

Finding Reproducible and Prognostic Radiomic Features in Variable Slice Thickness Contrast Enhanced CT of Colorectal Liver Metastases

Jacob J. Peoples <https://orcid.org/0000-0003-0191-7446>
School of Computing, Queen's University, Kingston, ON, Canada

jacob.peoples@queensu.ca

Mohammad Hamghalam <https://orcid.org/0000-0003-2543-0712>
School of Computing, Queen's University, Kingston, ON, Canada
Department of Electrical Engineering, Qazvin Branch, Islamic Azad University, Qazvin, Iran

Imani James
Department of Radiology, Memorial Sloan Kettering Cancer Center, New York, NY, USA

Maida Wasim
Department of Radiology, Memorial Sloan Kettering Cancer Center, New York, NY, USA

Natalie Gangai <https://orcid.org/0000-0003-0603-7165>
Department of Radiology, Memorial Sloan Kettering Cancer Center, New York, NY, USA

Hyunseon Christine Kang <https://orcid.org/0000-0001-5313-3526>
Department of Abdominal Imaging, The University of Texas MD Anderson Cancer Center, Houston, TX, USA

X. John Rong <https://orcid.org/0000-0002-1169-5042>
Department of Imaging Physics, The University of Texas MD Anderson Cancer Center, Houston, TX, USA

Yun Shin Chun <https://orcid.org/0000-0003-1384-8927>
Department of Surgical Oncology, The University of Texas MD Anderson Cancer Center, Houston, TX, USA

Richard K. G. Do <https://orcid.org/0000-0002-6554-0310>
Department of Radiology, Memorial Sloan Kettering Cancer Center, New York, NY, USA

Amber L. Simpson <https://orcid.org/0000-0002-4387-8417>
School of Computing, Queen's University, Kingston, ON, Canada
Department of Biomedical and Molecular Sciences, Queen's University, Kingston, ON, Canada

Abstract

Establishing the reproducibility of radiomic signatures is a critical step in the path to clinical adoption of quantitative imaging biomarkers; however, radiomic signatures must also be meaningfully related to an outcome of clinical importance to be of value for personalized medicine. In this study, we analyze both the reproducibility and prognostic value of radiomic features extracted from the liver parenchyma and largest liver metastases in contrast enhanced CT scans of patients with colorectal liver metastases (CRLM). A prospective cohort of 81 patients from two major US cancer centers was used to establish the reproducibility of radiomic features extracted from images reconstructed with different slice thicknesses. A publicly available, single-center cohort of 197 preoperative scans from patients who underwent hepatic resection for treatment of CRLM was used to evaluate the prognostic value of features and models to predict overall survival. A standard set of 93 features was extracted from all images, with a set of eight different extractor settings. The feature extraction settings producing the most reproducible, as well as the most prognostically discriminative feature values were highly dependent on both the region of interest and the specific feature in question. While the best overall predictive model was produced using features extracted with a particular setting, without accounting for reproducibility,

(C-index = 0.630 (0.603–0.649)) an equivalent-performing model (C-index = 0.629 (0.605–0.645)) was produced by pooling features from all extraction settings, and thresholding features with low reproducibility ($\text{CCC} \geq 0.85$), prior to feature selection. Our findings support a data-driven approach to feature extraction and selection, preferring the inclusion of many features, and narrowing feature selection based on reproducibility when relevant data is available.

Keywords: Radiomics, Texture Analysis, Reproducibility, Colorectal Liver Metastases, Quantitative Imaging Biomarkers, Computed Tomography, Prospective Studies, Reproducible Features

1. Introduction

Radiomic analysis as a field is predicated on the idea that radiological imaging contains meaningful biological information contained in the patterns of intensity values within regions of interest, which could contribute to a better understanding of patient health (Gillies et al., 2016). The typical approach in radiomic studies is to extract a large number of pre-defined quantitative imaging features from a region of interest, and then use machine learning methods to reduce the dimensionality of the feature set and build models of a biological correlate or a medical outcome of interest (Horvat et al., 2022). The promise of radiomics to develop quantitative imaging biomarkers is of broad interest because it poses a non-invasive means to characterize patient disease using routine clinical imaging. However, to be clinically deployed, radiomic models must be widely validated, and their robustness to variable imaging settings well-established.

In the present study, we are focused on contrast-enhanced abdominal computed tomography (CT) of the liver for patients with colorectal liver metastases (CRLM). These patients have an overall poor prognosis, which could potentially be improved by prognostic radiomic signatures that could better target patients for surgery or chemotherapy. In patients with liver metastases, radiomic models derived from contrast-enhanced CT have shown substantial prognostic capability in both survival modeling and prediction of chemotherapy response (Fiz et al., 2020). Furthermore, although most studies focus on radiomic features extracted from the metastases themselves, studies of CRLM have shown that features from the liver parenchyma also contain important information when predicting hepatic disease-free survival or overall survival after hepatic resection (Simpson et al., 2017), or progression-free survival after radiotherapy (Hu et al., 2022b). With these applications in mind, in this study we are seeking to gain a better understanding of the robustness of radiomic models in contrast-enhanced CT of CRLM.

One way to study the robustness of radiomic models is to study the reproducibility of the model inputs—that is, to understand how consistent the radiomic features are under real-world variations that occur in image acquisition and reconstruction. Treating the radiomic features as measurements drawn from radiological images, the reproducibility of the features can be studied from a metrological perspective (Raunig et al., 2014), using statistical measures such as the concordance correlation coefficient (CCC). Understanding the reproducibility of radiomic features, however, is difficult given the many factors affecting the results. Restricting the discussion to CT, scans are acquired from different medical centers, using different scanners, with different imaging acquisition parameters and protocols, after which the images are reconstructed with different algorithms, different res-

olution parameters, and different kernels, all of which have been shown to affect feature reproducibility (Zhao, 2021). Once the images are in hand, feature extraction itself is not without concerns regarding reproducibility; while standardization efforts such as the Image Biomarker Standardisation Initiative (IBSI) (Zwanenburg et al., 2020) have been instrumental in creating a common nomenclature and a standard set of well-defined features, they do not provide a final answer in how to set the many configurable parameters which can affect the final feature values computed by common software such as `pyradiomics` (van Griethuysen et al., 2017).

In this study, we are considering feature reproducibility particularly when the slice thickness used to reconstruct the CT scans is varied. When comparing features extracted from CT images with thinner or thicker slices, studies have found that features from thinner images are more reproducible across variations in segmentation (Hu et al., 2022a), or across repeat imaging (Zhao et al., 2016), and may also produce more accurate models (He et al., 2016; Li et al., 2018; Xu et al., 2022). In direct comparisons, many features had poor reproducibility when comparing features from images with different slice thicknesses in a variety of phantom studies (Zhao et al., 2014; Ger et al., 2018; Berenguer et al., 2018; Kim et al., 2019; Varghese et al., 2019; Ligerio et al., 2021; Ibrahim et al., 2022). Prospective studies producing multiple reconstructions for each patient have reproduced this result on patient images for lung cancer (Lu et al., 2016; Park et al., 2019; Erdal et al., 2020; Yang et al., 2021; Emaminejad et al., 2021), and liver metastases (Meyer et al., 2019). Poor feature reproducibility with respect to slice thickness is concerning because many real-world retrospective or multi-site data sets include images with a range of different slice thicknesses due to variations in local protocols (Ger et al., 2018). Although image interpolation to a common, isotropic voxel size is considered a best practice for preprocessing during feature extraction (Zwanenburg et al., 2020) in order to ensure the image features are comparable between images with different voxel sizes, the optimal choice of resolution and resampling algorithm is undecided. Furthermore, resampling to a common voxel size appears, on its own, to be insufficient to overcome the inconsistency of feature values due to slice thickness variation, except in a small subset of features (Shafiq-Ul-Hassan et al., 2017; Shafiq-ul-Hassan et al., 2018).

To further complicate matters, it has been shown that reproducibility is not necessarily consistent across cancer types, even for a single modality such as CT (van Timmeren et al., 2016). However, in a systematic review of radiomics reproducibility studies, Traverso et al. (2018) found that the literature was limited to a small number of cancer types, with the greatest number of studies addressing lung cancers, amongst which CT was the most common imaging modality. Across these studies there was not a clear consensus on the most reproducible features, although in CT they found agreement that first-order features tended to be more reproducible than higher-order texture features. Ultimately, it seems that the reproducibility of features is not easily generalized across different anatomies or cancers. Therefore, reproducibility studies for the region of interest (ROI) and disease under consideration is an important part of the validation any radiomics-based imaging biomarker.

In this study we present an analysis of the relationship between the reproducibility and prognostic value of radiomic features drawn from contrast enhanced CT of CRLM. We present a reproducibility analysis of radiomic features on a cohort of 81 prospectively enrolled patients from two major US cancer centers, who underwent contrast enhanced

abdominal imaging with a controlled and systematically varied protocol. Our analysis is primarily focused on the effects of slice thickness chosen at reconstruction time. Features were extracted from the largest liver metastasis and the liver parenchyma from each subject using a variety of different configurations, varying the level of resampling, and the method of aggregation used in computing the higher-order texture features. To investigate the relationship between reproducibility and prognostic value, we conducted an in depth univariate and multivariable survival modeling analysis on an independent, publicly available data set of 197 preoperative contrast enhanced CT scans of patients who underwent hepatic resection to treat CRLM (Simpson et al., 2023, 2024).

This paper is a significant expansion of our previous work (Peoples et al., 2023), with a greater focus on the relationship between the reproducibility and prognostic value of the radiomic features under consideration. In this paper, we are less concerned with finding the feature extraction settings that produce the most reproducible features. Instead, we focus on the integration of reproducibility information into the development of prognostic radiomic signatures. We present a joint, univariate analysis of both the reproducibility and prognostic discriminative ability of the features, taking a multi-objective optimization point of view. The multivariable analysis was revised to use a standard feature selection algorithm, and to conduct many iterations of the cross-validation to ensure our results were stable. Additional details are provided throughout the paper, giving greater context on the methodology, more visualisation and discussion of the results, and more interpretation of how our results fit into the greater context of the literature on the reproducibility of radiomics. Code for all analysis is available at github.com/jpeoples/melba2024.

2. Methods

2.1 CT Imaging and Segmentation

Contrast enhanced, portal venous phase CT scans were prospectively collected from a total of 81 patients with CRLM from two institutions, Memorial Sloan Kettering Cancer Center (New York, NY) (MSK) (n=44) and MD Anderson Cancer Center (Houston, TX) (MDA) (n=37), with institutional review board approval and informed consent.¹ Every scan was collected on a multi-detector CT scanner (Discovery CT750 HD; GE Healthcare, Madison, WI, USA) with 64 detector rows, and 0.625 mm detector width, for a total collimation width of 40 mm. The images were collected with a tube voltage of 120 kVp, and automated tube current modulation using GE Smart mA with a noise index of 14 (MSK) or 11 (MDA). The tube current range varied between centers, with MSK using range 220-380 mA, while MDA used range 275-650 mA. The gantry rotation time was 0.7 s with pitch factor 0.984 for MSK, and 0.5 s with pitch factor 0.516 for MDA. All images were reconstructed with the standard soft tissue convolution filter.

