# Machine Oriented Compression

Dan Jacobellis

The University of Texas at Austin

UT-SysML
Systems & Machine Learning

# Why do we need compression?

Modern sensor hardware provides incredible power efficiency



7-channel spatial audio array
**6 Mbits/sec**

2.5 Watt-hour battery

IMU, barometer, magnometer

**30 Hours of recording on a single charge**

8MP Front facing camera
**200 Mbits/sec**

2x high FOV scene cameras
**74 Mbits/sec**

projectaria.com

2x eye tracking cameras
**36 Mbits/sec**

High resolution signals are too large to process on device or transmit to the cloud

# Types of compression systems



**Lossless**

Huffman

ANS

Arithmetic

Golomb

**Lossy**

Resolution Reduction

Scalar quantization

Vector quantization

**Linear Transform coding**

JPEG 2000　　　　JPEG

WEBP　　　　AV1

**Learned Compression**
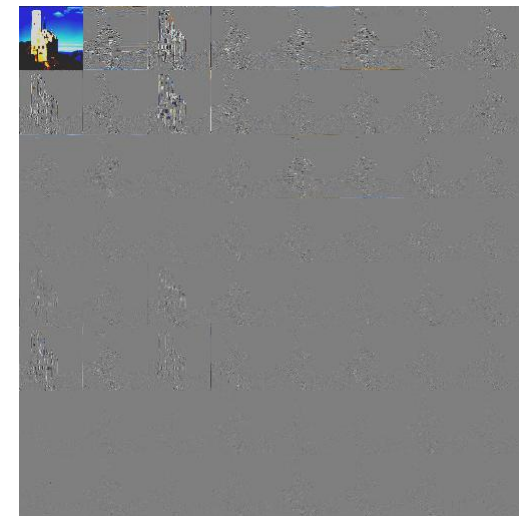
NNCP　　　　JPEG-AI

**Generative Compression**

Bits back coding　　　Stable Diffusion

Cosmos

# Linear transform coding

- Most of the bits that move across the internet today use linear transform coding (e.g. JPEG, AV1)

- These codecs use energy compacting transforms (e.g. DCT) to create a sparse representation

- Bits are allocated to different components using models of human perception

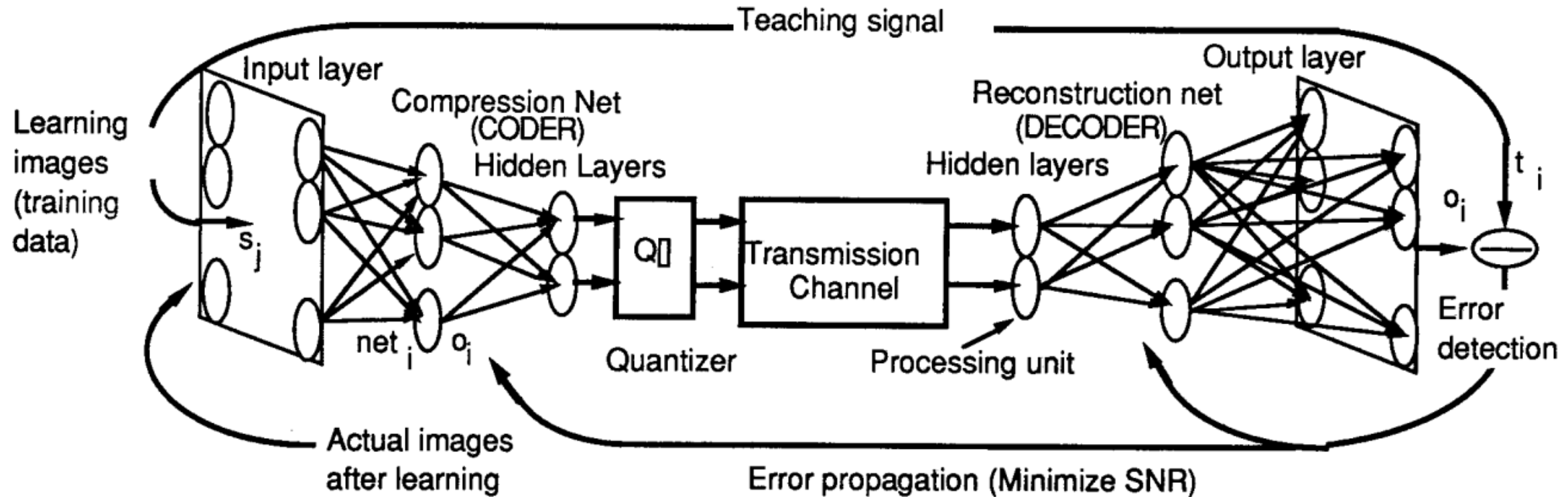- Exploit sparsity via entropy coding (RLE, huffman, etc)



$$\begin{bmatrix} 52 & 55 & 61 & 66 & 70 & 61 & 64 & 73 \\ 63 & 59 & 55 & 90 & 109 & 85 & 69 & 72 \\ 62 & 59 & 68 & 113 & 144 & 104 & 66 & 73 \\ 63 & 58 & 71 & 122 & 154 & 106 & 70 & 69 \\ 67 & 61 & 68 & 104 & 126 & 88 & 68 & 70 \\ 79 & 65 & 60 & 70 & 77 & 68 & 58 & 75 \\ 85 & 71 & 64 & 59 & 55 & 61 & 65 & 83 \\ 87 & 79 & 69 & 68 & 65 & 76 & 78 & 94 \end{bmatrix} \quad \begin{bmatrix} -26 & -3 & -6 & 2 & 2 & -1 & 0 & 0 \\ 0 & -2 & -4 & 1 & 1 & 0 & 0 & 0 \\ -3 & 1 & 5 & -1 & -1 & 0 & 0 & 0 \\ -3 & 1 & 2 & -1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

| 0 | -1 | 1 | -3 | 2 | -4 | -2 | 5 | -6 | -26 |
|---|-----|-----|------|------|-------|-------|-------|--------|--------|
| 1 | 000 | 001 | 0101 | 0110 | 01000 | 01001 | 01110 | 011110 | 011111 |

# Learned compression using neural networks



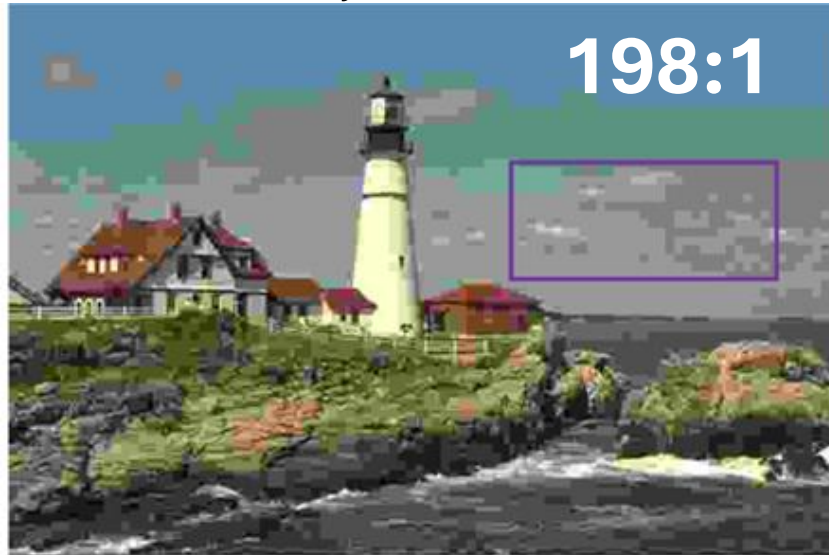Learn the transform and quantizer from representative data

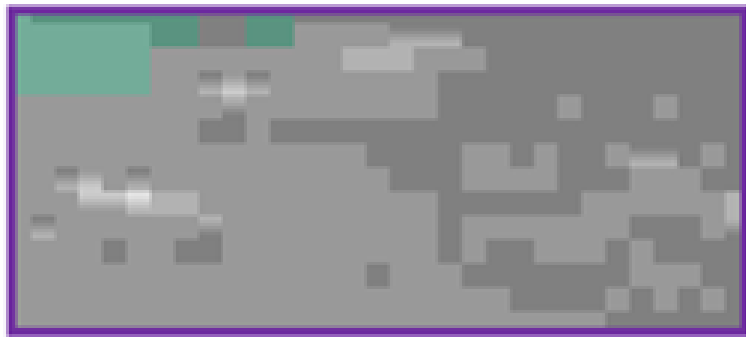**High compression efficiency**     **Poor computational efficiency**

Sonehara, et al. "Image data compression using a neural network model."
*International 1989 Joint Conference on Neural Networks*. IEEE, 1989.

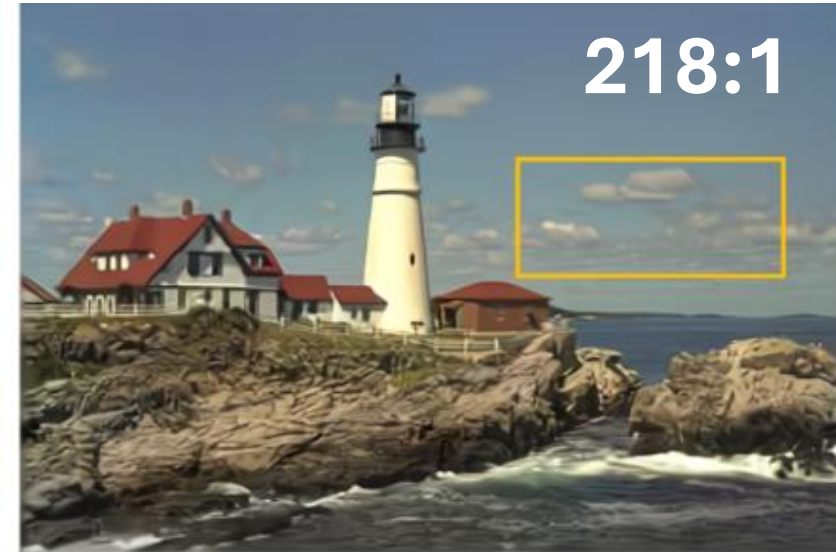# Compression efficiency vs computational efficiency

JPEG

Original

Neural Network

**198:1**

**218:1**

64 parameters

<500 MACs/pixel

Millions of Parameters

>100k MACs/pixel

# Generative compression

- Autoencoder will struggle to preserve **details, texture** and **high frequencies**

- Use a **generative model** to resynthesize the details

# Who is the perceiver?

We need **machine-oriented** compression systems
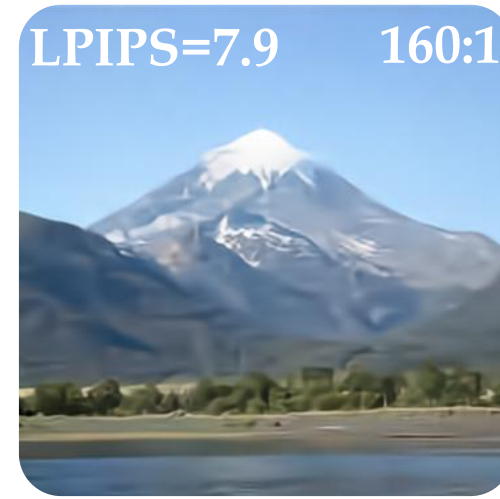
# Perceptual quality



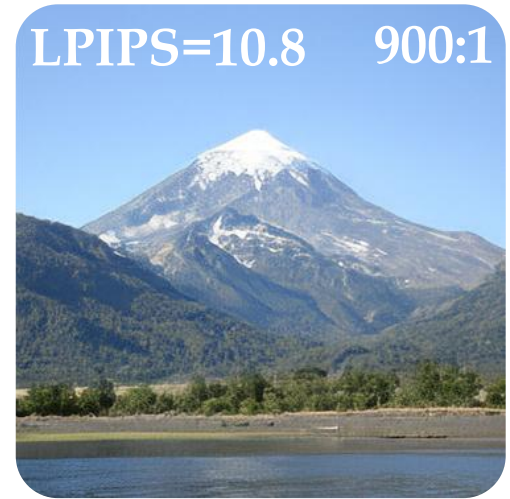Legacy transform codec — LPIPS=6.1  90:1

Modern transform codec — LPIPS=7.0  160:1

MSE Autoencoder — LPIPS=7.9  160:1

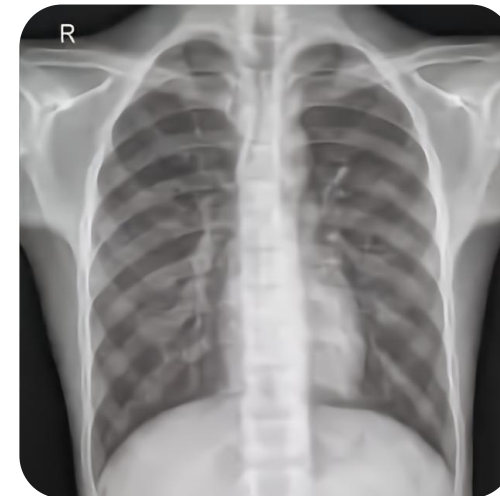Generative Autoencoder — LPIPS=10.8  900:1

Human <  Human <  Human <

Human perceptual quality is well modeled (e.g. LPIPS)
What about machine perception for specific applications?

?  ?  ?

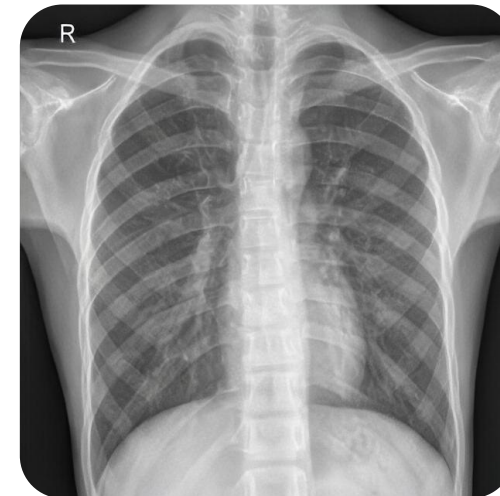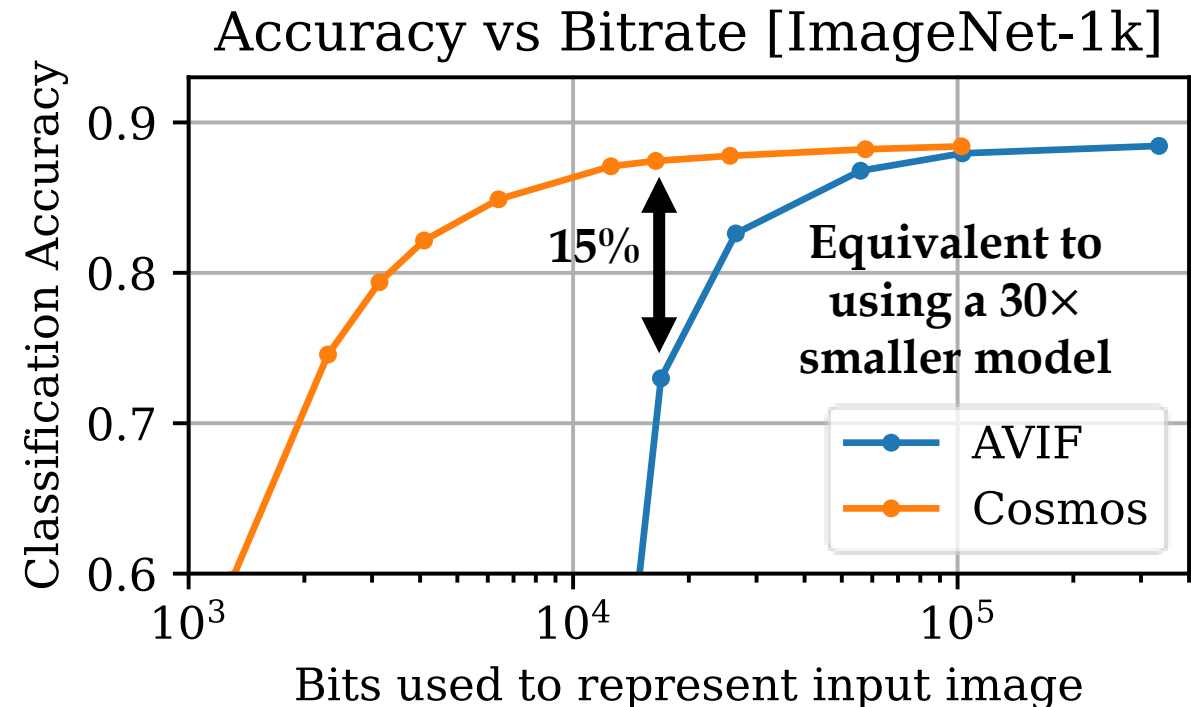# Machine perceptual quality

**Original (0.1 MP)**  **AVIF**  **Cosmos**



$2.5 \times 10^6$ bits     $3.1 \times 10^4$ bits     $2.5 \times 10^4$ bits

Does the high human perceptual quality of generative codecs translate to high machine perceptual quality?

