

Deep Image Generation Model Challenge in Surgery 2022: Structured description of the challenge design

CHALLENGE ORGANIZATION

Title

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

Deep Image Generation Model Challenge in Surgery 2022

Challenge acronym

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

AdaptOR 2022

Challenge abstract

Provide a summary of the challenge purpose. This should include a general introduction in the topic from both a biomedical as well as from a technical point of view and clearly state the envisioned technical and/or biomedical impact of the challenge.

The continuing AdaptOR Challenge aims to spark methodological developments in deep image generation models for the surgical domain. In this year's edition, we again focus on video-assisted mitral valve repair [1], which is becoming the novel state-of-the-art [2]. Especially the usage of 3D endoscopes, where left and right images captured from a stereo-camera are presented on a 3D compatible monitor have proven to be beneficial, since they enable better perception of the depth, and spatial relations of structures in the scene. Additionally, exploiting the stereo information enables quantitative analyses of downstream tasks for example in 3D pose estimation.

Towards this end and building upon the challenge in the previous year [7], the AdaptOR challenge 2022 proposes a task of novel view synthesis for endoscopic data. During training, the participants are provided the left and the right stereo camera images, and the test time task is to predict the corresponding image for a given image from the left camera. Clinically, this improves the perception of crucial structures in the surgical scene such as depth of chordae, relative position of papillary muscles, and size of mitral annulus. Furthermore, novel view synthesis is an integral sub-task of numerous existing cutting-edge depth estimation methods [5,6]. This enables not only real-time workflows, but also retrospective analyses on the existing 2D endoscopic data that would otherwise not be possible.

Intra-operative datasets in this challenge have varying camera angles, illumination, field of views and occlusions from tissues, tubes, and increase light reflections from surgical headlights. Especially demanding in these scenes, is the view-dependent appearance of the objects that are directly in front of the camera (eg. sutures). They render it difficult to train models that faithfully predict the missing image from the image pair or to define the correct correspondences between associated pixels. Therefore, the proposed task of novel view synthesis is difficult to solve.

To enhance the training split with data from a related domain, we additionally provide stereo frames captured

from a mitral valve surgical simulator. This data is captured from mitral valve repair performed on patient specific 3D printed silicone valves. They contain a comparably stable illumination (less varying reflections) and stereo relation and a more standardized view angle. Participants are invited to include approaches that can learn robust image features that can be potentially transferred to the intra-operative domain (e.g., [8]).

The dataset this year is an extension of our dataset we released in the previous year, where we now consider more surgeries and additional phases of mitral valve repair to significantly increase the sizes of the single splits at higher resolutions.

Challenge keywords

List the primary keywords that characterize the challenge.

Novel view synthesis, Image reconstruction, Image generation, Image synthesis, Generative models, Generative Adversarial Networks, Image transformation, Mitral Valve, Heart Valve, 3D Endoscopy, Surgical Training, Surgical Simulation

Year

The challenge will take place in ...

2022

FURTHER INFORMATION FOR MICCAI ORGANIZERS

Workshop

If the challenge is part of a workshop, please indicate the workshop.

Deep Generative Models for Medical Image Computing and Computer Assisted Intervention 2022 (DGM4MICCAI)

Duration

How long does the challenge take?

Half day.

Expected number of participants

Please explain the basis of your estimate (e.g. numbers from previous challenges) and/or provide a list of potential participants and indicate if they have already confirmed their willingness to contribute.

We expect about 5-10 teams to participate.

Our estimation is based on

- The number of registrations we received on our Synapse platform in the previous year for AdaptOR 2021 (6 teams). Due to multiple requests, we re-opened the challenge again in Dec 2021.
- Past EndoVis 19, 20, 21 which are challenges related to endoscopy (about 5-10 teams)

Publication and future plans

Please indicate if you plan to coordinate a publication of the challenge results.

Participants will firstly submit a half-page write up of the challenge explaining their approach (together with the docker container submission). Additionally, the participants have the option to submit an 8-page paper on their methods, adhering to the same schedule of the workshop, which will be published in LNCS. After the challenge is

concluded, a journal paper (preferable TMI or MedIA) will be submitted to summarize the challenge results. First and last author of each team will be invited as co-authors.

With this second edition of AdaptOR, we aim to establish a reoccurring event to support progress in this application field and in deep generative image model research.

Space and hardware requirements

Organizers of on-site challenges must provide a fair computing environment for all participants. For instance, algorithms should run on the same computing platform provided to all.

