Lecture 18: Image and Video Compression

Yuhao Zhu

http://yuhaozhu.com yzhu@rochester.edu CSC 259/459, Fall 2025 Computer Imaging & Graphics

The Roadmap

Theoretical Preliminaries

Human Visual Systems

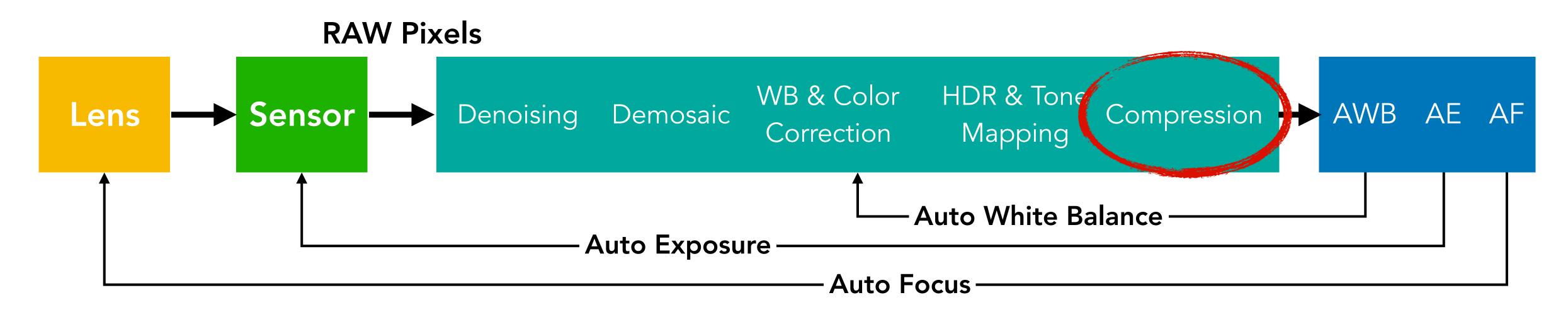
Display and Camera

Modeling and Rendering

Sources of Color

Display and Lighting
Photographic Optics
Image Sensor
Image Signal Processing
Image/Video Compression
Immersive Content

Compression



Most of cameras and smartphone allow you to shoot in RAW, i.e., exporting RAW images without any post-processing (in camera-native RGB space).

Otherwise, images/videos you get are compressed in the color space of the output device, mostly sRGB.



30-second video @ 1080p resolution (1920 x 1080 pixels per frame) @ 30 frames per second (FPS) 3 colors per pixel + 1 byte per color \rightarrow 6.2 MB/frame \rightarrow 6.2 MB x 30 s x 30 FPS = 5.2 GB total size Actual H.264 video file size: 65.4 MB (80-to-1 compression ratio).

Compression/encoding done in real-time without you even realizing it!





Basic Concepts and Ideas of (Any) Compression

Goal: reduce data size while maintaining high (visual) quality

Lossless compression vs. lossy compression

Two main techniques:

- Lossless: removing redundancies.
- Lossy: sacrificing "unimportant" details. In the context of image/video compression, "importance" is usually dictated by human perception.

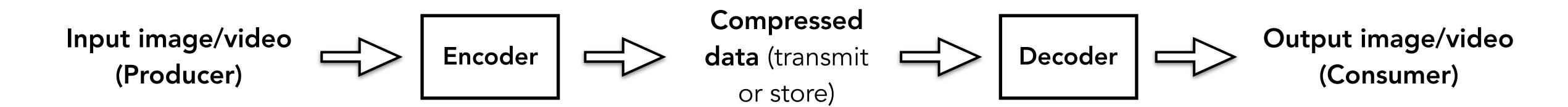


Image and Video Compression in Hardware

Nvidia
Xavier SoC

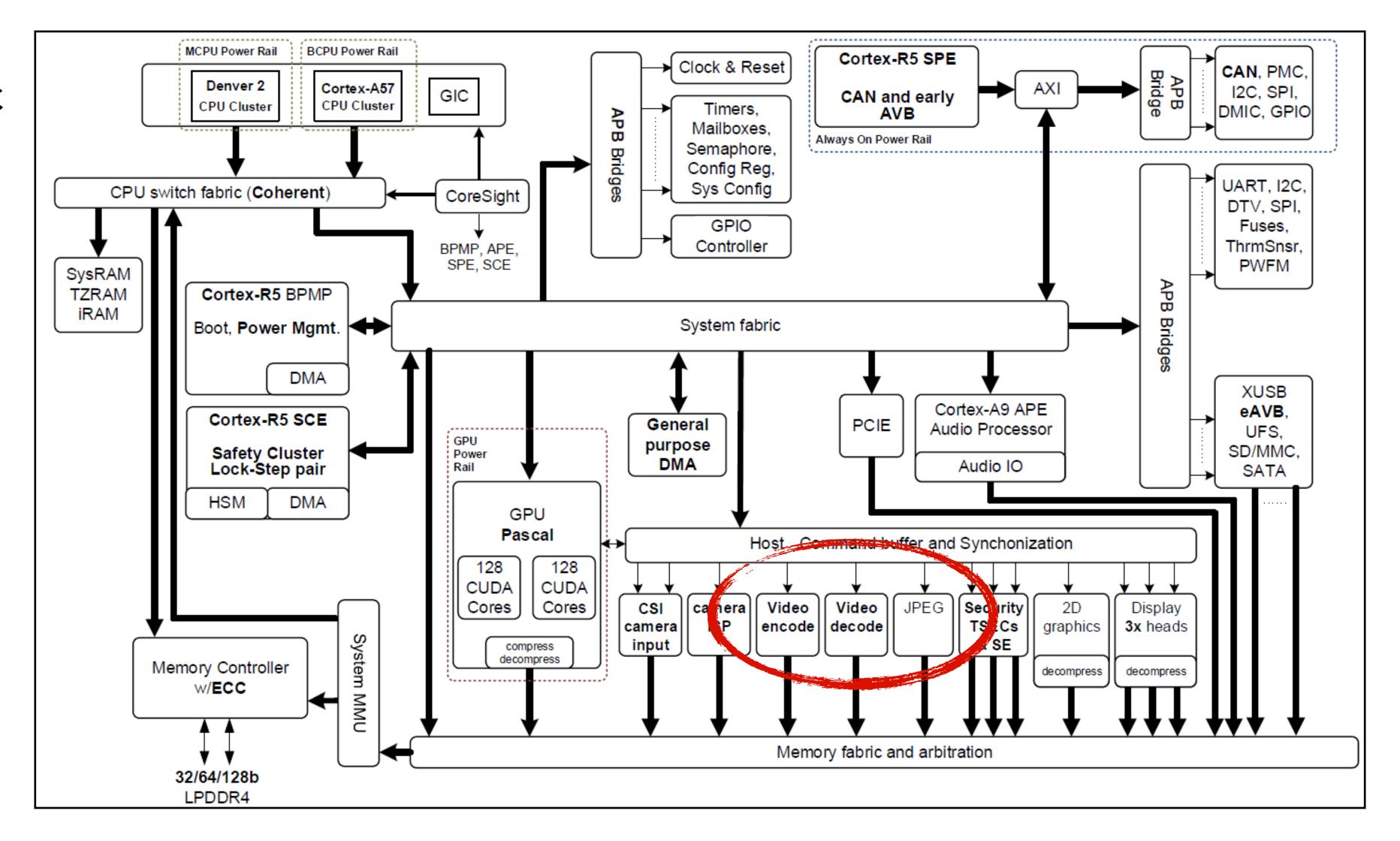
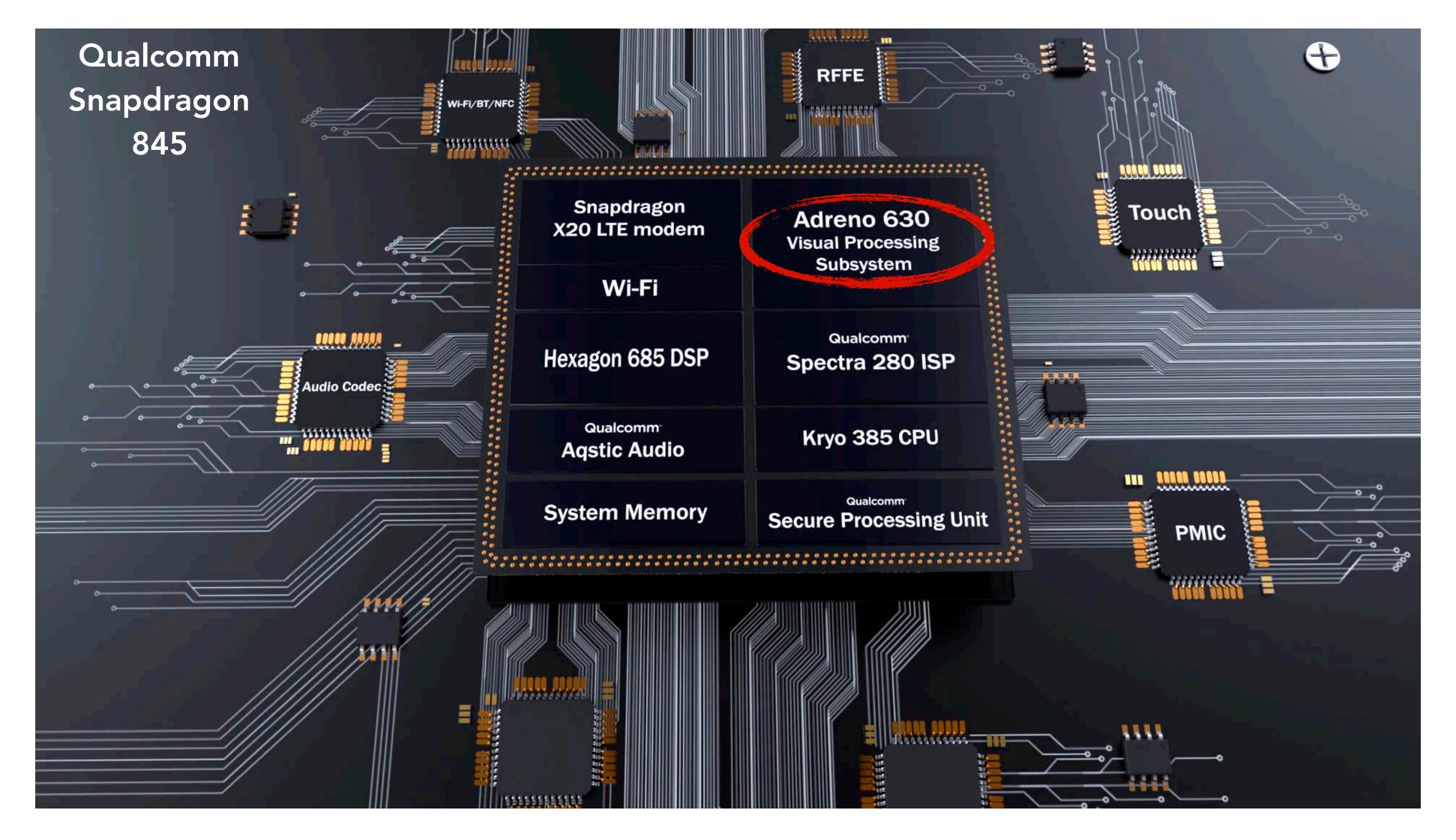


Image and Video Compression in Hardware



JPEG Image Compression

Two Big Ideas in JPEG Image Compression

JPEG is lossy, so must choose wisely what information to sacrifice.

Idea 1: retain luminance; subsample "colors".

- Luminance (brightness) encodes visual details.
- We are more sensitive to brightness differences than to chromatic differences; we can afford to lose information in colors, but not so much in luminance.

Idea 2: retain low frequency information; sacrifice high frequency information.

• Human visual system has a frequency threshold (photoreceptors are doing spatial sampling), and high-frequency information looks "busy"/"noisy" anyway, so a bit of inaccuracy/corruption in pixel values isn't very noticeable.

How Do We Separate Luminance and "Colors"

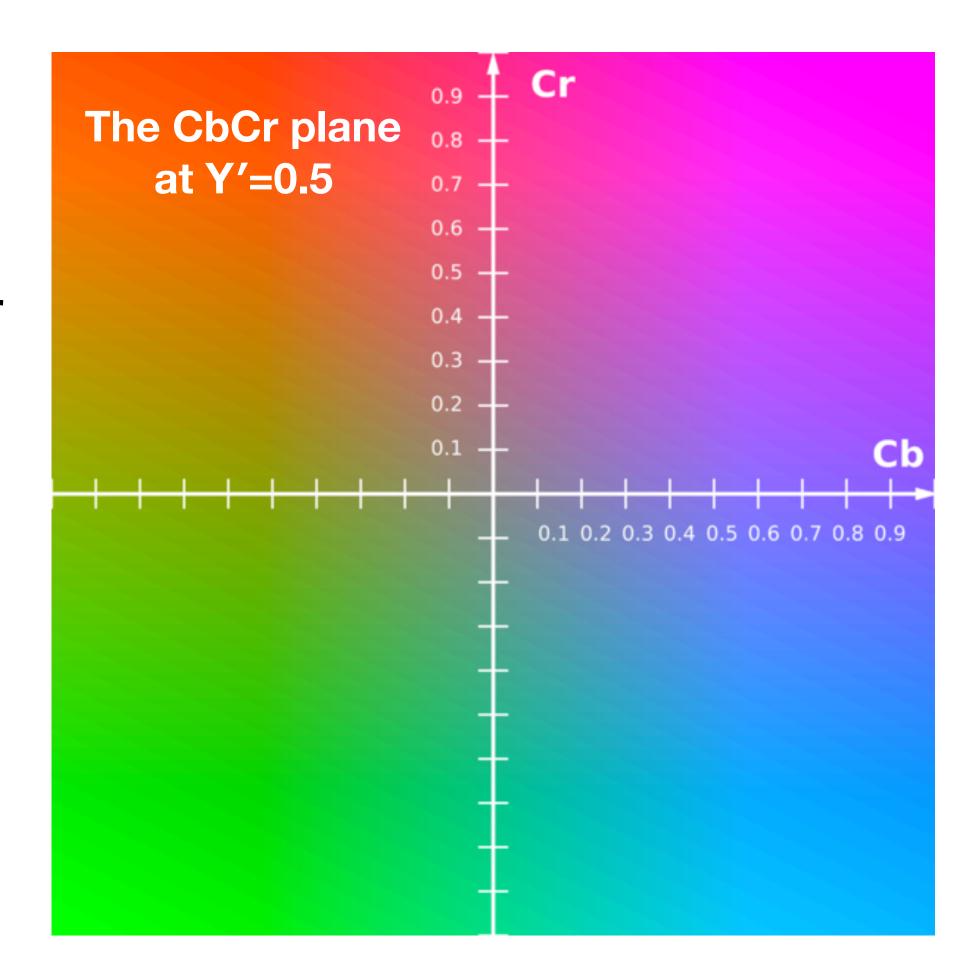
Recall Hering's Opponent Processes:

