

# 1 Radiomic signature accurately predicts the risk of 2 metastatic dissemination in late-stage non-small cell lung 3 cancer.

4 Agata Małgorzata Wilk<sup>1, 2, §</sup>, Emilia Kozłowska<sup>1, §,\*</sup>, Damian Borys<sup>1,3</sup>, Andrea  
5 D'Amico<sup>3</sup>, Krzysztof Fajarewicz<sup>1</sup>, Izabela Gorczewska<sup>3</sup>, Iwona Dębosz-  
6 Suwińska<sup>4</sup>, Rafał Suwiński<sup>5</sup>, Jarosław Śmieja<sup>1</sup>, Andrzej Swierniak<sup>1\*</sup>.

7 <sup>1</sup> Department of Systems Biology and Engineering, Silesian University of Tech-  
8 nology, Akademicka 16, 44-100, Gliwice, Poland; <sup>2</sup> Department of Biostatistics  
9 and Bioinformatics, Maria Skłodowska-Curie National Research Institute of On-  
10 cology, Gliwice Branch, Wybrzeże Armii Krajowej 15, 44-102, Gliwice, Poland;  
11 <sup>3</sup> Department of Nuclear Medicine and Endocrine Oncology, PET Diagnostics  
12 Unit, Maria Skłodowska-Curie National Research Institute of Oncology, Gliwice  
13 Branch, Wybrzeże Armii Krajowej 15, 44-102, Gliwice, Poland; <sup>4</sup> Department of  
14 Radiotherapy, Maria Skłodowska-Curie National Research Institute of Oncology,  
15 Gliwice Branch, Wybrzeże Armii Krajowej 15, 44-102, Gliwice, Poland; <sup>5</sup> II-nd  
16 Radiotherapy and Chemotherapy Clinic and Teaching Hospital, Maria Skłodow-  
17 ska-Curie National Research Institute of Oncology, Gliwice Branch, Wybrzeże  
18 Armii Krajowej 15, 44-102, Gliwice, Poland.

19 § These authors contributed equally

20 \* Corresponding author

## 22 Corresponding authors:

### 23 Prof. Andrzej Swierniak

24 Department of Systems Biology and Engineering

25 Silesian University of Technology

26 Akademicka 16, 44-100 Gliwice, Poland

27 E-mail: [andrzej.swierniak@polsl.pl](mailto:andrzej.swierniak@polsl.pl)

28 Phone number: 32 2372712

29

### 30 Ph.D. Emilia Kozłowska

31 Department of Systems Biology and Engineering

32 Silesian University of Technology

33 Akademicka 16, 44-100 Gliwice, Poland

34 E-mail: [emilia.kozlowska@polsl.pl](mailto:emilia.kozlowska@polsl.pl)

35 Phone number: +48 32 2372119

36

37 Running title: Predicting the risk of metastasis in lung cancer.

38 Word count: 3824

39 Number of Figures: 6

40 Number of Tables: 1

41

42

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

## 43 **Author Contributions**

- 44 (I) Conception and design: EK, AS.
- 45 (II) Administrative support: AS, KF, JS.
- 46 (III) Provision of study materials or patients: IDS, RS, AD.
- 47 (IV) Collection and assembly of data: IDS, DB, AD, IG, EK.
- 48 (V) Data analysis and interpretation: AMW, EK.
- 49 (VI) Manuscript writing: All authors.
- 50 (VII) Final approval of manuscript: All authors.

## 51 **Abstract**

52

### 53 **Background:**

54 Non-small cell lung cancer (NSCLC) is the most common type of lung cancer, and the median  
55 overall survival is approximately 2-3 years among patients with stage III disease. Furthermore,  
56 it is one of the deadliest types of cancer globally due to non-specific symptoms and the lack of  
57 a biomarker for early detection. The most important decision that clinicians need to make after  
58 a lung cancer diagnosis is the selection of a treatment schedule. This decision is based on,  
59 among others factors, the risk of developing metastasis.

60

### 61 **Methods:**

62 A cohort of 115 NSCLC patients treated using chemotherapy and radiotherapy with curative  
63 intent was retrospectively collated and included patients for whom positron emission tomogra-  
64 phy/computed tomography (PET/CT) images, acquired before radiotherapy, were available.  
65 The PET/CT images were used to compute radiomic features extracted from a region of interest,  
66 the primary tumor. Radiomic and clinical features were then classified to stratify the patients  
67 into short and long time to metastasis, and regression analysis was used to predict the risk of  
68 metastasis.

69

### 70 **Results:**

71 Classification based on binarized metastasis-free survival (MFS) was applied with moderate  
72 success. Indeed, an accuracy of 0.73 was obtained for the selection of features based on the  
73 Wilcoxon test and logistic regression model. However, the Cox regression model for metastasis  
74 risk prediction performed very well, with a concordance index (c-index) score equal to 0.84.

75

### 76 **Conclusions:**

77 It is possible to accurately predict the risk of metastasis in NSCLC patients based on radiomic  
78 features. The results demonstrate the potential use of features extracted from cancer imaging in  
79 predicting the risk of metastasis.

80

81 **Keywords:** NSCLC, metastasis, Cox regression, classification, radiomics

82

83

84

85

86

87

88

89

90

91

## 92 Introduction

93 Lung cancer is one of the most frequently diagnosed cancer types worldwide, constituting over  
94 11% of all cancer cases. With 2.2 million new diagnoses in 2020 alone, it was surpassed in  
95 incidence only by breast cancer, making the lung the most prevalent cancer site in men (with  
96 over 1.43 million diagnoses) and the third most prevalent in women after breast and colorectal  
97 cancers, which had 0.77 million diagnoses (1). While tobacco smoking is recognized as the  
98 primary cause of lung cancer, it can also be attributed to environmental factors such as air pol-  
99 lution, occupational exposure, and genetic predisposition (2–4). It is usually diagnosed at an  
100 advanced stage due to non-specific early-stage symptoms, which is reflected in the very high  
101 mortality rate. Indeed, the five-year survival rate for lung cancer does not exceed 20% (5–7),  
102 thus, it is the leading cause of cancer-related mortality and is responsible for 18% of all deaths  
103 from cancer (1).

104 Diagnosis of lung cancer involves medical imaging, including X-ray and positron emission to-  
105 mography/computed tomography (PET/CT), which allows for classification according to the  
106 tumor node metastasis (TNM) staging system. Detected lesions are sampled by endobronchial  
107 ultrasound (EBUS) guided bronchoscopy and undergo histopathological assessment. Manage-  
108 ment is stage-specific (7), with clinical guidelines divided into early-stage, locally advanced,  
109 and metastatic cancer (8,9). In early-stage lung cancer, lobectomy is the preferred treatment  
110 option. If the tumor is not initially resectable, neoadjuvant chemotherapy can be implemented  
111 to downgrade the tumor, which would eventually allow for surgery. For selected patients with  
112 comorbidities, stereotactic radiotherapy (SABR) may also be considered. For locally advanced  
113 cancer with lymph node involvement, platinum-based chemotherapy administered concurrently  
114 or sequentially with radiotherapy is the most commonly used curative therapeutic option, and  
115 it can be followed by maintenance immunotherapy. For advanced metastatic cancer, immune  
116 checkpoint inhibitors, with or without chemotherapy, are a viable therapeutic option. As mo-  
117 lecular diagnostics becomes routinely available, targeted therapies aimed at epidermal growth  
118 factor receptor (EGFR) (10), fibroblast growth factor receptor (FGFR) (11), anaplastic lym-  
119 phoma kinase (ALK) (12), or Kirsten rat sarcoma virus (KRAS) (13) are being used to treat  
120 mutation carriers.

121 One of the main reasons for the high mortality seen in lung cancer is its invasiveness, and most  
122 patients develop distant metastases. Unfortunately, metastatic tumors are often resistant to treat-  
123 ment, which leads to much shorter survival times for these patients. Although the exact mech-  
124 anisms of metastasis are still being investigated, it is known that cancer cells can spread by both  
125 blood and lymphatic vessels (14). Lung cancer metastases are most frequently observed in the  
126 brain, bones, liver, lung, and adrenal gland (15). Since the occurrence of distant metastasis is  
127 the turning point in the course of the disease, it might be considered an important endpoint in  
128 prognostic analysis, along with the standard endpoints. Furthermore, the ability to predict when  
129 lung cancer will metastasize could guide clinical decision-making and may be used to indicate  
130 the need for therapy intensification in high-risk patients.

