# Towards practical neural image compression: SlimCAE and DANICE
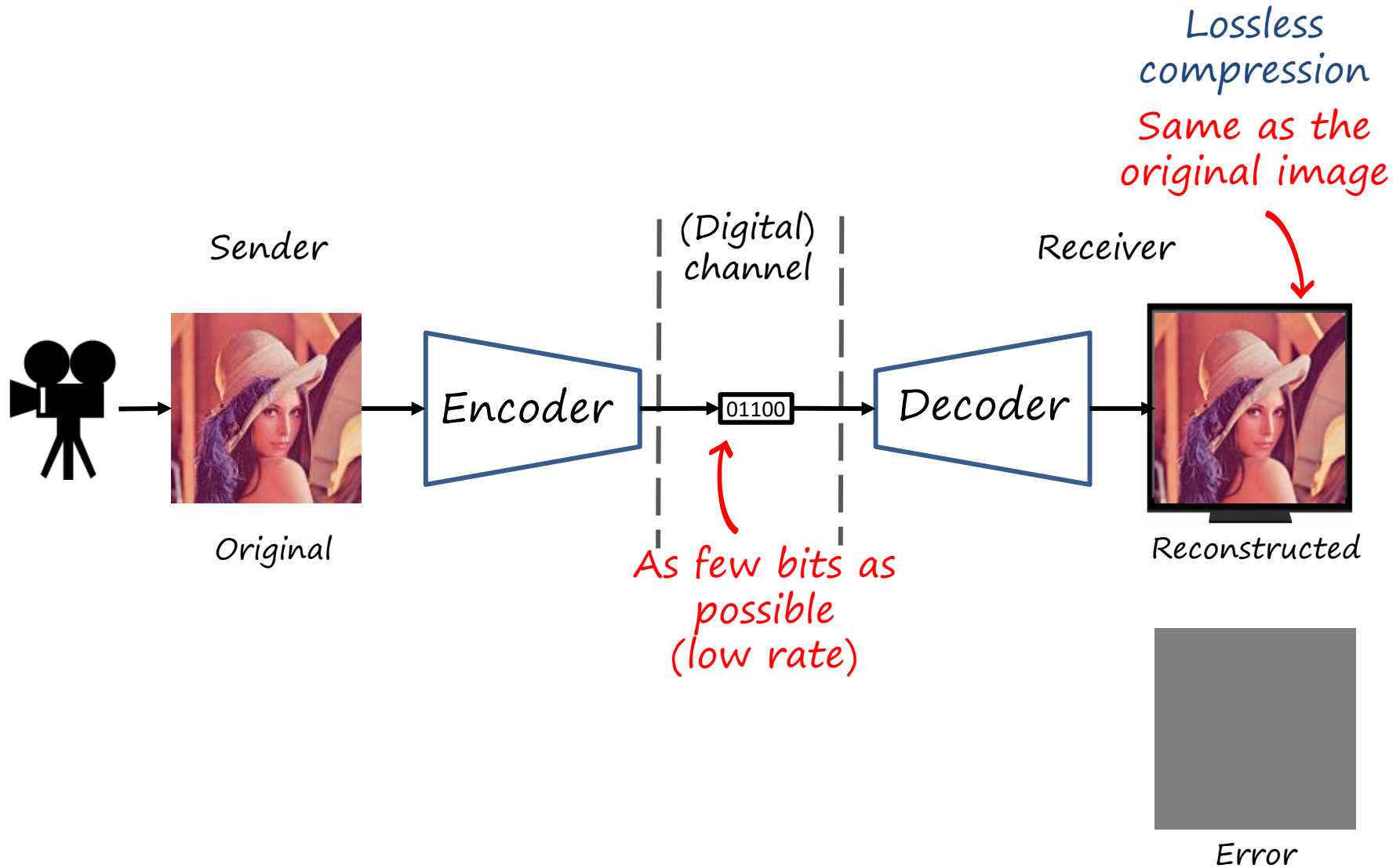
Fei Yang, Luis Herranz

CVPR 2021/CLIC2021
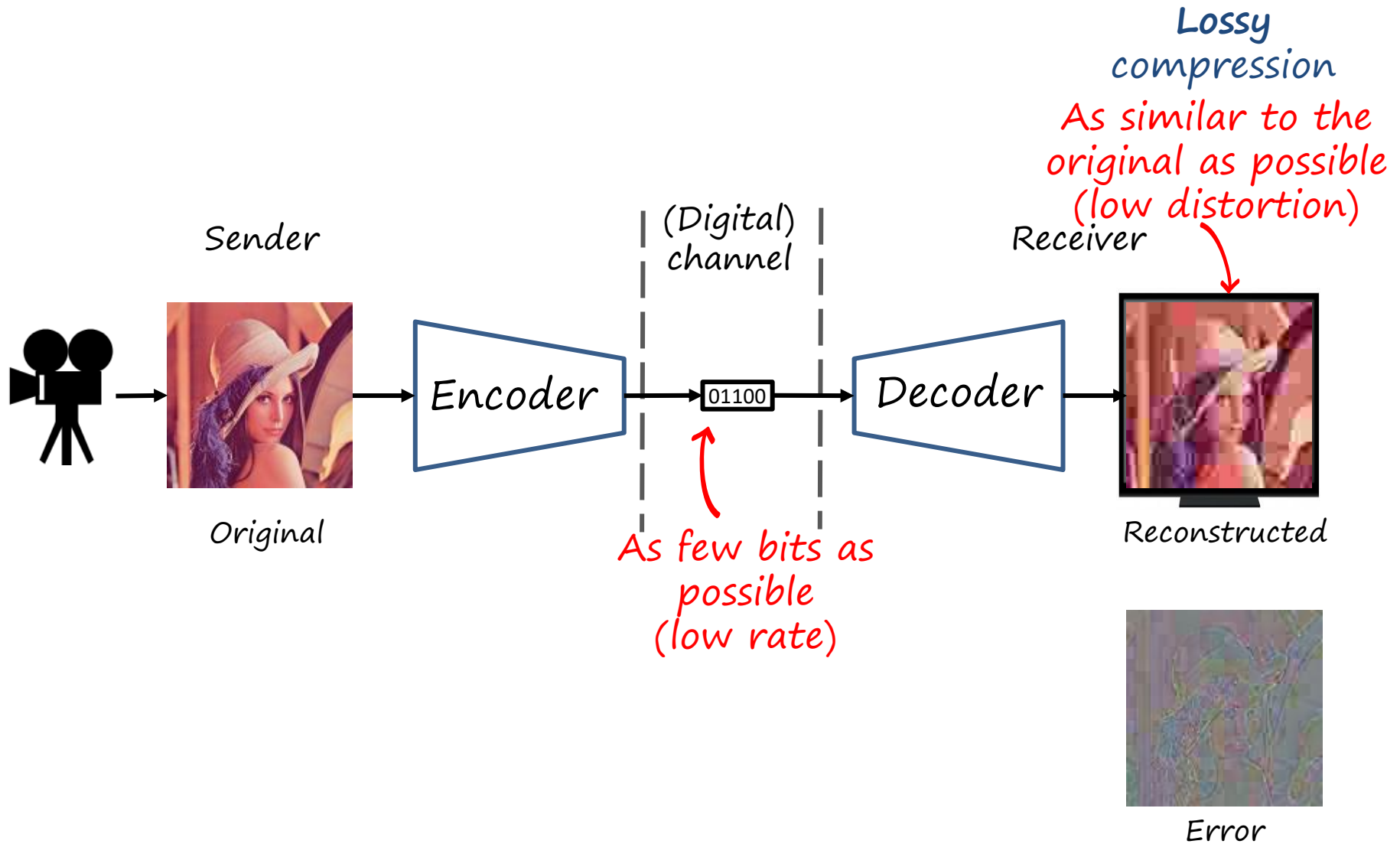
# Outline

- Introduction and motivation

- SlimCAE (CVPR 2021)

