

# Learned Disentangled Latent Representations for Scalable Image Coding for Humans and Machines

Ezgi Ozyilkan<sup>†,\*</sup>, Mateen Ulhaq<sup>‡,\*</sup>, Hyomin Choi<sup>\*</sup>, Fabien Racapé<sup>\*</sup>

¶ Joint first authors.

† Dept. of Electrical and Computer Engineering, New York University

‡ School of Engineering Science, Simon Fraser University

\* InterDigital – Emerging Technologies Lab

ezgi.ozyilkan@nyu.edu, mulhaq@sfu.ca,  
{hyomin.choi, fabien.racape}@interdigital.com

This work was done while E. Ozyilkan and M. Ulhaq were interns at InterDigital.



# Contents

1. Introduction to Scalable Image Compression
2. Related Work
3. Proposed Framework and Architecture
4. Experiments and Results
5. Information-Theoretic Insights into Information Flow
6. Final Remarks

# Traditional Transform Coding: JPEG in a nutshell

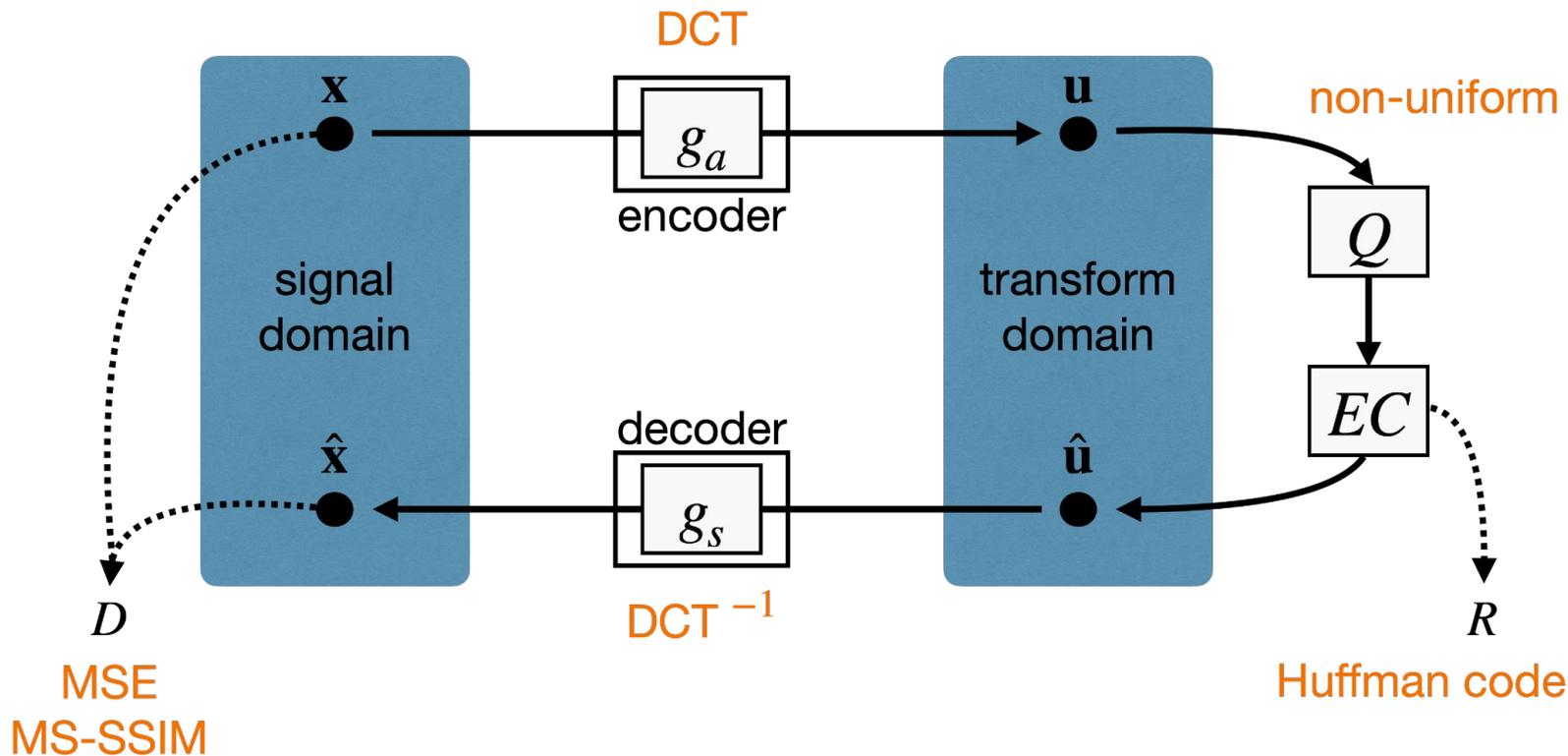


Figure adapted from [J. Ballé et al.], "End-to-end optimized image compression," ICLR, 2017.

