

# Learned transform compression with optimized entropy encoding

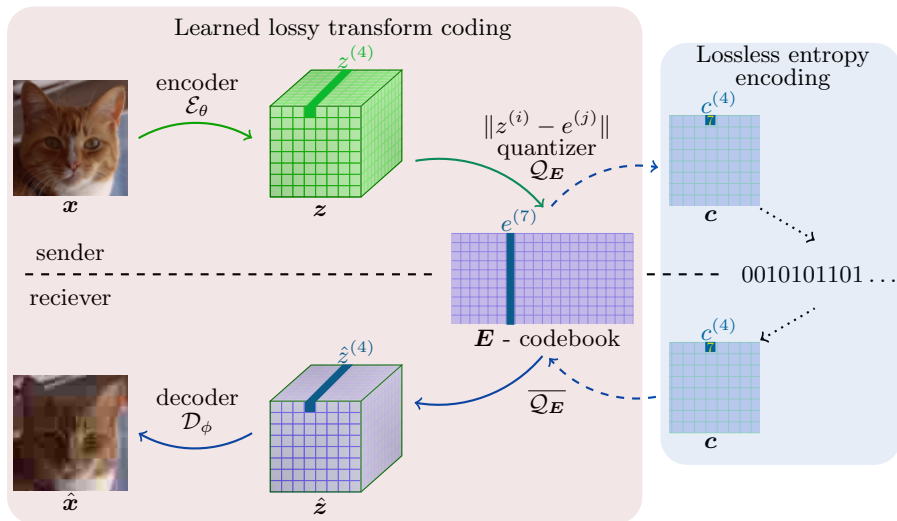
Magda Gregorová

DMML workshop 6 July 2021, Geneva

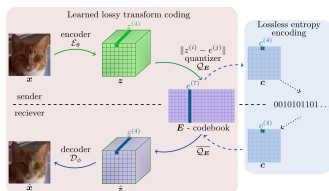
*In collaboration with:*

*Marc Desaulles & Alexandros Kalousis*

# Transform coding with vector quantization



# End-to-end optimized compression



Learn  $\mathcal{E}_\theta, \mathcal{D}_\phi, Q_E$  by minimizing

$$\mathcal{L} := \underbrace{\mathbb{E}_{\mu_x} d(\mathbf{x}, \hat{\mathbf{x}})}_{\text{distortion}} + \lambda \underbrace{\mathbb{E}_{\mu_c} l(\mathbf{c})}_{\text{rate}} \quad (\text{trade-off})$$

data  $\mathbf{x} \sim \mu_x$     symbols  $\mathbf{c} \sim \mu_c$   
(unknown probability measures)

*reconstruction error:*

$$d(\mathbf{x}, \hat{\mathbf{x}}) = \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2$$

(or  $\ell_1$ , MS-SSIM, ...)

*length of binary encoding:*

$$l(c) = -\log p_c(c) \quad (\text{Shannon})$$

$$\text{pmf: } \int_{\mathcal{A}} p_c \, d\# = \sum_{\mathbf{a} \in \mathcal{A}} p_c(\mathbf{a}) = \mu_c(\mathcal{A})$$

$$\text{rate} = \text{entropy:} \quad \mathbb{E}_{\mu_c} l(c) = -\mathbb{E}_{\mu_c} \log p_c(c) = \mathbb{H}_{\mu_c}(c)$$

unknown  $p_c \Rightarrow$  cannot evaluate  $\mathbb{H}_{\mu_c}(c) \Rightarrow$  replace by estimate  $q_c \approx p_c$

