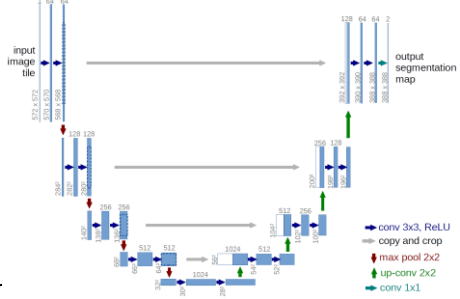# (Bench)mark: Pitfalls in AI Validation
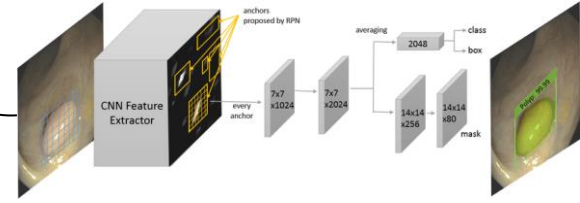
Annika Reinke

Div. Intelligent Medical Systems, German Cancer Research Center (DKFZ)

**dkfz.** GERMAN CANCER RESEARCH CENTER IN THE HELMHOLTZ ASSOCIATION

Research for a Life without Cancer

Ronneberger et al. U-net: Convolutional networks for biomedical image segmentation. MICCAI 2015

Qadir et al. Polyp Detection and Segmentation using Mask R-CNN: Does a Deeper Feature Extractor CNN Always Perform Better? ISMICT 2013
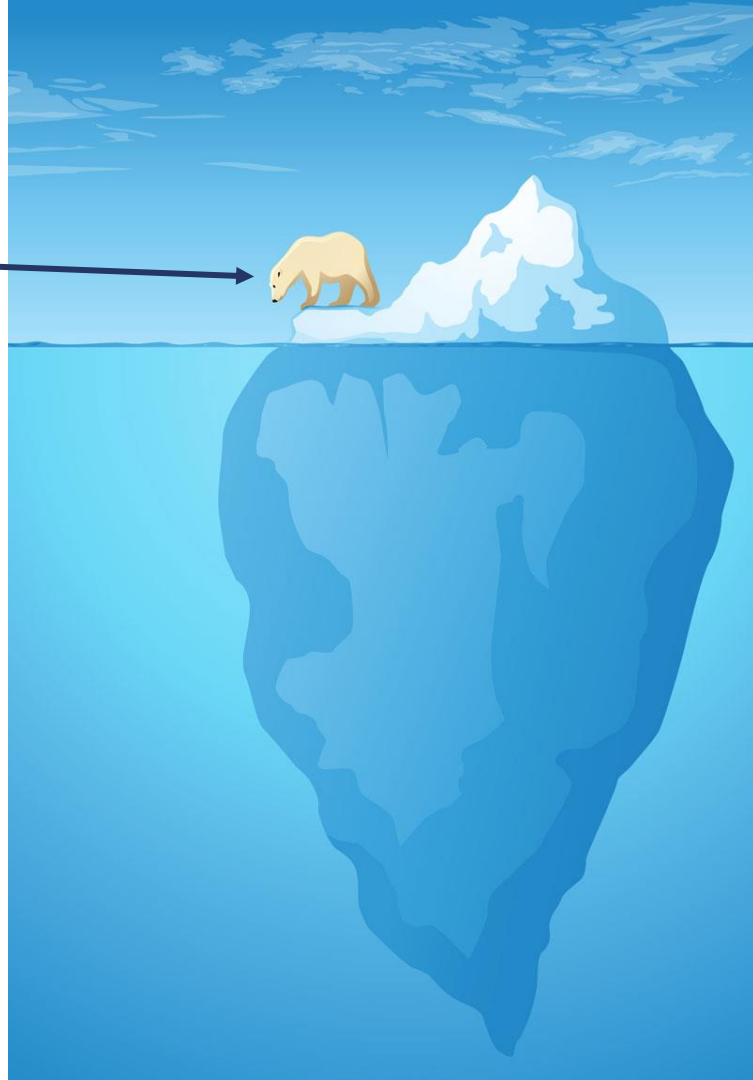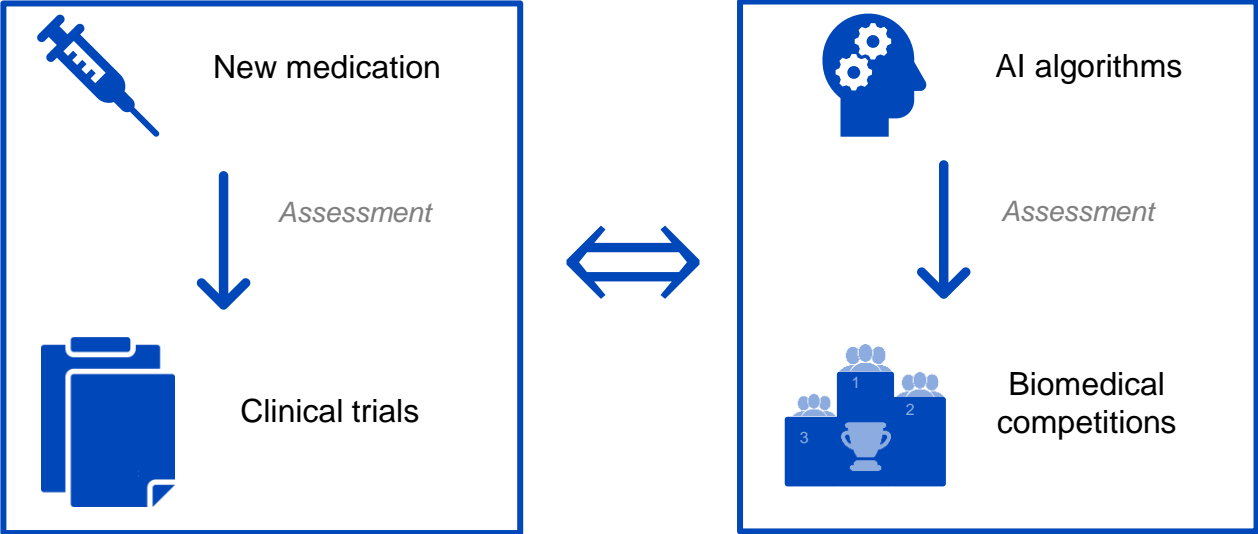
ML developer

Machine Learning (ML)

Image source:
Regulatory Affairs Professional Society

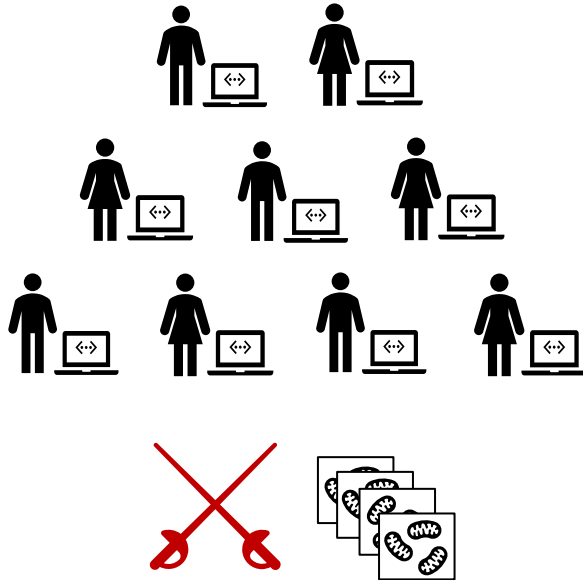ML developer

Machine Learning (ML)