The challenge will not be on-site. The online challenge will use the synapse platform.

TASK: Novel View Synthesis in Stereoendoscopy

SUMMARY

Abstract

Provide a summary of the challenge purpose. This should include a general introduction in the topic from both a biomedical as well as from a technical point of view and clearly state the envisioned technical and/or biomedical impact of the challenge.

described above

Keywords

List the primary keywords that characterize the task.

Novel View Synthesis, Stereo, Endoscopy, Mitral Valve, Heart, Surgery, Domain Adaptation, Gene Deep Learning, Machine Learning

ORGANIZATION

Organizers

a) Provide information on the organizing team (names and affiliations).

Jun. Prof. Dr. Sandy Engelhardt, University Hospital Heidelberg

Dr. Anirban Mukhopadhyay, Technical University Darmstadt

Prof. Dr. Raffaele De Simone, University Hospital Heidelberg

Lalith Sharan, University Hospital Heidelberg

Halvar Kelm, University Hospital Heidelberg

Henry Krumb, Technical University Darmstadt

b) Provide information on the primary contact person.

Sandy Engelhardt, University Hospital Heidelberg

Email: sandy.engelhardt@med.uni-heidelberg.de

Life cycle type

Define the intended submission cycle of the challenge. Include information on whether/how the challenge will be continued after the challenge has taken place. 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).

Examples:

- One-time event with fixed conference submission deadline
- Open call (challenge opens for new submissions after conference deadline)
- Repeated event with annual fixed conference submission deadline

One time event with fixed submission deadline.

Challenge venue and platform

a) Report the event (e.g. conference) that is associated with the challenge (if any).

MICCAI.

b) Report the platform (e.g. grand-challenge.org) used to run the challenge.

Docker submission: <https://docs.synapse.org/>

The challenge will be linked on grand-challenge.org (once the proposal is accepted).

c) Provide the URL for the challenge website (if any).

<https://adaptor2022.github.io/> (site under construction)

Participation policies

a) Define the allowed user interaction of the algorithms assessed (e.g. only (semi-) automatic methods allowed).

Fully automatic.

Additional points: Only fully automatic approaches are allowed

b) Define the 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.

No additional data and no models pre-trained on other datasets are allowed.

c) Define the participation policy for members of the organizers' institutes. For example, members of the organizers' institutes may participate in the challenge but are not eligible for awards.

Members of the organizers' institutes may participate in the challenge but are not eligible for awards.

d) Define the award policy. In particular, provide details with respect to challenge prizes.

Certificates will be provided for the top 3 performing teams. Upon acceptance of the challenge, we will seek for sponsorship of the winner(s) of the challenge from industry partners.

e) Define the policy for result announcement.

Examples:

- Top 3 performing methods will be announced publicly.
- Participating teams can choose whether the performance results will be made public.

All the results will be made available publicly. All teams will be invited to the half-day challenge event at DGM4MICCAI workshop to present their work in more detail.

f) Define the publication policy. In particular, provide details on ...

- ... who of the participating teams/the participating teams' members qualifies as author
- ... whether the participating teams may publish their own results separately, and (if so)
- ... whether an embargo time is defined (so that challenge organizers can publish a challenge paper first).

Challenge submission can be optionally accompanied by an 8-page LNCS format paper, describing individual methods by the participants in detail. The paper will be published in the „Deep Generative Model“ workshop proceedings after the workshop. There are no restrictions on the number of authors.

After the challenge, the challenge organizers will publish one challenge journal paper together with two participants of each challenge team summarizing the results. Each team should nominate two authors (typically it is the first and last author). An embargo period until the availability of this journal paper will be defined.

Submission method

a) Describe the method used for result submission. Preferably, provide a link to the submission instructions.

Examples:

- Docker container on the Synapse platform. Link to submission instructions: <URL>
- Algorithm output was sent to organizers via e-mail. Submission instructions were sent by e-mail.

The challenge cohort splits into two sets:

- 1) data acquired during simulating mitral valve repair on a surgical simulator ("Sim-Domain"),
- 2) intraoperative endoscopic data from mitral valve repair ("Intraop-Domain").

For the purpose of result verification and to encourage reproducibility and transparency, the submissions are evaluated with a docker container in the Synapse platform. More specifically:

When run a predefined command on unseen input data from the left camera of the intraoperative domain, the model should output the image for the corresponding right camera of the stereo-view.

We encourage participants to provide their code open source. The URL should be added in the half-page description and in the potential LNCS submission.

