# Domain-Specific Fusion Of Multiple Objective Quality Metrics

## Presenter: Yiannis Andreopoulos

VQEG meeting, May 2022
Rennes, France

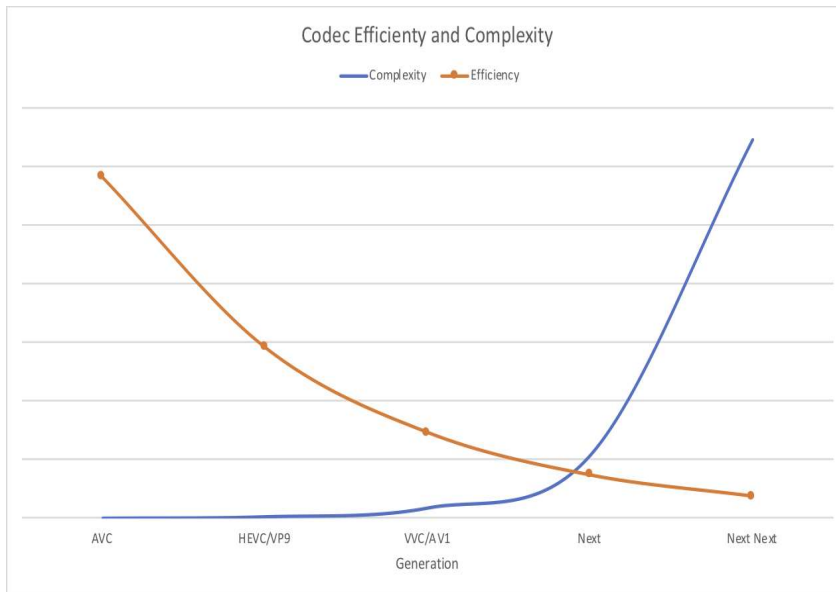Joint work of iSIZE, industrial collaborators and the Innovate UK SEQUOIA consortium, project: 96984

# iSIZE: What we do

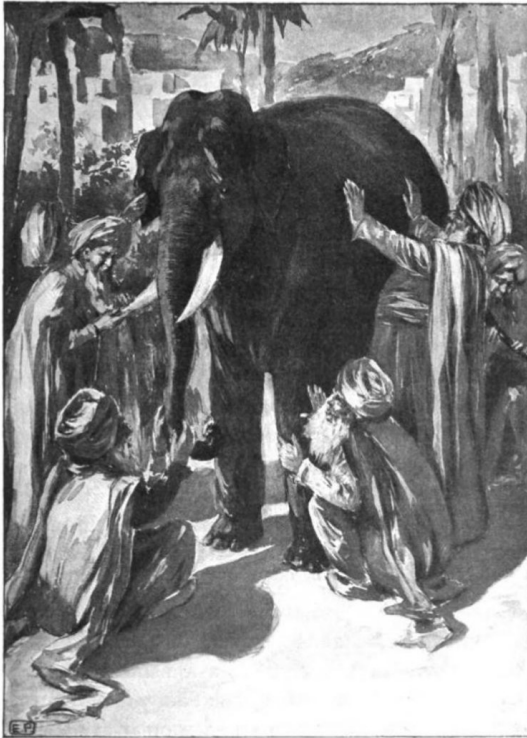| Problem we solve | Solution | Target market | Description |
|---|---|---|---|
| **Perceptual Quality** | Deep perceptual optimization<br><br>**BITSAVE**<br>demo.isize.co | • Entertainment / Media<br>• Video streaming<br>• Gaming<br>• Social media | • Deep psychovisual preprocessing for maximum bitrate savings. Significantly advance the development of AI-based quality metrics and quality scoring<br>• AI-based preprocessing that requires no change in encoding, delivery or decoding devices |
| **Noisy Video Content** | Deep perceptual denoising<br><br>**BITCLEAR**<br>http://bitclear.isize.co/ | • Social media/user uploads<br>• Post-decoder enhancement | • Remove compression noise from video content by addressing the problem across the quality-bitrate-complexity space<br>• Can work both as a server and as a client component (post-decoder) |
| **Low-Bitrate/Low-Latency Video Delivery** | Domain-specific generative video representations<br><br>**BITGEN** | • Conversational services<br>• Virtual reality/telepresence<br>• IoT/driverless technologies | • Extreme reduction in video bitrate, working in a compact latent space.<br>• Enable telepresence with near-zero latency. |

# iSIZE: Why we do it



Codec Efficienty and Complexity

https://www.linkedin.com/pulse/encoder-complexity-hits-wall-david-ronca/ (D. Ronca, Meta, 2019)



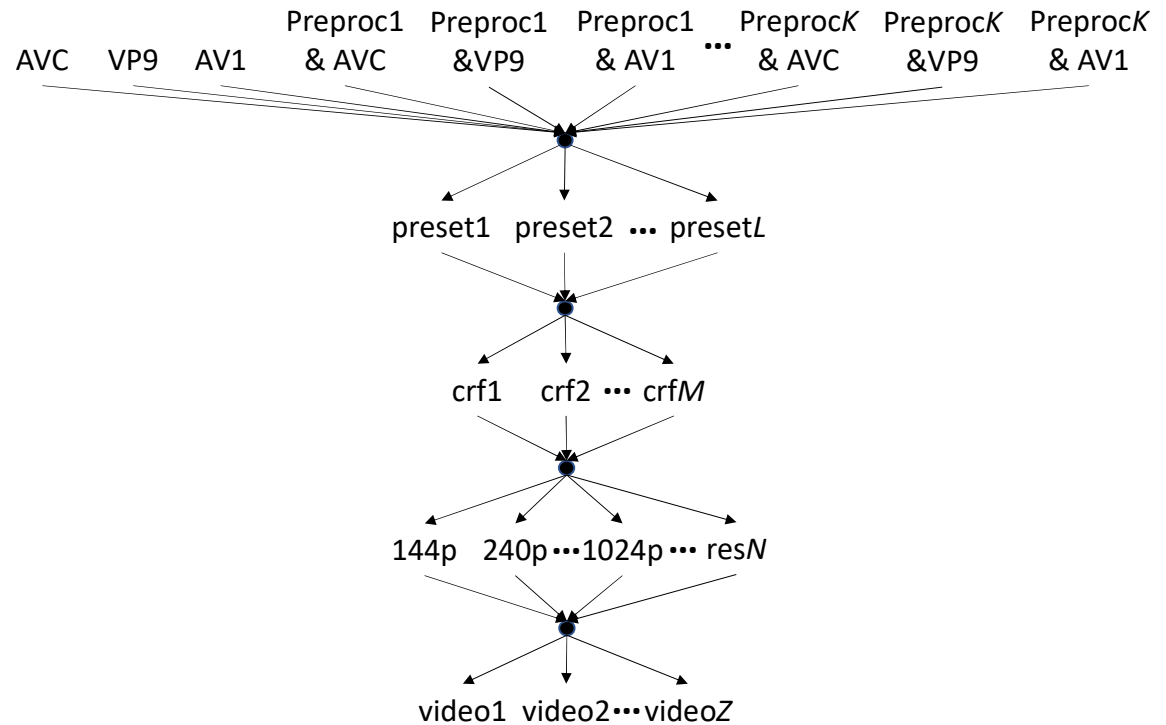Sikora, Proc. of the IEEE, 2005, https://doi.org/10.1109/JPROC.2004.839601

