

Updates on Maximum Likelihood Estimation (MLE) Methods for Subject Behavior Modeling

Zhi Li & Christos Bampis, *Netflix, USA*
Lucjan Janowski, *AGH, Poland*
Ioannis Katsavounidis, *Facebook, USA*

VQEG Meeting March 2020
3/10/2020

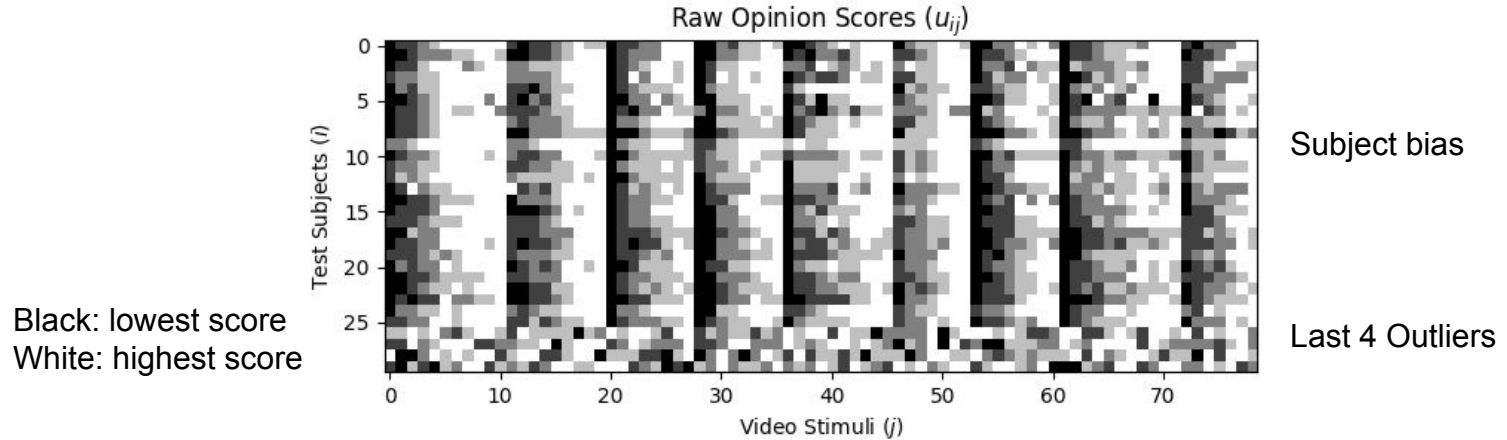
MLE Activities in SAM Study Group

- Motivation
 - Mission: to improve data quality coming from subjective experiment
 - Saw opportunity to improve data cleanup methods currently adopted in ITU-T/R recommendations via statistical methods (e.g. MLE)
- Progress since the last VQEG meeting
 - A comprehensive study with comparison to prior standards ITU-R BT.500 and ITU-T P.913
 - Bayesian Information Criterion (BIC) to validate model fitting with real-world data
 - Besides the original Newton-Raphson method, proposed an alternative solution based on projection, proven to be faster and more intuitive

Subjective Test



Raw opinion scores are noisy and unreliable



- Would MOS or DMOS be good enough?
- Corrective mechanisms
 - Subject outlier rejection
 - Subject bias removal

Prior Art: Subject Outlier Rejection (ITU-R BT.500)

For each test presentation, calculate the mean, \bar{u}_{jkr} , standard deviation, S_{jkr} , and kurtosis coefficient, β_{2jkr} , where β_{2jkr} is given by:

$$\beta_{2jkr} = \frac{m_4}{(m_2)^2} \quad \text{with} \quad m_x = \frac{\sum_{i=1}^N (u_{ijk} - \bar{u}_{jkr})^x}{N} \quad (5)$$

For each observer, i , find P_i and Q_i , i.e.:

for $j, k, r = 1, 1, 1$ to J, K, R

if $2 \leq \beta_{2jkr} \leq 4$, then:

if $u_{ijk} \geq \bar{u}_{jkr} + 2 S_{jkr}$ then $P_i = P_i + 1$

if $u_{ijk} \leq \bar{u}_{jkr} - 2 S_{jkr}$ then $Q_i = Q_i + 1$

else:

if $u_{ijk} \geq \bar{u}_{jkr} + \sqrt{20} S_{jkr}$ then $P_i = P_i + 1$

if $u_{ijk} \leq \bar{u}_{jkr} - \sqrt{20} S_{jkr}$ then $Q_i = Q_i + 1$

If $\frac{P_i + Q_i}{J \cdot K \cdot R} > 0.05$ and $\left| \frac{P_i - Q_i}{P_i + Q_i} \right| < 0.3$ then reject observer i

with:

N : number of observers

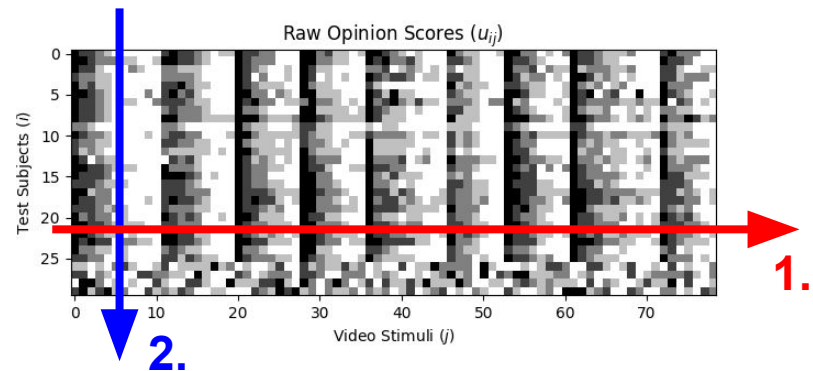
J : number of test conditions including the reference

K : number of test images or sequences

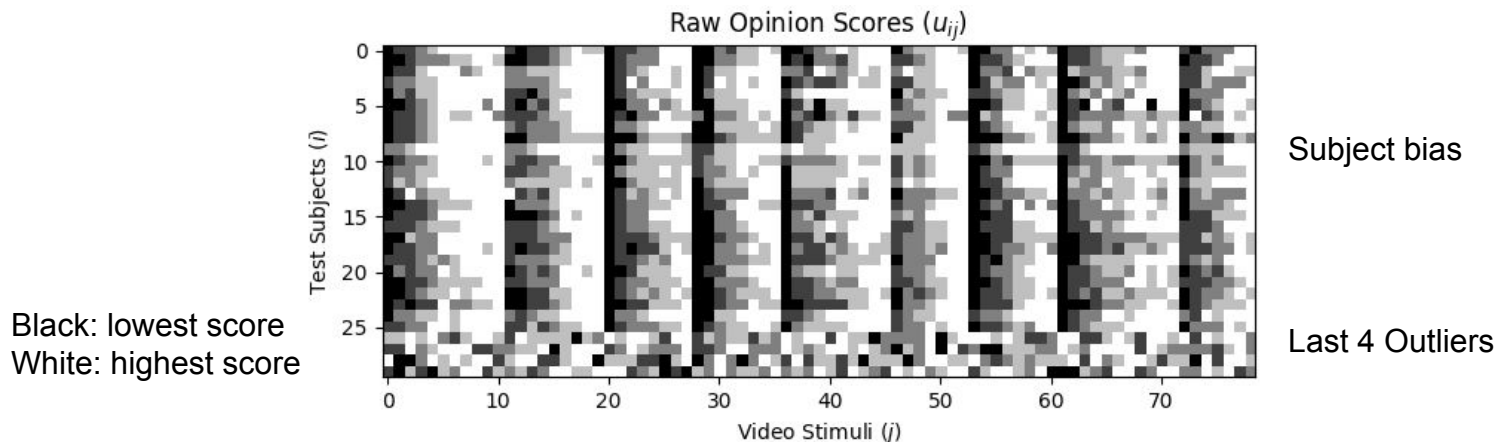
R : number of repetitions

L : number of test presentations (in most cases the number of presentations will be equal to $J \cdot K \cdot R$, however it is noted that some assessments may be conducted with unequal numbers of sequences for each test condition).

1. Video by video, the algorithm counts the number of instances when a subject's opinion score deviates by a few sigmas
2. Subject by subject, if the occurrences are more than a fraction, reject the subject



Limitations of BT.500-Style Subject Outlier Rejection



- All scores corresponding to rejected subjects are discarded - an overkill
- Often only identifies a subset of outliers
 - In the example above, only subjects #26, #28, #29 were rejected
- Hard-coded parameters and heuristic steps lack interpretability

Prior Art: Subject Bias Removal (ITU-T P.913)

First, estimate the MOS for each PVS:

$$\mu_{\psi_j} = \frac{1}{I_j} \sum_{i=1}^{I_j} o_{ij}$$

where:

o_{ij} is the observed rating for subject i and PVS j ;

I_j is the number of subjects that rated PVS j ;

μ_{ψ_j} estimates the MOS for PVS j , given the source stimuli and subjects in the experiment.

Second, estimate subject bias:

$$\mu_{\Delta_i} = \sum_{j=1}^{J_i} (o_{ij} - \mu_{\psi_j})$$

where:

μ_{Δ_i} estimates the overall shift between the i th subject's scores and the true values (i.e., opinion bias)

J_i is the number of PVSs rated by subject i .

Third, calculate the normalized ratings by removing subject bias from each rating:

$$r_{ij} = o_{ij} - \mu_{\Delta_i}$$

where:

r_{ij} is the normalized rating for subject i and PVS j .

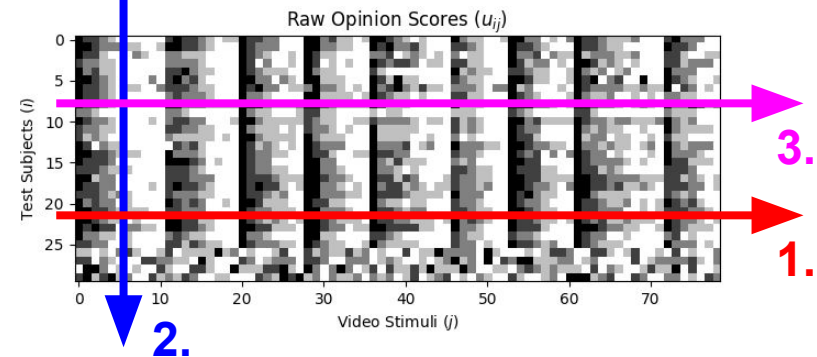
MOS and DMOS are then calculated normally. This normalization does not impact MOS:

$$\mu_{\psi_j} = \frac{1}{I_j} \sum_{i=1}^{I_j} r_{ij} = \frac{1}{I_j} \sum_{i=1}^{I_j} o_{ij}$$

where:

μ_{ψ_j} estimates the MOS of PVS j .

1. Video by video, estimate MOS by averaging over subjects
2. Subject by subject, estimate subject bias by comparing against MOS, and remove bias from opinion scores
3. Video by video, estimate MOS again based on bias-removed outlier-rejected opinion scores



Can we do better?

