# ELIC: Efficient Learned Image Compression with Unevenly Grouped Space-Channel Contextual Adaptive Coding

Dailan He[1*], Ziming Yang[1*], Weikun Peng[1], Rui Ma[1], Hongwei Qin[1], Yan Wang[12†]

SenseTime Research[1], Tsinghua University[2]

{hedailan, yangziming, pengweikun, marui, qinhongwei, wangyan1}@sensetime.com

wangyan@air.tsinghua.edu.cn

## Abstract

*Recently, learned image compression techniques have achieved remarkable performance, even surpassing the best manually designed lossy image coders. They are promising to be large-scale adopted. For the sake of practicality, a thorough investigation of the architecture design of learned image compression, regarding both compression performance and running speed, is essential. In this paper, we first propose uneven channel-conditional adaptive coding, motivated by the observation of energy compaction in learned image compression. Combining the proposed uneven grouping model with existing context models, we obtain a spatial-channel contextual adaptive model to improve the coding performance without damage to running speed. Then we study the structure of the main transform and propose an efficient model, ELIC, to achieve state-of-the-art speed and compression ability. With superior performance, the proposed model also supports extremely fast preview decoding and progressive decoding, which makes the coming application of learning-based image compression more promising.*

## 1. Introduction

In the past years, lossy image compression based on deep learning develops rapidly [4, 5, 15, 20, 22, 24, 29, 39, 40, 49, 50]. They have achieved remarkable progress on improving the rate-distortion performance, with usually much better MS-SSIM [47] than conventional image formats like JPEG [26] and BPG [8], which indicates better subjective quality. Some very recent works [18–20, 22, 49, 50] even outperform the still image coding of VVC [2], one of the best hand-crafted image and video coding standards at present, on both PSNR and MS-SSIM. These results are en-
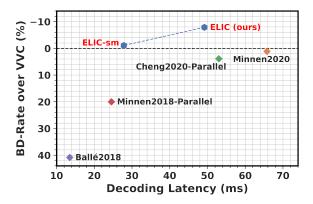


Figure 1. Rate-speed comparison on Kodak. Left-top is better.

couraging, as learned image compression has been proved as a strong candidate for next generation image compression techniques. In the near future, it is quite possible to deploy this line of image compression models in industrial applications. Yet, to make these approaches practical, we must carefully assess the running speed, especially the decoding speed of learned image compression.

One of the most important techniques in learned image compression is the joint backward-and-forward adaptive entropy modeling [15, 20, 22, 29, 39, 40, 48–50]. It helps convert the marginal probability model of coding-symbols to a joint model by introducing extra latent variables as prior [5, 39, 40], leading to less redundancy and lower bitrate. However, the backward-adaptive models along spatial dimension significantly break the parallelism, which inevitably slows down the decoding. To address the issue, He *et al*. [24] proposes to adopt checkerboard convolution as a parallel replacement to the serial autoregressive context model, which has a much better degree of parallelism with constant complexity. Minnen *et al*. [40] proposes to adopt a context model along channel dimension instead of the serial-decoded spatial ones, which also improves the parallelism. However, to achieve a non-trivial bit-saving with this channel-conditional model, the symbols are divided to

---

10 groups and coded progressively, which still slows down the overall inference. It is promising to delve into parallel multi-dimension contextual adaptive coding by combing these two models to achieve better coding ability [24], constituting one of the motivation of our work. In this paper, we investigate an uneven grouping scheme to speed up the channel-conditional method, and further combine it with a parallel spatial context model, to promote RD performance while keeping a fast running speed.

More and more complex transform networks also slow down the inference. As learned image compression is formulated as a sort of nonlinear transform coding [3, 21], another plot improving coding performance is the development of main transform. Prior works introduce larger networks [15, 20, 22, 32], attention modules [15, 22, 33, 35] or invertible structures [37, 50] to main analysis and synthesis networks. These heavy structures significantly improve the RD performance but hurt the speed. We notice that, with a relative strong and fast adaptive entropy estimation (*i.e.* the above mentioned adaptive coding approaches with hyperprior and context model), we can re-balance the computation between main transform and entropy estimation, to obtain low-latency compression models. This further motivates us to promote the contextual modeling technique.

Learned image compression is growing mature and tends to be widely used, but its lack of efficiency is still a critical issue. In this paper, we contribute to this field from following perspectives:

- We introduce information compaction property as an inductive bias to promote the expensive channel-conditional backward-adaptive entropy model. Combining it with the spatial context model, we propose a multi-dimension entropy estimation model named Space-Channel ConTeXt (SCCTX) model, which is fast and effective on reducing the bit-rate.

- With the proposed SCCTX model, we further propose ELIC model. The model adopts stacked residual blocks as nonlinear transform, instead of GDN layers [4]. It surpasses VVC on both PSNR and MS-SSIM, achieving state-of-the-art performance regarding coding performance and running speed (Figure 1 and Table 2).

- We propose an efficient approach to generate preview images from the compressed representation. To our knowledge, this is the first literature addressing the very-fast preview issue of learned image compression.

## 2. Related works

### 2.1. Learned lossy image compression

Learned lossy image compression [5, 15, 22, 24, 35, 39] aims at establishing a data-driven rate-distortion optimiza-

tion (RDO) approach. Given input image $\boldsymbol{x}$ and a pair of neural analyzer $g_a$ and neural synthesizer $g_s$, this learning-based RDO is formulated as:

$$\mathcal{L} = R(\hat{\boldsymbol{y}}) + \lambda D(\boldsymbol{x}, g_s(\hat{\boldsymbol{y}})) \tag{1}$$

where $\hat{\boldsymbol{y}} = \lceil g_a(\boldsymbol{x}) \rfloor$ represents the discrete coding-symbols to be saved and $\lceil \cdot \rfloor$ is the quantization operator. Balancing the estimated bit-rate $R$ and image reconstruction distortion $D$ with a rate-controlling hyper-parameter $\lambda$, we can train a set of neural networks $g_a, g_s$ to obtain various pairs of image en/de-coding models, producing a rate-distortion curve.

Ballé *et al.* [4] proposes to adopt a uniform noise estimator and a parametric entropy model to approximate the probability mass function $p_{\hat{\boldsymbol{y}}}$, so that its expected negative entropy $-\mathbb{E}[\log p_{\hat{\boldsymbol{y}}}(\hat{\boldsymbol{y}})]$ can be supervised as the $R(\hat{\boldsymbol{y}})$ term in eq. 1 in a differentiable manner with gradient-decent-based optimization. Later, the entropy model is further extended to a conditioned Gaussian form [5, 39]:

$$p_{\hat{\boldsymbol{y}}|\hat{\boldsymbol{z}}}(\hat{\boldsymbol{y}}|\hat{\boldsymbol{z}}) = \left[ \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2) * U(-0.5, 0.5) \right](\hat{\boldsymbol{y}}) \tag{2}$$

where the entropy parameters $\boldsymbol{\Theta} = (\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$ are calculated from extra computed or stored prior. Ballé *et al.* [5] adopts hyperprior $\hat{\boldsymbol{z}}$ to calculate the entropy parameters. $\hat{\boldsymbol{z}}$ is calculated from unquantized symbols $\boldsymbol{y}$ with a hyper analyzer $h_a$. It can be seen as side-information introduced to the neural coder, acting as the forward-adaptive method.

To painlessly improve the coding efficiency, several training, inference, and encoding-time optimizing approaches are proposed [23, 52, 54]. They can improve the RD performance without slow down the decoding, and can be used together with various coding architectures.

### 2.2. Backward-adaptive entropy models

Backward-adaptive coding is also introduced to learned image compression, including spatial context models [22, 24, 29, 39] and channel-conditional models [40]. Correlating current decoding symbols with already decoded symbols, this sort of approaches further save the bits.

