

Estimation of cost-optimal encoding ladders for tiled 360-degree videos in adaptive streaming systems

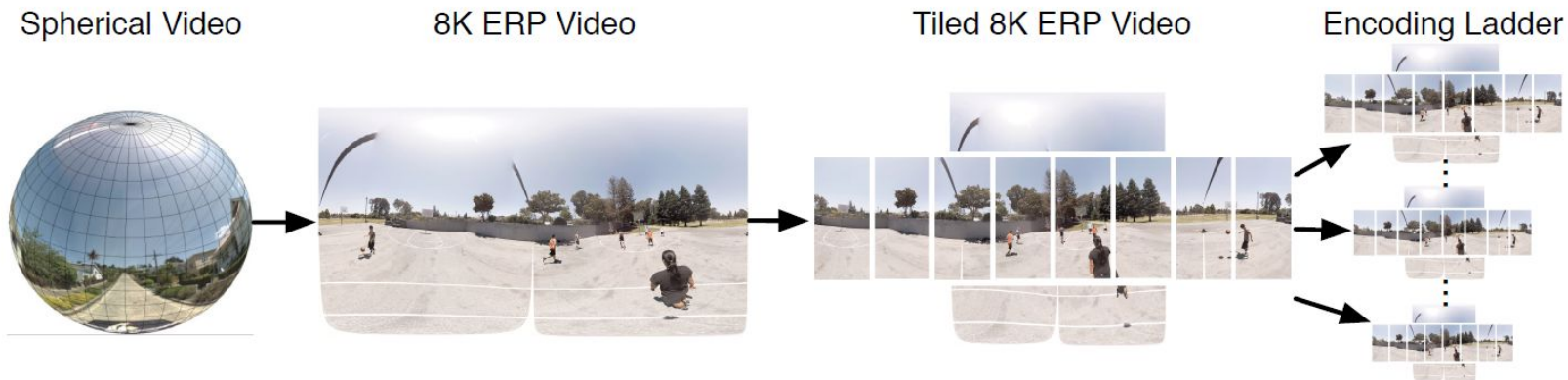
Cagri Ozcinar, Ana De Abreu, Sebastian Knorr, and Aljosa Smolic
Trinity College Dublin (TCD), Dublin 2, Ireland.

Outline

- Introduction
- Problem to be solved
- Proposed system model
 - Classification of the content type
 - Distortion modeling
 - Cost modeling
 - Problem formulation
- Evaluation
- Conclusion

