# Optimal Strategies for Live Video Streaming in the Low-latency Regime

Liyang Sun, Tongyu Zong, Yong Liu and Yao Wang, New York University

Haihong Zhu, Futurewei Technologies

# Video-over-Ideal-5G

24K, 360 Degree, Volumetric

AR/VR/MR

Live/Interactive

| Application |
| Transport |
| Network |
| Data Link |
| Physical |

Higher Throughput: Gbps

Lower Latency, 1ms

1

# Video-over-5G: real challenges

| Application Layer | User QoE Optimization with Realistic Network Assumptions |
|---|---|
| Transport | |
| Network | |
| Data Link | |
| Physical | |

❖ Users sensitive to video quality and temporal variation

❖ Video freezes/skips/black-screen detrimental to user QoE

❖ Long end-to-end video delay kills interactivity

❖ Users want mobile/wireless video

??? Consistently High-throughput/Low-delay from Lower Layers ???
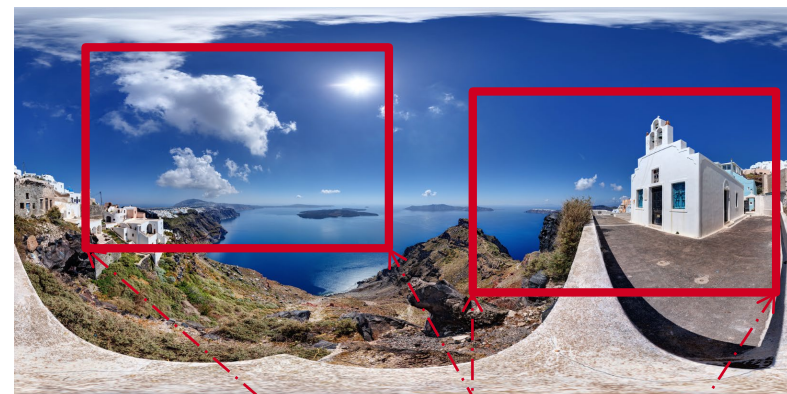
# 360-degree Video Streaming Projects (joint with Yao Wang)

1. **Two-tier on-demand 360° video streaming**
   - Field-of-View (FoV) streaming to reduce b.w. requirement (1/6)
   - two-tier segment coding/streaming to be robust against b.w. and FoV dynamics



Fraction of area within FoV: 120°/360° x 90°/180°=1/6
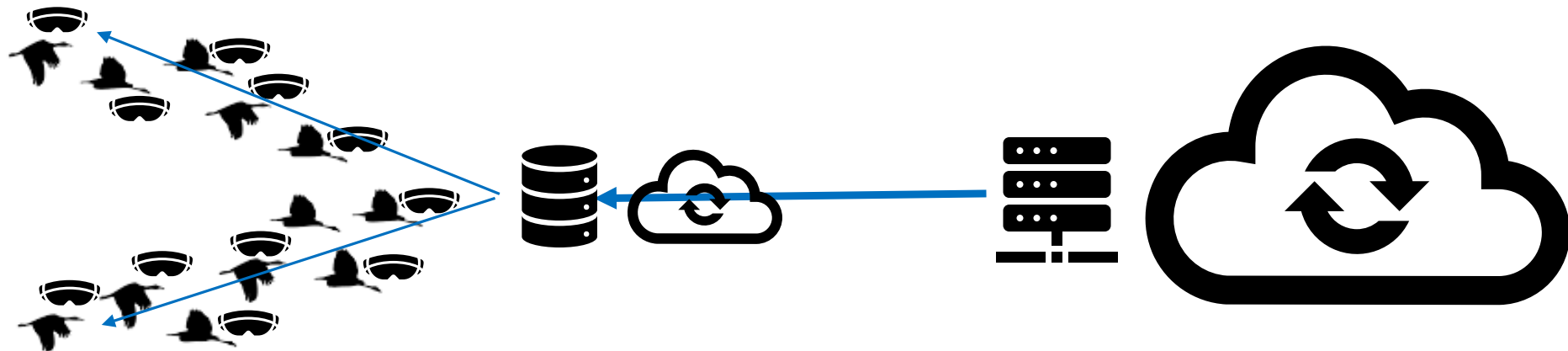
# 360-degree Video Streaming Projects (joint with Yao Wang)

1. **Two-tier on-demand 360° video streaming**
   - Field-of-View (FoV) streaming to reduce b.w. requirement (1/6)
   - two-tier segment coding/streaming to be robust against b.w. and FoV dynamics

2. **Flocking-based live 360° streaming from edge cloud**
   - users watching same live event form a "flock"
   - users with shorter video lag lead flock: populating cache, generating realtime "saliency" map

