

Quantifying Inter-Operator Variability and its Causes for Medical Semantic Segmentation

Anonymous ICCV submission

Paper ID 9

Abstract

001 *Uncertainty quantification is a vital component of the de-*
 002 *velopment of models for medical semantic segmentation, as*
 003 *we are increasingly faced with the question “when is ‘good’*
 004 *good enough?”. A necessary part of the answer to this*
 005 *will always involve understanding the variability of accept-*
 006 *able segmentation results, which in turn requires the inter-*
 007 *operator variability to be disclosed. This is a stumbling*
 008 *block for many, as quantifying inter-operator variability is*
 009 *expensive and time consuming - especially within medical*
 010 *imaging. In this work we not only highlight how critical*
 011 *understanding this variability is, but also provide a novel*
 012 *framework to better understand and predict regions of vari-*
 013 *ability by considering the texture of the scan. The source*
 014 *code will be released upon publication.*

015 1. Introduction

016 1.1. Semantic Segmentation

017 In computer vision, semantic segmentation is the task of
 018 determining which pixels in a given image belong to a tar-
 019 get class. Within medical imaging the output of semantic
 020 segmentation can be used to take measurements [12], de-
 021 tect abnormalities [16], and construct virtual copies of an
 022 object [11]. However, ground truth segmentations are of-
 023 ten time consuming (and therefore expensive) for a clini-
 024 cian to produce manually [13]. This has necessitated the
 025 development of segmentation models to reduce the amount
 026 of clinician input needed to generate segmentations [2]. A
 027 key difference between how a clinician produces a segmen-
 028 tation versus a standard segmentation model is the use of
 029 the object boundary. A clinician will define the object by its
 030 boundary; for them it is clear which spaces are in the object
 031 and which are not. In contrast, most models will define the
 032 object as a collection of pixels, with no regard for which
 033 ones form the boundary [17].

034 Inter-operator variability is an accepted issue across se-
 035 mantic segmentation, whilst also providing a helpful mar-

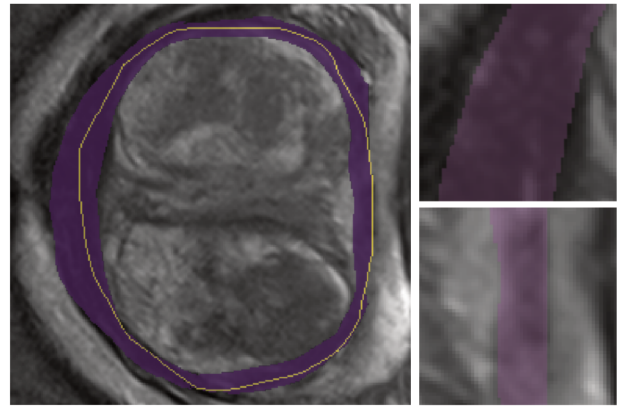


Figure 1. *Left:* An example segmentation (yellow) within the area of other clinician segmentations (purple). *Top Right:* A magnified region with high variability. *Bottom Right:* A magnified region with low variability.

gin of acceptable error for models. Despite the importance
 of annotator variability, relatively few publicly available
 datasets include variability statistics in their documentation.
 When available, variability is often stated as a single mea-
 surement (e.g. 0.05 DSC), with no further information on
 how certain parts of a segmentation may be prone to higher
 variability than others.

This paper has three aims:

1. To highlight how prevalent inter-operator variability is within semantic segmentation data.
2. To demonstrate how the distribution of this variability is often uneven.
3. To provide a framework for quantifying both the variability and its causes within a single image.

Predicting variability via models is a well-explored area; however, these methods often struggle when trying to validate their results against actual recorded uncertainty. The ideal method to validate such a model would be to compare it against the variability between a group of annotators; unfortunately, this is frequently unrealistic due to the cost of annotation. Instead, if we could understand the causes of

057 inter-operator variability then we may be able to assess the
058 performance of these models without the need for data fully
059 annotated by multiple clinicians.

060 1.2. Inter-Operator Variability

061 Inter-operator variability refers to the variability between
062 different annotators; in contrast to intra-operator variabil-
063 ity, which refers to the variability within the segmentations
064 done by the same annotator. Reporting inter-operator vari-
065 ability is particularly important if the dataset contains seg-
066 mentations drawn by multiple distinct annotators, as ‘er-
067 rors’ by the model may actually just be differences in an-
068 notator style.

069 As a basis for this work we considered the 18 seman-
070 tic segmentation challenges run in the MICCAI 2023 con-
071 ference, focusing on the number of annotators providing
072 the segmentations and the reported inter-operator variabil-
073 ity. This conference was chosen not only because it had a
074 large number of semantic segmentation challenges for med-
075 ical imaging, but the high-profile nature of the conference
076 meant that the datasets used received a lot of publicity, with
077 the prize money totaling in excess of \$40K.

078 All of the challenges had a single annotation per image,
079 but many had split these between multiple annotators. Of
080 the 18 challenges, only 2 had a single clinician providing
081 all annotations and more than half of the challenges had 5
082 or more clinicians providing annotations [3]. In the chal-
083 lenges where the exact number of annotators was not given,
084 it was made clear that there were still multiple annotators.
085 This means that we can expect inter-operator variability to
086 be present in at least 16 of the challenges.

