

Shape-based pose estimation for automatic standard views of the knee

Lisa Kausch^{1,2,3}, Sarina Thomas⁴, Holger Kunze⁵, Jan Siad El Barbari⁶, and Klaus Maier-Hein^{1,2,3}

¹ Division of Medical Image Computing, German Cancer Research Center (DKFZ), Heidelberg, Germany

l.kausch@dkfz-heidelberg.de

² National Center for Tumor Diseases (NCT) Heidelberg, Germany

³ Pattern Analysis and Learning Group, Department of Radiation Oncology, Heidelberg University Hospital, Germany

⁴ Department of Informatics, University of Oslo, Norway

⁵ Advanced Therapy Systems Division, Siemens Healthineers, Erlangen, Germany

⁶ MINTOS Research Group, Trauma Surgery Clinic Ludwigshafen, Germany

Abstract. Surgical treatment of complicated knee fractures is guided by real-time imaging using a mobile C-arm. Immediate and continuous control is achieved via 2D anatomy-specific standard views that correspond to a specific C-arm pose relative to the patient positioning, which is currently determined manually, following a trial-and-error approach at the cost of time and radiation dose. The characteristics of the standard views of the knee suggests that the shape information of individual bones could guide an automatic positioning procedure, reducing time and the amount of unnecessary radiation during C-arm positioning. To fully automate the C-arm positioning task during knee surgeries, we propose a complete framework that enables (1) automatic laterality and standard view classification and (2) automatic shape-based pose regression toward the desired standard view based on a single initial X-ray. A suitable shape representation is proposed to incorporate semantic information into the pose regression pipeline. The pipeline is designed to handle two distinct standard views simultaneously. Experiments were conducted to assess the performance of the proposed system on 3528 synthetic and 1386 real X-rays for the a.-p. and lateral standard. The view/laterality classifier resulted in an accuracy of 100%/98% on the simulated and 99%/98% on the real X-rays. The pose regression performance was $d\theta_{a.-p} = 5.8 \pm 3.3^\circ$, $d\theta_{lateral} = 3.7 \pm 2.0^\circ$ on the simulated data and $d\theta_{a.-p} = 7.4 \pm 5.0^\circ$, $d\theta_{lateral} = 8.4 \pm 5.4^\circ$ on the real data outperforming intensity-based pose regression.

Keywords: Shape-based pose estimation · Standard projections · Knee.

1 Introduction

Intraoperative imaging employing a mobile C-arm enables immediate and continuous control during orthopedic and trauma interventions. For optimal frac-

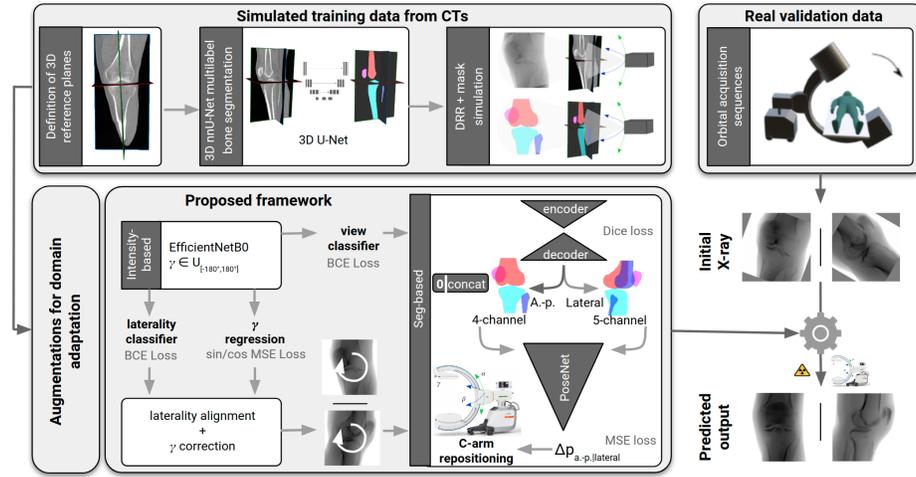


Fig. 1: Proposed shape-based pose estimation framework for automatic acquisition of standard views. A single architecture is trained for the representation of 2 distinct standard views simultaneously. The 2-step pipeline consists of a direct intensity-based combined view classification and in-plane rotation regression, followed by a segmentation-based pose regression focusing on the out-of-plane rotation. The pipeline is solely trained on synthetic data with automatically generated ground truth annotations and evaluated on real X-rays.

