NEURAL INFORMATION
PROCESSING SYSTEMS

# L-CAD: Language-based Colorization with Any-level Descriptions using Diffusion Priors

Zheng Chang[1#]   Shuchen Weng[2,3#]   Peixuan Zhang[1]   Yu Li[4]   Si Li[*1]   Boxin Shi[2,3]

[1]School of Artificial Intelligence, Beijing University of Posts and Telecommunications
[2]National Key Laboratory for Multimedia Information Processing, School of Computer Science, Peking University
[3]National Engineering Research Center of Visual Technology, School of Computer Science, Peking University
[4]International Digital Economy Academy

# Outline

# Background



1. A person wearing a spacesuit is mowing grass on the lawn.
2. A rabbit with sunglasses and a hat.
3. An astronaut piled up a pyramid with sand.
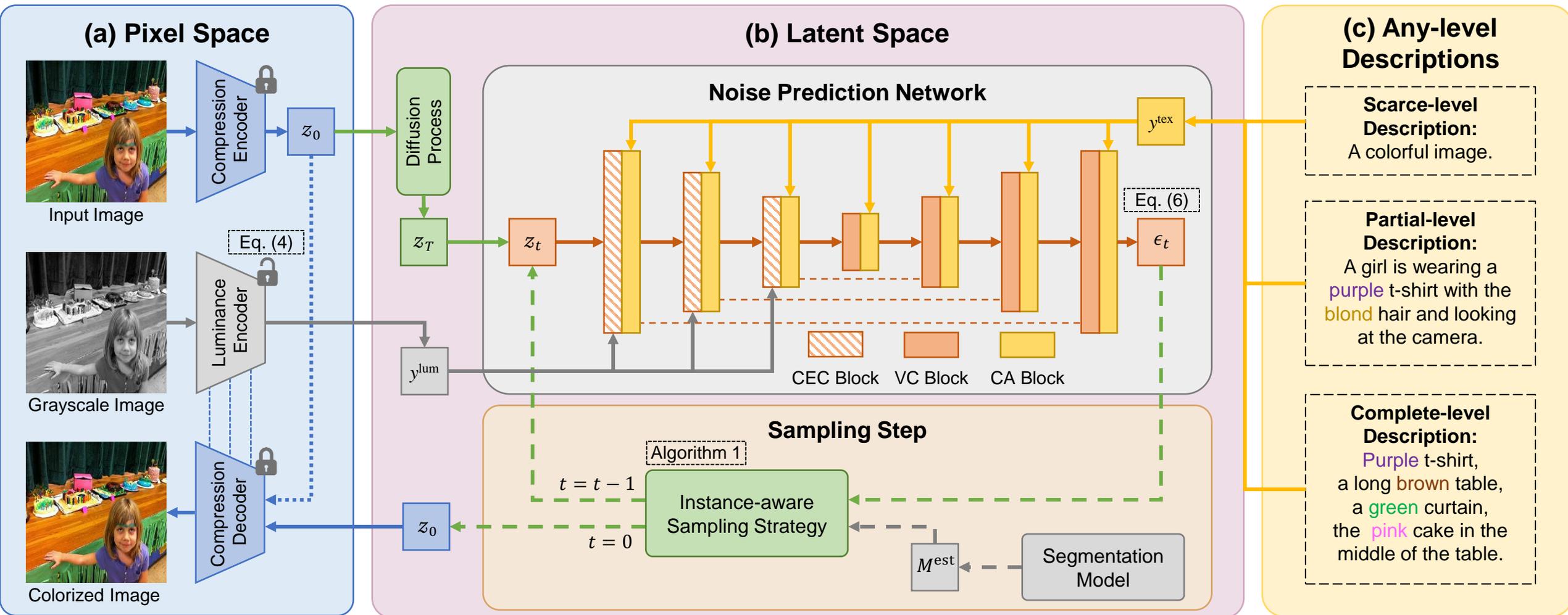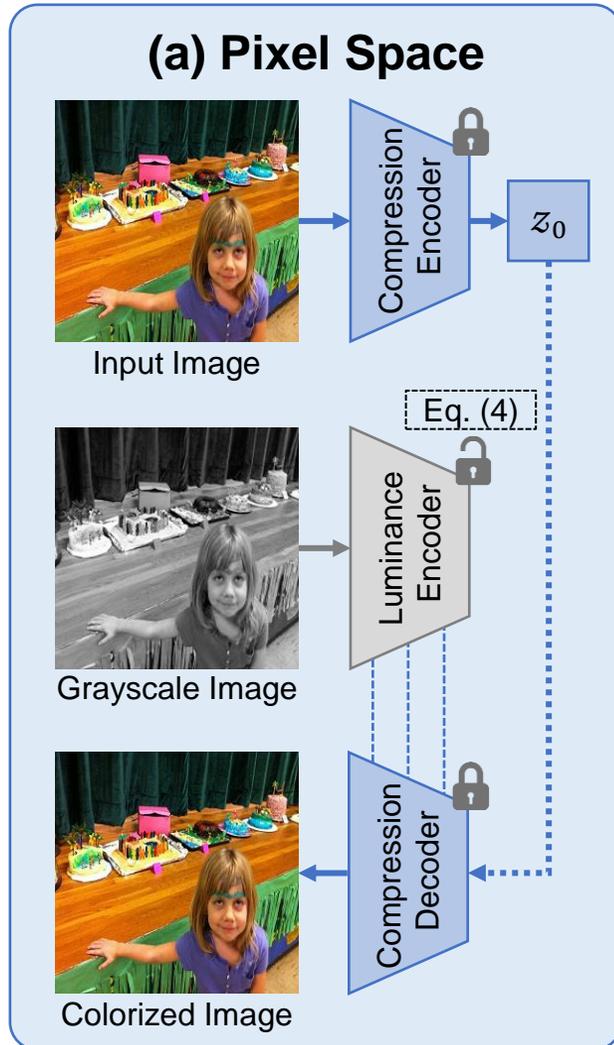4. The robot is pouring coffee……