087 Despite this, only 6 of these 16 challenges quantified
088 the annotator uncertainty in some way, and of those the
089 uncertainty was often described in a vague or unhelpful
090 manner. For example, the autoPET 2023 challenge only
091 describes the uncertainty in defining lesion boundaries by
092 saying “*The difference in segmentation volumes can range*
093 *from 5-30%*” [7].

094 The datasets used in these challenges will likely go on to
095 be used in many further studies, yet most provide no warn-
096 ing that multiple annotators were used or state the likely
097 variability in the ground truths. As it would be unrealistic
098 to expect each image to be annotated by multiple clinicians,
099 we will instead provide a method to allow single annota-
100 tions to reflect the uncertainty in the boundary.

101 2. Existing Work

102 2.1. Predicting Uncertainty

103 Previous work in predicting which areas of a ground truth
104 are likely to be the most uncertain often involves methods
105 such as training a model multiple times with different pa-
106 rameters, such as MC Dropout which uses different dropout

107 conditions [5]. The variance in the results on the same im-
108 age can then be used as a prediction for the uncertainty of
109 that model. Another example is CoraNet [15], in which the
110 model is trained with a conservative setting and a radical
111 setting to determine which regions are the most likely to be
112 uncertain.

113 Other methods consider that different annotators will
114 have different styles and that some annotators may be more
115 senior and therefore should be considered ‘more correct’.
116 An example of this can be seen in the work by Ji et al. [9],
117 which looked to calibrate multiple annotations with exper-
118 tise levels in order to produce a model that could emulate
119 an expert annotator.

120 2.2. Annotation Reliability

121 The standard way to increase the reliability of an annota-
122 tion is to have the annotation repeated, either by the same
123 clinician (to increase intra-operator reliability) or by mul-
124 tiple clinicians (to increase inter-operator reliability). The
125 annotations are then often combined into a single annota-
126 tion, by methods such as majority voting [6] or the STA-
127 PLE algorithm, which forms a probabilistic estimate of the
128 ‘true’ ground truth and weights the annotator segmentations
129 accordingly [18].

130 3. Method

131 3.1. Data

132 In order to properly assess the relationship between vari-
133 ability, noise, and segmentation regions, we searched for
134 publicly available semantic segmentation data with at least
135 three annotations per image. The resulting datasets were
136 then reduced to those of universal physiological features
137 (i.e. not lesions, tumors, or other features that have no set
138 location). The reason for this was that we wanted to focus
139 on tasks where the target segmentation had a fixed known
140 location, and therefore have neighboring structures in the
141 same place in each image.

142 This search resulted in just three datasets: the QUBIQ
143 Prostate dataset, the RIGA fundus dataset, and the Grey
144 Matter (GM) Spinal dataset. Each of these datasets has two
145 tasks per image, resulting in six unique tasks overall. Ex-
146 amples from each dataset can be seen in Figure 2.

147 3.1.1. QUBIQ Prostate

148 The QUBIQ Prostate data is a dataset of prostate MRI im-
149 ages, made available for the QUBIQ challenge [10]. It con-
150 tains 48 2D MRI slices, where each slice has been annotated
151 for two tasks by an expert clinician. The challenge does not
152 describe what the two tasks are segmenting, but it appears
153 that Task 1 is the total prostate volume and Task 2 is the
154 transitional zone volume [19].

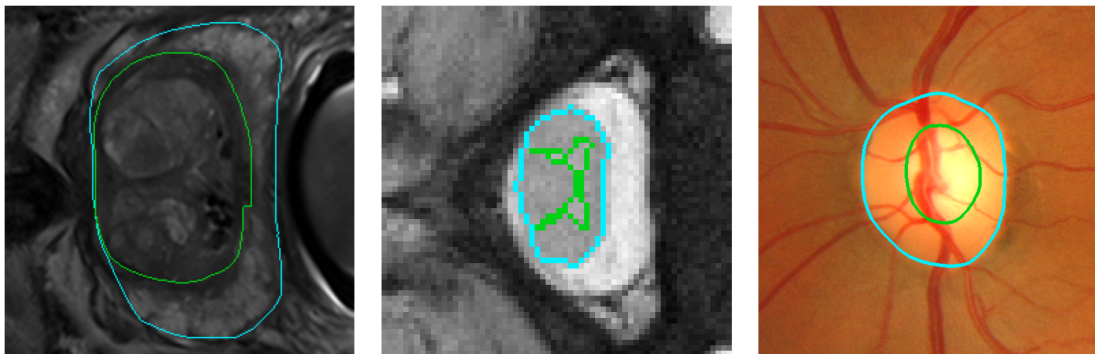


Figure 2. Example segmentations for both tasks for each of: QUBIQ Prostate (Left), GM Spinal (Centre), RIGA (Right)

