

**INTERNATIONAL ORGANISATION FOR STANDARDISATION  
ORGANISATION INTERNATIONALE DE NORMALISATION  
ISO/IEC JTC 1/SC 29/WG 4  
MPEG VIDEO CODING**

**ISO/IEC JTC 1/SC 29/WG 4 m62236**

**January 2023, Online**

**Title: [VCM] CE1.7 - Implementation of Poznan University of Technology tools from Proposal C in Video Coding for Machines Reference Software**

**Authors:**

**Marek Domański, Olgierd Stankiewicz, Sławomir Maćkowiak, Tomasz Grajek, Maciej Wawrzyniak, Jakub Stankowski, Sławomir Rózek, Dominik Cywiński, Jakub Szekięda, Jakub Siejak**

**Poznań University of Technology, Poznań, Poland**

**Abstract**

**1 Introduction**

This document described RoI-based simplification and retargeting method implemented in VCM-RS software [1]. The method originally has been implemented in Poznan University of Technology responses (proposal C) [2,3,4] to “Call for Proposals for Video Coding for Machines” [CFP].

The important part of the contribution is the detailed description of the bitsream syntax.

## 2 General description

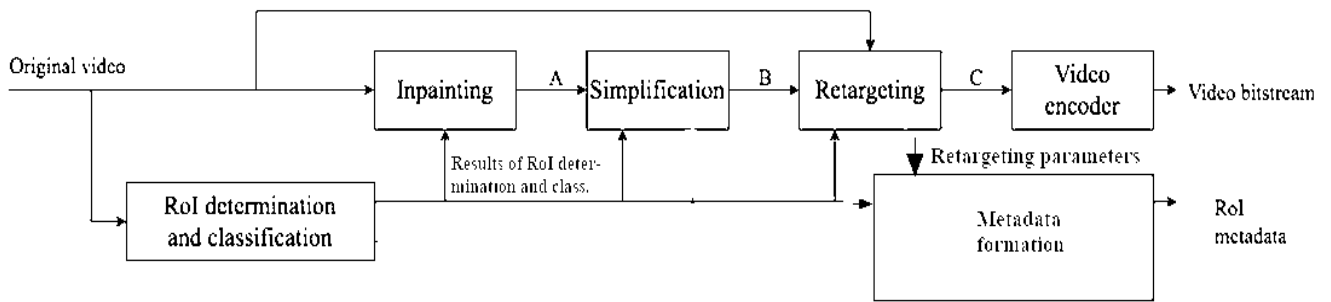


Fig. 1. The structure of the encoder

The block “Video encoder” is now called as inner encoder, here this is a VVC encoder.

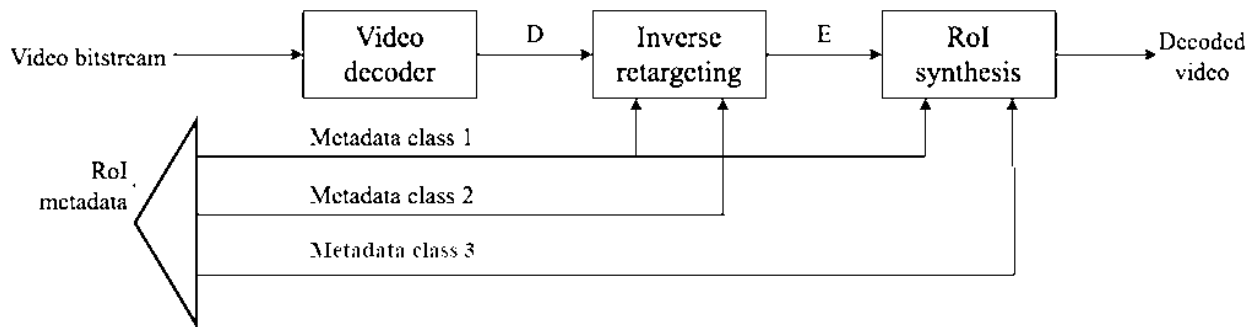


Fig. 2. The structure of the decoder.

