

The SFU logo consists of the letters 'SFU' in a white, bold, sans-serif font, centered within a solid red square.

SFU

SIMON FRASER UNIVERSITY
ENGAGING THE WORLD

OmniSense: Towards Edge-Assisted Online Analytics for 360-Degree Videos

Miao Zhang, Yifei Zhu, Linfeng Shen, Fangxin Wang, **Jiangchuan Liu**

BACKGROUND

Omnidirectional cameras have become increasingly affordable



Insta360 One X2

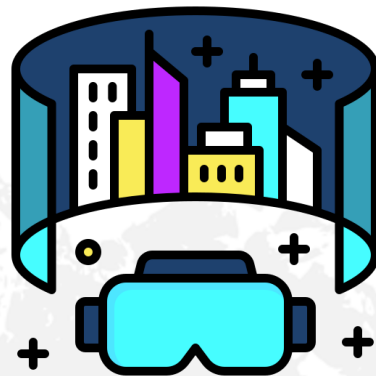


GoPro Max

Common Usage:

Entertainment:

Record videos for human viewers



Potential Usage:

Perception and Interaction:

Analyze videos for full situational awareness without blind spots



BACKGROUND

Video Analytics (VA): Evolving From Regular Videos to 360-Degree Videos

VA for regular videos:

Challenge 1: high-accuracy VA requires running **compute-intensive** deep neural networks (DNNs)

Solution: **offload** VA tasks to resource-rich edge or cloud servers

Challenge 2: offloading video data over networks requires **tremendous bandwidth**

Solution: **encoding configuration tuning, frame filtering, DNN model splitting**

Immersive VA for 360-degree videos:

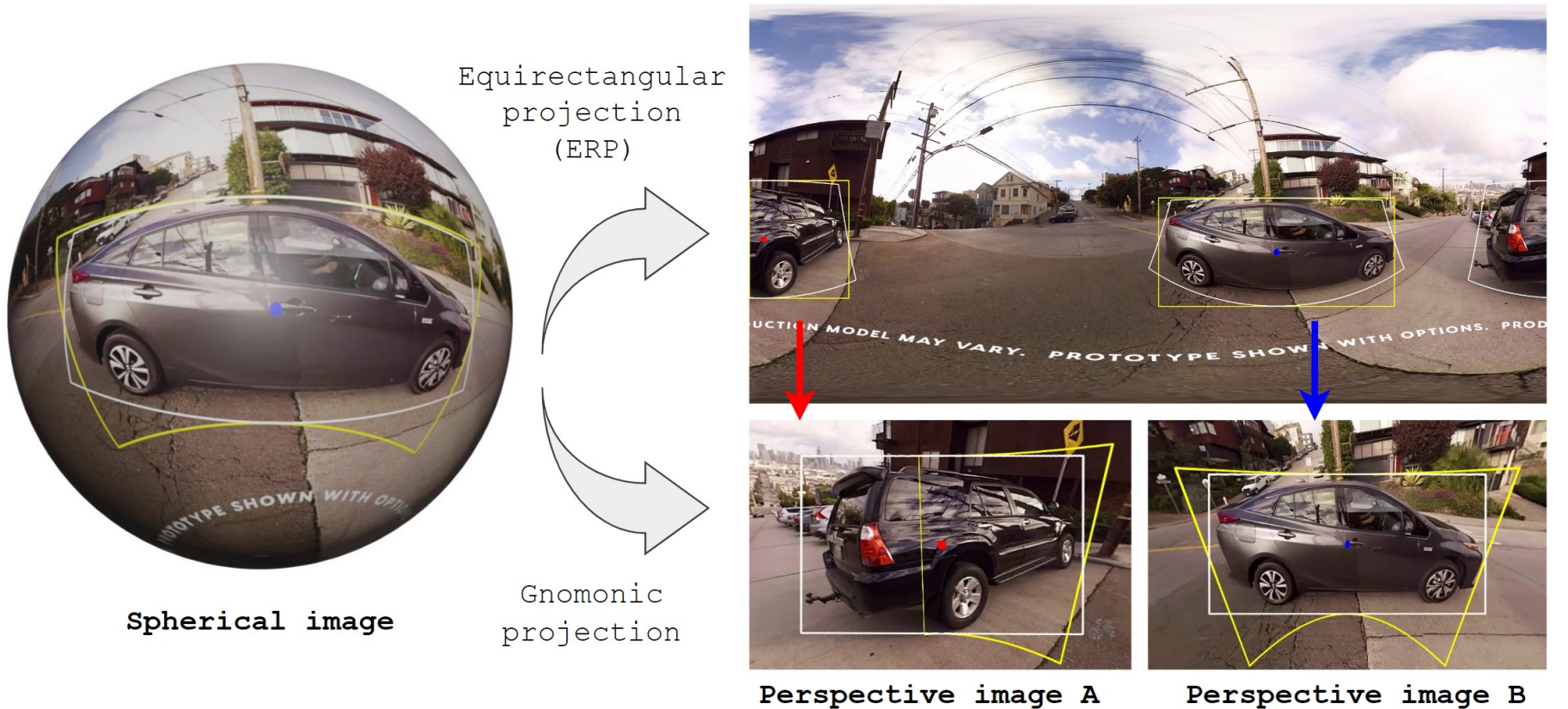
Challenge 1: 4 – 6 x large than regular videos. Much more **compute-intensive** and **bandwidth-intensive**

Challenge 2: particular geometry structure (spherical image)

BACKGROUND

Immersive VA – Potential Solution Discussion

- Direct inference on ERP images?
Distortion and discontinuity can hurt accuracy.
- Analyze many perspective images (PIs) to minimize distortion while covering the entire sphere?
Resource-intensive and time-consuming
- Design DNNs specialized for spherical geometry? Cannot directly benefit from advances in off-the-shelf vision models. Require extra efforts.