# 360-degree Video Streaming Projects (joint with Yao Wang)

1. ## Two-tier on-demand 360° video streaming
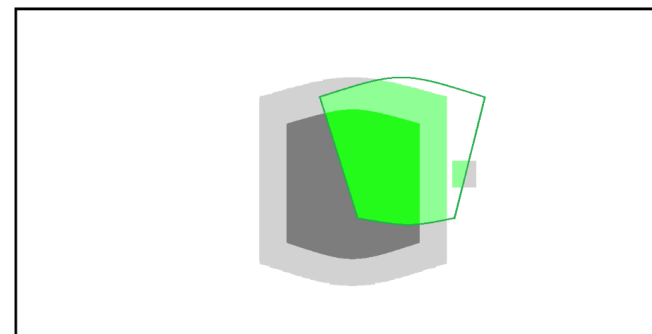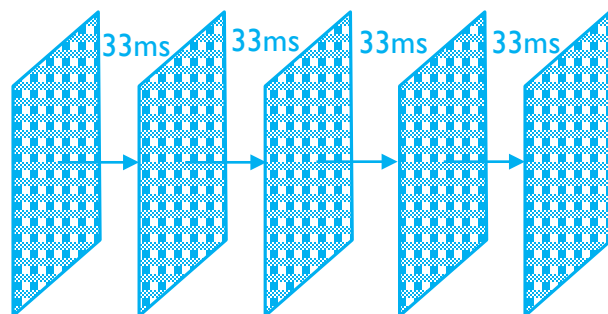   - Field-of-View (FoV) streaming to reduce b.w. requirement (1/6)
   - two-tier segment coding/streaming to be robust against b.w. and FoV dynamics

2. ## Flocking-based live 360° streaming from edge cloud
   - users watching same live event form a "flock"
   - users with shorter video lag lead flock: populating cache, generating realtime "saliency" map

3. ## Realtime coding and delivery for interactive streaming
   - frame-level coding and delivery to ensure tens milli-second latency
   - tile-based video rate allocation for FoV quality differentiation

33ms 33ms 33ms 33ms

5

# 360-degree Video Streaming Projects (joint with Yao Wang)

1. ## Two-tier on-demand 360° video streaming
   - Field-of-View (FoV) streaming to reduce b.w. requirement (1/6)
   - two-tier segment coding/streaming to be robust against b.w. and FoV dynamics
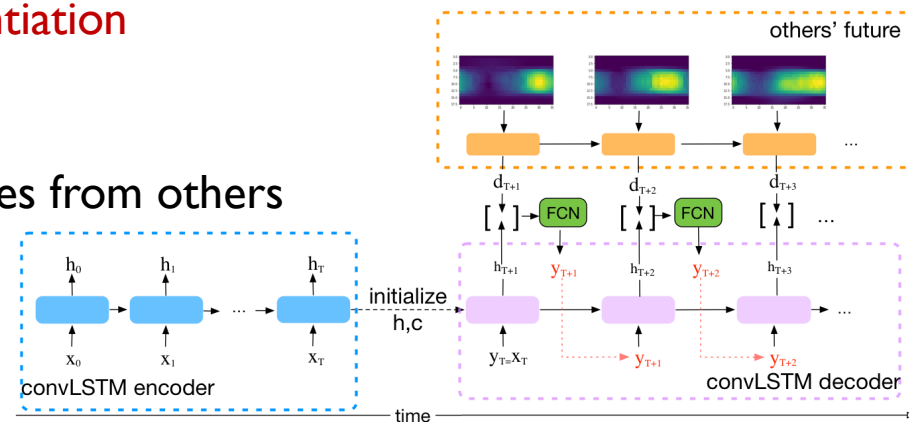
2. ## Flocking-based live 360° streaming from edge cloud
   - users watching same live event form a "flock"
   - users with shorter video lag lead flock: populating cache, generating realtime "saliency" map

3. ## Realtime coding and delivery for interactive streaming
   - frame-level coding and delivery to ensure tens milli-second latency
   - tile-based video rate allocation for FoV quality differentiation

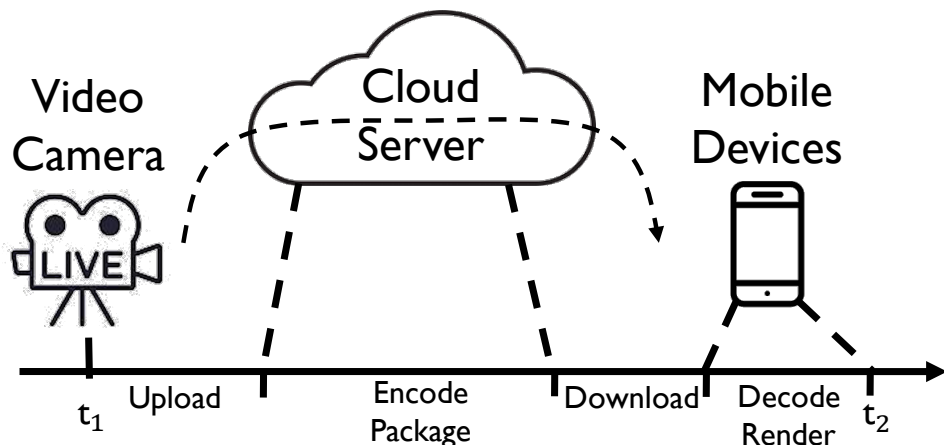4. ## Deep-learning based user FoV prediction
   - target user past FoV trajectory, and "future" trajectories from others
   - video content saliency map