A spatial context model refers to observable neighbors of each symbol vector $\hat{\boldsymbol{y}}_i$ at the $i$-th location:

$$\hat{\boldsymbol{y}}_{<i} = \{\hat{\boldsymbol{y}}_1, \ldots, \hat{\boldsymbol{y}}_{i-1}\} \tag{3}$$
$$\Phi_{\mathrm{sp},i} = g_{\mathrm{sp}}(\hat{\boldsymbol{y}}_{<i}) \tag{4}$$

where the context model $g_{\mathrm{sp}}(\cdot)$ is an autoregressive convolution [39, 42]. Each context representation $\Phi_{\mathrm{sp},i}$ is used to jointly predict entropy parameters accompanied by the hyperprior $\hat{\boldsymbol{z}}$. This approach demands symbol vectors $\hat{\boldsymbol{y}}_1, \ldots, \hat{\boldsymbol{y}}_{HW}$ to be decoded serially, which critically slows down the decoding [24, 35, 39]. He *et al.* [24] proposes to separate the symbols into anchors and non-anchors:

$$\hat{\boldsymbol{y}}_{<i}^{(\mathrm{anchor})} = \varnothing, \quad \hat{\boldsymbol{y}}_{<i}^{(\mathrm{nonanc})} = \hat{\boldsymbol{y}}^{(\mathrm{anchor})} \tag{5}$$

Figure 2. Visualization of sorted channels with the top-10 largest average energy. Left 1: the original image (`kodim19.png`). Lighter regions correspond to larger symbol magnitudes $|\hat{y}|$. It can been seen from the figures that most strong activation concentrates in the first channel (left 2) and the remained channels become sparse gradually.
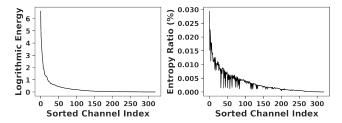


Figure 3. Energy and entropy distribution along channels. The results are evaluated with Minnen2018 [39] model, on Kodak. The channels are sorted by energy averaged over all 24 images.



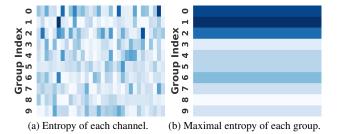(a) Entropy of each channel.  (b) Maximal entropy of each group.

Figure 4. A case study of adopting the 10-slice channel-conditional adaptive coding [40]. Deeper colors denote larger values. The entropy of each channel group is implicitly sorted. The beginning groups contain channels with the largest entropy. The results are from evenly grouped model trained with $\lambda = 0.045$, evaluated on `kodim08` from Kodak.

and adopts a checkerboard convolution as $g_{\mathrm{sp}}(\cdot)$, so the decoding of both anchors and non-anchors can be in parallel.

Another scheme to perform parallel backward adaption is to reduce the redundancy among channels. Minnen *et al.* [40] proposes to group the symbol channels to $K$ chunks as the channel-wise context:

$$\Phi_{\mathrm{ch}}^{(k)} = g_{\mathrm{ch}}^{(k)}(\hat{\boldsymbol{y}}^{<k}), \quad k = 2, \ldots, K \qquad (6)$$

where $\hat{\boldsymbol{y}}^{<k} = \{\hat{\boldsymbol{y}}^{(1)}, \ldots, \hat{\boldsymbol{y}}^{(k-1)}\}$ denotes already decoded channel chunks. Setting a proper $K$ is essential for this method to balance the compression performance and running speed. The larger $K$ is, the better the RD performance is [40], yet the slower the parameter estimation is (as the degree of parallelism decreases).

Multi-dimension adaptive coding approaches have been proposed but all of them still suffer from the slow-speed issue, to our knowledge. Liu *et al.* [33] proposes a 3D context model which performs a 3D convolution passing by all the channels. Li *et al.* [31] proposes a multi-dimension context model with non-constant complexity. Guo *et al.* [22] uses 2-chunk contextual modeling with serial global-and-local adaptive coding. Ma *et al.* [36] propose a cross-channel context model which uses a even denser referring scheme than 3D context models.

## 3. Parallel multi-dimension context modeling

### 3.1. Information compaction property

Energy compaction is an essential property of transform coding [21], *e.g.* DCT-based JPEG. With decomposed coef-

ficients extremely concentrated on lower frequencies, describing most structural and semantic information of the original image, higher frequencies can be compressed more heavily by using larger quantization steps to achieve a better rate-distortion trade-off.

We find this compaction also exists in learned analysis transform. We visualize each latent feature map $\hat{\boldsymbol{y}}^{(\ell)}$ of the mean-scale joint model, Minnen2018 [39], as its scaled magnitude (Figure 2) and draw the distribution of energy and entropy along channels (Figure 3). More strongly activated, the beginning channels have much larger average energy. Since the entropy distribution correlates to the energy distribution, it indicates an information compaction property. This phenomenon exists in all of the 5 models we test: Ballé2018 [5], Minnen2018 [39], Cheng2020 [15] and their parallel versions [24]. The information compaction in those models is orderless because the supervision conducted on the analyzer output channels is symmetric.

When adopting a channel-conditional approach, this property induces a group-level order. See Figure 4. Particular channels in the earlier encoded groups have much larger entropy, so they are allocated more bits. As the beginning channels are more frequently referred to by following channels, the major information implicitly concentrates on the beginning channels to help eliminate more channel-wise redundancy. The progressive coding results [40] also experimentally prove this, since we can reconstruct the main
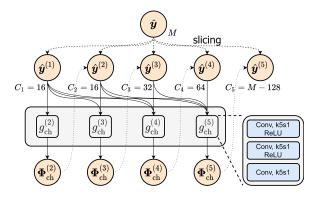
3

Figure 5. Proposed uneven grouping for channel-conditional (CC) adaptive coding. The $M$-channel coding-symbols $\hat{y}$ are grouped into 5 chunks with gradually increased number of channels $C_k$.

semantics of images only from the beginning channels.

We tend to understand this information compaction property of learned image compression from the perspective of sparse representation learning [41, 44, 45], yet the theoretical analysis and explanation of it are beyond the topic of this paper. We view it as a strong prior knowledge that helps us introduce inductive bias to improve the model design.

### 3.2. Unevenly grouped channel-wise context model

The visualization in Figure 4 shows that the later encoded channels contain less information, and they are less frequently used to predict following groups. Therefore, we can reduce the cross-group reference to speed up, by merging more later encoded channels into larger chunks. On the other hand, because of the information compaction, with less channel, the earlier encoded channel groups may still well help reduce the entropy of the following channels. Thus, a more elaborate channel grouping scheme may further improve this entropy estimation module by rebalancing the channel numbers of different groups. Yet, existing approaches often simply group the channels to chunks with the same size [22, 40] or adopt per-channel grouping [33, 36].

We propose an uneven grouping scheme, allocating finer granularity to the beginning chunks by using fewer channels and grow coarser gradually for the following chunks by using more channels. Thus, for symbols $\hat{y}$ with $M$ channels, we split them along the channel dimension to 5 chunks $\hat{y}^{(1)}, \ldots, \hat{y}^{(5)}$ with $16, 16, 32, 64, M-128$ channels respectively. Figure 5 shows the channel-conditional model with this uneven allocation. It only requires 5 times of parallel calculation to decode all the slices $\hat{y}^{(k)}$, which saves the running time.

### 3.3. SCCTX: space-channel context model

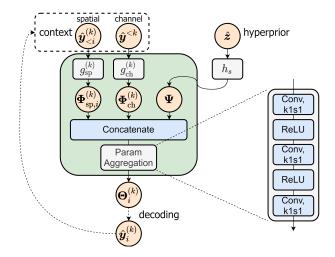Spatial context model and channel-conditional model eliminate redundancy along spatial and channel axes. As



Figure 6. Diagram of proposed space-channel context model.

| Method | Type | $G$ | $S$ ($M = 320$) |
|---|---|---|---|
| autoregressive [39] | sp. | $HW$ | 320 |
| checkerboard [24] | sp. | 2 | $160HW$ |
| 10-slice CC. [40] | ch. | 10 | $32HW$ |
| uneven CC. (ours) | ch. | 5 | $(16 \rightarrow 192)HW$ |
| SCCTX (ours) | mix | 10 | $(8 \rightarrow 96)HW$ |

Table 1. Comparison of different backward-adaptive coding approaches. $G$ denotes the number of groups, and $S$ is the group size. Symbols in the same group can be processed in parallel. The sizes are calculated for $M \times H \times W$ symbols where $M = 320$.

those dimensions are orthogonal, we assume the redundancy in those dimensions is orthogonal too. Thus, we combine the two models for a better backward-adaptive coding.

