

Learning for Image Compression with Deep Neural Networks



Ren Yang



Dr. Radu Timofte

Computer Vision Laboratory, D-ITET, ETH Zurich

ETH zürich

**Large amount of
high-resolution images/videos**



**Limited
bandwidth**



**Terminal
devices**

Limited storage



$$7296 \times 5472 = 39,923,712 \text{ pixels}$$

Uncompressed image: $39,923,712 \times 3 = 120 \text{ MB}$

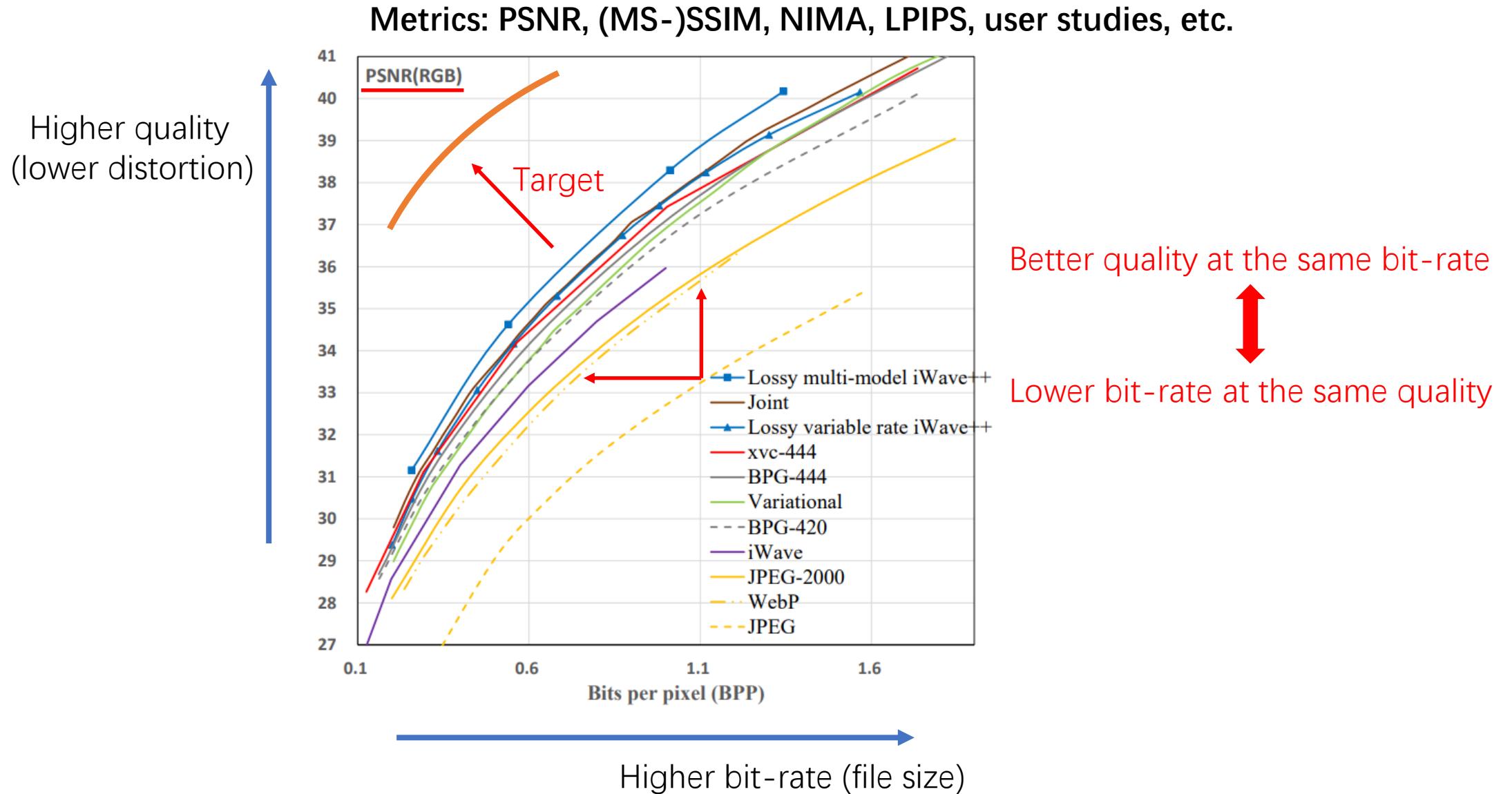
Uncompressed video (60 fps): $120 \text{ MB} \times 60 = 7.2 \text{ GBps}$ (18s needs 128 GB)

Lossless compression (.png): 44 MB

Lossy compression (.jpg): 9 MB

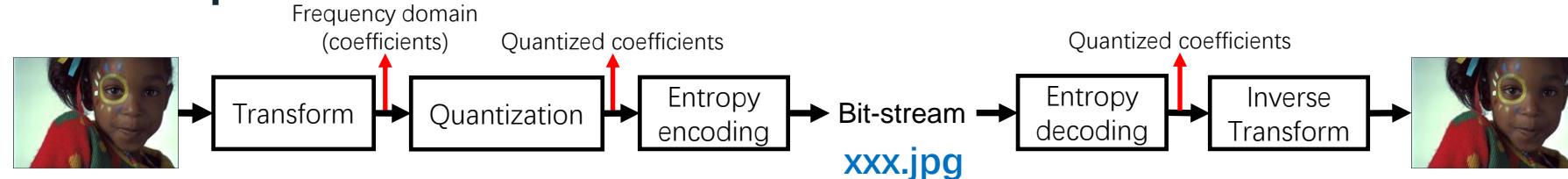
Image/video compression plays an important role in multimedia streaming, online conference, data storage, etc.

Rate-distortion trade-off

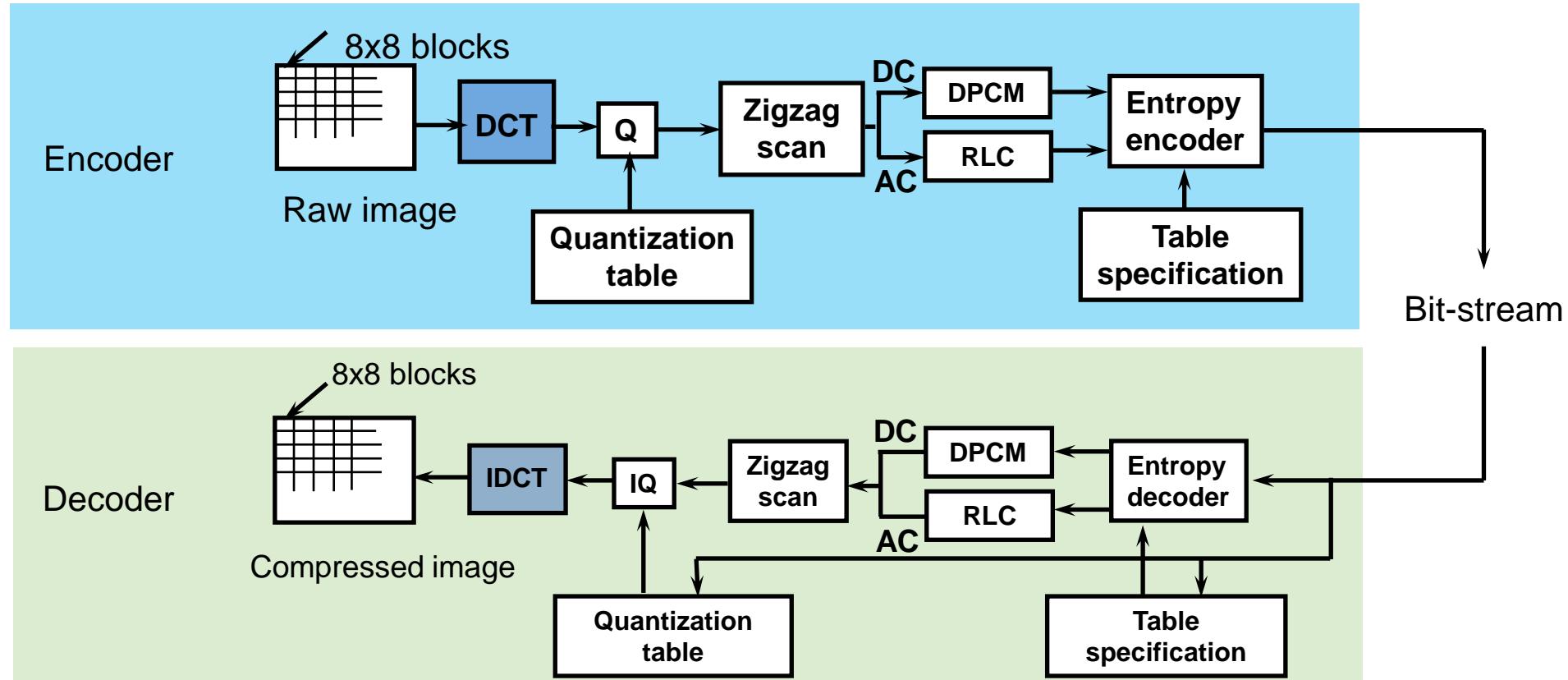


Traditional Image Compression

- Classical Architecture:



- Standards: JPEG (DCT + Huffman), JPEG2000 (DWT + Arithmetic coding), BPG (HEVC), ...
- Example: JPEG compression framework



Entropy coding

Entropy:

$$H(X) = E[I(X)] = E[-\log(P(X))]$$

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

Cross entropy:

$$H(p, q) = - \sum_{x \in \mathcal{X}} \frac{p(x)}{\text{real}} \log \frac{q(x)}{\text{estimated}} \quad (\text{Eq.1})$$

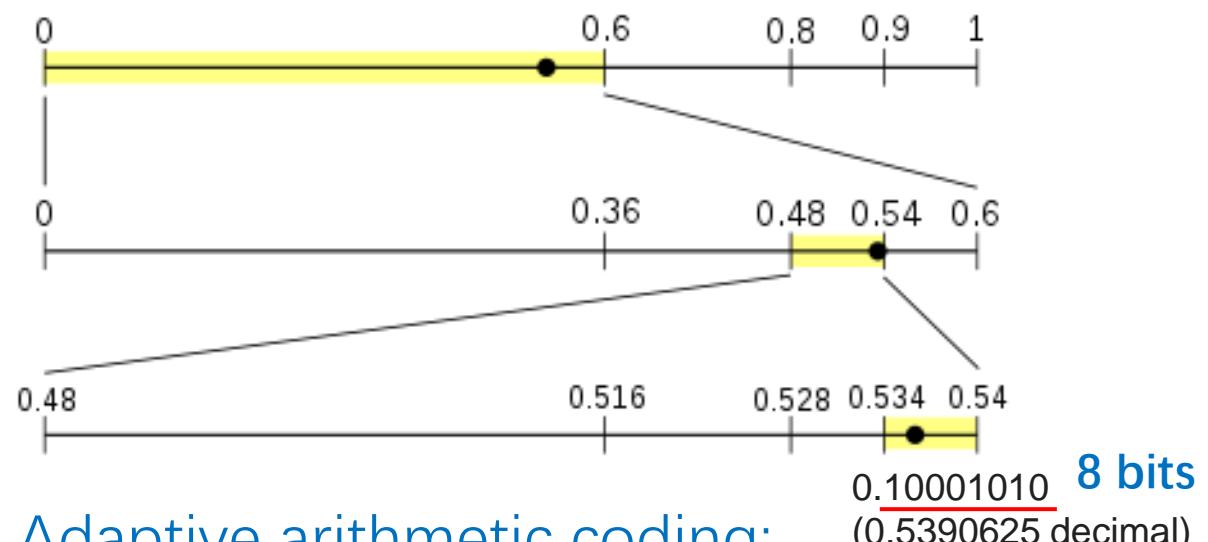
(Adaptive) arithmetic coding is theoretically able to losslessly compress data at

- bit-rate \cong cross entropy (with little overhead)

Arithmetic coding:

- 60% chance of symbol NEUTRAL
- 20% chance of symbol POSITIVE
- 10% chance of symbol NEGATIVE
- 10% chance of symbol END-OF-DATA.

NEUTRAL NEGATIVE END-OF-DATA message

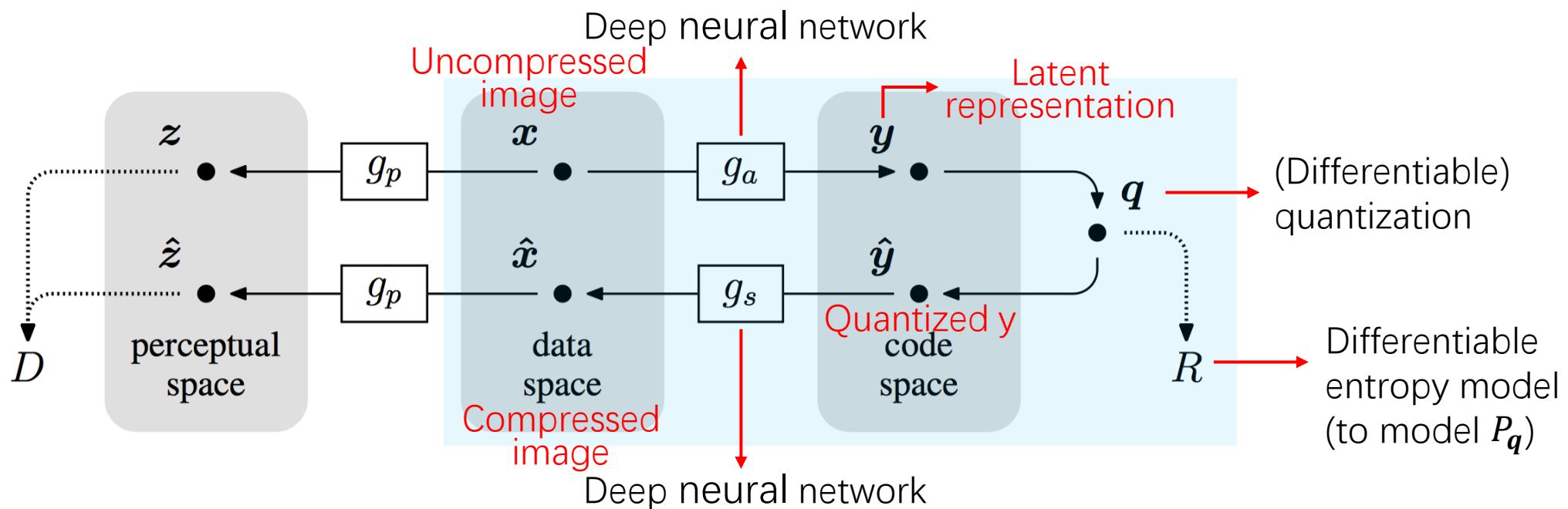


Adaptive arithmetic coding:

Changing the frequency (or probability) tables while processing the data.

Learned Image Compression

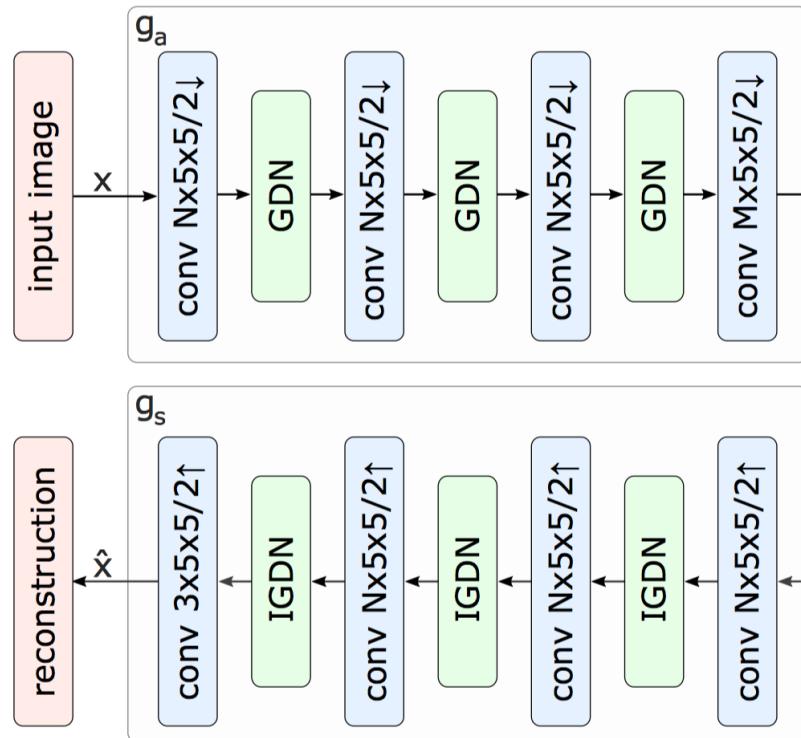
- Basic architecture [1]: **End-to-end trainable**



$$L[g_a, g_s, P_q] = \frac{-\mathbb{E}[\log_2 P_q] + \lambda \mathbb{E}[d(x, \hat{x})]}{R}$$

Learned Image Compression

- CNN transformer + **factorized** entropy model [1]



$$\tilde{\mathbf{y}} = \mathbf{y} + \Delta\mathbf{y} \sim \mathcal{U}(0, 1)$$

$$\hat{\mathbf{y}} = \text{round}(\mathbf{y})$$

Inference: quantization (not differentiable)

