

Variability and Repeatability of Quantitative Uptake Metrics in ^{18}F -FDG PET/CT of Non–Small Cell Lung Cancer: Impact of Segmentation Method, Uptake Interval, and Reconstruction Protocol

Mingzan Zhuang^{1,2}, David Vázquez García¹, Gerbrand M. Kramer³, Virginie Frings³, E.F. Smit⁴, Rudi Dierckx¹, Otto S. Hoekstra³, and Ronald Boellaard^{1,3}

¹Department of Nuclear Medicine and Molecular Imaging, University Medical Center Groningen, University of Groningen, Groningen, The Netherlands; ²The Key Laboratory of Digital Signal and Image Processing of Guangdong Province, Shantou University, Shantou, China; ³Department of Radiology and Nuclear Medicine, VU University Medical Center, Amsterdam, The Netherlands; and ⁴Department of Pulmonary Disease, VU University Medical Center, Amsterdam, The Netherlands

There is increased interest in various new quantitative uptake metrics beyond SUV in oncologic PET/CT studies. The purpose of this study was to investigate the variability and test–retest ratio (TRT) of metabolically active tumor volume (MATV) measurements and several other new quantitative metrics in non–small cell lung cancer using ^{18}F -FDG PET/CT with different segmentation methods, user interactions, uptake intervals, and reconstruction protocols. **Methods:** Ten patients with advanced non–small cell lung cancer received 2 series of 2 whole-body ^{18}F -FDG PET/CT scans at 60 min after injection and at 90 min after injection. PET data were reconstructed with 4 different protocols. Eight segmentation methods were applied to delineate lesions with and without a tumor mask. MATV, SUV_{max} , SUV_{mean} , total lesion glycolysis, and intraleSIONal heterogeneity features were derived. Variability and repeatability were evaluated using a generalized-estimating-equation statistical model with Bonferroni adjustment for multiple comparisons. The statistical model, including interaction between uptake interval and reconstruction protocol, was applied individually to the data obtained from each segmentation method. **Results:** Without masking, none of the segmentation methods could delineate all lesions correctly. MATV was affected by both uptake interval and reconstruction settings for most segmentation methods. Similar observations were obtained for the uptake metrics SUV_{max} , SUV_{mean} , total lesion glycolysis, homogeneity, entropy, and zone percentage. No effect of uptake interval was observed on TRT metrics, whereas the reconstruction protocol affected the TRT of SUV_{max} . Overall, segmentation methods showing poor quantitative performance in one condition showed better performance in other (combined) conditions. For some metrics, a clear statistical interaction was found between the segmentation method and both uptake interval and reconstruction protocol. **Conclusion:** All segmentation results need to be reviewed critically. MATV and other quantitative uptake metrics, as well as their TRT, depend on segmentation method, uptake interval, and reconstruction protocol. To obtain quantitative reliable metrics, with good TRT performance, the optimal segmentation method depends on local imaging procedure, the PET/CT system,

or reconstruction protocol. Rigid harmonization of imaging procedure and PET/CT performance will be helpful in mitigating this variability.

Key Words: variability; repeatability; segmentation method; non-small cell lung cancer; positron emission tomography imaging

J Nucl Med 2019; 60:600–607

DOI: 10.2967/jnumed.118.216028

PET imaging with ^{18}F -FDG is extensively used in oncology for diagnosis, staging, prognosis and response monitoring. Various quantitative metrics in PET imaging, such as metabolically active tumor volume (MATV), SUV, and intraleSIONal uptake heterogeneity, have been developed as indicators to quantify glucose metabolism in malignant tumors (1,2). However, the variability in segmentation techniques, user interaction during the segmentation, and imaging acquisition protocols presents particular challenges for consistently and accurately obtaining quantitative metrics.

Over the last 20 years, several segmentation methods have been developed and investigated in different tumor types, presenting large variability in terms of delineation accuracy and user interaction (3,4). As reported by the American Association of Physicists in Medicine (AAPM), validation for most published segmentation methods is either insufficient or inconsistent (5). Besides, although repeatability of quantitative metrics in PET imaging has been extensively explored (6,7), several recent studies have presented conflicting results. Tixier et al. (8) reported poor repeatability of various textural features in esophageal cancer, with only a few features being sufficiently reliable. However, van Velden et al. (9) found that most metrics had similar or better repeatability than SUV in non–small cell lung cancer (NSCLC). It is unclear whether these apparently conflicting results are caused by differences in tumor types, segmentation methodologies, applied imaging protocols, or a combination of these factors. The systematic comparison of the performances of a range of oncologic image–derived PET metrics obtained using different segmentation methods and imaging protocols is highly desirable.

Therefore, to understand the potential interactions among these aspects, we studied the variability of a representative set of frequently

Received Jun. 11, 2018; revision accepted Sep. 24, 2018.

For correspondence or reprints contact: Ronald Boellaard, Radiology and Nuclear Medicine, VU University Medical Centre, De Boelelaan 1117, 1081 HV Amsterdam, The Netherlands.

E-mail: r.boellaard@vumc.nl

Published online Nov. 2, 2018.