Figure 6 shows our space-channel context model (SCCTX). In the $k$-th unevenly grouped chunk, we apply a spatial context model $g_{\rm sp}^{(k)}$ to recognise spatial redundancy (eq. 4). It could be an autoregressive convolution [39] or its two-pass parallel adaption [24]. We introduce $g_{\rm ch}$ networks to model the channel-wise context $\Phi_{\rm ch}^{(k)}$ (Figure 5 and eq. 6). The output of spatial and channel branches at the $(k, i)$-th location, $\Phi_{{\rm sp},i}^{(k)}$ and $\Phi_{\rm ch}^{(k)}$, will be concatenated with hyperprior representation $\Psi$ and fed into a location-wise aggregation network to predict the entropy parameters $\Theta_i^{(k)} = (\mu_i^{(k)}, \sigma_i^{(k)})$ for the following en/de-coding of $\hat{y}_i^{(k)}$. Then the just obtained $\hat{y}_i^{(k)}$ will be used as context to compute $\Phi_{{\rm sp},(i+1)}^{(k)}$ or $\Phi_{\rm ch}^{(k+1)}$, till en/decoding the entire $\hat{y}$.

As shown in Table 1, by default we use the parallel checkerboard [24] model as the spatial context for SCCTX, which is only applied inside each channel chunk to be more efficient.
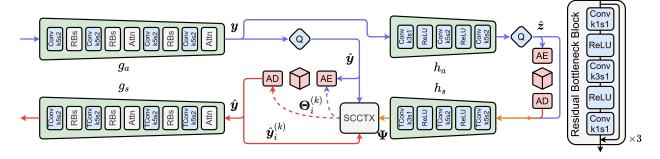
Figure 7. The overall architecture of ELIC. *RBs* denote stacks of residual bottleneck blocks [25] shown in the right. *Attn* blocks are attention modules proposed by Cheng *et al.* [15]. *AE* and *AD* are arithmetic en/de-coder, respectively. *TConv* denotes transposed convolution. The blue and red arrows denote the encoding and decoding data flow. The orange ones are shared by both encoding and decoding.

# 4. ELIC: efficient learned image compression with scalable residual nonlinearity

## 4.1. Stacking residual blocks for nonlinearity

For a long period, generalized divisive normalization (GDN) is one of the most frequently used techniques in learned image compression [5, 15, 22, 29, 50]. It introduces point-wise nonlinearity to the model [6], which aggregates information along the channel axis and scales feature vectors at each location. Different from linear affine normalization techniques like batch normalization or layer normalization, this nonlinear GDN performs more similarly to point-wise attention mechanism. Thus, we propose to investigate other nonlinear transform layers as alternatives of GDN. Note that this is different from existing investigations that view GDN as activation function [14, 15, 17, 27].

We replace the GDN/IGDN layers with stacks of residual bottleneck blocks [25]. Earlier works also investigate pure convolution networks for visual compression [11–14] and here we revisit it from the layer-level perspective. We observe that the performance further improves when the number of stacked blocks increases, because of the enhanced nonlinearity. Thus, a network with strong enough nonlinearity can express the intermediate features better for rate-distortion trade-off, even without GDN layers. A similar structure is also investigated by Chen *et al.* [11], as a part of non-local attention module. We experimentally prove that, even without attention mechanism, the RD performance can still be improved by simply stacking the residual blocks.

Stacking residual blocks allows us to better conduct scalable model profiling. It is also expected to benefit from modern training and boosting techniques like network architecture search [34, 53] and loss function search [30, 46], though they are out of the bound of this work. Residual block is also easier to be extended for dynamic or slimmable inference, while GDN requires special handling [51].
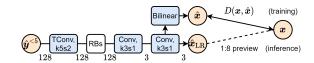


Figure 8. Structure of the proposed thumbnail synthesizer. *Bilinear* module denotes three bilinear upsampling layers by factor 2.

## 4.2. Architecture of ELIC

We summarize the above-mentioned techniques in our proposed model named ELIC (see Figure 7), to build an efficient learned image compression model with strong compression performance. We use the proposed SCCTX module to predict the entropy parameters $\Theta = (\mu, \sigma)$ of a mean-scale Gaussian entropy model [39]. Using this more powerful backward-adaptive coding allows adopting lighter transform networks compared with recent prior works [15, 20, 22, 50]. In the main transforms, we simply adopt stride-2 convolutions to de/in-crease the feature map sizes, following earlier settings [5, 39].

# 5. Quickly decoding thumbnail-preview

The major bottleneck of decoding process is the synthesis inference, which runs a heavy network to reconstruct the full-resolution image. When applying learned image compression, however, we do not always want to decode the full-resolution image. For instance, when looking through images saved on a server, we require the decoder to quickly generate thousands of thumbnail-preview images which have much lower resolution but keep the structure and semantics of the original images. Another case is progressive decoding preview, investigated in prior literature [40]. When the image is progressively and partially decoded step by step, frequently invoking the heavy synthesizer will critically slow down the overall decoding process.

5

On these occasions, image quality is far less important than decoding speed. Thus, directly decoding the full-resolution images using the heavy synthesizer is impractical.

We propose to train an additional tiny network, called thumbnail synthesizer, to reconstruct low-resolution images as thumbnail-preview. When adopting SCCTX, most semantic information is compacted in the earlier decoded channels. Hence we propose to generate the preview image only from the first 4 chunks (*i.e.* the first 128 channels).

The structure of our proposed thumbnail synthesizer is shown in Figure 8. After training the main models, we freeze all the learned parameters and change the main synthesizer to initialized thumbnail synthesizer. Then we restart the distortion optimization to train the model.

As the proposed thumbnail synthesizer is extremely light, its decoding only requires a few microseconds (w.r.t. $768 \times 512$ images). Compared with obtaining the preview image by down-sampling from the entirely reconstructed full-resolution image, using the proposed model to obtain the thumbnail-preview images is much more efficient.

## 6. Experiments

### 6.1. Settings

We train the models on the largest 8000 images picked from ImageNet [16] dataset. A noising-downsampling pre-processing is conducted following prior works [5, 24]. We use Kodak [28] and CLIC Professional [1] for evaluation.

The training settings are accordingly sketched from existing literature [5, 15, 24, 39]. For each architecture, we train models with various $\lambda$ standing for different quality presets. We set $\lambda = \{4, 8, 16, 32, 75, 150, 300, 450\} \times 10^{-4}$ for each model when optimizing MSE. Empirically, on Kodak, models trained with these $\lambda$ values will achieve average bits-per-pixel (BPP) ranging from 0.04 to 1.0. We set the number of channels $N = 192$ and $M = 320$ for all models. We train each model with an Adam optimizer with $\beta_1 = 0.9, \beta_2 = 0.999$. We set initial learning-rate to $10^{-4}$, batch-size to 16 and train each model for 2000 epochs (1M iterations, for ablation studies) or 4000 epochs (2M iterations, to finetune the reported ELIC models), and then decay the learning-rate to $10^{-5}$ for another 100-epoch training.

We use mixed quantization estimator to train channel-conditional models, following Minnen *et al.* [40]. It helps optimize the single Gaussian mean-scale entropy model [39], making it comparable with mixture models like GMM [15, 22, 29, 35]. Following prior works [24, 39] and community discussion[1], we encode each $\lceil y - \mu \rfloor$ to the bitstream instead of $\lceil y \rfloor$ and restore the coding-symbol as $\lceil y - \mu \rfloor + \mu$, which can further benefit the single Gaussian entropy model. When using the autoregressive con-
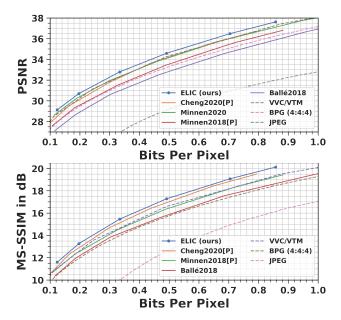
---

---



Figure 9. Rate-distortion curves of various image compression approaches. The results are evaluated on Kodak. All shown learned models are optimized for minimizing MSE.

---

text model with the two above-mentioned strategies together, the one-pass encoding [24] doesn't simply work for training because of the inconsistency in synthesizer input ($\lceil y - \mu \rfloor + \mu$ instead of $\lceil y \rfloor$). Therefore, we adopt a two-stage training. In the beginning 2000/3800 epochs, we sample uniform noise to estimate $\hat{y}$ input to SCCTX and use $\text{STE}(y)$ as the input of synthesizer. Thus, the faster one-pass encoding can be adopted. Then, we fed $\text{STE}(y-\mu)+\mu$ to both the synthesizer and SCCTX in the remained epochs.

To obtain a lightweight model, we do not adopt the LRP network [40], and the low-rate models ($\lambda \leq 0.0075$) perform worse if trained with target $\lambda$ values from scratch. Inspired by prior work [40], we train these models with $\lambda = 0.015$ at the beginning, and adapt them using target $\lambda$ values after decaying the learning-rate.

