

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

Reproducibility of tumor uptake heterogeneity characterization through textural feature analysis in ¹⁸F-FDG PET imaging

Florent TIXIER^{1,2}, Mathieu Hatt¹, Catherine Cheze Le Rest¹, Adrien Le Pogam¹,
Laurent Corcos², Dimitris Visvikis¹

¹ INSERM, U650, LaTIM, CHRU Morvan, Brest F-29200, France

² INSERM, U613, Faculty of Medicine, Brest F-29200, France

Corresponding author:

TIXIER Florent
LaTIM, INSERM U650,
CHRU Morvan, 5 avenue Foch,
29609 Brest Cedex, France
Tél: +33 (0)298-018-111
Fax: +33 (0)298-018-124
Email: florent.tixier@univ-brest.fr

Short running title: PET heterogeneities' reproducibility

31 **ABSTRACT**

32 ¹⁸F-FDG PET measurement of standardized uptake values (SUV) is increasingly used for
33 monitoring therapy response or predicting outcome. Alternative parameters computed
34 through textural analysis were recently proposed to quantify the tumor tracer uptake
35 heterogeneity as significant predictors of response. The primary objective of this study was
36 the evaluation of the reproducibility of these heterogeneity measurements. **Methods:**
37 Double-baseline ¹⁸F-FDG PET scans of 16 patients acquired within a period of 4 days prior to
38 any treatment were considered. A Bland-Altman analysis was carried out on six parameters
39 based on histogram measurements and 17 heterogeneity parameters based on textural
40 features obtained after discretization with values between 8 and 128. **Results:** SUV_{max} and
41 SUV_{mean} reproducibility were similar to previously reported studies with a mean percentage
42 difference of 4.7±19.5% and 5.5±21.2% respectively. By comparison better reproducibility
43 was measured for some of the textural features describing tumor tracer local heterogeneity,
44 such as entropy and homogeneity with a mean percentage difference of -2±5.4% and
45 1.8±11.5% respectively. Several of the tumor regional heterogeneity parameters such as the
46 variability in the intensity and size of homogeneous tumor activity distribution regions had
47 similar reproducibility to the SUV measurements with 95% confidence intervals of -22.5% to
48 3.1% and -1.1% to 23.5% respectively. These parameters were largely insensitive to the
49 discretization range values. **Conclusion:** Several of the parameters derived from textural
50 analysis describing tumor tracer heterogeneity at local and regional scales had similar or
51 better reproducibility as simple SUV measurements. These reproducibility results suggest
52 that these FDG PET image derived parameters which have already been shown to have a
53 predictive and prognostic value in certain cancer models, may be used within the context of
54 therapy response monitoring or predicting patient outcome.

55

56

57 Keywords: PET, ¹⁸F-FDG, texture analysis, reproducibility, oncology

58 INTRODUCTION

59 ^{18}F -FDG PET imaging is well established in clinical practice for diagnosis and staging. On the
60 other hand there is increasing interest in the use of this imaging modality within the context of
61 therapy response assessment or patient follow-up. For such applications, standardized
62 uptake value (SUVs) measurements are used, with the maximum of tumor activity
63 concentration (SUV_{max}) being the most popular since it is the easiest to obtain. The use of the
64 mean obtained in an 1cm^3 sphere centered on the voxel of maximum activity concentration
65 (SUV_{peak} (1)), has been proposed as an alternative since it should be more robust to noise
66 compared to SUV_{max} , remaining at the same time easy to derive. Additional PET image
67 derived parameters allowing a more complete lesion characterization include the mean SUV
68 (SUV_{mean}), the metabolically active tumor volume (MATV, defined as the tumor volume that
69 can be seen and delineated on a PET image) and the total lesion glycolysis (TLG, defined as
70 the product of MATV and its associated SUV_{mean}), although they all require an accurate
71 delineation of the functional tumor volume. Different studies have in the past explored the
72 role of such PET image derived parameters for assessing response to therapy (2-6). More
73 recently tracer uptake heterogeneity characterization based on textural analysis extracted
74 from PET images has been also proposed, allowing an improved predictive and prognostic
75 value to be derived from baseline PET scans (7,8).

76 Most frequently monitoring response to therapy involves a comparison of such PET image
77 derived parameters between a baseline PET scan and a second scan carried out early or
78 late during treatment, or after the end of treatment. In this case the variation of the
79 parameters between the two scans is used to characterize response (1). Whether
80 considering the % difference of PET image derived parameters between successive scans or
81 the absolute values on a baseline scan the definition of thresholds in order to identify
82 response or progressive disease requires, amongst others, an evaluation of the physiological
83 reproducibility that characterizes them. Such evaluations are performed on double baseline
84 scans acquired before any treatment within a few days interval from each other.

85 Until now only few studies have investigated the physiological reproducibility of such
86 measurements, almost exclusively focusing on SUVs (9-11), and more recently on the MATV
87 computed using different segmentation algorithms (12,13). Other authors have demonstrated
88 the sensitivity of several textural feature parameters to PET acquisition and reconstruction
89 settings (14), demonstrating the need for standardization in order for such image derived
90 parameters to be used in therapy response assessment studies. However, the physiological
91 reproducibility of these promising parameters extracted from the analysis of tumor activity
92 distributions has never been investigated. The objective of our study was therefore to
93 evaluate the reproducibility of textural features quantifying in a local, regional and global
94 fashion the tumor tracer uptake heterogeneities, thereby identifying the potential of these
95 parameters to be used for therapy response monitoring purposes. A comparison with the
96 physiological reproducibility of SUVs using the same patient datasets was also performed
97 since they are the most used parameters in current clinical practice and in order to facilitate a
98 direct comparison with previous reproducibility studies.