COPYRIGHT © 2019 by the Society of Nuclear Medicine and Molecular Imaging.

used quantitative metrics for NSCLC PET imaging as a function of segmentation method, user interaction, uptake interval, and reconstruction protocol, along with the repeatability of these metrics.

MATERIALS AND METHODS

Patients

We reanalyzed PET/CT scans from a prospective single-center study on 10 patients with advanced NSCLC who underwent double ^{18}F -FDG PET/CT scans at VU University Medical Center. Patient characteristics are listed in Table 1 and were previously described (9). All patients gave written informed consent before enrollment. This study was approved by the Medical Ethics Review Committee of the VU University Medical Center and was registered in the Dutch trial register (www.trialregister.nl, NTR3508).

Data Acquisition and Reconstruction

Patients fasted for at least 6 h before administration of ^{18}F -FDG. All scans were performed using an Ingenuity TF PET/CT scanner (Philips Healthcare). Two whole-body (i.e., skull vertex to mid-thigh) PET/CT scans were performed, one at 60 min after injection and another at 90 min. For each PET scan, a low-dose CT scan (120 kVp, 50 mAs) was also obtained. The same procedure was repeated within 3 d after the first examinations. For 2 patients, the 90-min PET scans were not collected because the patients could not comply with the long duration of the scan.

All PET images were reconstructed using 4 different protocols with necessary corrections (e.g., attenuation, scatter, random, and normalization), which included a vendor-provided body reconstruction protocol (ING), an EANM Research Ltd (EARL)-compliant reconstruction (10), a postreconstruction resolution model with 1 iteration (PSF1), and the same protocol with 2 iterations (PSF2). The matrix size of all reconstructed images was 144×144 with an isotropic voxel size of 4 mm (supplemental data, available at <http://jnm.snmjournals.org>).

Delineation Methods

Lesions were identified by a nuclear physician. For each lesion, 8 automated segmentation methods were applied (Supplemental Table 1): a method for automated segmentation using an active contour model (MASAC) (11), an affinity propagation algorithm (AP) (12), a contourlet-based active contour algorithm (CAC) (13), the contrast-oriented thresholding method (ST) of Schaefer et al. (14), segmentation using 41% of the maximum tumor value as a threshold (41MAX) (15), segmentation using 50% of the peak tumor value as a threshold, adapted for local background (A50P) (15), segmentation using an SUV of 2.5 as a threshold (SUV25), and segmentation using an SUV of 4.0 as a threshold (SUV40).

Each segmentation method was applied with and without a manually defined tumor mask, restricting the region growing to remain within the mask.

TABLE 1
Patient Characteristics

Characteristic	Median	Scan 1	Scan 2	<i>P</i>
Patients	10			
Men	6			
Age (y)	61 (45–66)			
Tumor type (histology)				
Adenocarcinoma	7			
Squamous cell carcinoma	3			
Tumor stage				
IIIb	3			
IV	7			
Tumor location: lung				
Lesions per patient	2 (1–13)			
Weight (kg)		76 (57–110)	75 (57–113)	0.781
Number of patients				
60-min uptake interval		10	10	
90-min uptake interval		10	8*	
Number of lesions				
60-min uptake interval		26	26	
90-min uptake interval		26	18*	
Injected activity (MBq)		248 (194–377)	238 (192–392)	0.800
Scan start time (min)				
60-min uptake interval		61 (59–67)	60 (60–63)	0.293
90-min uptake interval		92 (90–97)	90 (90–95)	0.219

*For 2 patients, 90-min PET scans were not collected because of patients' inability to comply with scan duration.

Qualitative data are expressed as number; continuous data are expressed as median followed by range in parentheses. *P* values are from Wilcoxon signed-rank test.

Performance Evaluations

The index “out-of-mask” (OM) was included as a metric of segmentation failure:

$$OM = 100 \times N_{\text{outside}} / N_{\text{total}}, \quad \text{Eq. 1}$$

where N_{outside} is the number of cases for which the segmentation method without a mask generated a segmentation expanding beyond the predefined tumor mask, and N_{total} is the total number of PET tumor segmentations. Thus, the out-of-mask index reflects the ability of a segmentation method to automatically segment the tumor without spatial constraints (i.e., without a mask). The lower the number, the more successful the method was to generate a tumor segmentation without the inclusion of nonlesioned ^{18}F -FDG-avid areas, or without mislocalization of the segmentation (e.g., jumping to a wrong location, such as a different tumor, kidney, bladder, myocardium, or liver).

Quantitative Uptake Metrics

The quantitative metrics evaluated in this study were MATV, SUV_{max} , SUV_{mean} , total lesion glycolysis (TLG), and several textural intratumor heterogeneity features. These features included a global heterogeneity indicator (i.e., area under the curve of the cumulative intensity histogram, CIH_{AUC}) (16), and some local heterogeneity features, such as homogeneity, entropy, dissimilarity, high-intensity emphasis (HIE), and zone percentage (ZP). These features were selected because of their reproducibility and robustness (8,16,17). MATV, SUV_{max} , SUV_{mean} , TLG, and CIH_{AUC} were calculated with in-house software, whereas local heterogeneity features were obtained with the Pyradiomics package (18). All features were extracted from the original images, without the application of any postprocessing (e.g., rebinning or filtering). Detailed information about the implementation of these metrics are presented as Supplemental Equations 1–5.

