**MeshMonk: Open-source large-scale intensive 3D phenotyping**

Julie D. White[1*], Alejandra Ortega-Castrillón[2,3], Harold Matthews[4,5,6], Arslan A. Zaidi[1,7], Omid Ekrami[8], Jonatan Snyders[9], Yi Fan[4,10], Tony Penington[4,5,6], Stefan Van Dongen[8], Mark D. Shriver[1], Peter Claes[2,3,4*]

[1]Department of Anthropology, The Pennsylvania State University, University Park, PA, USA.
[2]Department of Electrical Engineering, KU Leuven, Leuven, Belgium
[3]Medical Imaging Research Center, UZ Leuven, Leuven, Belgium
[4]Murdoch Children's Research Institute, Melbourne, Australia
[5]Royal Children's Hospital, Melbourne, Australia
[6]Department of Pediatrics, University of Melbourne, Melbourne, Australia
[7]Department of Biology, The Pennsylvania State University, University Park, PA, USA.
[8]Department of Biology, University of Antwerp, Antwerp, Belgium
[9]WebMonks, Hasselt, Belgium
[10]Melbourne Dental School, University of Melbourne, Melbourne, Australia

*Correspondence:
jdw345@psu.edu; peter.claes@kuleuven.be

1    **Abstract**
2
3    In the post-genomics era, an emphasis has been placed on disentangling 'genotype-phenotype'
4    connections so that the biological basis of complex phenotypes can be understood. However,
5    our ability to efficiently and comprehensively characterize phenotypes lags behind our ability to
6    characterize genomes. Here, we report a toolbox for fast and reproducible high-throughput
7    dense phenotyping of 3D images. Given a target image, a rigid registration is first used to orient
8    a template to the target surface, then the template is transformed further to fit the specific shape
9    of the target using a non-rigid transformation model. As validation, we used N = 41 3D facial
10   images registered with MeshMonk and manually landmarked at 19 locations. We demonstrate
11   that the MeshMonk registration is accurate, with 0.62 mm as the average root mean squared
12   error between the manual and automatic placements and no variation in landmark position or
13   centroid size significantly attributable to landmarking method used. Though validated using 19
14   landmarks for comparison with traditional methods, MeshMonk allows for automatic dense
15   phenotyping, thus facilitating more comprehensive investigations of 3D shape variation. This
16   expansion opens up an exciting avenue of study in assessing genomic and phenomic data to
17   better understand the genetic contributions to complex morphological traits.

1

## Introduction

The phenotypic complement to genomics is *phenomics*, which aims to obtain high-throughput and high-dimensional phenotyping in line with our ability to characterize genomes[1]. The paradigm shift is simple and similar to the one made in the Human Genome Project: instead of 'phenotyping as usual' or measuring a limited set of simplified features that seem relevant, why not measure it all? In contrast to genomic technologies, which successfully measure and characterize complete genomes, the scientific development of phenomics lags behind. However, with the advent of new technologies, hardware exists for extensively and intensively collecting quantitative phenotypic data. For example, 3D image surface and/or medical scanners provide the optimal means to capture information of biological morphology and appearance. Today, the challenge is to establish standardized and comprehensive phenotypic representations from large scale image data that can be used to study phenotypic variation in the context of genetic variation[2]. This is a challenge that we address with the development of the MeshMonk toolbox, which enables fast and reproducible high-throughput phenotyping of 3D images, or quasi-landmark indication, which can be applied to 3D facial images as well as 3D scans of other complex morphological structures.

Dense correspondence phenotyping is important beyond genomics and could be employed by anthropologists, biologists, and medical clinicians to accurately and reproducibly characterize anatomical structures such that underlying qualities about the structure can be understood. The study of variation and covariation in anatomy can provide insights into the genetic causes and evolution of the anatomical structure. In addition, comparing the anatomy of an individual patient to a control population can indicate pathology to a medical practitioner. Traditionally, this has been achieved using visual clinical assessment or by taking measurements between manually placed anatomical landmarks. Some examples include the endo- and exocanthi (the inner and outer corners of the eyes, respectively) and the pronasale (the tip of the nose).