99

100 **MATERIALS AND METHODS**

101 **Patients**

102 16 patients with newly diagnosed esophageal cancer were enrolled in this study. All of these
103 patients underwent two ^{18}F -FDG PET baseline scans before initiating any treatment. The two
104 scans were obtained within 2-7 days (median 4.2 days). PET images were acquired on a
105 PET/CT scanner (Gemini; Philips), with 2-min acquisitions per bed position, 60 min after the
106 injection of 6MBq/kg of ^{18}F -FDG. Data were reconstructed using a 3D row-action
107 maximization-likelihood algorithm (RAMLA (15)) with standard clinical protocol parameters (2
108 iterations, relaxations parameter of 0.05, and 5mm full width at half maximum 3D Gaussian
109 post-filtering). This analysis was carried out after obtaining the approval of the local
110 Institutional Ethics Review Board.

111

112 **Tumor Analysis**

113 The primary lesions of each patient were delineated with the Fuzzy Locally Adaptive
114 Bayesian (FLAB) algorithm which has been previously demonstrated to provide reproducible
115 MATV automatic delineations (mean difference between baseline scans of $5\pm 13\%$) (16).
116 SUV_{max} and mean SUV within the delineated tumor (SUV_{mean}) were extracted from the
117 primary tumor in each of the two baseline PET images for each patient. In addition, a number
118 of tumor heterogeneity parameters shown in table 1, whose value for prognosis and
119 prediction of outcome and treatment response on FDG PET images has been previously
120 investigated (7,8), were calculated based on the delineated 3D functional volumes.

121

122 **Textural Analysis**

123 We define texture as a spatial arrangement of a predefined number of voxels allowing the
124 extraction of complex image properties and we define a textural feature as a measurement
125 computed using a texture matrix (8). Given that these features quantify the spatial
126 relationship between voxels and their relative intensities, they can be associated to tracer
127 heterogeneity patterns within the functional volume of the tumor at different scales, namely
128 local and regional (using texture matrices) or global (using image-voxel-intensity histograms).
129 The first type of matrices is used to quantify local heterogeneity as they allow
130 characterization of the intensity variations between consecutive voxels. On the other hand,
131 the second type of matrices allows characterization of arrangements of larger homogeneous
132 areas (groups of voxels) within the tumors therefore providing information on tumor regional
133 heterogeneity.

134 Local heterogeneity parameters were derived using the co-occurrences matrices (17) and
135 were computed by considering a 26-connectivity (i.e. neighboring voxels in all 13 directions in

136 three dimensions) and a 1-distance (i.e. no gap) relationship between consecutive voxels.
137 On these matrices, 6 different parameters characterizing the local heterogeneity were
138 calculated by averaging the values on the 13 directions for each feature. The other type of
139 texture matrices is called intensity size-zone matrix (8, 18) and is constructed in two steps.
140 First, homogenous areas are identified within the tumor and a matrix linking the size of each
141 of these homogeneous areas to its intensity is constructed. 11 features characterizing the
142 regional heterogeneity were calculated from this matrix. For example, parameters can
143 quantify the presence of large areas with high intensity (HILAE) or small areas with a low
144 intensity (LISAE).

145 Other features characterizing regional heterogeneity include the variability in the size (SZV)
146 and the intensity (IV) of identified homogeneous tumor zones, as well as the ratio between
147 the number of homogeneous tumor zones and the overall tumor size (known as the zone-
148 percentage (ZP)). Regional heterogeneity formulae were summarized in table 2 and the
149 mathematical definition of all local features used in this study have been previously
150 summarized in Haralick et al (17). A complete list of texture matrices and their associated
151 features used in this work are included in table 1.

152 Building texture matrices on which the textural features are computed require a discretization
153 of the voxel values within the previously delineated MATV on a specific range of values. This
154 range has to be chosen as a power of two due to algorithmic constraints and in this study the
155 features were extracted by considering downsampling to ranges of 8, 16, 32, 64 and 128
156 distinct values. Figure 1 illustrates on a transaxial tumor slice the resulting resampled MATV
157 for each of these discretization ranges. This necessary downsampling step on the one hand
158 reduces image noise while on the other normalizes the tumor voxel intensities across
159 patients, subsequently facilitating the comparison of the extracted textural features. In a
160 previous study (8) there were no statistically significant differences shown in the extracted
161 textural feature values as a result of varying the number of discrete values in this resampling
162 normalization process. 64 discrete values were considered sufficient for a range of SUVs

163 between 4 and 20. In the present study the influence of this parameter in the physiological
164 reproducibility of the textural feature parameters was also assessed.

165

166 **Statistical Analysis**

167 The reproducibility of the quantitative values (q) for each parameter under investigation was
168 assessed by calculating the mean percentage difference relative to the mean of both
169 baseline scans using the following formula:

$$170 \quad \Delta = \frac{(q_1 - q_2)}{(q_1 + q_2)/2} \cdot 100 \quad Eq. 1$$

171 This analysis was performed for all parameters and in the case of the textural features for all
172 discretization values (from 8 to 128). A Kolmogorov-Smirnov test was first performed to verify
173 the normality of the distribution of Δ . Bland-Altman analysis (19) was subsequently used to
174 evaluate the differences for the image derived parameters considered. The mean and
175 standard deviation (SD) and the associated 95% confidence intervals (CI) were obtained.
176 Lower and upper reproducibility limits (LRL and URL), defining the reference range of
177 spontaneous changes, were calculated as $\pm 1.96 \times SD$ provided that the distribution were not
178 statistically different than a normal one. Intraclass correlation coefficients (ICC) were in
179 addition calculated providing an evaluation of the reliability of measurements, whereas their
180 reproducibility was estimated based on their precision (half the width of 95%CI * 100 %). The
181 differences in the calculated reproducibility of the textural feature parameters as a function of
182 the discretization values used in the normalization step was assessed using a paired student
183 t-test. P values of less than 0.05 were considered statistically significant.