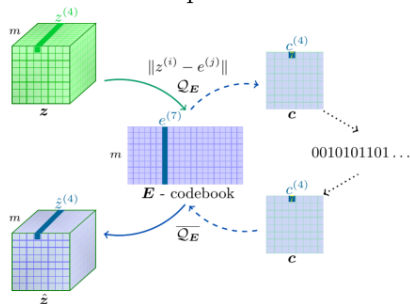
$$\text{rate} \approx \text{cross-entropy:} \quad \mathbb{E}_{\mu_c} l(c) \approx -\mathbb{E}_{\mu_c} \log q_c(c) = \mathbb{H}_{\mu_c | q_c}(c)$$

$\mathcal{E}_\theta, \mathcal{D}_\phi, Q_E, \mathcal{P}_\psi$

$$\mathcal{L} := \underbrace{\mathbb{E}_{\mu_x} d(\mathbf{x}, \hat{\mathbf{x}})}_{\text{distortion}} + \lambda \underbrace{\mathbb{H}_{\mu_c | q_c}(c)}_{\text{rate}}$$

# Vector quantization

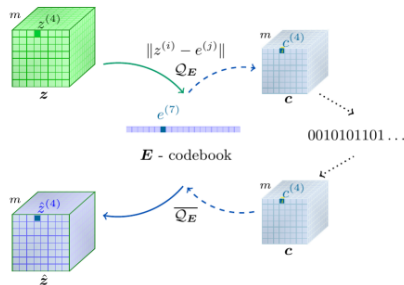
## Vector quantization



message length:  $d^2$

$$Q_E : \quad \hat{z}^{(i)} = \arg \min_{e^{(j)}} \|z^{(i)} - e^{(j)}\| \quad c^{(i)} = \{j : \hat{z}^{(i)} = e^{(j)}\}$$

## Scalar quantization



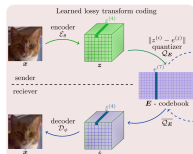
message length:  $d^2 m$

# Model learning - problems

## 1) Non-differentiability of quantization operation

$$\text{Forward: } \mathbf{x} \xrightarrow{\mathcal{E}_\theta} \mathbf{z} \xrightarrow{\mathcal{Q}_E} \hat{\mathbf{z}} \xrightarrow{\mathcal{D}_\phi} \hat{\mathbf{x}} \longrightarrow d(\mathbf{x}, \hat{\mathbf{x}})$$

$$\text{Backward: } \mathbf{x} \xleftarrow{\nabla_\theta} \nabla_{\mathbf{z}} \xleftarrow{\nabla_E} \nabla_{\hat{\mathbf{z}}} \xleftarrow{\nabla_\phi} \nabla_{\hat{\mathbf{x}}} \xleftarrow{\nabla} d(\mathbf{x}, \hat{\mathbf{x}})$$



## 2) Cross-entropy minimization does not minimize rate

$$\mathbb{H}_{\mu_c|q_c}(\mathbf{c}) = -\mathbb{E}_{\mu_c} \log q_c(\mathbf{c}) = \overbrace{D_{\text{KL}}(p_c||q_c)}^{\geq 0} + \mathbb{H}_{\mu_c}(\mathbf{c})$$

$$\min_{q_c} \mathbb{H}_{\mu_c|q_c}(\mathbf{c}) \Leftrightarrow \min_{q_c} D_{\text{KL}}(p_c||q_c)$$

$\mathbb{H}_{\mu_c}(\mathbf{c})$  not function of  $q_c$  so not optimized  $\Rightarrow$  rate not optimized

# Solutions - i), ii), iii)

## i) soft quantization for backward gradients

forward hard (non-differentiable):

$$p_z(\hat{\mathbf{z}} = \mathbf{e}^{(j)} | \mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{e}^{(j)} = \arg \min_{\mathbf{e}^{(i)} \in \mathbf{E}} \|\mathbf{z}(\mathbf{x}) - \mathbf{e}^{(i)}\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{\mathbf{z}}(\mathbf{x}) = \sum_{\mathbf{e}^{(j)} \in \mathbf{E}} p_z(\hat{\mathbf{z}} = \mathbf{e}^{(j)} | \mathbf{x}) \mathbf{e}^{(j)}$$

