



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Optimizing SUV Analysis: A Multicenter Study on Preclinical FDG-PET/CT Highlights the Impact of Standardization

Citation for published version:

Kuntner, C, Alcaide, C, Anestis, D, Bankstahl, JP, Boutin, H, Brasse, D, Elvas, F, Forster, D, Rouchota, MG, Tavares, A, Teuter, M, Wanek, T, Zachhuber, L & Mannheim, JG 2024, 'Optimizing SUV Analysis: A Multicenter Study on Preclinical FDG-PET/CT Highlights the Impact of Standardization', *Molecular Imaging and Biology*, vol. 26, no. 4, pp. 668-679. <https://doi.org/10.1007/s11307-024-01927-9>

Digital Object Identifier (DOI):

[10.1007/s11307-024-01927-9](https://doi.org/10.1007/s11307-024-01927-9)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Molecular Imaging and Biology

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



1 **Optimizing SUV Analysis: A Multicenter Study on Preclinical FDG-PET/CT Highlights**
2 **the Impact of Standardization**

3

4 Claudia Kuntner^{1, 2}, Carlos Alcaide³, Dimitris Anestis⁴, Jens P. Bankstahl⁵, Herve Boutin^{6,7},
5 David Brasse⁸, Filipe Elvas⁹, Duncan Forster¹⁰, Maritina G. Rouchota⁴, Adriana Tavares³, Mari
6 Teuter⁵, Thomas Wanek¹, Lena Zachhuber¹, Julia G. Mannheim^{11, 12}

7 ¹Department of Biomedical Imaging and Image-guided Therapy, Medical University of Vienna,
8 Vienna, Austria

9 ²Medical Imaging Cluster (MIC), Medical University of Vienna, Vienna, Austria

10 ³ University of Edinburgh, Edinburgh, United Kingdom

11 ⁴ BIOEMTECH, Athens, Greece

12 ⁵ Hannover Medical School, Hannover, Germany

13 ⁶ Division of Neuroscience & Experimental Psychology, Faculty of Biology, Medicine and
14 Health, University of Manchester, Manchester, UK

15 ⁷ INSERM, UMR 1253, iBrain, Université de Tours, Tours, France

16 ⁸ Institut Pluridisciplinaire Hubert Curien, Université de Strasbourg, CNRS, UMR7178,
17 Strasbourg, France

18 ⁹ Molecular Imaging Center Antwerp, University of Antwerpen, Antwerpen, Belgium

19 ¹⁰ Division of Informatics, Imaging and Data Sciences, Manchester Molecular Imaging Centre,
20 The University of Manchester, Manchester, UK.

21 ¹¹ Werner Siemens Imaging Center, Department of Preclinical Imaging and Radiopharmacy,
22 Eberhard-Karls University Tuebingen, Tuebingen, Germany

23 ¹² Cluster of Excellence iFIT (EXC 2180) "Image Guided and Functionally Instructed Tumor
24 Therapies", Tuebingen, Germany

25

26 Corresponding & first author: Claudia Kuntner

27 Address: Waehringer Guertel 18-20, 1090 Vienna, Austria

28 E-mail address: claudia.kuntner@meduniwien.ac.at

29 Telephone number: +43 (0)1 40160-64022

30

31 Short running title: Multicenter preclinical image analysis

32 Manuscript category: Original Article

33

34

35 **Abstract**

36 Purpose: Preclinical imaging, with translational potential, lacks a standardized method for
37 defining volumes of interest (VOIs), impacting data reproducibility. The aim of this study was
38 to determine the interobserver variability of VOI sizes and standard uptake values (SUV_{mean}
39 and SUV_{max}) of different organs using the same [¹⁸F]FDG-PET and PET/CT datasets analyzed
40 by multiple observers. In addition, the effect of a standardized analysis approach was
41 evaluated.

42 Procedures: In total, 12 observers (4 beginners and 8 experts) analyzed identical preclinical
43 [¹⁸F]FDG-PET-only and PET/CT datasets according to their local default image analysis
44 protocols for multiple organs. Furthermore, a standardized protocol was defined, including
45 detailed information on the respective VOI size and position for multiple organs, and all
46 observers reanalyzed the PET/CT datasets following this protocol.

47 Results: Without standardization, significant differences in the SUV_{mean} and SUV_{max} were
48 found among the observers. Coregistering CT images with PET images improved the
49 comparability to a limited extent. The introduction of a standardized protocol that details the
50 VOI size and position for multiple organs reduced interobserver variability and enhanced
51 comparability.

52 Conclusions: The protocol offered clear guidelines and was particularly beneficial for
53 beginners, resulting in improved comparability of SUV_{mean} and SUV_{max} values for various
54 organs. The study suggested that incorporating an additional VOI template could further
55 enhance the comparability of the findings in preclinical imaging analyses.

56

57 **Key words (3-5):** multicenter, image analysis, reproducibility, PET/CT, preclinical imaging

58 Introduction

59 Over the past few decades, preclinical molecular imaging, notably positron emission
60 tomography (PET) combined with computed tomography (CT), has become indispensable in
61 scientific medical research [1-2]. This approach offers multimodal imaging in preclinical
62 models that are highly translatable to clinical settings [3-4]. PET enables quantification of
63 biological processes in living subjects, achieved by defining regions or volumes of interest
64 (ROIs or VOIs) on the images to extract activity concentrations (typically given in kBq/cc).
65 Mathematical operations transform these activity concentrations into percent injected activity
66 or dose per volume of tissue (%IA/cc or %ID/cc) by normalizing them to the administered
67 activity or standardized uptake values (SUVs) by additionally normalizing to the body weight.
68 The SUV is used as a semiquantitative measurement of glucose uptake in tissue from a 2-
69 deoxy-2-[¹⁸F]fluoro-D-glucose ([¹⁸F]FDG) PET scan, especially in clinical practice [5]. The
70 SUV_{mean}, reflecting the mean voxel value within a VOI, is strongly influenced by the VOI
71 definition method and is susceptible to partial volume effects, resulting in greater variability.
72 Conversely, the SUV_{max}, which represents the voxel with the highest radioactivity
73 concentration, is less affected by observer variability but more affected by technical variations
74 [6].

75 A major limitation in preclinical imaging is the lack of standardized or fully automated methods
76 for defining VOIs. While some data-driven or semiautomatic segmentation methods exist, they
77 still require observer input to define or choose the proposed cluster. Anatomy-based automatic
78 segmentation methods rely heavily on annotated training images (magnetic resonance (MR)
79 and/or CT), but their effectiveness hinges on the quality and quantity of the database.
80 Currently, there is no widely accepted automated preclinical VOI delineation method.
81 Consequently, most preclinical image analysis is manual, with observers selecting regions for
82 analysis. Additionally, the availability of multiple software tools for preclinical PET/CT image
83 analysis, each with different features and pipelines, further complicates the issue.