Data set design,
Annotations,
Metrics,
Rankings,
Reporting,
Infrastructure,

…

# Assessment of AI algorithms

**dkfz.**

# Biomedical image analysis competitions



Up to €1 million price money

New state-of-the art method

Fame for researcher

…

- ✓ Challenges have led to common data sets used for validation

- ✓ Various fields of application covered

- ✓ Various modalities covered

# Algorithm benchmarking

**Table 12**
Comparison of existing methods.

| Methods (%) | Database | Classifier | Sensitivity (%) | Specificity (%) | Accuracy (%) | AUC (%) |
|---|---|---|---|---|---|---|
| Second technique (best) | DDSM | MLP | 96.87 | 95.94 | 96.47 | 95.10 |
| Second technique (average) | DDSM | MLP | 96.25 | 93.78 | 95.01 | 94.99 |
| Second technique (best) | MIAS | MLP | 92.70 | 90.54 | 90.16 | 95.58 |
| Second technique (average) | MIAS | MLP | 87.91 | 85.40 | 86.66 | 88.15 |
| Second technique (best) | DDSM | SCBDL | 80.70 | 79.00 | 80.00 | – |
| Wang and Yang et al., 2014; Wang, Li, & Gao, 2014 | DDSM | SVM | – | – | 92.74 | 96.50 |
| Liu and Tang, 2013 | DDSM | SVM | 92.00 | 93.00 | 93.00 | 94.39 |
| Saki et al., 2013 | MIAS | OWBPE | 90.10 | 88.06 | 89.28 | 92.80 |
| Zhang, Tomuro, Furst, & Raicu, 2012 | DDSM | SVM | – | – | 72.00 | – |
| Tahmasbi et al., 2011 | MIAS | MLP | 100 | 94.50 | 96.43 | 97.60 |
| Buciu and Gacsadi, 2011 | MIAS | PSVM | 84.61 | 80 | 82.30 | 78.00 |
| Tahmasbi et al., 2010 | MIAS | MLP | 90.10 | 94.44 | 92.80 | 98.00 |
| Verma et al., 2010 | DDSM | MLP | 85.00 | 92.50 | 88.75 | – |
| Verma et al., 2010 | DDSM | SCBDL | 97.50 | 97.50 | 97.50 | – |
| Verma et al., 2009 | DDSM | SCNN | 97.83 | 90.74 | 94.28 | – |
| Rojas-Domínguez and Nandi, 2009 | DDSM, MIAS | Bayesian, FLD | – | – | 81.00 | – |
| Mu et al., 2008 | MIAS | S2SP | – | – | – | 95.00 |
| Masotti, 2006 | DDSM | SVM | 90.00 | 95.50 | 92.75 | 97.80 |

Rouhi, et al. Benign and malignant breast tumors classification based on region growing and CNN segmentation. Expert Systems with Applications 2015.

| | Cats | CelebA | Cars | Chairs | Churches |
|---|---|---|---|---|---|
| 2D GAN [58] | 18 | 15 | **16** | 59 | 19 |
| Plat. GAN [32] | 318 | 321 | 299 | 199 | 242 |
| BlockGAN [64] | 47 | 69 | 41 | 41 | 28 |
| HoloGAN [63] | 27 | 25 | 17 | 59 | 31 |
| GRAF [77] | 26 | 25 | 39 | 34 | 38 |
| Ours | **8** | **6** | **16** | **20** | **17** |

**Table 1: Quantitative Comparison.** We report the FID score (↓) at $64^2$ pixels for baselines and our method.

| | CelebA-HQ | FFHQ | Cars | Churches | Clevr-2 |
|---|---|---|---|---|---|
| HoloGAN [63] | 61 | 192 | 34 | 58 | 241 |
| w/o 3D Conv | 33 | 70 | 49 | 66 | 273 |
| GRAF [77] | 49 | 59 | 95 | 87 | 106 |
| Ours | **21** | **32** | **26** | **30** | **31** |

**Table 2: Quantitative Comparison.** We report the FID score (↓) at $256^2$ pixels for the strongest 3D-aware baselines and our method.

Niemeyer and Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. IEEE/CVF 2021.

| Model | Top-1 Acc. | Top-5 Acc. | #Params | Ratio-to-EfficientNet | #FLOPs | Ratio-to-EfficientNet |
|---|---|---|---|---|---|---|
| **EfficientNet-B0** | **77.1%** | **93.3%** | **5.3M** | **1x** | **0.39B** | **1x** |
| ResNet-50 (He et al., 2016) | 76.0% | 93.0% | 26M | 4.9x | 4.1B | 11x |
| DenseNet-169 (Huang et al., 2017) | 76.2% | 93.2% | 14M | 2.6x | 3.5B | 8.9x |
| **EfficientNet-B1** | **79.1%** | **94.4%** | **7.8M** | **1x** | **0.70B** | **1x** |
| ResNet-152 (He et al., 2016) | 77.8% | 93.8% | 60M | 7.6x | 11B | 16x |
| DenseNet-264 (Huang et al., 2017) | 77.9% | 93.9% | 34M | 4.3x | 6.0B | 8.6x |
| Inception-v3 (Szegedy et al., 2016) | 78.8% | 94.4% | 24M | 3.0x | 5.7B | 8.1x |
| Xception (Chollet, 2017) | 79.0% | 94.5% | 23M | 3.0x | 8.4B | 12x |
| **EfficientNet-B2** | **80.1%** | **94.9%** | **9.2M** | **1x** | **1.0B** | **1x** |
| Inception-v4 (Szegedy et al., 2016) | 80.0% | 95.0% | 48M | 5.2x | 13B | 13x |
| Inception-resnet-v2 (Szegedy et al., 2017) | 80.1% | 95.1% | 56M | 6.1x | 13B | 13x |
| **EfficientNet-B3** | **81.6%** | **95.7%** | **12M** | **1x** | **1.8B** | **1x** |
| ResNeXt-101 (Xie et al., 2017) | 80.9% | 95.6% | 84M | 7.0x | 32B | 18x |
| PolyNet (Zhang et al., 2017) | 81.3% | 95.8% | 92M | 7.7x | 35B | 19x |
| **EfficientNet-B4** | **82.9%** | **96.4%** | **19M** | **1x** | **4.2B** | **1x** |
| SENet (Hu et al., 2018) | 82.7% | 96.2% | 146M | 7.7x | 42B | 10x |
| NASNet-A (Zoph et al., 2018) | 82.7% | 96.2% | 89M | 4.7x | 24B | 5.7x |
| AmoebaNet-A (Real et al., 2019) | 82.8% | 96.1% | 87M | 4.6x | 23B | 5.5x |
| PNASNet (Liu et al., 2018) | 82.9% | 96.2% | 86M | 4.5x | 23B | 6.0x |
| **EfficientNet-B5** | **83.6%** | **96.7%** | **30M** | **1x** | **9.9B** | **1x** |
| AmoebaNet-C (Cubuk et al., 2019) | 83.5% | 96.5% | 155M | 5.2x | 41B | 4.1x |
| **EfficientNet-B6** | **84.0%** | **96.8%** | **43M** | **1x** | **19B** | **1x** |
| **EfficientNet-B7** | **84.3%** | **97.0%** | **66M** | **1x** | **37B** | **1x** |
| GPipe (Huang et al., 2018) | 84.3% | 97.0% | 557M | 8.4x | - | - |