- Device power+heat dissipation and cloud-based scaling have both hit the wall
- Inflection point: quality metrics and neural network hardware now allow for AI-based pre- and post-processing
- Codecs are amazing SNR/SSIM-vs.-bitrate machines, but these loss functions have significant limitations

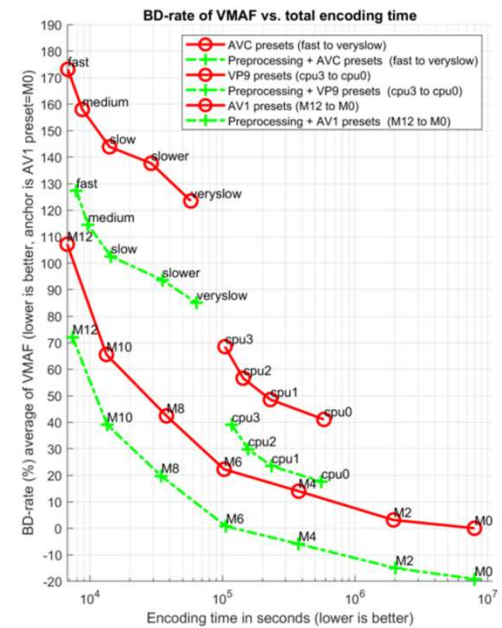# iSIZE The three challenges with visual quality assessment

1. Objective metrics (and humans) are myopic

AVC  VP9  AV1  Preproc1 & AVC  Preproc1 &VP9  Preproc1 & AV1  ···  PreprocK & AVC  PreprocK &VP9  PreprocK & AV1

preset1  preset2  ···  presetL

crf1  crf2  ···  crfM

144p  240p ··· 1024p ··· resN

video1  video2 ··· videoZ

2. The exploration space can surpass 1m tests for a 100-video library

# iSIZE The three challenges with visual quality assessment



3. Video processing algorithms are now increasingly optimized for perceptual quality metrics instead of signal distortion

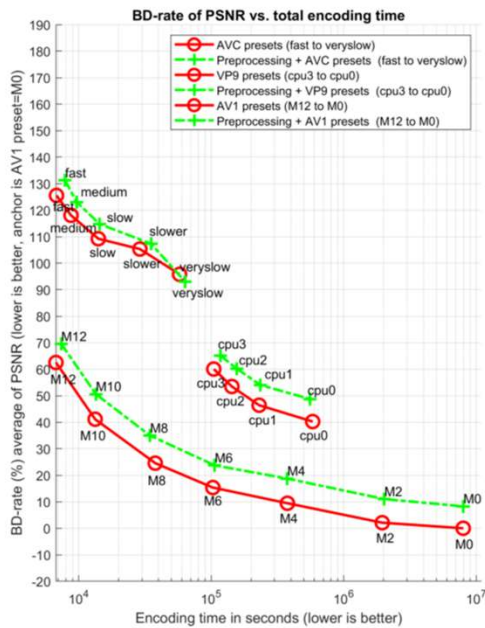→ This means that they may score well for metrics like VMAF, but this may be because of metric overfitting.

# Methods: **iSIZE BitSave preprocessing (mk3)**

Source → BitSave → Encoder → Delivery

AI-based pre-processing prior to encoding (AVC, HEVC, VP9, AV1)

One frame latency

Single pass processing per content for an entire ABR ladder

Improves encoding quality as measured by standard perceptual quality metrics (VMAF, SSIM, VIF), can also work in tandem with BitClear

Integrated within Intel OpenVINO, ONNX and Dolby Vision, easy to plug&play within any existing workflow

# Methods: iSIZE BitClear post-processing (mk3.5)

Delivery → Decoder → **BitClear** → Display

AI-based post-processing after decoding (AVC, HEVC, VP9, AV1)