### Description of the building blocks in a VCM encoder:

- a) ***RoI determination and classification***: In this block the original image/video is analyzed and potential regions of interest (RoI) are determined in the areas where interesting objects are located. This is done using a special deep neural network, specifically: detectron network. Then the RoIs are classified as corresponding to important objects, as corresponding to small objects and maybe into other classes.
- b) ***Simplification***: The simplification is aimed at production of the content that may be represented with reduced number of bits after compression without any significant deterioration of the efficiency of the machine vision tasks executed on decoded (decompressed) image/video. In principle the simplification reduces details in the background whereas the objects and their neighborhoods remain untouched or only slightly

simplified. In particular, the results of RoI determination and classification are used to identifying objects, the bands around objects and the background. Here, we do consider the objects already replaced by simple content by the inpainting. The background is simplified, and possibly some bands around selected objects or even some objects (usually the objects of big size – relatively in an image or in a video frame). The goal of simplification is to produce such image/video that can be easily represented by a very small number of bits (in compressed representation) whereas the object still can be well recognized (detected, tracked etc.). This may be optimized in context of particular class of object within given RoI.

c) **Retargeting**: The number of lines and columns of the input image or the video frame may be reduced according to the content identified in the process of RoI determination and classification. The lines and columns of samples are mostly removed when their do not comprise samples of objects. The lines and rows that comprise object samples may be reduced according to analysis of the objects. This reduction correspond to downsampling the object that must not deteriorate the ability to recognize the objects, therefore it is possible if the object is oversampled with respect of its texture and other features. The sample (line / column) removal must be preceded by the respective local lowpass (anti-aliasing) filtering. In classical retergating methods energy function has to be formulated which controls which rows and columns are primarily removed. Such would require signalization of the energy function (image) to the decoder. In our proposal a simple yet efficient retargeting technique is used which omits implicit energy image formulation, and instead rectangle-shaped RoI are used. Therefore, only coordinates of RoIs are transmitted.

d) **Object modeling using a dictionary** :

Some of the objects detected (of size lesser that 200x200 pixels) are classified as dictionary objects. Those objects are inpainted and not sent through inner codec (VVC). Instead, those objects are signalled as indices pointing to patches from the dictionary (known to the decoder). In the decoder those objects are inserted to the decoded image.

### 3 Implementation details

The tool has been implemented as ROI component in VCM-RS [1] software.

- `encoder_poznanroi.py` – main encoder component file
- `decoder_poznanroi.py` – main decoder component file
- `poznanroi_config_codec.py` – hard-coded parameters of the codec e.g. blurring parameters.
- `poznanroi_config_classes.py` – classes of objects used in the encoder
- `poznanroi_detectron.py` – interface to detectron network
- `poznanroi_blur.py` – implementation of simplification (blurring) algorithm
- `poznanroi_retarget.py` – implementation of retargeting algorithm
- `poznanroi_utils.py` – utilities related to general processing.
- `poznanroi_rois.py` - utilities related to RoIs
- `poznanroi_objs.py` - utilities related to objects
- `poznanroi_bitenc.py` – bitstream encoding utilities
- `poznanroi_bitdec.py` – bitstream decoding utilities

The current implementation of the tools has additional requirements:

- Detectron framework
- Tracker library
- Bitstring python library – bitstream encoding of RoIs

It can be noted that the current implementation does NOT require GPU for operation (CPU only). It however employs GPU if available to speed-up processing. We have notices negligible differences in encoder performance when CPU is used vs GPU.

## 4 Syntax

Table 1. VCM RoI information syntax.

Syntax element	C	Descriptor
VCM_info( payloadSize ) {		
bits_size_x	5	u
bits_size_y	5	u
output_image_size_x_minus_1	bits_size_x	u
output_image_size_y_minus_1	bits_size_y	u
if (codec_option_use_roi_retargetting) {		
roi_retargetting_info()		
}		
if (codec_option_use_dictionary_objects) {		
dictionary_objects_info()		
}		

bits\_size\_x - number of bits on which output\_image\_size\_x\_minus\_1 syntax element is signalled.

bits\_size\_y - number of bits on which output\_image\_size\_y\_minus\_1 syntax element is signalled.

output\_image\_size\_x\_minus\_1 - size of the original image to be outputted from the VCM decoder, minus 1.

output\_image\_size\_y\_minus\_1 - size of the original image to be outputted from the VCM decoder, minus 1.

codec\_option\_use\_roi\_retargetting - flag active in all Poznan University of Technology responses to CFP: A, B and C

codec\_option\_use\_dictionary\_objects - flag active only in Poznan University of Technology response C to CFP

Table 2. RoI retargeting information syntax.