However, manual landmarking is tedious to perform, difficult to standardize in practice, and prone to intra and inter-operator error[3–7]. Furthermore, sparse landmark configurations can only quantify form at defined landmarks that can be reliably identified and indicated by a human, and thus lack the resolution to fully characterize shape variation in between landmarks. An alternative is to automatically indicate quasi-landmarks across the entire surface of the structure. This is achieved by gradually warping a generic template composed of thousands of points into the shape of each target image through a non-rigid registration algorithm[8–12]. The coordinates of these warped templates, now in the shape of each target, can then be assessed in geometric morphometric analysis. An automatic approach like this is preferable for the analysis of large datasets, avoiding the problems of manual landmarking by multiple operators. Dense phenotype coordinates are also more suitable for applications that require synthesis of a recognizable instance of the actual structure, such as predicting a complete shape from DNA[13], synthetic growth and ageing of a face[14,15], constructing 3D facial composites for forensic applications[16], and characterization of dysmorphology for clinical diagnosis[17,18].

Surface registration, implemented in MeshMonk, defines a warping of the vertices from one (template) image to their corresponding locations on another (target) image and allows us to quantify and visualize both subtle and acute variation in surface form across a sample by finding the geometrical relationship (one-to-one correspondences) between 3D shapes[8–10,12,19]. The registration strategy is akin to fitting an elastic net onto a solid facial statue through a geometry-driven mapping of anatomically corresponding features. When the template is warped onto each target, the coordinates of any anatomical landmark, manually annotated on the template, can also be defined on each target, thus the complete quasi-landmark indication can also be

2

69 considered a method for automatic placement of sparse anatomical landmarks[20]. As a
70 validation, we compared a set of 19 sparse landmarks indicated manually by two observers and
71 automatically using MeshMonk.
72
73 **Results**
74
75 Accuracy
76 Direct comparison of manual and automatic landmark placements
77 As one measure of validation of the automatic landmark indications, we compared the raw
78 coordinate values of manual landmark indications with the raw coordinate values of automatic
79 landmark indications while considering the manual landmarks to be the "gold standard".
80 Because of the leave-one-out nature of our approach, we can compare the manual and
81 automatic landmark coordinates directly without fear of training bias. To compare landmark
82 indications, we calculated the root mean squared error between the *x, y*, and *z* coordinates for
83 manual and automatic indications (Table 1) and calculated the intraclass correlation coefficient
84 between the *x, y*, and *z* coordinates produced by the two methods. When comparing the
85 average of all six manual landmarking indications ($C_{ML}$) and the automatic landmarks trained
86 using this average ($C_{Auto}$), the highest difference after averaging standard deviation values
87 across all axes was 0.85 mm, for the right side exocanthion landmark (Table 1). Overall, the
88 average standard deviation between $C_{ML}$ and $C_{Auto}$ across all landmarks was 0.62 mm. Bland-
89 Altman comparisons showed that the 95% confidence intervals for the landmark indication
90 between methods are within 1.5 mm of a mean difference of 0 mm (Supplementary Fig. S1).
91 Most individuals fall within these confidence limits, with only a few comparisons from each axis
92 having differences greater than 3 mm. The intraclass correlation coefficients for each axis are
93 around 0.99, representing very high correlation and agreement between manual and automatic
94 landmark indications.
95
96 Centroid size comparison
97 We used estimates of centroid size (CS; the square root of the sum of squared distances from
98 each landmark to the geometric center of each landmark configuration) as an additional
99 assessment of the similarity between manual and automatic landmark placements, since
100 centroid sizes feature heavily in geometric morphometric assessments. The ICC of centroid
101 sizes calculated using the manual and automatic landmarks were all high ($ICC_A = 0.9589$, $ICC_B$
102 $= 0.9486$, $ICC_C = 0.9591$; Supplementary Fig. S2). Analysis of variance (ANOVA) by individual,
103 observer, and method shows that individual is the only significant factor in explaining variance in
104 centroid size ($F = 130.407$, $p < 2 \times 10^{-16}$; Table 2). Bland-Altman comparison showed that the
105 95% confidence intervals for the centroid size estimates between methods are 2 mm relative to
106 an average centroid size of about 165 mm (Supplementary Fig. S2).
107
108 Analysis of shape variance
109 A multivariate analysis of variance (MANOVA) on shape, based on the average of each
110 observer's manual landmark indications and automatic landmark configurations, separately, was
111 performed to determine if the variance explained by individual and observer factors was similar
112 in both methods (Supplementary Table S1). In both methods, individual variation contributed to
113 most of the variation in shape ($R^2_{ML} = 94\%$; $R^2_{Auto} = 97\%$). Differences in observer accounted for
114 1.9% of the variation in shape from manual landmarks and 2.6% of the variation in shape from
115 automatic landmarks. In total, 3.9% of the variation present in manual landmark shape
116 configurations was unexplained by our model while only 0.22% of the variation was unexplained
117 when testing the automatic landmark configurations. A MANOVA on Generalized Procrustes
118 Analysis (GPA) aligned manual and automatic configurations from each observer, with method,
119 individual, observer, and individual x observer as predictors showed that landmarking method