ture reduction and implant placement, correct acquisition of standard views that correspond to a specific C-arm pose relative to the patient is essential [17]. Incorrect standard views can exhibit superimposed anatomical structures, leading to overlooked errors that can result in malunion of fractures, functional impairment, or require revision surgeries. To enable deducing all three dimensions of the trauma case, at least two 2D fluoroscopic views are acquired in two distinct planes usually at right angles to each other. The current manual C-arm positioning procedure results in only 20% surgically relevant acquisitions while the remaining 80% are caused by the iterative positioning process, exposing patients and clinical staff to unnecessary radiation [15].

Recent developments towards robotic C-arms ask for automatic positioning methods. Many state of the art approaches require a patient-specific CT for intraoperative real-time simulation [4,5,8] or 2D-3D registration [2,7,16], external tracking equipment [6,15], manual landmark annotation [1] or do not estimate an optimal pose but reproduce intraoperatively recorded C-arm views employing augmented reality [6,19]. The inherent prior assumptions and severe inference with the clinical workflow hinder broad clinical applicability until today. In contrast to the majority of anatomical regions, the standard planes of the knee are not orthogonal to each other. The a.-p. standard view is characterized by symmetric projection of the joint gap, femoral and tibial condyles. The tibia surface projects line-shaped and the medial half of the fibula head is super-

imposed by the tibia. In the lateral standard view, both femoral condyles are aligned and the joint gap is maximized. Automatic deep learning-based positioning for standard views involves specific challenges for image understanding due to overlapping anatomical structures, the presence of surgical implants, and changing viewing directions and showed to benefit from extracting semantic information [10]. Inspired by that and considering that standard views of the knee anatomy are characterized by the shape information of the individual projected bones, we propose a complete framework to fully automate the C-arm positioning tasks during knee surgeries (Fig. 1). Our contribution is 4-fold: (1) We propose a novel framework that enables simultaneous automatic standard view classification, laterality classification, in-plane rotation correction, and subsequent view-independent shape-based C-arm positioning to the desired standard view while requiring only a single initial X-ray projection. One pose regression network can handle two distinct standard views of the knee anatomy. A suitable segmentation representation for the knee anatomy is proposed to recognize correct standard views, which explicitly incorporates semantic information to reflect on the actual clinical decision-making of surgeons. Since intraoperative X-rays with reference pose annotations do not exist, the proposed framework is solely trained on simulated data with automatically generated pose annotations. (2) We show that the proposed approach outperforms view-specific shape-based and intensity-based pose regression. (3) We show that the proposed shape representation and augmentation strategies aid generalization from simulated training data to real cadaver X-rays. (4) We investigate the importance of individual knee bones on the overall positioning performance for two distinct standard views.

2 Materials and Methods

An overview of the complete framework for fully automated C-arm positioning towards desired standard views during knee surgeries is given in Fig.1. The anterior-posterior (a.-p.) and the lateral standard view showed to be sufficient for various diagnostic entities [3].

2.1 Training data simulation

To address the interventional data scarcity problem, simulated training data was generated from a collection of CT and C-arm volumes using a realistic DRR simulation framework [20] complemented with corresponding 2D segmentations. Preprocessing involved the following steps: **(1) Field-of-view cropping:** Prevents superposition of the other laterality in the projection domain. **(2) Laterality alignment:** Prohibits ambiguities during pose estimation. **(3) Definition of 3D reference planes:** Two independent raters defined the 3D reference planes in the CT volumes utilizing a DRR preview integrated into the open-source Medical Imaging Toolkit [21] with interactive plane positioning. They serve as ground truth pose reference during simulation. **(4) 3D automatic bone segmentation:** To compute automatic 3D segmentations, a 3D