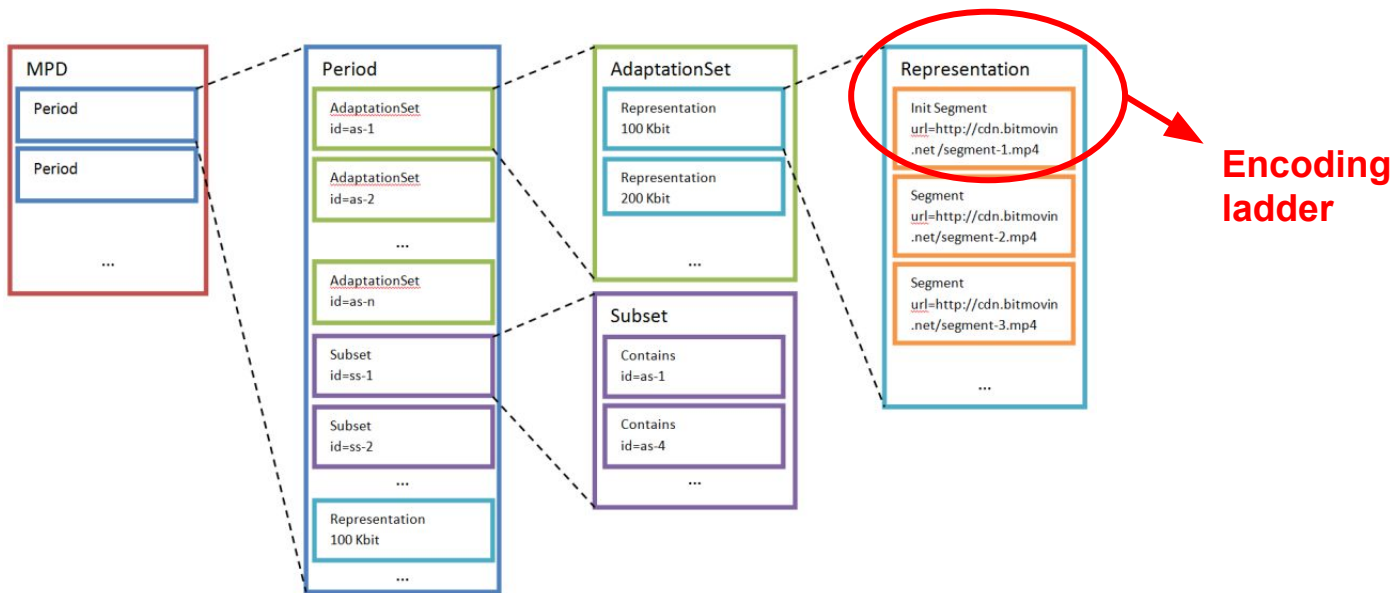
Introduction

- Streaming 360-degree video is challenging
- It is very high resolution, such as 4K, 8K equi-rectangular projection (**ERP**) or higher
- We only stream the user's field of view using tile-based encoding and adaptive streaming (**DASH**)




What do people usually do?

- ERP
- Tile-based encoding
- DASH with its extension, Spatial Relationship Description (SRD)



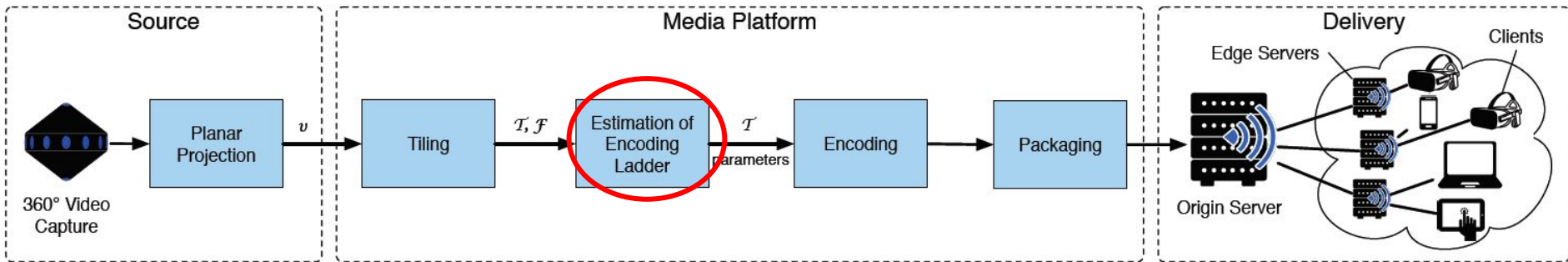
Problem to be solved

- Most recent work focused on the client's perspective without considering the service providers' perspective
- **Client's** perspective: end-users' latency, bandwidth, distortion, QoE
- **Provider's** perspective: computation cost and storage utilization

 Estimation method of cost-optimal encoding ladders in adaptive streaming systems by considering both the **provider's and client's perspective** for tiled 360 video streaming

Proposed system

- Estimation of Encoding ladder contains 4 major components:
 - Classification of the content type
 - Distortion modeling
 - Cost modeling
 - Problem formulation
- Minimize the service **provider's resource costs** while providing high quality 360° video streaming experience (**distortion**)



