# **Panoramic Vision Transformer for** Saliency Detection in 360° Videos

Heeseung Yun, Sehun Lee, Gunhee Kim





## 360° Video Saliency Detection

- One of the core problems in 360° video understanding
  - Captures whole surroundings of a scene instead of specific predetermined context
- Directly related to practical scenarios
  - 360° video summarization (*i.e.*, cinematography)
  - Dynamic rendering for VR



## 360° Video Saliency Detection

- A simple tricky question
  - "Which direction to watch if you were in the scene?"
- Problem 1. architecture design for panoramic videos
  Taking distortion & discontinuity into account
- Problem 2. modeling "saliency" in panoramic videos
  - Often ambiguous and subjective

- Ignoring geometric property
  - 😫 Simple
  - Dense NFoV projection is not scalable



- Using designated architecture
  - Well-tailored for spherical input format
  - **Not transferrable**



Esteves et al. Learning SO(3) Equivariant Representations with Spherical CNNs. In ECCV 2018. Lee et al. SpherePHD: Applying CNNs on a Spherical PolyHeDron Representation of 360 Images. In CVPR 2019.

- Additional modules & training for geometric adaptation
  - Transferrable architecture
  - Training overhead, layerwise error accumulation



- Previous approaches
  - Ignoring geometric property
  - Using designated architecture
  - Additional modules & training for geometric adaptation
- Solution: Local patch-based modeling
  - Geometry-aware approximation
  - Transferrable from most Vision Transformer (ViT) variants
  - Wo additional modules or training for geometric adaptation
  - ( Applicable to various 360° video formats

#### Modeling Saliency in Panoramic Videos

- Saliency itself is a longstanding question in CV
  - Saliency as self-information, anomaly, class activation, etc.
- "Which direction to watch if you were in the scene?"
  - Ambiguity and subjectivity
  - Supervised learning usually inapplicable
- Solution: leverage features from NFoV domain
  - Rich and readily available pretrained knowledge
  - Spatio-temporal feature consistency suffices for 360° video saliency

- Panoramic Vision Transformer for 360° videos
  - First to adopt ViT to encode omnidirectional imagery
  - No additional module, trivial overhead, video format-agnostic
  - Do not require class activation, optical flow for saliency
  - Competitive results in 360° video saliency detection & quality assessment



#### • Panoramic Vision Transformer: Encoder



• Panoramic Vision Transformer: Decoder



• Panoramic Vision Transformer: Decoder



• 360° video saliency detection (Wild360)



Cheng et al. Cube Padding for Weakly-Supervised Saliency Prediction in 360 Videos. In CVPR 2018.

• 360° video saliency detection (Wild360)



• 360° video saliency detection (Wild360)



Saliency score decomposition





Transfer from different pretrained knowledge



Applying to various 360° video formats



- Video quality assessment of 360° videos (VQA-ODV)
  - Relevant to rendering quality control in VR pipeline



Li et al. Bridge the Gap Between VQA and Human Behavior on Omnidirectional Video. In ACM MM 2018.

# **Concluding Remarks**

- Local patch-based architecture for 360° videos
  - First attempt to adopt ViT to encode omnidirectional imagery
  - No additional module, trivial overhead
- Spatio-temporal consistency of local patch features for measuring saliency
  - Leverage pretrained knowledge
  - Independent of class activation, optical flow, etc.
- Competitive results in 360° video saliency detection & quality assessment



