

#	Parameter name	Description	Representative instantiations
<b>Challenge organization</b>			
1	Challenge name	Full name of the challenge with year	Example: MICCAI Endoscopic Vision Challenge 2015
2	Challenge acronym	Acronym of the challenge (if any).	Example: EndoVis15
3	Challenge website	URL of challenge website (if any).	<ul style="list-style-type: none"> <li>- URL to challenge website</li> <li>- Private link to website under construction</li> <li>- No website</li> </ul>
4	Organizing institutions and contact person	Information on the organizing team including contact person and other team members.	Should include: <ul style="list-style-type: none"> <li>- Contact person with affiliation</li> <li>- Team members with affiliations</li> </ul>
5	Lifecycle type	Submission cycle of the challenge. Not every challenge closes after the submission deadline (one time event). Sometimes it is possible to submit results after the deadline (open call) or the challenge is repeated with some modifications (repeated event).  <i>Example 1 - Brain tumor segmentation: One time event</i>	<ul style="list-style-type: none"> <li>- One time event</li> <li>- Repeated event</li> <li>- Open call</li> </ul>
6	Challenge schedule	Timetable for the challenge which includes the release of training and test cases, the submission dates, possibly associated workshop days, release of results and other important dates.	Should include: <ul style="list-style-type: none"> <li>- Training data release(s)</li> <li>- Test data release(s)</li> <li>- Submission deadline</li> <li>- Conference day (if any)</li> </ul>
7	Ethics approval	Information on ethics approval, preferably Institutional Review Board, location, date and number of the ethics approval.  <i>Example 1 - Brain tumor segmentation: &lt;URL to ethics approval&gt;</i>	<ul style="list-style-type: none"> <li>- No ethics needed (due to in silico validation)</li> <li>- URL to ethics approval document</li> <li>- No ethics required (data downloaded from a public database)</li> </ul>
<b>Participation conditions</b>			
8	Interaction level policy	Allowed user interaction of the algorithms assessed.  <i>Example 1 - Brain tumor segmentation: Both automatic and semi-automatic algorithms can participate in the challenge.</i>	<ul style="list-style-type: none"> <li>- Fully interactive</li> <li>- Semi-automatic</li> <li>- Fully automatic</li> </ul>
9	Organizer participation policy	Participation policy for members of the organizers' institutes.  <i>Example 1 - Brain tumor segmentation: Members of the organizers' institutes may participate but they are not eligible for awards.</i>	<ul style="list-style-type: none"> <li>- Members of the organizers' institutes may participate but they are not eligible for awards and they will not be listed in the leaderboard.</li> <li>- Members of the organizers' institutes may not participate.</li> </ul>

10	Training data policy	<p>Policy on the usage of training data. The data used to train algorithms may, for example, be restricted to the data provided by the challenge or to publicly available data including (open) pre-trained nets.</p> <p><i>Example 1 - Brain tumor segmentation: The challenge training data may be complemented by other publicly available data.</i></p>	<ul style="list-style-type: none"> <li>- No policy as no training data is required</li> <li>- No additional data allowed</li> <li>- Publicly available data may be added</li> <li>- Private data may be added</li> <li>- Docker Container</li> </ul>
11	Submission format	<p>Method that is used for result submission.</p> <p><i>Example 1 - Brain tumor segmentation: Participants send the algorithm output to the organizers via email.</i></p>	<ul style="list-style-type: none"> <li>- Docker</li> <li>- Cloud</li> <li>- Upload whole code</li> <li>- Upload executable</li> <li>- Send algorithm output to organizers</li> <li>- API</li> <li>- Evaluation Platform</li> </ul>