Classification of the content type

- Spatial complexity and temporal complexity

$$\mathcal{F} = \{f_{spa}, f_{tmp}\}$$

- **2-pass constant rate factor (CRF)** encoding, which has the QPs slightly varied across the time based on the scene complexity, action, and motion
- The average size of I- and P- frames can be used to determine the complexity features
- Content type: o1, o2, and o3 (simple->complex)

Distortion modeling

- Spherical distortion (spherical content onto the planar surface)
- Weighted-to-spherically uniform mean square error (**WS-MSE**) [1]
- The noise power for the i -th representation of the j -th tile

$$d_{ij} = \frac{\sum_{x \in W} \sum_{y \in H} ((t_j(x, y) - \tilde{t}_{ij}(x, y))^2 q_j(x, y))}{\sum_{x \in W} \sum_{y \in H} q_j(x, y)} \quad q_j(x, y) = \cos \frac{(y + 0.5 - H/2)\pi}{H}$$

Reconstructed
Uncompressed
Weighting intensity
Pixel coordinate

Cost modeling - Encoding cost

- Encoding cost c_e can be described for the j -th tile of the i -th representation
- The **cost calculation model** used by the Amazon cloud service [2]

$$c_{ij}^e = \begin{cases} \mu_e, & r_{ij} \leq 720p \\ 2\mu_e, & 720p < r_{ij} \leq 1080p \\ 4\mu_e, & 1080p < r_{ij} \leq 4K \\ 8\mu_e, & 4K < r_{ij} \leq 8K \end{cases}$$

Constant term

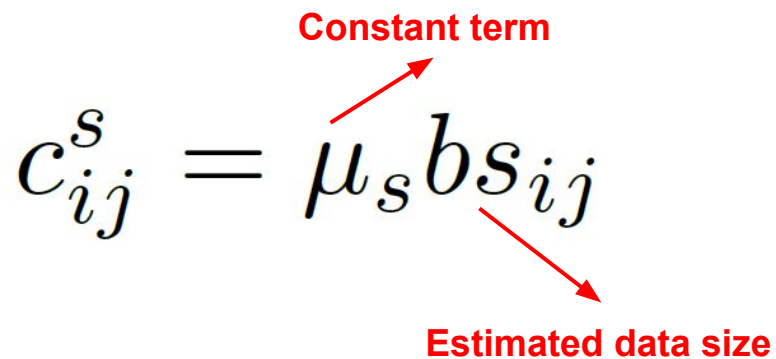
Cost modeling - Storage cost

- Linear cost model
- The **storage cost** for the j -th tile of the i -th representation

$$c_{ij}^s = \mu_s b s_{ij}$$

Constant term

Estimated data size



Problem formulation

- To minimize both the total spherical distortion and the total resource cost
- Constraints
 - **Bandwidth**, a set of given network bandwidth profiles $\{P\}$
 - **Computational and storage costs**, limitations for the encoding and storage costs
 - **Encoding rate**, the bitrate levels of the representations

$$\mathcal{L}^* : \operatorname{argmin}_{\mathcal{L}} \sum_{i \in \mathcal{L}} \sum_{p \in \mathcal{P}} (\gamma c_i + (1 - \gamma) d_i) a_{ip}$$

Pre-defined constant $[0,1]$

Decision variable $\{0,1\}$

Evaluation settings

Apple [15]		Axinom [21]		Netflix [16]	
Z (Mbps)	$W \times H$	Z (Mbps)	$W \times H$	Z (Mbps)	$W \times H$
45	8192×4096	45	8192×4096	43	8192×4096
30	8192×4096	30	8192×4096	30	4096×2048
20	4096×2048	21	4096×2048	23.5	4096×2048
11	3072×1536	12	3072×1536	17.5	3072×1536