Please refer to the supplementary material for more detailed experiment settings, including running environment and raw data accessing.

### 6.2. Quantitative results

To compare the RD performance of various models, we present BD-rate [9] computed from PSNR-BPP curves as the quantitative metric. The anchor RD performance is set to VVC or BPG accordingly (with an anchor BD-rate $0\%$). In the supplementary material we also show BD-rate results calculated on MS-SSIM.

| Model | Sp. Context | BD-Rate (%) | Inference Latency (ms) | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | Enc. Tot. | Dec. Tot. | YEnc. | YDec. | ZEnc. | ZDec. | Param. |
| ELIC (ours) | [P] | **-7.88** | 42.44 | 49.16 | 31.16 | 38.33 | 1.18 | 1.59 | 10.10 |
| ELIC-sm (ours) | [P] | -1.07 | 22.64 | 27.80 | 12.55 | 17.37 | 1.07 | 1.41 | 9.02 |
| Minnen2020 [40] | - | 1.11 | 60.76 | 65.75 | 10.60 | 15.41 | 1.36 | 1.54 | 48.80 |
| Cheng2020[P] [24] | [P] | 3.89 | 45.43 | 53.01 | 36.84 | 44.20 | 0.93 | 1.15 | 7.66 |
| Minnen2018[P] [24] | [P] | 20.00 | 18.84 | 24.52 | 11.63 | 17.32 | 1.47 | 1.46 | 5.74 |
| Ballé2018 [5] | - | 40.85 | 13.07 | 17.24 | 11.90 | 16.14 | 1.17 | 1.10 | - |
| Gao2021 [20] | [S] | **-10.94** | / | / | / | / | / | / | / |
| Guo2021 [22] | [S] | -7.02 | 190.41 | $> 10^3$ | 154.08 | 431.68 | 3.67 | 1.61 | $> 10^3$ |
| Xie2021 [50] | [S] | -0.54 | 61.20 | $> 10^3$ | 57.96 | 162.18 | 1.28 | 1.34 | $> 10^3$ |
| Cheng2020 [15] | [S] | 3.35 | 41.42 | $> 10^3$ | 36.99 | 47.35 | 0.96 | 1.12 | $> 10^3$ |
| Minnen2018 [39] | [S] | 14.92 | 18.14 | $> 10^3$ | 12.75 | 17.34 | 1.69 | 1.60 | $> 10^3$ |
| VVC (YUV 444) [2] | - | 0.00 | - | - | - | - | - | - | - |
| BPG [8] | - | 27.21 | - | - | - | - | - | - | - |
| JPEG [26] | - | 266.32 | - | - | - | - | - | - | - |

Table 2. RD and inference time of learned image compression models. The BD-Rate data is calculated relative to VVC (YUV 444) from PSNR-BPP curve on Kodak. The **Param.** column denotes inference latency of entropy parameter calculation during decoding (*e.g.* the running time of SCCTX in our ELIC model). **Enc. Tot.** and **Dec. Tot.** denote total network latency for encoding and decoding respectively. Models marked with [P] adopt parallel checkerboard context model and [S] denotes serial context model. We have not reproduced the very recent Gao2021 model, yet its speed reported by the authors [20] is slower than the referred Cheng2020 model (more than 1s/image).

### 6.2.1 Comparison of coding performance and speed

We compare the RD performance and coding speed of ELIC with exiting learning based approaches, as shown in Table 2. Apparently, our ELIC achieves remarkable improvement on compression quality and speed among learned image compression approaches. We also report BD-rates of manual-designed image compression approaches, but do not compare the speed of them because they cannot run on GPU. ELIC outperforms VVC (in YUV 4:4:4 colorspace) on RD performance regarding PSNR (and also MS-SSIM, see Figure 9). Without expensive RDO searching process, learned compression models including ELIC can achieve encoding speed as fast as decoding speed, which is potentially helpful to conduct high-quality real-time visual data transmission. This can be an advantage of learning based approaches, since currently best conventional image coders like VVC spend much longer time in encoding than decoding[2].

We present a slim architecture ELIC-sm, which is adapted from ELIC by removing attention modules and using fewer res-blocks (RB×1). It also achieves a remarkable RD performance while significantly reducing the latency.

To further present our approach, we highlight Figure 1 and Figure 9. Figure 1 shows the RD-latency relationship of various approaches. It can be seen from the figure that our ELIC model achieves state-of-the-art performance regarding Pareto optimum. We also draw all tested learned image compression approaches which can decode the Ko-

dak image in 100 microseconds in Figure 9, which further indicates the superiority of ELIC. For completeness, we present more RD curves and results tested on Kodak and other datasets in the supplementary material.

### 6.2.2 Ablation study

**Influence of uneven grouping and spatial adaption in SCCTX.** We evaluate the proposed uneven grouping model and compare it with the previous even grouping ap-

| Data | Grouping Scheme | BD-Rate | Param. Latency |
|---|---|---|---|
| Kodak | even (10-slice) | -3.81 | 13.46 |
| | + spatial[P] | -5.02 | 19.41 |
| | uneven | -2.78 | **6.45** |
| | + spatial[P] | **-5.63** | 10.10 |
| | + spatial[S] | -6.99 | $> 10^3$ |
| CLIC-P | even (10-slice) | -12.42 | 73.59 |
| | + spatial[P] | -14.11 | 85.44 |
| | uneven | -11.54 | **31.88** |
| | + spatial[P] | **-14.51** | 46.06 |
| | + spatial[S] | -15.51 | $> 10^3$ |
| VVC | - | 0.00 | - |

Table 3. Performance of different grouping schemes. Since image resolutions in CLIC Professional vary, we report the latency repeatedly tested on its largest 2048 × 1890 image.

| Model | Layer | BD-Rate | Latency | |
|---|---|---|---|---|
| | | | Enc. | Dec. |
| Ballé2018 | GDN | 8.23 | 12.27 | 17.30 |
| Ballé2018 | RB×1 | 5.68 | 12.45 | 17.51 |
| Ballé2018 | PWRB×3 | 4.59 | 18.15 | 22.84 |
| Ballé2018 | RB×3 | **-2.94** | 23.00 | 27.91 |
| Minnen2018 | GDN | -5.04 | 12.75 | 17.34 |
| Minnen2018 | RB×1 | -6.34 | 12.90 | 17.44 |
| Minnen2018 | PWRB×3 | -11.14 | 18.10 | 22.78 |
| Minnen2018 | RB×3 | **-13.56** | 23.01 | 28.21 |
| BPG | - | 0.00 | - | - |

Table 4. Comparison of different nonlinear layers, evaluated on Kodak. The columns marked as **Latency Enc./Dec.** denote inference time of main analysis/synthesis. $RB×n$ denotes a stack of $n$ residual bottleneck blocks [25] (the right in Figure 7). *PWRB* is point-wise residual block which removes the $3 \times 3$ convolution.

proach [40]. For a fair comparison, we train a line of models with the same architectures for main and hyper autoencoders as ELIC. The only difference between them is the grouping models for backward-adaptive coding. Table 3 shows the results evaluated on Kodak and CLIC Professional. With the proposed uneven grouping, the latency of adaptive entropy estimation decreases by half. Introducing the spatial context model to both even and uneven models gain RD promotion, with uneven+spatial[P] running speed still faster than the 10-slice even grouping models.

**Residual blocks versus GDN (Table 4).** After replacing all GDN/IGDN layers with residual bottleneck blocks, the RD performance improves on both Ballé2018 [5] and Minnen2018 [39] baselines while the inference latency is still on par with GDN version. When stacking more RB blocks as nonlinear transform, the BD-rates even reduce more, which proves the scalability. We also try to stack GDN layers but the training becomes very unstable and eventually fails.

Since the residual bottleneck module has a larger receptive field than GDN, we also evaluate a point-wise version of it (*PWRB*) to figure out the influence of a larger receptive field. It still improves BD-rates on both baselines. Note that residual blocks with larger receptive fields improve more on models without spatial context (Ballé2018), implying that nonlinear transform with larger receptive fields is also helpful to remove spatial redundancy.

### 6.3. Qualitative results

**Full reconstruction.** We present the high-resolution images coded by ELIC and prior approaches in the supplementary material to show the reconstruction performance.