184

185 **RESULTS**

186 For all considered features, Δ showed no significant differences from a normal distribution
187 according to the Kolmogorov-Smirnov test. Consequently, Bland-Altman analysis was
188 performed on all parameters. All of the reproducibility results using the Bland-Altman
189 analysis, including LRL and URL (and associated 95% CI), are provided in table 3 for both
190 intensity histogram parameters and textural features, whereas the ICCs and associated 95%
191 CI and precision are summarized in table 4. As figure 2A and Table 3 show SUV
192 measurements exhibited reproducibility levels in line with previously published studies. A
193 mean difference of $5\pm 20\%$ and associated LRL and URL of -34% and $+43\%$ were found for
194 SUV_{max} , and $6\pm 21\%$ mean difference, with -36% LRL and $+47\%$ URL for SUV_{mean} . ICC was
195 0.94 (95% CI: $0.82-0.98$; precision $\pm 8\%$) and 0.92 (95% CI: $0.78-0.97$; precision $\pm 10\%$) for
196 SUV_{max} and SUV_{mean} respectively. Amongst other global tumor heterogeneity characterization
197 parameters derived using the intensity histogram, kurtosis was found to have similar
198 reproducibility as SUV_{max} and SUV_{mean} but a lower ICC (0.80 with 95% CI between $0.44-0.93$;
199 precision $\pm 25\%$; figure 2B). COV (Mean/SD) was characterized by reproducibility limits
200 ranging between -43% and 51% and an ICC of 0.82 (95%CI: $0.49-0.94$; precision $\pm 23\%$).
201 Standard deviation, skewness and minimum intensity had the highest reproducibility limits
202 ranging between -45 and 60% .

203 Among the local heterogeneity parameters calculated on co-occurrence matrices, the
204 entropy, homogeneity and dissimilarity were characterized by reproducibility limits below
205 30% and an ICC precision below $\pm 16\%$, the most reproducible being the entropy, with LRL of
206 -13% and URL of 9% (figure 2C). The other local features (2^{nd} angular moment, contrast and
207 correlation) were characterized by lower reproducibility, with LRL and URL varying between -
208 40.9% and 62.7% , which is comparable with the reproducibility achieved for some of the
209 histogram based parameters such as skewness (LRL-URL between -54.2% and 53.6%) or
210 minimum intensity (LRL-URL between -45.6% and 58.2%). Both the intensity and the size
211 variability of uniform zones identified within the tumor, representing a measure of regional
212 tumor heterogeneity and previously shown as significant predictors of response to therapy,

213 have shown a better physiological reproducibility with LRL and URL of -56.7% to 37.3% and -
214 34.1% to 56.5% respectively (figure 2D). The respective ICCs for these measurements were
215 0.97 (95%CI: 0.93-0.99; precision $\pm 3\%$) and 0.97 (95%CI: 0.91-0.99; precision $\pm 4\%$). More
216 specifically the SD of the mean percentage difference was 23.1% and 24% for the textural
217 feature parameters related to the size and intensity variability of tumor uniform zones
218 compared to 19.5% and 21.2% in the case of the SUV_{max} and SUV_{mean} respectively. Other
219 regional heterogeneity features were not reproducible, as for example small area emphasis
220 (LRL and URL of -113% and +100%), low-intensity emphasis (LRL and URL of -112% to
221 +104%) and low-intensity small area emphasis (LRL and URL of -140% to +125%).

222 As illustrated in figure 3A, all of the textural parameters describing local tumor heterogeneity
223 were found to be insensitive to the chosen discretization values. Within this context no
224 statistically significant differences were found for the range of discretization values used (8 to
225 128) with a mean SD of 5% and 15% for 8 and 128 discretization values respectively.
226 Several of the regional heterogeneity parameters calculated on intensity size-zone matrices
227 were sensitive to the chosen discretization value, with statistically significant differences and
228 SD values twice as high or low with varying discretization, as shown in figure 3B. The large
229 area emphasis feature, for instance, was characterized by a mean difference of $29 \pm 79\%$ and
230 $4 \pm 30\%$ using 8 and 64 values respectively. On the other hand, the intensity and size
231 variability of uniform tumor areas as well as the high intensity emphasis zones were largely
232 independent (SD differences $< 20\%$) of the discretization values with non-statistically
233 significant differences.

234

235 **DISCUSSION**

236 Predicting and monitoring therapy response with PET imaging is one of the rising
237 applications of this modality. Characterizing intra-tumor heterogeneity of the radiotracer
238 uptake has been identified as a clinically relevant task and requires semi-automatic

239 validated, accurate, robust and reproducible tools (20). We have recently introduced the use
240 of textural features for the characterization of tumor heterogeneity within the context of
241 predicting tumor response to therapy using FDG PET imaging (8). It is clearly not
242 straightforward to associate each of these heterogeneity features with one specific
243 physiological process within the tumor, particularly in the case of FDG imaging. However,
244 since all these different parameters represent measurements of tumor local and regional
245 tracer uptake heterogeneity, a reasonable assumption is that their quantitation can be related
246 to underlying physiological processes, such as vascularization, perfusion, tumor
247 aggressiveness, or hypoxia (21, 22). All of these processes have been identified as
248 potentially contributing to the way the FDG uptake is spatially distributed within a tumor
249 volume.

250 A possible clinical significance of tumor uptake heterogeneity patterns can be related to the
251 efficiency of a given treatment regime. One example is in the case of combined chemo-
252 radiotherapy, where the delivery of a uniform radiation dose to a target tumor volume
253 independently of the actual tracer distribution within the tumor may be responsible for
254 possibly explaining failure of treatment (8, 20) Finer characterization of the heterogeneity as
255 obtained through textural features could therefore help identifying potential responders or
256 non responders before initiating treatment or early during treatment by characterizing the
257 evolution of uptake heterogeneity during treatment.

258 As the features are calculated within a delineated MATV, it is important to reduce the
259 potential variability that could arise from the reproducibility of the tumor volume delineation
260 step. There is indeed a large variability in the reproducibility results observed depending on
261 the segmentation algorithm used. It has been demonstrated that threshold-based delineation
262 may lead to poorly reproducible delineated MATV on double baseline scans (12,13). On the
263 other hand, the use of more sophisticated and robust segmentation algorithms (such as
264 FLAB) has been demonstrated to lead to satisfactory results with similar reproducibility as

265 SUV_{max} ($\pm 30\%$) (13). This delineation method was therefore used in this study in order to
266 minimize the impact of MATV delineation to the textural features reproducibility.