Tan and Le. EfficientNet: Rethinking model scaling for convolutional neural networks. International conference on machine learning 2019.

| | UKCF (Binary Targets) | | ADNI (Continuous Targets) | | MIMIC (Mixed Targets) | |
|---|---|---|---|---|---|---|
| | PRC(I) | PRC(C) | MSE(B) | MSE(C) | PRC | MSE |
| Base | 0.411±0.035* | 0.497±0.057* | 0.105±0.018* | 0.361±0.064 | 0.142±0.028* | 0.153±0.011 |
| REG | 0.415±0.030* | 0.518±0.052* | 0.096±0.014* | 0.360±0.066 | 0.143±0.019* | 0.152±0.010 |
| FEA | 0.410±0.033* | 0.521±0.054* | 0.092±0.012* | 0.356±0.068 | 0.144±0.030* | 0.152±0.012 |
| TEA | **0.483±0.045** | **0.583±0.072** | **0.063±0.010** | **0.330±0.066** | **0.239±0.039** | **0.150±0.012** |
| F/TEA | 0.457±0.037 | 0.576±0.071 | 0.073±0.010* | 0.338±0.067 | 0.166±0.023* | 0.154±0.011 |

Jarrett and van der Schaar. Target-embedding autoencoders for supervised representation learning. arXiv 2020.

dkfz.

Is the winner
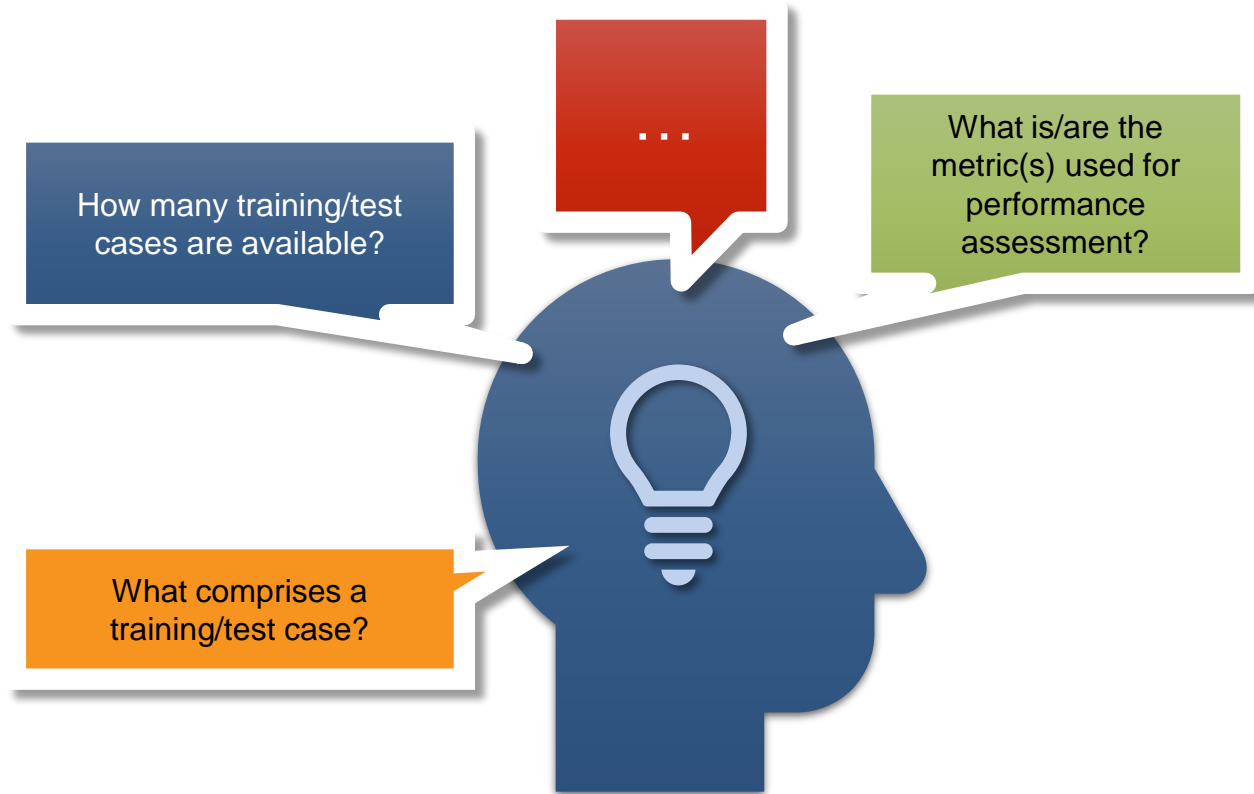really the best?

# Reporting

*"The one practice that can universally commended is the **transparent and complete reporting of all facets of a study**, allowing a critical reader to evaluate the work and fully understand its strengths and limitations*"

(Nature Neuroscience 2017, https://doi.org/10.1038/nn.4500)

**dkfz.**

# A lot of challenge parameters matter (the obvious)



How many training/test cases are available?

...

What is/are the metric(s) used for performance assessment?

What comprises a training/test case?

**dkfz.**

# A lot of challenge parameters matter (the "not so obvious")

Are challenge participants allowed to complement the training data with their own data?

...

How are missing values handled in the rankings?

Who annotated the data? How many observers? Based on which instructions?

**dkfz.**

# Analysis of > 500 competitions

- A median of **64%** of parameters were reported

- Only **6%** of parameters were reported by all challenges

- *Examples:*
    - **85%** of challenges did not give instructions on whether training data provided by challenge organizers may be complemented by other publicly available or private data
    - In **66%** of all tasks, there was no description on how the reference (i.e. gold standard) annotation was performed

Maier-Hein et al. Why rankings of biomedical image analysis competitions should be interpreted with care. **Nature Commun 2018**

**dkfz.**

# BIAS Reporting guideline

- BIAS (**B**iomedical **I**mage **A**nalysis **C**hallenge**S**) initiative: bring challenges to next level of quality

- Formed by MICCAI board challenge working group

- Developed **guideline for designing and reporting challenges**

- Registered BIAS with equator network



Maier-Hein, Reinke et al. BIAS: Transparent reporting of biomedical image analysis challenges, **Med Image Anal 2020**

