# Statistical Testing of Overlap-Based Performance Between an AI Segmentation Device and a Multi-Expert Human Panel Without Requiring a Reference Standard

**Tingting Hu,[a, *] Berkman Sahiner,[a] Shuyue Guan,[a] Mike Mikailov,[a] Kenny Cha,[a] Frank Samuelson,[a] Nicholas Petrick[a]**

[a]U.S. Food and Drug Administration, Silver Spring, Maryland, United States

## Abstract

Purpose: AI-based medical imaging devices often include lesion or organ segmentation capabilities. Existing methods for segmentation performance evaluation compare AI results to an aggregated reference standard using accuracy metrics like the Dice coefficient or Hausdorff Distance. However, these approaches are limited for lacking a gold standard and challenges in defining meaningful success criteria. To address this, we developed a statistical method to assess agreement between an AI device and multiple human experts without requiring a reference standard.

Approach: We propose a paired-testing method evaluating whether an AI device's segmentation performance significantly differs from multiple human experts'. The method compares device-to-expert dissimilarity with expert-to-expert dissimilarity, avoiding the need for a reference standard. We validated the method through: (1) statistical simulations where Dice coefficient performance is either shared ("overlap agreeable") or not shared ("overlap disagreeable") between the device and experts; (2) image-based simulations using 2D contours with shared or non-shared transformation parameters ("transformation agreeable or disagreeable"). We also applied the method to compare an AI segmentation algorithm to four radiologists using data from the Lung Image Database Consortium.

Results: Statistical simulations show the method controls type I error (~0.05) for overlap-agreeable and type II error (~0) for overlap-disagreeable scenarios. Image-based simulations show acceptable performance with mean type I error 0.07 (SD 0.03) for transformation-agreeable and mean type II error 0.07 (SD 0.18) for transformation-disagreeable cases.

Conclusions: The paired-testing method offers a new tool for assessing the agreement between an AI segmentation device and multiple human expert panelists without requiring a reference standard.

**Keywords**: segmentation assessment, multi-expert human panel, paired testing

**\*First Author**, E-mail: Tingting.Hu@fda.hhs.gov

## 1   Introduction

The segmentation of structures and lesions within medical images is an increasing focus of artificial intelligence (AI) based medical imaging devices. Examples include devices used for delineating lesions in disease detection and diagnosis, or for outlining organs in surgical planning or radiation therapy. The evaluation of segmentation algorithm performance is critical in the medical device domain to ensure patient safety and benefit, however, there is a relatively small amount of literature addressing the evaluation of segmentation algorithm performance even though there is a significant amount of research focused on AI segmentation development. Furthermore, existing literature on segmentation performance assessment primarily focuses on reviewing and comparing performance evaluation metrics (e.g., [1, 2]) or proposing new evaluation measures (e.g., [3, 4]). However, these proposed metrics and measures cannot be directly used for evaluating segmentation performance in the absence of a reference standard establishing the ground truth for segmentation task assessment. In many medical imaging tasks in the literature, reference contours for each object to be segmented are often defined by multiple human experts ('expert' may be referred to in the literature or submissions by various terms, such as observer, reader, reviewer, or truther). While using only one expert's reference contour to evaluate device performance simplifies the analysis, this approach fails to reflect the truth variability that exists even among high-level experts. Therefore, including multiple experts' reference contours better reflects the true nature of the problem, though this introduces complexity in evaluation and analysis. This complexity necessitates robust statistical methods to appropriately handle comparison of device segmentation against multi-expert references.

   A commonly used assessment approach in practice is, with a reference standard contour defined, the AI segmentation output is then compared, through an overlap metric, e.g., Dice [5],

67  Jaccard [6], or a distance-based metric, e.g., Hausdorff distance [1] with the reference. A limitation

68  of this approach is it requires prespecifying a meaningful cutoff for each metric. However,

69  selecting and justifying a clinically meaningful performance criterion can be difficult.

70     These two challenges with current segmentation assessment methods — the absence of a

71  definitive reference standard contour and the difficulty of establishing a clinically meaningful

72  performance goal — led us to explore a new statistical method that assesses interchangeability,

73  without requiring a reference standard, between an AI segmentation and segmentations from a

74  human expert panel. With the growing number of AI/ML segmentation devices, seeking for a

75  more generalizable segmentation assessment approach has become a pressing need in the

76  medical technology field.

77     The only relevant work we found is by Zou et al.[7], which addresses a similar problem of

78  evaluating a segmentation algorithm against multiple truthers (Example 2 in their paper).

79  However, their method aggregates the three human annotations into a single STAPLE-based

80  reference, then compares the algorithm to this composite. This would still require defining an

81  arbitrary success cutoff and leave the fundamental challenges unresolved. To address this gap,

82  we examined methodologies for assessing continuous estimation tasks (i.e., estimation tasks for

83  which the target quantity has a continuous value, such as area or volume measurements for

84  organs or lesions), where agreement measures and methods have been more widely discussed.

85  These include the Bland-Altman method [8], individual bioequivalence [9], the individual

86  equivalence index [10], [11], the agreement index [12], and the individual equivalence coefficient and

87  coefficient of individual agreement [13].

88     To address segmentation agreement, we adapted the interchangeability method proposed by

89  Obuchowski et al. [11] for numerical estimation to the segmentation context. The basic idea is to

90  condense reader-reader pairwise segmentation comparisons into a within-panel agreement score

91  for each image and, likewise, generate a corresponding device-panel agreement score. This

92  reformulates the segmentation comparison as a paired numerical score comparison. By tailoring

93  Obuchowski et al.'s method, we constructed a paired test statistic and apply the resulting

94  confidence interval to determine whether device-panel agreement significantly differs from

95  within-panel agreement.

96  In the remainder of the paper, we define the problem mathematically and present the

97  proposed methodology. We next present two simulation studies: one statistic-based and the other

98  image-based, and report the Type I and Type II errors for our proposed method.  Finally, we

99  conclude with a discussion of the findings and potential directions for future research.

## 2    Methodology

101  This section outlines the problem definition and proposed methodology. For illustration

102  simplicity, we focus on a single object to segment throughout this article. But the proposed

103  method can be applied to multi-object scenarios in a similar way as defined here.

### 2.1  Problem Formulation

105  Consider a testing dataset containing $n$ images, where each image is obtained from an

106  independent patient and contains a single object to segment.  On this dataset, an AI device,

107  denoted as $D$, segments the object of interest within each image. Concurrently, a human expert

108  panel denoted as $P$, comprised of $k$ experts, each independently performs manual annotation of

109  the object. Each segmentation, whether by the device or an expert, is represented as a binary

110  image.

111      The problem of interest is testing whether the segmentation performance of device $D$

112    significantly diverges from that of the human expert panel. Symbolically, the aim is to test

113    whether $\mu_D$ - $\mu_P$ significantly differs from zero, where $\mu_D$ represents the mean (dis)similarity

114    score between the device and panel across images, and $\mu_P$ denotes the mean (dis)similarity score

115    between experts within the panel across images, i.e., testing the null hypothesis

116    $$H_0: \mu_D - \mu_P = 0,$$

117    against the alternative hypothesis

118    $$Ha: \mu_D - \mu_P \neq 0.$$

119    Here, we adopt an equality-based null hypothesis formulation for simplicity, following the

120    fashion of hypothesis formulation in the FDA guidance [14] where the null hypothesis is set as no

121    treatment effect on the selected endpoint. This choice ensures alignment with the conventional

122    definitions of Type I and Type II errors, which will be utilized for method validation later in

123    Section 3.


124    *2.2  Related Work for Numeric Output*

125    Obuchowski et al. [11] proposed a metric called individual equivalence index to measure the

126    individual equivalence of imaging tests when the health outcome of interest is a numeric

127    variable. This metric is defined as below:

128    $$\gamma = E\left(Y_{jTik} - Y_{jRik\prime}\right)^2 - E\left(Y_{jRik} - Y_{jRik\prime}\right)^2 \quad (1)$$

129    where $Y_{jTik}$ denotes the result or measurement by the new test modality (T) by reader i for

130    subject j on occasion k, and $Y_{jRik}$ denotes the result or measurement by the existing reference

131    modality (R) by reader i for subject j on occasion k.

132    While the Obuchowski et al. [11] method was developed for an estimation task with a numeric

133    output, we adapt it to quantify segmentation agreement through suitable modifications.


134    *2.3 Proposed Method*


135    *2.3.1 Segmentation Interchangeability Metric*

136    From Eq. (1), it is evident that the central concept behind the interchangeability metric is to

137    compare the dissimilarity between the new test and reference test with the dissimilarity within

138    the reference test. Building on this idea, a natural extension of this approach to segmentation

139    outputs is replacing the dissimilarity metric for numeric outputs adopted by Obuchowski et al. [11]

140    (i.e., the mean squared difference) with an appropriate segmentation dissimilarity metric. One of

141    the most widely used similarity measures for segmentation is the Dice Similarity Coefficient

142    (DSC) [5]. We therefore use 1-DSC as the dissimilarity surrogate to tailor the original individual

143    equivalence index to segmentation. Based on this modification, we propose a segmentation

144    interchangeability metric denoted by δ, to evaluate segmentation agreement between an AI

145    device and a panel of human readers. The proposed metric is defined as follows:

146    $$\delta = E\{1 - DSC(device, reader\ panel)\} - E\{1 - DSC(within\ reader\ panel)\} \qquad (2)$$
147
148    where E denotes the expected value. Clearly, the closer $\delta$ is to zero, the more similar the device's

149    segmentation performance is to the human reader panel. The Dice coefficient for paired

150    segmentations on a single image is always positive and ranges between zero and one. These

151    properties make this segmentation interchangeability metric ($\delta$) well-suited for evaluating

152    segmentation performance at the individual image level.

153 *2.3.2 Point Estimate*

154 The point estimator for the proposed interchangeability metric (2) can be easily derived as

155 below.

156 $\hat{\delta} = \frac{1}{n}\sum_{j=1}^{n}\hat{\delta}(j) = \frac{1}{n}\sum_{j=1}^{n}\{\frac{1}{k}\sum_{i=1}^{k}[1 - DSC_{Di}(j)] - \frac{2}{k(k-1)}\sum_{i=1}^{k}\sum_{i'=i+1}^{k}[1 - DSC_{ii'}(j)]\}$ (3)

157 It can also alternatively be expressed as:

158 $\hat{\delta} = \frac{1}{nk}\sum_{j=1}^{n}\sum_{i=1}^{k}\hat{\delta}_i(j) = \frac{1}{nk}\sum_{j=1}^{n}\sum_{i=1}^{k}\{[1 - DSC_{Di}(j)] - \frac{1}{k-1}\sum_{\substack{i'=1 \\ i'\neq i}}^{k}[1 - DSC_{ii'}(j)]\}$ (4)

159 From the formulae (3) and (4), we can see $\hat{\delta}_i(j)$ is defined as the mean difference between

160 device and the $i^{th}$ individual reader on $j^{th}$ image, $\hat{\delta}(j)$ is the mean $\hat{\delta}_i(j)$ across all readers for $j^{th}$

161 image, and $\hat{\delta}$ is the mean of $\hat{\delta}(j)$ across all images.

