For this, the decoder shall locate the JPEG Extensions tone mapping marker whose index is given by the  $T_{mt}$  parameter of the JPEG Extensions table. The value  $R_0$  shall then be used as an index into the table  $D_k$  of the tone mapping table to retrieve a single precision floating point number  $\mu$ .
  - If the  $T_m$  parameter of the JPEG Extensions marker is 3, the scale factor  $\mu$  shall be derived from  $R_0$  and the JPEG Extensions parameters  $T_{lymax}$  and  $T_{lymin}$  as follows:

$$\mu = \exp(R_0 \times (T_{lymax} - T_{lymin}) / 255 + T_{lymin})$$

- Inverse tone mapping shall be applied to the residual components  $R_0$ . The tone mapping to be performed on residual component  $R_i$  depends on the  $T_{ri}$  parameter of the JPEG Extensions Specification marker.
  - If  $T_{ri}$  is zero, the output of the inverse tone mapping shall be identical to the input

$$O_i = R_i$$

- If  $T_{ri}$  is one, the decoder shall find the JPEG Extensions tone mapping curve whose  $M$  parameter is identical to the  $T_{rti}$  parameter of the JPEG Extensions table parameter and apply an inverse tone mapping following the specifications of Annex XXX where the input is clamped to the range  $[0, 255]$ .

$$O_i = D_k[\text{clamp}(R_i, 0, 255)]$$

- If  $T_{ri}$  is four, the output  $O_i$  shall be determined by a linear function from the input where the scaling parameters are given by  $T_{rmin_i}$  and  $T_{rmax_i}$  of the JPEG Extensions Specification marker:

$$O_i = R_0 \times (T_{rmax_i} - T_{rmin_i}) / 255 + T_{rmin_i}$$

- The output of the inverse residual tone mapping shall be multiplied with the luminance prediction value  $L_p$  computed from the low dynamic range and refinement data giving the scaled residual data  $P_i$ :

$$P_i = L_p \times O_i$$

- Inverse color decorrelation shall be applied to the  $P_i$  data. This shall be the identity transformation if the  $T_{rcc}$  parameter of the JPEG Extensions Specification marker is zero or if only one component is present. If  $T_{rcc}$  is one, the inverse ICT or inverse FCT specified in Annex XXX shall be used. If the FCT is used as a substitute to the ICT, samples must be converted to the fixpoint format and back, as required by the FCT. If  $T_{rcc}$  is four, the inverse RCT specified in Annex XXX shall be used. The output of this process are the inversely decorrelated prediction errors  $Q_i$ .
- The high dynamic range output  $F_i$ , i.e. the final output of the decoding process, is reconstructed from  $H_i$ , the predicted high dynamic range signal, the multiplier  $\mu$  and the inversely decorrelation prediction errors  $Q_i$  as follows:

$$F_i = \mu \times H_i + Q_i$$

Kommentar [tr15]: This is currently a limitation of the implementation. Does this make sense in general? Lossless decoding of subsampled data?

Kommentar [tr16]: Is this a good name? This is where the Dolby codestream goes encoding the quotient data.

Kommentar [tr17]: This need to include a clamping process to ensure table lookups are correct.

Kommentar [tr18]: Need to extend JPEG TONE to include IEEE float values.

Kommentar [tr19]: Note that  $\mu$  can be set to all-one (as required by Stuttgart for lossless coding) by setting  $T_{lymax}$  and  $T_{lymin}$  to both zero.

Kommentar [tr20]: This is what I find in the Dolby code. Please check.

Kommentar [tr21]: Do we actually need this? Dolby doesn't, I don't. Note that the residual data is signed, thus this only applies to the Dolby process naturally.

Kommentar [tr22]: This also includes the identity case required by Stuttgart and the zero case required by Dolby.

Kommentar [tr23]: Dolby uses here actually a slightly different transformation. Is this an error or intentional? Please check!

- Lossless coding restricts the choice of some of the parameters. The L transformation inversely decorrelating the low-dynamic range and refinement coding samples shall be the FCT. The Tu parameter shall be zero, disabling the generation of the luminance prediction value, and the noise floor value Tnl shall be 128, yielding a constant luminance prediction value Lp of one. The Tm parameter defining the generation of the multiplier  $\mu$  shall be three and the Tlymax and Tlymin parameters shall be zero, thus generating a constant multiplier of one. The residual tone mapping parameters Tr<sub>i</sub> shall be all four, and Trmax<sub>i</sub> = 255 and Trmin<sub>i</sub> = 0 shall hold for all i, thus defining an identity mapping between O<sub>i</sub> and R<sub>i</sub>. Finally, the H transformation inversely decorrelating the residual data shall be the RCT by setting the Trcc parameter to four.

Kommentar [tr24]: Should probably go into a separate section defining the profiles. There should also be a Dolby profile.

### 5.30 Splitting HDR data into LDR and residual data (informative)

Kommentar [tr25]: This should probably go into a separate "Recommendations and Guidelines" section.

This Recommendation/International Standard does not define a normative algorithm for computing a low dynamic range and a residual image from a given high dynamic range image. Nevertheless, some guidelines will be given that demonstrate how the coding mechanisms defined by this Recommendation/International Standard can be deployed to implement image coding in the high dynamic range regime.