1. Note that all scans used to generate the results in this paper are consistent with the standard clinical acquisitions at each site, that would have been acquired for these patients regardless of their involvement in the study. Informed consent was required because the full study protocol includes a second scan, taken just before or after the usual scan (within ± 15 s), to capture the test-retest repeatability of radiomic features in the presence of contrast. These results are not included in the present paper because the add-on scans were not processed at the time of preparation.

To study the reproducibility of radiomic features with respect to image reconstruction, each patient scan was retrospectively reconstructed with different slice thicknesses and levels of adaptive statistical iterative reconstruction (ASiR) after image acquisition. In particular, every combination of three different slice thicknesses—2.5 mm, 3.75 mm, 5 mm—and seven different levels of ASiR—from 0% (equivalent to filtered back-projection), to 60% in increments of 10%—were generated, giving a total of $3 \times 7 = 21$ reconstructions for every scan. In all cases, the slice thickness and slice interval were equal, such that the reconstructed slices were spatially contiguous, and non-overlapping.² Images were stored and transferred in Digital Imaging and Communications in Medicine (DICOM) format after deidentification from both MSK and MDA. All CT data was converted into Neuroimaging Informatics Technology Initiative (NIfTI) format for further processing.

2.2 Segmentation

A single reconstruction (slice thickness 5 mm and ASiR 20%) was chosen as the reference reconstruction for each patient, for manual segmentation verification and correction by an experienced radiologist (R. D.). The choice of 5 mm and 20% ASiR for the reference was made to more closely match the standard-of-care clinical imaging protocol at the radiologist’s home institution (MSK). The segmentation of these reference scans was completed in two phases: first, an automated segmentation was generated, after which the radiologist verified and corrected each mask in 3D Slicer (Kikinis et al., 2013). The initial segmentations were generated using an nnU-net model (Isensee et al., 2021) trained on a public database of 197 CT scans from patients with CRLM (Simpson et al., 2023, 2024), available from the Cancer Imaging Archive (TCIA) (Clark et al., 2013). Details on the development of this model can be found in several previous publications (Hamghalam et al., 2021; Mojtahedi et al., 2022; Hamghalam et al., 2023). After the radiologist corrected the segmentations for the reference reconstruction, segmentations for other slice thicknesses were generated by resampling the reference segmentation using nearest-neighbor interpolation with the Simple ITK software library. No changes were made to the segmentations for reconstructions with different ASiR, given that the tissue being imaged does not change in a given slice between ASiR levels.

2.3 Radiomic Feature Extraction

Features were extracted from every image using `pyradiomics` (van Griethuysen et al., 2017), which is one of several open-source packages implementing a large set of IBSI-compliant features. Two 3D ROIs were used: the largest tumor within the liver, and the liver parenchyma (with all tumors and vessels excluded). The `pyradiomics` library has seven default classes of features: shape, first order, gray level co-occurrence matrix (GLCM), gray level dependence matrix (GLDM), gray level size zone matrix (GLSZM), gray level run length matrix (GLRLM), neighboring gray-tone difference matrix (NGTDM). In this study we used

2. Because the CT scanners in this study (GE Discovery CT750) have an equal number of detectors and channels, the thinnest possible slices (corresponding to a single detector row width of 0.625 mm) are always acquired, and the slice thickness is applied only at reconstruction time, grouping data from multiple detector rows to produce thicker slices. In our case, 2.5 mm, 3.75 mm, and 5 mm slice thicknesses correspond to grouping 4, 6, or 8 individual 0.625 mm rows.

Table 1: Feature counts by class.

Feature Class	Count
First order	18
GLCM	24
NGTDM	5
GLDM	14
GLRLM	16
GLSZM	16
Total	93

all feature classes, with the exception of shape features, which were excluded because there was only one reference segmentation per patient, and therefore any variations across reconstructions would only reflect the effects of interpolation. From the remaining six classes, all default features were included. Although `pyradiomics` supports the extraction of features from a variety of derived images, in addition to the original intensity image, we did not include any analysis of these derived image features. The total number of features used in this study is broken down by feature class in Table 1, and a complete list is provided in Table 6 in the appendix.

2.3.1 TERMINOLOGY

Before continuing, we will establish some key terminology used throughout the remainder of this paper. We will frequently refer to two classes of radiomic features: first-order features, and texture features. For our purposes, first-order features are those that are computed strictly from the intensity histogram for all voxels in the ROI, and corresponds to the first order feature class in `pyradiomics`. Texture features, on the other hand, will be used strictly to refer to the remaining classes of higher-order features—GLCM, GLDM, GLSZM, GLRLM, and NGTDM—which share the common factor of accounting for relationships between the intensities of neighboring voxels. Note that this usage mirrors the terminology used by IBSI (Zwanenburg et al., 2020).

2.3.2 PREPROCESSING

A typical radiomic feature extraction pipeline includes a number of image preprocessing steps prior to feature computation. In our feature extraction, there were three preprocessing steps that every image underwent: image resampling, mask resegmentation, and intensity discretization.

Image resampling refers to the process of altering the input image resolution using an interpolation process. Given that texture features take account of relationships between neighboring voxels, resampling is important, because otherwise the neighbors being compared would not be an equal physical distance apart across images with different resolution, changing the meaning of the feature. To resolve this issue, the IBSI reference manual (Zwanenburg et al., 2016) recommends resampling images to a common resolution. Furthermore, the IBSI manual recommends a common *isotropic* resampling, in order to ensure the rota-

tional invariance of 3D texture features, which account for relationships between neighboring voxels in all 3D directions. Because radiomic features depend on both the underlying CT image, and a segmentation of the ROI, when applying resampling, both the image and segmentation mask must be interpolated.

Mask resegmentation refers to the removal of voxels outside of a preconfigured range when computing the first-order and texture features. Although Zwanenburg et al. (2016) do not give specific recommendations about the optimal settings for resegmentation, it is commonly applied to ensure that outlier intensity values (due to small errors in the segmentation mask, or due to artifacts) do not skew the resulting feature distributions.

Intensity discretization (or quantization) is a binning process used to reduce the number of unique intensity values in the image, which is used in the computation of texture features, as well as some first-order features which require a probability density based on the intensity histogram. Discretization is known to have a substantial effect on feature values (Shafiq-Ul-Hassan et al., 2017), although the effect on feature reproducibility may be limited (Larue et al., 2017).

In this study, all features were computed using a discretization level of 24 bins. The masks were resegmented using a window of $[-50, 350]$ Hounsfield units (HU), in order to exclude metal artifacts (stents, etc), as well as rare and implausible outlier intensities. All CT images were interpolated using B-splines, while the segmentation masks were interpolated using nearest-neighbor interpolation, which are the default algorithms in `pyradiomics`. We tested several resampling resolutions, which are described in greater detail below.

2.3.3 TEXTURE FEATURE AGGREGATION

The IBSI reference manual (Zwanenburg et al., 2016) breaks feature aggregation into three categories—2D, 2.5D, and 3D. The exact details of how the aggregation works varies across the classes of texture features, but the key factors are as follows.

3D vs 2D or 2.5D: All texture features consider the relationships between neighboring voxels. The first key factor in the aggregation methods is that while 3D aggregation includes neighbors from any direction in three dimensions, 2D and 2.5D aggregation only considers neighboring voxels within the same axial plane.

2.5D vs 2D: The second key factor differentiates 2D and 2.5D methods. Each class of texture feature is defined by an underlying matrix representing some aspect of the relationships between neighboring voxels. Individual features are then defined in terms of equations operating on this matrix. For 2D aggregation, the matrix is computed independently for each slice, and the resulting features are merged across slices by averaging. In the case of 2.5D aggregation, a single matrix is computed, which includes relationships from all slices, which is then used to compute the features in the usual way. This differs from 3D aggregation in that for 2.5D aggregation, only neighbors which share the same axial plane are considered, while in 3D aggregation, all neighbours in all three dimensions are included.

Directional vs non-directional: Another factor is whether the matrix underlying the feature class is defined directionally or not. The GLCM and GLRLM are both defined per direction. Taking symmetry into account, this means that there are four or thirteen unique GLCM/GLRLM in 2D or 3D respectively. Therefore, the matrices for each direction can

Table 2: Aggregation methods by feature class.

Type	Aggregation		Feature Classes
	2.5D	3D	
Directional	JJUI	ITBB	GLCM, GLRLM
Non-directional	62GR	KOBO	NGTDM, GLDM, GLSZM

be used to derive directional features, which can then be averaged, or the matrices can be merged via addition and used to compute single features. In the present work, we rely on the former approach for both 3D and 2.5D aggregated features, which corresponds to the IBSI-defined classes ITBB and JJUI, respectively.

The matrices underlying the remaining classes—GLDM, NGTDM, GLSZM—consider neighbors in all included directions in a single matrix. Therefore, no directional merging takes place. For these, 2D aggregation would compute a matrix per slice, which would then be used to derive per-slice features, which could be averaged; 2.5D aggregation would merge the 2D matrices across all slices, before computing features; 3D aggregation would compute a single matrix for the entire ROI, because neighbors in all directions are counted. In the present work we are using both 2.5D and 3D aggregation, which correspond to the IBSI-defined classes 62GR and KOBO, respectively.

A breakdown of the texture feature classes and the aggregation methods we are using in the present study is given in Table 2.

2.3.4 RESAMPLING AND AGGREGATION VARIANTS

Because our images are substantially anisotropic, with a larger axial slice thickness than in-plane pixel size, we wanted to investigate the optimality of resampling to an isotropic voxel size per the IBSI recommendations (Zwanenburg et al., 2016). Accordingly, we investigated three different levels of resampling: $1 \times 1 \times 1$ mm, $0.85 \times 0.85 \times 0.85$ mm, and $0.85 \times 0.85 \times 2.5$ mm, where 1 mm was chosen as a typical value from the radiomics literature, 0.85 mm was chosen as the median in-plane pixel spacing in the data set, and 2.5 mm was chosen as the 10th-percentile z-axis spacing, in analogy to the nnU-net resampling method for anisotropic imaging data sets (Isensee et al., 2021). The original distribution of in-plane pixel sizes in the data set prior to resampling is shown in Figure 1.

We also chose to extract features with both 3D and 2.5D feature aggregation at all resampling levels, to investigate whether 2.5D aggregation might be more suitable in anisotropic imaging. Because 2.5D aggregation never includes comparisons of voxels in neighboring axial planes, maintaining a common z-axis resolution is also not necessary for features to remain comparable. Therefore, in addition to the three previously described resampling levels, we also tried resampling the images to 0.85 mm or 1 mm in-plane resolutions, while preserving the original z-axis resolution, and extracting only 2.5D aggregated features.

All the aforementioned configurations led to a total of eight distinct feature extraction settings, which are summarized in Table 3. The table also introduces a naming scheme for these feature extraction settings, in order to enable the comparison of results across methods. In the naming scheme, the voxel resolution is indicated by the letters “L” (large, i.e. 1 mm),

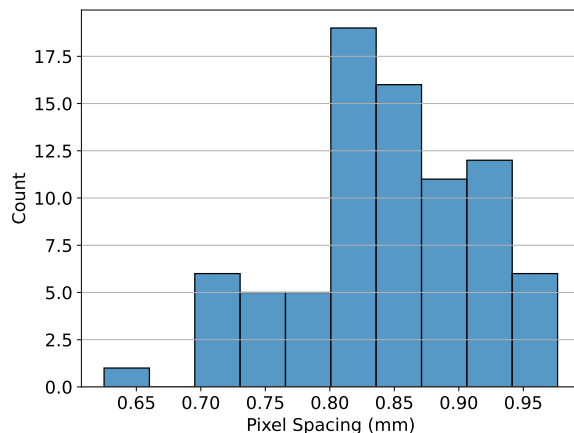


Figure 1: A histogram of the in-plane pixel spacing for images in our dataset. Pixel spacing was consistent across all reconstructions, so only counts for the reference reconstructions (5 mm slice thickness and 20% ASiR) are shown.