$\tilde{\mathbf{y}}$ or $\hat{\mathbf{y}}$

$$f_k(\underline{\mathbf{x}}) = g_k(\mathbf{H}^{(k)}\underline{\mathbf{x}} + \mathbf{b}^{(k)}) \quad 1 \leq k < K$$

$$f_K(\underline{\mathbf{x}}) = \text{sigmoid}(\mathbf{H}^{(K)}\underline{\mathbf{x}} + \mathbf{b}^{(K)})$$

$$g_k(\underline{\mathbf{x}}) = \underline{\mathbf{x}} + \mathbf{a}^{(k)} \odot \tanh(\underline{\mathbf{x}}) \quad \mathbf{H}^{(k)} = \text{softplus}(\hat{\mathbf{H}}^{(k)})$$

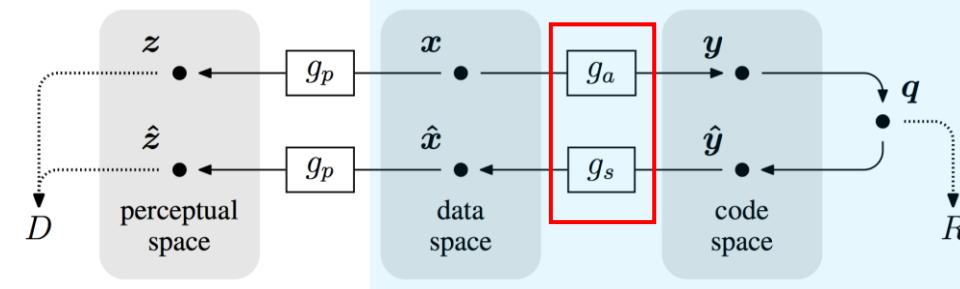
$$g'_k(\underline{\mathbf{x}}) = 1 + \mathbf{a}^{(k)} \odot \tanh'(\underline{\mathbf{x}}) \quad \mathbf{a}^{(k)} = \tanh(\hat{\mathbf{a}}^{(k)})$$

$$R = \mathbb{E}_{\underline{\mathbf{x}} \sim p_{\underline{\mathbf{x}}}} [-\log_2 p_{\hat{\mathbf{y}}} (Q(g_a(\underline{\mathbf{x}}; \phi_g)))] \quad \text{estimated bit-rate}$$

$$L(\theta, \phi) = \mathbb{E}_{\underline{\mathbf{x}}, \Delta\mathbf{y}} \left[- \sum_i \log_2 p_{\tilde{\mathbf{y}}_i} (g_a(\underline{\mathbf{x}}; \phi) + \Delta\mathbf{y}; \psi^{(i)}) + \lambda d(g_p(g_s(g_a(\underline{\mathbf{x}}; \phi) + \Delta\mathbf{y}; \theta)), g_p(\underline{\mathbf{x}})) \right]$$

bit-rate trade-off distortion

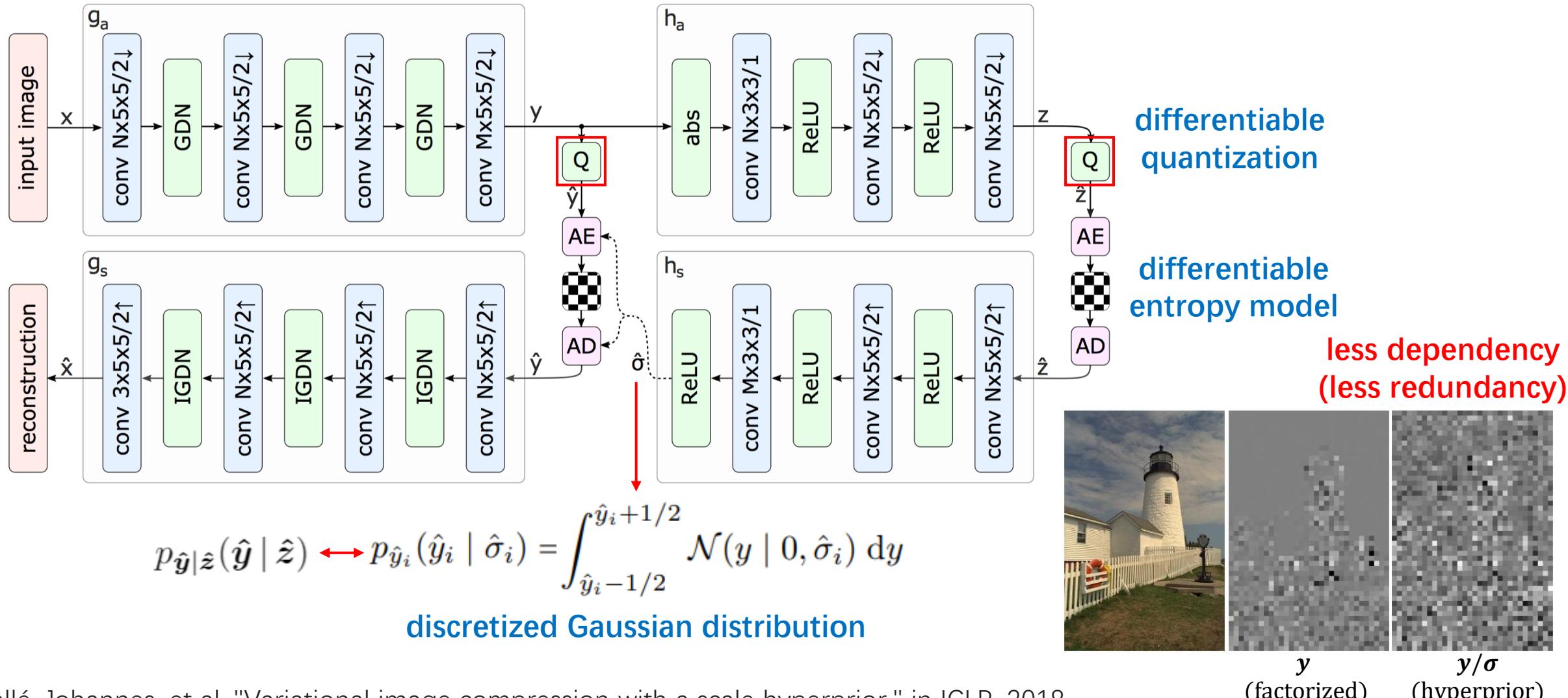
Optimized in an end-to-end manner



Training: differentiable quantization

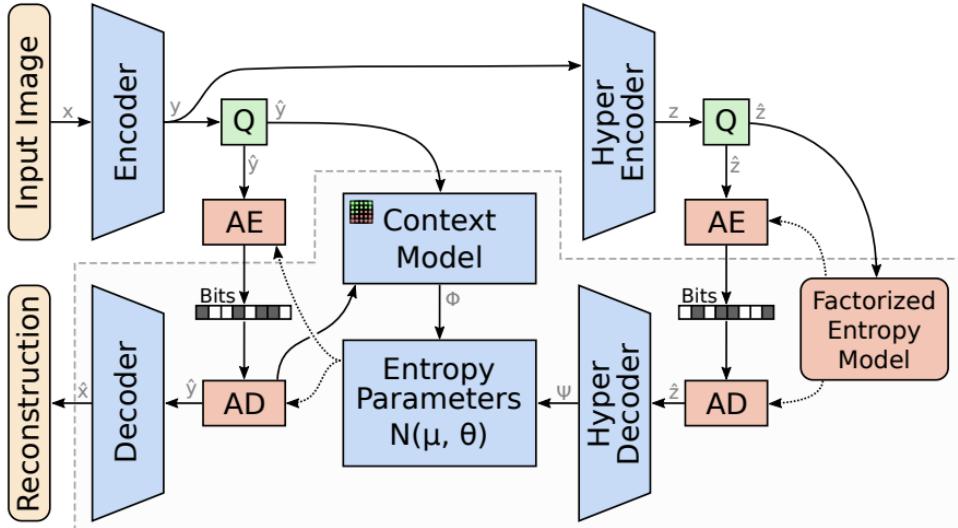
Learned Image Compression

- CNN transformer + **hyperprior** entropy model [2]



Learned Image Compression

- CNN transformer + **autoregressive** entropy model [3]

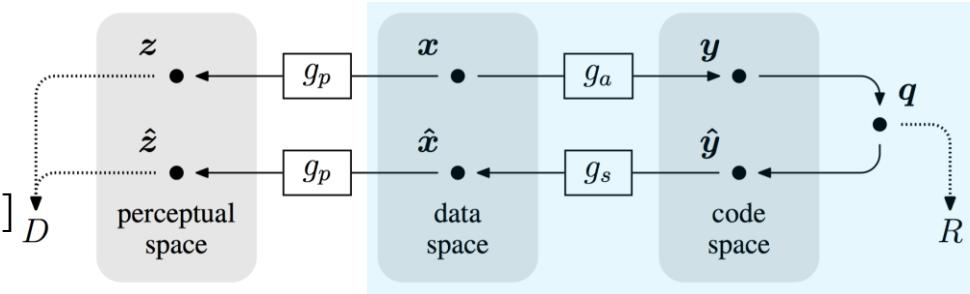


Due to the chain rule: $p(\mathbf{y}) = p(y_1) \cdot p(y_2|y_1) \cdot p(y_3|y_2, y_1) \dots p(y_N|y_{N-1})$

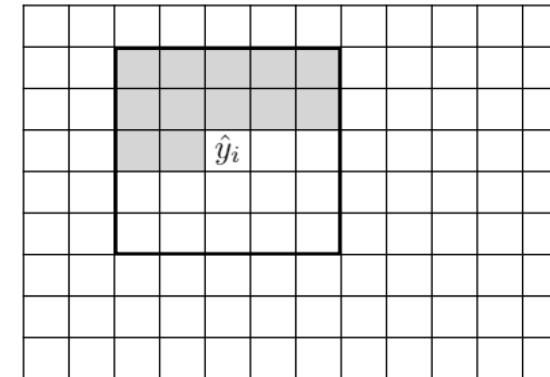
$$p_{\hat{\mathbf{y}}|\hat{\mathbf{z}}}(\hat{\mathbf{y}} | \hat{\mathbf{z}}) = \prod_{i=1}^N p_{\hat{y}_i|\hat{y}_{<i}, \hat{\mathbf{z}}}(\hat{y}_i | \hat{y}_{<i}, \hat{\mathbf{z}})$$

$$p_{\hat{\mathbf{y}}}(\hat{\mathbf{y}} | \hat{\mathbf{z}}, \theta_{hd}, \theta_{cm}, \theta_{ep}) = \prod_i \left(\mathcal{N}(\mu_i, \sigma_i^2) * \mathcal{U}\left(-\frac{1}{2}, \frac{1}{2}\right) \right)(\hat{y}_i)$$

with $\mu_i, \sigma_i = g_{ep}(\psi, \phi_i; \theta_{ep})$, $\psi = g_h(\hat{\mathbf{z}}; \theta_{hd})$, and $\phi_i = g_{cm}(\hat{y}_{<i}; \theta_{cm})$



Mask CNN [4]



$$\begin{matrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$$

first layer

$$\begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \end{matrix}$$

other layers

Algorithm 1 Constructing 3D Masks

```

1: central_idx  $\leftarrow \lceil (f_w \cdot f_H \cdot f_D)/2 \rceil$ 
2: current_idx  $\leftarrow 1$ 
3: mask  $\leftarrow f_w \times f_H \times f_D$ -dimensional matrix of zeros
4: for  $d \in \{1, \dots, f_D\}$  do
5:   for  $h \in \{1, \dots, f_H\}$  do
6:     for  $w \in \{1, \dots, f_w\}$  do
7:       if current_idx  $<$  central_idx then
8:         mask( $w, h, d$ )  $= 1$ 
9:       else
10:        mask( $w, h, d$ )  $= 0$ 
11:       current_idx  $\leftarrow$  current_idx + 1

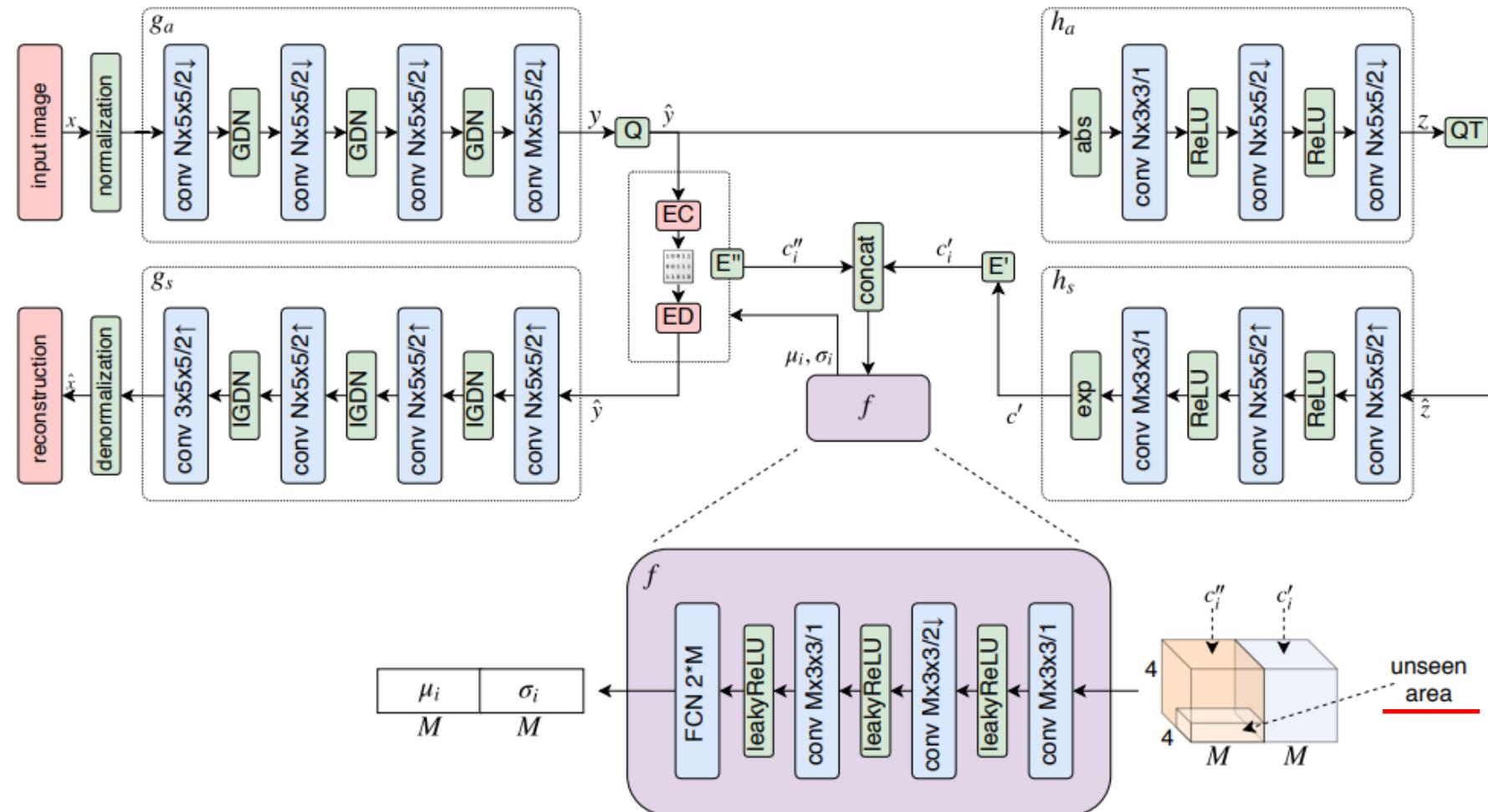
```

[3] Minnen, David, et al. "Joint autoregressive and hierarchical priors for learned image compression." in NeurIPS. 2018.

[4] Mentzer, Fabian, et al. "Conditional Probability Models for Deep Image Compression", in CVPR, 2018.