# Nonlinear Transform Coding

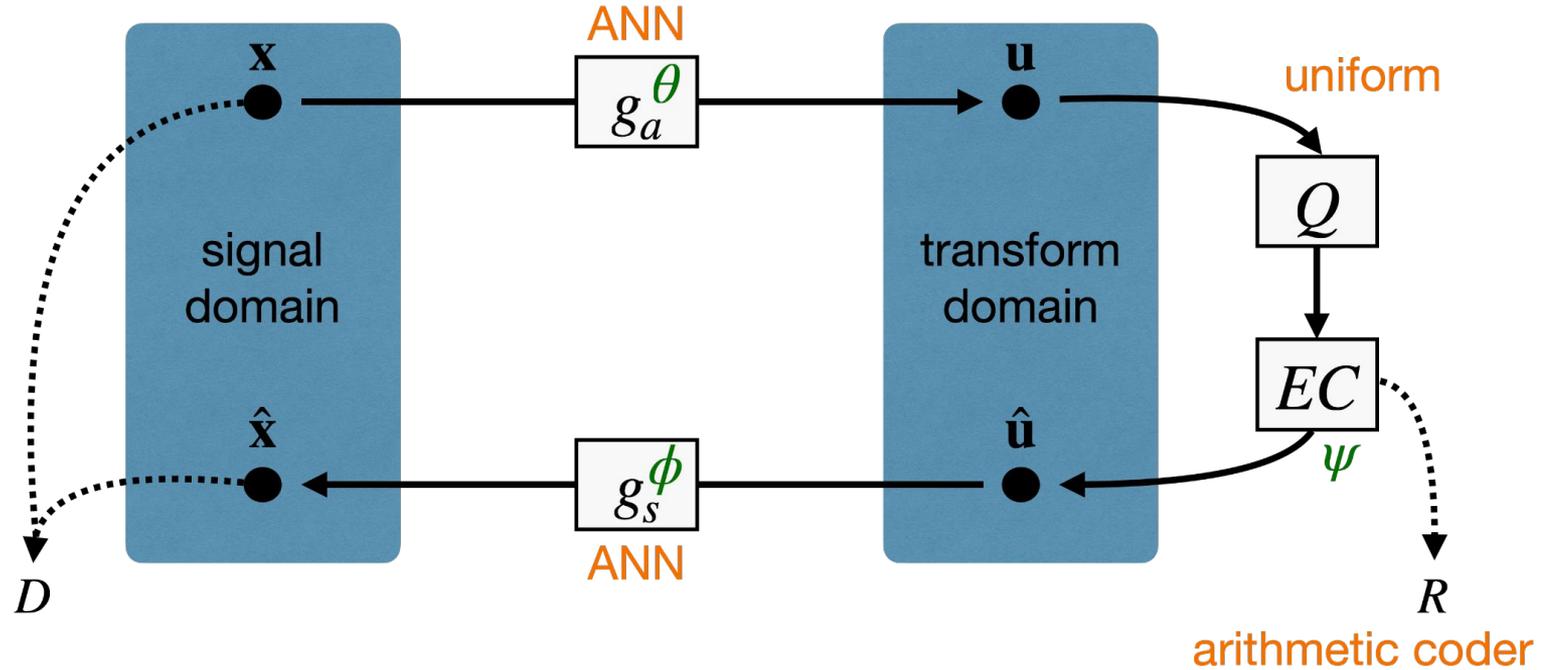


Figure adapted from [J. Ballé et al.], "End-to-end optimized image compression," ICLR, 2017.

# Multi-Task Image Coding

Split transform domain latent space for machine analytics. “Video Coding for Machines” (VCM).

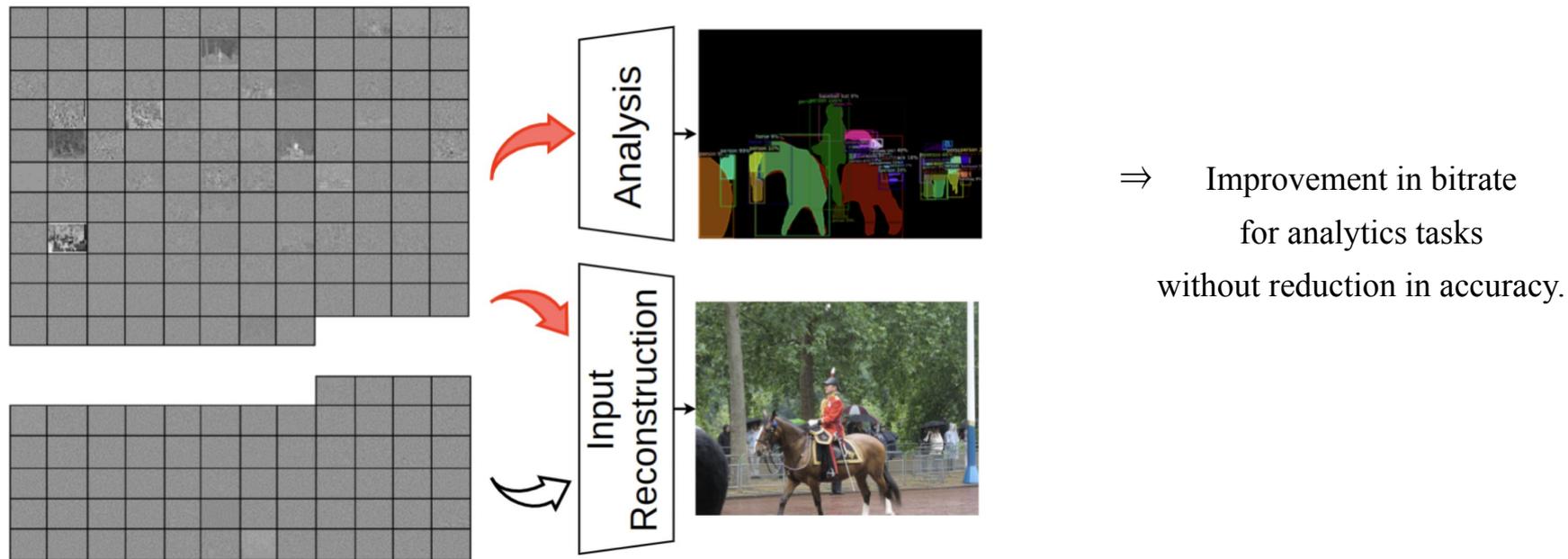


Figure courtesy of [H. Choi et al.], “Scalable Image Coding for Humans and Machines,” *IEEE Transactions on Image Processing*, 2022.