Red-Green, Blue-Yellow, Light-Dark

The perceptual opponent space is *not* a linear transformation from, say, cone space

• There are neural opponent spaces that are linear w.r.t. cone space (e.g., the DKL space)

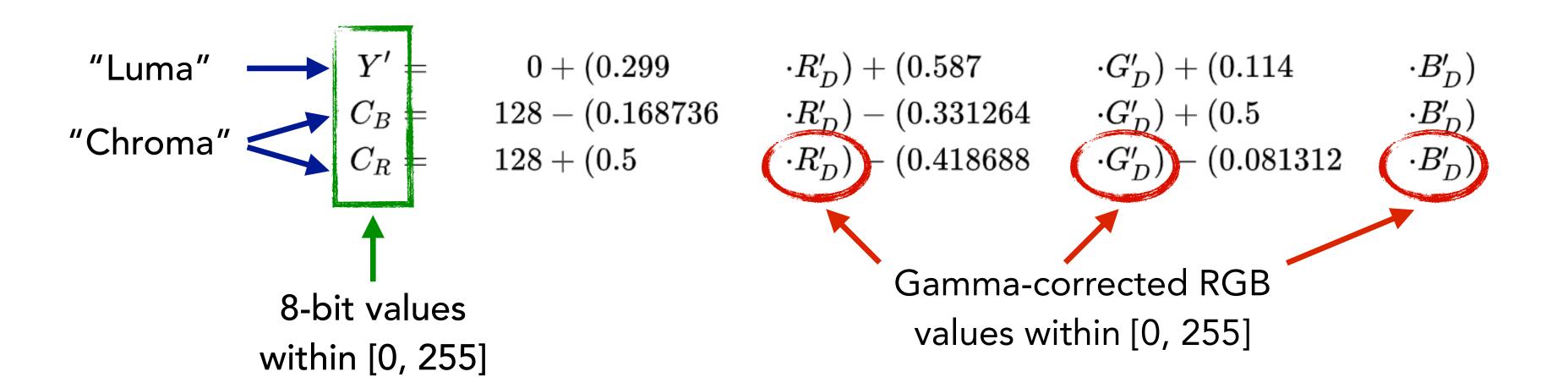
JPEG uses an empirical perceptual opponent space that is a linear transformation from sRGB



JPEG's Y'CbCr Color Space

RGB

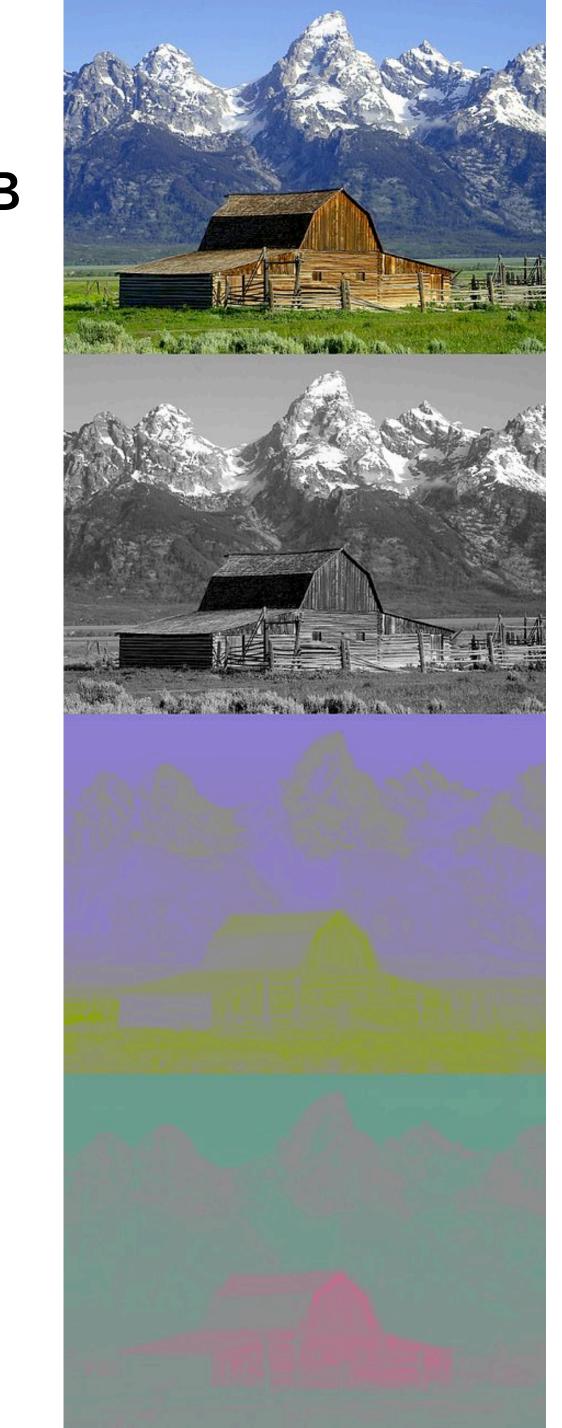




Cb

* Gamma correction here means the RGB values are not luminance-linear

Cr



Chroma Subsampling (Lossy)

Downsample the two chroma components but keep luma unchanged.

4:4:4 representation (no chroma subsampling, lossless)

Y'	Y'	Y'	Y'
Cb	Cb	Cb	Cb
Cr	Cr	Cr	Cr
Y'	Y'	Y'	Y'
Cb	Cb	Cb	Cb
Cr	Cr	Cr	Cr

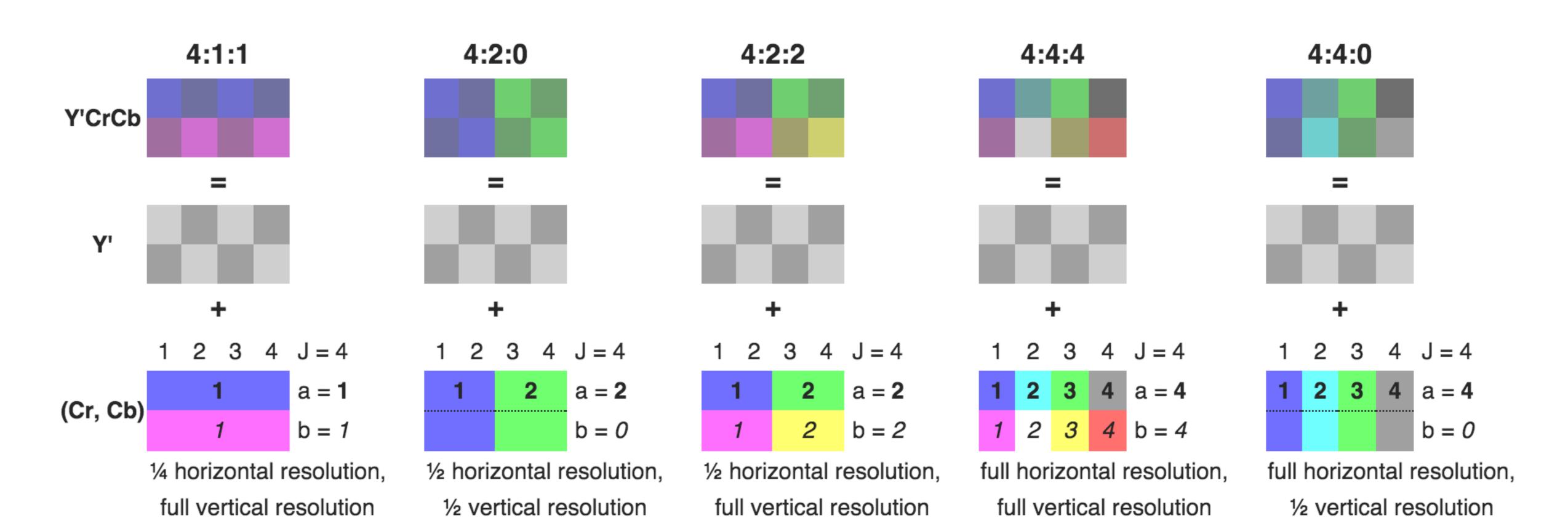
4:2:2 representation (Cb and Cr are subsampled 2X)

Y'	Y'	Y'	Y'
Cb		Cb	
Cr		Cr	
Y'	Y'	Y'	Y'
Y' Cb	Y'	Y' Cb	Y'

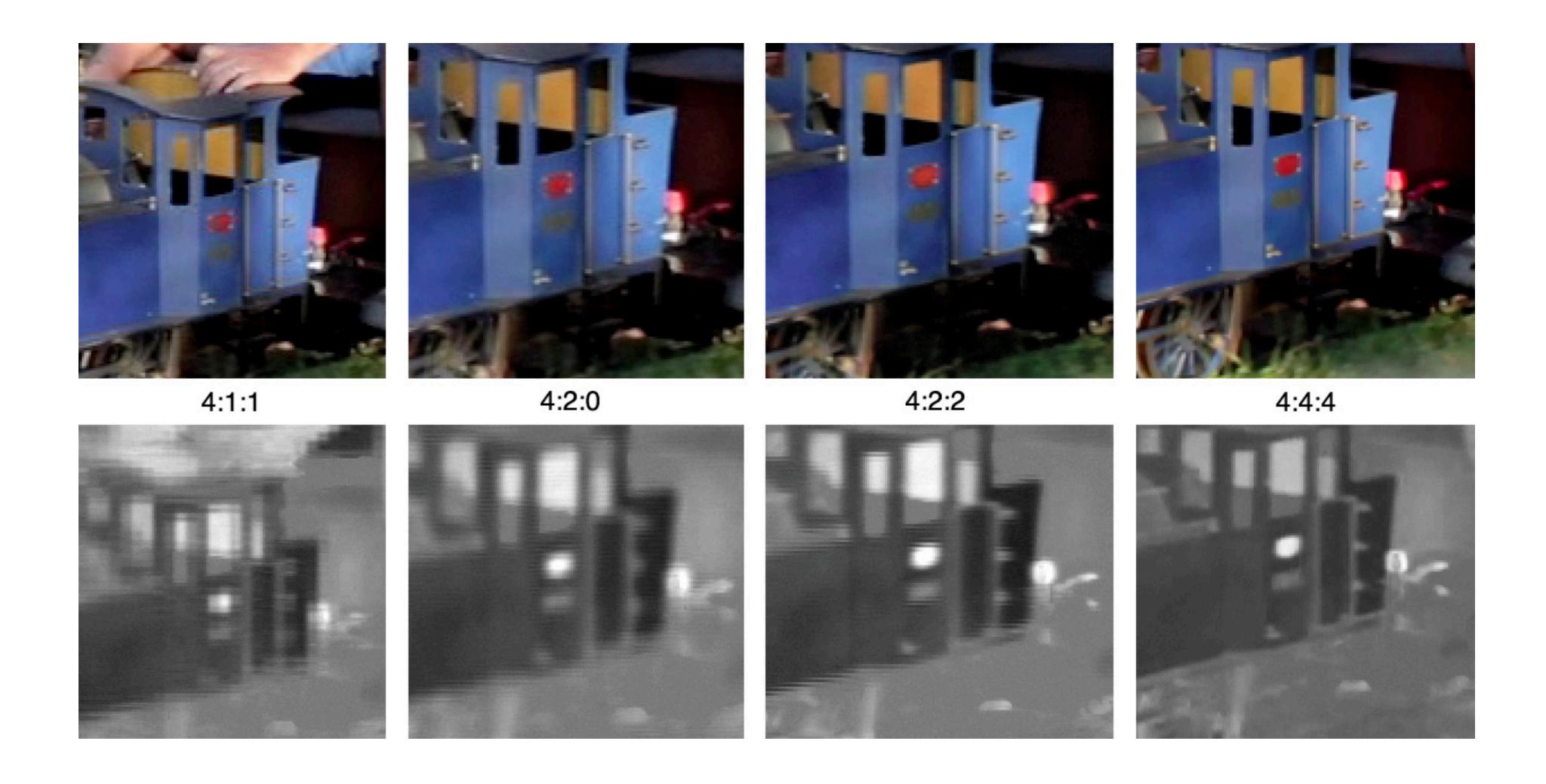
4:2:0 representation (Cb and Cr are subsampled 4X)

Y' Cb	Y'	Y' Cb	Y'		
Cr Y'	Y'	Cr Y'			

Chroma Subsampling (Lossy)