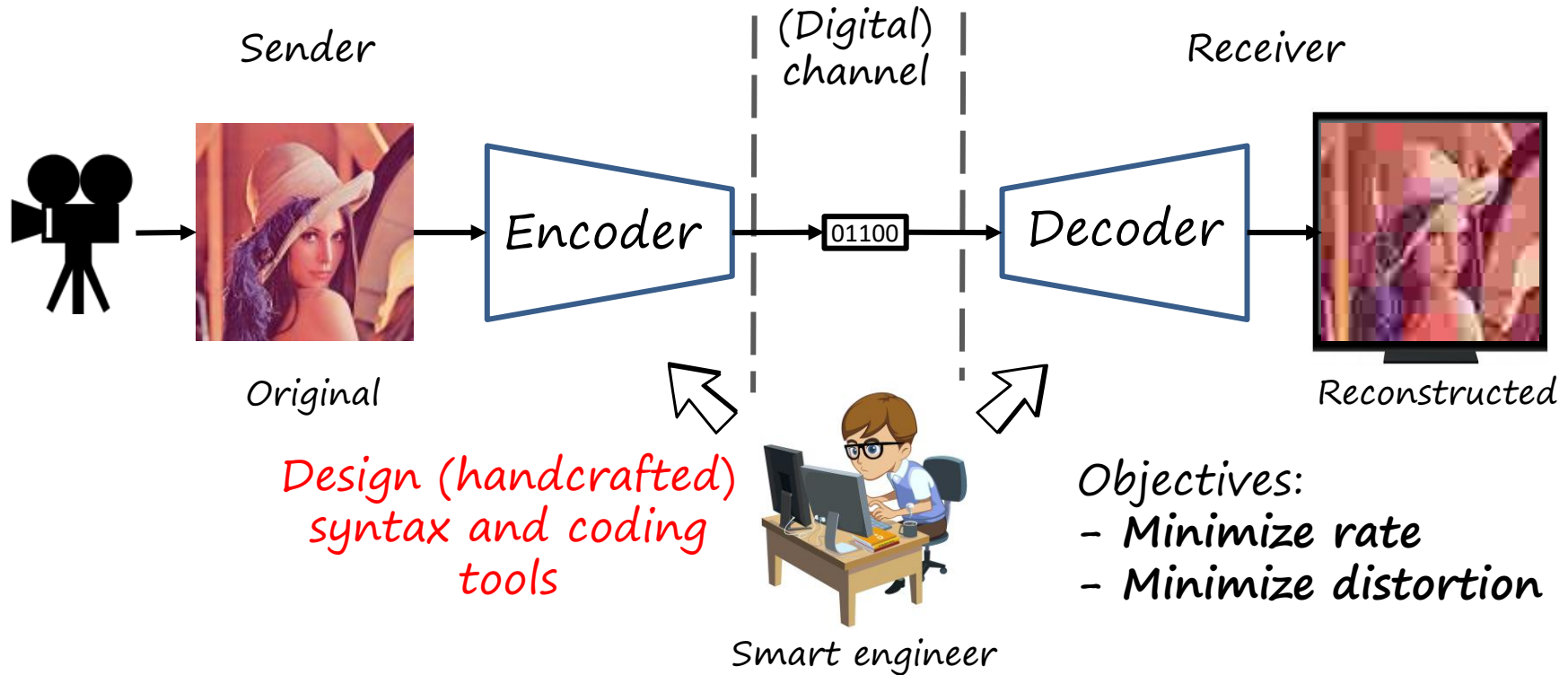- DANICE (CLIC workshop at CVPR 2021)
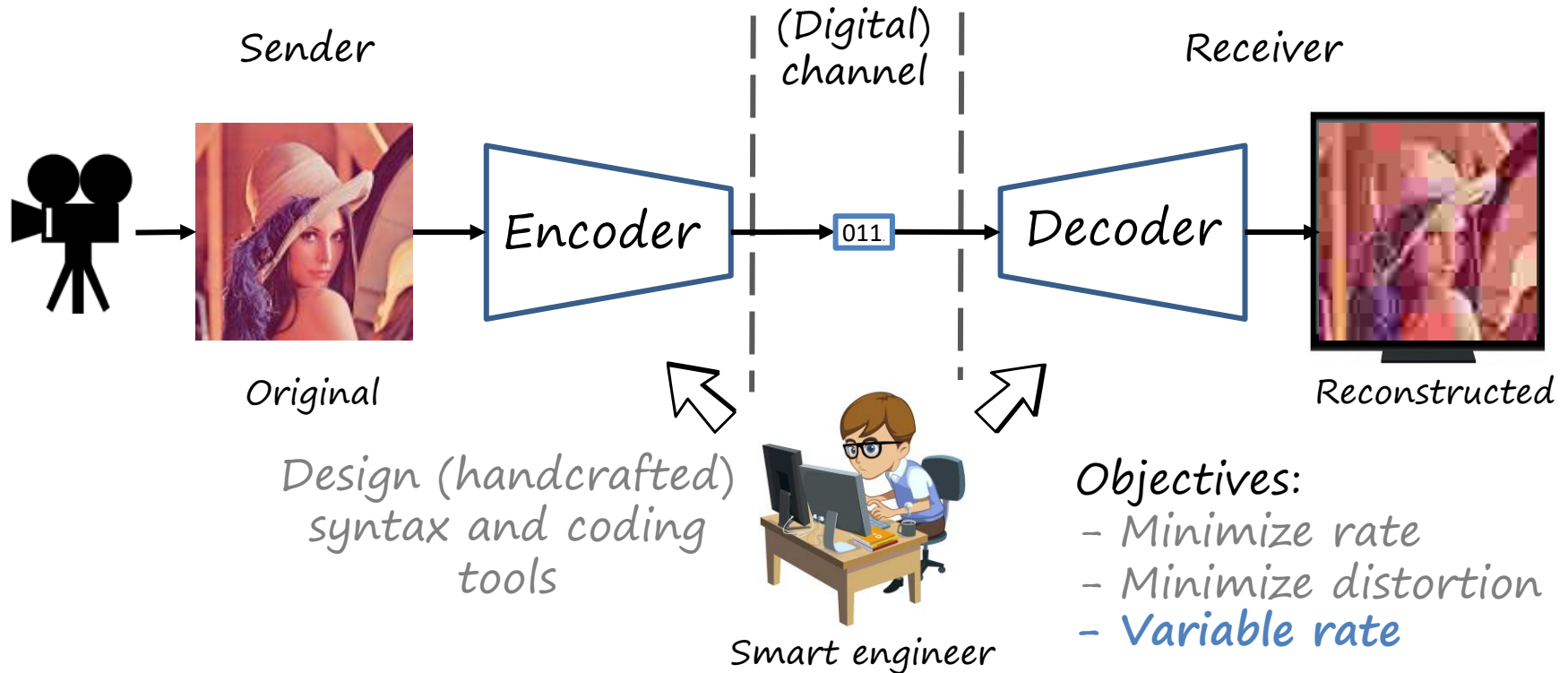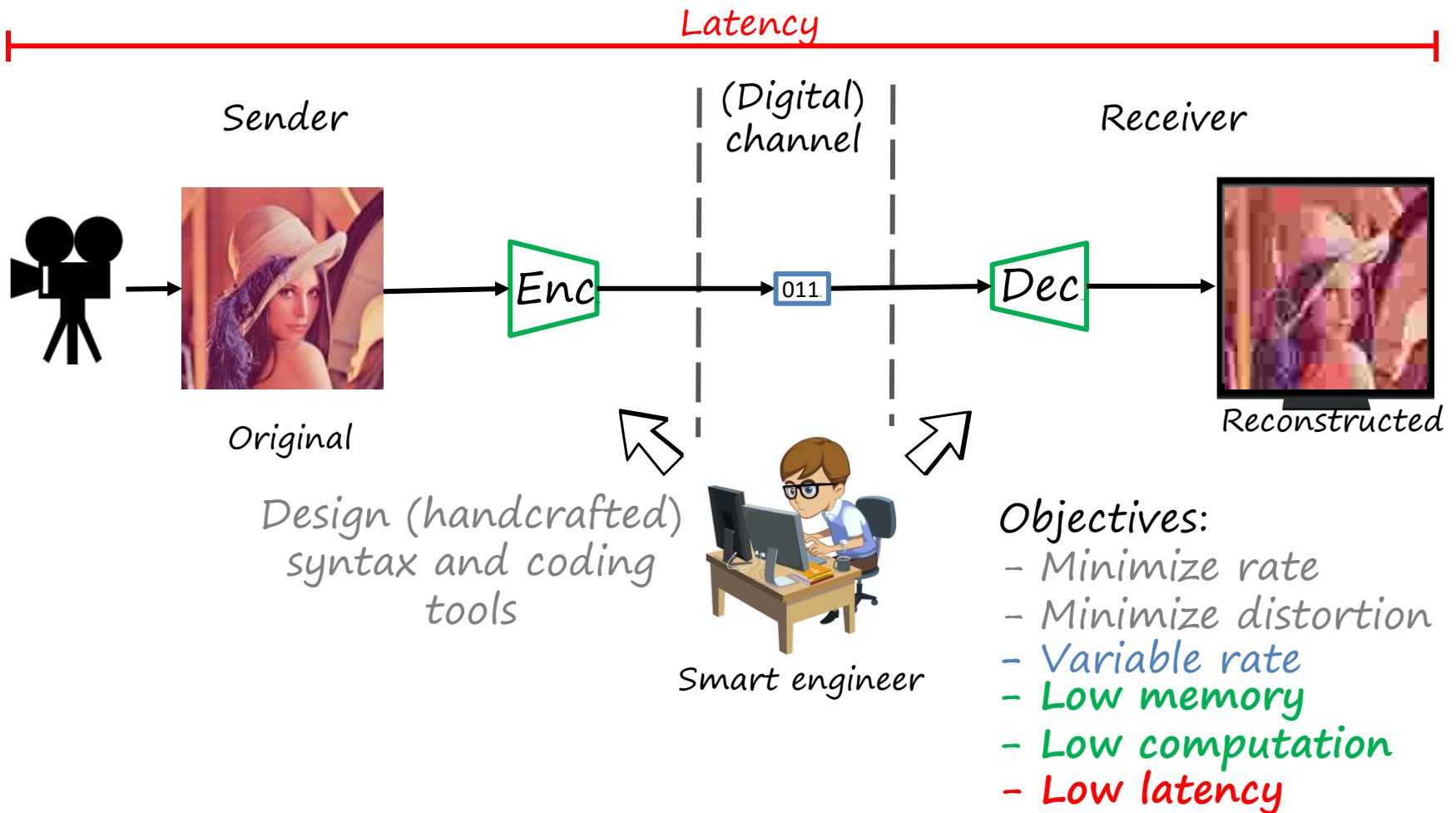
# The visual communication problem



Lossless compression

Same as the original image

Sender

(Digital) channel

Receiver

Encoder

01100

Decoder

Original

As few bits as possible (low rate)

Reconstructed

Error

# The visual communication problem



Sender

Original

Encoder

(Digital) channel

01100

As few bits as possible (low rate)

Decoder

Receiver

Reconstructed

Error

**Lossy** compression

As similar to the original as possible (low distortion)

# Developing traditional image/video codecs



Sender

(Digital) channel

Receiver

Original

Encoder

01100

Decoder

Reconstructed

Design (handcrafted) syntax and coding tools

Smart engineer

Objectives:
- **Minimize rate**
- **Minimize distortion**

# … for practical applications

Sender

(Digital) channel

Receiver

Encoder

011

Decoder

Original

Reconstructed

Design (handcrafted) syntax and coding tools

Smart engineer

Objectives:
- Minimize rate
- Minimize distortion
- **Variable rate**

# … for practical applications

Latency

Sender

(Digital) channel

Receiver

Enc

011

Dec

Original

Reconstructed

Design (handcrafted) syntax and coding tools

Smart engineer

**Objectives:**
- Minimize rate
- Minimize distortion
- Variable rate
- **Low memory**
- **Low computation**
- **Low latency**

# ... for practical applications

Latency

Sender

(Digital) channel

Receiver


Original

Enc

011

Dec


Reconstructed

Design (handcrafted) syntax and coding tools

Smart engineer

Objectives:
– Minimize rate
– Minimize distortion
– Variable rate
– Low memory
– Low computation
– Low latency
– **Compatibility**
– **Domain-specific**

# Basic pipeline



Feature encoder → **Quantization (lossy)** → Entropy encoder → `01100` → Entropy decoder → Feature decoder

**Entropy coding (lossless)**

Example: block-based transform coding (e.g. JPEG, MPEG-4)

Block partition → 8×8 DCT → Quantization → Entropy encoder → `01100` → Entropy decoder → 8×8 IDCT

**Per block**

# Neural image/video codecs

– Coding tools and syntax are **parametric** and **learned**
– Encoders/decoders and probability models are **deep neural networks**



Original

01100

Reconstructed

Train model

Collect data

Design network architecture

# Neural image compression

*Compressive autoencoder (CAE) [Theis2017, Balle2017]*
*(autoencoder+quantization+entropy coding)*

$$D\left(\;,\;\right) + \boldsymbol{\lambda} R\left(\boxed{01100}\right)$$



Encoder $\rightarrow$ 01100 $\rightarrow$ Decoder

*Optimize a weighted rate-distortion loss ($\lambda$ controls the **tradeoff**)*

# Neural image compression

# Architecture

Compressive autoencoder (CAE) [Theis2017, Balle2017]
(autoencoder+quantization+entropy coding)



$$D(\text{■},\text{■})$$
$$+\lambda R(\boxed{01100})$$

Feature encoder → Entropy encoder → $\boxed{01100}$ → Entropy decoder → Feature decoder

*Not differentiable!*

# Architecture (training)

Use differentiable proxies for end-to-end training



Model parameters      $\psi = (\theta, \phi, \nu)$

Loss      $J(\mathcal{X}^{\mathrm{tr}}, \psi; \lambda) = R(\mathcal{X}^{\mathrm{tr}}, \psi) + \lambda D(\mathcal{X}^{\mathrm{tr}}, \psi)$

Optimization problem    $\psi^* = \min_{\psi} J(\mathcal{X}^{\mathrm{tr}}, \psi; \lambda)$

# Autoencoder architecture

Balle et al.
[ICLR2017]

# Autoencoder architecture

Balle et al.
[ICLR2017]



Generalized divisive normalization (GDN) [Balle2016]

$$\hat{y}_i = \frac{y_i}{\left(\beta_i + \sum_j \gamma_{ij}\, y_j^2\right)^{1/2}}$$

$\gamma$     $\beta$

Learnable parameters

# Rate-distortion tradeoff λ



Each RD point is a different independent model (λ is fixed)

RD curve

Input

Decoded

Error

Low rate (λ=0.002)

High rate (λ=0.032)

PSNR= 31.1 dB
Rate= 0.08 bpp

PSNR= 36.2 dB
Rate= 0.41 bpp

# Is neural image compression practical?



$$D(\text{■},\text{■}) + \lambda R(\boxed{01100})$$

Encoder → 01100 → Decoder

**Limitations**

– $\lambda$ is fixed

– Heavy encoders/decoders

**Practical neural image compression?**
– Minimize rate ✓
– Minimize distortion ✓
– Variable rate ✗
– Low memory ✗
– Low computation ✗
– Low latency ✗