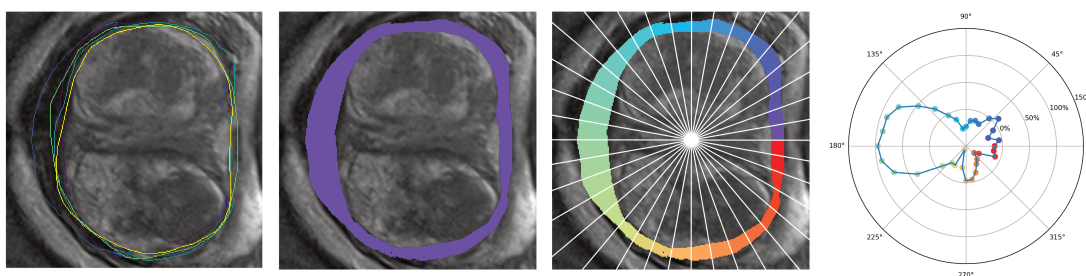


Figure 3. *Far Left*: A scan from the second QUBIQ Prostate task with multiple annotations. *Center Left*: The area disputed between annotators. *Center Right*: The disputed area split into 36 boundary regions. *Center Right*: The corresponding graph showing the change in disputed area width, as described in Section 3.4.1.

155 **3.1.2. RIGA**

156 The RIGA dataset contains 750 fundus images from three
157 different sources [1]. Each fundus image has the optic disc
158 (Task 1) and optic cup (Task 2) annotated by six experi-
159 enced ophthalmologists. Of the three sources (MESSIDOR,
160 Bin Rushed, and Magrabi), the MESSIDOR subset has 460
161 images, the Bin Rushed set has 195, and the Magrabi set
162 has 95. The Bin Rushed images are difficult to use due
163 to the ground truths being compressed on top of the fundus
164 images as single JPEG images, so we removed this subset
165 from the analysis.

166 **3.1.3. GM Spinal**

167 The GM Spinal Cord Segmentation Challenge data was re-
168 leased in 2017, containing 3D MRI data for 80 patients from
169 4 sites, resulting in a total of 407 2D slices [14]. Each
170 slice was annotated by four expert raters, who manually seg-
171 mented the white matter (Task 1) and gray matter (Task 2).
172 The software used to assist the annotations varied between
173 annotators, and is described in the GM Spinal publication.

174 **3.1.4. Pre Processing**

175 Each image was converted to 256 levels, with values nor-
176 malized within this range. No other pre-processing was ap-
177 plied.

3.2. Disputed and Agreed Areas

In order to assess how annotator variability changes across
a segmentation, we first need to establish which parts of the
segmentation are contributing towards the variability.

For a set of annotations A we consider the *agreed seg-
mentation*, $AS(A)$ and the *disputed segmentation* $DS(A)$.

The agreed area is the intersection of all annotations:

$$AS(A) = \bigcap_{a_i \in A} a_i$$

The disputed area is the area that is included by at least
one annotator, but not all annotators:

$$DS = \bigcup_{a_i \in A} a_i - AS(A)$$

Finally, we will also refer to the *total segmentation*
 $TS(A)$, which is the union of the disputed and agreed ar-
eas.

An example of the disputed area can be seen in Figure 3.

3.3. Boundary Regions

In order to compare variability between the same areas, we
split the disputed area into n non-overlapping regions. The

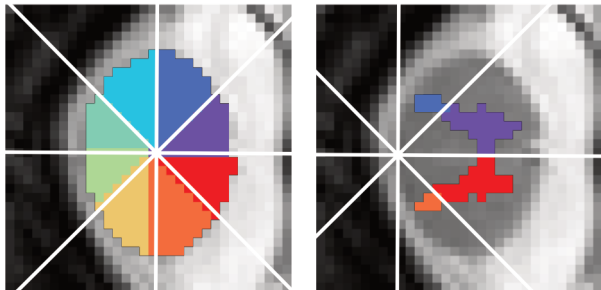


Figure 4. *Left*: An example of 8 regions for the GM Spinal Task 1. *Right*: An example of 8 regions for the GM Spinal Task 2, with centroid aligned to the left to allow more meaningful separation.

exact value of n is determined using the minimum boundary length across the entire dataset, with the constraint that each region must always contain at least one pixel.

The region each pixel p belongs to is determined by the angle between p and the centroid c of the agreed area for all tasks except Task 2 of the GM Spinal data. In this task, the concave shape of the gray matter meant that more meaningful regions were created by instead using a left aligned centroid, as seen in Figure 4.

For given n and c the regions of $\text{disp}(A)$ are then defined as:

$$R_{0 \leq i < n} = \left\{ p \in \text{disp}(A) \mid i \leq \frac{1}{n} \cdot \tan^{-1} \frac{p_y - c_y}{p_x - c_x} < i + 1 \right\}$$

An example of this for 36 regions can be seen in Figure 3.

This method has some weaknesses; regions with a longer medial axis (i.e. higher irregularity) are greater represented than those with a shorter medial axis, and if the medial axis is not convex then multiple disconnected segments of the axis could potentially be included in the same region. A more nuanced method would be an improvement for future work.

For an object with genus 0 the disputed area almost always appears as a band around the boundary of the object, as the topology of the object is generally not disputed. An exception to this is seen in the GM Spinal data, where the segmentations are so small that the disputed area is often only 1 pixel wide, and interrupted by sections with no disputed area.

In the simple case where the disputed area is an uninterrupted band, we can measure the width by iterating across the medial axis of the disputed area: the width is calculated as the shorted distance to the boundary of the disputed area. In cases where the medial axis is a set of disconnect segments, we instead join them with the corresponding sections

of the boundary of the agreed area (each of these pixels has width 0).