MOTIVATION

□ Measurement Setup

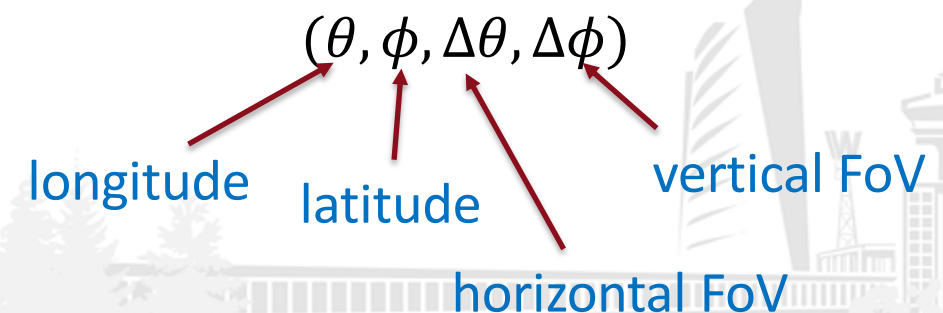
Vision Task (model): object detection (scaled-YOLOv4, which provides a set of model variants)

Video dataset: self-collected UHD 360-degree video dataset

Model Name (Index)	Model Size	Input Size
YOLOv4-Tiny-416 (1)	23 MB	416 x 416
YOLOv4-CSP-512 (2)	202 MB	512 x 512
YOLOv4-CSP-640 (3)	202 MB	640 x 640
YOLOv4-P5 (4)	271 MB	896 x 896
YOLOv4-P6 (5)	487 MB	1280 x 1280

Immersive Video Name	Video Source	Resolution	Frames
New-Orleans-drive	YouTube	7680 x 3840	2100
Expressway-drive	YouTube	5760 x 2880	2100
Chicago-drive	YouTube	7680 x 3840	2100
Sunny-walk1	Self-captured	5376 x 2688	2100
Sunny-walk2	Self-captured	5376 x 2688	2100
Cloudy-walk	Self-captured	5376 x 2688	2100

Immersive object detection criteria:
Spherical bounding box (SphBB): represented by a spherical region (SR)