84 For clinical PET/CT imaging, several studies have assessed inter- and intraobserver variability
85 and proposed methods to standardize image analysis [7-10]. Until now, there hasn't been any
86 study conducted on preclinical PET/CT imaging that includes a standardized image analysis.
87 Therefore, the present study assessed the variability in VOI size, SUV_{mean} , and SUV_{max}
88 measurements of multiple organs and tumors between different observers (grouped into
89 beginners and experts) when analyzing the same preclinical [^{18}F]FDG-PET-only and
90 [^{18}F]FDG-PET/CT datasets with free or commercially available image analysis software.
91 Furthermore, a standardized protocol was used, and all observers reanalyzed the PET/CT
92 datasets following this protocol; potential improvements in interobserver variability were
93 evaluated accordingly.

94

95 **Materials and Methods**

96 Imaging data

97 Twelve observers analyzed dynamic [^{18}F]FDG-PET-only (dynamic images 0-75 min, 25
98 frames; n=6) and [^{18}F]FDG-PET/CT (dynamic images 0-60 min, 19 frames; n=7) scans of
99 tumor-bearing mice. Two laboratories provided the datasets, which were acquired according
100 to local regulations. The images were provided in Bq/cc together with the injected activities
101 and weights of the mice in the scanner-specific and DICOM formats. Information regarding
102 the animal experiments and imaging protocols can be found in the Electronic Supplementary
103 Material (ESM). Co-registration of PET/CT data for part 2 and 3 was performed by one
104 observer to eliminate potential co-registration-induced influences.

105 Of the twelve observers, eight were experts in the analysis of preclinical images (> 4 years of
106 experience), whereas four were classified as beginners (< 1 year of experience). With the
107 exception of the dataset providers, all observers analyzed the images independently and
108 blinded to each other's assessments, utilizing their expertise and judgment.

109

110 Part 1: [¹⁸F]FDG-PET-only image analysis and reporting

111 The observers were asked to analyze the images according to their standard institutional
112 procedures, including the choice of image analysis software, the procedures for preparing the
113 images (e.g., adjustment of the animal's position), the radiation scale and time frames, and
114 the method of delineating VOIs. The observers were requested to delineate the following VOIs:
115 tumor, whole brain, muscle, heart (either whole heart or left ventricle), kidneys (left and right),
116 liver, and urinary bladder (short name bladder). An additional region covering the whole FOV
117 was delineated on the last time frame with a predefined size ($128 \times 128 \times 95$ voxels/ $51.2 \times$
118 51.2×75.62 mm³) to assess any software-related biases in image quantitation.

119 After analyzing the images, the observers completed a detailed report, including SUV_{mean} and
120 SUV_{max} (normalized to the body weight of the animals, respectively), VOI delineation method
121 (manual, thresholding, fixed objects, etc.), and volume (in mm³). They also specified how they
122 displayed the images (radiation scale, minimum and maximum values, kBq/cc, %IA/cc, or
123 SUV). As the datasets were dynamic, observers indicated the time frame (individual frame or
124 summed image) for VOI delineation. Time-activity curves (TACs) for all animals and organs
125 were plotted. Group differences (SUV_{mean} and SUV_{max}) were determined across observers and
126 animals based on the 10 min time frame from 55-65 min.

127 Part 2: [¹⁸F]FDG-PET/CT image analysis and reporting

128 The image analysis procedure for the PET/CT datasets was identical to that for the [¹⁸F]FDG-
129 PET-only datasets. Only the whole FOV region was adjusted ($256 \times 256 \times 159$ voxels/ 99.377
130 $\times 99.377 \times 126.564$ mm³) as a different PET scanner was used for these experiments. In
131 addition, the observers were asked to report on which dataset (PET or CT) each organ and
132 the tumor were delineated. Group differences (SUV_{mean} and SUV_{max}) were determined across
133 observers and animals based on the 5 min time frame from 55-60 min.

134

135 Part 3: Standardized [¹⁸F]FDG-PET/CT image analysis and reporting

136 The authors established a standardized tumor and organ VOI definition method based on
137 [¹⁸F]FDG-PET-only and [¹⁸F]FDG-PET/CT data analysis results. The protocol required to be
138 universally applicable across image analysis software tools. Consequently, data-driven
139 segmentation methods, such as multiclustering, were excluded from part 3, resulting in the
140 exclusion of observer E8. Observer B3's analysis was also omitted due to inability to meet the
141 standardized consensus specifications for VOI definition.

142 Observers unanimously opted to delineate organs and tumors using specific objects (ellipsoids
143 and boxes), with predefined VOI drawing on either PET or CT images. PET-related VOIs
144 adhered to a fixed radiation scale specified in SUV. VOIs for the brain, heart and tumor were
145 delineated on the CT images (and verified on the respective PET images), as the CT image
146 provided sufficient anatomical delineation to surrounding tissues. The VOIs for both muscle
147 regions, kidneys, liver and both bladder regions were delineated on the PET images (and
148 verified on the respective CT images) due to the fact that for most of these organs the [¹⁸F]FDG
149 uptake is very distinct and the low soft-tissue contrast of the CT doesn't enable a clear
150 delineation to surrounding tissues.

151 **Table 1** summarizes the objects and predefined VOI sizes and ranges. To explore VOI
152 position influence on quantitative analysis, two muscle regions (gluteus maximus and
153 biceps/triceps) and two urinary bladder regions (bottom and maximum fill) were included.

154

155 **Table 1** Details on the standardized VOI analysis. The PET-related VOIs were delineated at
156 the last time frame using the specified SUV radiation scale.

VOI	image used for VOI delineation	radiation scale (SUV)	shape	size	notes
tumor	CT	n.a.	ellipsoid	entire tumor	
brain	CT	n.a.	ellipsoid	7 x 5 x 10 mm ³	inside skull, control on PET that olfactory bulb and hardierian glands are excluded
heart	CT	n.a.	ellipsoid	>100 and <200 mm ³	
muscle	PET	0 - 2	box	2 x 2 x 3 mm ³	gluteus maximus, avoid spill in from bladder, control on CT that no bone is included
muscle	PET	0 - 2	box	2 x 2 x 3 mm ³	biceps/triceps, control on CT that no bone is included
kidney	PET	0 - 2	ellipsoid	~ 60 mm ³	definition of right and left side
liver	PET	0 - 2	box	4 x 4 x 4 mm ³	opposite to the stomach

bladder bottom	PET	0 - 10	box	2 x 2 x 2 mm ³	bottom of bladder
bladder maximum fill	PET	0 - 10	ellipsoid	entire bladder	draw on time frame with largest bladder fill

157

158 Statistical analysis

159 The mean or maximum radioactivity concentrations given as SUV_{mean} or SUV_{max} per animal
160 and organ over the 12 (part 1 and 2) and 10 (part 3) observers were used.