<b>Validation objective</b>			
12	Field(s) of application	<p>Medical or biological application that the algorithm was designed for.</p> <p><i>Example 1 - Brain tumor segmentation: Diagnosis</i></p>	<ul style="list-style-type: none"> <li>- Training</li> <li>- Intervention planning</li> <li>- Intervention follow-up</li> <li>- Diagnosis</li> <li>- Screening</li> <li>- Assistance (e.g. tracking tasks)</li> <li>- Research (e.g. cell tracking)</li> <li>- Cross-phase</li> <li>- Education</li> <li>- Prognosis</li> <li>- Prevention</li> <li>- Medical data management</li> </ul>
13	Task category(ies)	<p>Category(ies) of the algorithms assessed.</p> <p><i>Example 1 - Brain tumor segmentation: Segmentation</i></p>	<ul style="list-style-type: none"> <li>- Segmentation</li> <li>- Classification</li> <li>- Tracking</li> <li>- Retrieval</li> <li>- Detection</li> <li>- Localization</li> <li>- Registration</li> <li>- Reconstruction</li> <li>- Modeling</li> <li>- Simulation</li> <li>- Regression</li> <li>- Stitching</li> <li>- Restoration</li> <li>- Prediction</li> <li>- Denoising</li> </ul>

14	Target cohort	<p>Description of subjects/objects from whom the data would be acquired in the final application.</p> <p>Remark: A challenge could be designed around the task of medical instrument tracking in robotic kidney surgery. While the validation (see parameter validation cohort) could be performed ex vivo in a laparoscopic training environment with porcine organs, the final application (i.e. robotic kidney surgery) would be targeted on real patients with certain characteristics defined by inclusion criteria such as restrictions regarding gender or age.</p> <p><i>Example 1 - Brain tumor segmentation: Patients diagnosed with glioblastoma that got MRI scans for diagnosis including T1-weighted 3D acquisitions, T1-weighted contrast-enhanced (gadolinium contrast) 3D acquisitions and T2-weighted FLAIR 3D acquisitions.</i></p>	<ul style="list-style-type: none"> <li>- Healthy volunteers that undergo screening</li> <li>- Patients that undergo laparoscopic surgery</li> <li>- Patients that get an abdominal CT</li> <li>- Patients of a particular database</li> <li>- Patients referred for early Barrett's esophagus cancer without visible abnormalities</li> <li>- Patients attending a state-of-the-art cardiac MRI diagnostic center</li> <li>- Healthy volunteers that are recruited for a certain study</li> <li>- Patients that get chemotherapy</li> <li>- Men with clinical suspicion of having prostate cancer</li> <li>- Standardized cancer cell lines (such as HeLa)</li> <li>- Physicians that use a da Vinci Si for surgical training in an ex vivo setting</li> <li>- OR team (surgeons, nurses, ...) during liver transplantation</li> <li>- Specific journals with an oncology focus (for retrieval tasks)</li> </ul>
15	Algorithm target(s)	<p>Structure/subject/object component that the algorithm focuses on.</p> <p><i>Example 1 - Brain tumor segmentation: Glioblastoma</i></p>	<ul style="list-style-type: none"> <li>- Glioblastoma</li> <li>- Hepatocellular carcinoma (HCC)</li> <li>- Vessels</li> <li>- Liver</li> <li>- Tool tip</li> <li>- (Any) Tumor</li> <li>- Surgeon</li> <li>- Nurse</li> <li>- Specific cell type</li> <li>- Operating room</li> <li>- Specular reflections</li> <li>- Fiber pathway</li> </ul>
16	Data origin	<p>Region(s)/part(s) of subjects/objects from which the data would be acquired in the final application.</p> <p><i>Example 1 - Brain tumor segmentation: Brain</i></p>	<ul style="list-style-type: none"> <li>- Abdomen</li> <li>- Liver</li> <li>- Thorax</li> <li>- Whole body</li> <li>- Whole operating room</li> <li>- Cortical gray matter</li> <li>- Specific journal (for retrieval tasks)</li> <li>- Blood obtained from forearm</li> </ul>