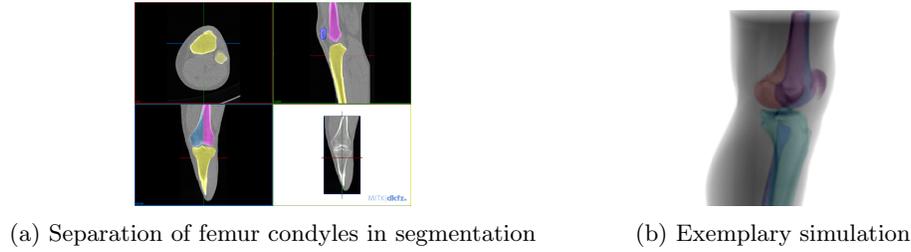


Fig. 2: Suitable segmentation representation to recognize the lateral standard view of the knee: In an ideal standard view the condyles’ overlap each other.

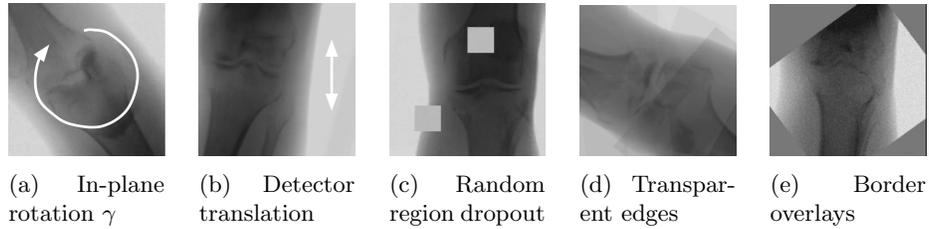


Fig. 3: Training data augmentations.

nnU-Net [9] is trained on a subset of 10 manually annotated CTs for the task of multilabel bone segmentation, segmenting the femur, tibia, fibula, and patella. **(5) Suitable segmentation representation:** The two femur condyles are not distinguishable in the shape-based representation, however, this is relevant for optimal lateral view recognition. Annotating the condyles as line features would result in an increased manual labeling effort in the projection domain. Alternatively, we propose to incorporate this information in the segmentation by separating the femur annotation symmetrically along the femoral shaft (Fig. 2a). This results in one additional segmentation label for the lateral standard to recognize condyles’ congruence and derive the directional pose offset. **(6) DRR and mask simulation:** DRRs are simulated for varying angulations of orbital and angular rotation $\alpha, \beta \in [-40^\circ, 40^\circ]$ around the defined reference standard. The DeepDRR simulation framework was extended to allow the forward projection of corresponding masks (Fig 2b). A set of augmentations is applied to the simulated dataset, to bridge the domain gap from synthetic data to real X-rays (Fig. 3). The in-plane rotation augmentation accounts for variable patient to C-arm alignment, the translation bridges the gap to the validation data where the joint gap is not centered like in the training data, random region dropout reflects superposition artifacts, transparent edge overlays reflect projection artifacts resulting from, e.g., the operating table, and border overlays account for the gamma correction interpolation artifacts. For simulation, a set of 24 CTs was considered, 15 CTs without metal and 9

C-arm volumes with metal. The data was divided 60 – 20 – 20% for training (15 CTs), validation (5 CTs), test (4 CTs).

2.2 Shape-based positioning framework:

The proposed shape-based positioning framework was trained jointly for both standard views (Fig. 1). It consists of two modules: The first is responsible for a view classification, in-plane rotation and laterality alignment, directly estimated from the image intensities. Thereby, pose ambiguities and data variation are addressed, simplifying the task for the subsequent module to estimate the optimal C-arm positioning for the desired standard view, employing shape features.

(1) Intensity-based multi-task classification and regression module:

For simultaneous in-plane rotation regression, view recognition, and laterality classification, an EfficientNet-B0 feature extractor [18,14] was extended with two binary classification heads with one output neuron, followed by sigmoid activation, and one regression head, with the same architecture, but 2 outputs, omitting the activation. The in-plane rotation γ is mapped to sin/cos-space to ensure a continuous Loss function during optimization. All training examples were aligned with the same laterality during data simulation which would otherwise result in pose ambiguities. Thus, to train the laterality classifier, the training examples were randomly flipped horizontally with $p = 0.5$ and the corresponding γ label was adapted accordingly. The weights were optimized using Binary Cross Entropy Loss for the classification heads, and Mean Squared Error for the regression head.