# Motivation

# Outline

- Introduction
  - Background
  - Motivation
- Method
  - Pipeline
  - Luminance-guided image compression
  - Semantic-aligned latent representation
  - Instance-aware sampling strategy
- Result
  - Comparison with Automatic Colorization
  - Comparison with Language-based Colorization
  - Ablation
  - Application

# Pipeline



**(a) Pixel Space**

Input Image

Grayscale Image

Colorized Image

Compression Encoder

Luminance Encoder

Eq. (4)

Compression Decoder

$z_0$

**(b) Latent Space**

Diffusion Process

$z_T$

$z_t$

**Noise Prediction Network**

$y^{\text{tex}}$

Eq. (6)

$\epsilon_t$

CEC Block  VC Block  CA Block

$y^{\text{lum}}$

**Sampling Step**

Algorithm 1

$t = t - 1$

Instance-aware Sampling Strategy

$t = 0$

$z_0$

$M^{\text{est}}$

Segmentation Model

**(c) Any-level Descriptions**

**Scarce-level Description:**
A colorful image.

**Partial-level Description:**
A girl is wearing a purple t-shirt with the blond hair and looking at the camera.

**Complete-level Description:**
Purple t-shirt,
a long brown table,
a green curtain,
the pink cake in the middle of the table.

# Luminance-guided image compression

## (a) Pixel Space



Input Image

Compression Encoder

$z_0$

Eq. (4)