6

# Low Latency Live Streaming



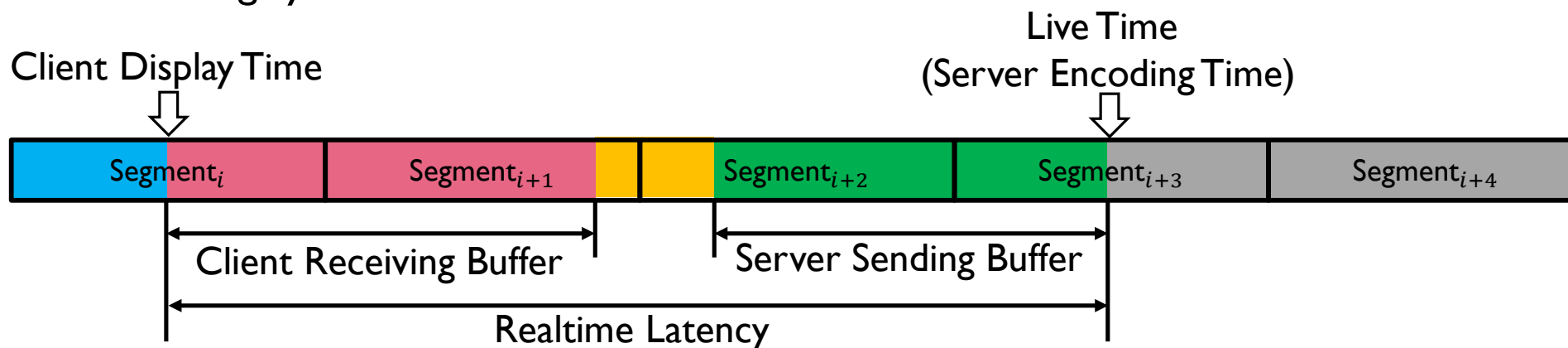| Application | Latency (s) |
|---|---|
| YouTube Live* | 7-11 |
| Facebook Live* | ~15 |
| Twitch* | ~15 |
| FOX, abc** | ~7 |

\*: Online Live Streaming
\*\*: TV broadcasting

- Sports, online gaming broadcast, social live UGC,

- Online live streaming still lags behind TV.

- User live/interactivity experience is ruined by long latency!

- Can we simply shorten latency in live streaming system? No!

