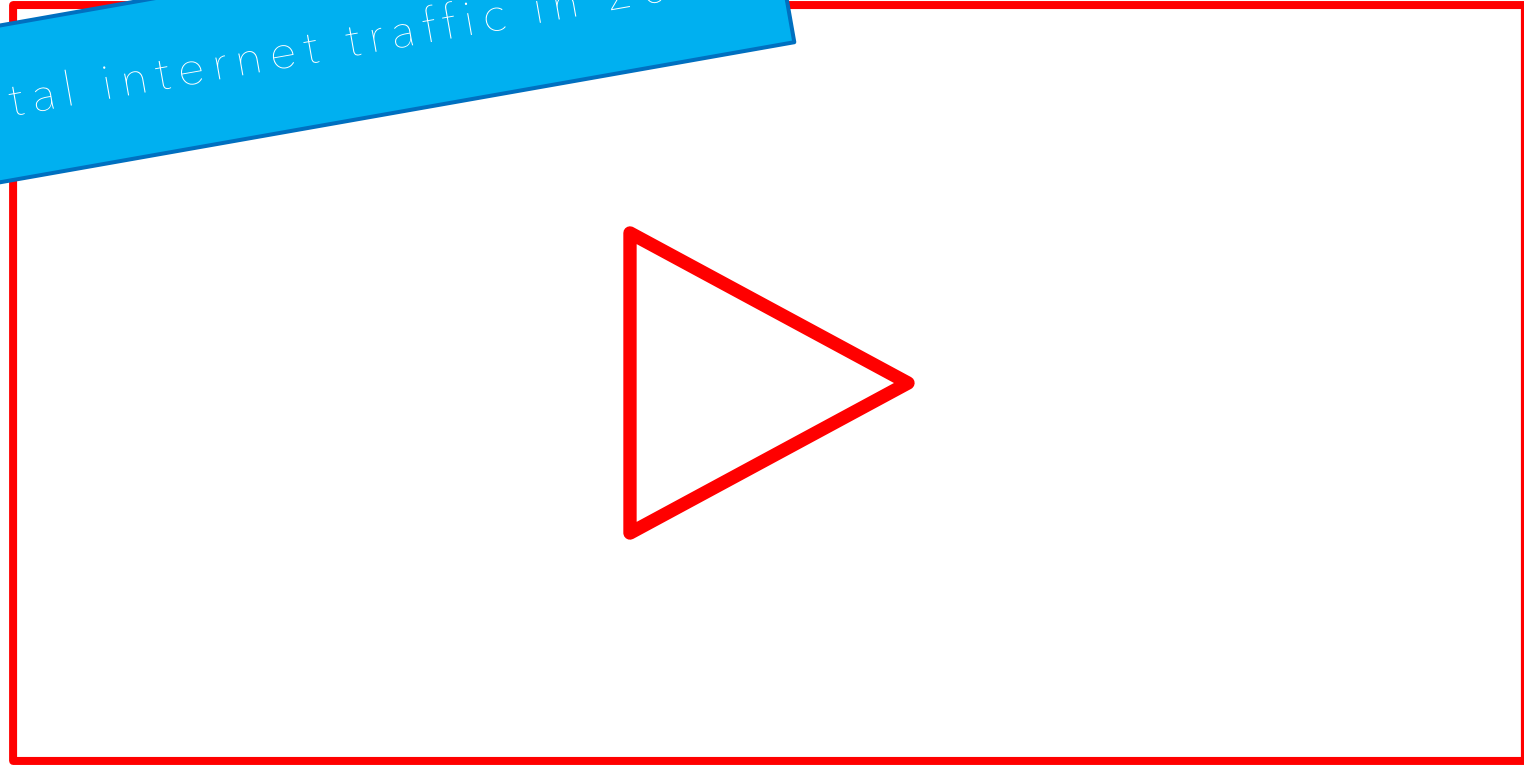# Neural Inter-Frame Compression for Video Coding

Presented by *Maxime Raafat*

*Abdelaziz Djelouah, Joaquim Campos, Simone Schaub-Meyer, Christopher Schroers (Disney Research)*

ICCV, October 27, 2019

75% of total internet traffic in 2017

# Previous work

- Image compression : JPEG, JPEG2000, BPG, WebP

- Neural Image Compression : deep learning for image compression

- Video compression : H.264 (AVC), H.265 (HEVC)