**Thumbnail-preview.** See Figure 10, where we present the thumbnail-preview of Kodak images generated from our



Figure 10. Quickly decoded thumbnail-preview images. For each pair of images, the left one is the downsampled ground-truth and the right is corresponding thumbnail-preview.
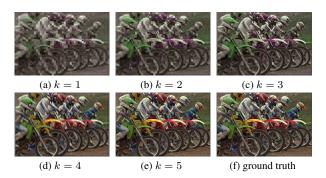


(a) $k = 1$    (b) $k = 2$    (c) $k = 3$

(d) $k = 4$    (e) $k = 5$    (f) ground truth

Figure 11. Progressive decoding. Each image is reconstructed from the first $k$ groups. When $k = 5$, the image is fully decoded.

thumbnail synthesizer. Since the thumbnail-preview images are in low-resolution by design, the relative low reconstruction quality (PSNR = $23.02$ dB evaluated on full-resolution Kodak) will not do harm to subjective feeling, as major semantic information remains and the artifacts are suppressed by image downsampling. The inference of thumbnail synthesizer takes about 3 microseconds, which is over 12 times faster than invoking the full synthesizer of ELIC.

**Progressive decoding** is an extended application of channel-conditional models [40]. We also present the progressive decoding results of ELIC. Different from Minnen *et al.* [40], at the $k$-th step, we directly fill the non-decoded channel chunks $\hat{\boldsymbol{y}}^{(k+1)}, \ldots, \hat{\boldsymbol{y}}^{(5)}$ with 0 before feeding them to the synthesizer. Therefore, extra calculation predicting non-decoded symbols is avoided. See Figure 11. With the beginning 16 channels, the structural information can be already reconstructed. The following channel groups further provide chrominance and high-frequency information. Since the progressive decoding is usually adopted for preview, we use the thumbnail synthesizer to quickly reconstruct the partially decoded ($k \leq 4$) images.

## 7. Discussion and conclusion

With the proposed SCCTX and the residual transform networks, we obtain state-of-the-art model ELIC, which better balances compression ability and running speed. Furthermore, we propose to train a thumbnail network for preview decoding, which also improves the utility. In the future, we will further investigate the information compaction phenomenon for improving the architecture.

Note that VVC is majorly designed for YUV 4:2:0 colorspace instead of YUV 4:4:4, as the former better reflects the sensitivity of human perception. We will also delve into models with both objective and subjective image quality remarkable, following exiting literature [10, 18, 38, 43].

# References

[1] Workshop and challenge on learned image compression (clic). http://www.compression.cc, 2020. 6

[2] Versatile video coding reference software version 12.1 (vtm-12.1). https://vcgit.hhi.fraunhofer.de/jvet/VVCSoftware_VTM/-/tags/VTM-12.1, 2021. 1, 7

[3] Johannes Ballé, Philip A Chou, David Minnen, Saurabh Singh, Nick Johnston, Eirikur Agustsson, Sung Jin Hwang, and George Toderici. Nonlinear transform coding. *IEEE Journal of Selected Topics in Signal Processing*, 15(2):339–353, 2020. 2

[4] Johannes Ballé, Valero Laparra, and Eero P Simoncelli. End-to-end optimized image compression. In *Int. Conf. on Learning Representations*, 2017. 1, 2

[5] Johannes Ballé, David Minnen, Saurabh Singh, Sung Jin Hwang, and Nick Johnston. Variational image compression with a scale hyperprior. In *Int. Conf. on Learning Representations*, 2018. 1, 2, 3, 5, 6, 7, 8

[6] J. Ballé. Efficient nonlinear transforms for lossy image compression. In *2018 Picture Coding Symposium (PCS)*, pages 248–252, June 2018. 5

[7] Jean Bégaint, Fabien Racapé, Simon Feltman, and Akshay Pushparaja. Compressai: a pytorch library and evaluation platform for end-to-end compression research. *arXiv preprint arXiv:2011.03029*, 2020. 7

[8] Fabrice Bellard. Bpg image format. https://bellard.org/bpg, 2015. 1, 7

[9] Gisle Bjontegaard. Calculation of average psnr differences between rd-curves. *VCEG-M33*, 2001. 6

[10] Yochai Blau and Tomer Michaeli. Rethinking lossy compression: The rate-distortion-perception tradeoff. In *International Conference on Machine Learning*, pages 675–685. PMLR, 2019. 9

[11] Tong Chen, Haojie Liu, Zhan Ma, Qiu Shen, Xun Cao, and Yao Wang. End-to-end learnt image compression via non-local attention optimization and improved context modeling. *IEEE Transactions on Image Processing*, 30:3179–3191, 2021. 5

[12] Tong Chen, Haojie Liu, Qiu Shen, Tao Yue, Xun Cao, and Zhan Ma. Deepcoder: A deep neural network based video compression. In *2017 IEEE Visual Communications and Image Processing (VCIP)*, pages 1–4. IEEE, 2017. 5

[13] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Deep convolutional autoencoder-based lossy image compression. In *2018 Picture Coding Symposium (PCS)*, pages 253–257. IEEE, 2018. 5

[14] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Deep residual learning for image compression. In *CVPR Workshops*, page 0, 2019. 5

[15] Zhengxue Cheng, Heming Sun, Masaru Takeuchi, and Jiro Katto. Learned image compression with discretized gaussian mixture likelihoods and attention modules. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1, 2, 3, 5, 6, 7

[16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE conference on computer vision and pattern recognition*, pages 248–255. IEEE, 2009. 6

[17] Qiang Duan, Xue Li, Qingshan Yin, Luoluo Feng, Jing Zhao, Yijin Teng, Xiaohui Duan, Yanhan Zhao, Ming Gao, Jianhua Wang, et al. A study on the generalized normalization transformation activation function in deep learning based image compression. In *Proceedings of Sixth International Congress on Information and Communication Technology*, pages 351–359. Springer, 2022. 5

[18] Hilmi E Egilmez, Ankitesh K Singh, Muhammed Coban, Marta Karczewicz, Yinhao Zhu, Yang Yang, Amir Said, and Taco S Cohen. Transform network architectures for deep learning based end-to-end image/video coding in subsampled color spaces. *IEEE Open Journal of Signal Processing*, 2:441–452, 2021. 1, 9

[19] Haisheng Fu, Feng Liang, Jianping Lin, Bing Li, Mohammad Akbari, Jie Liang, Guohe Zhang, Dong Liu, Chengjie Tu, and Jingning Han. Learned image compression with discretized gaussian-laplacian-logistic mixture model and concatenated residual modules. *arXiv preprint arXiv:2107.06463*, 2021. 1

[20] Ge Gao, Pei You, Rong Pan, Shunyuan Han, Yuanyuan Zhang, Yuchao Dai, and Hojae Lee. Neural image compression via attentional multi-scale back projection and frequency decomposition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14677–14686, 2021. 1, 2, 5, 7

[21] V. K. Goyal. Theoretical foundations of transform coding. *IEEE Signal Processing Magazine*, 18(5):9–21, Sep. 2001. 2, 3

[22] Zongyu Guo, Zhizheng Zhang, Runsen Feng, and Zhibo Chen. Causal contextual prediction for learned image compression. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021. 1, 2, 3, 4, 5, 6, 7

[23] Zongyu Guo, Zhizheng Zhang, Runsen Feng, and Zhibo Chen. Soft then hard: Rethinking the quantization in neural image compression. In *International Conference on Machine Learning*, pages 3920–3929. PMLR, 2021. 2

[24] Dailan He, Yaoyan Zheng, Baocheng Sun, Yan Wang, and Hongwei Qin. Checkerboard context model for efficient learned image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14771–14780, 2021. 1, 2, 3, 4, 6, 7

[25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5, 8

[26] ITU. Information technology - digital compression and coding of continuous - tone still images - requirements and guidelines. *CCITT, Recommendation*, 1992. 1, 7

[27] Ogun Kirmemis and A Murat Tekalp. Shrinkage as activation for learned image compression. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 1301–1305. IEEE, 2020. 5

[28] Eastman Kodak. Kodak lossless true color image suite (photocd pcd0992). http://r0k.us/graphics/kodak, 1993. 6

[29] Jooyoung Lee, Seunghyun Cho, and Munchurl Kim. An end-to-end joint learning scheme of image compression and quality enhancement with improved entropy minimization. *arXiv preprint arXiv:1912.12817*, 2019. 1, 2, 5, 6

[30] Chuming Li, Xin Yuan, Chen Lin, Minghao Guo, Wei Wu, Junjie Yan, and Wanli Ouyang. Am-lfs: Automl for loss function search. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8410–8419, 2019. 5