# Problem: Quality control after challenge acceptance

Direct quality control

No quality control



**STEP 1**
**Challenge submission**

**STEP 2**
**Challenge review**
(with modifications of document)

**STEP 2**
**Challenge accepted**

**STEP 3**
**Challenge execution**

**dkfz.**

# Solution: Challenge registration

Direct quality control

Indirect quality control

**STEP 1**
**Challenge submission**

**STEP 2**
**Challenge review**
(with modifications of document)

**STEP 2**
**Challenge accepted**

**STEP 3**
**Upload challenge document**

**STEP 4**
**Challenge execution**

**dkfz.**

# Challenge registration

| Challenge name | Acronym | DOI |
|---|---|---|
| 2nd Retinal Fundus Glaucoma Challenge | REFUGE2 | 10.5281/zenodo.3714946 |
| 3D Head and Neck Tumor Segmentation in PET/CT | HECKTOR | 10.5281/zenodo.3714956 |
| Anatomical Brain Barriers to Cancer Spread: Segmentation from CT and MR images | ABCs | 10.5281/zenodo.3714981 |
| Automated Segmentation of Coronary Arteries | ASOCA | 10.5281/zenodo.3714985 |
| Automatic Evaluation of Mycardial Infarction from Delayed-Enhancement Cardiac MRI | EMIDEC | 10.5281/zenodo.3714997 |
| Automatic Lung Cancer Detection and Classification in Whole-slide Histopathology | ACDC@LungHP | 10.5281/zenodo.3715000 |
| Automatic Structure Segmentation for Radiotherapy Planning Challenge 2020 **(Challenge withdrawn due to COVID-19 pandemic situation)** | StructSeg 2020 | 10.5281/zenodo.3718884 |
| Cerebral Aneurysm Detection and Analysis | CADA | 10.5281/zenodo.3715011 |
| Computational Precision Medicine Challenge on Brain Tumor Classification 2020 | CPM-RadPath | 10.5281/zenodo.3718893 |
| Diabetic Foot Ulcers Grand Challenge 2020 | DFUC 2020 | 10.5281/zenodo.3715015 |
| Endoscopic Vision Challenge 2020 | EndoVis | 10.5281/zenodo.3715645 |
| International Skin Imaging Collaboration Challenge: Using Dermoscopic Image Context to Diagnose Melanoma | ISIC 2020 | 10.5281/zenodo.3715749 |
| Intracranial Aneurysm Detection and Segmentation Challenge | ADAM | 10.5281/zenodo.3715847 |
| Large Scale Vertebrae Segmentation Challenge | VerSe'20 | 10.5281/zenodo.3715865 |
| Learn2Reg - The Challenge | L2R | 10.5281/zenodo.3715651 |
| Medical Out-of-Distribution Analysis Challenge | MOOD | 10.5281/zenodo.3715869 |
| MICCAI Brain Tumor Segmentation (BraTS) 2020 Benchmark: "Prediction of Survival and Pseudoprogression" | BraTS 2020 | 10.5281/zenodo.3718903 |
| Multi-Centre, Multi-Vendor & Multi-Disease Cardiac Image Segmentation Challenge | M&Ms | 10.5281/zenodo.3715889 |
| Multi-sequence CMR based Mycardial Pathology Segmentation Challenge | MyoPS 2020 | 10.5281/zenodo.3715931 |
| Quantification of Uncertainties in Biomedical Image Quantification | QUBIQ | 10.5281/zenodo.3718911 |
| Rib Fracture Detecion and Classification | RibFrac | 10.5281/zenodo.3715933 |

Preview

Page: 1 of 33 — Automatic Zoom

# Medical Out-of-Distribution Analysis Challenge: Structured description of the challenge design

Remark: This challenge have been slightly modified. All changes are highlighted in red.

## CHALLENGE ORGANIZATION

### Title

Use the title to convey the essential information on the challenge mission.

**Medical Out-of-Distribution Analysis Challenge**

### Challenge acronym

Preferable, provide a short acronym of the challenge (if any).

Files (144.9 kB)

| Name | Size | |
|---|---|---|
| MedicalOut-of-DistributionAnalysisChallenge_v2.pdf | 144.9 kB | Preview / Download |

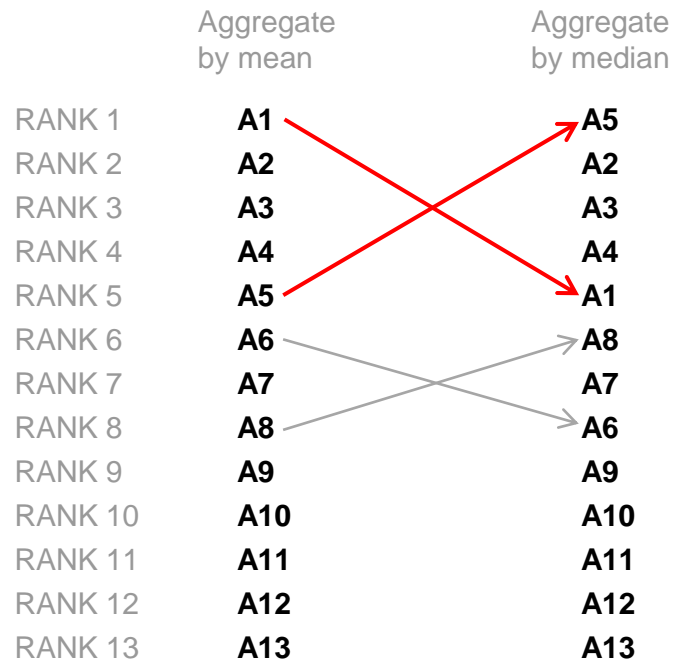md5:01c0625a7de75bfbf28497bf9dbc362d

# Rankings

# Rankings



|  | 👤💻 | 👤💻 | 👤💻 | ... |
|---|---|---|---|---|
| 🔬 | DSC: 0.87 | DSC: 0.68 | DSC: 0.94 | ... |
| 🔬 | DSC: 0.76 | DSC: 0.62 | DSC: 0.81 | ... |
| 🔬 | DSC: 0.90 | DSC: 0.71 | DSC: 0.86 | ... |
| 🔬 | DSC: 0.83 | DSC: 0.66 | DSC: 0.92 | ... |
| ... | ... | ... | ... | ... |

**dkfz.**

# Rankings

Data from MICCAI 2015 segmentation challenges

Challenge rankings are sensitive to a range of challenge design parameters:

- **Metric variant**
- Type of test case **aggregation**
- **Annotator**

| | Aggregate by mean | Aggregate by median |
|---|---|---|
| RANK 1 | A1 | A5 |
| RANK 2 | A2 | A2 |
| RANK 3 | A3 | A3 |
| RANK 4 | A4 | A4 |
| RANK 5 | A5 | A1 |
| RANK 6 | A6 | A8 |
| RANK 7 | A7 | A7 |
| RANK 8 | A8 | A6 |
| RANK 9 | A9 | A9 |
| RANK 10 | A10 | A10 |
| RANK 11 | A11 | A11 |
| RANK 12 | A12 | A12 |
| RANK 13 | A13 | A13 |

dkfz.