Yes; generative codecs often provide **better downstream performance** than conventional methods at **lower rates**

## Accuracy vs Bitrate [ImageNet-1k]



**15%**

**Equivalent to using a 30× smaller model**
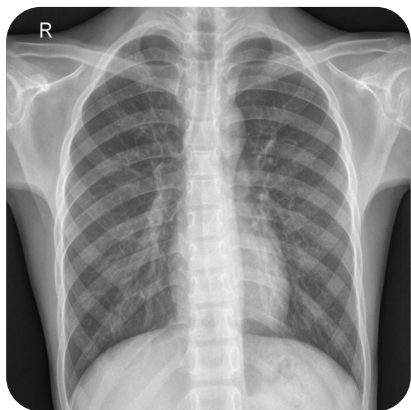
— AVIF
— Cosmos

Bits used to represent input image

# Does lossy compression always hurt accuracy?

How could lossy compression *increase* performance?

Datasets used to pre-train foundation models use legacy JPEG and MPEG compression at default settings

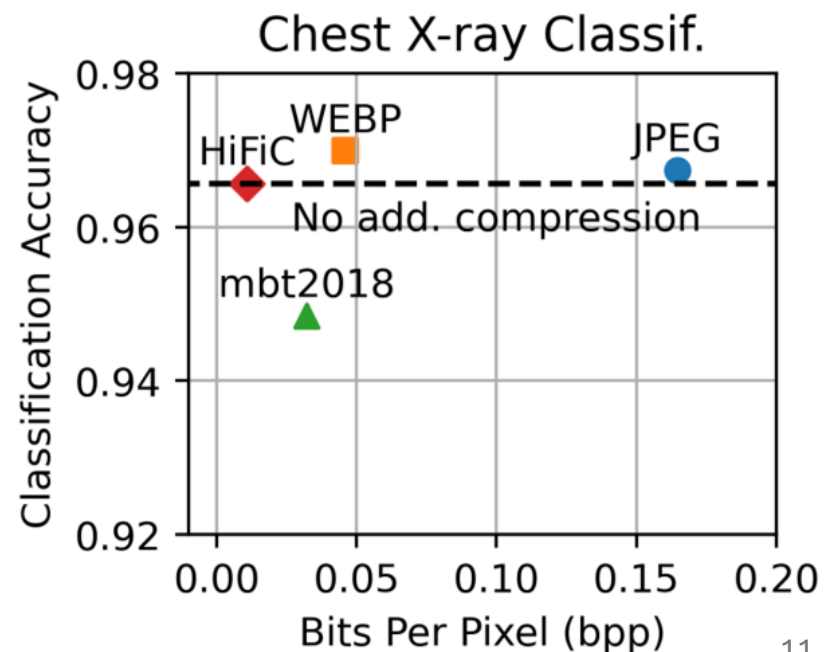High quality, lossless samples are **out of distribution!**
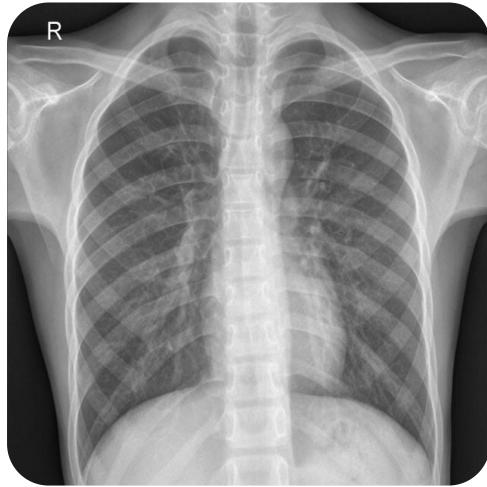
**Legacy transform coding**

**Modern transform coding**

**MSE-optimized autoencoder**

**Generative compression**



Pristine, high resolution inputs



Chest X-ray Classif.
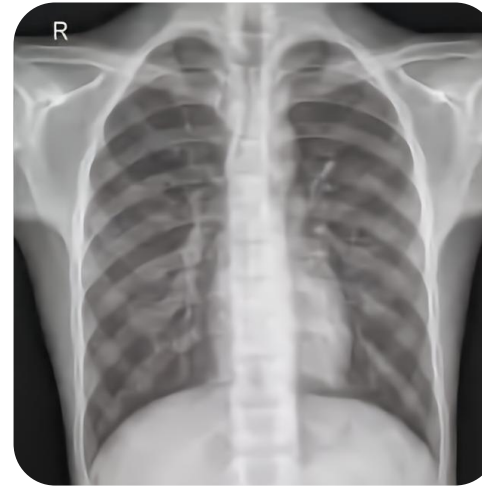
# Denoising effect of lossy compression



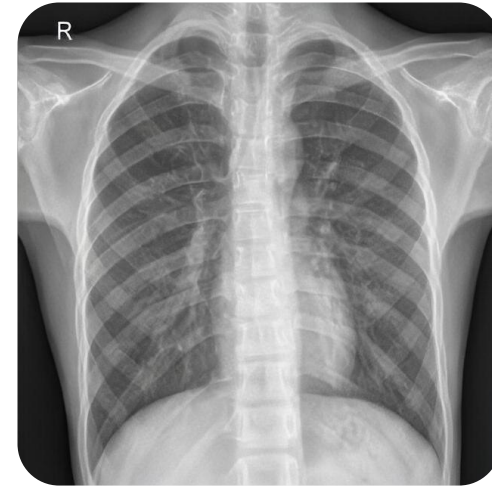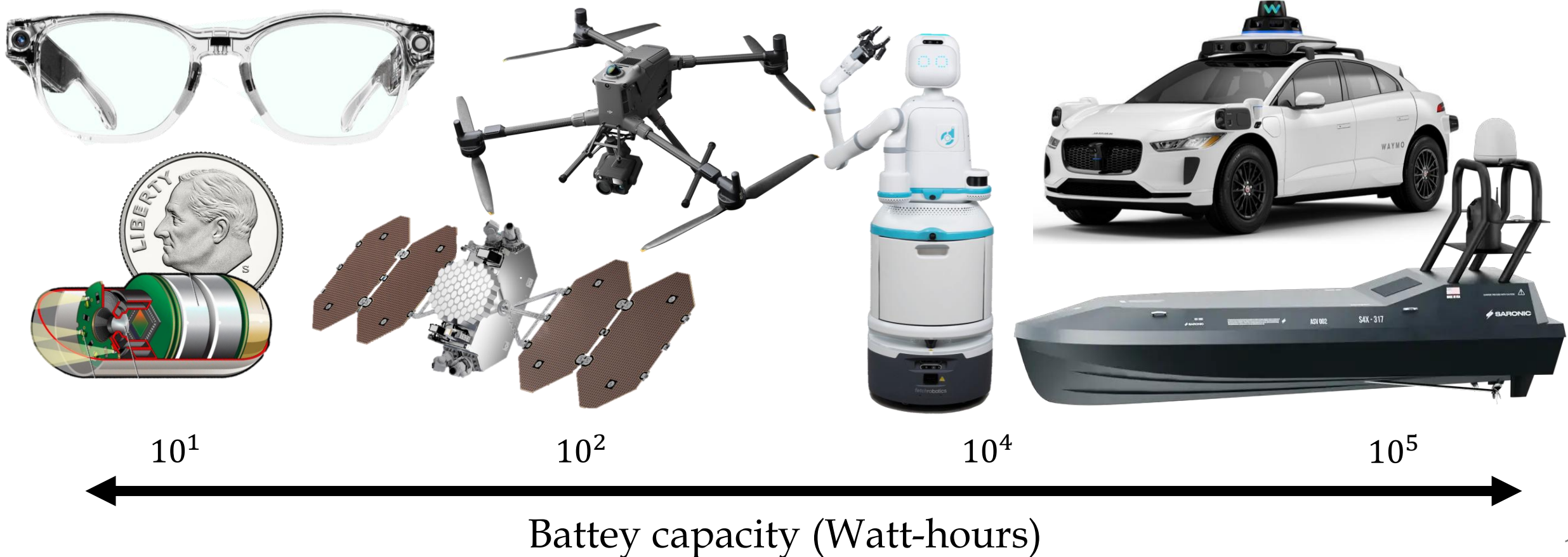| Original | Legacy transform coding (JPEG) | Modern transform coding (WEBP) | MSE-optimized Autoencoder | Generative model (HiFiC) |

# How much power is available for sensing?

Mobile, remote, and wearable sensors produce constant streams of **high resolution** signals

Sensor efficiency is increasing, while ML models get more expensive

**Solution:** divide computation between sensor and cloud



$10^1$      $10^2$      $10^4$      $10^5$

Battey capacity (Watt-hours)

# Compression for mobile and remote sensing

Mobile, remote, and wearable sensors produce constant streams of **high resolution** signals

Sensor efficiency is increasing, while ML models get more expensive

**Solution:** divide computation between sensor and cloud

**Sensor**

**Remote/Cloud**



Original Signal

Enc.

Dec.

Lossy reconstruction

ML Applications

Classification

Segmentation

Enhancement

**Demands high compression ratio**

**Degrades accuracy**

**Adds decoding overhead**

# Machine-oriented compression

Mobile, remote, and wearable sensors produce constant streams of **high resolution** signals

Sensor efficiency is increasing, while ML models get more expensive

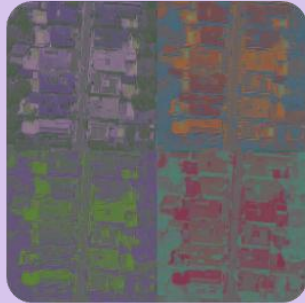**Solution:** divide computation between sensor and cloud



**Sensor**

Machine-interpretable features

**Remote/Cloud**

Enc.

Original Signal

ML Applications

Classification

Segmentation

Enhancement

Optional decoding

**Less bandwidth**     **Enhanced accuracy**     **More efficient ML**

# Machine-oriented compression

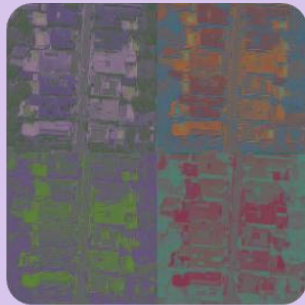> *What are ideal characteristics of the compression system?*

# Machine-oriented compression

*What are ideal characteristics of the compression system?*
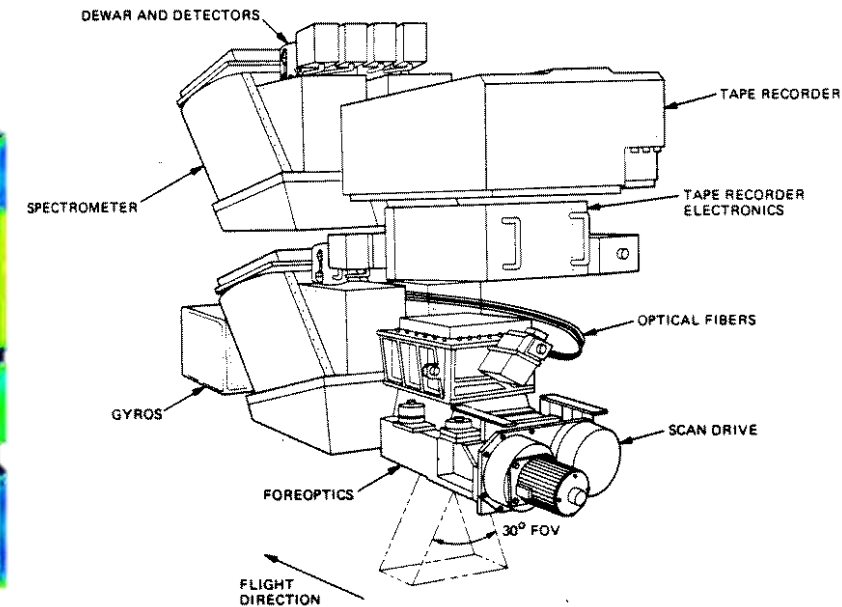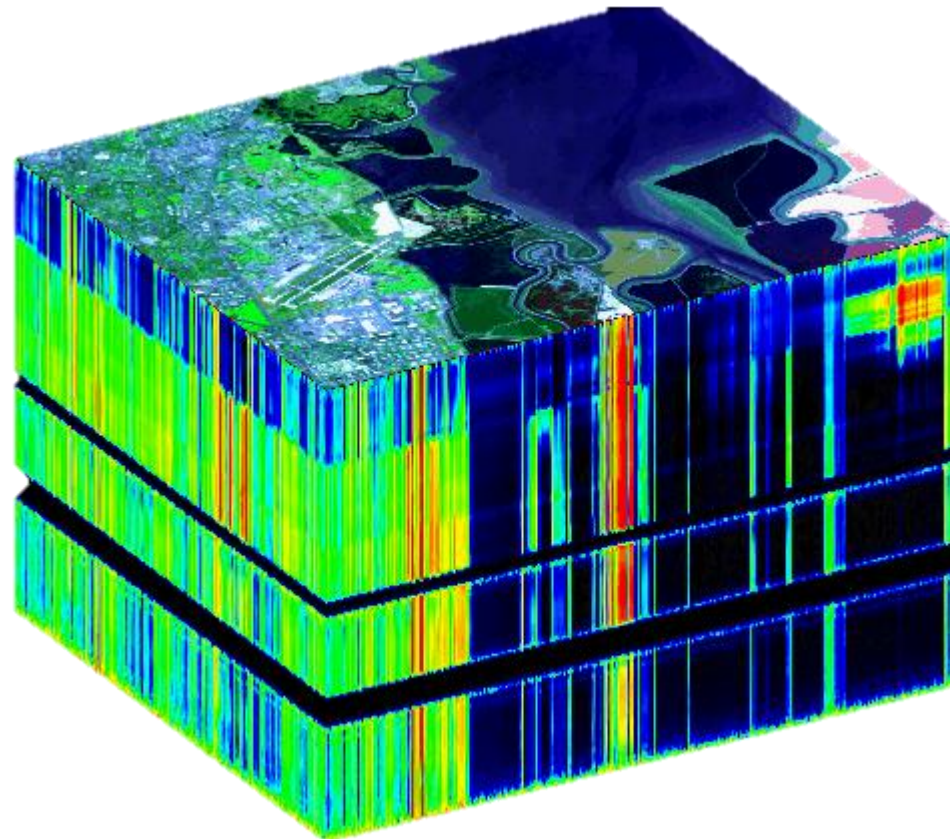
- **Support many modalities**
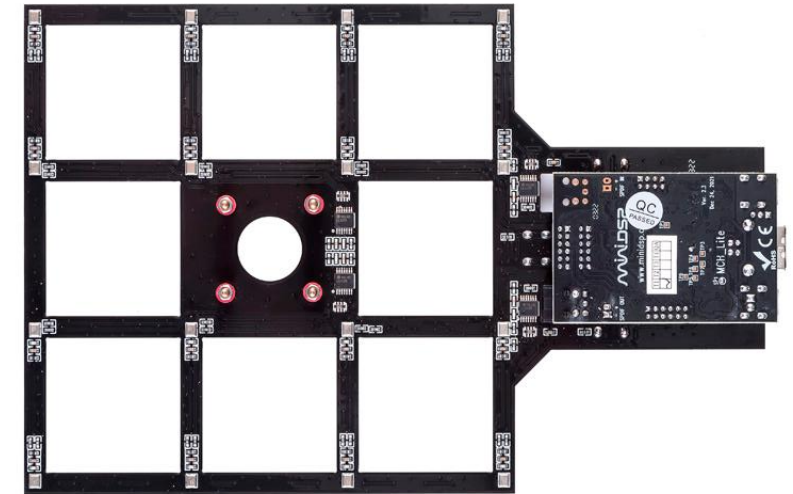
- **Hyperspectral**