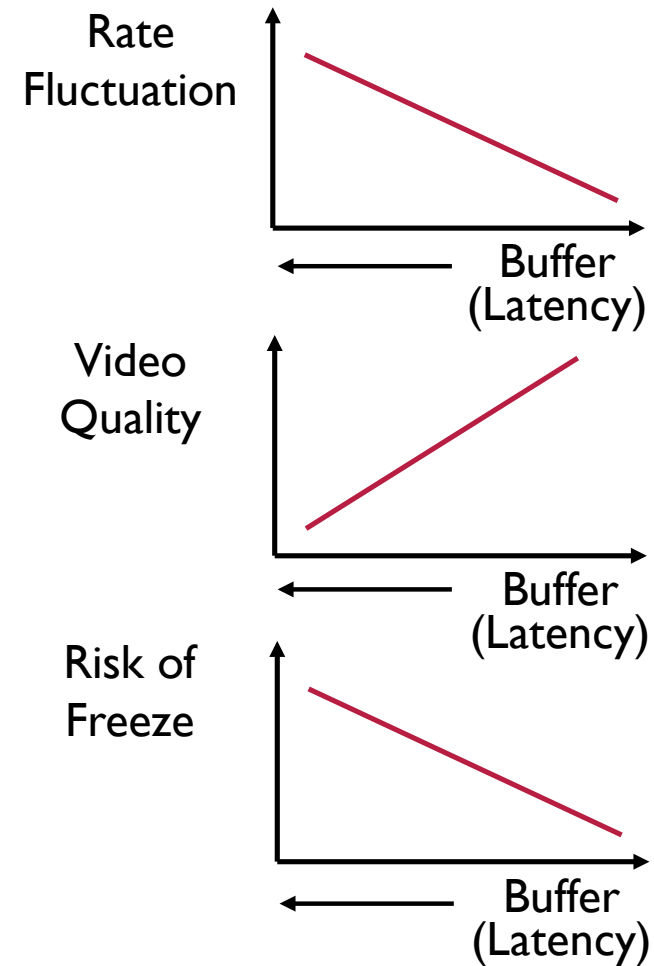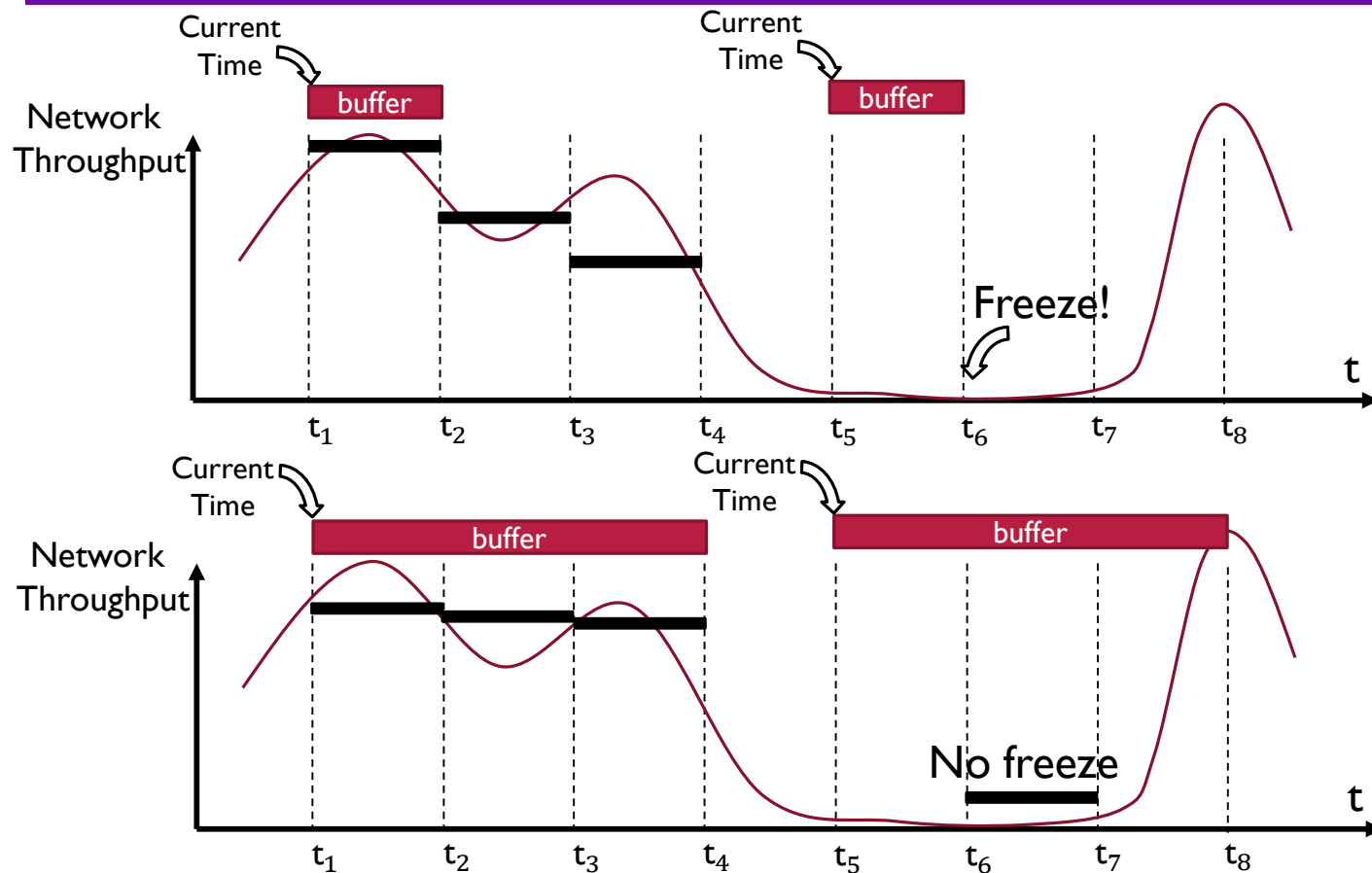# Low Latency and Buffer Length

- Live streaming system state at time T



- Realtime latency is the upper bound of client buffer length.
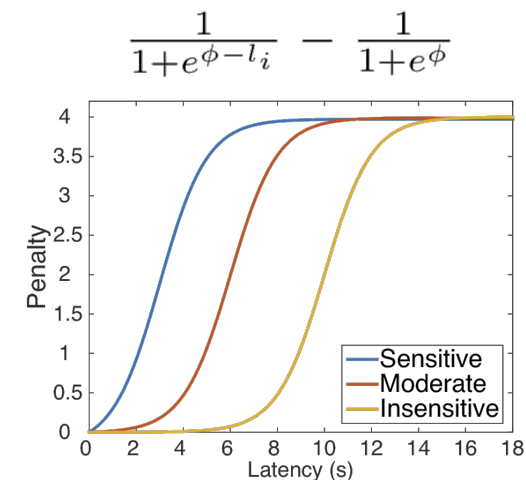
8

# Influence of Buffer Length (Latency)



- Shortening latency negatively impact all the other QoE metrics.

- Goal: Trade-off between latency and other metrics to maximize QoE.