# Example: Exchange ranking schemes

| | Default ranking scheme<br>Metric: DSC<br>Aggregate then rank with mean | Ranking scheme 01<br>Metric: DSC<br>Aggregate then rank with median | Ranking scheme 02<br>Metric: DSC<br>Rank then aggregate with mean | Ranking scheme 03<br>Metric: DSC<br>Rank then aggregate with median | Ranking scheme 04<br>Metric: HD<br>Aggregate then rank with mean | Ranking scheme 05<br>Metric: HD<br>Aggregate then rank with median | Ranking scheme 06<br>Metric: HD<br>Rank then aggregate with mean | Ranking scheme 07<br>Metric: HD<br>Rank then aggregate with median | Ranking scheme 08<br>Metric: HD95<br>Aggregate then rank with mean | Ranking scheme09<br>Metric: HD95<br>Aggregate then rank with median | Ranking scheme 10<br>Metric: HD95<br>Rank then aggregate with mean | Ranking scheme 11<br>Metric: HD95<br>Rank then aggregate with median |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RANK 1 | A1 | | | | | | | | | | | |
| RANK 2 | A2 | | | | | | | | | | | |
| RANK 3 | A3 | | | | | | | | | | | |
| RANK 4 | A4 | | | | | | | | | | | |
| RANK 5 | A5 | | | | | | | | | | | |
| RANK 6 | A6 | | | | | | | | | | | |
| RANK 7 | A7 | | | | | | | | | | | |
| RANK 8 | A8 | | | | | | | | | | | |
| RANK 9 | A9 | | | | | | | | | | | |
| RANK 10 | A10 | | | | | | | | | | | |
| RANK 11 | A11 | | | | | | | | | | | |
| RANK 12 | A12 | | | | | | | | | | | |
| RANK 13 | A13 | | | | | | | | | | | |

Reinke et al. How to Exploit Weaknesses in Biomedical Challenge Design and Organization. **MICCAI 2018**
Maier-Hein et al. Why rankings of biomedical image analysis competitions should be interpreted with care. **Nature Commun 2018**

dkfz.

# Analysis of Results

## 27%

of all reports are based solely on
ranking lists (without further visualization)

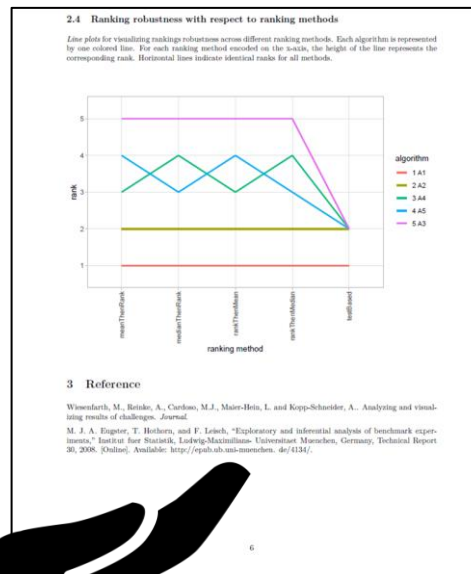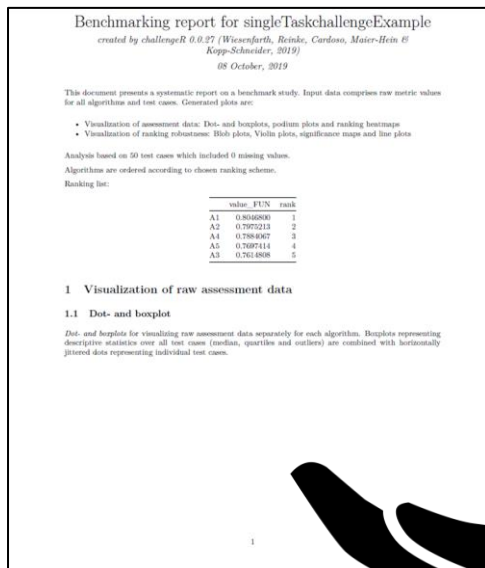# Why result analysis and visualization is critical: Example



Wiesenfarth, Reinke et al. Methods and open-source toolkit for analyzing and visualizing challenge results. **Scientific Reports 2021**

# Try it yourself: Metric values in, full PDF report out

**Input:**

Metric values in csv file

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | | task | alg_name | value | case |
| 2 | 251 | c_random | A1 | 0.705483 | 1 |
| 3 | 252 | c_random | A1 | 0.843386 | 2 |
| 4 | 253 | c_random | A1 | 0.660242 | 3 |
| 5 | 254 | c_random | A1 | 0.956698 | 4 |
| 6 | 255 | c_random | A1 | 0.861703 | 5 |
| 7 | 256 | c_random | A1 | 0.663634 | 6 |
| 8 | 257 | c_random | A1 | 0.879471 | 7 |
| 9 | 258 | c_random | A1 | 0.903639 | 8 |
| 10 | 259 | c_random | A1 | 0.888527 | 9 |
| 11 | 260 | c_random | A1 | 0.767565 | 10 |
| 12 | 261 | c_random | A1 | 0.953104 | 11 |
| 13 | 262 | c_random | A1 | 0.868738 | 12 |
| 14 | 263 | c_random | A1 | 0.706565 | 13 |
| 15 | 264 | c_random | A1 | 0.328561 | 14 |
| 16 | 265 | c_random | A1 | 0.932449 | 15 |
| | 266 | c_random | A1 | 0.810777 | 16 |
| | 7 | c_random | A1 | | 17 |
| | 8 | c_ra... | A1 | | |
| | 269 | c_random | A1 | 0.910619 | |

https://github.com/wiesenfa/challengeR

Icons created by the Noun Project

Wiesenfarth et al. Methods and open-source toolkit for analyzing and visualizing challenge results. **Scientific Reports 2021**

**dkfz.**

# Try it yourself: Metric values in, full PDF report out

**Output:**
Full PDF report

https://github.com/wiesenfa/challengeR

Wiesenfarth et al. Methods and open-source toolkit for analyzing and visualizing challenge results. **Scientific Reports 2021**

**dkfz.**

# Cheating

**You don't think people cheat?**

**20%**

of the MICCAI 2020 challenge organizers reported cheating!

# Example: Weaknesses in challenge design can be exploited



**MICCAI 2018 challenges:**
*No report* on
… full ranking scheme **(50%)**
… metrics **(15%)**

Legend:
- Differently reported as stated in proposal
- Not reported on the website
- Reported as stated in proposal

Reinke et al. How to Exploit Weaknesses in Biomedical Challenge Design and Organization. **MICCAI 2018**

**dkfz.**

# Example: Weaknesses in challenge design can be exploited

Ranking schemes are often not published before the challenge



*Tuning?*

Organizer's main competitor

| Ranking scheme 1 | Ranking scheme 2 |
|---|---|
| A1 | **A3** |
| A2 | A2 |
| A3 | **A1** |
| A4 | A4 |
| A5 | A5 |

Reinke et al. How to Exploit Weaknesses in Biomedical Challenge Design and Organization. **MICCAI 2018**

**dkfz.**

# Example: Missing value handling

**82%** of tasks provide no information about how missing data is handled



| Image | $I_1$ | $I_2$ | $I_3$ | $I_4$ | $I_5$ | $I_6$ |
|-------|-------|-------|-------|-------|-------|-------|
| **DSC** | 0.94 | **NaN** | 0.87 | 0.90 | **NaN** | 0.89 |

*Ignore NaNs*     *Set NaNs to worst possible value (here: 0)*

**Mean DSC:** 0.90     **Mean DSC:** 0.60

Reinke et al. How to Exploit Weaknesses in Biomedical Challenge Design and Organization. **MICCAI 2018**
Reinke et al. Common Limitations of Image Processing Metrics: A Picture Story. **arXiv 2021**

**dkfz.**

# Example: Missing value handling

What happens if algorithms systematically submit only the most plausible results?