Repeatability Evaluations

Repeatability of the metrics between the 2 scans was calculated as the test–retest ratio ($\text{TRT}_{\text{metric}}$):

$$\text{TRT}_{\text{metric}} = (\text{Metric}_{\text{scan1}} - \text{Metric}_{\text{scan2}}) / [(\text{Metric}_{\text{scan1}} + \text{Metric}_{\text{scan2}}) / 2], \quad \text{Eq. 2}$$

where $\text{Metric}_{\text{scan1}}$ and $\text{Metric}_{\text{scan2}}$ are the metric values obtained from the first and second scans, respectively.

Statistical Analysis

Statistical analysis was performed using SPSS Statistics 24.0 software (IBM). The generalized-estimating-equation model was used to account for repeated measurements and missing data. The independent working correlation matrix was selected for analysis, with an identity link function. The natural log transformation was applied

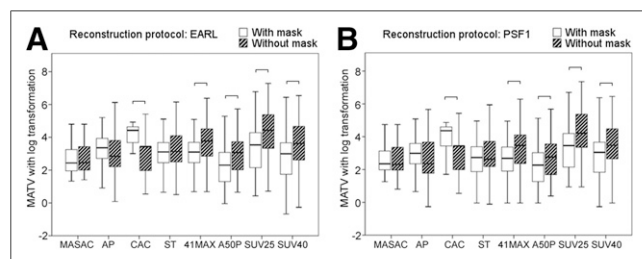


FIGURE 1. Box-and-whisker plots of MATV at 60 min for each segmentation method. For display purposes, outliers identified as $1.5 \times$ interquartile range were removed from plot (whiskers). Statistically significant differences are marked with horizontal line.

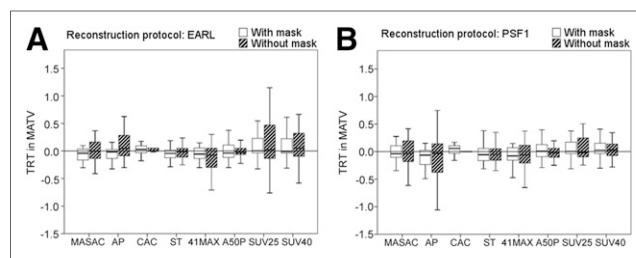


FIGURE 2. Box-and-whisker plots of TRT_{MATV} at 60 min for each segmentation method. For display purposes, outliers identified as $1.5 \times$ interquartile range were removed from plot (whiskers). No statistically significant differences were found.

to MATV, SUV_{max} , SUV_{mean} , TLG, and HIE to obtain normally distributed data.

To assess the influence of uptake interval and reconstruction protocol, the specific metric was selected as the dependent outcome in the generalized-estimating-equation model; the patient, scan, uptake interval, and reconstruction protocol were included as independent variables, along with the interaction effect between uptake interval and reconstruction. Similar settings were also used for the $\text{TRT}_{\text{metric}}$, excluding the “scan” variable. A post hoc pairwise comparison was performed when the test of model effect was shown to be significant, applying Bonferroni adjustment for multiple comparisons present in the test. P values of less than 0.05 were considered to be significant.

To explore the relationship between MATV and the other metrics, MATV was set as the dependent outcome, with each other metric included independently as the main effect in the generalized-estimating-equation model, and corrected for other factors such as patient, scan, uptake interval, reconstruction protocol, interaction of the uptake interval with the metric, interaction of the reconstruction protocol and the metric, and interaction of uptake interval, reconstruction protocol, and the metric. Similarly, the correlation between the TRT_{MATV} and TRT of the other metrics was also investigated. Moreover, scatterplots were also used to explore the relationships of TRT_{MATV} with MATV and SUV_{max} .

RESULTS

Tumor Mask

For 41MAX, A50P, SUV25, and SUV40, the use of a mask resulted in a significantly smaller (12%–22%) MATV, whereas CAC showed a significantly larger (35%) MATV with masking (Fig. 1). However, applying a tumor mask did not improve the MATV’s repeatability in most segmentation results (Fig. 2). Similar results were also found with the other reconstruction protocols.

As shown in Table 2, A50P displayed fewer incorrect segmentation results (30%) than the other segmentation methods. In general, CAC and SUV25 showed the worst out-of-mask index results (CAC, 77% at 60 min; SUV25, 86% at 90 min). Because no segmentation method correctly delineated all lesions without a mask, we used the results derived from the segmentation with a tumor mask for further analysis.