162 *2.3.3 Confidence Interval*

163 Various approaches can be used to construct confidence intervals (CIs) for $\hat{\delta}$. In this study,

164 we used both a parametric and non-parametric method to construct CIs. The parametric method

165 estimates the z-interval, given by $\hat{\delta} \pm Z_{\alpha/2} S$, where $S = \sqrt{\frac{1}{n-1}\sum_{j}^{n}(\hat{\delta}(j) - \hat{\delta})^2}$ is the sample

166 standard deviation of $\hat{\delta}(j)$, and the non-parametric method is a bootstrap approach that follows

167 the procedure outlined in [11].

168 A CI covering zero indicates that no significant difference in overlap-based segmentation

169 performance between device and panel. Note, this should not be interpreted as the two are the

170 same but only that a difference could not be established statistically [14]. A CI entirely below zero

171 indicates device-panel segmentation agreement is statistically higher than the within-panel

172 agreement, while a CI above zero indicates the device-panel agreement is significantly lower

173 than the within-panel agreement.

## 3 Simulation Studies

This section presents the design and results from our simulation studies conducted as part of the method validation process.

### 3.1 Overall Study Design

#### 3.1.1 Study Overview

In this article, we undertake two primary types of simulation studies to validate our proposed method. Simulation Study 1 is a statistics-based simulation study, which simulates Dice scores from a predefined statistical distribution. This approach defines the 'agreement/interchangeability' between readers and devices by controlling the Dice distribution characteristics for each. This simulation allows for the assessment of our segmentation interchangeability metric across a range of Dice distributions. A limitation of the Simulation Study 1 design is it may not perfectly capture relationships inherent in image-generated Dice scores. The relationship here refers to the correlations between pairwise Dice scores coming from segmentations of the same objects from different annotators. Simulating Dice scores simply by sampling from a statistical distribution may not fully capture these relationships.

To address this limitation, we conducted Simulation Study 2 where we directly compute Dice scores from contour masks rather than simulating Dice scores directly. To achieve this, we employed the Medical Image Segmentation Synthesis (MISS) Tool, developed by Guan et al. [15] to synthesize multiple segmentation masks. The MISS Tool has a set of adjustable parameters simulating six types of segmentation errors. Using the MISS Tool, we generate Dice scores by applying the following process: (1) simulate a set of truth contours, (2) generate reader and device annotation variations using the MISS tool based on the true contours, and (3) calculate the

Dice scores from the simulated contours. Study 2 establishes 'agreement/interchangeability' between readers and devices by either using the same or different transformation parameters within the MISS Tool. To distinguish between the two simulation studies, we refer to (dis)agreement in Study 1 as "overlap (dis)agreeable" and in Study 2 as "transformation (dis)agreeable".

*3.1.2 Performance Metrics*

For Study 1, we evaluate the performance of the proposed method using three metrics:

1. **Type I Error**: The probability of incorrectly rejecting the null hypothesis ($H_0$) when it is actually true.

2. **Type II Error**: The probability of failing to reject the null hypothesis ($H_0$) when it is actually false.

3. **Coverage Probability (CP)**: The proportion of times a confidence interval contains the true value of the parameter being estimated (in our case, the true value of $\delta$).

A desirable method is expected to have type I error near 0.05 for agreeable cases, type II error below 0.2 for disagreeable cases assuming a statistical power of 0.80, and CP near the confidence level (set as 0.95 in our study).

For Study 2 (image-based simulation), we assess performance using Type I and Type II errors. Coverage Probability (CP) is not applicable as the true value of $\delta$ is unknown because the mean Dice performance cannot be directly controlled or defined using the MISS tool.

215 *3.2 Study 1: Statistics-based Simulation*

216 *3.2.1 Parameter Configuration*

217 In Simulation Study 1, the number of readers in the panel, denoted as k, are either 2 or 3 in our

218 experiments. We define the following global parameters for each simulation conduction:

219     • $(\boldsymbol{\mu_0}, \boldsymbol{\sigma_0})$: mean and standard deviation for any reader-reader pair DSC.

220     • $(\boldsymbol{d_\mu}, \boldsymbol{d_\sigma})$, where $d_\mu = \mu_D - \mu_0$, $d_\sigma = \sigma_D - \sigma_0$, with $\mu_D$ *and* $\sigma_D$ representing the mean and

221          standard deviation of DSC for any device-reader pair. $d_\mu$ and $d_\sigma$ correspond to the

222          differences in the means and standard deviations for the device-reader DSC and the

223          reader-reader DSC.

224     • $(\boldsymbol{\rho_0}, \boldsymbol{\rho_D}, \boldsymbol{\rho_{D0}})$: pairwise reader-reader DSC correlation $\boldsymbol{\rho_0}$, pairwise device-reader DSC

225          correlation $\boldsymbol{\rho_D}$, and the between device-reader and reader-reader DSC correlation $\boldsymbol{\rho_{D0}}$,

226          respectively.

227 To mirror real-world data variations, where within-panel and device-panel dissimilarities

228 may have different means and standard deviations, we explore four scenarios:

229        I.    equal mean, equal variance (Overlap Agreeable): $\mu_D = \mu_0, \sigma_D = \sigma_0$;

230        II.    unequal mean, equal variance (Overlap Disagreeable): $\mu_D \neq \mu_0, \sigma_D = \sigma_0$;

231        III.    equal mean, unequal variance (Overlap Agreeable): $\mu_D = \mu_0, \sigma_D \neq \sigma_0$;

232        IV.    unequal mean, unequal variance (Overlap Disagreeable): $\mu_D \neq \mu_0, \sigma_D \neq \sigma_0$.

233 Parameter configurations for these scenarios are presented in Table 1.

234 Table 1: Parameter configuration scenarios for Study 1: statistic-based simulations.

| Scenario | Parameter | configuration |
|---|---|---|
| **I.** equal μ, equal σ (overlap agreeable): 32 | $(\mu_0, \sigma_0)$ | $\mu_0$ in {0.75, 0.8, 0.85, 0.9} $\sigma_0$ in {0.025, 0.05, 0.1, 0.15} |
| | $(d_\mu, d_\sigma)$ | $d_\mu = 0, d_\sigma = 0$ |

10

| | | |
|---|---|---|
| settings total | $(\rho_0, \rho_D, \rho_{D0})$ | $\rho_0 = \rho_D = \rho_{D0}$ in {*moderate, strong/very strong*} (correlation criteria is described above in Section 3.2.1) |
| **II.** unequal μ, equal σ (overlap disagreeable): 512 settings total | $(\mu_0, \sigma_0)$ | $\mu_0$ in {0.75, 0.8, 0.85, 0.9}<br>$\sigma_0$ in {0.025, 0.05, 0.1, 0.15} |
| | $(d_\mu, d_\sigma)$ | $d_\mu$ in {-0.05, -0.1, -0.25, -0.5}, $d_\sigma = 0$ |
| | $(\rho_0, \rho_D, \rho_{D0})$ | $\rho_0$ in {*moderate, strong/very strong*}<br>$\rho_D$ in {*moderate, strong/very strong*}<br>$\rho_{D0}$ in {*weak, very weak*} |
| **III.** equal μ, unequal σ (overlap Dice agreeable): 96 settings total | $(\mu_0, \sigma_0)$ | $\mu_0$ in {0.75, 0.8, 0.85, 0.9}<br>$\sigma_0$ in {0.025, 0.05, 0.1} |
| | $(d_\mu, d_\sigma)$ | $d_\mu = 0$, $d_\sigma$ in {0.02, 0.05, 0.1, 0.15} |
| | $(\rho_0, \rho_D, \rho_{D0})$ | $\rho_0 = \rho_D = \rho_{D0}$ in {*moderate, strong/very strong*} |
| **IV.** unequal μ, unequal σ (overlap disagreeable): 648 settings total | $(\mu_0, \sigma_0)$ | $\mu_0$ in {0.8, 0.85, 0.9}<br>$\sigma_0$ in {0.025, 0.05, 0.1} |
| | $(d_\mu, d_\sigma)$ | $d_\mu$ in {-0.05, -0.1, -0.25},<br>$d_\sigma$ in {0.02, 0.1, 0.2} |
| | $(\rho_0, \rho_D, \rho_{D0})$ | $\rho_0$ in {*moderate, strong/very strong*}<br>$\rho_D$ in {*moderate, strong/very strong*}<br>$\rho_{D0}$ in {*weak, very weak*} |

235

**$(\mu_0, \sigma_0)$** values in Table 1 were determined based on the mean and standard deviation values reported in the existing literature for human-observer-pair Dice scores in tasks such as CT-scan lung nodule/tumor segmentation and liver tumor/brain hematomas segmentation (e.g., [16, 17]). **$(d\mu, d\sigma)$** values in Table 1 were defined based on our exploration of LIDC-based synthesized data. In practice device performance is typically lower than human experts, thus **dμ** is set to be non-positive in our studies, and **dσ** is set to be non-negative across all settings. **$(\rho_0, \rho_D, \rho_{D0})$** configurations in Table 1 are based on the following assumptions: 1) Dice scores drawn from distributions with the same mean exhibit *moderate or higher* correlation, whereas Dice scores drawn from distributions with different means demonstrate *weak or lower* correlation (criteria for defining different correlation strength categories are described in Section 3.2.2); and 2) an equi-correlation-category assumption is used to define correlation structure for between annotator-pair-Dice (details see next section 3.2.2).

*3.2.2 Dataset Simulation*

249    For each setting configured in Table 1, we simulated a total of $N_{sim} = 1000$ datasets. Each

250    dataset contained k readers. Thus there is a total of k*(k-1)/2 reader-pair DSCs and k device-

251    reader-pair DSCs for each of *n* images. Dice scores were generated using a multivariate Beta

252    distribution as follows:

253    Step 1: Compute the Beta shape parameters $(a_0, b_0)$ and $(a_D, b_D)$ corresponding to the pre-

254    specified mean and standard deviation of $(\mu_0, \sigma_0)$ and $(\mu_D, \sigma_D)$ for the reader-reader Dice

255    distribution and the device-reader Dice distribution, respectively.