Learned Image Compression

- CNN transformer + **autoregressive** entropy model [5]



Learned Image Compression

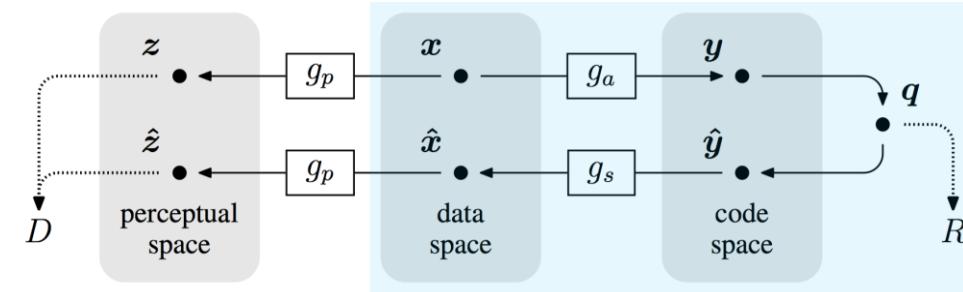
- Another differentiable quantization method [4]
given centers $\mathcal{C} = \{c_1, \dots, c_L\}$

$$\hat{z}_i = Q(z_i) := \arg \min_j \|z_i - c_j\|$$

Inference

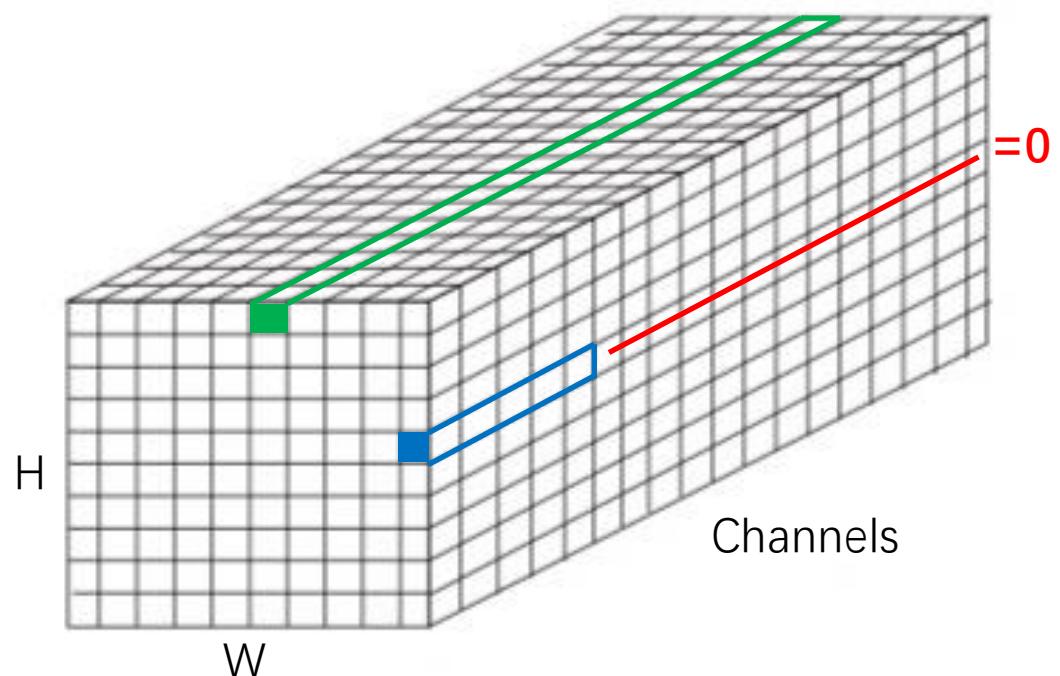
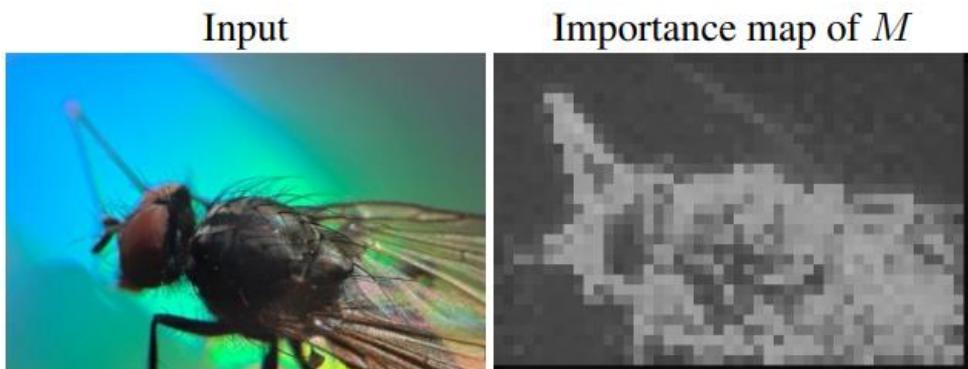
$$\tilde{z}_i = \sum_{j=1}^L \frac{\exp(-\sigma \|z_i - c_j\|)}{\sum_{l=1}^L \exp(-\sigma \|z_i - c_l\|)} c_j$$

Training: differentiable



$$\bar{z}_i = \text{tf.stopgradient}(\hat{z}_i - \tilde{z}_i) + \tilde{z}_i$$

- Importance map [4]



Learned Image Compression

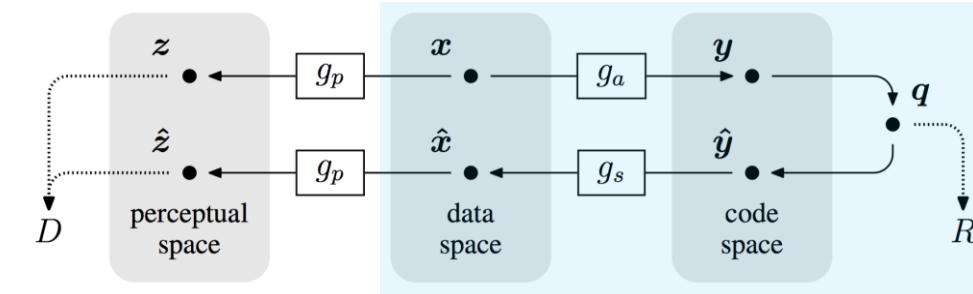
- Another differentiable quantization method [4]
given centers $\mathcal{C} = \{c_1, \dots, c_L\}$

$$\hat{z}_i = Q(z_i) := \arg \min_j \|z_i - c_j\|$$

Inference

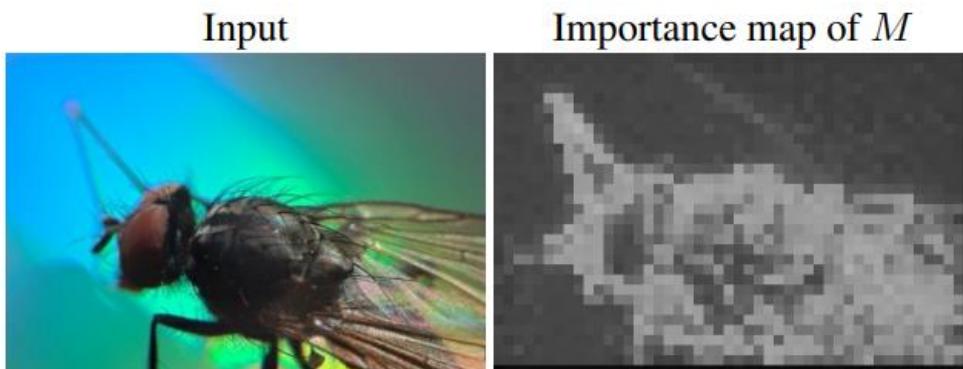
$$\tilde{z}_i = \sum_{j=1}^L \frac{\exp(-\sigma \|z_i - c_j\|)}{\sum_{l=1}^L \exp(-\sigma \|z_i - c_l\|)} c_j$$

Training: differentiable



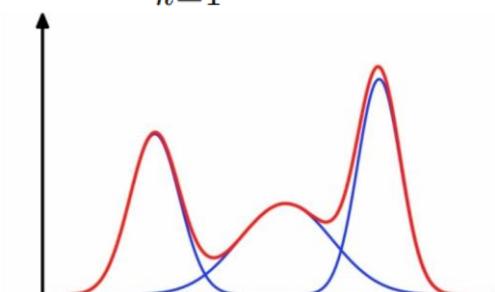
$$\bar{z}_i = \text{tf.stopgradient}(\hat{z}_i - \tilde{z}_i) + \tilde{z}_i$$

- Importance map [4]



- Gaussian Mixture Model (GMM) for entropy [6]

$$p_{\hat{y}|\hat{z}}(\hat{y}|\hat{z}) \sim \sum_{k=1}^K w^{(k)} \mathcal{N}(\mu^{(k)}, \sigma^{2(k)})$$

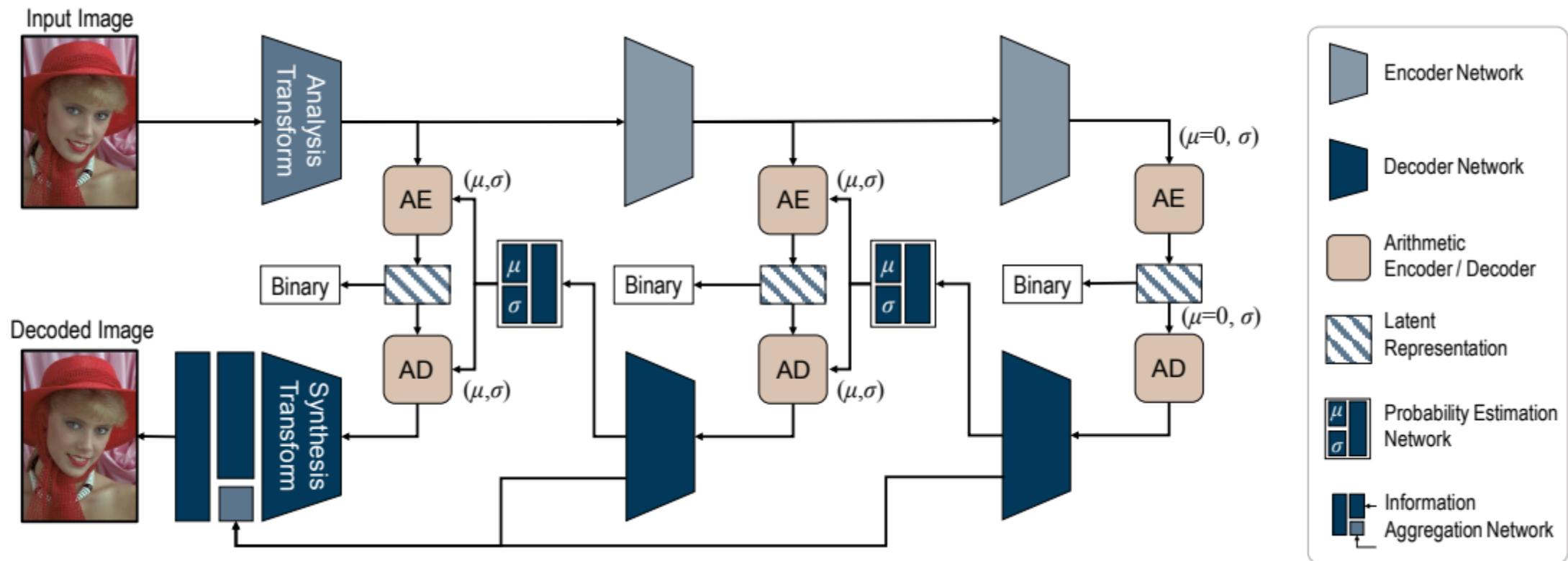
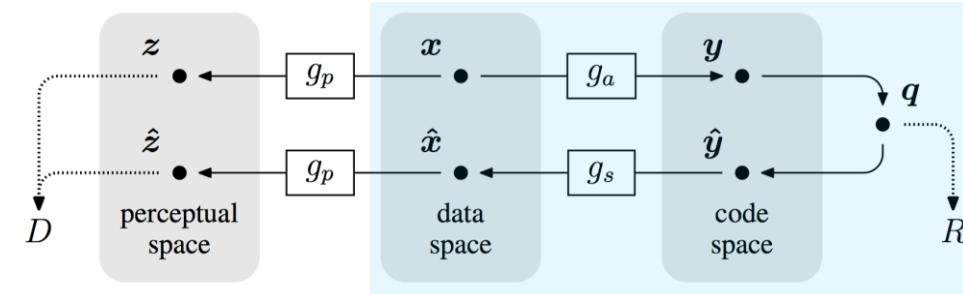


[4] Mentzer, Fabian, et al. "Conditional Probability Models for Deep Image Compression", in CVPR, 2018.

[6] Cheng et al. "Learned Image Compression with Discretized Gaussian Mixture Likelihoods and Attention Modules", in CVPR, 2020.

Learned Image Compression

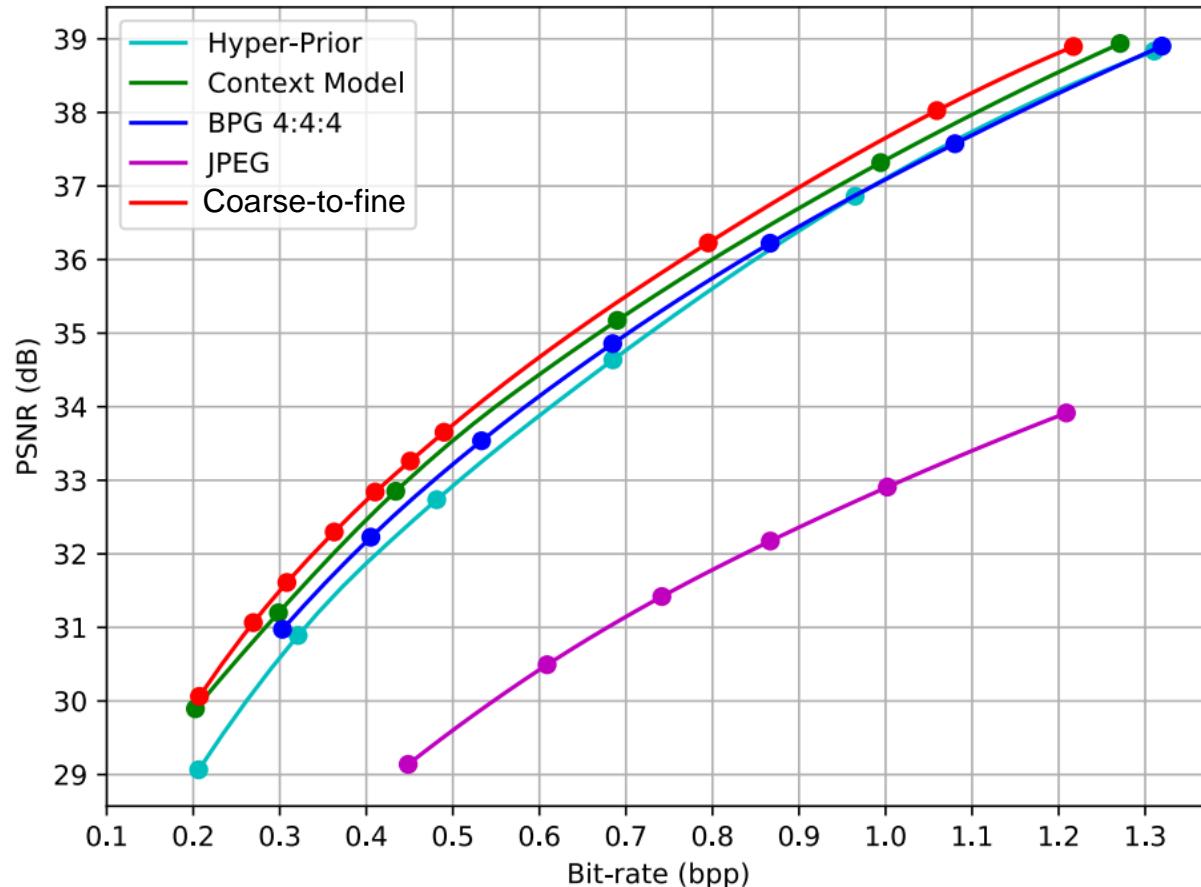
- CNN transformer + **coarse-to-fine** model [7]



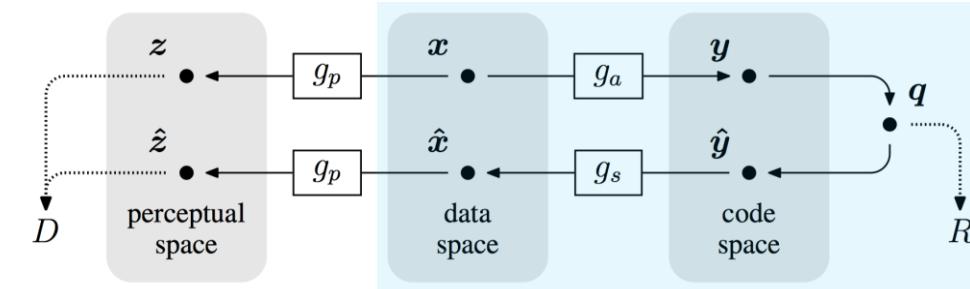
Learned Image Compression

- Performance

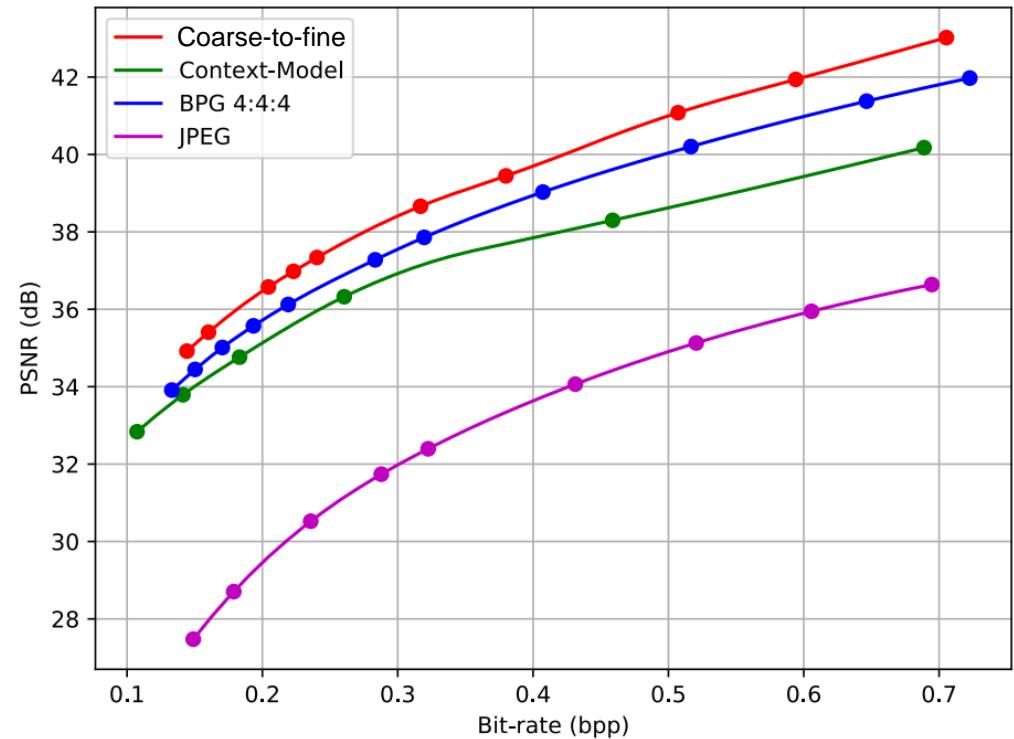
Comparison on Kodak image set



The context (autoregressive) and coarse-to-fine models outperform BPG 4:4:4 (latest traditional standard)



