

EXPLOITING SEGMENTATION LABELS AND REPRESENTATION LEARNING TO FORECAST THERAPY RESPONSE OF PDAC PATIENTS

Alexander Ziller^{1,*}, Ayhan Can Erdur^{1,*}, Friederike Jungmann¹
Daniel Rueckert^{1,2}, Rickmer Braren^{1,4}, Georgios Kaissis^{1,3}

¹ Klinikum Rechts der Isar, Technical University of Munich

² Imperial College London

³ Helmholtz Zentrum München

⁴ German Cancer Consortium DKTK, Partner Site Munich

ABSTRACT

The prediction of pancreatic ductal adenocarcinoma therapy response is a clinically challenging and important task in this high-mortality tumour entity. The training of neural networks able to tackle this challenge is impeded by a lack of large datasets and the difficult anatomical localisation of the pancreas. Here, we propose a hybrid deep neural network pipeline to predict tumour response to initial chemotherapy which is based on the *Response Evaluation Criteria in Solid Tumors* (RECIST) score, a standardised method for cancer response evaluation by clinicians as well as tumour markers, and clinical evaluation of the patients. We leverage a combination of representation transfer from segmentation to classification, as well as localisation and representation learning. Our approach yields a remarkably data-efficient method able to predict treatment response with a ROC-AUC of 63.7% using only 477 datasets in total.

Index Terms— personalised treatment, PDAC, representation learning, transfer learning

1. INTRODUCTION

Pancreatic ductal adenocarcinoma (PDAC) is among the tumour entities with the highest mortalities worldwide. Its diagnosis is challenging and relies fundamentally upon high-resolution imaging such as computed tomography (CT) imaging. Imaging likely yields a wealth of information about the tumour no imaging facts beyond the localisation and size of the lesion, its anatomical surroundings as well as the presence or absence of tumour spreading. However, this information is not considered for diagnosing PDAC and evaluating its treatment response according to the current guidelines [1]. With the rise of deep learning, models aiming to extract the aforementioned information and assist physicians in the assessment of such tumours have begun to appear [2]. For example, machine learning-assisted prediction of molecular

tumour subtype and patient prognosis under therapy have recently been presented [3, 4, 5]. As many PDAC are discovered at later stages where a primary resection is infeasible, chemotherapy is a first-line treatment for many patients, either to shrink the tumour prior to an operation or as a palliative indication. In this setting, the prediction of tumour response to initial treatment is a particularly interesting end-point. For example, deep neural networks trained to predict treatment response conditioned on the specific choice of chemotherapy can in the future be employed to select the personalised treatment with the highest success probability and/or best outcome. In clinical practice, standardised metrics of tumour response already exist: The *Response Evaluation Criteria in Solid Tumors* (RECIST) [6, 7] are a framework to classify whether a tumour is progressive, stable or regressive under therapy. Although RECIST can be used to derive labels for training neural networks, a number of challenges have so far inhibited the success of such machine learning approaches:

- As in most medical problems, sufficiently large datasets to train deep neural networks are difficult to acquire;
- Abdominal CT scans cover a large volume of the body, whereas the pancreas (and potential extrapancreatic tumour manifestations) represent(-s) a proportionally small area of interest. Thus, there exists a risk that models which operate on the entire CT volume miss relevant information in the scan;
- Compounding the two problems, methods which are capable of focusing on relevant areas within the image (e.g. attention-based architectures) are—in practice— even more data-hungry [8].

Our approach

To overcome the problem of data scarcity, we propose to leverage knowledge from related public datasets *but from an unrelated task*. Concretely, we propose a multi-step pipeline incorporating representations learned from a segmentation task into our classification model. We realise this through initially training a hybrid segmentation and object detection model to pre-crop the CT volume to the pancreatic region.

*equal contribution

Imitating the human approach to geometric pattern matching, we feed the segmentation masks into the classification model as assistive information. Moreover, our pipeline allows us to re-use the segmentation weights for the classification task, further increasing data efficiency. Finally, we use triplet loss-guided learning to enhance the quality of the classification model’s intermediate representations. Our contributions can be summarised as follows:

- We propose an integrated deep learning pipeline based on cascading a slice classifier and a segmentation model and complemented by transfer and representation learning. This results in –to our knowledge– the first successful approach to response prediction to chemotherapy from baseline CT scans in PDAC using fewer than 500 total datasets;
- We find that the incorporation of each of the intermediate steps outlined above in the pipeline yields substantial improvements in terms of performance and/or efficiency;
- We perform a detailed ablation study of the aforementioned intermediate components.