161 The coefficient of variation (CV, %) was calculated as the ratio of the standard deviation to the
162 mean to assess the extent of variability. Moreover, to account for the variability between
163 animals, the normalized difference was calculated for each animal and organ based on the 60
164 min values using the following equation:

165
$$\text{normalized difference} = \frac{\text{individual value} - \text{mean value}}{\text{mean value}}$$

166 The data are expressed as the mean ± standard deviation. Statistical analysis was performed
167 with Prism 9.5.0 Software (GraphPad, La Jolla, CA, USA) and SPSS Statistics (version 29.0,
168 IBM SPSS, IBM Corp., Armonk, NY, USA). Differences between the beginner and expert
169 groups were assessed by applying two-way ANOVA followed by a Bonferroni multiple
170 comparisons test, with an alpha level of 0.05 for each organ. Brown-Forsythe and Welch
171 ANOVA tests were performed to assess interobserver variability, followed by Dunnett's
172 multiple comparisons test, with individual variances computed for each comparison and organ.
173 The threshold of statistical significance was set to an adjusted p value ≤ 0.05.

174 Intraclass correlation coefficients (ICCs; single-measure, two-way random, absolute
175 agreement) were calculated based on the SUV_{mean} and SUV_{max} values to determine
176 interobserver reliability for the beginners, the experts, and all observers [11-12]. According to
177 Koo et al. [12], ICCs less than 0.5 can be classified as poor reliability, ICCs in the range of 0.5
178 to 0.75 as moderate reliability, ICCs between 0.75 and 0.8 as good reliability, and ICCs greater
179 than 0.9 as excellent reliability.

180

181 **Results**

182 **Selection of image analysis software programs and VOI definition methods**

183 Five different image analysis software programs were utilized in the present study. The
184 selected software and the typically used output units, radiation scales, and time frames are
185 summarized in the **Suppl. Tab. s1** (see ESM). One observer employed a data-driven
186 segmentation method (observer E8, BrainVISA/Anatomist) that used the local means analysis
187 method based exclusively on the dynamics (i.e., time-activity and level of uptake) of each
188 voxel in the PET images [13-14]. The VOIs of six of the remaining eleven observers were
189 defined in the last time frame. Some observers (3 out of 11) selected the time frame where
190 the respective organ was clearly visible for analysis. Seven out of the eleven observers applied
191 different radiation scales for specific organs (e.g., 0-2 SUV for muscle, 0-20 SUV for the heart),
192 whereas the rest used a fixed radiation scale for all organs. The whole FOV region evaluated
193 in parts 1 and 2 revealed no systematic software biases in image-based quantitation of the
194 mean and maximum activity values (**Suppl. Fig. s1**, see ESM). These small differences were
195 attributed to the VOI position in the whole FOV region.

196

197 **Parts 1 and 2: Individual [¹⁸F]FDG-PET-only and [¹⁸F]FDG-PET/CT image analysis**

198 VOI sizes

199 The VOI delineation methods vary from fixed objects (e.g., spheres for the whole brain and
200 heart) to manual drawings of VOIs on consecutive slices to those using thresholds (see **Fig.**
201 **1** for examples of VOI positions and shape for each software tool). Some observers applied
202 post-processing to re-orient the images according to the “standard” configuration in preclinical
203 imaging (head first, prone), whereas others analyzed the images in the orientation provided
204 by the scanner. The delineation methods used for each organ are summarized in the
205 supplementary methods (**Suppl. Fig. s2 and s3**, see ESM) for the PET-only and PET/CT
206 studies, respectively.

207 For the [¹⁸F]FDG-PET-only study, the tumor VOI was excluded from the analysis because
208 delineation was rather challenging due to the low uptake and small size of the tumors (most
209 of the observers could not identify the tumors).

210 The different delineation methods resulted in considerable variability in the VOI sizes, as
211 illustrated in **Fig. 2** ((**a**) [¹⁸F]FDG-PET-only; (**b**) [¹⁸F]FDG-PET/CT). The beginners delineated
212 significantly larger liver and heart VOIs than did the experts on the PET images (part 1). The
213 smallest variability in the VOI sizes in the beginner group was obtained for the heart (71%
214 CV), whereas in the expert group, the smallest variability was obtained for the kidneys (52%
215 CV). In contrast, the greatest variability was found in the muscle VOI (149% CV) for the
216 beginner group and in the liver VOI (210% CV) for the expert group.

217 On the [¹⁸F]FDG-PET-CT images (part 2), the beginners delineated significantly larger VOIs
218 than did the experts in the liver, heart, and brain. The smallest variability in VOI sizes was
219 obtained in the bladder for the beginners (37% CV) and in the tumor VOIs for the experts (40%
220 CV). The highest variability in VOI sizes was found in the muscle for the beginners (159% CV)
221 and in the liver for the experts (164% CV). In particular, the VOI drawn for the liver ranged
222 from 16 to 3619 mm³, which spans two orders of magnitude. Furthermore, the VOI position for
223 the muscle differed among the observers (e.g., for part 2, the lower left limb was delineated
224 by seven observers, the upper left limb was delineated by four observers, and the upper right
225 limb was delineated by one observer).

226

227 Organ-time activity curves

228 The organ TACs for part 1 [¹⁸F]FDG-PET-only images for a representative animal, subdivided
229 into beginner and expert groups, are shown in **Suppl. Fig. s4** (SUV_{mean}) and **Fig. s5** (SUV_{max})
230 in the ESM. The heart and kidney SUV_{mean} TACs exhibited greater interobserver variation in
231 the beginner group than in the expert group. The remaining organs revealed a similar pattern
232 between beginners and experts.

233 For the SUV_{max} of the TACs, the beginner group revealed greater interobserver variation for
234 the brain and muscle; interestingly, the experts showed greater variability than the beginners
235 for the liver and heart.

236 The inclusion of CT data (part 2) reduced the variability in the liver, brain, and muscle SUV_{mean}
237 TACs, as depicted in **Suppl. Fig. s6** and **Fig. s7** (see ESM). For the SUV_{max} of the TACs
238 (beginners: **Suppl. Fig. s8**; experts: **Suppl. Fig. s9**), reduced variability was detected mainly
239 for the muscle. The two groups of observers determined identical SUV_{max} TACs for the tumor,
240 kidney, and bladder.

241

242 Last time frame analysis

243 The SUV_{mean} and SUV_{max} values from the time frame covering 60 min were used to compare
244 the variability between groups (beginners and experts) and individual observers. For the PET-
245 only study, the calculated normalized difference based on the SUV_{mean} showed the greatest
246 deviation from 0 for the heart region (-0.25 ± 0.27 for beginners and 0.13 ± 0.18 for experts)
247 and the smallest deviation for the brain (0.01 ± 0.14 for beginners and -0.01 ± 0.14 for experts),
248 as displayed in the upper row of **Fig. 3 (a)**. In addition, statistically significant differences were
249 observed between the beginner and expert groups for the heart, muscle and bladder. The
250 ICCs revealed greater reliability within the expert groups for all organs except the brain,
251 although poor reliability was observed for the muscle and liver (ICCs<0.5).

