**Lecture 11:**

# The Present and Future of Video Conferencing Systems

**Visual Computing Systems**
**Stanford CS348K, Spring 2022**
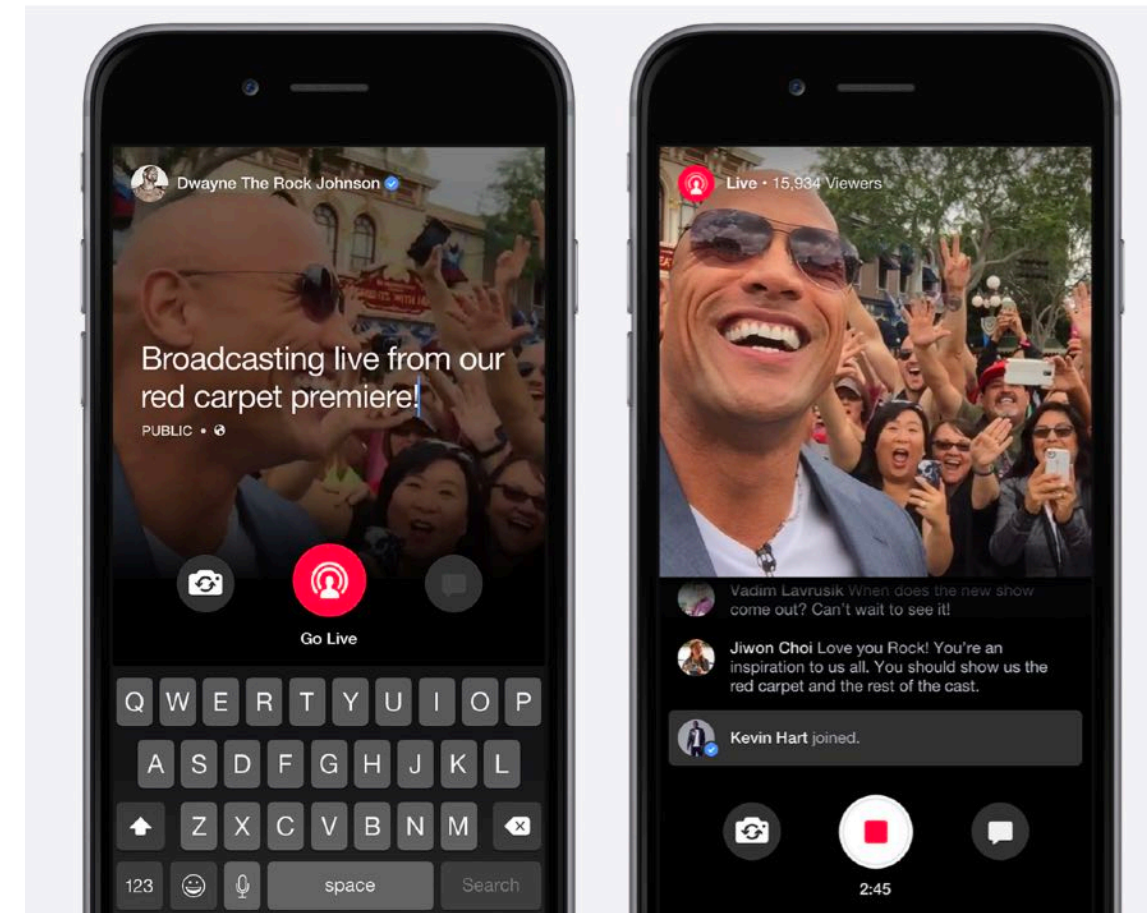
# Discussion

- **Google VCU paper (Ranganathan et al. 2021)**

# Types of video we watch on the internet
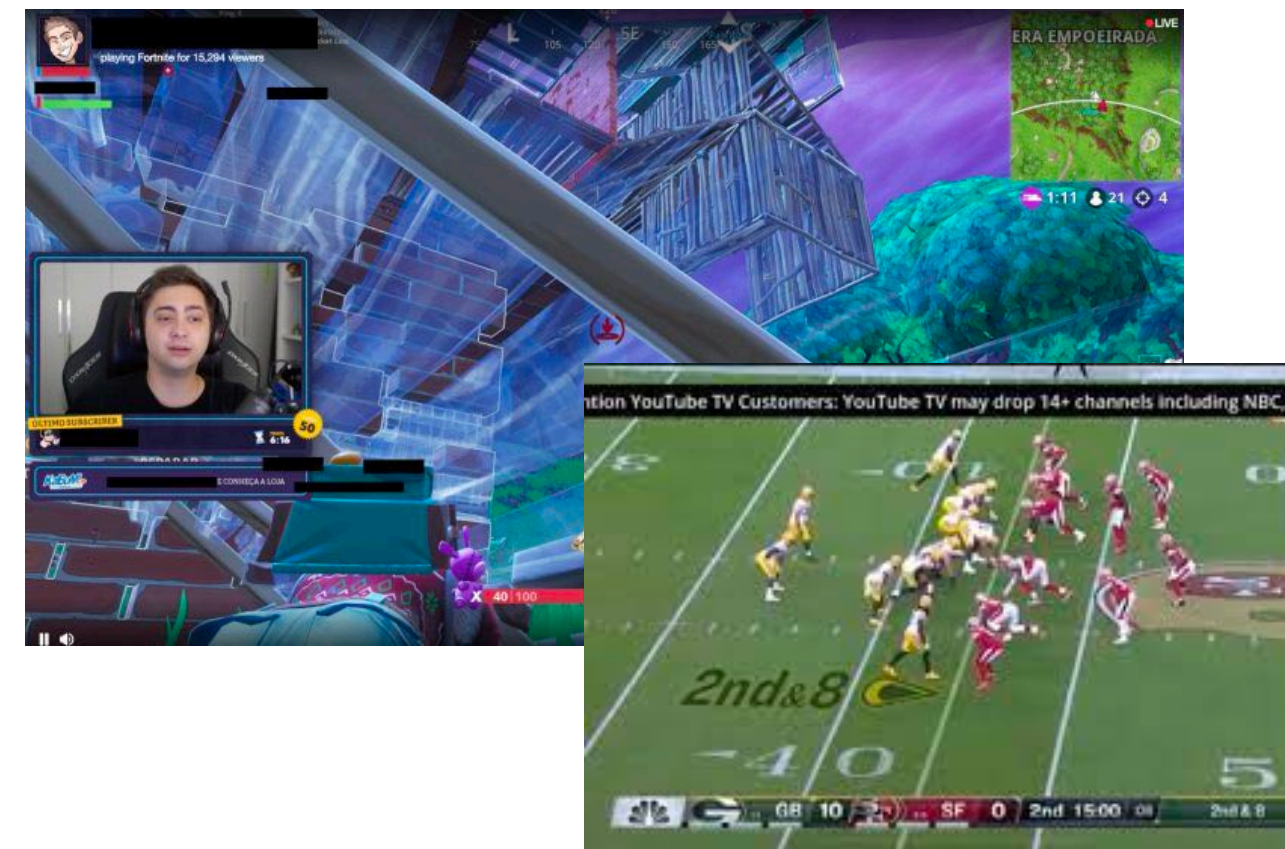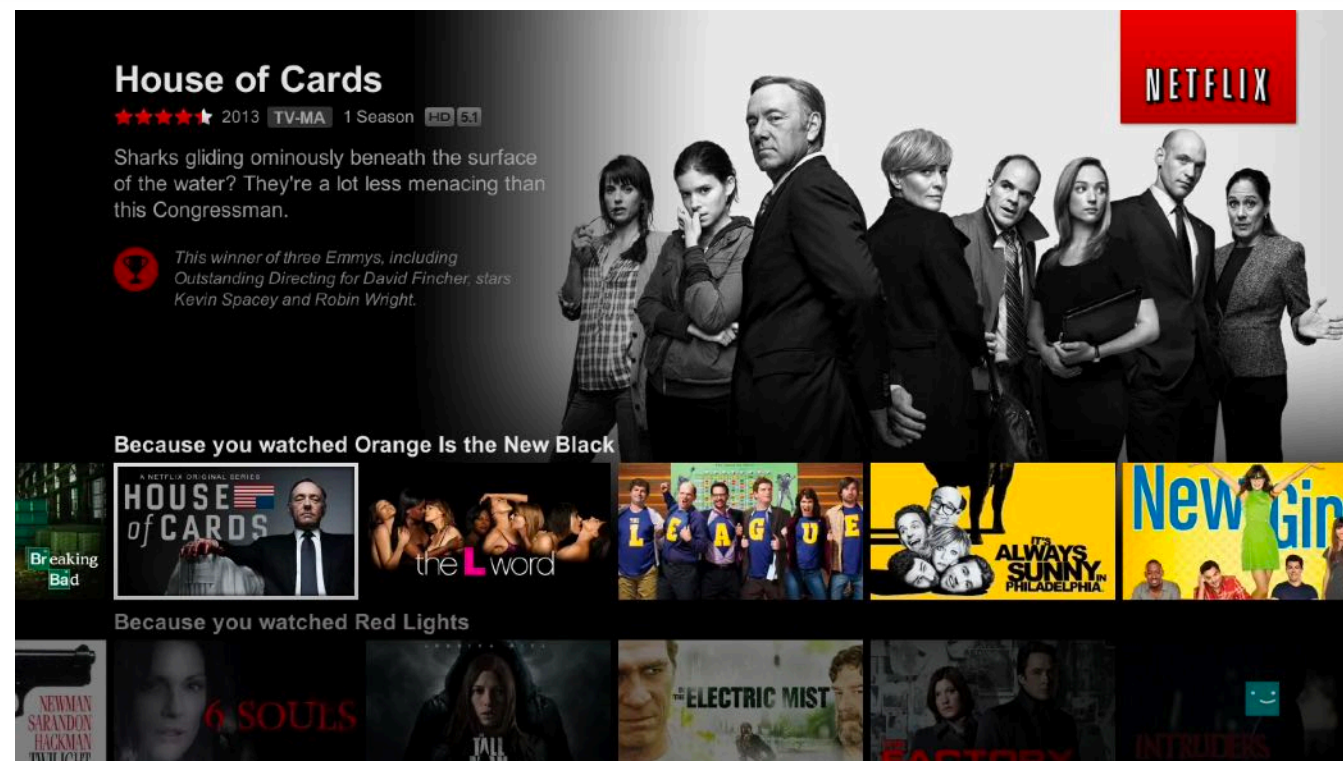
**Very different latency/bandwidth requirements…**