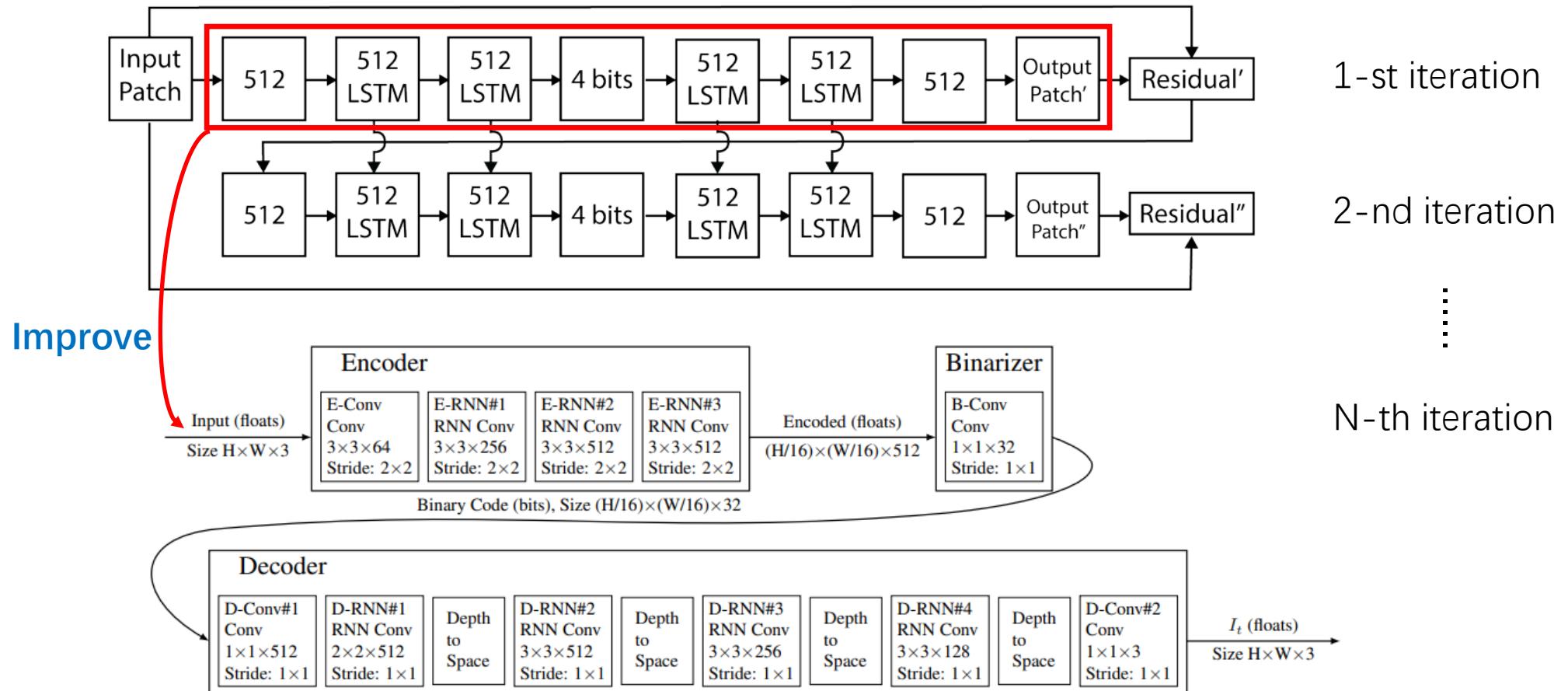
Comparison on Tecnick image set



The rank may vary on different datasets

Learned Image Compression

- Variable rate image compression: RNN-based methods [8, 9]

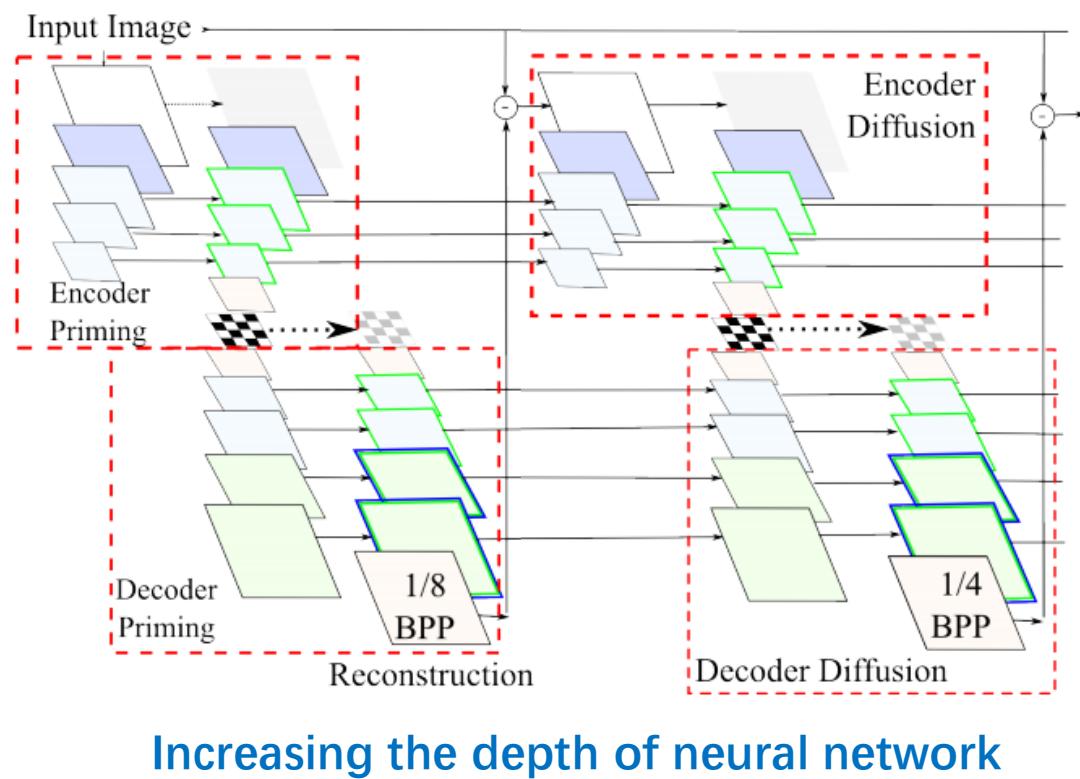
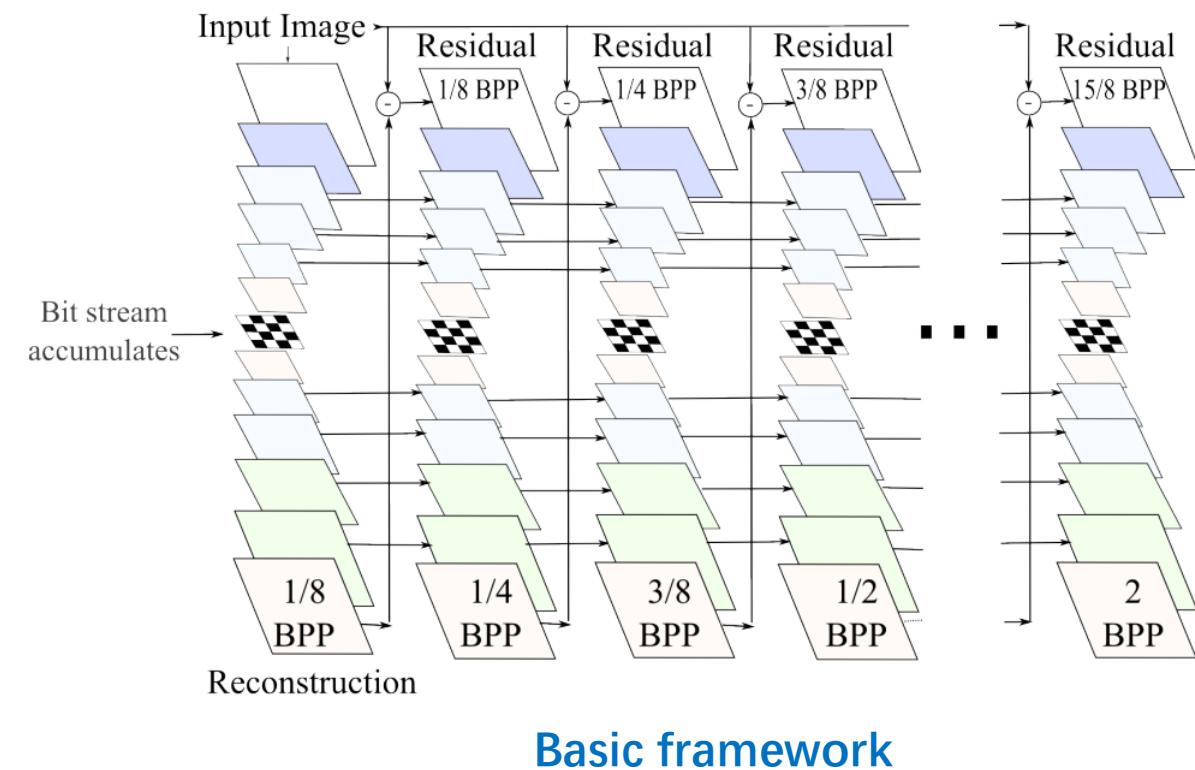


[8] Toderici, George, et al. "Variable Rate Image Compression with Recurrent Neural Networks." in ICLR. 2016.

[9] Toderici, George, et al. "Full Resolution Image Compression with Recurrent Neural Networks." in CVPR, 2017.

Learned Image Compression

- Variable rate image compression: RNN-based methods [10]



[10] Johnston, Nick, et al. "Improved Lossy Image Compression with Priming and Spatially Adaptive Bit Rates for Recurrent Networks." in CVPR. 2018.

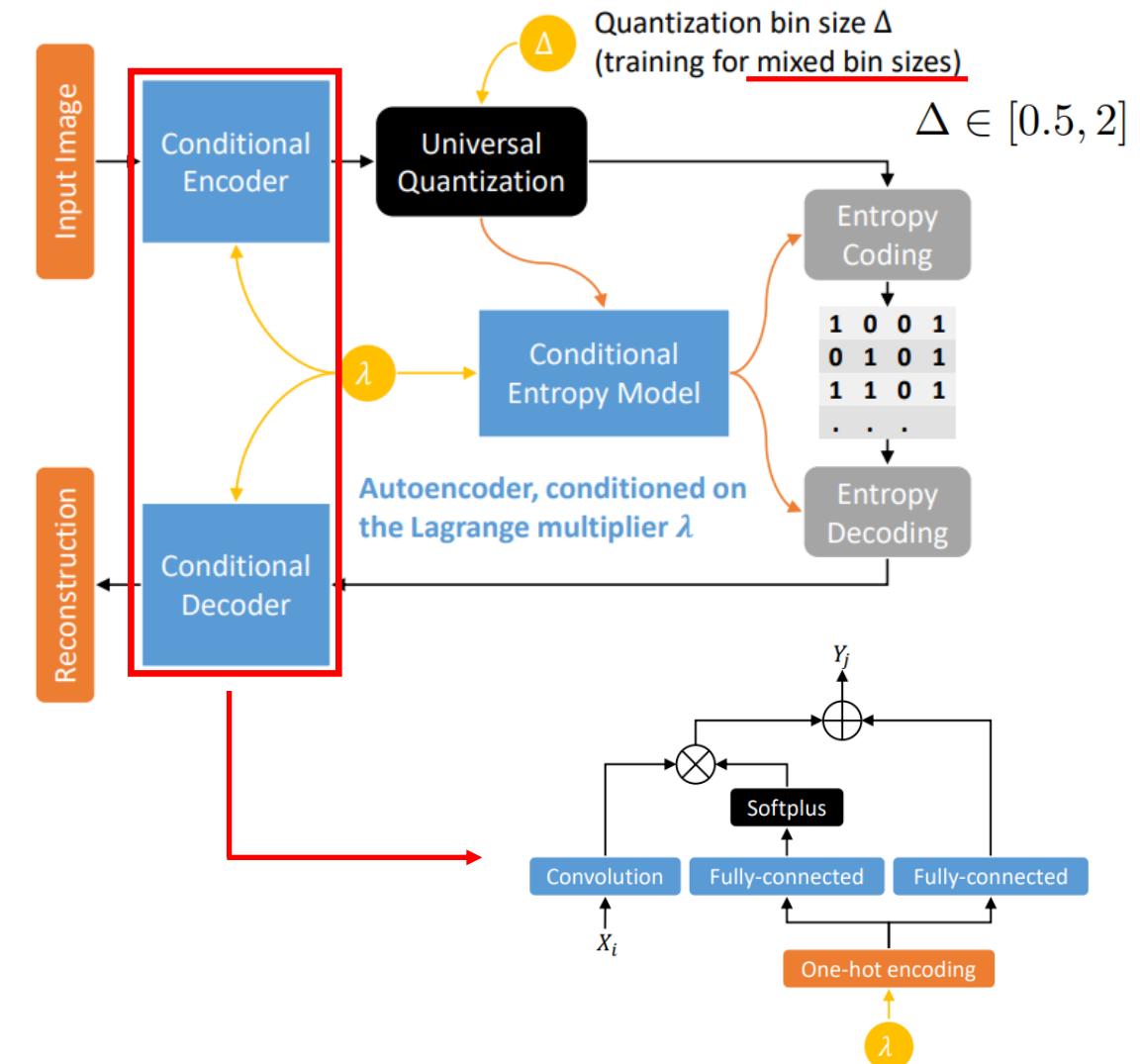
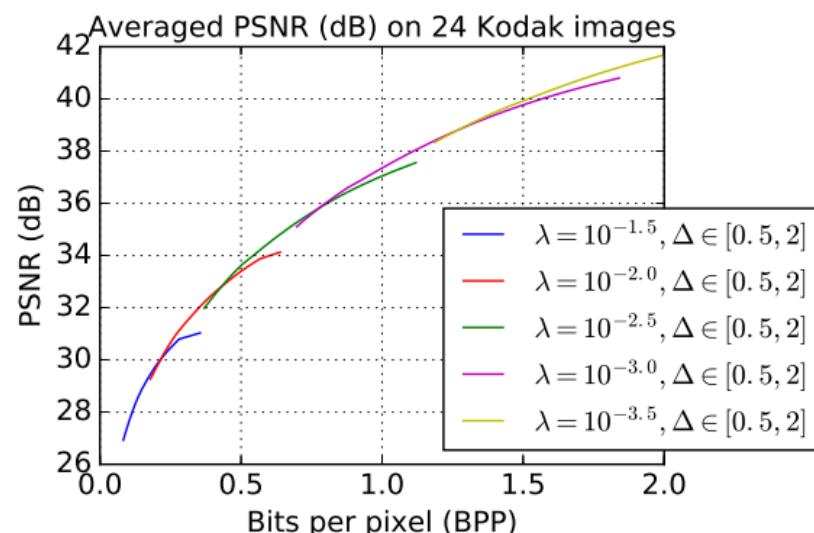
Learned Image Compression

- Variable rate image compression: Conditional autoencoder [11]

Loss function: $\min_{\phi, \theta} \{D_{\phi, \theta} + \underline{\lambda} R_{\phi}\}$