- Eight 8Kx4K resolution 360° ERP video test sequences
 - $V = \{\text{Train, Stitched left Dancing360 8K, Basketball, KiteFlite, ChairLift, SkateboardInLot}\}$
- Each video was split into $N = 10$ tiles
- Content type = $\{O\} = \{o_1, o_2, o_3\}$
- Resolution = $\{G\} = \{g_1, g_2, g_3\} = \{3072 \times 1536, 4096 \times 2048, 8192 \times 4096\}$
- Bandwidth = $\{P\} = \{p_1, p_2, p_3, p_4\}$

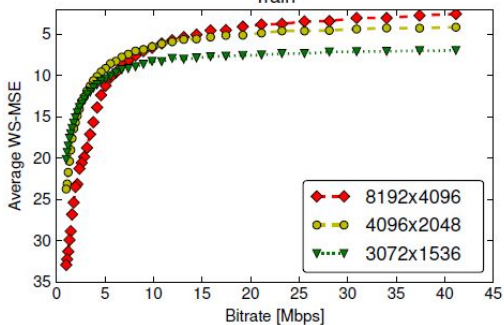
Sequence	f_{spa}	f_{tmp}	\mathcal{O}
<i>Train</i>	234	0.501	o_1
<i>Stitched_left_Dancing360_8K</i>	313	0.501	
<i>Basketball</i>	1167	0.502	o_2
<i>KiteFlite</i>	1547	0.502	
<i>ChairLift</i>	2842	0.502	o_3
<i>SkateboardInLot</i>	3977	0.502	

Evaluation results

- WS-MSE versus bitrate (in Mbps) performance
 - Each content type has various content dependencies for each encoding resolution and bitrate
 - The lowest complex encoding features, achieves a low distortion score
 - High-resolution version has a higher **sensitivity** for unpredictable motions, which requires further residuals to avoid visual distortions