Participants agree that the challenge organizers are allowed to use their submitted docker containers to run further meta-analysis.

b) Provide information on the possibility for participating teams to evaluate their algorithms before submitting final results. For example, many challenges allow submission of multiple results, and only the last run is officially counted to compute challenge results.

The challenge will be split into three phases: Training phase, Platform testing phase, Testing phase.

During training phase, the participating teams will be able to independently validate their results using cross-validation on the training data.

During the platform testing phase, they are allowed to use the official submission platform to resolve potential technical issues. We will use dummy datasets for sanity checks, e.g. to ensure the submission is in the correct format.

During the test phase, participants are allowed to make in total three submissions. The best result out of these three is selected as final result.

Furthermore, violation to the following rules will lead to disqualification:

- Each team is only allowed to register once and all submissions must be done from the same account.
- A single participant is only allowed to be part of one team.

Challenge schedule

Provide a timetable for the challenge. Preferably, this should include

- the release date(s) of the training cases (if any)
 - the registration date/period
 - the release date(s) of the test cases and validation cases (if any)
 - the submission date(s)
 - associated workshop days (if any)
 - the release date(s) of the results
- release date of the training cases: 15th April 2022
 - registration period: 15th April 2022 – 30th June 2022
 - platform testing: 01st July 2022 – 01st August 2022
 - submission of docker container: 15th July 2022 – 15th August 2022
 - LNCS paper submission date: Depending on the timeline of the workshop
 - associated workshop days: either 17th September 2022 or 23th September 2022
 - release date of the results: date of the workshop

Ethics approval

Indicate whether ethics approval is necessary for the data. If yes, provide details on the ethics approval, preferably institutional review board, location, date and number of the ethics approval (if applicable). Add the URL or a reference to the document of the ethics approval (if available).

We have received approval from the Local Ethics Committee from University Hospital Heidelberg to use the anonymized data. The registration numbers are S-658/2016 and S-777/2019.

Data usage agreement

Clarify how the data can be used and distributed by the teams that participate in the challenge and by others during and after the challenge. This should include the explicit listing of the license applied.

Examples:

- CC BY (Attribution)
- CC BY-SA (Attribution-ShareAlike)
- CC BY-ND (Attribution-NoDerivs)
- CC BY-NC (Attribution-NonCommercial)
- CC BY-NC-SA (Attribution-NonCommercial-ShareAlike)
- CC BY-NC-ND (Attribution-NonCommercial-NoDerivs)

CC BY NC SA.

Additional comments: By registering in the challenge, each team agrees (1) to use the data provided only in the scope of the challenge and (2) to neither pass it on to a third party nor to use it for any additional publication or for commercial use. After the challenge, the data will be made publicly available for non-commercial use.

Code availability

a) Provide information on the accessibility of the organizers' evaluation software (e.g. code to produce rankings). Preferably, provide a link to the code and add information on the supported platforms.

We will make the code available on the synapse platform that will be used to compute the metrics for ranking.

b) In an analogous manner, provide information on the accessibility of the participating teams' code.

Teams are encouraged to provide their code open source and to add the URL in the LNCS paper.

Conflicts of interest

Provide information related to conflicts of interest. In particular provide information related to sponsoring/funding of the challenge. Also, state explicitly who had/will have access to the test case labels and when.

No conflict of interest. Only challenge organizing team will have access to test case labels during the challenge.

MISSION OF THE CHALLENGE

Field(s) of application

State the main field(s) of application that the participating algorithms target.

Examples:

- Diagnosis
- Education
- Intervention assistance
- Intervention follow-up
- Intervention planning
- Prognosis
- Research
- Screening
- Training
- Cross-phase

Surgery.

Additional points: Intraoperative Support

Task category(ies)

State the task category(ies).

Examples:

- Classification
- Detection
- Localization
- Modeling
- Prediction
- Reconstruction
- Registration
- Retrieval
- Segmentation
- Tracking

Image Synthesis

Cohorts

We distinguish between the target cohort and the challenge cohort. For example, a challenge could be designed around the task of medical instrument tracking in robotic kidney surgery. While the challenge could be based on ex vivo data obtained from a laparoscopic training environment with porcine organs (challenge cohort), the final biomedical application (i.e. robotic kidney surgery) would be targeted on real patients with certain characteristics defined by inclusion criteria such as restrictions regarding sex or age (target cohort).

a) Describe the target cohort, i.e. the subjects/objects from whom/which the data would be acquired in the final biomedical application.