17	Assessment aim(s)	<p>Property(ies) of the algorithms aimed to be optimized.</p> <p>Remark: Ideally, the metrics used in the study assess the properties of the algorithm as defined by the parameter <i>assessment aim(s)</i>. For example, an assessment aim could be targeted on the accuracy of segmentation algorithms. Possible metrics to assess the accuracy include the Dice Similarity Coefficient (DSC) and the Hausdorff Distance (HD).</p> <p><i>Example 1 - Brain tumor segmentation: Accuracy of enhancing tumor/necrosis/edema segmentation</i></p>	<ul style="list-style-type: none"> <li>- Accuracy</li> <li>- Robustness</li> <li>- Reliability</li> <li>- Precision</li> <li>- Sensitivity</li> <li>- Specificity</li> <li>- Consistency</li> <li>- Runtime</li> <li>- Applicability</li> <li>- Feasibility</li> <li>- Complexity</li> <li>- Usability</li> <li>- User satisfaction</li> <li>- Criteria linked to ergonomics</li> <li>- Integration in (clinical) workflow</li> <li>- Hardware requirements</li> </ul>

## Study conditions

18	Validation cohort	<p>Subject(s)/object(s) from whom/which the data was acquired used to validate the algorithm.</p> <p>Remark: While a challenge is typically targeted on humans, validation may exclusively involve porcine models or phantoms.</p> <p><i>Example 1 - Brain tumor segmentation: Patients with glioblastoma (retrospective analysis)</i></p>	<ul style="list-style-type: none"> <li>- Specific mouse model</li> <li>- Porcine model</li> <li>- Physical phantom</li> <li>- Patients under controlled conditions</li> <li>- Patients in clinical routine</li> <li>- Porcine liver (<i>in vitro</i>)</li> <li>- <i>In silico</i> data</li> <li>- Healthy volunteers</li> </ul>
19	Center(s)	<p>Center(s) or institute(s) in which the data was acquired.</p> <p><i>Example 1 - Brain tumor segmentation: National Center for Tumor Diseases (NCT) Heidelberg</i></p>	<ul style="list-style-type: none"> <li>- Centers involved in the xy study</li> <li>- University Clinic xy</li> <li>- Centers that are part of the xy consortium</li> </ul>
20	Imaging modality(ies)	<p>Imaging technique(s) applied for training/test data acquisition.</p> <p><i>Example 1 - Brain tumor segmentation: MRI</i></p>	<ul style="list-style-type: none"> <li>- Magnetic Resonance Imaging (MRI)</li> <li>- Computed Tomography (CT)</li> <li>- Ultrasound (US)</li> <li>- 3D US</li> <li>- Intravascular US (IVUS)</li> <li>- Positron Emission Tomography (PET)</li> <li>- Light Microscopy (LM)</li> <li>- Electron Microscopy (EM)</li> <li>- X-ray</li> <li>- Optical Coherence Tomography (OCT)</li> <li>- Endomicroscopy (w/ or w/o dye)</li> <li>- SPECT</li> <li>- Video</li> <li>- Fluoroscopy</li> <li>- Thermography</li> </ul>
21	Context information	<p>Additional information given along with the images. The information may correspond directly to the image data (e.g. tumor volume), to the patient in general (e.g. gender, medical history) or to the acquisition process (e.g. medical device data during endoscopic surgery, calibration data for an image modality).</p> <p><i>Example 1 - Brain tumor segmentation: Clinical patient data: {age, gender, ...}</i></p>	<ul style="list-style-type: none"> <li>- No additional information</li> <li>- Genetic information</li> <li>- Age</li> <li>- Gender</li> <li>- Pathology</li> <li>- Clinical diagnoses</li> <li>- Patient number</li> <li>- Medical record</li> <li>- Weight</li> <li>- BMI</li> <li>- Race</li> <li>- Cancer (sub)type</li> <li>- Cancer/disease stage</li> <li>- Body weight/height</li> <li>- Smoking status</li> <li>- Clinical treatment details</li> <li>- Lab data</li> <li>- Clinical history</li> <li>- OR device data</li> <li>- Free text, such as the radiology report, the operation report or histopathology report</li> </ul>