Luminance Encoder

Grayscale Image

Compression Decoder
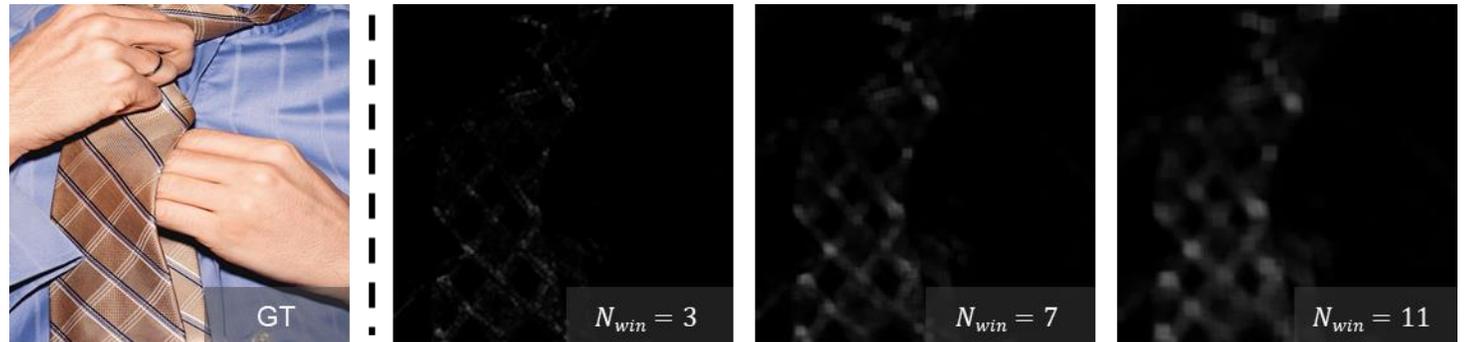
Colorized Image

**Training loss:**

$$\mathcal{L}_{\mathrm{pix}} = \mathcal{L}_{\mathrm{rec}} + \alpha \mathcal{L}_{\mathrm{dis}} + \beta \mathcal{L}_{\mathrm{per}}$$
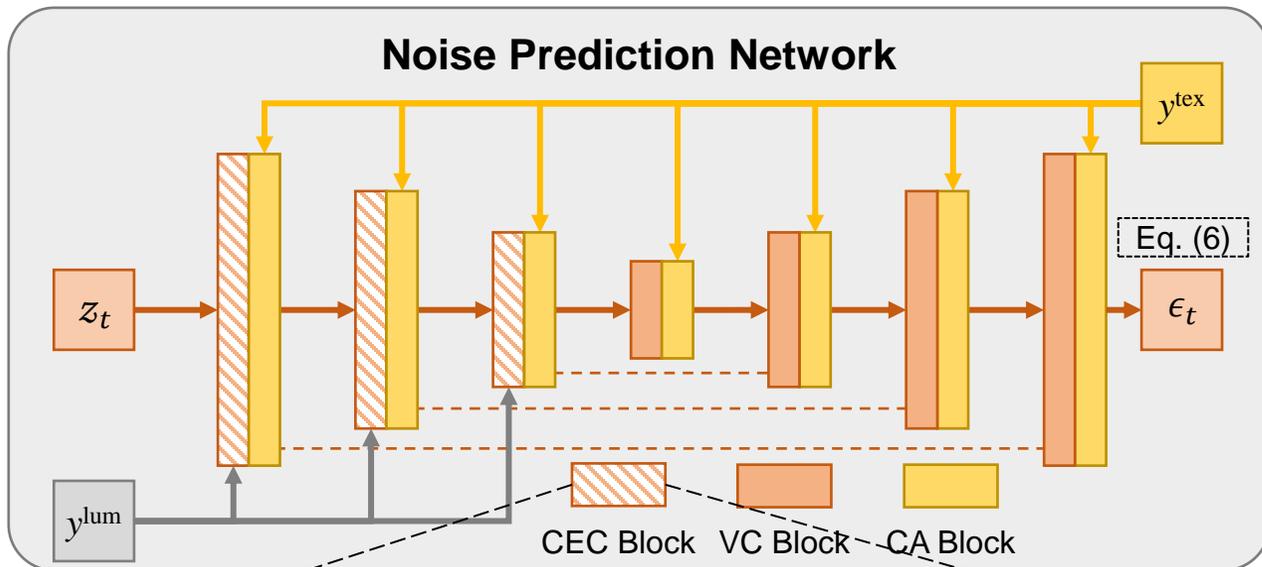
**New reconstruction loss:**

$$\mathcal{L}_{\mathrm{rec}} = \| M^{\mathrm{art}} \odot (x - \tilde{x}) \|_1 .$$

**Artifacts map:**

$$M^{\mathrm{art}}_{h,w} = \sum_{p \in \Omega_{(h,w)}} \left( \frac{\delta_p - \mu_p}{N_{\mathrm{win}}} \right)^2, \qquad \mu_p = \sum_{p \in \Omega_{(h,w)}} \frac{\delta_p}{N^2_{\mathrm{win}}},$$



GT      $N_{win} = 3$      $N_{win} = 7$      $N_{win} = 11$

# Semantic-aligned latent representation



**Noise Prediction Network**

CEC Block   VC Block   CA Block

Formula for extended convolution:

$$f'_{h,w} = \sum_{i=0}^{N_{\mathrm{k}}-1} \sum_{j=0}^{N_{\mathrm{k}}-1} \Big( \sum_{k=1}^{N_{\mathrm{fix}}} \omega^{\mathrm{fix}}_{i,j,k} f_{p,q,k} + \sum_{k=1}^{N_{\mathrm{ext}}} \omega^{\mathrm{ext}}_{i,j,N_{\mathrm{fix}}+k} \bar{y}^{\mathrm{lum}}_{p,q,k} \Big),$$