9

# Live Streaming QoE

- QoE Metrics:

$$QoE = \underbrace{a_1 Q(r_i)}_{\text{①}} - \underbrace{a_2 |Q(r_i) - Q(r_{i-1})|}_{\text{②}} - \underbrace{a_3 x_i}_{\text{③}} - \underbrace{a_4 n_i}_{\text{④}} - \underbrace{a_5 \boldsymbol{g}(l_i)}_{\text{⑤}}$$

Video Rate$^+$    Rate Fluctuation$^-$    Freeze$^-$    Skip$^-$    Latency$^-$



$$\frac{1}{1+e^{\phi - l_i}} - \frac{1}{1+e^{\phi}}$$

10

# Model of Live Streaming System

- System Evolution

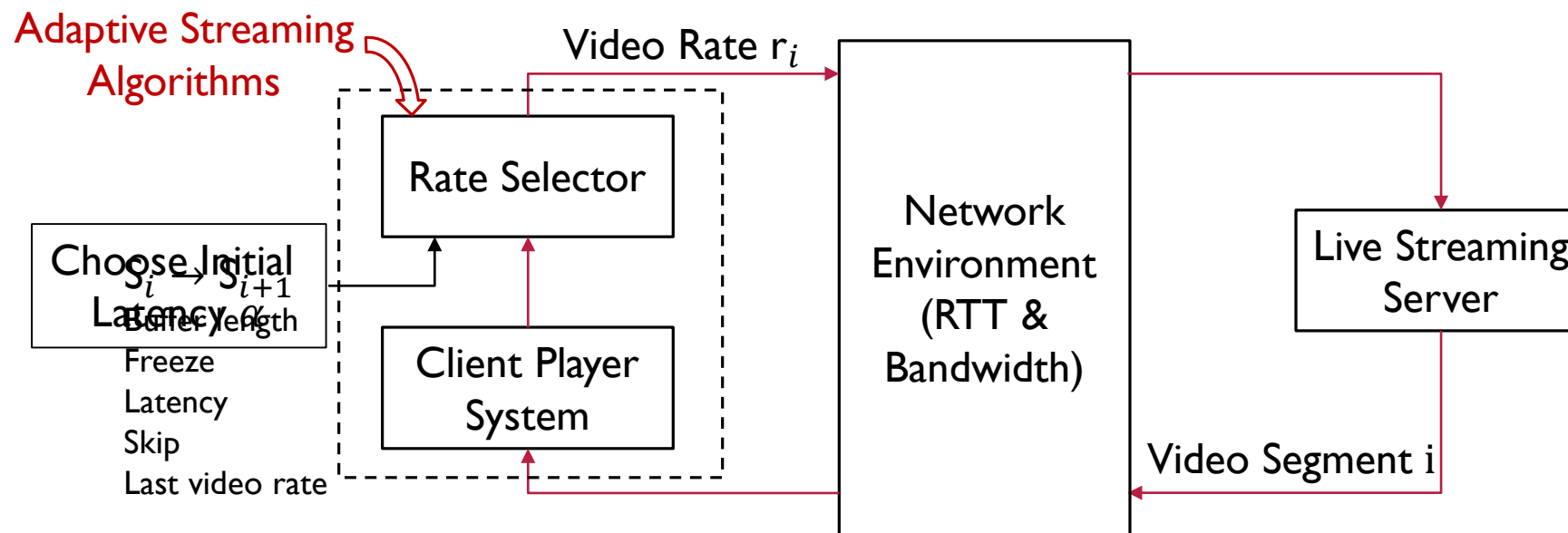Choose Initial Latency → Choose Rate → Download → Update System State → Choose Rate …

# Model of Live Streaming System

- System Evolution

Choose Initial Latency → Choose Rate → Download → Update System State → Choose Rate …

Adaptive Streaming Algorithms

Video Rate $r_i$

Rate Selector

Choose Initial Latency $\alpha$

$S_i \rightarrow S_{i+1}$
Buffer length
Freeze
Latency
Skip
Last video rate

Client Player System

Network Environment (RTT & Bandwidth)

Live Streaming Server

Video Segment i

# Optimal Streaming with Network Oracle

- Network condition for future m steps is available.

Algorithm Horizon-m

Initial



Optimal Solution: $\{ R_3 \ R_1 \ R_3 \dots \dots \}$

$m$

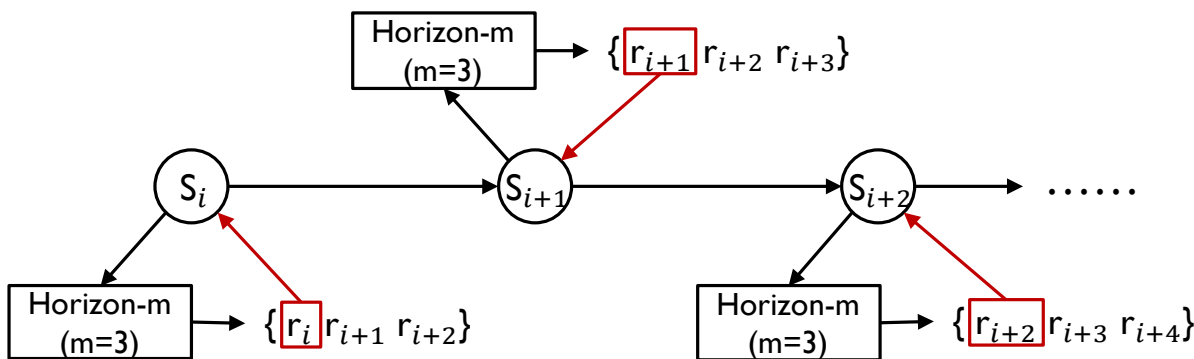- Large m causes state explosion.