- Two most dominant effects of test subject inaccuracy:
 - Subject bias
 - Picky viewers tend to be biased toward lower scores
 - Not every subject has “golden eyes” - their sensitivity to impairment varies
 - Subject inconsistency
 - Subjects may not rate consistently throughout a session
 - Outliers - a special case with very large inconsistency
- Our proposal:
 - A simple yet effective model to account for subjective bias and inconsistency
 - Jointly solve the model parameters via maximum likelihood estimation (MLE)
 - Incorporate implicit “subject outlier rejection” and “subject bias removal” during model solving

A Simple Model*

$$\begin{array}{ccccccc} \text{Raw} & & \text{True} & & \text{Subject} & & \text{Subject} \\ \text{Opinion} & & \text{Score} & & \text{Bias} & & \text{Inconsistency} \\ \text{Score} & & & & & & \\ \color{red}{\boxed{U_{ijr}}} & = & \color{green}{\boxed{\psi_j}} & + & \color{orange}{\boxed{\Delta_i}} & + & \color{purple}{\boxed{v_i}} X \end{array}$$

- U_{ijr} - Opinion score of subject i , stimulus j and repetition r
- ψ_j - true quality of stimulus j
- Δ_i - bias of subject i
- v_i - inconsistency (std) of subject i
- X - i.i.d. normal random variables, $X \sim N(0, 1)$

*The model is a simplified version of [Li&Bampis'17] without considering the ambiguity of content. Compared to the original, the simplified model is more efficient and stable.

Solving the Model via MLE

- Given observations $\{U_{ijr}\}$
- The task is to solve for free parameters $\theta = (\{\psi_j\}, \{\Delta_i\}, \{v_i\})$
- Define log-likelihood function $l(\theta)$

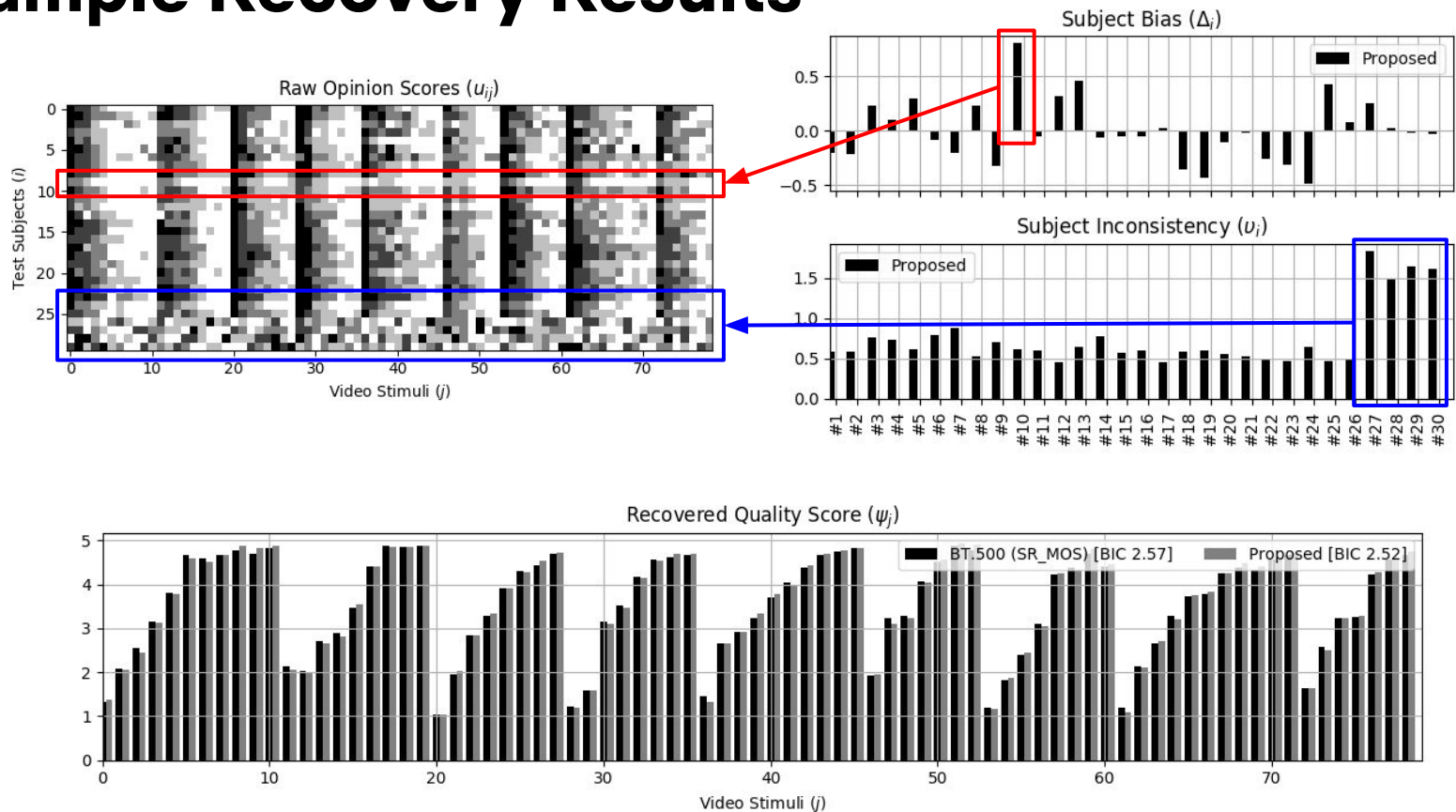
$$l(\theta) = \log P(U_{ijr} | \{\psi_j\}, \{\Delta_i\}, \{v_i\})$$

- Numerically solve to maximize the log-likelihood function

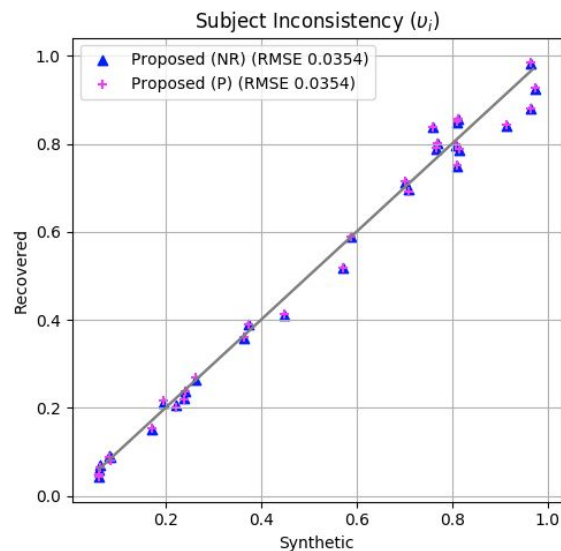
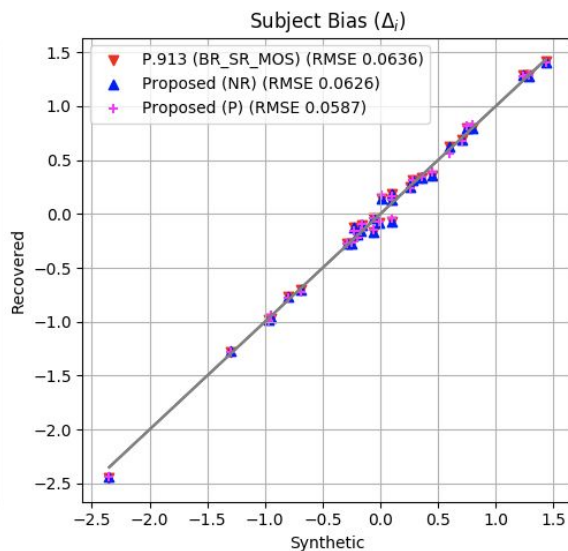
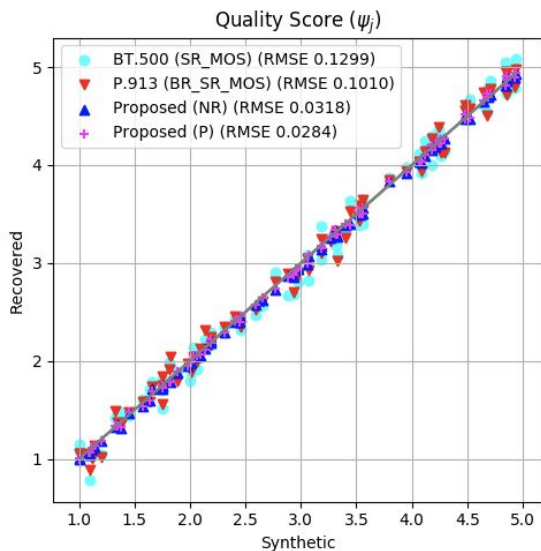
$$\hat{\theta} = \arg \max l(\theta)$$

- Proposed two numerical solutions
 - Newton-Raphson (NR) solution [Li&Bampis'17]
 - Projection-based (P) solution **NEW!** (thanks to Ioannis!)
 - Faster and strongly intuitive
 - Similar to ITU-T P.913, but 1) iterative 2) the projection is weighted by (sample count)/(residue variance)

Sample Recovery Results



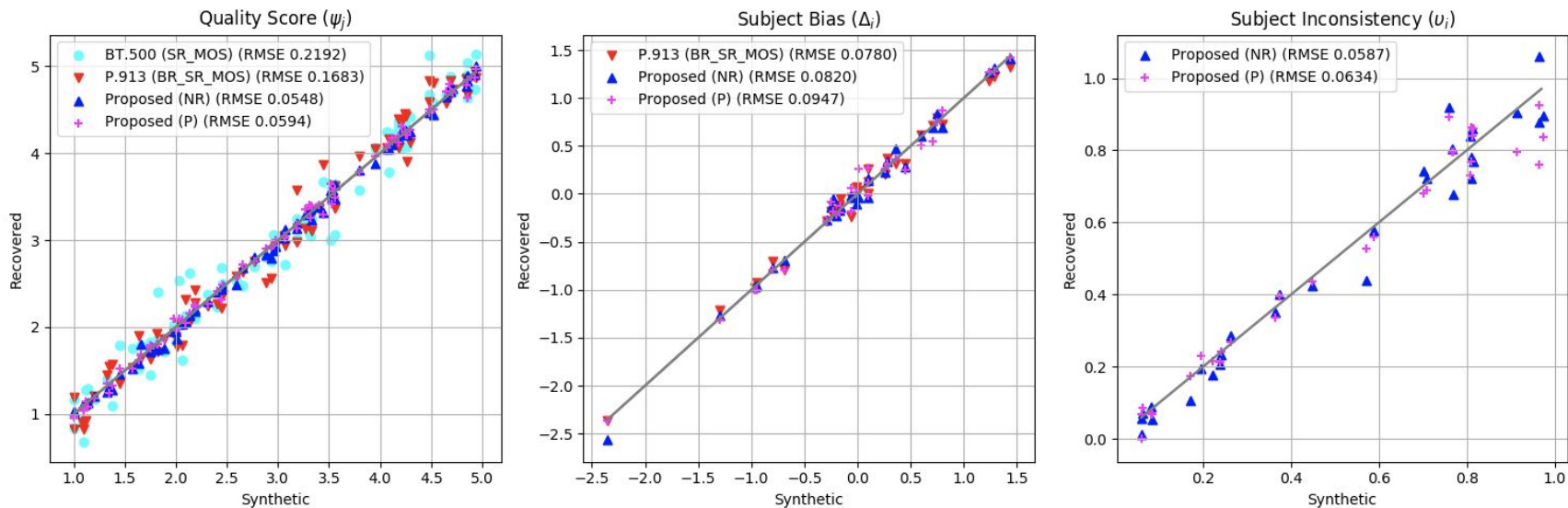
Validation Using Synthetic Data



- Synthetic data generation

- Randomly generate parameters according to $\psi_j \sim U[1, 5]$, $\Delta_i \sim N(0, 1)$, $v_i \sim U[0, 1]$
- Randomly generate observations according to parameters and model

Validation Using Synthetic Data (50% Missing Data)



- Synthetic data generation

- Randomly generate parameters according to $\psi_j \sim U[1, 5]$, $\Delta_i \sim N(0, 1)$, $\nu_i \sim U[0, 1]$
- Randomly generate observations according to parameters and model
- Data missing probability 0.5

SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error

Validation Using Bayesian Information Criterion

- BIC is a criterion for model fitting, balancing between:
 - The degree of freedom (number of parameters)
 - The goodness of fit (log-likelihood function)

$$\text{BIC} = \frac{\log(n) |\theta| - 2l(\theta)}{n}$$

- $|\theta|$ - the number of model parameters
- n - the number of observations (i.e. raw opinion scores)
- $l(\theta)$ - log-likelihood function

Bayesian Information Criterion (BIC)*

Dataset	MOS	BT.500 (SR_MOS)	P.913 (BR_SR_MOS)	Proposed (NR)	Proposed (P)
NFLX_dataset_public_raw_last4outliers	2.97	2.57	2.55	2.52	2.53
VQEGHD3_dataset_raw	2.75	2.74	2.39	2.30	2.31
HDTV_Phase_I_Experiment_1	2.45	2.46	2.38	2.20	2.22
HDTV_Phase_I_Experiment_2	2.72	2.72	2.52	2.32	2.33
HDTV_Phase_I_Experiment_3	2.72	2.71	2.37	2.29	2.29
HDTV_Phase_I_Experiment_4	2.96	2.96	2.51	2.27	2.27
HDTV_Phase_I_Experiment_5	2.77	2.77	2.47	2.33	2.33
HDTV_Phase_I_Experiment_6	2.51	2.49	2.32	2.16	2.16

*The model with the smallest BIC is preferred.

SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error

Bayesian Information Criterion (BIC)*

Dataset	MOS	BT.500 (SR_MOS)	P.913 (BR_SR_MOS)	Proposed (NR)	Proposed (P)
ITU-T_Supp_23_Experiment_1_BNR	2.91	2.91	2.35	2.31	2.31
MM2_1	2.80	2.78	2.83	2.74	2.75
MM2_2	3.89	3.89	3.52	3.13	3.13
MM2_3	2.48	2.47	2.45	2.41	2.42
MM2_4	2.74	2.73	2.62	2.47	2.47
MM2_5	2.90	2.82	2.67	2.64	2.64
MM2_6	2.81	2.74	2.74	2.72	2.74
MM2_7	2.73	2.72	2.76	2.67	2.71

*The model with the smallest BIC is preferred.

SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error

Bayesian Information Criterion (BIC)*

Dataset	MOS	BT.500 (SR_MOS)	P.913 (BR_SR_MOS)	Proposed (NR)	Proposed (P)
MM2_8	3.00	2.92	2.88	2.70	2.71
MM2_9	3.27	3.21	2.95	2.79	2.80
MM2_10	3.04	3.05	2.98	2.82	2.83
its4s2	3.63	3.63	2.96	2.59	2.59
its4s_AGH	3.14	3.04	2.76	2.63	2.63
its4s_NTIA	2.92	2.89	2.52	2.37	2.37

*The model with the smallest BIC is preferred.

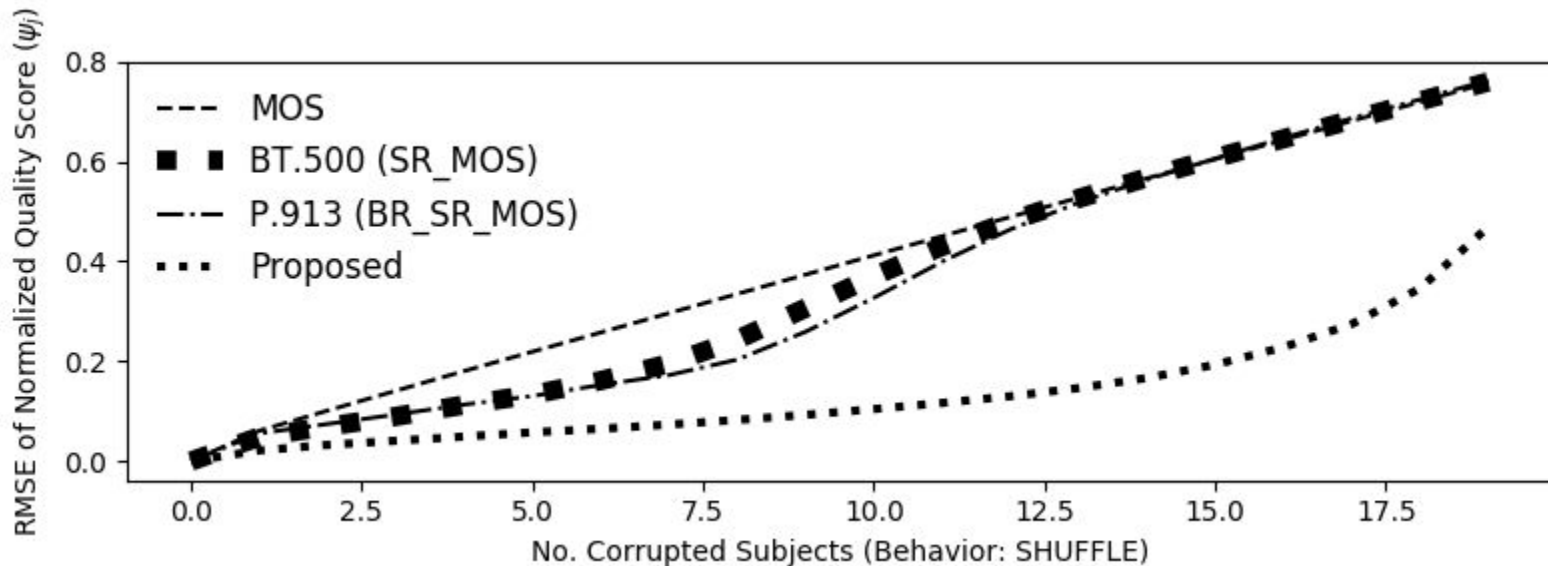
SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error

Robustness Against Subjects Giving Random Scores

Worse



Better



Random behavior: a subject's scores are shuffled among themselves

Y-axis: RMSE with respect to clean case

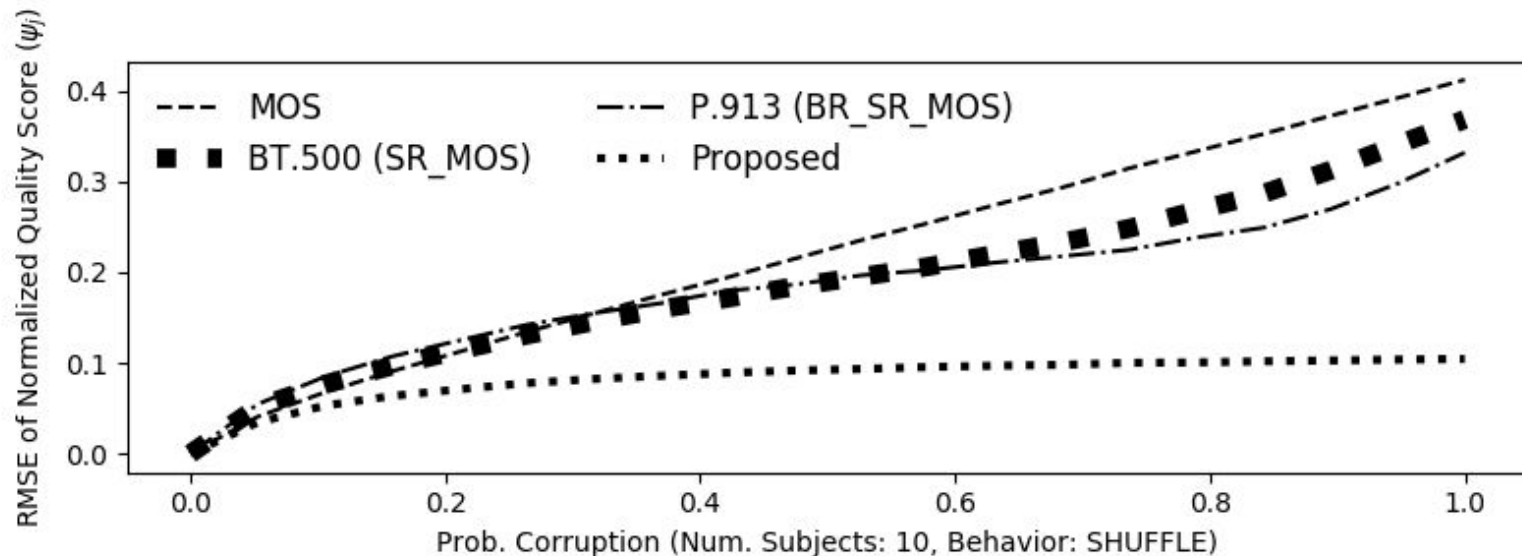
SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error

Robustness Against Increasing Corruption Probability

Worse



Better



10 random subjects are corrupted, with corruption probability varying from 0.0 to 1.0

Y-axis: RMSE with respect to clean case

SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error

Conclusions

- Recommendations such as ITU-R BT.500 and ITU-T P.913 standardize the procedure to clean up raw scores from subjective experiments through subject outlier rejection (SR) and subject bias removal (BR)
- We introduce a simple model and the corresponding parameter estimation procedure that implicitly takes into account both SR and BR, with the following advantages:
 - Better model fitting
 - Better robustness in the presence of outlier subjects in terms of recovery accuracy
 - Auxiliary information on test subjects on their bias and consistency, providing guides on subject selection
 - The projection-based solution provides strong intuition
- We propose to update ITU-R BT.500 and ITU-T P.910/P.913 with the new methodology



Backup Slides

The projection solution is faster and more intuitive

```
# video by video, estimate MOS by averaging over subjects
s_j = np.nanmean(x_ji, axis=1) # mean marginalized over i

# subject by subject, estimate subject bias by comparing with MOS
b_ji = x_ji - np.tile(s_j, (I, 1)).T
b_i = np.nanmean(b_ji, axis=0) # mean marginalized over j

MAX_ITR = 1000
DELTA_THR = 1e-8
EPSILON = 1e-8

itr = 0
while True:

    s_j_prev = s_j

    # calculate residue
    r_ji = x_ji - np.tile(s_j, (I, 1)).T - np.tile(b_i, (J, 1))

    # video by video, estimate MOS by averaging over subjects, inversely weighted by residue variance
    v_i = np.nanstd(r_ji, axis=0)
    s_ji = x_ji - np.tile(b_i, (J, 1))
    w_i = np.divide(cnt_i, v_i ** 2 + EPSILON)
    s_j = weighed_nanmean_2d(s_ji, weights=w_i, axis=1) # mean marginalized over i

    # subject by subject, estimate subject bias by comparing with MOS, inversely weighted by residue variance
    v_j = np.nanstd(r_ji, axis=1)
    b_ji = x_ji - np.tile(s_j, (I, 1)).T
    w_j = np.divide(cnt_j, v_j ** 2 + EPSILON)
    b_i = weighed_nanmean_2d(b_ji, weights=w_j, axis=0) # mean marginalized over j

    itr += 1

    delta_s_j = linalg.norm(s_j_prev - s_j)

    if delta_s_j < DELTA_THR:
        break

    if itr >= MAX_ITR:
        break
```