Patients with mitral insufficiency, which undergo minimally-invasive mitral valve repair.

b) Describe the challenge cohort, i.e. the subject(s)/object(s) from whom/which the challenge data was acquired.

The challenge cohort splits into two sets:

- 1) data acquired from mitral valve repair performed on a surgical simulator ("sim domain")
- 2) intraoperative endoscopic data from mitral valve repair ("intraop domain")

Imaging modality(ies)

Specify the imaging technique(s) applied in the challenge.

Endoscopy

Context information

Provide additional information given along with the images. The information may correspond ...

a) ... directly to the image data (e.g. tumor volume).

- Whether frame is captured from the left or the right camera
- To which (anonymized) patient and domain the frame belongs to

b) ... to the patient in general (e.g. sex, medical history).

No further information.

Target entity(ies)

a) Describe the data origin, i.e. the region(s)/part(s) of subject(s)/object(s) from whom/which the image data would be acquired in the final biomedical application (e.g. brain shown in computed tomography (CT) data, abdomen shown in laparoscopic video data, operating room shown in video data, thorax shown in fluoroscopy video). If necessary, differentiate between target and challenge cohort.

Mitral valve shown in endoscopy

b) Describe the algorithm target, i.e. the structure(s)/subject(s)/object(s)/component(s) that the participating algorithms have been designed to focus on (e.g. tumor in the brain, tip of a medical instrument, nurse in an operating theater, catheter in a fluoroscopy scan). If necessary, differentiate between target and challenge cohort.

Corresponding stereo view for a given image

Assessment aim(s)

Identify the property(ies) of the algorithms to be optimized to perform well in the challenge. If multiple properties are assessed, prioritize them (if appropriate). The properties should then be reflected in the metrics applied (see below, parameter metric(s)), and the priorities should be reflected in the ranking when combining multiple metrics that assess different properties.

- Example 1: Find highly accurate liver segmentation algorithm for CT images.
- Example 2: Find lung tumor detection algorithm with high sensitivity and specificity for mammography images.

Corresponding metrics are listed below (parameter metric(s)).

Sensitivity, Consistency, Precision.

Additional points: Generate images of the corresponding stereo-view with high similarity to the target stereo-image.

DATA SETS

Data source(s)

a) Specify the device(s) used to acquire the challenge data. This includes details on the device(s) used to acquire the imaging data (e.g. manufacturer) as well as information on additional devices used for performance assessment (e.g. tracking system used in a surgical setting).

Image1S 3D 30 degree optics (Karl Storz SE & CO KG).

• Imaging systems used:

Image1 Connect TC200, with resolution of 1080x1920, 25 fps

Image1S Connect TC200, 2160x3840, 25 fps

... for a stereo-pair saved in top-down format.

• Recorders used:

Karl Storz AIDA

DVI2PCIe capture card with Epiphan video capture software

b) Describe relevant details on the imaging process/data acquisition for each acquisition device (e.g. image acquisition protocol(s)).

Sim-Domain: Data was acquired on the minimally invasive training simulator (modified housing from (MICS MVR surgical simulator, Fehling Instruments GmbH & Co. KG, Karlstein, Germany). Camera angle is mainly from the upper left side. Light intensities 85-100%.

Intraop-Domain: Data was acquired during minimally invasive surgery (rightlateral thoracotomy). The distance and the camera orientation with respect to the mitral valve depends on the anatomical conditions of the patient. Light intensities varied between 85-100%.

c) Specify the center(s)/institute(s) in which the data was acquired and/or the data providing platform/source (e.g. previous challenge). If this information is not provided (e.g. for anonymization reasons), specify why.

The data was acquired at University Hospital Heidelberg (in MIC training lab and operating rooms).

d) Describe relevant characteristics (e.g. level of expertise) of the subjects (e.g. surgeon)/objects (e.g. robot) involved in the data acquisition process (if any).

The data sets are anonymized. No further characteristics are available.

Training and test case characteristics

a) State what is meant by one case in this challenge. A case encompasses all data that is processed to produce one result that is compared to the corresponding reference result (i.e. the desired algorithm output).

Examples:

- Training and test cases both represent a CT image of a human brain. Training cases have a weak annotation (tumor present or not and tumor volume (if any)) while the test cases are annotated with the tumor contour (if any).
- A case refers to all information that is available for one particular patient in a specific study. This information always includes the image information as specified in data source(s) (see above) and may include context information (see above). Both training and test cases are annotated with survival (binary) 5 years after (first) image was taken.