(2) Shape-based pose regression:

Following surgical characteristics for recognizing correct standard views of the knee, a view-independent shape-based pose regression framework was developed. The architecture is based on a 2D U-Net [12] with two view-specific segmentation heads, because the segmentation labels differ for both views. The extracted shape features are used as input for the pose regression network that outputs the necessary C-arm pose update $(\alpha, \beta, \gamma, \mathbf{t}) \in \mathbf{R}^6$ to acquire the desired standard view [11].

2.3 Validation data:

Real X-rays for validation were sampled from single Siemens Cios Spin[®] sequences generated during 3D acquisition of 6 knee cadavers. Preprocessing consisted of (1) definition of 3D standard reference planes, (2) laterality check, (3) sampling of X-rays around the defined reference standards in the interval $\alpha, \beta \in [-30^\circ, 30^\circ]$, and generation of ground truth pose labels. Since the Spin sequences are orbital acquisition sequences, only the orbital rotation is equidistantly covered in the validation set, while the angular rotation is constant for all X-rays sampled from the same sequence. The number of sampled X-rays per standard and view may differ, if the reference standard is located close to the edge of the orbital sequence (range: 102-124).

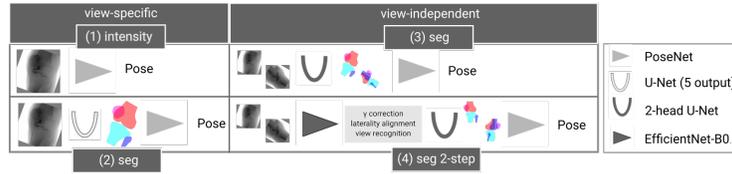


Fig. 4: Variants for performance comparison.

3 Experiments and Results

The proposed pipeline was evaluated considering the following research questions:

- (RQ1) Does the proposed shape-based pipeline outperform view-specific intensity-based and shape-based pose regression? How does it influence the generalization from synthetic to real data? (Sec. 3.1)
- (RQ2) How do individual bones influence the overall positioning performance? (Sec. 3.2)
- (RQ3) How accurate is the performance of the view and laterality classification? (Sec. 3.3)

Positioning performance was evaluated based on the angle $\theta = \arccos(\langle v_{pred}, v_{gt} \rangle)$ between the principal rays of the ground truth v_{gt} and predicted pose v_{pred} and the mean absolute error (AE) of in-plane rotation γ . The interrater variation of the reference standard planes defined by two independent raters serves as an upper bound for the reachable accuracy of a C-arm positioning approach trained on the reference annotations. It was assessed in terms of orientation differences θ ($\theta_{a.-p.} = 4.1 \pm 2.6^\circ$, $\theta_{lateral} = 1.8 \pm 1.3^\circ$).

The models were implemented using PyTorch 1.6.0, trained end-to-end with an 11 GB GeForce RTX 2080 Ti, and optimized with the Adam optimizer with a base learning rate of $\eta = 10^{-4}$ and batchsize 8 until convergence.

3.1 Importance of pipeline design choices (RQ1)

In an ablation study, the proposed shape-based view-independent pose regression was compared to view-specific direct intensity-based pose regression [11]. Further, the complete pipeline (2-step) is compared to a 1-step segmentation-based approach trained view-specific and view-independent (Fig. 4). Evaluation was performed on the simulated test DRRs and cadaveric X-rays (Fig. 5).

View-independent vs. view-specific networks: While the view-specific networks perform significantly better (lateral) or comparable (a.-p.) on the simulated data, the view-independent networks perform significantly better (a.-p.) or comparable (lateral) on the real data.