One frame latency

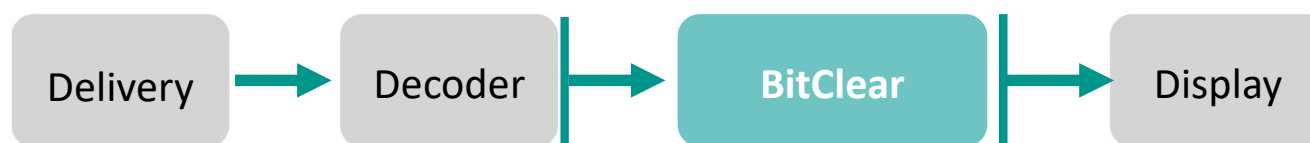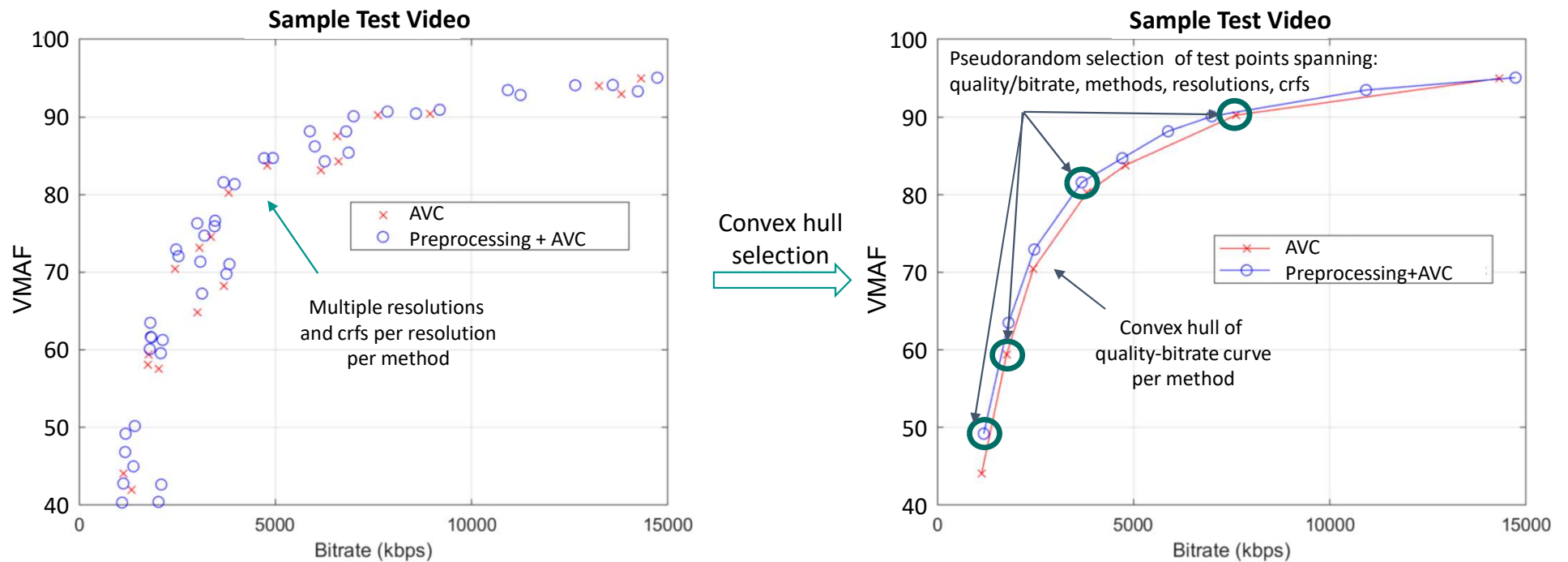Single pass processing per content for an entire ABR ladder

Improves decoding quality as measured by standard perceptual quality metrics (VMAF, SSIM, VIF), can also work in tandem with BitSave

Integrated within Intel OpenVINO and ONNX, easy to plug&play within any existing workflow

# Domain-specific fusion of multiple quality metrics



Sample Test Video
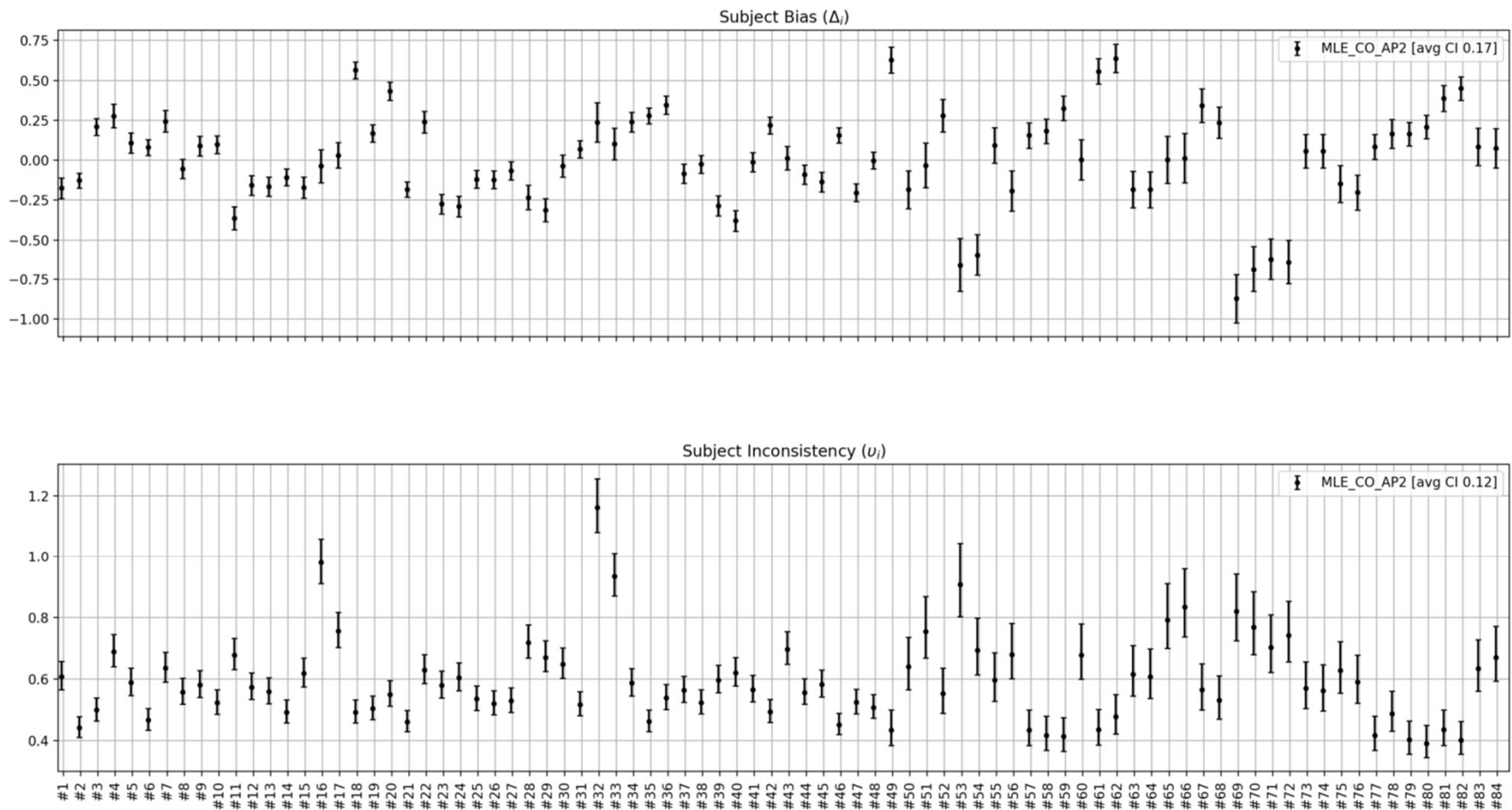
Convex hull selection

Sample Test Video

Three steps:

1. Convex-hull selection based on VMAF

2. Carry out P.910 ACR and post-processing

3. Fuse metrics to recovered quality scores using support vector regression
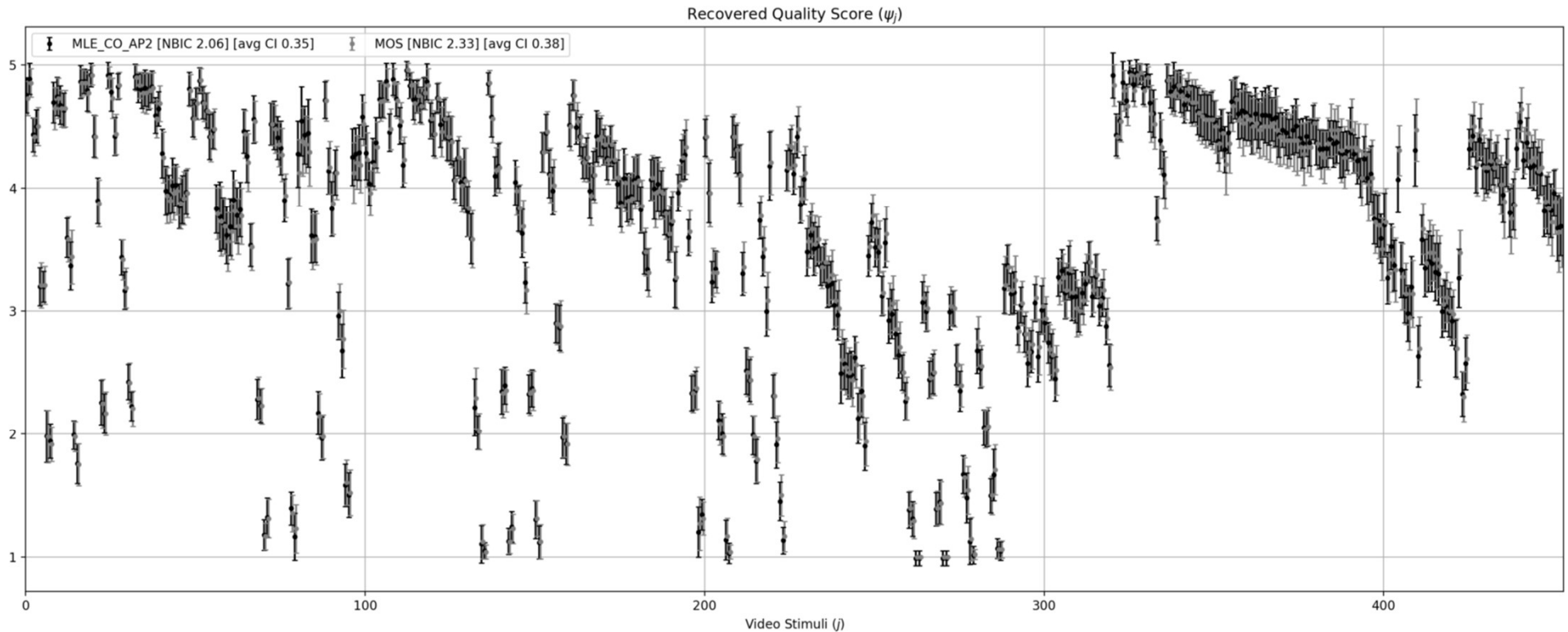
# P.910 ACR test setup and conditions

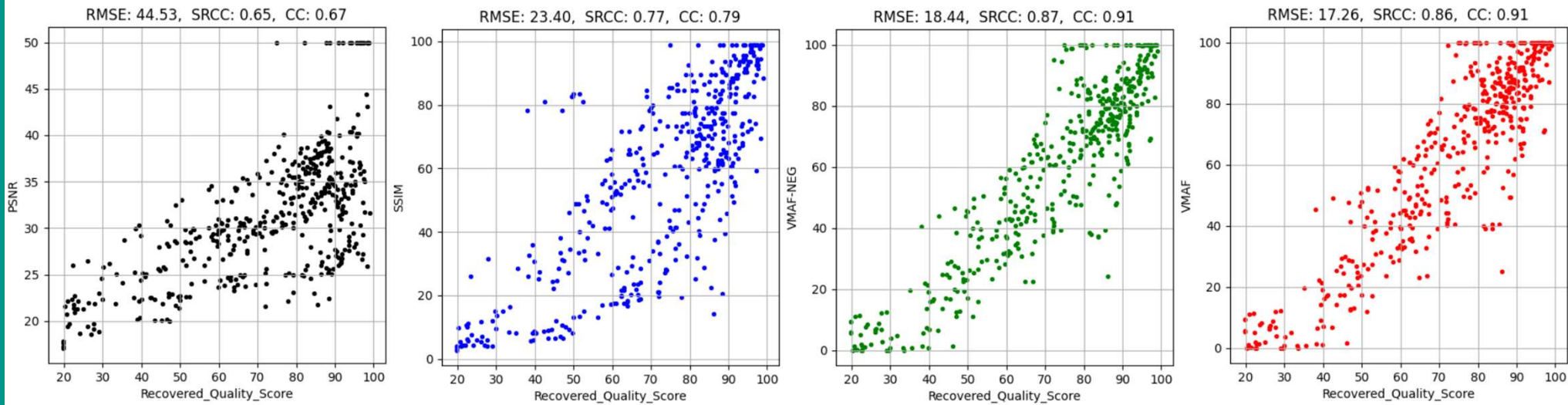| Setup Component | What was Used | Further Details on Settings | Comments |
|---|---|---|---|
| **Encoders** | AVC x264 (Lavc58.134.100 libx264)<br><br>WEBM VP9 (v1.10.0-48-g4ec84326c) | • 1080p, 720p, <u>540p</u>, <u>360p</u>, 216p (only underlined done for post-processing)<br>• Per resolution: AVC preset=veryslow, CRF={22,30,38,46} (medium used for post-processing)<br>• Per resolution: VP9 preset=0, CRF={32,38,44,<u>46</u>,<u>48</u>,<u>50</u>,52,<u>54</u>,<u>56</u>,<u>58</u>,<u>60</u>} (underlined CRFs done for 720p & 540p, preset=5 used for post-processing) | • The slowest preset of each encoder was used for preprocessing, faster presets for post-processing<br>• Constant-CRF encoding ensures quality remains consistent, no effects from rate control algorithms<br>• The range of CRFs ensures the full quality range of relevance to each resolution & application is sampled<br>• All lower resolutions were upscaled to 1080p for viewing using FFmpeg Lanczos-5 |
| **Content and test conditions** | AV2 CTC content https://media.xiph.org/video/aomctc/test_set/ P.910 ACR standard test conditions applied | • 3H distance, controlled lighting, same screen conditions for all tests<br>• Ratings from 1-5<br>• Raters were briefed on task and how to use the quality scaling | • All content replayed at 25fps, 1080p@50Hz TV screen, all TV filters were off<br>• 21 sequences at 1080p resolution (8bit) used, comprising a mixture of entertainment, sports, UGC, gaming, web browsing, and artistic content (16 sequences for post-processing) |
| **Raters and data processing** | • 48 raters for preprocessing (the underlined VP9 CRFs had 36 additional raters)<br>• 24 raters for post-processing<br>• The SUREAL package was used for post-processing | • All raters were screened for color blindness and good eyesight<br>• All 16368 ratings were used | • SUREAL: https://github.com/Netflix/sureal<br>• The full maximum likelihood estimation (MLE) model of SUREAL was used, which takes into account both subjects and contents<br>• For quality-bitrate plots per resolution and cross-resolution combined quality-bitrate plots, an MLE fit per codec was carried out and the recovered quality scores were used |