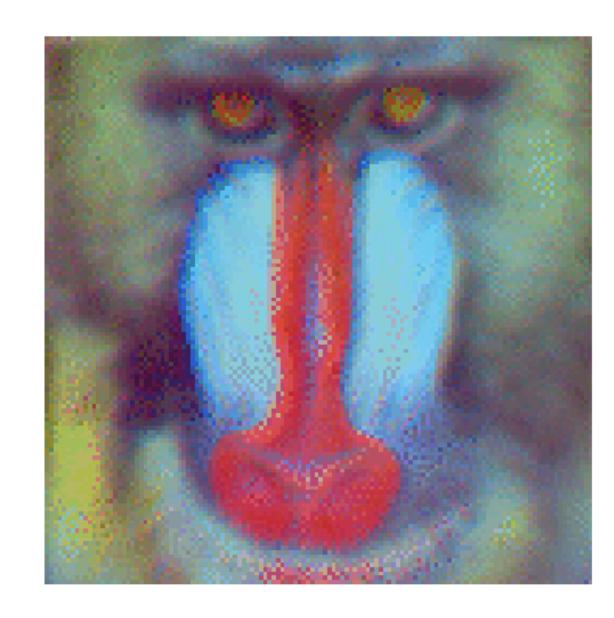


Examples

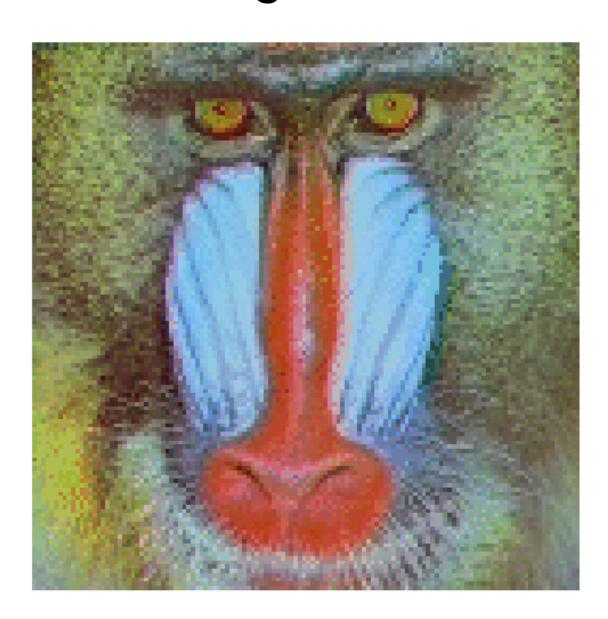


HVS is sensitive to blurring in dark-light channel

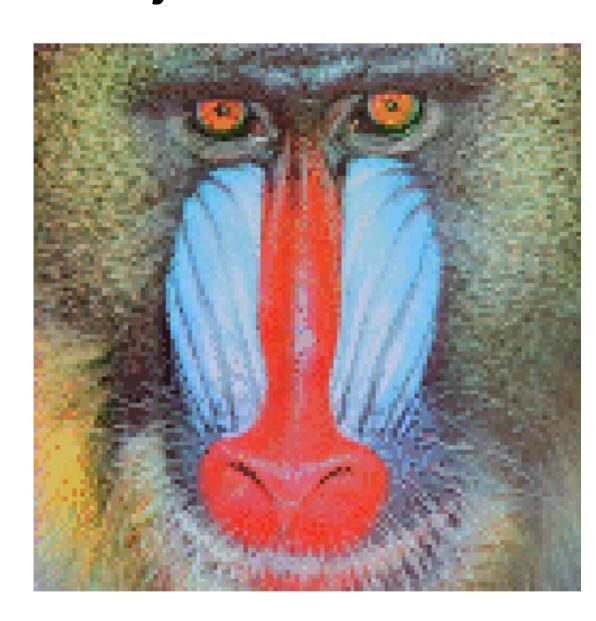
Blur dark-light channel



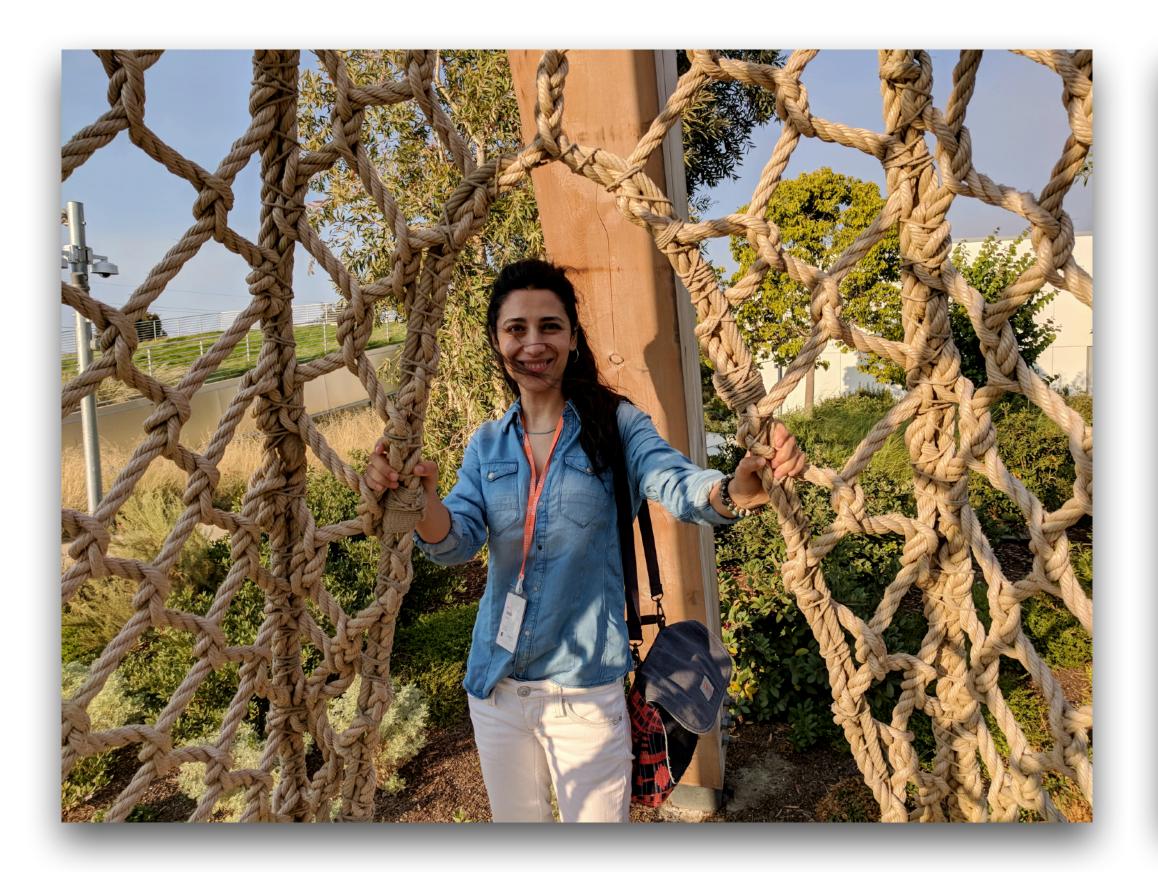
Blur red-green channel

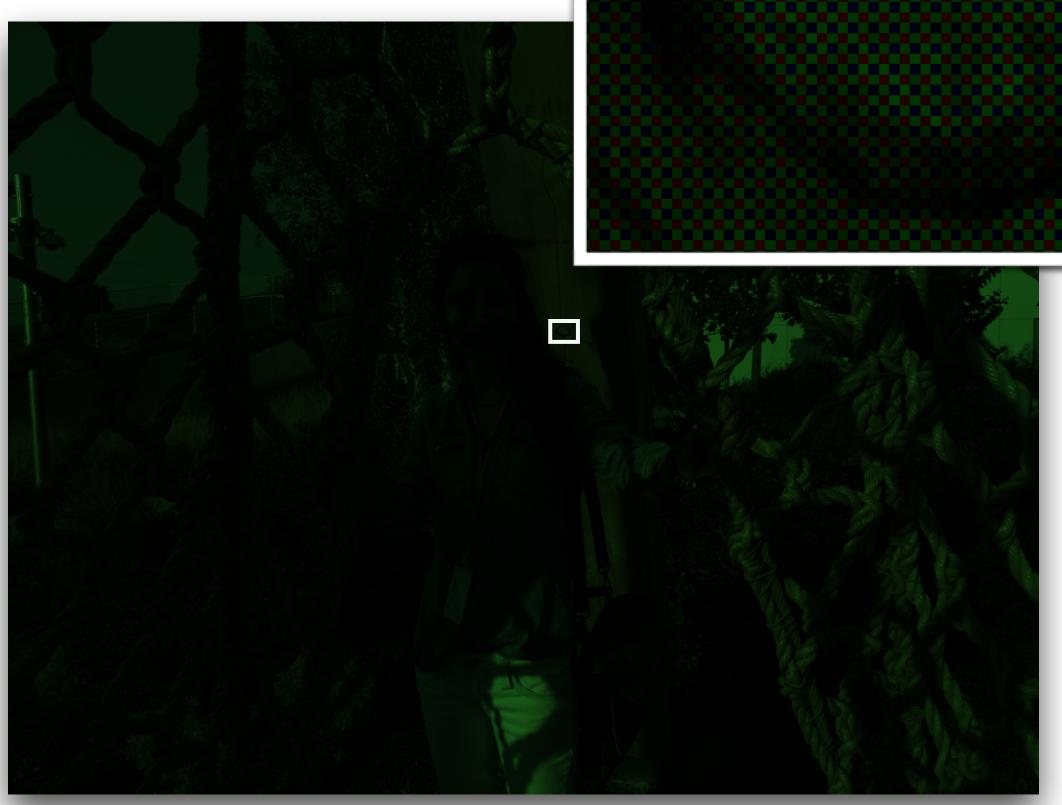


Blur yellow-blue channel



Subsampling in RAW RGB Space?





A bayer filter is a simple form of directly subsampling in the RGB color space. Details are visibly lost — that's why we need demosaicking!

High-Frequency Information Compression

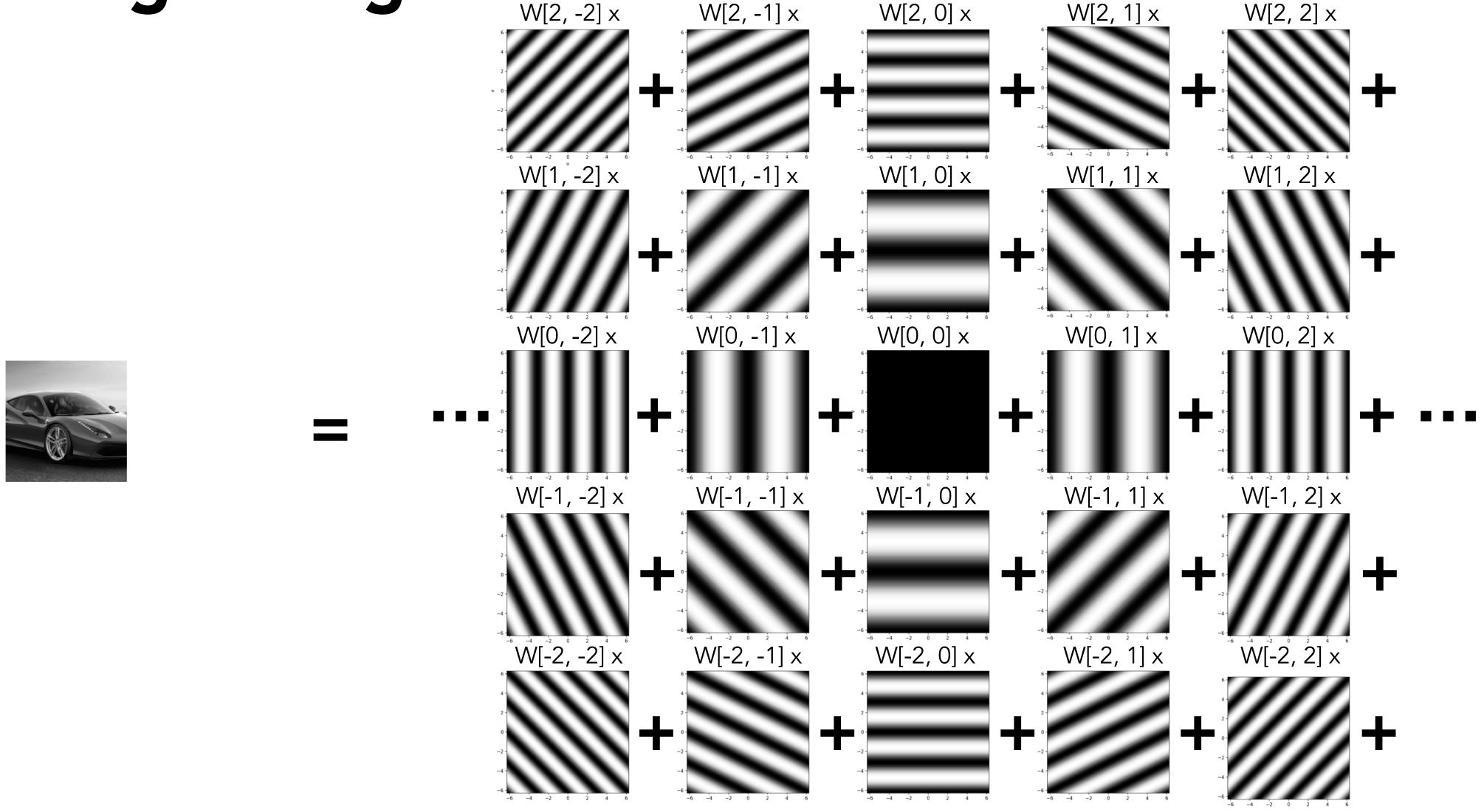
There is a highest frequency threshold beyond which human visual system can't see.

Because photoreceptors are performing a spatial sampling

RGB/Y'CbCr represent spatial domain information. Convert to frequency domain for compression.

Recall Fourier transform: any periodic function can be exactly represented as a weighted sum of simple sinusoids.

Decomposing 2D Signals W[2,-2] x

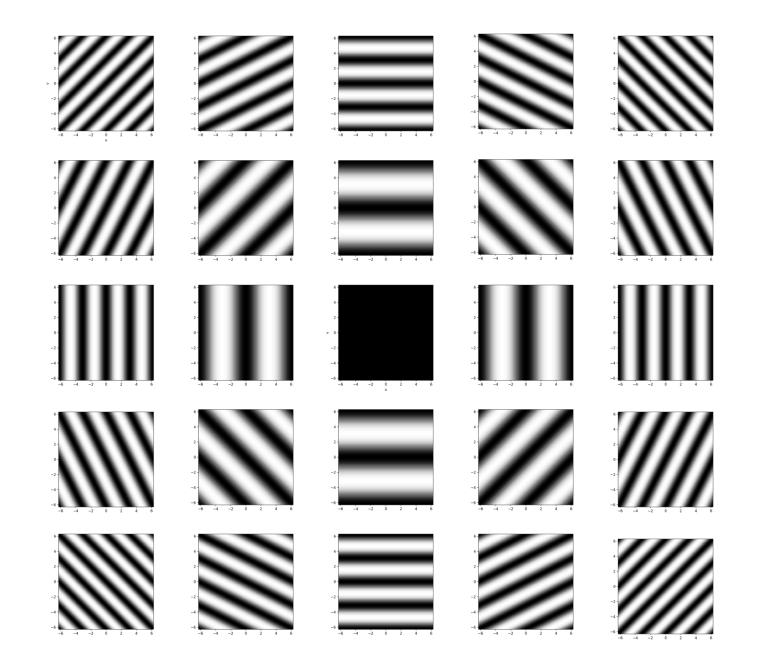


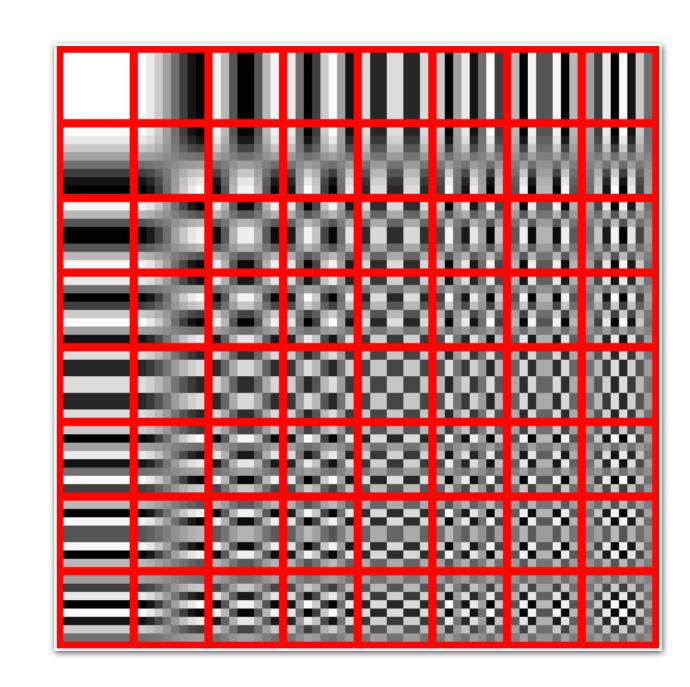
High-Frequency Information Compression

JPEG uses **Discrete Cosine Transformation** (DCT), just a different set of basis functions than what Discrete Fourier Transform uses.

• DCT (in image compression): any 8 ×8 zero-centered image can be represented by a weighted sum of the 64 8×8 images (basis functions).