256    Step 2: Generate the correlation matrix for pairwise Dices based on the pre-specified

257    correlation strength categories $(\rho_0, \rho_D, \rho_{D0})$ using the following: Set the diagonal elements in the

258    correlation matrix to 1. The upper triangle in the matrix was defined by randomly generating

259    correlation coefficients between pairwise reader-reader DSCs (DSC$_{ii'}$ vs. DSC$_{jj'}$, 1≤i<i'≤k,

260    1≤j<j'≤k, i<j) based on

$$Corr(DSCii', DSCjj') \sim \begin{cases} U(0, 0.2) & if \ \rho_0 \ = \ \text{``very weak''} \\ U[0.2, 0.4) & if \ \rho_0 \ = \ \text{``weak''} \\ U[0.4, 0.6) & if \ \rho_0 \ = \ \text{``moderate''} \\ U[0.6, 0.8) & if \ \rho_0 \ = \ \text{``strong''} \\ U[0.8, 1) & if \rho_0 = \ \text{``very strong''} \end{cases}$$

261    where U stands for the uniform distribution. The correlation categorization criteria here is taken

262    from the guidelines established by Evans [18], which are consistent with those used by LaMorte [19]

263    and Swinscow and Campbell [20]. Correlation between pairwise device-reader DSCs ($\rho_D$), and

264    between reader-reader DSCs and device-reader DSCs ($\rho_{D0}$) were determined in a similar fashion.

265    The correlation matrix, denoted as Σ and having dimensions k(k+1)/2 by k(k+1)/2, is symmetric;

266    therefore, the lower triangle can be filled by mirroring the upper triangle.

267    Step 3: Generate all annotator-pair DSCs via multivariate beta distribution as follows.

$$\left(DSC_{12}, \ldots, DSC_{(k-1)k}; DSC_{D1}, \ldots, DSC_{Dk}\right) \sim MultiBeta\{\vec{a}, \vec{b}, \Sigma\}$$

268    where $\vec{a} = \left(a_0\vec{1}^{(k-1)k/2}, a_D\vec{1}^k\right)$, $\vec{b} = \left(b_0\vec{1}^{(k-1)k/2}, b_D\vec{1}^k\right)$ $and$ $\vec{1}$ refers to an all-ones vector.

269    The beta shape parameters reflect the desired means and standard deviations from Step 1 and $\Sigma$

270    is the correlation matrix generated in Step 2. This step is implemented using modified code from

271    [21].

272        The statistical simulation was run using a non-synchronized parallelization techniques to

273    efficiently scale the large-scale simulations on a High-Performance Computing (HPC) clusters,

274    and the Pierre L'Ecuyer's [22] 'RngStreams' function was used to generate multiple independent

275    pseudo-random number streams ensuring no overlap across array job tasks.

276    *3.2.3 Results*

277        Below are selected results from the statistics-based simulation study. Note that in the rest of

278    the article, $n_{reader}$ and $n_{image}$ are used interchangeably with k and n, respectively, for ease of

279    reference. Performance metrics are computed for each setting based on $N_{Sim}$=1000, $n_{image}$=400,

280    and $n_{reader}$=2 and 3.

281        Figure 1 presents boxplots of the coverage probabilities (CP) and Type I/II error results for

282    both 2 and 3 readers across all conditions for the four scenarios. Tables 2–5 provide detailed

283    results for the Type I/II error and CP for the individual settings.  The tables include results for

284    the 3-reader scenario and include bootstrap CIs for all 32 settings in Scenario I and selected

285    settings from Scenarios II–IV. The 2-reader scenario and z-interval CIs are similar. The selected

286    settings from scenarios II and IV only include the  $d_\mu = -0.05$ setting as the type II error would

287    only be smaller for larger absolute differences $d_\mu < -0.05$.

288    From Figure 1 and Tables 2–5, it is evident that the proposed method performs well.

289    Specifically, Type I error remains close to 0.05 across all settings for overlap-agreeable cases

290    (Scenarios I and III), Type II error approaches 0 in all settings for overlap-disagreeable cases

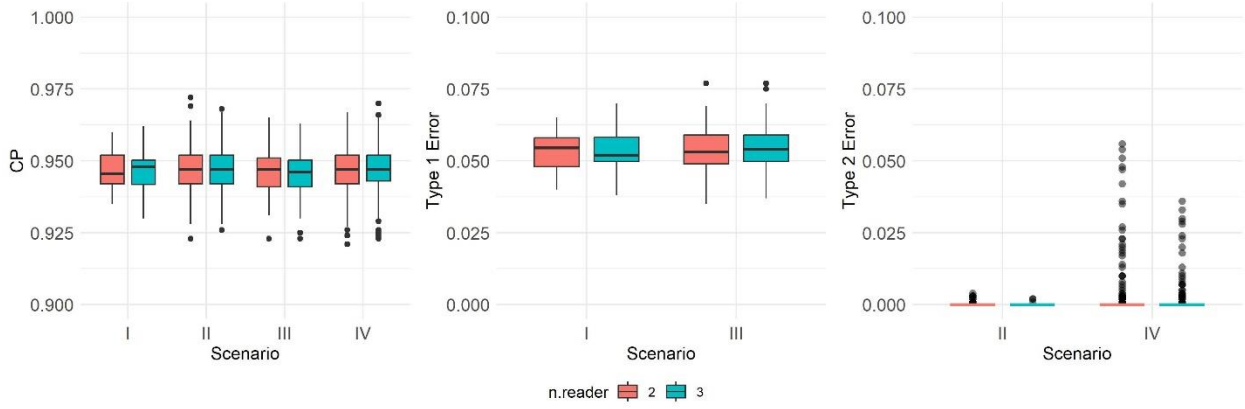291    (Scenarios II and IV), and CP consistently hovers around 0.95 across all the scenarios.



292

293    Figure 1. CP, Type I error and Type II error results for the statistics-based Simulation Study 1.  The figure includes

294        boxplot results aggregated across all the settings in Scenarios I – IV, and for both the 2-reader and 3-

295        reader scenarios.

296

297    Table 2. Results for all 32 settings in Scenario I: Equal Mean $\mu$ and Equal Standard Deviation $\sigma$ fromstatistics-based
298                                        simulation Study 1.

| $\mu_0$ | $\sigma_0$ | $\rho_0 = \rho_D = \rho_{D0}$ | CP | Type I Error |
|---|---|---|---|---|
| 0.75 | 0.025 | moderate | 0.937 | 0.063 |
| 0.8 | 0.025 | moderate | 0.955 | 0.045 |
| 0.85 | 0.025 | moderate | 0.947 | 0.053 |
| 0.9 | 0.025 | moderate | 0.943 | 0.057 |
| 0.75 | 0.05 | moderate | 0.949 | 0.051 |
| 0.8 | 0.05 | moderate | 0.93 | 0.07 |
| 0.85 | 0.05 | moderate | 0.951 | 0.049 |
| 0.9 | 0.05 | moderate | 0.939 | 0.061 |
| 0.75 | 0.1 | moderate | 0.942 | 0.058 |
| 0.8 | 0.1 | moderate | 0.948 | 0.052 |
| 0.85 | 0.1 | moderate | 0.949 | 0.051 |
| 0.9 | 0.1 | moderate | 0.953 | 0.047 |
| 0.75 | 0.15 | moderate | 0.95 | 0.05 |
| 0.8 | 0.15 | moderate | 0.946 | 0.054 |

14

| | | | | |
|---|---|---|---|---|
| 0.85 | 0.15 | moderate | 0.962 | 0.038 |
| 0.9 | 0.15 | moderate | 0.948 | 0.052 |
| 0.75 | 0.025 | strong/very strong | 0.957 | 0.043 |
| 0.8 | 0.025 | strong/very strong | 0.961 | 0.039 |
| 0.85 | 0.025 | strong/very strong | 0.95 | 0.05 |
| 0.9 | 0.025 | strong/very strong | 0.936 | 0.064 |
| 0.75 | 0.05 | strong/very strong | 0.957 | 0.043 |
| 0.8 | 0.05 | strong/very strong | 0.942 | 0.058 |
| 0.85 | 0.05 | strong/very strong | 0.941 | 0.059 |
| 0.9 | 0.05 | strong/very strong | 0.935 | 0.065 |
| 0.75 | 0.1 | strong/very strong | 0.948 | 0.052 |
| 0.8 | 0.1 | strong/very strong | 0.942 | 0.058 |
| 0.85 | 0.1 | strong/very strong | 0.932 | 0.068 |
| 0.9 | 0.1 | strong/very strong | 0.948 | 0.052 |
| 0.75 | 0.15 | strong/very strong | 0.945 | 0.055 |
| 0.8 | 0.15 | strong/very strong | 0.941 | 0.059 |
| 0.85 | 0.15 | strong/very strong | 0.961 | 0.039 |
| 0.9 | 0.15 | strong/very strong | 0.95 | 0.05 |

299

300 Table 3. Selected results for 8 Settings with $(\mu_0, d_\mu, \sigma_0) = (0.85, -0.05, 0.15)$ in Scenario II: Unequal Mean $\mu$ and
301 Equal Standard Deviation $\sigma$ from statistics-based Simulation Study 1.

302

| $\mu_0$ | $d_\mu$ | $\sigma_0$ | $\rho_0$ | $\rho_D$ | $\rho_{D0}$ | Type II error | CP |
|---|---|---|---|---|---|---|---|
| 0.85 | -0.05 | 0.15 | moderate | moderate | very weak | 0 | 0.955 |
| 0.85 | -0.05 | 0.15 | moderate | moderate | weak | 0 | 0.939 |
| 0.85 | -0.05 | 0.15 | moderate | strong/very strong | very weak | 0.002 | 0.94 |
| 0.85 | -0.05 | 0.15 | moderate | strong/very strong | weak | 0 | 0.944 |
| 0.85 | -0.05 | 0.15 | strong/very strong | moderate | very weak | 0 | 0.946 |
| 0.85 | -0.05 | 0.15 | strong/very strong | moderate | weak | 0 | 0.94 |
| 0.85 | -0.05 | 0.15 | strong/very strong | strong/very strong | very weak | 0.001 | 0.942 |
| 0.85 | -0.05 | 0.15 | strong/very strong | strong/very strong | weak | 0 | 0.941 |

303

304 Table 4. Selected results for 24 of 96 Settings with $\mu_0 = 0.75$ in Scenario III: Equal Mean $\mu$ and Unequal Standard
305 Deviation $\sigma$ from statistics-based Simulation Study 1.