Syntax element	C	Descriptor
roi_retargeting_info() {		
<b>hierarchy_power_base</b>	8	u
if (hierarchy_power_base != 0) {		
<b>bits_scale</b>	3	u
<b>bits_pos_x</b>	5	u
<b>bits_pos_y</b>	5	u
<b>bits_size_x</b>	5	u
<b>bits_size_y</b>	5	u
<b>bg_scale</b>	bits_scale	u
<b>bg_size_x_minus_1</b>	bits_pos_x	u
<b>bg_size_y_minus_1</b>	bits_pos_y	u
<b>bits_num_rois</b>	5	5
<b>num_rois</b>	bits_num_rois	u
current_scale = 0		
for (i = 0; i < num_rois; i++ ) {		
if (current_scale!=0) {		
update_scale	1	u
if (update_scale):		
<b>current_scale</b>	bits_scale	u
}		
roi_scale[i] = current_scale		
<b>roi_pos_x[i]</b>	bits_pos_x	u
<b>roi_pos_y[i]</b>	bits_pos_y	u
<b>roi_size_x_minus_1[i]</b>	bits_size_x	u
<b>roi_size_y_minus_1[i]</b>	bits_size_y	u

}		
}		

hierarchy\_power\_base - base of power function basing on which retargetting scaling is calculated

bits\_scale - number of bits on which bg\_scale and roi\_scale[i] syntax elements are signaled

bits\_pos\_x - number of bits on roi\_pos\_x[i] syntax element is signaled

bits\_pos\_y - number of bits on roi\_pos\_y[i] syntax element is signaled

bits\_size\_x - number of bits on roi\_size\_x\_minus\_1[i] syntax element is signaled

bits\_size\_y - number of bits on roi\_size\_y\_minus\_1[i] syntax element is signaled

bg\_scale - scale level of the background (highest hierarchy RoI element)

bg\_size\_x\_minus\_1 - size (minus 1) of the background (highest hierarchy RoI element)

bg\_size\_y\_minus\_1 - size (minus 1) of the background (highest hierarchy RoI element)

bits\_num\_rois - number of bits on which num\_rois syntax elements is signaled

num\_rois - number of RoIs (excluding background)

current\_scale - temporary decoding variable, used to update successive scales: roi\_scale[i]

roi\_scale[i] - scale of RoI index i

roi\_pos\_x[i] - position of RoI index i

roi\_pos\_y[i] - position of RoI index i

roi\_size\_x\_minus\_1[i] - width (minus 1) of RoI index i

roi\_size\_y\_minus\_1[i] - height(minus 1) of RoI index i

Table 3. Dictionary objects information syntax.

Syntax element	C	Descriptor
dictionary_objects_nfo() {		
<b>bits_num_dictobjs</b>	5	u
if (bits_num_dictobjs!=0) {		
<b>num_dictobjs</b>	bits_num_dictobj	u
if (num_dictobjs!=0) {		
<b>bits_class_id</b>	3	u
<b>bits_patch_id</b>	4	u
<b>bits_start_x_min</b>	5	u
<b>bits_start_y_min</b>	5	u
<b>bits_start_width</b>	5	u
<b>bits_start_height</b>	5	u
<b>is_sequence</b>	1	u
if (is_sequence) {		
<b>bits_start_frame</b>	4	u
<b>bits_num_frames</b>	4	u
<b>bits_v_x_min</b>	5	u
<b>bits_v_y_min</b>	5	u
<b>bits_v_width</b>	5	u
<b>bits_v_height</b>	5	u
}		
for (i = 0; i < num_dictobjs; i++) {		
<b>dictobj_class_id[i]</b>	bits_class_id	u
<b>dictobj_patch_id[i]</b>	bits_patch_id	u
<b>dictobj_flipped[i]</b>	1	
<b>dictobj_start_x_min[i]</b>	bits_start_x_min	u
<b>dictobj_start_y_min[i]</b>	bits_start_y_min	u
<b>dictobj_start_width[i]</b>	bits_start_width	u
<b>dictobj_start_height[i]</b>	bits_start_height	u
if is_sequence:		
<b>dictobj_start_frame[i]</b>	bits_start_frame	u
<b>dictobj_num_frames[i]</b>	bits_num_frames	u
<b>dictobj_v_x_min[i]</b>	bits_v_x_min	s
<b>dictobj_v_y_min[i]</b>	bits_v_y_min	s
<b>dictobj_v_width[i]</b>	bits_v_width	s
<b>dictobj_v_height[i]</b>	bits_v_height	s
}		
}		
}		