# Towards **practical** neural image compression



Main objectives
– Minimize rate
– Minimize distortion

Practical objectives
– **Variable rate**
– **Low memory**
– **Low computation**
– **Low latency**

MAE
[SPL2020]

SlimCAE
[CVPR2021]

Other practical considerations
– **Domain-specific codecs**
  (e.g. videoconference, screencast)
– **Backward/forward compatibility**
  (with legacy formats and encoders/decoders)

DANICE
[CLIC2021]

[SPL2020] Variable Rate Deep Image Compression with Modulated Autoencoder, Signal Processing Letters 2020
[CVPR2021] Slimmable compressive autoencoders for practical imaga compression, CVPR 2021
[CLIC2021] DANICE: Domain adaptation without forgetting in neural image compression, CLIC 2021 at CVPR 2021

# Slimmable Compressive Autoencoders for Practical Neural Image Compression

Fei Yang[1,2,3], Luis Herranz[1,2], Yongmei Cheng[3], Mikhail Mozerov[1,2]

[1]Computer Vision Center
[2]Universitat Autònoma de Barcelona
[3]Northwestern Politechnical University

CVPR 2021

# Variable rate neural image compression

Objective: one single model for multiple $\lambda$

Bottleneck scaling [Theis2017]          Feature modulation [MAE, cAE]



- Minimize rate ✓
- Minimize distortion ✓
- Variable rate ✓

- Low memory ✗
- Low computation ✗
- Low latency ✗

cAE: conditional autoencoder [Choi2019]
MAE: modulated autoencoder [Yang2020]

# Model capacity and rate-distortion

# Slimmable compressive autoencoder

Approach: slim the network to the minimal capacity for a given $\lambda$

Slimming [SlimCAE]



- Minimize rate ✓
- Minimize distortion ✓
- Variable rate ✓
- Lower memory ✓
- Lower computation ✓
- Lower latency ✓

(for low-mid rates)

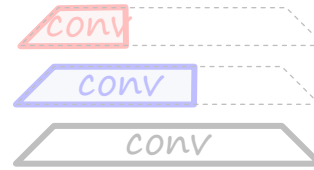Select w and slim

$\lambda$

# Slimmable layers in SlimCAE

SlimCAE
w$\in[$w$_1$,w$_2$, w$_3]$

SlimConv
SlimIGDN
SlimConv
SlimIGDN
SlimConv
SlimIGDN

SlimGDN
SlimConv
SlimGDN
SlimConv
SlimGDN
SlimConv

# Slimmable layers in SlimCAE



SlimCAE
$w \in [w_1, w_2, w_3]$

SlimConv
SlimIGDN
SlimConv
SlimIGDN
SlimConv
SlimIGDN

SlimGDN
SlimConv
SlimGDN
SlimConv
SlimGDN
SlimConv

Slimmable convolution [Yu2019]

conv
conv
conv

# Slimmable layers in SlimCAE
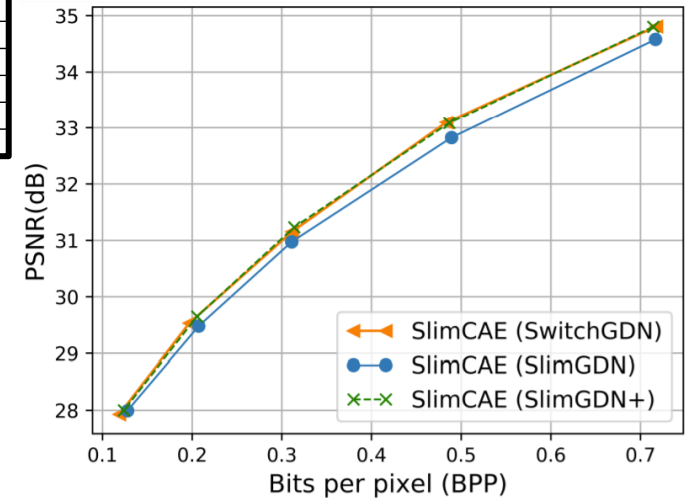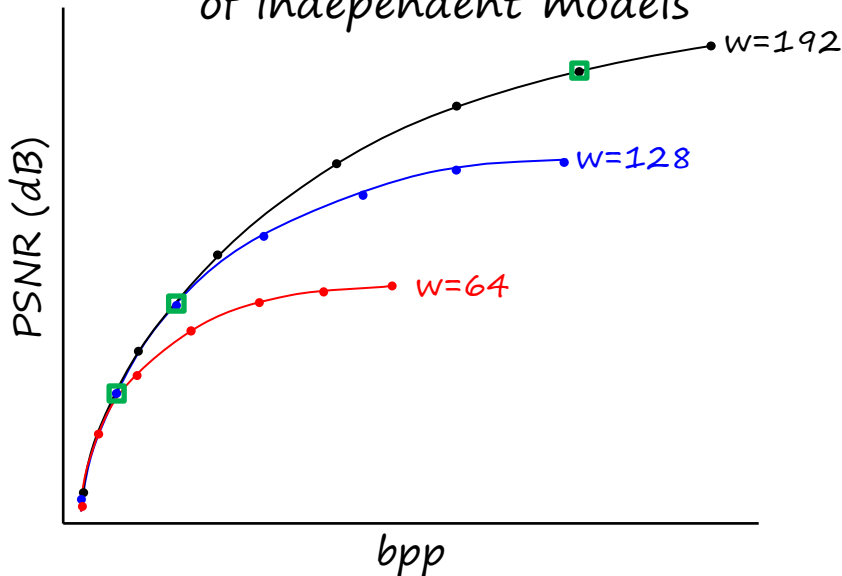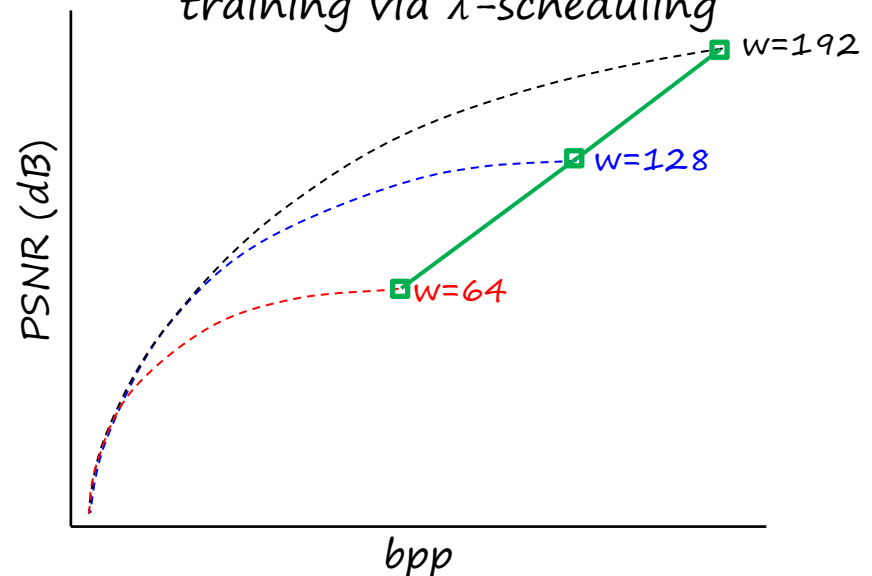
# Training SlimCAE

Problem: we need the optimal $\lambda$s to train the SlimCAE



Estimate from RD curves of independent models

Automatically estimate during training via $\lambda$-scheduling

1. Train several independent models for different w
2. Plot RD curves and find critical points
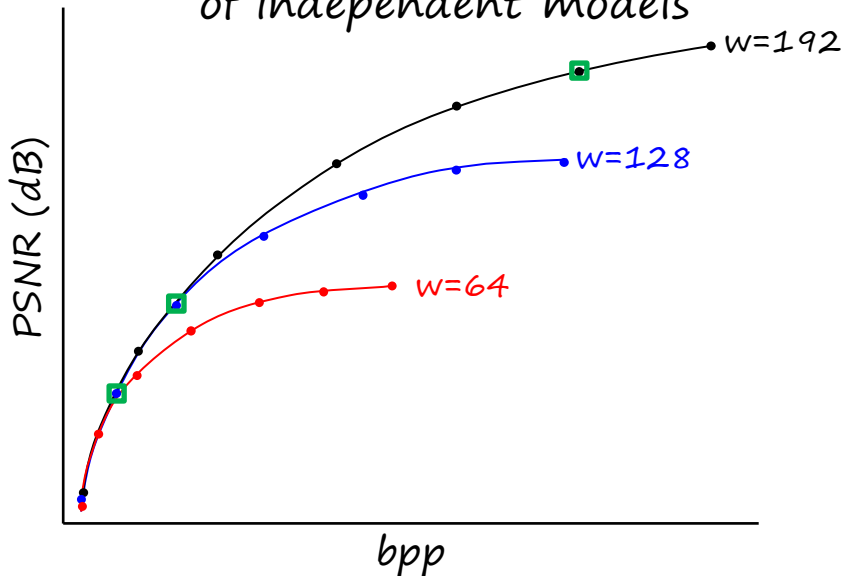3. Estimate optimal $\lambda$s from trained models

Problem: extremely expensive!

1. Train a SlimCAE with $\lambda_1 = \lambda_2 = \lambda_3$
2. While not converged do
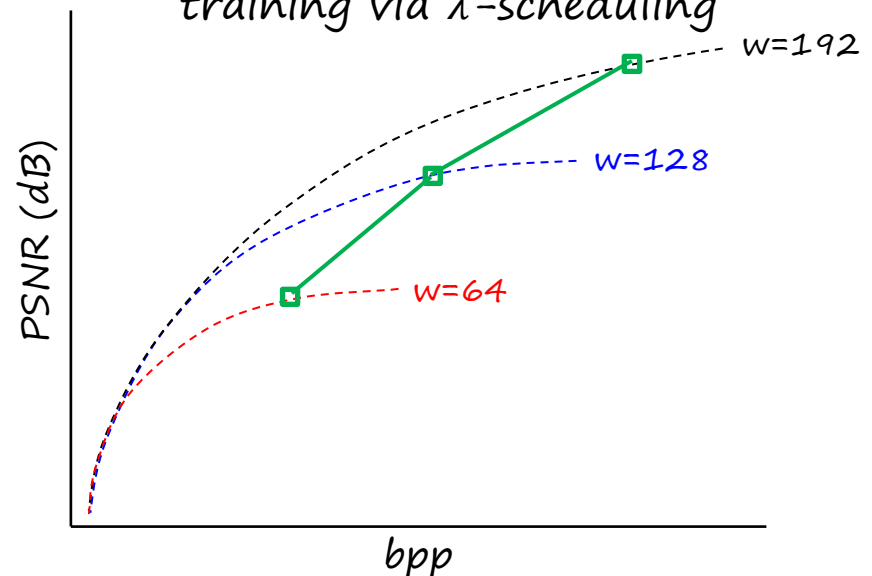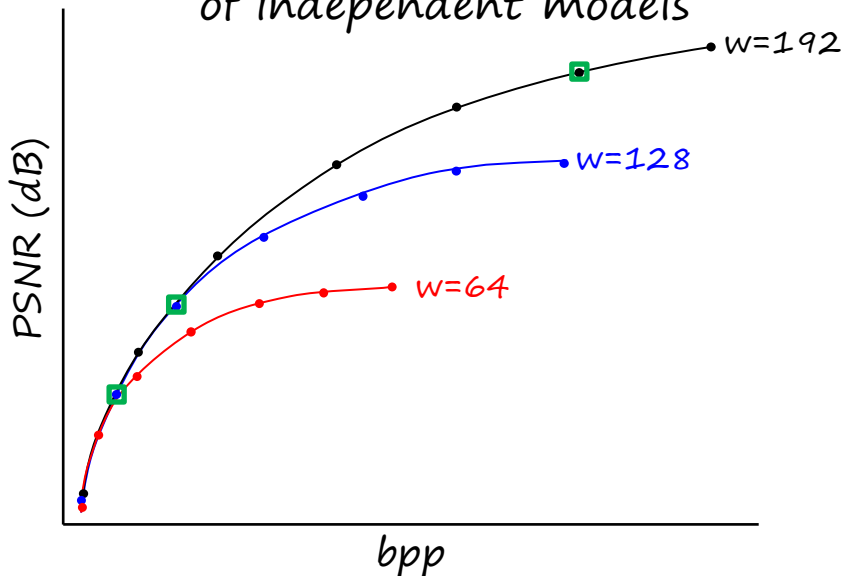   - Update $\lambda$s according to schedule
   - Optimize CAE

# Training SlimCAE

Problem: we need the optimal $\lambda$s to train the SlimCAE

### Estimate from RD curves of independent models



1. Train several independent models for different w
2. Plot RD curves and find critical points
3. Estimate optimal $\lambda$s from trained models

Problem: extremely expensive!