252 The calculated normalized difference based on the SUV_{max} (**Fig. 3(b)**) yielded the greatest
253 deviation from 0 for the muscle region among the beginners (0.24 ± 0.81) and for the bladder
254 among the experts (0.14 ± 0.95). The smallest deviation was found for the kidney region
255 (beginners: 0.01 ± 0.02 ; experts: -0.01 ± 0.07). Overall, no statistically significant differences
256 between the observer groups were observed. An overview of all the ICCs, including
257 confidence intervals (CIs), for each organ can be found in the supplementary materials (**Suppl.**
258 **Tab. s2**, see ESM).

259 Multiple statistically significant differences in the SUV_{mean} were detected between the
260 individual observers, especially for the heart and muscle VOIs, as shown in **Fig. 4(a)**. For the
261 SUV_{max} , the liver and muscle indices revealed multiple significant differences among the 12
262 observers (**Fig. 4(b)**). The individual p values are given in Suppl. Fig. s10 (see ESM).

263 For the PET/CT study, the normalized difference of the muscle for beginners and experts was
264 reduced (compare the middle row of **Fig. 3(a)**). However, statistically significant differences
265 between the observer groups were obtained for the heart, kidneys, bladder, and tumor. The
266 ICCs for the liver, muscle, and bladder showed improved reliability compared to those of part
267 1. Analyzing the normalized difference based on the SUV_{max} (**Fig. 3(b)**) yielded the largest
268 overall spread in the liver region (0.60 ± 1.67 for the beginners and -0.25 ± 0.73 for the experts,
269 $p < 0.0001$). No improvement in reliability was detected for the ICCs based on the SUV_{max} for
270 part 2 compared to part 1.

271 The interobserver SUV_{mean} and SUV_{max} variability are displayed in **Fig. 5(a)** and **6(a)**, revealing
272 multiple statistically significant differences in the heart and tumor regions (both SUV_{mean}) as
273 well as the liver and brain regions (both SUV_{max}). The individual p values between the
274 observers are given in Suppl. Fig. s11 and Fig. s12 (see ESM).

275

276 **Part 3: standardized [^{18}F]FDG-PET/CT image analysis**

277 The predefined VOI sizes reduced the variations, as shown in **Fig. 2(c)**. However, for the two
278 regions for which the entire structure was to be delineated, namely, the tumor and the bladder
279 at the maximum-fill level, significantly larger VOIs were determined by experts with great
280 variability (tumor: beginners: 41% CV; experts: 38% CV; bladder: beginners: 56% CV; experts:
281 45% CV).

282

283 Organ-time activity curves after standardization

284 The standardized image analysis method reduced the variation in the SUV_{mean} TACs of the
285 tumor, brain, liver, and kidney, as shown in panel B in the **Suppl. Fig. s6 and s7** (see ESM).
286 The muscle and bladder TACs exhibited different patterns depending on the VOI position. The
287 expert group obtained mostly congruent SUV_{max} TACs for the liver, heart, tumor, brain,
288 kidneys, and bladder maximum-fill VOIs (**Suppl. Fig. s9**), whereas the beginner group
289 obtained slightly greater variations (**Suppl. Fig. s8**, see ESM).

290

291 Last time frame analysis after standardization

292 The standardized analysis approach notably enhanced the normalized difference based on
293 SUV_{mean} for most organs, depicted in the lower row of **Fig. 3(a)**, correlating with higher ICCs
294 across most organs. Liver and brain index reliability significantly improved, achieving excellent
295 levels post-standardization. Initially poor heart and tumor reliability transformed into good and
296 moderate levels, respectively. Standardization notably elevated kidney index reliability from
297 moderate to excellent levels. However, statistically significant differences persisted between
298 observer groups for muscle gluteus maximus and urinary bladder maximum-fill regions.
299 Improvement in normalized difference based on SUV_{max} was inconsistent post-
300 standardization, with no improvement observed for tumor or urinary bladder (**Fig. 3(b)**).
301 Significant differences between observer groups were found for liver and gluteus maximus
302 region (SUV_{max}). Notably, liver and brain ICCs substantially improved in standardized analysis
303 (liver: part 2=0.08, part 3=0.43; brain: part 2=0.00, part 3=0.65).

304 The interobserver variability based on the SUV_{mean} values was markedly reduced using the
305 standardized image analysis approach. However, some statistically significant differences
306 between observers persisted in the tumor, biceps/triceps muscle, or maximum-fill urinary
307 bladder region (**Fig. 5(b)**). The individual p values between the observers are given in Suppl.
308 Fig. s13 (see ESM). For the SUV_{max} , no significant differences were found between the
309 observers for any of the organs (**Fig. 6(b)**).

310 **Discussion**

311 Quantifying radioactivity concentrations in small animal organs or tumors is standard in
312 preclinical imaging and relies on parameters such as the SUV_{mean} or SUV_{max} . However, the
313 variability and reproducibility of these parameters among different observers within a single
314 institution or across multiple centers remain poorly understood. Currently, each imaging lab
315 and often each observer within the same institution applies different workflows, experiences,
316 and judgments to analyze and segment PET images. These variations encompass factors
317 such as the position, size, and shape of VOIs; PET image display settings; and postprocessing
318 methods, potentially compromising comparability across observers and centers. Despite the
319 prevalence of preclinical [^{18}F]FDG-PET/CT studies, no multicenter consensus exists on a
320 reproducible image analysis method. This study represents the first comprehensive
321 multicenter [^{18}F]FDG-PET/(CT) investigation into the impact of image analysis methods on
322 results and the comparability of a standardized analysis approach. Our findings underscore
323 the significant influence of image analysis methods on [^{18}F]FDG-PET/(CT) study outcomes,
324 particularly regarding SUV_{mean} discrepancies attributed to regional position and size,
325 corroborating similar observations from prior studies [15].

326 Our first observation was that not all observers performed post-processing to re-orient the
327 images according to the “standard” configuration in preclinical imaging (head first, prone).
328 Some analyzed the images in the orientation provided by the scanner, which was for the
329 PET/CT study in feed first, prone. Thus, an agreement on the orientation of images to be used
330 (also with regard to future automatic segmentation applications) is therefore the first step
331 towards standardized image analysis. Without standardization, variations in VOI sizes were
332 observed between beginners and experts for multiple organs. These differences influenced
333 SUV_{mean} (e.g., heart) and SUV_{max} (e.g., liver in PET/CT) analyses, suggesting that VOI size
334 impacts uptake. However, for certain organs (e.g., the liver in PET-only and the brain in
335 PET/CT), despite significant differences in VOI size, SUV analysis was unaffected by
336 homogeneous [^{18}F]FDG uptake.