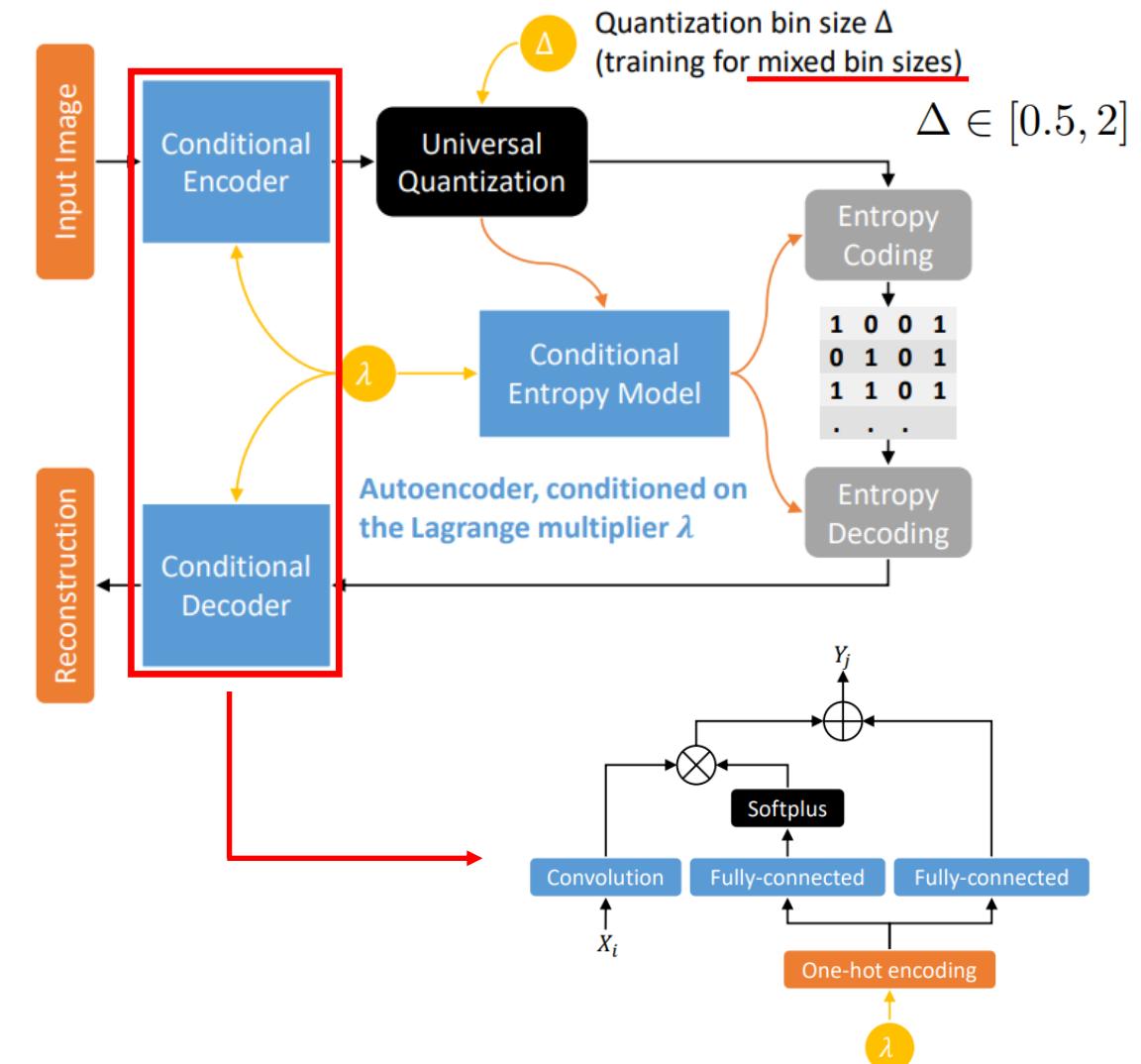
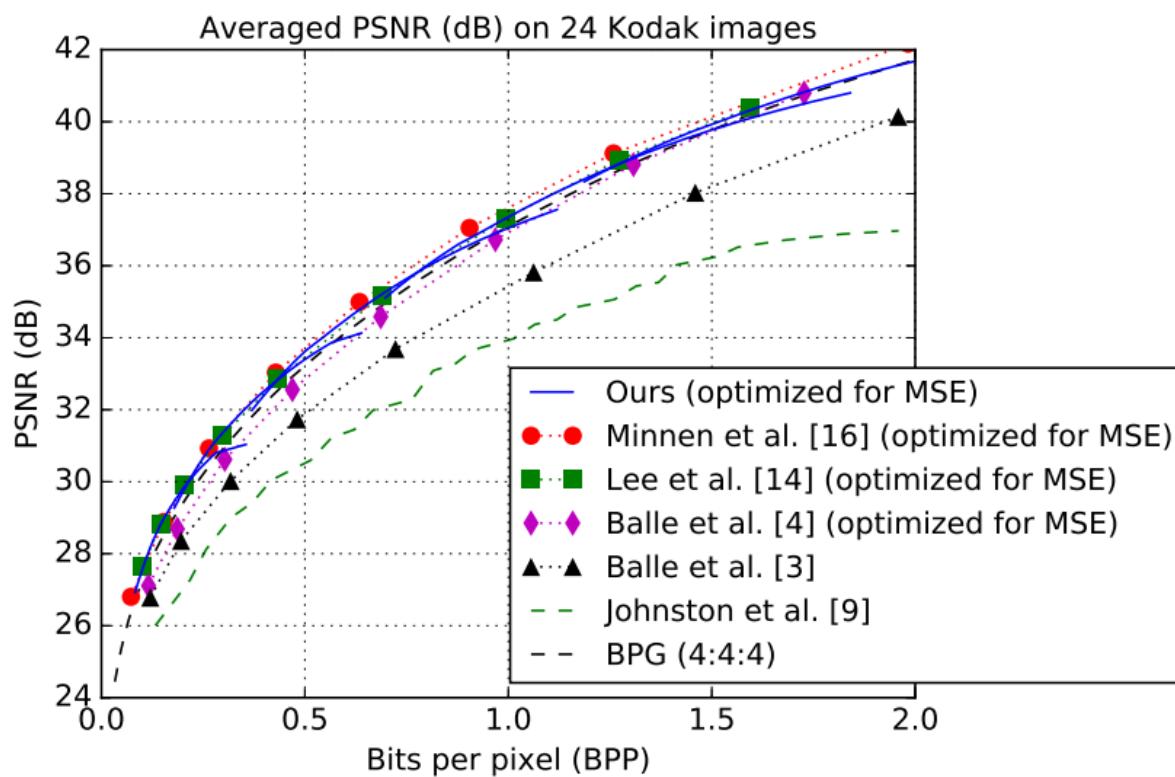
$$\min_{\phi, \theta} \sum_{\lambda \in \Lambda} (D_{\phi, \theta}(\lambda) + \lambda R_{\phi, \theta}(\lambda))$$

$$\min_{\phi, \theta} \sum_{\lambda \in \Lambda} \mathbb{E}_{p(\Delta)} [D_{\phi, \theta}(\lambda, \Delta) + \lambda R_{\phi, \theta}(\lambda, \Delta)]$$



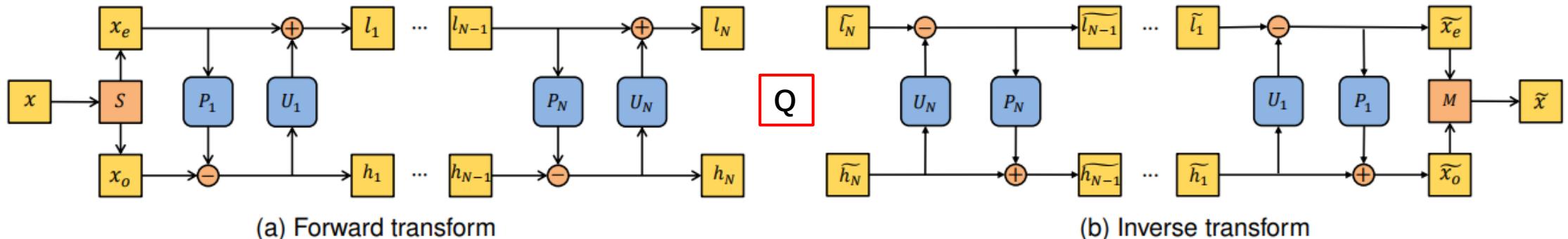
Learned Image Compression

- Variable rate image compression: Conditional autoencoder [11]

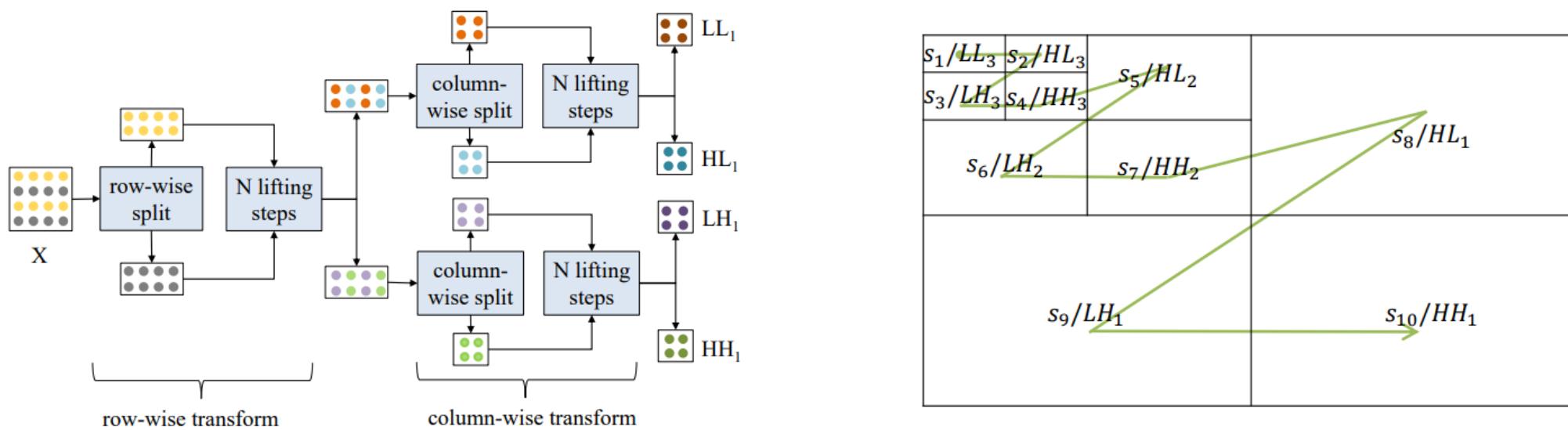


Learned Image Compression

- Variable rate image compression: Wavelet-like transformer [12]

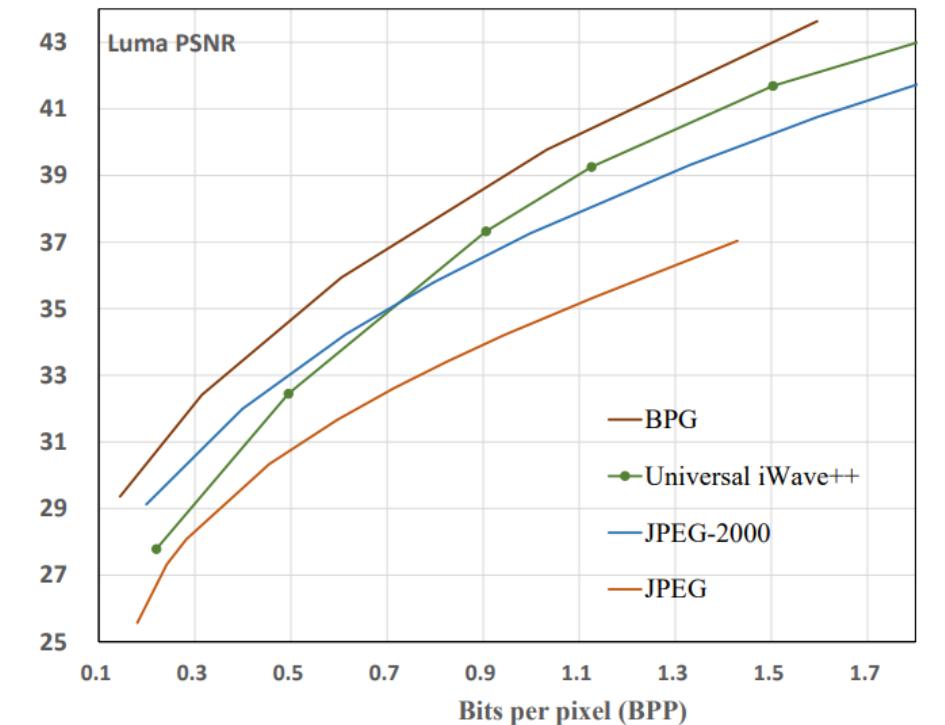
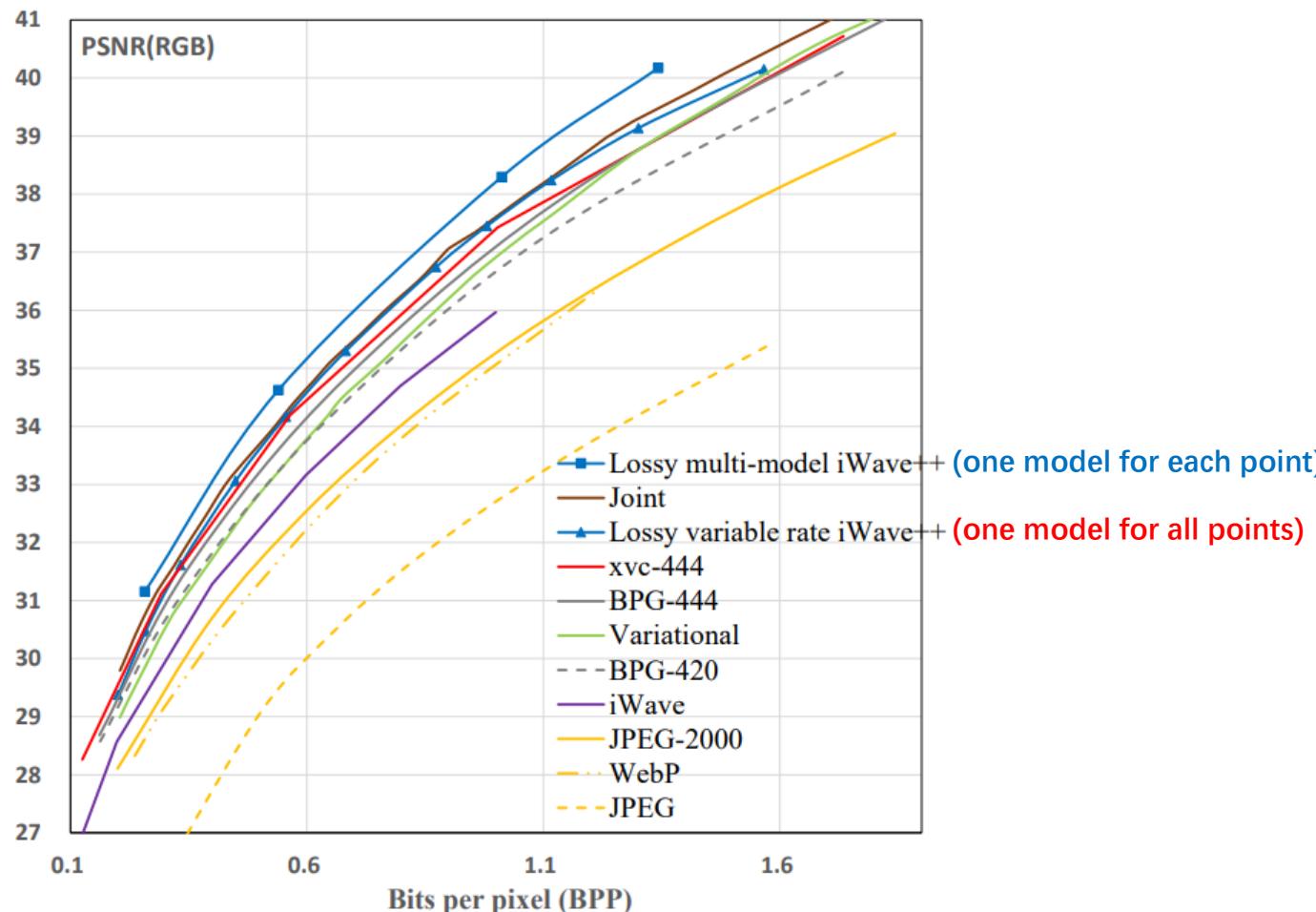


Invertible: achieving lossy and lossless compression by the same framework



Learned Image Compression

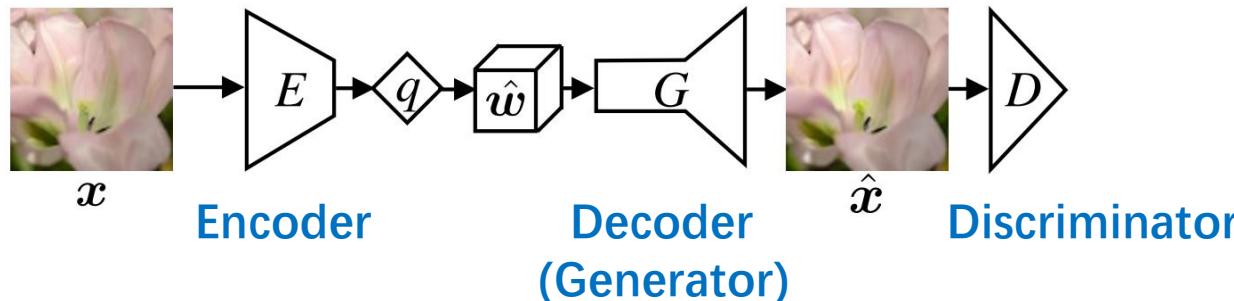
- Variable rate image compression: Wavelet-like transformer [12]



