

Journal of Cachexia, Sarcopenia and Muscle

Validation of a deep learning model for automatic segmentation of skeletal muscle and adipose tissue on L3 abdominal CT images

--Manuscript Draft--

Manuscript Number:	
Full Title:	Validation of a deep learning model for automatic segmentation of skeletal muscle and adipose tissue on L3 abdominal CT images
Article Type:	Original Article
Corresponding Author:	David P J van Dijk, MD, MSc, PhD Maastricht University Medical Centre+: Maastricht Universitair Medisch Centrum+ Maastricht, Limburg NETHERLANDS
Corresponding Author Secondary Information:	
Corresponding Author's Institution:	Maastricht University Medical Centre+: Maastricht Universitair Medisch Centrum+
Corresponding Author's Secondary Institution:	
Corresponding Author E-Mail:	david.van.dijk@mumc.nl
First Author:	David P J van Dijk, MD, MSc, PhD
First Author Secondary Information:	
Order of Authors:	David P J van Dijk, MD, MSc, PhD Leroy F. Volmer Ralph Brecheisen Ross D. Dolan Adam S. Bryce David K. Chang Donald C. McMillan Jan H.M.B. Stoot Malcolm A. West Sander S. Rensen Andre Dekker Leonard Wee Steven W.M. Olde Damink Body Composition Collaborative
Order of Authors Secondary Information:	
Abstract:	<p>Background Body composition assessment using abdominal computed tomography (CT) images is increasingly applied in clinical and translational research. Manual segmentation of body compartments on L3 CT images is time-consuming and requires significant expertise. Robust high-throughput automated segmentation is key to assess large patient cohorts and ultimately, to support implementation into routine clinical practice. By training a deep learning neural network (DLNN) with several large trial cohorts and performing external validation on a large independent cohort, we aim to demonstrate the robust performance of our automatic body composition segmentation tool for future use in patients.</p> <p>Methods L3 CT images and expert-drawn segmentations of skeletal muscle, visceral adipose</p>

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

	<p>tissue, and subcutaneous adipose tissue of patients undergoing abdominal surgery were pooled (n = 3,187) to train a DLNN. The trained DLNN was then externally validated in a cohort with L3 CT images of patients with abdominal cancer (n = 2,535). Geometric agreement between automatic and manual segmentations was evaluated by computing two-dimensional Dice Similarity (DS). Agreement between manual and automatic annotations were quantitatively evaluated in the test set using Lin's Concordance Correlation Coefficient (CCC) and Bland-Altman's Limits of Agreement (LoA).</p> <p>Results</p> <p>The DLNN showed rapid improvement within the first 10,000 training steps and stopped improving after 38,000 steps. There was a strong concordance between automatic and manual segmentations with median DS for skeletal muscle, visceral adipose tissue, and subcutaneous adipose tissue of 0.97 (interquartile range, IQR: 0.95-0.98), 0.98 (IQR: 0.95-0.98), and 0.95 (IQR: 0.92-0.97), respectively. Concordance correlations were excellent: skeletal muscle 0.964 (0.959-0.968), visceral adipose tissue 0.998 (0.998-0.998), and subcutaneous adipose tissue 0.992 (0.991-0.993). Bland-Altman metrics (relative to approximate median values in parentheses) indicated only small and clinically insignificant systematic offsets : 0.23 HU (0.5%), 1.26 cm².m⁻² (2.8%), -1.02 cm².m⁻² (1.7%), and 3.24 cm².m⁻² (4.6%) for skeletal muscle average radiodensity, skeletal muscle index, visceral adipose tissue index, and subcutaneous adipose tissue index, respectively. Assuming the decision thresholds by Martin et al. for sarcopenia and low muscle radiation attenuation, results for sensitivity (0.99 and 0.98 respectively), specificity (0.87 and 0.98 respectively), and overall accuracy (0.93) were all excellent.</p> <p>Conclusion</p> <p>We developed and validated a deep learning model for automated analysis of body composition of patients with cancer. Due to the design of the DLNN, it can be easily implemented in various clinical infrastructures and used by other research groups to assess cancer patient cohorts or develop new models in other fields.</p>
<p>Suggested Reviewers:</p>	<p>Richard Skipworth, MD, PhD University of Edinburgh rskipworth@hotmail.com Expert in cachexia and body composition</p> <p>Oliver Bathe, MD University of Calgary Oliver.Bathe@albertahealthservices.ca Surgeon, expert on cachexia and body composition</p> <p>Vera Mazurak, PhD Professor, University of Alberta Oliver.Bathe@albertahealthservices.ca Expert on body composition and adipose tissue.</p> <p>Joost Klaase, MD, PhD Groningen University j.m.klaase@umcg.nl Surgeon, professor in body composition and prehabilitation.</p>
<p>Author Comments:</p>	<p>Re MS# JCSM-D-23-00160 Validation of a deep learning model for automatic segmentation of skeletal muscle and adipose tissue on L3 abdominal CT images.</p> <p>Dear Dr. Skipworth, Dear Dr. Anker and Dr. von Haehling,</p> <p>We would like to thank you and the reviewers for providing a critique of our work and granting us the opportunity to resubmit our work. We believe that the changes we have made (highlighted in bold and red text) will address the concerns raised. Below this letter are our specific responses to the critique by the reviewers.</p> <p>Our manuscript presents a deep learning neural network which has been trained for automatic body composition segmentation on L3 CT-images, using several large patient cohorts (n=3187). We then externally validated our model on a large cohort (n=2535) consisting of pancreatic and colon cancer patients. The algorithm performed excellent</p>

compared with manual segmentations.

In the field of AI tools for potential clinical use, we think that robust, fully inter-institutional and large-scale external testing with unseen datasets remain a rarity. Therefore the quality and robustness of study in this work is rare and remains highly deserving of attention. We believe work such as this will be highly cited as a new benchmark for this topic. What stands out is that the algorithm was able to successfully segment "challenging" scans (e.g. with anatomic variations) due to the large heterogeneous patient cohorts it had been trained on. The algorithm is an ideal tool to study large patient cohorts with relative ease, which would otherwise be impossible to assess manually. In addition, due to the design of the algorithm and integration of automatic L3 selection, it can be easily implemented in various clinical infrastructures and used by other research groups to assess cancer patient cohorts or develop new models in other fields.

I hope that you will agree that this study is of sufficient quality and novelty to merit publication in the Journal of Cachexia, Sarcopenia and Muscle.

Sincerely,

David P.J. van Dijk, MD, PhD
On behalf of the team

Response to reviewers' comments:

Reviewer #2:

1.Comment: These results and conclusion are replicating existing work and already old news. There are numerous papers that have not only shown L3 segmentation for SKM, VAT, SAT but also IMAT (which wasn't covered here) and go a step beyond to finding the L3 automatically. Line 58 states that "Scientifically, our L3 segmentation tool enables assessment of large (incl. historical) cohorts that would be unfeasible to segment manually". This is also not accurate as the task of extracting the L3 manually hasn't been automated, which would present a significant bottleneck in the analysis of large cohorts.

Response (1/3): Thank you for your comment, however in the field of AI tools for potential clinical use, we think that robust, fully inter-institutional and large-scale external testing with unseen datasets remain a rarity. Therefore the quality and robustness of study in this work is rare and remains highly deserving of attention. We believe work such as this will be highly cited as a new benchmark for this topic. Automatic L3 extraction was not the focus of this algorithm nor is it novel of itself, because fully stand-alone automatic L3 selection tools already existed for a long time and is sometimes part of vendor-provided standard CT-scanner software packages. However, we do agree with the reviewer that some method of automatic L3 vertebra extraction remains important for using the tool in large cohorts or for clinical implementation. To this end, adopted a highly modular software design (i.e. the chaining together of highly specialized "narrow AI" components) and thus incorporated the Total Segmentator externally validated whole-body segmentation deep learning model into our workflow (TotalSegmentator).

Added: Methods, p10, lines 15-22: "For use on large cohorts and for ease of future clinical implementation, automatic vertebra localization is necessary. We have integrated a state-of-the-art externally validated and open-source tool known as TotalSegmentator (<https://github.com/wasserth/TotalSegmentator>)^{14,15}. In keeping with the "narrow AI" paradigm, we have chained together highly specialized AI tools for each task. TotalSegmentator was first used for automated segmentation of all visible vertebrae in a volumetric CT study. The resulting labelled masks were used to locate all the slices intersecting L3, and then we selected the CT slices closest to the centre of the segmented object (see Figure S2)."

AND

Results, p14, lines 7-10: "We tested the accuracy of L3 mid-slice localization from TotalSegmentator using a small independent test cohort of 30 subjects. The tool correctly extracted the CT-slice at L3 in 30 out of 30 cases (100%)."