Uptake Interval and Reconstruction Protocol

Overall, MATV at a 90-min uptake interval was larger than at 60 min for CAC, A50P, SUV25, and SUV40 but smaller for MASAC, AP, ST, and 41MAX (Fig. 3), specially affecting those protocols with lower spatial resolution (EARL and ING). These observed differences were statistically significant for all methods, with the exception of MASAC, CAC, and A50P (Table 3; Supplemental Table 2). For example, direct comparison (i.e., without log

TABLE 2
Comparison of Out-of-Mask Index for Each Segmentation Method with Different Uptake Intervals and Reconstruction Protocols

Uptake interval	Reconstruction	MASAC	AP	CAC	ST	41MAX	A50P	SUV25	SUV40
60 min (%)	EARL	54	65	77	64	50	31	67	42
	ING	54	64	75	62	44	29	71	42
	PSF1	54	52	77	56	40	27	75	40
	PSF2	54	50	77	56	35	29	73	42
90 min (%)	EARL	53	64	72	56	42	31	81	53
	ING	53	64	78	56	42	33	81	56
	PSF1	53	56	78	53	47	31	86	61
	PSF2	53	50	72	50	33	28	86	67
Total (%)		53	58	76	56	42	30	78	50

transformation) of MATV in EARL reconstructed data showed a median increase of 7% (interquartile range, 1%–13%) for A50P, SUV25, and SUV40 delineations versus a median decrease of –4% (interquartile range, –8% to 2%) for MASAC, AP, CAC, ST, and 41MAX. In addition, except for SUV40, most segmentation methods showed a slightly smaller MATV with reconstruction protocols that provided higher spatial resolutions (in ascending order: EARL, ING, PSF1, and PSF2) at both uptake intervals.

With each segmentation method, SUV_{max} , SUV_{mean} , TLG, entropy, and ZP increased significantly ($P < 0.001$) from 1% to 6% at the 90-min uptake interval as compared with the 60-min interval. In

CAC, homogeneity was independent of uptake interval, whereas all other segmentation methods showed significantly lower homogeneity (2%, $P < 0.001$) at 90 min than at 60 min (Table 3; Supplemental Table 2).

SUV_{max} , SUV_{mean} , TLG, homogeneity, entropy, and ZP were significantly affected by the reconstruction protocol regardless of the segmentation method. For most segmentation methods, SUV_{max} , SUV_{mean} , entropy, and ZP increased from 1% to 6% at reconstruction protocols with higher spatial resolution, whereas CIH_{AUC} and homogeneity decreased slightly (1%) in these cases (Supplemental Table 3). Compared with other metrics, dissimilarity and HIE were hardly affected by either uptake interval or reconstruction protocol.

There were significant interaction effects for homogeneity; that is, it correlated not only with uptake interval and reconstruction protocol but also with their combinations. For most segmentation methods, the PSF2 reconstruction protocol at 90 min of uptake showed the lowest homogeneity, whereas the EARL protocol showed the highest homogeneity at 60 min, except for ST and 41MAX.

Repeatability: Effect of Uptake Interval and Reconstruction Protocol

Uptake interval had no effect on TRT for any metric or segmentation method, whereas the used reconstruction protocol affected TRT in SUV_{max} for each segmentation method (Fig. 4; Table 4). In general, the TRT in SUV_{max} was worse with higher-spatial-resolution reconstruction protocols. There were no evident interactions between uptake interval and reconstruction protocol for any of the metrics or segmentation methods.

Relationship Between MATV and Other Metrics as Well as Their Repeatability

There was a significant relationship between MATV and other metrics (Supplemental Table 7). Similarly, TRT_{MATV} strongly correlated with TRT for SUV_{mean} , TLG, CIH_{AUC} , and ZP (Supplemental Table 8). These relationships were also affected by the different uptake intervals and the reconstruction protocols.

Relationship Between TRT_{MATV} and MATV or SUV_{max}

For most segmentation methods, the repeatability of MATV was better at larger MATVs and higher SUVs (Figs. 5 and 6, respectively). A similar trend for the relationship between TRT_{MATV} and

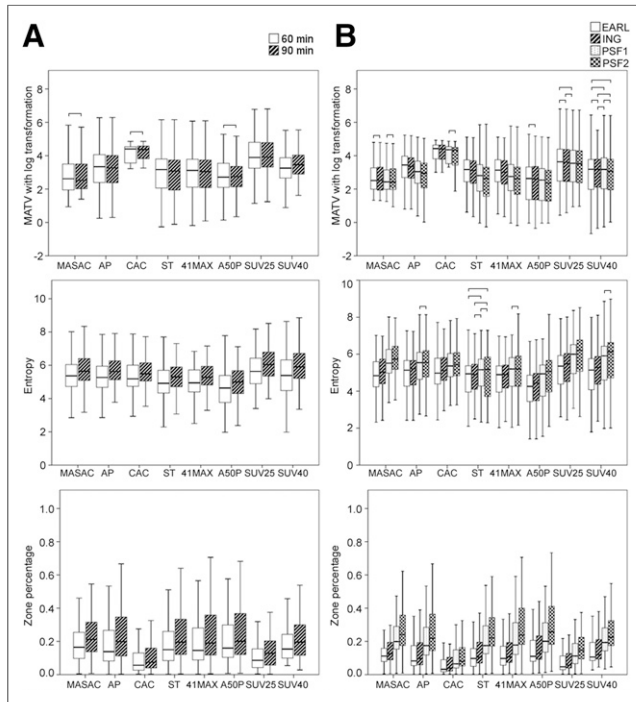


FIGURE 3. Box-and-whisker plots for various metrics as function of uptake interval (A) and reconstruction protocol (B). For display purposes, outliers identified as $1.5 \times$ interquartile range were removed from plot (whiskers). Comparisons without statistically significant differences are marked with horizontal line.