One case refers to one stereo frame extracted from the simulation sessions and the intra-operative videos.

Stereo-frames were saved in top-down format (left image top, right image bottom) and split.

Relevant scenes of mitral annuloplasty were identified before extracting the frames from the entire video recordings. Every 120th frame was extracted. In scenes with rapid changes, every 10th frame was extracted and in scenes with only few changes, every 240th frame was extracted.

b) State the total number of training, validation and test cases.

We do not release an explicit validation data set. Participants can split the training data accordingly during the training phase.

The current estimate of the training data in both the domains, and the test data is given below.

Training surgical simulator domain:

1600 stereo frames from 12 simulations with an average of 133 frames per session.

Training, Intra-operative domain:

6200 stereo frames from 10 surgeries/patients with an average 620 frames per surgery/patient.

Testing, Intra-operative domain:

450 stereo frames from 5 surgeries/patients with an average of 90 frames per surgery/patient.

c) Explain why a total number of cases and the specific proportion of training, validation and test cases was chosen.

To reflect a real-world scenario, the split must be conducted on the level of the patients/simulations. The training-test ratio is 14:5, which means that 26% are used for testing. According to [3], the median ratio of training cases to test cases in past challenges is 0.75.

The idea behind the challenge is to keep the number of intraoperative patients lower to encourage participants to incorporate the frames from Sim-Domain in the training process to achieve better generalization performance.

d) Mention further important characteristics of the training, validation and test cases (e.g. class distribution in classification tasks chosen according to real-world distribution vs. equal class distribution) and justify the choice.

Additionally, data is collected from multiple phases of the surgery such as suturing, ring implantation, ring knotting, inspection, saline test, etc. Each surgery or patient contains data from multiple phases of the surgery. In addition, the test set is curated to avoid over-dependence on data from one patient, i.e. the test data comprises of a small number of frames but from diverse surgeries with different camera angles, illumination and field of view.

Annotation characteristics

a) Describe the method for determining the reference annotation, i.e. the desired algorithm output. Provide the information separately for the training, validation and test cases if necessary. Possible methods include manual image annotation, in silico ground truth generation and annotation by automatic methods.

If human annotation was involved, state the number of annotators.

No human annotation involved

b) Provide the instructions given to the annotators (if any) prior to the annotation. This may include description of a training phase with the software. Provide the information separately for the training, validation and test cases if necessary. Preferably, provide a link to the annotation protocol.

-

c) Provide details on the subject(s)/algorithm(s) that annotated the cases (e.g. information on level of expertise such as number of years of professional experience, medically-trained or not). Provide the information separately for the training, validation and test cases if necessary.

-

d) Describe the method(s) used to merge multiple annotations for one case (if any). Provide the information separately for the training, validation and test cases if necessary.

-

Data pre-processing method(s)

Describe the method(s) used for pre-processing the raw training data before it is provided to the participating teams. Provide the information separately for the training, validation and test cases if necessary.

Stereo-frames were acquired in top-down format (left image on the top, right image at the bottom), and split and named with the corresponding file names. The split images were resized to 540x960, which maintains the aspect ratio of the original image.

Sources of error

a) Describe the most relevant possible error sources related to the image annotation. If possible, estimate the magnitude (range) of these errors, using inter-and intra-annotator variability, for example. Provide the information separately for the training, validation and test cases, if necessary.

-

b) In an analogous manner, describe and quantify other relevant sources of error.

Some scenes have a poor view of the mitral valve owing to endoscope artefacts. For example, in some scenes the valve is occluded by tissue, or is marred by heavy fogging, or is affected by motion artefact due to camera motion. We remove these scenes from the dataset.

ASSESSMENT METHODS

Metric(s)

a) Define the metric(s) to assess a property of an algorithm. These metrics should reflect the desired algorithm properties described in assessment aim(s) (see above). State which metric(s) were used to compute the ranking(s) (if any).

- Example 1: Dice Similarity Coefficient (DSC)
- Example 2: Area under curve (AUC)

The metric used for evaluation is a combination of local distance-based metric (L1) and a perceptual metric (SSIM).

b) Justify why the metric(s) was/were chosen, preferably with reference to the biomedical application.

A combination of local and perceptual metrics have been found to be a more reliable indicator of image similarity than only one of the metrics [4]. Local metrics such as L1 assume independence of pixels and ignore local context. Using SSIM metrics take into account the neighbourhood information of a pixel, but have issues at image boundaries unlike the local image metrics.