- Neural Video Compression : deep learning for video compression

# Pipeline

Interpolation based video compression technique, compatible with neural image compression methods, while expressing residual information in latent space.
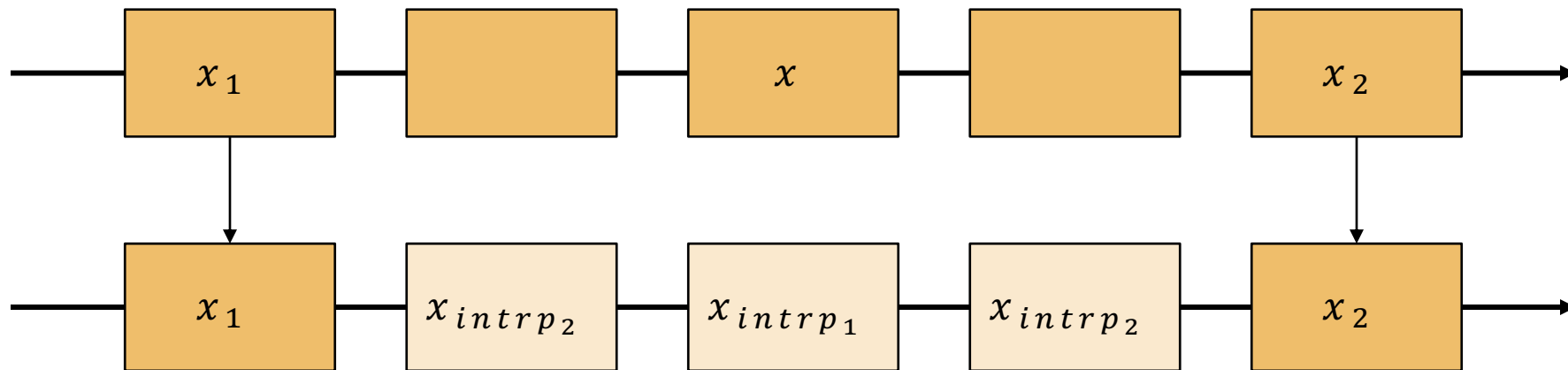
- Interpolation with compression constraints