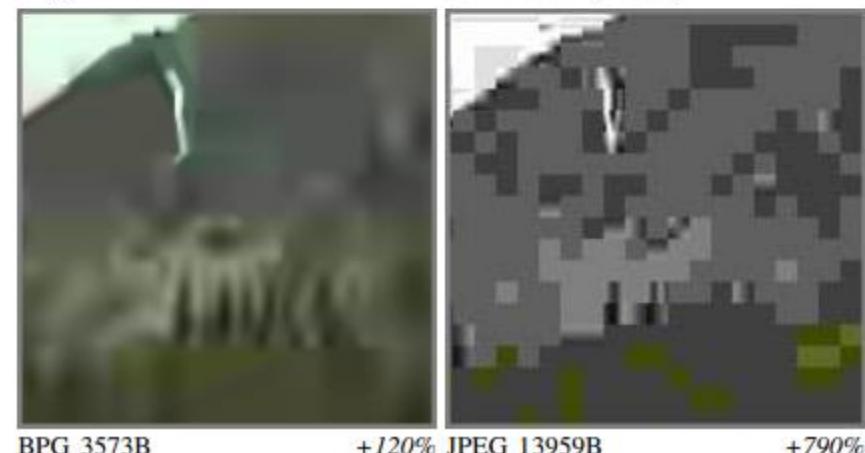
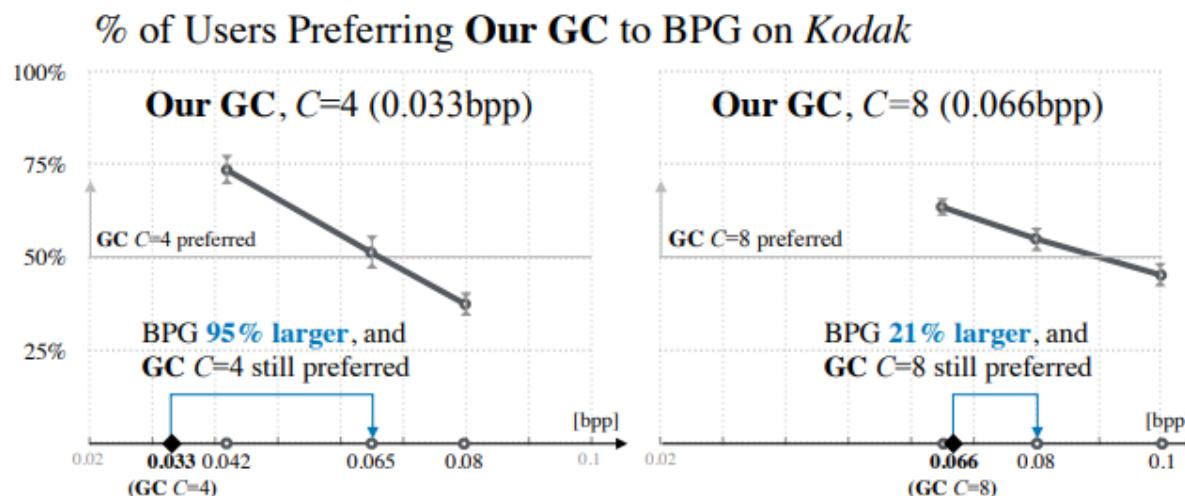
one model for both lossy and lossless compression

Learned Image Compression

- Generative image compression: GAN-based methods [13]

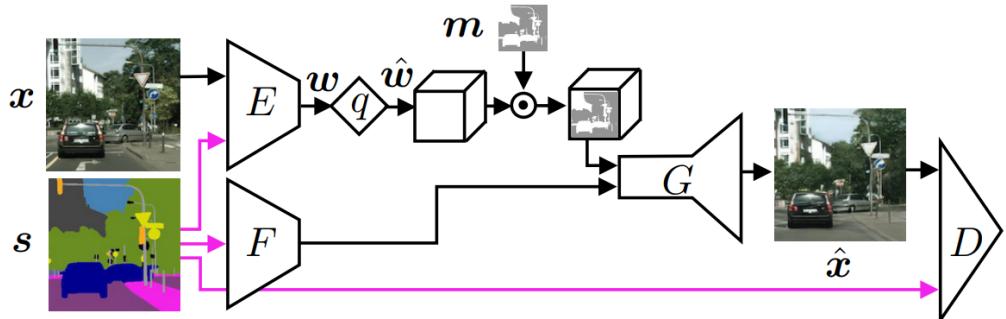


$$\min_{E,G} \max_D \frac{\mathbb{E}[f(D(\hat{w}))] + \mathbb{E}[g(D(G(\hat{w})))]}{+\lambda\mathbb{E}[d(x, G(\hat{w}))] + \beta H(\hat{w})}, \text{ RD loss}$$



Learned Image Compression

- Generative image compression: GAN-based methods [13]



Conditional GAN: $\mathcal{L}_{\text{cGAN}} := \max_D \mathbb{E}[f(D(x, s))] + \mathbb{E}[g(D(G(z, s), s))]$

Selective generative compression (SC): binary heatmap m



road (0.146bpp, -55%)



car (0.227bpp, -15%)



all synth. (0.035bpp, -89%)



people (0.219bpp, -33%)



building (0.199bpp, -39%)

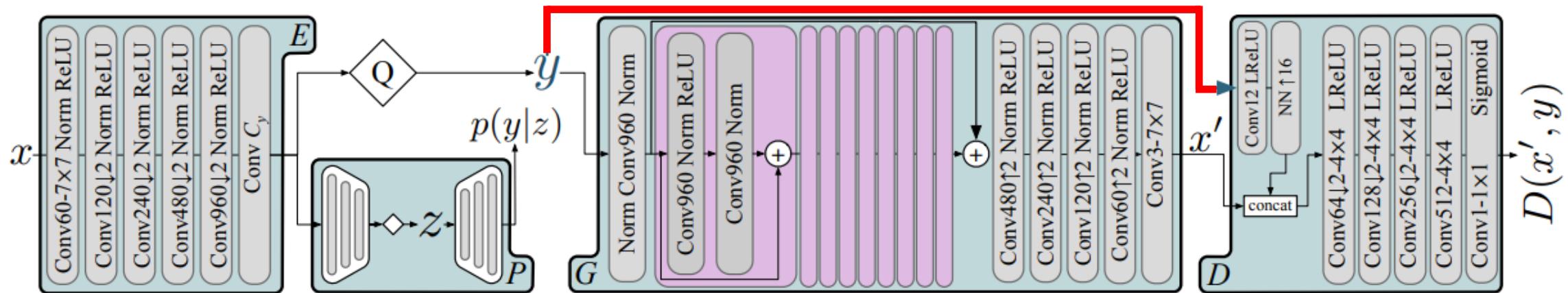


no synth. (0.326bpp, -0%)

Learned Image Compression

- Generative image compression: GAN-based methods [14]

High-Fidelity Generative Image Compression



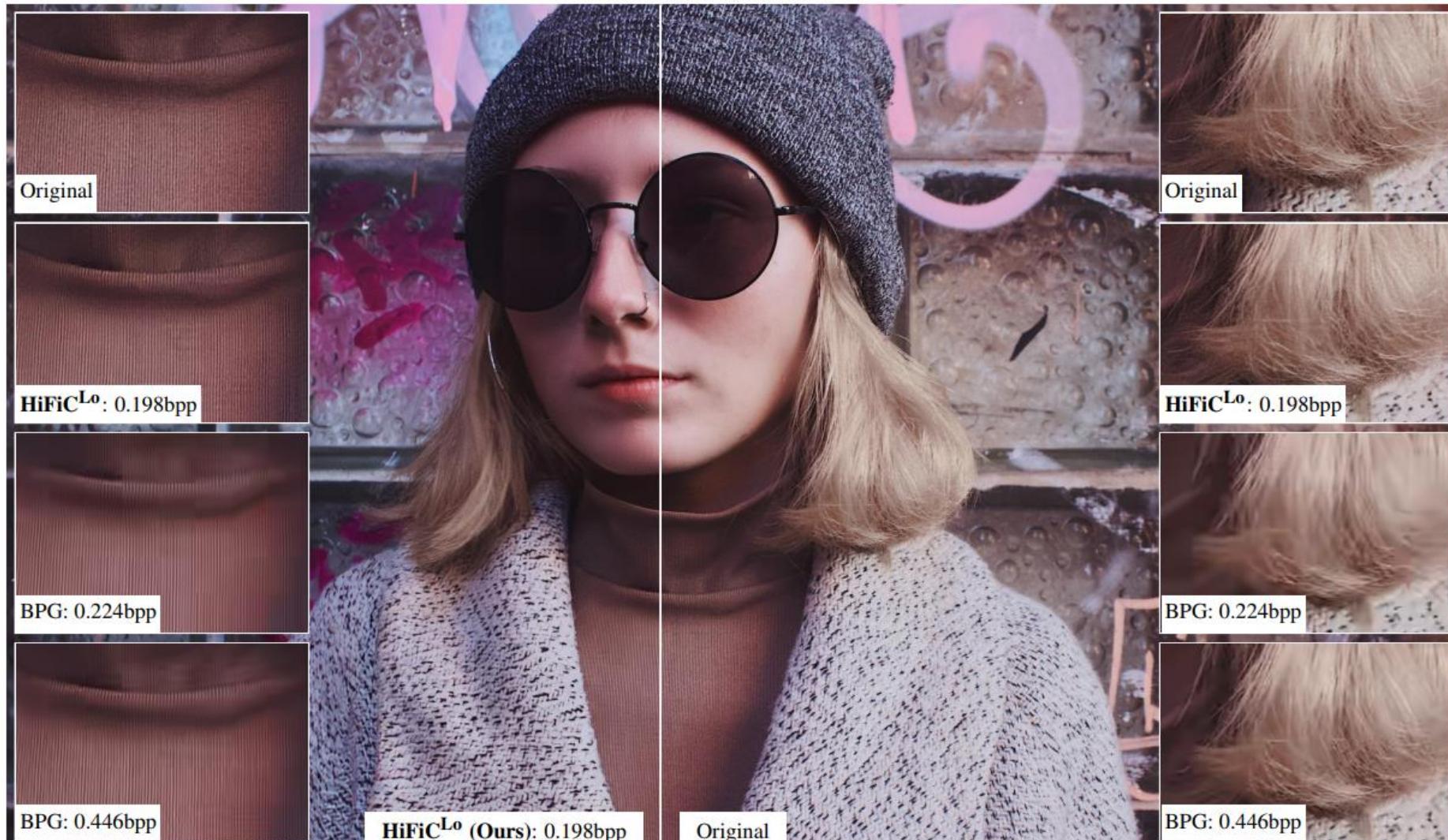
Conditional discriminator:

$$\mathcal{L}_{EGP} = \mathbb{E}_{x \sim p_X} [\lambda r(y) + d(x, x') - \beta \log(D(x', y))],$$

$$\mathcal{L}_D = \mathbb{E}_{x \sim p_X} [-\log(1 - D(x', y))] + \mathbb{E}_{x \sim p_X} [-\log(D(x, y))].$$

Learned Image Compression

- Generative image compression: GAN-based methods [14]



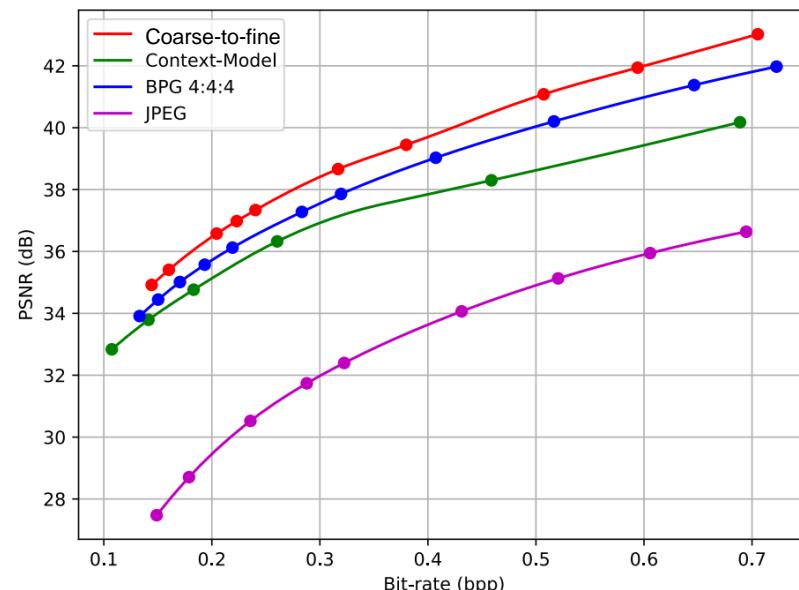
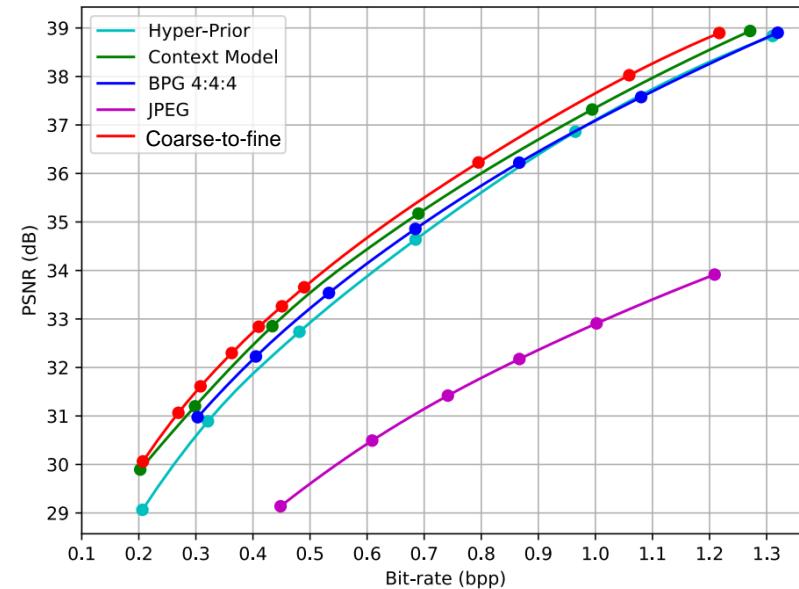
[14] Mentzer, Fabian, et al., "High-Fidelity Generative Image Compression." in NeurIPS. 2020.

Learned Image Compression

Conclusion:

- CNN-based methods
 - Factorized entropy model
 - Hyperprior entropy model
 - Autoregressive entropy model
 - Coarse-to-fine entropy model
 - Conditional auto-encoder (variable bit-rates)
 - Invertible auto-encoder (lossy and lossless by one framework)
- RNN-based methods
 - Variable bit-rate
- GAN-based methods
 - Photo-realistic compressed image with low bit-rate

The state-of-the-art learned image compression methods successfully outperform the latest traditional compression standard BPG 4:4:4



Learned Image Compression

- Will learning-based compression be standardized?

JPEG initiates standardisation of image compression based on AI

The 89th JPEG meeting was held online from 5 to 9 October 2020.

During this meeting multiple JPEG standardisation activities and explorations were discussed and progressed. Notably, the call for evidence on learning-based image coding was successfully completed and evidence was found that this technology promises several new functionalities while offering at the same time superior compression efficiency, beyond the state of the art.

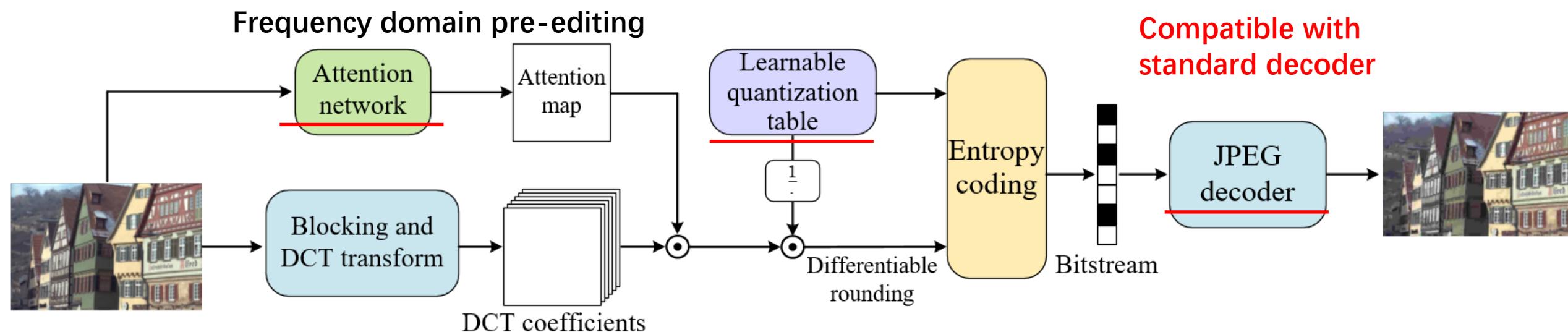
JPEG AI

At the 89th meeting the submissions to the Call for Evidence on learning-based image coding were presented and discussed. Four submissions were received in response to the Call for Evidence. The results of the subjective evaluation of the submissions to the Call for Evidence were reported and discussed in detail by experts. It was agreed that there is strong evidence that learning-based image coding solutions can outperform the already defined anchors in terms of compression efficiency, when compared to state-of-the-art conventional image coding architecture. Thus, it was decided to create a new standardisation activity for a JPEG AI on learning-based image coding system, that applies machine learning tools to achieve substantially better compression efficiency compared to current image coding systems, while offering unique features desirable for an efficient distribution and consumption of images. This type of approach should allow to obtain an efficient compressed domain representation not only for visualisation, but also for machine learning based image processing and computer vision. JPEG AI releases to the public the results of the objective and subjective evaluations as well as a first version of common test conditions for assessing the performance of leaning-based image coding systems.