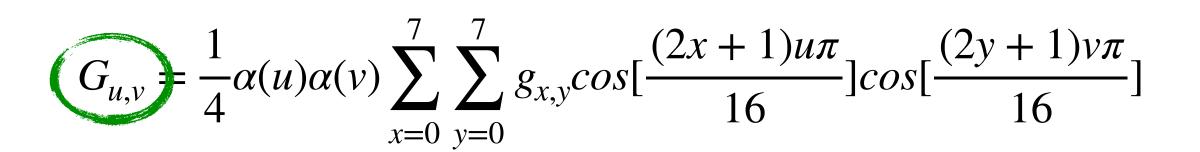
Basis
Functions
used in DFT





Basis
Functions
used in DCT

DCT (Lossless)

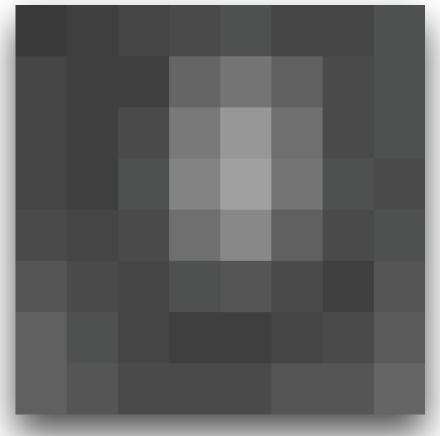


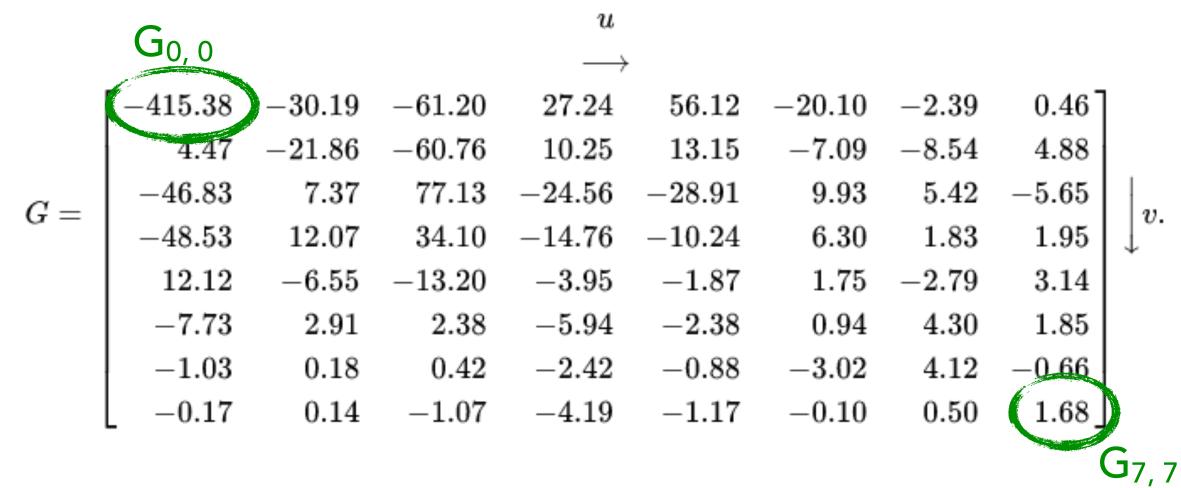
Zero-centered

$$g = \begin{bmatrix} -76 & -73 & -67 & -62 & -58 & -67 & -64 & -55 \\ -65 & -69 & -73 & -38 & -19 & -43 & -59 & -56 \\ -66 & -69 & -60 & -15 & 16 & -24 & -62 & -55 \\ -65 & -70 & -57 & -6 & 26 & -22 & -58 & -59 \\ -61 & -67 & -60 & -24 & -2 & -40 & -60 & -58 \\ -49 & -63 & -68 & -58 & -51 & -60 & -70 & -53 \\ -43 & -57 & -64 & -69 & -73 & -67 & -63 & -45 \\ -41 & -49 & -59 & -60 & -63 & -52 & -50 & -34 \end{bmatrix}$$









Weights (coefficients)

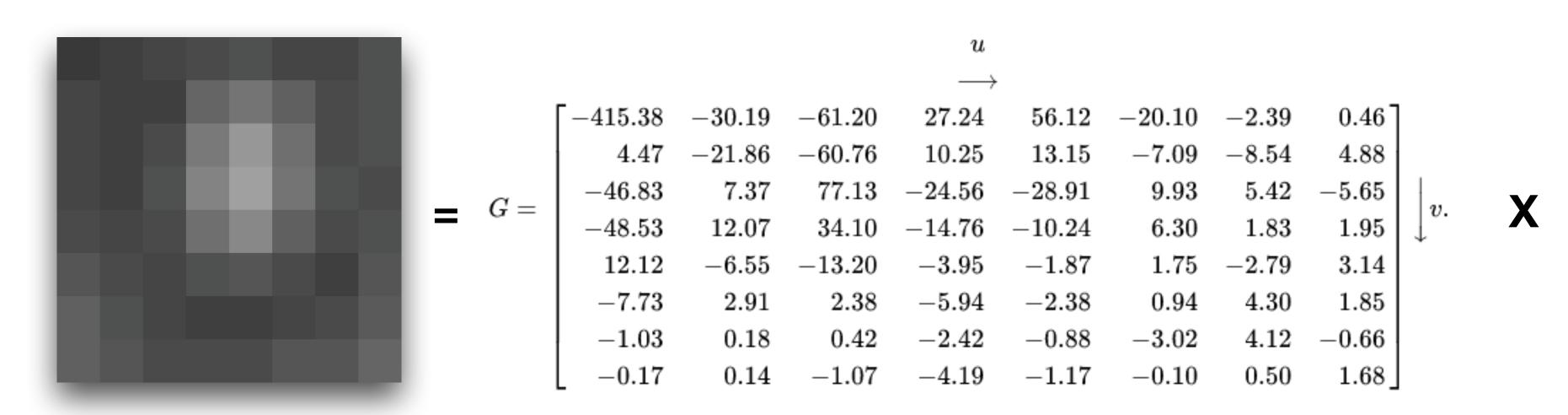
DCT (Lossless)

Calculating the DCT coefficients is mathematically lossless, but could be lossy if the computation isn't done using enough precision.

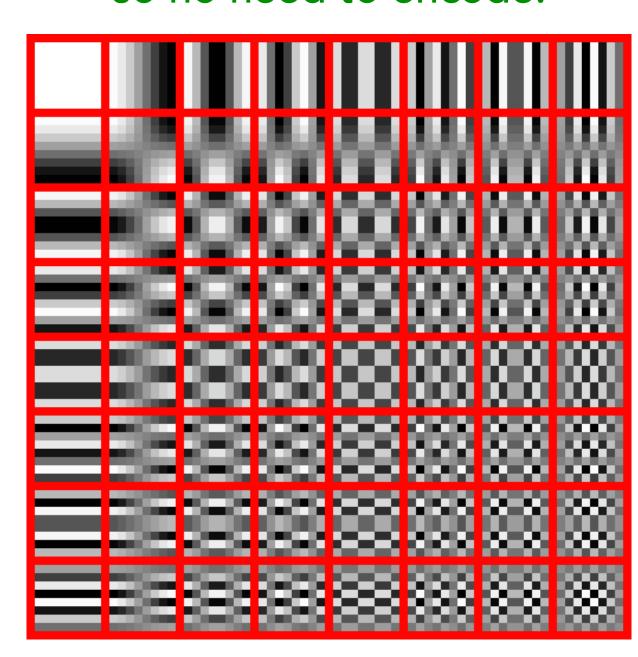
• Often 16-bit computation is needed even for 8-bit images.

An 8x8 block.

Weights (coefficients)

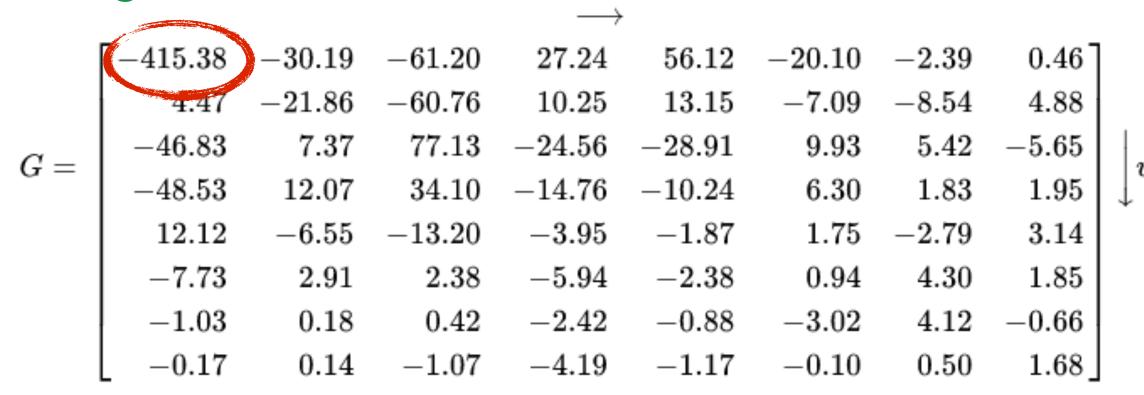


The 8x8 basis functions. Fixed, so no need to encode.



Coefficient Quantization (Lossy)

Weights (coefficients)



B and G aren't equivalent (lossy), but B is small fix-point numbers and thus requires fewer bits to encode.

Quantization matrix (fixed, no need to encode)

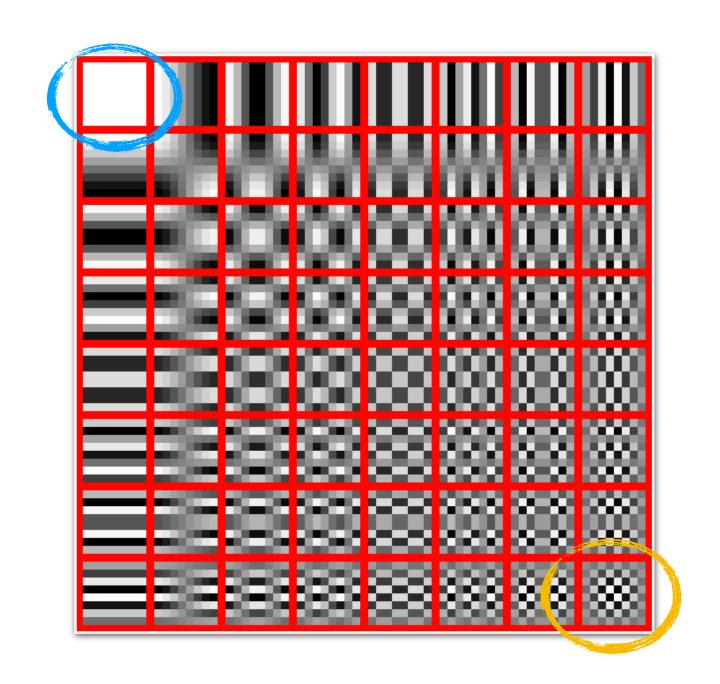
$$Q = \begin{bmatrix} 16 & 11 & 10 & 16 & 24 & 40 & 51 & 61 \\ 12 & 12 & 14 & 19 & 26 & 58 & 60 & 55 \\ 14 & 13 & 16 & 24 & 40 & 57 & 69 & 56 \\ 14 & 17 & 22 & 29 & 51 & 87 & 80 & 62 \\ 18 & 22 & 37 & 56 & 68 & 109 & 103 & 77 \\ 24 & 35 & 55 & 64 & 81 & 104 & 113 & 92 \\ 49 & 64 & 78 & 87 & 103 & 121 & 120 & 101 \\ 72 & 92 & 95 & 98 & 112 & 100 & 103 & 99 \end{bmatrix}$$

$$B_{u,v} = round(\frac{G_{u,v}}{Q_{u,v}})$$

Quantization Matrix

JPEG standard specifies different quantization matrices at different quality levels. Lower quality means larger magnitudes in the matrix (quantize more).

• The exact values in the matrix are derived from modeling human perception to different frequencies.

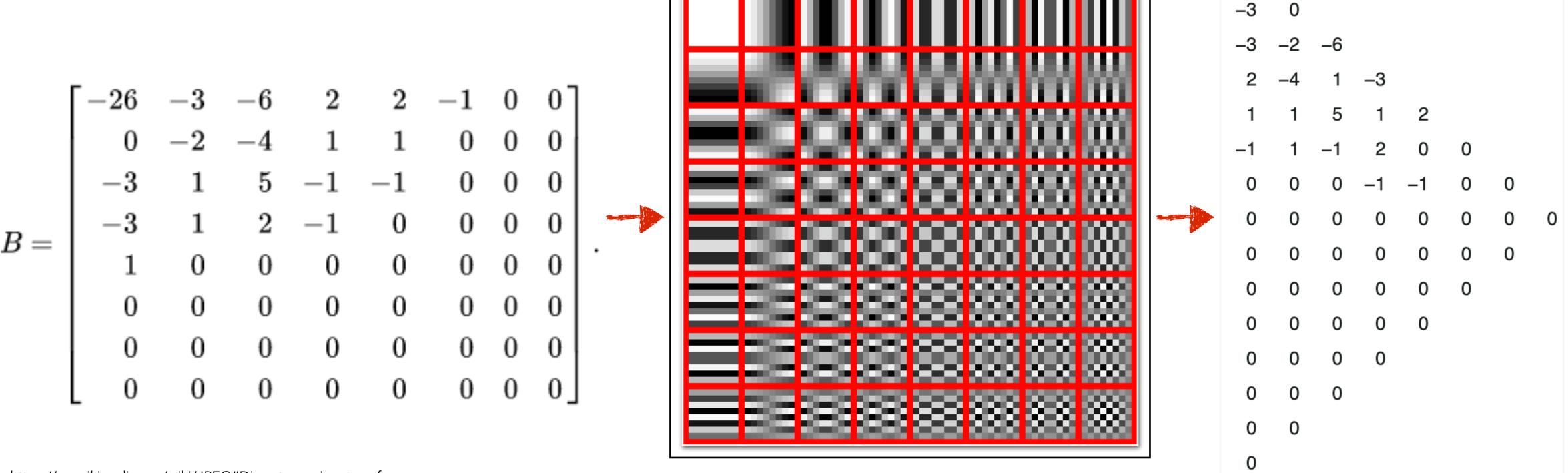


	_ •	_ •							
Qua	ation	ma		۲		• 1		11 11 • 1	
	T 16	11	10	_OW-	trequ	ency	weigh	nts are	e smaller, and high-
	12	12	14	fre	quen	cy we	eights	are la	arger. As a result,
		12	1-7	αı	ıantiz	ation	zeros	out r	nost of the high-
	14	13	16	-					
0 -	$_{O} = \begin{bmatrix} 14 & 17 \end{bmatrix}$			tr	equer	ncy co	peffici	ents -	— a good thing!
<i>₩</i> −	18	22	37	56	68	109	103	77	•
	24	35	55	64	81	104	113	92	
	49	64	78	87	103	121	120	101	
	72	92	95	98	112	100	103	99	

https://en.wikipedia.org/wiki/JPEG#Discrete_cosine_transform

Entropy Encoding (Lossless)

- 1. Reorder quantized coefficients in zig-zag order, which orders the frequencies in the ascending order.
 - Why zig-zag? Look at the pattern in the basis function.