120 did not significantly account for variation in landmark placement ($F = 0.3463$; $p = 0.987$; Table
121 3).
122
123 Reliability
124 Intra- and inter-observer error of manual landmarks
125 The quantitative study of morphology using 3D coordinates requires specific attention to
126 measurement error and has a robust presence in the literature. For each observer, we
127 calculated the intra-observer error of the manual landmarks as the standard deviation between
128 the *x*, *y*, and *z* coordinates of each observer's three landmarking iterations. Supplementary
129 Table S2 reports intra-observer standard deviations for the manual landmark indications along
130 each axis, averaged across images. The average standard deviation of observer A across all
131 landmarks was 0.58 mm while the average standard deviation of observer B across all
132 landmarks was 0.44 mm. The average inter-observer error, measured as the standard deviation
133 between the average *x*, *y*, and *z* coordinates of each observer's landmarking iterations was 0.40
134 mm. This range of deviation is considered highly precise and is similar to previously reported
135 measures of landmark error[6,21].
136
137 The analysis of measurement and observer error for the manual landmarks alone, assessed
138 using a MANOVA for shape, with individual, observer, observer x individual, and nested
139 observer x landmarking iteration as factors showed that non-individual factors contributed
140 significantly to variation in shape (Supplementary Table S3). Individual variation contributed to
141 most of the variation in shape (85%), as expected. Simple measurement error accounted for
142 3.5% of the total variation in shape. Additional to this, differences in observer accounted for
143 1.8% of shape variation, and deviation across landmarking iterations contributed an additional
144 1.5% of the total variation in shape. In total, non-individual effects contributed to 15% of the total
145 shape variation, with 8.3% of this variation unexplained by the model.
146
147 Comparison of manual and automatic inter-observer errors
148 By treating the automatic landmark indications as if they were performed by a third observer, we
149 calculated "inter-observer" errors to compare the variation of automatic and manual
150 landmarking. In this assessment, we compared inter-observer errors calculated using only the
151 manual landmarks ($A_{ML}$ vs. $B_{ML}$) with error estimates calculated by replacing one of the
152 observer's manual landmark indications with the automatic indications trained using that
153 observer's average. This resulted in two extra estimations of inter-observer error ($A_{ML}$ vs. $B_{Auto}$
154 and $A_{Auto}$ vs. $B_{ML}$), calculated as the standard deviation between *x*, *y*, and *z* coordinates
155 (Supplementary Table S4; Supplementary Fig. S3). The mean manual landmarking inter-
156 observer error was 0.40 mm while both manual-automatic comparisons had mean standard
157 deviation values of 0.53 mm (Supplementary Table S4). A paired t-test between the manual
158 landmark error values and each of the manual-automatic comparison showed that the landmark
159 indications that were significantly different between the two methods tended to be those where
160 facial texture likely assisted in the placement of the manual landmarks (e.g. localizing the crista
161 philtra by looking at the differences in color between the lips and the skin; Supplementary Table
162 S5). This result indicates that automatic sparse landmarking using MeshMonk will likely produce
163 more robust results when given input data that has a strong anatomical orientation (e.g. the
164 nasion and pogonion). Even given these differences in variance, the manual-automatic
165 comparisons did not produce errors that were completely outside the range of inter-observer
166 errors, a sign of the reliability of the MeshMonk registration.
167
168 As an illustration of the low errors between automatic landmark indications trained using
169 different observers, we calculated the standard deviation between automatic landmark
170 indications trained using the average of observer A's three landmark indications and the