AND

Discussion, p16 lines 15-22: "For volumetric CTs as input, an important consideration is how to select the slice intersecting the middle of the L3 vertebra, and more generally in case the user arbitrarily wishes to select some other vertebra. In keeping with the "narrow AI" paradigm, we have elected to implement a modular software design such that highly specialized DLNN are joined up sequentially in a workflow to accomplish a meaningful task. For the present, we integrated the state-of-the-art and validated TotalSegmentator tool to automatically localize spinal vertebrae. If a superior vertebrae segmentation tool should emerge in future, we could relatively easily adapt our workflow to incorporate the new tool, compared to "all-in-one" monolithic software design."

Response (2/3): Regarding IMAT: this was purposefully not included in the algorithm because IMAT is very heterogeneously distributed within and among skeletal muscles (Bhullar et al. J Cachexia Sarcopenia Muscle. 2020 Jun;11(3):735-747), invalidating/complicating accurate quantification of its volume based on single-slice analysis.

Response (3/3): With regard to the novelty of this study: Our deep learning neural network was trained on CT-slices of >3000 individual patients, with different types of cancer/disease, from different centres, that were previously manually segmented by different researchers trained in body composition analysis. This heterogeneity created a highly robust algorithm which can handle challenging CT-slices (i.e. suboptimal patient positioning, anatomical variations; which is the clinical reality). To our knowledge, we are the first group to successfully validate an algorithm on a completely different external patient cohort (of >2500 slices), again with different cancer types and manual segmentation (as ground truth) performed by external researchers trained in body composition analysis. In the discussion section, we discussed the differences and novelties of our algorithm compared to the other published algorithms in more detail (page 16, lines 1-9): "Some other automated segmentation tools have been developed. The largest cohort (n=12,128) was used for development of the AI tool published by Magudia et al.¹⁴ Their tool performed well with similar dice scores to our algorithm. Their training cohort only included 604 pancreatic cancer patients while the large (n=12,128) hospital dataset was used to derive reference curves. However, the large hospital dataset only included patients without cancer and cardiovascular disease, making it less applicable to a clinical population of subjects with cancer who frequently display body composition alterations. In addition, analysis of CT-scans of cancer patients can be more challenging due to anatomic abnormalities and suboptimal patient positioning. As patients with cancer were excluded, the tool by Magudia et al. could perform worse in cancer cohorts. Our analyses did not exclude patients with anatomical variations or unconventional patient positioning, likely resulting in a more robust segmentation tool. Dabiri et al. published an automated segmentation tool which was trained on two cohorts of patients with cancer (n=2529).²¹ Their segmentation tool performed similarly well compared with our segmentation tool. However, in contrast to our study, they did not perform external validation, making it uncertain how their AI performs in other cohorts."

2.Comment: Conclusion is stated as following: "Due to the design of the DLNN, it can be easily implemented in various clinical infrastructures and used by other research groups to assess cancer patient cohorts or develop new models in other fields." However, there is no evidence to support this conclusion. Implementation in various clinical infrastructures is anything but easy and no trained model weights are offered to the community (DLNN code offering is useless as it is not so difficult a project to build a DLNN in a few lines of code now).

Response: We agree with the reviewer that implementation in clinical infrastructure can be challenging. However the comment of the reviewer, that the main thing is to write a few lines of code to make a DLNN, is also highly inaccurate. The first true challenge of AI implementation is to have sufficiently massive dataset for training and then once again for robust interstitial validation, which we have achieved and is relevant with regards to writing a few lines of code. The second challenge is to embed the DLNN tools within a workflow that is useable by radiographers and clinicians. To that end, we have embedded the complexity of the workflow into a web browser-based graphical user interface. Additionally, the code for the DLNN architecture and untrained model is

made publicly available and open access, thus it is even unnecessary for future workers to write any lines of code. A non-profit service is also provided for external institutional users; access to the trained model can be requested via the public website www.mosamatic.com.

Added: Data availability statement: p18 line 8: "The trained model is available upon request through www.mosamatic.com."

Modified: Conclusion, p 17 lines 16-20: "To simplify future use and potential integration of the DLNN-based automated segmentation workflow, we have incorporated the steps into a web browser-based graphical user interface. Clinical implementation within our own institution is not within the scope of this study, but is the subject of a future study. For external institutional users who wish to access the trained model for research, please see data availability statement."

3.Comment: Ultimately, all efforts that seek to improve the state of art in body composition deserve support and on that basis, this paper deserves to be published somewhere and will likely get published somewhere. Maybe if this project actually did some level of analysis of some interesting cancer cachexia/sarcopenia cohort linking to clinical outcomes, as part of that larger effort, a new method development can be explained and published in JCSM. If presenting methodological development replicating already known knowledge as has been submitted, then, in my humble opinion, since there is nothing exciting or novel that this paper contributes to the knowledge base in the sarcopenia and cachexia community, it does not merit a publication in JCSM.

Response: See answer 1 for novelty and impact. The objective of this project was to provide an externally validated automatic deep learning body composition segmentation algorithm to the scientific community. A clean methodological paper facilitates future use in large clinical cohorts by different research groups with different research topics. In our opinion, adding a clinical outcome to this paper would make the paper less readable and less approachable by other research groups.

Reviewer #3:

4.Comment: Its not entirely clear if the L3 images are still having to be individually, by the researcher, extracted, saved as a DICOM and then analysed by the DLNN or if the program finds the appropriate image given the whole scan? It is this process which is the most time consuming and would be best automated.

Response: In this study, all L3 images provided by the external validation institution (Glasgow) had been extracted manually from their respective CT studies. Therefore, automated selection was not strictly needed. However, for ease of future use and potential clinical studies, we have already integrated a standalone deep-learning based whole-body segmentation tool to automatically locate all visible spinal vertebrae in volumetric CTs into our processing workflow (please see our response above to Reviewer #2).

5.Comment: Following on from this point it would be good to include in the discussions some of the limitations of the program and how it can be developed going forward.

Response: We agree with the reviewer that it is important to state limitations of the DLNN and potential improvements. We therefore added the following section to the discussion section.

Added: Discussion, p16-17, lines 23-6: "While the DLNN showed excellent performance, even with challenging CT-scans, it has its limitations. In particular, analysis of CT-scans of patients with anatomical abnormalities (e.g. large abdominal hernia, colostomy, profound edema) or of patients with abnormal/non-standard positioning in the CT-scanner can lead to (partially) incorrect segmentations. Such challenging CT-images should then be manually corrected and stored prospectively. In due time, this cohort of "challenging CT-images" can be used to retrain and improve the DLNN. In addition, different deep learning segmentation algorithms will have different limitations depending on the cohort. A comparative study using both healthy individuals and different patient groups could provide insight into how these different

algorithms perform and if one algorithm is preferred over the other in specific cohorts.”

6.Comment: Page 8 Line 9 'Subjects' I do think this would benefit from a better title to acknowledge these are patients. As with Line 11 'abdominal surgical subjects' - I think should be replaced with something more appropriate like 'A total of 3,187 patients requiring abdominal surgery who had undergone a CT scan prior to surgery contributed..' (see 'patient' characteristics).

Response: We agree with the reviewer and changed the text accordingly.

Changed: Methods: p8 lines 3-6: “Patients, A total of 3,187 patients requiring abdominal surgery who had undergone a CT scan prior to surgery contributed by 32 distinct centres were used for DLNN development (see general patient characteristics in Table 1).”

Reviewer #4:

7.Comment: There are several published NN for body composition analysis, and the authors comment on page 6 that other methods may be underpowered. It would therefore be interesting to see how other NN perform on the cohorts used in this manuscript. Can the authors run published algorithms, such as by Magudia et al. or the DAFS platform from Voronoi, and compare their performance (including accuracy and speed) to their NN? This would be a valuable demonstration of the utility of their method - or at least its performance compared to others.

Response: We agree with the reviewer that a comparison with other published DLNNs would offer much insight, particularly to see if some algorithms are preferred over others in certain (patient) cohorts. Such a comparative study has been initiated and will be performed and published in due course. We added the importance of a comparative study in the discussion section.

Added: Discussion, p17 lines 2-6: “In addition, different deep learning segmentation algorithms will have different limitations depending on the selected cohort. A comparative study using both healthy individuals and different patient groups could provide insight in how these different algorithms perform and if one algorithm is preferred over the other in specific cohorts.”

8.Comment: The authors should specify the hardware environment/requirements on which this analysis is run, and also provide information on computation time, since the authors comment on the significant reduction in labour between manual and automated - would be nice to see that exact comparison exactly for their method and therefore the obvious utility if integrated into clinical practice at some point.

Response: Indeed, the time saved by automatic segmentation vs manual segmentation is enormous, especially for large cohorts. We added a comparison between manual and automatic segmentation time, as well as technical requirements needed to run the algorithm.