Table 3: Feature extraction methods.

Name	Resampling (mm)		Aggregation
	In-plane	z-axis	
L2i	1	None	2.5D
L2	1	1	2.5D
L3	1	1	3D
S2i	0.85	None	2.5D
S2	0.85	0.85	2.5D
S3	0.85	0.85	3D
A2	0.85	2.5	2.5D
A3	0.85	2.5	3D

“S” (small, i.e., 0.85 mm), and “A” (anisotropic, i.e. $0.85 \times 0.85 \times 2.5$ mm). The feature aggregation is indicated by a number “2” or “3”, for 2.5D or 3D aggregation, respectively. Finally, a lower-case “i” is appended to the name to indicate that the resampling was restricted to in-plane only, preserving the z-axis resolution.

2.4 Reproducibility Analysis

The CCC (Lin, 1989) was used to measure reproducibility of radiomic features. The first phase of the analysis was restricted to the reference ASiR level of 20%. The standard pairwise CCC was computed for every feature across all pairs of slice thicknesses (2.5 mm vs. 3.75 mm, 2.5 mm vs. 5 mm, 3.75 mm vs 5 mm). The paired Wilcoxon sign-rank test (Wilcoxon, 1945) was used to test the statistical significance of the change in CCC between slice thicknesses.

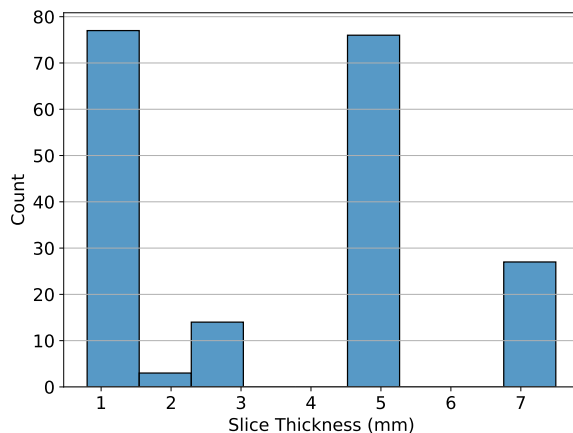


Figure 2: A histogram of the slice thickness for images in the survival data set. Slice thicknesses fell in the range $[0.8, 7.5]$ mm.

The second phase of the analysis used a linear mixed model (LMM) for each feature in order to compute a generalized CCC using the data from all three slice thicknesses (Carrasco and Jover, 2003). In this model, the reconstructions with different ASiR levels were also included, and controlled for as a fixed-effect when computing the CCC. In brief, using this approach, an LMM is computed for each feature, which is used to estimate the variance due to subject, σ_s^2 , the variance due to both slice thickness, σ_t^2 and ASiR, σ_a^2 , along with an error term, σ_e^2 . Following Carrasco and Jover (2003), the generalized CCC is then

$$\text{CCC} = \frac{\sigma_s^2}{\sigma_s^2 + \sigma_t^2 + \sigma_e^2}. \quad (1)$$

For each feature, we used this method to compute a spectrum of CCCs, across all feature extraction settings, and ROIs.

2.5 Survival Analysis on Independent Data Set

The reproducibility of a feature is an independent consideration from its value as a predictor in a given context. Ultimately, in radiomics the goal is to model a given outcome, and therefore, the predictive or prognostic values of features can not be sacrificed in order to create reproducible radiomic signature. To address this concern we performed survival analysis on the aforementioned public data set of 197 CRLM patients (Simpson et al., 2023, 2024). In this data set, the reproducibility of features across slice thickness is an important consideration, because thickness varied widely across scans (range $[0.8, 7.5]$ mm, see Figure 2). All scans in this data set were acquired prior to a hepatic resection to treat CRLM, and the repository includes right-censored data on overall survival time post-operation. For each of these pre-treatment scans, we extracted features using all eight of the different feature extractor settings previously described. In each case, the settings used to configure pyradiomics were identical to those used on the reproducibility dataset.

We completed the survival analysis in two phases, the first looking at univariate relationships of the radiomic features with overall survival, and the second considering multivariable models. In the first phase we computed the discriminative ability of each individual feature with respect to overall survival, as measured by Harrell's C-index (Harrell et al., 1996). Features with negative discriminative ability (C-index < 0.5) were negated, ensuring all C-index values fell in the range $[0.5, 1]$.

In the second phase of the analysis, we did a repeated 10-fold cross-validation of a multivariable Cox proportional hazards (CPH) model building procedure. Features from both ROIs were combined into one larger set for feature selection. We first restricted the candidate features to those with $CCC \geq CCC_t$ in the reproducibility analysis, where CCC_t was a predetermined threshold. The model building process in each iteration then consisted of a feature selection, followed by CPH modeling. Features were selected by computing the univariate C-index for each feature, as described above, and removing any with C-index < 0.55 . After this, univariate CPH models were constructed for each remaining feature, and removed whenever the feature was not significant with $p < 0.1$.³ Finally, the remaining features were reduced down to a predetermined number of features using the minimum redundancy, maximum relevancy (mRMR) feature selection algorithm (Ding and Peng, 2005). The 10-fold cross validation was repeated 100 times in order to get a stable estimate of the performance of the resulting multivariable models. In this case, we did not compute a final model, because our goal was only to compare the performance across the models.

The entire multivariable modeling process was conducted for every feature extraction setting individually, as well as for the case of all features from all extraction settings combined into one large set. For each feature set, the process was repeated for each CCC threshold in the set $CCC_t \in \{0, 0.8, 0.85, 0.9, 0.95\}$. Finally, for each feature set, and each CCC threshold, the process was repeated for every feature count in the set $\{1, 2, 4, 8, 16, 32, 64\}$. This resulted in a total of 9 feature sets \times 5 CCC thresholds \times 7 feature counts = 315 different multivariable model cross-validation experiments.

2.6 Statistical Analysis

2.6.1 HIERARCHICAL CLUSTERING

To find patterns in the reproducibility and univariate predictive value of features, we used a hierarchical clustering and dendrogram visualization (Müllner, 2011; Bar-Joseph et al., 2001). For the reproducibility analysis, the rows of the matrix corresponded to the features, while the columns corresponded to each combination of ROI and feature extractor, and the values were the CCC. By including both ROIs in the analysis, we were able to see clusters of patterns across and between ROIs. For the univariate survival analysis, a similar clustering approach was used, although the two ROIs were split, because the same feature in either ROI may have a different biological relevance to survival. Therefore, in the univariate survival case, the rows corresponded to each feature, while the columns corresponded to the feature extractors, and the analysis was repeated for both ROIs. In all cases, the hierarchical clustering was done for both the rows and columns, using Ward linkage (Ward, 1963).

3. We used a weak threshold for significance, because this is simply a filtering step.

2.6.2 PARETO EFFICIENCY AND PARETO FRONT

In multi-objective optimization, a potential solution A *Pareto dominates* solution B if for every objective, A is better than B , or equally good. A solution A is *Pareto efficient* if there exists no other solution that Pareto dominates it. The *Pareto front* refers to the set of all Pareto efficient solutions. In other words, the Pareto front is the set of all solutions for which no improvement on any objective is possible, without a deterioration of some other objective. In order to better understand the relationship between the reproducibility and prognostic value of the features, we considered the set of features that are Pareto efficient over both the CCC and univariate C-index when features were grouped across all extractors. To understand the relationships between extractors, we also computed a Pareto front across feature extractors for each individual feature.

3. Results

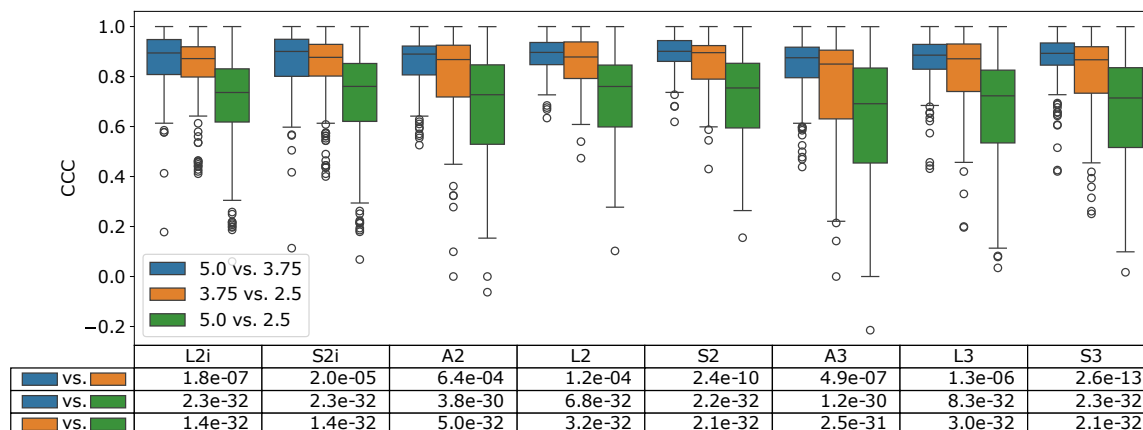


Figure 3: Box plots of the features compared pairwise across slice thicknesses, broken down by feature extraction setting. The statistical significance of the change in CCC value between different pairs of slice thicknesses is listed below the plot, based on the results of a Wilcoxon sign-rank test.

The results of the first phase of the reproducibility analysis are summarized in Figure 3. Here, the distribution of CCCs across all feature values for each extractor is compared across each pair of slice thicknesses. Intuitively, the largest difference in slice thickness (2.5 mm vs. 5 mm) shows a significantly lower distribution of CCCs than the other two pairs, which compare slice thicknesses which are closer together (2.5 mm vs. 3.75 mm, and 3.75 mm vs 5 mm). These differences were highly statistically significant based on a Wilcoxon sign-rank test between all CCC pairs ($p < 3.8 \times 10^{-30}$ in the bottom two rows of the table in Figure 3). Interestingly, the 3.75 mm vs. 5 mm pairs also seem to have slightly higher CCC than the 2.5 mm vs. 3.75 mm pairs, at a lower level of significance ($p < 6.4 \times 10^{-4}$ in the top row of the table in Figure 3).