An image coding process in which a low dynamic range image is extended to high dynamic range by means of an additive residual error can be realized by disabling the S and R tone mapping processes, setting the S-transformation to zero and selecting a noise floor of 1.0. This process has the advantage that the absolute error in the high-dynamic range process is easily controllable and may even be zero, i.e. lossless coding is possible. The following steps should be taken to implement this coding process:

- The Tu parameter of the JPEG Extensions marker is set to zero, indicating that the S transformation will be disabled.
- The Tnl parameter is set to 128 indicating a noise floor of 1.0.
- The Tm parameter is set to three, indicating an exponential curve for the S-curve. Additionally, Tlymax and Tlymin are set to zero such that the exponential curve is identical to one. Thus, the  $\mu$  factor is always one and the output scaling of the low-dynamic range image is disabled.
- All Tr<sub>i</sub> parameters are set to four, indicating a linear relation between the residuals R<sub>i</sub> and O<sub>i</sub>. Additionally, Trmax<sub>i</sub> is set to 255 and Trmin<sub>i</sub> to zero, such that the R-lookup tables implement the identity transformation.
- The L-lookup tables are chosen as the inverses of a suitable tone mapping process. For example, the encoder might apply the Reinhard tone mapping operator [xxx] to generate a low-dynamic range version of the input HDR image, and compute a numerical inverse of the tone mapping operation which is then stored as lookup table in the JPEG Extensions Inverse Tone Mapping marker segment. For this, the Tt<sub>i</sub> parameters are zero.
- The low dynamic range image generated by forwards tone mapper is directly compressed by a legacy ITU.T Rec.T.81 | ISO/IEC 10918-1 encoder.
- The residual image is generated by inversely dequantizing the low-dynamic range DCT coefficients C<sub>ij</sub> and applying an inverse DCT transformation to them. The resulting components are then upsampled and inversely decorrelated. The difference between these samples and the original image define then the residual signal Q<sub>i</sub>. This essentially requires reproducing parts of the decoding process at the encoder.
- The difference signal Q<sub>i</sub> is decorrelated by the ICT or RCT implementing the (approximate) inverse of the H transformation, generating the R<sub>i</sub> signal.
- The R<sub>i</sub> signal is quantized and encoded by the residual coding process defined in Annex XXX. While coding this signal by the legacy ITU.T Rec.T.81|ISO/IEC 10918-1 specifications is possible, it is not recommended since the signal consists of already decorrelated noise which does not compress well by the legacy coding mechanism.
- The output rate can be controlled by adjusting the quantization matrix of the LDR coding process and the quantization parameter of the residual coding process. For low qualities, it is typically most efficient to disable the residual signal completely and defining the quality of the HDR image entirely by the legacy quantization matrix. For higher qualities, the legacy quantization matrix should be left constant and should define relatively fine quantization, i.e. the entries in the quantization matrices should be small. Quality is then controlled by the high-dynamic range quantization parameter.

Lossless coding can be implemented by a refinement of the above coding process, namely

- The H transformation must be the RCT, i.e. the Trcc parameter is set to four.
- The DCT to compute the decoder generated LDR data at encoder side is the fixpoint DCT



- The inverse decorrelation transformation is the FCT.
- The quantization value for the HDR data is set to zero.
- The residual coding process defined in Annex XXX is selected to encode the error residuals.

NOTE – the choice of the forwards DCT, the L-curves or the quantization matrix of the legacy codestream is irrelevant as far as lossless decoding is concerned. It influences only the compression performance of the encoder, though not the ability to reconstruct the original image samples without loss.

By disabling the residual coding completely and additionally setting the L-lookup tables to identity, an extension of the legacy DCT process for bit depths up to 12 bits per sample is possible. This allows a lightweight, minimum effort coder for images of up to 12 bits per sample. Coding performance of this extension is not quite as good as the performance of the residual process defined above. For this, only refinement coding is used.

EDITOR'S NOTE: Description of Dolby's encoding process should go here, additionally with its advantages and when to pick it.

Kommentar [tr26]: TBD by Dolby.  
Need to understand how coding works here.

## Annex F

### Lossy Coding of Residual Data

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

TO BE SUPPLIED BY DOLBY. This is likely only a 10918-1 codestream encapsulated in an APP11-marker. Can we change the marker?

## Annex G

### Lossless and Near Lossless Coding of Residual Data

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

This Annex describes one method for coding and decoding residual data that works most effective on additive residual data, the other one is the DCT method specified in Annex YYY. Since the residual coding process specified here can be implemented without loss, it shall be used for lossless coding of still images. It is also quite effective for lossy coding of high-dynamic range images and is then combined with the quantization of the residual data specified in Annex XXX and, optionally, noise shaping, specified in the same Annex. The lossy residual coding method specified in Annex YYY is then another coding option.

Entropy coded residual data is encapsulated in JPEG Extensions Residual Data markers. Whenever such markers are present, additional residual scans over the image provide a correction signal that enhances the high-dynamic range image or enables lossless coding by encoding the residual error caused by the non-ideally invertible FDCT and FCT decorrelation transformation. Unlike the legacy coding process or the refinement coding process, residual coding does not transform its coefficients by a DCT. First of all, the DCT is not exactly invertible and is thus unsuitable for lossless coding, and second because the residual signal is mostly noise and additional transformations do not provide compression gains. Instead, entropy coding and decoding operates directly in the spatial domain on the residual signal, and so does the quantization and noise shaping specified in Annex XXX. Nevertheless, the residual signal is also separated into 8×8 blocks that are aligned with the blocks the DCT transformation is run on for creating the legacy and refinement data, and these blocks are scanned in the order defined in ITU.T Rec.T.81|ISO/IEC 10918-1.