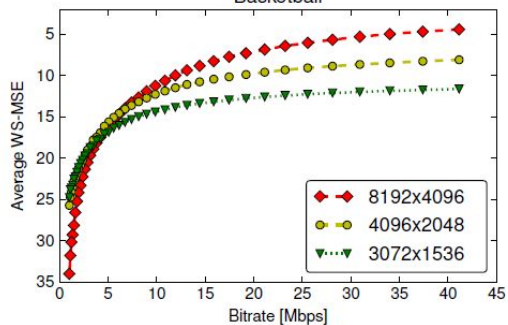
o1

Train



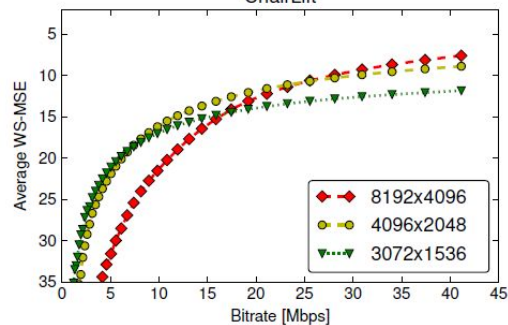
o2

Basketball



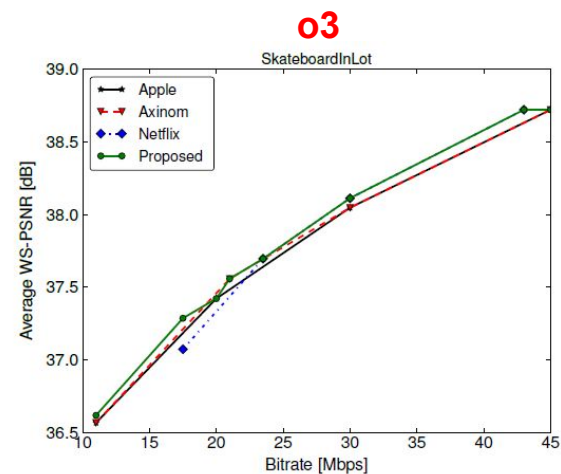
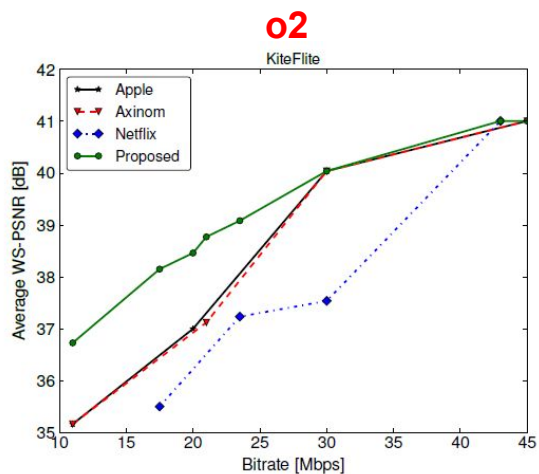
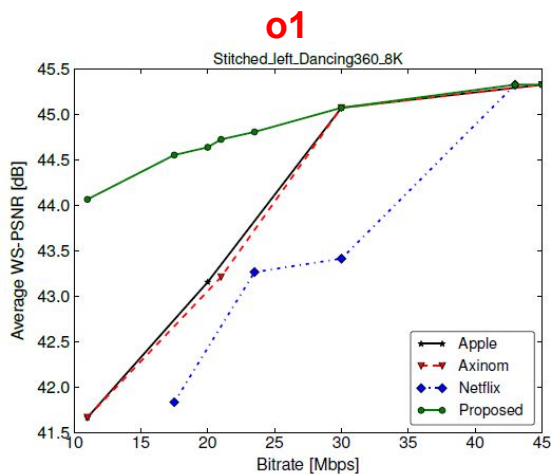
o3

ChairLift



Evaluation results

- RD performance gain
 - Proposed method considerably increases the objective video quality (i.e., **WS-PSNR**) [1]
 - High bitrate savings between 10-30 Mbps bandwidth ranges for the content types o1 and o2



Evaluation results (cont.)

- Bjøntegaard metric [1] (**BD-rate**)
- A negative BD-rate indicates a decrease of bitrate at the same quality
- Proposed method provides considerable bitrate savings compared to the recommended encoding ladders at the same bitrates.

Sequence v	Streaming vendor		
	Apple	Axinom	Netflix
<i>Stitched_left_Dancing360_8K</i>	-4.838	-7.070	-1.102
<i>KiteFlite</i>	-13.937	-20.395	-68.299
<i>SkateboardInLot</i>	-1.659	-1.094	-1.144

Table 5: BD-rate saving (%) of the proposed method.

$$\mathcal{L}^* : \operatorname{argmin}_{\mathcal{L}} \sum_{i \in \mathcal{L}} \sum_{p \in \mathcal{P}} (\gamma c_i + (1 - \gamma) d_i) a_{ip}$$

Evaluation results (cont.)

- Content type o1 increases its encoding resolution and decreases its target encoding rate
- Content type o3, decreases both its encoding resolution and target encoding rate
- Resolution = {G} = {g1,g2,g3} = {3072x1536, 4096x2048, 8192x4096}