MOTIVATION

□ Measurement Setup

Ground-truth annotation: self-developed automated annotation method



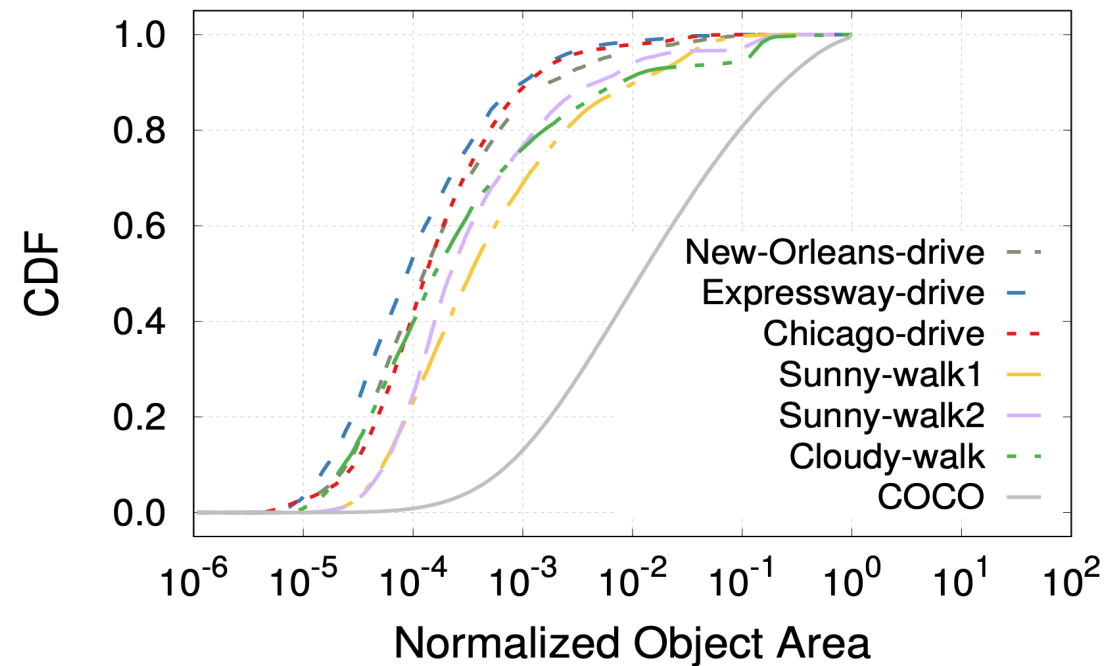
Frame source: Sunny-walk2



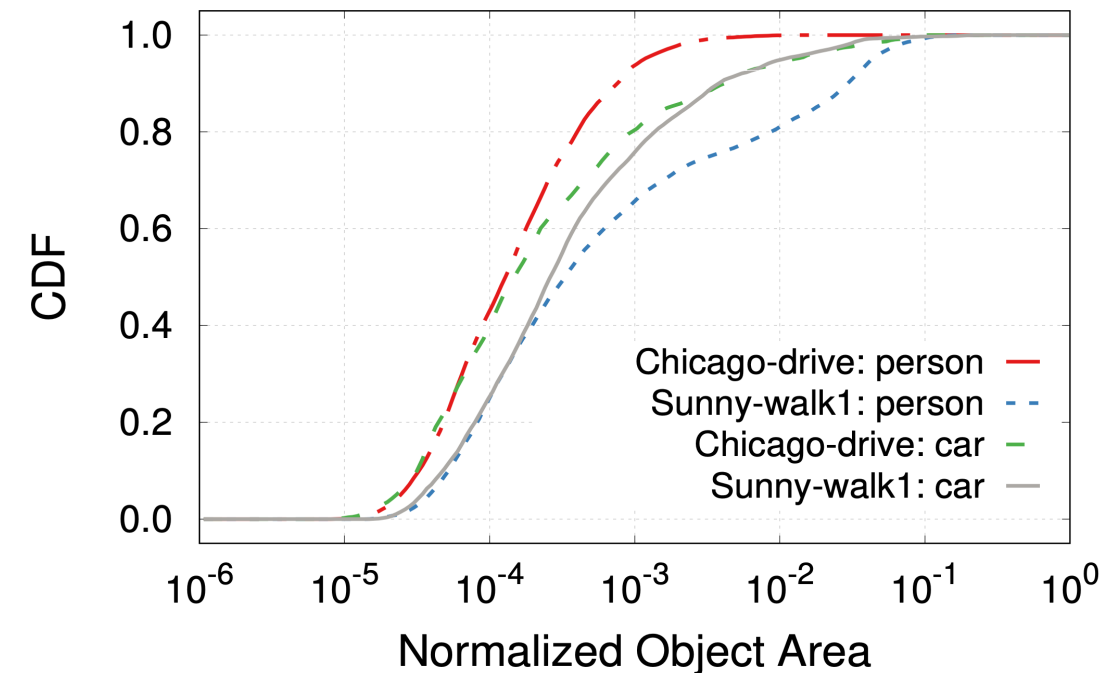
Frame source: Chicago-drive

MOTIVATION

Measurement Insights



CDFs of NOA for 360-degree videos and 2D images



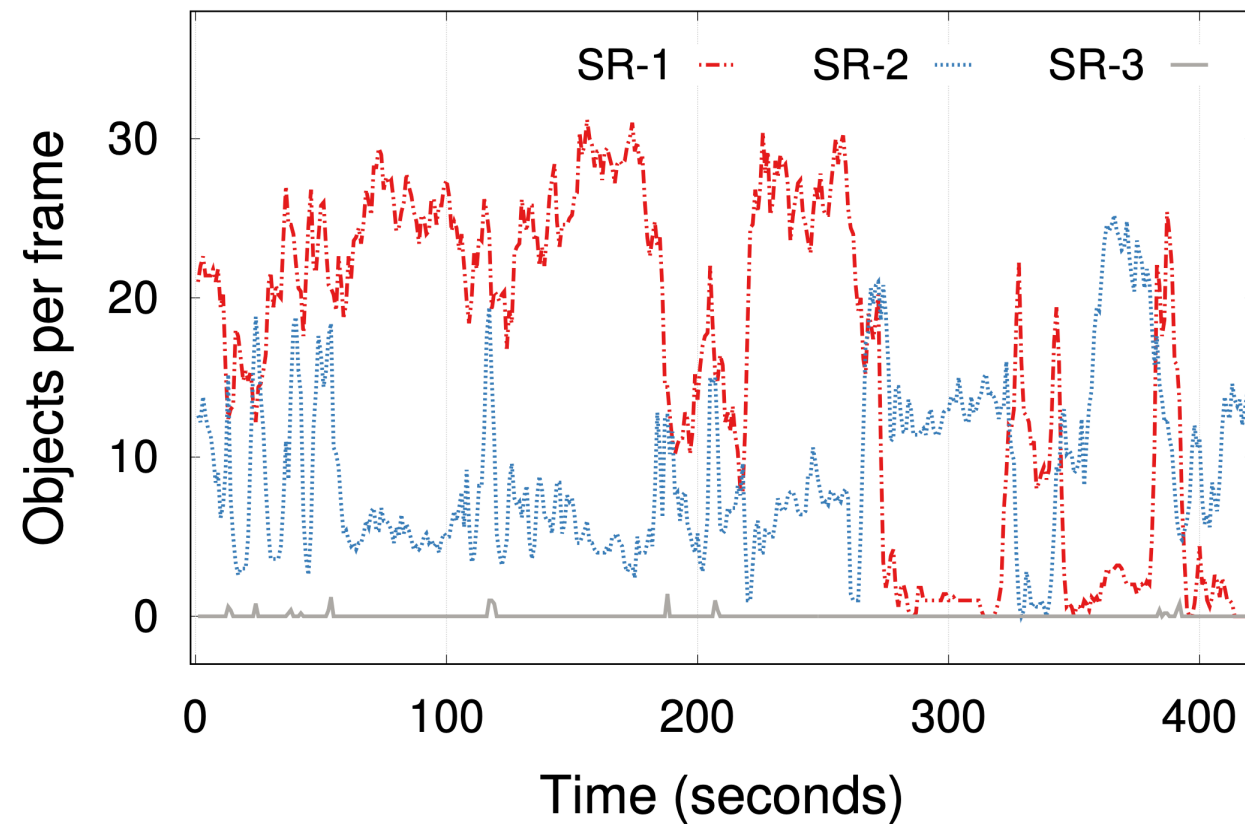
CDFs of NOA for two specific object categories

Most objects in 360-degree videos only occupy a **tiny** area of a frame, which **cannot** be handled by off-the-shelf vision models.

Both **object size** and **category** are crucial reference factors in characterizing video content.

MOTIVATION

□ Measurement Insights



Object number variations in different spherical regions (video: New-Orleans-drive).