Training loss:

$$\mathcal{L}_{\mathrm{lat}} = \mathbb{E}_{t,z_0,\epsilon \sim \mathcal{N}(0,1)} \Big[ \| \epsilon_t - \epsilon_\theta(z_t, t, y^{\mathrm{tex}}, y^{\mathrm{lum}}) \|^2 \Big],$$

# Instance-aware sampling strategy

**Algorithm 1:** Instance-aware sampling strategy

**input** : Roughly estimated object contour $M^{\text{est}}$

**output** : Colorized latent representation $z_0$

**for** $t = T \dots 1$ **do**

    $\_, M^{\text{att}}_* = \epsilon_\theta(z_t, t, y^{\text{lum}}, y^{\text{tex}})$

    **for** $l = 1 \cdots L$ **do**

        $\hat{M}^{\text{est}}_l \leftarrow \text{Downsampling}(M^{\text{est}}, l)$

        $\mathcal{M} \leftarrow \text{Sigmoid}(M^{\text{att}}_l)$

        $\hat{M}^{\text{att}}_l \leftarrow M^{\text{att}}_l - \lambda \nabla_{M^{\text{att}}_l} \mathcal{L}_{\text{BCE}}(\mathcal{M}, \hat{M}^{\text{est}}_l)$

    **end**

    $\hat{\epsilon}_{t,\_} = \epsilon_\theta(z_t, t, y^{\text{lum}}, y^{\text{tex}})\{M^{\text{att}}_* \leftarrow \hat{M}^{\text{att}}_*\}$

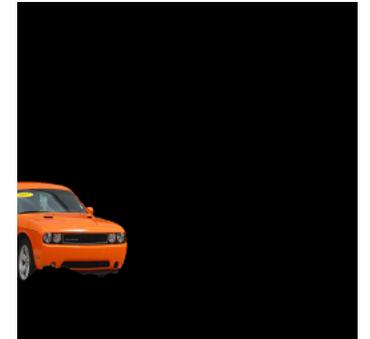    $z_{t-1} = \text{DDIM}(z_t, \hat{\epsilon}_t, t)$