- Experimental Results

# Interpolation with compression constraints

# Interpolation model

Reference frames (keyframes) $\mathcal{K}_x = \{x_1, x_2, \ldots, x_k\}$

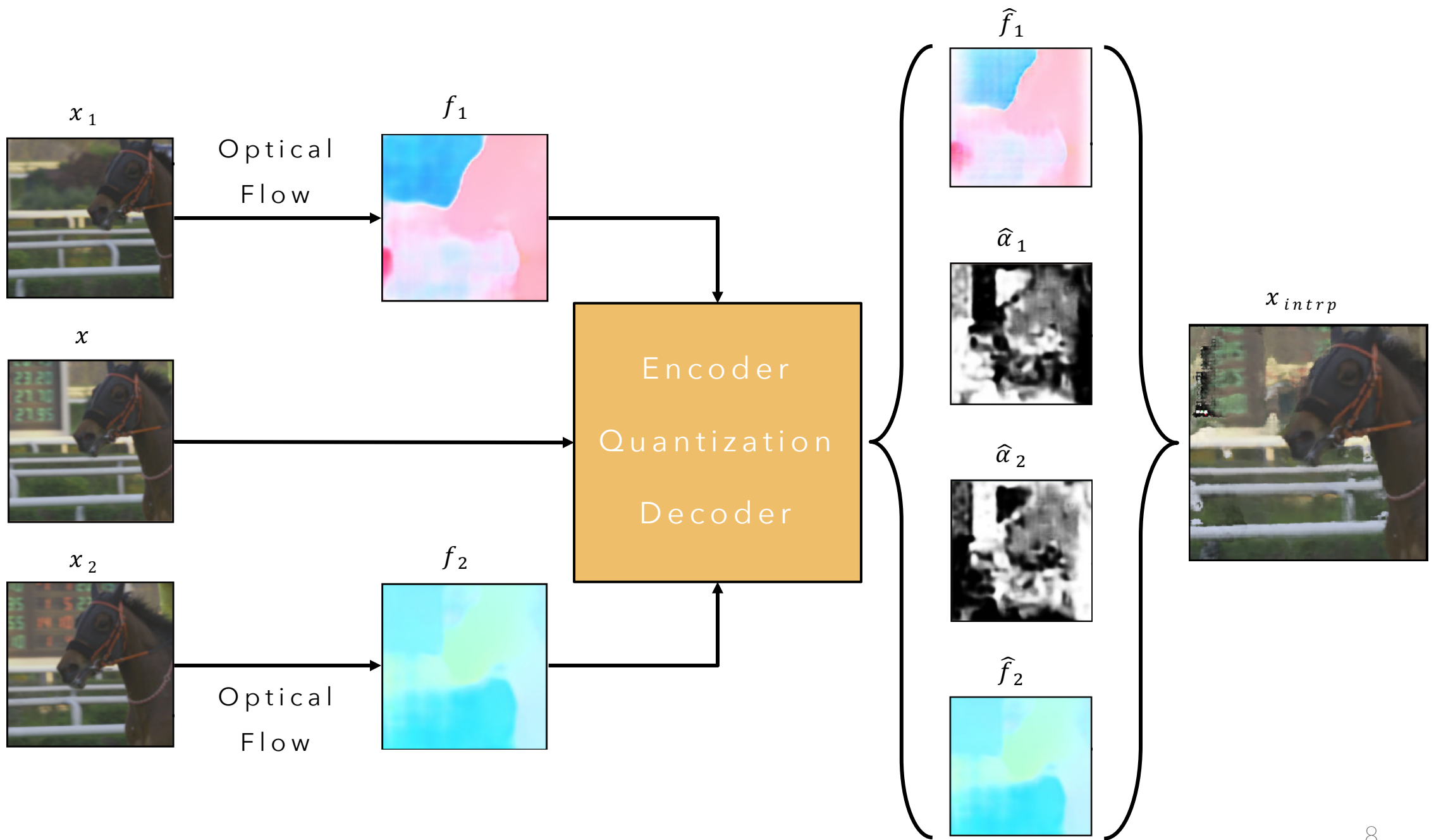Predict intermediate frame $x$

# Interpolation model

Reference frames (keyframes) $\mathcal{K}_x = \{x_1, x_2, \ldots, x_k\}$

Predict intermediate frame $x$

$$x_{intrp} = \sum_{i=1}^{k} \widehat{\alpha}_i \, \omega(x_i, \widehat{f}_i), \text{ with } \sum_{i=1}^{k} \widehat{\alpha}_i = 1$$

$\widehat{f}_i$ = quantized displacement map of $x_i$ w.r.t. to $x$

$\widehat{\alpha}_i$ = quantized blending coefficient of $x_i$

$x_1$

$x$

$x_2$

Optical Flow

Optical Flow

$f_1$

$f_2$

Encoder

Quantization

Decoder

$\hat{f}_1$

$\hat{\alpha}_1$

$\hat{\alpha}_2$

$\hat{f}_2$

$x_{intrp}$

# Compression constraints

Quantized latent representation $\widehat{q}$ should occupy as little storage as possible, while minimizing distortion on the interpolation result.

$$L\left(\phi, \phi', p_{\widehat{q}}\right) = \mathbb{E}_{x \sim p_x}[-log_2 p_{\widehat{q}}(\widehat{q}) + \lambda * d(x, x_{intrp})]$$

$\phi$ and $\phi'$ = encoder-decoder network parameters

$p_{\widehat{q}}$ = entropy model ; $\lambda$ = compression-rate vs distortion regularizer

# Latent space residuals

Minimize the transmitted residual information between $x_{intrp}$ and $x$.

$$r = y - y_{intrp} = g_\phi(x) - g_\phi(x_{intrp})$$

Quantize $r \rightarrow \hat{r}$

$$\hat{x} = g_{\phi'}(y_{intrp} + \hat{r})$$

$g_\phi$ and $g_{\phi'}$ = encoder and decoder

# Latent space residuals

Minimize the transmitted residual information between $x_{intrp}$ and $x$.

$$L(\phi, \phi', p_{\hat{q}}) \leftarrow L(\phi, \phi', p_{\hat{q}}) + \mathbb{E}_{x \sim p_x}[-log_2 p_{\hat{r}}(\hat{r}) + \lambda * d(x, \hat{x})]$$

$\phi$ and $\phi'$ = encoder-decoder network parameters

$p_{\hat{r}}$ = entropy model for residual values

$\lambda$ = regularizer

# Network architectures

- Encoder $g_\phi$ : 5 blocks (each one convolutional* and one Generalized Normalization Transformation layer).