### Automatically estimate during training via $\lambda$-scheduling



1. Train a SlimCAE with $\lambda_1 = \lambda_2 = \lambda_3$
2. While not converged do
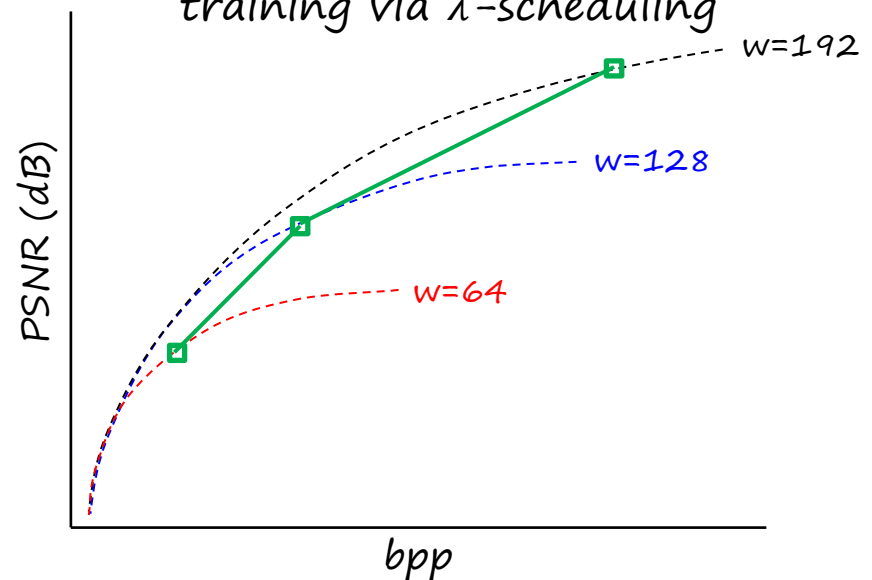   - Update $\lambda$s according to schedule
   - Optimize CAE

# Training SlimCAE

Problem: we need the optimal $\lambda$s to train the SlimCAE

Estimate from RD curves of independent models



Automatically estimate during training via $\lambda$-scheduling



1. Train several independent models for different w
2. Plot RD curves and find critical points
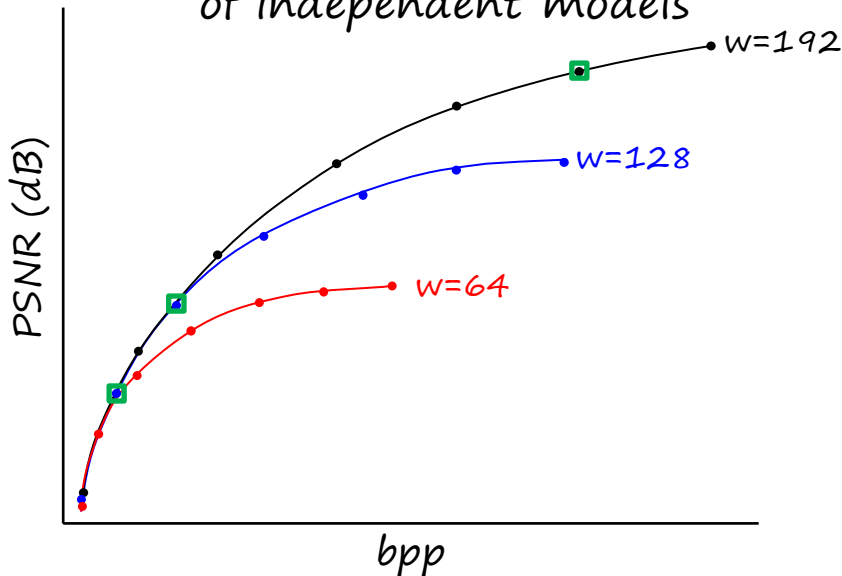3. Estimate optimal $\lambda$s from trained models

   Problem: extremely expensive!

1. Train a SlimCAE with $\lambda_1 = \lambda_2 = \lambda_3$
2. While not converged do
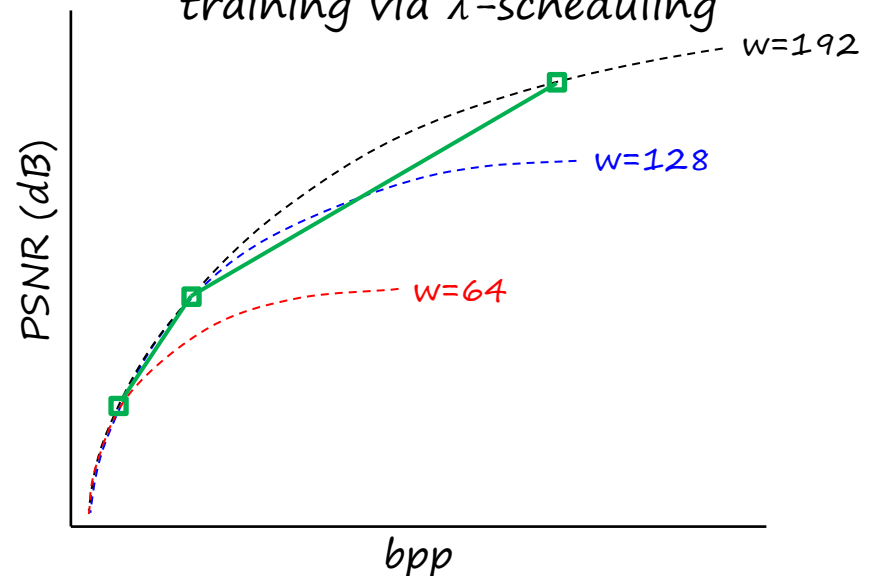   - Update $\lambda$s according to schedule
   - Optimize CAE

# Training SlimCAE

Problem: we need the optimal $\lambda$s to train the SlimCAE

Estimate from RD curves
of independent models



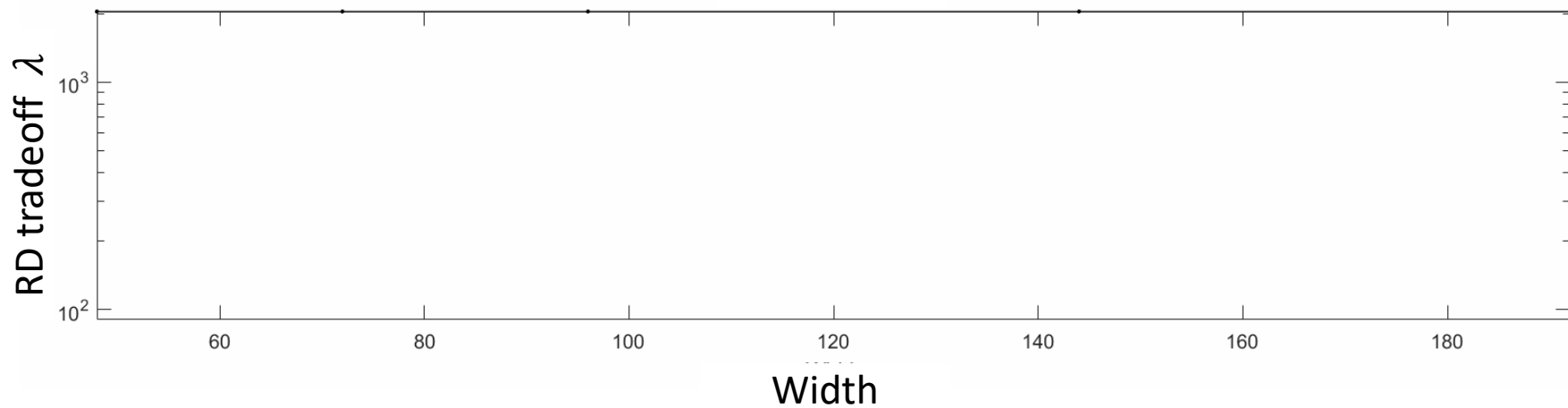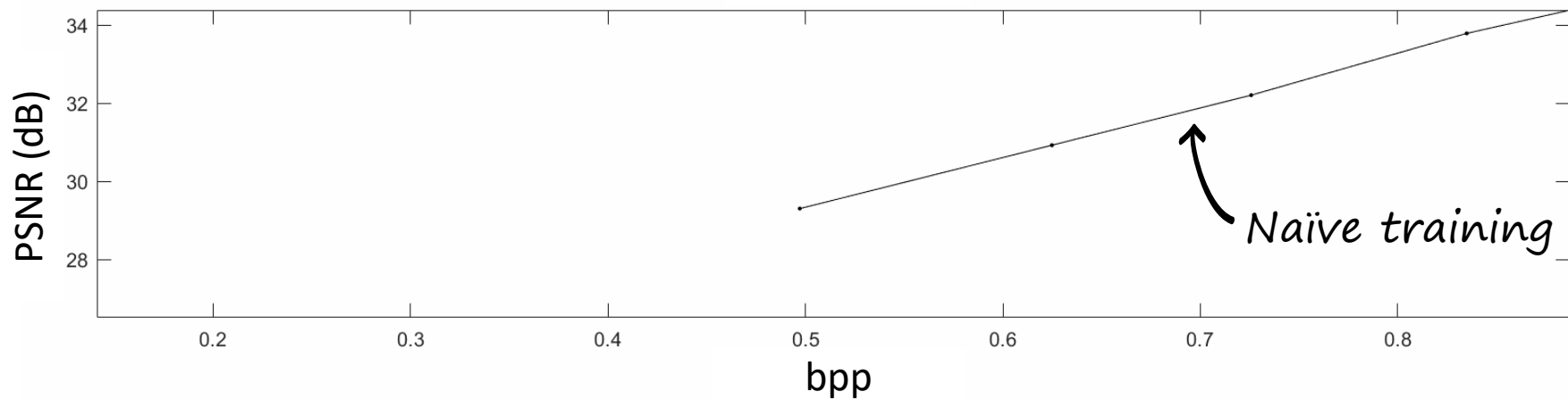Automatically estimate during
training via $\lambda$-scheduling