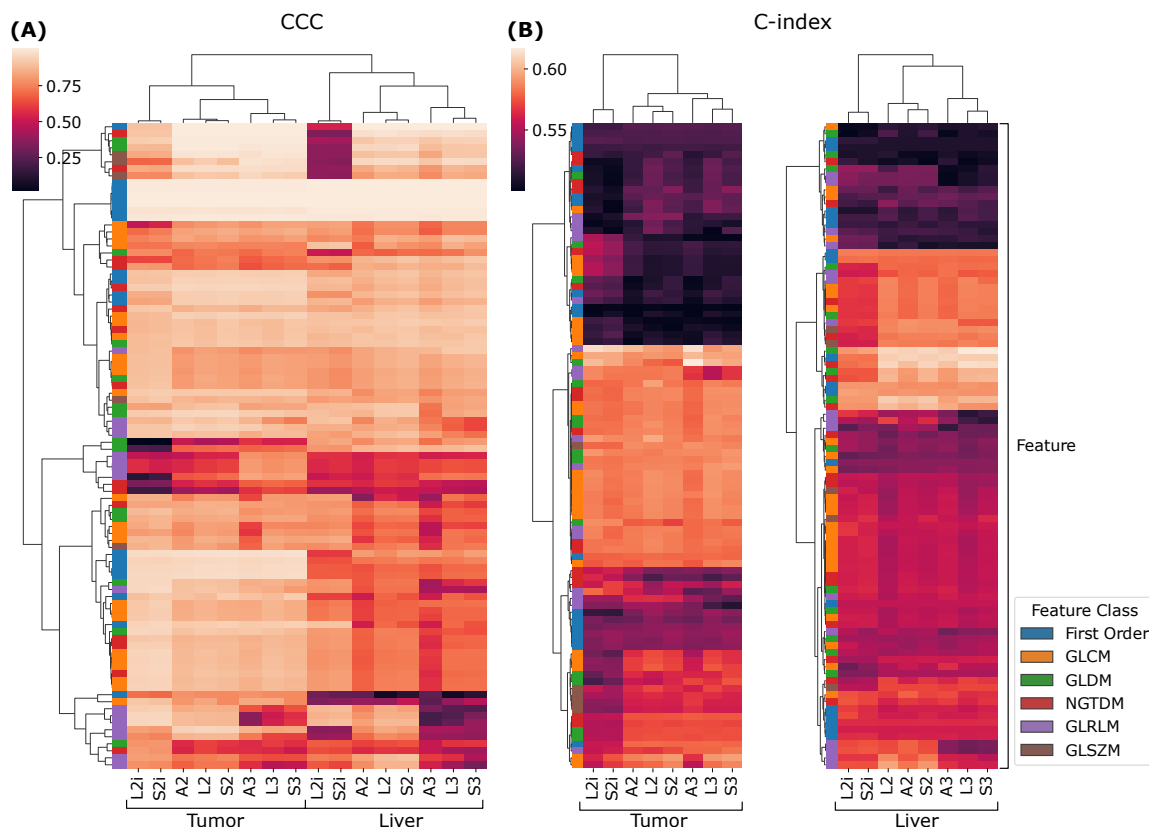


Figure 4: Cluster maps of the CCC (A) and C-index (B) of all features. Each row corresponds to a unique feature, while each column is an extractor setting. For the CCCs, the liver and tumor extractor results are joined, and clustered together, to emphasize the patterns of reproducibility across ROIs and extractor. For the C-index, the liver and tumor results are clustered separately. The feature class for each row is indicated by the left-most column of each heat map.

We now turn to the second phase of the reproducibility analysis, based on the generalized CCCs computed across all slice thicknesses using the LMM model. Figure 4 (A) summarizes the results across all features, extraction settings, and ROIs. The results are visualized as a heat map, where each row is a feature, and each column corresponds to a combination of ROI and extractor setting, and the values are CCCs. Both the rows, and columns underwent hierarchical clustering based on Euclidean distance and Ward linkage. We can see, near the top of the map, a cluster of features that are highly reproducible across all extraction settings. As can be seen in the left-most column of the map, these highly reproducible features all belong to the first-order feature class. With some exceptions, within both ROIs, the CCCs across extractors tend to follow similar patterns; however, overall, the liver features visually appear to tend toward lower CCCs. This tendency is confirmed by a Wilcoxon sign-rank test, which shows that feature CCCs between the two ROIs are statistically significantly lower in the liver parenchyma, with $p \in [8.8 \times 10^{-11}, 3.6 \times 10^{-5}]$

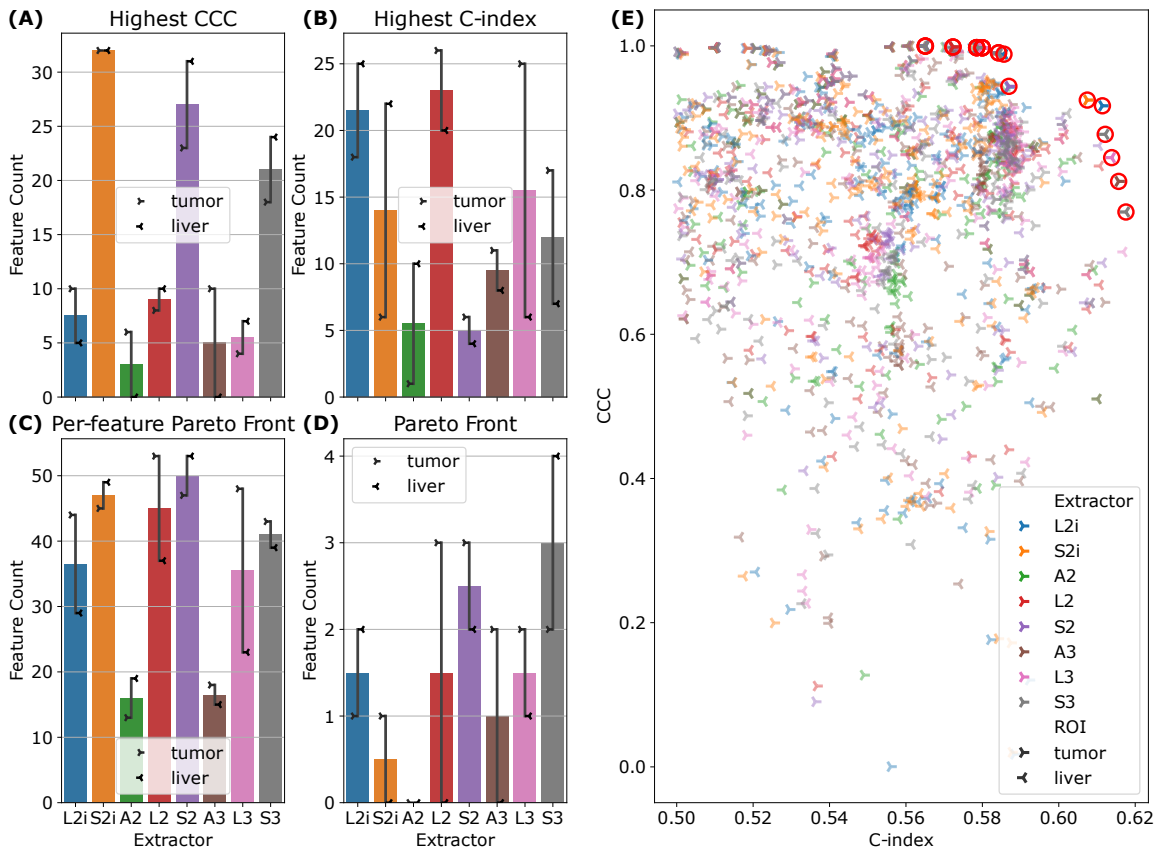


Figure 5: (A)–(C): Bar plots counting for how many features each extractor produces the highest CCC (A), highest C-index (B), or a pair of CCC and C-index that is Pareto efficient for that feature (C). (D): A bar plot of how many features from each extractor are on the Pareto front for *all* features across all extractors. In (A)-(D), the line indicators show the number of features broken down by ROI (tumor or liver for left and right, respectively), while the bar height corresponds to the average across the two ROIs. (E): A scatter plot of the C-index and CCC for all features, color coded by ROI, with points on the Pareto front rendered with full opacity and circled in red.

across all extractors. Interestingly, the CCCs between ROIs and extractor settings are sufficiently different that they cluster logically: the two ROIs form two higher level clusters; and within each ROI, L2i and S2i form their own cluster, and within the other cluster, the 2.5D and 3D settings cluster together.

Figure 4 (B) and (C) show similar cluster maps of the univariate C-indexes, computed against overall survival in the publicly available cohort, for each feature, from the tumor and liver ROIs, respectively. The ROIs were clustered separately because the ROIs are biologically different, and therefore the discriminative ability of a given feature may be completely different in a different ROI. Both maps cluster features into three broad groups—

Table 4: The performance of the top 10 parameter combinations in terms of test-set C-index and 95% confidence interval, in the multivariable CPH experiments. The average proportion of features drawn from the liver across all runs is listed in the final column.

Extractor	Feature Count	CCC _t	Harrel's C-index	Prop. Liver Features
L2i	4	0.00	0.630 (0.603–0.649)	0.52
all	4	0.85	0.629 (0.605–0.645)	0.53
L2i	8	0.00	0.627 (0.607–0.648)	0.57
all	8	0.80	0.626 (0.606–0.645)	0.55
L2i	4	0.90	0.625 (0.603–0.643)	0.28
all	4	0.80	0.624 (0.600–0.645)	0.53
L2i	8	0.90	0.624 (0.602–0.645)	0.24
all	8	0.85	0.622 (0.605–0.644)	0.56
all	8	0.00	0.622 (0.599–0.643)	0.50
S2i	8	0.00	0.620 (0.592–0.645)	0.61

a low C-index group at the top, a high C-index group in the middle, and a middle group at the bottom. Both the high C-index and low C-index groups are larger in the tumor ROI than in the liver parenchyma. Broadly speaking, the features tend to show similar trends in terms of predictive value across the extractor settings. On the other hand, there are some exceptions: at the bottom of the tumor cluster map (B), there is a cluster of features where the L2i and S2i features tend to be less predictive than the others. Similarly, in the high C-index cluster of the liver ROI, L2i and S2i again tend to be less predictive. The column clusters in both maps reproduce the finding for CCC in Figure 4 (A): three high-level clusters corresponding to (L2i and S2i), (A2, L2, and S2), and (A3, L3, and S3).

To examine the relationship between feature reproducibility and discriminative ability, we considered the relationships between the CCC and C-index for each feature. For every feature, we checked which extractors produced the highest CCC, the highest C-index, and which extractors were Pareto efficient for both CCC and C-index in that feature. The results are summarized in Figure 5 (A)–(C). Comparing Figure 5 (A) and (B), we see that the results are inconsistent—there is no set of extractors that produce both the highest CCC and the highest C-index. With the exceptions of A2 and A3, Figure 5 (C) shows that all feature extraction settings are well-represented on the Pareto fronts across features. Finally, grouping all features from all extractors together, and computing the set of Pareto efficient features for CCC and C-index shows that all extractors except A2 contribute (see Figure 5 (D)), with S3 contributing the greatest number of features (6). A scatter plot of all features is shown, color coded by extractor setting, with the Pareto front highlighted, in Figure 5 (E). Finally, the set of all 23 Pareto efficient features is listed in Table 5 in the appendix. It is worth noting that several features that are listed are actually equal, and therefore over-counted in Figure 5 (D). Furthermore, note that first-order features across S2 and S3, or L2 and L3 are also equal, and therefore some features that appear in both sets are still identical. These groups of equal features are marked in Table 5. Figure 5 (A)–(D)

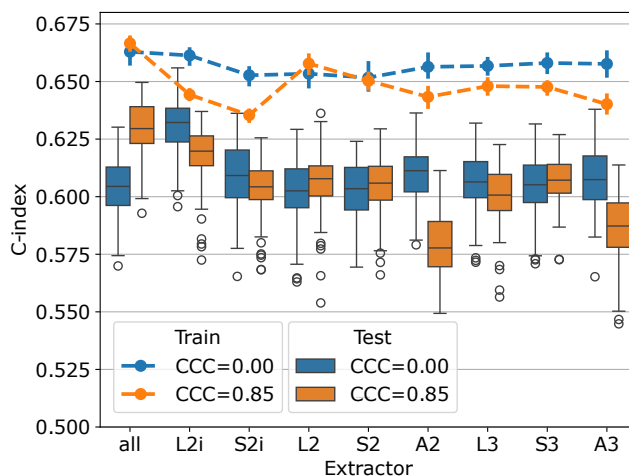


Figure 6: The cross-validation performance of the 4-feature models across extractor settings is plotted for $CCC_t = 0$ (blue) and $CCC_t = 0.85$ (orange). The test set results are summarized in box plots, while the average performance on the train set across folds is summarized by lines. The train set error bars correspond to 90% confidence intervals for the performance across repeated cross-validations.

also indicates that the optimal feature extraction method for a given criterion (C-index, CCC, or Pareto front) appears to vary not only across specific features, but across ROIs. Indeed, most of the extractors, and both ROIs are represented, at least to some extent, on the overall Pareto front across all features and extractors (see Figure 5 (D)).