Ranking method(s)

a) Describe the method used to compute a performance rank for all submitted algorithms based on the generated metric results on the test cases. Typically the text will describe how results obtained per case and metric are aggregated to arrive at a final score/ranking.

The evaluation is performed on the intra-operative test dataset using weighted L1 and SSIM.

In case the metrics are tied between two teams, the time stamp of the submission will be used for the ranking, and the earlier submission will have a higher ranking.

b) Describe the method(s) used to manage submissions with missing results on test cases.

In the case where no predictions are provided, the similarity between the target image and a zero image is lower compared to when a prediction is provided. Similarly, the local metric determined by L1 yields a higher error compared to when no predictions are provided.

c) Justify why the described ranking scheme(s) was/were used.

Both metrics are commonly used for novel view synthesis.

Statistical analyses

a) Provide details for the statistical methods used in the scope of the challenge analysis. This may include

- description of the missing data handling,
- details about the assessment of variability of rankings,
- description of any method used to assess whether the data met the assumptions, required for the particular statistical approach, or
- indication of any software product that was used for all data analysis methods.

We will assess the ranking variability with Kendall's tau analysis. In particular, we will investigate whether only using L1 as metric or SSIM will lead to different challenge rankings.

The results during the challenge event will be reported with transparency.

b) Justify why the described statistical method(s) was/were used.

Kendall's tau may quantify differences between rankings (1: identical ranking; -1: inverse ranking). However, even for high values of Kendall's tau, critical changes in the ranking may occur [3]. Therefore, we will additionally provide the complete alternative ranking lists.

Further analyses

Present further analyses to be performed (if applicable), e.g. related to

- combining algorithms via ensembling,
- inter-algorithm variability,
- common problems/biases of the submitted methods, or
- ranking variability.

In a surgical training scenario, an additional idea is to transform not-so-realistic phantom data into more realistic surgical images [8]. Therefore, we encourage the participants to use image-to-image translation approaches to transform simulated data into more realistic appearance, however, this is not mandatory.

ADDITIONAL POINTS

References

Please include any reference important for the challenge design, for example publications on the data, the annotation process or the chosen metrics as well as DOIs referring to data or code.

[1] Carpentier, A., Deloche, A., Dauptain, J., Soyer, R., Blondeau, P., Piwnica, A., Dubost, C., McGoon, D.C.: A new reconstructive operation for correction of mitral and tricuspid insufficiency. *The Journal of Thoracic and Cardiovascular Surgery* 61 (1), 1–13 (1971)

[2] Casselman Filip P., Van Slycke Sam, Wellens Francis, De Geest Raphael, Degrieck Ivan, Van Praet Frank, Vermeulen Yvette, Vanermen Hugo: Mitral Valve Surgery Can Now Routinely Be Performed Endoscopically. *Circulation* 108 (10 suppl 1), II-48 (2003). DOI <https://doi.org/10.1161/01.cir.0000087391.49121.ce>

[3] Maier-Hein, L., Eisenmann, M., Reinke, A. et al. Why rankings of biomedical image analysis competitions should be interpreted with care. *Nat Commun* 9, 5217 (2018). <https://doi.org/10.1038/s41467-018-07619-7>

[4] Zhao, H., Gallo, O., Frosio, I., & Kautz, J. (2018). Loss Functions for Neural Networks for Image Processing. ArXiv:1511.08861 [Cs]. <http://arxiv.org/abs/1511.08861>

[5] Watson, J., Mac Aodha, O., Turmukhambetov, D., Brostow, G. J., & Firman, M. (2020). Learning Stereo from Single Images. ArXiv:2008.01484 [Cs]. <http://arxiv.org/abs/2008.01484>

[6] Hou, Y., Solin, A., & Kannala, J. (2021). Novel View Synthesis via Depth-guided Skip Connections. ArXiv:2101.01619 [Cs]. <http://arxiv.org/abs/2101.01619>

[7] <https://adaptor2021.github.io/> doi: 10.5281/zenodo.4646979

[8] Engelhardt S., De Simone R., Full P.M., Karck M., Wolf I. (2018) Improving Surgical Training Phantoms by Hyperrealism: Deep Unpaired Image-to-Image Translation from Real Surgeries. In: Frangi A., Schnabel J., Davatzikos C., Alberola-López C., Fichtinger G. (eds) Medical Image Computing and Computer Assisted Intervention – MICCAI 2018. MICCAI 2018. Lecture Notes in Computer Science, vol 11070. Springer, Cham, doi: 10.1007/978-3-030-00928-1_