[31] Mu Li, Kede Ma, Jane You, David Zhang, and Wangmeng Zuo. Efficient and effective context-based convolutional entropy modeling for image compression. *IEEE Transactions on Image Processing*, 29:5900–5911, 2020. 3

[32] Chaoyi Lin, Jiabao Yao, Fangdong Chen, and Li Wang. A spatial rnn codec for end-to-end image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13269–13277, 2020. 2

[33] Haojie Liu, Tong Chen, Peiyao Guo, Qiu Shen, Xun Cao, Yao Wang, and Zhan Ma. Non-local attention optimized deep image compression. *arXiv preprint arXiv:1904.09757*, 2019. 2, 3, 4

[34] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *International Conference on Learning Representations*, 2018. 5

[35] Jiaheng Liu, Guo Lu, Zhihao Hu, and Dong Xu. A unified end-to-end framework for efficient deep image compression. *arXiv preprint arXiv:2002.03370*, 2020. 2, 6

[36] Changyue Ma, Zhao Wang, Ruling Liao, and Yan Ye. A cross channel context model for latents in deep image compression. *arXiv preprint arXiv:2103.02884*, 2021. 3, 4

[37] Haichuan Ma, Dong Liu, Ning Yan, Houqiang Li, and Feng Wu. End-to-end optimized versatile image compression with wavelet-like transform. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2

[38] Fabian Mentzer, George D Toderici, Michael Tschannen, and Eirikur Agustsson. High-fidelity generative image compression. *Advances in Neural Information Processing Systems*, 33:11913–11924, 2020. 9

[39] David Minnen, Johannes Ballé, and George D Toderici. Joint autoregressive and hierarchical priors for learned image compression. In *Advances in Neural Information Processing Systems*, pages 10771–10780, 2018. 1, 2, 3, 4, 5, 6, 7, 8

[40] David Minnen and Saurabh Singh. Channel-wise autoregressive entropy models for learned image compression. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 3339–3343. IEEE, 2020. 1, 2, 3, 4, 5, 6, 7, 8

[41] Andrew Ng et al. Sparse autoencoder. *CS294A Lecture notes*, 72(2011):1–19, 2011. 4

[42] Aäron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, and Koray Kavukcuoglu. Condi-

tional image generation with pixelcnn decoders. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 4797–4805, 2016. 2

[43] Yash Patel, Srikar Appalaraju, and R Manmatha. Saliency driven perceptual image compression. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 227–236, 2021. 9

[44] Marc Ranzato, Y-Lan Boureau, Yann LeCun, et al. Sparse feature learning for deep belief networks. *Advances in neural information processing systems*, 20:1185–1192, 2008. 4

[45] Marc Ranzato, Christopher Poultney, Sumit Chopra, Yann LeCun, et al. Efficient learning of sparse representations with an energy-based model. *Advances in neural information processing systems*, 19:1137, 2007. 4

[46] Baocheng Sun, Meng Gu, Dailan He, Tongda Xu, Yan Wang, and Hongwei Qin. Hlic: Harmonizing optimization metrics in learned image compression by reinforcement learning. *arXiv preprint arXiv:2109.14863*, 2021. 5

[47] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003. 1

[48] Andreas Weinlich, Peter Amon, Andreas Hutter, and Andre Kaup. Probability distribution estimation for autoregressive pixel-predictive image coding. *IEEE Transactions on Image Processing*, 25(3):1382–1395, 2016. 1

[49] Yaojun Wu, Xin Li, Zhizheng Zhang, Xin Jin, and Zhibo Chen. Learned block-based hybrid image compression. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021. 1

[50] Yueqi Xie, Ka Leong Cheng, and Qifeng Chen. Enhanced invertible encoding for learned image compression. In *Proceedings of the ACM International Conference on Multimedia*, 2021. 1, 2, 5, 7

[51] Fei Yang, Luis Herranz, Yongmei Cheng, and Mikhail G Mozerov. Slimmable compressive autoencoders for practical neural image compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4998–5007, 2021. 5

[52] Yibo Yang, Robert Bamler, and Stephan Mandt. Improving inference for neural image compression. *Advances in Neural Information Processing Systems*, 33:573–584, 2020. 2

[53] Jiahui Yu, Pengchong Jin, Hanxiao Liu, Gabriel Bender, Pieter-Jan Kindermans, Mingxing Tan, Thomas Huang, Xiaodan Song, Ruoming Pang, and Quoc Le. Bignas: Scaling up neural architecture search with big single-stage models. In *European Conference on Computer Vision*, pages 702–717. Springer, 2020. 5

[54] Jing Zhao, Bin Li, Jiahao Li, Ruiqin Xiong, and Yan Lu. A universal encoder rate distortion optimization framework for learned compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1880–1884, 2021. 2

# Supplementary Material

**ELIC: Efficient Learned Image Compression with Unevenly Grouped Space-Channel Contextual Adaptive Coding**

## 1. Detailed network architecture

### 1.1. Architecture of transform networks

| Analyzer $g_a$ | Synthesizer $g_s$ |
|---|---|
| in: 3-channel image | in: $M$-channel symbols |
| Conv $5 \times 5$, s2, $N$ | Attention |
| ResBottleneck$\times 3$ | TConv $5 \times 5$, s2, $N$ |
| Conv $5 \times 5$, s2, $N$ | ResBottleneck$\times 3$ |
| ResBottleneck$\times 3$ | TConv $5 \times 5$, s2, $N$ |
| Attention | Attention |
| Conv $5 \times 5$, s2, $N$ | ResBottleneck$\times 3$ |
| ResBottleneck$\times 3$ | TConv $5 \times 5$, s2, $N$ |
| Conv $5 \times 5$, s2, $M$ | ResBottleneck$\times 3$ |
| Attention | TConv $5 \times 5$, s2, 3 |

Table 1. Architecture of ELIC main transform networks.

| Analyzer $g_a$ | Synthesizer $g_s$ |
|---|---|
| in: 3-channel image | in: $M$-channel symbols |
| Conv $5 \times 5$, s2, $N$ | TConv $5 \times 5$, s2, $N$ |
| ResBottleneck$\times 1$ | ResBottleneck$\times 1$ |
| Conv $5 \times 5$, s2, $N$ | TConv $5 \times 5$, s2, $N$ |
| ResBottleneck$\times 1$ | ResBottleneck$\times 1$ |
| Conv $5 \times 5$, s2, $N$ | TConv $5 \times 5$, s2, $N$ |
| ResBottleneck$\times 1$ | ResBottleneck$\times 1$ |
| Conv $5 \times 5$, s2, $M$ | TConv $5 \times 5$, s2, 3 |

Table 2. Architecture of ELIC-sm main transform networks.

As shown in Table 1 and Table 2, our main networks are adapted from those previously used by Ballé *et al*. (2018), Minnen *et al*. (2018) and Minnen *et al*. (2020). We replace GDN/IGDN layers by residual blocks, as mentioned in the main text. We follow previous works to introduce a variable $N$ representing the channel number of intermediate features, and lift the output channel number of analyzer to $M$. Thus, the coding-symbol $\hat{y}$ is an $H \times W \times M$ tensor.

We adopt residual bottleneck structures as nonlinear transform, which have a dimensional bottleneck which can both provide larger receptive field and keep the volume of calculation acceptable. The channel number of intermedi-

ate features before and after $3 \times 3$ convolution is $\frac{N}{2}$. We sequentially stack 3 such residual blocks after each down/up-sampling convolution in ELIC (Table 1), and instead use only one at each position in ELIC-sm (Table 2). We also use the attention modules adopted by Cheng *et al*. in the larger ELIC models to further enhance the nonlinearity.

We simply adopt the three-layer hyper analyzer and synthesizer, following previous works (Minnen *et al*., 2018; Minnen *et al*., 2020; He *et al*., 2021). The hyper synthesizer output is an $H \times W \times (2M)$ tensor $\boldsymbol{\Psi}$ as shown in Figure 7 in the main text.

### 1.2. Architecture of SCCTX networks

Following Minnen *et al*. (2020), we use $5 \times 5$ convolutions to analyze cross-channel redundancy. The architecture of $g_{\mathrm{ch}}$ network is frankly sketched from Minnen *et al*. (2020). We follow the suggestions of Minnen *et al*. (2018) and He *et al*. (2021) and use single-layer $5 \times 5$ autoregressive/checkerboard-masked convolution as spatial context model $g_{\mathrm{sp}}$. Both $g_{\mathrm{sp}}^{(k)}$ and $g_{\mathrm{ch}}^{(k)}$ output $H \times W \times (2M^{(k)})$ features, where $M^{(k)}$ is the channel number of the $k$-th channel group.