TABLE 3
Generalized-Estimating-Equation Model: Significance Results (*P* Values) for All Metrics Tested

Variable	Metric	MASAC	AP	CAC	ST	41MAX	A50P	SUV25	SUV40
Uptake interval	MATV	0.215	0.002	0.386	0.003	0.024	0.213	0.001	0.005
	SUV _{max}	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	SUV _{mean}	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	TLG	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	CIH _{AUC}	0.433	0.569	0.001	0.065	0.143	0.115	<0.001	<0.001
	Homogeneity	<0.001	<0.001	0.134	<0.001	<0.001	<0.001	<0.001	<0.001
	Entropy	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	Dissimilarity	0.282	0.134	0.529	0.549	0.641	0.666	0.396	0.165
	HIE	0.086	0.201	0.011	0.441	0.239	0.390	<0.001	0.085
	ZP	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
Reconstruction	MATV	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	0.122
	SUV _{max}	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	SUV _{mean}	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	TLG	<0.001	<0.001	0.006	<0.001	<0.001	<0.001	<0.001	0.001
	CIH _{AUC}	<0.001	0.156	<0.001	0.401	0.002	<0.001	<0.001	<0.001
	Homogeneity	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	Entropy	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	Dissimilarity	0.365	<0.001	0.005	0.002	0.009	0.143	0.716	0.138
	HIE	0.066	0.082	<0.001	0.927	0.524	<0.001	<0.001	0.297
	ZP	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
Uptake interval × reconstruction	MATV	0.624	0.210	0.373	0.233	0.138	0.158	0.017	0.223
	SUV _{max}	0.865	0.966	0.966	0.708	0.775	0.775	0.775	0.412
	SUV _{mean}	0.730	0.722	0.530	0.104	0.689	0.243	0.042	0.086
	TLG	0.552	0.102	0.620	0.294	0.071	0.230	0.025	0.241
	CIH _{AUC}	0.535	0.192	0.351	0.117	0.122	0.324	0.031	0.491
	Homogeneity	0.008	<0.001	0.018	0.050	0.225	0.021	0.002	0.021
	Entropy	0.107	0.030	0.022	0.090	0.268	0.004	0.002	0.339
	Dissimilarity	0.364	0.161	0.505	0.715	0.951	0.632	0.852	0.956
	HIE	0.793	0.349	0.374	0.040	0.099	0.207	0.429	0.550
	ZP	0.218	0.078	0.007	0.040	0.215	0.041	0.001	0.119

Statistically significant results (*P* < 0.05) are presented in bold.

MATV or SUV_{max} was also observed at other uptake intervals and for other reconstructions.

DISCUSSION

Our study showed that segmentation methods are influenced by different user interactions, uptake intervals, and reconstruction protocols, suggesting that all segmentation results need to be reviewed critically. User interaction during the segmentation process is often required in medical imaging (19,20). In our study, no segmentation method could delineate all lesions correctly without a tumor mask, indicating the necessity of manually defining a tumor mask, especially for tumors adjacent to high-activity areas.

In our study, MASAC, CAC, and A50P were statistically independent of the uptake interval (i.e., 60 vs. 90 min) in MATV, whereas SUV25 and SUV40 showed larger MATVs at 90 min than at 60 min after ¹⁸F-FDG administration (6% and 10%,

respectively). Because lesional uptake was higher at 90 than 60 min, these 2 segmentation methods, taking the absolute SUVs as threshold values, tended to generate larger MATVs at 90 min, as occurred in our experiments. However, 41MAX, using relative thresholds, show a slight decrease in MATV at an increased uptake interval. Moreover, MATV obtained with most segmentation methods depends on the reconstruction protocol, and thus, these methods require careful consideration when used in different clinical scenarios.

It has been reported that intralesional heterogeneity correlates with treatment outcome (21). However, regardless of the uncertainties in segmentation methods, differences in acquisition protocols also result in changes in image quality, thus influencing the results for these extracted features (22,23). We found that intralesional heterogeneity increases with uptake interval or spatial image resolution (in ascending order: EARL, ING, PSF1, and PSF2),