#### 5.31 Decoding of Additive Residual Error Signals (normative)

Decoding the additive residual error signal consists of the following steps: First, creation of a continuous codestream from the data encapsulated in the JPEG Extensions Residual Data markers, and second entropy decoding of this bitstream. The output of this process is a residual signal  $R_{ij}$  organized in 8×8 blocks that is then inversely quantized, optionally noise-shaped by an adaptive low-pass filter before entering the upsampling, inverse tone mapping and inverse multi component decorrelation transformation, see Figure XXX in clause YYY. This annex only describes the first two steps, creation of a continuous codestream and its decoding.



In the first step, the payload data of all JPEG Extensions Residual data markers, i.e. the marker contents without the marker size, the Common Identifier and the JPEG Extensions Type Identifier shall be concatenated in the order they appear in the codestream. The output of this step is a continuous codestream defining one or several scans over the residual data, consisting optionally of table definitions defining the Huffman table for the entropy decoder, followed by scan headers in the form of SOS marker segments defining the components participating in the scan, followed by the entropy coded data itself. A frame header shall not be present in this codestream. Figure H-1 displays the organization of the residual codestream resulting from concatenating the payload data of the JPEG Extensions Residual data markers.

[Tables/ misc]	Scan header	ECS <sub>0</sub>	[RST <sub>0</sub> ]	...	ECS <sub>last</sub>	[RST <sub>last</sub> ]	[Tables /misc]	Scan header	ECS <sub>0</sub>	[RST <sub>0</sub> ]	...	ECS <sub>last</sub>	[RST <sub>last</sub> ]
-------------------	----------------	------------------	---------------------	-----	---------------------	------------------------	-------------------	----------------	------------------	---------------------	-----	---------------------	------------------------

**Figure G-1: Syntax of the Contents of the Residual Data Continuous Codestream**

Data in the entropy coded segments ECS<sub>0</sub> to ECS<sub>last</sub> shall be decoded by a modified progressive or successive approximation scan algorithm as defined in Annex G of ITU.T Rec.T.81|ISO/IEC 10918-1:

- The DC decoding step shall be skipped, i.e. subclause G.1.2.1 of ITU.T Rec.T.81|ISO/IEC 10918-1 is not effective and shall be skipped.
- The AC decoding algorithms defined in subclauses G.1.2.2 and G.1.2.3, specifically figures G.3 and G.7 also apply if Ss = 0, i.e. if the start of the spectral selection is zero. For example, if Ss = 0, the flow charts G.3 and G.7 start with K=-1 and the first coefficient decoded is ZZ(0), the top left block edge.
- Consequently, only the AC Huffman tables shall be used for decoding the entropy coded data segment. The DC Huffman table indices in the scan header shall be ignored. If no Huffman table definition is found in the residual codestream, the Huffman tables defined in the frame or scan header of the legacy codestream apply.
- The dimensions of the frame and the subsampling factors are defined in the frame header of the legacy codestream.

NOTE – Since the scan type defined in this Annex decodes residual data that has not been transformed by a DCT, no specific DC coefficient that would require special treatment is present; instead, all residual error coefficients follow an AC statistics that allows a consistent coding mechanism. The progressive encoding mechanism of AC coefficients defined in ITU.T Rec.T.81 | ISO/IEC 10918-1, especially the block-skip mode, allows efficient encoding of data that contains many zeros. A similar block skip mode is not available for DC coefficients which was removed for this reason. Note that the modified progressive coding mode is used for residual data regardless of the frame type indicated by the legacy codestream.

### 5.32 Coding of Additive Residual Error Signals (informative)

Encoding of non-DCT transformed additive residual data consists of two steps: First, the data is encoded by one or several modified progressive or subsequent approximation scans. The resulting codestream, complete with Huffman table definitions and scan headers, but without the frame header, is split into one or several data chunks encapsulated into JPEG Extensions Residual Data markers which are then embedded into the legacy codestream.

In the first step, the residual data is first separated into 8×8 blocks that are aligned to the DCT blocks of the LDR image. Note that the residual data is untransformed, but potentially quantized and noise-shaped. This data is then encoded by one or many modified progressive scans, and optionally additional subsequent approximation scans, forming one or many entropy coded segments. The modifications made to the entropy coding algorithm are identical to those discussed in subclause XXX, namely DC coding is skipped and AC coding also applies to the residual coefficient at block position 0,0. Optionally, the modified progressive and subsequent approximation scans may redefine Huffman tables to improve coding efficiency of the residual data.

The output of the first step is a continuous codestream that optionally includes one or several DHT marker segments, and one or many scans each consisting of a SOS marker followed by an entropy coded segment. The layout of this codestream follows then Figure H-1 in the previous subclause.

In the second step, this continuous codestream is separated into chunks, each of which no longer than 2<sup>16</sup>-9 bytes. These chunks are inserted as payload data into JPEG Extensions Residual data marker segments which are injected into the legacy codestream.

NOTE – It is recommended, though not required, to define Huffman tables that are specific to the residual coding scans as this helps to improve the coding efficiency. Since spectral selection does not select any part of a spectrum, but rather selects the spatial position of the residual data within an 8×8 block in the image, using this feature of progressive or subsequent approximation scans is not particularly useful, though valid. Successive approximation scans may provide a subsequent refinement of the residual data, though.