# Scalable Image Compression

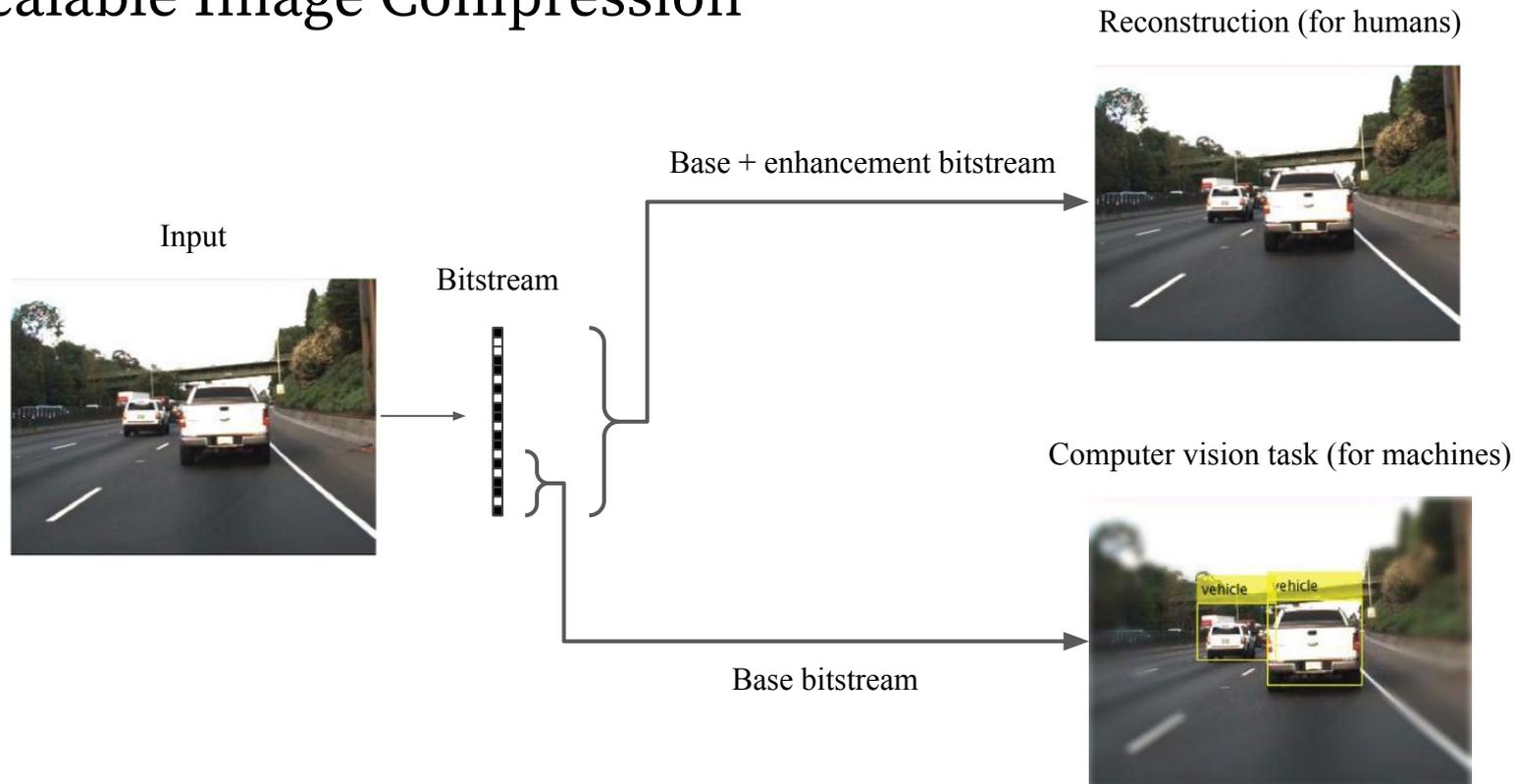
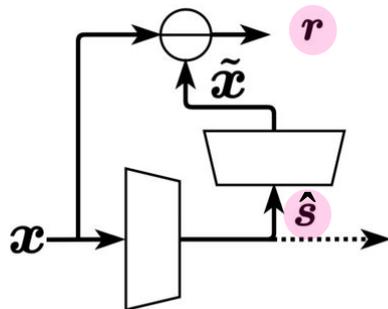


Image taken from <https://www.mathworks.com/discovery/object-detection.html>

# Prior Work

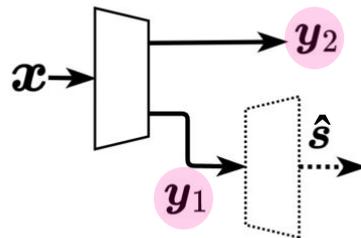
Chamain et al.



(a)

“Enhancement” is residual error in reconstructing from “base”.

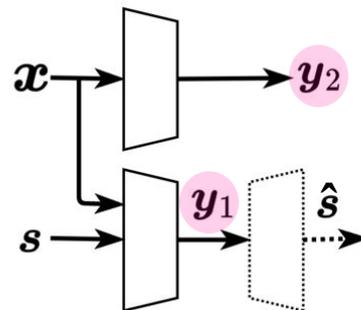
Choi et al.



(b)

“Base” and “enhancement” obtained from same transform on  $x$ .

Proposed



(c)

“Base” from  $x, s$  and “enhancement” only from  $x$ .

*Transmitted bitstreams are highlighted.*

[Yan et al.], “End-to-end optimized image compression for machines, a study,” *DCC*, 2021.

[Choi et al.], “Scalable Image Coding for Humans and Machines,” *IEEE Transactions on Image Processing*, 2022.