Using these regions, we can investigate whether particular regions of a segmentation have consistently higher inter-operator variability or changes in texture.

3.4. Inter-Operator Variability

3.4.1. Disputed Area Width

To measure the average width of the disputed area for a region we consider the width for each point $p \in \text{RM}_i$, where RM_i is the intersection of region R_i and the medial axis. The width at p is the shortest distance to the boundary of the disputed area.

Rather than the absolute width, we instead measure the percentage change between the width w_p at point p and the mean width \bar{w} . The mean change in width ΔW for region R_i is then:

$$\Delta W_{R_i} = \frac{100}{|\text{RM}_i|} \sum_{p \in \text{RM}_i} w_p - \bar{w}$$

A region with a high change in width will have a higher inter-operator variability than a region in the same image with a lower change in width.

3.4.2. Dice Similarity Coefficient

The second method we will use to measure the inter-operator variability between regions is the change in the Dice Similarity Coefficient (DSC) [4]. DSC is commonly used in medical imaging semantic segmentation as a way to measure the similarity of two segmentations, and as such is also often used to record inter-operator variability.

The value of DSC for a region R_i measures the similarity between the total segmentation TS and the agreed segmentation AS :

$$\text{TS}_{R_i} = R_i \cap \text{TS} \tag{1}$$

$$\text{AS}_{R_i} = R_i \cap \text{AS} \tag{2}$$

$$DSC(R_i) = \frac{2 |\text{AS}_{R_i}|}{|\text{AS}_{R_i}| + |\text{TS}_{R_i}|} \tag{3}$$

We then evaluate ΔDSC as the percentage difference against the value of DSC for all regions ($DSC(\cup R)$):

$$\Delta DSC(R_i) = \frac{100 DSC(R_i)}{DSC(\cup R)}$$

3.5. Texture

The aim for this paper is to compare inter-operator variability to a textures measures that could be obtained from a dataset with just one annotation per image. We therefore



Figure 5. *Left*: A GM Spinal image with patch in red. *Center*: The magnified patch. *Right*: The corresponding patch GLCM.

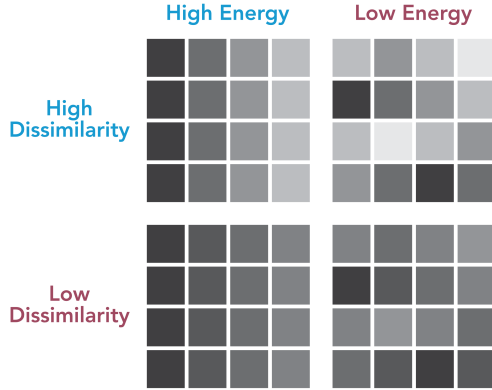


Figure 6. Example patches showing combinations of high and low energy and dissimilarity, for a GLCM with comparison distance of 1 pixel and horizontal direction. Within each patch the difference in intensity between each pixel and its horizontal neighbour is the same.

273 measure the texture for each point p along the boundary of
274 each individual segmentation for the six tasks.

275 For each point p on the segmentation boundaries we assess
276 the texture using a Gray Level Co-occurrence Matrix
277 (GLCM) for a patch P of size m centered around p [8].
278 The value of m is dependent on the smallest segmentation
279 area in the dataset for that task. The GLCM identifies the
280 frequency that each intensity level occurs within the patch
281 with a distance of d and along an axis of θ degrees. This
282 matrix is then used to compute the energy and dissimilarity
283 values.

284 The value of θ for a point p is decided using the angle
285 of p from the centroid c , as described in Section 3.3. This
286 approximately gives a value of θ that is perpendicular to the
287 boundary, in order to identify any clear boundaries in the
288 original scan. An example patch and corresponding GLCM
289 matrix can be seen in Figure 5.

290 3.5.1. Energy

291 The energy of the GLCM measures how orderly the pixel
292 value differences are within the patch [8]. A GLCM with
293 high energy will have pixels of similar intensities close to

294 each other, whereas one with low energy will have pixels of
295 different intensities ‘mixed’ with each other. For a point p
296 with GLCM matrix P , the energy is then:

$$297 \quad e(p) = \sum_{i,j=0}^{levels-1} P_{i,j}(-\log(P_{i,j}))$$

298 For a segmentation boundary region R_i , the mean change
299 in energy is calculated against the mean segmentation en-
300 ergy \bar{e} :

$$301 \quad \Delta E_{R_i} = \frac{100}{|R_i|} \sum_{p \in R_i} e_p - \bar{e}$$

302 Examples of patches with low versus high energy can be
303 seen in the rows of Figure 6.

304 3.5.2. Dissimilarity

305 Dissimilarity is a measure of contrast within the patch [8].
306 A GLCM with high dissimilarity will have low contrast, re-
307 gardless of the order of the pixels or the absolute values
308 compared to the rest of the image. For a point p with GLCM
309 matrix P , the dissimilarity is calculated as:

$$310 \quad d(p) = \sum_{i,j=0}^{N-1} P_{i,j} |i - j|$$

311 Similarly, for a segmentation boundary region R_i the mean
312 change in dissimilarity is calculated against the mean seg-
313 mentation dissimilarity \bar{d} :