## Annex H

### Entropy Coding of Refinement Data

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

This Annex defines the procedures for decoding and encoding of images using an extended bitdepths in the DCT domain. If this coding method is used, the entropy coded data visible to legacy decoders describes an image with a precision of eight bits per sample, which is extended by a refinement coding procedure described in this Annex to up to 12 bits. Refinement coding can be combined with other coding methods, for example the residual coding method defined in Annex G, then allowing even higher bit depths of up to 16 bits per sample.

Refinement coding is signaled by a non-zero value of the  $h$  parameter in the JPEG Extensions Specification marker; the value of  $h$  describes the additional bit precision becoming available due to the refinement scan described in this Annex. It may take values from zero – refinement scan disabled – to four, resulting in 12 bits precision for the output of the DCT process.

If the refinement scan is enabled by setting  $h \neq 0$ , the quantization and inverse quantization process defined in ITU.T Rec. T.81 has to be modified to include the additional bits, and the entropy coding and decoding procedures includes additional scans over the image to represent such bits. The DCT transformation and multi component decorrelation transformation remain unchanged except that the DC level shift has to be modified to include the additional bits, and care must be taken to avoid overflow conditions as the range of the input and output data is extended. These considerations were included in the fixpoint DCT transformation defined in Annex XXX which shall be used for lossless reconstruction; it may also be used for lossy reconstruction of images and then provides a good compromise between speed and precision.

#### 5.33 Modifications of the inverse quantization process for residual decoding (normative)

The regular dequantization process in the absence of refinement scans multiplies the entropy decoded DCT coefficient value  $S_{i,j}$  at block position  $i,j$  with the quantization matrix entry  $q_{i,j}$  to get the DCT coefficient  $C_{i,j}$ . In the presence of a residual scan,  $h$  additional bits  $T_{i,j}$  are decoded by the refinement scan that shall be included as  $h$  least significant bits in the reconstruction of the DCT coefficient  $C_{i,j}$  as follows:

$$C_{i,j} = (2^h \times S_{i,j} + T_{i,j}) \times q_{i,j}$$

NOTE – An equivalent implementation strategy for the above inverse quantization procedure is to decode the data from the legacy scans represented by the entropy coded data following the SOS markers as if the point transformation parameter would be set to  $AI+h$  (or  $0+h$ ) instead of  $AI$  (or zero).  $AI$  is the point transformation parameter of the scan header, see subclause A.4 and B.2.3 of ITU.T Rec.T.81|ISO/IEC 10918-1. Data decoded by the refinement scan uses an unmodified decoding procedure and is then automatically placed in the least significant bits of  $S_{i,j}$ .

#### 5.34 Modifications of the quantization process for residual coding (informative)

The quantization procedure generates in the presence of refinement scans from the DCT coefficients  $C_{i,j}$  two data sets, the legacy quantized data  $S_{i,j}$  which will be encoded by the entropy coding methods provided by ITU.T Rec.T.81|ISO/IEC 10918-1, and the refinement data  $R_{i,j}$  which undergoes refinement coding.  $S_{i,j}$  and  $T_{i,j}$  are computed from the DCT coefficients  $C_{i,j}$  and the quantization matrix  $q_{i,j}$  as follows:

$$\begin{aligned} X_{i,j} &= \lfloor C_{i,j}/q_{i,j} + 0.5 \rfloor \\ S_{0,0} &= \lfloor X_{0,0}/2^h \rfloor & T_{0,0} &= X_{0,0} - 2^h \times S_{0,0} \\ S_{i,j} &= \text{sign}(X_{i,j}) \times \lfloor |X_{i,j}|/2^h \rfloor & T_{i,j} &= X_{i,j} - 2^h \times S_{i,j} \quad \text{for } (i,j) \neq (0,0) \end{aligned}$$

NOTE – The division and rounding algorithm for computing  $S_{i,j}$  generates a "fat zero" for AC coefficients, i.e. a dead zone of twice the size of a regular quantization bucket. DC coding does not follow this convention to avoid drift errors when encoding. The rounding conventions picked here are intentionally identical to the rounding procedure used when encoding data with the successive approximation of the progressive scan defined in ITU.T Rec.T.81|ISO/IEC 10918-1, and this is, in fact, one possible implementation strategy: The entropy coded segment is created by a legacy encoder as if the successive approximation parameter would be  $AI+h$  (or  $h$ ) instead of  $AI$  (or zero), the remaining bits are then encoded by the refinement scan defined in this Annex.

### 5.35 Decoding process for entropy coded refinement data (normative)

This subclause defines the decoding procedure for the refinement data  $T_{ij}$  that extends the precision of the DCT coefficients by  $h$  bits. Entropy coded refinement data is present whenever the  $h$  parameter of the JPEG Extensions Specification marker is nonzero.

In a first step, the decoder shall concatenate the payload data of all JPEG Extensions Refinement ('FINE') markers in the order they appear in the bitstream to form a single consecutive bitstream. The payload data is defined by the marker body without the marker segment length, the JPEG Extensions Common Identifier and the JPEG Extensions Type Identifier.

The resulting stream consists of a series of scan headers and entropy coded data following the syntax of ITU.T Rec.T.81|ISO/IEC 10918-1, except that neither a frame header nor table definitions shall be present. Restart markers, if indicated in the legacy codestream, shall then also be present here. The scans included in this stream are decoded by a modification of the successive approximation decoding algorithm using Huffman decoding defined in Annex G of ITU.T Rec.T.81|ISO/IEC 10918-1, regardless of the frame type indicated in the legacy codestream (i.e. regardless of whether the frame type is baseline, extended sequential or progressive).