Same as ITU-T P.913

Similar to ITU-T P.913,
but weighted

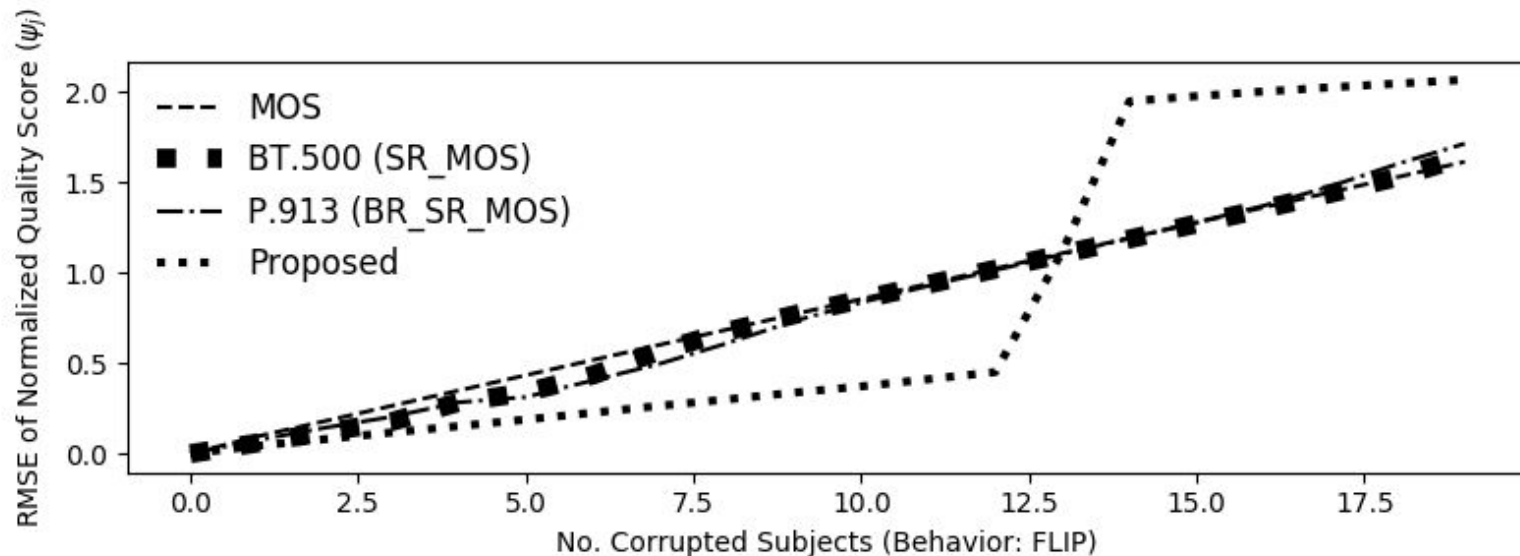
The weight is proportional
to the sample count, and
inversely proportional to the
residue variance

Robustness Against Subjects Giving “Flipped” Scores

Worse



Better



Malicious behavior: scores are “flipped”, for example, 1 for 5, 2 for 4, 2.5 for 3.5, and so on
Y-axis: RMSE with respect to clean case

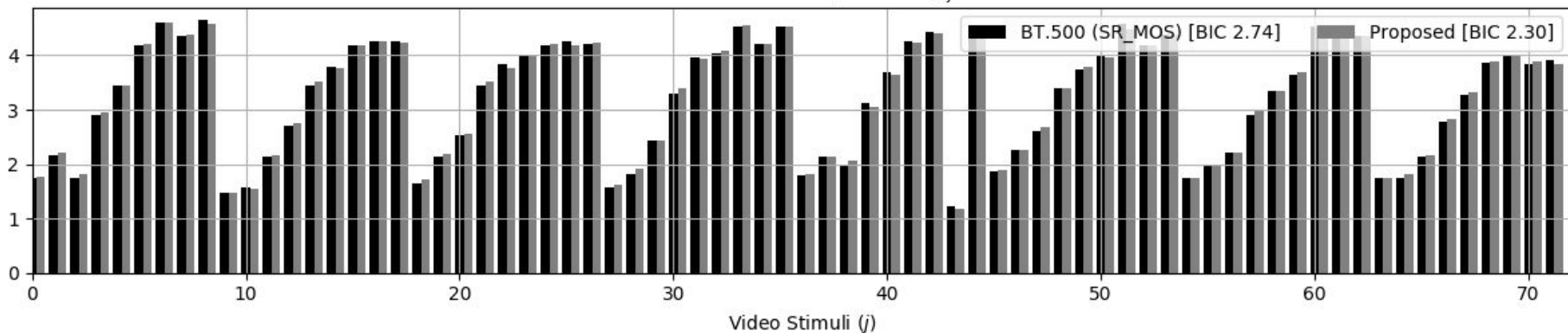
SR: subject rejection; BR: bias removal; MOS: mean opinion score; RMSE: root mean squared error



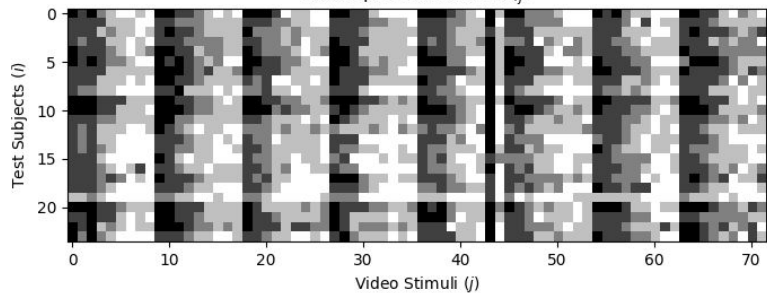
More Datasets

VQEGHD3_dataset_raw

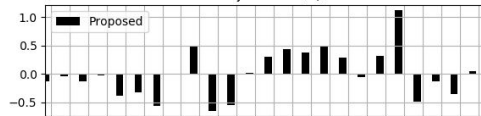
Recovered Quality Score (ψ_j)



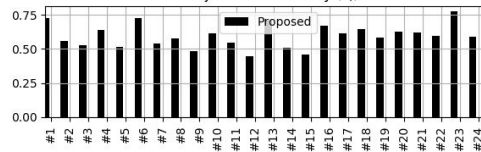
Raw Opinion Scores (u_{ij})



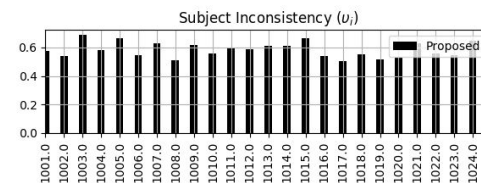
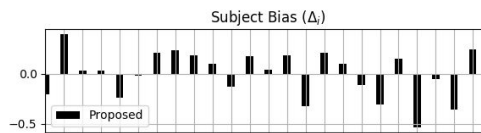
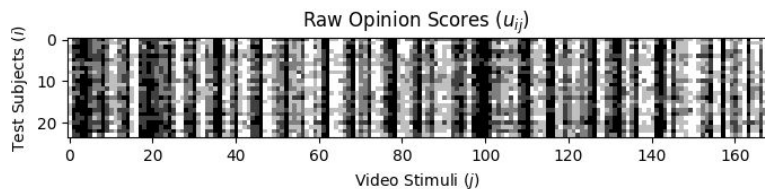
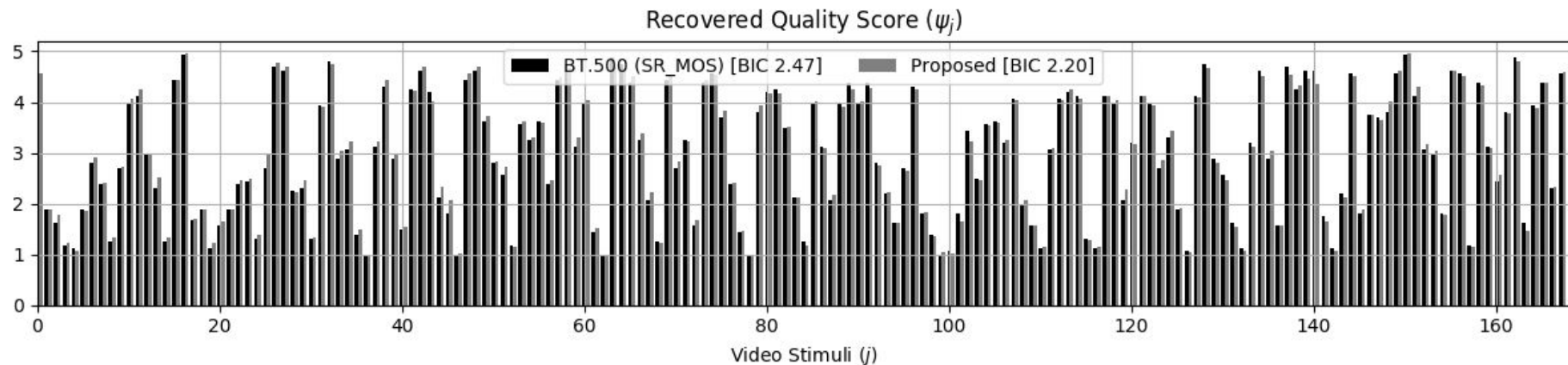
Subject Bias (Δ_i)



Subject Inconsistency (ν_i)

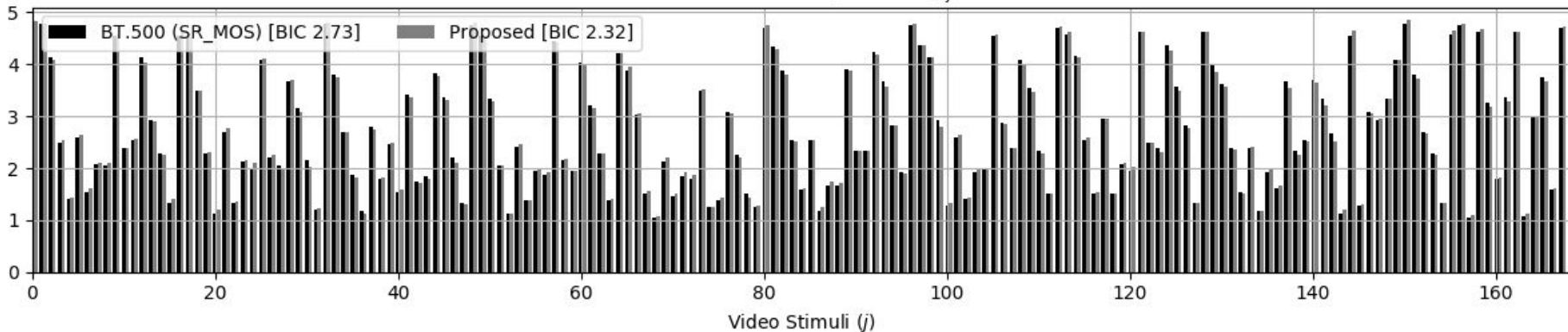


HDTV_Phase_I_Experiment_1

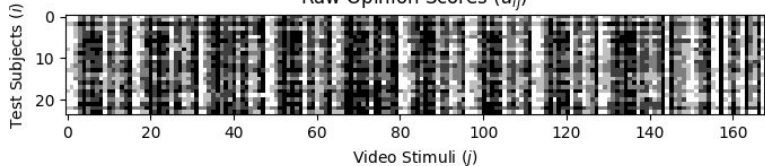


HDTV_Phase_I_Experiment_2

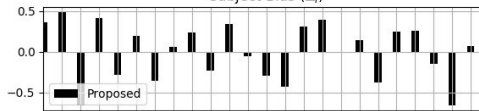
Recovered Quality Score (ψ_j)



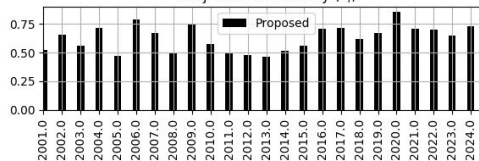
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