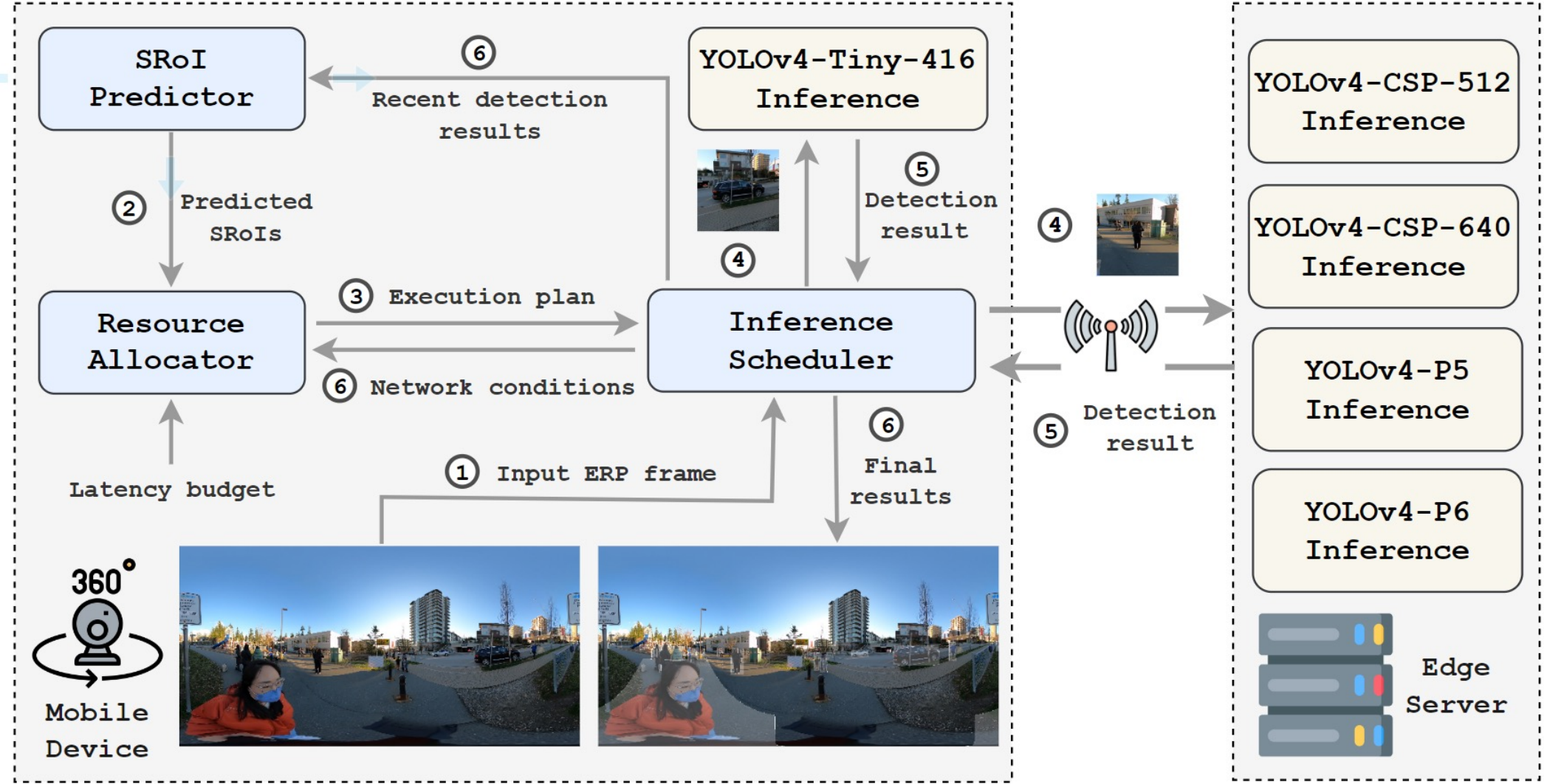
The spatial distribution of objects is **biased** and **highly dynamic**.



- Prune useless pixels before offloading.
- Apply **different** models to different SRs.
- Adapt** to content variations.

SYSTEM DESIGN

Overview of OmniSense



□ Lightweight SRoI Prediction

Motivation: consecutive frames have the **smallest** content differences.

Solution: Propose a **lightweight** spherical RoI (SRoI) prediction algorithm based on **the most recent** detection results. Design a **spherical object discovery** mechanism to discover new objects.

Algorithm 1: SRoI Prediction Algorithm

Input: $f; \gamma; O$ (detected objects of the most recent δ frames)
Output: A set of predicted SRoIs \mathcal{R}

```
1 Initialize SRoI sets  $\mathcal{S} \leftarrow \emptyset, \mathcal{S}' \leftarrow \emptyset$  ;
2 Get the number of all historical objects  $N \leftarrow |O|$  ;
3 foreach object  $o \in O$  do
4   if o can be covered by an  $f \times f$  FoV then
5     merged  $\leftarrow \text{False}$  ;
6     foreach SRoI  $s \in \mathcal{S}$  do
7       hFoV, vFoV  $\leftarrow$  merged horizontal and vertical
8         FoVs for the set  $s.objects \cup \{o\}$ ;
9       if hFoV  $< f$  and vFoV  $< f$  then
10        s.objects  $\leftarrow s.objects \cup \{o\}$  ;
11        s.FoV  $\leftarrow (hFoV, vFoV)$  ;
12        merged  $\leftarrow \text{True}$ ; break ;
13   if not merged then
14     new_s  $\leftarrow$  create a new SRoI with o ;
15      $\mathcal{S} \leftarrow \mathcal{S} \cup \{new\_s\}$  ;
16   else
17     Create a new special SRoI  $s'$  with o ;
18     s'.center  $\leftarrow o.center$ ; s'.FoV  $\leftarrow \gamma \times o.FoV$  ;
19     Calculate content characteristics s'.ccv based on o ;
20     s'.alpha  $\leftarrow 1 / N$  ;
21      $\mathcal{S}' \leftarrow \mathcal{S}' \cup \{s'\}$  ;
22 foreach SRoI  $s \in \mathcal{S}$  do
23   Calculate SRoI center s.center according to s.FoV ;
24   Calculate s.ccv based on s.objects;
25   s.alpha  $\leftarrow |s.objects| / N$  ;
26   s.FoV  $\leftarrow (f, f)$  ;
27  $\mathcal{R} \leftarrow \mathcal{S}' \cup \mathcal{S}$ ;
28 return  $\mathcal{R}$  ;
```

□ Content-Specific Model Performance Estimation

Model inference latency : offline profiling

Model accuracy estimation: Accuracy varies with the analyzed image content.

Quantify the detection capability of a model. Define the **general accuracy vector (gav)**:

$$\mathbf{A}_i = [a_i^{s1}, \dots, a_i^{sn}, a_i^{m1}, \dots, a_i^{mn}, a_i^{l1}, \dots, a_i^{ln}]$$