**Watching videos**



**Live streaming
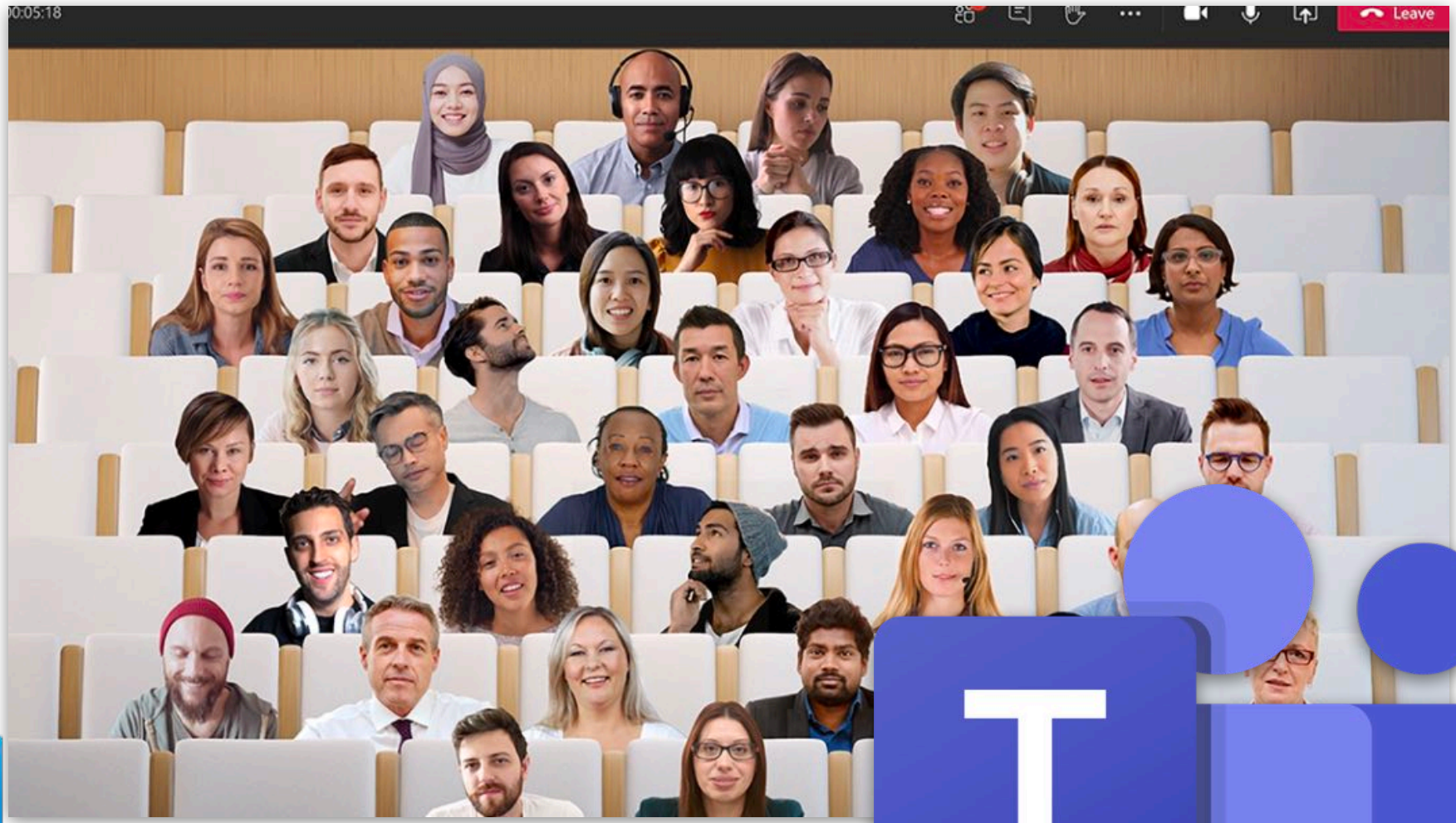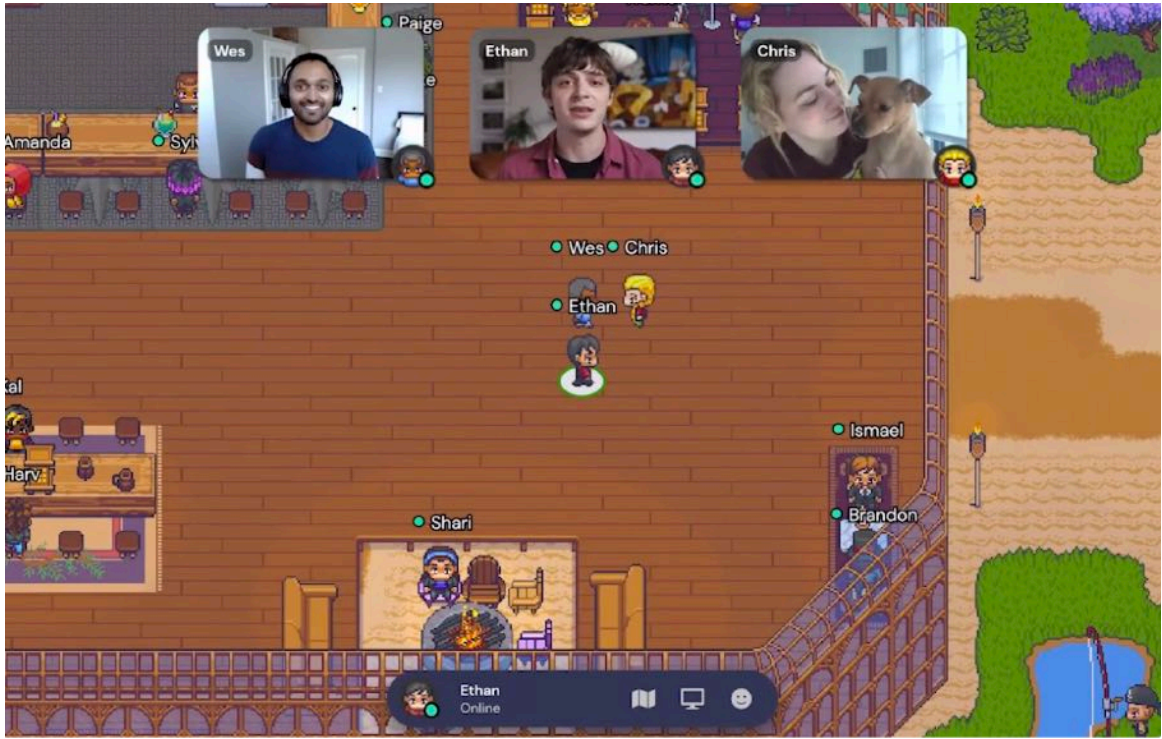(Live TV, twitch, personal live streams)**