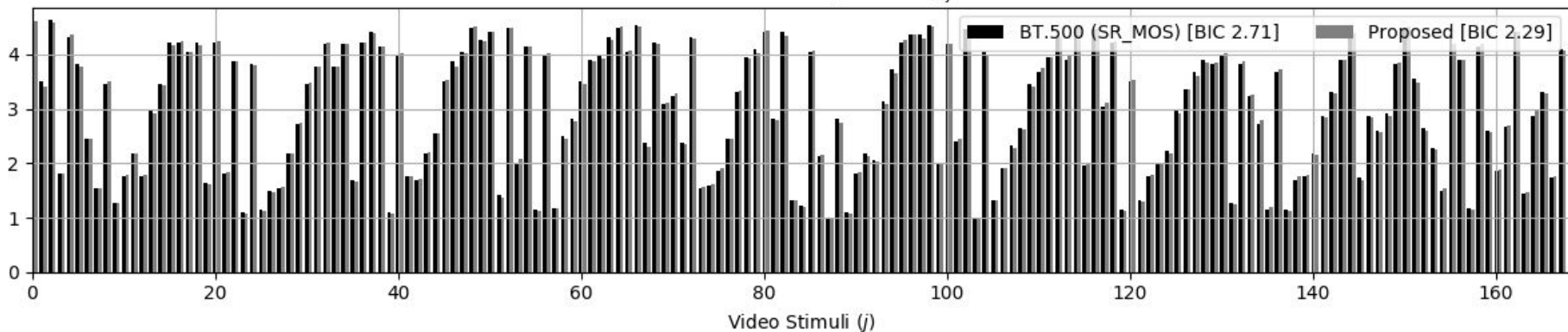


Subject Inconsistency (ν_i)

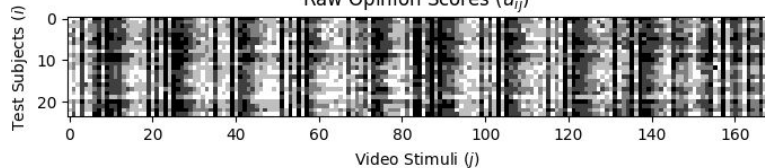


HDTV_Phase_I_Experiment_3

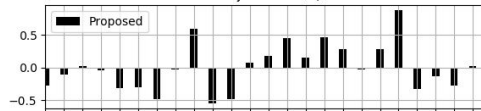
Recovered Quality Score (ψ_j)



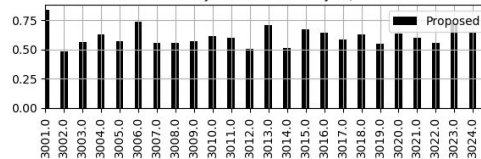
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

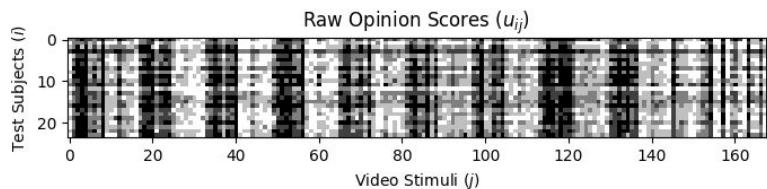
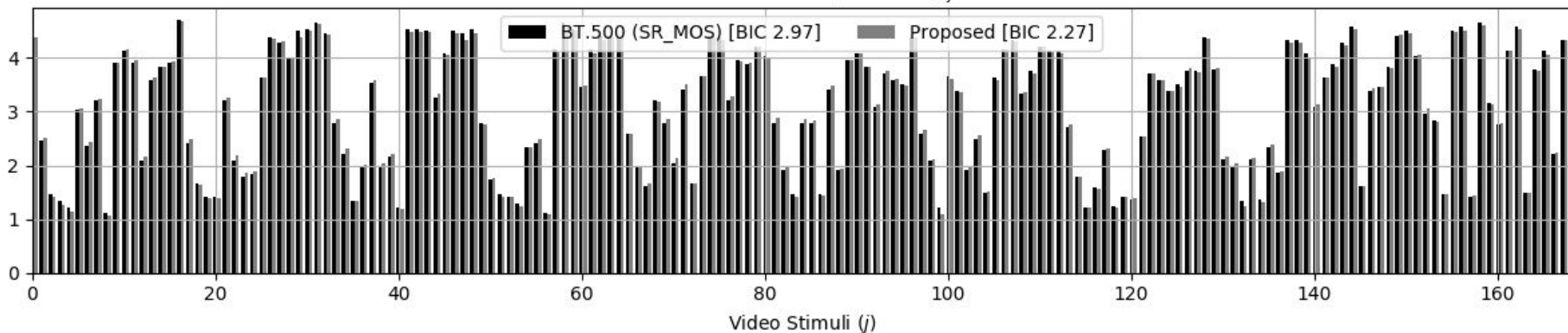


Subject Inconsistency (v_i)

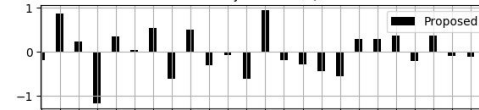


HDTV_Phase_I_Experiment_4

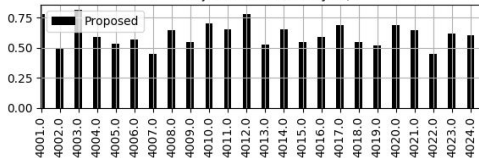
Recovered Quality Score (ψ_j)



Subject Bias (Δ_i)

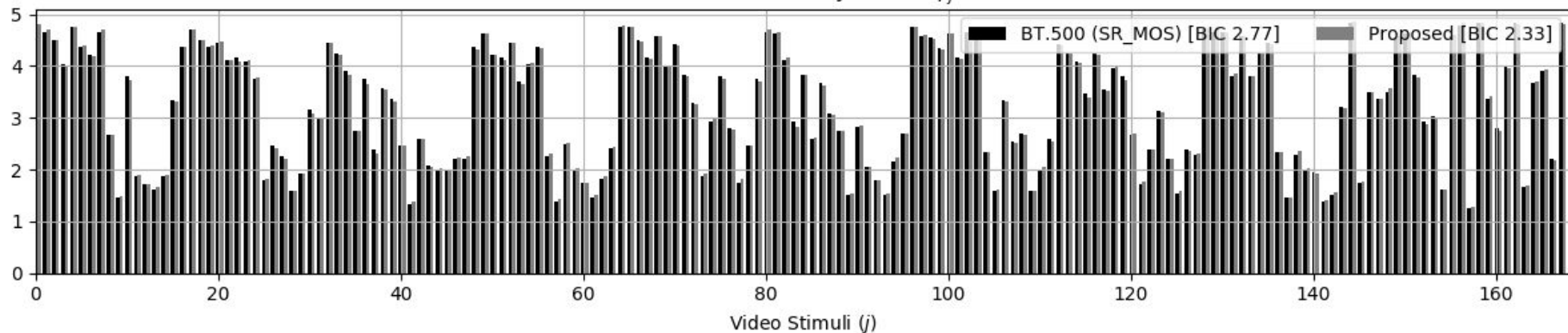


Subject Inconsistency (v_i)

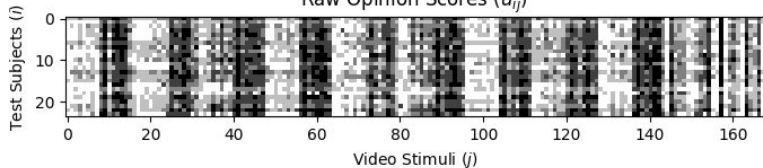


HDTV_Phase_I_Experiment_5

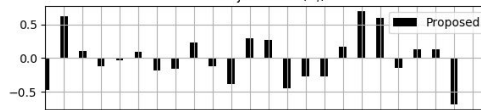
Recovered Quality Score (ψ_j)



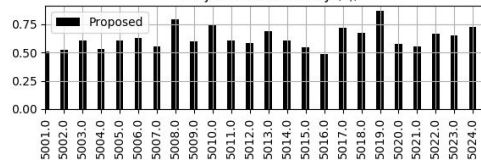
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

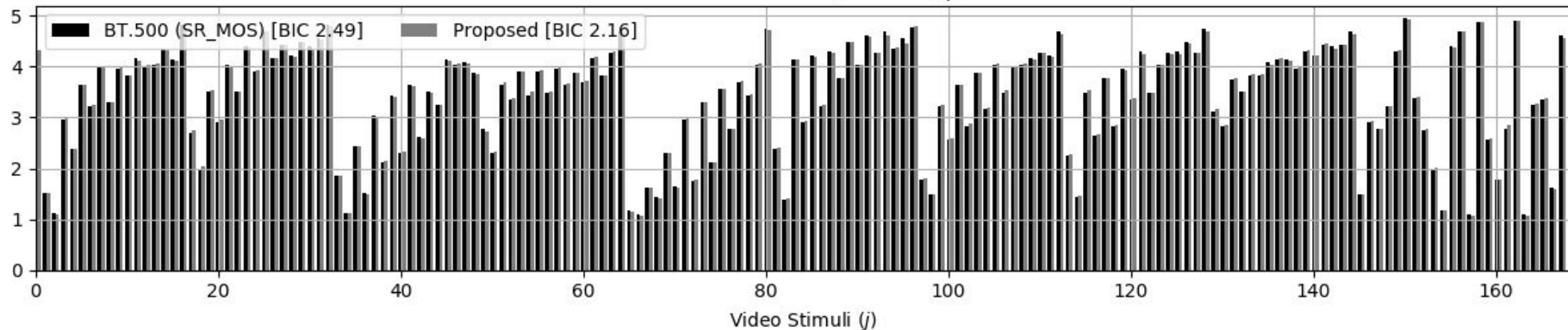


Subject Inconsistency (v_i)

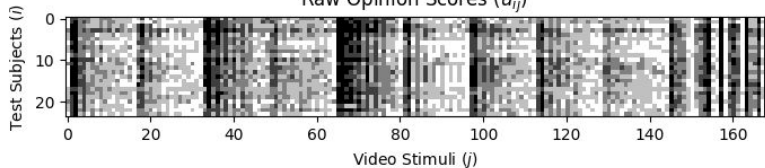


HDTV_Phase_I_Experiment_6

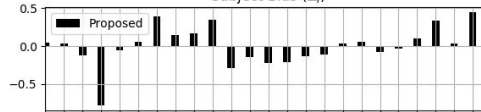
Recovered Quality Score (ψ_j)



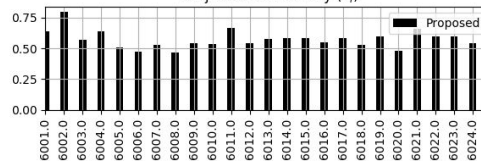
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

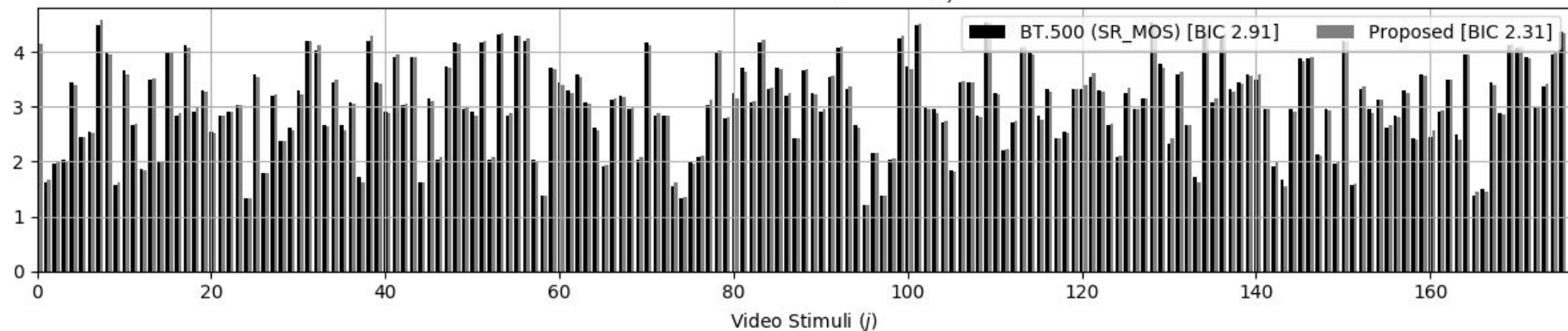


Subject Inconsistency (v_i)