The multivariable CPH cross-validation experiments produced results for 315 combinations of three parameters (feature extraction setting (8 extractors plus 1 for all), feature count (1, 2, 4, 8, 16, 32, 64), CCC threshold (0, 0.8, 0.85, 0.9, 0.95)). The test-set C-index and 95% confidence intervals across all 100 repetitions are listed for the ten highest performing parameter combinations in Table 4. Models based on the features from extractor L2i attained the highest C-index of 0.630 (0.603–0.649), with 4 features being selected, and no CCC thresholding. The second highest performance of 0.629 (0.605–0.645) was attained using features from all extractors, with 4 features being selected, and a threshold of $CCC \geq 0.85$. Based on these top two performing cases, the performance of the models using 4 features, across all extractors, and $CCC_t = 0$ and $CCC_t = 0.85$ are visualized in Figure 6. For most extractors, the thresholding of features to those with CCC over 0.85 appears to have little effect on the performance of the resulting models. In a few cases, L2i, A2, and A3, the performance is reduced after thresholding. When including features from all extractors (left-most in the figure), the performance increases after CCC thresholding. The performance on the training set tends to be lower after thresholding, suggesting that the models using reproducible features may be less prone to overfitting.

To visualize the sensitivity of these results with respect to the CCC threshold and the number of features selected, we have visualized the results from all parameter settings in Figure 7. Each plot shows the C-index for both the training and test sets, plotted against the

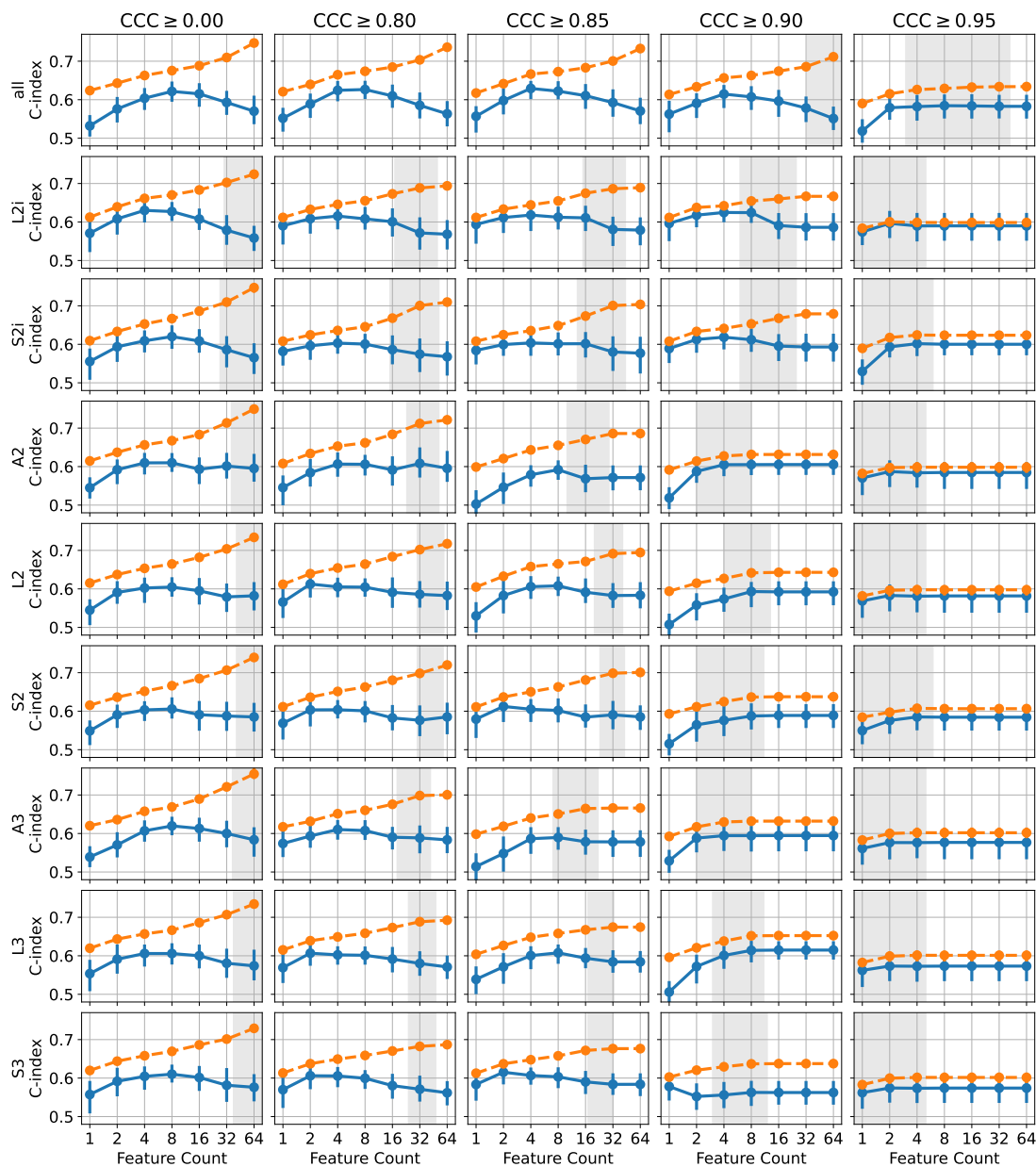


Figure 7: The C-index for the train set (blue) and test set (orange) is plotted against feature count, for every extractor setting (rows) and every CCC threshold (columns). The 90% confidence interval for the number of available features after CCC thresholding and univariate filtering is indicated by the gray regions.

number of features selected in the model. In some cases, after thresholding, and applying the univariate feature filter (removing features with C-index below 0.55 or with p -value over 0.1), there may have been fewer than the desired number of features remaining for

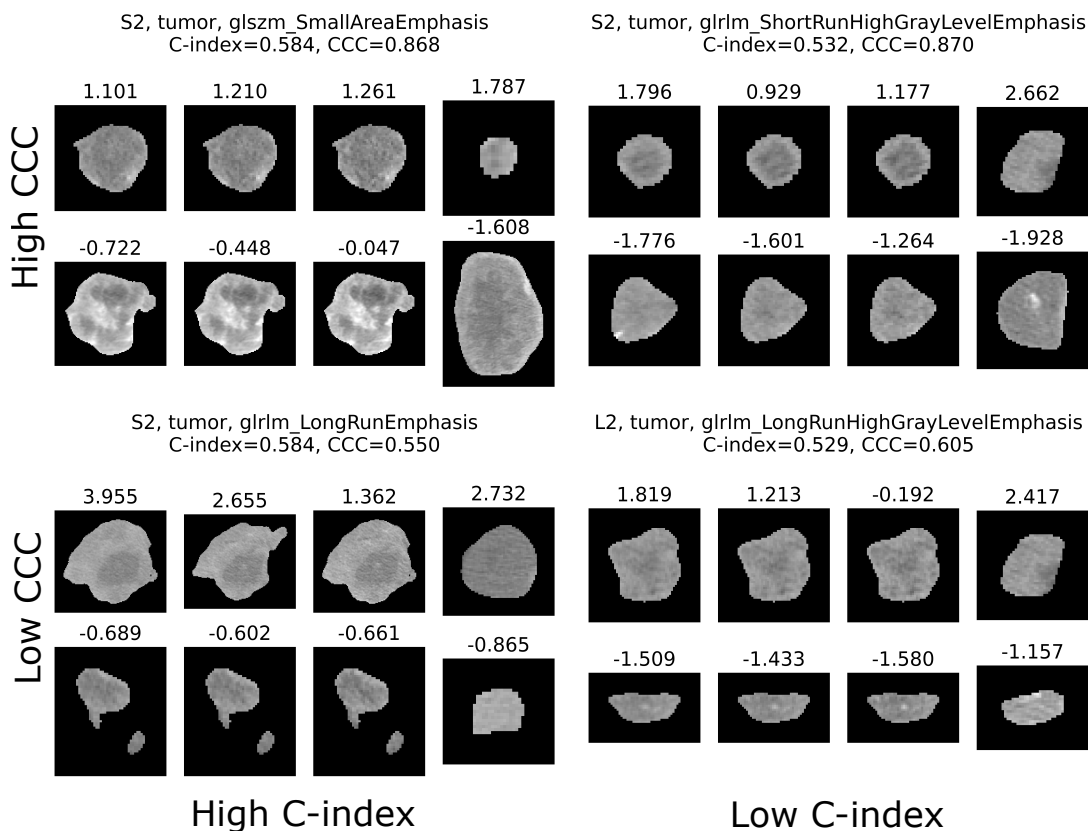


Figure 8: Features with high/low CCC (top/bottom) and high/low C-index (left/right) are visualized. For each, an example tumor with a high value for the selected feature, and a low value are illustrated in the two rows. The columns display the same tumor for an example patient from the reproducibility data set (first three columns, corresponding to slice thickness 5 mm, 3.75 mm and 2.5 mm), and the survival data set (fourth column).

feature selection. In these cases, all features passing the threshold and univariate filter were used. The 95% confidence intervals for the number of available features at the time of applying mRMR for the different threshold levels is indicated by the gray regions, where either limit was less than 64. Typically, beyond this point, requesting more features does not change the performance of either the train or test set, because the number of features is saturated. This effect is clearly visible in the plots for $CCC \geq 0.9$ and $CCC \geq 0.95$, where fewer features tend to remain after thresholding, due to the stringent reproducibility requirement. Broadly, we observe that the best performance tends to be obtained using 4 or 8 features, insofar as there are sufficiently many features available after thresholding. Again we can observe that, although the effect appears small, the training set performance is slightly reduced when the features are thresholded on CCC, indicating a potential reduction in overfitting. For higher thresholds, this effect is primarily due to the number of features

being saturated, but for $CCC \geq 0.8$ and $CCC \geq 0.85$ the number of available features is high enough to observe the difference even before feature count saturation.