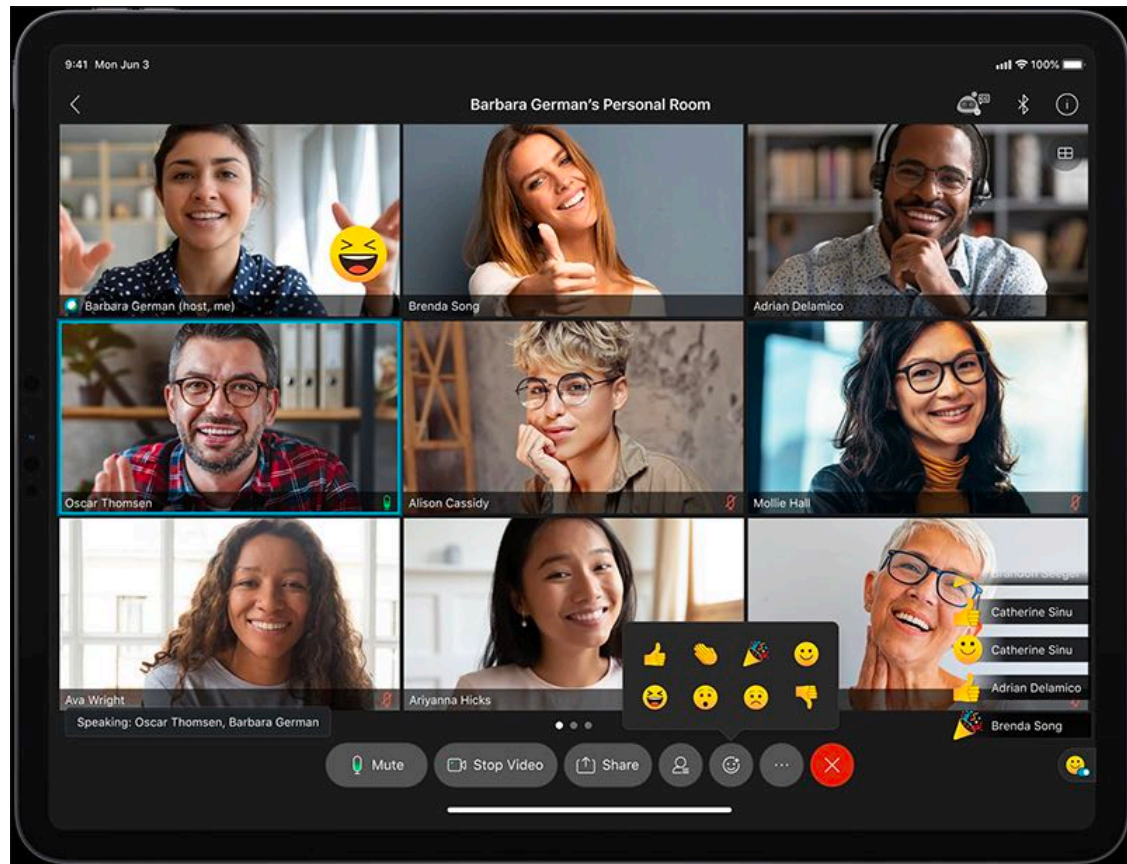**Videoconferencing,
Streaming gaming**

# Videoconferencing systems

# As you can imagine, a lot of modern interest in video conferencing(big and small!)

# Let's design a video conferencing system

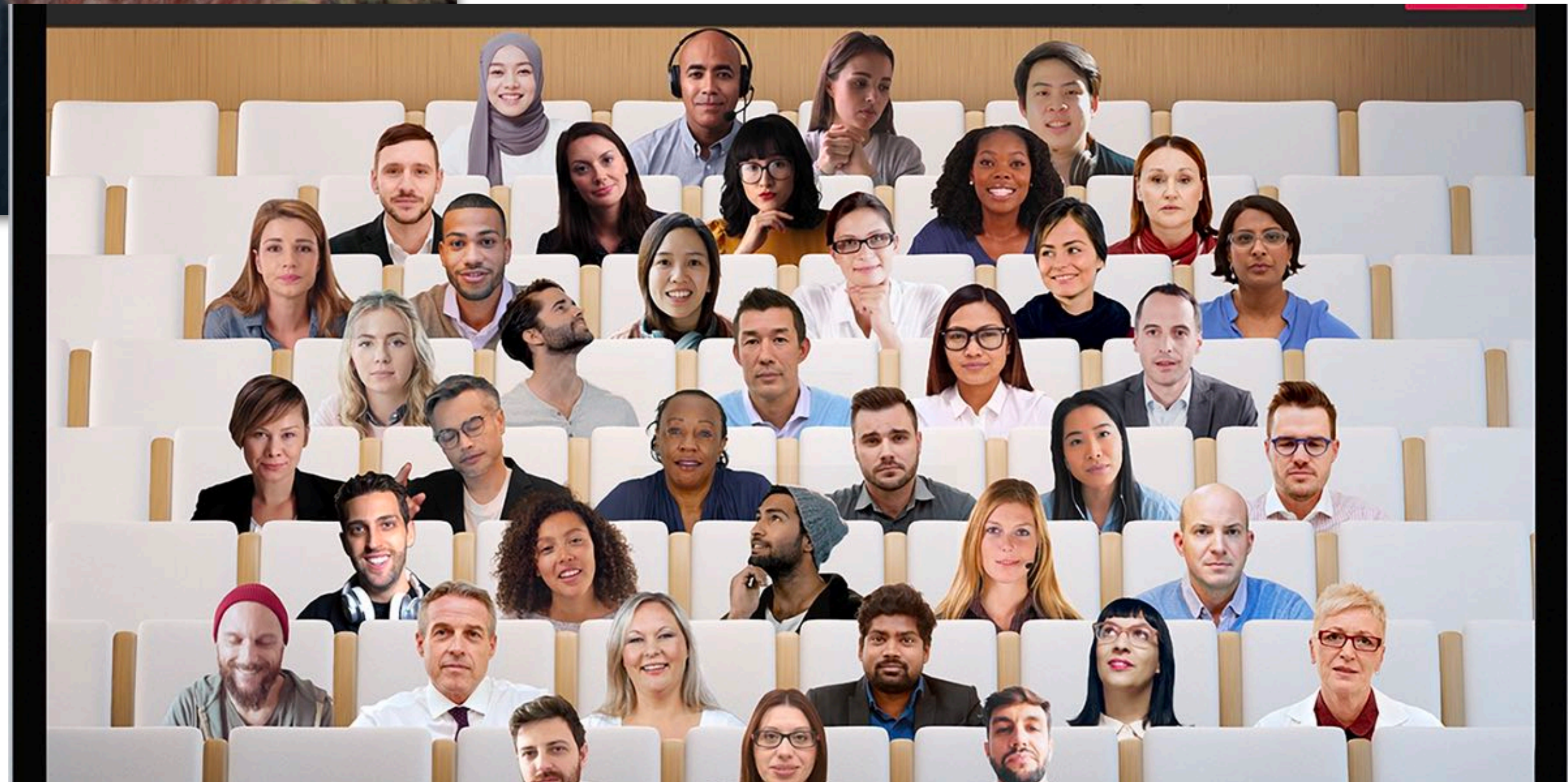## We want to deliver a visually rich experience similar to features of modern platforms

# Let's design a video conferencing system



Kayvon Fatahalian

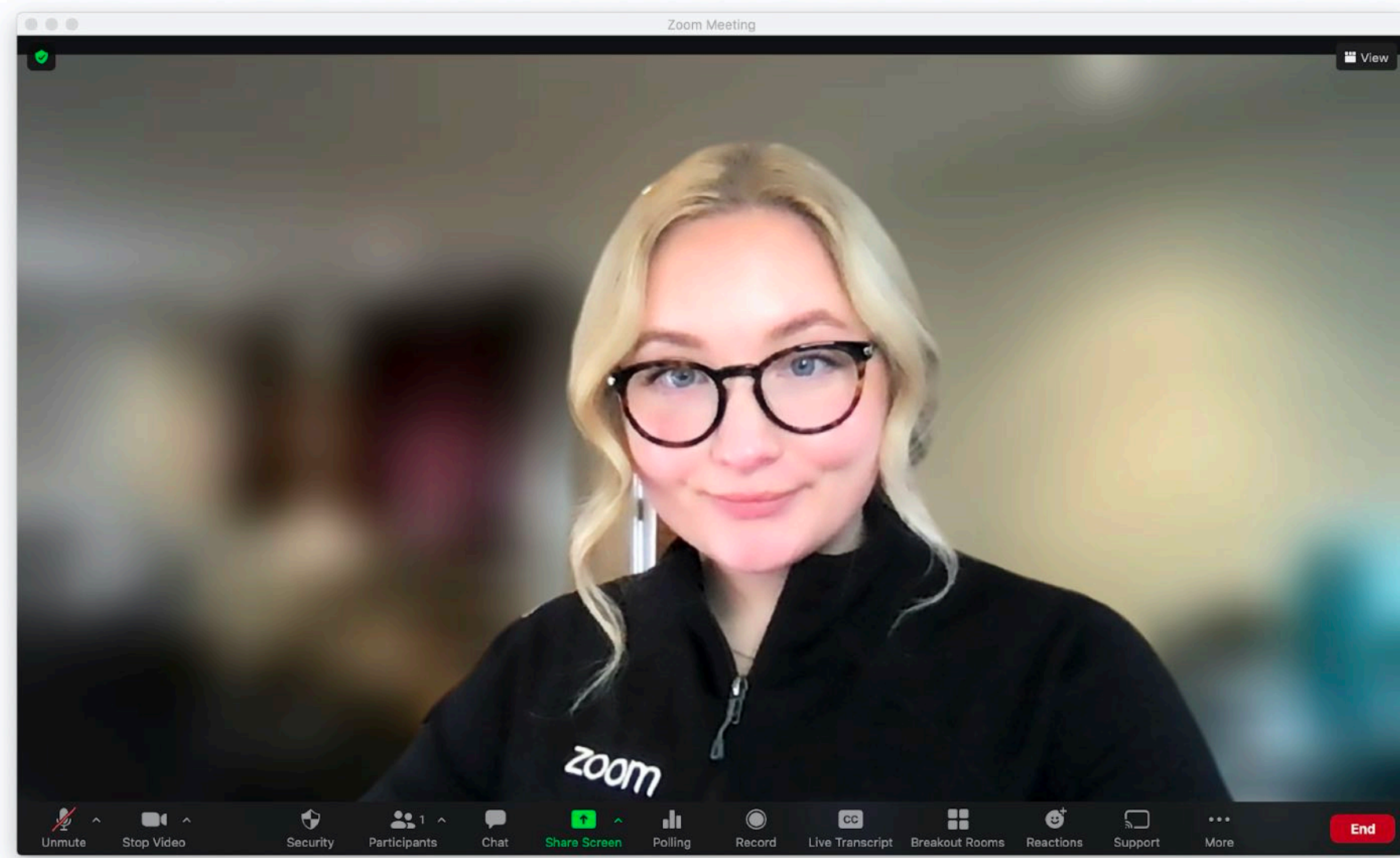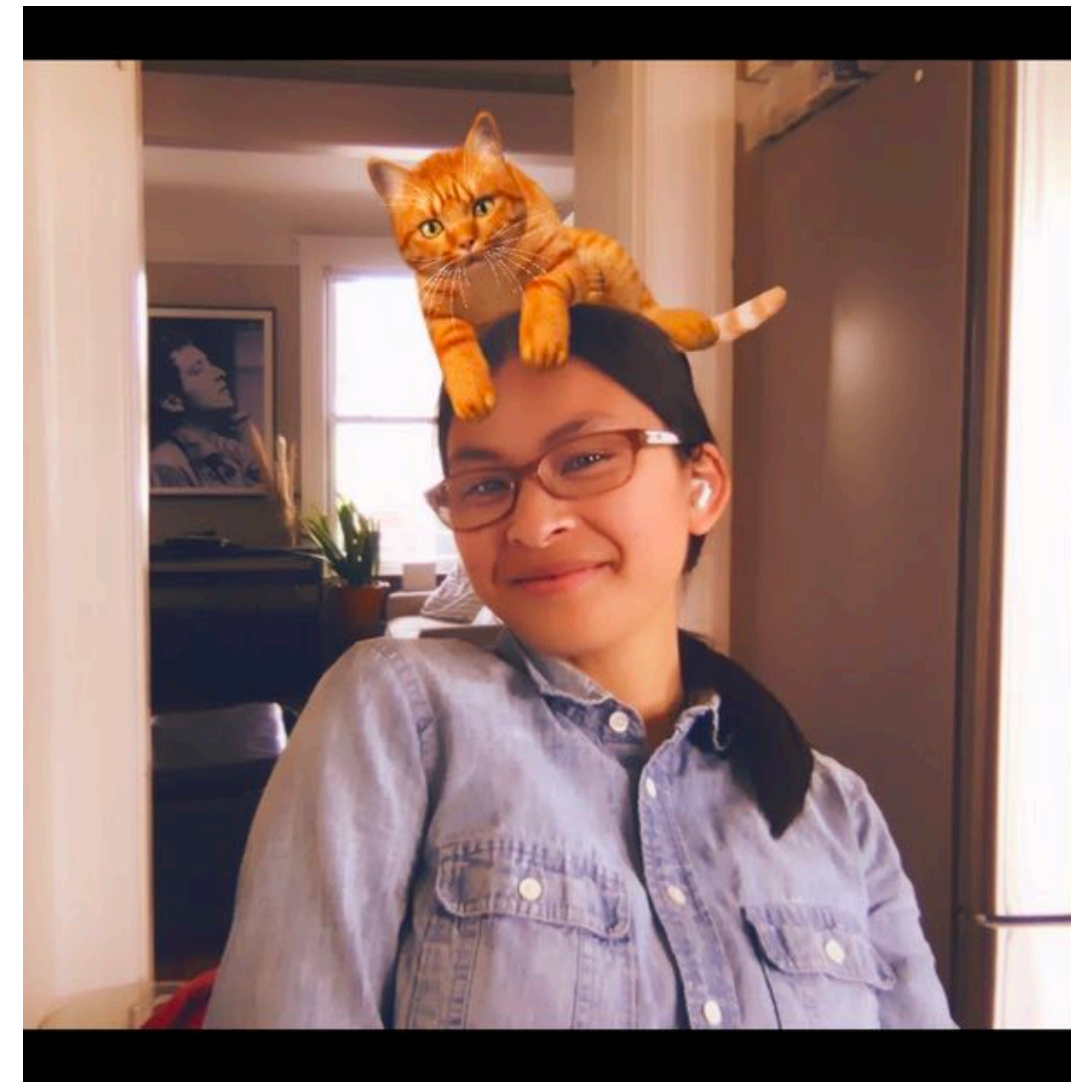**Segment participant from background**