- Decoder $g_{\phi'}$ : 3 RGB output channels for the image $\hat{x}$ and 5 output channels for $x_{intrp}$ ($\hat{f}_1$ and $\hat{f}_2$, as well as $\hat{\alpha}_1$ and $\hat{\alpha}_2 = 1 - \hat{\alpha}_1$).

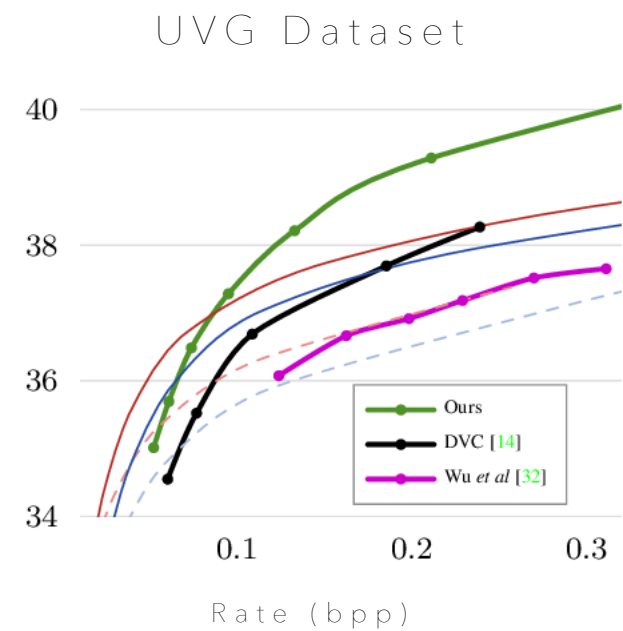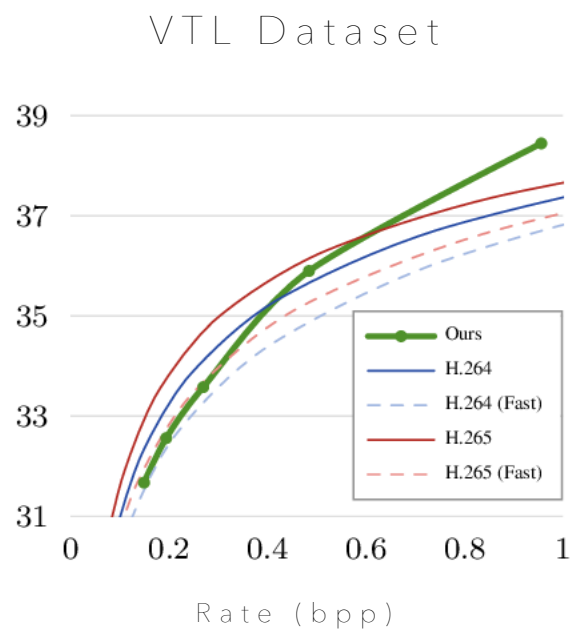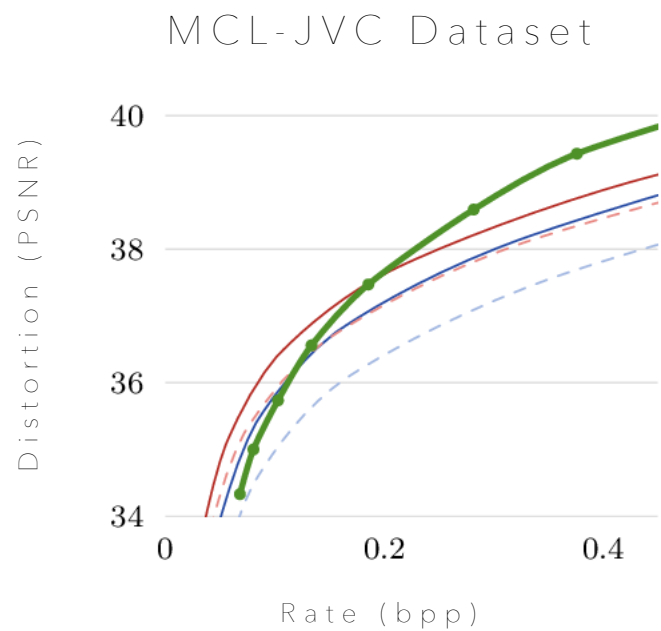- Training performed with the mean squared error (MSE) for the distortion function $d$.