[Tables/ misc]	Scan header	ECS <sub>0</sub>	[RST <sub>0</sub> ]	...	ECS <sub>last</sub>	[RST <sub>last</sub> ]	[Tables /misc]	Scan header	ECS <sub>0</sub>	[RST <sub>0</sub> ]	...	ECS <sub>last</sub>	[RST <sub>last</sub> ]
-------------------	----------------	------------------	---------------------	-----	---------------------	------------------------	-------------------	----------------	------------------	---------------------	-----	---------------------	------------------------

**Figure G-1: Syntax of the Contents of the Refinement Data Continuous Codestream**

The entropy coding algorithm used for refinement decoding is identical to the successive approximation decoding defined in Annex G of ITU.T Rec.T.81|ISO/IEC 10918-1 except that the scan pattern  $ZZ(K)$  has to be initialized before starting the first refinement scan by the entropy decoded data of the legacy codestream shifted by  $2^h$  bits, i.e.

$$ZZ(K) = 2^h \times S_{(i,j)(k)}$$

where  $(i,j)(k)$  is the  $k$ -th position of the zig-zag scan defined in Figure 5 of ITU.T Rec.T.81|ISO/IEC 10918-1. The refinement data  $T_{ij}$  after decoding is then defined by the output  $ZZ(K)$  of all successive approximation scans included in the JPEG Extensions refinement markers:

$$T_{(i,j)(k)} = ZZ(K) - 2^h \times S_{(i,j)(k)}$$

NOTE – An alternative implementation strategy would first decode the legacy portion of the codestream using either the baseline, sequential or progressive decoding procedure, but with a point transformation/successive approximation parameter of  $AI+h$  (or  $h$ ) for each scan included in the legacy codestream. The refinement scan extends these scans on the same scan buffer by using a successive approximation scan on the same scan buffer while using the successive approximation parameter exactly as indicated in the scan headers included in the JPEG Extensions Refinement marker. Note further that the refinement scan always uses a successive approximation scan pattern, regardless of the frame type of the legacy codestream.

### 5.36 Encoding of refinement data (informative)

This Annex specifies an encoding process of refinement data that is compatible to the decoding process defined in the above subclause. In the first step, the scan pattern  $ZZ(K)$  is initialized with the quantized data encoded in the legacy codestream,  $S_{ij}$ , plus the refinement data  $T_{ij}$ :

$$ZZ(K) = 2^h \times S_{(i,j)(k)} + T_{(i,j)(k)}$$

where  $(i,j)(k)$  is the  $k$ -th position in the zig-zag scan. Then, one or several successive approximation scans following the specifications in Annex G of ITU.T Rec.T.81|ISO/IEC 10918-1 encode the  $h$  least significant bits of the data, i.e. the  $AI$  parameter of the scans generated will vary between  $h-1$  and zero.

The resulting entropy coded data, together with the scan headers defining the scan pattern, but without any frame header or table definitions, are split into chunks of at most  $2^{16}-9$  bytes each and inserted as payload data of the JPEG Extensions Refinement markers into the legacy bitstream.

NOTE – An alternative implementation would first use a legacy scan process with a point transformation of  $AI+h$  (or  $h$ ) bits to encode the LDR data, and would then run a successive approximation scan on the same data buffer with an unmodified parameter of  $AI$ . Note further that the refinement encoding procedure always uses the successive approximation scan of the progressive coding mode, regardless of the frame type of the legacy codestream.

## Annex I

### Discrete Cosine Transformation

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

This Annex defines a fixpoint version of the inverse DCT transformation implementations of the lossless profile shall follow. If lossless reconstruction is not required, alternative implementations of the DCT process are possible; nevertheless, the algorithm described here define a DCT implementation that is conforming to ITU.T Rec.T.81|ISO/IEC 10918-1 and that may be used even for the lossy decoding process.

### 5.37 Overview on the Inverse Fixpoint DCT Transformation (normative)

The Inverse Fixpoint DCT Transformation takes integer samples  $Q_{i,j}$  in the form of an  $8 \times 8$  block, and generates from that fixpoint samples  $Y_{i,j}$  also organized in an  $8 \times 8$  block. The input samples  $Q_{i,j}$  are integers generated by the dequantization procedure described in ITU.T Rec.T.81|ISO/IEC 10918-1, the output samples  $Y_{i,j}$  are fixpoint numbers with four fractional bits, i.e. integers premultiplied by 16, that are suitable for the Inverse Fixpoint Decorrelation Transformation (FCT) defined in Annex C. The steps of the algorithms are as follows:

- Level-shift the DC coefficient  $Q_{0,0}$  and scale all coefficients to the precision of the FCT:
 
$$X_{0,0} = 2^4 \times Q_{0,0} + 2^{7+h+3+4}$$
 where  $h$  is the number of hidden DCT bits specified by the  $h$  parameter of the JPEG Extensions marker.
 
$$X_{i,j} = 2^4 \times Q_{i,j} \text{ for } (i,j) \neq (0,0)$$
- Run a one dimensional inverse fixpoint DCT process over all columns of the  $8 \times 8$  block
  - by setting  $A_k = X_{k,l}$  where  $k$  is the row and  $l$  is the column index
  - following the procedure of the one-dimensional DCT process taking the vector  $A$  as input and generating the vector  $B$  as output
  - backwards scaling  $B$  and generating one column of the  $8 \times 8$  block  $Z$  from  $B$ 

$$Z_{k,l} = \lfloor (B_k + 2^8) / 2^9 \rfloor$$
  - repeating the above steps for all columns  $l \in [0,7]$
- Run a one dimensional inverse fixpoint DCT process over all rows of the  $8 \times 8$  block
  - By setting  $A_k = Z_{j,k}$  where  $j$  is the row and  $k$  is the column index
  - Following the procedure of the one-dimensional DCT process taking the vector  $A$  as input and generating the vector  $B$  as output
  - Backwards scaling the output  $B$  and inserting this as one row into the output matrix  $Y$ 

$$Y_{j,k} = \lfloor (B_k + 2^{8+3}) / 2^{9+3} \rfloor$$
  - Repeating the above steps for all rows  $j \in [0,7]$

Kommentar [tr27]: Do we actually need to include the DC shift here? Check with the old standard.