# Prior Work: Choi et al.

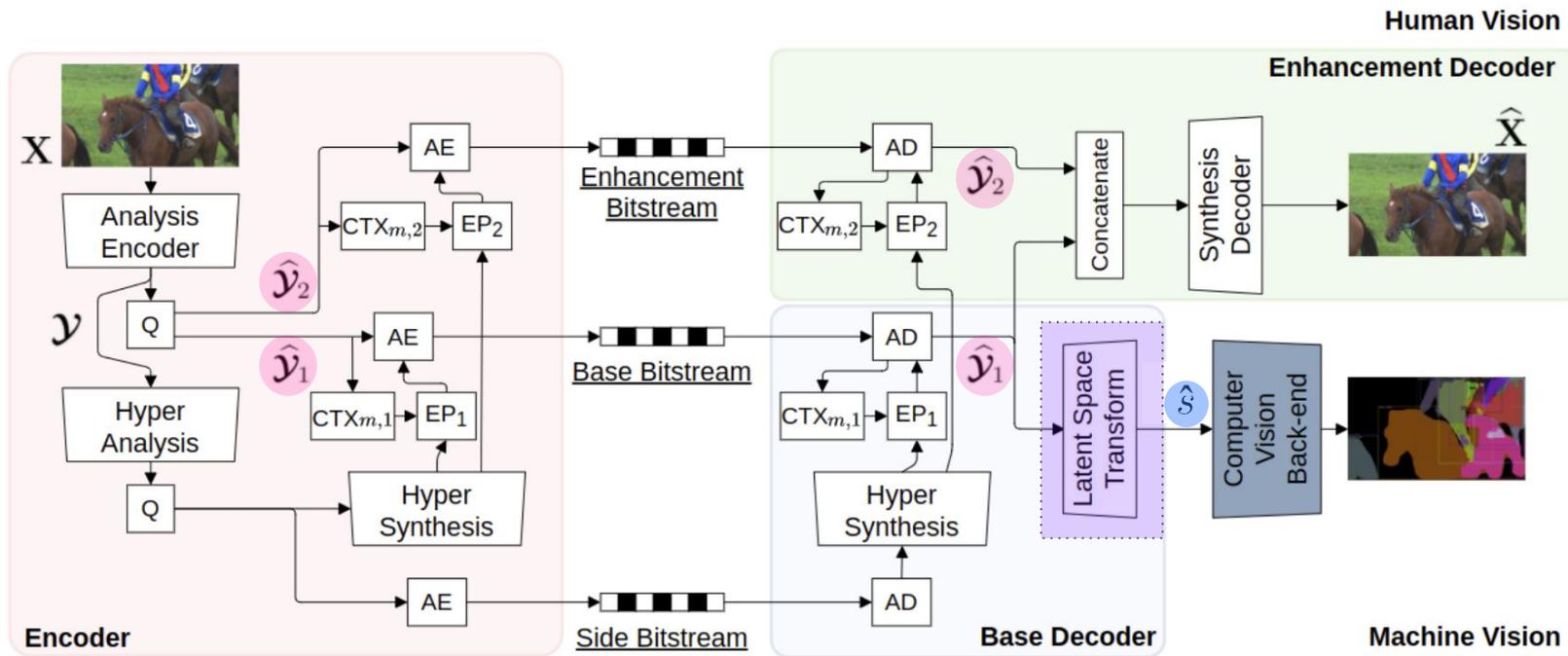


Figure adapted from [H. Choi et al.], "Scalable Image Coding for Humans and Machines," *IEEE Transactions on Image Processing*, 2022.

# Idea: Learned Disentangled Latent Spaces

Motivation is to have little (or none!) excess rate:  $I(\mathbf{y}_1; \mathbf{y}_2) \approx 0$ .

Proposed approach is based on variational inference.

$$\begin{aligned} p_{\theta}(\mathbf{x}, \mathbf{s}, \mathbf{y}_1, \mathbf{y}_2) &= p(\mathbf{y}_1) p(\mathbf{y}_2 \mid \mathbf{y}_1) p_{\theta}(\mathbf{x} \mid \mathbf{y}_1, \mathbf{y}_2) p_{\theta}(\mathbf{s} \mid \mathbf{y}_1, \mathbf{y}_2, \mathbf{x}) \\ &= p(\mathbf{y}_1) p(\mathbf{y}_2) p_{\theta}(\mathbf{x} \mid \mathbf{y}_1, \mathbf{y}_2) p_{\theta}(\mathbf{s} \mid \mathbf{y}_1) \end{aligned}$$

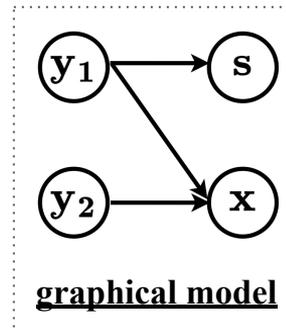
by chain rule

since  $\mathbf{y}_1 \perp\!\!\!\perp \mathbf{y}_2$   
and  $(\mathbf{s} \perp\!\!\!\perp \mathbf{y}_2) \mid \mathbf{y}_1$

The data likelihood is given by integrating:

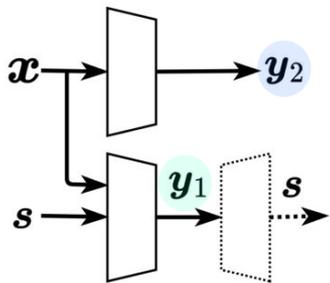
$$p_{\theta}(\mathbf{x}, \mathbf{s}) = \iint p_{\theta}(\mathbf{x}, \mathbf{s}, \mathbf{y}_1, \mathbf{y}_2) d\mathbf{y}_1 d\mathbf{y}_2$$

Unfortunately, intractable!



# Overcoming Intractability

Introduce variational posterior.

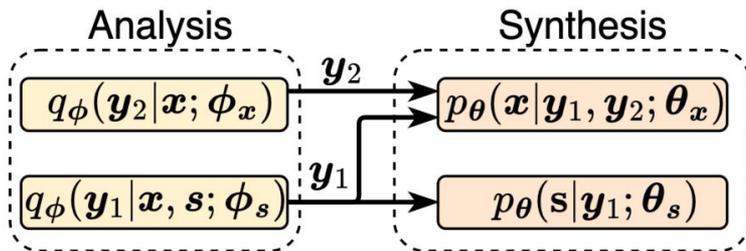


$$q_{\phi}(\mathbf{y}_1, \mathbf{y}_2 \mid \mathbf{x}, \mathbf{s}) = \underbrace{q_{\phi}(\mathbf{y}_1 \mid \mathbf{x}, \mathbf{s})}_{\mathbf{y}_1 \text{ derived from } \mathbf{x}, \mathbf{s}} \underbrace{q_{\phi}(\mathbf{y}_2 \mid \mathbf{x})}_{\mathbf{y}_2 \text{ derived from } \mathbf{x}}$$

Impose above factorization by system model.

Loss function construction turns out to be very similar to Ballé et al. (2018).

We seek to minimize Kullback-Leibler (KL) divergence between  $q_{\phi}, p_{\theta}$ .