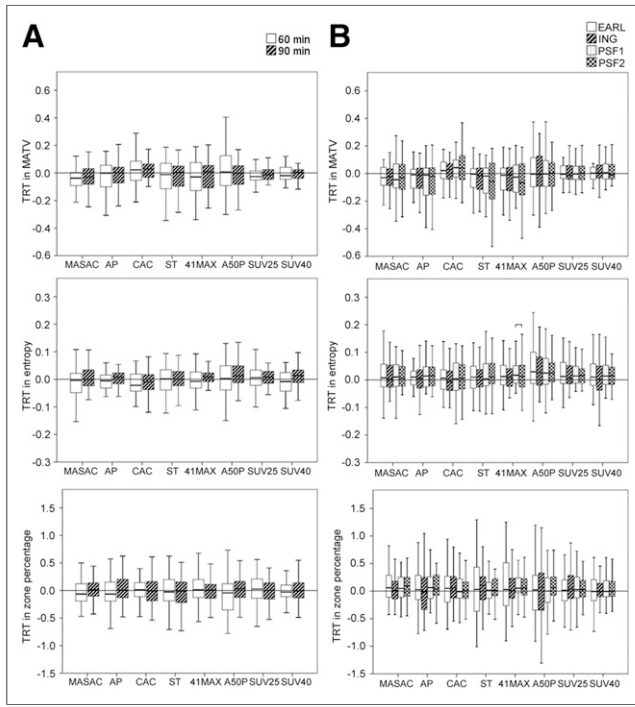


FIGURE 4. Box-and-whisker plots for TRT of various metrics as function of uptake interval (A) and reconstruction protocol (B). For display purposes, outliers identified as $1.5 \times$ interquartile range were removed from plot (whiskers). Statistically significant differences are marked with horizontal line.

presented as the decrease in CIH_{AUC} and homogeneity and the increase in entropy and ZP, although dissimilarity and HIE showed less association with uptake interval or reconstruction. The lower the CIH_{AUC} or homogeneity, the higher the heterogeneity of the image, whereas the higher the entropy or ZP, the more details an image carries and the more heterogeneous are the tumor features in the image. Similar results were also found by Lasnon et al. (1), who showed that PSF images resulted in higher heterogeneity than EARL-compliant images.

We found that the repeatability of most metrics was independent of the tracer uptake interval and reconstruction protocol, for each segmentation method evaluated. Moreover, MATV and other metrics were highly correlated, as well as their TRTs. This finding may seem to be inconsistent with the results of Hatt et al. (24), but we believe it can be explained by the use of different segmentation procedures and acquisition protocols. Moreover, to identify predictors of repeatability in MATV, the correlations of TRT_{MATV} with MATV and SUV_{max} were also investigated. We found that, in general, the repeatability of MATV was better with high values of MATV or SUV_{max} , suggesting that small lesions are more likely to be affected by variation in imaging procedures, consistent with our previous study (25).

As proposed by AAPM report 211, accuracy evaluation of segmentation methods is required for each PET scanning condition (5). Our study confirms and further supports this recommendation. We observed that MATV, as well as most of the other metrics, depends not only on the segmentation method but more specifically on its specific combination with uptake interval and

TABLE 4
Generalized-Estimating-Equation Model: Significance Results (*P* Values) for All TRT Metrics Tested

Variable	TRT	MASAC	AP	CAC	ST	41MAX	A50P	SUV25	SUV40
Reconstruction	MATV	0.742	0.216	0.256	0.246	0.182	0.901	0.179	0.792
	SUV_{max}	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001	<0.001
	SUV_{mean}	0.942	0.018	0.409	0.136	0.016	0.520	0.081	0.200
	TLG	0.614	0.602	0.166	0.371	0.652	0.928	0.288	0.737
	CIH_{AUC}	0.197	0.297	0.474	0.095	0.609	0.114	0.007	0.022
	Homogeneity	0.668	0.195	0.378	0.310	0.568	0.121	0.109	0.761
	Entropy	0.483	0.030	0.041	0.272	0.040	0.457	0.345	0.089
	Dissimilarity	0.493	0.388	0.605	0.261	0.550	0.221	0.562	0.417
	HIE	0.125	0.498	0.157	0.562	0.549	0.783	0.036	0.809
	ZP	0.197	0.452	0.775	0.340	0.602	0.235	0.069	0.731
Uptake interval \times reconstruction	MATV	0.307	0.082	0.579	0.524	0.924	0.121	0.058	0.908
	SUV_{max}	0.896	0.372	0.372	0.156	0.367	0.367	0.367	0.351
	SUV_{mean}	0.166	0.022	0.697	0.798	0.346	0.196	0.035	0.663
	TLG	0.309	0.121	0.661	0.301	0.872	0.119	0.158	0.913
	CIH_{AUC}	0.174	0.048	0.332	0.358	0.733	0.113	0.313	0.602
	Homogeneity	0.965	0.250	0.938	0.038	0.009	0.248	0.077	0.367
	Entropy	0.939	0.422	0.198	0.331	0.023	0.840	0.226	0.540
	Dissimilarity	0.534	0.926	0.902	0.362	0.712	0.572	0.421	0.567
	HIE	0.045	0.122	0.342	0.759	0.635	0.225	0.047	0.415
	ZP	0.530	0.585	0.948	0.467	0.523	0.403	0.795	0.347

Statistically significant results ($P < 0.05$) are presented in bold.