1. Train several independent models for different w
2. Plot RD curves and find critical points
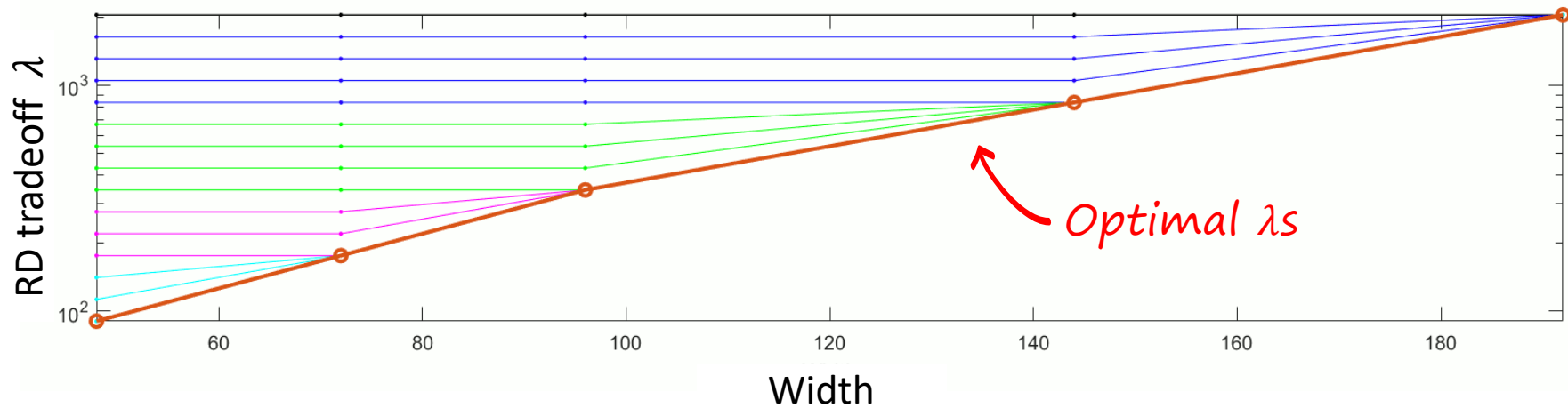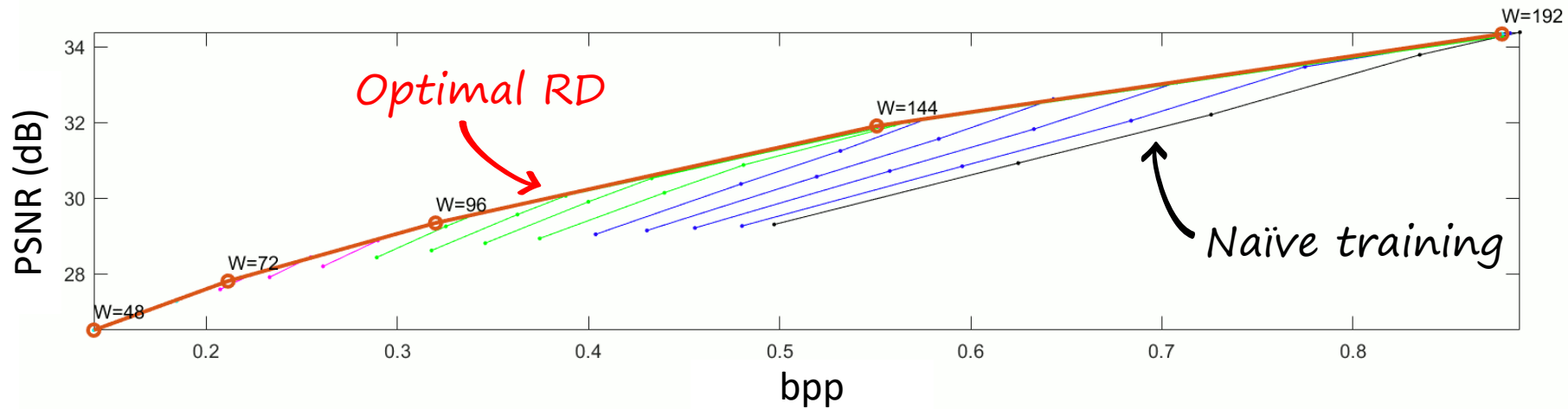3. Estimate optimal $\lambda$s from trained models

Problem: extremely expensive!

1. Train a SlimCAE with $\lambda_1 = \lambda_2 = \lambda_3$
2. While not converged do
   • Update $\lambda$s according to schedule
   • Optimize CAE

Directly train one model!

# $\lambda$-scheduling. Example



*Naïve training*

PSNR (dB) vs bpp

RD tradeoff $\lambda$ vs Width

# λ-scheduling

# Performance comparison



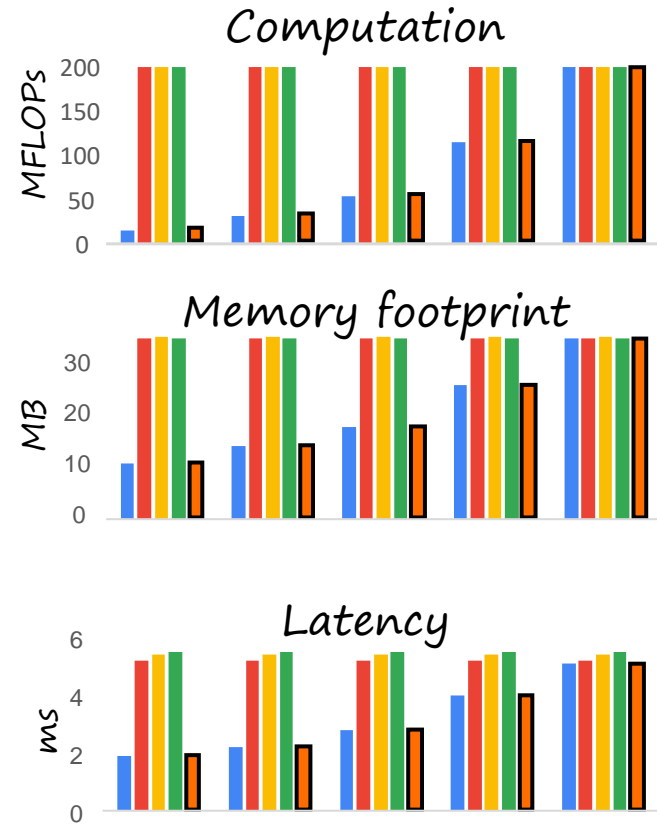Independent CAEs (each with minimal capacity) — Scaling [Theis2017] — MAE [Yang2020] — cAE [Choi2019] — SlimCAE (ours)

Rate-distortion

Encoder

Computation

Memory footprint

Latency

# Thanks!

https://arxiv.org/abs/2103.15726

https://github.com/FireFYF/SlimCAE



Fei Yang      Luis Herranz      Yongmei Cheng      Mikhail Mozerov

# DANICE: Domain adaptation without forgetting in neural image compression

Sudeep Katakol[1,(2)], Luis Herranz[2,3], Fei Yang[2,3,4], Marta Mrak[5]

[1]University of Michigan, Ann Arbor, [2]Computer Vision Center, [3]Universitat Autònoma de Barcelona, [4]Northwestern Politechnical University, [5]BBC R&D
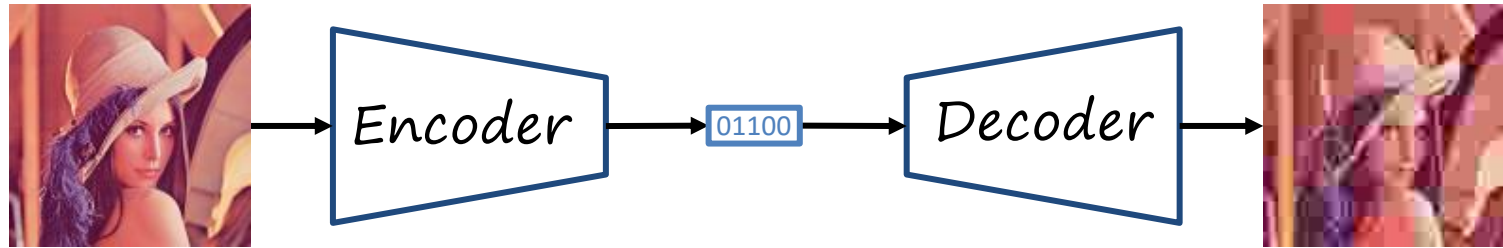
CLIC 2021 (@CVPR 2021)

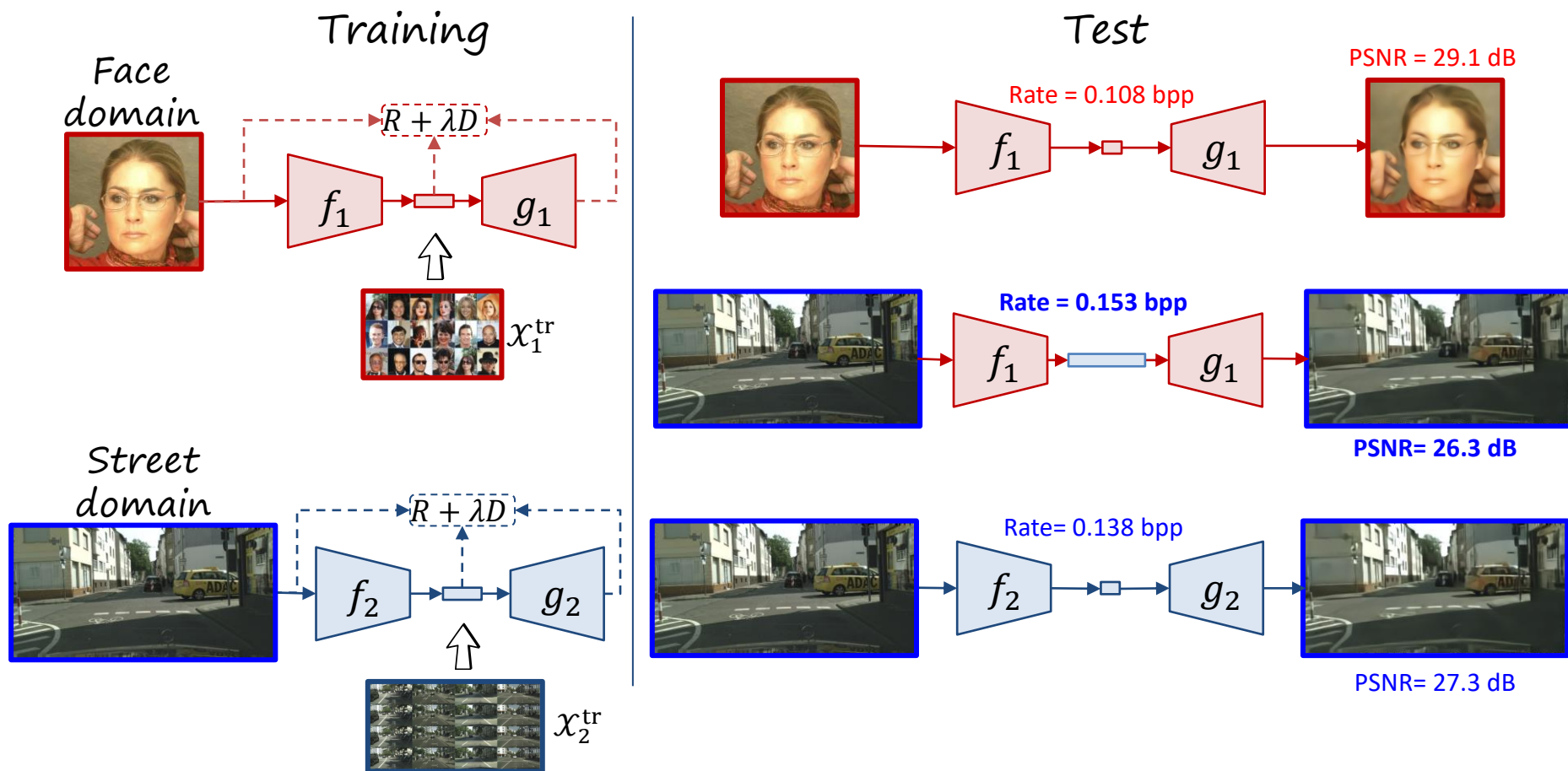# Towards **practical** neural image compression