267 The parameters extracted from the intensity histogram characterize the distribution of the
268 voxel intensities without taking into consideration spatial relationships between the voxels.
269 For this reason, the features extracted from the histogram can be denoted as global. The
270 maximum intensity of the histogram, corresponding to the SUV_{max}, had the best
271 reproducibility along with kurtosis and mean SUV with a SD of the mean percentage
272 difference of 19.5%, 18% and 21.2% with an ICC of 0.94, 0.80 and 0.92 respectively. These
273 reproducibility results are similar to these reported on previous reproducibility studies
274 concerning the SUVs measurements. The reproducibility for the other tumor global features,
275 namely the minimum intensity, standard deviation and skewness, was worse with LRL and
276 URL at -54% to 58%, which may compromise their potential for clinical use in order to
277 characterize tumor response or progression.

278 The local heterogeneity features derived from co-occurrence matrices provide far more
279 complex information than the intensity histogram as they are focusing on the relationship
280 between voxels and their neighbors at a local scale. Despite this characteristic of being very
281 specific and local parameters, some of these features (entropy, local homogeneity) exhibited
282 even better reproducibility than the SUV_{max}. These tumor local heterogeneity features were
283 previously identified amongst other tumor heterogeneity characteristics as being capable of
284 classifying esophageal cancer patients with high specificity and sensitivity regarding
285 response to combined radiochemotherapy. On the other hand, other local heterogeneity
286 features such as contrast, 2nd angular moment or correlation were characterized by larger
287 reproducibility limits between -40% and 63% (ICC ≥ 0.94). Finally most of the local
288 heterogeneity parameters were found to be robust versus changes in the discretization
289 value.

290 Regarding regional heterogeneity features, several parameters (SAE, LAE, LIE, LISAE,
291 LILAE, HILAE and ZP) were found to be sensitive to the choice of the discretization value.
292 Some of them (particularly SAE, LIE and LISAE) were also found to have poor
293 reproducibility. All of these parameters are focusing on the smaller homogenous and lower
294 intensity regions, which on the one hand are expected to be less reproducible and on the
295 other hand not of the highest interest in terms of characterizing regional tumor FDG uptake
296 heterogeneities. Other regional heterogeneity parameters such as the features characterizing
297 large homogeneous and high intensity tumor regions (LAE, HIE, HILAE) may be more
298 interesting for predicting response to therapy. The high intensity areas, corresponding to high
299 radiotracer uptake regions, are associated to the more aggressive tumor parts. On the other
300 hand, the large homogeneous areas represent more robust tumor characteristics since they
301 are less likely to result from statistical noise or partial volume effects. Among these regional
302 heterogeneity parameters, only the high intensity regions feature exhibit a reproducibility
303 similar to the SUV_{max} (LRL -36% to URL +44%, ICC 0.82), and therefore sufficient to be
304 considered as a parameter of interest for characterizing patient response.

305 Finally, the parameters corresponding to the variability in the size or intensity (SZV and IV
306 respectively) of the homogeneous areas are also good indicators of the regional tumor
307 heterogeneity having already shown potential for patient differentiation in terms of response
308 to therapy. These parameters highlight the repartition of the intensity values or region sizes
309 within the tumor (high tumor heterogeneity corresponding to high variability of the radiotracer
310 distribution, corresponding in turn to high intensity variability). A good reproducibility with a
311 SD of the mean percentage difference of 24% and an ICC of 0.97 (compared to 19.5% for
312 the SUV_{max}) was measured for these regional heterogeneity features.

313 Our study suggests that a careful selection of the parameters to quantify local and regional
314 heterogeneity may provide both a complete and reproducible characterization of the tracer
315 uptake spatial heterogeneity within tumors in FDG PET images. It should be emphasized that
316 these parameters exhibiting the highest reproducibility in this study were also the ones that

317 were found to be significant predictors of patient response in a previous study (local
318 homogeneity and entropy, intensity variability and size-zone variability) (8).

319 One of the limitations of the current study is the small sample of patients, which is however of
320 the same size and in line with previously published reproducibility studies (9-11). On the
321 other hand, although our reproducibility results were established on FDG PET images of
322 esophageal cancer lesions, these lesions displayed a large range of sizes and tracer uptake
323 heterogeneity patterns. These results obviously require confirmation for other cancer models
324 and/or radiotracers. Partial volume effects (PVE) were not specifically investigated in this
325 work, although since tumors were all larger than 10cm^3 and in the same body region, PVE is
326 expected to have a low impact on an inter-patient basis for this dataset as far as the
327 reproducibility evaluation is concerned. On the other hand, PVE correction can be expected
328 to have a potentially more important role on the absolute quantification of the heterogeneity
329 parameters, and therefore the impact of partial volume effects correction within this context
330 will be the focus of further investigations.

331 Finally, in this study we assumed that a satisfactory reproducibility range for textural features
332 could be considered as $\sim\pm 30\text{-}40\%$ (SD of $15\text{-}20\%$) upper and lower limits. This was chosen
333 accordingly to what was previously defined as reproducibility limits for the use of SUV and
334 tumor metabolic volume measurements. This means that in order to be used for response
335 monitoring purposes, a given parameter has to exhibit higher changes during treatment than
336 its reproducibility range observed in double baseline scans. However, no study has yet to
337 investigate the evolution of textural features on sequential PET scans and the correlation of
338 these changes with therapy response. Such a study will provide an estimation of the range of
339 changes for these parameters between a pre- and post- or early into treatment scans. This
340 range of values, in comparison with the reproducibility limits of the same parameters as
341 established in the present study, would allow evaluating the potential of using these
342 heterogeneity measures within the context of assessing response to therapy with serial FDG
343 PET scans.

344 **CONCLUSIONS**

345 The physiological reproducibility varied significantly among the various tumor heterogeneity
346 features under investigation, only a few of them being identified as reproducible. Based on
347 our results, heterogeneity parameters that should be preferentially considered for tumor
348 heterogeneity characterization since they are the most reproducible include entropy,
349 homogeneity and dissimilarity for local characterization, and variability in the size and
350 intensity of homogeneous tumor areas for regional characterization.