Added: Results, p12 11-15: Segmentation speed, The DLNN was able to segment a single CT-image in around 2 seconds and the whole external validation cohort in around 90 minutes only using the CPU. Considering that an experienced clinical researcher train in body composition analysis needs a minimum of 2-5 minutes to segment a single CT-image, the use of automatic segmentation can potentially save months of work when assessing large cohorts.

AND

Methods, p10 lines 12-13: “The trained algorithm can run easily on a conventional office laptop with standard specifications.”

9.Comment: It is not completely clear how the NN manages hands, arm and other extraneous objects - were these recognised and filtered out? There is mention of other methods not taking this into account, but I can't tell if it was considered here.

Response: Generally, patients included in the training and validation cohorts did not have their arms on the CT-image. However, we have tested our algorithm before in

trauma patients (in which many patients had one or both arms displayed in the CT-image) and the algorithm performed well. This was mentioned in the introduction (Introduction, p7 lines 5-10): “In previous work, the DLNN that is the subject of this paper was independently validated using a large polytrauma patient cohort extracted from the same university hospital, albeit at a different department and for a clinically distinct setting. This was nonetheless considered a challenging validation attempt due to the large variation in patient positioning (including arms and hands appearing inside the field of view) as well as radiation artifacts (e.g., from metal devices attached to the patient). Even with this challenging cohort, the present DLNN model performed very well.”

10.Comment: Even though the original model was published in a separate paper previously, it would be useful to have a description of the training cohort characteristics underlying the model, in order to understand how this validation cohort compares to it, and the degree of comparable clinical demographics. Can the authors also comment on any concerns, if any, in the 3187 cohort of patients with different, as opposed to one, types of cancers?

Response: As the training set consisted of several international cohorts, we do not have access to the full set of patient characteristics. Their main characteristics are summarized in table 1. We deliberately chose cohorts from a variety of (international) centres, as well as different disease and cancer types. This prevents overfitting of the AI-model (e.g. to a single disease type or specific medical centre).

Changed: Discussion, P 16, lines 9-11: “Our analyses did not exclude patients with anatomical variations or unconventional patient positioning, which prevents overfitting the model to a specific patient group and will likely result in a more robust segmentation tool.”

11.Comment: In the discussion, the authors state that data “supports the use of body composition analysis in the standard diagnostic work-up”. These statements are controversial, since any diagnostic measures should have relevant clinical implications – given the uncertainties regarding the independent prognostic effect of body composition alterations across different cancer types and histologies within a single cancer type, and the fact that there are no tailored medical interventions to address the impact of e.g. low SAT with proven patient benefit, I would recommend to reword this statement and highlight the caveats.

Response: We agree with the reviewer that it is still too early to provide treatment advice based on body composition. Automatic body composition segmentation will greatly facilitate data collection and creation of larger cohorts for specific cancers and their subtypes, potentially enabling clinical implementation in the future. We have changed our phrasing in the discussion section:

Changed: Discussion p15, lines 5-7: Larger cohorts are needed for each cancer type, as these could support the use of body composition analysis in the standard diagnostic work-up, and potentially aid in clinical treatment decision-making.

12.Comment: Also in the discussion, and an extension of my concerns as above re: emphasis of clinical utility without mentions of specific clinical scenarios or potential caveats, the messaging would benefit from some consideration as to exactly how clinical treatment decision making would be influenced - more intense follow up? more frequent imaging? additional treatment if deemed high risk? and so on. That said, I do like the idea that body composition combined with other risk factors has the potential to be more powerful in terms of its clinical risk predictive ability.

Response: We agree with the reviewer that it is important to pay attention to/discuss the potential implications of body composition analysis for clinical practice. We therefore suggested several hypothetical uses and added potential scenarios.

Added: Discussion p15, lines 11-16: “In the end, integrating body composition data with established prognostic factors such as tumour stage may improve prediction of a patient’s prognosis. A combined tumour and host focused approach would provide a basis for clinical trials aimed at exploring whether body composition-based prognostic

information can be used as a basis for treatment decision making (e.g. palliative intent instead of curative intent, or indication for/selection of (neo)adjuvant therapy).”

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Validation of a deep learning model for automatic segmentation of skeletal
2 muscle and adipose tissue on L3 abdominal CT images

3 David P.J. van Dijk^{1,2,*§}, Leroy F. Volmer^{3,*}, Ralph Brecheisen^{1,2}, Ross D. Dolan⁴, Adam S. Bryce^{5,6}, David
4 K. Chang^{5,6}, Donald C. McMillan⁴, Jan H.M.B. Stoot⁷, Malcolm A. West⁸, Sander S. Rensen^{1,2}, Andre
5 Dekker², Leonard Wee^{2,**}, Steven W.M. Olde Damink^{1,2,9,**}, Body Composition Collaborative

6 *Authors contributed equally

7 **Authors contributed equally

9 ¹Department of Surgery, Maastricht University Medical Centre, Maastricht, The Netherlands

10 ²NUTRIM School of Nutrition and Translational Research in Metabolism, Maastricht University, The
11 Netherlands

12 ³Department of Radiotherapy (MAASTRO), School of Oncology and Reproduction, Maastricht University,
13 Maastricht, The Netherlands.

14 ⁴Academic Unit of Surgery, School of Medicine, University of Glasgow, Glasgow Royal Infirmary,
15 Glasgow, United Kingdom

16 ⁵Wolfson Wohl Cancer Research Centre, School of Cancer Sciences, University of Glasgow, Glasgow,
17 United Kingdom

18 ⁶West of Scotland Pancreatic Unit, Glasgow Royal Infirmary, Glasgow, United Kingdom

19 ⁷Department of Surgery, Zuyderland Medical Centre, Sittard-Geleen, The Netherlands

20 ⁸Academic Unit of Cancer Sciences, Faculty of Medicine, University of Southampton, Southampton, UK

21 ⁹Department of General, Visceral and Transplant Surgery, University Hospital Aachen, Aachen, Germany

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Collaborators:**

2 The Body Composition Collaborative: Thais T.T. Tweed^{1,2}, Stan Tummers², Gregory van der Kroft¹,
3 Marjolein A.P. Ligthart¹, Merel R. Aberle¹, Lubbers Tim¹, Bart C. Bongers³, Jorne Ubachs⁴, Roy F.P.M.
4 Kruitwagen⁴, Siân Pugh⁵, John N. Primrose⁶, John A. Bridgewater⁷, Philip H. Pucher⁸, Nathan J. Curtis⁹,
5 Stephan B. Dreyer^{10,11}, Michael Kazmierski¹²

6
7 ¹Department of Surgery, Maastricht University Medical Centre, Maastricht, The Netherlands

8 ²Department of Surgery, Zuyderland Medical Centre, Sittard-Geleen, The Netherlands

9 ³Department of Nutrition and Movement Sciences, School of Nutrition and Translational Research in
10 Metabolism (NUTRIM), Faculty of Health, Medicine and Life Sciences, Maastricht University, Maastricht,
11 The Netherlands

12 ⁴Department of Obstetrics and Gynaecology, Maastricht University Medical Center, Maastricht, The
13 Netherlands

14 ⁵Department of Medical Oncology, University of Southampton, Southampton, UK

15 ⁶Department of Surgery, University of Southampton, Southampton, UK

16 ⁷UCL Cancer Institute, University College London, London, UK

17 ⁸Department of General Surgery, Portsmouth Hospitals University NHS Trust, Portsmouth, UK

18 ⁹Academic Unit of Cancer Sciences, Faculty of Medicine, University of Southampton, Southampton, UK

19 ¹⁰Wolfson Wohl Cancer Research Centre, School of Cancer Sciences, University of Glasgow, Glasgow,
20 United Kingdom

21 ¹¹West of Scotland Pancreatic Unit, Glasgow Royal Infirmary, Glasgow, United Kingdom

22 ¹²Department of Radiotherapy (MAASTRO), School of Oncology and Reproduction, Maastricht
23 University, Maastricht, The Netherlands.

24
25 **Keywords:** body composition, deep learning, convolutional neural networks, image segmentation, cancer
26 cachexia, computed tomography

1
2
3
4 1
5
6 2 §Corresponding author:
7
8 3 David P.J. van Dijk
9
10 4 Maastricht University, Department of Surgery
11
12
13 5 P.O. box 616
14
15 6 6200 MD, Maastricht
16
17 7 The Netherlands
18
19
20 8 Tel: +31433881526
21
22 9 Fax: +31433884154
23
24 10 Email: david.van.dijk@mumc.nl
25
26 11
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Abstract**

2

3 **Background**

4 Body composition assessment using abdominal computed tomography (CT) images is increasingly applied
5 in clinical and translational research. Manual segmentation of body compartments on L3 CT images is time-
6 consuming and requires significant expertise. Robust high-throughput automated segmentation is key to
7 assess large patient cohorts and ultimately, to support implementation into routine clinical practice. By
8 training a deep learning neural network (DLNN) with several large trial cohorts and performing external
9 validation on a large independent cohort, we aim to demonstrate the robust performance of our automatic
10 body composition segmentation tool for future use in patients.