131 The search for accurate prognostic biomarkers in lung cancer is hindered by its high heteroge-  
132 neity and complexity. Nonetheless, clinical and molecular characteristics have shown some  
133 promise in predicting metastasis. Metastasis-associated lung adenocarcinoma transcript 1 (*MA-  
134 LAT-1*), a long non-coding ribonucleic acid (RNA), was demonstrated to be significantly asso-  
135 ciated with metastasis in non-small cell lung cancer (NSCLC) (16). Meanwhile, cancer antigen  
136 125 (CA125) and neuron-specific enolase (NSE) were found to be indicative of liver metastasis  
137 (17). NSE, histological type, number of metastatic lymph nodes, and tumor grade were used to

138 construct a nomogram for use in brain metastasis prediction (18). Vimentin expression was also  
139 identified as a potential predictor of brain metastasis in *EGFR*-mutant NSCLC patients (19).

140 Recently, medical imaging has gained attention as an alternative source of biomarkers (20–22).  
141 It has distinct advantages over molecular markers in that it is non-invasive, requires no addi-  
142 tional assays, and utilizes information acquired during a routine diagnostic procedure. As such,  
143 imaging biomarkers are also the fastest to obtain, making them perfect for therapy planning.  
144 Two main strategies can be used to acquire biomarkers. The first is to directly analyze raw  
145 images. Another solution is radiomics, in which segmented images are subjected to feature ex-  
146 traction. This method provides numerical variables that describe the shape and texture of the  
147 region of interest (ROI), which can then be used in statistical or machine-learning models.

148 The radiomics-based approach has been successfully applied for different endpoints in lung  
149 cancer, including overall survival (OS) and progression-free survival (PFS) (23). It has also  
150 shown promising results for the prediction of distant metastases. Coroller et al. (24) selected a  
151 radiomic signature based on CT images to predict distant metastasis in lung adenocarcinoma.  
152 Wu et al. constructed and validated a Cox proportional hazards model using <sup>18</sup>F-fluorodeoxy-  
153 glucose PET (<sup>18</sup>F-FDG PET) imaging to predict freedom of distant metastasis in early-stage  
154 NSCLC patients (25). Fave et al. (26) demonstrated that adding pre-treatment radiomic features  
155 extracted from CT images could improve the ability of clinical prognostic models to predict  
156 distant metastasis (26). Meanwhile, Dou et al. (27) focused on locally advanced lung adenocar-  
157 cinoma and investigated radiomic features from the primary tumor and peritumoral region (27).

158 In this study, 115 NSCLC patients with various histological subtypes were retrospectively an-  
159 alyzed. The prognostic value of standard clinical features, and radiomic features extracted from  
160 PET/CT images acquired for radiotherapy planning, were evaluated by determining if they  
161 could be used to predict time to distant metastasis. To answer this question, machine learning  
162 models were constructed for continuous and categorical metastasis-free survival (MFS) predic-  
163 tion.

164

## 165 **Materials and Methods**

### 166 **Study design**

167 A cohort of NSCLC patients was collated to investigate if PET/CT imaging routinely performed  
168 for radiotherapy planning could help in planning the future treatment strategy, with a focus on  
169 predicting the risk and time of relapse with distant metastases. MFS was defined as the time  
170 elapsed between diagnosis and the detection of distant metastasis or the time of death/last fol-  
171 low-up if distant metastases did not emerge. In addition, classification algorithms were used to  
172 predict if MFS would be short or long.

173 As the prediction of metastasis risk was the focus of the study, the primary lung cancer tumor  
174 was the ROI. Using the available PET/CT scans, radiomic features were extracted from the ROI  
175 and assessed.

176 The specific clinical question considered in this work was whether or not a radiomic signature  
177 could be extracted that would help discriminate between a primary tumor that has the potential  
178 to metastasize early from one that metastasizes late or not at all.

## 179 **Study population**

180 Data were collected retrospectively at the Maria Sklodowska-Curie National Research Institute  
181 of Oncology, Gliwice Branch (NRIO). The cohort consisted of 115 patients with NSCLC who  
182 were treated with curative intent at the Institute between 2009 and 2017. All patients in the  
183 cohort had been treated with a combination of chemotherapy and radiotherapy. Most of the  
184 patients received a platinum-based doublet with vinorelbine. Patients received between one and  
185 six cycles (median four), followed by radiotherapy (RT) with a total dose between 60 and 70  
186 Gray (Gy) in two Gy fractions. The study was approved by the Local Bioethical Committee of  
187 the NRIO in accordance with national regulations. Formal written consent was obtained from  
188 all participants of the study. The clinical data were anonymized before the computational anal-  
189 ysis.

190 All patients underwent PET/CT imaging for radiotherapy planning. Only patients with non-  
191 detectable distant tumors at the onset of treatment were assessed. However, most patients had  
192 locally disseminated tumors to the lymph nodes, as they were diagnosed late due to non-specific  
193 symptoms.

194 In the cohort, 72% of patients were males, and 28% were female. This is consistent with popu-  
195 lation data showing that most lung cancer patients are male. The median age of patients in the  
196 cohort was 61 years, and over half of the patients had tumors located in the left lung. The most  
197 prevalent cancer subtype was squamous cell carcinoma, which constituted two-thirds of all  
198 cases, followed by large cell carcinoma (24.3%) and adenocarcinoma (7.0%). Detailed charac-  
199 teristics of the cohort are presented in Table 1.

200 The median time-to-metastasis was 2.77 years, with a secondary tumor observed most fre-  
201 quently in the second lung, brain, bones, and liver.

202

## 203 **Positron emission tomography/computed tomography data acquisition and** 204 **segmentation**

205 The PET/CT images were acquired at the NRIO using Philips GeminiGXL 16 (Philips, Am-  
206 sterdam, Netherlands) (24 patients) and Siemens Biograph mCT 131 (Siemens AG, Munich,  
207 Germany) (88 patients) PET/CT scanners. For each patient, the ROI was contoured by the same  
208 experienced nuclear medicine specialist using Medical Image Merge (MIM) 7.0.1 software and  
209 the PET Edge™ tool (both MIM Software Inc., OH, USA).

## 210 **Extraction of radiomic features**

211 Feature extraction was performed with PyRadiomics version 3.0.1, a Python package designed  
212 to increase the reproducibility of radiomic studies (28). Using the PET dataset, 105 standard  
213 features were calculated. Radiomic features belong to one of three classes, including first-order  
214 statistics such as energy, entropy, and minimum, as well as shape features such as volume,  
215 surface area, and sphericity, and texture features including Gray Level Co-occurrence Matrix  
216 (GLCM), Gray Level Dependence Matrix (GLDM), Gray Level Run Length Matrix (GLRLM),  
217 Gray Level Size Zone Matrix (GLSZM), and Neighboring-Gray Tone Difference Matrix  
218 (NGTDM).



## 219 **Metastasis-free survival categorization**

220 For the classification, a threshold of one year was used to create two classes, which included  
221 patients with MFS below and over this threshold. Due to the presence of censored observations  
222 (in the cohort this primarily signified the patient's death), such stratification divided patients  
223 into a group who suffered either metastasis or death within a year (66 patients), and those who  
224 did not (49 patients). To create subgroups that were more related to the research question, the  
225 binary MFS was defined as “short” if the patient developed metastasis within a year (25 pa-  
226 tients) and “long” if the patient developed metastasis or was censored after longer than a year  
227 (49 patients).

## 228 **Statistical analysis**

229 Statistical analysis was performed using the R environment (version 4.1.3). For survival analy-  
230 sis, survival (version 3.2-13) was used, caret (version 6.0.93) and RandomForest (version 4.7-  
231 1.1) were used for classification, and randomForestSRC (version 3.1.1) was used to perform  
232 random survival forest. A heatmap of the radiomic features was created with ComplexHeatmap  
233 (version 2.10.0).

234 Filtering of the radiomic features was applied based on the Pearson correlation coefficient to  
235 avoid redundancy, with a cutoff threshold equal to 0.9 (see Supplementary Table 1). Since the  
236 PET images were acquired using two scanners, principal component analysis was applied to  
237 determine if there was any grouping of samples due to the scanner used (see Supplementary  
238 Figure 1).

239 The clinical and radiomic features with potential for event-free survival (EFS) and MFS pre-  
240 diction were assessed (see Supplementary Table 2). In addition, differences in the values of  
241 radiomic and clinical features between ‘short’ and ‘long’ MFS patient subgroups were investi-  
242 gated statistically. Fisher's exact test was performed for categorical variables, while the Mann-  
243 Whitney U test was used for continuous variables (see Supplementary Table 3). A log-rank test  
244 was also conducted for both categorical and continuous features (see Supplementary Table 4).  
245 As the log-rank test assesses if there is a significant difference between two or more survival  
246 curves, continuous features were binarized with respect to the median value.