## Validation Datasets

22	Distribution of training and test cases	<p>Describes how training and test data were split and for what reason this division was chosen. This should include information (1) on why a specific proportion of training/test data was chosen, (2) why a certain total amount of cases was chosen and (3) why certain characteristics were chosen for the training/test set (e.g. class distribution according real-world distribution vs equal class distribution).</p> <p><i>Example 1 - Brain tumor segmentation: 80% training data and 20% test data according to common practice in machine learning.</i></p>	<ul style="list-style-type: none"> <li>- Not applicable as no training data is provided</li> <li>- Randomly distributed</li> <li>- Balanced false and negative cases</li> <li>- 80% training data, 20% test data as recommended by [ref]</li> </ul>
23	Category of training data generation method	<p>Method for determining the desired algorithm output for the training data. Possible methods include manual image annotation, in silico ground truth generation and annotation by automatic methods, and no training data generated.</p> <p><i>Example 1 - Brain tumor segmentation: Hybrid: Initiation by algorithm and refinement/correction by expert physician</i></p>	<ul style="list-style-type: none"> <li>- Ground truth from simulation (exact)</li> <li>- Reference from algorithm</li> <li>- Reference from single human rater</li> <li>- Reference from multiple human raters</li> <li>- Hybrid: Initiation by algorithm, refinement by expert physician</li> <li>- Reference derived from clinical practice (diagnosis/disease code etc.)</li> <li>- Crowdsourced annotations</li> </ul>
24	Number of training cases	<p>Number of cases that can be used for algorithm training and parameter optimization. A case encompasses all data that is processed to produce one result (e.g. one segmentation) as well as the corresponding reference result.</p> <p><i>Example 1 - Brain tumor segmentation: 400</i></p>	<ul style="list-style-type: none"> <li>- No training data provided</li> <li>- 100 images</li> <li>- 100 raw endoscopic video sequences with a total of 1,000 fully annotated frames</li> </ul>
25	Characteristic of training cases	<p>Additional information on the training cases describing their nature, such as the level of detail of the annotations (e.g. fully vs weakly annotated).</p> <p><i>Example 1 - Brain tumor segmentation: Pixel-level segmentation of the structures of interest and additional clinical information as described in context information.</i></p>	<ul style="list-style-type: none"> <li>- No training data provided</li> <li>- Full annotation (pixel level)</li> <li>- Weak annotation (image level): tumor volume, disease stage</li> <li>- Mixed annotation: 1,000 fully annotated images, 100 weakly annotated images</li> <li>- 100 endoscopic video images with 10 fully annotated training images</li> </ul>
26	Annotation policy for the training cases	<p>Instructions given to the annotators prior to training case annotation. This may include description of a training phase with the software.</p> <p><i>Example 1 - Brain tumor segmentation: The annotator was instructed to segment the edema using the T2 and FLAIR images. The enhancing tumor was subsequently to be segmented on the T1 contrast-enhanced modality. Finally, the necrotic core was to be outlined using the T1 and contrast-enhanced T1 image. The annotations were to be performed in axial slices. The undergraduate student received training on 5 cases (by the radiologist) to extract the weak labels (see parameter: context information).</i></p>	<ul style="list-style-type: none"> <li>- Challenge-specific detailed instructions - e.g. should an annotation be performed along a tumor boundary or including a safety zone? Is it allowed to guess a boundary if not clearly visible</li> <li>- URL to annotation instructions</li> <li>- What tissue would you resect?</li> <li>- Where would you take a (small) biopsy?</li> </ul>