# Machine-oriented compression

*What are ideal characteristics of the compression system?*

- **Support many modalities**

- Hyperspectral

- **Spatial Audio**

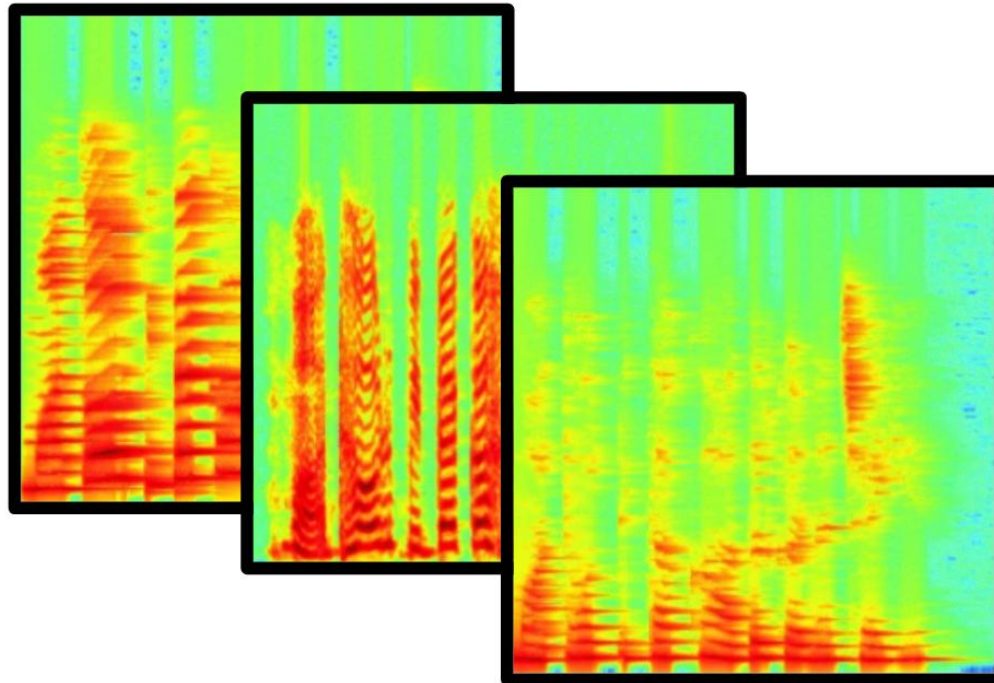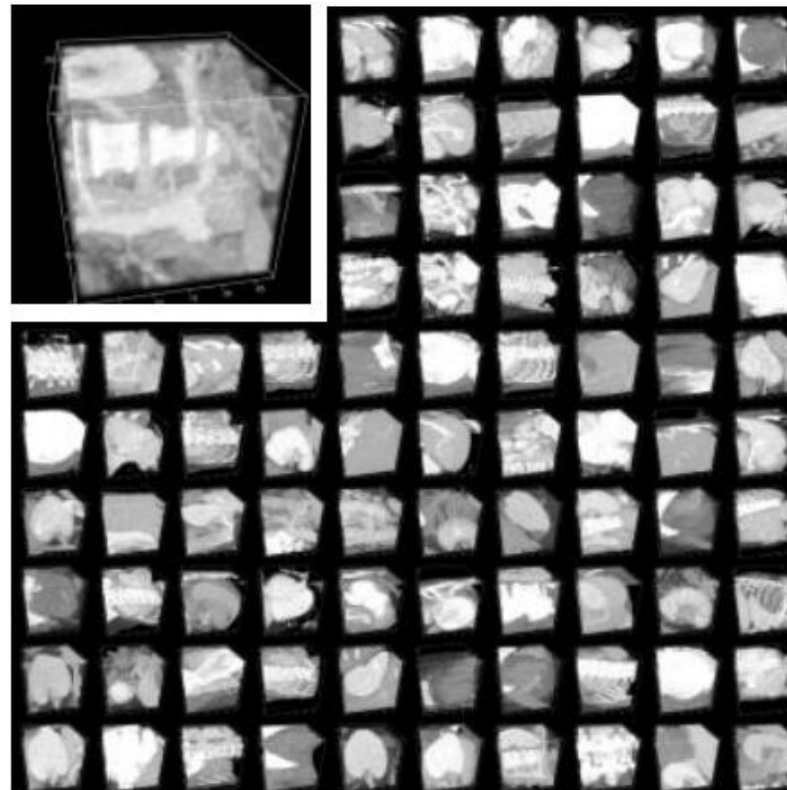# Machine-oriented compression

> *What are ideal characteristics of the compression system?*
>
> - **Support many modalities**

- Hyperspectral

- Spatial Audio

- **3D volumes, medical images**

# Machine-oriented compression

*What are ideal characteristics of the compression system?*

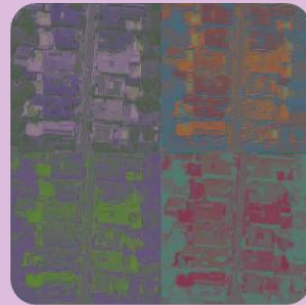- Support many modalities
- **Allow efficient encoding**



**Sensor**

Original Signal

Enc.

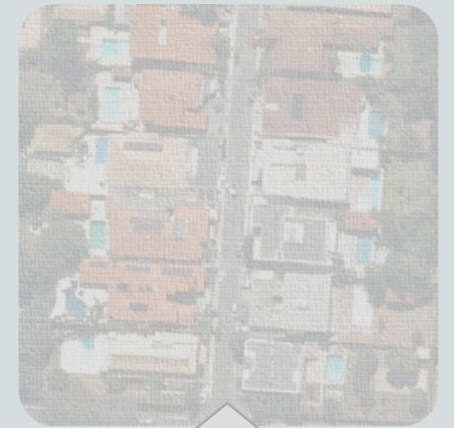Machine-interpretable features

Remote/Cloud

ML Applications

Classification

Segmentation

Enhancement

Optional decoding

**Less bandwidth**     **Enhanced accuracy**     **More efficient ML**

# Machine-oriented compression

**What are ideal characteristics of the compression system?**

- Support many modalities
- Allow efficient encoding
- **Preserve details**

Generative models synthesize details
For recognition, we must **preserve** details

Original        Stable Diff. VAE

# Machine-oriented compression

*What are ideal characteristics of the compression system?*

- Support many modalities • Allow efficient encoding • Preserve details
- **Achieve high compression rate**



**Sensor**

Machine-interpretable features

**Remote/Cloud**

Enc.

ML Applications

Classification

Segmentation

Enhancement

⋮

Original Signal

Optional decoding

**Less bandwidth**   **Enhanced accuracy**   **More efficient ML**

# Machine-oriented compression

> *What are ideal characteristics of the compression system?*
>
> - Support many modalities   • Allow efficient encoding   • Preserve details
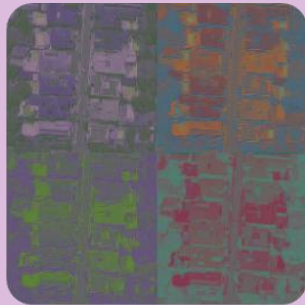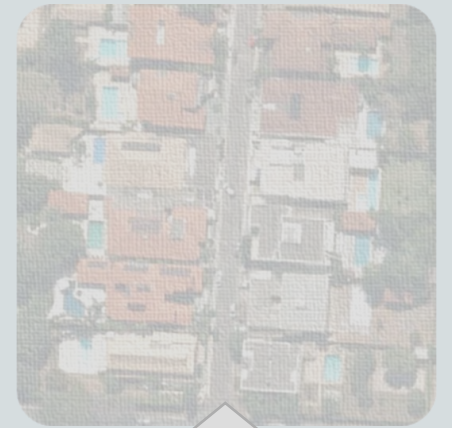> - Achieve high compression rate   • **Accelerate downstream ML models**



Sensor

Machine-interpretable features

Remote/Cloud

Original Signal

ML Applications
- Classification
- Segmentation
- Enhancement
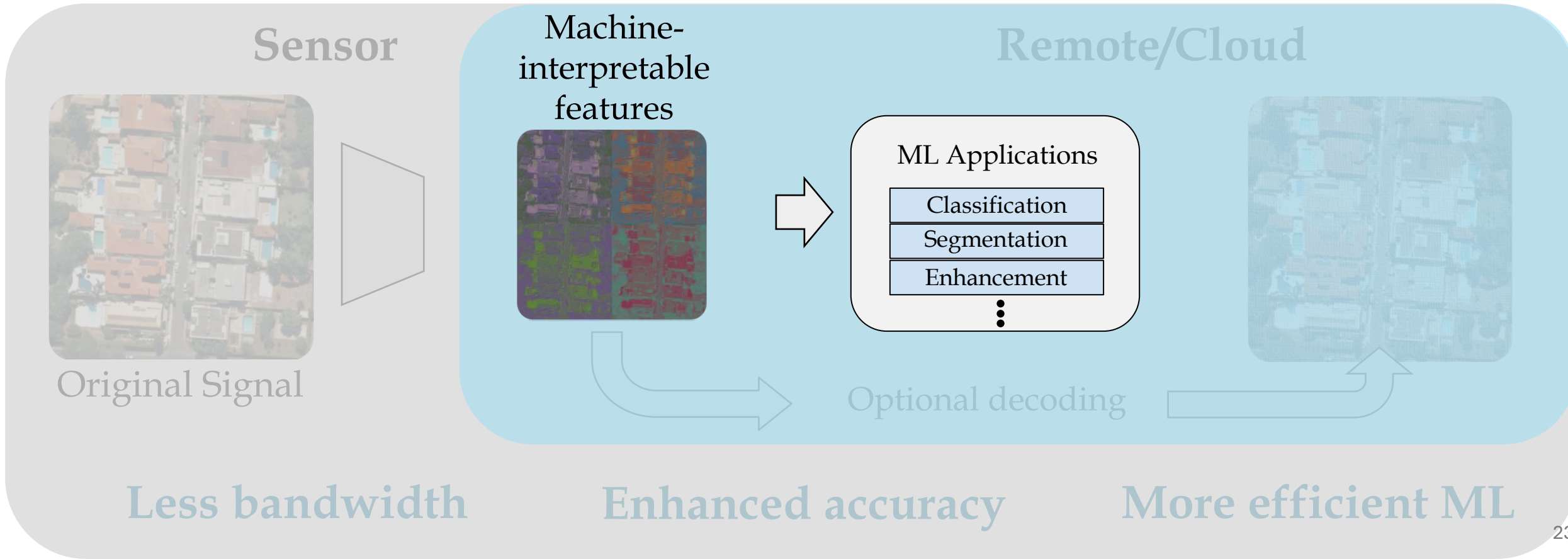
Optional decoding

**Less bandwidth**       **Enhanced accuracy**       **More efficient ML**

# Machine-oriented compression

What are ideal characteristics of the compression system?

- Support many modalities
- Allow efficient encoding
- Preserve details
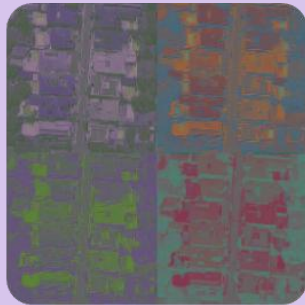- Achieve high compression rate
- Accelerate downstream ML models



**Sensor**

Original Signal

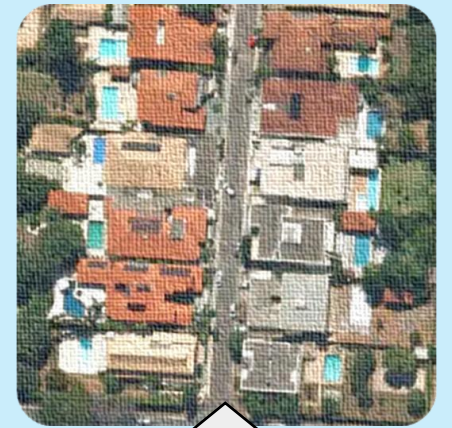Enc.