## 247 **Cross-validation**

248 The value of any type of predictive model lies in its applicability to unknown data, and not just  
249 its ability to fit the training data. Cross-validation enables evaluation of the model's ability to  
250 generalize by removing part of the data from the cohort and applying them in the estimation of  
251 model performance. In addition, data partitioning at the beginning of each iteration prevents  
252 information leakage.

253 For a more consistent comparison between the regression and classification results, modified  
254 k-fold data partitioning was applied. Firstly, the data was ordered according to (continuous)  
255 MFS values. Then, the observations were assigned consecutive numbers, from one to five,  
256 which were used as cross-validation folds. Such partitioning ensures proper stratification of  
257 both continuous and binarized MFS.

## 258 **Classification algorithms**

259 The observed relationships between binary MFS and binary EFS and extracted features (both  
260 clinical radiomic) were verified by employing classification models. Firstly, three main feature  
261 selection methods were applied, including Student's t-test, Wilcoxon test, and a mutual infor-  
262 mation test. To investigate the impact of a varying number of features on classification quality,  
263 between 1 and 10 features were tested. Since only the mutual information method handles both  
264 categorical and continuous variables, a hybrid selection was used for the other two methods by  
265 applying the main method for continuous variables and Fisher's exact test for categorical vari-  
266 ables. The categorical variables that passed the significance threshold equal to 0.1 were added  
267 to the model.

268 The following classification methods were tested: K Nearest Neighbor (KNN) with different K  
269 values (for clarity, only the best one, K=5, is presented), random forest, support vector machines  
270 (SVM) with linear and radial kernels, and logistic regression (LogReg). Considering the inconsis-  
271 tent orders of magnitude for radiomic features, a z-score transformation was used to scale  
272 the data. In each k-fold iteration, the scaling parameters (mean and standard deviation) were  
273 determined from the training set and applied to both the training and test sets. Classification  
274 accuracy was then used to assess model performance.

## 275 **Regression algorithms**

276 For the prediction of continuous MFS, Cox proportional hazards regression (using survival R  
277 package) and random survival forest (using randomForestSRC R package) were applied. Vari-  
278 able selection was performed based on univariate analysis, with the Harrell Concordance index  
279 (C-index) adopted as a ranking metric. The model performance was validated using the k-fold  
280 partitioning described above. Again, models containing between 1 and 10 features were tested.

281

## 282 **Radiomic-based risk score**

283 Although cross-validation facilitates the estimation of prediction quality, the results and se-  
284 lected features can be different in each iteration due to subsampling. Therefore, all selections  
285 were repeated on the entire dataset to obtain conclusive feature rankings. To demonstrate the  
286 validity of the obtained signature, the Cox model was chosen, which is the classic approach to  
287 survival data analysis with known interpretation. The patients were then divided into high-risk  
288 and low-risk groups based on the calculated median risk score, and MFS was compared using  
289 Kaplan-Meier curves.

## 290 **Results**

### 291 **Patient characteristics**

292 The cohort included only NSCLC patients, as it is the most common type of lung cancer. Most  
293 patients (67%) had squamous histopathological subtypes, and almost two-thirds had an ad-  
294 vanced stage of the primary tumor (T3 or T4). In total, 37 patients eventually developed distant  
295 metastases. Figure 2A shows a Kaplan-Meier plot for MFS probability in the entire patient  
296 cohort.

297 None of the clinical features of the cohort were informative in relation to the time to metastasis  
298 onset (Supplementary Table 2). This means that clinicians are unable to predict if a particular  
299 patient will develop metastatic cancer, based only on clinical variables at diagnosis. On the  
300 other hand, 34 radiomic features were statistically significant against continuous MFS, 36 fea-  
301 tures against the binarized MFS, and 18 against EFS.

302

### 303 **Integration of clinical and radiomic data**

304 High correlations were observed between the radiomic features, which resulted in only 65 of  
305 105 features passing the initial correlation filtering. The highest redundancy was found for the  
306 first-order features (6 out of were 18 kept) and the lowest for the GLSZM features (15 out of 16  
307 were kept) (see Supplementary Table 1).

308 Correlations between radiomic and clinical features were mostly low, signifying that both da-  
309 tasetes carried independent information. Also, the hierarchical clustering of radiomic features  
310 did not correspond to any discernible grouping of clinical features (see Figure 2).

311 Figure 2B shows the normalized z-score values of radiomic features for each patient. The pa-  
312 tients were divided into short and long EFS groups. As can be seen from the results, the hierar-  
313 chical clustering correctly divided patients into these two groups. Furthermore, it was observed  
314 that the radiomic feature spectrum varied between patients with short and long EFS. This  
315 demonstrates that there is potential for the use of radiomic features in predicting EFS.

### 316 **Classification of advanced non-small cell lung cancer**

317 As expected, no clinical features were selected by the models. The feature rankings obtained  
318 for EFS and MFS prediction differed, which aligns with the different interpretations of these  
319 endpoints. While the rankings varied with respect to feature selection and classification meth-  
320 ods, there was some consistency among the top features. Indeed, TotalEnergy, ZoneEntropy,  
321 and RootMeanSquared favored EFS prediction, while Variance, TotalEnergy, RunLength-  
322 NonUniformity (GLRLM), SizeZoneNonUniformityNormalized, and Maximum2DDiameter-  
323 Column favored MFS prediction.

324 The highest accuracy for EFS prediction (approx. 0.65) was achieved using the SVM classifier  
325 with linear kernels for the mutual information selection of eight features. The highest accuracy  
326 for MFS prediction (approx. 0.73) was achieved using the LogReg classifier for the five features  
327 selected using the Wilcoxon test. Due to the imbalanced classes, with “long” (treated as the  
328 negative class) being the predominant group, the models for MFS tended to yield high speci-  
329 ficity and relatively low sensitivity. Most models performed better for a small number of fea-  
330 tures.

### 331 **Prediction of risk of metastasis**

332 For regression-based models, the tendency was similar, with the highest predictive ability ob-  
333 served for a small number of features. The highest median C-index across folds was reached  
334 for two features (GLRLM and NGTDM Business) in Cox regression and one feature (shapeMi-  
335 norAxisLength) in the random survival forest.



336 The mean C-index for the best set of features using Cox regression was 0.84, whereas the C-  
337 index for the random survival forest was 0.8. The inclusion of more features in the model re-  
338 sulted in a loss of prediction quality due to overfitting. No clinical features were selected for  
339 the best models, which is consistent with the preliminary patient cohort analysis.

340 Feature selection on the entire dataset revealed that the two top features for Cox regression,  
341 SmallAreaLowGrayLevelEmphasis (GLSZM) and GLRLM, also held high-ranking positions  
342 in the classification approach. Therefore, the Cox model was constructed using those two fea-  
343 tures. The high-risk and low-risk groups (Figure 6) had significantly different MFS, with the  
344 log-rank test  $P < 0.001$ .

345

## 346 Discussion

347 Lung cancer is the leading cause of cancer-related death worldwide, claiming over 1.7 million  
348 lives yearly. It is characterized by high invasiveness, and the occurrence of distant spread sig-  
349 nificantly influences survival and treatment options. This necessitates the search for prognostic  
350 biomarkers that could help determine the time to metastasis onset. With the rapid development  
351 of the radiomics field, researchers have turned to medical imaging, which is routinely per-  
352 formed and non-invasive, as a source of information that could shed some light on the tumor  
353 dissemination process and aid clinicians in therapy planning.

354 A cohort of NSCLC patients with different subtypes and stages of the disease was collated. It  
355 was concluded that the standard clinical data available for the patients, except for higher meta-  
356 static potential exhibited by the squamous subtype, were largely uninformative regarding me-  
357 tastasis occurrence. To assess the potential of radiomics for MFS prediction, we extracted 105  
358 radiomic features from PET/CT scans, using the primary tumor as the ROI. Regression and  
359 machine learning methods were then used to select radiomic signatures that could predict the  
360 risk of metastasis and achieved a C-index of 0.84 for the Cox proportional hazards model and  
361 0.8 for the random survival forest, and an accuracy of 0.72 for the KNN classifier. These results  
362 confirm that medical images contain information that could be successfully applied to MFS  
363 prediction.