backward soft (differentiable):

$$\hat{p}_z(\hat{\mathbf{z}} = \mathbf{e}^{(j)} | \mathbf{x}) = \frac{\exp(-\sigma \|\mathbf{z} - \mathbf{e}^{(j)}\|)}{\sum_i^k \exp(-\sigma \|\mathbf{z} - \mathbf{e}^{(i)}\|)}$$

$$\tilde{\mathbf{z}}(\mathbf{x}) = \sum_{\mathbf{e}^{(j)} \in \mathbf{E}} \hat{p}_z(\hat{\mathbf{z}} = \mathbf{e}^{(j)} | \mathbf{x}) \mathbf{e}^{(j)}$$

## Solutions - i), ii), iii)

## ii) pushforward measure

$\mu_x$  - unknown & fixed,  $\mu_c$  pushforward  $f_*(\mu_x)$ ,  $f = \mathcal{Q}_E \circ \mathcal{E}_\theta$  - unknown & not fixed  
change  $\mathcal{E}_\theta$  and  $\mathcal{Q}_E$  to change  $\mu_c$  and hence rate  $\mathbb{H}_{\mu_c}$

## iii) soft cross-entropy

$$\mathbb{H}_{\mu_c|q_c}(c) = -\mathbb{E}_{\mu_c} \log q_c(c) = -\int_{\mathbf{x}} \sum_j p_c(c = j|\mathbf{x}) \log q_c(c = j) d\mu_x$$

hard cross-entropy (no gradients to  $\mathcal{E}_\theta$  and  $\mathcal{Q}_E$ , no rate effect):

$$\mathbb{H}_{\mu_c|q_c}(c) = h(c) \approx -\frac{1}{n} \sum_i^n \sum_j p_c(c = j|\mathbf{x}_i) \log q_c(c = j), \quad p_c(c = j|\mathbf{x}) = p_z(\hat{\mathbf{z}} = \mathbf{e}^{(j)}|\mathbf{x})$$

soft cross-entropy (gradients to  $\mathcal{E}_\theta$  and  $\mathcal{Q}_E$ , rate effect):

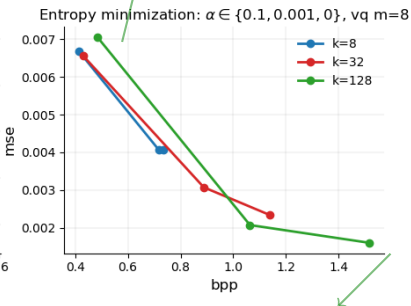
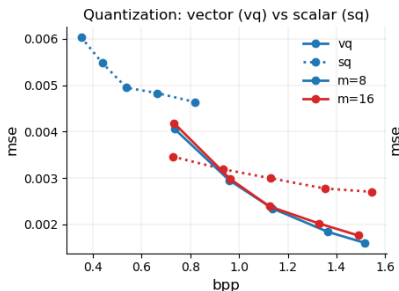
$$\mathbb{H}_{\mu_c|q_c}(c) = s(c) \approx -\frac{1}{n} \sum_i^n \sum_j \hat{p}_c(c = j|\mathbf{x}_i) \log \text{sg} q_c(c = j), \quad \hat{p}_c(c = j|\mathbf{x}) = \hat{p}_z(\hat{\mathbf{z}} = \mathbf{e}^{(j)}|\mathbf{x})$$

$$\mathcal{L}(\mathcal{E}_\theta, \mathcal{D}_\phi, \mathcal{Q}_E, \mathcal{P}_\psi) := \sum_i^n d(\mathbf{x}, \mathcal{D}_\phi[\text{sg}(\hat{\mathbf{z}}_i - \tilde{\mathbf{z}}_i) + \tilde{\mathbf{z}}_i]) + \alpha s(\mathbf{c}_i) + \beta h(\mathbf{c}_i)$$