# Let's design a video conferencing system

**Perform image processing to enhance look of video feed**



**Blur background**



**Render additional content**

**Studio Effects**

☑ Apply to all future meetings

⌄ **Eyebrows**

⌄ **Moustache & Beard**

⌄ **Lip Color**

**Adjust lighting**

# Other forms of augmentation



**Real-time translation and captioning**

# Let's design a video conferencing system

**Large gallery views: companies raced to provide 7x7 gallery in 2020**



Maximum participants displayed per screen in Gallery View:

○ 25 participants    ● 49 participants

# Deliver to wide range of clients and network settings

# Setup…

**Consider issues like latency…**



Cloud

West Coast
Servers

East Coast
Servers

Personal
computer

Personal
computer

# Q. Should we transcode/process video on our cloud servers?

- **What are advantages (to users? To us the provider)?**

- **What are disadvantages?**

# Implementing gallery view

Cloud routes compressed video bitstreams to users
(Does not manipulate bits)

. . .

Zoom calls this "multimedia routing"

Receiving client "renders" all streams into appropriate display

# One drawback of this design

- If each client is providing a single compressed video stream, that means each person on the video call must receive the same bits right? (What if they are on different network connections?)

# Scalable video codec (SVC)

- **"Scalable" compressed video bitstream: subsets of the bitstream encode valid video streams for a decoder**

  - **Implication: if packets get lost, the remaining packets form a valid H.264 bitstream, albeit at lower resolution or quality**

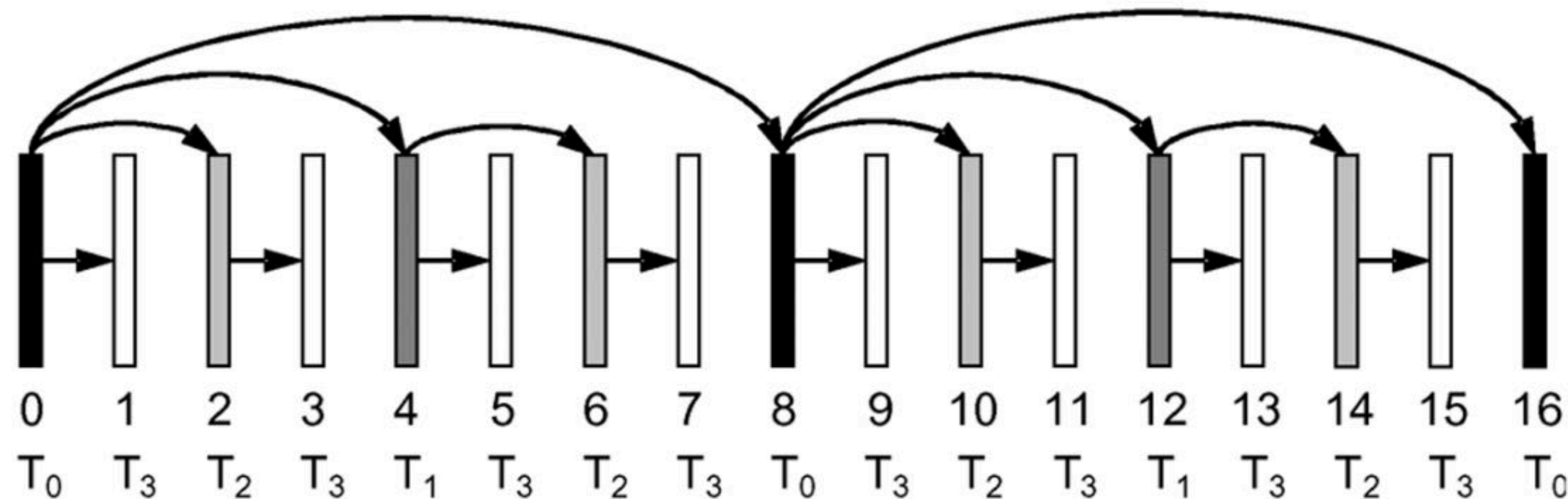**Example: temporal scalability**



| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|
| $T_0$ | $T_3$ | $T_2$ | $T_3$ | $T_1$ | $T_3$ | $T_2$ | $T_3$ | $T_0$ | $T_3$ | $T_2$ | $T_3$ | $T_1$ | $T_3$ | $T_2$ | $T_3$ | $T_0$ |