364 Several studies have shown the potential of radiomic features in predicting distant metastasis  
365 in lung cancer, with most of them focusing on either a particular subtype or stage. Coroller et  
366 al. (24) investigated radiomic features extracted from CT images for predicting distant metas-  
367 tasis in lung adenocarcinoma, which had a C-index of 0.61 on an independent validation set.  
368 Fave et al. (26) demonstrated that combining pre-treatment radiomic features with clinical in-  
369 formation improved the ability of prognostic models to predict distant metastasis in stage III  
370 NSCLC patients, reporting a C-index of 0.63 (24). Wu et al. used features extracted from PET  
371 images to predict the freedom of distant metastasis, with a high C-index of 0.71 in independent  
372 validation (25). However, this work only focused on early-stage lung cancer. Dou et al. (27)  
373 presented an interesting approach, extracting features from both the tumor and tumor rim and  
374 achieving a C-index of 0.64 in a cohort of patients with locally advanced lung adenocarcinoma  
375 (27). In the current work, significantly better model quality was achieved in a cohort including  
376 patients with varying subtypes (squamous cell carcinoma, adenocarcinoma, large cell carci-  
377 noma) and stages.

378 While a regression approach, such as a Cox proportional hazards model, is typically used for  
379 survival-type analysis, the risk score it yields does not directly translate to the time of event

380 occurrence. The C-index only compares pairs of observations, resulting in a global assessment  
381 of whether a higher risk is related to a shorter time-to-event. Therefore, classification was also  
382 performed and achieved an accuracy of 0.72 in cross-validation.

383 After testing several methods and approaches to variable selection, it was observed that similar  
384 predictive ability could be achieved for different feature sets, which indicates that even unre-  
385 lated radiomic features carry equivalent information. Interestingly, the quality dropped drasti-  
386 cally with increased feature numbers in all models. This suggests that features with high pre-  
387 dictive potential perform much worse when combined than when used in isolation, and empha-  
388 sizes the importance of selecting algorithms that are sensitive to feature interactions.

389 Certain variables retained high positions across different selections. These included GLRLM,  
390 NGTDM Strength, and NGTDM Business. This demonstrates that these radiomic features are  
391 important for predicting if and when metastasis will occur in a lung cancer patient.

392 This analysis was not without limitations. While the study design ensured all images were con-  
393 toured by one expert, which prevented bias, this did not allow for an assessment of the repro-  
394 ducibility of radiomic feature extraction. In addition, plans are in place to collect an independent  
395 patient cohort to validate the signature. Future work will also investigate tumor growth and  
396 dissemination dynamics, to achieve more clinically meaningful predictions.

397

## 398 **Conclusions**

399 Based on a cohort comprised of 115 NSCLC patients, clinical features routinely collected dur-  
400 ing diagnostic procedures are not sufficient for the prediction of the risk of metastasis. Medical  
401 images (PET/CT scans) were investigated as a potential source of prognostic markers by as-  
402 sessing radiomic features in various classes of predictive models. A model based on two texture  
403 features (GLSZM and GLRLM) was constructed, which divided the patient cohort into low-  
404 risk and high-risk groups that significantly differed in MFS. The findings of this study have the  
405 potential to help clinicians make adjustments to therapy and create a rational basis for the in-  
406 tensification of systemic treatment in high-risk lung cancer patients.

407

## 408 **Acknowledgments**

409 This work was supported by the Polish National Science Centre, Grant Number: UMO-  
410 2020/37/B/ST6/01959, and Silesian University of Technology statutory research funds. Calcu-  
411 lations were performed on the Ziemowit computer cluster in the Laboratory of Bioinformatics  
412 and Computational Biology, created in the EU Innovative Economy Programme  
413 POIG.02.01.00-00-166/08 and expanded in the POIG.02.03.01-00-040/13 project. Data analy-  
414 sis was partially carried out using the Biotest Platform developed within project  
415 PBS3/B3/32/2015, which was financed by the Polish National Centre of Research and Devel-  
416 opment (NCBiR). This work was carried out in part by the Silesian University of Technology  
417 internal research funding (A.M.W., E.K., D.B., K.F., J.S., and A.S.). The founders have no role  
418 in designing the study and writing the manuscript.

419

## 420 **Footnotes**

### 421 **Reporting checklist**

422 The authors have completed the “**Prediction Model Development and Validation**” reporting  
423 checklist.

### 424 **Data Sharing Statement**

425 The authors submitted a data-sharing statement along with the manuscript.

### 426 **Conflict of interest**

427 The authors declare no competing interests.

### 428 **Ethical statement**

429 The authors are accountable for all aspects of the work and will ensure that questions related to  
430 the accuracy or integrity of any part of the work are appropriately investigated and resolved.  
431 The study was conducted in accordance with the Declaration of Helsinki (as revised in 2013).  
432 The study was approved by the institutional board of Maria Skłodowska-Curie National Re-  
433 search Institute of Oncology (Gliwice Branch), and individual consent for this retrospective  
434 analysis was waived.

435

436

437

438

## 439 **References**

- 440 1. Ferlay J, Colombet M, Soerjomataram I, Parkin DM, Piñeros M, Znaor A, et al. Cancer  
441 statistics for the year 2020: An overview. *Int J Cancer* [Internet]. 2021 Aug 15 [cited  
442 2022 Nov 27];149(4):778–89. Available from: [https://pub-  
443 med.ncbi.nlm.nih.gov/33818764/](https://pub-med.ncbi.nlm.nih.gov/33818764/)
- 444 2. Alberg AJ, Samet JM. Epidemiology of lung cancer [Internet]. Vol. 123, *Chest*. Chest;  
445 2003 [cited 2022 Nov 27]. p. 21S-49S. Available from: [https://pub-  
446 med.ncbi.nlm.nih.gov/12527563/](https://pub-med.ncbi.nlm.nih.gov/12527563/)
- 447 3. dela Cruz CS, Tanoue LT, Matthay RA. Lung Cancer: Epidemiology, Etiology, and  
448 Prevention. Vol. 32, *Clinics in Chest Medicine*. 2011. p. 605–44.
- 449 4. Barta JA, Powell CA, Wisnivesky JP. Global epidemiology of lung cancer [Internet].  
450 Vol. 85, *Annals of Global Health*. *Ann Glob Health*; 2019 [cited 2022 Nov 27]. Availa-  
451 ble from: <https://pubmed.ncbi.nlm.nih.gov/30741509/>
- 452 5. Lu T, Yang X, Huang Y, Zhao M, Li M, Ma K, et al. Trends in the incidence, treat-  
453 ment, and survival of patients with lung cancer in the last four decades. *Cancer Manag  
454 Res* [Internet]. 2019 [cited 2022 Nov 27];11:943–53. Available from: [https://pub-  
455 med.ncbi.nlm.nih.gov/30718965/](https://pub-med.ncbi.nlm.nih.gov/30718965/)