# Proof of concept - experiments

$\mathcal{E}_\theta, \mathcal{D}_\phi$ : CNN, stride-2 down-/up-sampling, 64 kernels size 3-4, 10 residual blocks with skip connections

$$\mathcal{P}_\psi : q_c(c) = \prod_i^{d^2} q_{c_i}(c_i), \quad q_{c_i} = q_{c_j}$$



ADAM, one cycle cosine schedule  
 $\sigma = 1, \beta = 1$  Imagenet32





## Future work ideas

- i) technical improvements
- ii) instance specific dictionary
- iii) squeeze more from entropy
- iv) BB-ANS for VQVAE
- v) other random ideas

## i) Technical improvements

### *Better probability model:*

Current  $q_c(\mathbf{c}) = \prod_i^d q_i(c_i), q_i = q_j \Rightarrow$  **more complex model e.g. AR or IDF**

Challenge:  $q_c \approx p_c$  not fixed but evolving during training; training stability?

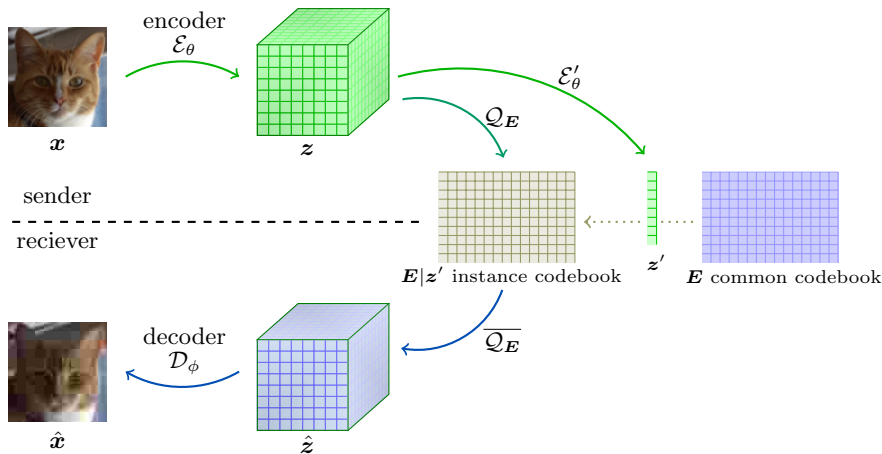
### *Backward gradient info:*

soft relaxation vs streight-through vs. soft relaxation with annealing

### *Initialization of $\mathbf{E}$ :*

random uniform vs k-means++

## ii) Instance specific dictionary



Idea:  $E|z'$  better for specific instance  $x$  than generic  $E$

transmitt:  $c, z' \rightarrow$  trade-off size of  $z'$  vs  $E$

$C_\xi(z', E) = E|z'$ : architecture so that not ignoring  $z'$ , complex vs constrained transformation (e.g. completely free vs only shuffle columns to improve entropy)

### iii) Squeeze more from entropy

link to MAP, ELBO (VAE), more info theory?



## iv) BB-ANS for VQVAE

Townsend (2019) BB-ANS: efficient lossless compression using VAE

open question - need to discretize latent  $\mathbf{z}$  before encoding via ANS

van den Oord (2017) VQVAE: learned discretization via vector quantization

$$p(\mathbf{x}) = \sum_{\mathbf{c}} p(\mathbf{x}|\mathbf{c})p(\mathbf{c}), \quad \mathbf{c} \in \{0, 1, \dots, K\}$$

$$\text{deterministic: } q(\mathbf{c} = k|\mathbf{x}) = \begin{cases} 1 & \text{if } k = \arg \min_i \|\mathbf{z}(\mathbf{x}) - \mathbf{e}^{(i)}\|_2^2 \\ 0 & \text{otherwise} \end{cases}$$