11 **Methods**

12 L3 CT images and expert-drawn segmentations of skeletal muscle, visceral adipose tissue, and
13 subcutaneous adipose tissue of patients undergoing abdominal surgery were pooled (n = 3,187) to train a
14 DLNN. The trained DLNN was then externally validated in a cohort with L3 CT images of patients with
15 abdominal cancer (n = 2,535). Geometric agreement between automatic and manual segmentations was
16 evaluated by computing two-dimensional Dice Similarity (DS). Agreement between manual and automatic
17 annotations were quantitatively evaluated in the test set using Lin's Concordance Correlation Coefficient
18 (CCC) and Bland-Altman's Limits of Agreement (LoA).

19 **Results**

20 The DLNN showed rapid improvement within the first 10,000 training steps and stopped improving after
21 38,000 steps. There was a strong concordance between automatic and manual segmentations with median
22 DS for skeletal muscle, visceral adipose tissue, and subcutaneous adipose tissue of 0.97 (interquartile range,
23 IQR: 0.95-0.98), 0.98 (IQR: 0.95-0.98), and 0.95 (IQR: 0.92-0.97), respectively. Concordance correlations
24 were excellent: skeletal muscle 0.964 (0.959-0.968), visceral adipose tissue 0.998 (0.998-0.998), and
25 subcutaneous adipose tissue 0.992 (0.991-0.993). Bland-Altman metrics (relative to approximate median
26 values in parentheses) indicated only small and clinically insignificant systematic offsets : 0.23 HU (0.5%),

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 1.26 cm².m⁻² (2.8%), -1.02 cm².m⁻² (1.7%), and 3.24 cm².m⁻² (4.6%) for skeletal muscle average
2 radiodensity, skeletal muscle index, visceral adipose tissue index, and subcutaneous adipose tissue index,
3 respectively. Assuming the decision thresholds by Martin et al. for sarcopenia and low muscle radiation
4 attenuation, results for sensitivity (0.99 and 0.98 respectively), specificity (0.87 and 0.98 respectively), and
5 overall accuracy (0.93) were all excellent.

6 **Conclusion**

7 We developed and validated a deep learning model for automated analysis of body composition of patients
8 with cancer. Due to the design of the DLNN, it can be easily implemented in various clinical infrastructures
9 and used by other research groups to assess cancer patient cohorts or develop new models in other fields.

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Introduction

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Body composition assessment using routine abdominal computed tomography (CT) images is increasingly applied in clinical and translational research. By measuring the tissue area at the level of the third lumbar vertebra (L3) and scaling for subject height, precise assessments of total body mass of skeletal muscle (SM), visceral adipose tissue (VAT), and subcutaneous adipose tissue (SAT) can be made.¹ Body composition has been found to be highly independently predictive of survival, especially among cancer patients. In particular, low skeletal muscle mass (i.e., sarcopenia), low adipose tissue mass, and decreased skeletal muscle radiodensity (i.e., myosteatorsis) have been shown to be associated with shorter overall survival in various cancer types.²⁻⁴

Body composition exhibits substantial heterogeneity among people due to natural variation in age, sex, race, and build.⁵ These intrinsic inter-personal differences are unrelated to disease and may therefore obscure disease related body composition effects, necessitating large population-based data cohorts to adjust for them.

Manual segmentation of body compartments on L3 CT images is time-consuming and requires significant expertise. Therefore, robust high-throughput automated segmentation is key to body composition assessment in large patient cohorts and ultimately, to support implementation of body composition assessment into routine clinical practice. A deep learning neural network (DLNN) can be an essential part of such an automated workflow.

One challenge for developing a robust DLNN is that patients do not always have the ideal CT scans for body composition assessment, such that variable orientation of the patient, degradation of image quality due to radiation artefacts, and individual-specific anatomical attributes may result in poor performance of an automated segmentation algorithm.⁶ A systematic review revealed that one in three DLNN studies of body composition segmentation have been developed with less than 100 unique human subjects, and more than half of the reviewed studies used exclusively single-institutional datasets.⁷ Robust DLNNs need to be trained on datasets that are large enough to incorporate the heterogeneity created by a variety of scanners,

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 image acquisition settings, image reconstruction kernels, patient positioning protocols, and sufficiently high
2 heterogeneity of subject clinical presentations. Additionally, the quantitative performance of DLNNs need
3 to be comprehensively evaluated with external test datasets sourced from a wholly independent clinical
4 workflow and a separate clinical setting from the one used to train the DLNN.⁸

5 In previous work, the DLNN that is the subject of this paper had been independently validated using a large
6 polytrauma patient cohort extracted from the same university hospital, albeit at a different department and
7 for a clinically distinct setting.⁹ This was nonetheless considered a challenging validation attempt due to
8 the large variation in patient positioning (including arms and hands appearing inside the field of view) as
9 well as radiation artifacts (e.g., from metal devices attached to the patient). Even with this challenging
10 cohort, the present DLNN model performed very well.

11 A robust, fully inter-institutional and large-scale external testing with unseen datasets is needed for
12 developing a quality AI tool for potential clinical use. This paper presents the first validation of the
13 Mosamatic DLNN in a surgical oncology cohort using data from a separate hospital, using previously
14 unseen scanners, with independent radiology scan protocols, and with reference delineations provided by
15 independent clinicians.

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Patients and methods**

2

3 *Patients*

4 **A total of 3,187 patients requiring abdominal surgery who had undergone a CT scan prior to surgery**

5 **contributed by 32 distinct centres were used for DLNN development (see general patient**

6 **characteristics in Table 1).** These comprised of de-identified data abstracted from previously ethics board-

7 approved clinical studies; permission for secondary analysis was obtained via the principal investigators of

8 the respective studies. We used L3 CT slices from: three colorectal liver metastases trials - two from

9 multiple sites across the UK and a single-institution study in The Netherlands; two ovarian cancer trials

10 among five participating Dutch centers; and one pancreatic cancer trial of patients operated either in

11 Aachen, Germany, or in Maastricht, the Netherlands.

12 An independent external validation set comprised 2,535 L3 CT slices at different time intervals taken from

13 1,054 unique subjects diagnosed with either resectable colorectal or pancreatic cancer (see Table 1).^{10,11}

14 Ethical approval was granted by the West of Scotland Research Ethics Committee, Glasgow.

15

16 *Image acquisition and reference segmentations*

17 The aforementioned datasets comprised CT scans from a broad range of equipment vendors and image

18 acquisition settings. Images were archived in DICOM (Digital Imaging and Communications in Medicine)

19 format. Table S1 (see online supplementary materials) summarizes the diverse imaging settings as recorded

20 in DICOM metadata.

21 All human-made segmentations in this study were created with *Slice-o-matic* (Tomovision, Quebec,

22 Canada). Regions of interest (ROIs) were defined using standardized Hounsfield Unit (HU) ranges (SM: -

23 29 to +150, VAT: -150 to -50, SAT: -190 to -30). Absolute areas were normalized by physical height

24 squared to derive skeletal muscle index (SMI), visceral adipose tissue index (VATI), and subcutaneous

25 adipose tissue index (SATI). Mean HU in SM at L3 was used as the skeletal muscle radiation attenuation

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 (SMRA). All human reference segmentations were made by clinical researchers trained to perform body
2 composition analysis in Slice-o-matic.

3 Previously published analyses on the external validation dataset had been made with ImageJ (National
4 Institutes of Health, v1.47, <http://rsbweb.nih.gov/ij/>), but this method was shown to overestimate adipose
5 tissue areas relative to other software.¹² Every validation subject in this study was therefore independently
6 re-annotated in Slice-o-matic by the original data owners. To ensure consistency for direct comparison, we
7 re-computed areas and mean HU for all subjects with independent Python code, and confirmed equivalent
8 values with each version of Slice-o-matic used to 2 decimal places or better.

9
10 *Deep learning neural network (DLNN)*

11 A DLNN for multi-label segmentation of SM, VAT, and SAT was built from a canonical 2D U-Net,¹³ with
12 minor change in the input layer to match the dimensions of a CT slice (512x512). An essential development
13 for this work was to chain two independently-trained U-Net networks; the first U-Net was developed to
14 segment the whole abdomen, whilst ignoring hands, arms, CT mattress and extraneous medical devices that
15 sometimes appeared in the CT field of view. The second U-Net was specialized for segmenting SM, VAT,
16 and SAT within the abdominal outline detected by the first U-Net (see online supplementary materials
17 Figure S1 and its accompanying text).