* kernel size $k$ = 5, stride size $s$ = 2

# Experimental Results

# Experimental Setting

- Keyframes positioned every 12 frames

- Peak Signal to Noise Ratio (for distortion measures)

- Training on 3 datasets (max length = 300 frames)
    - MCL-JVC (resolution : 1920 x 1080)
    - VTL : Video Trace Library (resolution : 352 x 288)
    - UVG : Ultra Video Group (resolution: 1920 x 1080)

# Video codec comparisons



MCL-JVC Dataset

VTL Dataset

UVG Dataset

# Advantages of the proposed interpolation
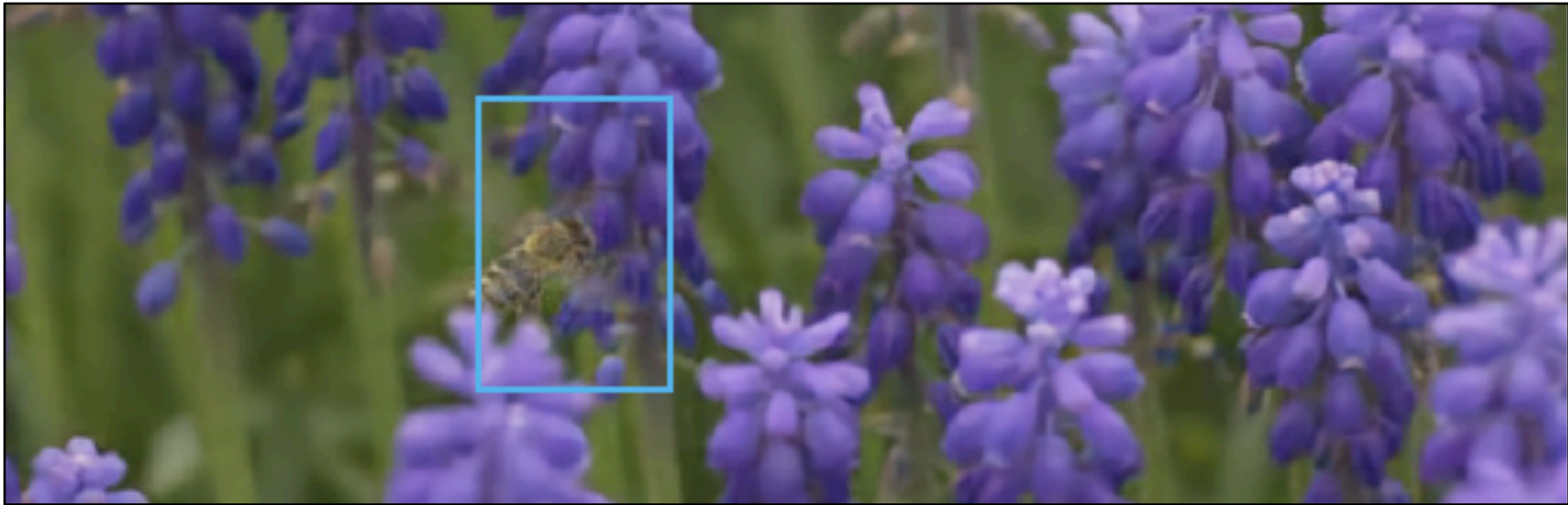
Flow + Interpolation

Ours (lower bit-rate)

Ours (higher bit-rate)



0.028bpp

0.027bpp

0.24bpp

0.024bpp

0.021bpp
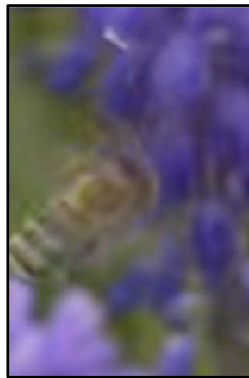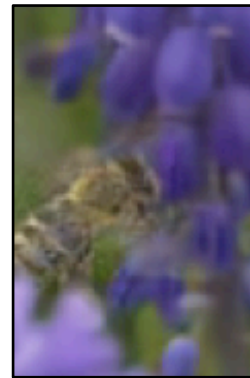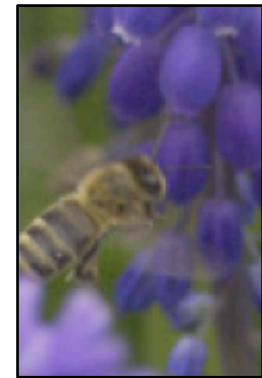
0.27bpp

H.264
0.02bpp

H.265
0.02bpp

Ours
0.02bpp

Ground
Truth

# Questions?