NOTE – the DCT process defined in this subclause is designed such that 32 bit wide variables are sufficient to implement all steps for  $h \leq 4$ , i.e. at most 12 bits sample precision without causing overflows.

### 5.38 One dimensional inverse fixpoint DCT transformation (normative)

This subclause defines a one dimensional inverse fixpoint process that is a building block of the two-dimensional fixpoint DCT transformation defined in the subclause above. It takes an eight-dimensional vector  $A_k$  of fixpoint numbers prescaled by 4 bits as input, and generates an eight-dimensional vector  $B_k$  of fixpoint numbers prescaled by  $9+4=13$  bits as output.

- Set  $Z_1 = (A_2 + A_6) \times 277$
- Set  $T_2 = Z_1 - A_6 \times 946$
- Set  $T_3 = Z_1 + A_2 \times 392$
- Set  $T_0 = (A_0 + A_4) \times 512$
- Set  $T_1 = (A_0 - A_4) \times 512$
- Set  $T_{10} = T_0 + T_3$
- Set  $T_{13} = T_0 - T_3$

- Set  $T_{11} = T_1 + T_2$
- Set  $T_{12} = T_1 - T_2$
- Set  $Z_4 = A_7 + A_3$
- Set  $Z_5 = A_5 + A_1$
- Set  $Z_6 = (Z_4 + Z_5) \times 602$
- Set  $Z_7 = (A_7 + A_1) \times -461$
- Set  $Z_8 = (A_5 + A_3) \times -1312$
- Set  $Z_9 = Z_4 \times -1004 + Z_6$
- Set  $Z_{10} = Z_5 \times 200 + Z_6$
- Set  $T_{30} = A_7 \times 153 + Z_7 + Z_9$
- Set  $T_{31} = A_5 \times 1051 + Z_8 + Z_{10}$
- Set  $T_{32} = A_3 \times 1573 + Z_8 + Z_9$
- Set  $T_{33} = A_1 \times 769 + Z_7 + Z_{10}$
- Set  $B_0 = T_{10} + T_{33}$
- Set  $B_7 = T_{10} - T_{33}$
- Set  $B_1 = T_{11} + T_{32}$
- Set  $B_6 = T_{11} - T_{32}$
- Set  $B_2 = T_{12} + T_{31}$
- Set  $B_5 = T_{12} - T_{31}$
- Set  $B_3 = T_{13} + T_{30}$
- Set  $B_4 = T_{13} - T_{30}$

NOTE – The above algorithm implements an (unscaled) one dimensional DCT transformation that is, by itself, not exactly invertible. The lossless process will correct errors from the forwards transformation by encoding residuals between this DCT and the ideal transformation in a residual scan.

### 5.39 Overview on the Fixpoint DCT Transformation (informative)

This subclause provides an overview and implementation guideline for a forwards reversible fixpoint DCT. Many other implementation choices are possible, even for the lossless profile, provided that the residual is computed according to inverse fixpoint DCT specified above. The forwards fixpoint DCT transform takes an  $8 \times 8$  block of sample values  $Y_{i,j}$  as input and generates from that a DCT transformed signal  $X_{i,j}$  which is also organized in an  $8 \times 8$  block. In the implementation presented here, the input samples are fixpoint numbers premultiplied by 16, i.e. have 4 fractional bits. The output coefficients  $Q_{i,j}$  are also represented in fixpoint and have  $9+4+3=16$  fractional bits and form the input for the quantization procedure defined in ITU-T Rec.T.81|ISO/IEC 10918-1. The number of fractional bits in both the input and output samples have been selected such that the DCT transformation can be implemented on 32 bit architectures without causing overflow conditions for input samples of up to 12 bits precision. This corresponds to the maximum allowable choice of  $h = 4$ .

The forwards fixpoint DCT consists of the following steps:

- Extract row  $k$  from the source signal  $Y_{k,l}$  and denote this row-vector by  $B_l$ .
- Run the one dimensional forwards fixpoint DCT algorithm on  $B_k$  described in the next subclause resulting in an output vector  $A_k$ .
- Scale the output vector  $A_k$  back into a fixpoint representation and insert it as  $k^{\text{th}}$  row into the intermediate block  $Z_{k,l}$ :

$$Z_{k,l} = \lfloor (A_k + 2^8) / 2^9 \rfloor$$

- Repeat the above steps for all rows  $l \in [0,7]$
- Extract column  $j$  from the intermediate block  $Z_{i,j}$  and use this as input vector  $B_i$  for the one dimensional fixpoint DCT

$$B_i = Z_{i,j}$$

- Run the one dimensional forwards fixpoint algorithm on  $B_i$  giving  $A_i$ .
- Insert the output vector as column  $j$  into the intermediate output DCT block  $X_{i,j}$

$$X_{i,j} = A_i$$

- Repeat the above steps for all columns  $i \in [0,7]$  to fill the entire block  $X_{i,j}$ .
- Generate the final output by level shifting the DC coefficient and copying all other coefficients to  $Q_{i,j}$

$Q_{0,0} = [X_{0,0} - 2^{7+h+9+4+3+3}]$  where  $h$  is the number of hidden DCT bits defined by the  $h$  parameter of the JPEG Extensions specification marker.

$$Q_{i,j} = X_{i,j} \quad \text{for } (i,j) \neq (0,0)$$

Kommentar [tr28]: Do we need to include the DC shift here? What about the "point transformation"?