Encoder → 01100 → Decoder

Other practical considerations
- **Domain-specific codecs**
  (e.g. videoconference, screencast)
- **Backward/forward compatibility**
  (with legacy formats and encoders/decoders)
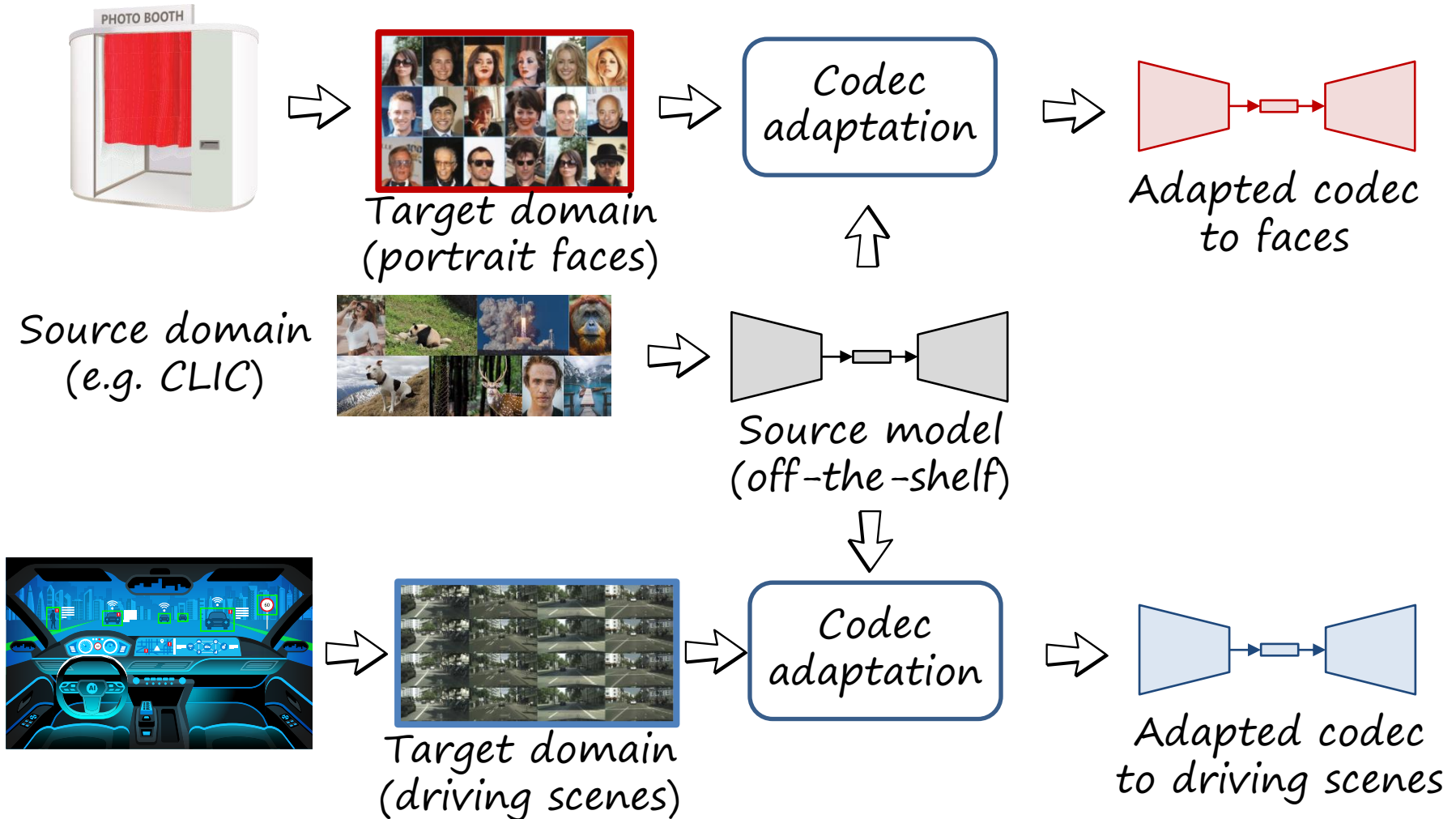
DANICE
[CLIC2021]

# Rate-distortion optimality of learned codecs

*Learned codecs are only optimal in the domain of the training data*



**Training**

Face domain

$R + \lambda D$

$f_1$   $g_1$

$x_1^{tr}$

Street domain

$R + \lambda D$

$f_2$   $g_2$

$x_2^{tr}$

**Test**

PSNR = 29.1 dB

Rate = 0.108 bpp

$f_1$   $g_1$

Rate = 0.153 bpp

$f_1$   $g_1$

PSNR= 26.3 dB

Rate= 0.138 bpp
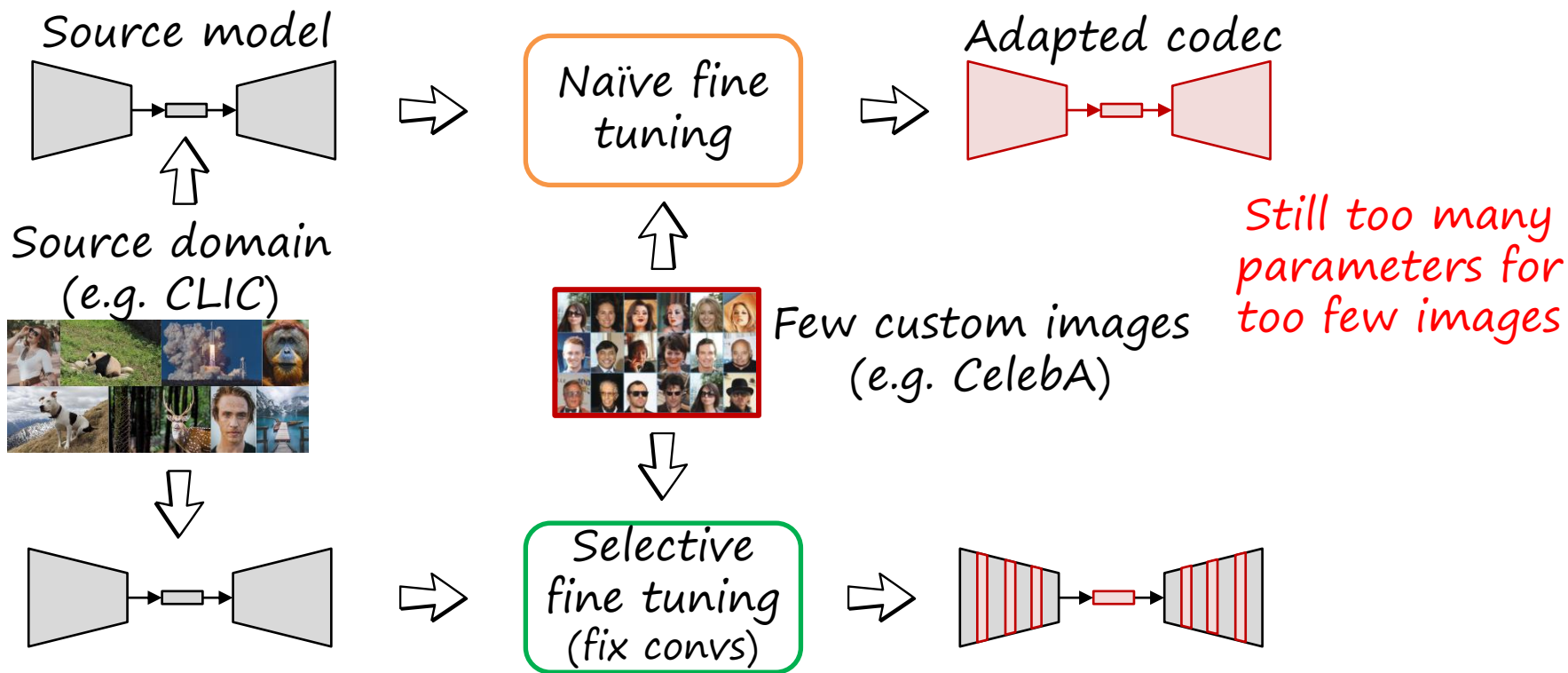
$f_2$   $g_2$

PSNR= 27.3 dB

# Domain Adaptation in Neural Image ComprEssion (DANICE)

Learned codecs can be customized with user content to specific domains
Problem: usually we don't have enough custom data; training is expensive
Solution: transfer pre-trained codecs
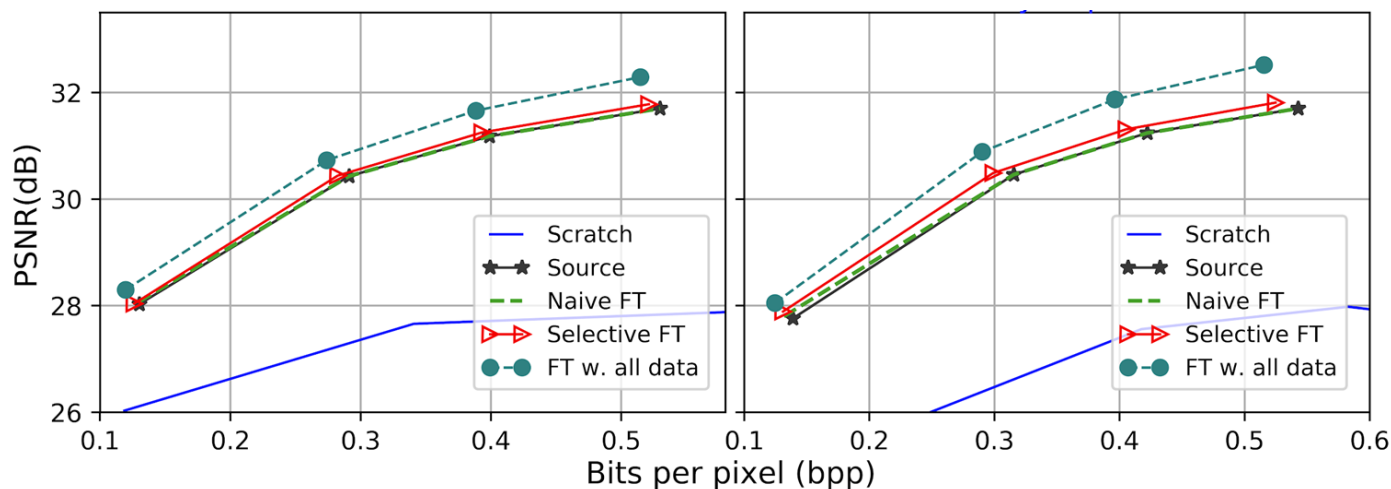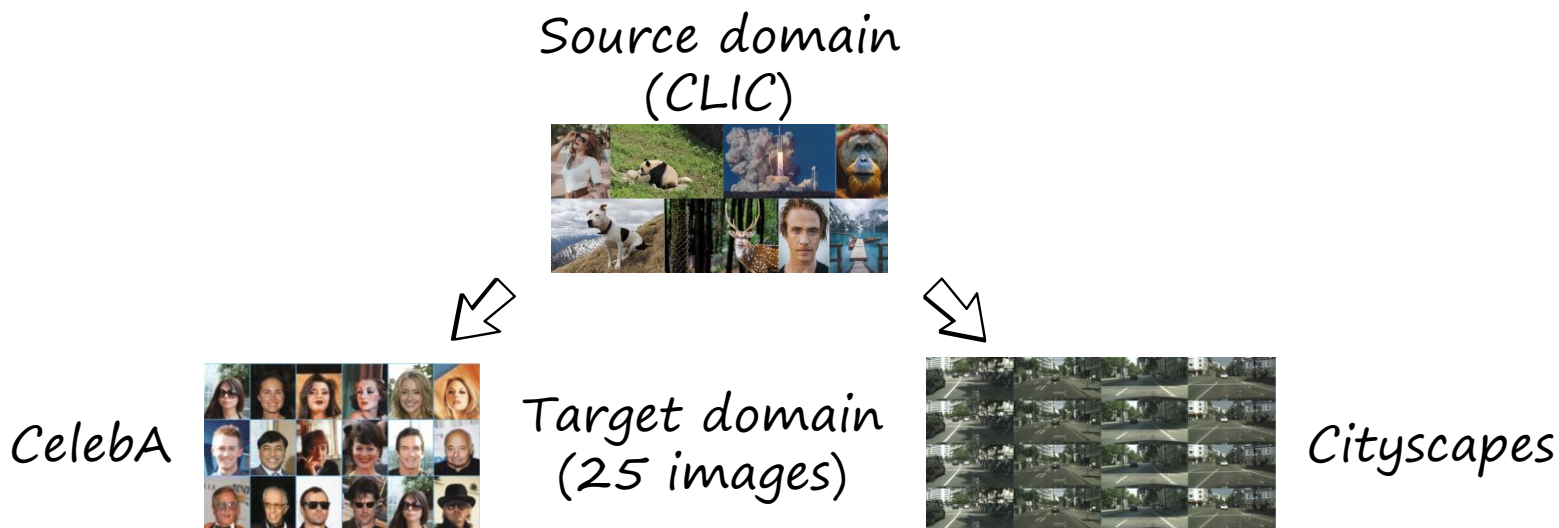
# Domain adaptation via fine tuning