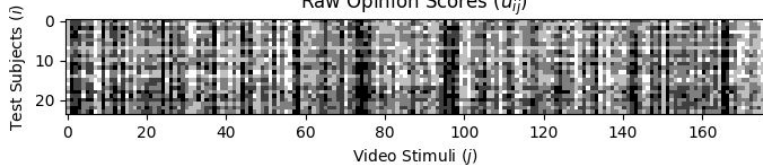


ITU-T_Supp_23_Experiment_1_BNR

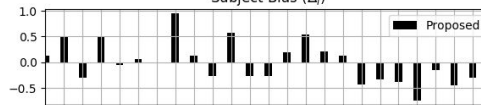
Recovered Quality Score (ψ_j)



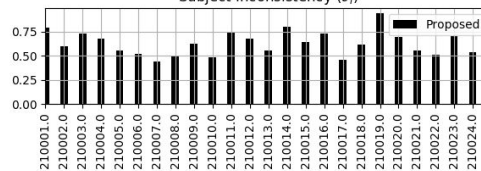
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

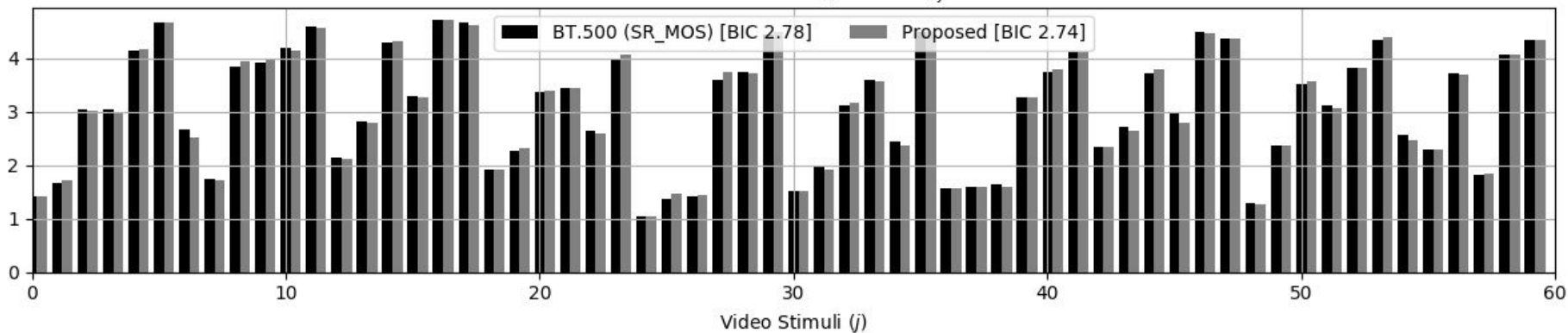


Subject Inconsistency (ν_i)

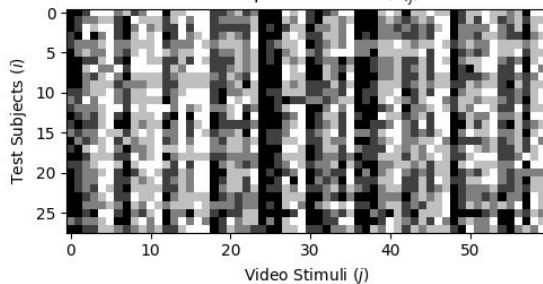


MM2_1

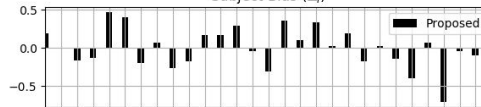
Recovered Quality Score (ψ_j)



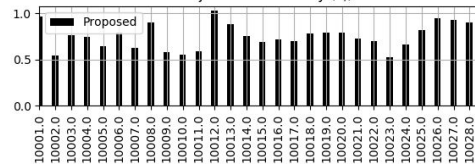
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

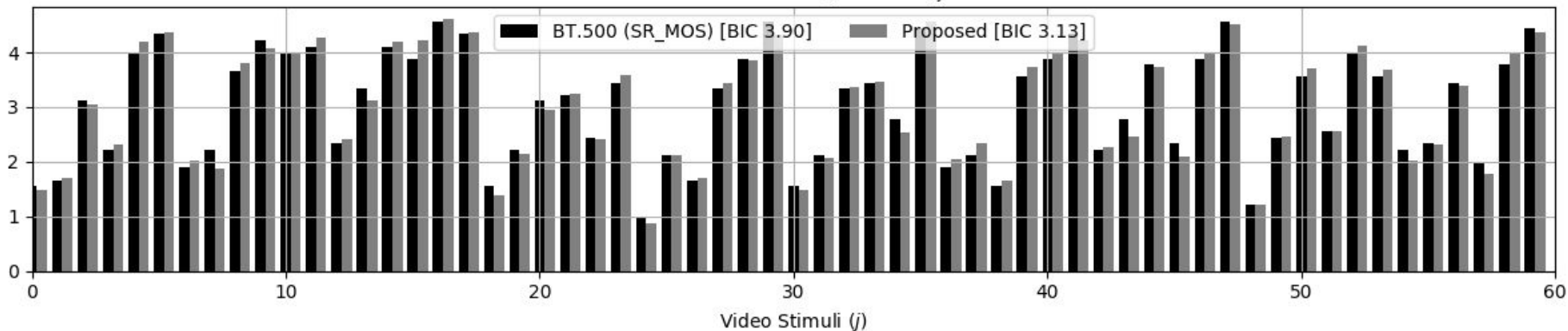


Subject Inconsistency (ν_i)

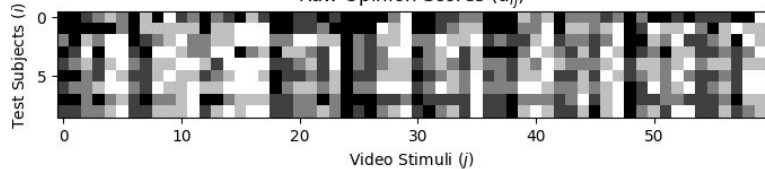


MM2_2

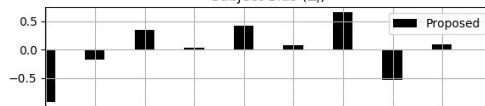
Recovered Quality Score (ψ_j)



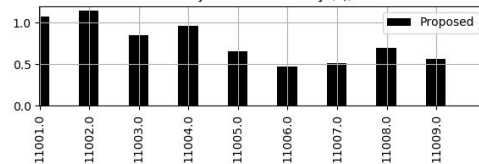
Raw Opinion Scores (u_{ij})



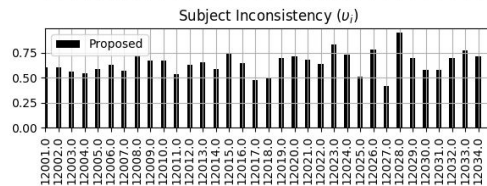
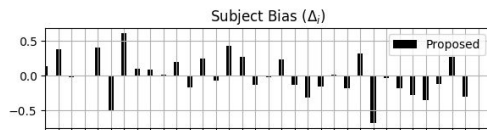
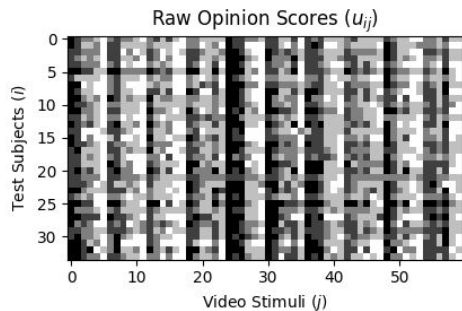
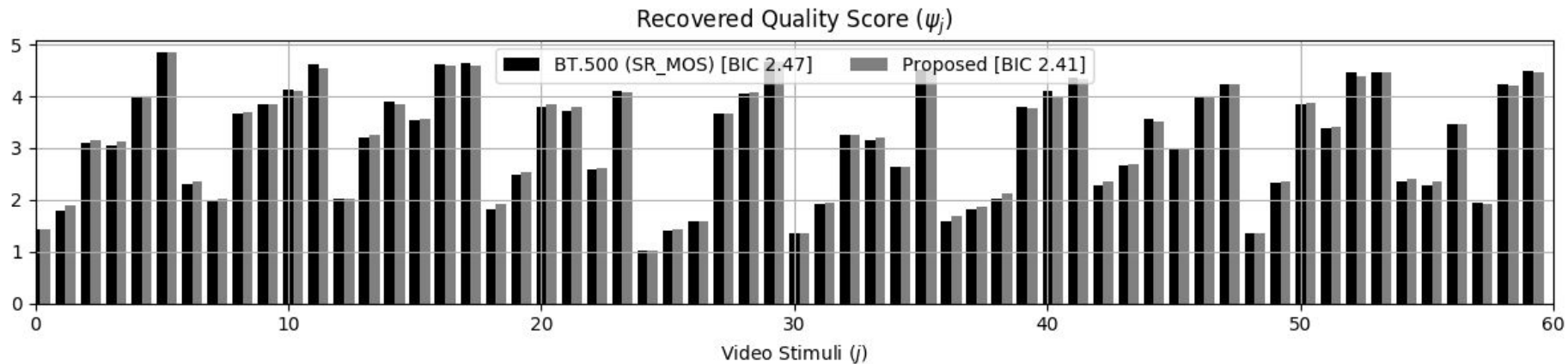
Subject Bias (Δ_i)



Subject Inconsistency (θ_i)

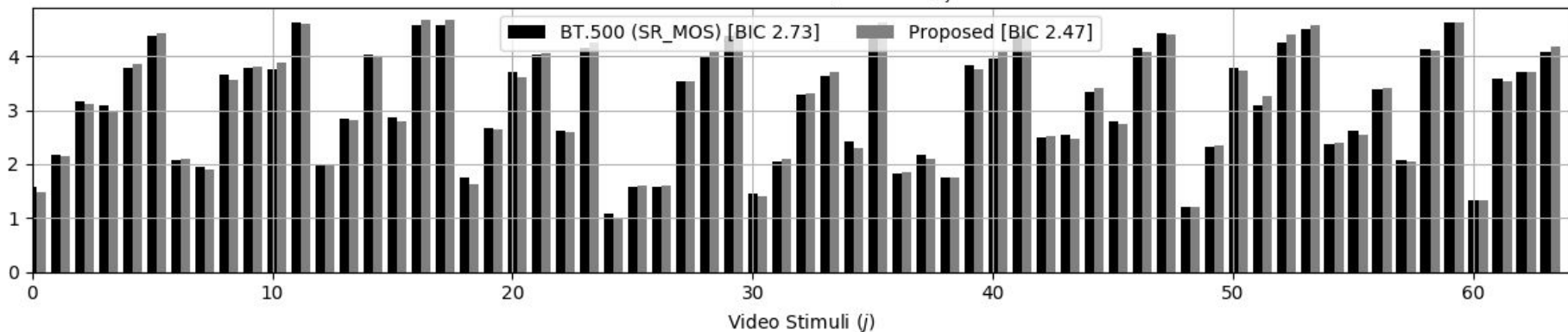


MM2_3

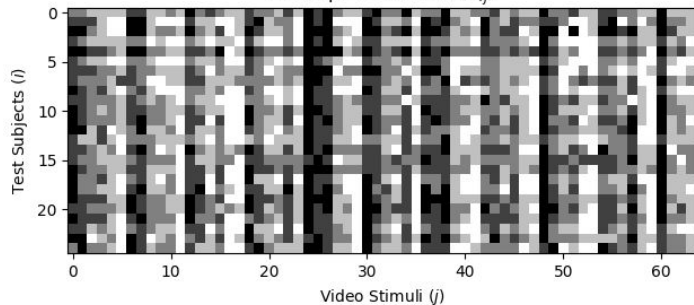


MM2_4

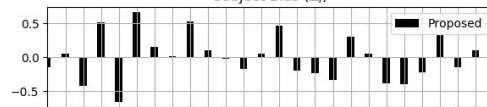
Recovered Quality Score (ψ_j)