Prior work

So far, very few works have tackled the challenging task of treatment response prediction in PDAC. For example, [9] break the task down to a comparison instead of a forecasting task. Other relevant previous approaches estimate diameters of lesions, which are an important part of the RECIST criteria [10]. In terms of methodology, the closest work to our approach are [11], who use lung segmentation models to classify chest X-rays. Similarly to the second stage of our approach this work uses an upstream model to predict a segmentation, which is used to generate a cropping bounding box over a region of interest. Moreover, [12, 13] simultaneously learn segmentation and classification tasks. However, these works have both labels available for both tasks, whereas we learn the tasks disjointly with only one label being available per dataset and task. Finally, our approach to utilise autoencoder-like architectures for representation learning has also been explored in [14, 15], however –unlike us– not with segmentation architectures.

2. METHODS

As discussed above, we propose a multi-stage pipeline consisting of three sequential convolutional neural networks (CNNs) to increase the chemotherapy response classification performance. Furthermore, we evaluate the effect of each individual component in our pipeline.

2.1. Dataset and Tasks

Our approach is motivated by the fact that we only have access to two datasets with disjoint labels. For the first two stages, we are dependent on segmentation labels. As these are not present in the classification dataset, we use the pub-

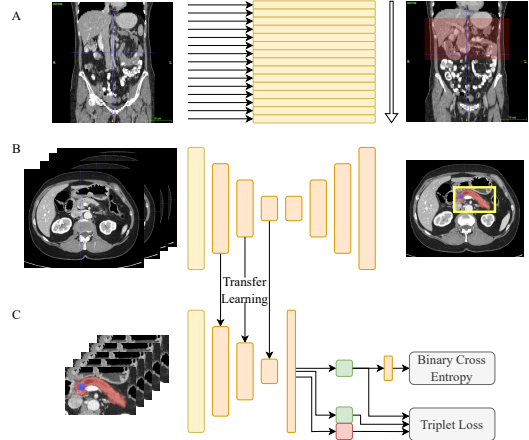


Fig. 1. Overview of our approach. A: We classify each slice along the z -direction whether it contains pancreatic tissue in a 2.5-D approach. B: We use a segmentation network, to crop in x and y -direction, as well as feeding segmentation masks to the classification model. We re-use the network weights of the feature extractor part for the classification stage. C: We train a classifier, first using a triplet loss and afterwards using binary cross entropy. Images are extracted from the Medical Segmentation Decathlon dataset.

lic *Medical Segmentation Decathlon* [16] dataset. It contains $N = 281$ abdominal CT scans alongside per-voxel labels whether a point is part of the pancreas, cancerous tissue or background.

We use an in-house dataset of $N = 477$ CT scans for the classification task, where we predict treatment response labels, grouped into a binary prediction of progressive disease vs. stable and regressive diseases. As the clinically important distinction is whether patients experienced tumour progression or not, we perform a binary prediction on progressive cases (PD) against all other stable and regressive cases. This yields a class distribution of $N = 171$ progressive vs $N = 306$ stable or regressive CT scans. We split the dataset into a training subset of $N = 420$ (151PD/269 other) samples and a testing subset of $N = 57$ (20PD/37 other) samples.

In the following, we describe the general model we train, as well as additional pre-processing by segmentation models and representation learning.

2.2. Classification

Our overall goal is to train a network that predicts the treatment response in form of a RECIST label classification. Acting as the baseline of our work, we train a voxel-based 3-D ResNet50 CNN to predict the treatment response label of the scan. We then evaluate the effect of the subsequent components in our integrated model cascade on the classification performance, as discussed in Section 3.

2.3. Stage I: Slice Classification

In order to limit the search space along the z -Dimension of a CT scan, we train a binary classifier which predicts for each slice whether it contains at least one pancreas or tumour pixel. Compared to a full segmentation model, this has the advantage of being very robust, and the model can be trained very efficiently and evaluated on large datasets while re-using the segmentation labels. To account for the three-dimensional structure of the task, we apply a Long Short-Term Memory (LSTM) cell after a 2-D encoder. The LSTM consecutively receives extracted image features and thus holds information about previous slices, turning the model into a 2.5-D classifier. As the 2-D encoder, we employ a ResNet50 pre-trained on ImageNet [17]. The single LSTM cell following the encoder has 1024 as the hidden size. Our 2.5-D approach is based on the notion that we can leverage prior knowledge of the classification target, i.e. that it is a single continuous object within the scan. Moreover, we post-process the predictions and automatically include slices between positively classified slices, as we know that these contain the organ of interest.