Source model

Source domain
(e.g. CLIC)

Naïve fine tuning

Adapted codec

Still too many parameters for too few images

Few custom images
(e.g. CelebA)

Selective fine tuning
(fix convs)

Experiments

|  | CLIC → CelebA | | CLIC → Cityscapes | |
|---|---|---|---|---|
| Source model | 19.24 | | 23.93 | |
| Number of target images | Naïve fine tuning | Selective fine tuning | Naïve fine tuning | Selective fine tuning |
| 10 | 19.24 | **16.46** | 22.96 | **17.54** |
| 25 | 18.76 | **14.93** | 18.44 | **15.79** |
| 50 | 15.59 | **13.73** | 16.29 | **15.33** |

BD-rate
(reference: training with all target data)

# Domain adaptation via fine tuning



Source domain (CLIC)

Target domain (25 images)

CelebA

Cityscapes
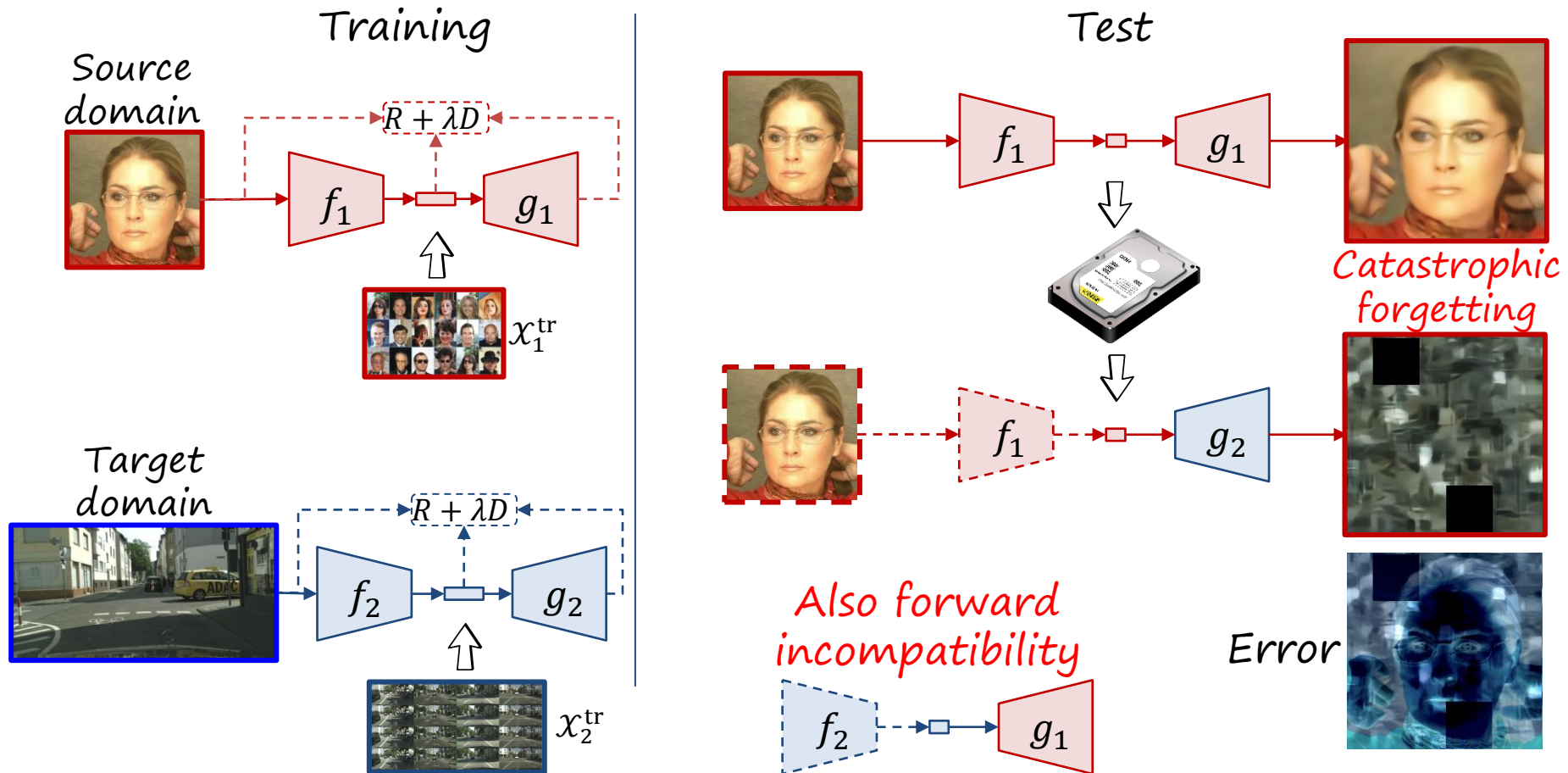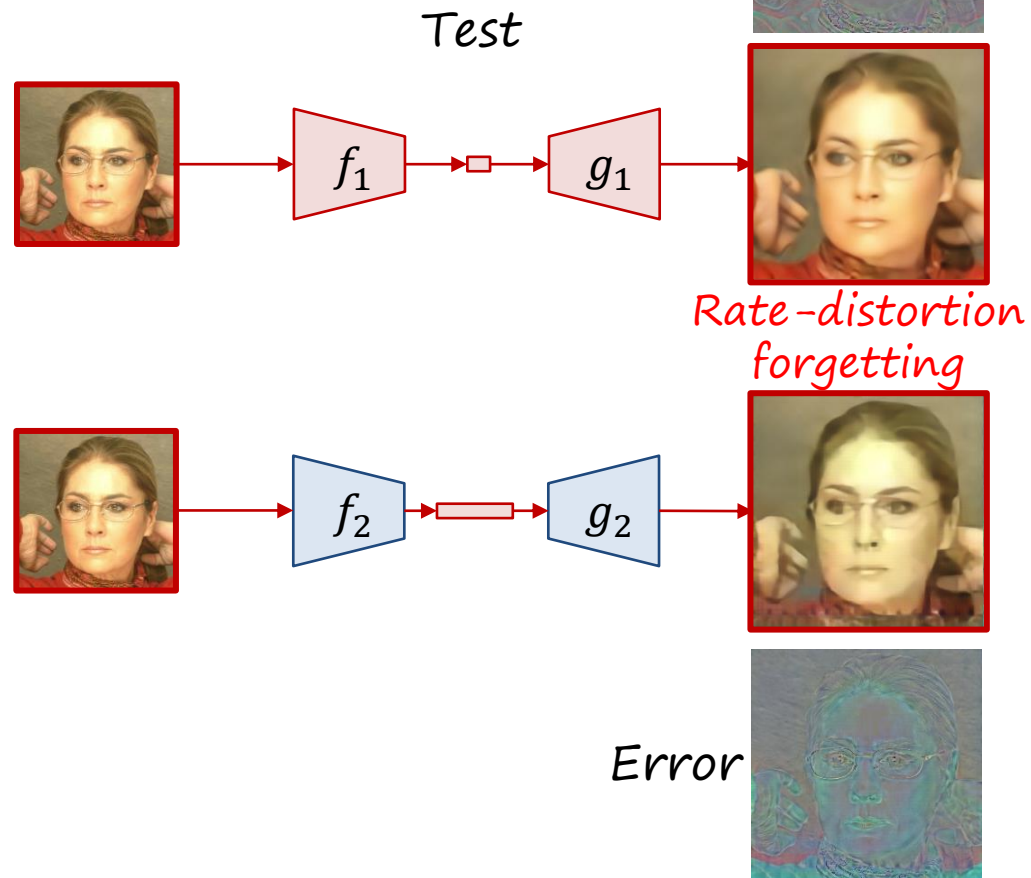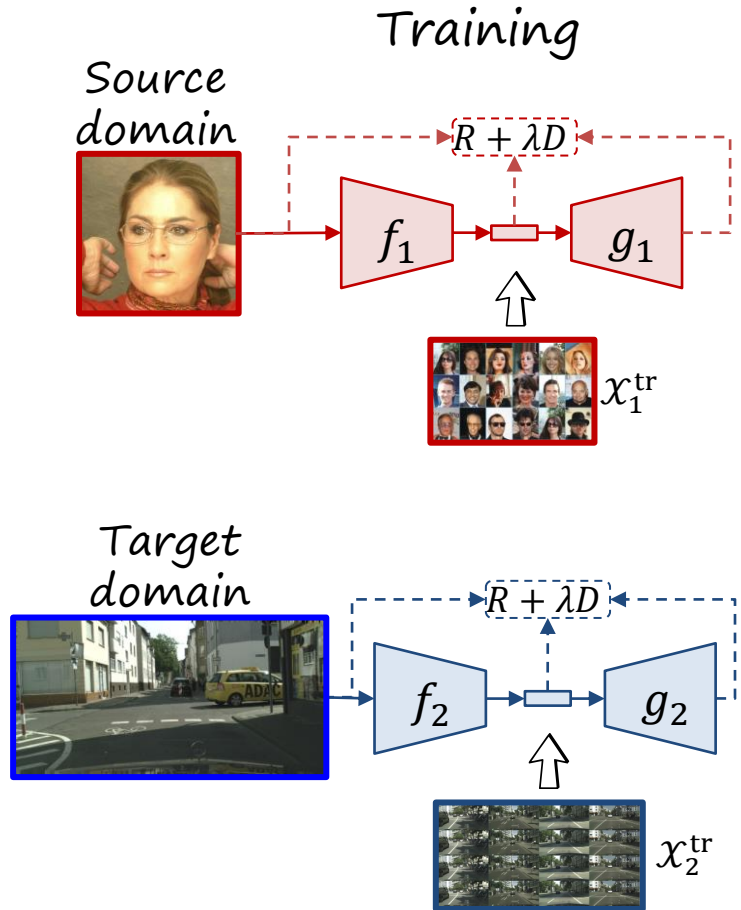
# Backward incompatibility with legacy bitstreams: catastrophic forgetting