New Trends in Learned Image Compression

- Learning-based compression compatible with traditional standards (e.g., JPEG)

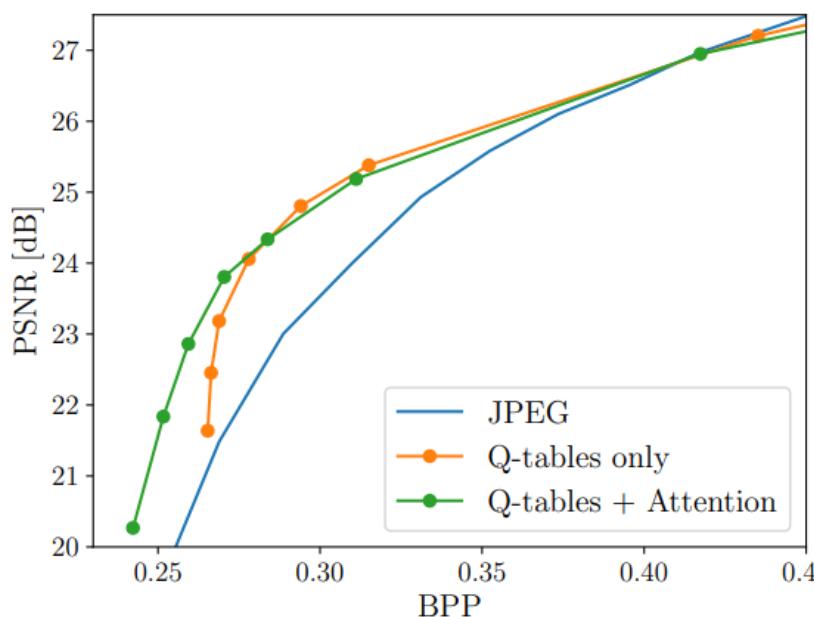
We made an attempt: [15]



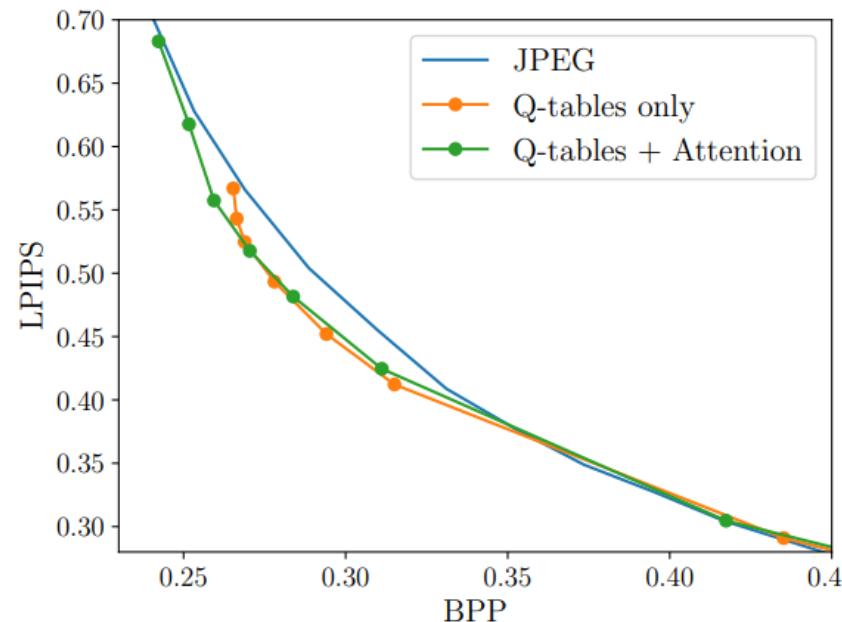
New Trends in Learned Image Compression

- Learning-based compression compatible with traditional standards (e.g., JPEG)

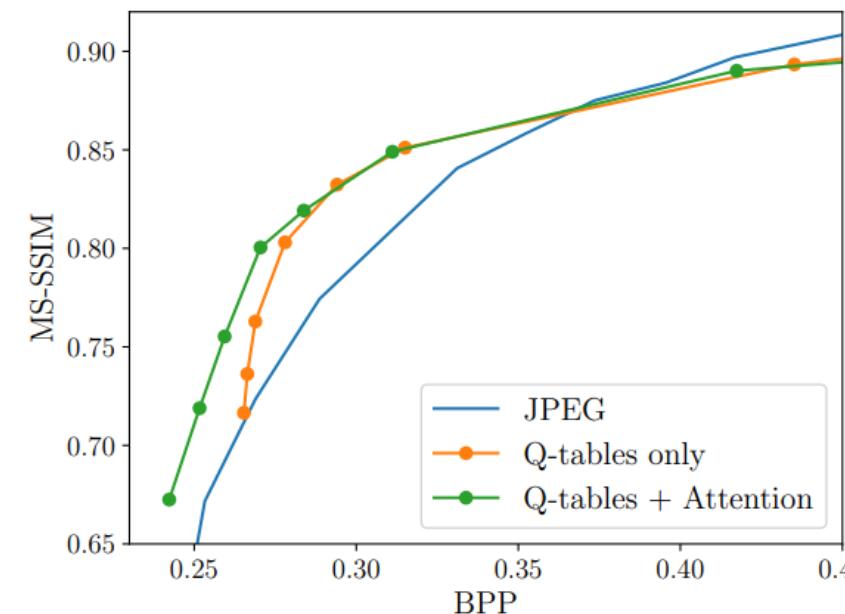
We made an attempt: [15]



PSNR on Kodak



LPIPS on Kodak



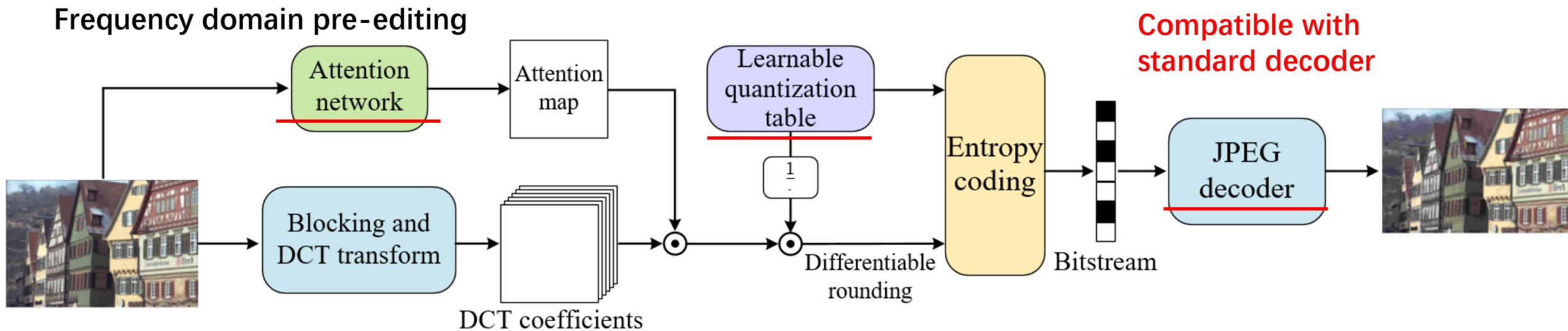
MS-SSIM on Kodak

New Trends in Learned Image Compression

- Learning-based compression compatible with traditional standards (e.g., JPEG)

We made an attempt: [15]

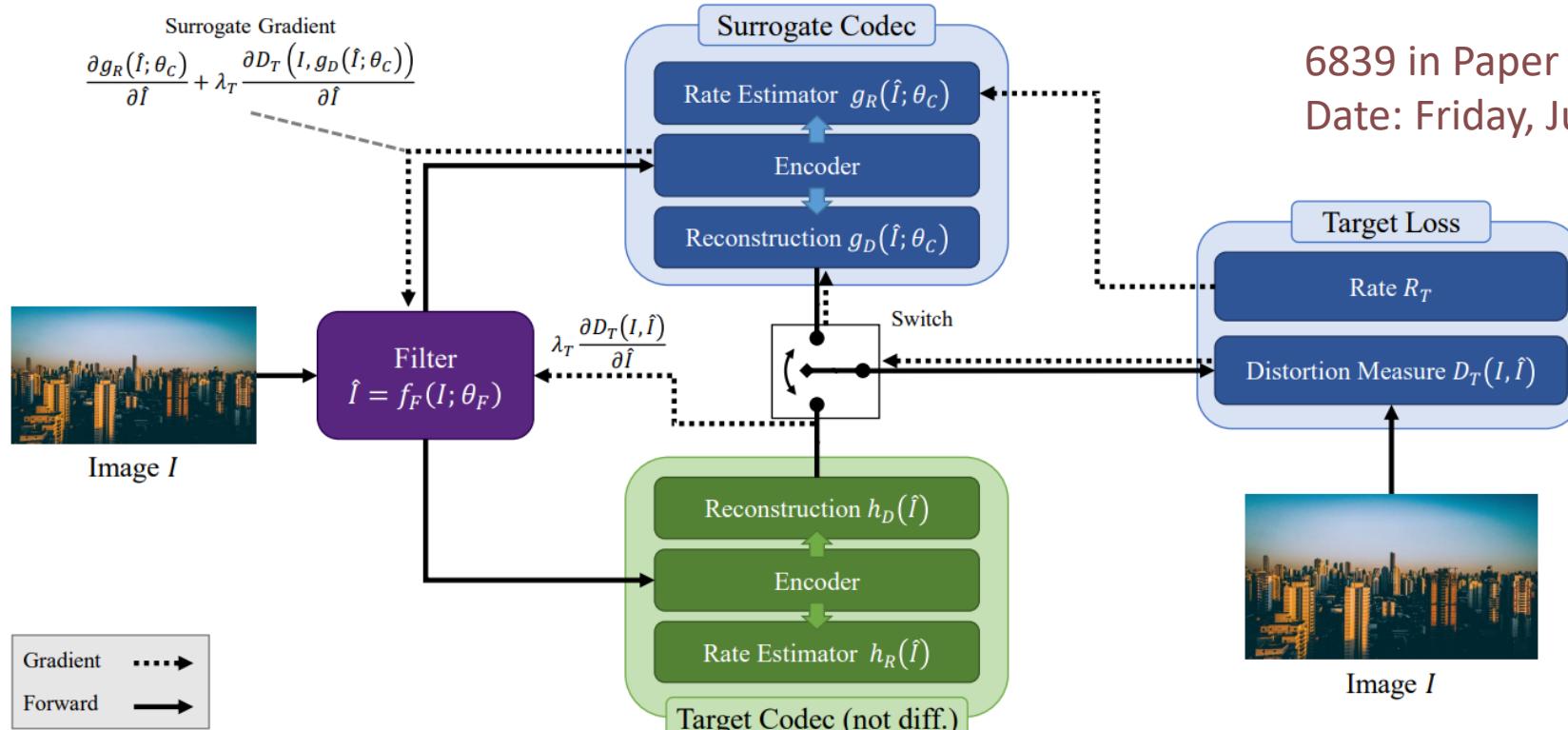
Frequency domain pre-editing



- We achieve better rate-distortion performance **without changing the standard decoder**
- The compressed image can be decoded (viewed) on **any common device**, e.g., mobile, IPad, PC, etc.

New Trends in Learned Image Compression

- Learning-based compression compatible with traditional standards (e.g., JPEG)

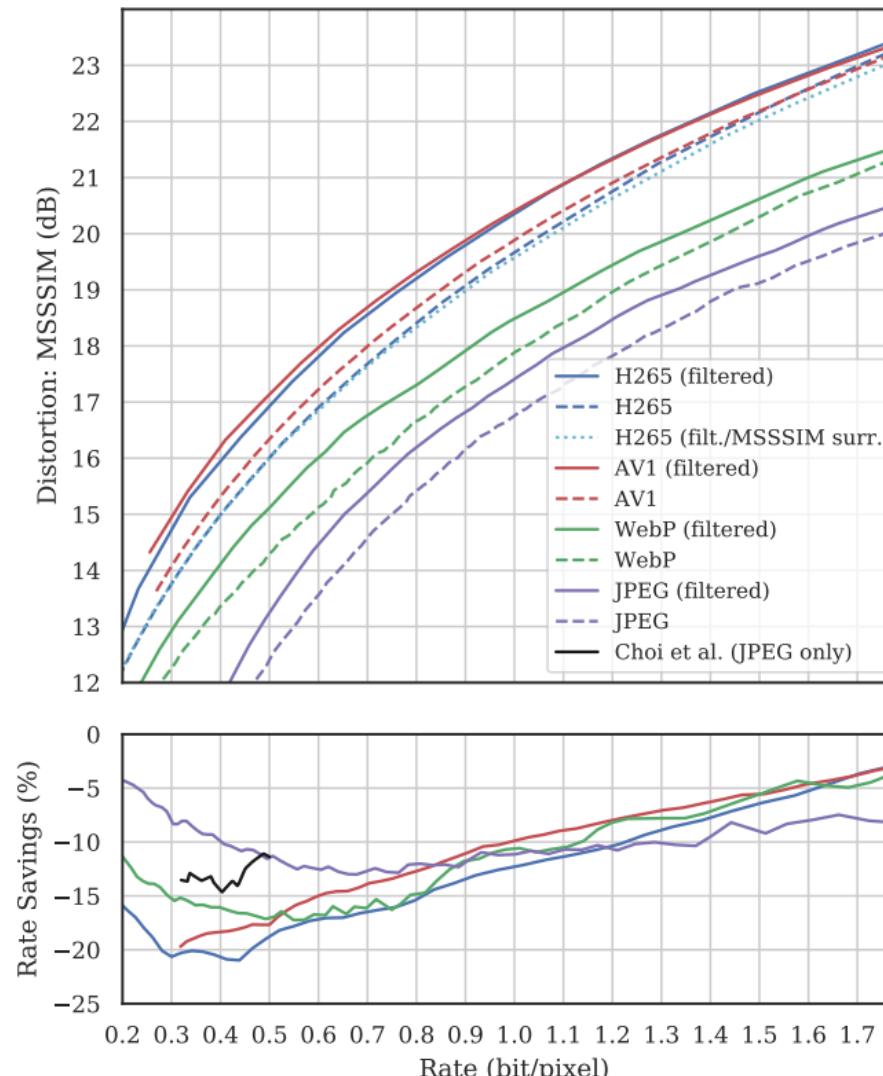


6839 in Paper Session Twelve
Date: Friday, June 25, 2021 6:00 – 8:30

Figure 1. **Structural overview of our method.** The goal is to obtain a trained filter $f_F(I; \theta_F)$ to optimise the input image I for encoding by a target codec. This target is typically not differentiable. A surrogate codec is used instead. It provides a differentiable rate estimate. For the reconstruction there are two options as indicated by the switch. The first is to take the surrogate's reconstruction, the second to invoke the target during the forward pass. The gradient flows back accordingly either through the surrogate's reconstruction or directly into the filter. The target distortion measure D_T can be chosen freely. The surrogate codec is pre-trained with a distortion measure similar to the one of the target codec so as to imitate its behaviour. At testing time, the filter is applied to the image before it's encoded by the target codec.

New Trends in Learned Image Compression

- Learning-based compression compatible with traditional standards (e.g., JPEG)



6839 in Paper Session Twelve
Date: Friday, June 25, 2021 6:00 – 8:30

New Trends in Learned Image Compression

- Learned Image Compression with Implicit Neural Network [17]

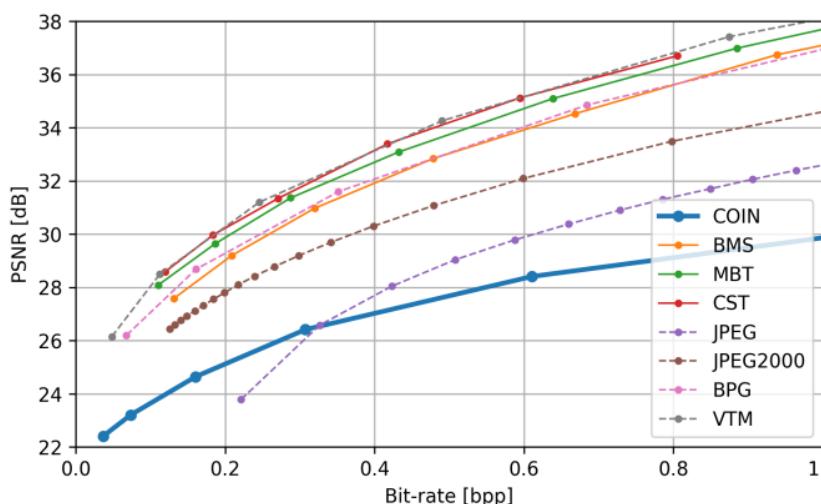
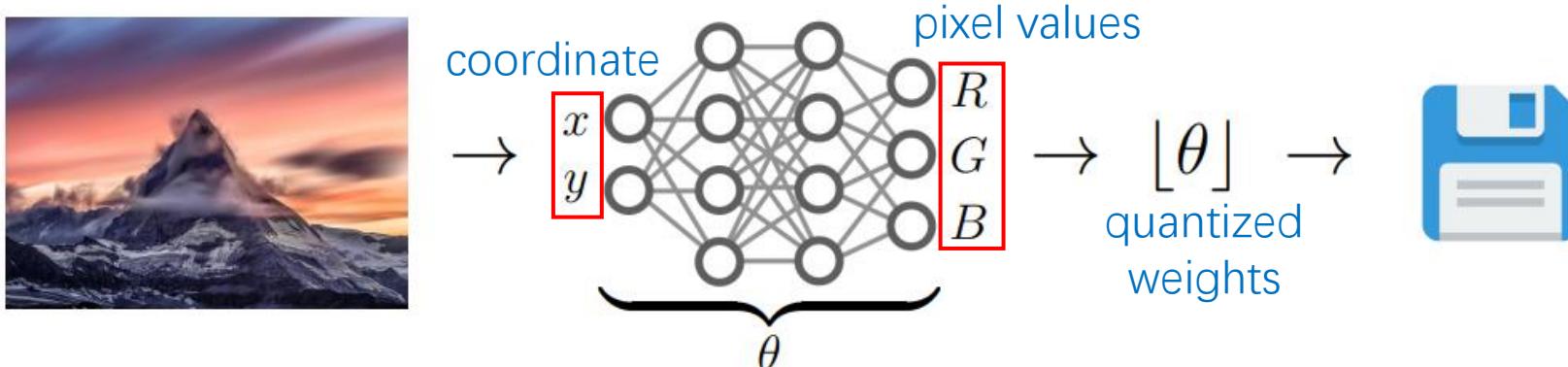


Figure 2: Rate distortion plots on the Kodak dataset.