**Layer 0: ($T_0$) defines valid video at frame rate R**

**Layer 1 ($T_1$) defines bumps frame rate to 2R**

**. . .**

**Note how layer 0 information is used to predict higher layer information**

# Scalable video codec (SVC)

**Example: spatial scalability**
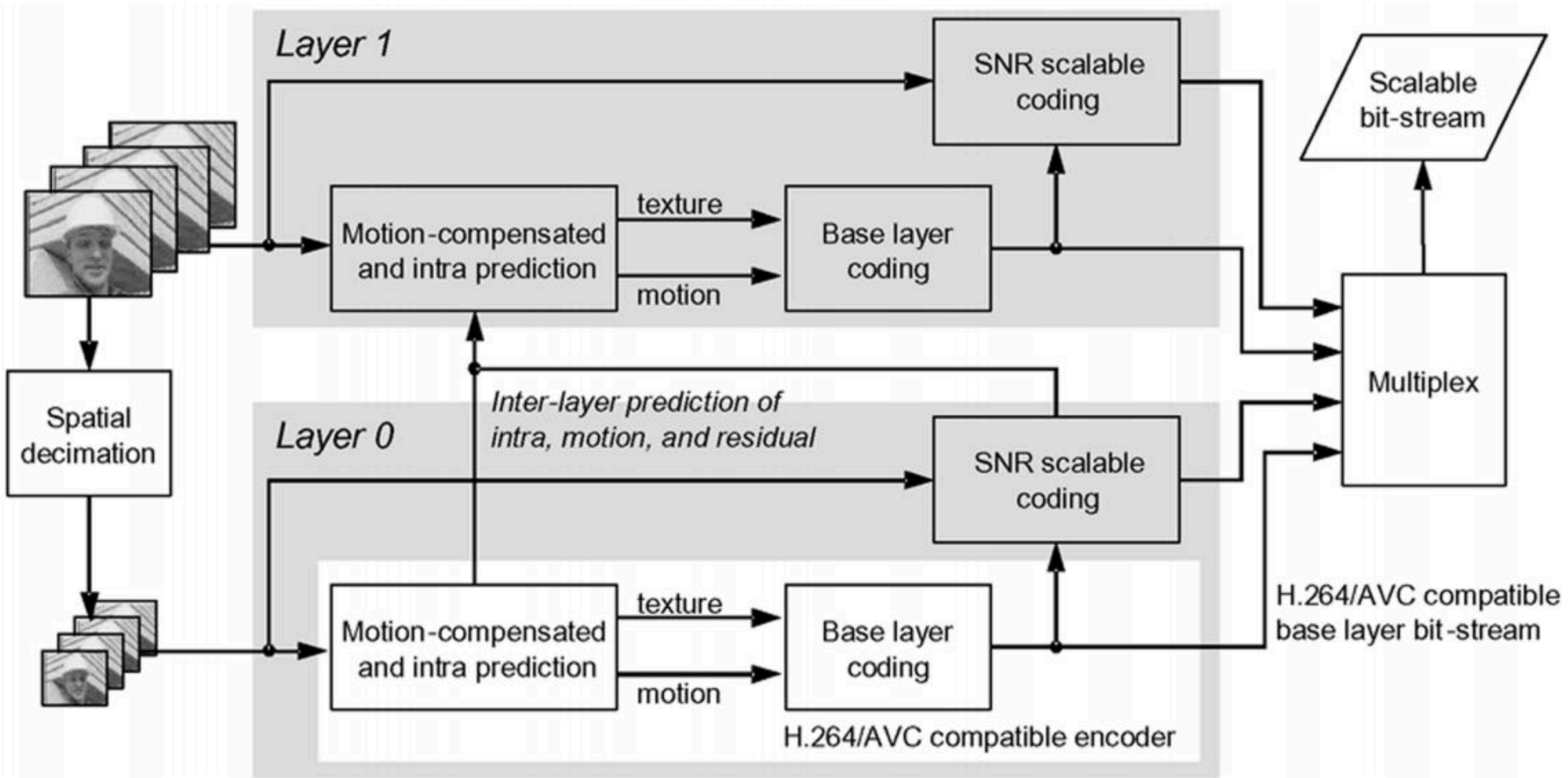


**Layer 1:
(Higher res)**

**Layer 0:
(Low res)**

Again, note how layer 0 information is used to predict higher layer information
(Higher efficiency than independently encoding two video streams)

Layer 0: defines valid video at low resolution (and low frame rate)
Layer 1: provides additional information for higher resolution (and higher frame rate) video
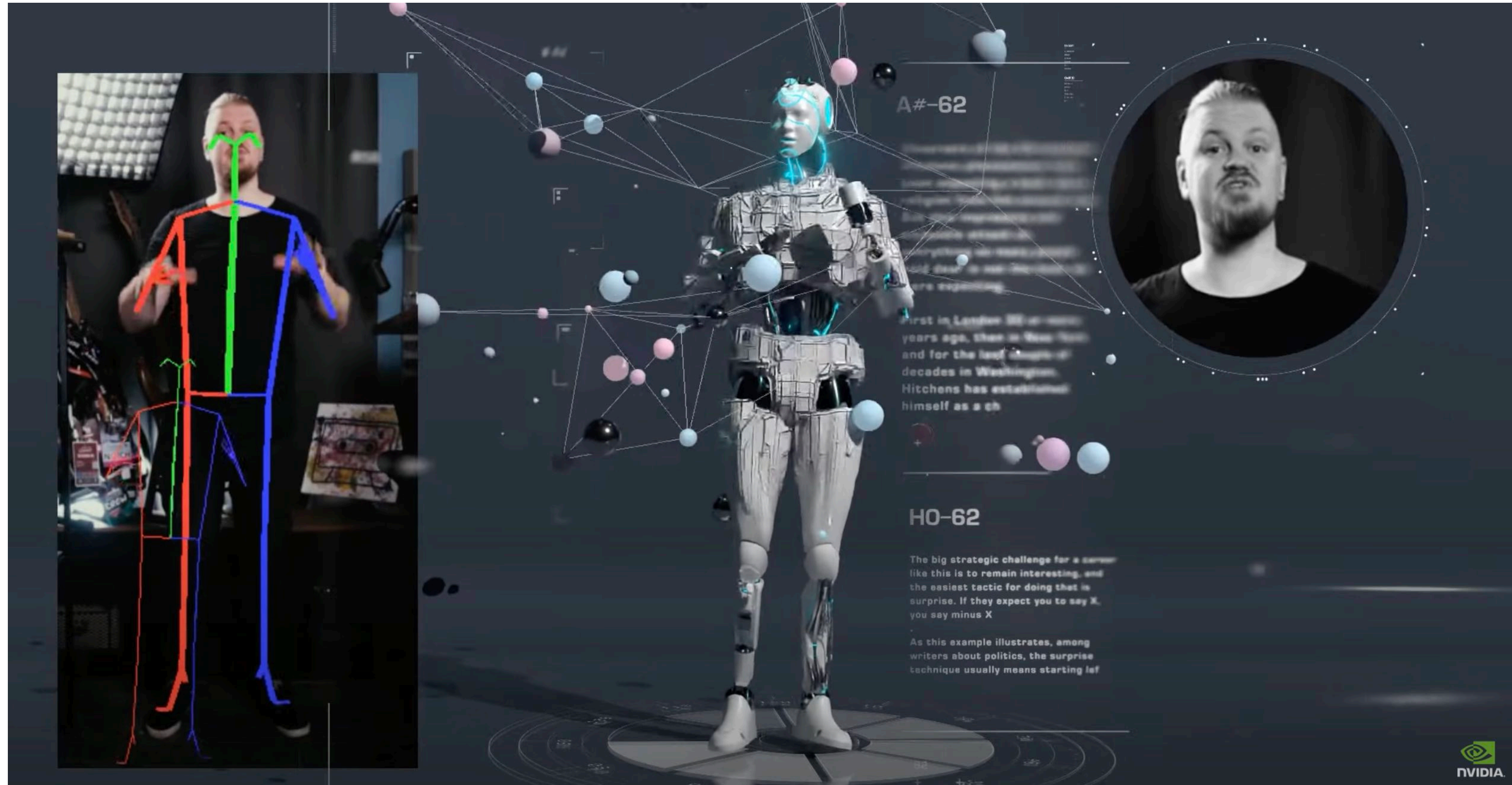
# Scalable video codec (SVC) encoder



**Costs: higher encoding/decoding costs**
**(But possible on modern clients as SVC is supported in hardware)**

# NVIDIA Maxine

**GPU-accelerated video processing for video conferencing applications**



**Examples: avatar control, video superresolution, advanced background segmentation**

# What do people *really need*?

# "Zoom fatigue" is very real



[Bailenson 2021]

FEBRUARY 23, 2021

## Stanford researchers identify four causes for 'Zoom fatigue' and their simple fixes

*It's not just Zoom. Popular video chat platforms have design flaws that exhaust the human mind and body. But there are easy ways to mitigate their effects.*

BY VIGNESH RAMACHANDRAN

Even as more people are logging onto popular video chat platforms to connect with colleagues, family and friends during the COVID-19 pandemic, Stanford researchers have a warning for you: Those video calls are likely tiring you out.

Prompted by the recent boom in videoconferencing, communication Professor Jeremy Bailenson, founding director of the Stanford Virtual Human Interaction Lab (VHIL), examined the psychological

1) Excessive amounts of close-up eye contact is highly intense.

2) Seeing yourself during video chats constantly in real-time is fatiguing.

3) Video chats dramatically reduce our usual mobility.

4) The cognitive load is much higher in video chats.

# The best camera is the one that's off?

## Best funny Zoom background trick: Put yourself in a looping video so you can skip the meeting

Now you can duck out on those hourlong conference calls.