- **25%** of non-winning algorithms would have been ranked first

- In **9%** of tasks, every single participating algorithm could have been ranked first

Reinke et al. How to Exploit Weaknesses in Biomedical Challenge Design and Organization. **MICCAI 2018**

# Metrics

Hamming loss
Dice similarity coefficient
Specificity
Volumetric overlap error
Jaccard
Average precision
Rand error
Area under curve
Average mesh distance
False positive rate
F1 score
Jaccard
Volume
ROC
Adjusted rand index
Recall
Absolute surface distance
False positive rate
Spearmans correlation coefficients
Volumetric overlap error
Root mean square error
Hausdorff distance
Hausdorff distance
F1 score
Average perpendicular distance
Mean average precision
Accuracy
Precision
Kappa
Average mesh distance
Averaged confusion matrix diagonal
Error rate
Average perpendicular distance
Rand error
Error rate
Root mean square error
Completeness
Volume
Dice similarity coefficient
Kappa
Area under curve
Homogeneity
Dice similarity coefficient
Average recall
Accuracy
Absolute surface distance
Specificity
Completeness
Area under curve
Hausdorff distance 95
Hausdorff distance 95
Error rate
Precision
Precision
F1 score
False negative rate
Mean average precision
False positive rate
Hausdorff distance
ROC
Average precision
Hamming loss
Adjusted rand index
Maximum surface distance
Average symmetric surface distance
Accuracy
Interclass correlation
Interclass correlation
Kappa
Average recall
Average mesh distance
Average symmetric surface distance
Adjusted rand index
Average recall
Jaccard
Average symmetric surface distance
Spearmans correlation coefficient
Maximum surface distance
Absolute surface distance
Volume
Euclidean distance
Root mean square error
Homogeneity
Specificity
False negative rate
Sweet spot coverage
Euclidean distance
Recall
Rand error
Hausdorff distance 95
Recall
Sweet spot coverage
Euclidean distance
Average precision
Maximum surface distance
Average perpendicular distance

# How metrics are currently selected

The two metrics are **widely used** to measure segmentation accuracy.

One of the **commonly used** metrics to validate scene segmentation

It is **widely used** in the field.

Is **commonly used** for the validation of reconstruction data.

**Standard** computer vision segmentation metrics.

Is a **common metric** in clinical diagnosis validation.

**Count**

| | 0 | 100 | 200 | 300 | 400 | 500 |

Dice similarity coefficient
Hausdorff distance
Average surface distance
Adjusted rand index
Interclass correlation
Recall
Average symmetric surface distance
Precision
Specificity
Accuracy
Euclidean distance
Volume
Jaccard index
Area under curve
F1 score
Hausdorff distance 95
Kappa
Number

**dkfz.**

# Class imbalance

**Goal:** Classify patients into sick (positive class) and healthy (negative class)



Sick class: 97 patients

Healthy class: 3 patients



Reality

Accuracy = 97%

**dkfz.**

# Class imbalance

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} = \frac{97 + 0}{97 + 0 + 3 + 0} = 0.97$$



Accuracy = 97%

Specificity = 0%

**dkfz.**

# Most common metric: Dice Similarity Coefficient (DSC)



$$DSC(A,B) = \frac{\blacksquare + \blacksquare}{\blacksquare + \blacksquare}$$

$$= \frac{2\,|A \cap B|}{|A| + |B|}$$

$$IoU(A,B) = \frac{\blacksquare}{\blacksquare + \blacksquare - \blacksquare}$$

$$= \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

$$= \frac{|A \cap B|}{|A \cup B|}$$

A    B    A ∩ B

Reinke et al. A discovery dive into the world of evaluation — Do's, don'ts and other considerations. **Medium Blogpost 2021**

**dkfz.**

# Shape unawareness

# Inappropriate phrasing of the problem:
Object detection vs. segmentation



Reference        Prediction 1        Prediction 2

1 object detected ❌      3 objects detected ✅

Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022
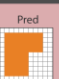Reinke et al. Common Limitations of Image Processing Metrics: A Picture Story. arXiv 2021

**dkfz.**

# Metric aggregation

Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022
Reinke et al. Common Limitations of Image Processing Metrics: A Picture Story. arXiv 2021

# Uncovering problems is good…

## Common Limitations of Image Processing Metrics: A Picture Story

ANNIKA REINKE*, German Cancer Research Center (DKFZ), Germany and Heidelberg University, Germany
MINU D. TIZABI, German Cancer Research Center (DKFZ), Germany
CAROLE H. SUDRE, University College London, UK and King's College London, UK
MATTHIAS EISENMANN, German Cancer Research Center (DKFZ), Germany
TIM RÄDSCH, German Cancer Research Center (DKFZ), Germany and understandAI GmbH, Germany
MICHAEL BAUMGARTNER, German Cancer Research Center (DKFZ), Germany
LAURA ACION, CONICET – Universidad de Buenos Aires, Argentina and University of Iowa, USA
MICHELA ANTONELLI, King's College London, UK and University College London, UK
TAL ARBEL, McGill University, Canada

Solving them is even better!

# Problem-aware metric recommendation framework



Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022

# Problem fingerprint



| Image processing category identified by category mapping | | Semantic segmentation (SS): assignment of one or multiple category labels to each pixel. |
|---|---|---|

**Domain interest-related properties (selection)**

| Particular importance of structure boundaries | | The application requires exact structure boundaries. |
|---|---|---|
| Particular importance of structure center (e.g. in cells, vessels) | | The application requires accurate knowledge of structure centers. |
| Compensation for annotation imprecisions requested | | The reference annotation is typically only an approximation of the (forever unknown) ground truth. It may be desirable to compensate for known uncertainties, such as intra-rater or inter-rater variability, by configuring the metric accordingly. This is only possible for some metrics. |
| ... | ... ... | ... |

**Target structure-related properties (selection)**

| Small size of structures relative to pixel size | | Structures of the provided class are consistently small relative to the grid size in such a way that a single pixel makes up at least several percentage points of the structure volume. |
|---|---|---|
| High variability of structure sizes (within one image, across images) | | The target structures vary substantially in size, such that some structures are several times the sizes of others. |
| ... | ... ... | ... |

**Data set-related properties (selection)**

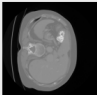| Presence of class imbalance | | The class prevalences differ substantially. |
|---|---|---|
| Non-independence of test cases | | The test cases are hierarchically structured, indicating non-independence of test cases. |
| ... | ... ... | ... |