- 456 6. Torre LA, Siegel RL, Jemal A. Lung Cancer Statistics. *Adv Exp Med Biol* [Internet].  
457 2016 [cited 2022 Nov 27];893:1–19. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/26667336/)  
458 [med.ncbi.nlm.nih.gov/26667336/](https://pub-med.ncbi.nlm.nih.gov/26667336/)
- 459 7. Molina JR, Yang P, Cassivi SD, Schild SE, Adjei AA. Non-small cell lung cancer: Epi-  
460 demiology, risk factors, treatment, and survivorship. In: *Mayo Clinic Proceedings* [In-  
461 ternet]. Elsevier Ltd; 2008 [cited 2021 Mar 20]. p. 584–94. Available from:  
462 <http://www.mayoclinicproceedings.org/article/S0025619611607350/fulltext>
- 463 8. Planchard D, Popat S, Kerr K, Novello S, Smit EF, Faivre-Finn C, et al. Metastatic  
464 non-small cell lung cancer: ESMO Clinical Practice Guidelines for diagnosis, treatment  
465 and follow-up. *Ann Oncol* [Internet]. 2018 Oct 1 [cited 2022 Nov 27];29(Suppl  
466 4):iv192–237. Available from: <https://pubmed.ncbi.nlm.nih.gov/30285222/>
- 467 9. Postmus PE, Kerr KM, Oudkerk M, Senan S, Waller DA, Vansteenkiste J, et al. Early  
468 and locally advanced non-small-cell lung cancer (NSCLC): ESMO Clinical Practice  
469 Guidelines for diagnosis, treatment and follow-up. *Annals of Oncology* [Internet]. 2017  
470 [cited 2022 Nov 27];28(suppl\_4):iv1–21. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/28881918/)  
471 [med.ncbi.nlm.nih.gov/28881918/](https://pub-med.ncbi.nlm.nih.gov/28881918/)
- 472 10. Le T, Gerber DE. Newer-Generation EGFR Inhibitors in Lung Cancer: How Are They  
473 Best Used? *Cancers (Basel)* [Internet]. 2019 Mar 1 [cited 2023 Jan 19];11(3). Available  
474 from: <https://pubmed.ncbi.nlm.nih.gov/30875928/>
- 475 11. Zhou Z, Liu Z, Ou Q, Wu X, Wang X, Shao Y, et al. Targeting FGFR in non-small cell  
476 lung cancer: implications from the landscape of clinically actionable aberrations of  
477 FGFR kinases. *Cancer Biol Med* [Internet]. 2021 May 5 [cited 2023 Jan 19];18(2):490.  
478 Available from: [/pmc/articles/PMC8185861/](https://pubmed.ncbi.nlm.nih.gov/35185861/)
- 479 12. Yuan M, Huang LL, Chen JH, Wu J, Xu Q. The emerging treatment landscape of tar-  
480 geted therapy in non-small-cell lung cancer. *Signal Transduct Target Ther* [Internet].  
481 2019 [cited 2022 Nov 27];4(1). Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/31871778/)  
482 [med.ncbi.nlm.nih.gov/31871778/](https://pub-med.ncbi.nlm.nih.gov/31871778/)
- 483 13. Salgia R, Pharaon R, Mambetsariev I, Nam A, Sattler M. The improbable targeted ther-  
484 apy: KRAS as an emerging target in non-small cell lung cancer (NSCLC). *Cell Rep*  
485 *Med* [Internet]. 2021 Jan 19 [cited 2022 Nov 27];2(1). Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/33521700/)  
486 [med.ncbi.nlm.nih.gov/33521700/](https://pub-med.ncbi.nlm.nih.gov/33521700/)
- 487 14. Popper HH. Progression and metastasis of lung cancer. *Cancer Metastasis Rev* [Inter-  
488 net]. 2016 Mar 1 [cited 2022 Nov 27];35(1):75–91. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/27018053/)  
489 [med.ncbi.nlm.nih.gov/27018053/](https://pub-med.ncbi.nlm.nih.gov/27018053/)
- 490 15. Riihimäki M, Hemminki A, Fallah M, Thomsen H, Sundquist K, Sundquist J, et al.  
491 Metastatic sites and survival in lung cancer. *Lung Cancer* [Internet]. 2014 Oct 1 [cited  
492 2022 Nov 27];86(1):78–84. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/25130083/)  
493 [med.ncbi.nlm.nih.gov/25130083/](https://pub-med.ncbi.nlm.nih.gov/25130083/)
- 494 16. Molina JR, Yang P, Cassivi SD, Schild SE, Adjei AA. Non-small cell lung cancer: epi-  
495 demiology, risk factors, treatment, and survivorship. *Mayo Clin Proc* [Internet]. 2008  
496 [cited 2022 Nov 27];83(5):584–94. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/18452692/)  
497 [med.ncbi.nlm.nih.gov/18452692/](https://pub-med.ncbi.nlm.nih.gov/18452692/)
- 498 17. Wang CF, Peng SJ, Liu RQ, Yu YJ, Ge QM, Liang R bin, et al. The Combination of  
499 CA125 and NSE Is Useful for Predicting Liver Metastasis of Lung Cancer. *Dis Mark-*  
500 *ers* [Internet]. 2020 [cited 2022 Nov 27];2020. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/33376560/)  
501 [med.ncbi.nlm.nih.gov/33376560/](https://pub-med.ncbi.nlm.nih.gov/33376560/)
- 502 18. Zhang F, Zheng W, Ying L, Wu J, Wu S, Ma S, et al. A Nomogram to Predict Brain  
503 Metastases of Resected Non-Small Cell Lung Cancer Patients. *Ann Surg Oncol* [Inter-  
504 net]. 2016 Sep 1 [cited 2022 Nov 27];23(9):3033–9. Available from: [https://pub-](https://pub-med.ncbi.nlm.nih.gov/27090794/)  
505 [med.ncbi.nlm.nih.gov/27090794/](https://pub-med.ncbi.nlm.nih.gov/27090794/)

- 506 19. Teocharoen R, Ruangritchankul K, Vinayanuwattikun C, Sriuranpong V, Sitthide-  
507 atphaiboon P. Vimentin expression status is a potential biomarker for brain metastasis  
508 development in EGFR-mutant NSCLC patients. *Transl Lung Cancer Res* [Internet].  
509 2021 Feb 1 [cited 2022 Nov 27];10(2):790–801. Available from: [https://pub-  
511 med.ncbi.nlm.nih.gov/33718022/](https://pub-<br/>510 med.ncbi.nlm.nih.gov/33718022/)
- 512 20. Lee G, Lee HY, Park H, Schiebler ML, van Beek EJR, Ohno Y, et al. Radiomics and  
513 its emerging role in lung cancer research, imaging biomarkers and clinical manage-  
514 ment: State of the art. *Eur J Radiol* [Internet]. 2017 Jan 1 [cited 2022 Nov 27];86:297–  
515 307. Available from: <https://pubmed.ncbi.nlm.nih.gov/27638103/>
- 516 21. Avanzo M, Stancanello J, Pirrone G, Sartor G. Radiomics and deep learning in lung  
517 cancer. *Strahlenther Onkol* [Internet]. 2020 Oct 1 [cited 2022 Nov 27];196(10):879–87.  
518 Available from: <https://pubmed.ncbi.nlm.nih.gov/32367456/>
- 519 22. Stieb S, McDonald B, Gronberg M, Engeseth GM, He R, Fuller CD. Imaging for Tar-  
520 get Delineation and Treatment Planning in Radiation Oncology: Current and Emerging  
521 Techniques. *Hematol Oncol Clin North Am* [Internet]. 2019 Dec 1 [cited 2023 Jan  
522 19];33(6):963–75. Available from: <https://pubmed.ncbi.nlm.nih.gov/31668214/>
- 523 23. Zhang Y, Oikonomou A, Wong A, Haider MA, Khalvati F. Radiomics-based Prognosis  
524 Analysis for Non-Small Cell Lung Cancer. *Sci Rep* [Internet]. 2017 Apr 18 [cited 2022  
525 Nov 27];7. Available from: <https://pubmed.ncbi.nlm.nih.gov/28418006/>
- 526 24. Coroller TP, Grossmann P, Hou Y, Rios Velazquez E, Leijenaar RTH, Hermann G, et  
527 al. CT-based radiomic signature predicts distant metastasis in lung adenocarcinoma.  
528 *Radiother Oncol* [Internet]. 2015 Mar 1 [cited 2022 Nov 27];114(3):345–50. Available  
529 from: <https://pubmed.ncbi.nlm.nih.gov/25746350/>
- 530 25. Wu J, Aguilera T, Shultz D, Gudur M, Rubin DL, Loo BW, et al. Early-Stage Non-  
531 Small Cell Lung Cancer: Quantitative Imaging Characteristics of (18)F Fluorodeoxy-  
532 glucose PET/CT Allow Prediction of Distant Metastasis. *Radiology* [Internet]. 2016  
533 Oct 1 [cited 2022 Nov 27];281(1):270–8. Available from: [https://pub-  
535 med.ncbi.nlm.nih.gov/27046074/](https://pub-<br/>534 med.ncbi.nlm.nih.gov/27046074/)
- 536 26. Fave X, Zhang L, Yang J, MacKin D, Balter P, Gomez D, et al. Delta-radiomics fea-  
537 tures for the prediction of patient outcomes in non-small cell lung cancer. *Sci Rep* [In-  
538 ternet]. 2017 Dec 1 [cited 2022 Nov 27];7(1). Available from: [https://pub-  
540 med.ncbi.nlm.nih.gov/28373718/](https://pub-<br/>539 med.ncbi.nlm.nih.gov/28373718/)
- 541 27. Dou TH, Coroller TP, van Griethuysen JJM, Mak RH, Aerts HJWL. Peritumoral radi-  
542 omics features predict distant metastasis in locally advanced NSCLC. *PLoS One* [Inter-  
543 net]. 2018 Nov 1 [cited 2022 Nov 27];13(11). Available from: [https://pub-  
545 med.ncbi.nlm.nih.gov/30388114/](https://pub-<br/>544 med.ncbi.nlm.nih.gov/30388114/)
- 546 28. van Griethuysen JJM, Fedorov A, Parmar C, Hosny A, Aucoin N, Narayan V, et al.  
547 Computational radiomics system to decode the radiographic phenotype. *Cancer Res*  
548 [Internet]. 2017 Nov 1 [cited 2023 Jan 19];77(21):e104–7. Available from: [https://pub-  
550 med.ncbi.nlm.nih.gov/29092951/](https://pub-<br/>549 med.ncbi.nlm.nih.gov/29092951/)

## 547 Tables

548 **Table 1.** Patient characteristics. For continuous variables, median and quartiles are listed.