By **Gordon Ung**
Executive Editor, PCWorld | APR 13, 2020 3:30 AM PDT

### Yes, you can make a Zoom background of yourself pretending to pay attention

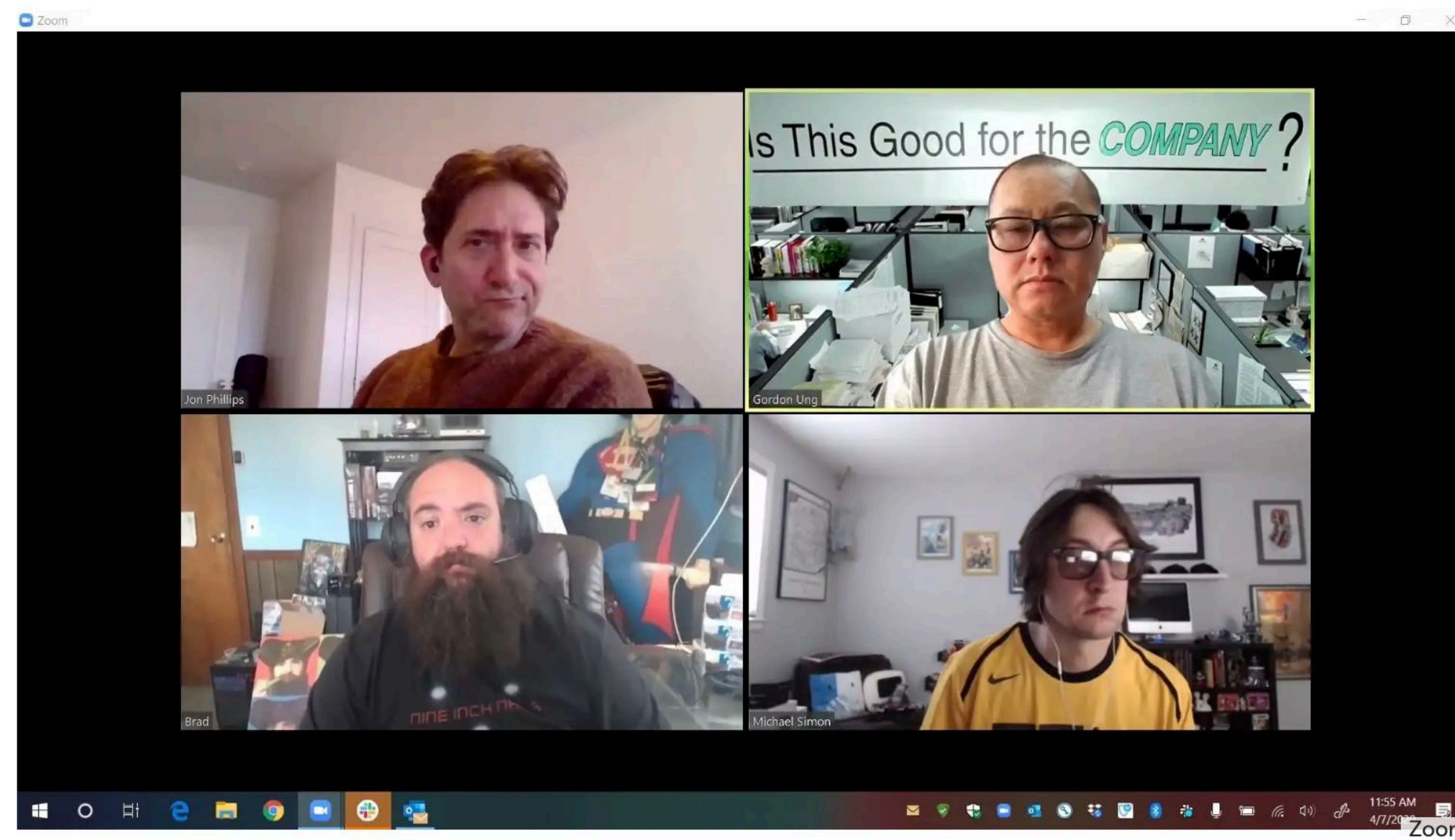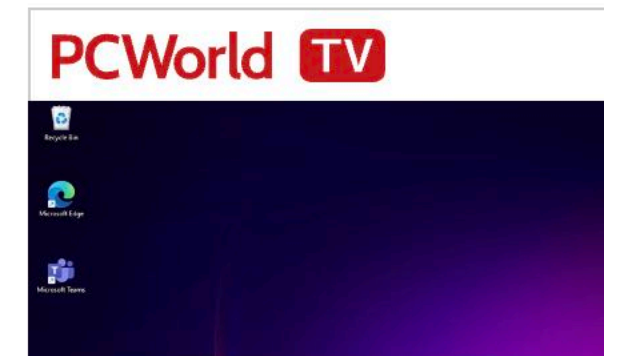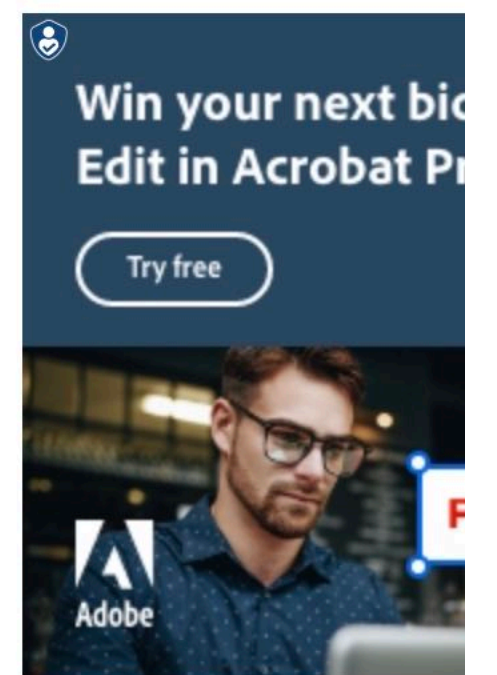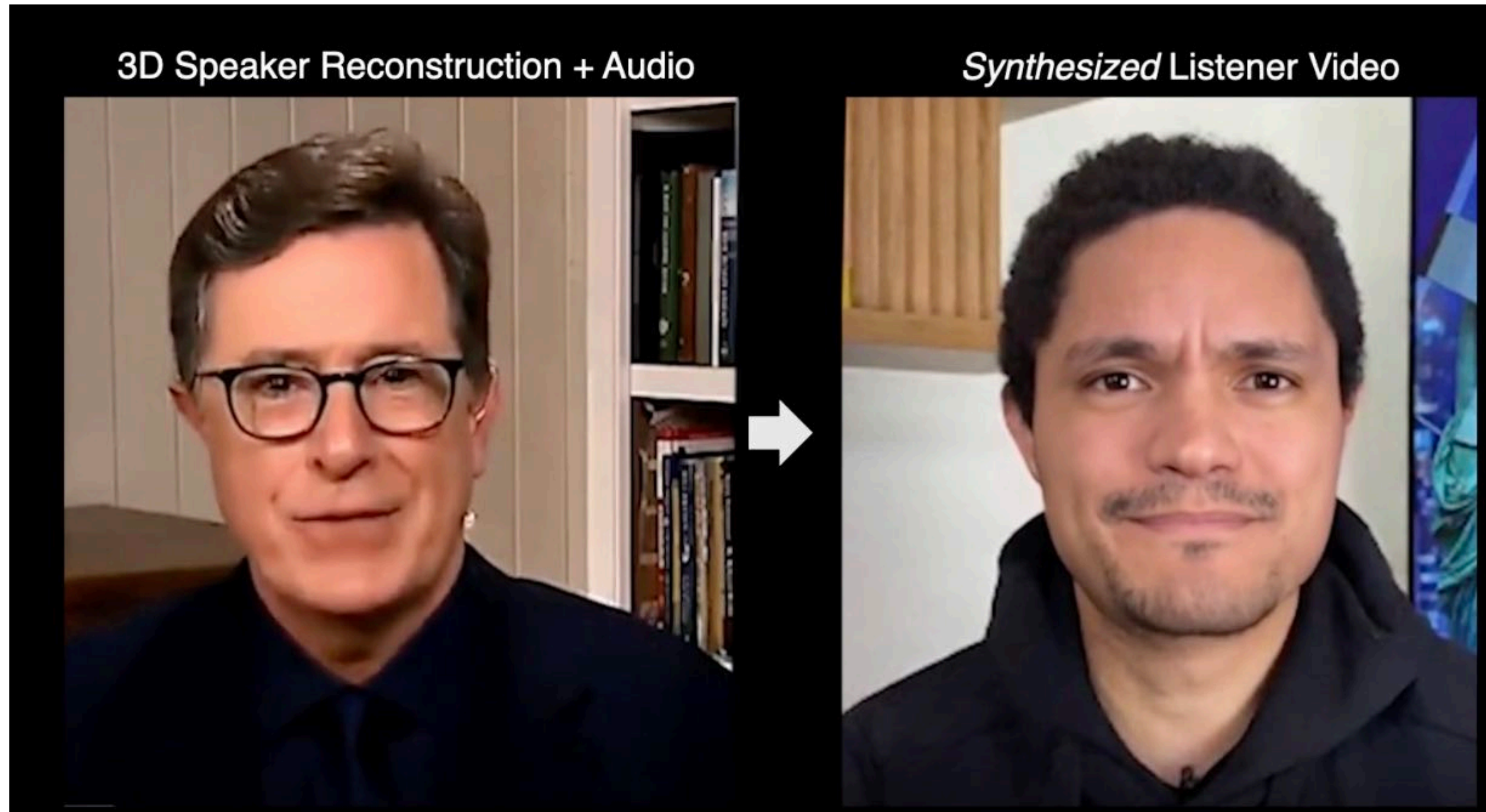And it's surprisingly easy to do, too.