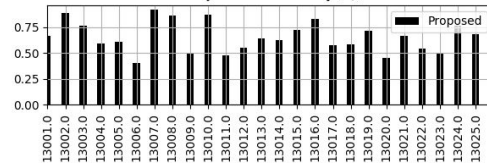
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

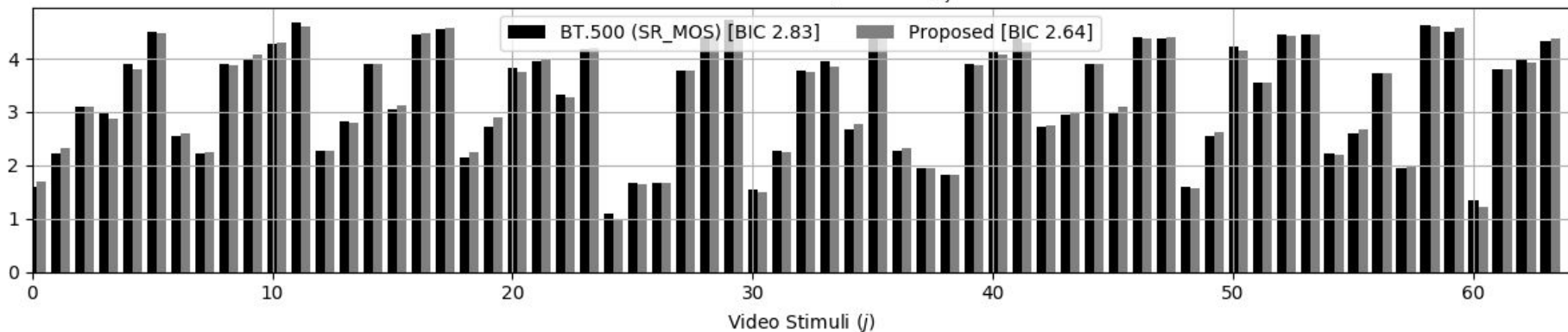


Subject Inconsistency (θ_i)

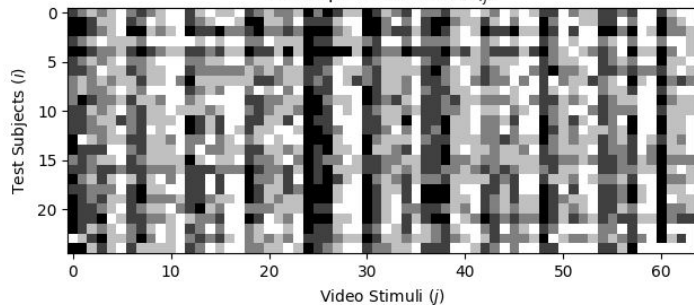


MM2_5

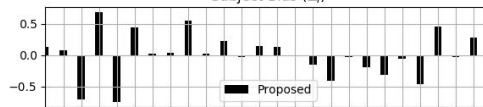
Recovered Quality Score (ψ_j)



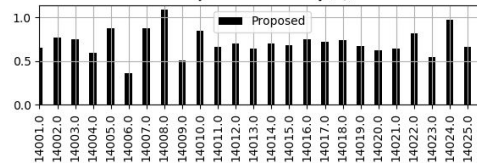
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

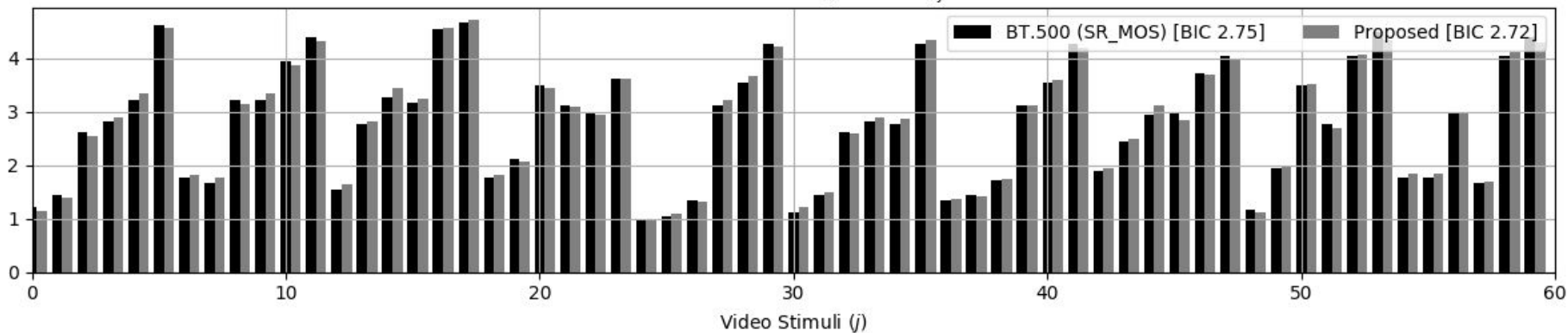


Subject Inconsistency (θ_i)

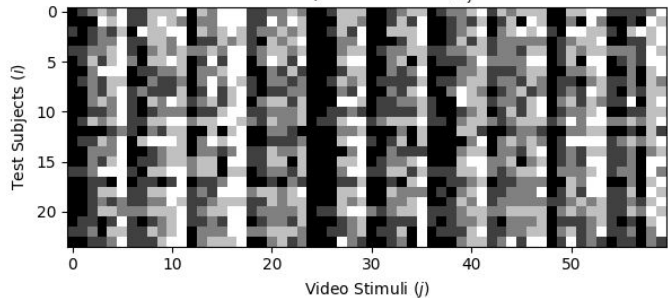


MM2_6

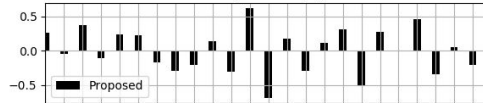
Recovered Quality Score (ψ_j)



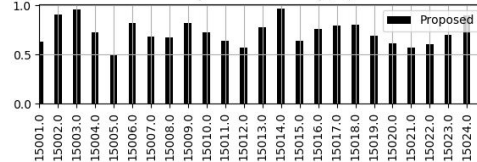
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

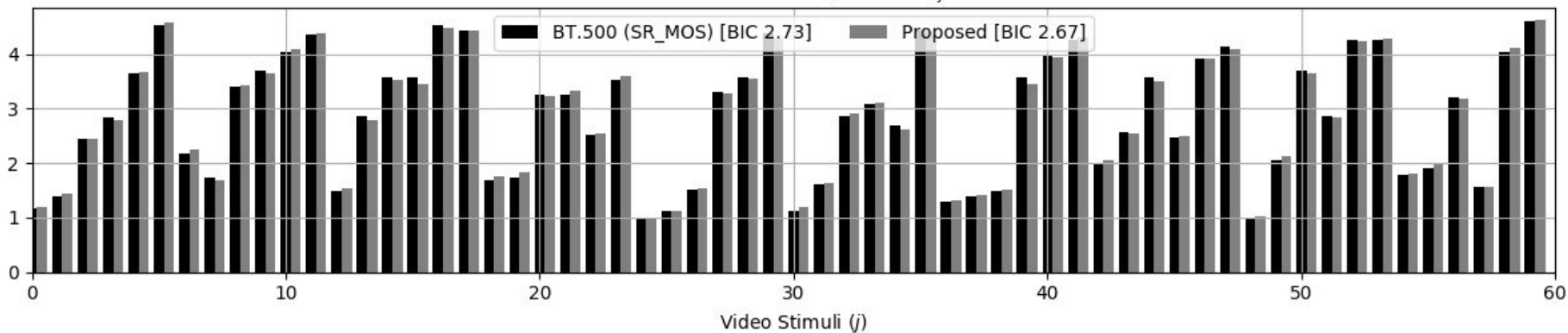


Subject Inconsistency (ν_i)

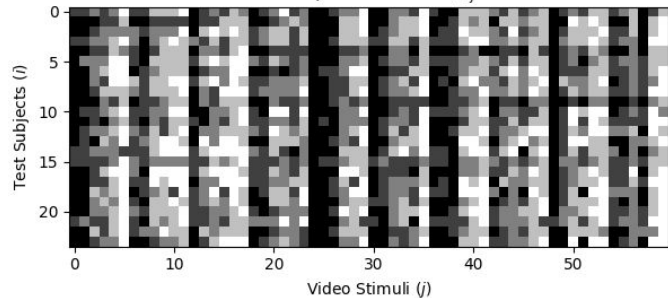


MM2_7

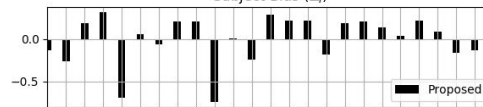
Recovered Quality Score (ψ_j)



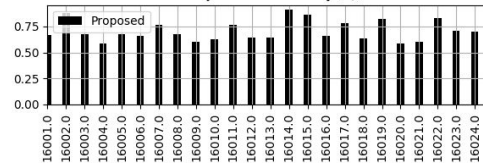
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

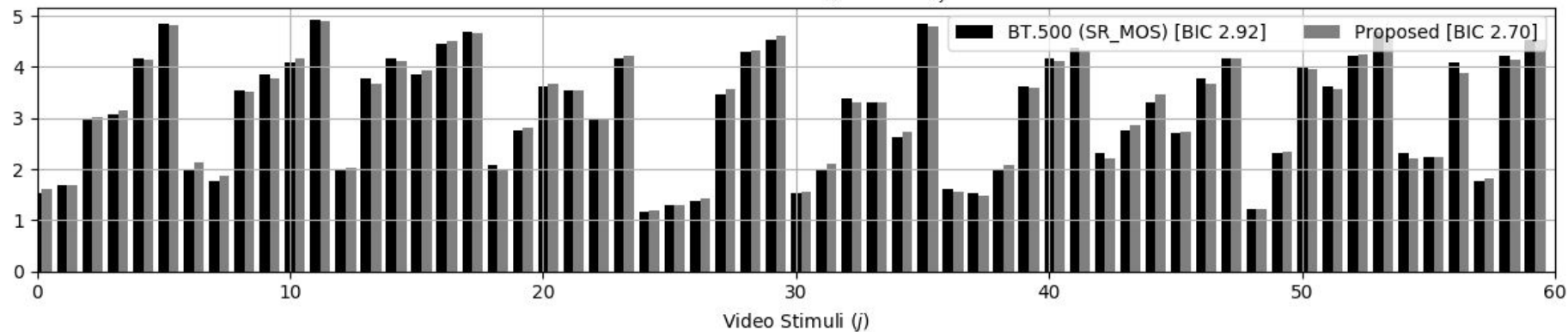


Subject Inconsistency (ν_i)