**Algorithm 1** Optimal Streaming for Horizon-$m$

**Input:** $\mathcal{S}_1$: the initial state; $m$: look-ahead horizon; $\{w_i, rtt_i, i \in [1, m]\}$: future available bandwidth and rtt; $\mathcal{R}$: available rates;
**Output:** $\{r_i^*, i \in [1, m]\}$: optimal rate sequence.
**Initialization:** The possible states at stage 1: $\Omega_1 = \{\mathcal{S}_1\}$.
1: *Branch-and-Bound State Expansion*
2: **for** each segment $i \in [1, m]$ **do**
3:     $\Omega_i = \emptyset$
4:     **for** each state $\mathcal{S}$ in $\Omega_{i-1}$ **do**
5:         **for** each $R_j \in \mathcal{R}$ **do**
6:             $\mathcal{S}_i' = \boldsymbol{f}(\mathcal{S}, R_j, \{w_i, rtt_i\})$
7:             **if** $\mathcal{S}_i'$ could be part of the overall optimal solution **then**
8:                 $\Omega_i \leftarrow \Omega_i \bigcup \mathcal{S}_i'$
9:             **end if**
10:         **end for**
11:     **end for**
12: **end for**

13: Find Optimal Transition $\mathcal{S}_1 \xrightarrow{r_1^*} \mathcal{S}_2^* \in \Omega_2 \cdots \xrightarrow{r_{m+1}^*} \mathcal{S}_{m+1}^* \in \Omega_{m+1}$ to maximize accumulated QoE $\sum_{i=1}^{m} QoE(\mathcal{S}_i, r_i)$ through DP.
14: **return** $r_{[1,\cdots,m]}^*$
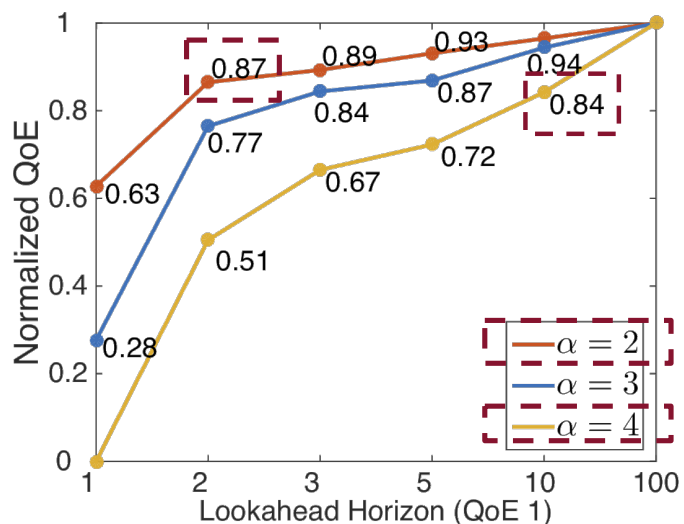
12