Therefore, we have the concatenated adaptive representation tensor $[\boldsymbol{\Phi}_{\mathrm{sp},i}^{(k)}, \boldsymbol{\Phi}_{\mathrm{ch}}^{(k)}, \boldsymbol{\Psi}] \in \mathbb{R}^{H \times W \times (4M^{(k)}+2M)}$. The parameter aggregation network linearly reduces the dimensions to $2M^{(k)}$.

## 2. Detailed experimental settings

We implement, train, and evaluate all learning-based models on PyTorch 1.8.1. We use NVIDIA TITANXP to test both RD performance and inference speed. To test the speeds, we reproduce previously proposed models and evaluate them under the same running conditions for fair comparison. Since most of the models adopt reparameterization techniques, we fix the reparameterized weights before testing the speed. We follow a common protocol to test the latency with GPU synchronization. When testing each model, we drop the latency results (we do not drop them when evaluating RD performance) of the first 6 images to get rid of the influence of device warm-up, and average the running time of remained images to get the precious speed results.

We do not enable the deterministic inference mode (*e.g.* `torch.backends.cudnn.deterministic`) when testing the model speeds for two reasons. First, we tend

to believe that the deterministic issue can be well solved with engineering efforts, such as using integer-only inference. Thus, the deterministic floating-point inference is unnecessary. Second, the deterministic mode extremely slows down the speed of specific operators, like transposed convolutions which are adopted by ELIC and earlier baseline models (Ballé *et al.* and Minnen *et al.*), making the comparison somewhat unfair.

We obtain the RD results of prior works by asking the authors via emails or accessing released data directly.

### 2.1. Visualization explanation

We visualize the coding-symbols to investigate the information compaction property. The visualization results of the lighthouse image is shown in Figure 2 of the main text. The $\ell$-th gray scale image is the rescaled magnitude $\boldsymbol{f}^\ell$ of the $\ell$-th symbol channel $\hat{\boldsymbol{y}}^{(\ell)}$:

$$\boldsymbol{f}^{(\ell)} = \frac{|\hat{\boldsymbol{y}}^{(\ell)}|}{\max \hat{\boldsymbol{y}}}, \quad \ell = 1, 2 \ldots, M \quad (1)$$

Inspired by the concept of energy in the signal processing community, we introduce the term energy to describe the average square value of each symbol channel:

$$e^{(\ell)} = \frac{1}{HW} \sum_{\hat{y}_i^{(\ell)} \in \hat{\boldsymbol{y}}^{(\ell)}} \hat{y}_i^{(\ell)2} \quad (2)$$

In the main text, Figure 2 presents the channels with the largest energy values. Figure 3-left further shows the logarithmic energy $\log(e)$ of each channel.

## 3. More rate-distortion results

We also calculate the BD-rate over VVC when adopting MS-SSIM as the distortion metric. The results are shown in Table 3. Several baselines are skipped since the original authors do not provide the MS-SSIM results evaluated on the MSE-optimized models.

For completeness, we evaluate our ELIC model at larger BPP range (from 1.0 to 1.5 on Kodak). We train these models with $\lambda = \{0.08, 0.16\}$. The high-rate models still perform well. We draw the results on larger figures (Figure 1 for PSNR and Figure 2 for MS-SSIM) to more clearly show the superiority of the proposed approach. We also evaluate the models on CLIC-Professional and report Figure 3.

We further evaluate our model performance when optimizing for MS-SSIM, following prior works. Thus, the distortion term of the loss function is replaced by $1 -$ MSSSIM($\boldsymbol{x}$). We use $\lambda = \{3, 12, 40, 120\}$ to train such models, following Cheng *et al.* (2020). Figure 4 shows the results. On Kodak, when achieving the same MS-SSIM, ELIC optimized for MS-SSIM saves about half of the bitrate compared with VVC, which is encouraging.

| Model | BD-Rate (%) (PSNR) | BD-Rate (%) (MS-SSIM in dB) |
|---|---|---|
| ELIC (ours) | **-7.88** | **-12.73** |
| ELIC-sm (ours) | -1.07 | -5.59 |
| Minnen2020 | 1.11 | - |
| Cheng2020[P] | 3.89 | -7.19 |
| Minnen2018[P] | 20.00 | 4.23 |
| Ballé2018 | 40.85 | 16.41 |
| Gao2021 | -10.94 | - |
| Guo2021 | -7.02 | - |
| Wu2021 [43] | -5.71 | - |
| Xie2021 | -0.54 | -7.82 |
| Cheng2020 | 3.35 | -3.17 |
| Minnen2018 | 14.92 | -0.68 |
| VVC | 0.00 | 0.00 |

Table 3. BD-rates over VVC for learned image compression models. The BD-Rate data is calculated from rate-distortion curves over data points with BPP < 1 on Kodak. We present results of two distortion metrics, PSNR and MS-SSIM in dB. Models marked with [P] adopt parallel checkerboard context model and [S] denotes serial context model.

| Layer | YEnc. (s) | YDec. (s) |
|---|---|---|
| GDN | 0.500 | 1.108 |
| RB×1 | 0.448 | 1.165 |
| PWRB×3 | 0.583 | 1.191 |
| RB×3 | 0.800 | 1.448 |

Table 4. Ballé2018 running speed on CPU with different nonlinear layers, evaluated on Kodak.

| Layer | YEnc. | | YDec. | |
| | Param. (M) | MACs (G) | Param. (M) | MACs (G) |
|---|---|---|---|---|
| GDN | 3.51 | 36.89 | 3.51 | 127.63 |
| RB×1 | 3.76 | 47.64 | 3.75 | 138.38 |
| PWRB×3 | 3.73 | 46.53 | 3.73 | 137.26 |
| RB×3 | 4.48 | 78.71 | 4.48 | 169.44 |

Table 5. Ballé2018 parameter volumes and MACs with different nonlinear layers, evaluated on Kodak.

## 4. Res-blocks v.s. GDN on CPU

For completeness, we also evaluate the latency of models with various nonlinear blocks on Intel Core i7-7700 (Table 4) as well as their parameter volumes and numbers of multiply-accumulate operations (MACs) (Table 5). To test parameter volumes and MACs, we use the FLOPs profiler

of DeepSpeed [1]. After replacing GDN blocks with residual blocks (RB×1), the adapted Ballé2018 models have slightly larger MACs yet their en/de-coding speeds are on par with the original model with GDNs. This is similar to the GPU results in Table 4 of the main body.

## 5. Image reconstruction results

See Figures 5,6, we compare the reconstruction results of the proposed ELIC and Cheng *et al*. (2020), since their synthesis latency is close. We also present the decoded image of VTM12.1 for reference.

## 6. Progressive decoding

In the main text, we present progressive decoding results with thumbnail synthesizer and zero filling. In Figures 7,8, We further show progressive decoding results with full synthesizer and/or mean filling.
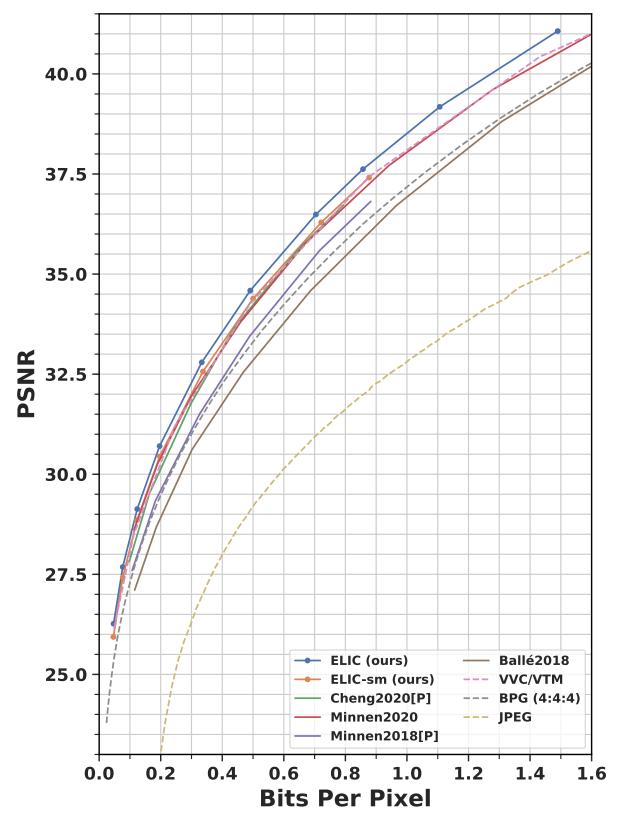
---

