# **Deep Video Codec Control for Vision Models**

 $\bf Christoph Reich^{1,2,3,4}$  $\bf Christoph Reich^{1,2,3,4}$  $\bf Christoph Reich^{1,2,3,4}$ , Biplob Debnath<sup>1</sup>, Deep Patel<sup>1</sup>, Tim Prangemeier $^3$ , Daniel Cremers<sup>2,4</sup>, and Srimat Chakradhar<sup>1</sup>  $^1$ NEC Laboratories America, Inc., <sup>2</sup> TU Munich,  $^3$  TU Darmstadt,  $^4$  Munich Center for Machine Learning (MCML)

# *tl;dr:*

Standard codecs are not optimized for current deep vision models. We present the first end-to-end learnable **Deep Video Codec Control**

# to optimize standard codecs for vision models w/o breaking standardization.

# **Introduction**

- Learn control network  $C_{\theta}$  to predict macroblock-wise quantization parameters QP
- Stay within the available bandwidth budget (rate control)
- Maximize downstream performance of a deep vision model DNN (e.g., DeepLabV3)
- Standard lossy video codecs are part of almost all real-world video processing pipelines
- Standardization is key to ensuring interoperability & low costs in real-world applications
- Existing standard codecs are *not optimized* for current deep vision models
- **We aim to optimize standard video codecs (e.g., H.264) for deep vision models**

**Tab. 1.** High-level comparison of our Deep Video Codec Control with existing approaches.

- We use a loss between the downstream prediction on the *coded* and the *original video* to **maximize vision performance**
- For **rate control**, we penalize the control network using a bandwidth loss
- **Motivation** We regularize the control network to generate a bandwidth close to the target bandwidth **and the surrougate** Control variates theory used for learning the surrogate [\[2\]](#page-0-2)



# **Problem Formulation**

**Goal:** Optimize deep vision performance for a given rate *w/o breaking standardization*

<span id="page-0-1"></span>
$$
\max_{\mathbf{Q} \in \mathcal{D}} \mathbf{M}(\mathbf{DNN}(\mathbf{H}.264(\mathbf{V}, \mathbf{C}_{\theta}(\mathbf{V}, b))))
$$
\n
$$
\mathbf{s}.\mathbf{t}. \tilde{b} \leq b. \tag{1}
$$

- **Fig. 3** Qualitative surrogate model results.
- Our surrogate approximates H.264 coding well
- Relative file size error typically below 5%





TECHNISCHE<br>UNIVERSITÄT<br>DARMSTADT

• Control problem can be formulated as a constrained optimization problem (*cf.* Eq. [\(1\)](#page-0-1))

# **Deep Video Codec Control**

**We train our Deep Video Codec Control end-to-end using a downstream model and our differentiable codec surrogate.**



Tab. 2 Semantic segmentation validation results. BW (acc<sub>bw</sub>) & segmentation accuracies (acc<sub>seg</sub>) for difference BW tolerances reported. Metrics averaged over ten BW conditions.

**Fig. 1** The control network predicts high-dimensional codec parameters for an input clip and a given dynamic bandwidth condition.

# **Differentiable Codec Surrogate**

**We learn a differentiable H.264 surrogate predicting both the coded clip and generated file size/bandwidth.**



- **Fig. 2** Differentiable H.264 codec surrogate model.
- 



How can we learn our Deep Video Codec Control such that performance (Eq. [\(1\)](#page-0-1)) is maximized?

- E Video encoding-decoding is not differentiable
- $\ell$  Reinforcement learning does not scale to large action spaces (high-dimensional QP)
- ∼ Learning on a proxy task is suboptimal [\[3\]](#page-0-3)
- ✓ **Learn a differentiable surrogate model of the video codec**
- **Learn control by using end-to-end learning, optimizing Lagrange function of Eq. [\(1\)](#page-0-1)**

# **Surrogate Results**



# **Conclusion**

- **We demonstrate that learning an end-to-end deep codec control is feasible**
- Our Deep Video Codec Control outperforms traditional rate control modules

## *Future research questions:*

- How to facilitate multiple downstream models and tasks [\[4\]](#page-0-4)
- How to generalize our Deep Video Codec Control to other standard codecs (e.g., H.265 [\[5\]](#page-0-5))







# **Codec Control Results**

We demonstrated the effectiveness of our Deep Video Codec Control on the tasks of semantic segmentation and optical flow estimation (see paper).



# *<b>P* Our Deep Codec Control consistently outperformed 2-pass ABR

 $\mathscr{D}$  We are able to preserve up to 20% in semantic segmentation accuracy

We also analyzed the control performance when transferred between vision tasks

**Tab. 3** Transfer results of our codec control from optical flow estimation to semantic segmentation on Cityscapes. We also report results when directly trained on semantic segmentation.



• Transferring between tasks leads to a drop in downstream performance

**This demonstrates that out control learns a task-specific behavior** 

# **References**

- <span id="page-0-0"></span>[1] T. Wiegand *et al.*, "Overview of the H.264/AVC video coding standard," *IEEE Trans. Circ. Syst. Video Tech.*, vol. 13, no. 7, pp. 560–576, 2003.
- <span id="page-0-2"></span>[2] W. Grathwohl *et al.*, "Backpropagation through the void: Optimizing control variates for black-box gradient estimation," in *ICLR*, 2018.
- <span id="page-0-3"></span>[3] K. Du *et al.*, "AccMPEG: Optimizing video encoding for accurate video analytics," in *MLSys*, vol. 4, 2022, pp. 450–466.
- <span id="page-0-4"></span>[4] Y.-H. Chen *et al.*, "TransTIC: Transferring transformer-based image compression from human perception to machine perception," in *ICCV*, 2023.
- <span id="page-0-5"></span>[5] G. J. Sullivan *et al.*, "Overview of the high efficiency video coding (HEVC) standard," *IEEE Trans. Circ. Syst. Video Tech.*, vol. 22, no. 12, pp. 1649–1668, 2012.