306

| $\mu_0$ | $\sigma_0$ | $d_\sigma$ | $\rho_0 = \rho_D = \rho_{D0}$ | Type I error | CP |
|---|---|---|---|---|---|
| 0.75 | 0.025 | 0.15 | strong or very strong | 0.054 | 0.946 |
| 0.75 | 0.1 | 0.02 | moderate | 0.048 | 0.952 |
| 0.75 | 0.05 | 0.02 | moderate | 0.058 | 0.942 |
| 0.75 | 0.025 | 0.02 | moderate | 0.05 | 0.95 |
| 0.75 | 0.1 | 0.05 | moderate | 0.049 | 0.951 |
| 0.75 | 0.05 | 0.05 | moderate | 0.054 | 0.946 |
| 0.75 | 0.025 | 0.05 | moderate | 0.048 | 0.952 |
| 0.75 | 0.1 | 0.1 | moderate | 0.067 | 0.933 |
| 0.75 | 0.05 | 0.1 | moderate | 0.056 | 0.944 |
| 0.75 | 0.025 | 0.1 | moderate | 0.061 | 0.939 |
| 0.75 | 0.1 | 0.15 | moderate | 0.055 | 0.945 |
| 0.75 | 0.05 | 0.15 | moderate | 0.058 | 0.942 |
| 0.75 | 0.025 | 0.15 | moderate | 0.054 | 0.946 |
| 0.75 | 0.1 | 0.02 | strong or very strong | 0.066 | 0.934 |
| 0.75 | 0.05 | 0.02 | strong or very strong | 0.048 | 0.952 |
| 0.75 | 0.025 | 0.02 | strong or very strong | 0.052 | 0.948 |
| 0.75 | 0.1 | 0.05 | strong or very strong | 0.051 | 0.949 |
| 0.75 | 0.05 | 0.05 | strong or very strong | 0.05 | 0.95 |
| 0.75 | 0.025 | 0.05 | strong or very strong | 0.055 | 0.945 |
| 0.75 | 0.1 | 0.1 | strong or very strong | 0.042 | 0.958 |
| 0.75 | 0.05 | 0.1 | strong or very strong | 0.055 | 0.945 |
| 0.75 | 0.025 | 0.1 | strong or very strong | 0.059 | 0.941 |
| 0.75 | 0.1 | 0.15 | strong or very strong | 0.046 | 0.954 |
| 0.75 | 0.05 | 0.15 | strong or very strong | 0.061 | 0.939 |

307

308 Table 5. Selected results for 24 of 648 Settings with $(\mu_0, d_\mu, \sigma_0) = (0.8, -0.05, 0.1)$ in Scenario IV: Unequal Mean
309 $\mu$ and Unequal Standard Deviation $\sigma$ from statistics-based Simulation Study 1.

| $\mu_0$ | $d_\mu$ | $\sigma_0$ | $d_\sigma$ | $\rho_0$ | $\rho_D$ | $\rho_{D0}$ | Type II error | CP |
|---|---|---|---|---|---|---|---|---|
| 0.8 | -0.05 | 0.1 | 0.2 | moderate | moderate | very weak | 0.007 | 0.956 |
| 0.8 | -0.05 | 0.1 | 0.2 | moderate | moderate | weak | 0.002 | 0.956 |
| 0.8 | -0.05 | 0.1 | 0.2 | moderate | strong/very strong | very weak | 0.033 | 0.952 |
| 0.8 | -0.05 | 0.1 | 0.2 | moderate | strong/very strong | weak | 0.02 | 0.932 |
| 0.8 | -0.05 | 0.1 | 0.2 | strong/very strong | moderate | very weak | 0.009 | 0.952 |
| 0.8 | -0.05 | 0.1 | 0.2 | strong/very strong | moderate | weak | 0.003 | 0.947 |
| 0.8 | -0.05 | 0.1 | 0.2 | strong/very strong | strong/very strong | very weak | 0.036 | 0.945 |
| 0.8 | -0.05 | 0.1 | 0.2 | strong/very strong | strong/very strong | weak | 0.018 | 0.95 |
| 0.8 | -0.05 | 0.1 | 0.1 | moderate | moderate | very weak | 0 | 0.959 |
| 0.8 | -0.05 | 0.1 | 0.1 | moderate | moderate | weak | 0 | 0.947 |
| 0.8 | -0.05 | 0.1 | 0.1 | moderate | strong/very strong | very weak | 0 | 0.948 |
| 0.8 | -0.05 | 0.1 | 0.1 | moderate | strong/very strong | weak | 0 | 0.948 |
| 0.8 | -0.05 | 0.1 | 0.1 | strong/very strong | moderate | very weak | 0 | 0.937 |

| 0.8 | -0.05 | 0.1 | 0.1 | strong/very strong | moderate | weak | 0 | 0.945 |
|---|---|---|---|---|---|---|---|---|
| 0.8 | -0.05 | 0.1 | 0.1 | strong/very strong | strong/very strong | very weak | 0.001 | 0.946 |
| 0.8 | -0.05 | 0.1 | 0.1 | strong/very strong | strong/very strong | weak | 0 | 0.935 |
| 0.8 | -0.05 | 0.1 | 0.02 | moderate | moderate | very weak | 0 | 0.951 |
| 0.8 | -0.05 | 0.1 | 0.02 | moderate | moderate | weak | 0 | 0.943 |
| 0.8 | -0.05 | 0.1 | 0.02 | moderate | strong/very strong | very weak | 0 | 0.954 |
| 0.8 | -0.05 | 0.1 | 0.02 | moderate | strong/very strong | weak | 0 | 0.941 |
| 0.8 | -0.05 | 0.1 | 0.02 | strong/very strong | moderate | very weak | 0 | 0.938 |
| 0.8 | -0.05 | 0.1 | 0.02 | strong/very strong | moderate | weak | 0 | 0.951 |
| 0.8 | -0.05 | 0.1 | 0.02 | strong/very strong | strong/very strong | very weak | 0 | 0.952 |
| 0.8 | -0.05 | 0.1 | 0.02 | strong/very strong | strong/very strong | weak | 0 | 0.946 |

310

### 3.3  Simulation Study 2: Image-based Simulation Study

312  In Study 2, we directly compute Dice scores from contour masks rather than simulating Dice

313  scores from statistical distributions. To achieve this, we employed the Medical Image

314  Segmentation Synthesis (MISS) Tool, developed by Guan et al.[15] to synthesize segmentation

315  masks with multiple controlled types of segmentation errors, such as spiculations and

316  shape/alignment changes.

**MISS-Tool Overview**

318  The MISS-tool emulates segmentation errors by modifying truth masks of anatomical objects

319  through a set of adjustable parameters that simulate six typical segmentation errors: boundary

320  spiculation, under/over-sizing, centroid location errors, overlap variations, shape/alignment

321  details, and the introduction of satellite structures. These segmentation error types are

322  implemented through four primary image processing methods:

323  1.  **Affine Transformation**: Modifies segmentation contours through resizing (changing

324     height and width ratios), shifting (location changes in x, y coordinates), and rotation

17

325      (angle parameter).

326      2.  **Spiculation**: Adjusts contours in polar coordinates by adding Gaussian functions to

327          create spike-like protrusions or indentations on boundaries, representing boundary

328          irregularities.

329      3.  **Fourier Descriptor (FD) Modification**: Modifies contours in the spatial frequency

330          domain by keeping low-frequency components (basic shape) while allowing changes to

331          middle-frequency components and removing high-frequency components (fine details of

332          the shape).

333      4.  **Satellite Structure Synthesis**: Adds separate small objects in nearby regions of the true

334          object, simulating disconnected components that are incorrectly included in segmentation

335          results.

336  For our simulation study, we employed three of the four image processing methods (affine

337  transformation, spiculation, and Fourier descriptor modification) to generate synthetic

338  segmentation variations that approximate potential segmentation difference between different

339  annotators.

340  *3.3.1 Parameter Configuration*

341  In Study 2, we consider cases containing synthetic contours with 3, 5, 7 and 9 readers. Contours

342  made by all annotators on the same image are provided as 2D binary images of identical

343  dimensions. We define the following parameters to control the synthetic contours used in Study

344  2 based on the MISS-tool methodology:

18

- **Affine transformation parameters: Resizing ratios (Rx, Ry):** control the width and height scaling of the segmentation boundaries, where ratios equal to 1 mean no change. **Location shift (Sx, Sy):** controls spatial displacement in x and y coordinates measured in pixels. **Rotation angle (φ):** controls the rotational variation of the segmented contours.

- **Fourier transformation parameters:** (as shown in the Supplementary Figure S1) **Detail:** number of non-zero Fourier descriptors from low to high frequences that will be kept or modified. This controls the level of boundary detail preservation, where higher values retain more fine-grained boundary features. **Range:** number of Fourier descriptors with middle frequency components to be modified. This determines the frequency bandwidth of modifications applied to the contour. **Magnitude:** strength value controlling how the middle-ranged Fourier descriptors are modified. This determines the extent of shape changes applied through frequency domain manipulation.

- **Spiculation parameters:** for adding Gaussian functions to create spike-like protrusions in polar coordinates (as shown in the Supplementary Figure S2), **the following parameters are used: Center (c):** angular position (in degrees, 0-360°) where spiculation is added to the contour boundary. **Height (h):** magnitude of spike protrusion or indentation, where positive values create convex spiculations and negative values create concave indentations. **Width (w):** angular spread of the spiculation feature, determining how broad the spike modification appears on the boundary.

We then explored two scenarios in Study 2:

I. device segmentation is **Transformation Agreeable** with the reader panel, and

II. device segmentation is **Transformation Disagreeable** with the reader panel.

367      Each setting is characterized by the 11 parameters defined above. Four affine parameters

368    (Rx, Ry, Sx, Sy) are used as "*tunable*" parameters, which introduce variability across settings.

369    The remaining seven parameters (including all Fourier and spike transformation parameters) are

370    designated as "*default*" parameters, maintaining uniformity across all settings. Table 6 and

371    Appendix A.1 present the configurations for the *tunable* and *default* parameters for the two

372    scenarios being studied. This design allows us to approximate the range of DSC distributions

373    reported in the literature, as confirmed in results later.

374              Table 6. Configuration for *tunable* parameters in image-based Simulation Study 2.

| scenario | parameters | configuration |
|---|---|---|
| **I.** reader panel vs. device: Transformation Agreeable | $(n_{image}, n_{reader})$ | $n_{image}$ in (100, 250, 750) <br> $n_{reader}$ in (3, 5, 7, 9) |
| | $(R_x, R_y)$ for both panel and device | $R_{max} = 1.15, 1.1,$ or $1.05$ <br> $R_x \sim U[\,2 - R_{max}, R_{max}]$ <br> $R_y \sim U[\,2 - R_{max}, R_{max}]$ |
| | $(S_x, S_y)$ for both panel and device | $S_{max} = 0, 2,$ or $3$ pixels <br> $S_x \sim U[-S_{max}, S_{max}]$ <br> $S_y \sim U[-S_{max}, S_{max}]$ |
| **II.** reader panel vs. device: Transformation Disagreeable | $(n_{image}, n_{reader})$ | $n_{image}$ in (100, 250, 750) <br> $n_{reader}$ in (3, 5, 7, 9) |
| | $(R_x, R_y, S_x, S_y)$ for panel | Same as in Scenario I above |
| | $(R_x, R_y, S_x, S_y)$ for device | For each of the 9 reader panel settings, the device $(R_x, R_y, S_x, S_y)$ parameter set is selected from the remaining 8 settings. |

375

376    *3.3.2 Dataset Simulation*

377    We simulate each Study 2 image dataset and corresponding annotations using the following

378    steps.