**Misalignment between encoding-decoding latent spaces (i.e. bitstream syntax incompatible)**
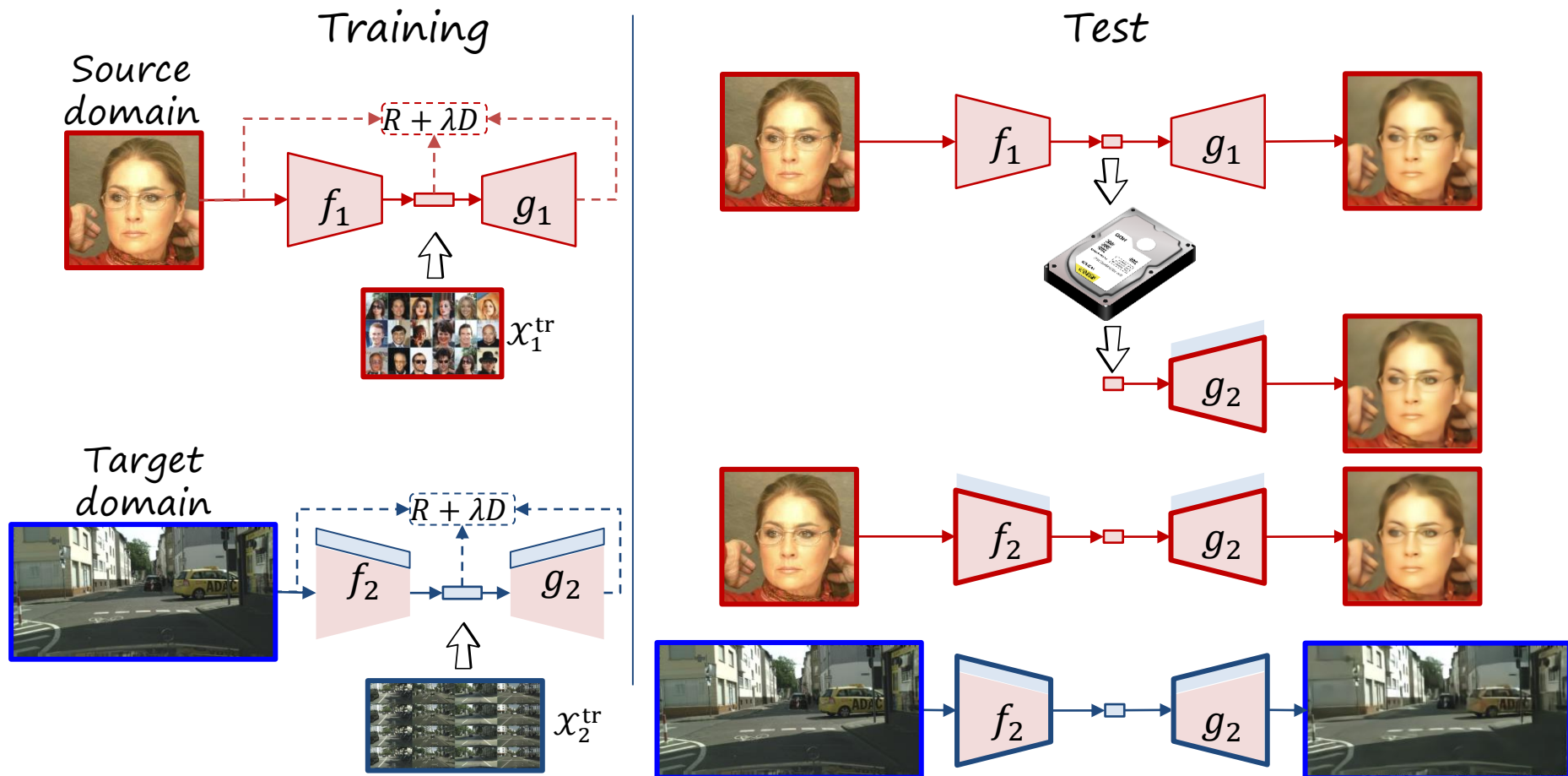
# Rate-distortion forgetting

Encoding-decoding latent spaces aligned, but suboptimal
(i.e. bitstream syntax compatible, yet degraded)

Training

Test

Source
domain

$R + \lambda D$

$f_1$    $g_1$

$x_1^{tr}$

$f_1$    $g_1$

Rate-distortion
forgetting

Target
domain

$R + \lambda D$

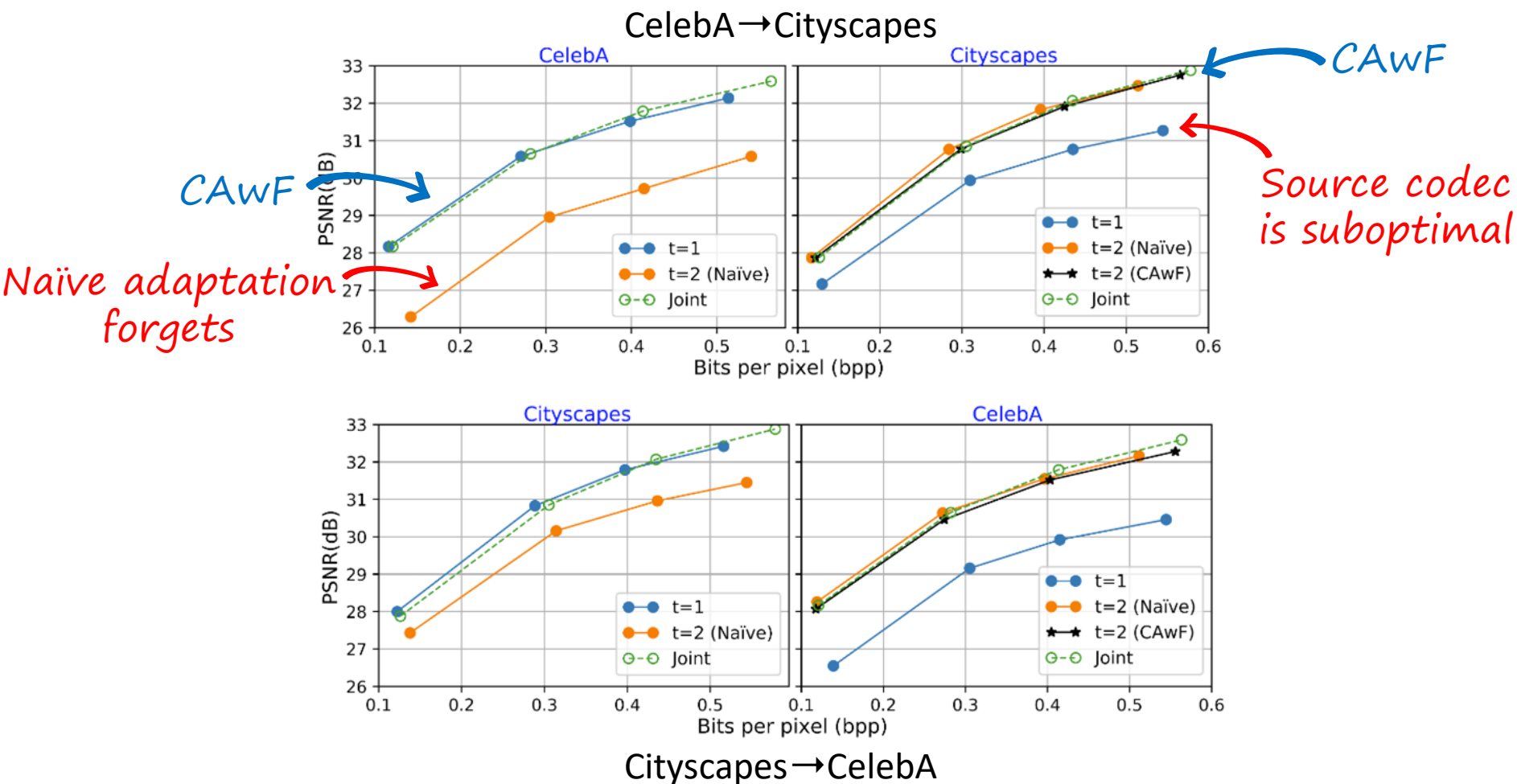$f_2$    $g_2$

$x_2^{tr}$

$f_2$    $g_2$

Error

# Codec adaptation without forgetting (CAwF)

Freeze source codec, and learn target codec as an enhancement layer
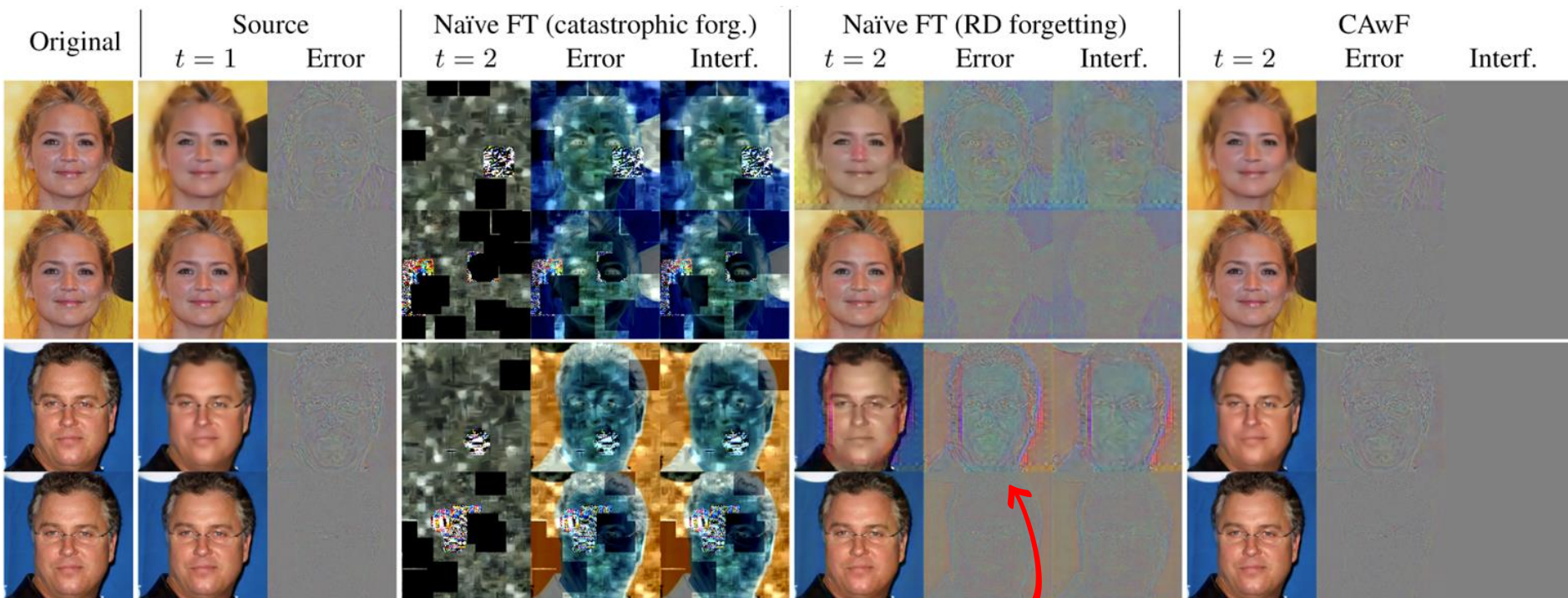Drawback: adds additional parameters

# Codec adaptation without forgetting (CAwF)

# Codec adaptation without forgetting (CAwF)



CelebA→Cityscapes
(source domain)

Codec adaptation
artifacts

# Thanks!

https://arxiv.org/abs/2103.15726



Sudeep Katakol      Luis Herranz      Fei Yang      Marta Mrak