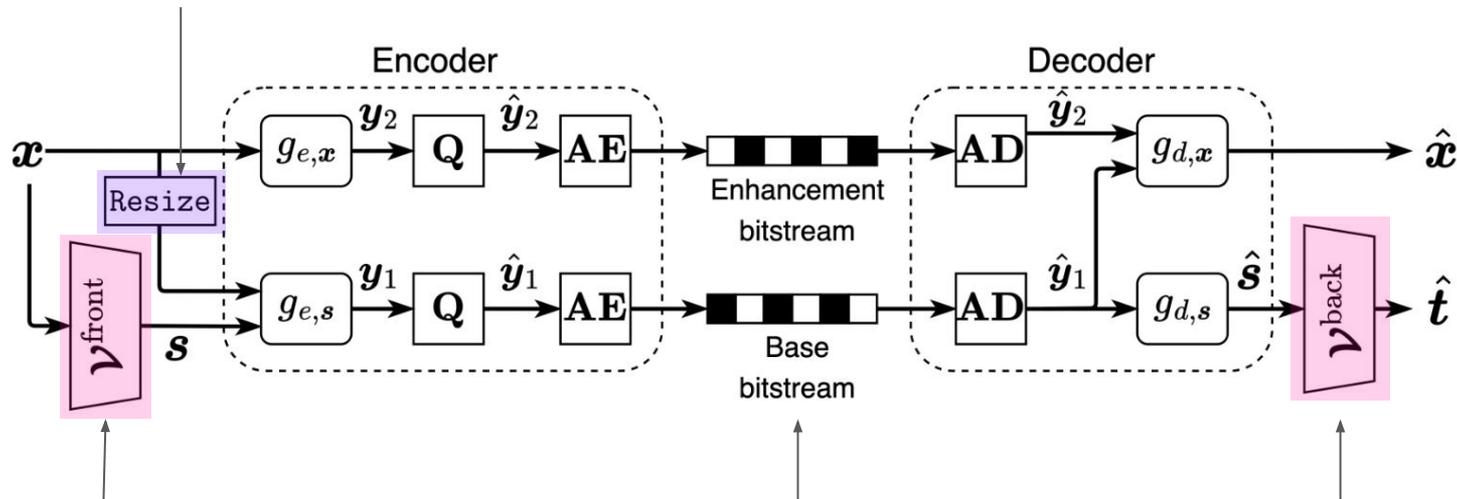
Minimize KL between  $q_\phi, p_\theta$  over dataset of  $\mathbf{x}, \mathbf{s}$ :

$$\begin{aligned}
 \mathcal{L} &= \mathbb{E}_{\mathbf{x}, \mathbf{s} \sim p(\mathbf{x}, \mathbf{s})} \left[ D_{\text{KL}}(q_\phi(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2 | \mathbf{x}, \mathbf{s}) \parallel p_\theta(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2 | \mathbf{x}, \mathbf{s})) \right] \\
 &= \mathbb{E}_{\mathbf{x}, \mathbf{s} \sim p(\mathbf{x}, \mathbf{s})} \mathbb{E}_{\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2 \sim q_\phi} \left[ \left( \overbrace{\log q_\phi(\tilde{\mathbf{y}}_1 | \mathbf{x}, \mathbf{s}; \phi_s)}^0 + \overbrace{\log q_\phi(\tilde{\mathbf{y}}_2 | \mathbf{x}; \phi_x)}^0 \right) \right. \\
 &\quad \left. - \left( \underbrace{\log p_\theta(\mathbf{x} | \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2; \theta_x)}_{D_x} + \underbrace{\log p_\theta(\mathbf{s} | \tilde{\mathbf{y}}_1; \theta_s)}_{D_s} + \underbrace{\log p(\tilde{\mathbf{y}}_1)}_{R_{y_1}} + \underbrace{\log p(\tilde{\mathbf{y}}_2)}_{R_{y_2}} \right) \right] + \text{const.}
 \end{aligned}$$

$\mathcal{L} = R_{y_1} + R_{y_2} + \lambda \cdot D_x + \gamma \cdot D_s$

# Proposed Architecture

“Resize” to match latent dimensions.

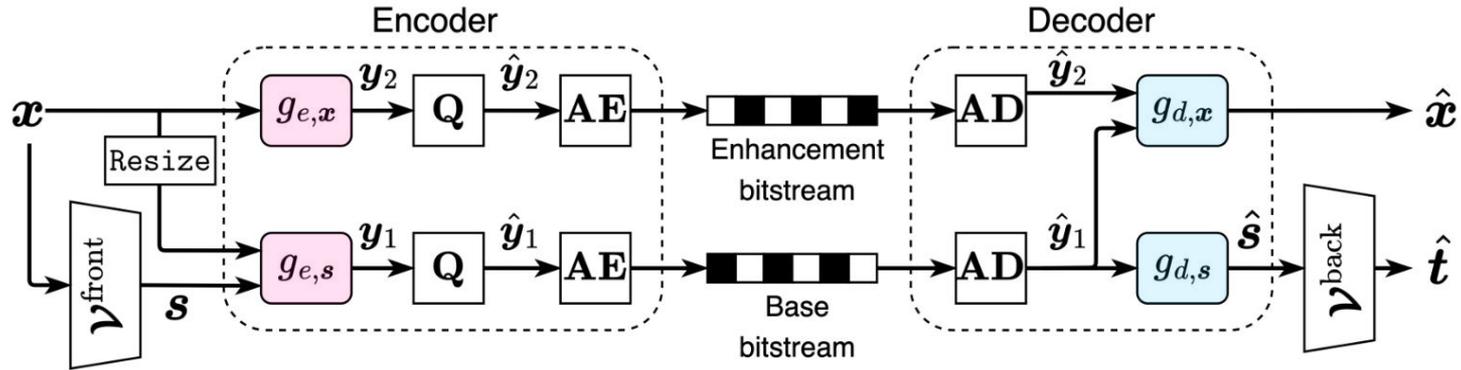


Features are generated from “front” half of task model.

Model is able to create a more task-optimized bitstream.

Features are fed into “back” half of task model.

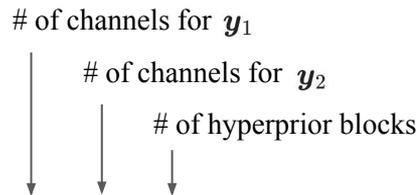
# Proposed Architecture



No.	Encoder				Decoder			
	Layer	In/Out	Layer	In/Out	Layer	In/Out	Layer	In/Out
1	conv5s1	$C_s + 3/N$	conv5s2	$3/N$	deconv5s1	$M_1/N$	deconv5s2	$M/N$
2	conv5s1	$N/N$	conv5s2	$N/N$	deconv5s1	$N/N$	deconv5s2	$N/N$
3	conv5s2	$N/M_1$	conv5s2	$N/N$	deconv5s2	$N/C_s$	deconv5s2	$N/N$
4			conv5s2	$N/M_2$			deconv5s2	$N/3$



# Experimental Setup



- Various architecture configurations for the tuple  $(M_1, M_2, H)$ .
- Train on Vimeo-90K dataset with “distortion” computed using mean-squared error (MSE).
- Evaluate object detection on COCO 2014 validation dataset using mAP (IoU=0.5).
- Evaluate input reconstruction on Kodak dataset using MSE and MS-SSIM.
- Benchmark performance in comparison with:
  - Standard codecs such as HEVC, VVC  $\Rightarrow$  **do not support task-scalability!**
  - Comparative model (*without* PixelCNN-style autoregression) from Choi et al.