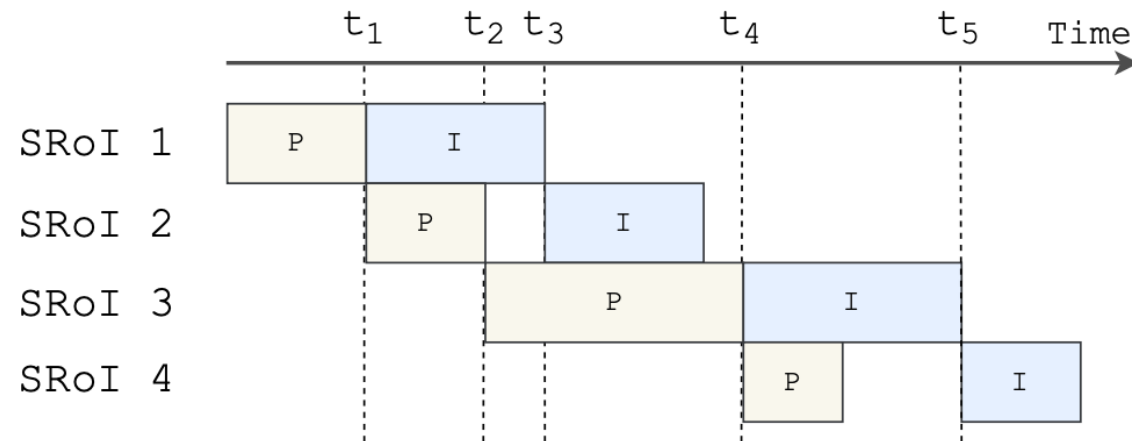
Quantify the content characteristics of an SRol. Define the **content characteristics vector (ccv)**:

$$\mathbf{P}_j = [p_j^{s1}, \dots, p_j^{sn}, p_j^{m1}, \dots, p_j^{mn}, p_j^{l1}, \dots, p_j^{ln}]$$

Estimate detection accuracy of model i on SRol j : $\mathbf{A}_i \cdot \mathbf{P}_j$

□ Latency-Constrained Model Allocation

$$\begin{aligned} & \max_{x_{i,j}} \sum_{j \in \mathcal{R}} \sum_{i \in \mathcal{M}} A_{i,j} \cdot x_{i,j} \\ & s.t. \begin{cases} \mathcal{L}(\mathcal{X}) \leq T, & \mathcal{X} = \{x_{i,j} \mid i \in \mathcal{M}, j \in \mathcal{R}\} \\ \sum_{i \in \mathcal{M}} x_{i,j} = 1, & \forall j \in \mathcal{R} \\ x_{i,j} \in \{0, 1\}, & \forall i \in \mathcal{M}, \forall j \in \mathcal{R} \end{cases} \end{aligned}$$



Algorithm 2: Dynamic Programming Algorithm

Input: $\{A_{i,j}\}; \{d_{i,j}\}; \{d_{i,j}^P\}; \{d_{i,j}^I\}; T$
Output: The optimal execution plan

- 1 $\mathcal{S}(1) \leftarrow \emptyset$;
- 2 **foreach** *model* $i \in \mathcal{M}$ **do**
- 3 **if** $d_{i,1} \leq T$ **then**
- 4 $\mathcal{S}(1) \leftarrow \mathcal{S}(1) \cup \{(A_{i,1}, d_{i,1}^P, d_{i,1}^I, [i])\}$;
- 5 **for** $j = 1$ **to** $r - 1$ **do**
- 6 $\mathcal{S}(j+1) \leftarrow \emptyset$;
- 7 **foreach** *quaternion* $(v, t^P, t, m_list) \in \mathcal{S}(j)$ **do**
- 8 **foreach** *model* $i \in \mathcal{M}$ **do**
- 9 $cur_t \leftarrow \max(t^P + d_{i,j+1}^P, t + d_{i,j+1}^I)$;
- 10 **if** $cur_t \leq T$ **then**
- 11 $cur_v \leftarrow v + A_{i,j+1}$;
- 12 $cur_t^P \leftarrow t^P + d_{i,j+1}^P$;
- 13 $m_list.append(i)$;
- 14 $\mathcal{S}(j+1) \leftarrow$
 $\mathcal{S}(j+1) \cup \{(cur_v, cur_t^P, cur_t, m_list)\}$;
- 15 Remove dominated execution plans from $\mathcal{S}(j+1)$;
- 16 Return the execution plan with the highest v in $\mathcal{S}(r)$;

EVALUATION

□ System Implementation

Mobile device: Nvidia Jetson TX2; Edge server: a desktop with an Nvidia GeForce GTX 1080Ti

The mobile device and server are connected by ASUS AC 1900 router.

OpenCV for Image and video operations; ZeroMQ for data transmission

□ Experimental Setup

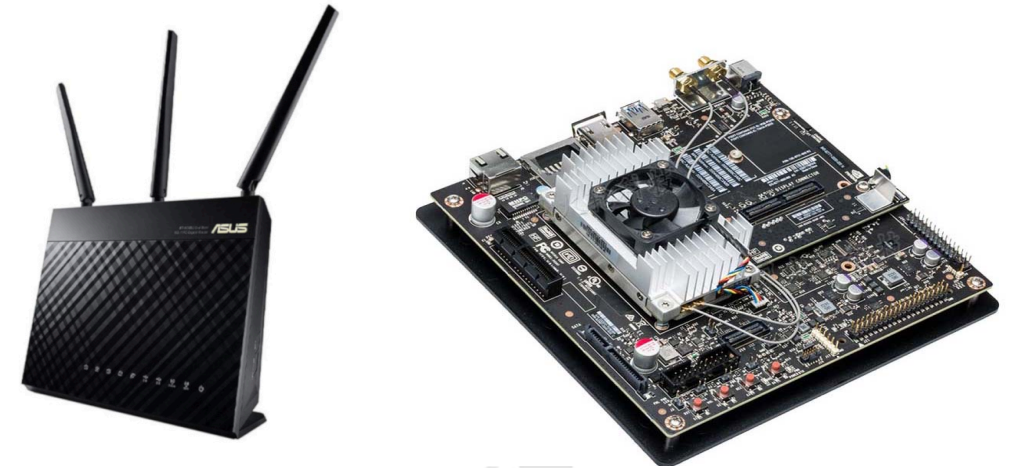
Videos and models: same as in the measurement study