Entropy Encoding (Lossless)

- Reorder quantized coefficients in zig-zag order, which orders the frequencies in the ascending order.
- 2. Apply run-length encoding. Compresses zeros.
- 3. Apply Huffman coding (or arithmetic coding) to compress the rest.

JPEG Algorithm Recap

A combination of lossless and lossy techniques.

1. Gamma encode RGB values.

Evenly encode perceived brightness, not luminance. Lossy due to quantization errors, but it's not unique to JPEG compression.

- 2. Convert image to Y'CbCr color space.
- 3. Chroma subsampling.



Human eyes are sensitive to luminance difference but not chrominance.

- 4. For each channel (Y', Cb, Cr)
 - A. For each 8x8 pixel block
 - 1. Compute DCT coefficients
 - 2. Quantize DCT coefficients
 - 3. Entropy encoding

Human eyes are insensitive to defects in high-frequency information, but not low-frequency information.

Compression Quality

Start seeing quantization artifact (blocking)



Quality level (Q) = 100Compression ratio: 2.7 : 1



Quality level (Q) = 25Compression ratio: 23:1



Severe blocking artifacts



Quality level (Q) = 1Compression ratio: 144:1

Image Compression Summary

JPEG is lossy (chroma subsampling and coefficient quantization).

• PNG and TIFF are lossless compression standards.

Quantization artifacts are most visible across edges and sharp corners (high frequency areas), where information loss is the greatest.

• For this reason, JPEG isn't good for compressing drawings and graphics. Co-designing graphics algorithms with compression algorithms?

The compressing (encoding) and uncompressing (decoding) system (software and hardware) is called the "codec", which are most definitely done specialized hardware. Your smartphone has one.

Video Compression

Video Compression

Compressing a sequence of continuous frames together.

What about compressing each image using JPEG?

That's what Motion-JPEG (M-JPEG) does.

- Support by most video codecs, but not widely used.
- Not well-standardized.
- Not very efficient.
- Fundamentally, compressing image individually ignores the temporal dimension in a video.

Common Video Compression Standards

MPEG-2 part 2, a.k.a., H.262.

Old but still widely used.

MPEG-4 part 10, a.k.a., H.264 or Advanced Video Coding (AVC).

• Most widely used; 91% of video industry developers (9/2019); mandated by blue-ray.

MPEG-H part 2, a.k.a., H.265 or High Efficiency Video Coding (HEVC).

Gradually taking over

Common Video Compression Standards

There are usually different "profiles" and "levels" within each MPEG standard.

• Profiles refer to different algorithmic choices; levels refer to different parameter choices

MPEG standards don't define the encoding process, but defines the format of a coded stream and the decoding process.

• You can implement your own H264 encoder/compression algorithm as long as the compressed format complies with the standard.

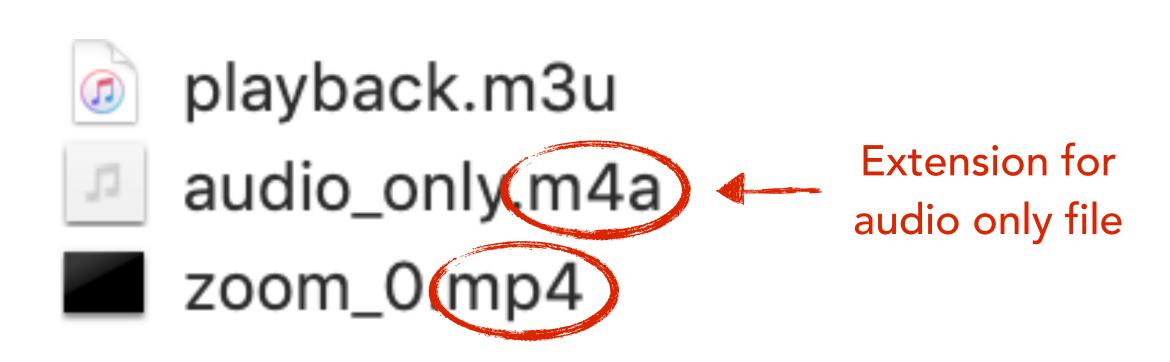
VP8/VP9, from Google, open-sourced; compete with H.264 and H.265.

Container Format

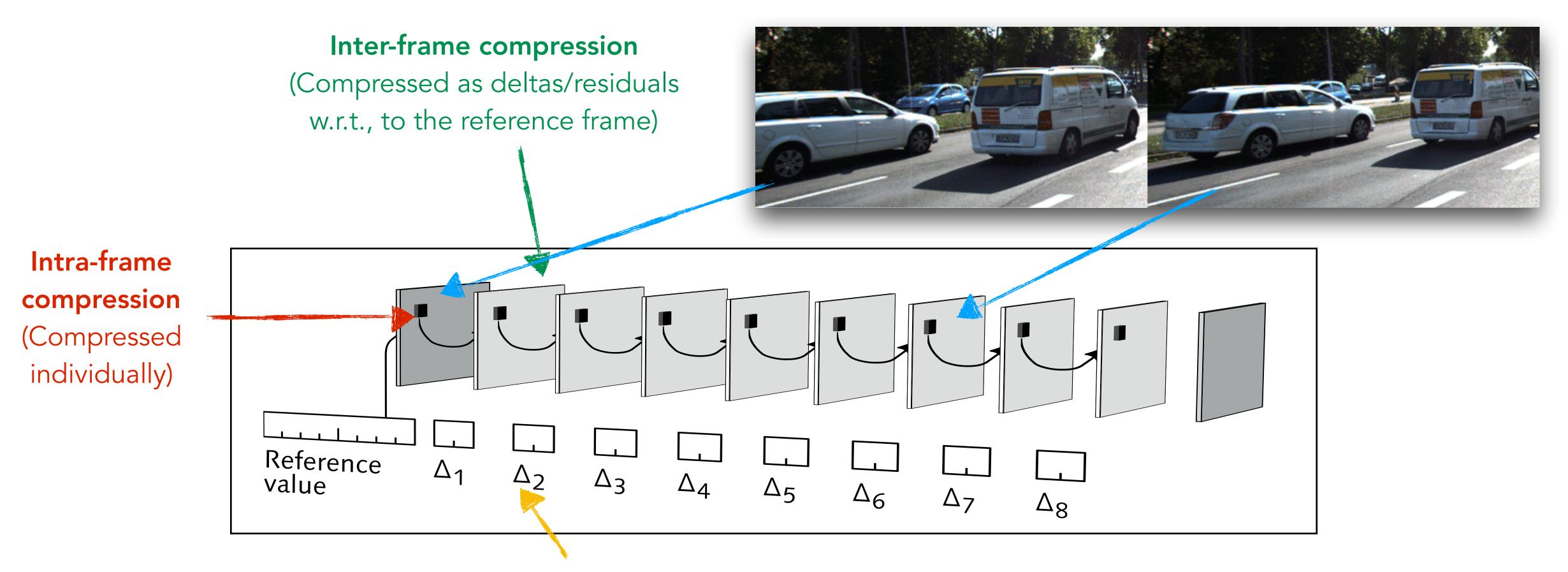
A container format bundles different data types, e.g., video, audio, subtitles.

Perhaps the most commonly used container format is .mp4, which is defined in MPEG-4 Part 14.

- Supports H.264 and H.265 video encoding and MPEG-4 Part 3 audio encoding.
- How many of you owned a "MP4 player"?! Can you imagine we use to buy devices that do nothing but playing .mp4 files?

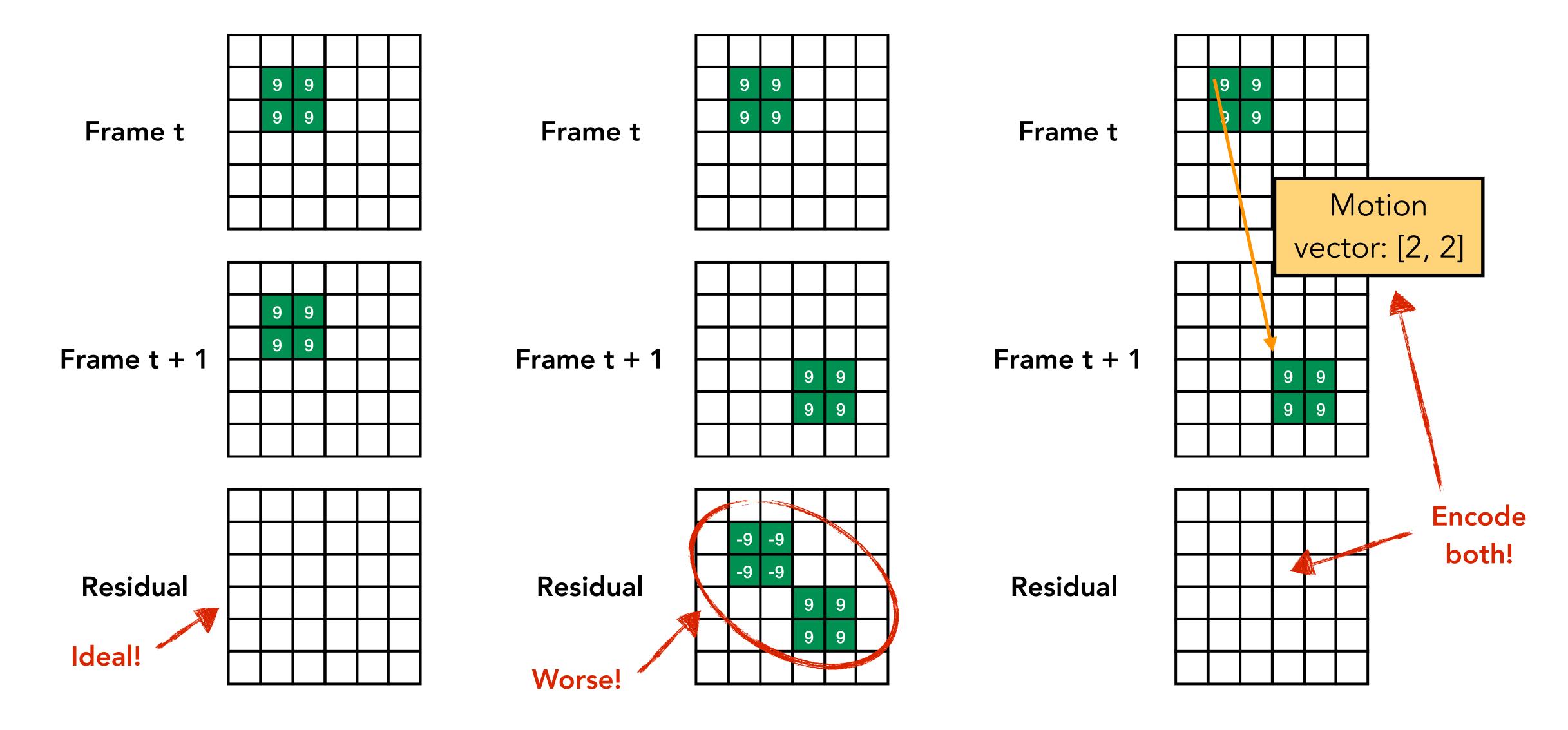


Video Compression: The Big Ideas



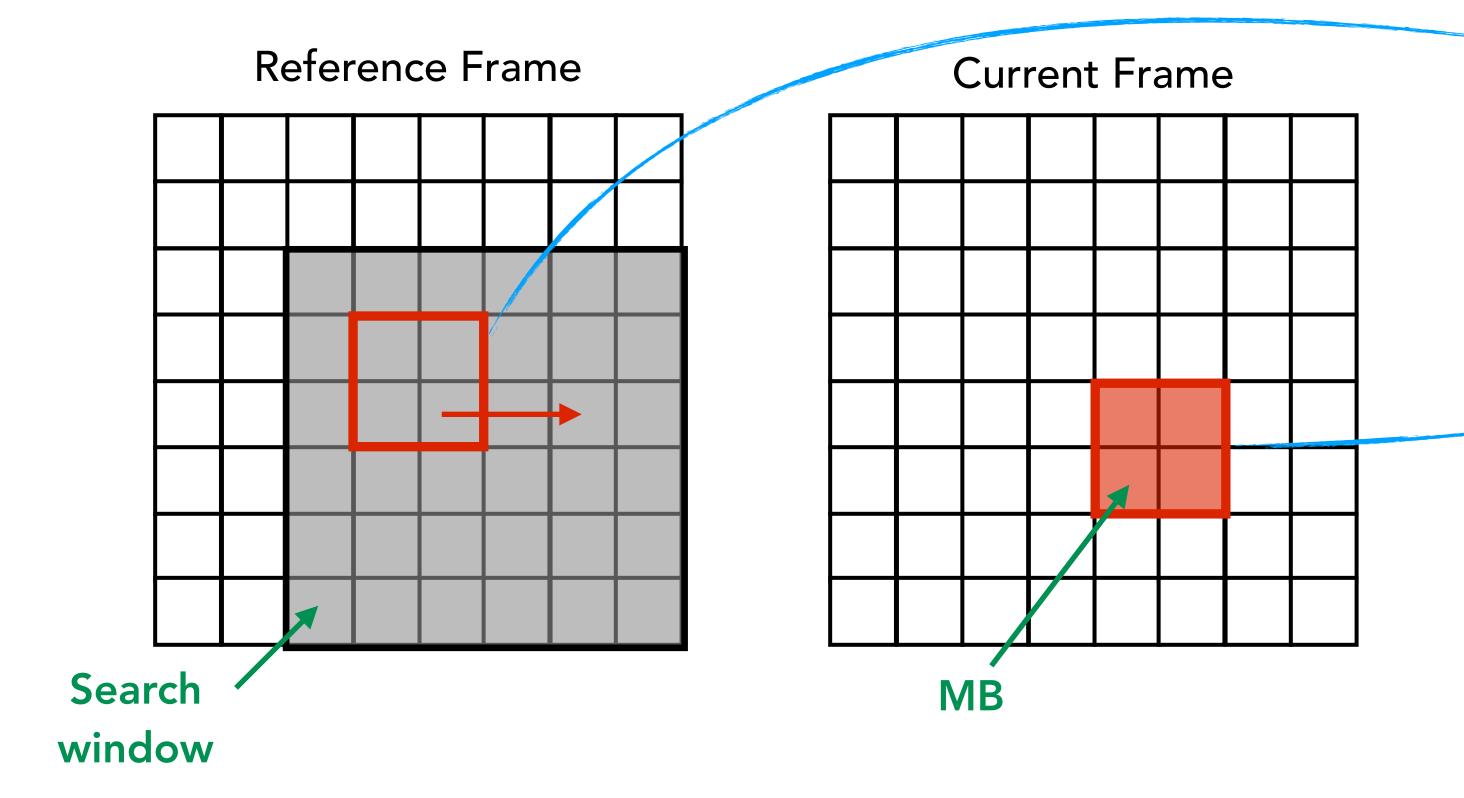
Provided that residuals can be encoded more compactly than the image itself.