[HEVC] [http://hevc.hhi.fraunhofer.de/svn/svn\\_HEVCSoftware/tags/HM-16.20+SCM-8.8/](http://hevc.hhi.fraunhofer.de/svn/svn_HEVCSoftware/tags/HM-16.20+SCM-8.8/)

[VVC] [https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware\\_VTM/-/tags/VTM-12.3/](https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware_VTM/-/tags/VTM-12.3/)

[Vimeo-90K] Xue et al. “Video Enhancement with Task-Oriented Flow,” *IJCV*, 2019.

[COCO 2014] T.-Y. Lin et al., “Microsoft COCO: Common objects in context,” 2014.

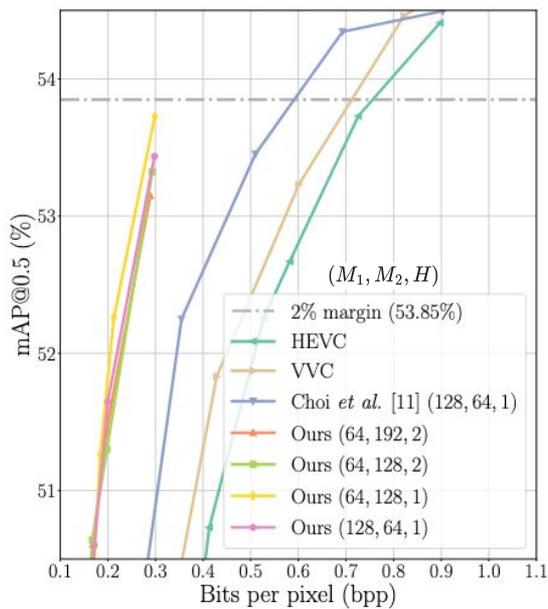
[Kodak] <http://r0k.us/graphics/kodak/>

[MS-SSIM] Z. Wang et al., “Multiscale structural similarity for image quality assessment,” *Asilomar Conf. Signals, Systems, and Computers*, 2003.

[H. Choi et al.] “Scalable Image Coding for Humans and Machines,” *IEEE Transactions on Image Processing*, 2022.

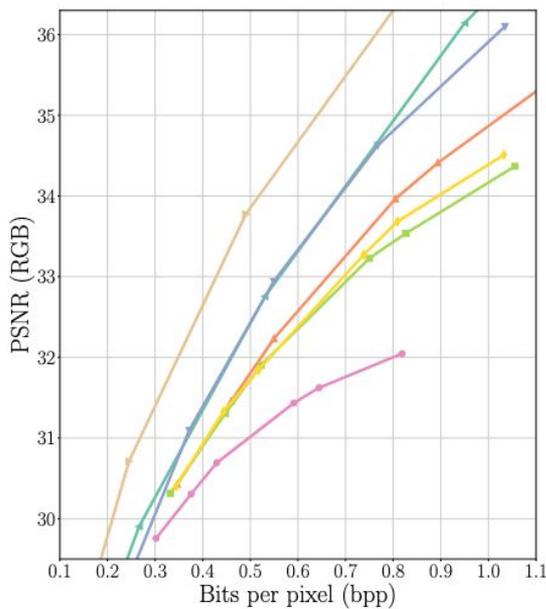
[PixelCNN] Oord et al., “Pixel Recurrent Neural Networks,” *PMLR*, 2016.

# Performance Across Various Metrics



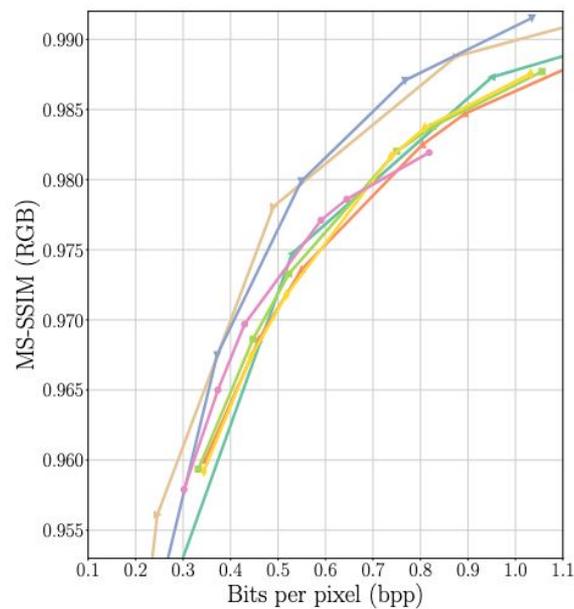
(a)

Task accuracy on COCO 2014 val  
mAP (IoU=0.5) vs. bpp



(b)