The output block  $Q_{i,j}$  contains now DCT coefficients preshifted by  $9+4+3=16$  bits. Quantization can now be performed by dividing  $Q_{i,j}$  by  $2^{16}q_{i,j}$  and rounding to the nearest integer, where  $q_{i,j}$  is the quantization matrix for the component currently transformed.

#### 5.40 One dimensional fixpoint DCT transformation (informative)

This subclause provides an implementation example for an unscaled one dimensional fixpoint DCT that is used as a building block of the two-dimensional fixpoint DCT specified in the above subclause. Alternative implementations are possible, even for the lossless compression path. The input to this algorithm is an eight-dimensional vector  $B_i$  of fixpoint coefficients preshifted by four bits, the output a vector  $A_i$  of eight fixpoint numbers preshifted by  $4+9=13$  bits.

- Set  $T_0 = B_0 + B_7$
- Set  $T_1 = B_1 + B_6$
- Set  $T_2 = B_2 + B_5$
- Set  $T_3 = B_3 + B_4$
- Set  $T_{10} = T_0 + T_3$
- Set  $T_{12} = T_0 - T_3$
- Set  $T_{11} = T_1 + T_2$
- Set  $T_{13} = T_1 - T_2$
- Set  $A_0 = (T_{10} + T_{11}) \times 512$
- Set  $A_4 = (T_{10} - T_{11}) \times 512$
- Set  $Z_1 = (T_{12} + T_{13}) \times 277$
- Set  $A_2 = Z_1 + T_{12} \times 392$
- Set  $A_6 = Z_1 - T_{12} \times 946$
- Set  $T_{20} = B_0 - B_7$
- Set  $T_{21} = B_1 - B_6$
- Set  $T_{22} = B_2 - B_5$
- Set  $T_{23} = B_3 - B_4$
- Set  $T_{32} = T_{20} + T_{22}$
- Set  $T_{33} = T_{21} + T_{23}$
- Set  $Z_2 = (T_{32} + T_{33}) \times 602$
- Set  $Z_3 = (T_{20} + T_{23}) \times -461$

- Set  $Z_4 = (T_{21} + T_{22}) \times -1312$
- Set  $Z_5 = T_{32} \times -200 + Z_2$
- Set  $Z_6 = T_{33} \times 1004 + Z_2$
- Set  $A_1 = T_{20} \times 769 + Z_3 + Z_5$
- Set  $A_3 = T_{21} \times 1573 + Z_4 + Z_6$
- Set  $A_5 = T_{22} \times 1051 + Z_4 + Z_5$
- Set  $A_7 = T_{23} \times 153 + Z_3 + Z_6$

NOTE – The above transformation is, by itself, not exactly invertible, neither that nor the choice of the forwards DCT matters for lossless reconstruction as long as the inverse fixpoint DCT is implemented exactly as shown.

## Annex J

### Quantization and Noise Shaping of the Residual Image

(This Annex forms a normative and integral part of this Recommendation | International Standard.)

In this Annex and all of its subclauses, the flow charts and tables are normative only in the sense that they are defining an output that alternative implementations shall duplicate.

This Annex and its subclauses describe the procedure by which residual data coefficients decoded by the lossless process defined in Annex XXX are reconstructed before merging the reconstructed data with the LDR data following the procedures of Annex A. The reconstruction after entropy decoding consists of two steps: An optional adaptive low-pass filter removing the high-frequency noise generated by the optional noise-shaping filter at encoder side, and an inverse quantization step that reconstructs the residuals from the encoded coefficients. If noise shaping is turned off and the quantization bucket size is set to one, lossless reconstruction of the residuals is possible, enabling lossless compression in combination with the fixed point DCT and the FCT decorrelation transformation.

#### 5.41 Adaptive Low-Pass Filtering and Inverse Quantization (normative)

This subclause specifies the adaptive low-pass filtering and inverse quantization algorithm decoders shall follow if the S bit in the JPEG Extensions Specification Marker is set to one. If it is set to zero, only the inverse quantization procedure defined in the following subclause applies.

The input data to the adaptive low-pass filtering and inverse quantization is an  $8 \times 8$  block of reconstructed integer residual coefficients  $C_{n,m}$  coming from the lossless residual data entropy decoding process defined in Annex XXX. This data block is separated again into  $2 \times 2$  subblocks  $S_{ij}^k$  where  $k$  is the subblock index ( $k=0..15$ ) and  $i,j$  the position within the subblock ( $i,j=0..1$ ). One has:

$$S_{ij}^k = C_{(i+k \% 4, j+k / 4)}$$

The quantization step size  $q_R$  that applies to inverse quantization of the residual coefficients is taken to be the entry  $q_{63}$  of one of two possible quantization matrices whose indices are given in the JPEG Extensions Specification marker. For component 0, and for images consisting of a single component only, the quantization matrix given by the index  $Lt$  in the JPEG Extensions Specification marker shall be used. For components 1 and 2, the quantization matrix with index  $Ct$ , also found in the JPEG Extensions Specification marker, shall be used to define the parameter  $q_R$ .