1-step vs. 2-step: The proposed 2-step approach performs significantly better or comparable than a 1-step shape-based pose regression approach on most validation cases (8/12) in viewing direction θ . Regarding the γ rotation, the 2-step

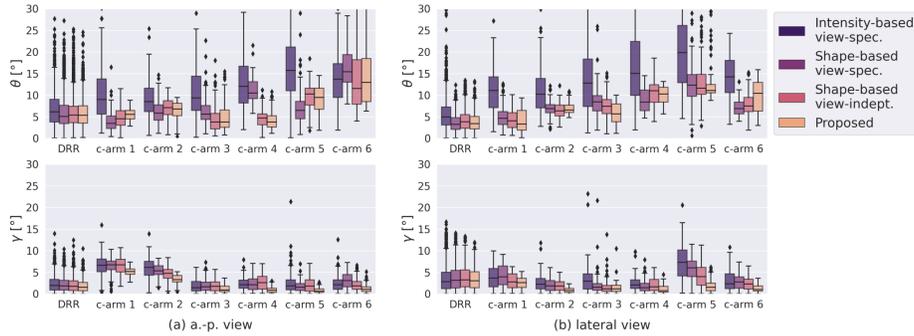


Fig. 5: Pose regression performance on simulated and real X-rays compared for different pipeline variants.

approach improves performance across all validation cases.

Generalization from DRR to X-ray: The shape-based pose regression network combined with joint view-independent training clearly boosts the performance compared to direct intensity-based pose regression from $d\theta_{a.-p}^{X-ray} = 12.2 \pm 6.8^\circ$, $d\theta_{lateral}^{X-ray} = 14.4 \pm 7.6^\circ$ to $d\theta_{a.-p}^{X-ray} = 7.4 \pm 5.0^\circ$, $d\theta_{lateral}^{X-ray} = 8.4 \pm 5.4^\circ$.

3.2 Importance of individual bones on overall performance (RQ2)

Fig. 7 shows the importance of individual segmented bone classes on the overall positioning performance evaluated on the test DRRs (3528 DRRs). The fibula has very little influence on the positioning for both views. The patella is only important for the a.-p. view, while tibia and femur are relevant for both views. The condyle assignment for the lateral view determines the rotation direction for the orbital and angular rotation (α , β). Inverting the assignment of left and right femur condyle results in a sign flip in α , β .

3.3 Accuracy of view and laterality classification (RQ3)

The classifier performances were assessed on the synthetic (3528 DRRs, 4 CTs) and real data (1386 X-rays, 6 C-arm scans). The view classifier (a.-p. / lateral) achieved an accuracy of 100% on the test DRRs and 99% on the X-rays. The laterality classifier (left / right) resulted in an accuracy of 98% on the test DRRs and 98% on the X-rays.

4 Discussion and Conclusion

A complete framework for automatic acquisition of standard views of the knee is proposed that can handle several standard views simultaneously. The complete pipeline is trained on simulated data with automatically generated annotations

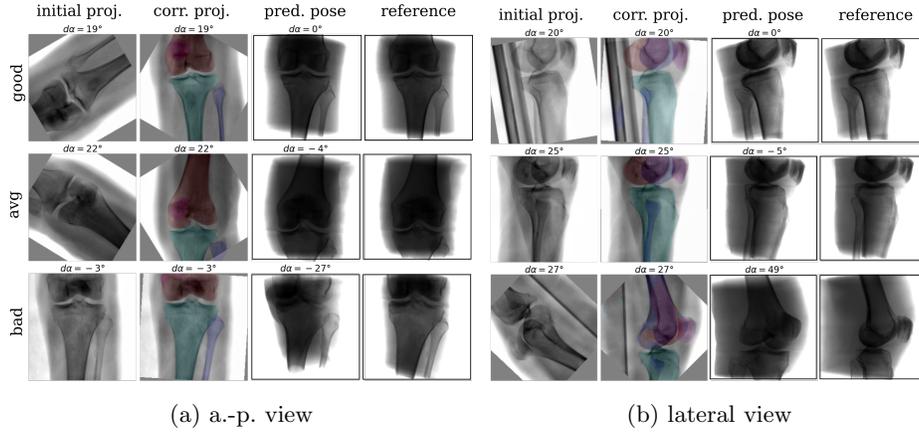


Fig. 6: Exemplary visual results for the a.-p. and lateral standard positioning with good, average, and bad performance measured with respect to the orbital rotation offset α to the reference standard pose. The initial real X-ray projection is visualized along with the in-plane corrected projection with the overlaid segmentation mask, and the predicted standard view simulated after C-arm repositioning according to the predicted correction pose side-by-side with the desired reference standard.