Networks: shape the traffic to typical 5G mobile throughputs

Performance metric: Spherical mAP; Mean end-to-end (E2E) latency

Baselines: ERP; CubeMap

Latency budgets: from 500 ms to 4,500 ms based on baselines



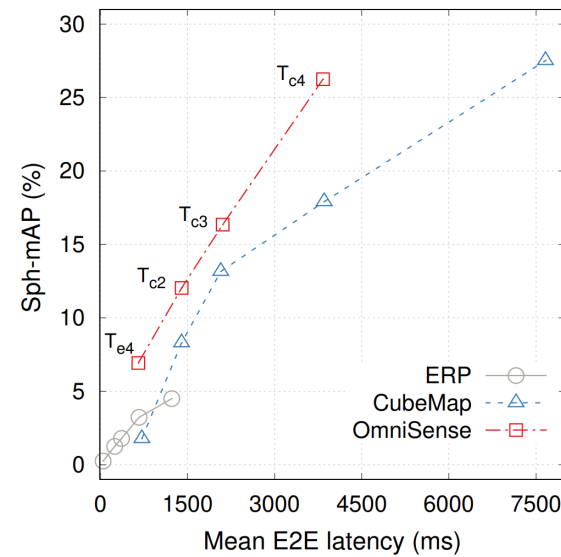
EVALUATION

□ Evaluation Results

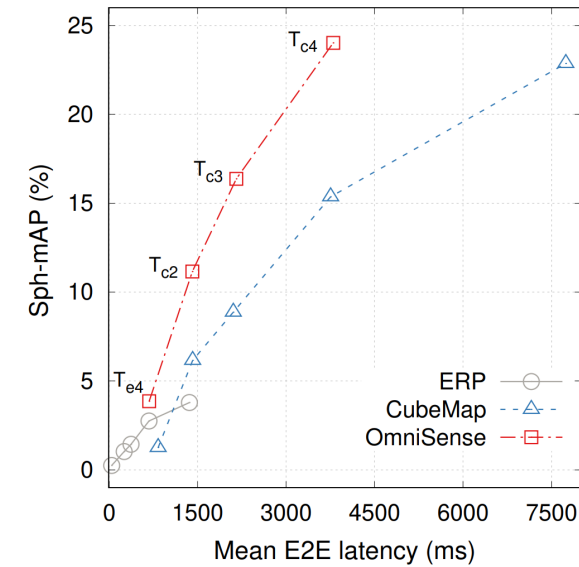
Overall performance comparisons of various methods on different videos