# Sliding Window with Horizon-(Small m)



**Algorithm 2** Sliding Horizon-$m$ Streaming

**Input:** $\mathcal{S}_1$: initial state; $\alpha$ and $\beta$: startup parameters; $m$: look-ahead horizon; $N$: live streaming duration; $\{w_i, rtt_i, i \in [1, N]\}$: available bandwidth and rtt; $\mathcal{R}$: available rates.
**Output:** $\{r_i, i \in [1, N]\}$: rate sequence for all segments
1: Download the first $\beta$ segments using predefined rate selection strategy $r_{[1,\cdots,\beta]}$, obtain $\mathcal{S}_{\beta+1}$
2: **for** each segment $i \in [\beta + 1, N]$ **do**
3: $\quad rr_i^{(m)} = \text{Horizon-m}(\mathcal{S}_i, m, \{w_{[i,i+m-1]}, rtt_{[i,i+m-1]}\}, \mathcal{R})$
4: $\quad r_i = rr_i^{(m)}[1]$
5: $\quad \mathcal{S}_{i+1} = f(\mathcal{S}_i, r_i, \{w_i, rtt_i\})$
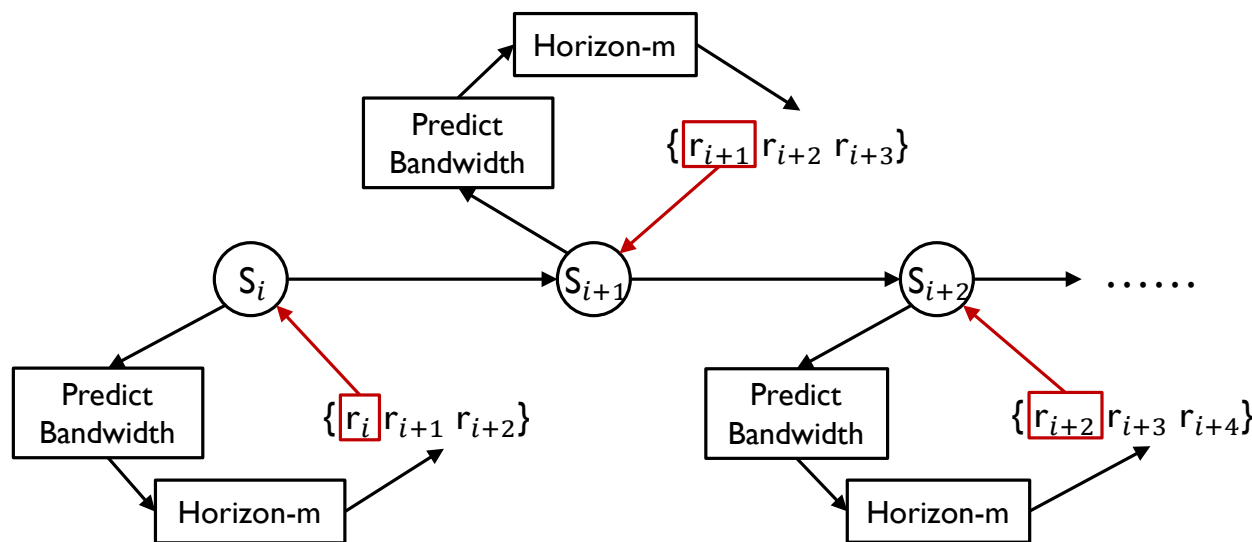6: **end for**
7: **return** $r_{[1,\cdots,N]}$



- When latency $\alpha$=2, small lookahead horizon (m=2) is needed to get high QoE.

- If $\alpha$=4, similar normalized QoE is achieved when m=10.

If latency is short, future information of short lookahead horizon is needed to achieve close-to-optimal QoE.

13

# Model Predictive Control (MPC) for Live Streaming

- Future network information is NOT available.

- Bandwidth predictions

  - Harmonic Mean, Hidden Markov Model (HMM), Recursive Least Squared (RLS) and LSTM.
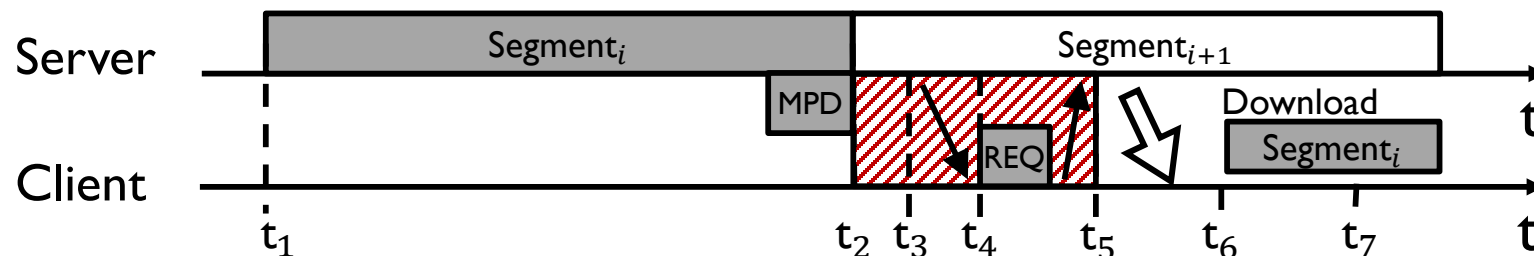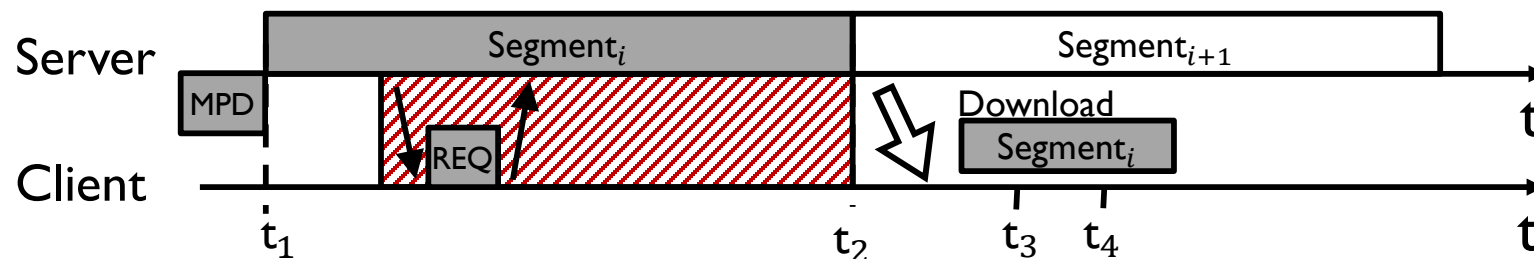