4

171  average of observer B's three landmark indications ($A_{Auto}$ vs. $B_{Auto}$; Supplementary Table S6,
172  Supplementary Fig. S3). The variance of the average standard deviation values were
173  significantly different for all landmarks except labiale superius, where we could not reject the null
174  hypothesis that the variances of the two standard deviation distributions were equal ($F = 2.4213$,
175  $p = 0.1236$). Supplementary Fig. S3 shows that the variance between automatic landmarking
176  indications ($A_{Auto}$ vs. $B_{Auto}$) is easily identified as being smaller than the manual landmark inter-
177  observer error ($A_{ML}$ vs. $B_{ML}$).
178
179  **Discussion**
180
181  Through studies utilizing manually placed sparse landmarks, we have begun to understand the
182  biological basis and evolution of complex phenotypes, both normative and clinical. However,
183  there is still much to be learned. One avenue for improvement is to expand and speed up the
184  production and analysis of data using methods derived from engineering and computer vision,
185  which allow for the description of shapes as "big data" structures instead of sparse sets of
186  landmarks or linear distances, thus matching our ability to describe phenotypes with our ability
187  to describe genomes. To this end, we introduce the MeshMonk registration framework, giving
188  researchers the opportunity to quickly and reliably establish a homologous set of positions
189  across entire samples. We have validated this framework using a sparse set of landmarks,
190  though the registration framework produces thousands of landmarks to finely characterize the
191  structure.
192
193  MeshMonk represents a step forward in our ability to describe complex structures, like the
194  human face, for clinical and non-clinical purposes. Consider Figure 1, showing the starting
195  template for facial image registration (left) as well as three example faces (right). Each point on
196  the images represents a quasi-landmark data point that is homologous and can be compared
197  across faces. Researchers are no longer limited to a few homologous points, chosen because
198  they can be reliably indicated over hundreds of hours of work. Instead, minute details of the face
199  can be identified and compared across thousands of images in a few hours, and additional
200  images can be incorporated just as easily, regardless of the camera system with which they
201  were captured, allowing for the incorporation of images from different sources and databases
202  (e.g. Facebase.org).
203
204  Because of the relative newness of dense correspondence phenotyping, few studies have
205  focused on the accuracy and reliability of the resulting registrations. Previous studies using
206  versions of the MeshMonk framework have shown that the error associated with the registration
207  of the template onto facial images is 0.2 mm[22] and parameters of the toolbox have been fine-
208  tuned, as discussed elsewhere[8] and in the Supplemental Methods. To provide some validation
209  regarding the ability of the registration process to accurately identify anatomical positions of
210  interest, we used a set of 40 faces with manual landmark indications to "train" positions of
211  interest on the template, then automatically indicate these positions on a face that was not
212  present in the training dataset. In the comparison of manual and automatic landmark
213  indications, the positions of the manual landmarks were considered to be the gold standard, as
214  they have a long history of use and validation in morphological studies[5,21]. By limiting ourselves
215  to a set number of sparse landmarks, we cannot necessarily speak to the accuracy of structures
216  not involved in our validation (i.e. the cheeks), but we argue that the results for our comparison
217  speak highly of the fidelity with which the MeshMonk registration framework aligns to underlying
218  anatomical structures.
219
220  In the direct comparison of sparse landmarks placed manually and using the MeshMonk
221  toolbox, the average difference between the manual and automatic placements was low