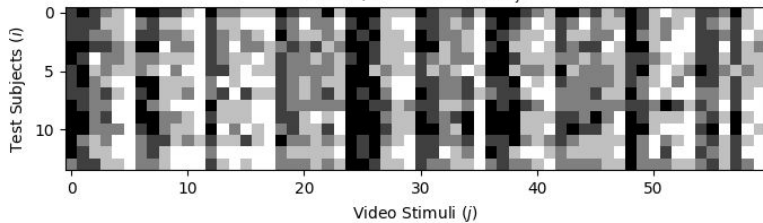


MM2_8

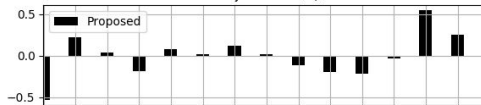
Recovered Quality Score (ψ_j)



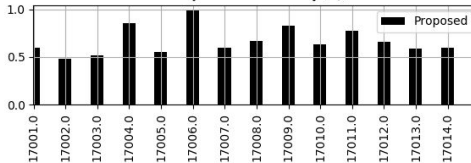
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

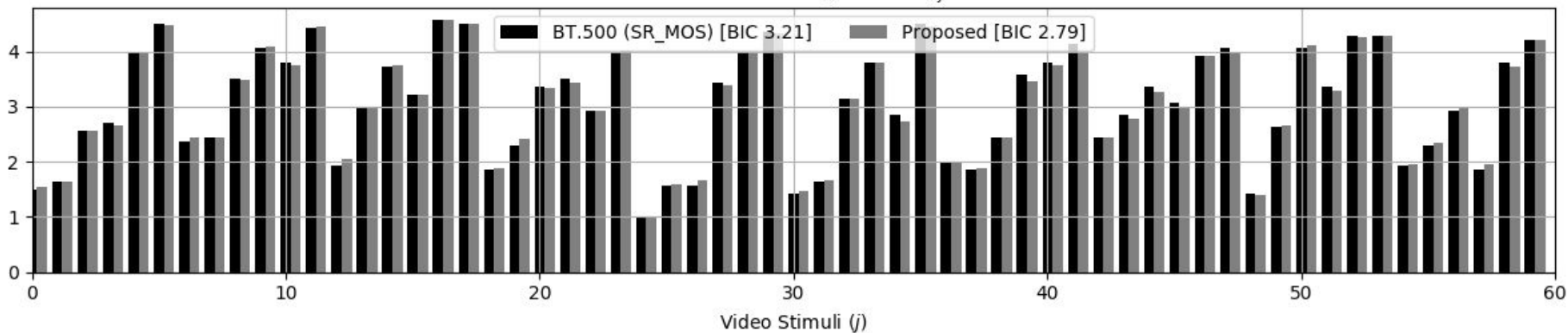


Subject Inconsistency (ν_i)

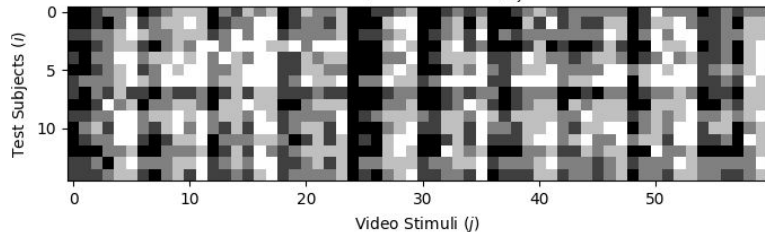


MM2_9

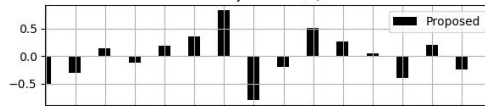
Recovered Quality Score (ψ_j)



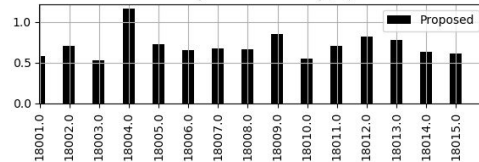
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

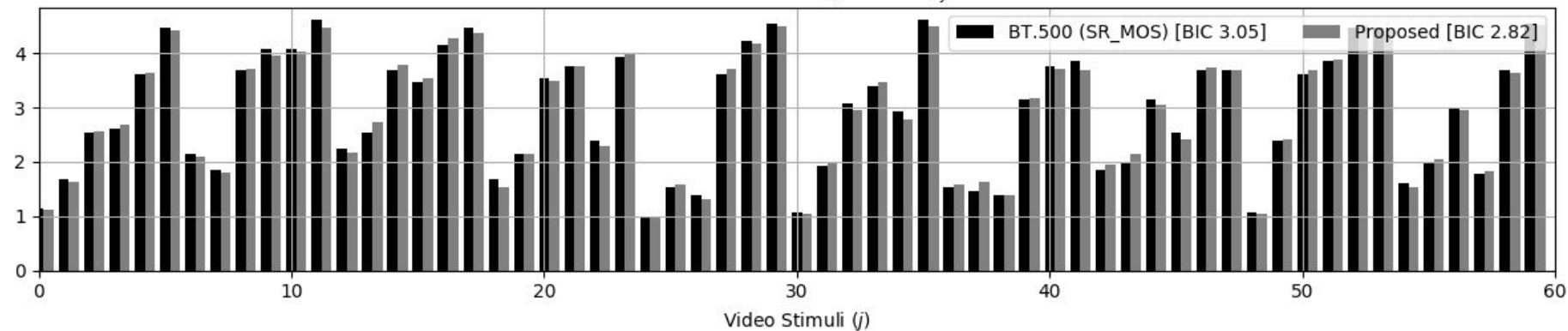


Subject Inconsistency (ν_i)

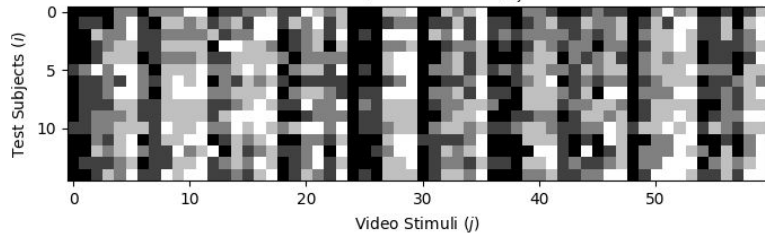


MM2_10

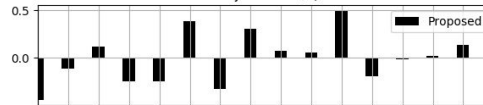
Recovered Quality Score (ψ_j)



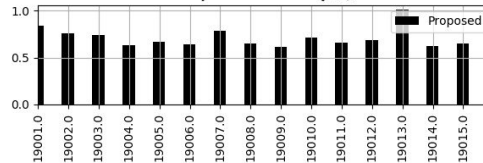
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

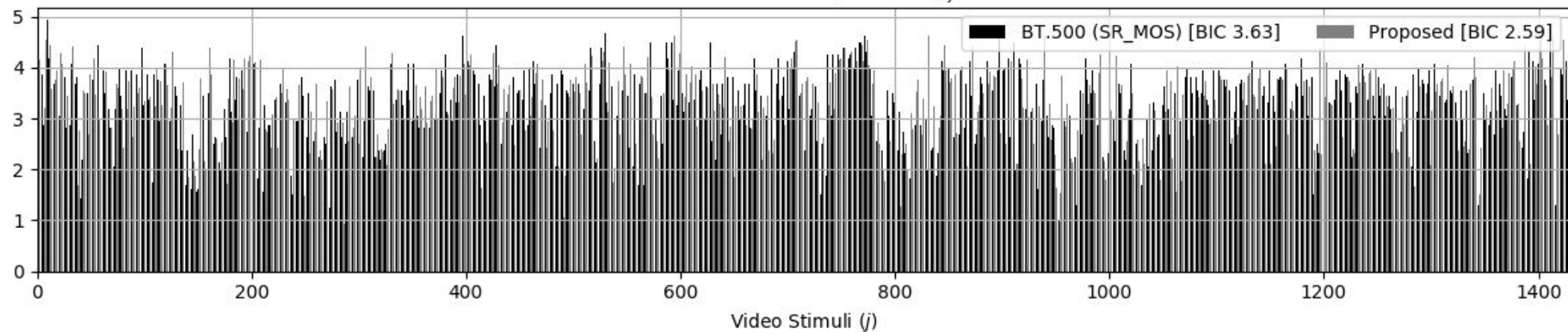


Subject Inconsistency (v_i)



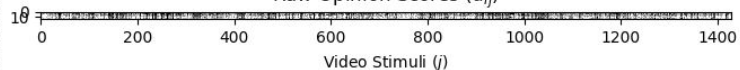
its4s2

Recovered Quality Score (ψ_j)

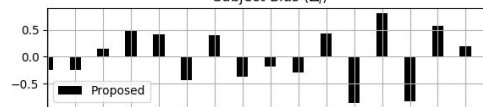


Test Subjects (i)

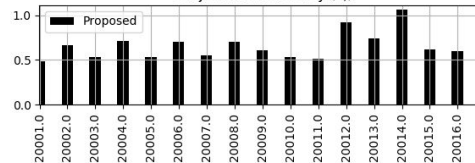
Raw Opinion Scores (u_{ij})



Subject Bias (Δ_i)

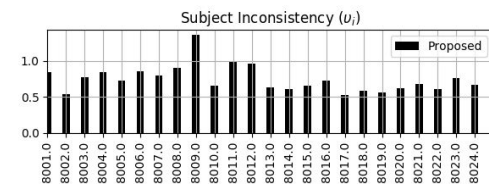
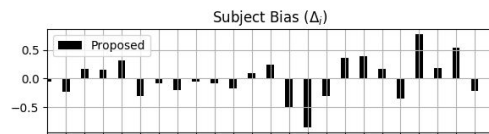
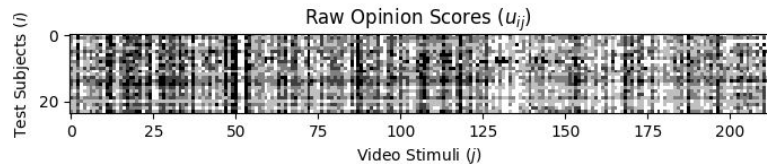
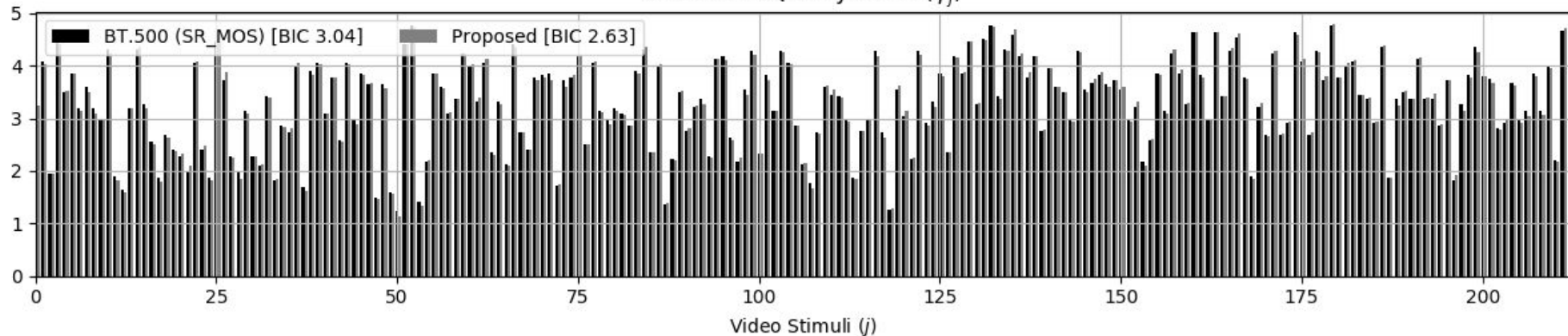


Subject Inconsistency (ν_i)



its4s_AGH

Recovered Quality Score (ψ_j)



its4s_NTIA

Recovered Quality Score (ψ_j)