Machine-interpretable features

**Remote/Cloud**
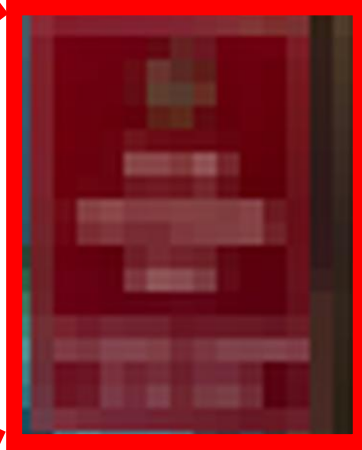
ML Applications
- Classification
- Segmentation
- Enhancement
⋮

Optional decoding

**Less bandwidth**   **Enhanced accuracy**   **More efficient ML**

# Comparison of existing codec designs



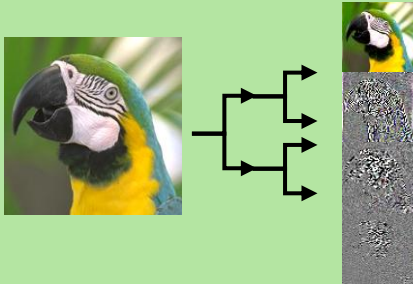Resample     WEBP     DGML (Cheng2020)     Stable Diff. VAE

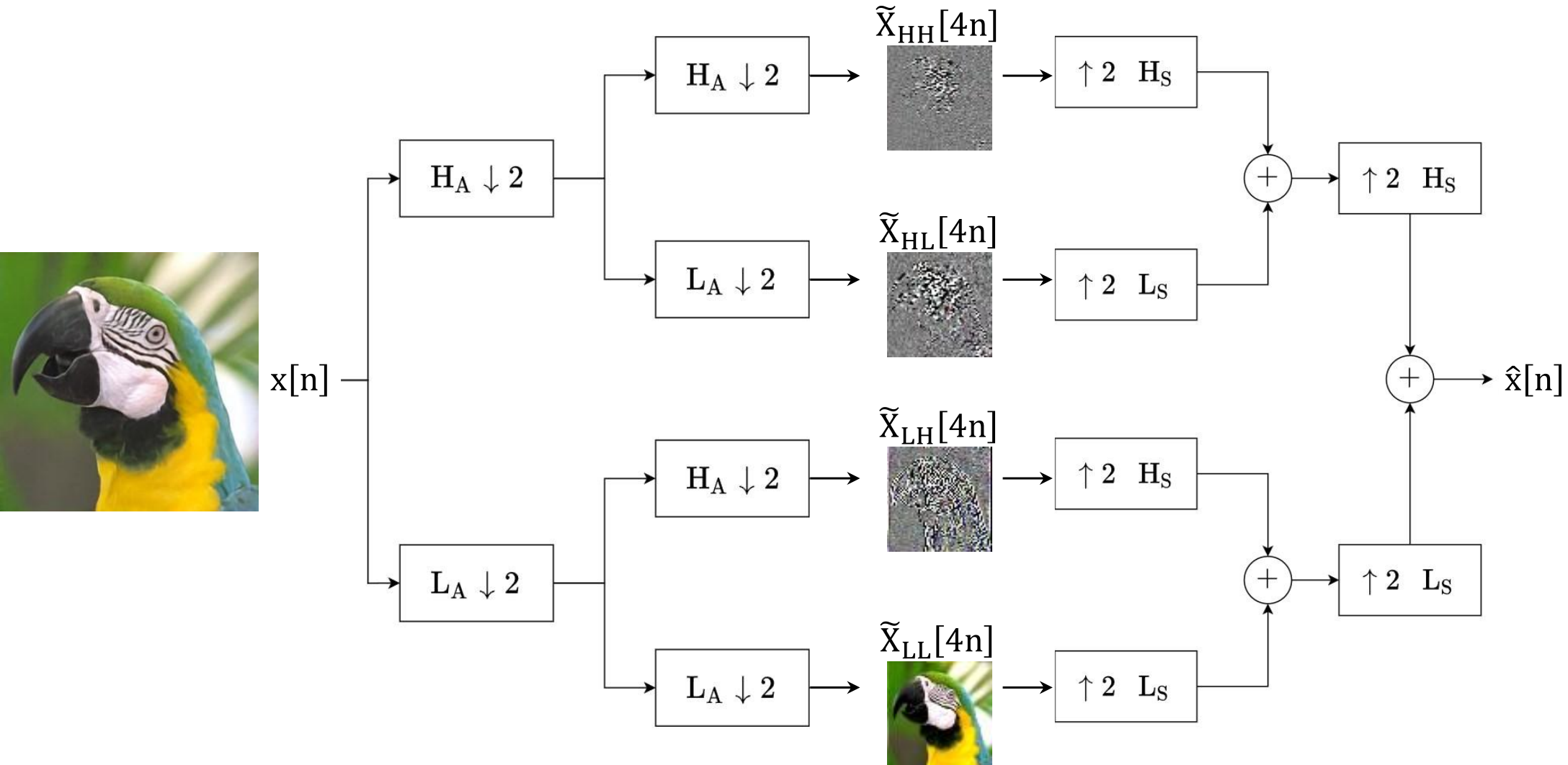| | RR | LTC | E2ELC | GenAE | Goal |
|---|---|---|---|---|---|
| **Allow efficient encoding** | ✅ | ✅ | ❌ | ❌ | ✅ |
| **Accelerate downstream ML** | ✅ | ❌ | ❌ | ✅ | ✅ |
| **Achieve high compression rate** | ❌ | ✅ | ✅ | ❌ | ✅ |
| **Preserve details** | ❌ | ✅ | ✅ | ❌ | ✅ |
| **Support many modalities** | ✅ | ❌ | ✅ | ❌ | ✅ |

# Proposed design

## Encoding efficiency

### Inspired by linear transform coding

Forgo expensive DNN-based analysis transform; leverage efficient, separable transform for energy compaction instead (wavelet packet decomposition)

# Wavelet packet transform

# WPT exchanges spatial resolution with channels



No information loss

Energy compaction

# Proposed design

## Encoding efficiency

### Inspired by linear transform coding

Forgo expensive DNN-based analysis transform; leverage efficient, separable transform for energy compaction instead (wavelet packet decomposition)
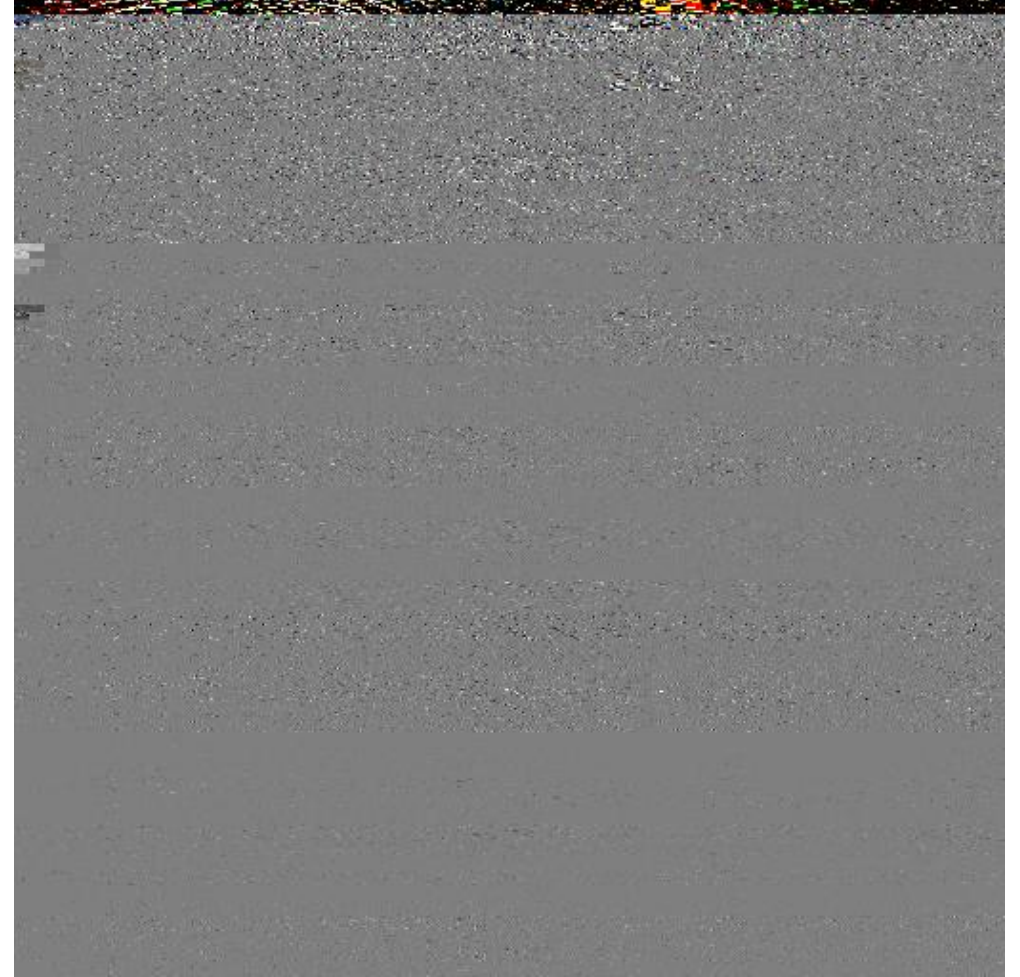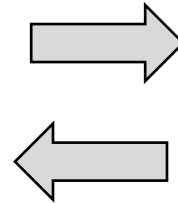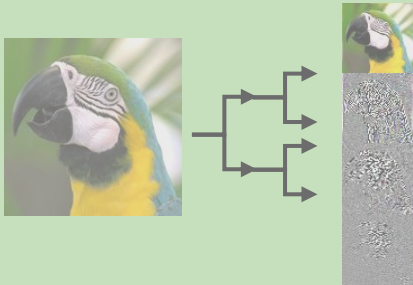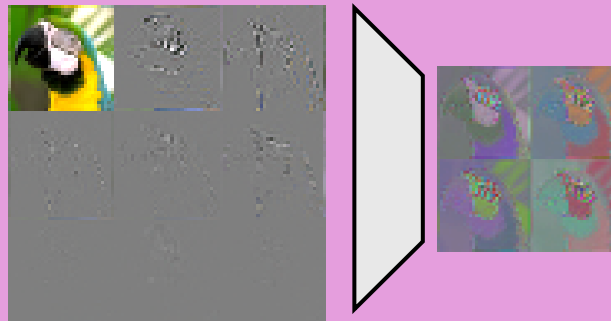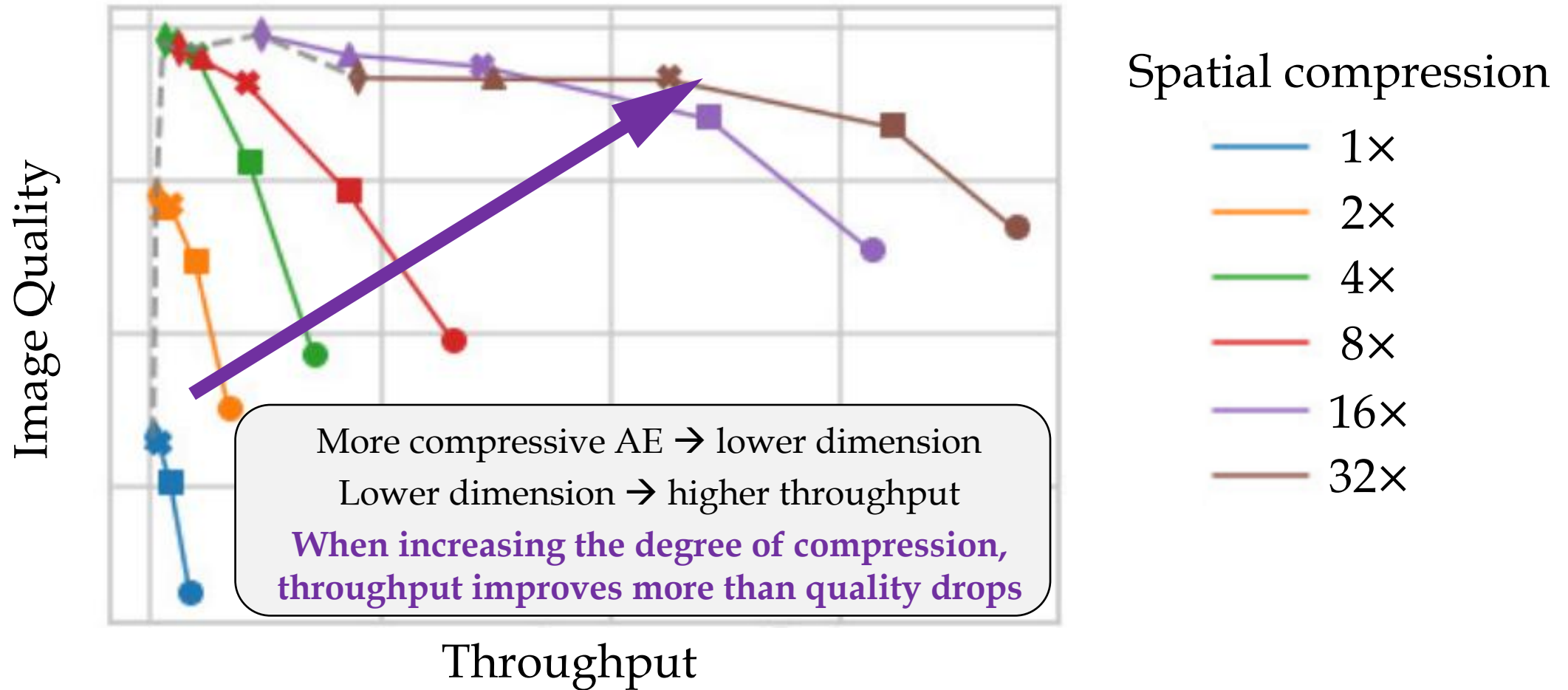


## Dimension reduction

### Inspired by generative autoencoders

Don't rely exclusively on sparsity; use channel bottleneck to provide guaranteed, uniform dimensionality reduction to accelerate downstream models

# Autoencoder for dimension reduction



Spatial compression

— 1×
— 2×
— 4×
— 8×
— 16×
— 32×

Image Quality

Throughput

More compressive AE → lower dimension
Lower dimension → higher throughput
**When increasing the degree of compression, throughput improves more than quality drops**

"High-Resolution Image Synthesis with Latent Diffusion Models"
(aka "Stable Diffusion")  Rombach et al. 2021

# Autoencoder for dimension reduction



113× more expensive than WEBP
34M param. DNN

Encoder (Lossy)

Decoder (Generative)

Original Image
(3×512×512)

48× lower dimension

"Latent" Image
(4×64×64)

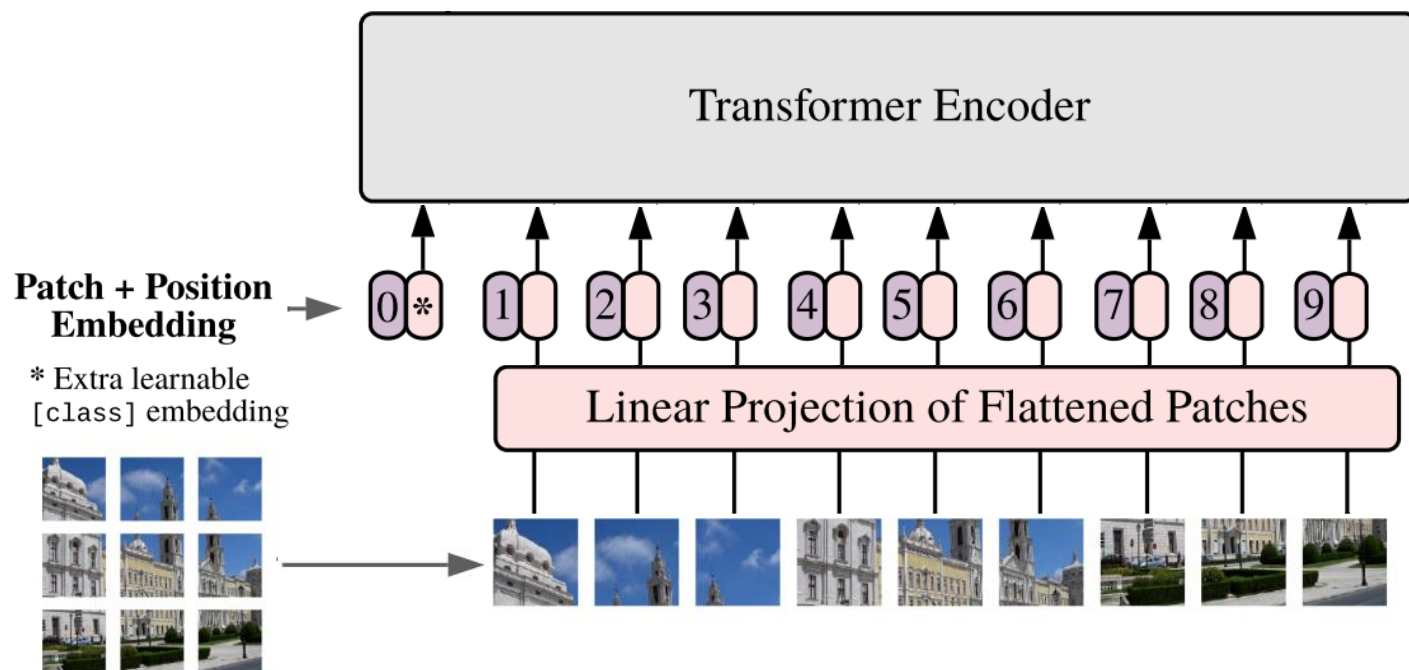Decoded Image
(3×512×512)