A visualization of tumors with different values for features with high/low C-index (prognostic value), and high/low CCC (reproducibility) is given in Figure 8.

4. Discussion

Although in the present study we have restricted our analysis to the reproducibility of radiomic features when slice thickness is varied at reconstruction time, the ultimate purpose of the data collection effort from which this study draws is to look at the effects of contrast timing and reconstruction parameters using a test-retest paradigm with two portal venous phase images collected within 15s of each other. As part of a multicenter prospective study systematically varying contrast timing and reconstruction parameters, this paper presents early results from a research effort that will be a substantial step forward in the understanding of the reproducibility of radiomics for contrast-enhanced CT of CRLM.

Our results show that when comparing the 20% ASiR images pairwise at different slice thicknesses, the CCCs for features from the 5 mm and 2.5 mm images are significantly lower than when we compare 5 mm and 3.75 mm, or 2.5 and 3.75 mm pairs. The degradation in consistency of features as the change in slice thickness increases is consistent with findings in the literature from phantom studies (Zhao et al., 2014; Kim et al., 2019; Ligerio et al., 2021), as well as lung cancer CT (Lu et al., 2016; Park et al., 2019; Erdal et al., 2020). The consistency of this finding across cancers is noteworthy given that feature reproducibility is anatomy/disease-specific, even within CT (van Timmeren et al., 2016). We also found, at a lower level of statistical significance, that the features had higher CCC in the 5 mm and 3.75 mm pairs, compared to the 2.5 mm and 3.75 mm pairs. We hypothesize that this result could be due to the lower noise level in thicker sliced images, compared to thinner slices, or due to resampling issues when upsampling the ROI segmentations to 2.5 mm.

Moving to liver cancer more specifically, Perrin et al. (2018) examined a database of consecutive patients from MSK diagnosed with liver malignancy, with the additional inclusion requirement that they had two contrast-enhanced abdominal CT scans taken within no more than 14 days of each other. By including two scans per patient with a small separation in time, this data set served as an approximation of a test-retest study, allowing the study of the reproducibility of radiomic features across two consecutive scans. The database included patients with multiple liver cancers including liver metastases (n=22), intrahepatic cholangiocarcinoma (n=10), and hepatocellular carcinoma (n=6). Because image acquisition and reconstruction parameters naturally varied between the scans for a given patient, this data set allowed the exploration of the effects of these different parameters on the overall concordance of the radiomic features between the consecutive scans, as measured by CCC. One parameter that Perrin et al. (2018) reported on was pixel spacing, which is closely related to slice thickness, given that both variables directly affect the resulting voxel volume of the 3D image. They found that scan pairs with a greater difference in pixel spacing had lower feature agreement for features extracted from both tumor and liver parenchyma ROIs, highlighting along with our own results that the feature reproducibility across variations of any parameters affecting voxel volume, including slice thickness and field of view, is an important consideration when drawing inference from radiomic features. For all the variables

they considered, including pixel spacing, Perrin et al. (2018) also found that the tumor tended to have a larger number of reproducible features than the liver parenchyma. Similarly, our reproducibility analysis showed that features extracted from the liver parenchyma were overall less reproducible across slice thickness variation than those drawn from the largest metastasis, regardless of the extractor setting under consideration.

Despite the lower reproducibility of the features drawn from the liver parenchyma, our univariate survival analysis showed that the liver features still contained a relevant signal with regard to overall patient survival in the independent retrospective data set. Indeed, the highest overall univariate C-index was attained by a feature from the liver parenchyma—GLDM SmallDependenceLowGrayLevelEmphasis (C-index=0.6176)—and features from the liver parenchyma were well-represented on the Pareto front for C-index and CCC (see Table 5). Similarly, in the top-performing multivariable models, across all cross-validation runs, the average proportion of selected features coming from the liver was usually over 50% (see Table 4). The presence of important information relevant to patient outcomes in the radiological texture of the liver parenchyma is biologically feasible, and consistent with previous analyses in the literature (Simpson et al., 2017; Hu et al., 2022b). Given the apparent importance of the liver parenchyma features, future studies should investigate methods of tailoring the feature extraction to improve the reproducibility of the features extracted from the liver parenchyma, without degrading their prognostic value.

The IBSI reference manual (Zwanenburg et al., 2016) recommends resampling to a common voxel size to ensure the features are comparable between images with different resolutions. However, several empirical studies have shown that resampling has a limited ability to mitigate the effects of slice thickness, or resolution generally, on radiomic features. Ligero et al. (2021) showed in a phantom study that resampling could improve feature concordance across varying slice thickness to a limited extent, but found that harmonization of features treating slice thickness as a batch effect had a stronger effect. In lung CT, studies have found limited improvement when resampling images to the same voxel volume with linear interpolation (Yang et al., 2021), although deep-learning-based super-resolution algorithms have shown promise (Park et al., 2019). In a study varying pixel spacing by varying field of view during reconstruction, Mackin et al. (2017) found that resampling on its own actually worsened reproducibility of features across pixel sizes; however, combining resampling with a low-pass Butterworth filter improved feature concordance. Shafiq-Ul-Hassan et al. (2017) found in a phantom study that resampling improved feature concordance across images acquired with different voxel sizes for only a small subset of features. The majority of these features could also be corrected without resampling, by including a correction factor based on the voxel volume and/or ROI volume—a result which they subsequently validated in a lung cancer cohort (Shafiq-ul Hassan et al., 2018). Using a similar approach to feature correction, Escudero Sanchez et al. (2021), in a study of contrast enhanced CT imaging of liver cancer, found that while resampling to the mean spacing in all directions improved the number of features that were reproducible across changes in slice thickness, after feature correction, the result was modest in terms of absolute improvement over extracting features from the original images. In aggregate, these studies show that isotropic resampling is, on its own, not sufficient to mitigate the effects of slice thickness variation on the resulting radiomic feature values.

Our results also indicate limitations of resampling to mitigate voxel resolution effects. In the present study, we were further motivated to question the use of isotropic resampling—as recommended by IBSI (Zwanenburg et al., 2016) to preserve rotational invariance of feature definitions—due to the inherently anisotropic nature of our abdominal images. In our previous study on this topic (Peoples et al., 2023), we showed that aggregating features in 2.5D with no z-axis resampling (i.e. L2i and S2i) tended to produce more reproducible texture features, and still produced good prognostic models. In the present work, we conducted a more in-depth analysis of the relationship between feature extraction method, reproducibility, and prognostic value. From Figure 4, we can see that no one extraction method produces more reproducible or more prognostic features in every case. Indeed, both Figure 4 and Figure 5 illustrate that there is no consistently “best” feature extraction approach across all features and ROIs, when we consider either C-index or CCC in isolation, or when we consider both by looking at the Pareto front. Furthermore, while the best overall performance in the multivariable survival models was obtained by using features from L2i, the combination of all extractors was able to rival this performance when we removed features with $CCC < 0.85$ across slice thicknesses. These results suggest that when feature reproducibility across important variations in imaging protocol is known a priori (such as from a test-retest study), including reproducibility into the feature selection process can improve model performance. In this case, highly reproducible and prognostic models can be achieved without optimizing the feature extraction process, by including features extracted using a variety of settings, and allowing the best features to be selected in a data-driven manner accounting for both reproducibility and discriminative ability. This data-driven approach is reminiscent of a method used by Vallières et al. (2017), wherein features extracted with multiple settings, such as resampling resolutions, were pooled into a larger table, and allowed to be considered as separate features during model building. In the absence of reproducibility scores for each feature, choosing settings that can produce features that are more robust to protocol variations present in the data set may be more important, but it is unclear how to choose these settings without a reproducibility study, given the study-specific nature of feature CCCs (van Timmeren et al., 2016).

One limitation of this study is the use of only one fixed bin count of 24 when computing the texture features. The bin count is a key parameter for computing texture features, and has a large effect on the results; though, for many features this effect is at least somewhat predictable (Shafiq-Ul-Hassan et al., 2017; Shafiq-ul Hassan et al., 2018). Furthermore, phantom studies have suggested that variation of the discretization, although it affects the feature values, does not have a large effect on feature reproducibility (Larue et al., 2017). Due to the large number of modeling experiments, and already large number of dimensions under consideration in the present study, we chose to avoid adding bin count as an additional variable at this time. In a study of contrast enhanced CT of hepatocellular carcinoma, Escudero Sanchez et al. (2021) found that the optimal bin count for feature reproducibility across slice thicknesses was in the range of (32–64), which given the step-size in their analysis, is nearly encompassing our chosen value of 24. Despite this, future work should consider the effect of bin count on feature reproducibility and prognostic value.

In conclusion, our results demonstrate the strong effect of slice thickness on feature reproducibility for contrast enhanced CT imaging of CRLM. Although some methods of feature extraction may mitigate the effects of slice thickness on some features, overall, we

found that the extractor producing the most reproducible value for a given feature can vary across features and ROIs. Similarly, the greatest discriminative ability for a given feature may be attained by different extractors, dependent on feature, and ROI. Given this, our results support a data-driven approach, where features from a variety of extractor settings are all considered, and selected in a manner accounting for reproducibility across relevant variations in protocol. Where disease-specific reproducibility metrics are not available, some methods of feature extraction may perform better due to improvements in reproducibility across certain prognostic features (such as the L2i features in the present study), however it is unclear how to determine this without a reproducibility study. Overall, our results demonstrate that we can find radiomic features that are both reproducible across slice thickness variation, and prognostic in patients undergoing hepatic resection, in the context of contrast enhanced CT of colorectal liver metastases.

Acknowledgments

This work was supported by National Institutes of Health grant R01 CA233888.

Ethical Standards

The work follows appropriate ethical standards in conducting research and writing the manuscript, following all applicable laws and regulations regarding treatment of animals or human subjects. The study was conducted in accordance with the principles described in the declaration of Helsinki, with approval of our local IRB and informed consent of all patients.

Conflicts of Interest

We declare that we have no conflicts of interest.

Data availability

The retrospective data set used for survival analysis is publicly available from TCIA (Simpson et al., 2023). The prospective data set used for the reproducibility analysis is still being actively collected and prepared, and is therefore not available, though our intent is ultimately to release it on TCIA.

References

Ziv Bar-Joseph, David K Gifford, and Tommi S Jaakkola. Fast optimal leaf ordering for hierarchical clustering. *Bioinformatics*, 17(suppl_1):S22–S29, June 2001.