Brian Lloyd
2 years ago

Share

We've all been in Zoom video conference meetings that drag on longer than a bad

Win your next bio
Edit in Acrobat Pr

Try free

Adobe

PCWorld TV

# Synthesizing reactions?

**Input: audio of speaker**

**Output: video of listener's reaction**



3D Speaker Reconstruction + Audio → *Synthesized* Listener Video

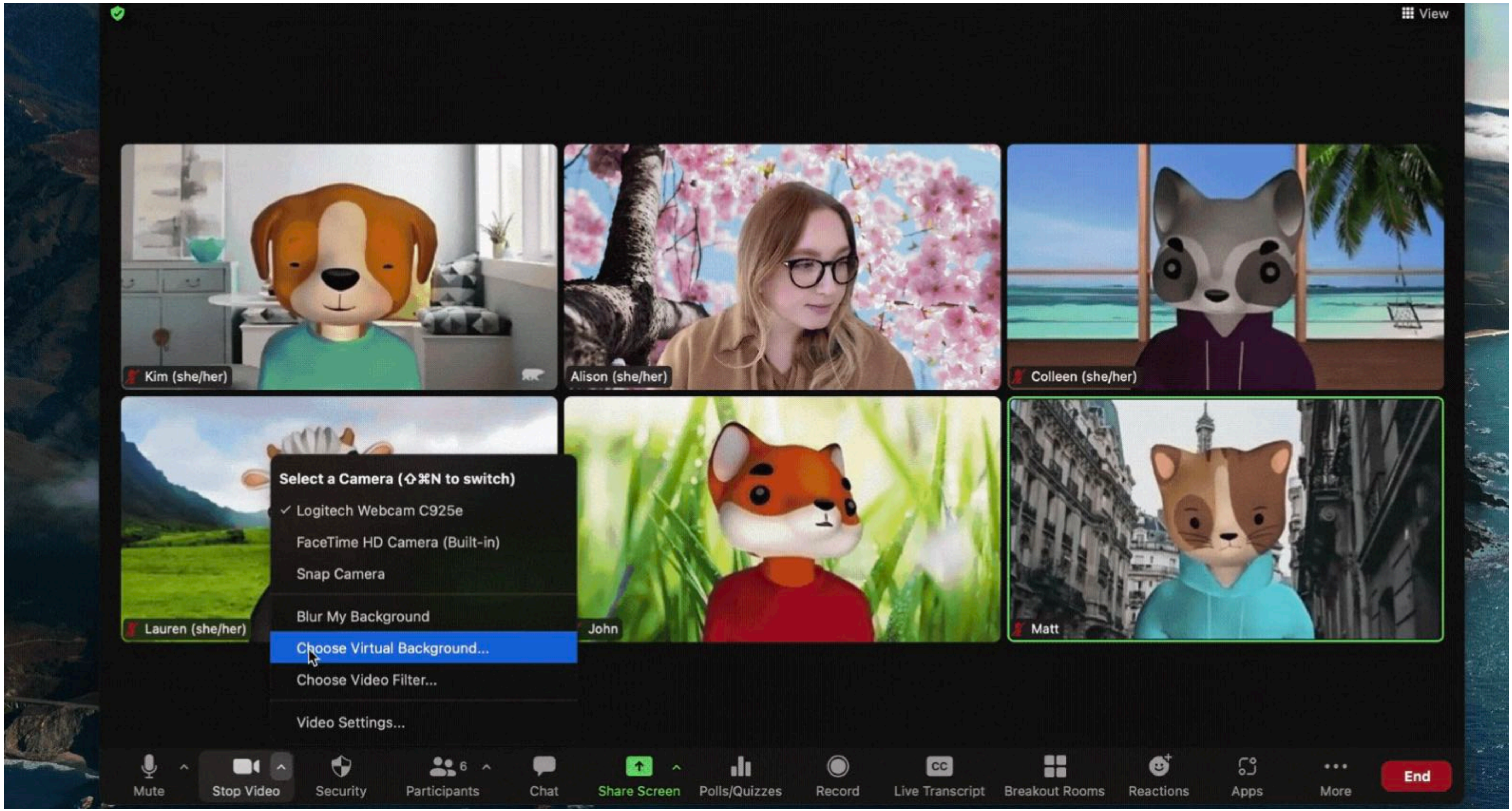# User-triggered effects (examples: audio clips, "reactions")

# Neural volumes

- **Learn to encode multiple views of a person into a latency code (z) that is decoded into a volume than can be rendered with conventional graphics techniques *from any viewpoint***



- **Motivated by VR applications**

# Zoom avatars / Snapcam Lenses

# More examples (demo)

# Discussion:

Where is the ethical line between "augmenting or abstracting what's real" and "fake"?

How can technology strike a balance between facilitating better forms of communication + a sense of presence (e.g., working from home) vs. ensuring privacy and personal space?

Can you think of ways where widespread use of near-photorealistic digital personals (e.g., for work calls) might lead to unexpected harms?

Do virtual meetings promote more diverse representation? Set it back?

# (If time)
# Co-designing video compressor and network transport

# Status quo

- **Video encoder proceeds to compress video frames, targeting a bit rate (on average) provided by the network protocol**

- **But any one frame may be too large or small (some may be hard to predict)**



target bit rate

**video codec**

*24 frames/s*

compressed frames

**transport protocol**

*300 packets/s*

**Encoder: targets an average bit rate (bits/second)**

**Protocol: attempts to determine and use the available capacity of the network**

**But generates individual frames
(which individually may or may not exceed network capacity)**

- **If the encoder overshoots, packet loss occurs. As a result, frames get dropped**

# Consider challenges



**Sender realizes packet carrying frame 2 has been dropped (e.g, it was too big)**

**But sender cannot re-encode frame at lower size because it's moved on and has different internal state**

# Stateless (functional) video encoder

```
// prob model: tables representing encoding of values in video stream
// reference_images contains three prior images
state:= (prob_model, reference_images[3]);


// just a full image
keyframe := image pixels for entire frame


// prediction_modes and motion vectors define how to predict current
// frame given decoder state
// residue is correction to this prediction
interframe := (prediction_modes, motion_vectors, residue)


// decoding a frame generates one image of pixels, and
// an updated decoder state
decode(state, compressed_frame) -> (new_state, image)


// generate an interframe approximating image given the current
// decoder state.  This operation requires expensive motion estimation.
encode-given-state(state, image, quality_param) -> interframe
```
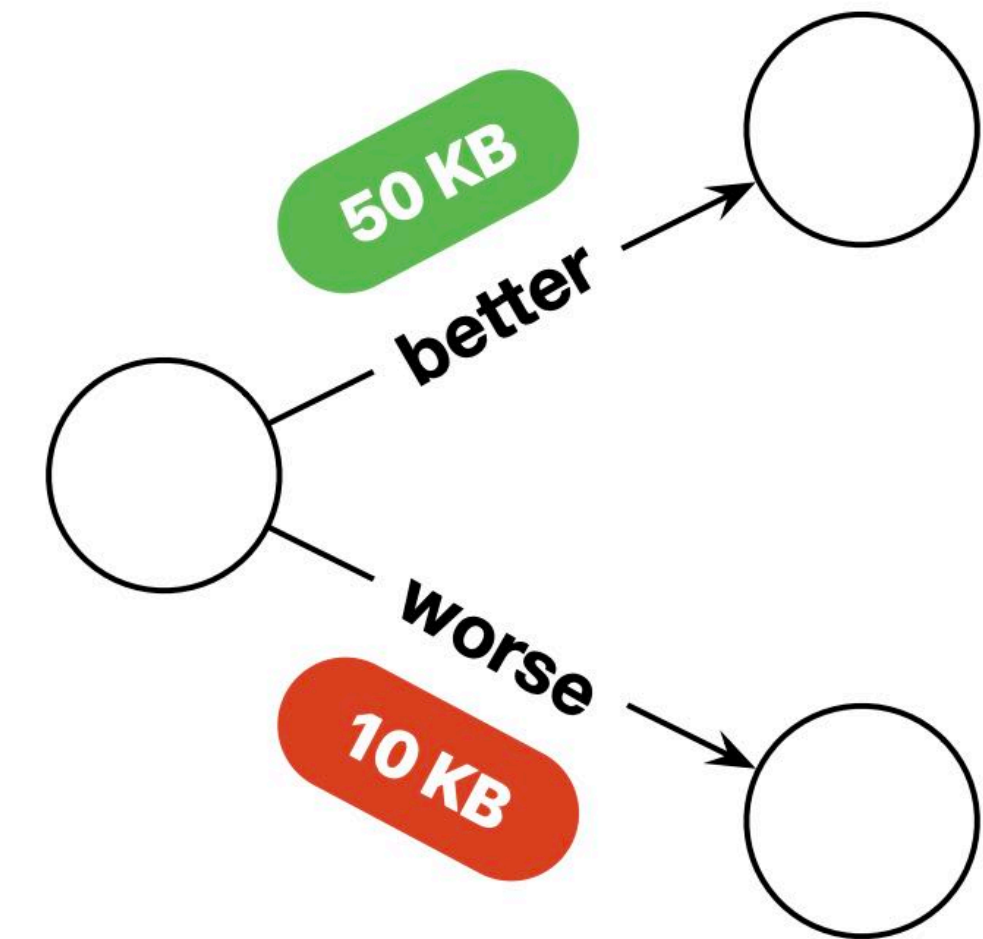
[Fouladi et al. 2017]

# Salsify: codec presents network three options

For each frame, codec presents the transport with *three* options:

🔺 A slightly-higher-quality version,

🔻 A slightly-lower-quality version,

✖ Discarding the frame.