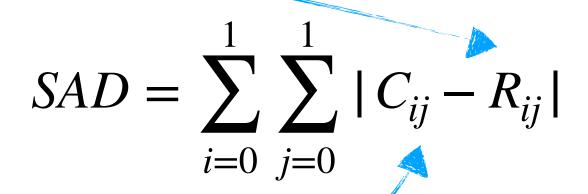
Residuals



Motion Estimation

In MPEG video compression, motion estimation is done using the **block matching** algorithm, which operates at the granularity of **macroblocks** (MB), and uses some form of **absolute difference** as the cost function.





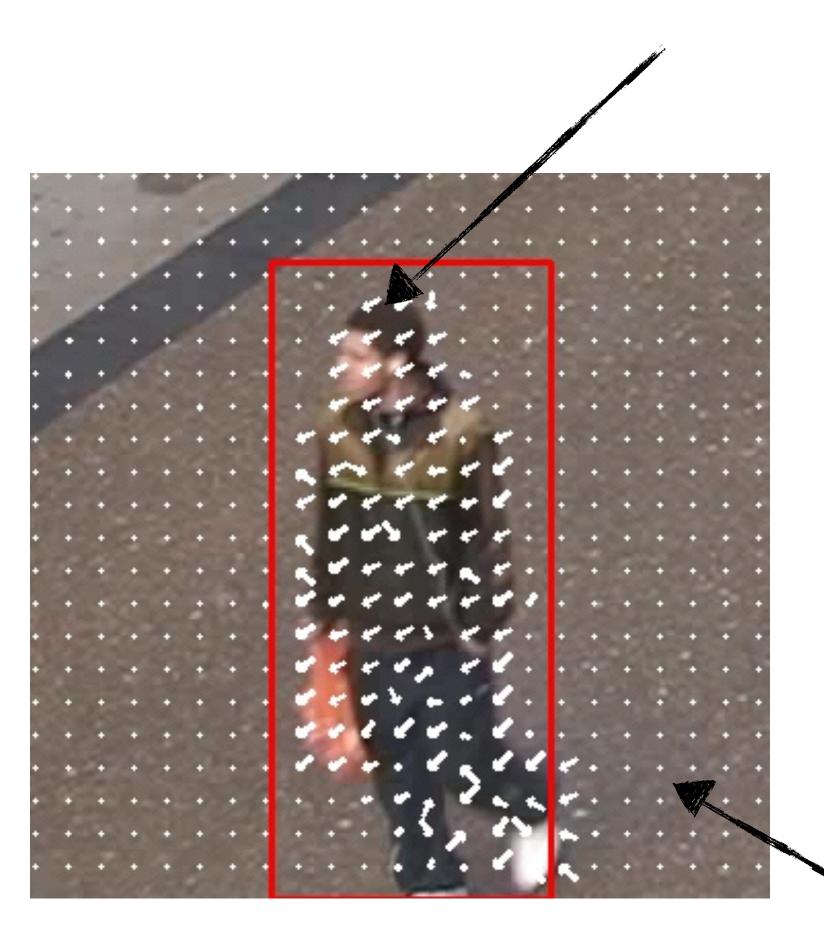
The MB in the search window that has the lowest SAD wins. Each MB is associated with a motion vector (MV).

Motion Vectors

A frame from a still camera capturing a moving person. Each arrow represents the MV of a 8x8 MB.

Motion vectors can be used to inferred the scene: object detection, tracking, etc.

Looking at just the MVs (not the image itself), can you guess where the object is moving?



Why do these pixels have zero motion?

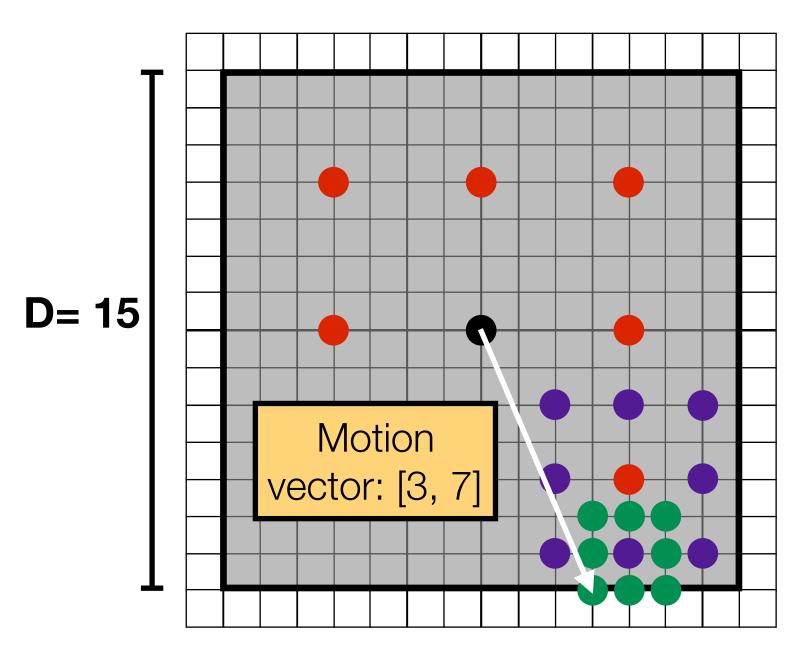
Motion Estimation Heuristics

Exhaustive search (ES) is costly. Lots of heuristics are proposed. Always a trade-off between accuracy vs. compute efficiency.

Three step search (TSS) is a classic heuristics that adaptively narrows the

search range (coarse-to-fine search).

In ES, 255 MBs are to be searched if the search window size is 15x15. With TSS, only 25 MBs are searched. 25X speed-up.

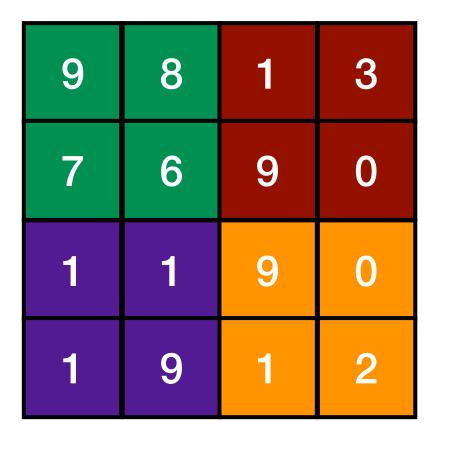


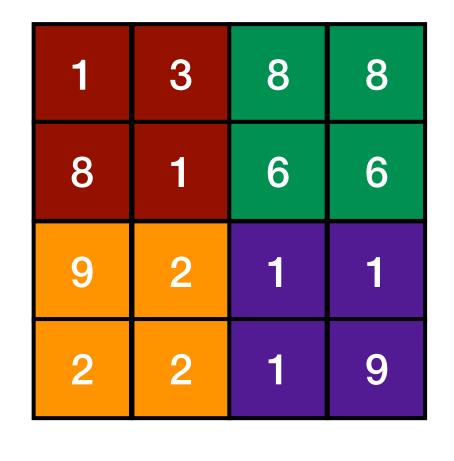
Encoding Residuals

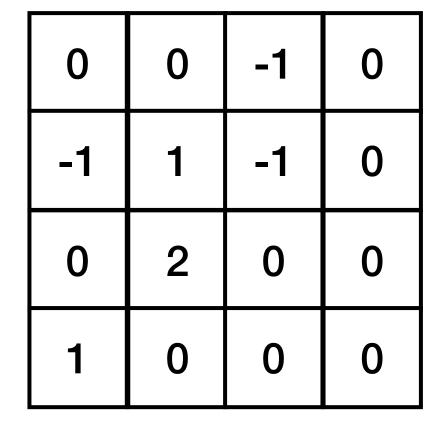
Motion estimation will most likely not find a perfect match, i.e., the cost function will not be 0. Therefore, the residuals must still be encoded.

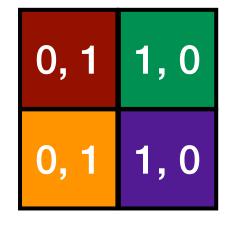
The residuals form a residual image, which is encoded in a JPEG-like fashion.

• The quantization matrices are usually different from the ones used in encoding images.









Frame t

Frame t+1

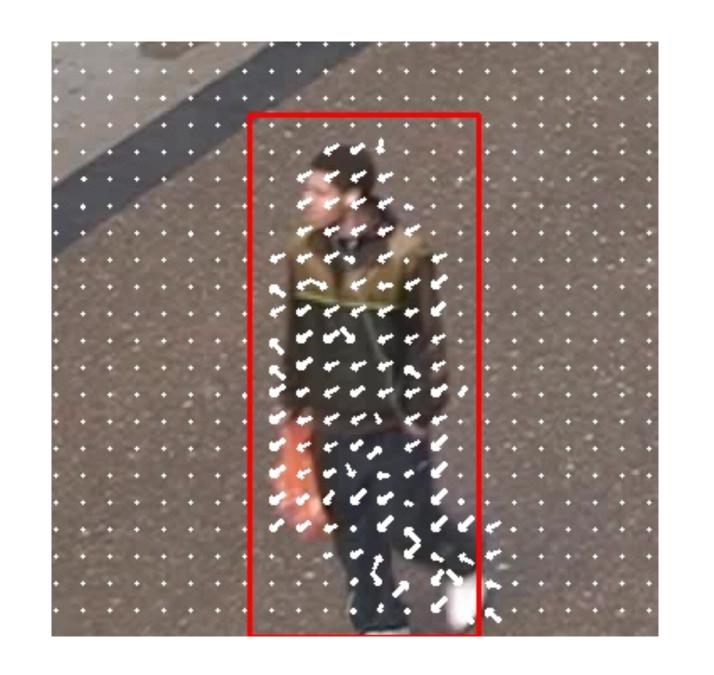
Residual

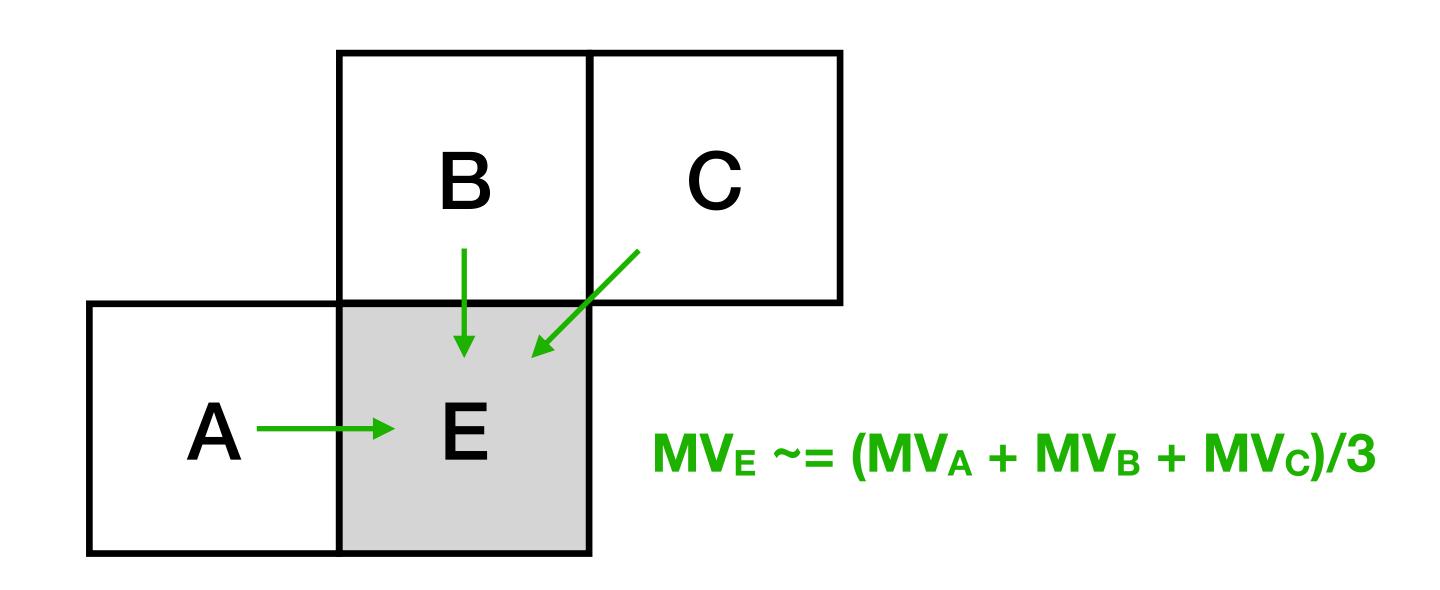
Motion vectors

Encoding Motion Vectors

Motion vectors of nearby MBs are often correlated (rigid objects).

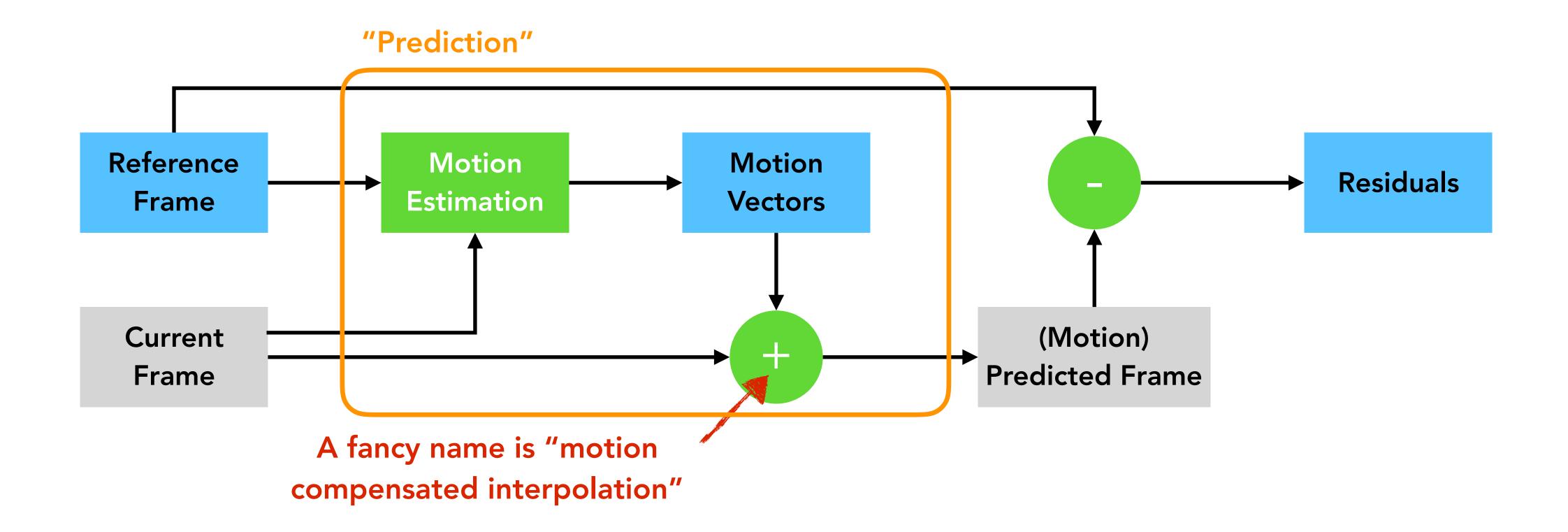
Idea: use nearby MVs to estimate/predict the current MV, compute the residuals, and encode the MV residuals.