$$314 \quad \Delta D_{R_i} = \frac{100}{|R_i|} \sum_{p \in R_i} d_p - \bar{d}$$

315 Examples of patches with low versus high dissimilarity can
316 be seen in the columns of Figure 6. In these examples we
317 can also see how the combination of energy and dissimilarity
318 affects how easily a boundary can be identified. Patches
319 with a continuous edge that runs perpendicular to θ will
320 have a higher energy, and this edge will be easier to detect
321 if the dissimilarity is higher. In areas with low dissimilarity
322 any boundaries will be much harder to make out – an indi-
323 cator that the patch is monotone, with no clear boundaries,
324 is high energy with low dissimilarity.

325 4. Results

326 4.1. Overall Variability

327 The mean size and standard deviation for the total segmen-
328 tation area can be seen in Table 1, which shows the sub-
329 stantial difference in size between the QUBIQ and RIGA
330 tasks versus the GM Spinal tasks. The RIGA tasks also
331 have a much larger standard deviation of the total segmen-
332 tation area, corroborating the need for variation measures
333 that do not use absolute pixel values.

	QUBIQ Prostate		RIGA		GM Spinal	
	Task 1	Task 2	Task 1	Task 2	Task 1	Task 2
Mean	44075.4	28901.3	52452.7	18523.8	990.5	527.2
S.D.	14339.8	14185.5	66009.5	25454.6	194.7	93.5

Table 1. Total segmentation area in pixels (mean and standard deviation) for all tasks.

	QUBIQ Prostate		RIGA		GM Spinal	
	Task 1	Task 2	Task 1	Task 2	Task 1	Task 2
Mean	20.0%	39.0%	17.3%	53.6%	13.8%	45.7%
S.D.	4.4%	13.8%	5.5%	15.5%	3.7%	6.6%

Table 2. Proportion of disputed area (mean and standard deviation) for all tasks.

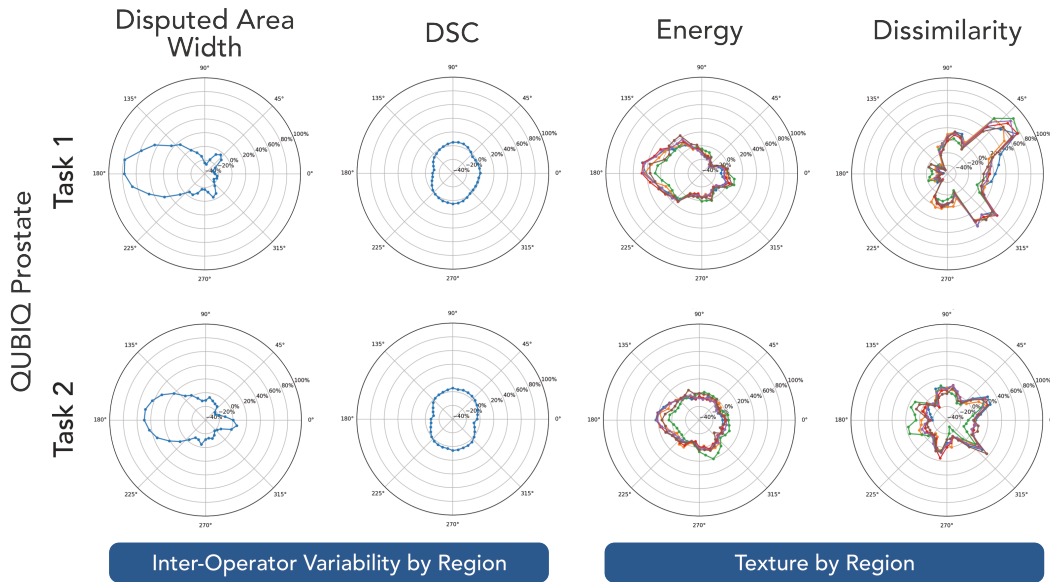


Figure 7. Inter-operator variability and texture measures by region for the QUBIQ Prostate Tasks

334 Table 2 then shows the overall size of the disputed area
 335 as a proportion of the total segmentation for all six tasks.
 336 Although this measurement does not show how variability
 337 changes across regions of a segmentation, we can see that
 338 across all six tasks the disputed area was a significant pro-
 339 portion of the total segmentation.

340 An interesting observation is that the second task for
 341 each of the three datasets had a notable increase in the pro-
 342 portion of the disputed area. The second feature for these
 343 datasets was always entirely contained within the first task
 344 (and therefore smaller). The cause for this increase in vari-
 345 ability could theoretically be because annotators made seg-
 346 mentations for both tasks at the same time, and did not zoom
 347 in on the second smaller task.

348 Finally, Table 3 shows the DSC values between the
 349 agreed and total segmentation areas for all tasks. The ob-
 350 servations here are similar to that of table 2, with the second

tasks having notably lower DSC scores.