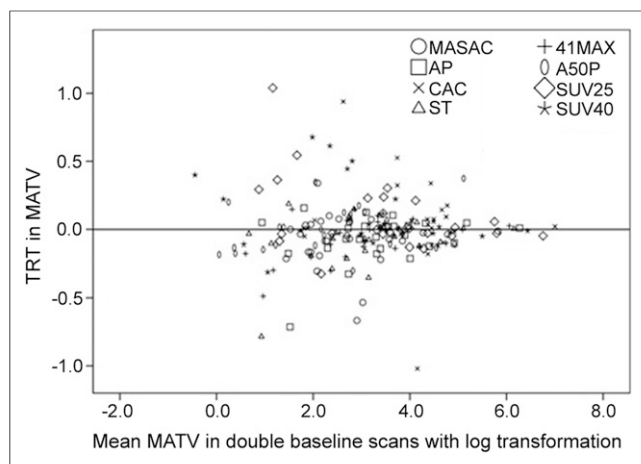


FIGURE 5. TRT_{MATV} as function of MATV with EARL reconstruction protocol at 60 min for each segmentation method.

reconstruction protocol. In other words, methods and procedures that may work well under one condition may be outperformed by other methods under different conditions. Therefore, it seems that the selection of the best segmentation method is highly dependent on the imaging procedures and conditions at hand, confirming the AAPM recommendation to evaluate performance for each scanning condition. Despite the publication of strict imaging guidelines (10,26), there remains considerable variability in imaging procedures. To some extent, these are mitigated by scanner accreditation programs (27), but residual variability will likely remain and require implementation of the AAPM report 211 recommendations.

The absence of ground truth in our study does not allow the accuracy of measured values to be assessed. In addition, although numerous data were included in our study to explore their interactions, these data were derived from ^{18}F -FDG PET images from only 10 NSCLC patients, which may not be sufficient to fully demonstrate their relationships in other clinical scenarios. Therefore, further studies are needed to establish a benchmark to evaluate their accuracy under different conditions.

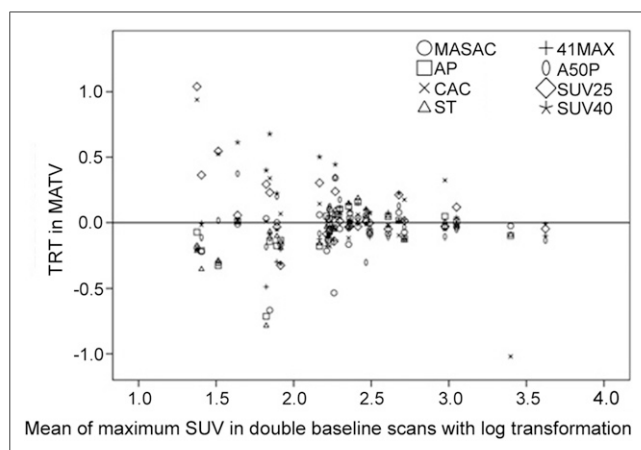


FIGURE 6. TRT_{MATV} as function of SUV_{max} with EARL reconstruction protocol at 60 min for each segmentation method.

CONCLUSION

Quantitative results derived from ^{18}F -FDG PET/CT studies on NSCLC patients show that all segmentation results need to be critically reviewed and that MATV, and other quantitative metrics, depend on segmentation method, uptake interval, and reconstruction protocol. Methods that perform well under one condition may not be suitable under different circumstances or studies. These interactions also suggest that to obtain reliable quantitative metrics with a good TRT performance, the optimal segmentation method depends on the local imaging procedures, PET/CT systems, or reconstruction protocols used. Rigid harmonization of imaging procedures and PET/CT performance will be helpful in mitigating this variability (28–30).

DISCLOSURE

This work was supported by an Open Grant (2014GDDSIPL-06) from the Key Laboratory of Digital Signal and Image Processing of Guangdong Province, Shantou University. No other potential conflict of interest relevant to this article was reported.

ACKNOWLEDGMENTS

We thank Prof. Habib Zaidi, Prof. Qingchun Qiu, and Zemian Chen for their assistance during the research.