# Preprocessing results: Subject bias & inconsistency



Subject Bias ($\Delta_i$)

Subject Inconsistency ($\upsilon_i$)

# Preprocessing results: Recovered quality scores



Recovered Quality Score ($\psi_j$)

MLE_CO_AP2 [NBIC 2.06] [avg CI 0.35]   MOS [NBIC 2.33] [avg CI 0.38]
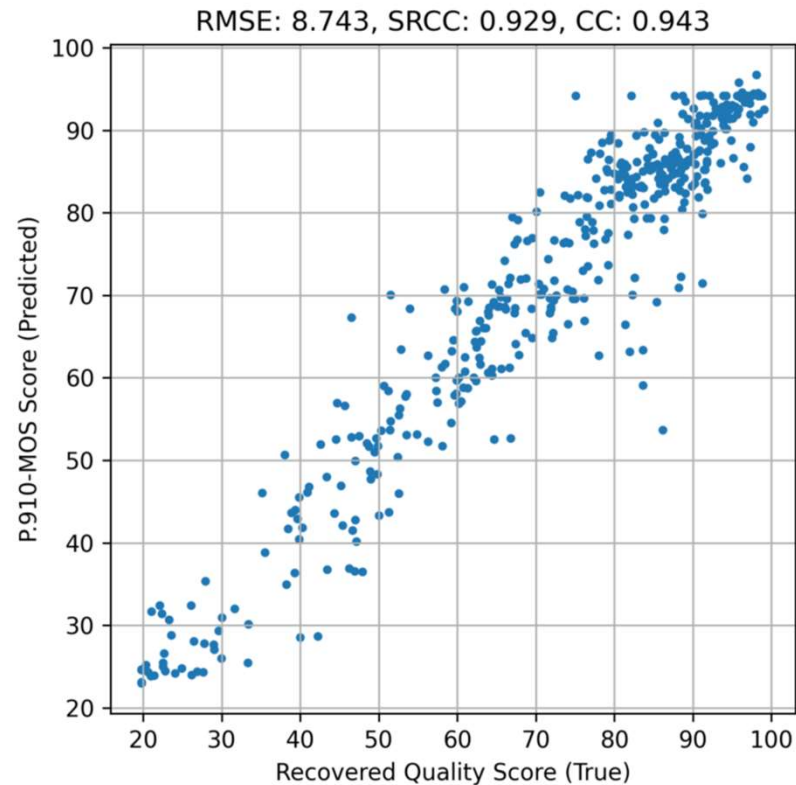
Video Stimuli ($j$)

- The Recovered Quality Scores (RQS) span the entire quality range and are adjusted according to bias, uncertainty and inconsistency based on SUREAL's methodology

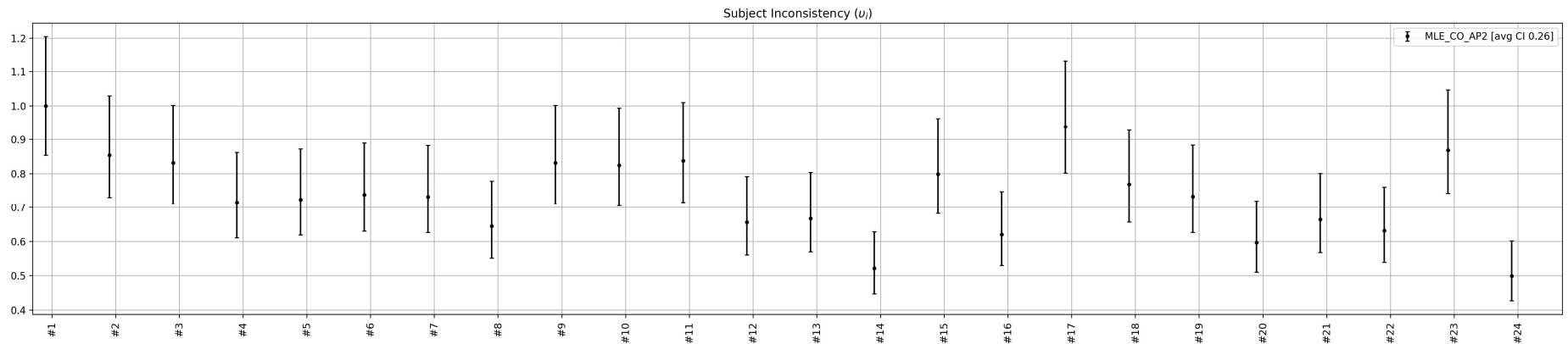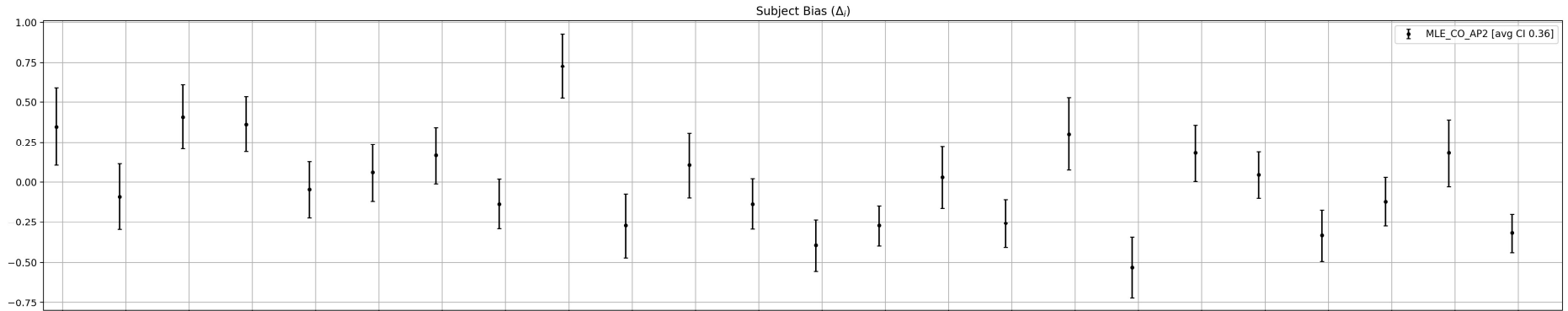# Preprocessing results: Metrics vs. RQS scatter plot