and evaluated on real intraoperative X-rays. To bridge the domain gap, different augmentation strategies are suggested that address intraoperative confounding factors, e.g., the OR table. View-independent training and multi-label shape features improve the generalization from simulated training to real X-rays and outperform direct intensity-based approaches. View-independent networks result in more training data which showed to improve the generalization from simulated training to real X-rays. The 2-step approach increases robustness and simultaneously automates necessary preprocessing tasks like laterality and standard view recognition, which can be performed with very high accuracy on simulated (100%, 98%) and real data (99%, 98%). The approach is fast and easy to translate into the operating room as it does not require any additional technical equipment. Assuming that the surgeon acquires the initial X-ray with a pose offset within the capture range of $[-30^\circ, 30^\circ]$, it has the potential to reduce time and unnecessary radiation during manual C-arm positioning. Furthermore, the segmentation features can serve as a sanity check and indicate the reliability of the pose regression result. Further experiments with a larger training set covering more anatomical variation, e.g., patella baja and different flexion angles [13], can potentially address observed failure cases.

Data use declaration: The data was obtained retrospectively from anonymized databases and not generated intentionally for the study. The acquisition of data from living patients had a medical indication and informed consent was not required. The corresponding consent for body donation for these purposes has been obtained.

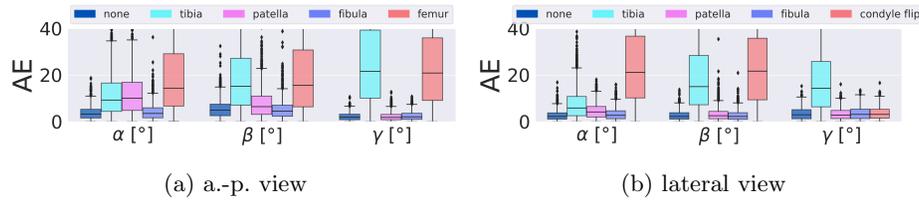


Fig. 7: Importance of individual bones on positioning performance (3528 test DRRs). One segmentation channel was set to zero at a time during inference, 'none' corresponds to the reference performance utilizing all channels.

References

1. Binder, N., Bodensteiner, C., Matthäus, L., Burgkart, R., Schweikard, A.: Image guided positioning for an interactive c-arm fluoroscope. In: *Computer Assisted Radiology and Surgery*. pp. 5–7 (2006)
2. Bott, O.J., Dresing, K., Wagner, M., Raab, B.W., Teistler, M.: Use of a C-arm fluoroscopy simulator to support training in intraoperative radiography. *Radiographics* **31**(3), E31–E41 (2011)
3. Cockshott, W.P., Racoveanu, N., Burrows, D., Ferrier, M.: Use of radiographic projections of knee. *Skeletal radiology* **13**(2), 131–133 (1985)
4. De Silva, T., Punnoose, J., Uneri, A., Goerres, J., Jacobson, M., Ketcha, M.D., Manbachi, A., Vogt, S., Kleinszig, G., Khanna, A.J., et al.: C-arm positioning using virtual fluoroscopy for image-guided surgery. In: *Medical Imaging: Image-guided Procedures, Robotic Interventions, and Modeling*. vol. 10135, p. 101352K. International Society for Optics and Photonics (2017)
5. Fallavollita, P., Winkler, A., Habert, S., Wucherer, P., Stefan, P., Mansour, R., Ghotbi, R., Navab, N.: Desired-view controlled positioning of angiographic C-arms. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 659–666. Springer (2014). https://doi.org/10.1007/978-3-319-10470-6_82
6. Fotouhi, J., Unberath, M., Song, T., Gu, W., Johnson, A., Osgood, G., Armand, M., Navab, N.: Interactive flying frustums (IFFs): Spatially aware surgical data visualization. *International Journal of Computer Assisted Radiology and Surgery* **14**(6), 913–922 (2019). <https://doi.org/10.1007/s11548-019-01943-z>
7. Gong, R.H., Jenkins, B., Sze, R.W., Yaniv, Z.: A cost effective and high fidelity fluoroscopy simulator using the image-guided surgery toolkit (IGSTK). In: *Medical Imaging: Image-Guided Procedures, Robotic Interventions, and Modeling*. vol. 9036, p. 903618. International Society for Optics and Photonics (2014). <https://doi.org/10.1117/12.2044112>
8. Haiderbhai, M., Turrubiates, J.G., Gutta, V., Fallavollita, P.: Automatic C-arm positioning using multi-functional user interface. *Canadian Medical and Biological Engineering Society Proceedings* **42** (2019)
9. Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods* **18**(2), 203–211 (2021). <https://doi.org/10.1038/s41592-020-01008-z>
10. Kausch, L., Thomas, S., Kunze, H., Norajitra, T., Klein, A., El Barbari, J., Privalov, M., Vetter, S., Mahnken, A.H., Maier-Hein, L., Maier-Hein, K.: C-arm po-