4.2. Variability and Texture by Region

351 For each of the three datasets, we report the disputed area
 352 width, DSC, energy, and dissimilarity as a mean value for
 353 each region. These results can be seen in Figures 7, 8, and
 354 9 as polar graphs, with the angular coordinate representing
 355 the angle of each region from the centroid and the radial
 356 coordinate representing the measured value.
 357
 358

4.2.1. QUBIQ Prostate Data

359 The results for the QUBIQ Prostate data can be seen in Fig-
 360 ure 7. For both tasks there is a clear increase in variability
 361 in the center left regions, with an average increase in width
 362 of up to 80% in Task 1 and 50% in Task 2. This corresponds
 363 to a decrease in region DSC of up to -10% in both cases (as
 364 Task 1 is of a larger volume).
 365

	QUBIQ Prostate		RIGA		GM Spinal	
	Task 1	Task 2	Task 1	Task 2	Task 1	Task 2
Mean	0.882	0.792	0.899	0.608	0.925	0.705
S.D.	0.029	0.136	0.036	0.154	0.021	0.055

Table 3. DSC between agreed area and total segmentation area (mean and standard deviation) for all tasks.

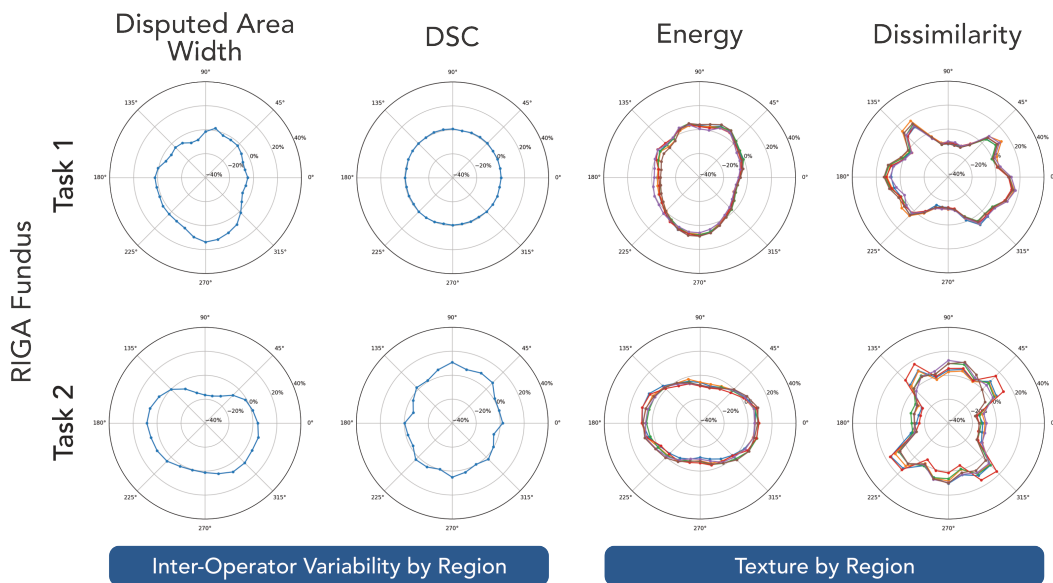


Figure 8. Inter-operator variability and texture measures by region for the RIGA Fundus Tasks

366 For both tasks, the increase in width in the center right
 367 regions corresponds to a similar increase in energy for the
 368 same regions, as well as a decrease in dissimilarity. As men-
 369 tioned in Section 3.5.2 the combination of high energy with
 370 low dissimilarity indicates a monotone region, which ap-
 371 pears to be the case here.

372 **4.2.2. RIGA Fundus Data**

373 The results for the RIGA fundus data can be seen in Figure
 374 8. As with the QUBIQ data a similar magnitude of differ-
 375 ence in width (+15% for Task 1 and -18% in Task 2) results
 376 in a much greater effect in region DSC for the task with a
 377 smaller overall volume.

378 A contributing factor to variability that is not considered
 379 here is blood vessels obscuring the segmentation bound-
 380 aries. Though the blood vessels often run perpendicular
 381 to the boundary, they may still affect the texture measure-
 382 ments.

383 **Task 1** The optic disc task has a region of high energy
 384 and low dissimilarity at the center bottom region, corre-
 385 sponding to a higher disputed area width. Though this area
 386 is less monotone than in the QUBIQ data, it once again cor-
 387 relates to an increase in variability.

388 **Task 2** Due to the smaller volume of the optic cup, any

389 blood vessels present will take up a larger proportion of
 390 the volume than with the optic disc. Despite the symmetry
 391 of the texture measurements the variability is much more
 392 asymmetric, with a notable decrease in variability in the top
 393 center segment.

394 **4.2.3. GM Spinal Data**

395 The results for the GM Spinal data can be seen in Figure
 396 9. This data was the most difficult to analyze due to the rel-
 397 ative size of the segmentations and the complex ‘butterfly’
 398 shape of the second task. The width of the disputed area
 399 was rarely greater than one pixel for both tasks, causing any
 400 slight variations to have a disproportionate affect on the re-
 401 sults.

402 **Task 1** The first task had relatively low variance per
 403 region in both width and DSC, with a slight increase in
 404 width for the right hand side and spikes of around +10%
 405 at the very top and bottom. The energy of the boundary
 406 was consistent for all annotators; however there were no-
 407 table changes in dissimilarity in the top and bottom quad-
 408 rants. These changes (ranging from -20% to +10%) oc-
 409 curred where the boundary becomes the most convex, which
 410 may have caused a lack in visible features for annotators to
 411 follow.