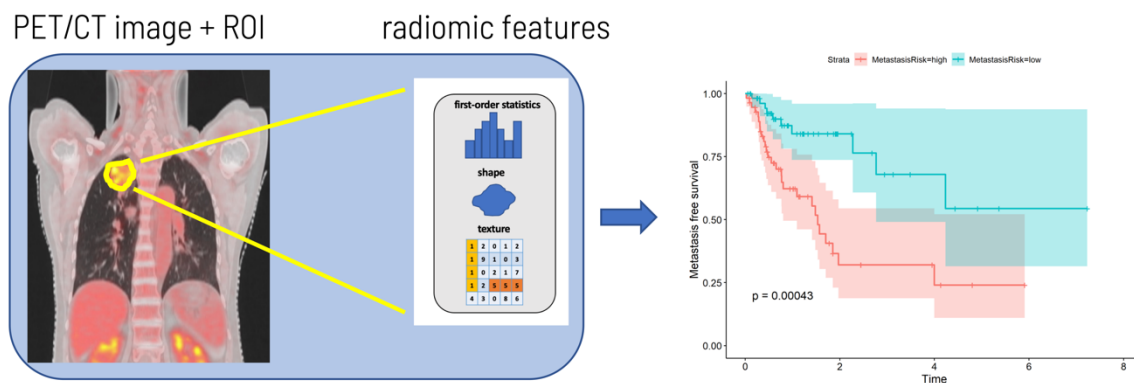
		<b>N = 115</b>
<b>Sex</b>	Male	83 (72.2%)
	Female	32 (27.8%)



<b>Age</b>		61 (57-67)
<b>Histopathology</b>	Squamous	77 (67.0%)
	Adenocarcinoma	8 (7.0%)
	Large cell	28 (24.3%)
	Other	2 (1.7%)
<b>Location</b>	Left	65 (56.5%)
	Right	50 (43.5%)
<b>T</b>	1	4 (3.5%)
	2	37 (32.2%)
	3	37 (32.2%)
	4	37 (32.2%)
<b>N</b>	0	19 (16.5%)
	1	6 (5.2%)
	2	83 (72.2%)
	3	7 (6.1%)
<b>M</b>	0	115 (100%)
	1	0 (0%)
<b>Zubrod score</b>	0	34 (29.6%)
	1	80 (69.6%)
	2	1 (0.9%)

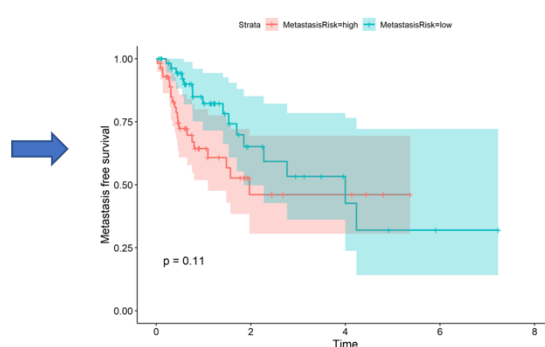
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564

565 **Figures**



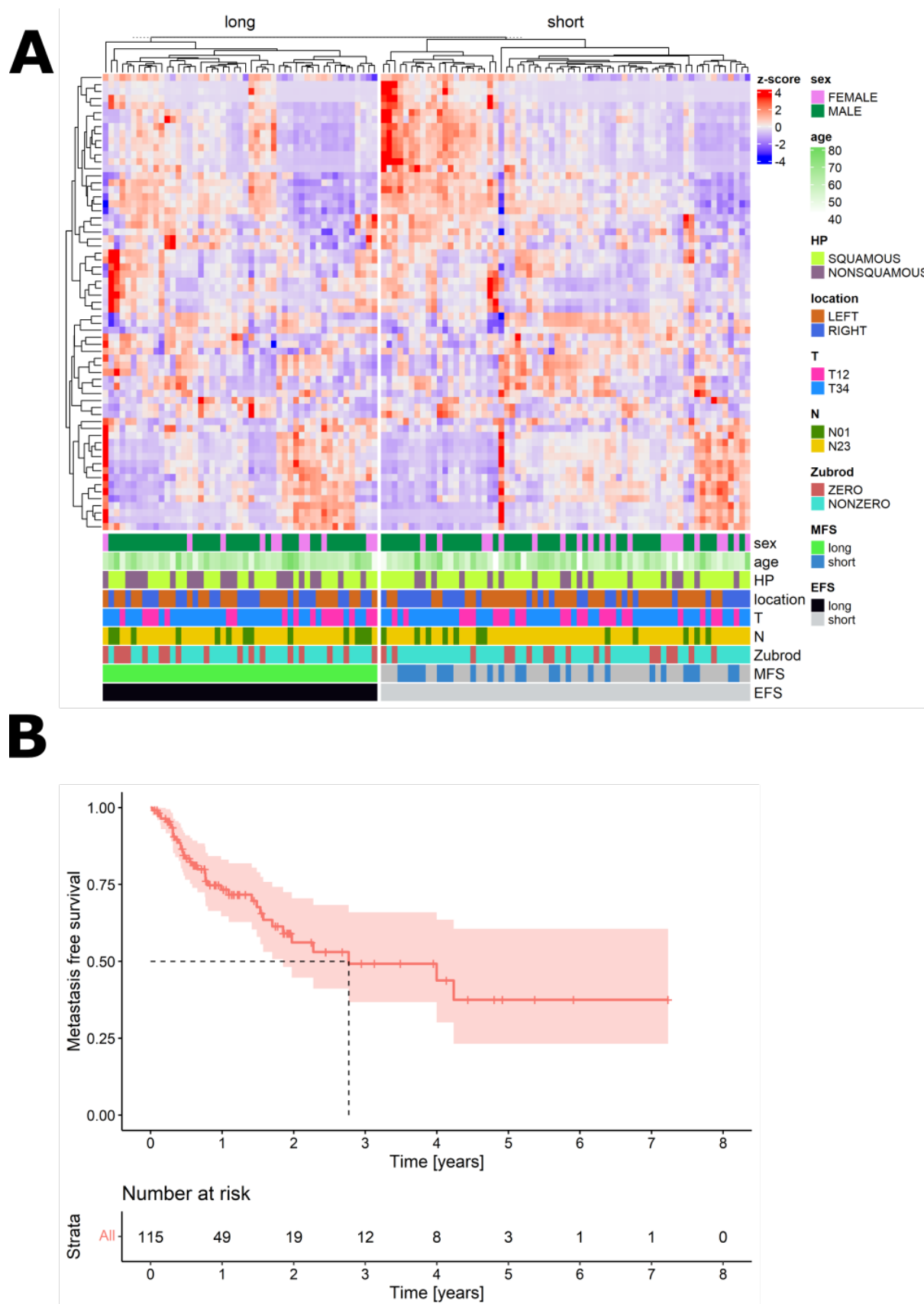
clinical features

Patient ID	Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Feature 6	Feature 7	Feature 8	Feature 9
ID1	0.01	12	11	23.4	12	0.01	12	0.01	0.01
ID2	0.01	23	1	12.6	23	0.01	23	0.01	0.01
ID3	0.03	21	56	56.2	21	0.03	21	0.03	0.03
ID4	0.00	11	100	11.1	11	0.00	11	0.00	0.00
...	0.04	65	1	67.0	65	0.04	65	0.04	0.04
IDM	0.08	11	0	12.4	11	0.08	11	0.08	0.08