RMSE: 44.53, SRCC: 0.65, CC: 0.67

RMSE: 23.40, SRCC: 0.77, CC: 0.79

RMSE: 18.44, SRCC: 0.87, CC: 0.91

RMSE: 17.26, SRCC: 0.86, CC: 0.91

- VMAF-NEG and VMAF are well aligned to Recovered Quality Scores, with correlation of 91%
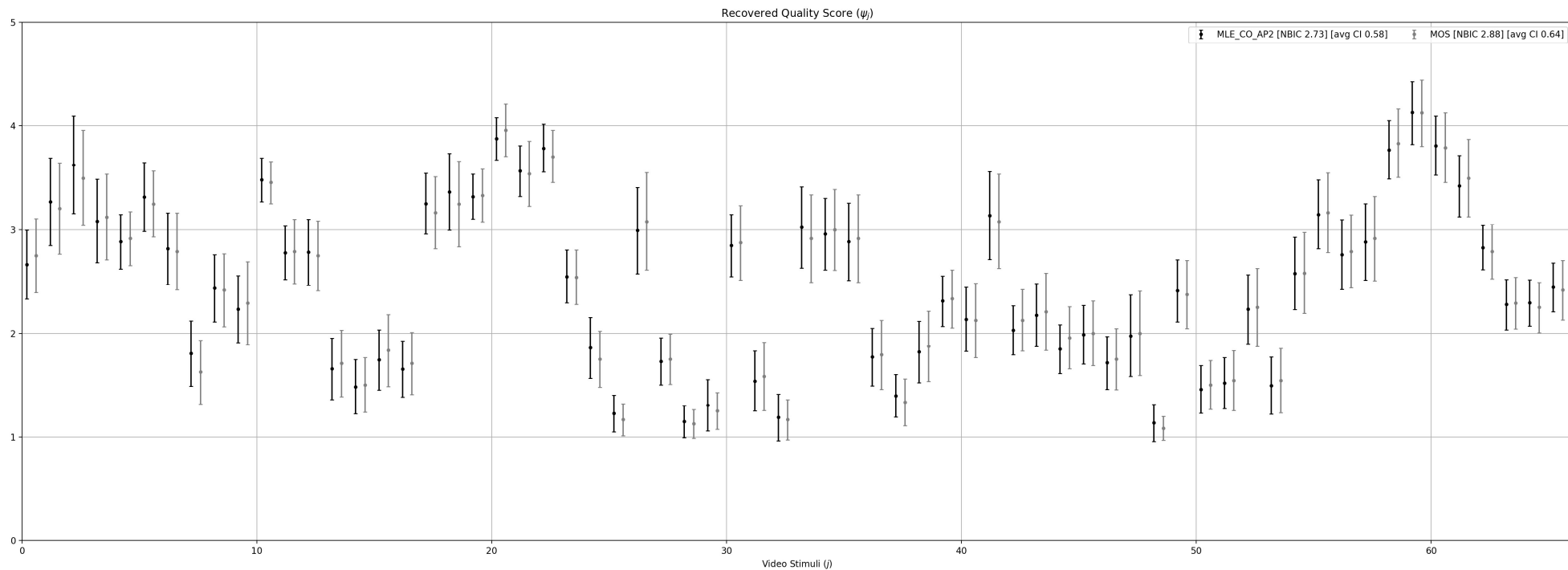
# **Preprocessing results: P.910 SVR model**



- Scatter plot of SVR with $v$=0.5 (proportion of support vectors vs. total samples), $\gamma$=0.85 (radius of RDF), $C$=1 (regularization term) predicted scores vs recovered quality scores

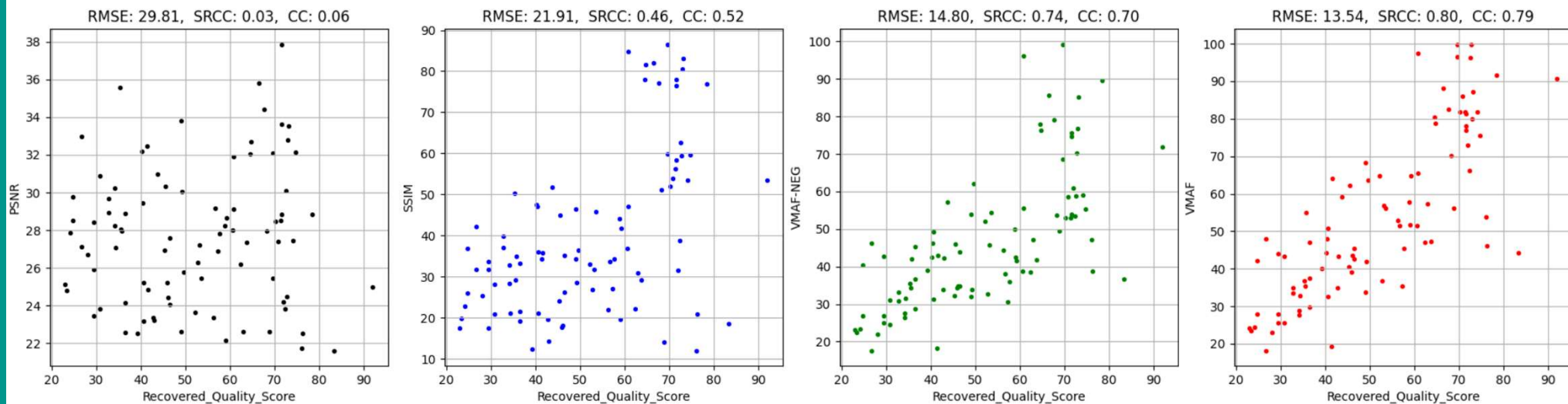# Post-processing results: Subject bias & inconsistency



Subject Bias ($\Delta_i$)

MLE_CO_AP2 [avg CI 0.36]

Subject Inconsistency ($\upsilon_i$)

MLE_CO_AP2 [avg CI 0.26]