5

222   (Supplementary Fig. S1), with the average root mean squared error across all landmarks
223   ranging from 0.62 to 0.68 mm (Table 1), which is well within the range of acceptable error for
224   manual landmarks[5,6,21] and similar or below errors reported in other comparisons of manual and
225   automatic landmarking methods[23–26]. When assessing landmarking methods separately, the
226   variance in landmark configuration attributable to individual and observer factors is similar, with
227   considerably less variation left unexplained by a MANOVA model using automatic landmark
228   configuration as the response (Supplementary Table S1). When assessing manual and
229   automatic landmark configurations in a single MANOVA, the landmarking method is a
230   nonsignificant factor, indicating that variation in scans is not attributable to variation in
231   landmarking method (Table 3). This result was also reproduced when comparing centroid sizes
232   calculated using manual and automatically placed landmarks (Table 2), speaking to the high
233   correspondence between landmark indications placed by human observers and those indicated
234   by the MeshMonk toolbox.
235
236   The validation results together suggest that the MeshMonk toolbox is able to reliably reproduce
237   information given by manual landmarking. Though the larger contribution of the MeshMonk
238   toolbox is the ability to quickly and densely characterize entire 3D surfaces, our illustration using
239   a small number of manually placed landmarks as a training set could be useful for studies
240   seeking specifically to study a sparse set of landmarks, perhaps to add more images to a
241   dataset that is already manually landmarked or to add additional landmarks to an analysis.
242   Utilization of the MeshMonk toolbox also gives the opportunity to minimize variation due to
243   different observers. Take, for example, datasets with manual landmarks indicated by two
244   different observers. During the course of analysis, the inter-observer error of these observers
245   would have to be calculated and taken into account when interpreting results. From our own
246   study, the inter-observer error of the manual landmarks placed by two different observers was
247   0.40 mm (Supplementary Table S2). With the automatic landmarking framework implemented
248   during this study, we can minimize both intra-observer variance for a single scan (by averaging
249   together all indications of that scan by a single observer) and intra-observer variance across
250   scans by placing all indications from the training dataset on the template mesh and averaging
251   the entire training set before using MeshMonk to place them in an automatic fashion on the
252   target image. This process finely tunes the position of the landmark, such that even if the
253   training sets were indicated by two different observers, the variation in automatic landmark
254   indication is much smaller than the variation in manual landmark indication, averaging 0.27 mm
255   in our study (Supplementary Table S6; Supplementary Fig. S3).
256
257   A visual hallmark of the ability of spatially dense surface registration to reliably represent
258   anatomical structures is found in the crispness of "average shapes," constructed by averaging
259   together all registered surfaces in a study sample. Because the MeshMonk registration aligns
260   closely with the underlying anatomical structure, averages across the study samples continue to
261   cleanly resemble the structure and detail is not lost in the averaging process. As depicted in
262   Figure 2, consider the sample average of the 41 faces in this work and 100 mandible scans. In
263   the rigid-only averages, details are overly smoothed compared to the level of detail present in
264   the rigid plus non-rigid registration averages. For example, it is obvious to the naked eye that
265   the sharpness of the eyes, nose, philtrum and mouth for the facial average, and the alveolar
266   crest, mental foramen, and coronoid and condylar processes for the mandible, are clearly better
267   represented with the rigid plus non-rigid registration. Thus, non-expert readers can easily
268   evaluate the quality of dense-correspondence morphometrics research by looking at the
269   average surfaces, which are typically used in manuscript figures, with the understanding that
270   high quality registration leads to sharp average scans where anatomical positions of interest are
271   clearly defined.
272

6

273  In this study, we present MeshMonk, an open-source resource for intensive 3D phenotyping on
274  a large scale. Through dense-correspondence registration algorithms, like MeshMonk, we can
275  advance our ability to integrate genomic and phenomic data to explore variation in complex
276  morphological traits and answer evolutionary and clinical questions about normal-range
277  variation, growth and development, dysmorphology, and taxonomic classification.
278
279  **Materials and Methods**
280
281  Explanation of Meshmonk registration
282  The core functionality of the MeshMonk toolbox is implemented in C++, with a focus on
283  computational speed and memory to enable the processing of large datasets of 3D images.
284  Interaction with the toolbox is provided using MATLAB$^{TM}$, enabling an easy to use
285  implementation and visualization environment for the user. A schematic of the complete surface
286  registration algorithm is presented in Figure 3 and a short video of the registration on this
287  example face is also available at the following GitHub account
288  (https://github.com/juliedwhite/MeshMonkValidation/).
289
290  To initiate the process, a rigid registration based on the iterative closest point algorithm[27] is
291  performed to better align the template to the target surface. During the rigid registration, the
292  transformation model is constrained to changing the position (translation), orientation (rotation),
293  and scale of the template only. Subsequently, a non-rigid registration is done that will alter the
294  shape of the template to match the shape of the target surface. During the non-rigid registration,
295  a visco-elastic model is enforced, ensuring that points that lie close to each other move
296  coherently[8]. At any iteration during the registration, for both the rigid and non-rigid registration
297  steps, correspondences are updated by using pull-and-push forces (symmetrical
298  correspondences)[28] and a weighted k-neighbor approach (Supplementary Fig. S4)[8].
299
300  3D surface images typically contain artifacts such as holes and large triangles indicating badly
301  captured or missing parts. Any correspondence to such artifacts is meaningless and are
302  indicated as correspondence outliers, not to be considered when updating the transformation
303  model, though they are consistently transformed along with the inliers. The MeshMonk toolbox
304  allows for the identification of outliers either deterministically or stochastically, or a combination
305  of both. In each iteration, correspondences are updated and outliers are identified, then an
306  updated transformation model is used. The smoothness of the transformation model is
307  parametrized by convolving the displacement vectors between corresponding points with a
308  Gaussian[29]. The amount of smoothing is high (multiple Gaussian convolution runs) at the
309  beginning iterations, when correspondences are still noisy and hard to define, and reduces
310  gradually towards the later iterations, when correspondences are more accurately defined.
311
312  Parameters and tuning
313  Given a dataset of 3D images of interest, the entire MeshMonk procedure can be optimized by
314  setting a variety of parameters in the toolbox, and a parameter tuning can be done based on
315  two "quality" measures. First, a quality of "shape fit" is defined as the root mean squared
316  distance of all template points to the target surface after registration. This essentially measures
317  how well the shape of the template was adapted to the target shape and can be measured over
318  multiple images to deduct an overall quality of shape fit from the dataset. Second, an indication
319  of the consistency of point indications across the same dataset is obtained following the
320  principle of minimum description length in shape modelling[30]. Given two models explaining the
321  same amount of variance, the model requiring fewer parameters is favored, or given two models
322  with the same number of parameters, the one explaining more variance in the data is favored.
323  To this end, a principal component analysis (PCA) is used to assess registration quality