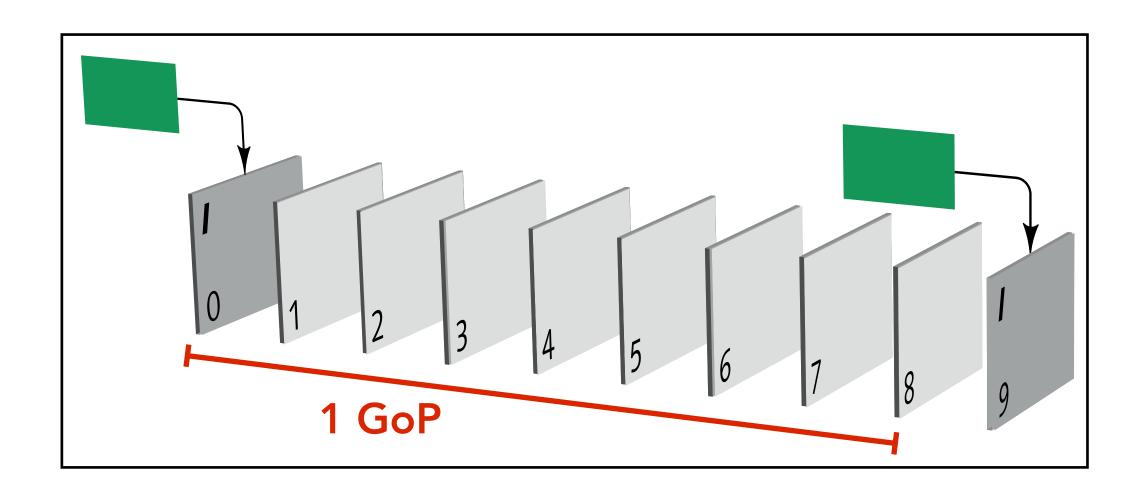
Inter-Frame Encoding Flow

Data that needs to be encoded



Reference Frame Choice

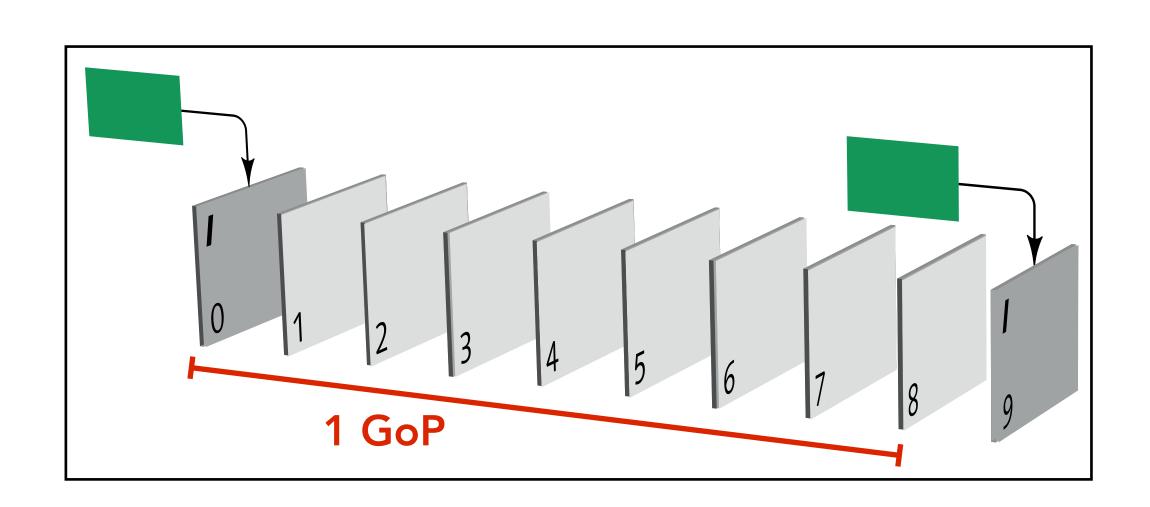
In MPEG, a video is partitioned into successive groups of pictures (GoPs), each of which contains a sequence of successive frames.

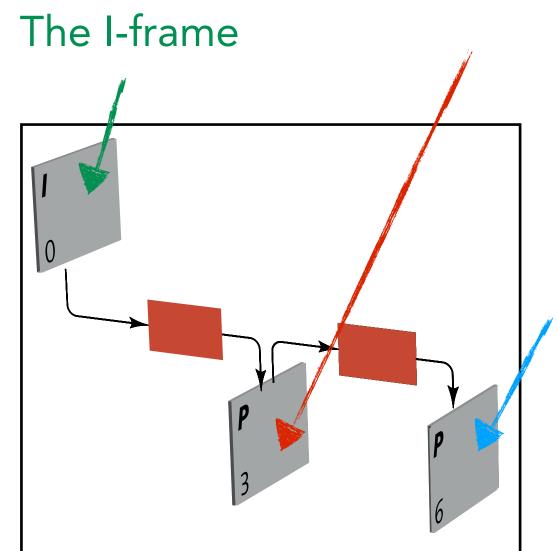


Reference Frame Choice

Each GoP has three (mutually exclusive) kinds of frames:

- An **I-frame** is intra-compressed and is a reference frame.
- A **P-frame** is a frame that is predicted from a reference frame, and itself can be used as a reference frame.





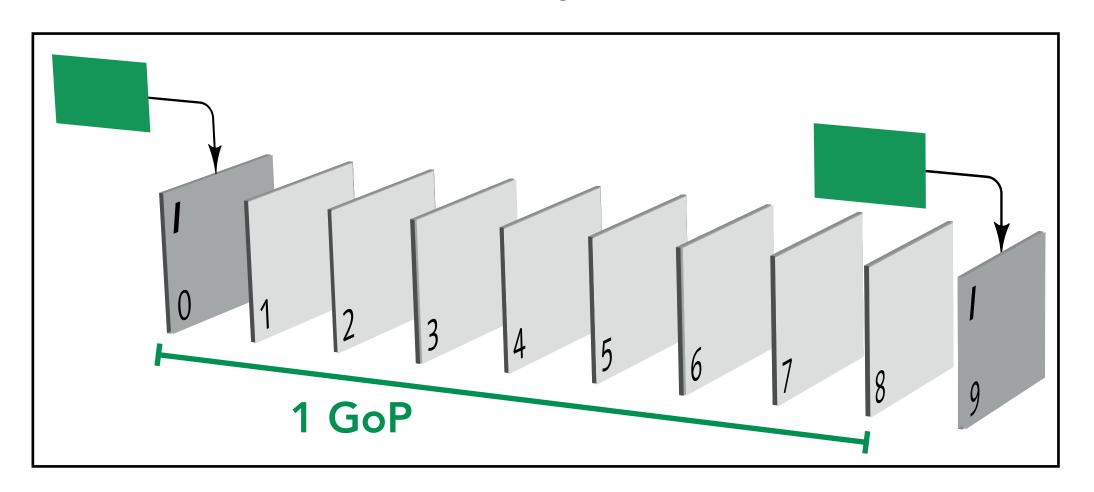
A P-frame, predicted from the I-frame, and becomes a reference frame itself

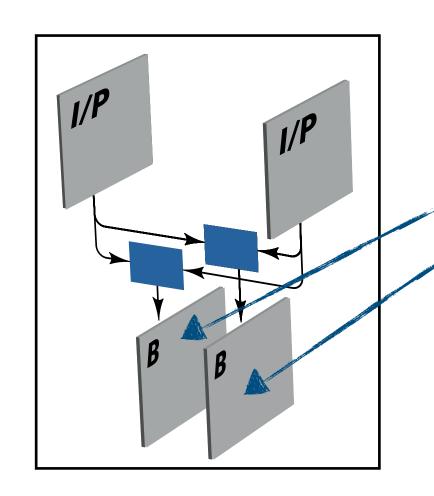
Another P-frame, predicted from P-frame 3, and is not used as a reference frame.

Reference Frame Choice

Each GoP has three (mutually exclusive) kinds of frames:

- An I-frame is intra-compressed and is a reference frame.
- A **P-frame** is a frame that is predicted from a reference frame, and itself can be used as a reference frame.
- A **B-frame** is a frame that is predicted from two reference frames (because the residuals are smaller that way), but itself can't be used as a reference frame.





Two B-frames. Usually a B-frame is predicted from the average of the two reference frames (but could also take a weighted average).

Reordering Frames in a GoP

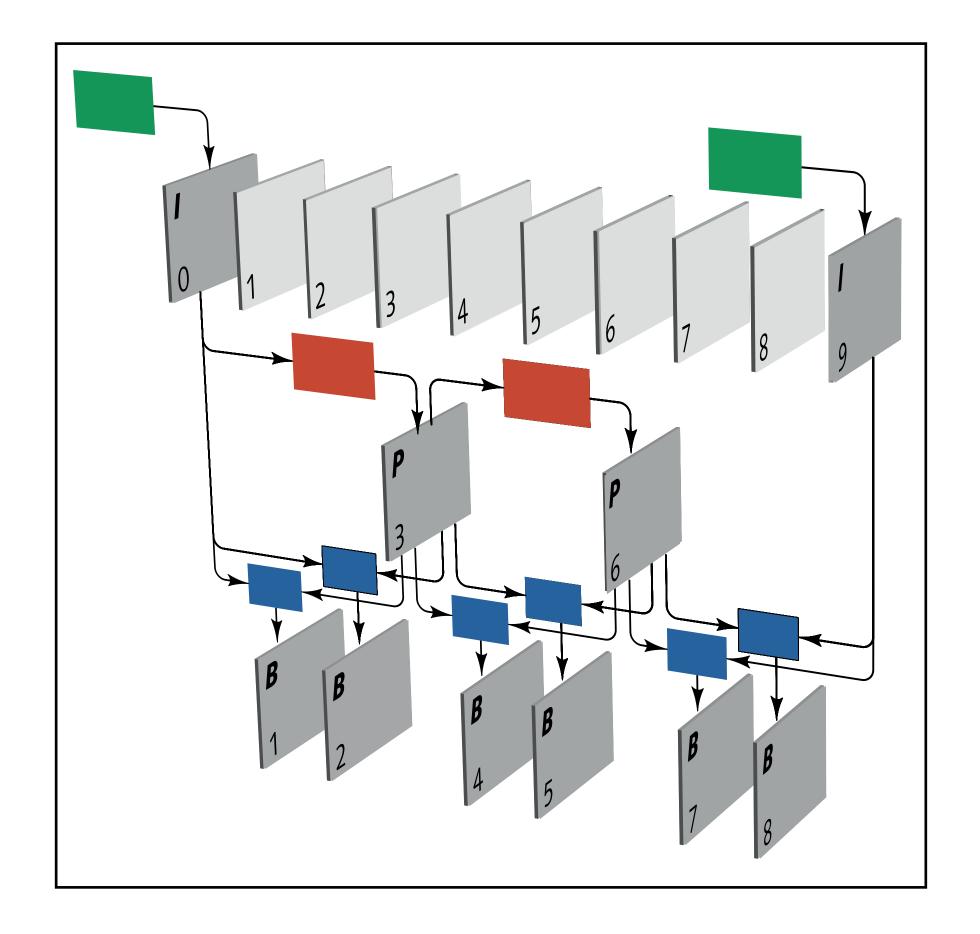
I₀B₁B₂P₃B₄B₅P₆B₇B₈

Bad for live streaming/teleconferencing!!!

- 1. When encoding, need to buffer B_1 and B_2 before P_3 arrives (e.g., from camera).
- 2. When decoding, need to buffer B_1 and B_2 before P_3 arrives (e.g., from Internet).

$$I_0P_3B_1B_2P_6B_4B_5(I_9)B_7B_8$$

Solution: reorder frames during transmission!

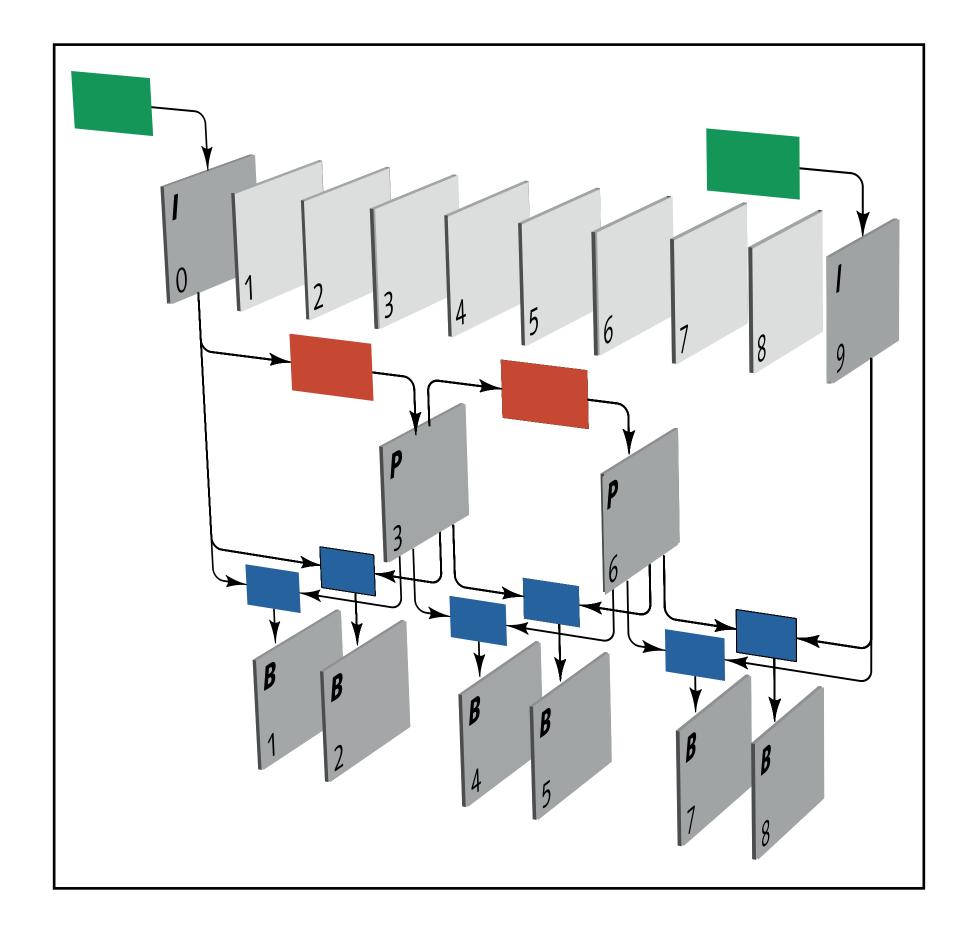


Incorporating the Decoder in the Encoder

 B_1 is generated from I_0 and P_3 , but both I_0 and P_3 themselves are compressed in a lossy fashion. What the decoder sees are the lossy versions I_0 and P_3 .