27	Annotator(s) of training data	<p>Details on the subjects/algorithms who/which annotated the training data.</p> <p><i>Example 1 - Brain tumor segmentation: Weak annotation (parameter context information) extracted from medical reports by undergraduate medical student; full image annotation performed by radiologist.</i></p>	<ul style="list-style-type: none"> <li>- No training data provided</li> <li>- Surgeon who has done &gt;100 cases of a specific type of surgery</li> <li>- Undergraduate physician (third year)</li> <li>- Engineer who developed the software</li> <li>- Physician with no prior experience in usage of the software</li> <li>- Crowd</li> <li>- Algorithm xy</li> </ul>
28	Annotation aggregation method(s) for training cases	<p>Method(s) used to merge multiple annotations for one case.</p> <p><i>Example 1 - Brain tumor segmentation: (only one observer)</i></p>	<ul style="list-style-type: none"> <li>- No aggregation</li> <li>- Simultaneous Truth and Performance Level Estimation (STAPLE)</li> <li>- Majority vote</li> <li>- An additional annotator resolves conflicts</li> <li>- Average</li> <li>- Selective and Iterative Method for Performance Level Estimation (SIMPLE)</li> <li>- Level-set based approach maximizing the a posteriori probability (LSML)</li> <li>- Strict combination (positive if and only if all annotators agree)</li> <li>- No training data is required</li> </ul>
29	Category of reference generation method	<p>Method for determining the reference (i.e. the desired algorithm result, also referred to as gold standard) which is used for assessing the participants' algorithms. Possible methods include manual image annotation, in silico ground truth generation and annotation by automatic methods.</p> <p><i>Example 1 - Brain tumor segmentation: Manual annotation</i></p>	<ul style="list-style-type: none"> <li>- Ground truth from simulation (exact)</li> <li>- Reference from algorithm</li> <li>- Reference from single human rater</li> <li>- Reference from multiple human raters</li> <li>- Acquired through previously validated methods according to [ref]</li> <li>- Reference derived from clinical practice (diagnosis/disease code etc.)</li> <li>- Crowdsourced annotations</li> <li>- Hybrid methods (e.g. initiation by algorithm, refinement by expert physician)</li> </ul>
30	Number of test cases	<p>Number of cases used to assess the performance of an algorithm. A case encompasses all data that is processed to produce one result as well as the corresponding reference result (typically not provided to the participants).</p> <p><i>Example 1 - Brain tumor segmentation: 100</i></p>	<ul style="list-style-type: none"> <li>- 100 images</li> <li>- 100 raw endoscopic video sequences with a total of 1,000 fully annotated frames</li> </ul>
31	Characteristic of test cases	<p>Additional information on the test cases describing their nature, such as the level of detail of the annotations (e.g. fully vs weakly annotated).</p> <p><i>Example 1 - Brain tumor segmentation: Pixel-level segmentation of the structures of interest and additional clinical information as described in context information</i></p>	<ul style="list-style-type: none"> <li>- Full annotation (pixel level)</li> <li>- Weak annotation (image level): tumor volume, disease stage</li> <li>- Mixed annotation: 1,000 fully annotated images, 100 weakly annotated images</li> <li>- 100 endoscopic video images with 10 fully annotated test images</li> </ul>

32	Annotation policy for test cases	<p>Instructions given to the annotators prior to test case annotation. This may include description of a training phase with the software.</p> <p><i>Example 1 - Brain tumor segmentation: The annotator was instructed to segment the edema using the T2 and FLAIR images. The enhancing tumor was subsequently to be segmented on the T1 contrast-enhanced modality. Finally, the necrotic core was to be outlined using the T1 and contrast-enhanced T1 image. The annotations were to be performed in axial slices. The undergraduate student had received training on extracting the weak labels when annotating the training images.</i></p>	<ul style="list-style-type: none"> <li>- Challenge-specific detailed instructions - e.g. should an annotation be performed along a tumor boundary or including a safety zone? Is it allowed to guess a boundary if not clearly visible?</li> <li>- URL to annotation instructions</li> <li>- What tissue would you resect?</li> <li>- Where would you take a (small) biopsy?</li> </ul>