18 Pixel intensities were clipped to the range [-500, +500] HU for the abdomen segmentation network. The
19 reference abdominal region was generated by computing the outermost continuous contour of the human
20 expert's SAT region before morphologically filling in every pixel inside. The range of intensities was
21 further clipped to [-200, +200] HU to train the multi-label segmentation of muscle and fat. In each network,
22 clipped intensities were scaled between [0,1] via standard min-max normalization. Pre-processed CT
23 images were stored and handled in DICOM format. Human expert segmentations were extracted from
24 Slice-o-matic in its proprietary TAG format and converted to Python (NumPy) array objects before training
25 the deep learning model.

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Hands, arms, and other extraneous objects were rare within the training set, thus we synthetically over-
2 sampled images with extraneous objects outside the abdomen until they comprised 50% of each training
3 batch while developing the abdomen U-Net. To train the muscle and fat multi-level segmentation network,
4 all available 3,187 subjects were randomly shuffled and split into 80% for training and 20% for validation.
5 Given the relatively large sample size, a (non-overlapping) 80-20 split is superior to alternative methods
6 like K-fold cross-validation where each validation block ultimately ends up being “seen” by the training
7 algorithm, potentially introducing bias due to data leakage. More details of DLNN construction have been
8 provided in online supplementary materials.

9 CT slices and human-drawn (reference) annotations for the external validation were not revealed until the
10 final DLNN model had been selected and all its model weights permanently fixed. Pre-processing of the
11 test set followed the same steps as aforementioned. The full DLNN code (stripped of all trained models and
12 patient data) is made open access (see data availability statement). **The trained algorithm can run easily
13 on a conventional office laptop with standard specifications.**

14
15 *Automatic L3-selection*

16 **For use on large cohorts and for ease of future clinical implementation, automatic vertebra
17 localization is necessary. We have integrated a state-of-the-art externally validated and open-source
18 tool known as TotalSegmentator (<https://github.com/wasserth/TotalSegmentator>).^{14,15} In keeping
19 with the “narrow AI” paradigm, we have chained together highly specialized AI tools for each task.
20 TotalSegmentator was first used for automated segmentation of all visible vertebrae in a volumetric
21 CT study. The resulting labelled masks were used to locate all the slices intersecting L3, and then we
22 selected the CT slices closest to the centre of the segmented object (see Figure S2).**

23
24 *Analysis*

25 Geometric agreement was evaluated by using 2D Dice Similarity (DS) comparing the DLNN segmentations
26 of SM, SAT, and VAT against the corresponding annotation made by human experts. DS computes the area

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 of the intersection between human and DLNN segmentations as a fraction of half the summated area
2 (human-drawn area plus DLNN-drawn area). Perfect geometric agreement implies $DS = 1$, and if the
3 intersection area is zero then $DS = 0$. Agreement of SMI, VATI, SATI, and SMRA between manual and
4 automatic annotations were quantitatively evaluated in the test set using Lin’s Concordance Correlation
5 Coefficient (CCC) and Bland-Altman’s Limits of Agreement (LoA) (with and without repeated
6 measurements). By using the human-drawn annotations in the test set as reference and then applying the
7 risk classification supplied by Martin et al,² we computed the diagnostic performance (sensitivity,
8 specificity, balanced accuracy, and agreement kappa) of the DLNN results. Statistical analyses were
9 performed in R (version 4.2.0).

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Results

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **Results**

2

3 *Model training*

4 Total loss and DS curves in the training dataset show DLNN model convergence within about 40,000 steps

5 (see Figure S3). There was rapid improvement within the first 10,000 steps but DS was largely stable

6 thereafter. Total (Dice+L2) loss continued to decrease gradually but we stopped model training after 38,000

7 steps, since there was very little to gain with further training. The DLNN weights after the last training step

8 were thus fixed as the “final model” for subsequent testing. The established segmentation tool was named

9 MosaMatic.

11 *Segmentation speed*

12 **The DLNN was able to segment a single CT-image in around 2 seconds and the whole external**

13 **validation cohort in around 90 minutes. Considering that an experienced clinical researcher trained**

14 **in body composition analysis needs a minimum of 2-5 minutes to segment a single CT-image, the use**

15 **of automatic segmentation can potentially save months of work when assessing large cohorts.**

17 *Concordance between manual and DLNN segmentations*

18 The overall distribution of DS for SM, VAT, and SAT in the quarantined validation dataset are summarized

19 in the box-whisker plot shown in Figure 1(a). The median DS for SM was 0.97 (interquartile range, IQR:

20 0.95-0.98), with a tail of outliers down to a minimum DS of 0.45. The distributions of DS for VAT (median:

21 0.98, IQR: 0.95-0.98) and SAT (median: 0.95, IQR: 0.92-0.97) were highly skewed, with extreme outliers

22 landing near zero (these were patients with very small amounts of total adipose tissue). The DS is known

23 to be overly sensitive for small volumes, and this can also be seen in our results – Figure 1(b, c, and d).

24 Lin’s CCC evaluation of SMRA, SMI, VATI, and SATI comparing expert segmentations (as reference)

25 and DLNN results (as test) was excellent, as shown in Figure 2 (a-d). Numerical measures of the

26 concordance correlation coefficient (CCC), bias correction factor for slope of agreement, and finally the

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Bland-Altman intervals of agreement without repeated scans are provided in Table 2. The CCC ranges from
2 0.964 (for SMI) up to 0.998 (for VATI). The errors in the agreement slope, as indicated by deviation from
3 the dotted line in Figure 2, were all close to unity, indicating no major deviations from the ideal, which is
4 supported by bias correction multipliers being better than 0.991 (i.e. no correction implies 1.00). Based on
5 our large cohort, median *in vivo* values (which are in reality age- and sex-dependent) of SMRA, SMI, VATI
6 and SATI roughly fall in the vicinity of 50 HU, 45 cm².m⁻², 60 cm².m⁻² and 70 cm².m⁻². The Bland-Altman
7 metrics (with percentages in parentheses) indicate only small systematic offsets of 0.23 HU (1.0%), 1.26
8 cm².m⁻² (2.9%), 1.02 cm².m⁻² (2.5%), and 3.24 cm².m⁻² (4.9%) for SMRA, SMI, VATI, and SATI,
9 respectively. The upper and lower limits of the Bland-Altman tests indicate SATI had the widest random
10 variation component (-6.7 to 13 cm².m⁻²). Most importantly for risk stratification by muscle fat content, the
11 random noise component of SMRA was estimated at about 2 to 3 HU in magnitude, and correspondingly
12 for SMI about 3 to 5 cm².m⁻² in magnitude.

13
14 *Consistent concordance for repeated measurements*

15 In 449 subjects, we obtained a repeated CT image at varying time intervals ranging from within a month
16 up to 12 months. Whereas the scope of this study was not to objectively quantify longitudinal precision, we
17 can already derive some preliminary insight into stability with repeated imaging over time using this data.
18 The concordance plots for SMRA, SMI, VATI, and SATI for *repeated scans* are equivalent to Figure 2 (see
19 Figure S4). There was no evidence of divergence from the high concordance observed in the agreement on
20 primary CTs. According to CCC metrics and Bland-Altman limits with repeated measures, there are no
21 notable changes between agreement of body composition indices between primary (top half of Table 2) and
22 repeat scans (bottom half of Table 2).

23
24 *Accuracy*

25 We tested the clinical significance of using the DLNN segmentations with respect to a change in
26 stratification for sarcopenia and low SMRA using the widely used thresholds defined by Martin et al.²

1
2
3
4 1 Overall accuracy of stratification was 0.93 for sarcopenia (sensitivity: 0.99, specificity: 0.87) and 0.98 for
5
6 2 low SMRA (sensitivity: 0.98, specificity: 0.98). The discretized agreement (Cohen’s inter-rater kappa) was
7
8 3 0.85 for sarcopenia and 0.96 for low SMRA, which is generally considered as being excellent. For
9
10 4 completeness, a 2x2 confusion matrix for sarcopenia and low SMRA is included in the online supplemental
11
12 5 materials as Figure S5.
13
14

15 6
16
17 7 ***Automatic L3-selection***
18
19 8 **We tested the accuracy of L3 mid-slice localization from TotalSegmentator using a small**
20
21 9 **independent test cohort of 30 subjects. The tool correctly extracted the CT-slice at L3 in 30 out of**
22
23 10 **30 cases (100%).**
24
25
26

27 11
28 12 **Discussion**
29
30

31 13 In this study, we present our high performing and externally validated deep learning model for automated
32
33 14 segmentation of CT-based L3 slices. Due to its excellent performance in both internal and external
34
35 15 validation cohorts, the DLNN-generated segmentation can reliably replace manual segmentation when
36
37 16 performing body composition assessment. This opens up new possibilities both in clinical and scientific
38
39 17 settings, such as cost- and time-effective clinical implementation and large cohort/population studies.
40