351

352

353 **References**

- 354 1. Wahl RL, Jacene H, Kasamon Y, Lodge MA. From RECIST to PERCIST: evolving
355 considerations for PET response criteria in solid tumors. *J Nucl Med.* 2009; 50(suppl
356 1):122S-150S.
- 357 2. Cazaentre T, Morschhauser F, Vermandel M et al. Pre-therapy 18F PET quantitative
358 parameters help in predicting the response to radioimmunotherapy in non-Hodgkin
359 lymphoma. *Eur J Nucl Med Mol Imaging.* 2010;37:494-504.
- 360 3. Rizk NP, Tang L, Adusumilli PS, et al. Predictive value of initial PET SUVmax in
361 patients with locally advanced esophageal and gastroesophageal junction adenocarcinoma.
362 *J Thoracic Oncol.* 2009;4:875-879.
- 363 4. Leibold T, Akhurst TJ, Chessin DB, et al. Evaluation of (18)F-FDG-PET for early
364 detection of suboptimal response of rectal cancer to preoperative chemoradiotherapy : a
365 prospective analysis. *Ann Surg Oncol.* 2011;18:2783-2789.
- 366 5. Shamim SA, Kumar R, Shandal V, et al. FDG PET/CT evaluation of treatment
367 response in patients with recurrent colorectal cancer. *Clin Nucl Med.* 2011;36:11-16.
- 368 6. Hatt M, Visvikis D, Albarghach NM, et al. Prognostic value of 18F-FDG PET image-
369 based parameters in oesophageal cancer and Impact of tumour delineation methodology *Eur*
370 *J Nucl Med Mol Imaging.* 2011;38:1191-1202.
- 371 7. El Naqa I, Grigsby P, Apte A, et al. Exploring feature-based approaches in PET images
372 for predicting cancer treatment outcomes. *Pattern Recognit.* 2009;42:1162-1171.
- 373 8. Tixier F, Cheze-Le-Rest C, Hatt M, et al. Intratumor heterogeneity characterized by
374 textural features on baseline 18F-FDG PET images predicts response to concomitant
375 radiochemotherapy in esophageal cancer. *J Nucl Med.* 2011;52:369-378.
- 376 9. Weber WA, Ziegler SI, Thodtmann R, Hanauske AR, Schwaiger M. Reproducibility of
377 metabolic measurements in malignant tumors using FDG PET. *J Nucl Med.* 1999;40:1771-
378 1777.

- 379 10. Nahmias C, Wahl LM. Reproducibility of standardized uptake value measurement
380 determined by 18F-FDG PET in malignant tumors. *J Nucl Med.* 2008;49:1804-1808.
- 381 11. Paquet N, Albert A, Foidart J, Hustinx R. Within patient variability of FDG
382 standardized uptake values in normal tissues. *J Nucl Med.* 2004;45:784-788.
- 383 12. Frings V, de Langen AJ, Smit EF, et al. Repeatability of metabolically active volume
384 measurements with 18F-FDG and 18F-FLT PET in non-small cell lung cancer. *J Nucl Med.*
385 2010;51:1870-1877.
- 386 13. Hatt M, Cheze-Le Rest C, Aboagye EO, et al., Reproducibility of 18F-FDG and 3'-
387 deoxy-3'-18F-fluorothymidine PET tumor volume measurements. *J Nucl Med.* 2010;51:1368-
388 1376.
- 389 14. Galavis PE, Hollensen C, Jallow N, Paliwal B, Jeraj R. Variability of textural features
390 in FDG PET images due to different acquisition modes and reconstruction parameters. *Acta*
391 *Oncol.* 2010;49:1012-1016.
- 392 15. Browne J, de Pierro AB. A row-action alternative to the EM algorithm for maximizing
393 likelihood in emission tomography. *IEEE Trans Med Imaging.* 1996; 15:687-699.
- 394 16. Hatt M, Cheze le Rest C, Turzo A, Roux C, Visvikis D. A fuzzy locally adaptive
395 Bayesian segmentation approach for volume determination in PET. *IEEE Trans Med*
396 *Imaging.* 2009; 28(6):881-893.
- 397 17. Haralick RM, Shanmugam K, Dinstein I. Textural features for image classification.
398 *IEEE Trans Syst Man Cybern.* 1973;3:610-621.
- 399 18. Thibault G, Fertil B, Navarro C, Pereira S. Texture indexes and gray level size zone
400 matrix: application to cell nuclei classification. *Pattern Recognition Inf Process.* 2009;140-
401 145.
- 402 19. Bland JM, Altman DG. Statistical methods for assessing agreement between two
403 methods of clinical measurement. *Lancet.* 1986; 327:307-310.
- 404 20. Basu S, Kwee TC, Gatenby R, et al. Evolving role of molecular imaging with PET in
405 detecting and characterizing heterogeneity of cancer tissue at the primary and metastatic

406 sites, a plausible explanation for failed attempts to cure malignant disorders. *Eur J Nucl Med*
407 *Mol Imaging*. 2011; 38:987-991.

408 21. Rajendran JG, Schwartz DL, O'Sullivan J, et al. Tumour hypoxia imaging with 18F
409 fluoromisonidazole positron emission tomography in head and neck cancer. *Clin Cancer*
410 *Res*. 2006; 12:5435–5441.

411 22. Kunkel M, Reichert TE, Benz P, et al. Overexpression of Glut-1 and increased glucose
412 metabolism in tumours are associated with a poor prognosis in patients with oral squamous
413 cell carcinoma. *Cancer*. 2003; 97:1015–1024.