337 Introducing anatomical references in part 2 reduced variability in heart and muscle regions but
338 had no effect on liver or brain regions. However, overall reliability and comparability did not
339 improve universally. Comparing parts 1 and 2 is challenging due to the different image sets
340 analyzed. However, this design showcases variability between studies (e.g., small vs. large
341 tumors with necrotic areas), mitigating potential biases from part 1 to part 2.

342 Based on the results from these two studies, the participants in this study reached a consensus
343 on the standardized VOI delineation method utilized in part 3.

344 Standardization improved the consistency and shape of SUV_{mean} TACs in the liver, brain, and
345 kidney, while nearly identical SUV_{max} TACs were obtained in the liver, heart, tumor, brain,
346 kidneys, and urinary bladder. Reduced interobserver variability poststandardization was
347 evidenced by reduced deviation and improved ICCs across organs, except for muscle and
348 urinary bladder regions. Muscle VOIs are small and prone to spill over from adjacent bone
349 regions, making muscle-fat differentiation challenging despite the use of anatomical
350 information from CT scans. Intensive training and visual aids are recommended for
351 comparability improvement. For maximum-fill bladder VOIs, inconsistent time frame choices
352 hindered comparisons between parts 2 and 3. Nevertheless, considering its importance in
353 dosimetric studies, assessing bladder necessity and employing frame-by-frame analysis for
354 volumetric changes are advised.

355 Furthermore, the significant differences between beginners and experts found by the
356 normalized difference analysis in the heart, kidneys, and tumor diminished after
357 standardization (**Fig. 3(b)** and **3(c)**). We concluded that the use of a standardized approach
358 reduced the interobserver variability in the SUV analysis. In addition, we propose to create a
359 VOI template for each preclinical PET/CT and PET/MR study that includes a standardized VOI
360 positioning and size as well as detailed information on the segmentation method. For
361 multicenter studies, we recommend reaching a consensus on the use of single analysis
362 software for evaluating and providing VOI template files. For single-center studies, a VOI

363 template from the first animal analyzed will ensure reproducibility for the remaining animals
364 and help train new personnel.

365 In general, the SUV_{max} revealed a lower interobserver variability than the SUV_{mean} in our study.
366 However, as the SUV_{max} represents only a single voxel within a region, the SUV_{mean} might be
367 a more stable marker for underlying tissue uptake. Therefore, both measures can be valuable
368 in multicenter studies.

369 Despite its strengths, our study has several limitations. First, mid-level observers were not
370 included, potentially biasing the results, as experience levels were subjectively categorized as
371 beginners or experts. Additionally, the varied backgrounds of the participating observers (e.g.,
372 physics, chemistry, biology, etc.) may have influenced interpretation. Secondly, validation
373 using gamma-counter data was not available. Third, the use of different image analysis
374 software led to the use of various segmentation tools, hindering detailed discrepancy
375 identification within segmented VOIs. Finally, the standardized protocol lacked optimization,
376 notably omitting a VOI template for precise location visualization. Addressing these limitations
377 in future studies could enhance the accuracy and reproducibility of the findings.

378 It has to be noted that depending on the specific tracer used, standardized image analysis
379 protocols need to be re-defined to address tracer-specific factors that might impact the
380 reproducibility of image analysis. This also applies for the acquisition of the imaging data, for
381 which standardized protocols – depending on the used tracer – can also significantly enhance
382 reproducibility [16].

383 The 12 observers in this study represent 8 different preclinical imaging facilities in Europe and
384 all observers were asked to use their default image analysis method and software tool to
385 analyze the provided PET(/CT) data. Only 1 observer analyzed the data using an automated
386 segmentation tool. Automatic organ segmentation has been an active field of research for
387 decades [17-22], and current research in this field includes the development of artificial
388 intelligence (AI)-assisted solutions [23]. Nevertheless, manual delineation will still be the

389 standard method for image analysis until these tools are applicable to a broader community
390 with sufficient training databases and a variety of VOI templates. The variety of chosen
391 software tools and methods utilized in this study represents in our opinion the current standard
392 in preclinical imaging. However, the transition to AI-guided automatic segmentation will
393 certainly be a strong focus within the next decade and thus will potentially improve the
394 comparability and reliability of preclinical multicenter image analysis.

395

396 **Conclusion**

397 For the first time, the present study demonstrated the significant influence of image analysis
398 on the obtained quantitative data; this work is intended as the basis for a discussion of further
399 standardization approaches in preclinical imaging. Moreover, the authors aim to raise
400 awareness of potential pitfalls when preclinical data are analyzed by multiple observers with
401 different levels of experience. Our study verified that the comparability of image analysis
402 significantly improves when detailed standardized image analysis protocols are used. This
403 approach will be of particular interest not only for preclinical multicenter studies but also for
404 studies performed over a long period within the same institution, where the observers might
405 vary.

406

407 **Acknowledgment**

408 For this work, the methodological advice of the Institute of Clinical Epidemiology and Applied
409 Biometry of the University of Tübingen was applied. We would like to express our sincere
410 thanks to Mr. Blumenstock for his support.

411

412 **Author contributions**

413 CK and JGM designed the study. CK, HB and DF provided the data. CK, CA, DA, JB, HB, BD,
414 FE, DF, MT, TW, LZ and JGM analyzed the image data. AT and MGR interpreted the data.
415 CK and JGM performed the comparability analysis of all observer analyses. All the authors
416 were involved in critically revising the manuscript. All the authors have read and approved the
417 final version of the manuscript.

418

419 **Funding**

420 This work was supported by the COST Action "Correlated Multimodal Imaging in Life
421 Sciences" (COMULIS, CA17121) and by the Chan Zuckerberg Initiative, Advancing Imaging
422 through Collaborative Projects (COMULISglobe, 2023-321161).