Sequence v	γ	Representation i											
		1	2	3	4	5	6	7	8	9	10	11	12
<i>Stitched_Left_Dancing360_8K</i>	0.0	(g ₁ ,1.47)	(g ₁ ,1.78)	(g ₁ ,2.15)	(g ₁ ,3.8)	(g ₁ ,4.6)	(g ₁ ,5.6)	(g ₂ ,10.84)	(g ₂ ,13.11)	(g ₂ ,15.87)	(g ₂ ,28.11)	(g ₃ ,34.01)	(g ₃ ,41.15)
	0.1	(g ₂ ,1.34)	(g ₂ ,1.61)	(g ₂ ,1.95)	(g ₂ ,2.60)	(g ₃ ,3.14)	(g ₃ ,3.80)	(g ₃ ,6.12)	(g ₃ ,7.40)	(g ₃ ,8.96)	(g ₃ ,17.45)	(g ₃ ,21.12)	(g ₃ ,25.55)
	0.5	(g ₂ ,1.00)	(g ₂ ,1.21)	(g ₂ ,1.47)	(g ₂ ,2.36)	(g ₃ ,2.86)	(g ₃ ,3.46)	(g ₃ ,6.12)	(g ₃ ,7.40)	(g ₃ ,8.96)	(g ₃ ,17.45)	(g ₃ ,21.12)	(g ₃ ,25.55)
<i>KiteFlite</i>	0.0	(g ₁ ,1.47)	(g ₁ ,1.78)	(g ₂ ,2.15)	(g ₂ ,3.80)	(g ₂ ,4.60)	(g ₃ ,5.56)	(g ₃ ,10.84)	(g ₃ ,13.11)	(g ₃ ,15.87)	(g ₃ ,28.11)	(g ₃ ,34.01)	(g ₃ ,41.15)
	0.1	(g ₁ ,1.47)	(g ₁ ,1.78)	(g ₂ ,2.15)	(g ₂ ,3.80)	(g ₂ ,4.60)	(g ₃ ,5.56)	(g ₃ ,6.73)	(g ₃ ,8.14)	(g ₃ ,9.85)	(g ₃ ,17.45)	(g ₃ ,21.12)	(g ₃ ,25.55)
	0.5	(g ₁ ,1.00)	(g ₁ ,1.21)	(g ₁ ,1.47)	(g ₂ ,2.36)	(g ₂ ,2.86)	(g ₂ ,3.46)	(g ₃ ,6.12)	(g ₃ ,7.40)	(g ₃ ,8.96)	(g ₃ ,17.45)	(g ₃ ,21.12)	(g ₃ ,25.55)
<i>SkateboardInLot</i>	0.0	(g ₁ ,1.47)	(g ₁ ,1.78)	(g ₁ ,2.15)	(g ₁ ,3.80)	(g ₁ ,4.60)	(g ₁ ,5.56)	(g ₂ ,10.84)	(g ₂ ,13.11)	(g ₂ ,15.87)	(g ₂ ,28.11)	(g ₃ ,34.01)	(g ₃ ,41.15)
	0.1	(g ₁ ,1.47)	(g ₁ ,1.78)	(g ₁ ,2.15)	(g ₁ ,2.86)	(g ₁ ,3.46)	(g ₁ ,4.18)	(g ₁ ,6.12)	(g ₁ ,7.40)	(g ₁ ,8.96)	(g ₁ ,17.45)	(g ₂ ,21.12)	(g ₂ ,25.55)
	0.5	(g ₁ ,1.21)	(g ₁ ,1.47)	(g ₁ ,1.78)	(g ₁ ,2.36)	(g ₁ ,2.86)	(g ₁ ,3.46)	(g ₁ ,6.12)	(g ₁ ,7.40)	(g ₁ ,8.96)	(g ₂ ,17.45)	(g ₂ ,21.12)	(g ₂ ,25.55)

Table 6: Results of the proposed encoding ladder estimation for $\gamma = 0$, $\gamma = 0.1$, and $\gamma = 0.5$.

Conclusion

- A novel encoding ladder estimation method for tiled 360 video streaming systems, considering both the provider's and client's perspectives
- Proposed method provides **cost-optimal** and enhanced video streaming **experiences** for VR end-users

Sequence v	$\Delta\text{cost} (\%)$		$\Delta\text{distortion} (\%)$	
	$\gamma = 0.1$	$\gamma = 0.5$	$\gamma = 0.1$	$\gamma = 0.5$
<i>Stitched_left_Dancing360_8K</i>	37.463	39.683	-13.628	-42.914
<i>KiteFlite</i>	33.165	39.206	-9.564	-25.326
<i>SkateboardInLot</i>	37.214	38.884	-8.977	-15.26

Table 7: Total cost saving and distortion gain with respect to $\gamma=0.0$.