414

415 **Figure captions**

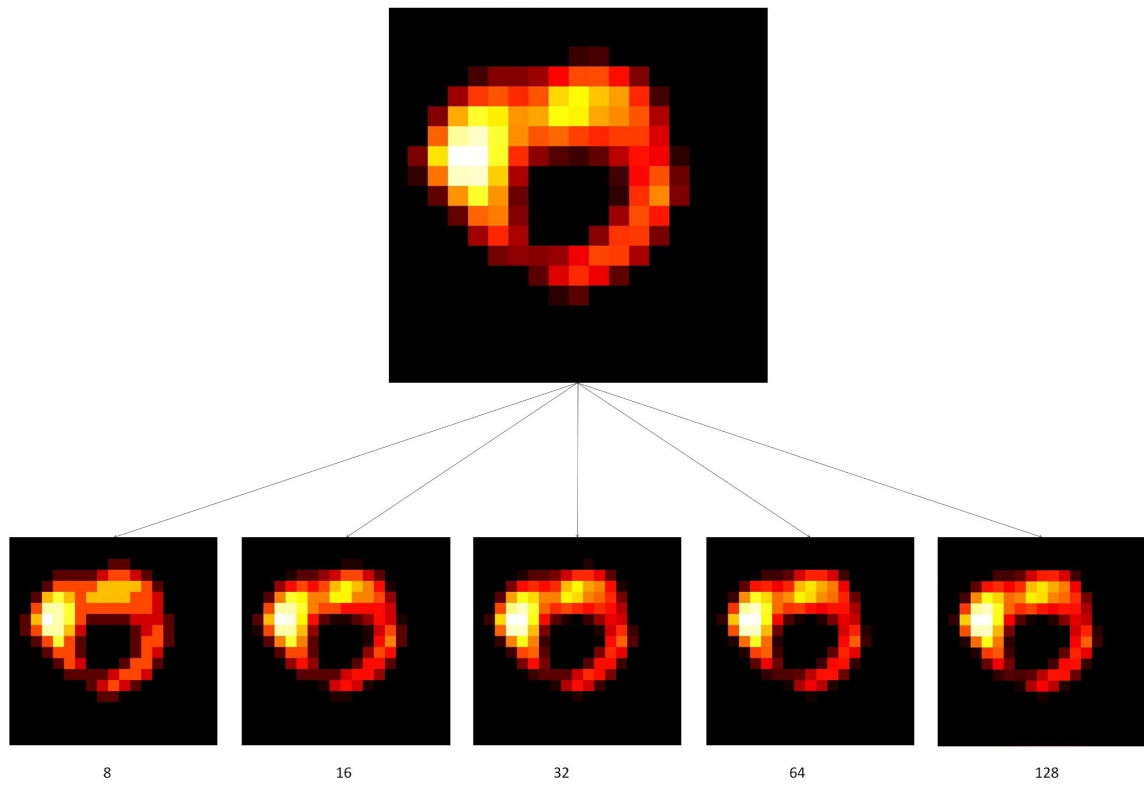
416

417 **Figure 1:** Illustration of one of the tumors considered in this study (sagittal slice) for varying
418 discretization values (from 8 to 128 distinct values).

419 **Figure 2:** Bland-Altman plots of intensity histogram parameters: SUV_{max} (A) and kurtosis (B);
420 as well as textural features heterogeneity parameters: entropy (C) and size-zone variability
421 (D). Lines show combined mean, 95%CI, as well as upper and lower reproducibility limits

422 **Figure 3:** Plots showing the standard deviation of the mean percentage difference as a
423 function of the discretization value for parameters derived from co-occurrences matrices
424 (entropy, dissimilarity, contrast) (A) and intensity size-zone matrices (LISAE: Low-intensity
425 small-area emphasis, SZV: Size-zone variability, ZP: zone percentage) (B).