50 KB

better

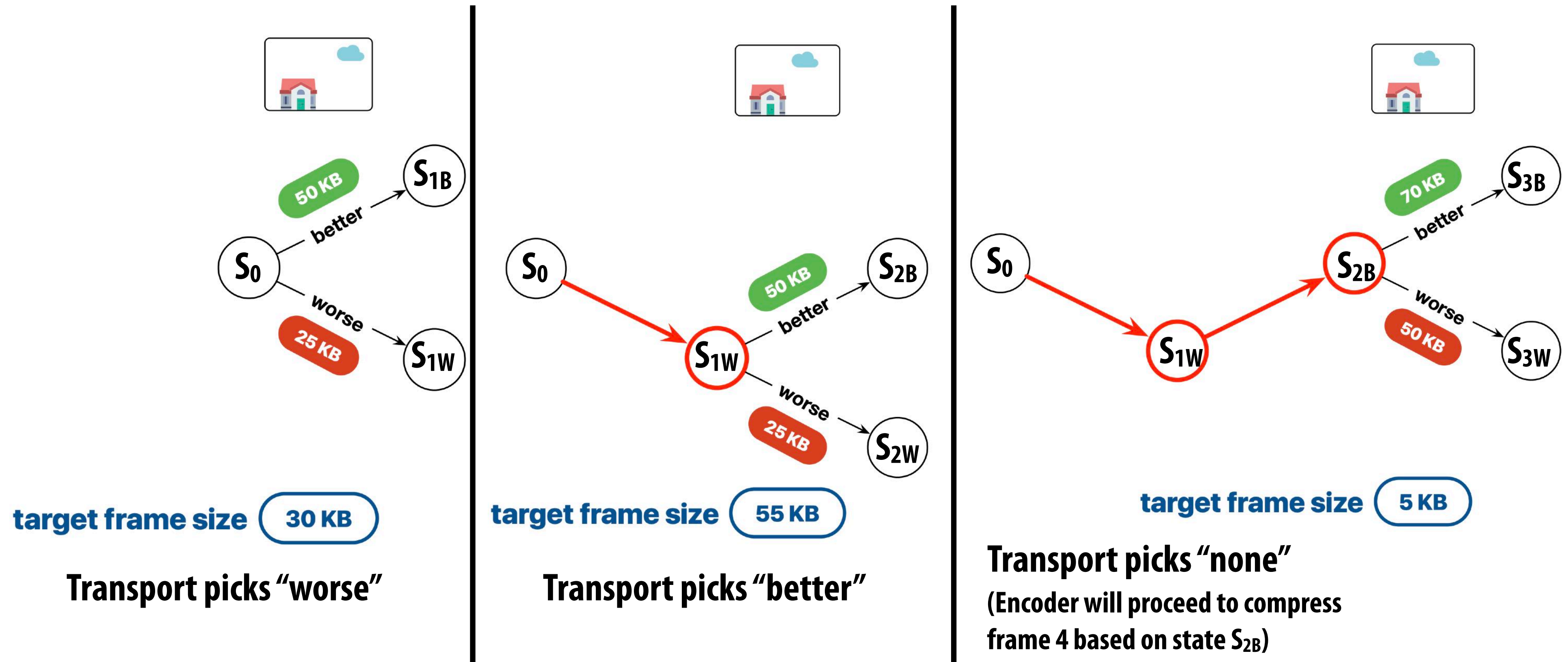worse

10 KB

**Notice roll of functional encoder.**

**Can encode "better", reset to previous state, and then encode "worse".**

# Salsify's "video aware transport protocol: network determines what to transmit based on size of compressed frames
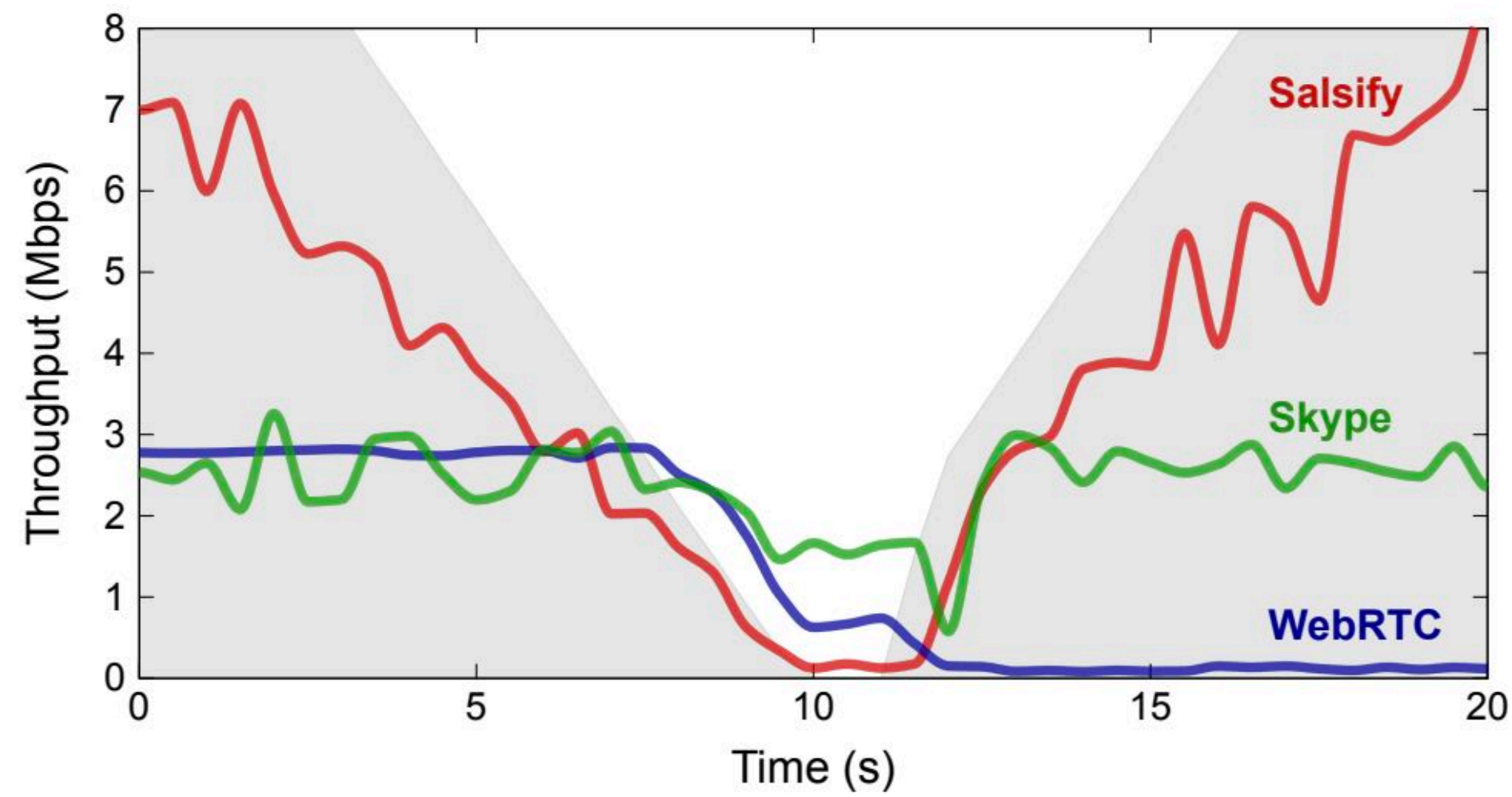
Before: network tried to send whatever the compressor generated.

Notice roll of functional encoder.

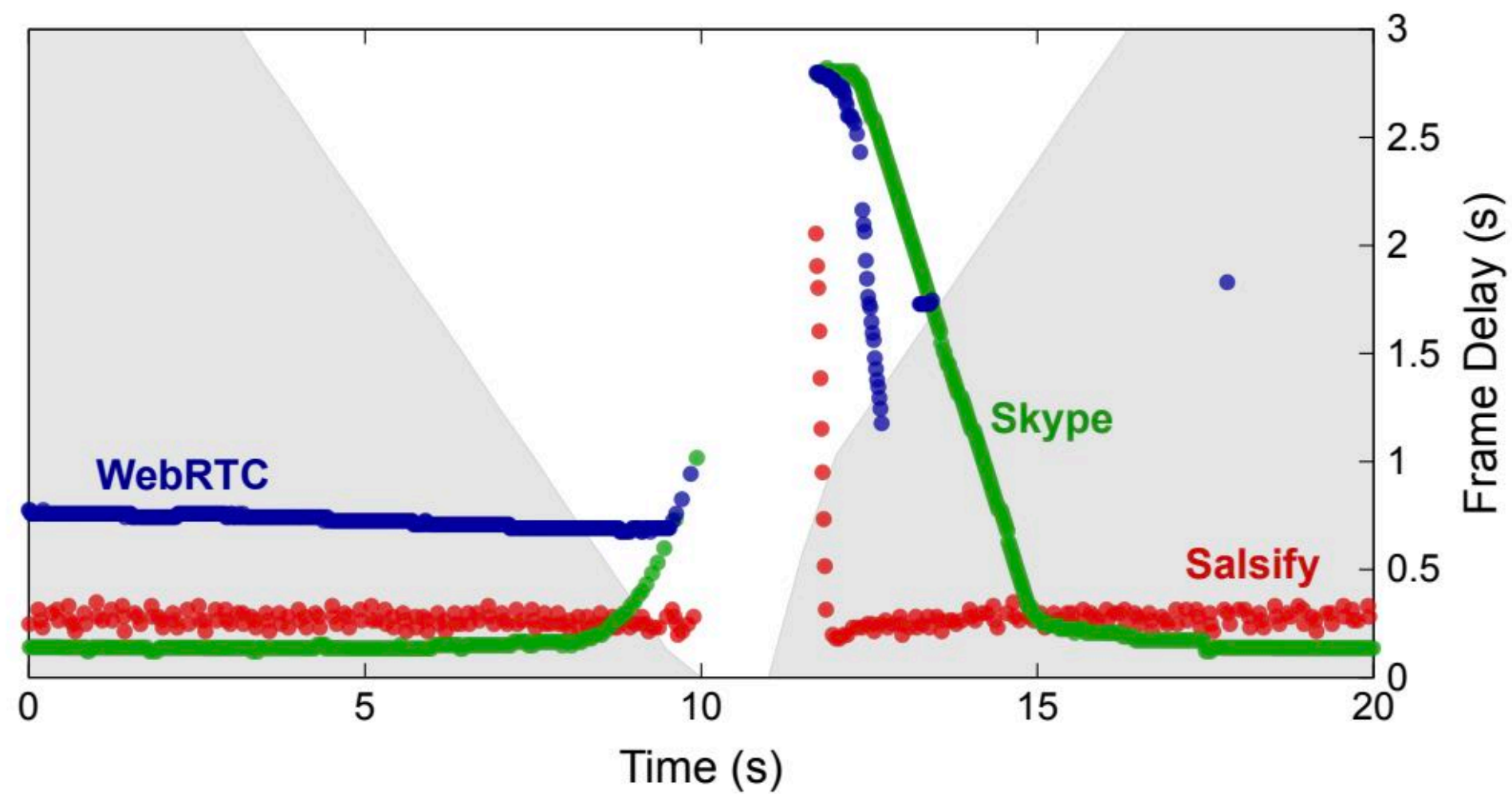Can resume encoding from state that results from transport's choice.



target frame size ( 30 KB )

**Transport picks "worse"**

target frame size ( 55 KB )

**Transport picks "better"**

target frame size ( 5 KB )

**Transport picks "none"**
(Encoder will proceed to compress
frame 4 based on state $S_{2B}$)

# Much faster recovery from network changes



(a) Throughput

Gray region shows capacity of network:
(Simulating an outage at 10 seconds)



(b) Frame delay