# Post-processing results: Recovered quality scores



Recovered Quality Score ($\psi_j$)

MLE_CO_AP2 [NBIC 2.73] [avg CI 0.58]     MOS [NBIC 2.88] [avg CI 0.64]
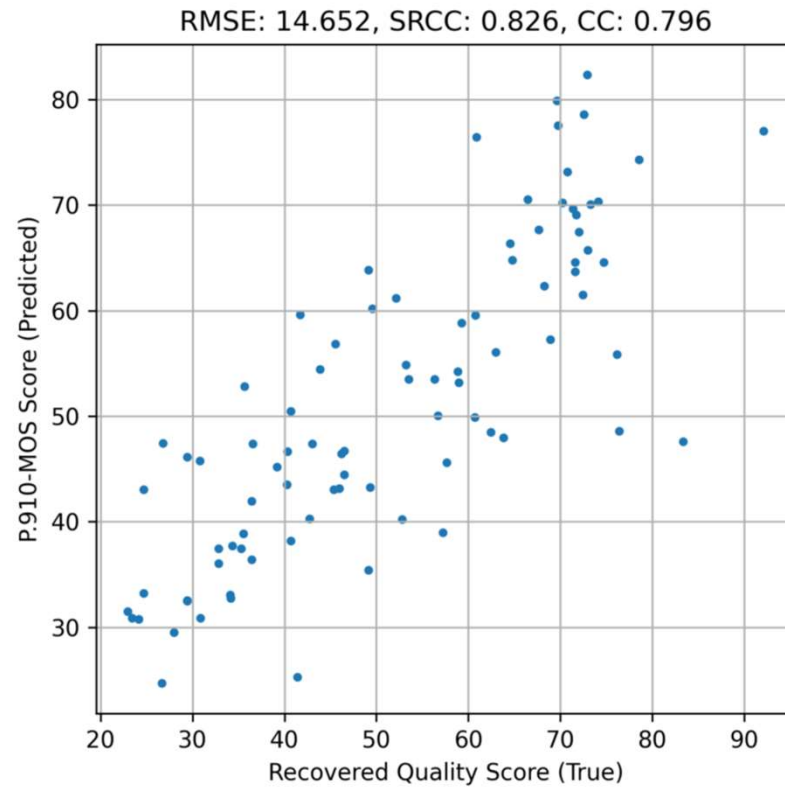
Video Stimuli ($j$)

- The Recovered Quality Scores (RQS) span the entire quality range and are adjusted according to bias, uncertainty and inconsistency based on SUREAL's methodology

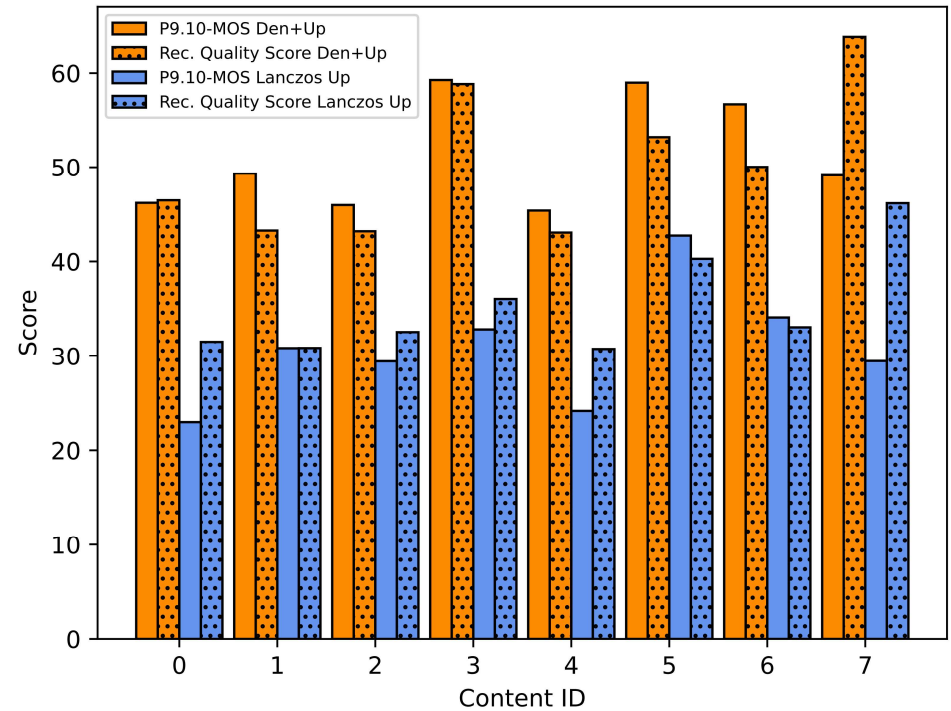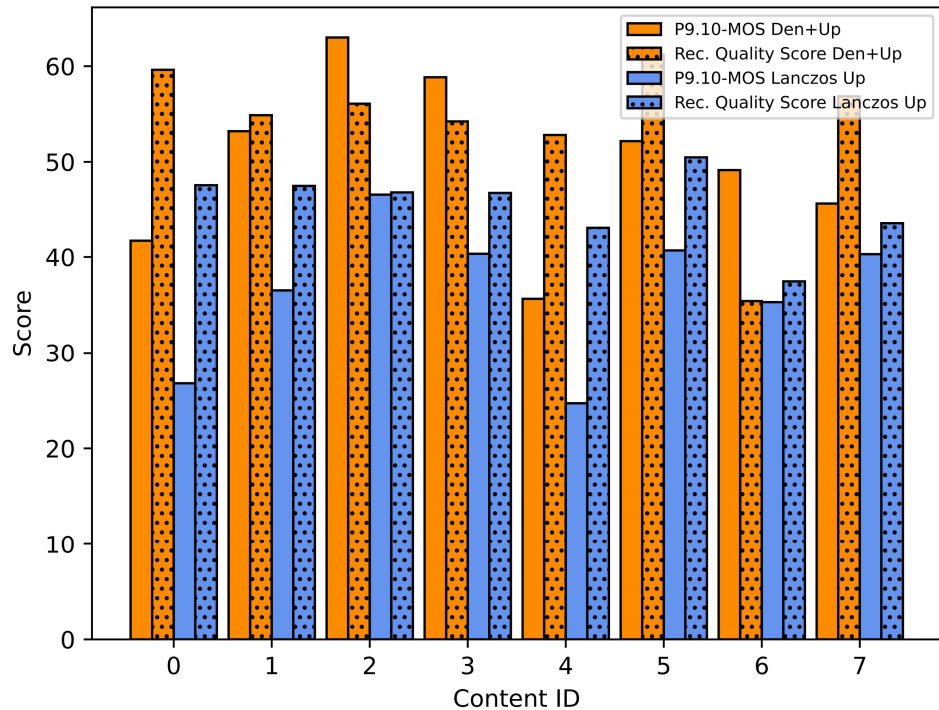# Post-processing results: Metrics vs. RQS scatter plot