Input reconstruction on Kodak  
PSNR vs. bpp



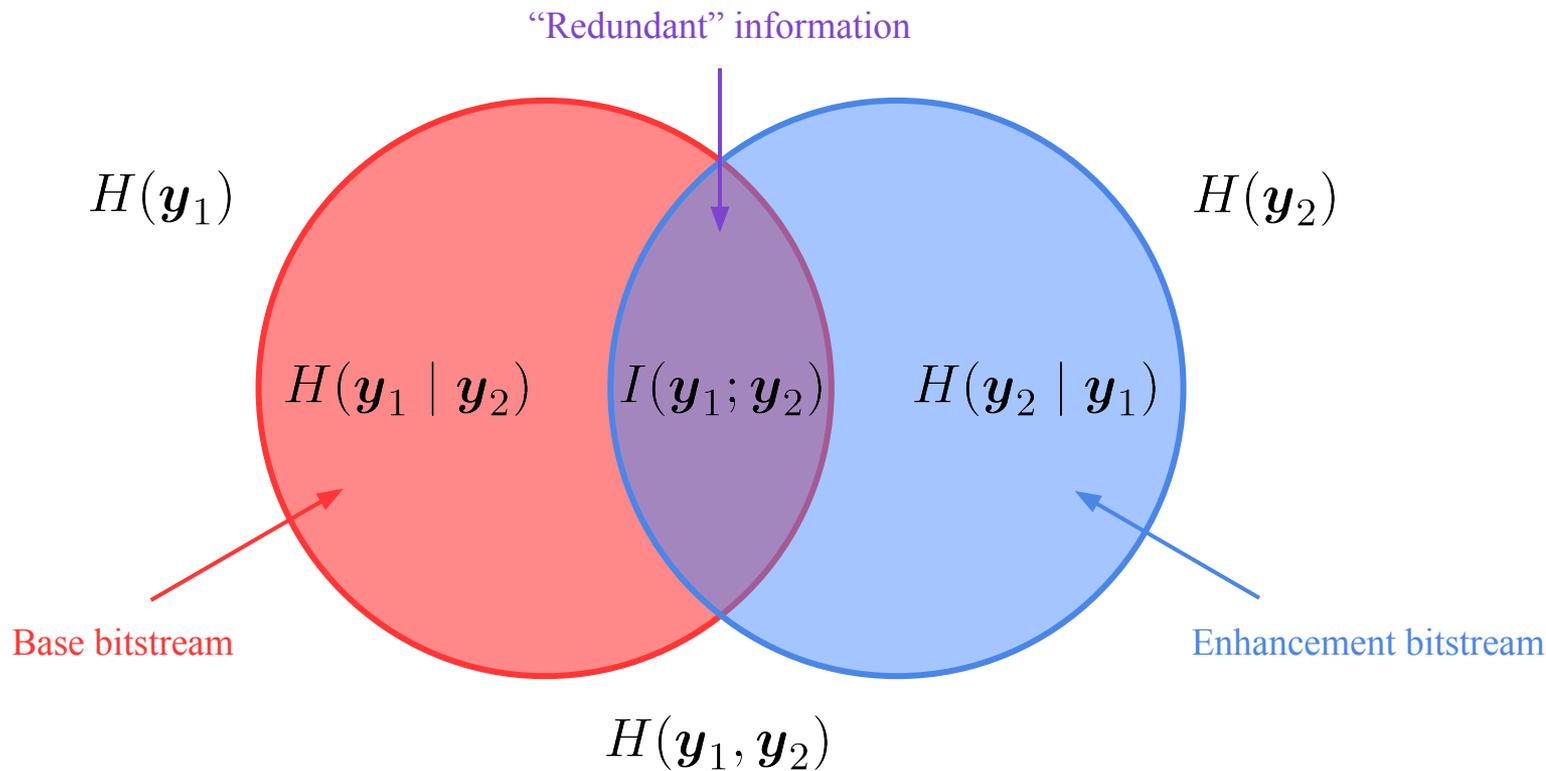
(c)

Input reconstruction on Kodak  
MS-SSIM vs. bpp

Baseline accuracy of YOLOv3 on COCO 2014 val, including JPEG-compressed images, is 55.85% mAP at 4.80 bpp.



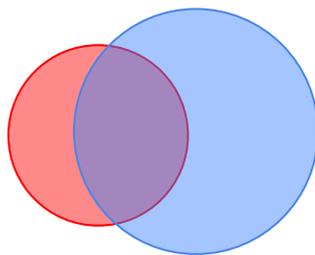
# Quick Recap of Entropy and Mutual Information



# Disentanglement

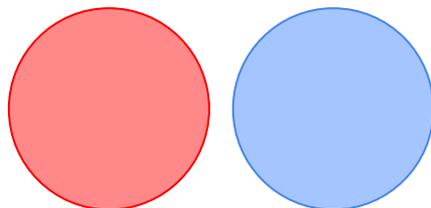


(i)



Redundancy in bitstreams  
Shaded area = 1  
Practical rate cost = 1.5

(ii)



Fully disentangled  
Shaded area = 1  
Practical rate cost = 1

$$\text{Redundancy} \propto I(\mathbf{y}_1; \mathbf{y}_2)$$

$$\text{Shaded area} = H(\mathbf{y}_1, \mathbf{y}_2)$$

$$\text{Practical rate cost} = H(\mathbf{y}_1) + H(\mathbf{y}_2)$$

# Redundancy

Definition: 
$$\text{Rdn}(\mathbf{y}_i | \mathbf{y}_j) \triangleq \frac{I(\mathbf{y}_i; \mathbf{y}_j)}{H(\mathbf{y}_i)} = 1 - \frac{H(\mathbf{y}_i | \mathbf{y}_j)}{H(\mathbf{y}_i)}$$

Codec	Base	Enhancement
	$H(\mathbf{y}_1)$	$H(\mathbf{y}_2)$
Ours	0.3	0.7
Choi et al.	0.6	0.05



$$0 \leq \text{Rdn}(\mathbf{y}_2 | \mathbf{y}_1) \leq 0.4$$



$$0 \leq \text{Rdn}(\mathbf{y}_2 | \mathbf{y}_1) \leq 1.0$$

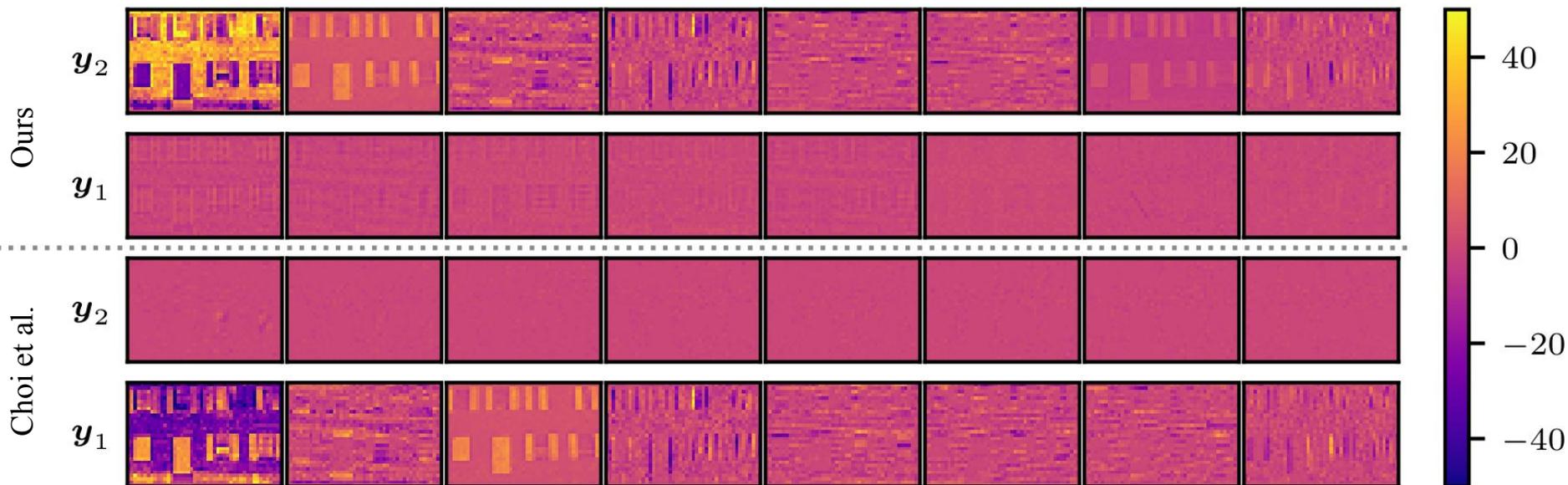
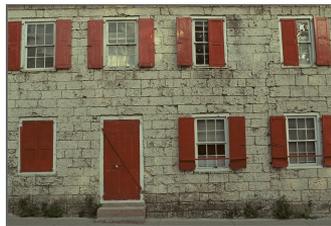