379    Step 1: Generate $n$ true contours, each comprising a single randomly located and pixelized

380    circle. These true contours serve as the foundation for generating both reader and device

381    contours.

382    Step 2: Simulate $k + 1$ contours for each true contour representing annotations by $k$ readers

383    and the device, respectively. This is achieved by applying the MISS tool transformations using

384    the parameter setting described in Table 6 to the true contours.

385    Two sets of contours generated using the same set of MISS parameters are defined as

386    Transformation Agreeable. Otherwise, they are deemed Transformation Disagreeable. Figure 2

387    provides an illustrative example of a true contour and various contours synthesized using the

388    MISS tool transformations for three experts and a Transformation Disagreeable device.

389



390    a)  truth       b) expert 1       c) expert 2       d) expert 3       e) device

391    Figure 2.  Example of synthetic expert and device contours based on an initial true contour.  Note, this is a

392           Transformation Disagreeable example where the device contour is based on a different set of MISS

393           parameters compared to that of the experts.  a) true contour, b)-d) synthetic reader contours, e)

394           transformation disagreeable device contour.

395    For each setting in Table 6, we simulate a total of $N_{Sim}$ imaging datasets. For each image

396    dataset, we compute image-level Dice scores for each pair of annotators. Subsequently, based on

397    the Dice scores computed for each image, we calculate the performance metrics as defined in

398    Section 3.1. The results are summarized in next subsection.

399     *3.3.3 Results*

400     Below are results for image-based simulation study. Performance metrics are computed for

401     each setting based on $N_{Sim}$=100, $n_{image}$=100, 250, 750, and $n_{reader}$=3, 5, 7, 9. All results

402     presented here are based on the z-interval CI approach. The results for the bootstrapping

403     approach are omitted as they are similar.

404     *3.3.3.1 Scenario I: Expert panel vs. device transformation-agreeable*
405     Table 7 summarizes the mean DSC performance of the expert panel, the mean DSC

406     performance difference between the device and panel, and the Type I error of the proposed

407     method for each Transformation Agreeable setting, as configured in Table 6. The results in Table

408     7 indicate that the proposed method performs well across all Transformation Agreeable settings

409     simulated, with Type I error close to 0.05 on average and small standard deviations (ranging

410     from 0.02 to 0.04) across all of the settings evaluated. These findings suggest that when the

411     device and expert panel share the same transformation pattern, their segmentation performance is

412     similar (reflected by a mean DSC difference close to 0 in Table 7), and the proposed method

413     achieves reasonable Type I error (ranging between 0.03 – 0.10 in Table 7, with a mean of 0.065

414     across all settings). Figure 3 illustrates the Type I error as a function of $n_{image}$ (i.e., sample size)

415     and $n_{reader}$ (i.e., panel size) for each of the nine individual settings under Scenario I:

416     Transformation Agreeable.

417     Table 7. Mean (standard deviation) for Within-Panel DSC, Device-Panel DSC, and Type 1 Error by panel
418             transformation pattern from Image-Simulation Study 2, Scenario I: Transformation Agreeable.
419

| | | Expert Panel Transformation parameters | | |
| --- | --- | --- | --- | --- |
| | | $S_{max}$=0 | $S_{max}$=2 | $S_{max}$=3 |
| Reader-pair DSC $\hat{\mu}_0$ | $R_{max}$=1.05 | 0.95 (0.04) | 0.88 (0.09) | 0.84 (0.13) |
| | $R_{max}$=1.10 | 0.93 (0.04) | 0.87 (0.09) | 0.83 (0.12) |
| | $R_{max}$=1.15 | 0.90 (0.05) | 0.85 (0.09) | 0.81 (0.12) |
| Device-Reader DSC | $R_{max}$=1.05 | 0.95 (0.04) | 0.88 (0.09) | 0.84 (0.13) |

| | | | | |
|---|---|---|---|---|
| $\hat{\mu}_D$ | $R_{max}=1.10$ | 0.93 (0.04) | 0.87 (0.09) | 0.83 (0.12) |
| | $R_{max}=1.15$ | 0.90 (0.05) | 0.85 (0.09) | 0.81 (0.12) |
| Type 1 error | $R_{max}=1.05$ | 0.10 (0.04) | 0.05 (0.02) | 0.03 (0.03) |
| $\alpha$ | $R_{max}=1.10$ | 0.09 (0.02) | 0.06 (0.02) | 0.06 (0.02) |
| | $R_{max}=1.15$ | 0.07 (0.02) | 0.07 (0.03) | 0.05 (0.02) |

420

421



422

423 Figure 3. Plots of the Type 1 error results by transformation pattern ($R_{max}$, $S_{max}$) for Image Simulation Study 2,
424     Scenario I: Transformation Agreeable. The black dashed horizontal line marks the expected type 1 error
425     level of 0.05.

426
427 *3.3.3.2 Scenario II: expert panel vs. device transformation-disagreeable*

428
429 Scenario II (Transformation Disagreeable) includes 72 settings. Table 8 summarizes the mean

430 within-panel DSC, the mean device-panel DSC, and the Type II error of the proposed method by

431 expert panel transformation pattern defined by defined by (Rmax, Smax). Each cell's summary

432 statistics are aggregated over the eight Transformation Disagreeable device settings

433 corresponding to that panel pattern.

434 Table 8. Mean (standard deviation) of Within-Panel DSC, Device-Panel DSC, and Type 2 Error by Panel
435 Transformation Pattern, Image-Simulation Study 2, Scenario II: Transformation Disagreeable.

436

| | | Expert Panel Transformation Pattern | | |
| | | $S_{max}=0$ | $S_{max}=2$ | $S_{max}=3$ |
|---|---|---|---|---|
| Reader-pair DSC $\hat{\mu}_0$ ($\hat{\sigma}_0$) | $R_{max}=1.05$ | 0.95 (0.04) | 0.88 (0.09) | 0.84 (0.13) |
| | $R_{max}=1.10$ | 0.93 (0.04) | 0.87 (0.09) | 0.83 (0.12) |
| | $R_{max}=1.15$ | 0.90 (0.05) | 0.85 (0.09) | 0.81 (0.12) |
| Device-Reader DSC: $\hat{\mu}_D$ ($\hat{\sigma}_D$) | $R_{max}=1.05$ | 0.89 (0.07) | 0.87 (0.09) | 0.85 (0.1) |
| | $R_{max}=1.10$ | 0.89 (0.07) | 0.87 (0.09) | 0.85 (0.1) |
| | $R_{max}=1.15$ | 0.88 (0.07) | 0.86 (0.08) | 0.84 (0.1) |
| Type 2 error $\beta$ | $R_{max}=1.05$ | 0.00 (0.00) | 0.08 (0.18) | 0.16 (0.28) |
| | $R_{max}=1.10$ | 0.01 (0.03) | 0.07 (0.16) | 0.14 (0.26) |
| | $R_{max}=1.15$ | 0.01 (0.04) | 0.09 (0.18) | 0.07 (0.17) |

437
438
439 Table 8 reveals that the average Type II error for both interval approaches is below the

440 threshold of 0.20, demonstrating the proposed method performs well in transformation-

441 disagreeable settings overall. Across all Transformation Disagreeable settings, the mean

442 (standard deviation) Type II error is 0.072 (0.182) for the z-interval approach. While not

443 explicitly detailed, the mean (standard deviation) of the Type II error based on bootstrap

444 confidence intervals is similar, at 0.070 (0.177).

445 The results highlight an interesting contrast in error fluctuations. Unlike Type I error, which

446 shows relatively small variability for both Overlap Agreeable (Figure 1) and Transformation

447      Agreeable scenarios (std = 0.03, Table 7), the Type II error exhibits low variability in Overlap

448      Disagreeable scenarios (Figure 1) but high variability in Transformation-Disagreeable scenarios

449      (std = 0.18, Table 8).  This disparity arises because Transformation Disagreeable scenarios do

450      not necessarily equate to Overlap Disagreeable. Certain Transformation Disagreeable patterns

451      can still yield similar overlap-based segmentation performance. For example, the transformation

452      patterns $(R_{max}, S_{max}) = (1.05, 2)$ and $(1.10, 2)$ produce similar mean (standard deviation) DSC

453      values of 0.88 (0.09) and 0.87 (0.09), respectively. The inclusion of both Overlap Agreeable and

454      Overlap Disagreeable cases within the transformation-disagreeable scenario contributes to higher

455      variability in Type II error.

456      Table 8 also demonstrates that as the absolute mean DSC difference $\hat{d}_\mu = \hat{\mu}_D - \hat{\mu}_0$  increases,

457      the mean Type II error approaches 0. This is as expected, as the method is expected to perform

458      better when the DSC performance difference between the device and the panel becomes more

459      pronounced.

460      An analysis of subgroup trends reveals that as the maximum reader shift $(R_{max})$ increases

461      from 0 to 2 to 3, the mean Type II error rises from 0.01 to 0.08 to 0.13, respectively. This

462      increase corresponds to a greater uncertainty in DSC performance as well because  higher pixel

463      shift counts are allowed. In Table 8, the DSC standard deviation increases from 0.04 to 0.09 to

464      0.13 as $R_{max}$ increases from 0 to 2 to 3. This trend suggests that statistical power increases when

465      expert transformation patterns exhibit less variability. This would be akin to a situation where all

466      panel members have substantial experience and expertise, resulting is less variation among their

467      segmentation contours. However, it is important to note that deliberately selecting experts to

468      ensure near-identical patterns (e.g., all from one institution with identical training) is

469      inappropriate, as this level of agreement is unrealistic in the broader clinical context.

470        To further visualize the method's performance for individual settings, Figure 4 illustrates the

471        Type II error for each setting with fixed panel transformation parameters $(R_{max}, S_{max}) = (1.10, 2)$.

472        Results for other (Rmax, Smax) settings are omitted here and provided instead in Supplementary

473        Figures S3–S10, as their patterns are similar to the trends shown in Figure 4.  Each subplot in

474        Figure 4 displays the Type II error of the proposed method when the device adopts a

475        transformation pattern disagreeable to the panel across different numbers of images and experts.