41 18 Clinically and subject to clinical implementation study to follow this work, our automated L3 body
42
43 19 composition segmentation tool is intended to be easily implemented in standard practice for all routine CT-
44
45 20 scans, which clinicians can then use for prognostic risk assessment and treatment decision making. Changes
46
47 21 in body composition over time can be detected during oncologic follow-up, which might provide early
48
49 22 indications of treatment effect or disease progression/recurrence. Going from a prognostic tool to a
50
51 23 predictive tool – in which the tool is used for treatment decisions - still remains a large step to take as large
52
53 24 international data-sets are needed to provide clinical reference values.
54
55

56 25 Body composition is highly variable among sex, age, race, and cancer types.^{3,4,16-18} For this reason,
57
58 26 developed clinical cut-offs vary greatly among different patient cohorts and prognostic models of outcome
59
60
61

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 (e.g. survival) are likely to fail during external validation.^{3,19} In addition, body composition can be
2 dependent on other clinical parameters and may have stronger prognostic effects when combined with
3 parameters such as systemic inflammation and weight loss.^{10,20,21} We have previously demonstrated that
4 such combinations or “host phenotypes” are more predictive of overall survival than tumor-based
5 prognostic scores in patients with colorectal liver metastases.²⁰ **Larger cohorts are needed for each
6 cancer type, as these could support the use of body composition analysis in the standard diagnostic
7 work-up, and potentially aid in clinical treatment decision-making.** Automated body composition
8 analysis is the only way of acquiring sufficient data for adequate Z-scoring and accounting for the
9 aforementioned patient characteristics. While cut-offs are necessary for clinical use, we advocate the
10 development of a clinical risk calculator, as the prognostic effect of body composition variables are
11 incremental⁴ and should therefore not be arbitrarily forced into dichotomic cut-offs. **In the end, integrating
12 body composition data with established prognostic factors such as tumour stage may improve
13 prediction of a patient’s prognosis. A combined tumour and host focused approach would provide a
14 basis for clinical trials aimed at exploring whether body composition-based prognostic information
15 can be used as a basis for treatment decision making (e.g. palliative intent instead of curative intent,
16 or indication for/selection of (neo)adjuvant therapy).**
17 Scientifically, our L3 segmentation tool enables assessment of large (incl. historical) cohorts that would be
18 unfeasible to segment manually. In addition, as the AI has learned from multiple observers, it has not
19 learned an expert’s specific signature, ensuring a more stable output. However, the true value of automated
20 segmentation is that it facilitates the inclusion of body composition as a study parameter in RCTs, as the
21 time and effort of analysis is reduced from a couple of months to a few minutes. This enables stratification
22 and selection of patients with different body compositions, creating either homogenous or heterogeneous
23 cohorts as required. Including body composition is particularly important in oncology as it is related to
24 chemotherapy effectiveness and toxicity.²² Ideally, chemotherapy dosing should be based on lean mass to
25 prevent dose-limiting toxicities for which DLNN would be a logical application in the future.

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 Some other automated segmentation tools have been developed. The largest cohort (n=12,128) was used
2 for development of the AI tool published by Magudia et al.¹⁶ Their tool performed well with similar dice
3 scores to our algorithm. Their training cohort only included 604 pancreatic cancer patients while the large
4 (n= 12,128) hospital dataset was used to derive reference curves. However, the large hospital dataset only
5 included patients without cancer and cardiovascular disease, making it less applicable to a clinical
6 population of subjects with cancer who frequently display body composition alterations. In addition,
7 analysis of CT-scans of cancer patients can be more challenging due to anatomic abnormalities and
8 suboptimal patient positioning. As patients with cancer were excluded, the tool by Magudia et al. could
9 perform worse in cancer cohorts. **Our analyses did not exclude patients with anatomical variations or
10 unconventional patient positioning, which prevents overfitting the model to a specific patient group
11 and will likely result in a more robust segmentation tool.** Dabiri et al. published an automated
12 segmentation tool which was trained on two cohorts of patients with cancer (n=2529).²³ Their segmentation
13 tool performed similarly well compared with our segmentation tool. However, in contrast to our study, they
14 did not perform external validation, making it uncertain how their AI performs in other cohorts.
15 For volumetric CTs as input, an important consideration is how to select the slice intersecting the middle
16 of the L3 vertebra, and more generally in case the user arbitrarily wishes to select some other vertebra. In
17 keeping with the “narrow AI” paradigm, we have elected to implement a modular software design such that
18 highly specialized DLNN are joined up sequentially in a workflow to accomplish a meaningful task. For
19 the present, we integrated the state-of-the-art and validated TotalSegmentator tool to automatically localize
20 spinal vertebrae. If a superior vertebrae segmentation tool should emerge in future, we could relatively
21 easily adapt our workflow to incorporate the new tool, compared to “all-in-one” monolithic software design.
22 **While the DLNN showed excellent performance, even with challenging CT-scans, it has its
23 limitations. In particular, analysis of CT-scans of patients with anatomical abnormalities (e.g. large
24 abdominal hernia, colostomy, profound edema) or of patients with abnormal/non-standard
25 positioning in the CT-scanner can lead to (partially) incorrect segmentations. Such challenging CT-
26 images should then be manually corrected and stored prospectively. In due time, this cohort of**

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 **“challenging CT-images” can be used to retrain and improve the DLNN. In addition, different deep**
2 **learning segmentation algorithms will have different limitations depending on the cohort. A**
3 **comparative study using both healthy individuals and different patient groups could provide insight**
4 **into how these different algorithms perform and if one algorithm is preferred over the other in**
5 **specific cohorts.**

6 The key step forward will be implementing automated segmentation into clinical practice and making it
7 easily accessible for new research initiatives. Our tool was created in such a way that it can be easily
8 integrated in clinical imaging software or work independent alongside existing imaging infrastructure. To
9 ensure easy access for research purposes, the untrained AI will be freely available for scientific use and the
10 trained AI can be used under license through a web-app or docker by other research groups. This enables
11 rapid implementation and much needed data collection to develop clinical prediction tools.

12
13 **Conclusion**

14 In this study, we developed a reliable deep learning model that was externally validated for automated
15 analysis of body composition of patients with cancer. **To simplify future use and potential integration**
16 **of the DLNN-based automated segmentation workflow, we have incorporated the steps into a web**
17 **browser-based graphical user interface. Clinical implementation within our own institution is not**
18 **within the scope of this study, but is the subject of a future study. For external institutional users who**
19 **wish to access the trained model for research, please see data availability statement.**

20
21
22 **Supplemental material**

23 Please see the attached document for online access.

24 **Data availability statement**

25 This work concerns only secondary re-use of clinical study data of patients, which were obtained in de-
26 identified form with permission from the original principal investigators. Each study had previously been

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
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 reviewed by a competent ethics body. Data may be obtained from the aforementioned principal
2 investigators upon reasonable request. Source code for data preparation of CT slices and human reference
3 annotations, along with the DLNN model architecture, are publicly available here under a Creative
4 Commons 4.0 CC-BY-NC License:
5 https://github.com/MaastrichtU-CDS/BodyCompL3_DLNN_Open_Code

6 **The trained model is available upon request through www.mosamatic.com.**

7 **Acknowledgments**

8 The New EPOC study was supported by Cancer Research UK.

9 **Conflicts of interest**

10 None.

11

References:

1. Mourtzakis M, Prado CMM, Lieffers JR, Reiman T, McCargar LJ, Baracos VE. A practical and precise approach to quantification of body composition in cancer patients using computed tomography images acquired during routine care. *Appl Physiol Nutr Metab* 2008; **33**(5): 997-1006.
2. Martin L, Birdsell L, Macdonald N, Reiman T, Clandinin MT, McCargar LJ, et al. Cancer cachexia in the age of obesity: skeletal muscle depletion is a powerful prognostic factor, independent of body mass index. *J Clin Oncol* 2013; **31**(12): 1539-47.
3. van Dijk DP, Bakens MJ, Coolsen MMM, Rensen SS, van Dam RM, Bours MJ, et al. Low skeletal muscle radiation attenuation and visceral adiposity are associated with overall survival and surgical site infections in patients with pancreatic cancer. *Journal of cachexia, sarcopenia and muscle* 2017; **8**(2): 317-26.
4. van Dijk DPJ, Zhao J, Kemter K, Baracos VE, Dejong CHC, Rensen SS, et al. Ectopic fat in liver and skeletal muscle is associated with shorter overall survival in patients with colorectal liver metastases. *J Cachexia Sarcopenia Muscle* 2021; **12**(4): 983-92.
5. Heymsfield SB, Gonzalez MC, Lu J, Jia G, Zheng J. Skeletal muscle mass and quality: evolution of modern measurement concepts in the context of sarcopenia. *Proceedings of the Nutrition Society* 2015; **74**(4): 355-66-66.
6. Ha J, Park T, Kim HK, Shin Y, Ko Y, Kim DW, et al. Development of a fully automatic deep learning system for L3 selection and body composition assessment on computed tomography. *Scientific reports* 2021; **11**(1): 21656.
7. Bedrikovetski S, Seow W, Kroon HM, Traeger L, Moore JW, Sammour T. Artificial intelligence for body composition and sarcopenia evaluation on computed tomography: A systematic review and meta-analysis. *European journal of radiology* 2022; **149**: 110218.
8. Norgeot B, Quer G, Beaulieu-Jones BK, Torkamani A, Dias R, Gianfrancesco M, et al. Minimum information about clinical artificial intelligence modeling: the MI-CLAIM checklist. *Nat Med* 2020; **26**(9): 1320-4.
9. Ackermans L, Volmer L, Wee L, Brecheisen R, Sánchez-González P, Seiffert AP, et al. Deep Learning Automated Segmentation for Muscle and Adipose Tissue from Abdominal Computed Tomography in Polytrauma Patients. *Sensors (Basel, Switzerland)* 2021; **21**(6).
10. Dolan RD, Almasaudi AS, Dieu LB, Horgan PG, McSorley ST, McMillan DC. The relationship between computed tomography-derived body composition, systemic inflammatory response, and survival in patients undergoing surgery for colorectal cancer. *Journal of Cachexia, Sarcopenia and Muscle* 2018.
11. Dolan RD, Abbass T, Sim WMJ, Almasaudi AS, Dieu LB, Horgan PG, et al. Longitudinal Changes in CT Body Composition in Patients Undergoing Surgery for Colorectal Cancer and Associations With Peri-Operative Clinicopathological Characteristics. *Frontiers in nutrition* 2021; **8**: 678410.
12. Dolan RD, Tien Y-T, Horgan PG, Edwards CA, McMillan DC. The relationship between computed tomography-derived body composition and survival in colorectal cancer: the effect of image software. *JCSM Rapid Communications* 2020; **3**(2): 81-90.
13. Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image Segmentation. 2015; Cham: Springer International Publishing; 2015. p. 234-41.
14. Wasserthal J, Meyer M, Breit H, Cyriac J, Yang S, Segeroth M. TotalSegmentator: robust segmentation of 104 anatomical structures in CT images. 2022. URL: <https://arxiv.org/abs/2208.05868>. arXiv: 2208.05868.
15. Isensee F, Jaeger PF, Kohl SAA, Petersen J, Maier-Hein KH. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* 2021; **18**(2): 203-11.

1
2
3
4 1 16. Magudia K, Bridge CP, Bay CP, Babic A, Fintelmann FJ, Troschel FM, et al. Population-
5 2 Scale CT-based Body Composition Analysis of a Large Outpatient Population Using Deep
6 3 Learning to Derive Age-, Sex-, and Race-specific Reference Curves. *Radiology* 2021; **298**(2):
7 4 319-29.
8 5 17. Tweed TTT, van der Veen A, Tummers S, van Dijk DPJ, Luyer MDP, Ruurda JP, et al.
9 6 Body Composition Is a Predictor for Postoperative Complications After Gastrectomy for Gastric
10 7 Cancer: a Prospective Side Study of the LOGICA Trial. *J Gastrointest Surg* 2022; **26**(7): 1373-
11 8 87.
12 9 18. Rutten IJG, Van Dijk DPJ, Kruitwagen RFPM, Beets-Tan RGH, Olde Damink SWM, Van
13 10 Gorp T. Changes in skeletal muscle mass during neoadjuvant chemotherapy are related to
14 11 survival in ovarian cancer. *Journal of Cachexia, Sarcopenia and Muscle* 2016; **7**(4): 458-66.
15 12 19. Petermann-Rocha F, Balntzi V, Gray SR, Lara J, Ho FK, Pell JP, et al. Global prevalence
16 13 of sarcopenia and severe sarcopenia: a systematic review and meta-analysis. *J Cachexia*
17 14 *Sarcopenia Muscle* 2022; **13**(1): 86-99.
18 15 20. Van Dijk DP, Krill M, Farshidfar F, Li T, Rensen SS, Olde Damink SW, et al. Host
19 16 phenotype is associated with reduced survival independent of tumor biology in patients with
20 17 colorectal liver metastases. *Journal of Cachexia, Sarcopenia and Muscle* 2018.
21 18 21. Martin L, Senesse P, Gioulbasanis I, Antoun S, Bozzetti F, Deans C, et al. Diagnostic
22 19 criteria for the classification of cancer-associated weight loss. *Journal of clinical oncology : official*
23 20 *journal of the American Society of Clinical Oncology* 2015; **33**(1): 90-9.
24 21 22. Hopkins JJ, Sawyer MB. A review of body composition and pharmacokinetics in oncology.
25 22 *Expert review of clinical pharmacology* 2017; **10**(9): 947-56.
26 23 23. Dabiri S, Popuri K, Ma C, Chow V, Feliciano EMC, Caan BJ, et al. Deep learning method
27 24 for localization and segmentation of abdominal CT. *Computerized medical imaging and graphics*
28 25 *: the official journal of the Computerized Medical Imaging Society* 2020; **85**: 101776.
29 26
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

medRxiv preprint doi: <https://doi.org/10.1101/2023.04.23.23288981>; this version posted January 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

Table 1. General patient characteristics for the deep learning neural network development sets and the external test set.

Study ID	MODEL DEVELOPMENT SETS					TEST SET
	FROGS*	New EPOC*	Zuyd ⁺	MUMC**	MUMC / Aachen ^x	UG [^]
Diagnosis	Emergency laparotomy	Colorectal liver metastases		Ovarian cancer	Pancreatic cancer	Pancreatic cancer + colorectal cancer
Time interval	2017-2019	2007-2012	2013-2017	2002-2015	2015-2019	2008-2019
Sample size	804	153	1587	339	304	1054 (147 pancreatic, 907 colorectal)
No. male (%)	374 (47%)	-	883 (56%)	0 (0%)	161 (53%)	567 (54%)
No. female (%)	430 (53%)	-	704 (44%)	339 (100%)	143 (47%)	487 (46%)
Ages (median)	25-95 (68)	-	32-98 (70)	30-101 -	10-88 (74)	23-93 (69)
Range BMI in kg.m ⁻² (median)	14-58 (26)	-	15-53 (26)	- -	- (25.4)	14-59 (27)

*Bristol, Poole, Bournemouth, Royal Marsden, Surrey, Portsmouth, Velindre, Sheffield, Imperial Charing Cross, Imperial St. Mary, Christie, Southend, Yeovil, North Middlesex, Southampton, Guys, Aintree, Winchester, Cambridge, Princess Alexandra, Bedford, Salisbury, University College London, Basingstoke, Pennine (UK). ⁺Zuyderland Medical Centre Geleen/Heerlen (The Netherlands). ^{**}Maastricht University Medical Centre, Radboud University Medical Centre Nijmegen, Bernhoven Medical Centre Uden, St. Jansdal Medical Centre Ede (The Netherlands). ^x Maastricht University Medical Centre (Netherlands), RWTH Uniklinik Aachen (Germany). [^]Glasgow Royal Infirmary (UK). - No individual values extracted.

medRxiv preprint doi: <https://doi.org/10.1101/2023.04.23.23288981>; this version posted January 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC 4.0 International license](#).

Table 2. Concordance correlation, bias correction factor, and Bland-Altman agreement without repeated measures (n = 1054)

Bland-Altman estimates of agreement for primary scan only (n = 1054)			
	Concordance correlation (95% confidence interval)	Bias correction factor	Bland-Altman agreement (95% lower-upper limits)
SMRA	0.991 (0.990 – 0.992)	0.999	0.23 (-2.06 – 2.52) HU
SMI	0.964 (0.959 – 0.968)	0.991	1.26 (-3.11 – 5.63) cm ² ·m ⁻²
VATI	0.998 (0.998 – 0.998)	0.999	-1.02 (-4.55 – 2.50) cm ² ·m ⁻²
SATI	0.992 (0.991 – 0.993)	0.997	3.24 (-6.69 – 13.2) cm ² ·m ⁻²
Bland-Altman estimates of agreement for repeated scans only (n = 449)			
	Concordance correlation (95% confidence interval)	Bias correction factor	Bland-Altman agreement (95% lower-upper limits)
SMRA	0.991 (0.990 – 0.992)	0.999	0.18 (-2.08 – 2.45) HU
SMI	0.973 (0.969 – 0.976)	0.997	0.75 (-3.56 – 5.06) cm ² ·m ⁻²
VATI	0.998 (0.998 – 0.998)	0.999	-1.07 (-4.55 – 2.41) cm ² ·m ⁻²
SATI	0.992 (0.991 – 0.993)	0.998	2.55 (-8.36 – 13.4) cm ² ·m ⁻²

medRxiv preprint doi: <https://doi.org/10.1101/2023.04.23.23288981>; this version posted January 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC 4.0 International license](https://creativecommons.org/licenses/by-nc/4.0/).