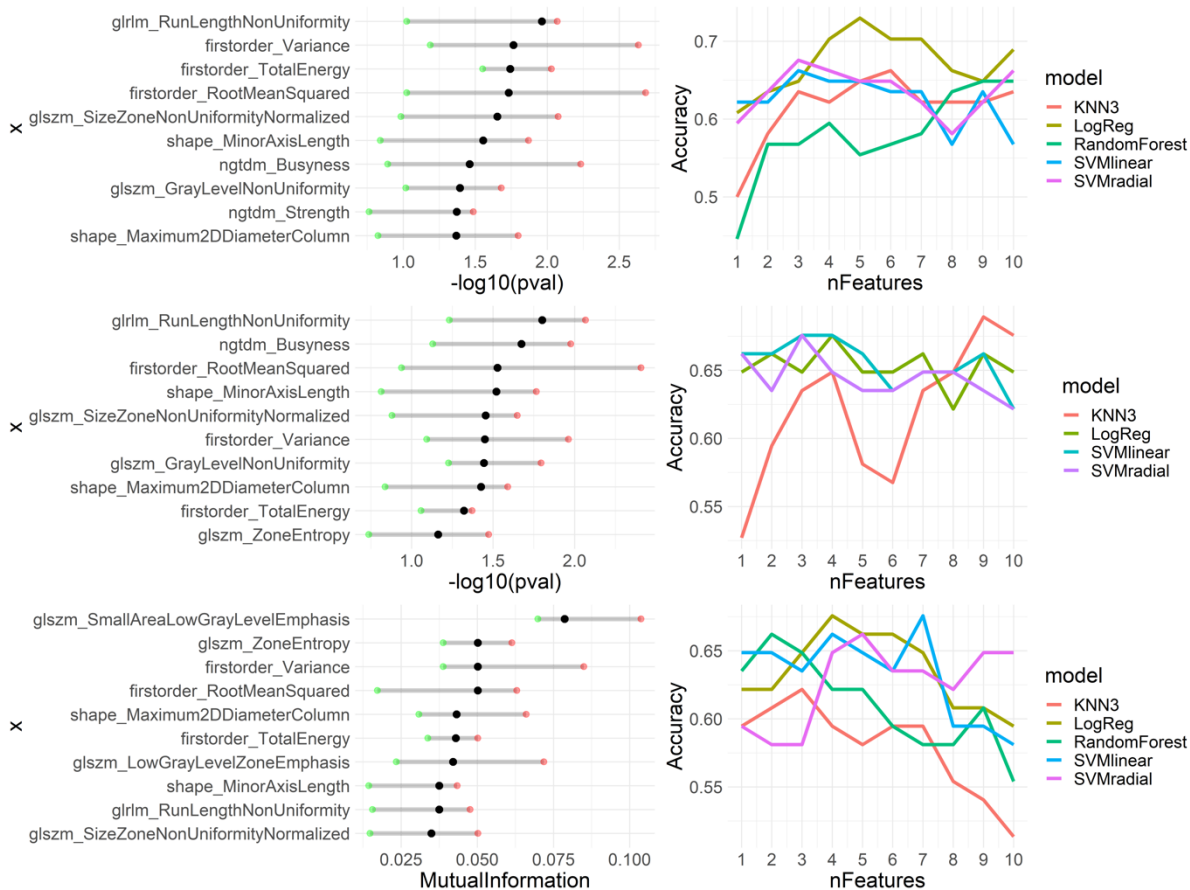
566

567 **Figure 1:** Project workflow. Positron emission tomography/computed tomography (PET/CT)  
 568 images were acquired and radiomic features were extracted from regions of interest (ROI). In-  
 569 tegration of clinical and radiomic data led to the prediction of short-term and long-term metas-  
 570 tasis-free survival (MFS) and the risk of metastasis. The output from the workflow was a radi-  
 571 omic signature, which could be used for the prediction of metastasis risk in newly diagnosed  
 572 NSCLC patients being treated with platinum-based chemotherapy.



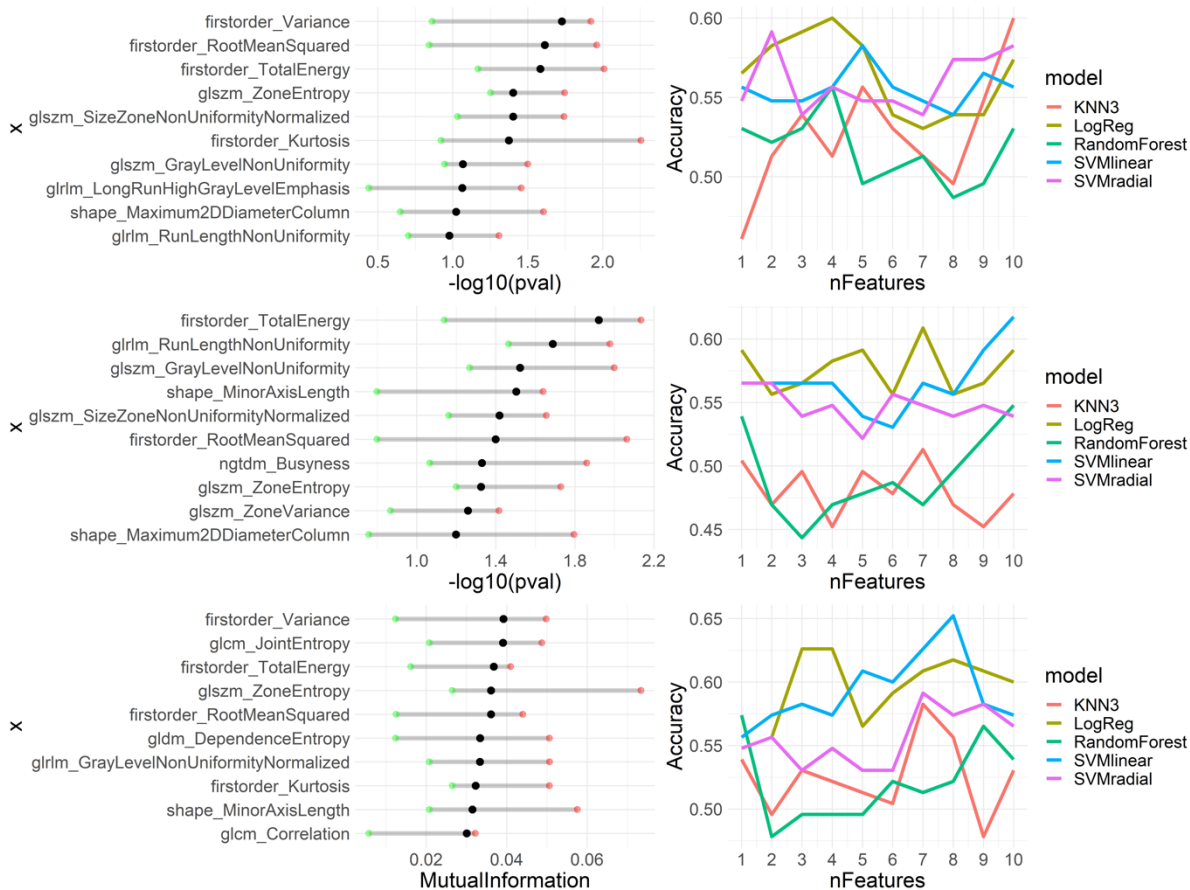
573

574 **Figure 2:** The integration of clinical and radiomic data. A. Kaplan-Meier plot of metastasis-  
 575 free survival (MFS) for the entire population. B. Integration of clinical and radiomic data. Pa-  
 576 tients were split by the binary event-free survival (EFS).