In the first step, the decoder shall compute the 16 averages over the inversely quantized subblocks by:

$$Avg^k = \left\lfloor (\sum_{i,j=0..1} S_{ij}^k \times q_R + 2) / 4 \right\rfloor$$

In a second step, coefficients are inversely quantized by multiplying the entropy decoded coefficient value  $C_{n,m}$  with the quantization bucket size  $q_R$ , but those coefficients that differ from the average in their subblock by not more than one quantization bucket are replaced by the average. One has:

$$R_{(i+k \% 4, j+k / 4)} = \begin{cases} Avg^k & \text{if } |Avg^k - q_R \times S_{ij}^k| < q_R \\ q_R \times S_{ij}^k & \text{otherwise} \end{cases}$$

NOTE – Noise shaping at the encoder side moves the quantization noise into higher frequencies. The adaptive averaging step described here removes the noise pattern introduced by the encoder whenever the noise level is below an expected noise threshold that is due to the noise shaping process. This averaging reduces the spatial resolution of the image but improves on average the amplitude resolution of the residual data and thus reduces potential visual artifacts due to quantization in the spatial domain.

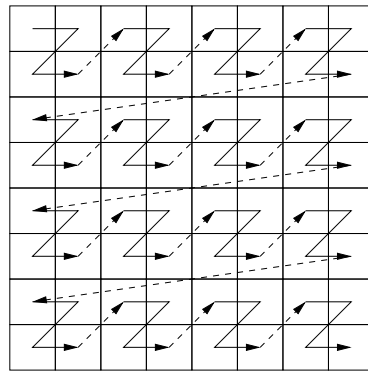
## 5.42 Inverse Quantization (normative)

In case noise shaping is turned off by setting the S parameter of the JPEG Extensions Specification marker to zero, inverse dequantization of the residual data shall be implemented by multiplying the decoded coefficient value  $C_{ij}$  with the quantization bucket size  $q_R$ . The parameter  $q_R$  shall be found as follows: For the component 0, find the quantization matrix whose index is given by the Lt parameter of the JPEG Extensions specification marker; for components 1 and 2, find the quantization matrix with index Ct. Then set  $q_R$  to the entry  $q_{63}$  of the found matrix; all other entries are irrelevant for inverse quantization of residual data. The residual data  $R_{ij}$  shall then be reconstructed by  $R_{ij} = q_R \times C_{ij}$  where  $C_{ij}$  is the decoded symbol generated by the entropy decoder.

NOTE – Quantization matrices are encoded in DQT markers that can be found in the tables of the legacy codestream, or the tables of the continuous codestream of the JPEG Extensions Residual data markers.

## 5.43 Quantization and Noise Shaping of Residual Data (informative)

If the S bit in the JPEG Extensions Specification marker is set, the encoder will include a noise shaping process to reduce the quantization noise at the cost of the spatial resolution of the residual data signal. This subclause describes a noise shaping and quantization process that is compatible to the normative adaptive low-pass filter defined in subclause XXX. For that, let  $R_{ij}$  be the residual signal, organized in an 8×8 block that is to be encoded by the near lossless coding method defined in Annex XXX. Let  $SO(k)$  be the scan order over the block defined in figure J-1. At the start of the scan, set the cumulative quantization error, E, to zero.



**Figure J-1: Scan Pattern for the Noise Shaping of the Residual Data**

Then, for each position  $k = 0..63$  in the scan pattern, compute the following:

- Set  $V = R_{SO(k)} + E$ , i.e. V is the source signal plus the cumulative error.
- Set the quantized signal to  $C_{SO(k)} = \text{sign}(V) \times \lfloor |V|/q_R \rfloor$
- Update the cumulative error  $E = E + R_{SO(k)} - q_R \times C_{SO(k)}$
- Repeat the above steps for  $k=0..63$  in the order of the scan pattern of Figure J-1.

The output of this process is the quantized signal  $C_{ij}$  which is entropy encoded by the algorithm specified in Annex XXX.

NOTE – The rounding convention selected here creates a "fat zero" (deadzone) around the zero bucket that is twice the size of all other buckets. This is intentional and reduces the amount of noise encoded in the residual scan. Other rounding conventions may be tried, and may, depending on the image, also provide better results.

## 5.44 Quantization of Residual Data (informative)

This subclause describes the quantization process of residual data that is used if the S parameter of the JPEG Extensions Specification marker is zero, and hence noise shaping is disabled. This process is suitable for lossless coding if the quantization bucket size  $q_R$  is equal to one. The parameter  $q_R$  is found as follows: For the component 0, find the quantization matrix whose index is given by the Lt parameter of the JPEG Extensions specification marker; for components 1 and 2, find the quantization matrix with index Ct. Then set  $q_R$  to the entry  $q_{63}$  of the found matrix; all other entries are irrelevant for quantization of residual data. The symbols  $C_{ij}$  to be coded by the lossless and near lossless entropy coding process for residual data defined in Annex XXX are then found by dividing the residual signal  $R_{ij}$  by  $q_R$  and rounding to zero:



$$C_{ij} = \text{sign}(R_{ij}) \times \lfloor |R_{ij}| / q_R \rfloor$$

NOTE – Other rounding conventions than the one shown here may be possible and may, depending on the source image, even provide better results.