423

424 **Conflict of interest**

425 Author DA and MGR are employees of the company BIOEMTECH.

426

427 **Ethical approval**

428 All applicable institutional and/or national guidelines for the care and use of animals were
429 followed.

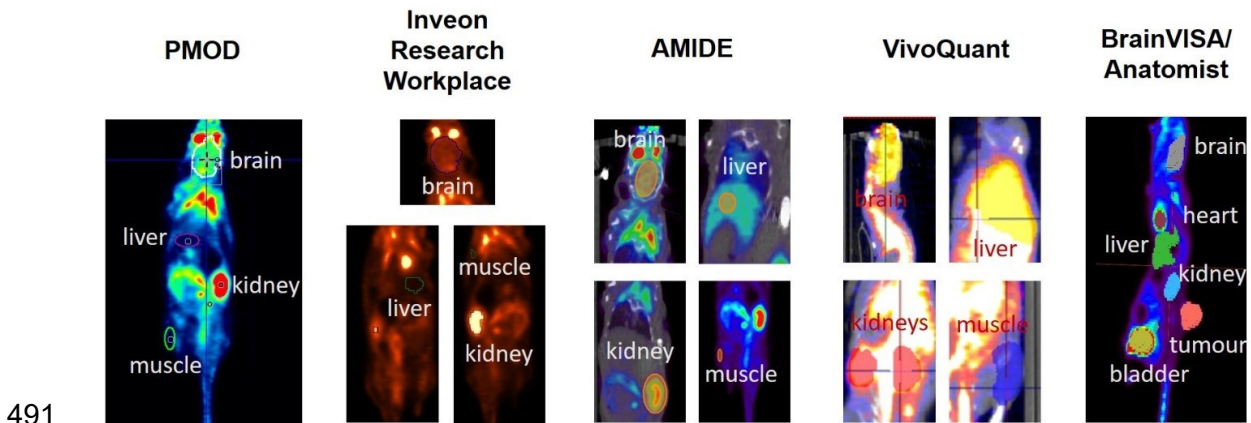
430

431 **References**

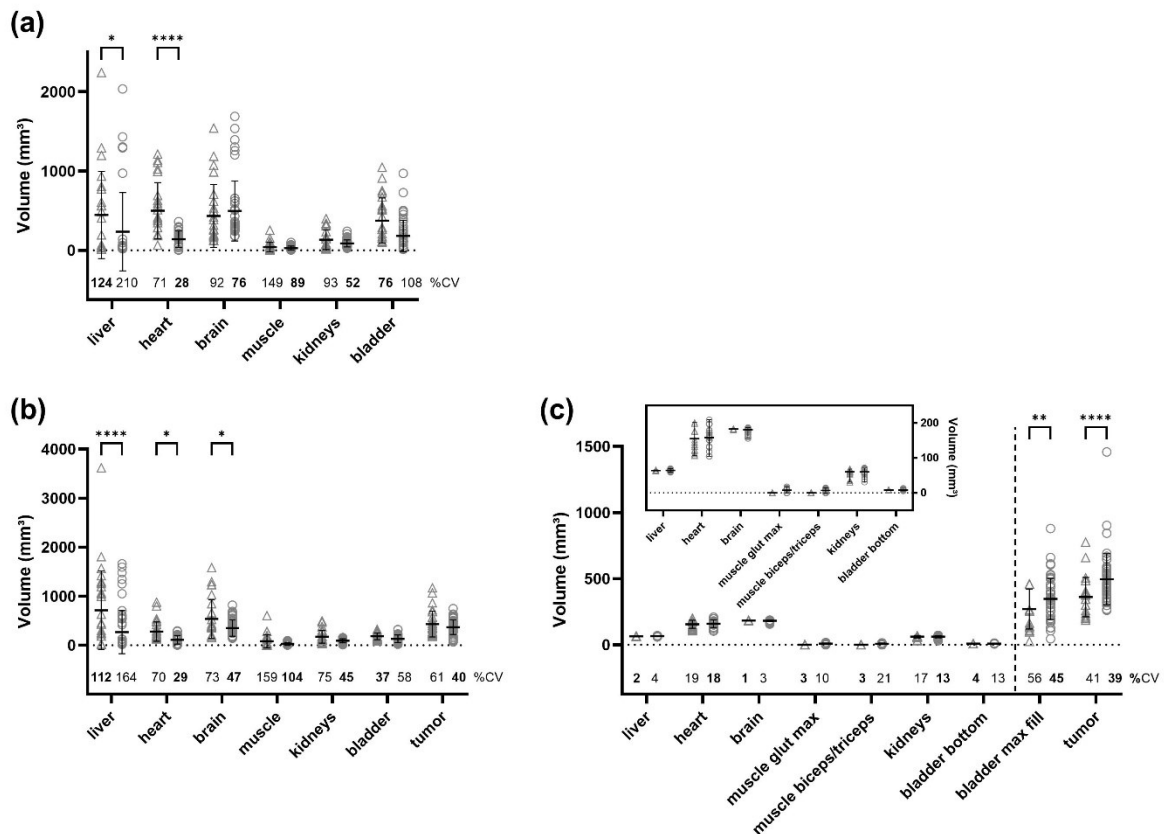
- 432 1. Lewis JS, Achilefu S, Garbow JR, Laforest R, Welch MJ (2002) Small animal imaging.
433 current technology and perspectives for oncological imaging. *Eur J Cancer* 38:2173-2188.
- 434 2. Kiessling F, Pichler BJ, Hauff P (2017) Small Animal Imaging.
- 435 3. Cherry SR, Gambhir SS (2001) Use of positron emission tomography in animal
436 research. *ILAR J* 42:219-232.