577

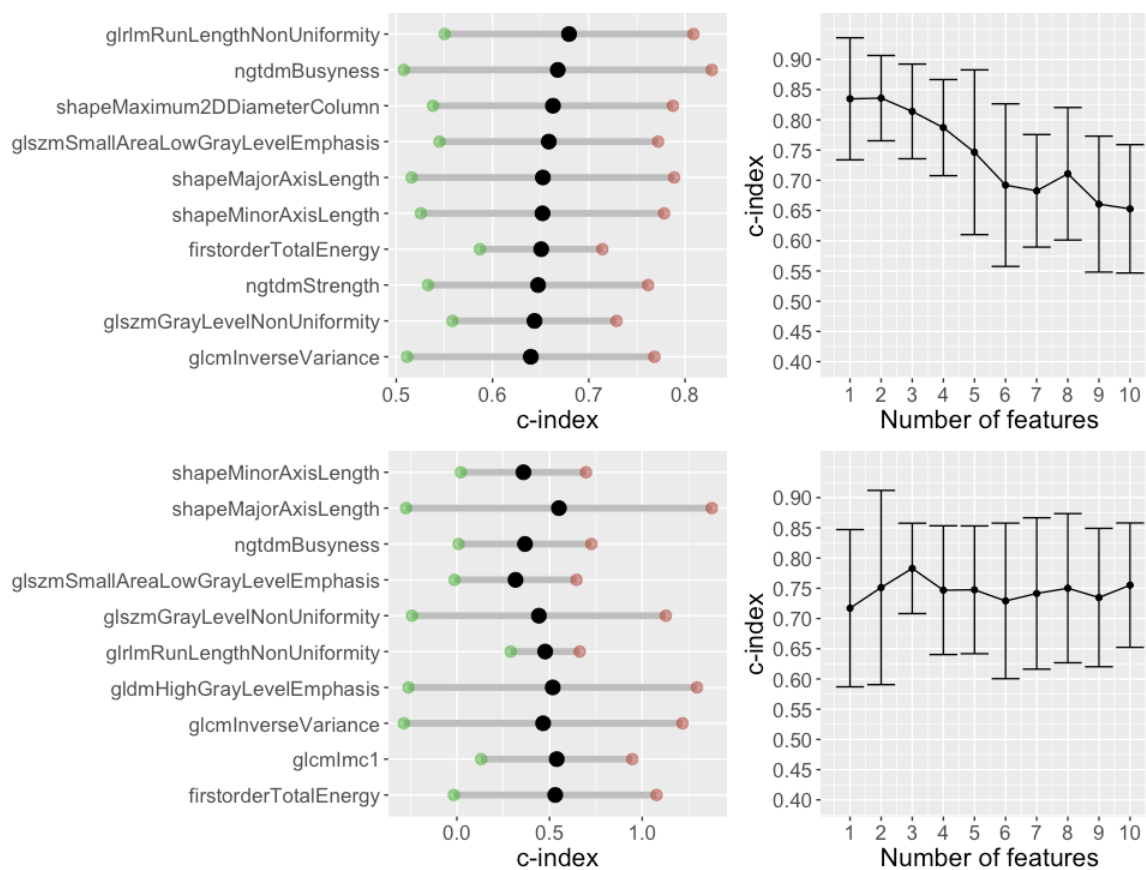
578 **Figure 3:** Metastasis-free survival (MFS) prediction using the classification approach. Top  
 579 row: Wilcoxon test; middle row: Student's t-test; bottom row: mutual information test. Left  
 580 column: feature selection in a 5-fold cross-validation. Features were ranked according to the  
 581  $-\log_{10}(p\text{-value})$  for the Wilcoxon test and Student's t-test selections, and mutual information  
 582 score for mutual information selection. Black dots indicate the median value across folds, green  
 583 dots indicate the lowest value across folds, and red dots indicate the highest value. Right col-  
 584 umn: classification results for the test set in a 5-fold cross-validation for different models, de-  
 585 pending on the number of features.



586

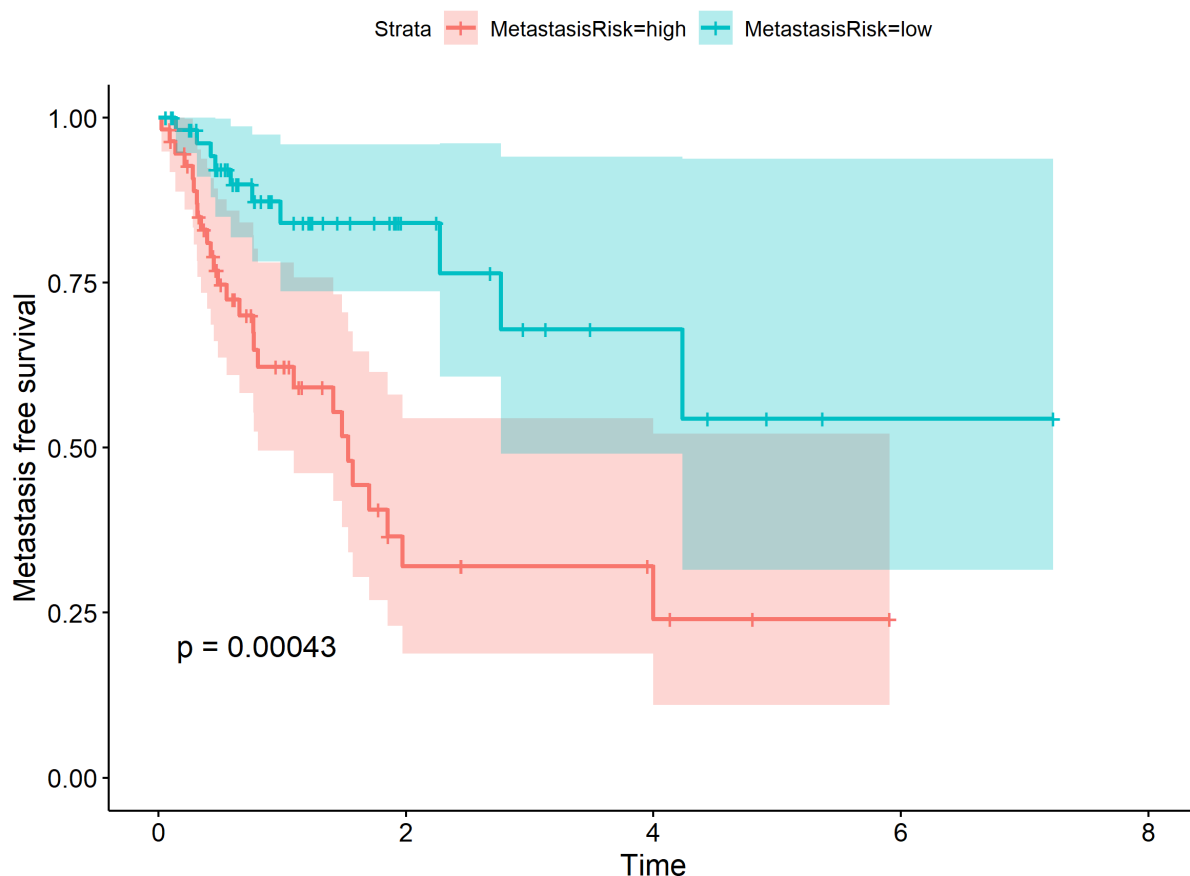
587 **Figure 4:** Event-free survival (EFS) prediction using the classification approach. Top row: Wil-  
 588 coxon test; middle row: Student's t-test; bottom row: mutual information. Left column: feature  
 589 selection in a 5-fold cross-validation. Features are ranked according to the  $-\log_{10}(p\text{-value})$  for  
 590 the Wilcoxon test and Student's t-test selections, and mutual information score for mutual in-  
 591 formation selection. Black dots indicate the median value across folds, green dots indicate the  
 592 lowest value across folds, and red dots indicate the highest value across folds. Right column:  
 593 classification results for the test set in a 5-fold cross-validation for different models, depending  
 594 on the number of features.





595

596 **Figure 5:** Metastasis-free survival (MFS) prediction using a regression approach. Top row: Cox  
 597 regression, bottom row: random survival forest. Left column: feature selection in a 5-fold  
 598 cross-validation. Features were ranked according to the concordance index value for the univariate  
 599 model. Right column: Prediction results for the test set in a 5-fold cross-validation, depending  
 600 on the number of features.



601  
602 **Figure 6:** Kaplan-Meier plot of metastasis-free survival for the whole cohort. The patients were  
603 divided into high-risk and low-risk groups according to a Cox model constructed using the  
604 feature selection and number for which the cross-validation accuracy was highest.

## 605 **Supplementary appendix**

606 **Supplementary Table 1.** Summary of the correlation-based feature filtering.

607  
608 **Supplementary Table 2.** Potential of the data to predict EFS (clinical variables). Fisher's ex-  
609 act test was applied for p-value estimation.

610  
611 **Supplementary Table 3.** Potential of the data to predict MFS and EFS. For binary MFS and  
612 EFS, the Mann-Whitney U test was used. For the continuous MFS log-rank test. Variables  
613 statistically significant against continuous MFS are highlighted in bold.

614  
615 **Supplementary Table 4.** Log-rank test for continuous MFS.

616  
617 **Supplementary Figure 1.** Principal component analysis of the radiomic features (after corre-  
618 lation filtering). Colors correspond to the PET/CT scanner. There is no visible grouping of  
619 samples according to the scanner.

620  
621 **Supplementary Figure 2.** Correlation between radiomic features.

622

623 **Supplementary Figure 3.** Kaplan-Meier plot for metastatic-free survival with high/low  
624 SmallAreaLowGrayLevelEmphasis value

625

626 **Supplementary Figure 4.** Kaplan-Meier plot for metastatic-free survival with high/low Run-  
627 LengthNonUniformity value

628