- VMAF-NEG and VMAF are better aligned to Recovered Quality Scores, with correlation of 70% to 79%

# Post-processing results: P.910 SVR model
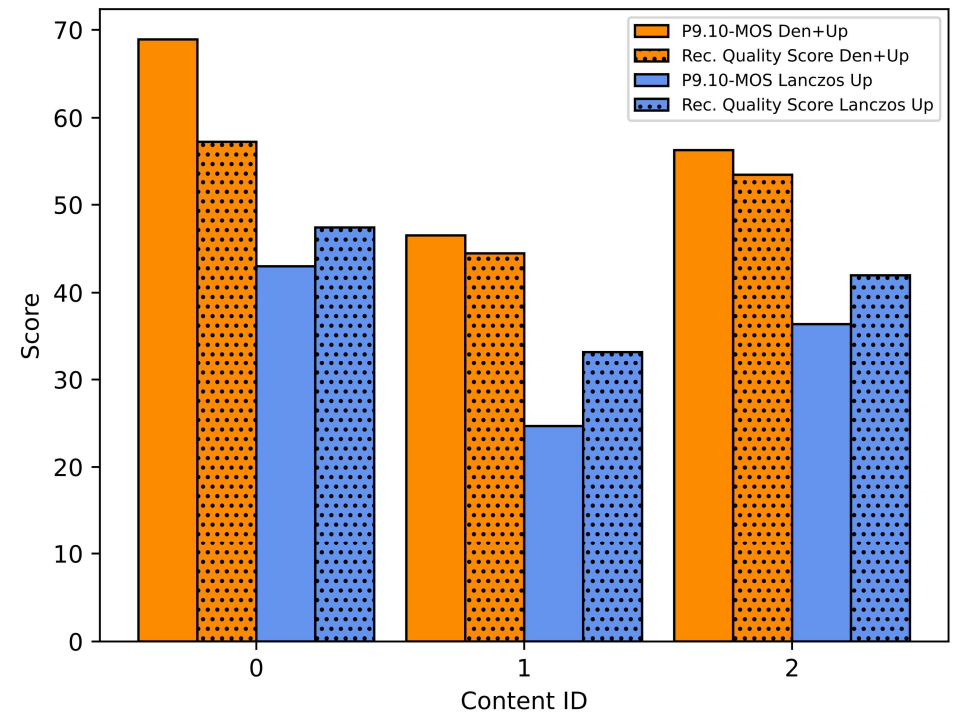


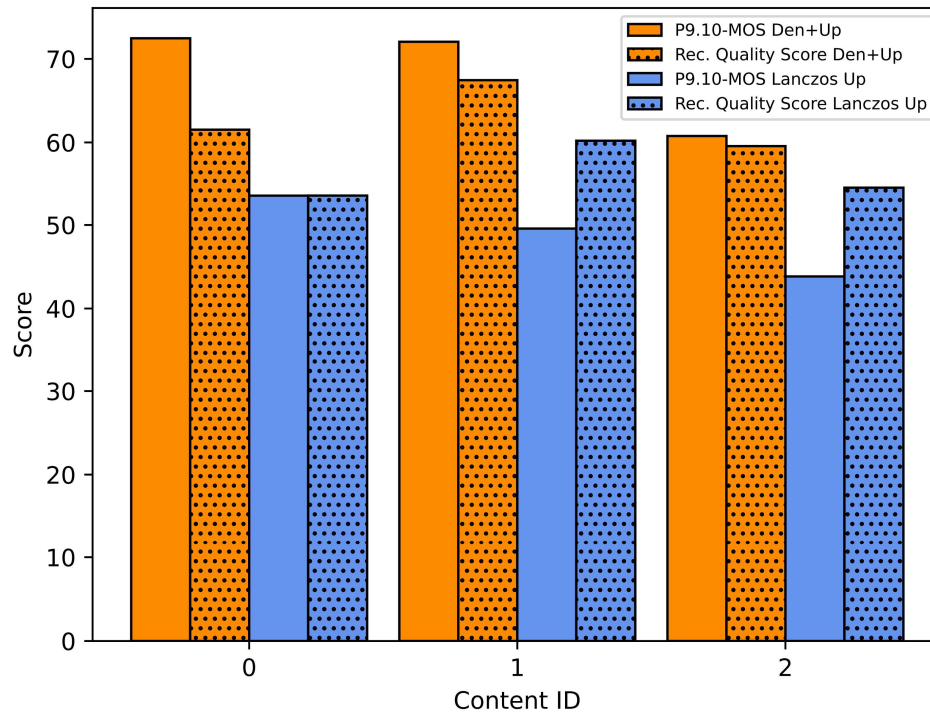RMSE: 14.652, SRCC: 0.826, CC: 0.796

- In this case, it is mainly VMAF-NEG and VMAF that contribute to the SVR fit

# Post-processing results: Bar plots for AVC medium+high CRF



- In this case, P.910-MOS of post-processing ("Den+Up") and Lanczos comes close to the results of VMAF

# Post-processing results: Bar plots for VP9 medium+high CRF



- In this case, P.910-MOS of post-processing ("Den+Up") and Lanczos comes close to the results of VMAF

# Conclusion: Some key take-aways

- Domain-specific fusion of metrics can help get closer to the true ACR recovered quality scores if a single metric does not dominate

- The presented methodology is easy to apply and allows for quick testing (and re-testing!) as versions improve

- In the case of iSIZE preprocessing, we found that P.910 ACR results come between VMAF-NEG and VMAF (i.e., VMAF-NEG with some allowance for gain limit may suffice (e.g., 2%-5%)

- In the case of iSIZE post-processing, due to the use of GAN losses, only VMAF-NEG and VMAF remain relevant; the overall average gains of post-processing were 1 point in the 5-point ACR scale or 14 VMAF points

- Pseudo-random sampling of the convex hull of points helps (100-fold reduction in sampling), there are probably further ways of optimizing the distribution sampling that we have not considered

- It would probably be interesting to add other metrics (LPIPS, no-reference metrics) to our tests