"High-Resolution Image Synthesis with Latent Diffusion Models"
(aka "Stable Diffusion")  Rombach et al. 2021

# Does the encoder need to be so expensive?

Synthesizing details is hard
Discarding details is easy
→Use a simple encoder (e.g. linear projection)



Patch + Position Embedding
* Extra learnable [class] embedding

"An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale" (aka "ViT") Beyer et al. 2021

## ViT-B/16

| Patch size | 3×16×16 |
|---|---|
| Sequence Len | 196 |
| Embedding Dim | 768 |
| Compression | 1:1 |
| Accuracy | 86.1 |

## ViT-B/32

| Patch size | 3×32×32 |
|---|---|
| Sequence Len. | 49 |
| Embedding Dim | 768 |
| Compression | 4:1 |
| Accuracy | 83.3 |

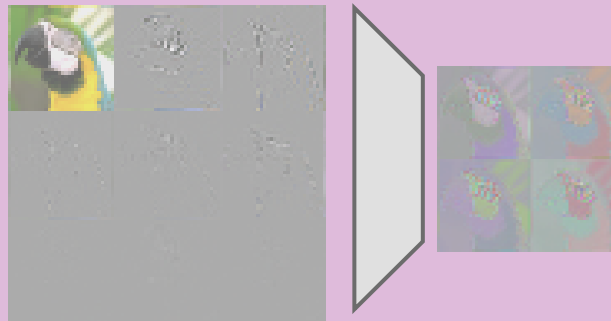# Proposed design

## Inspired by linear transform coding

Forgo expensive DNN-based analysis transform; leverage efficient, separable transform for energy compaction instead (wavelet packet decomposition)
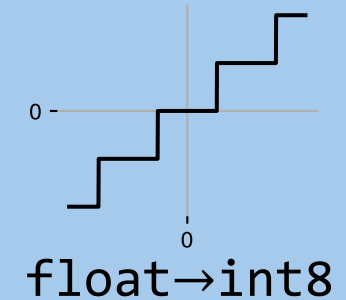


## Inspired by generative AEs

Don't rely exclusively on sparsity; use channel bottleneck to provide guaranteed, uniform dimensionality reduction to accelerate downstream models
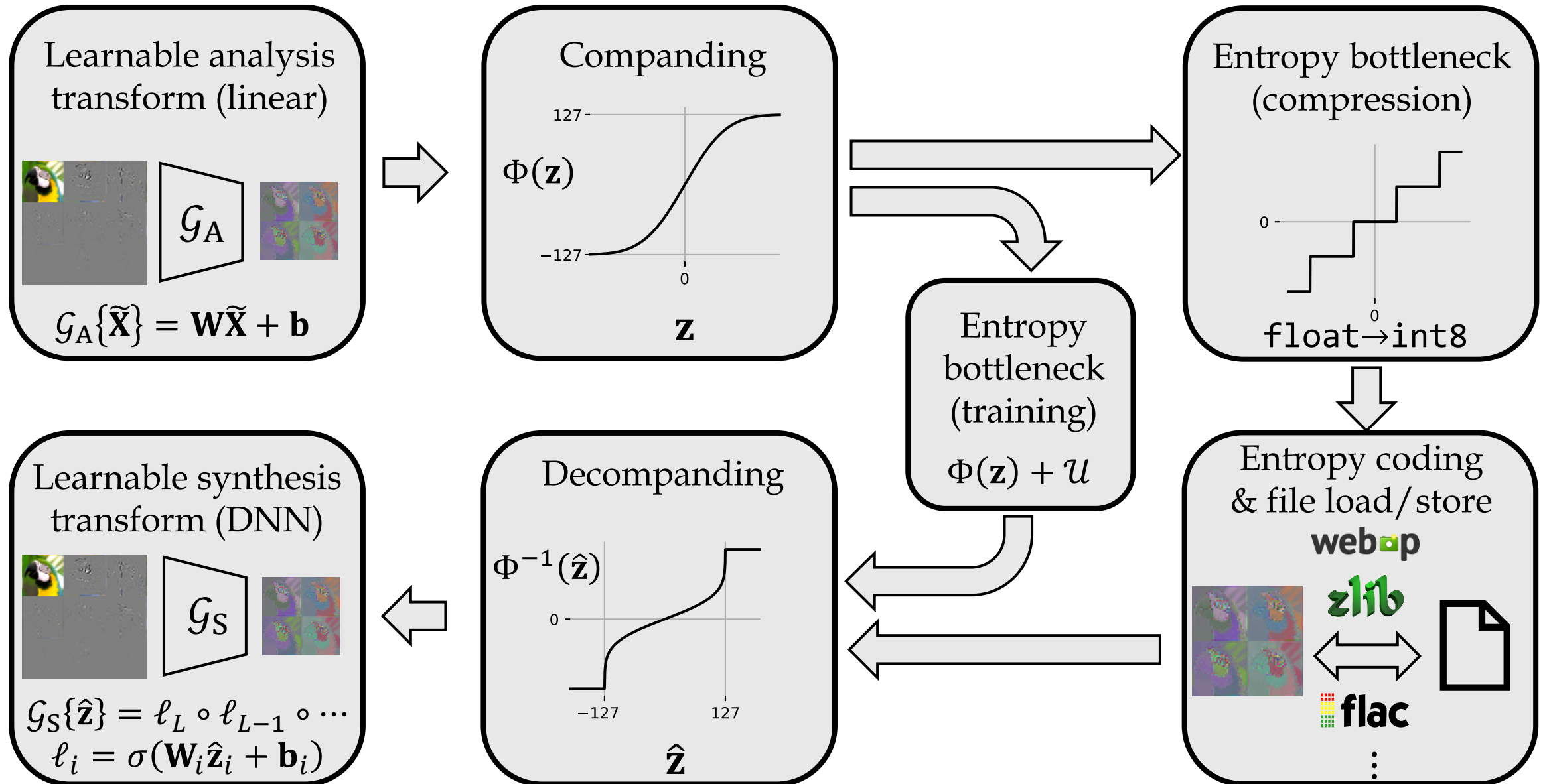


## Inspired by E2E learned compression

Guarantee resilience to quantization via additive noise during training. Leverage existing lossless codecs as a compression multiplier.



float→int8
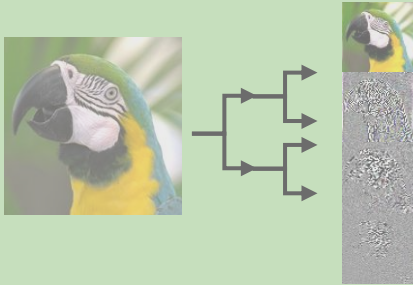
33

# E2E learned compression: quantization and entropy coding

**Learnable analysis transform (linear)**

$\mathcal{G}_A\{\widetilde{\mathbf{X}}\} = \mathbf{W}\widetilde{\mathbf{X}} + \mathbf{b}$

**Companding**

$\Phi(\mathbf{z})$

127

−127

0

$\mathbf{z}$

**Entropy bottleneck (compression)**

0

float→int8

**Entropy bottleneck (training)**

$\Phi(\mathbf{z}) + \mathcal{U}$

**Learnable synthesis transform (DNN)**

$\mathcal{G}_S\{\hat{\mathbf{z}}\} = \ell_L \circ \ell_{L-1} \circ \cdots$
$\ell_i = \sigma(\mathbf{W}_i\hat{\mathbf{z}}_i + \mathbf{b}_i)$

**Decompanding**

$\Phi^{-1}(\hat{\mathbf{z}})$

0

−127       127

$\hat{\mathbf{z}}$

**Entropy coding & file load/store**

webp

zlib

flac

⋮

# Proposed design

## Encoding efficiency

### Inspired by linear transform coding

Forgo e... ...WHT... ... analysi... ... efficien... ... energy compaction instead (wavelet packet decomposition)

## Dimension reduction

### Inspired by generative AEs

Do... ... ... dimensionality reduction to accelerate downstream models

## Compression ratio

### Inspired by E2E learned compression

Conv... ... ...ili... ... e noise ... ...ge existing lossless codecs as a compression multiplier.

float→int8

**WaLLoC**: **Wa**velet **L**earned **Lo**ssy **C**ompression

35

# WaLLoC workflow

**Encoding efficiency**

**Dimension reduction**

**Compression ratio**



Wavelet packet analysis

$$\widetilde{\mathbf{X}} = \text{WPT}\{\mathbf{x}\}$$

Wavelet packet synthesis

$$\hat{\mathbf{x}} = \text{IWPT}\{\widehat{\widetilde{\mathbf{X}}}\}$$

Learnable analysis transform (linear)

$$\mathcal{G}_A\{\widetilde{\mathbf{X}}\} = \mathbf{W}\widetilde{\mathbf{X}} + \mathbf{b}$$

Learnable synthesis transform (DNN)

$$\mathcal{G}_S\{\hat{\mathbf{z}}\} = \ell_L \circ \ell_{L-1} \cdots$$
$$\ell_i = \sigma(\mathbf{W}_i\hat{\mathbf{z}}_i + \mathbf{b}_i)$$

Entropy bottleneck (training)

$$\Phi(\mathbf{z}) + \mathcal{U}$$

Entropy bottleneck (compression)

$$0$$

`float→int8`

Entropy coding & file load/store

web**o**p

z**l**ib

flac

# How to avoid the pitfalls of generative autoencoders?

Resample

WEBP

DGML (Cheng2020)

Stable Diff. VAE

|  | RR | LTC | E2ELC | GenAE | Goal |
|---|---|---|---|---|---|
| Allow efficient encoding | ✅ | ✅ | ❌ | ❌ | ✅ |
| Accelerate downstream ML | ✅ | ❌ | ❌ | ✅ | ✅ |
| Achieve high compression rate | ❌ | ✅ | ✅ | ❌ | ✅ |
| **Preserve details** | ❌ | ✅ | ✅ | ❌ | ✅ |
| **Support many modalities** | ✅ | ❌ | ✅ | ❌ | ✅ |

# Loss function

$$\mathcal{L}(x, \hat{x}) = \underbrace{\text{MSE}(\text{LPF}\{x\}, \text{LPF}\{\hat{x}\})}_{\substack{\text{Pooled MSE} \\ \text{(does not penalize high frequencies)}}} + \underbrace{\mathcal{L}_{\text{LPIPS}}(x, \hat{x})}_{\substack{\text{Lerned perceptual} \\ \text{patch similarity}}} + \underbrace{\mathcal{L}_{\text{GAN}}(x, \hat{x})}_{\substack{\text{Adversarial loss using} \\ \text{VGG16 discriminator}}}$$

**Only preserves low frequency details**

**Requires pre-trained models specific to RGB images**

$$\mathcal{L}(x, \hat{x}) = \text{MSE}(x, \hat{x})$$

$$\hat{x} = \text{decode}(\text{encode}(x) + \mathcal{U})$$

**Better preservation of high frequency details**

**Supports a wide range of modalities**

"Training VQGAN and VAE, with detailed explanation"
S. Ryu, 2024. github.com/cloneofsimo/vqgan-training

# Comparison of autoencoder designs (RGB image)

# Comparison of autoencoder designs (stereo audio)

# How does it perform on downstream applications?

## Image Classification



⇨ Cat

## Document Understanding



Q: What is the date mentioned in the second table?

A: [ "05-12-92" ]

## Colorization



## Source Separation

# Comparison vs. resolution reduction

Reduced resolution & reduced computation
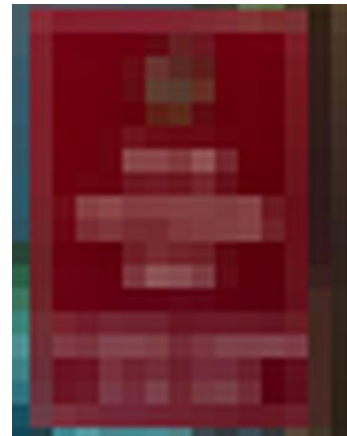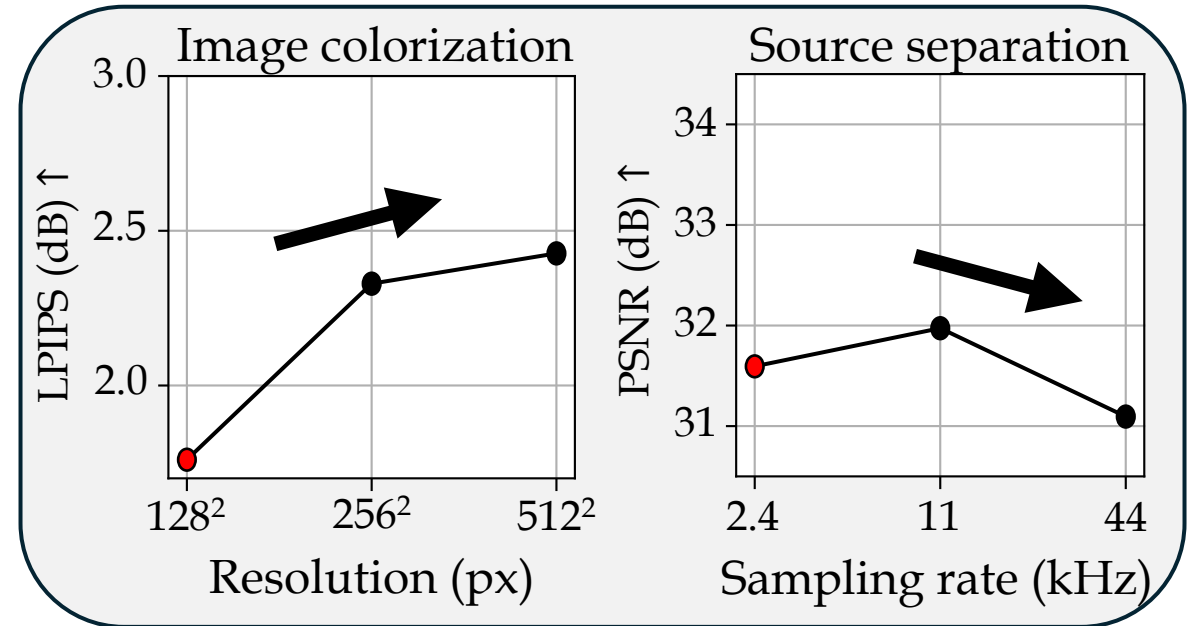


**4× lower latency**

**21GB→8GB GPU Mem**

**85% → 44% Accuracy**

Baseline ●

Resample ●

# Comparison vs. resolution reduction

Increased resolution & fixed computation



**Diminishing return of larger patches / filters**

Image colorization

Source separation

Baseline 🔴

Larger patches ●

# Comparison vs. resolution reduction

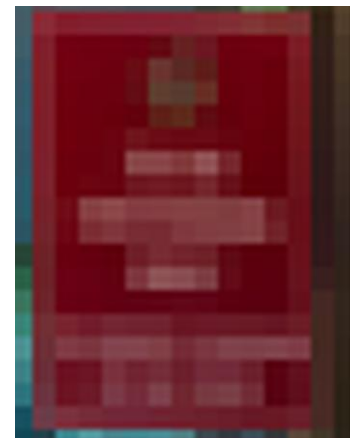Reduced resolution & reduced computation

Increased resolution & fixed computation



Baseline 🔴    Pixels ──●──    Ours ──★──

# Visual Comparison



Original
786 KB

JPEG
6 KB

# Visual Comparison



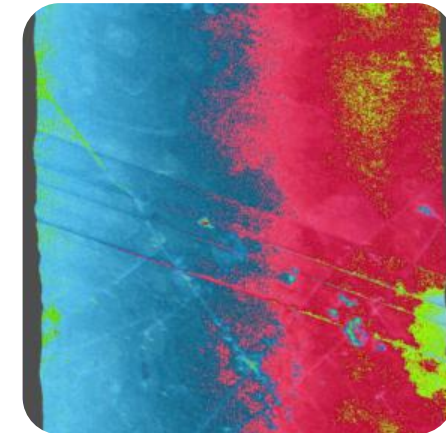Original — 786.43 KB

WaLLoC — 5.5 KB

# Areas for improvement

- Can we make it competitive in terms of the rate-distortion-complexity trade-off?