- 437 4. Phelps ME (2004) PET: Molecular Imaging and Its Biological Applications. 1st ed.
438 Softcover of orig. ed. 2004 ed., Berlin: Springer.
- 439 5. Kinahan PE, Fletcher JW (2010) Positron emission tomography-computed
440 tomography standardized uptake values in clinical practice and assessing response to
441 therapy. *Semin Ultrasound CT MR* 31:496-505.
- 442 6. De Luca GMR, Habraken JBA (2022) Method to determine the statistical technical
443 variability of SUV metrics. *EJNMMI Phys* 9:40.
- 444 7. Suzuki A, Nakamoto Y, Terauchi T, et al. (2007) Inter-observer Variations in FDG-PET
445 Interpretation for Cancer Screening. *Japanese Journal of Clinical Oncology* 37:615-622.
- 446 8. Büyükdereli G, Güler M, Şeydaoğlu G (2016) Interobserver and Intraobserver
447 Variability among Measurements of FDG PET/CT Parameters in Pulmonary Tumors. *Balkan*
448 *Med J* 33:308-315.
- 449 9. Gerke O, Vilstrup MH, Segtnan EA, Halekoh U, Høilund-Carlson PF (2016) How to
450 assess intra- and inter-observer agreement with quantitative PET using variance component
451 analysis: a proposal for standardisation. *BMC Med Imaging* 16:54.
- 452 10. Guezennec C, Bourhis D, Orhac F, et al. (2019) Inter-observer and segmentation
453 method variability of textural analysis in pre-therapeutic FDG PET/CT in head and neck
454 cancer. *PLoS One* 14:e0214299.
- 455 11. Fisher RA (1992) Statistical Methods for Research Workers. In *Breakthroughs in*
456 *Statistics: Methodology and Distribution*, Eds. Kotz S, Johnson NL. New York, NY: Springer
457 New York, pp 66-70.
- 458 12. Koo TK, Li MY (2016) A Guideline of Selecting and Reporting Intraclass Correlation
459 Coefficients for Reliability Research. *J Chiropr Med* 15:155-163.
- 460 13. Maroy R, Boisgard R, Comtat C, et al. (2008) Segmentation of rodent whole-body
461 dynamic PET images: an unsupervised method based on voxel dynamics. *IEEE Trans Med*
462 *Imaging* 27:342-354.
- 463 14. Maroy R, Boisgard R, Comtat C, et al. (2010) Quantitative organ time activity curve
464 extraction from rodent PET images without anatomical prior. *Med Phys* 37:1507-1517.
- 465 15. Habte F, Budhiraja S, Keren S, Doyle TC, Levin CS, Paik DS (2013) In situ study of
466 the impact of inter- and intra-reader variability on region of interest (ROI) analysis in preclinical
467 molecular imaging. *Am J Nucl Med Mol Imaging* 3:175-181.
- 468 16. Mannheim JG, Mamach M, Reder S, et al. (2019) Reproducibility and Comparability of
469 Preclinical PET Imaging Data: A Multicenter Small-Animal PET Study. *J Nucl Med* 60:1483-
470 1491.
- 471 17. Baiker M, Milles J, Dijkstra J, et al. (2010) Atlas-based whole-body segmentation of
472 mice from low-contrast Micro-CT data. *Med Image Anal* 14:723-737.
- 473 18. Khmelinskii A, Baiker M, Kaijzel EL, Chen J, Reiber JH, Lelieveldt BP (2011)
474 Articulated whole-body atlases for small animal image analysis: construction and applications.
475 *Mol Imaging Biol* 13:898-910.
- 476 19. Wang H, Stout DB, Chatziioannou AF (2012) Estimation of mouse organ locations
477 through registration of a statistical mouse atlas with micro-CT images. *IEEE Trans Med*
478 *Imaging* 31:88-102.
- 479 20. Akselrod-Ballin A, Dafni H, Addadi Y, et al. (2016) Multimodal Correlative Preclinical
480 Whole Body Imaging and Segmentation. *Sci Rep* 6:27940.
- 481 21. Yan D, Zhang Z, Luo Q, Yang X (2017) A Novel Mouse Segmentation Method Based
482 on Dynamic Contrast Enhanced Micro-CT Images. *PLoS One* 12:e0169424.
- 483 22. Wang H, Han Y, Chen Z, Hu R, Chatziioannou AF, Zhang B (2019) Prediction of major
484 torso organs in low-contrast micro-CT images of mice using a two-stage deeply supervised
485 fully convolutional network. *Phys Med Biol* 64:245014.
- 486 23. Schoppe O, Pan C, Coronel J, et al. (2020) Deep learning-enabled multi-organ
487 segmentation in whole-body mouse scans. *Nature Communications* 11:5626.
- 488

490 **Figures**



492 **Fig. 1** Representative images of multiple VOI positions for the individual software tools utilized
 493 for analysis. With the BrainVISA software, a 3D rendering of the VOIs is displayed.



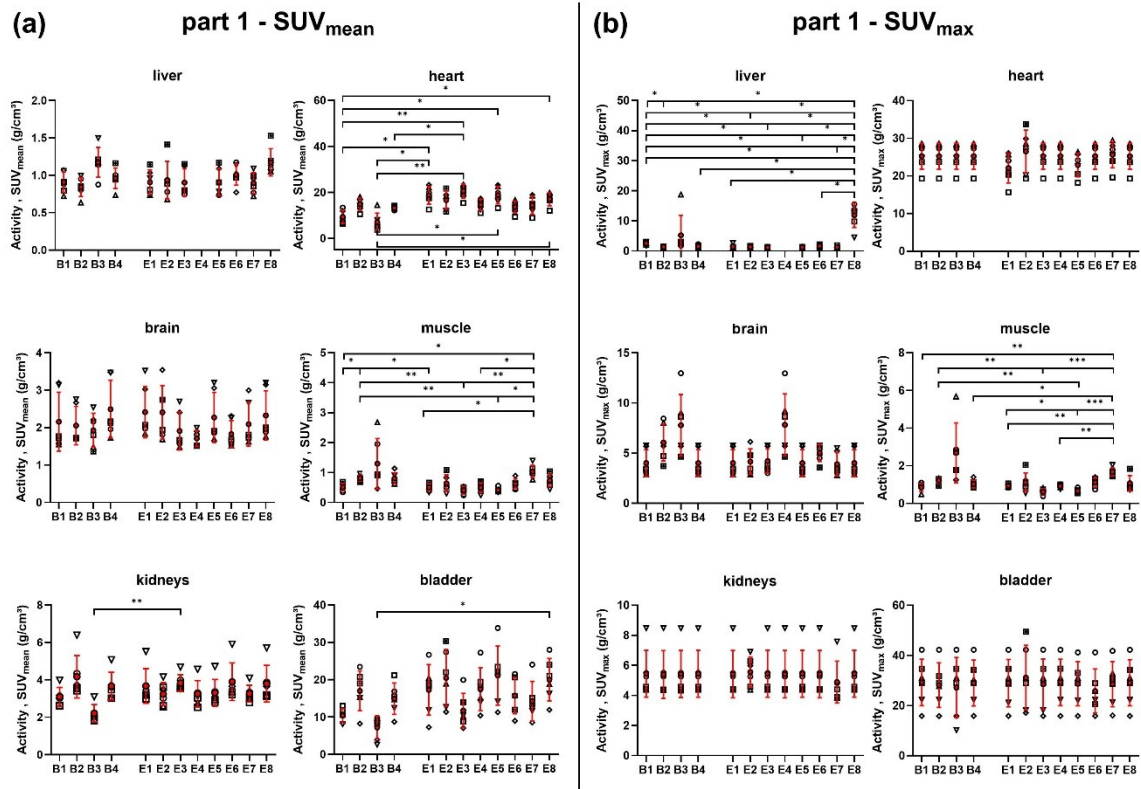
494
 495 **Fig. 2** VOI sizes delineated by the beginner (n=4, open triangle) or expert (n=8, open circle)
 496 group on the (a) $[^{18}\text{F}]\text{FDG}$ -PET-only (n=6) and (b) $[^{18}\text{F}]\text{FDG}$ -PET/CT (n=7) images. In (c), the

497 VOI sizes after the standardization procedure are shown. The mean values \pm standard
498 deviations are displayed. (* $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$; **** $p < 0.0001$; two-way ANOVA
499 followed by Bonferroni multiple comparisons test). The coefficient of variation (%CV) values
500 for each organ are provided separately for beginners and experts. The bold text marks lower
501 %CV values for beginners or experts. (Abbreviations used: bladder – urinary bladder, muscle
502 glut max – muscle gluteus maximus, bladder bottom – bottom of the urinary bladder, bladder
503 max fill – urinary bladder at maximum fill).

504

514 maximus, bladder bottom – bottom of the urinary bladder, bladder max fill – urinary bladder at
 515 maximum fill).