This means we should really predict/motion-estimate B_1 from I_0 and P_3 .

How? Incorporate the decoder in the encoder so that we use I_0 ' and P_3 ' exactly as how decoder will see them.

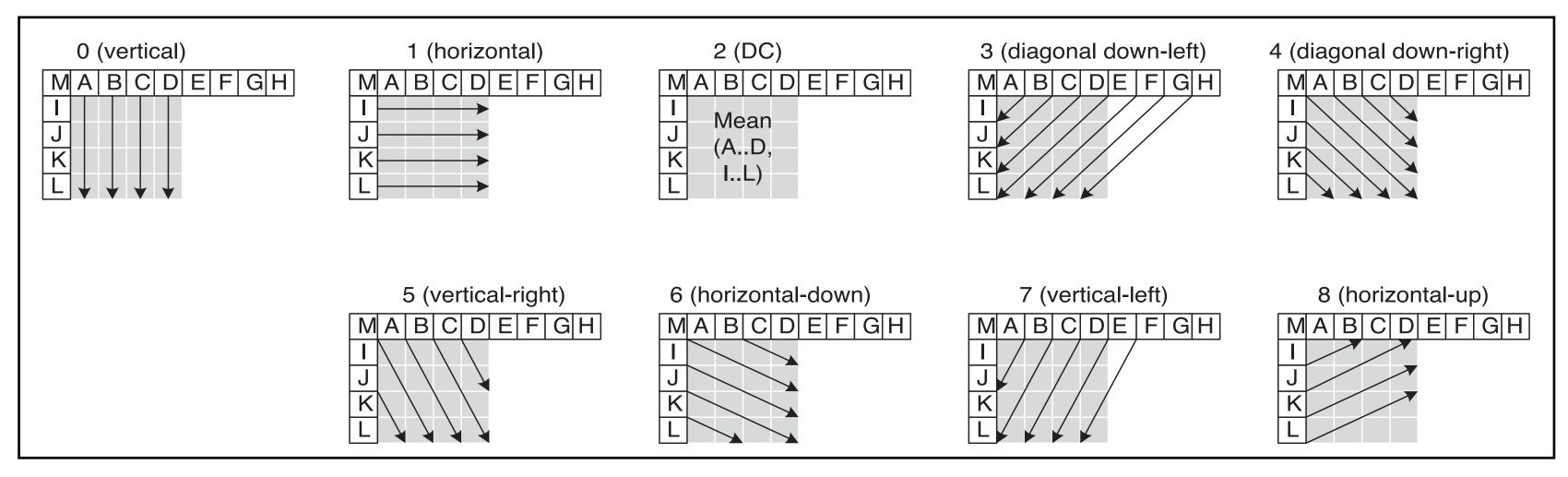


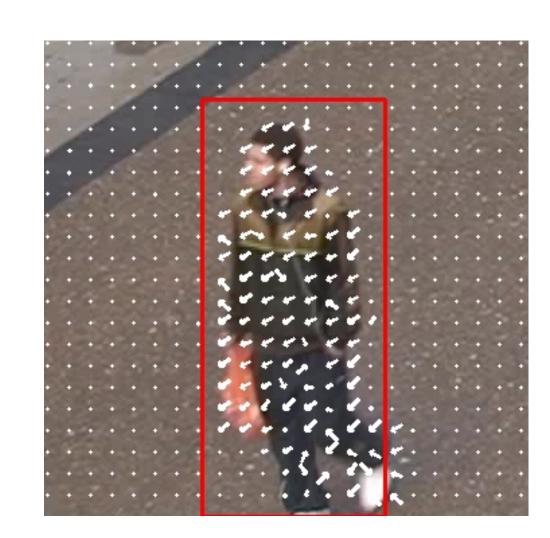
Intra-Frame Encoding (Spatial Prediction)

Prior to H.264, I-frames are simply compressed using JPEG-like techniques.

Starting from H.264, I-frames macroblocks are spatially predicted, exploring the observations that spatially adjacent pixels are strongly correlated.

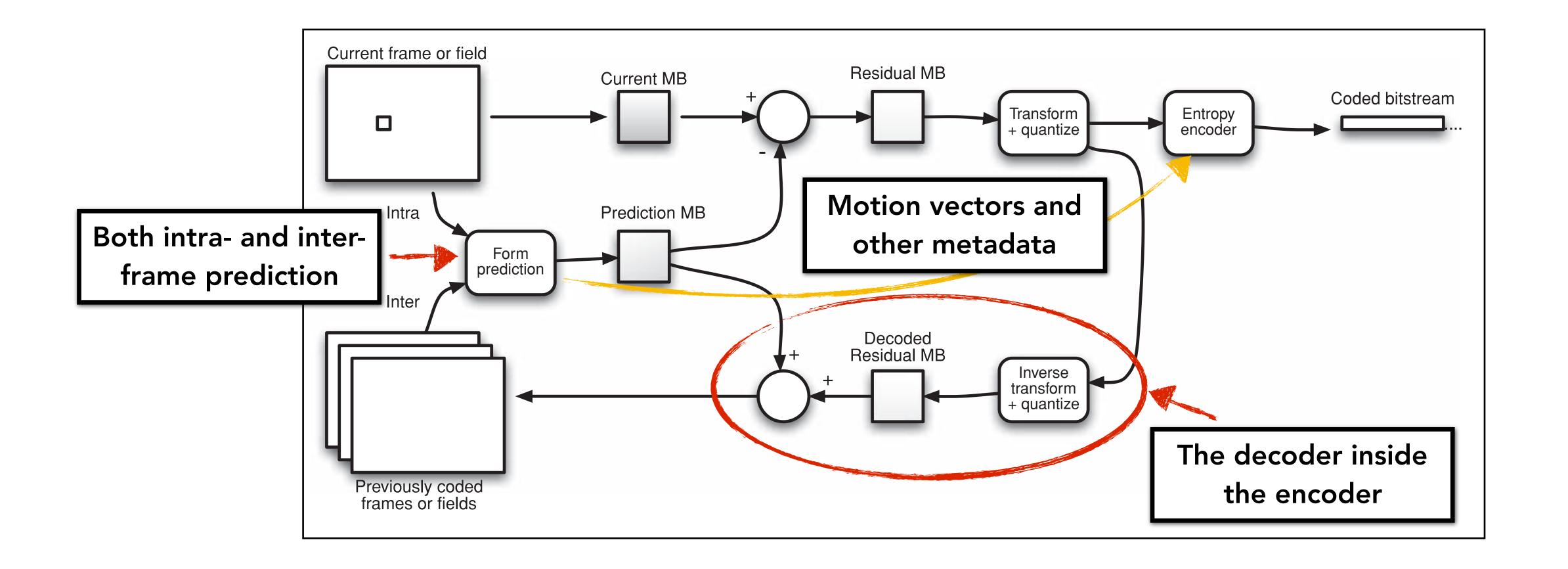
• Again residuals are then encoded (same as in predicted frames).





Modes for spatially-predicting a 4x4 tile

H.264 Encoding Architecture



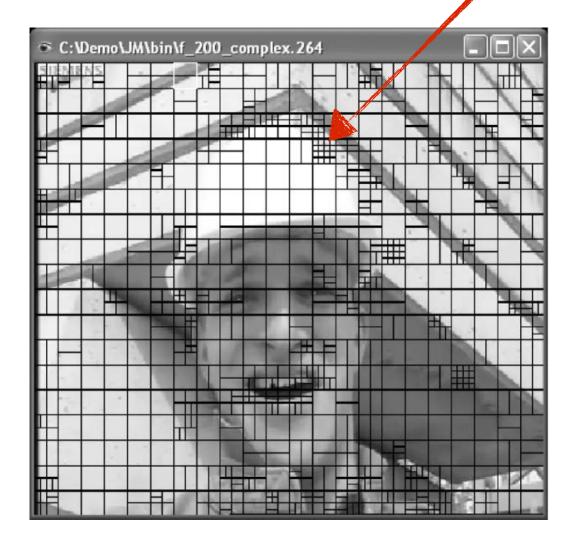
Other Details/Optimizations in Video Codecs

Luma and chroma components are dealt with separately, both in intracompressed and inter-compressed schemes.

Allow sub-pixel motion (smaller residuals).

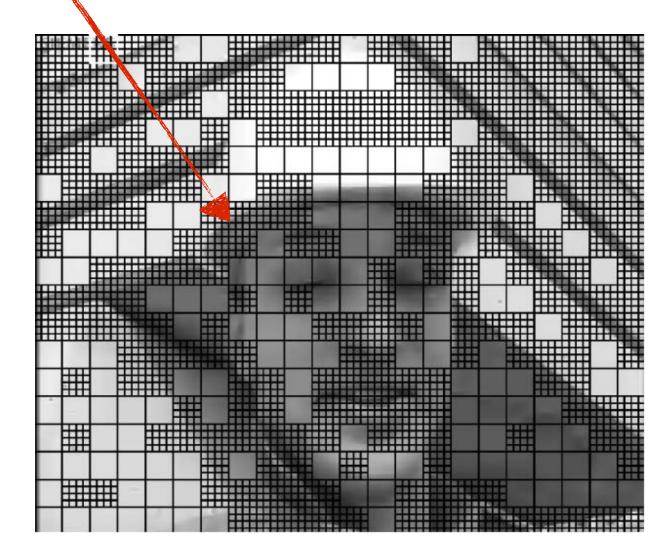
Allow different MB sizes.

 Different MB sizes in a P-frame.



"Busy" areas in a frames are better predicted using smaller MBs.

Different MB sizes in an I-frame.



Video Compression Recap

Spatial redundancy

• Pixels in a small neighborhood have strong correlations. I-frames are spatially predicted.

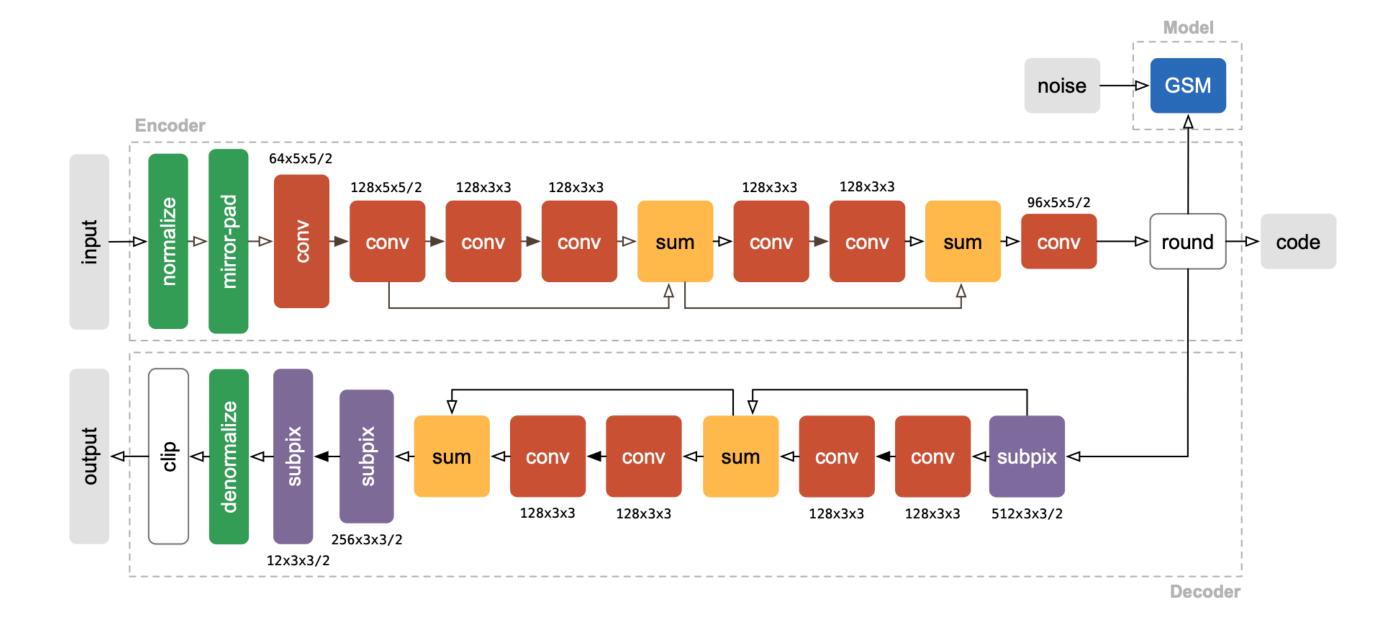
Temporal redundancy

• Consecutive frames have strong correlation. P/B-frames are temporally predicted.

Encoding time is dominated by motion estimation. Encoding and decoding are widely support by fixed-function hardware (ASICs).

Video-on-demand Video Transcoding (VOD). Speed is key. 1080p Universal Live or VOD Universal Play Original transcode transcode format 144p Measure INPUT PROTOCOLS INPUT CODECS OUTPUT CODECS **OUTPUT PROTOCOLS** H.265/HEVC, H.264/AVC, VP9. H.265/HEVC Adobe RTMP, Apple HLS, Adobe HDS, VP8 MPEG4 Part 2, MPEG2 H.264/AVC, RTSP/RTP, MPEG-TS, MPEG-DASH, Microsoft Smooth H.263 (v2), VP9 ICY (SHOUTcast/Icecast) Streaming, Adobe RTMP, RTSP/RTP, MPEG-TS MP3, AAC, AAC-LC, AAC, AAC-LC, 1080p HE-AAC+ v1 & v2, HE-AAC+ v1 & v2, Opus, G.711 MPEG1 Part 1/2, Speex, **Popular** G.711, Opus, Vorbis Play transcode 144p Spend more time optimizing 720p/2Mbps INPUT STREAM popular videos. TRANSCODING 480p/1Mbps ----Transrate A huge design space (different resolutions, codecs, compute capabilities, 240p/400Kbps bandwidths) Must make trade-offs!

Al For Compression and Compressing for Al



Immediate goal: remove empirical decisions and heuristics from compression. Learn the best compressed representation automatically.



A better question: how to design the compression algorithm (network) for computer vision algorithms? After all, many videos will be consumed by not humans, but computers.

- 1. Vision algorithms might require a different quality measure from human
- 2. Different vision algorithms might require a different compression scheme