REFERENCES

- Lasnon C, Majdoub M, Lavigne B, et al. F-18-FDG PET/CT heterogeneity quantification through textural features in the era of harmonisation programs: a focus on lung cancer. *Eur J Nucl Med Mol Imaging*. 2016;43:2324–2335.
- Vanhove K, Liesbet M, Louis E, Thomeer M, Adriaensens P, Boellaard R. Prognostic value of quantitative FDG PET/CT uptake metrics in NSCLC [abstract]. *J Nucl Med*. 2015;56(suppl 3):177.
- Zaidi H, El Naqa I. PET-guided delineation of radiation therapy treatment volumes: a survey of image segmentation techniques. *Eur J Nucl Med Mol Imaging*. 2010;37:2165–2187.
- Foster B, Bagci U, Mansoor A, Xu Z, Mollura DJ. A review on segmentation of positron emission tomography images. *Comput Biol Med*. 2014;50:76–96.
- Hatt M, Lee J, Schmidlein CR, et al. Classification and evaluation strategies of auto-segmentation approaches for PET: report of AAPM Task Group No. 211. *Med Phys*. 2017;44:e1–e42.
- Lodge MA. Repeatability of SUV in oncologic F-18-FDG PET. *J Nucl Med*. 2017;58:523–532.
- Desseroit MC, Tixier F, Weber WA, et al. Reliability of PET/CT shape and heterogeneity features in functional and morphologic components of non-small cell lung cancer tumors: a repeatability analysis in a prospective multicenter cohort. *J Nucl Med*. 2017;58:406–411.
- Tixier F, Hatt M, Le Rest CC, Le Pogam A, Corcos L, Visvikis D. Reproducibility of tumor uptake heterogeneity characterization through textural feature analysis in F-18-FDG PET. *J Nucl Med*. 2012;53:693–700.
- van Velden FH, Kramer GM, Frings V, et al. Repeatability of radiomic features in non-small-cell lung cancer [^{18}F]FDG-PET/CT studies: impact of reconstruction and delineation. *Mol Imaging Biol*. 2016;18:788–795.
- Boellaard R, Delgado-Bolton R, Oyen WJ, et al. FDG PET/CT: EANM procedure guidelines for tumour imaging: version 2.0. *Eur J Nucl Med Mol Imaging*. 2015;42:328–354.
- Zhuang M, Dierckx RA, Zaidi H. Generic and robust method for automatic segmentation of PET images using an active contour model. *Med Phys*. 2016;43:4483.
- Foster B, Bagci U, Ziyue X, et al. Segmentation of PET images for computer-aided functional quantification of tuberculosis in small animal models. *IEEE Trans Biomed Eng*. 2014;61:711–724.
- Abdoli M, Dierckx RAJO, Zaidi H. Contourlet-based active contour model for PET image segmentation. *Med Phys*. 2013;40:082507.
- Schaefer A, Kremp S, Hellwig D, Rube C, Kirsch CM, Nestle U. A contrast-oriented algorithm for FDG-PET-based delineation of tumour volumes for the radiotherapy of lung cancer: derivation from phantom measurements

- and validation in patient data. *Eur J Nucl Med Mol Imaging*. 2008;35:1989–1999.
15. Frings V, van Velden FH, Velasquez LM, et al. Repeatability of metabolically active tumor volume measurements with FDG PET/CT in advanced gastrointestinal malignancies: a multicenter study. *Radiology*. 2014;273:539–548.
 16. Tixier F, Vriens D, Cheze-Le Rest C, et al. Comparison of tumor uptake heterogeneity characterization between static and parametric ^{18}F -FDG PET images in non-small cell lung cancer. *J Nucl Med*. 2016;57:1033–1039.
 17. Hatt M, Tixier F, Cheze Le Rest C, Pradier O, Visvikis D. Robustness of intratumour ^{18}F -FDG PET uptake heterogeneity quantification for therapy response prediction in oesophageal carcinoma. *Eur J Nucl Med Mol Imaging*. 2013;40:1662–1671.
 18. van Griethuysen JJM, Fedorov A, Parmar C, et al. Computational radiomics system to decode the radiographic phenotype. *Cancer Res*. 2017.
 19. Olabarriaga SD, Smeulders AWM. Interaction in the segmentation of medical images: a survey. *Med Image Anal*. 2001;5:127–142.
 20. Ramkumar A, Dolz J, Kirisli HA, et al. User interaction in semi-automatic segmentation of organs at risk: a case study in radiotherapy. *J Digit Imaging*. 2016;29:264–277.
 21. Bashir U, Siddique MM, Mclean E, Goh V, Cook GJ. Imaging heterogeneity in lung cancer: techniques, applications, and challenges. *AJR*. 2016;207:534–543.
 22. Galavis PE, Hollensen C, Jallow N, Paliwal B, Jeraj R. Variability of textural features in FDG PET images due to different acquisition modes and reconstruction parameters. *Acta Oncol*. 2010;49:1012–1016.
 23. Yan J, Chu-Shern JL, Loi HY, et al. Impact of image reconstruction settings on texture features in ^{18}F -FDG PET. *J Nucl Med*. 2015;56:1667–1673.
 24. Hatt M, Majdoub M, Vallieres M, et al. F-18-FDG PET uptake characterization through texture analysis: investigating the complementary nature of heterogeneity and functional tumor volume in a multi-cancer site patient cohort. *J Nucl Med*. 2015;56:38–44.
 25. Frings V, de Langen AJ, Smit EF, et al. Repeatability of metabolically active volume measurements with ^{18}F -FDG and ^{18}F -FLT PET in non-small cell lung cancer. *J Nucl Med*. 2010;51:1870–1877.
 26. Harkness BA, Fahey FH. The current state of nuclear medicine physics training: findings of the AAPM/SNMMI Task Force. *J Nucl Med*. 2016;57:1146–1147.
 27. Graham MM. The Clinical Trials Network of the Society of Nuclear Medicine. *Semin Nucl Med*. 2010;40:327–331.
 28. Beyer T, Antoch G, Muller S, et al. Acquisition protocol considerations for combined PET/CT imaging. *J Nucl Med*. 2004;45(suppl):25S–35S.
 29. Graham MM, Badawi RD, Wahl RL. Variations in PET/CT methodology for oncologic imaging at U.S. academic medical centers: an imaging response assessment team survey. *J Nucl Med*. 2011;52:311–317.
 30. Velasquez LM, Boellaard R, Kollia G, et al. Repeatability of ^{18}F -FDG PET in a multicenter phase I study of patients with advanced gastrointestinal malignancies. *J Nucl Med*. 2009;50:1646–1654.