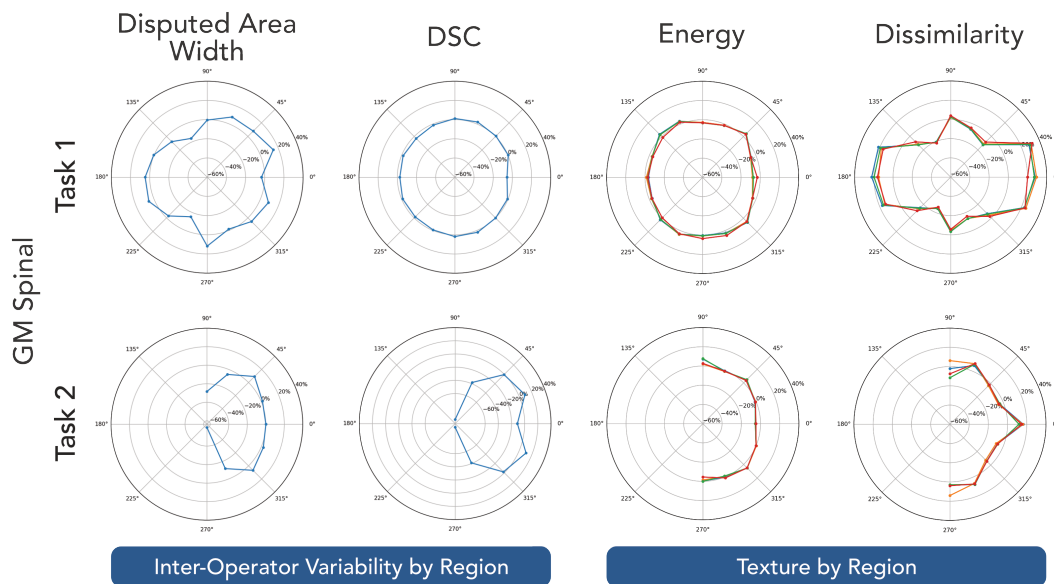


Figure 9. Inter-operator variability and texture measures by region for the GM Spinal Tasks

412 **Task 2** Due to its complex shape, the centroid of this
 413 task was instead placed to the left of the actual centroid.
 414 The width, DSC, energy, and dissimilarity were all relatively
 415 constant for the majority of the segmentation, with the
 416 exception of the long and thin dorsal horns. Counter-
 417 intuitively both the width and the DSC were significantly
 418 lower here, though the DSC much more so. The lack of
 419 variance of energy or dissimilarity in these regions suggests
 420 that the inter-operator variability here was not caused by
 421 the texture of the image; instead, this could potentially be
 422 due to the geometric complexity and thinness of the object.

423 **5. Conclusion**

424 When evaluating a model against a dataset it is vital to
 425 know the possible variability – with disputed areas being
 426 up to 50% of the total segmentation, the choice of annotator
 427 could have a significant impact on the measured accuracy
 428 of the model. For medical imaging this variability can also
 429 help provide the range of potential measurements that can
 430 be taken from a segmentation.

431 This work has compared multiple measures of inter-
 432 operator variability for segmentation regions for six tasks
 433 across three datasets. From this analysis we can conclude
 434 that certain regions may have consistently higher variabil-
 435 ity, and that for circular segmentations disputed area
 436 width is a more sensitive method than DSC for identifying
 437 these variations.

438 We have also compared these measures of inter-operator
 439 variability to texture measures of energy and dissimilarity
 440 for individual annotations, with the aim of providing a sub-

stantiated explanation for the uncertainty.

441 From this we have found that regions of high energy
 442 and low dissimilarity frequently correlate with an increase
 443 in inter-operator variability. These regions will be more
 444 monotonic than other regions in the image, providing no
 445 clear boundary for the annotator to follow. Similarly regions
 446 with low energy and high dissimilarity, which appear
 447 disordered and noisy, also correlate with regions with higher
 448 inter-operator variability.
 449

450 A limitation of this work is the lack of publicly available
 451 datasets to evaluate. While inter-operator availability is
 452 described for privately available datasets it is often done so
 453 as a single value, with no distinction for regions of the
 454 feature. By reporting variability with an approach such as
 455 the one described here much more nuance is obtained, without
 456 the need to make the dataset publicly available.

457 An area that would benefit from further study is the
 458 relationship between the geometry and the variability of a
 459 task. We theorized that this could be a contributor to the
 460 uncertainty in the second GM Spinal task, however further
 461 analysis would be needed to confirm this.

462 **6. Compliance with Ethical Standards**

463 This study was conducted retrospectively using ethically
 464 acquired publicly available human subject data. The authors
 465 have no interests to disclose.