33	Annotator(s) of test cases	<p>Details on the subjects/algorithms who/which annotated the test data.</p> <p><i>Example 1 - Brain tumor segmentation: Radiologist with 5 years of experience</i></p>	<ul style="list-style-type: none"> <li>- Surgeon who has done &gt;100 cases of a specific type of surgery</li> <li>- Undergraduate physician (third year)</li> <li>- Engineer who developed the software</li> <li>- Physician with no prior experience in usage of the software</li> <li>- Crowd</li> <li>- Algorithm xy</li> </ul>
34	Annotation aggregation method(s) for test cases	<p>Method(s) used to merge multiple annotations for one case (if any).</p> <p><i>Example 1 - Brain tumor segmentation: No merging</i></p>	<ul style="list-style-type: none"> <li>- No aggregation (ranking provided for each annotator)</li> <li>- STAPLE</li> <li>- Majority vote</li> <li>- An additional annotator resolves conflicts</li> <li>- Average</li> <li>- SIMPLE</li> <li>- LSML</li> <li>- Strict combination (positive if and only if all annotators agree)</li> </ul>
<b>Assessment method</b>			
35	Metric(s)	<p>Function(s) to assess a property of an algorithm. These functions should reflect the validation objective (see parameter <i>assessment aim(s)</i>).</p> <p><i>Example 1 - Brain tumor segmentation: 95% HD and precision applied separately to necrosis, enhancing tumor and edema.</i></p>	<ul style="list-style-type: none"> <li>- Hausdorff distance (HD)</li> <li>- Dice Similarity Coefficient (DSC)</li> <li>- Jaccard index</li> <li>- Computation time</li> <li>- Recall</li> <li>- Precision</li> <li>- Area under curve (AUC)</li> <li>- Root mean square error (RMSE)</li> <li>- Absolute volume difference</li> <li>- True positive rate</li> <li>- Computational complexity</li> <li>- Average Symmetric Surface Distance (ASSD)</li> <li>- F1-Score</li> <li>- Specificity</li> <li>- Intraclass correlation coefficient</li> <li>- Concordance index</li> <li>- MAP</li> </ul>



36	Justification of metrics	<p>Justification why the metric(s) were chosen, preferably with reference to the clinical application.</p> <p><i>Example 1 - Brain tumor segmentation: 95% Hausdorff distance as opposed to standard HD: Try to avoid that outliers have too much weight. All other metrics are commonly used in segmentation assessment (cf. ref. xy).</i></p>	<ul style="list-style-type: none"> <li>- According to best practice recommendations [ref]</li> <li>- According to paper [ref]</li> </ul>
37	Rank computation method	<p>Method used to compute a rank for all participants based on the generated results on the test data. It may include methods for aggregating over all test cases and/or for determining a final rank from multiple single metric-based ranks. It also includes the ranking order for tied positions.</p> <p><i>Example 1 - Brain tumor segmentation: For each participant <math>p_i</math> and each test case <math>c_j</math>: Compute the metric values for the 95% Hausdorff distance and precision. For each participant <math>p_i</math> and each test case <math>c_j</math>, determine the rank corresponding to both metrics (i.e. <math>R(\text{precision}, p_i, c_j)</math>: descending order for precision, <math>R(\text{HD}, p_i, c_j)</math> ascending for 95% HD). For each participant <math>p_i</math> and each test case <math>c_j</math>, compute the average rank <math>R(p_i, c_j)</math> over both metric ranks. Finally, compute the average over all case-specific ranks to get one final rank for each participant <math>p_i</math>.