Codec entropy rates (in bits per pixel) measured at 2% loss threshold in mAP.

Bounds on redundancy in enhancement bitstream under respective entropy models.

# Feature maps

Input



top-8 channels ordered by rate

$y_1$  = base (for machine vision)

$y_2$  = enhancement (for humans)



# Conclusion and Future Work

- DNN-based image codec with a new variational formulation.
  - Offers latent-space scalability for human and machine tasks.
  - New way of disentangling the learned latent representations.
- Significant bit reductions at the base layer.
- Needs further investigation about improving reconstruction quality while maintaining the analytics performance.



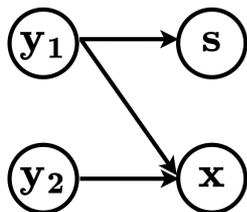
Thank you



# Idea: Learned Disentangled Latent Spaces

Motivation is to have little (or none!) excess rate:  $I(\mathbf{y}_1; \mathbf{y}_2) \approx 0$ .

Proposed approach is based on variational inference.



Introduce variables:  $y_1 \in \mathbb{R}^{d_1}$ ,  $y_2 \in \mathbb{R}^{d_2}$

Under this model, data likelihood is given by:

$$p_{\theta}(x, s, y_1, y_2) = \underbrace{p(y_1)p(y_2)}_{\text{assumption 1: } Y_1 \perp Y_2} p_{\theta}(x | y_1, y_2) \underbrace{p_{\theta}(s | y_1)}_{\text{assumption 2: } Y_2 \perp Y_1 - S}.$$

assumption 1:  $Y_1 \perp Y_2$

assumption 2:  $Y_2 \perp Y_1 - S$

$$\Rightarrow p_{\theta}(x, s) = \iint p_{\theta}(x, s, y_1, y_2) dy_1 dy_2.$$

intractable integral !!



# Overcoming Intractability

Introduce factored variational posterior.

$$q_{\phi}(y_1, y_2 | x, s) = \underbrace{q_{\phi}(y_1 | x, s)}_{\text{red}} \underbrace{q_{\phi}(y_2 | x)}_{\text{blue}}$$

*extract 'common' latents from both (X, S)* *extract 'enhancement' ones only from X*

Impose above factorization by system model.

Loss function construction turns out to be very similar to Ballé et al. (2018).

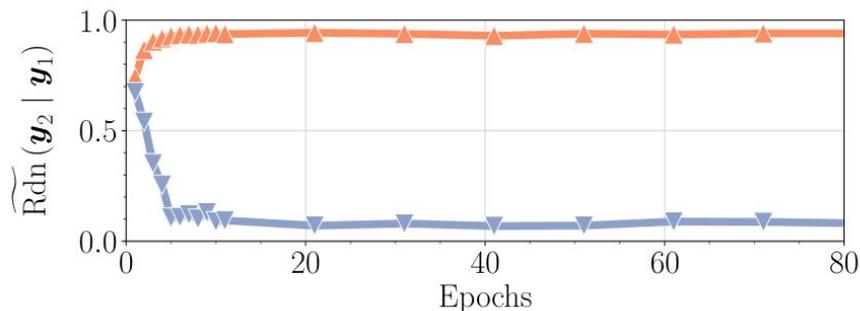
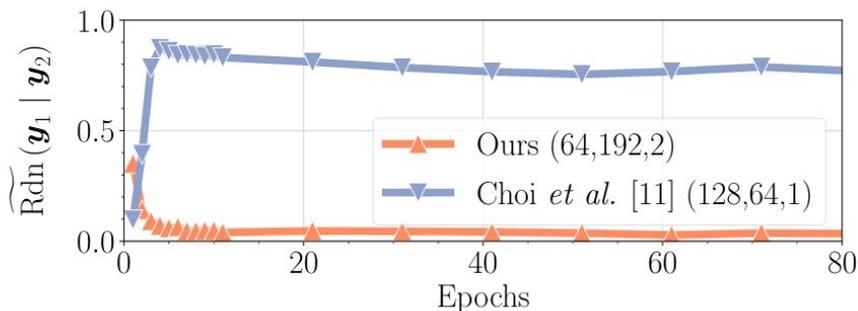
We seek to minimize Kullback-Leibler (KL) Divergence between  $q_{\phi}, p_{\theta}$ .

# Redundancy during training

Definition: 
$$\text{Rdn}(\mathbf{y}_i | \mathbf{y}_j) \triangleq \frac{I(\mathbf{y}_i; \mathbf{y}_j)}{H(\mathbf{y}_i)} = 1 - \frac{H(\mathbf{y}_i | \mathbf{y}_j)}{H(\mathbf{y}_i)}$$

Estimate: 
$$\widetilde{\text{Rdn}}(\mathbf{y}_i | \mathbf{y}_j) = 1 - \frac{1}{H(\mathbf{y}_i)} \sum_{k \in \{1, \dots, K\}} p(k) H(\bar{\mathbf{y}}_i | c(\bar{\mathbf{y}}_j) = k)$$

↑  
Clustering function



Evolution of the redundancy metrics during training.