476

Figure 4.  Type 2 error by device transformation pattern $(R^D_{max}, S^D_{max})$, given panel transformation $(R_{max}, S_{max}) =$
$(1.10, 2)$, with mean(std) within-panel Dice 0.87(0.09) for Image Simulation Study 2, Scenario II:
Transformation Disagreeable. The black dashed horizontal line marks type 2 error level at 0.2.

480        Figure 4 shows a decreasing Type II error as sample size or panel size increases. These

481        trends indicate that increasing the sample size and the panel size (assuming all experts have

482        similar DSC performance) enhances the power of the proposed method. The results presented in

483        Figure 4 further corroborate the findings in Table 8, demonstrating that the proposed method

484    achieves consistently low Type II error in settings with larger difference in Dice performance,

485    except in transformation-disagreeable but with smaller difference in the Dice performance. But

486    even in these challenging cases, increasing the sample size or the number of experts in the panel

487    effectively reduces the Type II error. Notably, the instances of Type II error exceeding 0.2 in

488    Figure 4 are unlikely to pose practical concerns, as the DSC performance differences in these

489    settings are minimal (e.g., as small as 0.01) and are shown to be mitigated to below 0.2 by

490    increasing the sample size.

491    **4   Application**

492    This section illustrates a real-world example to demonstrate the practical application of the

493    segmentation interchangeability metric for comparison of an AI segmentation and multiple

494    human experts.

495    *4.1   Database and Data Preparation*

496    The data are taken from the LIDC-IDRI database [23] developed by the Lung Image Database

497    Consortium. This is one of the largest publicly available datasets for lung nodule detection and

498    segmentation. It contains data from 1010 patients (1018 studies) and 2660 nodules, with slice

499    thickness varying from 0.45 mm to 5.0 mm. Each study includes clinical thoracic CT images

500    accompanied by an XML file documenting the annotations made by four experienced thoracic

501    radiologists. This dataset is anonymized, and all protected health information (PHI) is removed

502    [23].

503       A U-Net segmentation model with a ResNeXt encoder and ImageNet transfer learning was

504    trained using a subset of LIDC-IDRI data. This U-Net model was then compared to the

505    segmentation performance of 4 radiologist annotators using a testing dataset based on 124

506    independent LIDC-IDRI images. Details of training and testing data preparation as well as model

507    training are provided in the Appendix A2.

508    *4.2    Results*

509    Figure 5 shows the distribution of DSC values for pairwise annotators. The results suggest the AI

510    model may not agree with the panel, as evidenced by a lower DSC between the AI model and the

511    readers compared to the DSC between pairwise readers. However, this visual comparison does

512    not substantiate a statistical assessment of the difference.



513

514                    Figure 5. Distribution of DSC values for each pair of annotators.

515        We now apply our proposed statistical testing method to assess the interchangeability of the

516    AI model and four radiologists in the panel. Table 10 reports the mean (standard deviation) of

517    within-panel DSC $\hat{\mu}_P(\hat{\sigma}_P)$, the mean (standard deviation) for the AI-panel comparison $\hat{\mu}_A(\hat{\sigma}_A)$,

518    the difference between the two $\hat{d} = \hat{\mu}_P - \hat{\mu}_A$ and 95% CI derived from bootstrapping and z-

519    interval approaches derived using the proposed method.

520    Table 9. Agreement assessment results for the U-Net model (device) and the expert panel on

521    the LIDC-IDRI test dataset.

| Annotator vs. Reference Panel | $\hat{\mu}_P$ $(\hat{\sigma}_P)$ | $\hat{\mu}_A$ $(\hat{\sigma}_A)$ | $\hat{d} = \hat{\mu}_P - \hat{\mu}_A$ (95% CI) bootstrapping approach | $\hat{d} = \hat{\mu}_P - \hat{\mu}_A$ (95% CI) z-interval approach |
|---|---|---|---|---|
| U-Net vs. 4 readers | 0.7619 (0.0937) | 0.6896 (0.1172) | 0.0723 (0.0584, 0.087) | 0.0723 (0.0583, 0.0863) |

522

523    The results in Table 9 indicate that the trained U-Net model vs. 4-reader panel comparison

524    yields a 95% z-interval CI of (0.0583, 0.0856) (with a similar result from the bootstrapping CI,

525    as shown in Table 9). This indicates that the U-Net model has significantly lower agreement with

526    the panel than the within-panel agreement, and thus it is not interchangeable with the reader

527    panel. We note that the trained AI model presented here is used exclusively for illustrative

528    purposes, and the results of this specific use case should ***not*** be interpreted as evidence that U-

529    Net segmentation models more generally underperform relative to human experts.


530    **5    Discussion**

531    In this work, we developed a segmentation interchangeability metric and statistical method for

532    evaluating agreement between an AI device and a panel of human experts. Through a statistical

533    and an image-based simulation studies, we demonstrated that the proposed method exhibits well-

534    controlled Type I error and good Type II error behavior. The novelty of this method lies in its

535    ability to directly assess the interchangeability of a segmentation AI device with multiple human

536    experts without requiring a reference standard. This distinguishes our method from the

537    traditional approach that do require defining reference standard contours. Additionally, setting an

538    acceptable performance goal with conventional methods is challenging due to the lack of a

539    widely accepted clinical cutoff for the Dice score. In contrast, our method compares device

29

540   performance directly with an expert panel as a control, eliminating the need for a surrogate

541   ground truth. Setting a performance goal based on the mean DSC difference is also simpler, as

542   the target value is typically close to 0.

543   A limitation of the proposed method is it treats reader effect as fixed. This may not be a

544   major concern, as treating a small reader panel as fixed is not uncommon. In the context of

545   multi-reader studies, this approach has been used by Bandos et al. [25] and discussed by Hillis and

546   Schartz [26]. It would be an interesting future direction to incorporate methods accounting for

547   truther panel variability to refine our proposed method as an extension. Another limitation of the

548   method is it is designed to assess differences in overlap-based segmentation performance not

549   distance-based performance. One can easily substitute other metric into the proposed method,

550   such as a distance-based metric, but we have not determined how well the method's assumptions

551   hold.

552   The simulation study results indicate that increasing the sample size or the expert panel size

553   can enhance the power of a study utilizing the proposed method. In practice, it may be

554   challenging to expand the panel size while maintaining a high within-panel agreement level. As

555   such, increasing the sample size may be a most practical and feasible approach compared to

556   enlarging the panel.

557

558   **Appendix A**

559   *A1 Table. Configuration for Default parameters in Image-based Simulation*

| transform | parameter | configuration |
|---|---|---|
| fourier | (detail, range, magnitude) | detail = # of all pixels in original contour |
| | | range = round(detail * 0.2) + 1 |
| | | magnitude = 2 |

| | | center ~ U[0°, 360°] |
|---|---|---|
| spike | (center, height, width) | height ~ estimate lesion diameter * U[0.01, 0.2] (randomly assign ± sign) |
| | | width ~ estimate lesion diameter * U[0.01-0.2] / 2 |
| affine - rotation | φ | φ ~ U[0°, 360°] |

560  *A2  Training/Testing Data Preparation and Model Training*

561  Based on LIDC-IDRI database, we prepare the final data to be used for method demonstration

562  using following steps:

563  1. For each nodule, we draw the central slice from each scan (if there is even number of

564  slices (say, 2m slices) in a scan, we take $(m+1)^{th}$ slice as the central slice for this scan).

565  2. From the imaging dataset created in Step 1, we randomly draw a 70% sample (at the

566  patient level) and use these cases to train an AI segmentation model (a **U-net** architecture

567  with ResNeXt Encoder that is pretrained on ImageNet Database) using an IoU metric.

568  The reference standard mask for the training data, which is required for AI model training

569  purposes, is created by applying Majority Vote (MV) rule to the 4 radiologists' manual

570  annotations. For illustration, the figure below presents an example of the original slice,

571  alongside annotations from four radiologists, the aggregated consensus derived using the

572  Majority Vote (MV) criterion, and the annotation extracted by the U-Net algorithm.

31

Figure A1: Radiologists annotations of an LIDC-IDRI lesion along with the majority vote (MV) consensus, and the contour produced by our U-Net segemantation algorithm.

3. From the remaining 30% of patients not used in training, we collected cases where the 4 radiologists agreed on the lesion (defined here as having at least a 1 pixel overlap between annotations from any pair of radiologists in panel). One image per patient is then randomly chosen to form the final testing dataset, resulting in 124 independent images.

4. Using contours generated by the algorithm and those annotated by the four radiologists on the testing data, compute paired annotator Dice scores. And then the computed Dice scores will be fed to the proposed method for annotator vs. reference panel interchangeability assessment.

**Declaration of interest statement**

The authors declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

## Disclaimer

The mention of commercial products, their sources, or their use in connection with material reported herein is not to be construed as either an actual or implied endorsement of such products by the Department of Health and Human Services. This is a contribution of the U.S. Food and Drug Administration and is not subject to copyright.

## Acknowledgments

## Biographies

Tingting Hu, PhD, is a visiting scientist in the Division of Imaging, Diagnostics, and Software Reliability, Center for Devices and Radiological Health, U.S. Food and Drug Administration. She received her PhD in statistics from Florida State University. Her research interests include statistics.

Berkman Sahiner is a senior biomedical research scientist with the Division of Imaging, Diagnostics and Software Reliability, Center for Devices and Radiological Health, U.S. Food and Drug Administration. He has a PhD in electrical engineering and computer science from the University of Michigan, Ann Arbor. His research is focused on the evaluation of medical imaging and computer-assisted diagnosis devices, including devices that incorporate machine learning and artificial intelligence. He is a fellow of SPIE, AAPM and AIMBE.

Shuyue Guan, PhD, a current staff fellow in the Division of Imaging, Diagnostics, and Software Reliability, Center for Devices and Radiological Health, U.S. Food and Drug Administration. PhD received in Biomedical Engineering from the George Washington University. Primary research interests are image processing and machine learning.

Nicholas Petrick, Ph.D. has expertise in medical artificial intelligence development and assessment. He is currently the Deputy Director for the Division of Imaging, Diagnostics, and Software Reliability in the Center for Devices and Radiological Health, U.S. Food and Drug Administration and is a member of FDA's Senior Biomedical Research and Biomedical Product Assessment Service.

## Code, Data, and Materials Availability

The data used in the application section—LIDC-IDRI images—are publicly available at:

## Figure Caption List

Figure 1.CP, Type I error and Type II error results for the statistics-based Simulation Study 1. The figure includes boxplot results aggregated across all the settings in Scenarios I – IV, and for both the 2-reader and 3-reader scenarios.