</i></p>	<p>Example 1:</p> <ol style="list-style-type: none"> <li>0. Initialization: For each participant <math>p_i</math> and each test case <math>c_j</math>: compute metric values <math>M1(p_i, c_j)</math> and <math>M2(p_i, c_j)</math> for metrics <math>M1</math> and <math>M2</math>.</li> <li>1. Metric-based aggregation: For each participant <math>p_i</math> compute the median over all cases <math>c_j</math> for each metric <math>M1(p_i)</math> and <math>M2(p_i)</math>.</li> <li>2. For each participant <math>p_i</math>, compute the sum over the two metrics as <math>M1(p_i) + M2(p_i)</math>.</li> <li>3. Build rank for each participant by sorting the values <math>M1(p_i) + M2(p_i)</math> for each participant.</li> </ol> <p>Example 2:</p> <ol style="list-style-type: none"> <li>0. Initialization: For each participant <math>p_i</math> and each test case <math>c_j</math>: compute metric values <math>M1(p_i, c_j)</math>, <math>M2(p_i, c_j)</math> and <math>M3(p_i, c_j)</math> for metrics <math>M1</math>-<math>M3</math>.</li> <li>1. Case-based aggregation: For each participant <math>p_i</math> and each case <math>c_j</math>, determine the performance score <math>s_j</math> on case <math>c_j</math>: <math>s_j := 1/3 (M1(p_i, c_j) + M2(p_i, c_j) + M3(p_i, c_j))</math>.</li> <li>2. For each participant <math>p_i</math> and each case <math>c_j</math>, determine the rank <math>R(p_i, c_j)</math> for case <math>c_j</math> according to score <math>s_j</math>.</li> <li>3. Compute the average over all case-specific ranks for each participant <math>p_i</math> <math>s_i := 1/N * \sum_j (R(p_i, c_j))</math> to obtain the final rank.</li> </ol>
38	Interaction level handling	<p>Methods to handle any diversity in the level of user interaction when generating the performance ranking.</p> <p><i>Example 1 - Brain tumor segmentation: Weighting function (automatic methods are ranked higher)</i></p>	<ul style="list-style-type: none"> <li>- Indication in ranking</li> <li>- Separate ranking for fully-automatic methods</li> <li>- Only automatic methods allowed</li> </ul>
39	Missing data handling	<p>Methods used to manage submissions with missing results on test cases.</p> <p><i>Example 1 - Brain tumor segmentation: In case of missing data for participant <math>p_i</math> and case <math>c_j</math>, the case-based ranks for all metrics <math>m</math> <math>R(m, p_i, c_j)</math> are set to the maximum.</i></p>	<ul style="list-style-type: none"> <li>- Missing data not allowed (incomplete submissions not evaluated)</li> <li>- Missing data ignored</li> <li>- Missing data handled as in [ref]</li> </ul>



40	Uncertainty handling	<p>Method(s) used to make uncertainties in ranking explicit.</p> <p><i>Example 1 - Brain tumor segmentation: Test the sensitivity of the ranking with bootstrapping according to [ref].</i></p>	<ul style="list-style-type: none"> <li>- Test sensitivity of the ranking by <ul style="list-style-type: none"> <li>* Leaving out test data</li> <li>* Bootstrapping approaches</li> <li>* Changes in rank computation details</li> <li>* Changes in reference annotation</li> </ul> </li> </ul>
41	Statistical test(s)	<p>Statistical test(s) used to compare the results of challenge participants.</p> <p><i>Example 1 - Brain tumor segmentation: T-test used to test the stability of the first three ranks as described in [ref].</i></p>	<p>Quantities on which the hypothesis is taken:</p> <ul style="list-style-type: none"> <li>- Stability of the ranking</li> <li>- See whether the best results have statistically significant differences</li> </ul> <p>Tests:</p> <ul style="list-style-type: none"> <li>- Wilcoxon–Mann–Whitney test</li> <li>- t-test (paired, unpaired, one-sided, two-sided)</li> <li>- Saphiro-Wilk test</li> </ul>