**end**
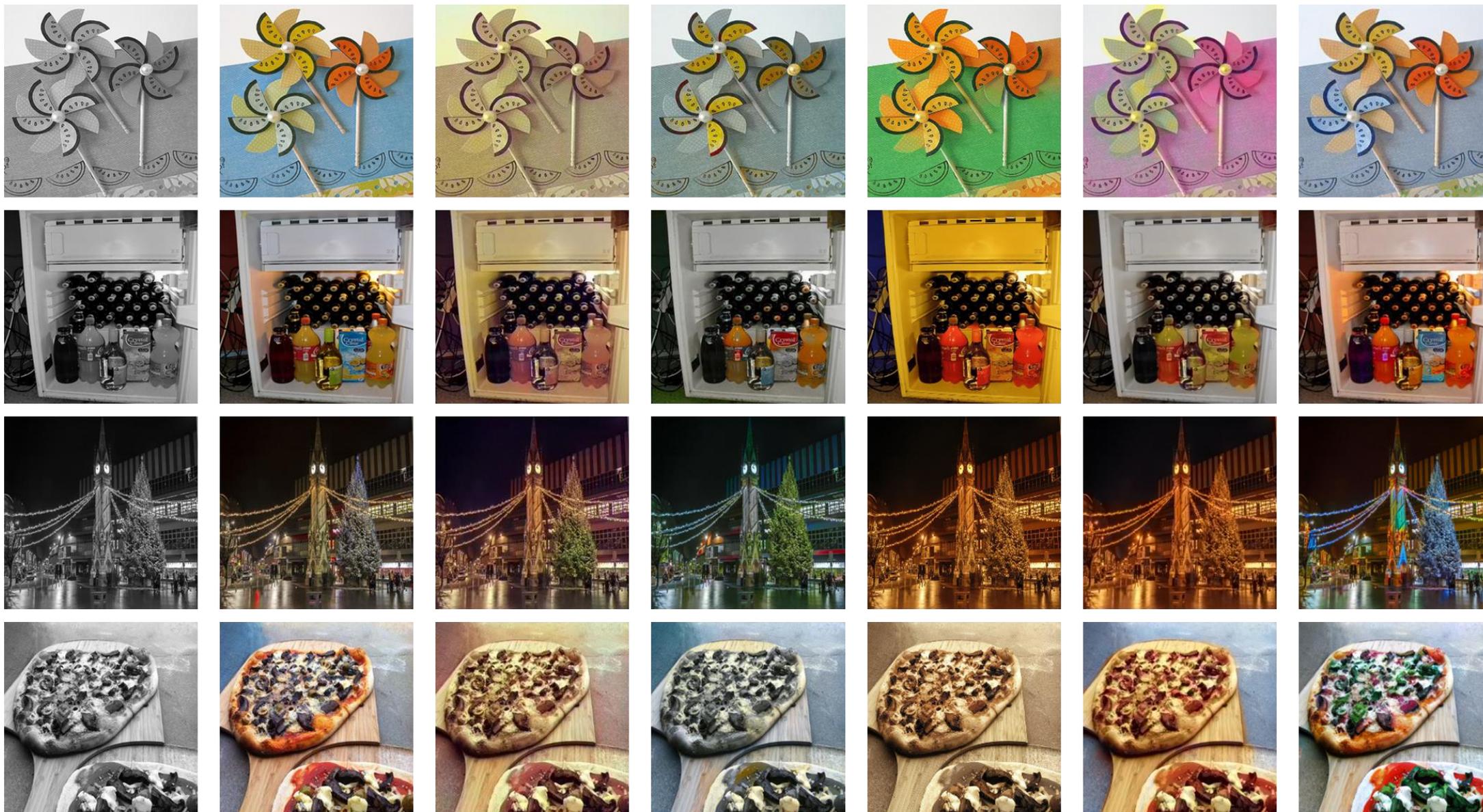


Ground-truth



$M^{est}$ for "left orange car"

$M^{est}$ for "right yellow car"

# Outline

- Introduction
  - Background
  - Problems
- Method
  - Pipeline
  - Luminance-guided image compression
  - Semantic-aligned latent representation
  - Instance-aware sampling strategy
- Result
  - Comparison with Automatic Colorization
  - Comparison with Language-based Colorization
  - Ablation
  - Application

# Comparison with Automatic Colorization



Grayscale     Ground Truth     CIC     BigColor     DISCO     CT$^2$     Ours

# Comparison with Language-based Colorization



| | Description | Grayscale | Ground Truth | ML2018 | L-CoDe | L-CoIns | Ours |
|---|---|---|---|---|---|---|---|

**Complete level**

The two chairs on the left are green, and the ones on the right are blue and pink.

The man on the left is wearing a **black** suit and red tie, and the man on the right is wearing a red suit and **black** tie.

**Partial level**

Orange carrots and some vegetable on a white plate ready to be cut.

A double decker red white and purple bus.

# Ablation



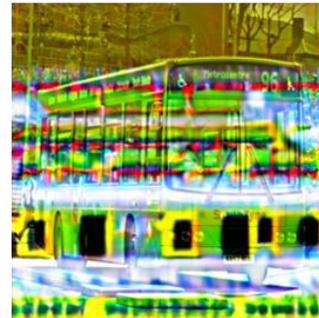| | Description | Grayscale | Ground Truth | *W/o* LIC | *W/o* SLR | *W/o* ISS | Ours |
|---|---|---|---|---|---|---|---|
| Complete level | The boy is wearing a yellow T-shirt, a white hat, and blue pants, holding a pizza. | | | | | | |
| Partial level | A purple bus is parked on the roadside. | | | | | | |
| Scarce level | A colorful image. | | | | | | |

# Application



*1960. "Judge Myles Paige. Columbus, Georgia."*

In front of the green vegetable stand stood a man in yellow clothes and green trousers.

*1957. "Plymouth vs. Ford on the streets of Oakland circa."*

There was a cyan car in the middle of the road.

There was a yellow car in the middle of the road.

*1940. " The family of Mr. Timothy Levy Crouch at their Thanksgiving Day dinner "*

A colorful image.

*1897. "Palm walk on Lake Worth, Palm Beach."*

There are many green trees along the beige path.

*1925. "Washington, D.C. Judge Geo. H. MacDonald & Geo. G. Adams."*

The man on the left was wearing an orange suit and the man on the right is wearing a pink suit.

The man on the left was wearing a gray suit and the man on the right is wearing a khaki suit.

*1940. "Stephen A. Lynch Jr. residence, Sunset Island, Miami Beach."*

A colorful image.