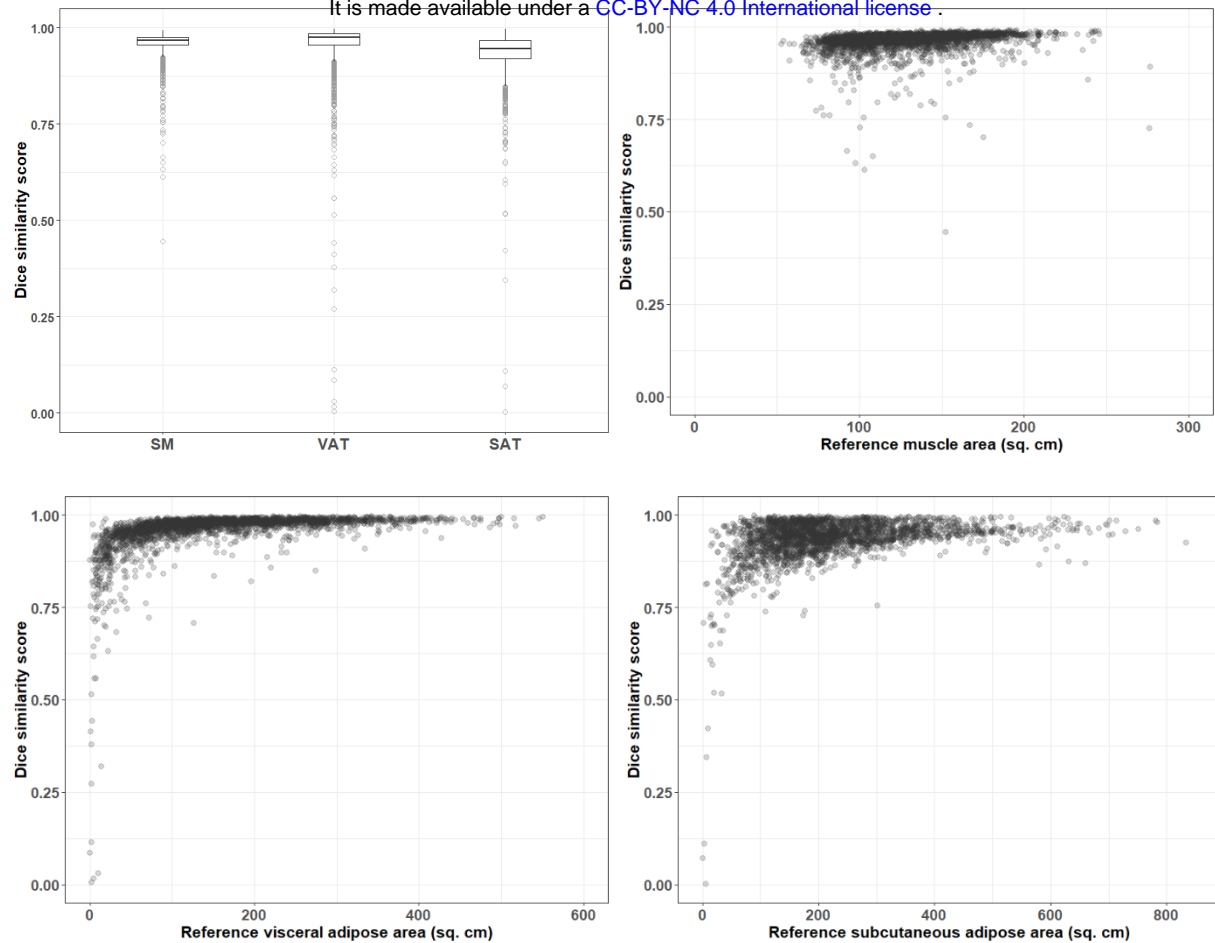


Figure 1: Distribution of geometric DS on L3 slice for skeletal muscle (SM), subcutaneous fat (SAT), and visceral fat (VAT)

(a) Box-whisker plot showing the median DS as the solid horizontal line and the interquartile range as the upper and lower limits of the box. The vertical line ends indicate 1%-tile and 99%-tile, and outliers outside this range are shown as individual dots. (b) – (d) show the distribution of DS as a function of SM area, VAT area, and SAT area, respectively.

medRxiv preprint doi: <https://doi.org/10.1101/2023.04.23.23288981>; this version posted January 22, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC 4.0 International license](#).

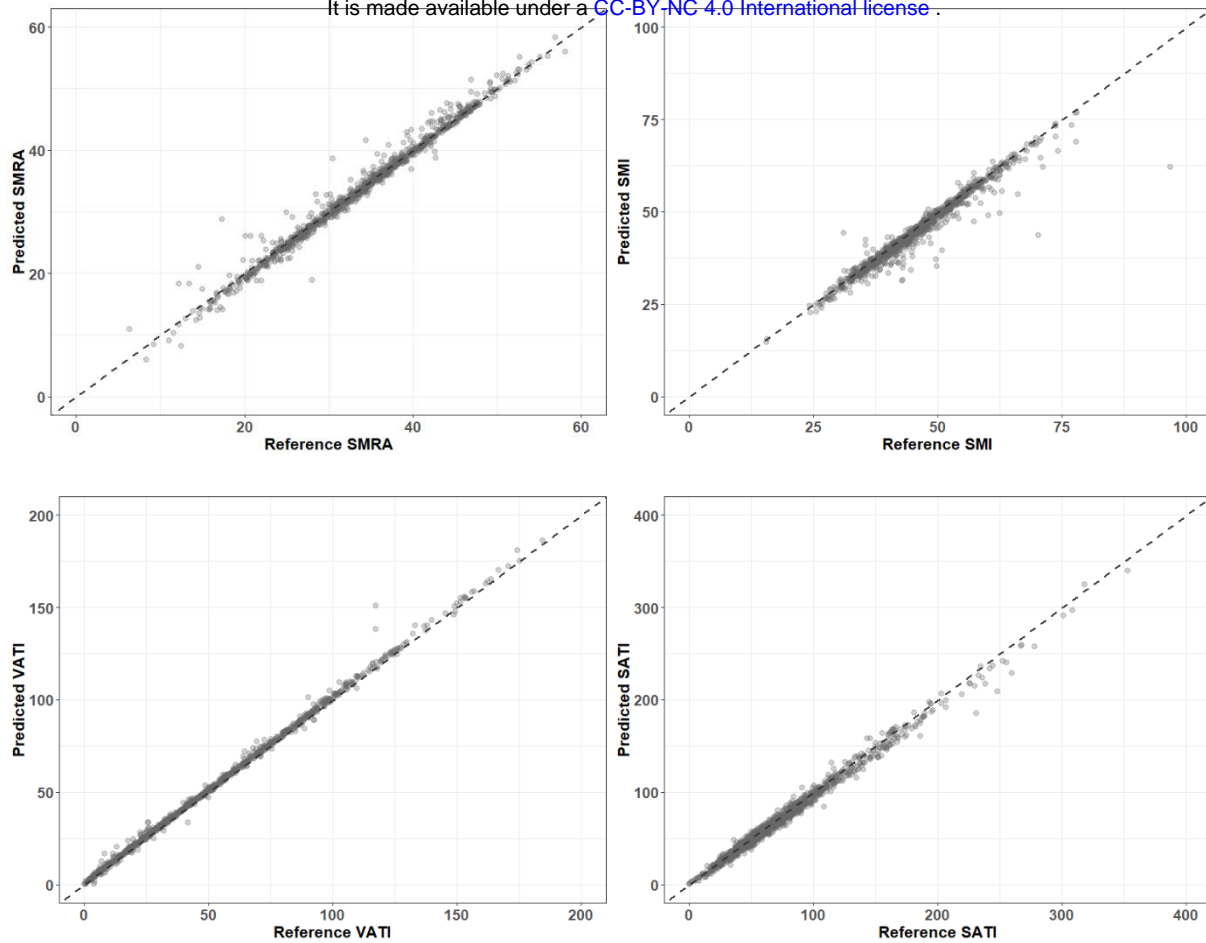
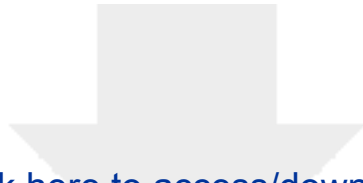


Figure 2: Lin's concordance correlation (CCC) plots

(a) skeletal muscle attenuation (SMRA), (b) skeletal muscle index (SMI), (c) visceral fat index (VATI) and (d) subcutaneous fat index (SATI). The units of SMRA are HU. The units of SMI, VATI, and SATI are all $\text{cm}^2 \cdot \text{m}^{-2}$. Reference values were defined as those extracted from human-drawn segmentations. Predicted values were extracted from DLNN-made segmentations.



Click here to access/download
Supplementary Material
19-7-2023 Algorithm paper supplemental files.docx