- How can we support a wider range of specialized signals types and modalities?

- Can we decouple the "generative" part of the decoding process to make it optional?
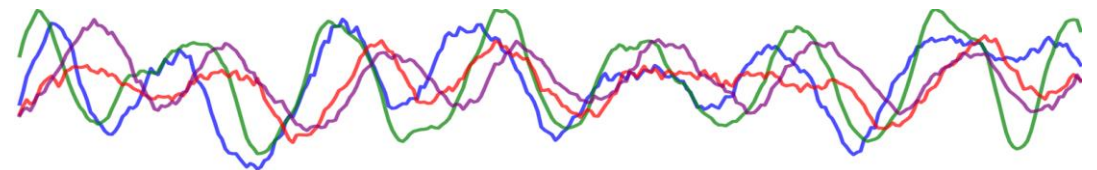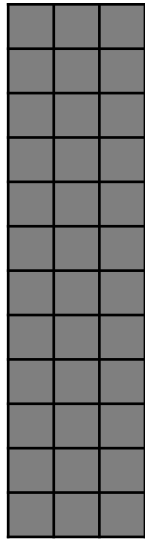
**Medical Images and 3D volumes**

**Hyperspectral & HDR**

**Video**

**Multi-channel & spatial audio**

# Even a single linear projection can be expensive
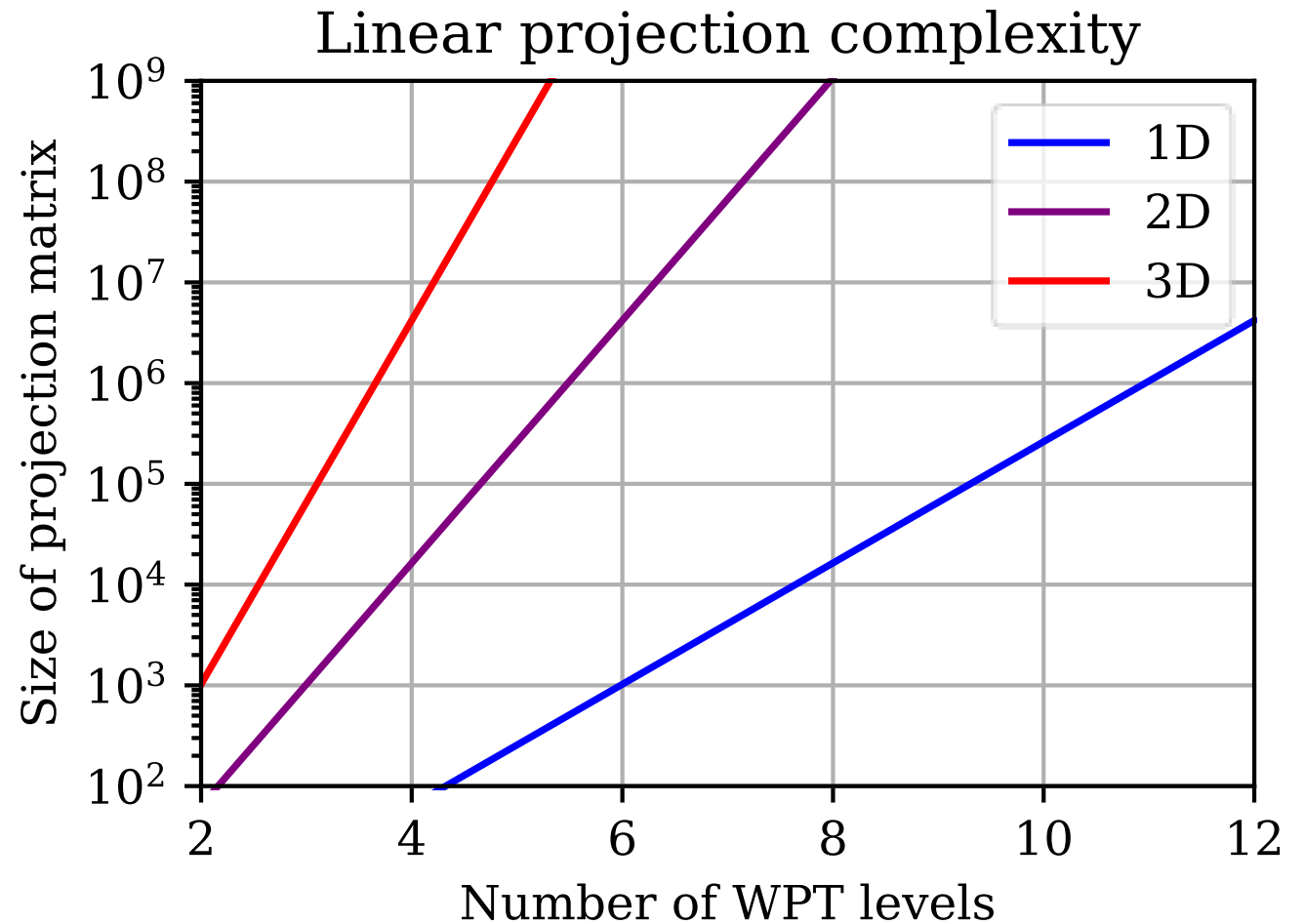
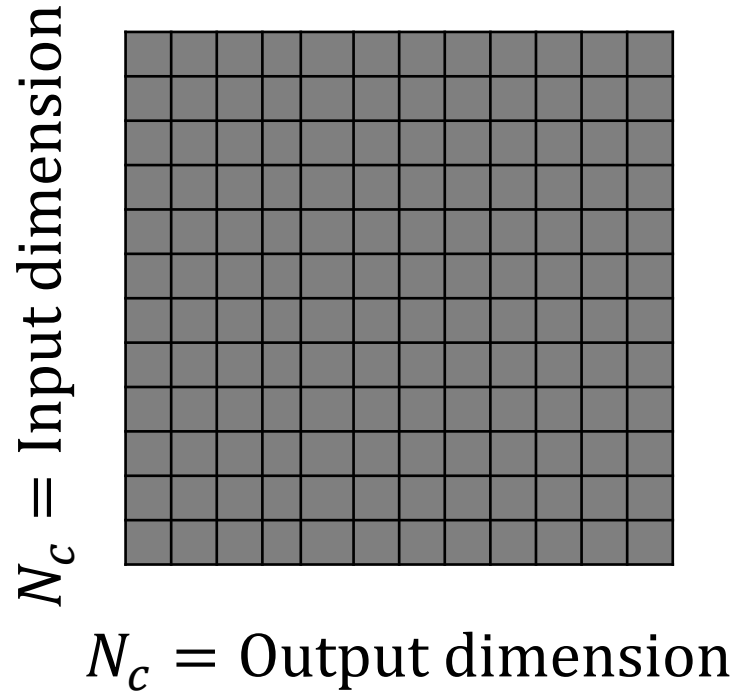**Dense $1 \times 1$ conv. layer**

$N_{\text{in}} = (2^J)^D$
input dimension

$N_{\text{out}} = \dfrac{N_{\text{in}}}{R}$
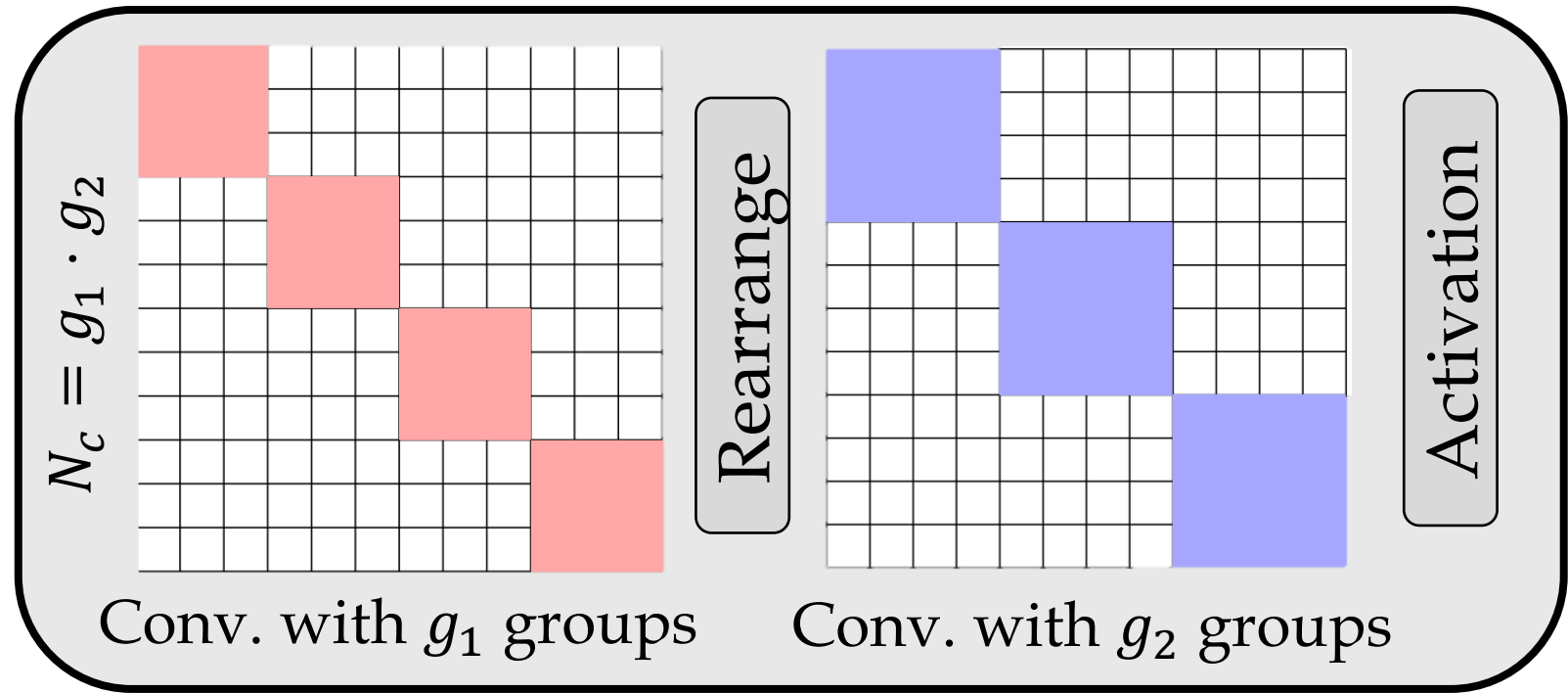
Latent dimension