**Algorithm output-related properties (selection)**

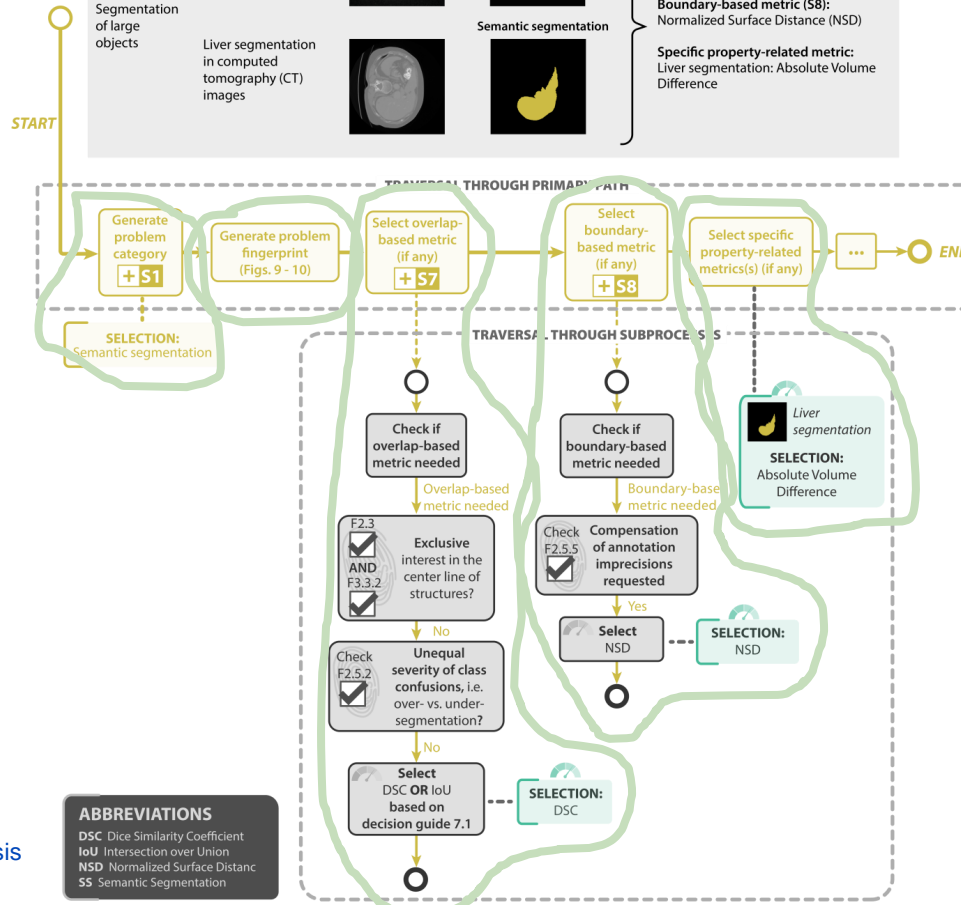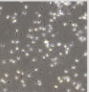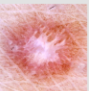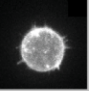| Possibility of algorithm output not containing the target structure(s) | | The algorithm may yield output images only comprising the background class. |
|---|---|---|
| ... | ... ... | ... |

Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. **arXiv 2022**

Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022

Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022
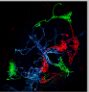
| DESCRIPTION OF PROBLEM | SCENARIO | SAMPLE INPUT IMAGE | RECOMMENDED OUPUT IMAGE | RECOMMENDATION |
|---|---|---|---|---|
| Segmentation of large objects | Lung cancer cell segmentation from microscopy images | | | **Overlap-based metric (S7):** Dice Similarity Coefficient (DSC) |
| | | | Semantic segmentation | **Boundary-based metric (S8):** Normalized Surface Distance (NSD) |
| | Liver segmentation in computed tomography (CT) images | | | **Specific property-related metric:** Liver segmentation: Absolute Volume Difference |

**TRAVERSAL THROUGH PRIMARY PATH**

START

Generate problem category + S1 → Generate problem fingerprint (Figs. 9 - 10) → Select overlap-based metric (if any) + S7 → Select boundary-based metric (if any) + S8 → Select specific property-related metrics(s) (if any) → ... → END

**SELECTION:** Semantic segmentation

Maier-Hein/Reinke et al.
Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022

| DESCRIPTION OF PROBLEM | SCENARIO | SAMPLE INPUT IMAGE | RECOMMENDED OUPUT IMAGE | RECOMMENDATION |
|---|---|---|---|---|
| Segmentation of large objects | Lung cancer cell segmentation from microscopy images | | | **Overlap-based metric (S7):** Dice Similarity Coefficient (DSC) |
| | | | Semantic segmentation | **Boundary-based metric (S8):** Normalized Surface Distance (NSD) |
| | Liver segmentation in computed tomography (CT) images | | | **Specific property-related metric:** Liver segmentation: Absolute Volume Difference |

**START**

**TRAVERSAL THROUGH PRIMARY PATH**

Generate problem category  **+ S1**

Generate problem fingerprint (Figs. 9 - 10)

Select overlap-based metric (if any)  **+ S7**

Select boundary-based metric (if any)  **+ S8**

Select specific property-related metrics(s) (if any)

··· **END**

**SELECTION:** Semantic segmentation

**TRAVERSAL THROUGH SUBPROCESSES**

Check if overlap-based metric needed

Overlap-based metric needed

F2.3 **AND** F3.3.2 — Exclusive interest in the center line of structures?

No

Check F2.5.2 — Unequal severity of class confusions, i.e. over- vs. under-segmentation?

No

Select DSC **OR** IoU based on decision guide 7.1

**SELECTION:** DSC

Check if boundary-based metric needed

Boundary-based metric needed

Check F2.5.5 **Compensation of annotation imprecisions requested**

Yes

Select NSD

**SELECTION:** NSD

Liver segmentation

**SELECTION:** Absolute Volume Difference

**ABBREVIATIONS**
**DSC** Dice Similarity Coefficient
**IoU** Intersection over Union
**NSD** Normalized Surface Distanc
**SS** Semantic Segmentation

Maier-Hein/Reinke et al.
Metrics reloaded: Pitfalls and recommendations for image analysis validation. arXiv 2022

# Instantiation for common biomedical use cases



Maier-Hein/Reinke et al. Metrics reloaded: Pitfalls and recommendations for image analysis validation. **arXiv 2022**