2.4. Stage II: Segmentation

As the second pipeline component, we train a 3-D segmentation network. For this, we use a voxel-based architecture, *DynU-Net* from MonAI [18] model zoo. The segmentation step is more expensive in terms of training and inference time, as well as more prone to inaccuracies compared to our slice classification model. However, it contributes considerable utility to the overall pipeline. It allows us to reduce the search space not only in the z -direction as above, but also in the x and y directions by converting the segmentation to a bounding box around the organ. Moreover, we later feed the predicted segmentations as additional image channels to the next stage and thus introduce additional (assistive) information. This approach is motivated by the human geometric pattern-matching approach, where the detection of specific shapes facilitates overall recognition. Lastly, we can reuse the network weights of the encoder for the next stage of the pipeline and profit from task-specific pre-trained weights as a *transfer learning* scheme. The re-utilisation of weights allows us to reduce the size of the RECIST classifier network and thus the computational power and training and inference time decrease.

2.5. Representation learning

As the final component of our pipeline, we evaluate the benefit of a triplet loss [19], such that the intermediate representations of samples of the same class are close, whereas the representations of different classes are further apart in the latent space. For this, we always sample an *anchor* image, a *positive* image from the same class and a *negative* image from the negative class. We calculate the loss in terms of

L_2 -distance of the intermediate representations after the feature extractor stage of the model (i.e. immediately before the linear layers). We use this loss term to condition the network on this task in a first training stage, and only in a second stage trained the classification part using a binary cross-entropy loss.

3. EXPERIMENTS

3.1. Slice Classification and Segmentation

Both models are trained separately on the Medical Segmentation Decathlon [16] dataset. We set a maximum of 200 epochs but monitor appropriate metrics and stop the training at convergence. On the held-out validation set, our slice classification network reaches **90.7% accuracy** on predicting the presence of pancreatic tissue in a slice. The DynU-Net segmentation network succeeds to segment the pancreas with a **Dice score of 77%** and the tumour with **40%**.

3.2. Classification Results

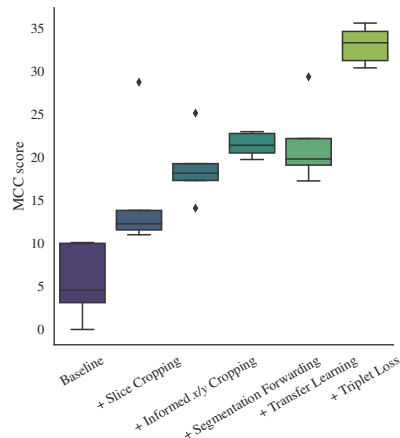


Fig. 2. Box-and-whisker plots of the MCC scores for each integrated pipeline component averaged over $N = 5$ runs.

In the following, we report the classification performance of the individual pipeline components as well as the overall pipeline in our experimental evaluation. We use the Matthew’s correlation coefficient (MCC) [20] as the main metric of our evaluation due to its clear advantages over other metrics [21], most prominently, its robustness to class imbalance. For the sake of completeness, we also report the binary accuracy (unweighted) as well as the area under the receiver-operator characteristic curve (AUC-ROC), however we stress that these results should be interpreted with caution as they are susceptible to class imbalance. We perform each experiment with 5 random seeds and report the standard deviation of the runs. As a baseline, we train our classifier model without modifications on the input CT-scan. We use

Method	MCC		Accuracy		AUC-ROC	
	μ	σ	μ	σ	μ	σ
Baseline	5.6%	4.4%	41.0%	4.5%	45.5%	7.9%
+ slice (z) cropping	15.5%	7.5%	54.2%	9.8%	57.5%	4.4%
+ informed x/y cropping	18.8%	4.0%	61.1%	4.4%	59.2%	3.1%
+ segmentation forwarding	21.5%	1.4%	53.3%	4.9%	59.3%	3.5%
+ transfer learning	21.6%	4.7%	64.1%	3.1%	59.8%	3.6%
+ triplet loss	33.1%	2.2%	67.2%	3.6%	63.7%	2.6%