$\text{MACs} = 2^{2JD}/R$

## Linear projection complexity

Size of projection matrix vs. Number of WPT levels

Legend: 1D, 2D, 3D

# Lightweight, FFT-inspired structured operations

**Dense $1 \times 1$ conv. layer**

$N_C$ = Input dimension

$N_C$ = Output dimension

$$\text{MACs} = N_C^2$$

**Structured $1 \times 1$ conv. layer**



$N_c = g_1 \cdot g_2$

Conv. with $g_1$ groups

Rearrange

Conv. with $g_2$ groups

Activation
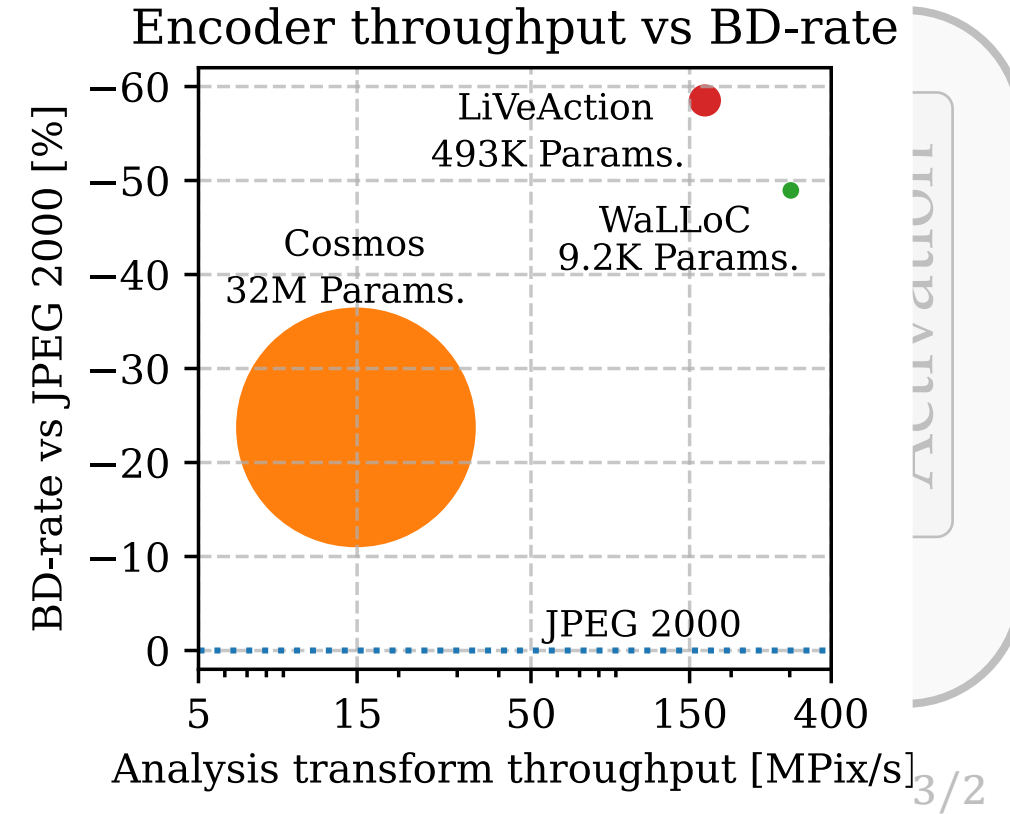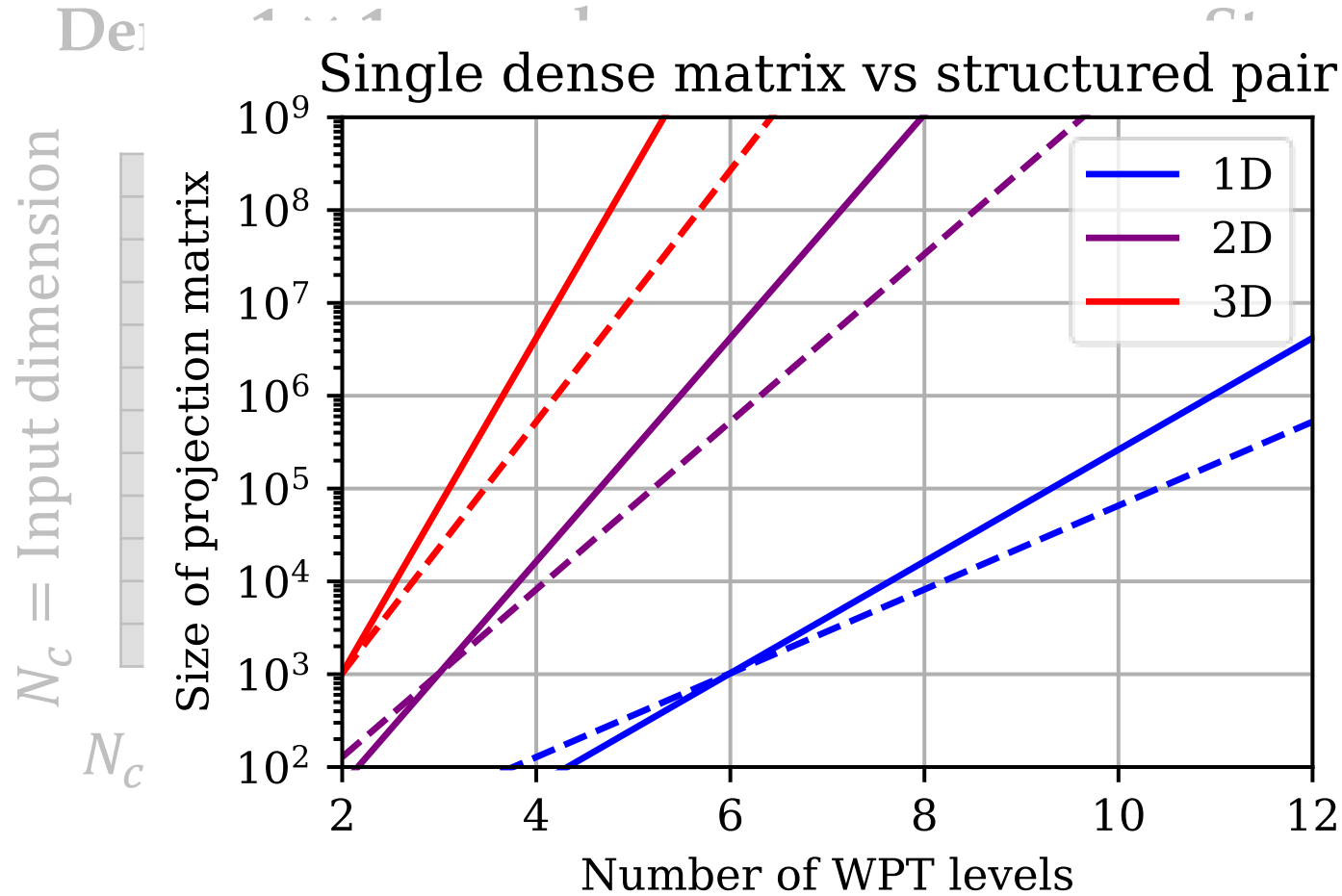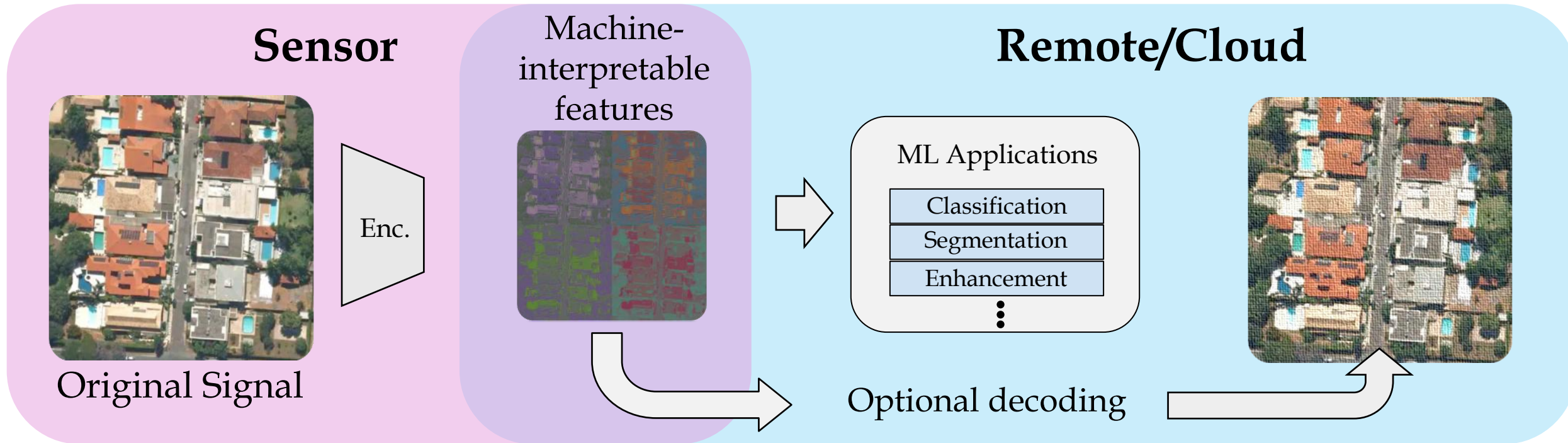
$$g_1 \approx g_2 \approx \sqrt{C_{\text{input}} \cdot 2^{J \cdot D}}$$

$$N_C \log N_C \leq \text{MACs} \leq N_C^{3/2}$$

# Lightweight, FFT-inspired structured operations



Single dense matrix vs structured pair

Encoder throughput vs BD-rate
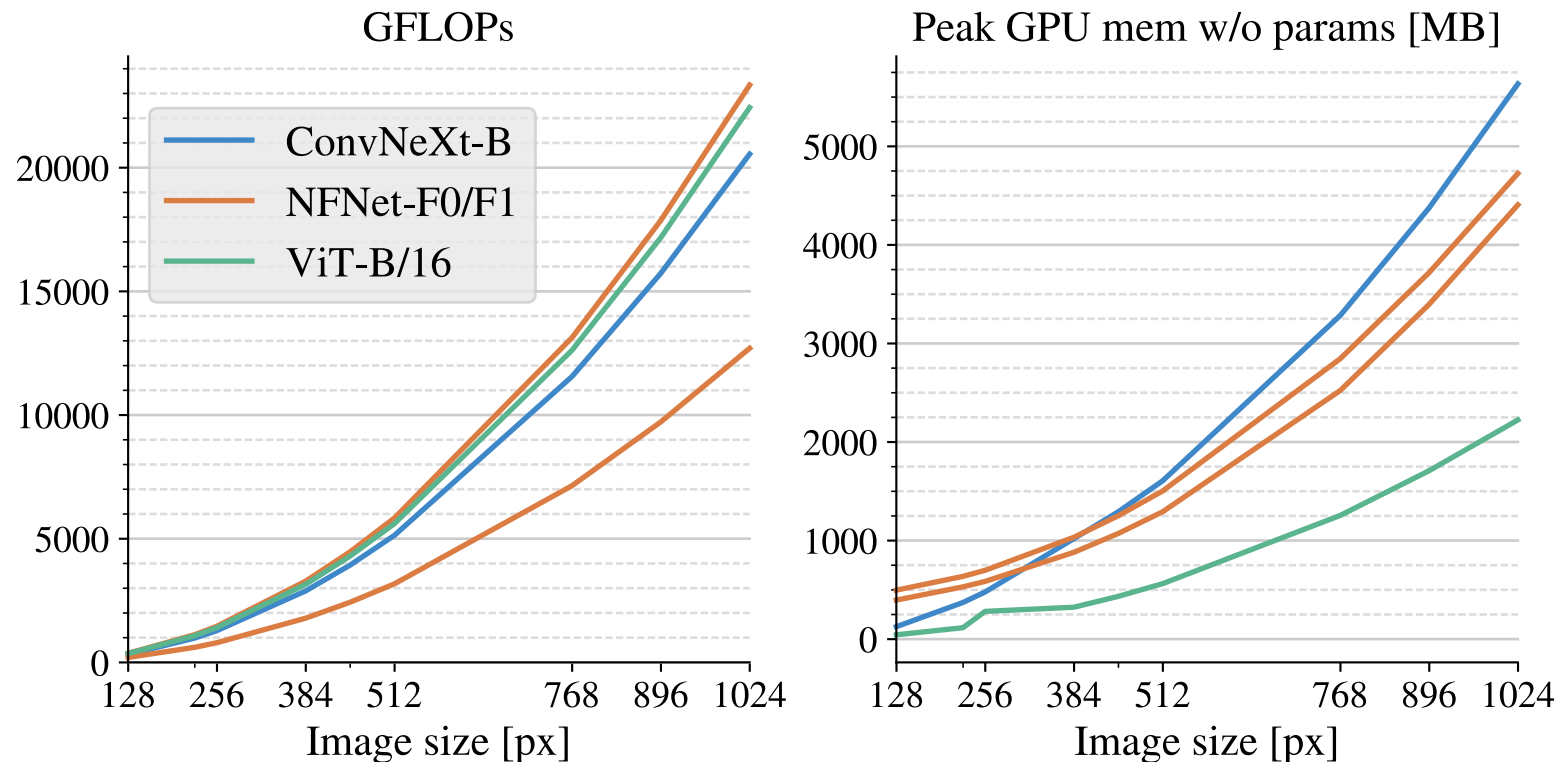
# Asymmetric Design



**Encoding efficiency is very important**

**Decoder can run on cloud AI supercomputers (or throw it away completely)**
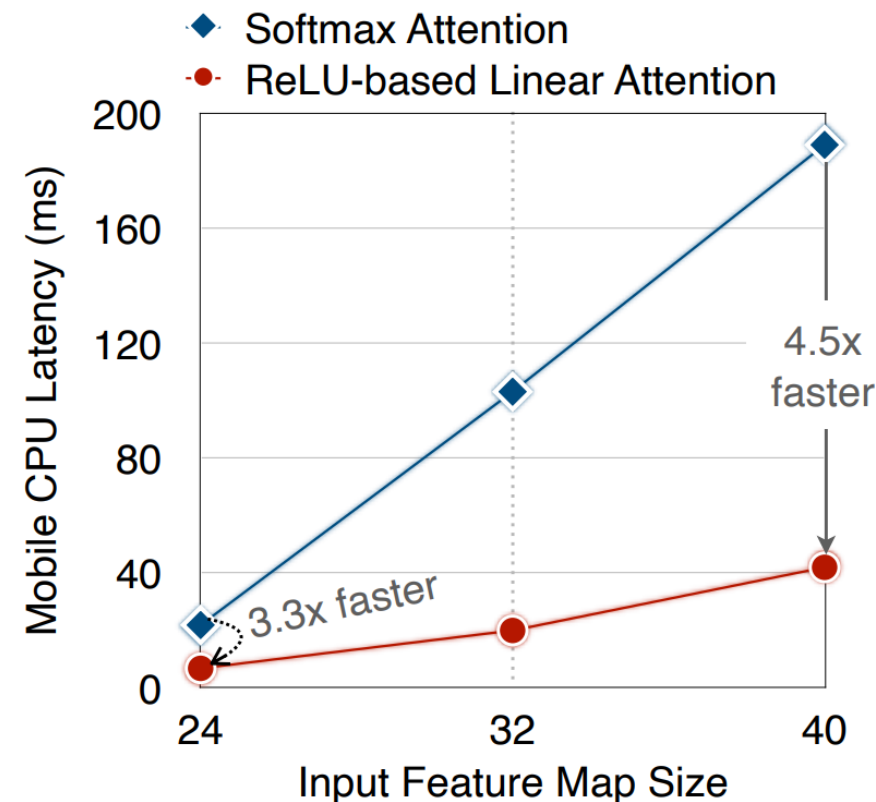
# Improving the synthesis transform

- High complexity synthesis transform is tolerable at runtime

- Training should still be possible in reasonable number of GPU hours

- Support modalities with high spatial and/or temporal resolution



Beyer, Lucas. "On the speed of ViTs and CNNs." (2024).

# ND-generalized ViT decoder with linear attention

- Global receptive field

→Exploit non-local redundancies

- No excessive compute requirements

→Train at high resolution on single GPU

- No position encoding or batch norm

→ Works for any modality



Cai, Han, et al. "Efficientvit: Lightweight multi-scale attention for high-resolution dense prediction." *Proceedings of the IEEE/CVF international conference on computer vision.* 2023.
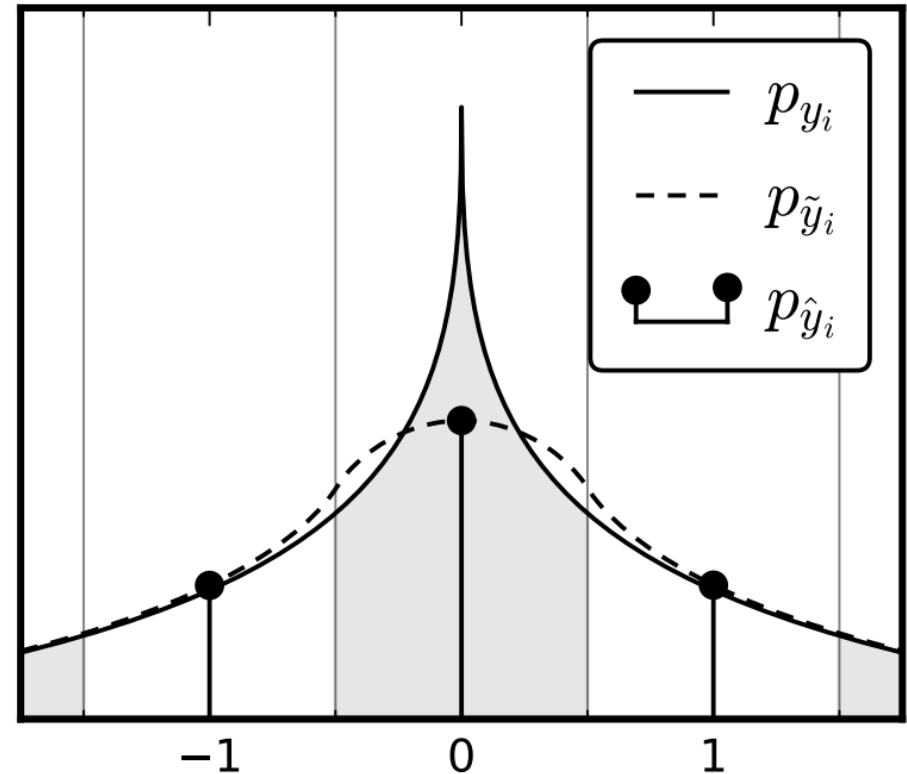
# LiVeAction Overview

- Can we make it competitive in terms of the rate-distortion-complexity trade-off?