- sitioning for spinal standard projections in different intra-operative setting. International Conference on Medical Image Computing and Computer-Assisted Intervention pp. 352–362 (2021). https://doi.org/10.1007/978-3-030-87202-1_34
11. Kausch, L., Thomas, S., Kunze, H., Privalov, M., Vetter, S., Franke, J., Mahnken, A.H., Maier-Hein, L., Maier-Hein, K.: Toward automatic C-arm positioning for standard projections in orthopedic surgery. *International Journal of Computer Assisted Radiology and Surgery* **15**(7), 1095–1105 (2020). <https://doi.org/10.1007/s11548-020-02204-0>
 12. Klein, A., Wasserthal, J., Greiner, M., Zimmerer, D., Maier-Hein, K.H.: `basic_unet_example` (v2019.01) (2019). <https://doi.org/10.5281/zenodo.2552439>
 13. Krönke, S., von Berg, J., Brueck, M., Bystrov, D., Gooßen, A., Harder, T., Lundt, B., May, J.M., Wieberneit, N., Wissel, T., et al.: Cnn-based pose estimation for assessing quality of ankle-joint x-ray images. In: *Medical Imaging 2022: Image Processing*. vol. 12032, pp. 344–352. SPIE (2022)
 14. Mairhöfer, D., Laufer, M., Simon, P.M., Sieren, M., Bischof, A., Käster, T., Barth, E., Barkhausen, J., Martinetz, T.: An ai-based framework for diagnostic quality assessment of ankle radiographs. In: *Medical Imaging with Deep Learning* (2021)
 15. Matthews, F., Hoigne, D.J., Weiser, M., Wanner, G.A., Regazzoni, P., Suhm, N., Messmer, P.: Navigating the fluoroscope’s C-arm back into position: An accurate and practicable solution to cut radiation and optimize intra-operative workflow. *Journal of Orthopaedic Trauma* **21**(10), 687–692 (2007). <https://doi.org/10.1097/BOT.0b013e318158fd42>
 16. Miao, S., Piat, S., Fischer, P., Tuysuzoglu, A., Mewes, P., Mansi, T., Liao, R.: Dilated FCN for multi-agent 2D/3D medical image registration. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 32 (2018)
 17. Norris, B.L., Hahn, D.H., Bosse, M.J., Kellam, J.F., Sims, S.H.: Intraoperative fluoroscopy to evaluate fracture reduction and hardware placement during acetabular surgery. *Journal of Orthopaedic Trauma* **13**(6), 414–417 (1999). <https://doi.org/10.1097/00005131-199908000-00004>
 18. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: *International conference on machine learning*. pp. 6105–6114. PMLR (2019)
 19. Unberath, M., Fotouhi, J., Hajek, J., Maier, A., Osgood, G., Taylor, R., Armand, M., Navab, N.: Augmented reality-based feedback for technician-in-the-loop C-arm repositioning. *Healthcare Technology Letters* **5**(5), 143–147 (2018). <https://doi.org/10.1049/htl.2018.5066>
 20. Unberath, M., Zaech, J.N., Lee, S.C., Bier, B., Fotouhi, J., Armand, M., Navab, N.: DeepDRR – a catalyst for machine learning in fluoroscopy-guided procedures. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 98–106. Springer (2018)
 21. Wolf, I., Vetter, M., Wegner, I., Böttger, T., Nolden, M., Schöbinger, M., Hasenteufel, M., Kunert, T., Meinzer, H.P.: The medical imaging interaction toolkit. *Medical image analysis* **9**(6), 594–604 (2005)