Table 1. Classification metrics on our validation data. Here, MCC stands for the *Matthew’s Correlation Coefficient*, accuracy is the binary classification accuracy, *AUC-ROC* is the area under the receiver-operator characteristic curve. μ denotes the mean result, σ denotes the standard deviation of the results for 5 runs with different random seeds.

an input shape of $256 \times 256 \times 256$ for this baseline. In the following configurations we can reduce this shape without losing the image resolution thanks to learned cropping. We compare this set of experiments to each incremental change introduced above. Cropping not only reduces the input to relevant features for the classification network, but also allows to lift the resolution limit. In the case where we only use a slice classifier to crop in the z -dimension we perform centre cropping in the x and y direction, whereby we exploit the fact that the pancreas is located centrally in the abdomen. Our findings are summarised in Table 1. The baseline model predicts almost random outputs, with an MCC score close to zero and a high standard deviation between the different runs. A simple pre-processing step to classify each slice whether it contains the pancreas and thus allowing a cropping operation of the scan along the z -direction already improves the results considerably (**slice cropping**). We suspect that an important reason to this is that by cropping, we can train on higher resolution input images without reaching hardware limits and reduce the scan to its relevant anatomical regions. However, the deviation between runs within this setting is still remarkably high. This variation is reduced by using a segmentation model to additionally crop in the x and y directions (**informed x/y -cropping**) instead of naïve centre cropping, and even further by adding the segmentations as input to the classifier (**segmentation forwarding**). Reusing the feature extractor of the segmentation *DynU-Net* model and **transfer learning** it as a classification network did not yield a notable improvement in terms of test metrics. However, this model is much smaller and thus less computational intensive, which allowed a $2.5\times$ average training speed-up as well as a substantially decreased energy consumption. Lastly, the introduction of an additional **triplet loss** term on the intermediate representations of our classification network further improves the metrics. We reach a final MCC of 33.1% on average with a standard deviation of 2.2%. MCC is a more pessimistic metric than accuracy, which is at 67.2% (36 – 41/57 patients over 5 runs), but very robust to class imbalances, which is important in our dataset. These results are visually summarised in Figure 2.

4. CONCLUSION

In this work we analyse the effects of iterative changes to the problem of treatment response classification of PDAC on baseline CT scans. We showed that exploiting available segmentation labels, as well as using representation learning can yield a large improvement in classification results. We hope that this work is a first step towards accurately predicting the chemotherapy response of PDAC patients and thus leading to an improved personalised treatment scheme.

5. ACKNOWLEDGMENTS

This work was supported by the European Research Council (Deep4MI - 884622), German Research Foundation, Priority Programme SPP2177, Radiomics: Next Generation of Biomedical Imaging, German Cancer Consortium Joint Funding UPGRADE Programme: Subtyping of Pancreatic Cancer based on radiographic and pathological Features. The funders played no role in the design or execution of the study, nor on the decision to prepare or submit the manuscript. We used Jakob Richter’s Punkreas as basis for parts of our code (<https://gitlab.com/sanddorn/punkreas>).

6. REFERENCES

- [1] “Leitlinienprogramm Onkologie (Deutsche Krebsgesellschaft, Deutsche Krebshilfe, AWMF): S3-Leitlinie Exokrines Pankreaskarzinom, Langversion 2.0, AWMF Registernummer: 032-010OL,” 2021.
- [2] Diego Ardila, Atilla P Kiraly, Sujeeth Bharadwaj, Bokyung Choi, Joshua J Reicher, Lily Peng, Daniel Tse, Mozziyar Etemadi, Wenxing Ye, Greg Corrado, et al., “End-to-end lung cancer screening with three-dimensional deep learning on low-dose chest computed tomography,” *Nature medicine*, vol. 25, no. 6, pp. 954–961, 2019.
- [3] Friederike Jungmann, Georgios A Kaissis, Sebastian Ziegelmayr, Felix Harder, Clara Schilling, Hsi-Yu Yen,