- How can we support a wider range of specialized signals types and modalities?

- Can we decouple the "generative" part of the decoding process to make it optional?

- FFT-like structured matrix operations in encoder

- ND-generalized vision transformer decoder with linear attention

- Simplified rate penalty

# Rate objective

$$\min_{\mathcal{E},\mathcal{D}} \underbrace{\|x - \mathcal{D}((\mathcal{E}(x))\|^2}_{\text{MSE Distortion}} + \underbrace{H(\mathcal{E}(x))}_{\text{latent rate}}$$

- Standard approach to optimize $H(\mathcal{E}(x))$ involves fitting a continuous proxy distribution $p_{y_i}$

- Requires an auxiliary optimizer with additional hyperparameters and separate learning rate



Ballé, Jona, Valero Laparra, and Eero P. Simoncelli. "End-to-end Optimized Image Compression." *International Conference on Learning Representations*. 2017.
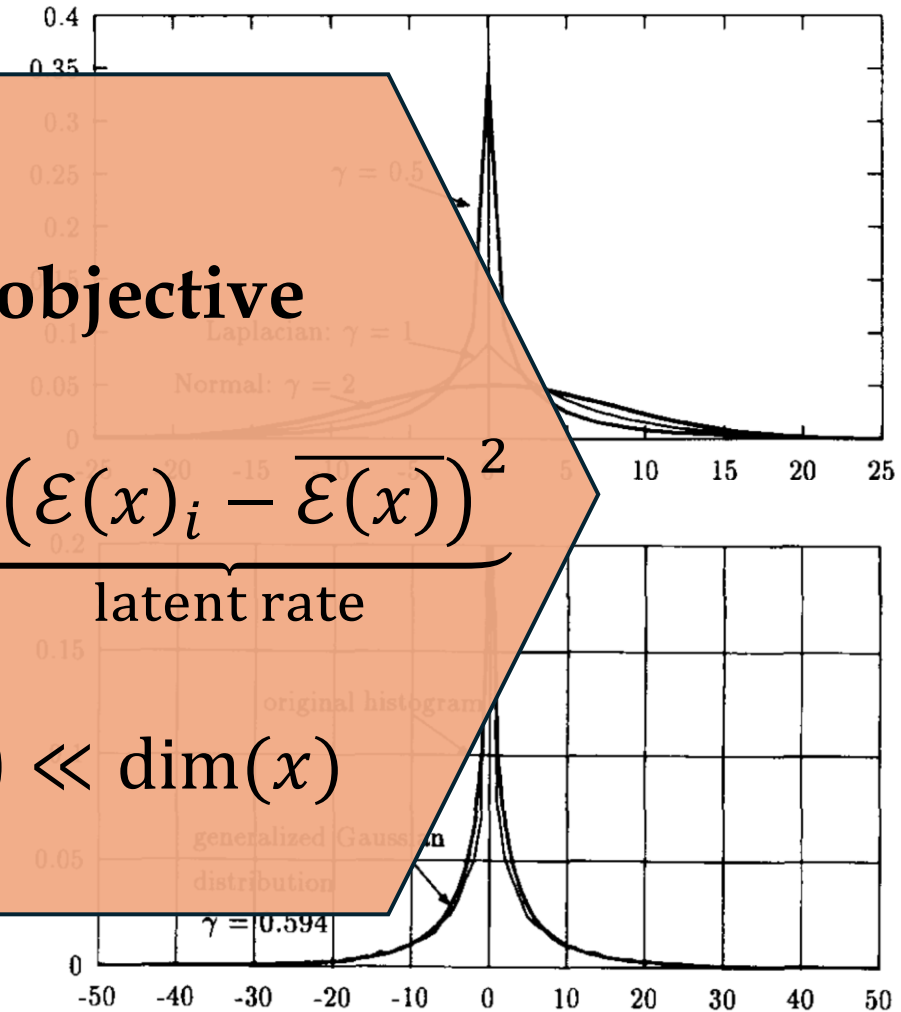
# Simplified rate penalty

- Intensity of sub-band filter outputs follow generalized gaussian distribution for natural signals

- Empirically, analysis transform also follow GGD

- For exponential family, minimizing variance is equivalent to minimizing rate
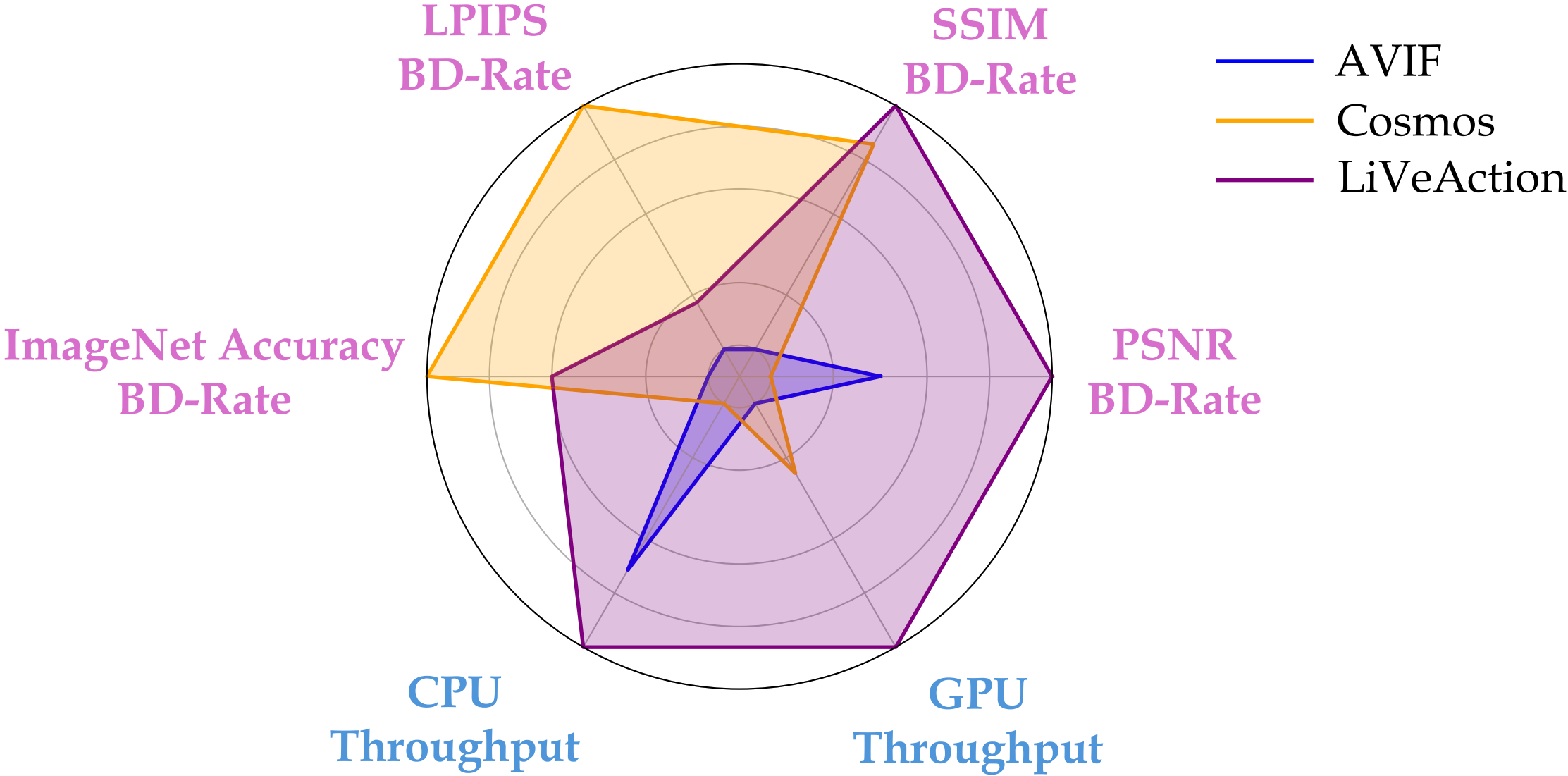
**Simplied training objective**

$$\min_{\mathcal{E},\mathcal{D}} \underbrace{\|x - \mathcal{D}((\mathcal{E}(x))\|^2}_{\text{MSE Distortion}} + \underbrace{\sum(\mathcal{E}(x)_i - \overline{\mathcal{E}(x)})^2}_{\text{latent rate}}$$

$$\text{Subject to } \dim(\mathcal{E}(x)) \ll \dim(x)$$

Sharifi, Karnran, and Alberto Leon-Garcia. "Estimation of shape parameter for generalized Gaussian distributions in subband decompositions of video." *IEEE Transactions on Circuits and Systems for Video Technology* 5.1 (1995)
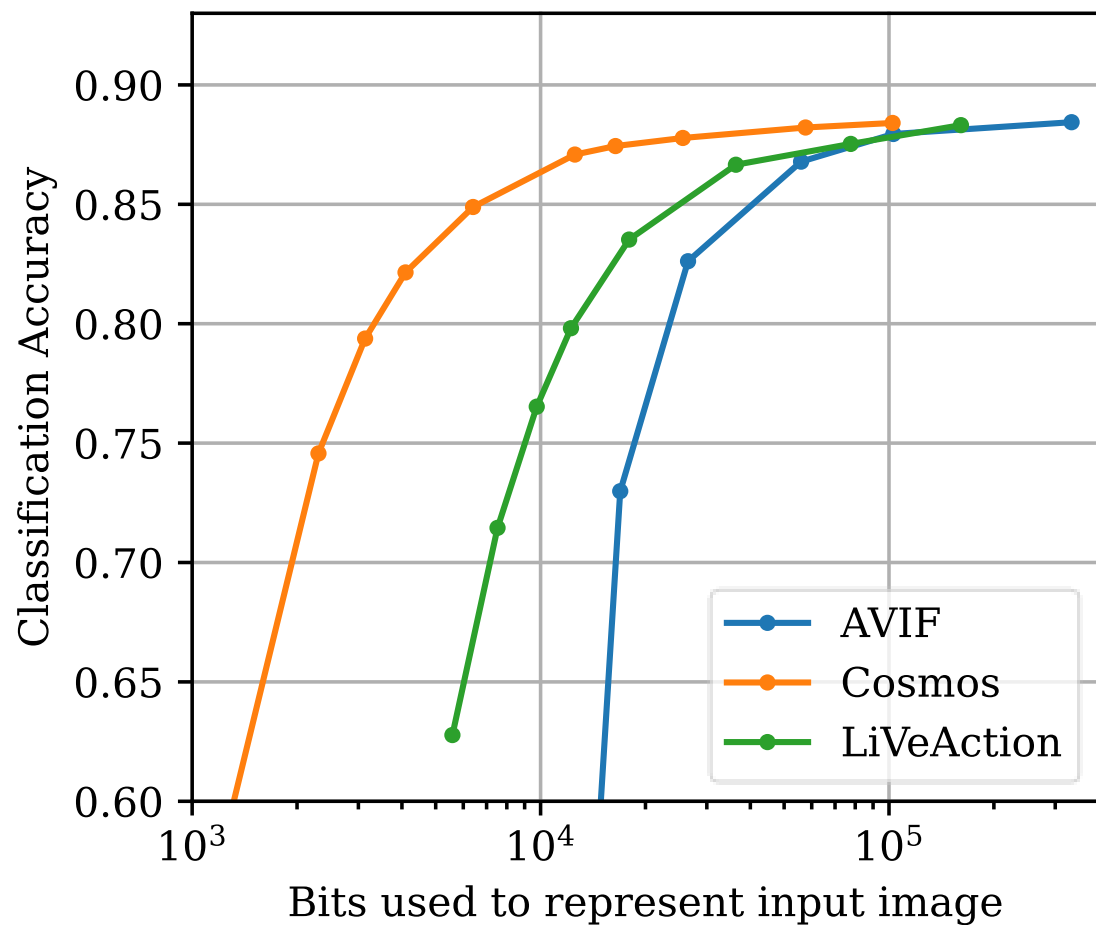
# Compression efficiency and computational efficiency

# Machine Perception

# Other modalities
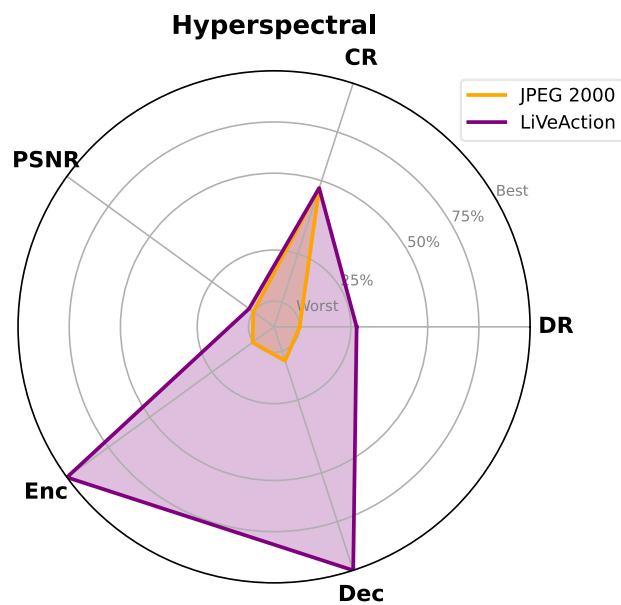
# Generative Enhancement



(a) Original

(b) Cosmos

(c) LiVeAction

(d) LiVeAction + FLUX

# Publications

[1] D. Jacobellis, D. Cummings, and N.J. Yadwadkar. "Machine Perceptual Quality: Evaluating the Impact of Severe Lossy Compression on Audio and Image Models." *Data Compression Conference*. IEEE, 2024.

[2] D. Jacobellis and N.J. Yadwadkar."Learned Compression for Compressed Learning." *Data Compression Conference*. IEEE, 2025.

[3] D. Jacobellis and N.J. Yadwadkar. "LiVeAction: a Lightweight, Versatile, and Asymmetric Neural Codec Design for Real-time Operation." Under Review.

[4] D. Jacobellis, M. Ulhaq, F. Racapé, H. Choi, and N.J. Yadwadkar." Dedelayed: Deleting remote inference delay via on-device correction." Under Review.

# Software releases

Installation → `pip install walloc`

Audio→ [Pre-trained codec](#)

Images→ [Pre-trained codec](#)

Training (1D) → [Tutorial](#)

Training (2D) → [Tutorial](#)

More details available:
[https://ut-sysml.org/walloc/](https://ut-sysml.org/walloc/)

Contact: [danjacobellis@utexas.edu](mailto:danjacobellis@utexas.edu)

# Software releases

Installation → `pip install livecodec`

Dozen+ pre-trained codecs available on Hugging Face
https://hf.co/danjacobellis/liveaction

Training → Tutorial

More details available:
https://ut-sysml.org/liveaction/

Contact: danjacobellis@utexas.edu