19.8 % - 114.6 % relative accuracy improvement

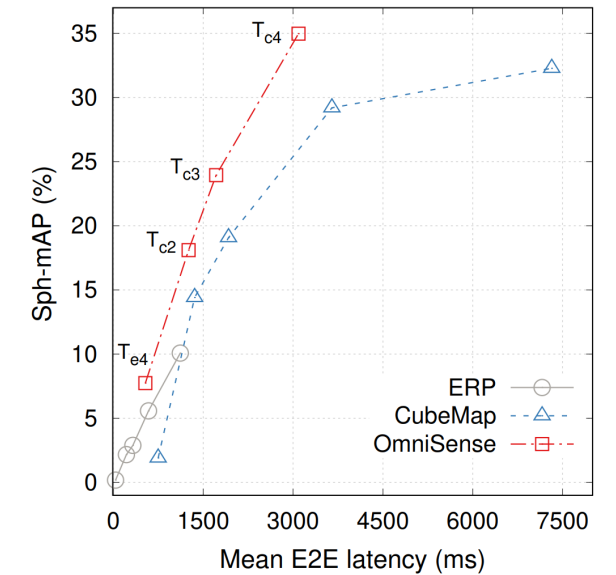
2.0 x – 2.4 x speedups



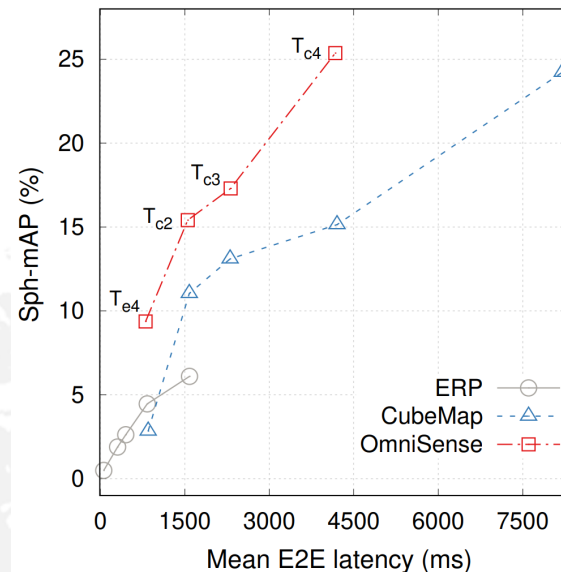
(a) Video: New-Orleans-drive



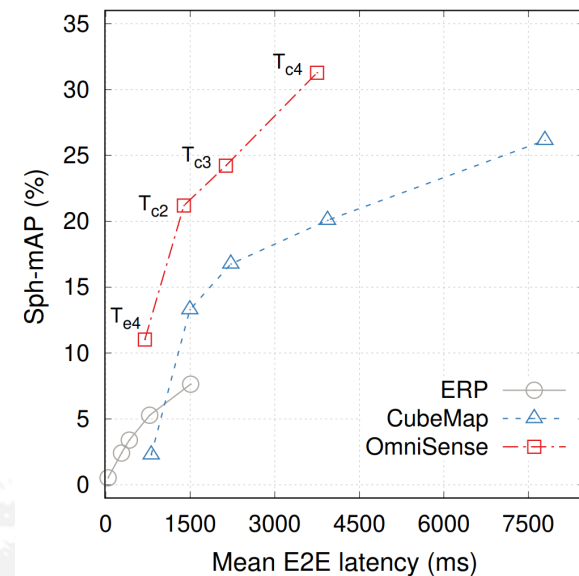
(b) Video: Chicago-drive



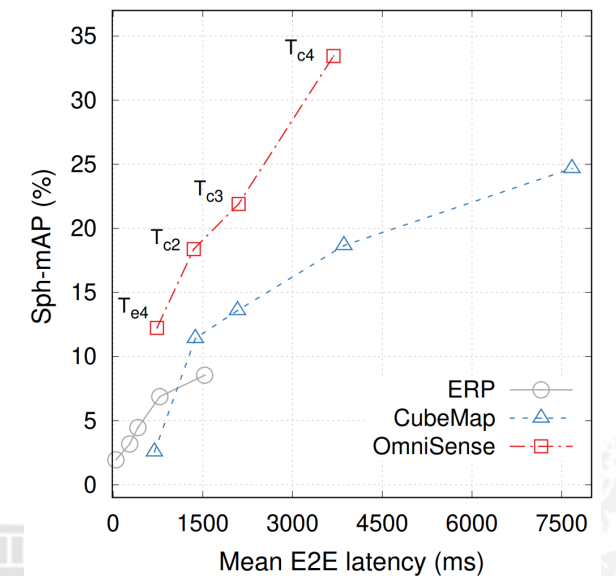
(c) Video: Expressway-drive



(d) Video: Sunny-walk1



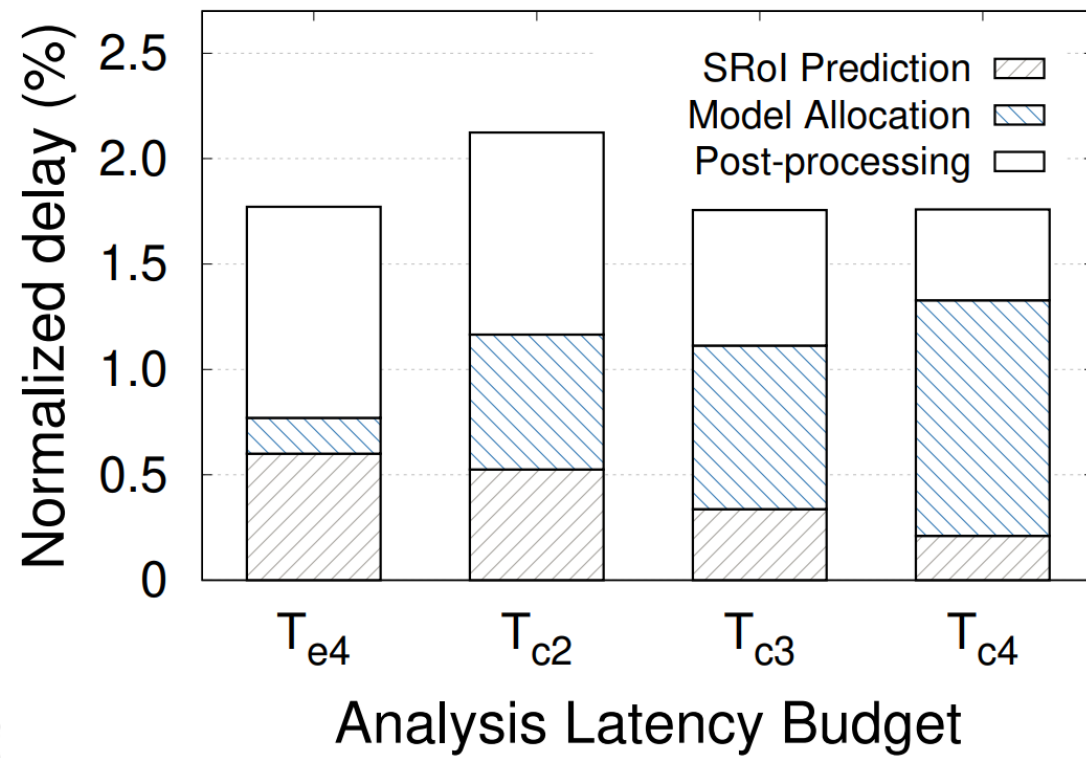
(e) Video: Sunny-walk2



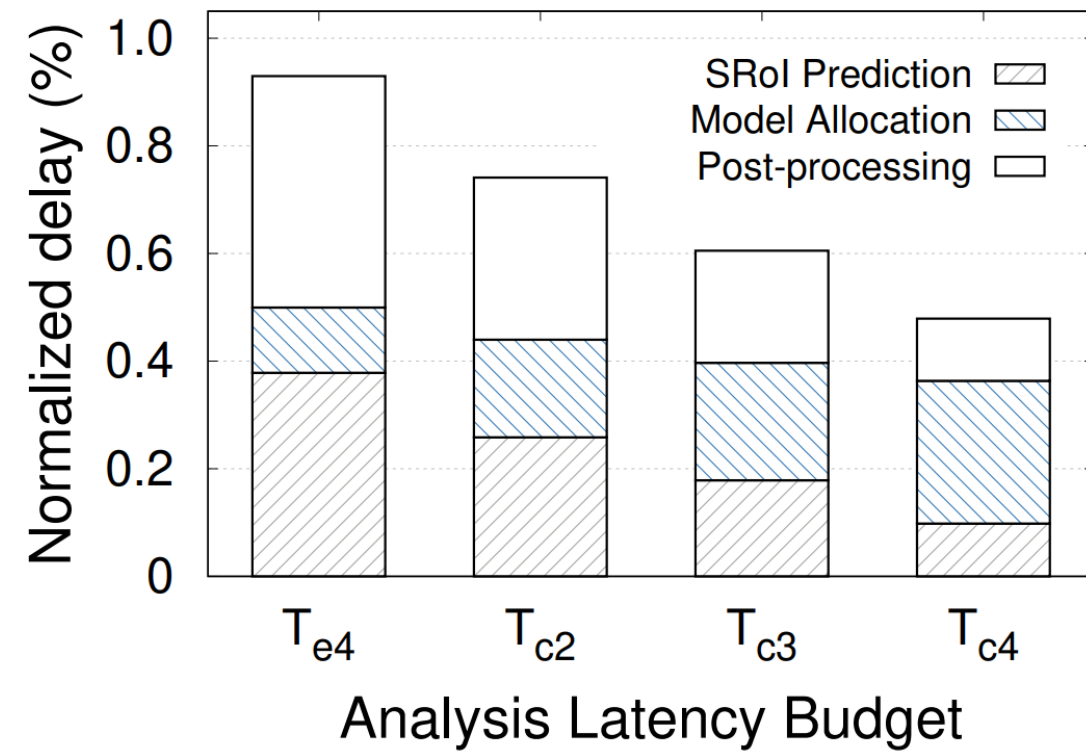
(f) Video: Cloudy-walk

EVALUATION

□ Evaluation Results



(a) video: Chicago-drive

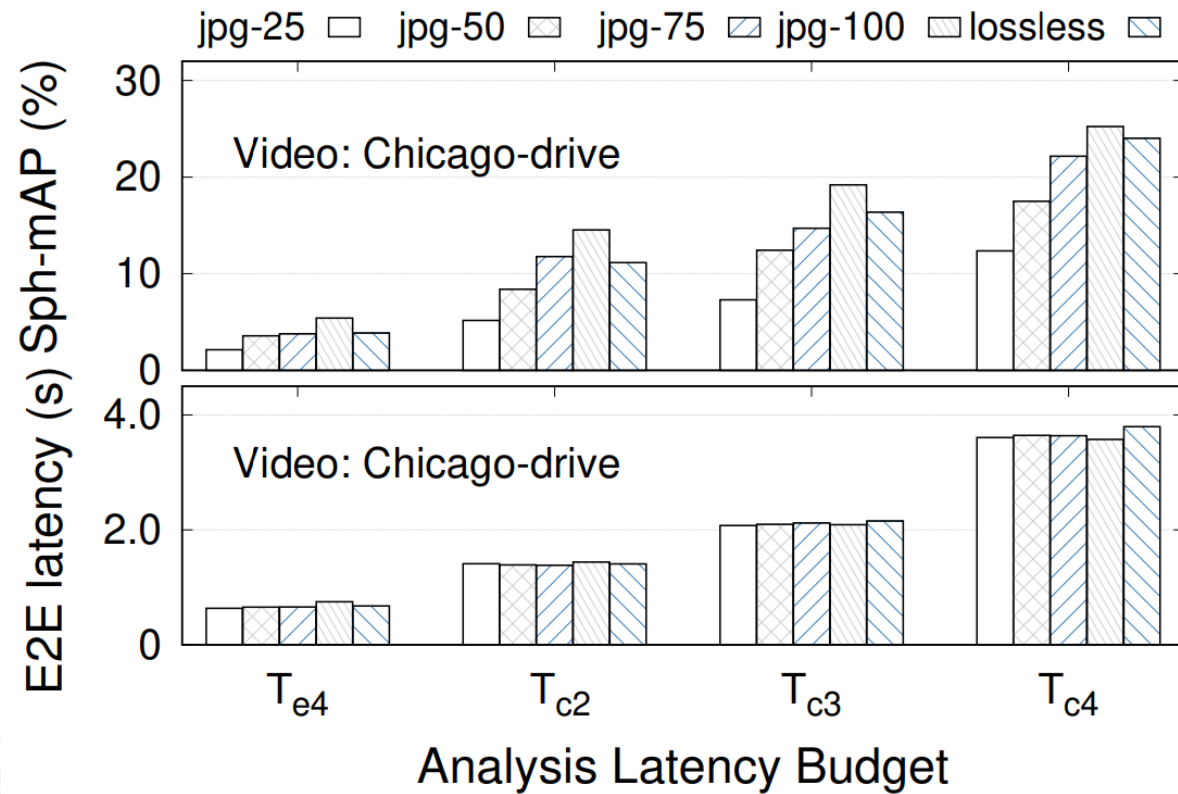


(b) video: Sunny-walk2

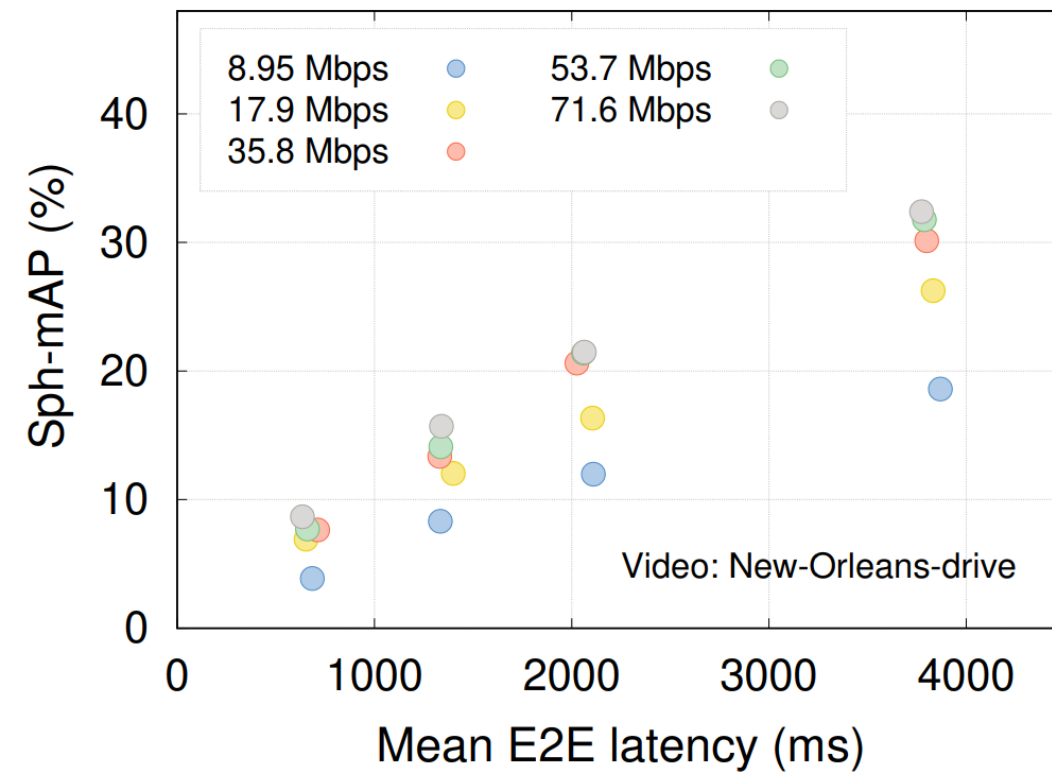
System overhead on the mobile device. The delays are normalized by the corresponding mean E2E latencies.

EVALUATION

□ Evaluation Results



(a) Perspective image compression quality



(b) Network bandwidth

Sensitivity to the compression quality of perspective images and network throughput

SUMMARY

- Immersive video analytics will be essential in unlocking the full potential of 360-degree videos.
- Our analysis of 360-degree content characteristics reveals new resource-saving opportunities in online analytics.
- OmniSense achieves low-latency and high-accuracy immersive video analytics by fine-grained content-aware resource adaptation.

THANK YOU

Q & A