$\Rightarrow D_{\text{KL}}(q(\mathbf{c} = k|\mathbf{x})\|p(\mathbf{c})) = \log K \Rightarrow$  can be dropped from loss  
 streight-through to backprob to  $\mathbf{z} \Rightarrow$  needs k-means loss for  $\mathbf{E}$  updates

fully deterministic scheme not amenable to BB

*Proposed method:*

$$\text{stochastic: } q(\mathbf{c} = k|\mathbf{x}) = \frac{\exp(\|\mathbf{z}(\mathbf{x}) - \mathbf{e}^{(k)}\|_2^2)}{\sum_i^K \exp(\|\mathbf{z}(\mathbf{x}) - \mathbf{e}^{(i)}\|_2^2)}$$

$$c \sim q(\mathbf{c} = k|\mathbf{x}) \quad \text{Gumbel soft-max etc.} \quad D_{\text{KL}}(q(\mathbf{c} = k|\mathbf{x})\|p(\mathbf{c})) = -\mathbb{H}_{q(\mathbf{c}|\mathbf{x})} + \log K$$

stochastic amenable to BB, What does it bring compared to conti latent?

note:  $\min -\mathbb{H}_{q(\mathbf{c}|\mathbf{x})}$  good for BB

## v) Other random ideas

### *Side info:*

Use side info (e.g. ABB meta-data to generate data and compress only the differences

⇒ major patterns covered by generations, compress only the irregularities (surprises)

### *Variable representation power:*

sender has access to full model so can evaluate the transmission error → if too big, compress less and vice versa

⇒ train multiple models (hierarchical, composable) with different rates and apply these selectively to different instances (e.g. driven by  $||\mathbf{z}_i - \mathcal{Q}(\mathbf{z}_i)||$ )

### *Autoregressive dictionary:*

current quantizer  $\mathcal{Q}$  uses single  $\mathbf{E}$  for quantizing all latent vectors  $\mathbf{z}$

⇒ learn  $\mathbf{E}^{(1)}$  to be used for  $\mathbf{z}^{(1)}$  and  $f : (\mathbf{E}^{(i)}, \mathcal{Q}(\mathbf{z}^{(i)}) \rightarrow \mathbf{E}^{(i+1)}$  to be used for  $\mathbf{z}^{(i)}$

## References

- Agustsson, E., Mentzer, F., Tschannen, M., Cavigelli, L., Timofte, R., Benini, L., & Van Gool, L. (2017). “*Soft-to-Hard Vector Quantization for End-to-End Learning Compressible Representations.*” arXiv:1704.00648.
- Balle, J., Laparra, V. & Simoncelli, E. P. (2017). “*End-to-end Optimized Image Compression.*” ICLR.
- Cover, T. M. & Thomas, T. M. (2006). “*Elements of Information Theory.*” Wiley.
- Habibian, A., van Rozendaal, T. Tomczak, J. M., & Cohen, T. S. (2019). “*Video Compression With Rate-Distortion Autoencoders.*” ICCV.
- Mentzer, F., Agustsson, F., Tschannen, M., Timofte, R., Van Gool, L. (2018). “*Conditional Probability Models for Deep Image Compression.*” CVPR.
- Sayood, K. (2012). “*Introduction to Data Compression.*” Elsevier
- Theis, L., Shi, W., Cunningham, A. & Huszar, F. (2017). “*Lossy Image Compression with Compressive Autoencoders.*” ICLR.
- van den Oord, A., Vinyals, O. & Kavukcuoglu, K. (2017). “*Neural Discrete Representation Learning.*” NeurIPS.
- Williams, W., Ringer, S., Ash, T., Hughes, J., MacLeod, D. & Dougherty, J. (2020). “*Hierarchical Quantized Autoencoders.*” NeurIPS.
- Townsend, J., Bird, T. & Barber, D. (2019) “*Practical Lossless Compression with Latent Variables using Bits Back Coding*” ICLR.