426



427

8

16

32

64

128

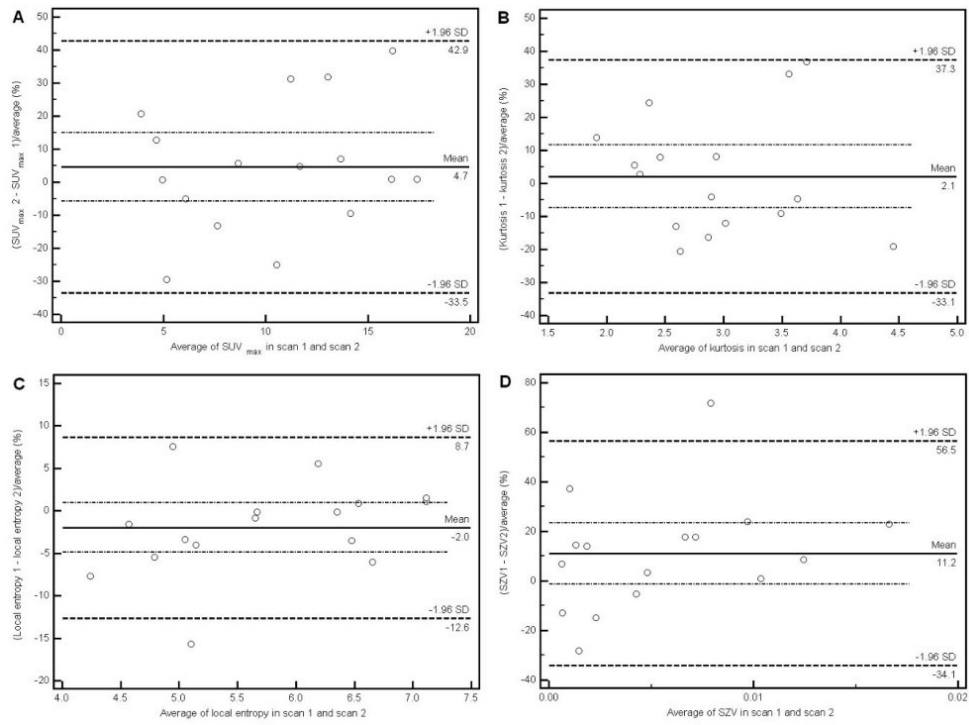
428

429

430

431

Figure 1



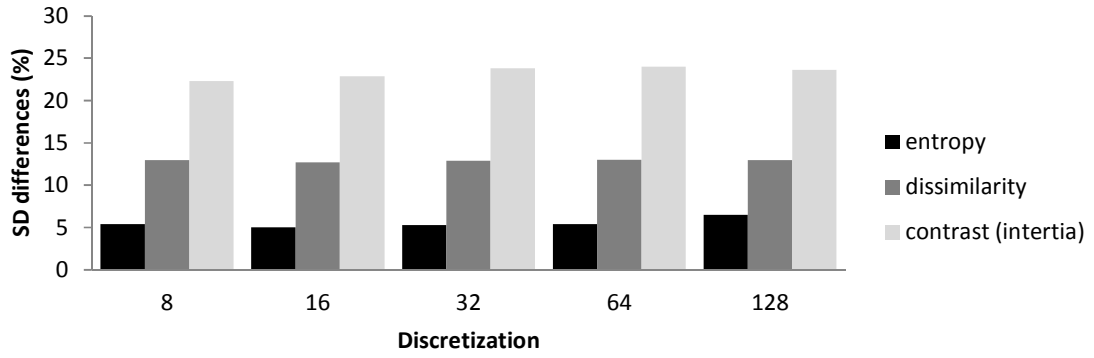
432

433

434

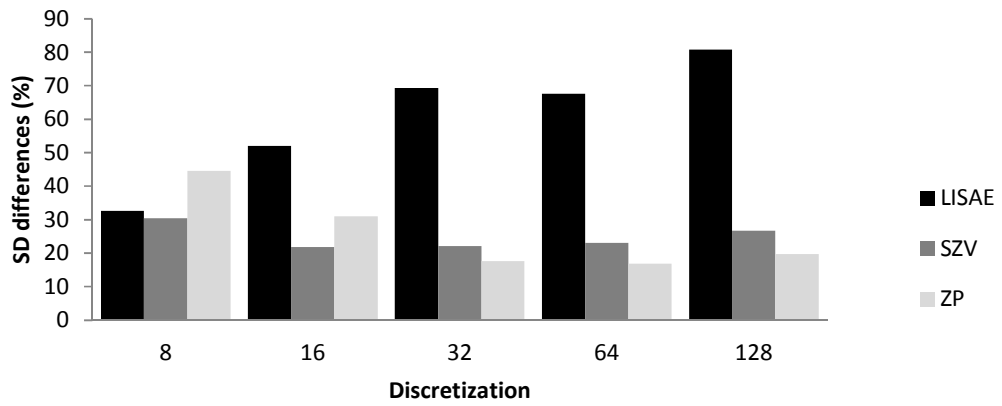
Figure 2

435 **A**



436

437 **B**



438

439

Figure 3

440

| Type | Feature | Scale |
|--|--|----------|
| Features based on intensity histogram | Minimum intensity | Global |
| | Maximum intensity (SUV_{max}) | |
| | Mean intensity (SUV_{mean}) | |
| | Variance | |
| | SD | |
| | Skewness | |
| | Kurtosis Mean/SD | |
| Features based on intensity-size-zone matrix | Small-area emphasis (SAE) | Regional |
| | Large-area emphasis (LAE) | |
| | Intensity variability (IV) | |
| | Size-zone variability (SZV) | |
| | Zone percentage (ZP) | |
| | Low-intensity emphasis (LIE) | |
| | High-intensity emphasis (HIE) | |
| | Low-intensity small-area emphasis (LISAE) | |
| | High-intensity small-area emphasis (HISAE) | |
| | Low-intensity large-area emphasis (LILAE) | |
| High-intensity large-area emphasis (HILAE) | | |
| Features based on co-occurrence matrices | Second angular moment | Local |
| | Contrast (inertia) | |
| | Entropy | |
| | Correlation | |
| | Homogeneity Dissimilarity | |

442

443

Table 1

444

| | |
|---|---|
| $SAE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N \frac{z(i,j)}{j^2}$ $LAE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N j^2 \cdot z(i,j)$ $IV = \frac{1}{\Omega} \sum_{i=1}^M \left[\sum_{j=1}^N \frac{z(i,j)}{i^2} \right]^2$ $SZV = \frac{1}{\Omega} \sum_{i=1}^N \left[\sum_{j=1}^M \frac{z(i,j)}{i^2} \right]^2$ $ZP = \Omega / \sum_{i=1}^M \sum_{j=1}^N j^2 \cdot z(i,j)$ $LIE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N \frac{z(i,j)}{i^2}$ | $HIE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N i^2 \cdot z(i,j)$ $LISAE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N \frac{z(i,j)}{i^2 \cdot j^2}$ $HILAE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N i^2 \cdot j^2 \cdot z(i,j)$ $LILAE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N \frac{j^2 \cdot z(i,j)}{i^2}$ $HILAE = \frac{1}{\Omega} \sum_{i=1}^M \sum_{j=1}^N \frac{i^2 \cdot z(i,j)}{j^2}$ |
|---|---|

Ω : number of homogeneous areas within the tumor

z : intensity size-zone matrix

M : used discretization value

N : size of the largest homogeneous area within the tumor

$z(i,j)$ represents the number of areas with an intensity I and a size j

446

447

448

Table 2

449

450

451 Reproducibility results for all considered image derived parameters, including SUVs and
 452 textural features (calculated using a downsampling range of 64 values).