Roberto Berenguer, María Del Rosario Pastor-Juan, Jesús Canales-Vázquez, Miguel Castro-García, María Victoria Villas, Francisco Mansilla Legorburo, and Sebastià Sabater. Ra-

- diomics of CT features may be nonreproducible and redundant: Influence of CT acquisition parameters. *Radiology*, 288(2):407–415, August 2018.
- Josep L Carrasco and Lluís Jover. Estimating the generalized concordance correlation coefficient through variance components. *Biometrics*, 59(4):849–858, December 2003.
- Kenneth Clark, Bruce Vendt, Kirk Smith, John Freymann, Justin Kirby, Paul Koppel, Stephen Moore, Stanley Phillips, David Maffitt, Michael Pringle, Lawrence Tarbox, and Fred Prior. The cancer imaging archive (TCIA): Maintaining and operating a public information repository. *J. Digit. Imaging*, 26(6):1045–1057, 2013.
- Chris Ding and Hanchuan Peng. MINIMUM REDUNDANCY FEATURE SELECTION FROM MICROARRAY GENE EXPRESSION DATA. *J. Bioinform. Comput. Biol.*, 03(02):185–205, 2005.
- Nastaran Emaminejad, Muhammad Wasil Wahi-Anwar, Grace Hyun J Kim, William Hsu, Matthew Brown, and Michael McNitt-Gray. Reproducibility of lung nodule radiomic features: Multivariable and univariable investigations that account for interactions between CT acquisition and reconstruction parameters. *Med. Phys.*, 48(6):2906–2919, June 2021.
- Barbaros S Erdal, Mutlu Demirer, Kevin J Little, Chiemezie C Amadi, Gehan F M Ibrahim, Thomas P O’Donnell, Rainer Grimmer, Vikash Gupta, Luciano M Prevedello, and Richard D White. Are quantitative features of lung nodules reproducible at different CT acquisition and reconstruction parameters? *PLoS One*, 15(10):e0240184, October 2020.
- Lorena Escudero Sanchez, Leonardo Rundo, Andrew B Gill, Matthew Hoare, Eva Mendes Serrao, and Evis Sala. Robustness of radiomic features in CT images with different slice thickness, comparing liver tumour and muscle. *Sci. Rep.*, 11(1):8262, April 2021.
- Francesco Fiz, Luca Viganò, Nicolò Gennaro, Guido Costa, Ludovico La Bella, Alexandra Boichuk, Lara Cavinato, Martina Sollini, Letterio S Politi, Arturo Chiti, and Guido Torzilli. Radiomics of liver metastases: A systematic review. *Cancers*, 12(10):2881, October 2020.
- Rachel B Ger, Shouhao Zhou, Pai-Chun Melinda Chi, Hannah J Lee, Rick R Layman, A Kyle Jones, David L Goff, Clifton D Fuller, Rebecca M Howell, Heng Li, R Jason Stafford, Laurence E Court, and Dennis S Mackin. Comprehensive investigation on controlling for CT imaging variabilities in radiomics studies. *Sci. Rep.*, 8(1):13047, August 2018.
- Robert J Gillies, Paul E Kinahan, and Hedvig Hricak. Radiomics: Images are more than pictures, they are data. *Radiology*, 278(2):563–577, February 2016. ppublish.
- Mohammad Hamghalam, Alejandro F Frangi, Baiying Lei, and Amber L Simpson. Modality completion via gaussian process prior variational autoencoders for multi-modal glioma segmentation. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: Part VII*, pages 442–452, 2021.

- Mohammad Hamghalam, Richard K G Do, and Amber L Simpson. Attention-based CT scan interpolation for lesion segmentation of colorectal liver metastases. In *Proc. SPIE 12468, Medical Imaging 2023: Biomedical Applications in Molecular, Structural, and Functional Imaging*, volume 12468, pages 186–193, 2023.
- F E Harrell, Jr, K L Lee, and D B Mark. Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat. Med.*, 15(4):361–387, February 1996. ppublish.
- Lan He, Yanqi Huang, Zelan Ma, Cuishan Liang, Changhong Liang, and Zaiyi Liu. Effects of contrast-enhancement, reconstruction slice thickness and convolution kernel on the diagnostic performance of radiomics signature in solitary pulmonary nodule. *Sci. Rep.*, 6(1):34921, October 2016.
- Natally Horvat, Joao Miranda, Maria El Homsy, Jacob J Peoples, Niamh M Long, Amber L Simpson, and Richard K G Do. A primer on texture analysis in abdominal radiology. *Abdom. Radiol. (NY)*, 47(9):2972–2985, September 2022.
- Peng Hu, Liye Chen, Yaoying Zhong, Yudong Lin, Xiaojing Yu, Xi Hu, Xinwei Tao, Shushen Lin, Tianye Niu, Ran Chen, Xia Wu, and Jihong Sun. Effects of slice thickness on CT radiomics features and models for staging liver fibrosis caused by chronic liver disease. *Jpn. J. Radiol.*, 40(10):1061–1068, October 2022a.
- Ricky Hu, Ishita Chen, Jacob Peoples, Jean-Paul Salameh, Mithat Gönen, Paul B Romesser, Amber L Simpson, and Marsha Reyngold. Radiomics artificial intelligence modelling for prediction of local control for colorectal liver metastases treated with radiotherapy. *Phys. Imaging Radiat. Oncol.*, 24:36–42, October 2022b.
- Abdalla Ibrahim, Bruno Barufaldi, Turkey Refaee, Telmo M Silva Filho, Raymond J Acciavatti, Zohaib Salahuddin, Roland Hustinx, Felix M Mottaghy, Andrew D A Maidment, and Philippe Lambin. MaasPenn radiomics reproducibility score: A novel quantitative measure for evaluating the reproducibility of CT-based handcrafted radiomic features. *Cancers*, 14(7):1599, March 2022.
- Fabian Isensee, Paul F Jaeger, Simon A A Kohl, Jens Petersen, and Klaus H Maier-Hein. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nat. Methods*, 18(2):203–211, February 2021.
- Ron Kikinis, Steve D Pieper, and Kirby G Vosburgh. 3D Slicer: A platform for subject-specific image analysis, visualization, and clinical support. In *Intraoperative Imaging and Image-Guided Therapy*, pages 277–289. Springer New York, 2013.
- Young Jae Kim, Hyun-Ju Lee, Kwang Gi Kim, and Seung Hyun Lee. The effect of CT scan parameters on the measurement of CT radiomic features: A lung nodule phantom study. *Comput. Math. Methods Med.*, 2019:8790694, February 2019.
- Ruben T H M Larue, Janna E van Timmeren, Evelyn E C de Jong, Giacomo Feliciani, Ralph T H Leijenaar, Wendy M J Schreurs, Meindert N Sosef, Frank H P J Raat, Frans H R van der Zande, Marco Das, Wouter van Elmpt, and Philippe Lambin. Influence

- of gray level discretization on radiomic feature stability for different CT scanners, tube currents and slice thicknesses: a comprehensive phantom study. *Acta Oncologica*, 56(11):1544–1553, 2017.
- Yajun Li, Lin Lu, Manjun Xiao, Laurent Dercle, Yue Huang, Zishu Zhang, Lawrence H Schwartz, Daiqiang Li, and Binsheng Zhao. CT slice thickness and convolution kernel affect performance of a radiomic model for predicting EGFR status in non-small cell lung cancer: A preliminary study. *Sci. Rep.*, 8(1):17913, December 2018.
- Marta Ligeró, Olivia Jordi-Ollero, Kinga Bernatowicz, Alonso Garcia-Ruiz, Eric Delgado-Muñoz, David Leiva, Richard Mast, Cristina Suarez, Roser Sala-Llonch, Nahum Calvo, Manuel Escobar, Arturo Navarro-Martin, Guillermo Villacampa, Rodrigo Dienstmann, and Raquel Perez-Lopez. Minimizing acquisition-related radiomics variability by image resampling and batch effect correction to allow for large-scale data analysis. *Eur. Radiol.*, 31(3):1460–1470, March 2021.
- Lawrence I-Kuei Lin. A concordance correlation coefficient to evaluate reproducibility. *Biometrics*, 45(1):255–268, March 1989.
- Lin Lu, Ross C Ehmke, Lawrence H Schwartz, and Binsheng Zhao. Assessing agreement between radiomic features computed for multiple CT imaging settings. *PLoS One*, 11(12):e0166550, December 2016.
- Dennis Mackin, Xenia Fave, Lifei Zhang, Jinzhong Yang, A Kyle Jones, Chaan S Ng, and Laurence Court. Harmonizing the pixel size in retrospective computed tomography radiomics studies. *PLoS One*, 12(9):e0178524, September 2017.
- Mathias Meyer, James Ronald, Federica Vernuccio, Rendon C Nelson, Juan Carlos Ramirez-Giraldo, Justin Solomon, Bhavik N Patel, Ehsan Samei, and Daniele Marin. Reproducibility of CT radiomic features within the same patient: Influence of radiation dose and CT reconstruction settings. *Radiology*, 293(3):583–591, December 2019.
- Ramtin Mojtahedi, Mohammad Hamghalam, Richard K G Do, and Amber L Simpson. Towards optimal patch size in vision transformers for tumor segmentation. In *Multiscale Multimodal Medical Imaging. MMMI 2022*, pages 110–120, 2022.
- Daniel Müllner. Modern hierarchical, agglomerative clustering algorithms. *arXiv:1109.2378 [stat. ML]*, 2011.
- Sohee Park, Sang Min Lee, Kyung Hyun Do, June Goo Lee, Woong Bae, Hyunho Park, Kyu Hwan Jung, and Joon Beom Seo. Deep learning algorithm for reducing CT slice thickness: Effect on reproducibility of radiomic features in lung cancer. *Korean J. Radiol.*, 20(10):1431–1440, October 2019.
- Jacob J Peoples, Mohammad Hamghalam, Imani James, Maida Wasim, Natalie Gangai, Hyunseon Christine Kang, Xiujiang John Rong, Yun Shin Chun, Richard K G Do, and Amber L Simpson. Examining the effects of slice thickness on the reproducibility of CT radiomics for patients with colorectal liver metastases. In *Uncertainty for Safe Utilization of Machine Learning in Medical Imaging. UNSURE 2023*, pages 42–52, 2023.