bits\_num\_dictobjs - number of bits on which num\_dictobjs syntax element is signaled

num\_dictobjs - number of dictionary objects

bits\_class\_id - number of bits on dictobj\_class\_id[i] syntax element is signaled

bits\_patch\_id - number of bits on dictobj\_patch\_id[i] syntax element is signaled

bits\_start\_x\_min - number of bits on dictobj\_start\_x\_min[i] syntax element is signaled

bits\_start\_y\_min - number of bits on dictobj\_start\_y\_min[i] syntax element is signaled

bits\_start\_width - number of bits on dictobj\_start\_width[i] syntax element is signaled

bits\_start\_height - number of bits on dictobj\_start\_height[i] syntax element is signaled

is\_sequence - flag specifying whether velocity parameters of dictionary objects are transmitted (or whether they are static across GOP).

bits\_start\_frame - number of bits on dictobj\_start\_frame[i] syntax element is signaled

bits\_num\_frames - number of bits on dictobj\_num\_frames[i] syntax element is signaled

bits\_v\_x\_min - number of bits on dictobj\_v\_x\_min[i] syntax element is signaled

bits\_v\_y\_min - number of bits on dictobj\_v\_y\_min[i] syntax element is signaled

bits\_v\_width - number of bits on dictobj\_v\_width[i] syntax element is signaled

bits\_v\_height - number of bits on dictobj\_v\_height[i] syntax element is signaled

dictobj\_class\_id[i] - index of class from dictionary for object i

dictobj\_patch\_id[i] - index of texture patch from dictionary for class dictobj\_class\_id[i], for object i

dictobj\_flipped[i] – value 1 means object is horizontally flipped

dictobj\_start\_x\_min[i] - starting position of object i

dictobj\_start\_y\_min[i] - starting position of object i

dictobj\_start\_width[i] - starting width object i

dictobj\_start\_height[i] - starting height object i

dictobj\_start\_frame[i] - index of first frame in GOP, when object i appears

dictobj\_num\_frames[i] - number of frames in GOP, in which object i appears

dictobj\_v\_x\_min[i] - velocity (differential change) of starting position of object i

dictobj\_v\_y\_min[i] - velocity (differential change) of starting position of object i

dictobj\_v\_width[i] - velocity (differential change) of width of object i

dictobj\_v\_height[i] - velocity (differential change) of height of object i

## 5 References

- [1] Honglei Zhang , "Introduction to the VCM reference software (VCM-RS) ",ISO/IEC JTC 1/SC 29/WG 4 m62003, January 2023, Online
- [2] Marek Domański, Olgierd Stankiewicz, Sławomir Maćkowiak, Sławomir Rózek, Tomasz Grajek, Jakub Szekięda, Dominik Cywiński, Jakub Siejak, "[VCM] Poznań University of Technology Proposal A in response to CfP on Video Coding for Machines", ISO/IEC JTC 1/SC 29/WG 2 m60727, Online – October 2022.
- [3] Marek Domański, Olgierd Stankiewicz, Sławomir Maćkowiak, Sławomir Rózek, Tomasz Grajek, Jakub Szekięda, Dominik Cywiński, Jakub Siejak, "[VCM] Poznań University of Technology Proposal B in response to CfP on Video Coding for Machines", ISO/IEC JTC 1/SC 29/WG 2 m60728, Online – October 2022.
- [4] Marek Domański, Olgierd Stankiewicz, Sławomir Maćkowiak, Sławomir Rózek, Tomasz Grajek, Jakub Szekięda, Dominik Cywiński, Jakub Siejak, "[VCM] Poznań University of Technology Proposal C in response to CfP on Video Coding for Machines", ISO/IEC JTC 1/SC 29/WG 2 m60729, Online – October 2022.
- [5] "Call for Proposals for Video Coding for Machines" wg2n00191
- [6] ISO/IEC JTC1/SC29/WG5, "VVC Reference Model (VTM)"