- Katja Steiger, Wilko Weichert, Rebekka Schirren, Is-han Ekin Demir, et al., “Prediction of tumor cellularity in resectable pdac from preoperative computed tomography imaging,” *Cancers*, vol. 13, no. 9, pp. 2069, 2021.
- [4] Georgios Kaissis, Sebastian Ziegelmayer, Fabian Lohöfer, Hana Algül, Matthias Eiber, Wilko Weichert, Roland Schmid, Helmut Friess, Ernst Rummeny, Donna Ankerst, et al., “A machine learning model for the prediction of survival and tumor subtype in pancreatic ductal adenocarcinoma from preoperative diffusion-weighted imaging,” *European radiology experimental*, vol. 3, no. 1, pp. 1–9, 2019.
- [5] Georgios Kaissis, Sebastian Ziegelmayer, Fabian Lohöfer, Katja Steiger, Hana Algül, Alexander Muckenhuber, Hsi-Yu Yen, Ernst Rummeny, Helmut Friess, Roland Schmid, et al., “A machine learning algorithm predicts molecular subtypes in pancreatic ductal adenocarcinoma with differential response to gemcitabine-based versus folfirinnox chemotherapy,” *PloS one*, vol. 14, no. 10, pp. e0218642, 2019.
- [6] Lawrence H Schwartz, Lesley Seymour, Saskia Litière, Robert Ford, Stephen Gwyther, Sumithra Mandrekar, Lalitha Shankar, Jan Bogaerts, Alice Chen, Janet Dancey, et al., “Recist 1.1—standardisation and disease-specific adaptations: Perspectives from the recist working group,” *European journal of cancer*, vol. 62, pp. 138–145, 2016.
- [7] Saskia Litiere, Sandra Collette, Elisabeth GE de Vries, Lesley Seymour, and Jan Bogaerts, “Recist—learning from the past to build the future,” *Nature reviews Clinical oncology*, vol. 14, no. 3, pp. 187–192, 2017.
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” *ICLR*, 2021.
- [9] Lin Lu, Laurent Dercle, Binsheng Zhao, and Lawrence H Schwartz, “Deep learning for the prediction of early on-treatment response in metastatic colorectal cancer from serial medical imaging,” *Nature communications*, vol. 12, no. 1, pp. 1–11, 2021.
- [10] Youbao Tang, Ke Yan, Jinzheng Cai, Lingyun Huang, Guotong Xie, Jing Xiao, Jingjing Lu, Gigin Lin, and Le Lu, “Lesion segmentation and recist diameter prediction via click-driven attention and dual-path connection,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 341–351.
- [11] Hilda Azimi, Jianxing Zhang, Pengcheng Xi, Hala Asad, Ashkan Ebadi, Stephane Tremblay, and Alexander Wong, “Improving classification model performance on chest x-rays through lung segmentation,” *arXiv preprint arXiv:2202.10971*, 2022.
- [12] Saba Saleem, Javeria Amin, Muhammad Sharif, Muhammad Almas Anjum, Muhammad Iqbal, and Shui-Hua Wang, “A deep network designed for segmentation and classification of leukemia using fusion of the transfer learning models,” *Complex & Intelligent Systems*, vol. 8, no. 4, pp. 3105–3120, 2022.
- [13] Yue Zhou, Houjin Chen, Yanfeng Li, Qin Liu, Xuanang Xu, Shu Wang, Pew-Thian Yap, and Dinggang Shen, “Multi-task learning for segmentation and classification of tumors in 3d automated breast ultrasound images,” *Medical Image Analysis*, vol. 70, pp. 101918, 2021.
- [14] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle, “Greedy layer-wise training of deep networks,” *Advances in neural information processing systems*, vol. 19, 2006.
- [15] Maximilian Kohlbrenner, Russell Hofmann, Sabbir Ahmed, and Youssef Kashef, “Pre-training cnns using convolutional autoencoders,” *TU Berlin, Berlin, Germany*, 2017.
- [16] Michela Antonelli, Annika Reinke, Spyridon Bakas, Keyvan Farahani, Annette Kopp-Schneider, Bennett A Landman, Geert Litjens, Bjoern Menze, Olaf Ronneberger, Ronald M Summers, et al., “The medical segmentation decathlon,” *Nature communications*, vol. 13, no. 1, pp. 1–13, 2022.
- [17] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [18] MONAI Consortium, “MONAI: Medical Open Network for AI,” 2022.
- [19] Florian Schroff, Dmitry Kalenichenko, and James Philbin, “Facenet: A unified embedding for face recognition and clustering,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [20] Brian W Matthews, “Comparison of the predicted and observed secondary structure of t4 phage lysozyme,” *Biochimica et Biophysica Acta (BBA)-Protein Structure*, vol. 405, no. 2, pp. 442–451, 1975.
- [21] Davide Chicco and Giuseppe Jurman, “The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation,” *BMC genomics*, vol. 21, no. 1, pp. 1–13, 2020.