- Thomas Perrin, Abhishek Midya, Rikiya Yamashita, Jayasree Chakraborty, Tome Saidon, William R Jarnagin, Mithat Gonen, Amber L Simpson, and Richard K G Do. Short-term reproducibility of radiomic features in liver parenchyma and liver malignancies on contrast-enhanced CT imaging. *Abdom. Radiol. (NY)*, 43(12):3271–3278, 2018.
- David L Raunig, Lisa M McShane, Gene Pennello, Constantine Gatsonis, Paul L Carson, James T Voyvodic, Richard L Wahl, Brenda F Kurland, Adam J Schwarz, Mithat Gonen, Gudrun Zahlmann, Marina V Kondratovich, Kevin O’Donnell, Nicholas Petrick, Patricia E Cole, Brian Garra, Daniel C Sullivan, and QIBA Technical Performance Working Group. Quantitative imaging biomarkers: A review of statistical methods for technical performance assessment. *Stat. Methods Med. Res.*, 24(1):27–67, June 2014.
- Muhammad Shafiq-Ul-Hassan, Geoffrey G Zhang, Kujtim Latifi, Ghanim Ullah, Dylan C Hunt, Yoganand Balagurunathan, Mahmoud Abraham Abdalah, Matthew B Schabath, Dmitry G Goldgof, Dennis Mackin, Laurence Edward Court, Robert James Gillies, and Eduardo Gerardo Moros. Intrinsic dependencies of CT radiomic features on voxel size and number of gray levels. *Med. Phys.*, 44(3):1050–1062, March 2017.
- Muhammad Shafiq-ul Hassan, Kujtim Latifi, Geoffrey Zhang, Ghanim Ullah, Robert Gillies, and Eduardo Moros. Voxel size and gray level normalization of CT radiomic features in lung cancer. *Sci. Rep.*, 8(1):10545, July 2018.
- Amber L Simpson, Alexandre Doussot, John M Creasy, Lauryn B Adams, Peter J Allen, Ronald P DeMatteo, Mithat Gonen, Nancy E Kemeny, T Peter Kingham, Jinru Shia, William R Jarnagin, Richard K G Do, and Michael I D’Angelica. Computed tomography image texture: A noninvasive prognostic marker of hepatic recurrence after hepatectomy for metastatic colorectal cancer. *Ann. Surg. Oncol.*, 24(9):2482–2490, 2017.
- Amber L Simpson, Jacob Peoples, John M Creasy, Gabor Fichtinger, Natalie Gangai, Andras Lasso, Krishna Nand Keshava Murthy, Jinru Shia, Michael I D’Angelica, and Richard K G Do. Preoperative CT and survival data for patients undergoing resection of colorectal liver metastases (Colorectal-Liver-Metastases) (Version 2) [Data set], 2023.
- Amber L Simpson, Jacob Peoples, John M Creasy, Gabor Fichtinger, Natalie Gangai, Krishna N Keshavamurthy, Andras Lasso, Jinru Shia, Michael I D’Angelica, and Richard K G Do. Preoperative CT and survival data for patients undergoing resection of colorectal liver metastases. *Sci. Data*, 11(1):172, February 2024.
- Alberto Traverso, Leonard Wee, Andre Dekker, and Robert Gillies. Repeatability and reproducibility of radiomic features: A systematic review. *Int. J. Radiat. Oncol. Biol. Phys.*, 102(4):1143–1158, November 2018.
- Martin Vallières, Emily Kay-Rivest, Léo Jean Perrin, Xavier Liem, Christophe Furstoss, Hugo J W L Aerts, Nader Khaouam, Phuc Felix Nguyen-Tan, Chang-Shu Wang, Khalil Sultanem, Jan Seuntjens, and Issam El Naqa. Radiomics strategies for risk assessment of tumour failure in head-and-neck cancer. *Sci. Rep.*, 7(1):10117, August 2017.
- Joost J M van Griethuysen, Andriy Fedorov, Chintan Parmar, Ahmed Hosny, Nicole Aucoin, Vivek Narayan, Regina G H Beets-Tan, Jean-Christophe Fillion-Robin, Steve Pieper,

- and Hugo J W L Aerts. Computational radiomics system to decode the radiographic phenotype. *Cancer Res.*, 77(21):e104–e107, November 2017.
- Janna E van Timmeren, Ralph T H Leijenaar, Wouter van Elmpt, Jiazhou Wang, Zhen Zhang, André Dekker, and Philippe Lambin. Test-retest data for radiomics feature stability analysis: Generalizable or study-specific? *Tomography*, 2(4):361–365, December 2016.
- Bino A Varghese, Darryl Hwang, Steven Y Cen, Joshua Levy, Derek Liu, Christopher Lau, Marielena Rivas, Bhushan Desai, David J Goodenough, and Vinay A Duddalwar. Reliability of CT-based texture features: Phantom study. *J. Appl. Clin. Med. Phys.*, 20(8):155–163, August 2019.
- Joe H Ward. Hierarchical grouping to optimize an objective function. *J. Am. Stat. Assoc.*, 58(301):236–244, March 1963.
- Frank Wilcoxon. Individual comparisons by ranking methods. *Biom. Bull.*, 1(6):80, December 1945.
- Yan Xu, Lin Lu, Shawn H Sun, Lin-Ning E, Wei Lian, Hao Yang, Lawrence H Schwartz, Zheng-Han Yang, and Binsheng Zhao. Effect of CT image acquisition parameters on diagnostic performance of radiomics in predicting malignancy of pulmonary nodules of different sizes. *Eur. Radiol.*, 32(3):1517–1527, March 2022.
- Shouxin Yang, Ning Wu, Li Zhang, and Meng Li. Evaluation of the linear interpolation method in correcting the influence of slice thicknesses on radiomic feature values in solid pulmonary nodules: a prospective patient study. *Ann. Transl. Med.*, 9(4):279, February 2021.
- Binsheng Zhao. Understanding sources of variation to improve the reproducibility of radiomics. *Front. Oncol.*, 11, 2021.
- Binsheng Zhao, Yongqiang Tan, Wei Yann Tsai, Lawrence H Schwartz, and Lin Lu. Exploring variability in CT characterization of tumors: A preliminary phantom study. *Transl. Oncol.*, 7(1):88–93, February 2014.
- Binsheng Zhao, Yongqiang Tan, Wei-Yann Tsai, Jing Qi, Chuanmiao Xie, Lin Lu, and Lawrence H Schwartz. Reproducibility of radiomics for deciphering tumor phenotype with imaging. *Scientific Reports*, 6(1), 2016.
- Alex Zwanenburg, Stefan Leger, Martin Vallières, and Steffen Löck. Image biomarker standardisation initiative: Reference manual. *arXiv:1612.07003 [cs. CV]*, 2016.
- Alex Zwanenburg, Martin Vallières, Mahmoud A Abdalah, Hugo J W L Aerts, Vincent Andrearczyk, Aditya Apte, Saeed Ashrafinia, Spyridon Bakas, Roelof J Beukinga, Ronald Boellaard, Marta Bogowicz, Luca Boldrini, Irène Buvat, Gary J R Cook, Christos Davatzikos, Adrien Depeursinge, Marie-Charlotte Desseroit, Nicola Dinapoli, Cuong Viet Dinh, Sebastian Echegaray, Issam El Naqa, Andriy Y Fedorov, Roberto Gatta, Robert J Gillies, Vicky Goh, Michael Götz, Matthias Guckenberger, Sung Min Ha, Mathieu Hatt,

Fabian Isensee, Philippe Lambin, Stefan Leger, Ralph T H Leijenaar, Jacopo Lenkowitz, Fiona Lippert, Are Losnegård, Klaus H Maier-Hein, Olivier Morin, Henning Müller, Sandy Napel, Christophe Nioche, Fanny Orlhac, Sarthak Pati, Elisabeth A G Pfaehler, Arman Rahmim, Arvind U K Rao, Jonas Scherer, Muhammad Musib Siddique, Nanna M Sijtsema, Jairo Socarras Fernandez, Emiliano Spezi, Roel J H M Steenbakkens, Stephanie Tanadini-Lang, Daniela Thorwarth, Esther G C Troost, Taman Upadhaya, Vincenzo Valentini, Lisanne V van Dijk, Joost van Griethuysen, Floris H P van Velden, Philip Whybra, Christian Richter, and Steffen Löck. The image biomarker standardization initiative: Standardized quantitative radiomics for high-throughput image-based phenotyping. *Radiology*, 295(2):328–338, 2020.

Appendix A. Additional Results

Table 5 lists all features on the Pareto front for CCC and C-index when combining all feature extractor settings and ROIs.

Table 5: Features on the Pareto front for CCC and C-index when combining all feature extractor settings and ROIs.

Extractor	ROI	Class	Name	CCC	C-index
L2i	Tumor	GLSZM	ZoneEntropy	0.9169	0.6114
	Liver	First order	10Percentile	0.9885	0.5857
	Liver	First order	Mean	0.9998	0.5650
S2i	Tumor	GLSZM	ZoneEntropy	0.9248	0.6075
L2	Tumor	First order	Energy	0.9974	0.5800
	Tumor	First order	TotalEnergy	0.9974	0.5800
	Tumor	GLRLM	GrayLevelNonUniformity	0.9986	0.5723
S2	Tumor	First order	Energy	0.9976	0.5785
	Tumor	First order	TotalEnergy	0.9976	0.5785
	Tumor	GLDM	SmallDependenceHighGrayLevelEmphasis	0.9439	0.5869
	Liver	First order	10Percentile	0.9904	0.5842
	Liver	First order	Mean	0.9997	0.5651
A3	Tumor	GLDM	DependenceEntropy	0.8120	0.6157
	Tumor	GLRLM	GrayLevelNonUniformity	0.9985	0.5724
L3	Tumor	First order	Energy	0.9974	0.5800
	Tumor	First order	TotalEnergy	0.9974	0.5800
	Liver	GLRLM	ShortRunLowGrayLevelEmphasis	0.8450	0.6138
S3	Tumor	First order	Energy	0.9976	0.5785
	Tumor	First order	TotalEnergy	0.9976	0.5785
	Liver	First order	10Percentile	0.9904	0.5842
	Liver	First order	Mean	0.9997	0.5651
	Liver	GLDM	SmallDependenceLowGrayLevelEmphasis	0.7696	0.6176
	Liver	GLRLM	ShortRunLowGrayLevelEmphasis	0.8773	0.6122

Table 6 lists all features used in this study.

Table 6: All features from all feature classes. Detailed feature definitions can be found in the `pyradiomics` documentation. <https://pyradiomics.readthedocs.io/en/latest/features.html>

Feature Class	Name
First order	10Percentile

Continued on next page

Table 6: All features from all feature classes.

Feature Class	Name
	90Percentile
	Energy
	Entropy
	InterquartileRange
	Kurtosis
	Maximum
	Mean
	MeanAbsoluteDeviation
	Median
	Minimum
	Range
	RobustMeanAbsoluteDeviation
	RootMeanSquared
	Skewness
	TotalEnergy
	Uniformity
	Variance
GLCM	Autocorrelation
	ClusterProminence
	ClusterShade
	ClusterTendency
	Contrast
	Correlation
	DifferenceAverage
	DifferenceEntropy
	DifferenceVariance
	Id
	Idm
	Idmn
	Idn
	Imc1
	Imc2
	InverseVariance
	JointAverage
	JointEnergy
	JointEntropy
	MCC
	MaximumProbability
	SumAverage
	SumEntropy

Continued on next page

Table 6: All features from all feature classes.

Feature Class	Name
	SumSquares
GLDM	DependenceEntropy DependenceNonUniformity DependenceNonUniformityNormalized DependenceVariance GrayLevelNonUniformity GrayLevelVariance HighGrayLevelEmphasis LargeDependenceEmphasis LargeDependenceHighGrayLevelEmphasis LargeDependenceLowGrayLevelEmphasis LowGrayLevelEmphasis SmallDependenceEmphasis SmallDependenceHighGrayLevelEmphasis SmallDependenceLowGrayLevelEmphasis
GLRLM	GrayLevelNonUniformity GrayLevelNonUniformityNormalized GrayLevelVariance HighGrayLevelRunEmphasis LongRunEmphasis LongRunHighGrayLevelEmphasis LongRunLowGrayLevelEmphasis LowGrayLevelRunEmphasis RunEntropy RunLengthNonUniformity RunLengthNonUniformityNormalized RunPercentage RunVariance ShortRunEmphasis ShortRunHighGrayLevelEmphasis ShortRunLowGrayLevelEmphasis
GLSZM	GrayLevelNonUniformity GrayLevelNonUniformityNormalized GrayLevelVariance HighGrayLevelZoneEmphasis LargeAreaEmphasis LargeAreaHighGrayLevelEmphasis LargeAreaLowGrayLevelEmphasis LowGrayLevelZoneEmphasis SizeZoneNonUniformity

Continued on next page

Table 6: All features from all feature classes.

Feature Class	Name
	SizeZoneNonUniformityNormalized
	SmallAreaEmphasis
	SmallAreaHighGrayLevelEmphasis
	SmallAreaLowGrayLevelEmphasis
	ZoneEntropy
	ZonePercentage
	ZoneVariance
NGTDM	Busyness
	Coarseness
	Complexity
	Contrast
	Strength