466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522

References

[1] Ahmed Almazroa, Sami Alodhayb, Essameldin Osman, Es- lam Ramadan, Mohammed Hummadi, Mohammed Dlaim, Muhannad Alkatee, Kaamran Raahemifar, and Vasudevan Lakshminarayanan. Retinal fundus images for glaucoma analysis: the riga dataset. In *Medical Imaging 2018: Imag- ing Informatics for Healthcare, Research, and Applications*, pages 55–62. SPIE, 2018. 3

[2] Reza Azad, Ehsan Khodapanah Aghdam, Amelie Rauland, Yiwei Jia, Atlas Haddadi Avval, Afshin Bozorgpour, Sanaz Karimijafarbigloo, Joseph Paul Cohen, Ehsan Adeli, and Dorit Merhof. Medical image segmentation review: The suc- cess of u-net. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. 1

[3] The Medical Image Computing and Computer As- sisted Intervention Society. Miccai registered challenges. [https://miccai.org/index.php/special- interest - groups / challenges / miccai - registered - challenges/](https://miccai.org/index.php/special-interest-groups/challenges/miccai-registered-challenges/), 2025. Accessed: 2025-01-20. 2

[4] Lee R Dice. Measures of the amount of ecologic association between species. *Ecology*, 26(3):297–302, 1945. 4

[5] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016. 2

[6] Mudasir A Ganaie, Minghui Hu, Ashwani Kumar Malik, Muhammad Tanveer, and Ponnuthurai N Suganthan. Ensemble deep learning: A review. *Engineering Applications of Artificial Intelligence*, 115:105151, 2022. 2

[7] Sergios Gatidis, Thomas Kustner, Michael Ingrisch, Clemens Cyranand, and Jens Kleesiek. Automated lesion segmentation in whole-body fdg-pet/ct - domain generaliza- tion. <https://zenodo.org/records/7845727>, 2023. Accessed: 2025-02-20. 2

[8] Mryka Hall-Beyer. Glcm texture: A tutorial v. 3.0 march 2017. *University of Calgary*, 2017. 5

[9] Wei Ji, Shuang Yu, Junde Wu, Kai Ma, Cheng Bian, Qi Bi, Jingjing Li, Hanruo Liu, Li Cheng, and Yefeng Zheng. Learning calibrated medical image segmentation via multi- rater agreement modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12341–12351, 2021. 2

[10] Hongwei Bran Li, Fernando Navarro, Ivan Ezhov, Amirhos- sein Bayat, Dhritiman Das, Florian Kofler, Suprosanna Shit, Diana Waldmannstetter, Johannes C Paetzold, Xi- aobin Hu, et al. Qubiq: Uncertainty quantification for biomedical image segmentation challenge. *arXiv preprint arXiv:2405.18435*, 2024. 2

[11] Marco Mandolini, Agnese Brunzini, Giulia Facco, Alida Mazzoli, Archimede Forcellese, and Antonio Gigante. Com- parison of three 3d segmentation software tools for hip sur- gical planning. *Sensors*, 22(14):5242, 2022. 1

[12] Caroline Petitjean, Maria A Zuluaga, Wenjia Bai, Jean- Nicolas Dacher, Damien Grosgeorge, Jérôme Caudron, Su Ruan, Ismail Ben Ayed, M Jorge Cardoso, Hsiang-Chou Chen, et al. Right ventricle segmentation from cardiac mri: a collation study. *Medical image analysis*, 19(1):187–202, 2015. 1

[13] Dzung L Pham, Chenyang Xu, and Jerry L Prince. Current methods in medical image segmentation. *Annual review of biomedical engineering*, 2(1):315–337, 2000. 1

[14] Ferran Prados, John Ashburner, Claudia Blaiotta, Tom Brosch, Julio Carballido-Gamio, Manuel Jorge Cardoso, Benjamin N Conrad, Esha Datta, Gergely Dávid, Benjamin De Leener, et al. Spinal cord grey matter segmentation chal- lenge. *Neuroimage*, 152:312–329, 2017. 3

[15] Yinghuan Shi, Jian Zhang, Tong Ling, Jiwen Lu, Yefeng Zheng, Qian Yu, Lei Qi, and Yang Gao. Inconsistency-aware uncertainty estimation for semi-supervised medical image segmentation. *IEEE transactions on medical imaging*, 41 (3):608–620, 2021. 2

[16] Toufique A Soomro, Lihong Zheng, Ahmed J Afifi, Ahmed Ali, Shafiullah Soomro, Ming Yin, and Junbin Gao. Image segmentation for mr brain tumor detection using machine learning: a review. *IEEE Reviews in Biomedical Engineer- ing*, 16:70–90, 2022. 1

[17] Michael J Trimpl, Sergey Primakov, Philippe Lambin, Eleanor PJ Stride, Katherine A Vallis, and Mark J Gooding. Beyond automatic medical image segmentation—the spec- trum between fully manual and fully automatic delineation. *Physics in Medicine & Biology*, 67(12):12TR01, 2022. 1

[18] Simon K Warfield, Kelly H Zou, and William M Wells. Sim- ultaneous truth and performance level estimation (staple): an algorithm for the validation of image segmentation. *IEEE transactions on medical imaging*, 23(7):903–921, 2004. 2

[19] Jianli Yang, Qiaozhi Ma, Jiqiang Liu, Haiping Zu, Siqing Dong, Ying Liu, Gang Guo, Binbin Nie, and Xuetao Mu. Multiparametric magnetic resonance imaging with compre- hensive assessment of prostate volume, morphology, and composition better reflects the correlation with international prostate symptom score. *Urology*, 177:134–141, 2023. 2