Figure 2. Example of synthetic expert and device contours based on an initial true contour. Note, this is a Transformation Disagreeable example where the device contour is based on a different set of MISS parameters compared to that of the experts. a) true contour, b)-d) synthetic reader contours, e) Transformation Disagreeable device contour.

Figure 3. Plots of the Type 1 error results by transformation pattern (Rmax, Smax) for Image Simulation Study 2, Scenario I: Transformation Agreeable. The black dashed horizontal line marks the expected type 1 error level of 0.05.

Figure 4.Type 2 error by device transformation pattern ($R_{max}^D, S_{max}^D$), given panel transformation ($R_{max}, S_{max}$) = $(1.10, 2)$, with mean(std) within-panel Dice 0.87(0.09) for image simulation study Scenario II: Transformation Disagreeable. The black dashed horizontal line marks desirable type 2 error level at 0.2.

Figure 5. Distribution of DSC values for each pair of annotators.

Figure A1: Radiologists annotations of an LIDC-IDRI lesion along with the majority vote (MV) consensus, and the contour produced by our U-Net segmentation algorithm.

## Table Caption List

Table 1: Parameter configuration scenarios for Study 1: statistic-based simulations.

Table 2. Results for all 32 settings in Scenario I: Equal Mean μ and Equal Standard Deviation σ from statistics-based simulation Study 1.

Table 3. Selected results for 8 Settings with $(\mu_0, d_\mu, \sigma_0) = (0.85, -0.05, 0.15)$ in Scenario II: Unequal Mean μ and Equal Standard Deviation σ from statistics-based Simulation Study 1.

Table 4. Selected results for 24 of 96 Settings with $\mu_0 = 0.75$ in Scenario III: Equal Mean μ and Unequal Standard Deviation σ from statistics-based Simulation Study 1.

Table 5. Selected results for 24 of 648 Settings with $(\mu_0, d_\mu, \sigma_0) = (0.8, -0.05, 0.1)$ in Scenario IV: Unequal Mean μ and Unequal Standard Deviation σ from statistics-based Simulation Study 1.

Table 6. Configuration for tunable parameters in image-based Simulation Study 2.

646 Table 7. Mean (standard deviation) for Within-Panel DSC, Device-Panel DSC, and Type 1 Error by panel
647 transformation pattern from Image-Simulation Study 2, Scenario I: Transformation Agreeable.

648 Table 8. Mean (standard deviation) of Within-Panel DSC, Device-Panel DSC, and Type 2 Error by Panel
649 Transformation Pattern, Image-Simulation Study 2, Scenario II: Transformation Disagreeable.

650 Table 9. Agreement assessment results for the U-Net model (device) and the expert panel using the LIDC-IDRI test
651 dataset.

## References

653 1. Taha, A. A., & Hanbury, A. (2015). Metrics for evaluating 3D medical image segmentation: analysis,

654    selection, and tool. BMC medical imaging, 15(1), 1-28.

655 2. Wang, Z., Wang, E., & Zhu, Y. (2020). Image segmentation evaluation: a survey of

656    methods. Artificial Intelligence Review, 53(8), 5637-5674.

657 3. Udupa, J. K., LeBlanc, V. R., Zhuge, Y., Imielinska, C., Schmidt, H., Currie, L. M., ... & Woodburn,

658    J. (2006). A framework for evaluating image segmentation algorithms. Computerized medical

659    imaging and graphics, 30(2), 75-87.

660 4. Cárdenes, R., de Luis-Garcia, R., & Bach-Cuadra, M. (2009). A multidimensional segmentation

661    evaluation for medical image data. Computer methods and programs in biomedicine, 96(2), 108-124.

662 5. Dice, L. R. (1945). Measures of the amount of ecologic association between species. Ecology, 26(3),

663    297-302.

664 6. Jaccard, P. (1912). The distribution of the flora in the alpine zone. 1. New phytologist, 11(2), 37-50.

665 7. Zou, K. H., Warfield, S. K., Bharatha, A., Tempany, C. M., Kaus, M. R., Haker, S. J., ... & Kikinis,

666    R. (2004). Statistical validation of image segmentation quality based on a spatial overlap index1:

667    scientific reports. Academic radiology, 11(2), 178-189.

668 8. Bland, J. M., & Altman, D. (1986). Statistical methods for assessing agreement between two methods

669    of clinical measurement. The lancet, 327(8476), 307-310.

670 9. FDA (2001). Guidance for industry: Statistical approaches to establishing bioequivalence, Food and

671    Drug Administration, Center for Drug Evaluation and Research (CDER).

672     10. Obuchowski, N. A. (2001). Can electronic medical images replace hard-copy film? Defining and

673        testing the equivalence of diagnostic tests. Statistics in medicine, 20(19), 2845-2863.

674     11. Obuchowski, N. A., Subhas, N., & Schoenhagen, P. (2014). Testing for interchangeability of imaging

675        tests. Academic Radiology, 21(11), 1483-1489.

676     12. Shao, J., & Zhong, B. (2004). Assessing the agreement between two quantitative assays with repeated

677        measurements. Journal of Biopharmaceutical Statistics, 14(1), 201-212.

678     13. Barnhart, H. X., Kosinski, A. S., & Haber, M. J. (2007). Assessing individual agreement. Journal of

679        biopharmaceutical statistics, 17(4), 697-719.

680     14. FDA CDER (2022). Guidance for industry: Multiple Endpoints in Clinical Trials.

681        https://www.fda.gov/regulatory-information/search-fda-guidance-documents/multiple-endpoints-

682        clinical-trials

683     15. Guan, S., Samala, R. K., Arab, A., & Chen, W. (2023, April). MISS-tool: medical image

684        segmentation synthesis tool to emulate segmentation errors. In Medical Imaging 2023: Computer-

685        Aided Diagnosis (Vol. 12465, pp. 273-281). SPIE.

686     16. Wang, S., Zhou, M., Liu, Z., Liu, Z., Gu, D., Zang, Y., ... & Tian, J. (2017). Central focused

687        convolutional neural networks: Developing a data-driven model for lung nodule

688        segmentation. Medical image analysis, 40, 172-183.

689     17. Joskowicz, L., Cohen, D., Caplan, N., & Sosna, J. (2019). Inter-observer variability of manual

690        contour delineation of structures in CT. European radiology, 29, 1391-1399.

691     18. Evans, J. D. (1996). Straightforward statistics for the behavioral sciences. Thomson Brooks/Cole

692        Publishing Co.

693     19. LaMorte W. W. The Correlation Coefficient (r) [Internet]. [Accessed January 2024]. Boston

694        University School of Public Health. Available at: https://sphweb.bumc.bu.edu/otlt/MPH-

695        Modules/PH717-QuantCore/PH717-Module9-Correlation-Regression/PH717-Module9-Correlation-

696        Regression4.html.

697    20. Swinscow, T. D. V., & Campbell, M. J. (2002). Statistics at square one (pp. 111-25). London: Bmj.

698         Available at: https://www.bmj.com/about-bmj/resources-readers/publications/statistics-square-

699         one/11-correlation-and-regression.

700    21. Whuber (https://stats.stackexchange.com/users/919/whuber), How to construct a multivariate Beta

701         distribution?, URL (version: 2021-10-21): https://stats.stackexchange.com/q/549262

702    22. L'Ecuyer, P. (2015, December). Random number generation with multiple streams for sequential and

703         parallel computing. In 2015 Winter Simulation Conference (WSC) (pp. 31-44). IEEE.

704    23. Armato III, S. G., McLennan, G., Bidaut, L., McNitt-Gray, M. F., Meyer, C. R., Reeves, A. P., ... &

705         Clarke, L. P. (2011). The lung image database consortium (LIDC) and image database resource

706         initiative (IDRI): a completed reference database of lung nodules on CT scans. Medical

707         physics, 38(2), 915-931.

708    24. FDA CDRH (2022). Guidance for Industry and Food and Drug Administration Staff Computer-

709         Assisted Detection Devices Applied to Radiology Images and Radiology Device Data - Premarket

710         Notification [510(k)] Submissions. https://www.fda.gov/media/77635/download

711    25. Bandos, A. I., Rockette, H. E., & Gur, D. (2006). A permutation test for comparing ROC curves in

712         multireader studies: a multi-reader ROC, permutation test. Academic radiology, 13(4), 414-420.

713    26. Hillis, S. L., & Schartz, K. M. (2018). Multireader sample size program for diagnostic studies:

714         demonstration and methodology. Journal of Medical Imaging, 5(4), 045503-045503.
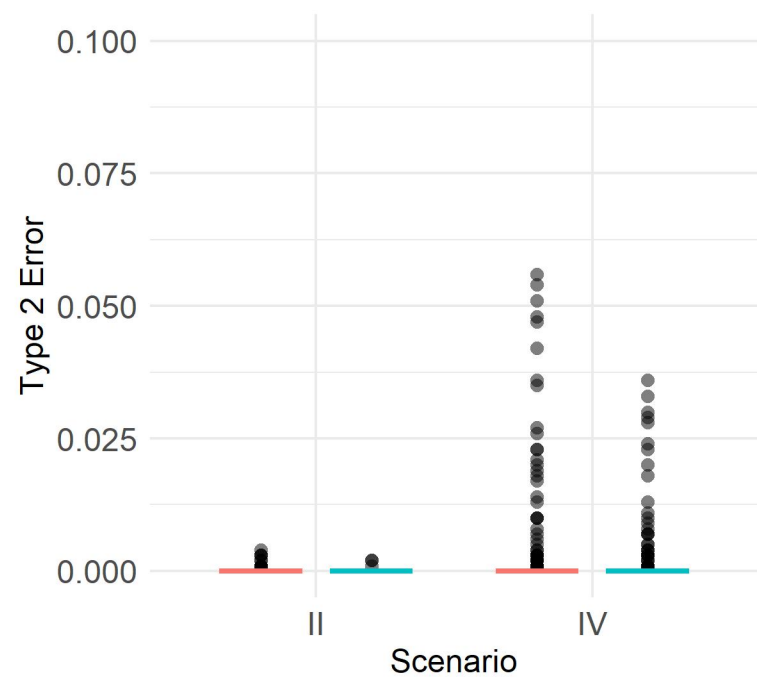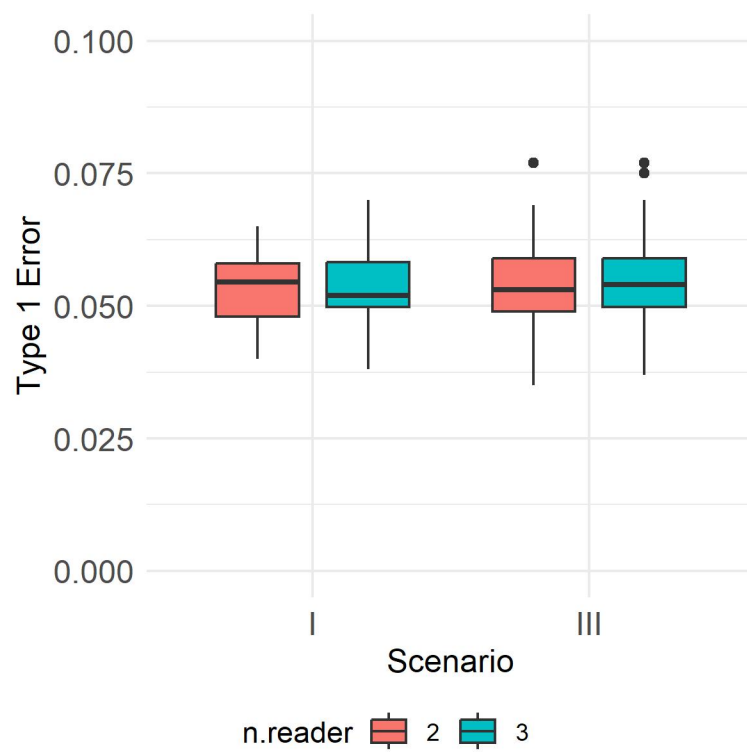
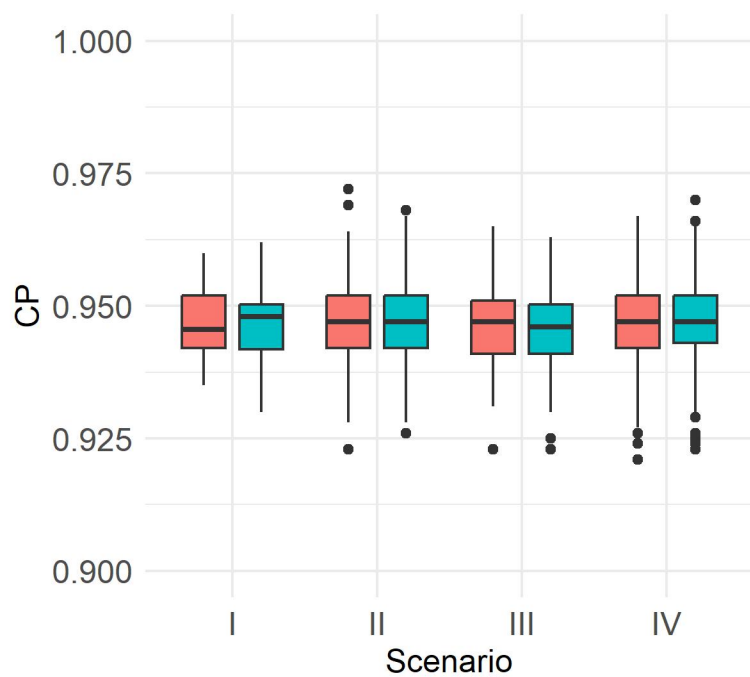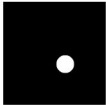Expert 0    Expert 1    Expert 2    Expert 3

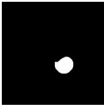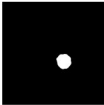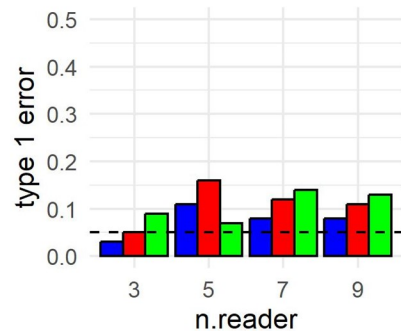Input Slice    MV Consensus    U-Net Extractor
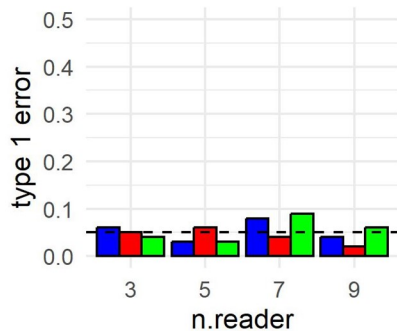
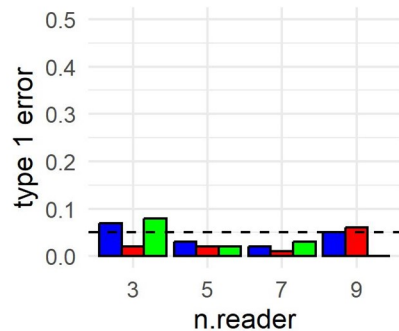a) truth     b) expert 1     c) expert 2     d) expert 3     e) device

Top row of panels:
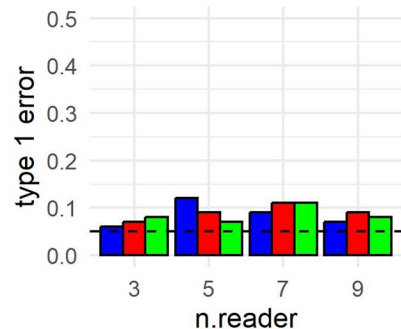- $R_{max} = 1.05$ $S_{max} = 0$ $\hat{\mu}_0 = 0.95$ $\hat{\sigma}_0 = 0.04$
- $R_{max} = 1.05$ $S_{max} = 2$ $\hat{\mu}_0 = 0.88$ $\hat{\sigma}_0 = 0.09$
- $R_{max} = 1.05$ $S_{max} = 3$ $\hat{\mu}_0 = 0.84$ $\hat{\sigma}_0 = 0.13$

Middle row of panels:
- $R_{max} = 1.1$ $S_{max} = 0$ $\hat{\mu}_0 = 0.93$ $\hat{\sigma}_0 = 0.04$
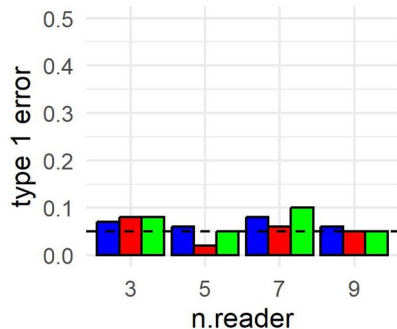- $R_{max} = 1.1$ $S_{max} = 2$ $\hat{\mu}_0 = 0.87$ $\hat{\sigma}_0 = 0.09$
- $R_{max} = 1.1$ $S_{max} = 3$ $\hat{\mu}_0 = 0.83$ $\hat{\sigma}_0 = 0.12$
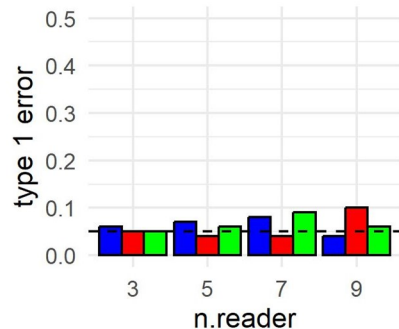
Bottom row of panels:
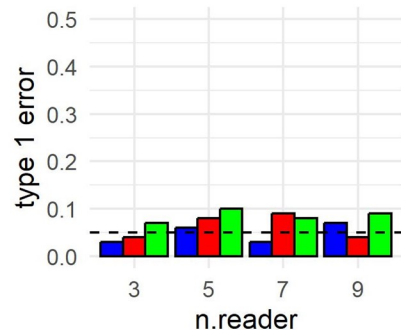- $R_{max} = 1.15$ $S_{max} = 0$ $\hat{\mu}_0 = 0.9$ $\hat{\sigma}_0 = 0.05$
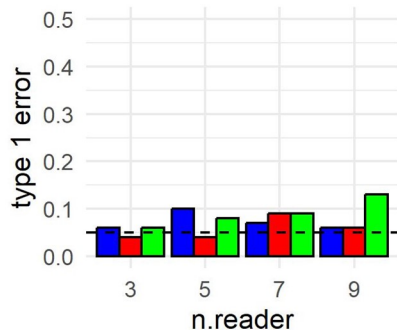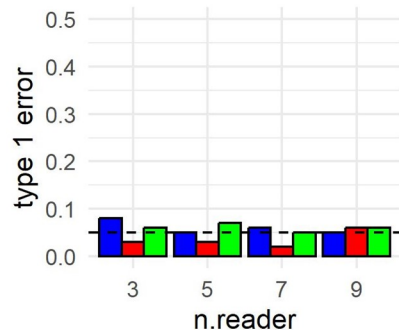- $R_{max} = 1.15$ $S_{max} = 2$ $\hat{\mu}_0 = 0.85$ $\hat{\sigma}_0 = 0.09$
- $R_{max} = 1.15$ $S_{max} = 3$ $\hat{\mu}_0 = 0.81$ $\hat{\sigma}_0 = 0.12$

Axis labels: type 1 error (y-axis); n.reader (x-axis)

n.image: 100, 250, 750

Top row, left to right:
$R^D_{max} = 1.05$ $S^D_{max} = 0$ $\hat{\mu}_D = 0.9$ $\hat{\sigma}_D = 0.06$

$R^D_{max} = 1.05$ $S^D_{max} = 2$ $\hat{\mu}_D = 0.87$ $\hat{\sigma}_D = 0.09$

$R^D_{max} = 1.05$ $S^D_{max} = 3$ $\hat{\mu}_D = 0.85$ $\hat{\sigma}_D = 0.11$

$R^D_{max} = 1.1$ $S^D_{max} = 0$ $\hat{\mu}_D = 0.89$ $\hat{\sigma}_D = 0.06$

Bottom row, left to right:
$R^D_{max} = 1.1$ $S^D_{max} = 3$ $\hat{\mu}_D = 0.85$ $\hat{\sigma}_D = 0.1$

$R^D_{max} = 1.15$ $S^D_{max} = 0$ $\hat{\mu}_D = 0.88$ $\hat{\sigma}_D = 0.06$

$R^D_{max} = 1.15$ $S^D_{max} = 2$ $\hat{\mu}_D = 0.86$ $\hat{\sigma}_D = 0.09$

$R^D_{max} = 1.15$ $S^D_{max} = 3$ $\hat{\mu}_D = 0.84$ $\hat{\sigma}_D = 0.1$

Axis labels: type 2 error (y-axis), n.reader (x-axis)

n.image ■ 100 ■ 250 ■ 750

Dice

AI vs. reader 0, AI vs. reader 1, AI vs. reader 2, AI vs. reader 3, reader 0 vs. 1, reader 0 vs. 2, reader 0 vs. 3, reader 1 vs. 2, reader 1 vs. 3, reader 2 vs. 3

■ AI vs. reader   ■ reader vs. reader