MPC based Practical Live Streaming Algorithm
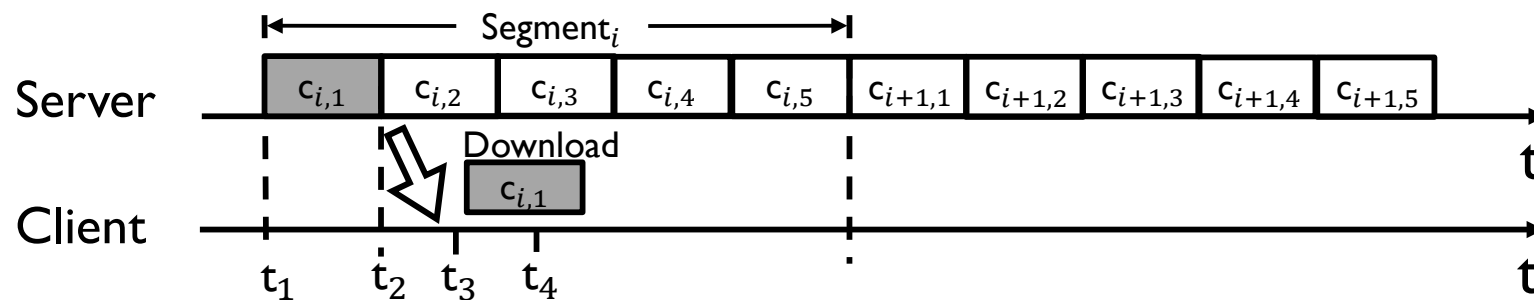
14

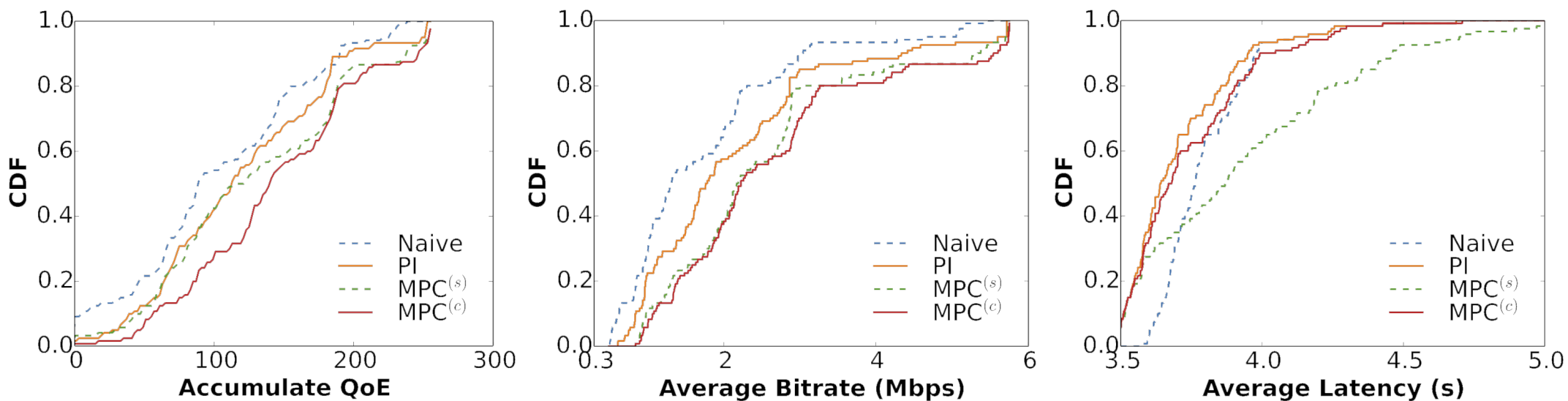# Segment & Chunk based Streaming



- Segment Mode
- Server-wait Mode
- Chunk Mode

# Trace-driven Experiments and Evaluation

- 4G cellular bandwidth dataset with 150 traces collected in NYC.

- Naïve ($\gamma \widehat{\omega}$), PI-Controller ($\gamma_p \widehat{\omega}$) and MPC (segment and chunk mode).



- MPC based algorithms outperform Naïve and PI-Controller.

- MPC$^s$ suffers more latency (caused by freeze) than MPC$^c$.

- MPC$^c$ achieves highest QoE with highest bitrate and lowest latency in most cases.

16

# Conclusions & Ongoing Work

- Low latency is crucial, balance between latency and other QoE metrics

- MPC based streaming algorithms can improve the QoE performance with low latency.

- Chunk-based delivery is helpful to support low latency live video streaming.

- Optimal Streaming Policy from Deep Reinforcement Learning (DRL) vs. Model-based RL

- Optimal Playback Pace Adaption

17

# THANK YOU!