**Metrics Reloaded –**
A new recommendation framework for biomedical image analysis validation

https://arxiv.org/abs/2206.01653

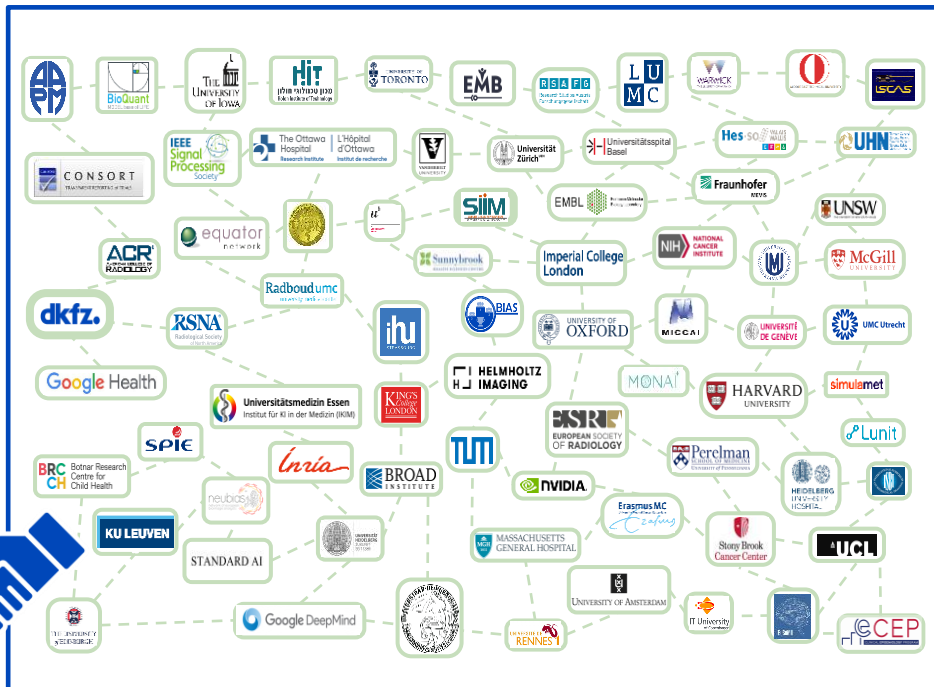Working on web-based tool for guiding the user through the process!

**Intelligent Medical Systems**
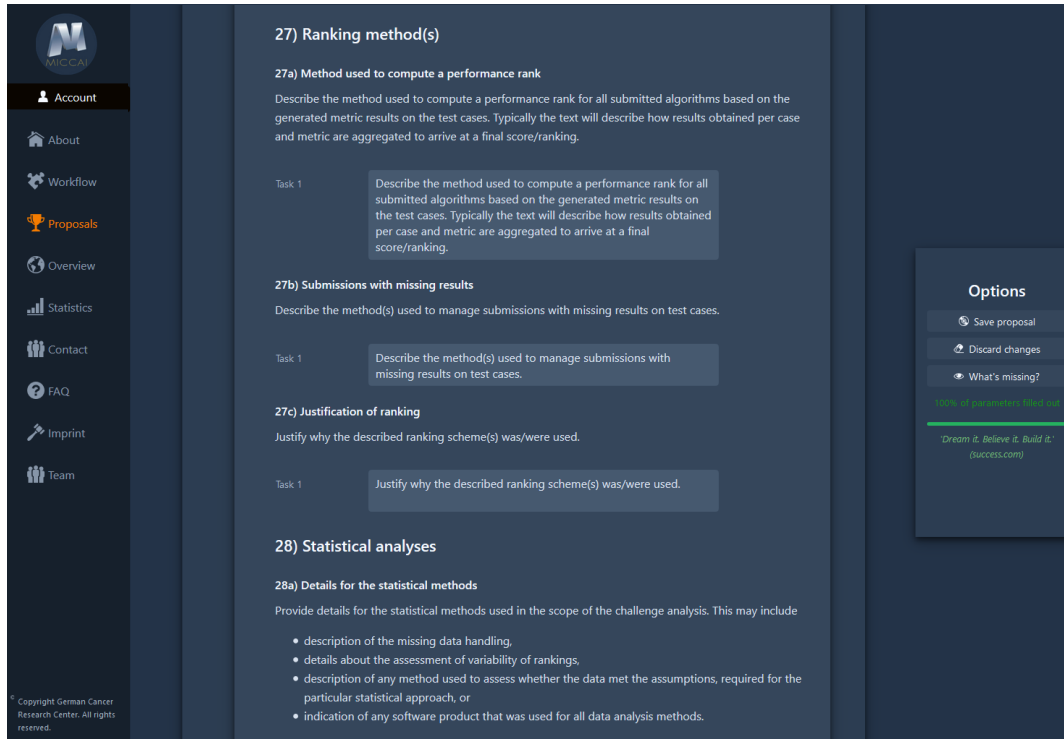**DKFZ**

@DKFZ_CAMI_lab
#BiomedicalChallenges
#benchmarking

HELMHOLTZ IMAGING

erc European Research Council
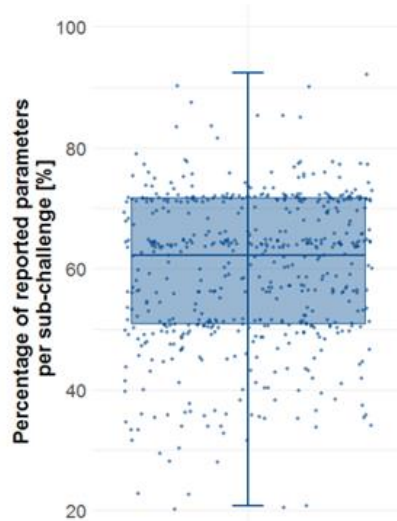
dkfz.

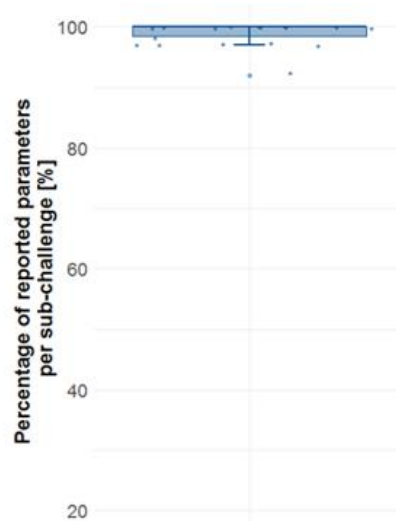# New: Structured challenge submission system



**DEVELOPERS:**
Annika Reinke, Sinan Onogur, Matthias Eisenmann, Keno März, Sebastian Pirmann
Div. Computer Assisted Medical Interventions (CAMI), German Cancer Research Center, DKFZ)

MICCAI 2018 — 2018

MICCAI 2019 — 2019

MICCAI 2020
ISBI 2020
MIDL 2020 — 2020

MICCAI 2021
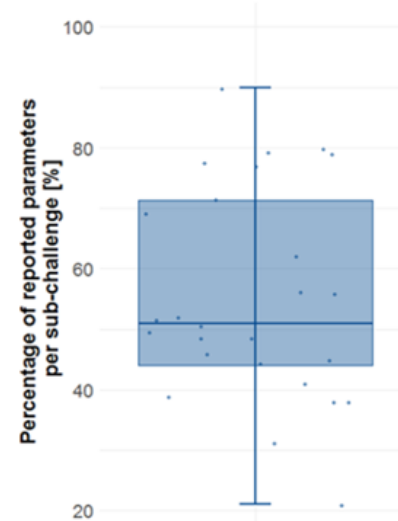ISBI 2021 — 2021

MICCAI 2022 — 2022

# Problem: Quality control after challenge acceptance



2007 - 2016

2018 Proposals

2018 Websites
Captured: July 2018

**dkfz.**