[1] https://www.deepspeed.ai/tutorials/flops-profiler/

Figure 1. **PSNR-BPP curve on Kodak.** All models are optimized for minimizing MSE. All presented learning-based models can decode the $512 \times 768$ Kodak image in 100 microseconds.
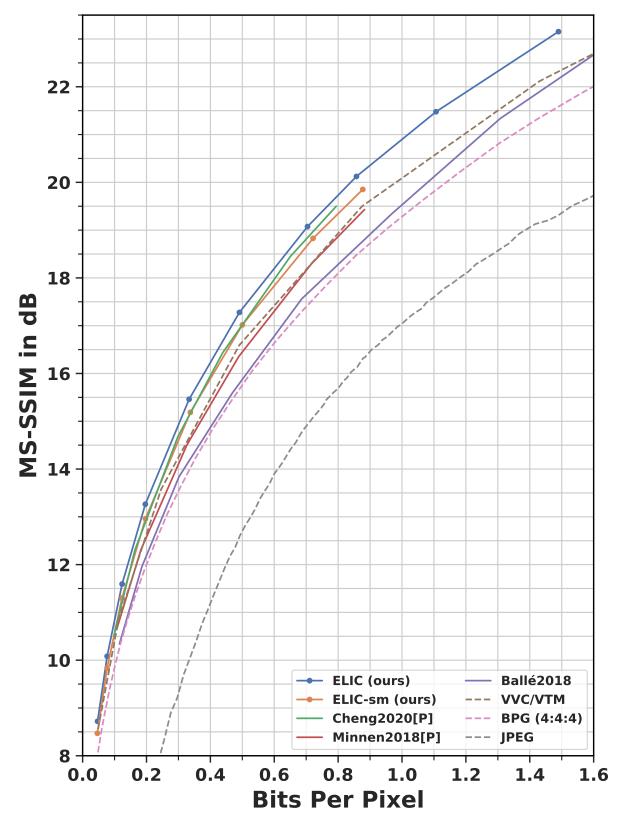
Figure 2. **MS-SSIM RD curve on Kodak.** All models are optimized for minimizing MSE. All presented learning-based models can decode the $512 \times 768$ Kodak image in 100 microseconds.
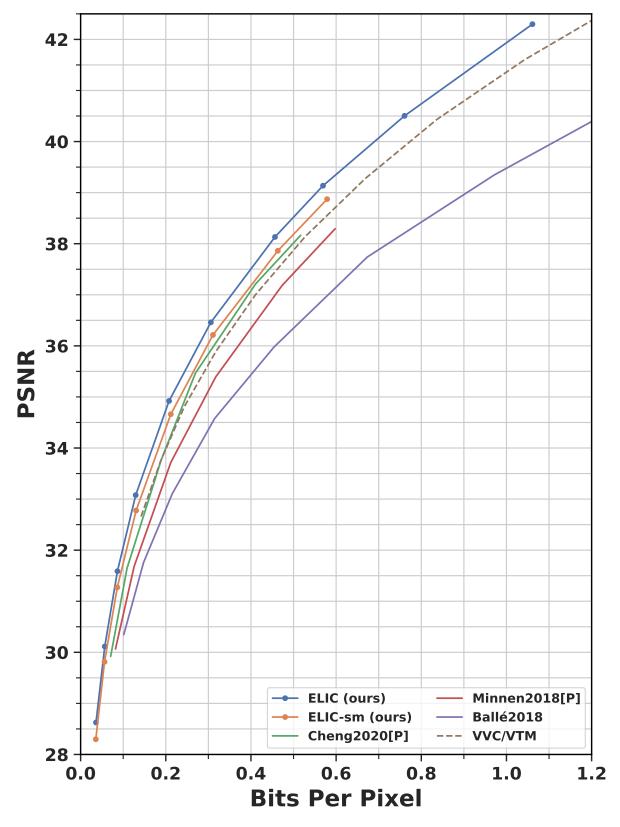
Figure 3. **PSNR-BPP curve on CLIC-Professional.** All models are optimized for minimizing MSE. All presented learning-based models can decode the $512 \times 768$ Kodak image in 100 microseconds.
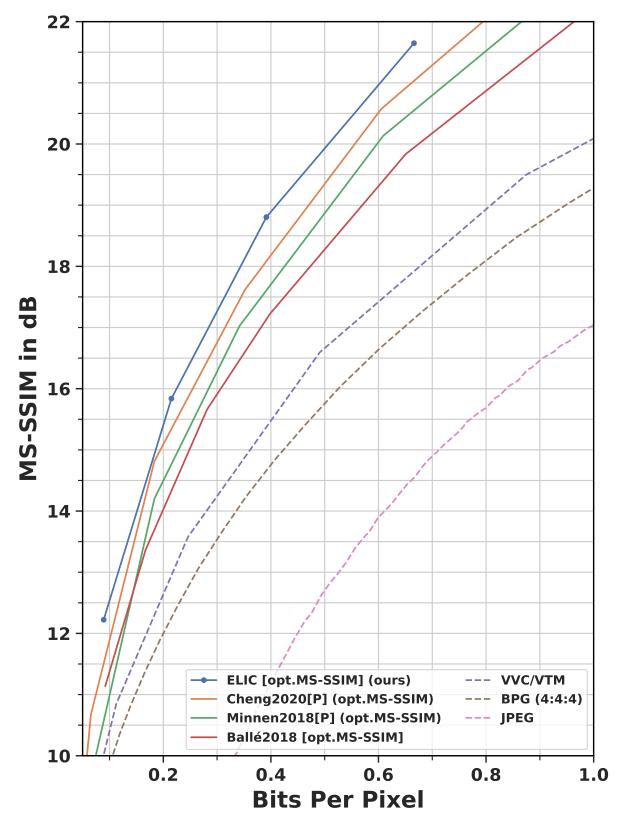
Figure 4. **MS-SSIM RD curve on Kodak (opt. MS-SSIM).** All models are optimized for maximizing MS-SSIM. All presented learning-based models can decode the $512 \times 768$ Kodak image in 100 microseconds.

(a) Cheng *et al.* (2020). BPP= 0.129. PSNR= 29.36      (b) Ours. BPP= 0.144. PSNR= 30.42

(c) Ground-truth.      (d) VTM12.1. BPP= 0.151. PSNR= 30.28

Figure 5. Qualitative comparison on reconstructed lighthouse (`kodim19`) image.

(a) Cheng *et al*. (2020). BPP= 0.155. PSNR= 29.85

(b) Ours. BPP= 0.127. PSNR= 31.66

(c) Ground-truth.

(d) VTM12.1. BPP= 0.151. PSNR= 31.65

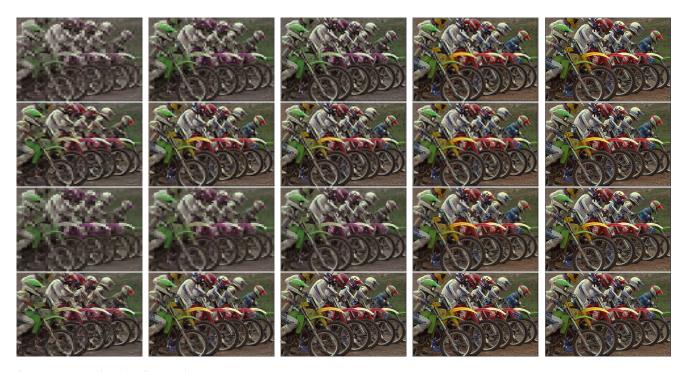Figure 6. Qualitative comparison on reconstructed sculpture (`kodim17`) image.

Figure 7. Progressive decoding results (kodim05).
Line 1 (thumb-zero): the first 4 groups are decoded with thumbnail decoder and zero filling (the same setting as the main text).
Line 2 (thumb-mean): the thumbnail-decoding results with mean filling.
Line 3 (full-zero): all the 5 groups are decoded with full synthesizer, using zero filling.
Line 4 (full-mean): the full-decoding results with mean filling.



Figure 8. Progressive decoding results of the house image (kodim24). The image order is the same as Figure 7.