516



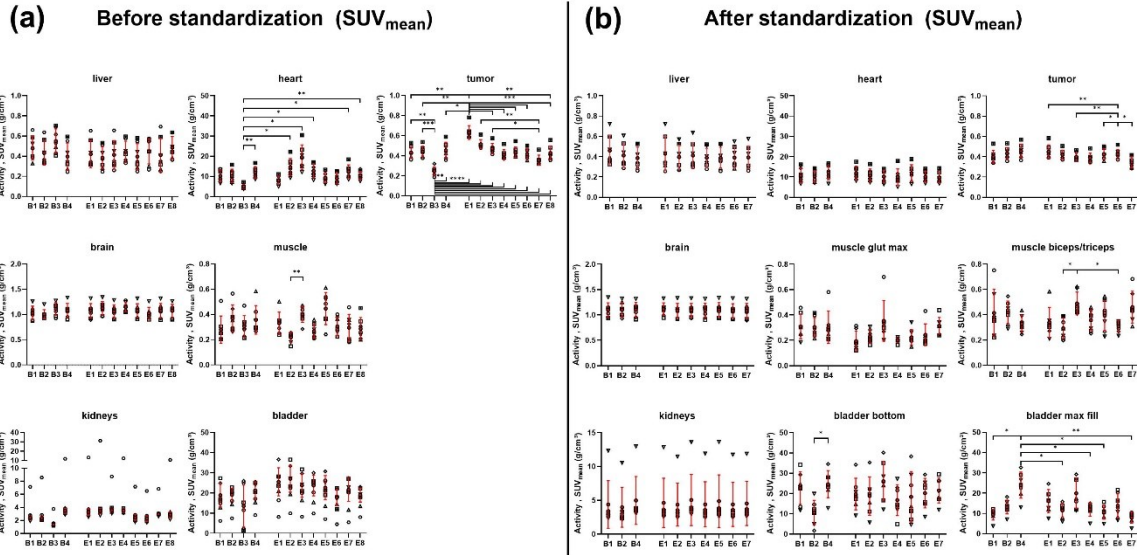
517

518 **Fig. 4** (a) SUV_{mean} and (b) SUV_{max} analysis as a function of beginner or expert observers for
 519 $[^{18}F]FDG$ -PET-only data from the liver, heart, brain, muscle, mean kidney, and urinary bladder.
 520 Individual values, as well as the mean \pm standard deviation, are displayed. B1-4: beginners 1
 521 to 4; E1-8: experts 1 to 8. Differences between individual observers were assessed by Brown-
 522 Forsythe and Welch ANOVA followed by Dunnett's T3 multiple comparisons test (* $p < 0.05$;
 523 ** $p < 0.01$; *** $p < 0.001$; **** $p < 0.0001$). Expert 4 did not analyze the liver. (Abbreviations used:
 524 bladder – urinary bladder).

525

526

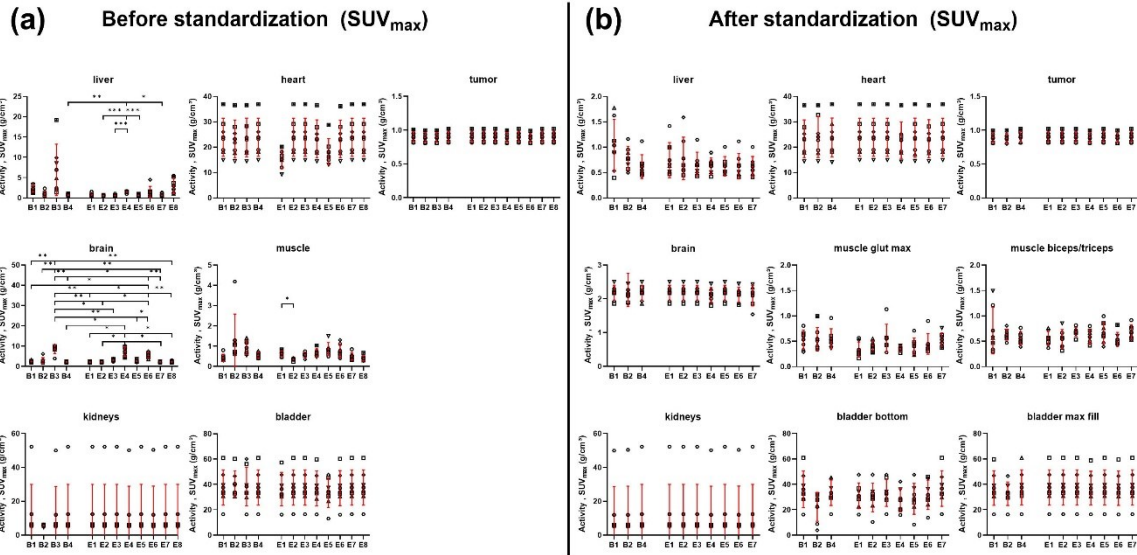
527



528

529 **Fig. 5** SUV_{mean} analysis as a function of beginner or expert observers from [¹⁸F]FDG-PET/CT
 530 data for the selected regions **(a)** before and **(b)** after standardization. Individual values, as well
 531 as the mean ± standard deviation, are displayed. B1-4: beginners 1 to 4; E1-8: experts 1 to 8.
 532 Differences between individual observers were assessed by Brown-Forsythe and Welch
 533 ANOVA followed by Dunnett's T3 multiple comparisons test (*p<0.05; **p<0.01; ***p<0.001;
 534 ****p<0.0001). The analyses of observers B3 and E8 were not included in the standardized
 535 [¹⁸F]FDG-PET/CT analysis because they were not applicable for the standardized protocol.
 536 (Abbreviations used: bladder – urinary bladder, muscle glut max – muscle gluteus maximus,
 537 bladder bottom – bottom of the urinary bladder, bladder max fill – urinary bladder at maximum
 538 fill).

539



540

541 **Fig. 6** SUV_{max} analysis as a function of beginner or expert observers from $[^{18}F]$ FDG-PET/CT
 542 data for the selected regions (a) before and (b) after standardization. Individual values, as well
 543 as the mean \pm standard deviation, are displayed. B1-4: beginners 1 to 4; E1-8: experts 1 to 8.
 544 Differences between individual observers were assessed by Brown-Forsythe and Welch
 545 ANOVA followed by Dunnett's T3 multiple comparisons test (* $p < 0.05$; ** $p < 0.01$; *** $p < 0.001$;
 546 **** $p < 0.0001$). The analyses of observers B3 and E8 were not included in the standardized
 547 $[^{18}F]$ FDG-PET/CT analysis because they were not applicable for the standardized protocol.
 548 (Abbreviations used: bladder – urinary bladder, muscle glut max – muscle gluteus maximus,
 549 bladder bottom – bottom of the urinary bladder, bladder max fill – urinary bladder at maximum
 550 fill).