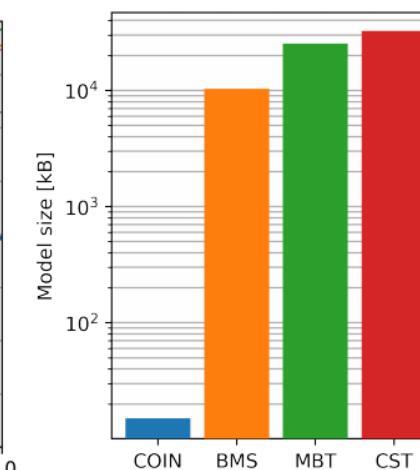


Figure 3: Model sizes at 0.3bpp.

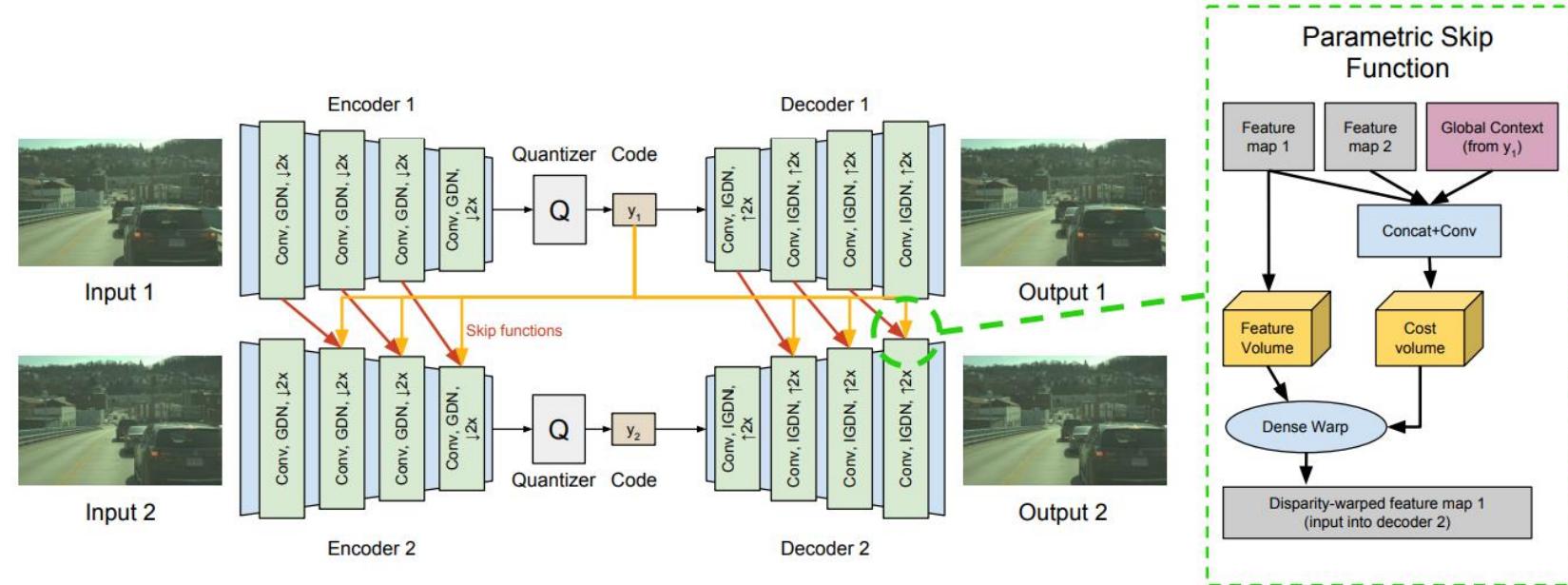
New Trends in Learned Image Compression

- Learning for Stereo Image Compression

DSIC: Deep Stereo Image Compression [18]



Many applications such as autonomous vehicles and 3D movies involve the use of stereo camera pairs. Stereo image compression can be seen as in-between the work of image and video compression. [15]

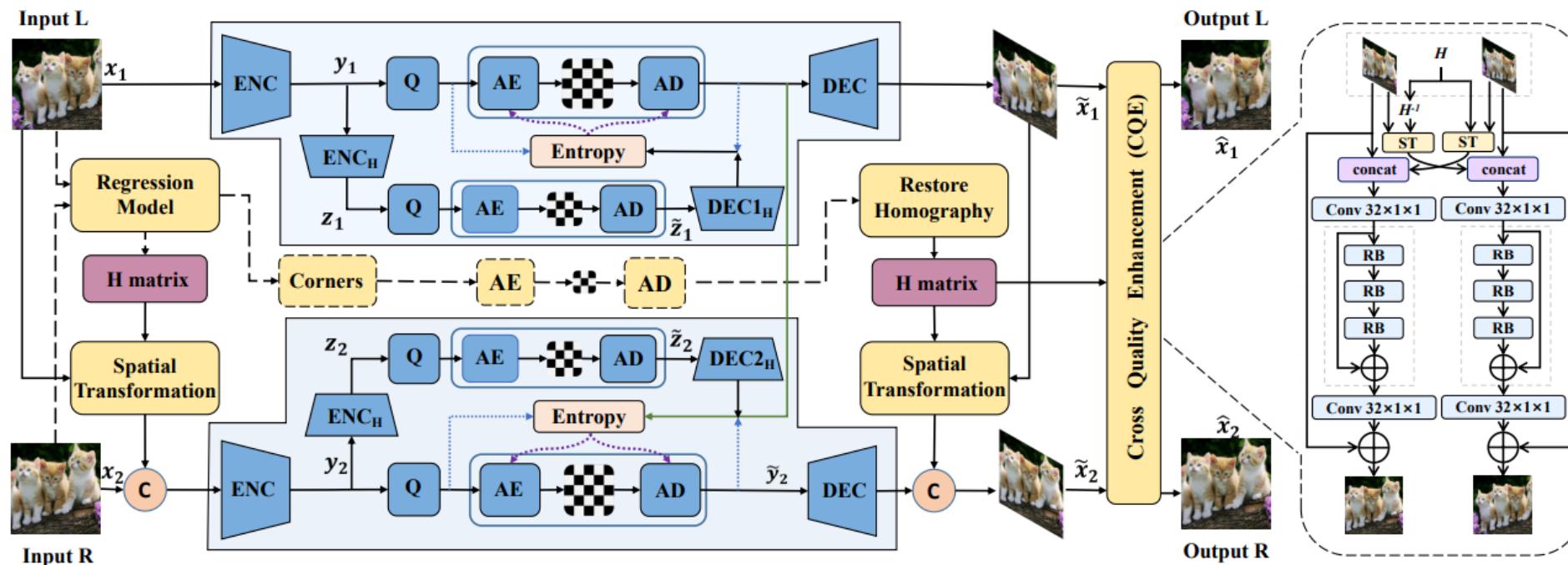


New Trends in Learned Image Compression

- Learning for Stereo Image Compression

Deep Homography for Efficient Stereo Image Compression [19]

6772 in Paper Session Two
Date: Monday, June 21, 2021 22:00 – 24:30



Stereo images are captured at the same time → only angle change w/o temporal motion

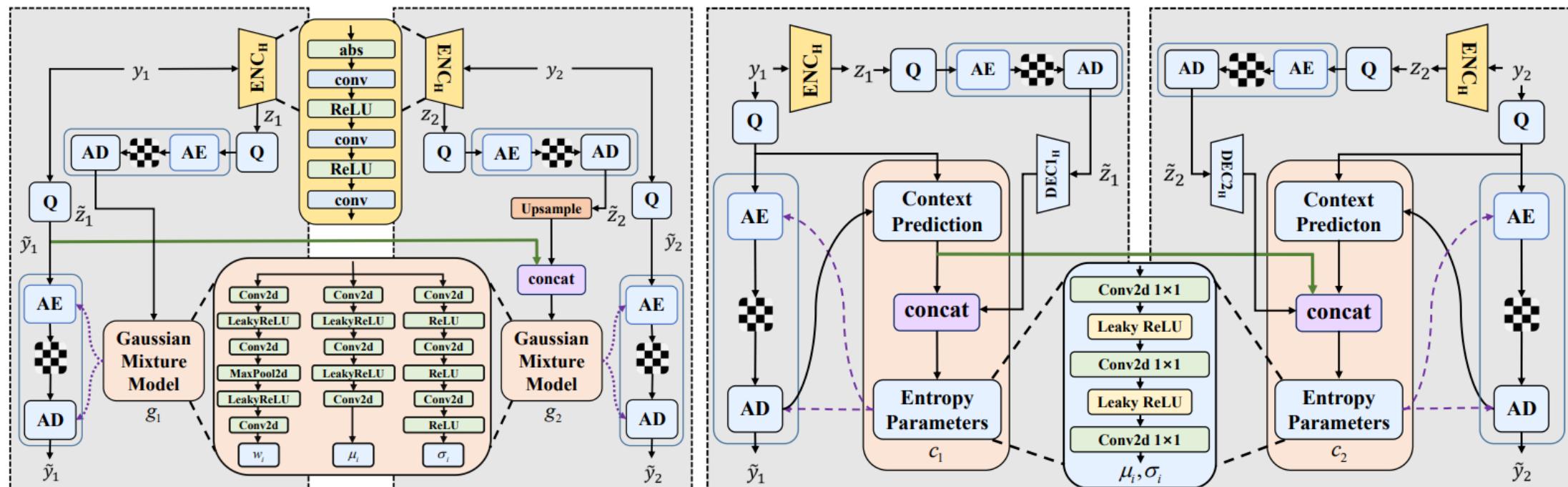
$$\begin{bmatrix} u' \\ v' \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{H} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$

New Trends in Learned Image Compression

- Learning for Stereo Image Compression

6772 in Paper Session Two

Date: Monday, June 21, 2021 22:00 – 24:30

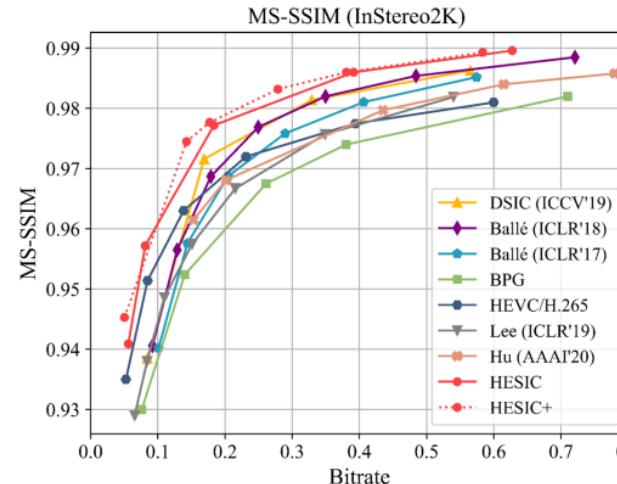
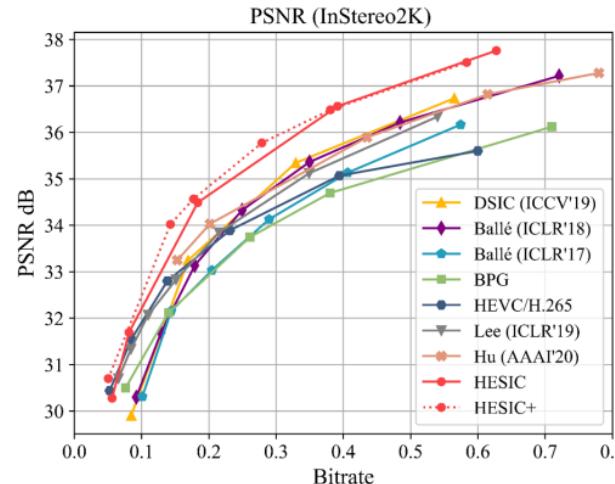


(a) GMM-based stereo entropy model.

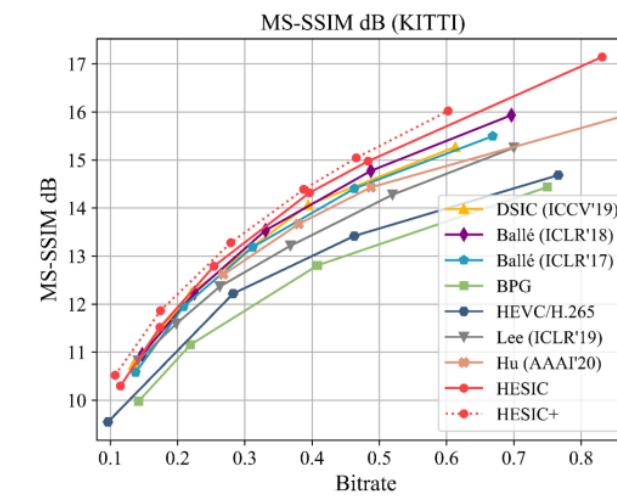
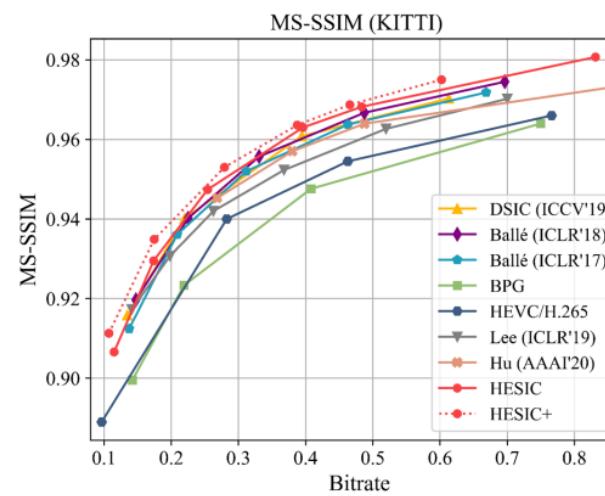
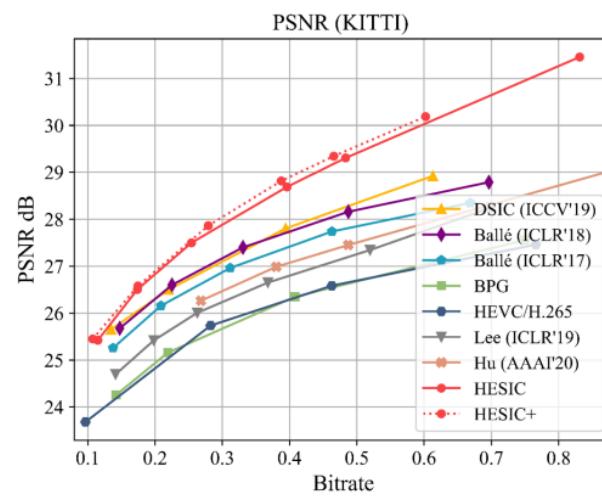
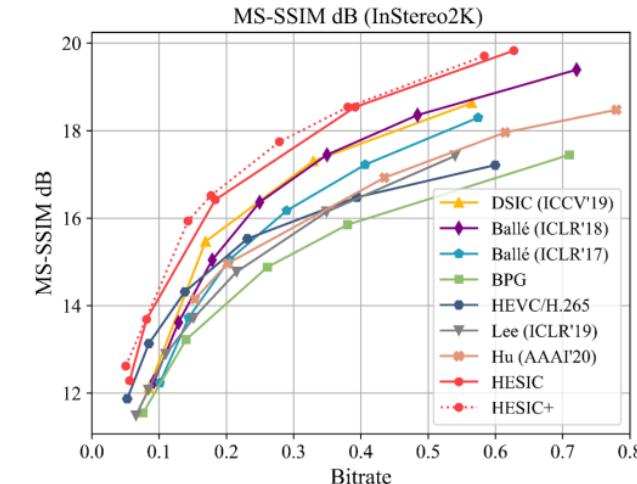
(b) Context-based stereo entropy model.

New Trends in Learned Image Compression

- Learning for Stereo Image Compression



6772 in Paper Session Two
Date: Monday, June 21, 2021 22:00 – 24:30



Learned Image Compression

- Open source codes:

- Ballé et al., (factorized), Ballé et al., (hyperprior):

<https://github.com/tensorflow/compression> (**TensorFlow**)

- Ballé et al., (factorized), Ballé et al., (hyperprior), Minnen et al., (autoregressive):

<https://interdigitalinc.github.io/CompressAI/index.html> (**PyTorch**)

- Lee et al., (context-adaptive):

https://github.com/JooyoungLeeETRI/CA_Entropy_Model

- Mentzer et al., (autoregressive + importance map):

<https://github.com/fab-jul/imgcomp-cvpr>

- Cheng et al., (GMM entropy model):

<https://github.com/ZhengxueCheng/Learned-Image-Compression-with-GMM-and-Attention>

- Hu et al., (coarse-to-fine):

<https://github.com/huzi96/Coarse2Fine-ImaComp>

- Ma et al., (wavelet-like transformer):

<https://github.com/mahaichuan/Versatile-Image-Compression>

- Mentzer et al., (generative compression):

<https://github.com/tensorflow/compression/tree/master/models/hifc>

- Dupont et al., (compression with implicit neural network):

<https://github.com/EmilienDupont/coin>

- Deng et al., (stereo image compression):

<https://github.com/ywz978020607/HESIC>

Learning for Visual Data Compression

CVPR 2021 Tutorial



Thanks for attention

Q & A



Ren Yang
ETH Zurich, Switzerland

ren.yang@vision.ee.ethz.ch



Radu Timofte
ETH Zurich, Switzerland

radu.timofte@vision.ee.ethz.ch