| Texture | Feature | Mean±SD | 95% CI | LRL | 95% CI for LRL | URL | 95% CI for URL |
|--|--|---------------|---------------|-----------------|-----------------|---------------|-----------------------------------|
| Global | Minimum intensity | 6.3 ± 26.5 | -7.8 to 20.4 | -45.6 | -70.2 to -20.9 | 58.2 | 33.6 to 82.8 |
| | Maximum intensity (SUV _{max}) | 4.7 ± 19.5 | -5.7 to 15.0 | -33.5 | -51.7 to -15.4 | 42.9 | 24.7 to 61.0 |
| | Mean intensity (SUV _{mean}) | 5.5 ± 21.2 | -5.8 to 16.8 | -36.1 | -55.8 to 16.4 | 47.1 | 27.3 to 66.8 |
| | SD | 1.2 ± 23.2 | -11.1 to 13.6 | -44.18 | -65.7 to -22.6 | 46.6 | 25.1 to 68.2 |
| | Skewness | -0.3 ± 27.5 | -15.0 to 14.3 | -54.2 | -79.8 to -28.6 | 53.6 | 28.0 to 79.2 |
| | Kurtosis | 2.1 ± 18.0 | -7.4 to 11.7 | -33.1 | -49.8 to -16.4 | 37.3 | 20.6 to 54.0 |
| | Mean/SD | 4.1 ± 24.1 | -8.8 to 16.9 | -43.2 | -65.6 to -20.7 | 51.3 | 28.9 to 73.7 |
| Local | 2nd ang moment | 10.9 ± 26.4 | -3.2 to 25.0 | -40.9 | -65.5 to -16.3 | 62.7 | 38.1 to 87.3 |
| | Contrast (inertia) | 5.4 ± 24.0 | -18.1 to 7.4 | -52.3 | -74.6 to -30.0 | 41.6 | 19.3 to 63.9 |
| | Entropy | -2.0 ± 5.4 | -4.9 to 0.9 | -12.6 | -17.7 to -7.6 | 8.7 | 3.6 to 13.8 |
| | Correlation | -0.6 ± 27.7 | -15.3 to 14.1 | -54.8 | -15.3 to 14.1 | 53.6 | 27.9 to 79.3 |
| | Homogeneity | 1.8 ± 11.5 | -4.4 to 7.9 | -20.8 | -31.5 to -10.1 | 24.4 | 13.6 to 35.1 |
| | Dissimilarity | -2.1 ± 13.0 | -9.0 to 4.9 | -27.6 | -39.7 to -15.5 | 23.5 | 11.4 to 35.6 |
| Regional | Small Area Emphasis (SAE) | -6.0 ± 54.3 | -35.0 to 22.9 | -112.5 | -163.0 to -62.0 | 100.4 | 49.9 to 150.9 |
| | Large Area Emphasis (LAE) | 3.6 ± 30.0 | -12.4 to 19.6 | -55.2 | -83.1 to -27.3 | 62.4 | 34.5 to 90.3 |
| | Intensity Variability (IV) | -9.7 ± 24.0 | -22.5 to 3.1 | -56.7 | -79.0 to -34.4 | 37.3 | 15.0 to 59.6 |
| | Size-Zone Variability (SZV) | 11.2 ± 23.1 | -1.1 to 23.5 | -34.1 | -55.6 to -12.6 | 56.5 | 35.0 to 78.0 |
| | Zone Percentage (ZP) | -2.7 ± 16.9 | -11.7 to 6.2 | -35.8 | -51.5 to -20.1 | 30.3 | 14.6 to 46.0 525.9 to 155.8 |
| | Low-Intensity Emphasis (LIE) | -4.0 ± 55.3 | -33.5 to 25.4 | -112.4 | -163.9 to -61.0 | 104.4 | |
| | High-Intensity Emphasis (HIE) | 3.9 ± 20.4 | -7.0 to 14.8 | -36.1 | -55.1 to -17.1 | 44.0 | 24.9 to 63.0 |
| | Low-Intensity Small Area Emphasis (LISAE) | -7.0 ± 67.6 | -43.1 to 29.0 | -139.5 | -202.4 to -76.6 | 125.4 | 62.5 to 188.3 |
| | High-Intensity Small Area Emphasis (HISAE) | 1.0 ± 31.2 | -15.6 to 17.6 | -60.1 | -89.1 to -31.1 | 62.0 | 33.0 to 91.0 |
| | Low-Intensity Large Area Emphasis (LILAE) | 1.8 ± 28.9 | -13.6 to 17.2 | -54.9 | -81.8 to 28.0 | 58.5 | 31.6 to 85.4 |
| High-Intensity Large Area Emphasis (HILAE) | 3.5 ± 35.8 | -15.6 to 22.6 | -66.7 | -100.1 to -33.4 | 73.7 | 40.4 to 107.1 | |

453

454

Table 3

455

456
457

Reliability of measurements using ICCs (calculated using a downsampling range of 64 values).

| Texture | Feature | ICC | 95% CI | Precision |
|--|--|--------------|---------------|-----------|
| Global | Minimum intensity | 0.99 | 0.92 to 0.99 | ± 4% |
| | Maximum intensity (SUV _{max}) | 0.94 | 0.82 to 0.98 | ± 8% |
| | Mean intensity (SUV _{mean}) | 0.92 | 0.78 to 0.97 | ± 10% |
| | SD | 0.99 | 0.96 to 0.99 | ± 2% |
| | Skewness | 0.82 | 0.49 to 0.94 | ± 23% |
| | Kurtosis | 0.80 | 0.44 to 0.93 | ± 25% |
| | Mean/SD | 0.82 | 0.49 to 0.94 | ± 23% |
| Local | 2nd ang moment | 0.95 | 0.85 to 0.98 | ± 7% |
| | contrast (inertia) | 0.94 | 0.82 to 0.98 | ± 8% |
| | Entropy | 0.98 | 0.93 to 0.99 | ± 3% |
| | correlation | 0.98 | 0.94 to 0.99 | ± 3% |
| | homogeneity | 0.88 | 0.64 to 0.96 | ± 16% |
| | dissimilarity | 0.93 | 0.81 to 0.98 | ± 9% |
| Regional | Small Area Emphasis (SAE) | 0.61 | -0.11 to 0.86 | ± 38% |
| | Large Area Emphasis (LAE) | 0.89 | 0.70 to 0.96 | ± 13% |
| | Intensity Variability (IV) | 0.97 | 0.93 to 0.99 | ± 3% |
| | Size-Zone Variability (SZV) | 0.97 | 0.91 to 0.99 | ± 4% |
| | Zone Percentage (ZP) | 0.84 | 0.55 to 0.95 | ± 20% |
| | Low-Intensity Emphasis (LIE) | 0.68 | 0.08 to 0.89 | ± 41% |
| | High-Intensity Emphasis (HIE) | 0.82 | 0.48 to 0.94 | ± 23% |
| | Low-Intensity Small Area Emphasis (LISAE) | 0.59 | -16 to 0.86 | ± 35% |
| | High-Intensity Small Area Emphasis (HISAE) | 0.83 | 0.52 to 0.94 | ± 21% |
| | Low-Intensity Large Area Emphasis (LILAE) | 0.93 | 0.80 to 0.98 | ± 9% |
| High-Intensity Large Area Emphasis (HILAE) | 0.78 | 0.36 to 0.92 | ± 28% | |

458

459

Table 4

460