324  because, if the point indications were performed consistently, few PCs are required to explain
325  variation in the registration results. A parameter tuning was done for the facial data in this work
326  prior to the validation and is described in the Supplemental Methods.
327
328  Validation sample and data curation
329  Our collaborative group has recruited participants through several studies at Pennsylvania State
330  University, recruited at the following locations: State College, PA (IRB 44929 and 4320); New
331  York, NY (IRB 45727); Urbana-Champaign, IL (IRB 13103); Dublin, Ireland; Rome, Italy;
332  Warsaw, Poland; and Porto, Portugal (IRB 32341). All procedures were performed in
333  accordance with the relevant guidelines and regulations from the Pennsylvania State University
334  Institutional Review Board and all participants signed a written informed consent form before
335  participation. Participants additionally gave optional informed consent for publication of their
336  images in a variety of formats, including online open-access publications.
337
338  Stereo photogrammetry was used to capture 3D facial surfaces of N~6,000 participants using
339  the 3dMD Face 2-pod and 3-pod systems (3dMD, Atlanta, GA). This well-established method
340  generates a dense 3D point cloud representing the surface geometry of the face from multiple
341  2D images with overlapping fields of view. During photo capture, participants were asked to
342  adopt a neutral facial expression with their mouth closed and to gaze forward, following
343  standard facial image acquisition protocols[31].
344
345  Manual placement of validation landmarks
346  Of the larger sample, N=41 surface images were chosen at random for validation, excluding
347  participants who reported facial surgery or injury. These images were diverse with respect to
348  sex, age, height, weight, and 3D camera system used (Supplementary Table S7). 3dMDpatient
349  was used to record the 3D coordinates of 19 standard landmarks (7 midline and 12 bilateral)
350  from each unaltered surface in wavefront.obj format (Supplementary Fig. S5; Supplementary
351  Table S8). Two independent observers placed landmarks three times each, with at least 24
352  hours in-between landmarking sessions, resulting in six total landmark indications for each facial
353  image. For each individual, we checked for gross landmark coordinate errors before analysis. In
354  the subsequent analysis, $A_{ML}$ represents the average manual landmarks from observer A, $B_{ML}$
355  represents the average manual landmarks from observer B, while the combined average of all
356  six manual landmark indications is denoted as $C_{ML}$.
357
358  Automatic placement of validation landmarks
359  To obtain automatic indications of the 19 validation landmarks, each of the validation faces was
360  registered using MeshMonk and the manual landmark placements were transferred to the
361  registered face by coordinate conversion (Figure 4A)[32]. Because the registered faces are now in
362  the same coordinate system as the original template, we can subsequently transfer the manual
363  landmark indications to the original pre-registration template, giving a set of 41 x 2 observers x 3
364  indications = 246 manual landmark positions on the template scan (Figure 4B). One by one,
365  each face was left out while averaging the other 40 landmark placements to "train" the
366  automatic landmarks (Figure 4C). These averages were then transferred onto the left-out
367  (target) face, resulting in the automatic placement of the validation landmarks using a "training"
368  set that did not include the target face (Figure 4D). Further detail on this process can be found
369  in the Supplemental Methods.
370
371  The placement of automatic landmarks was performed three times: once using the average of
372  observer A's manual landmark indications as input ($A_{Auto}$), again using the average of observer
373  B's manual landmark indications ($B_{Auto}$), and a final time using the combined average of all six

374     manual landmark indications from both observers ($C_{Auto}$). This process resulted in three
375     placements of automatic landmarks for comparison.
376
377     <u>Accuracy</u>
378     We assessed the accuracy of the MeshMonk automatic landmark placements by calculating the
379     root mean squared error (RMSE) between manual and automatic coordinates. We also
380     calculated Bland-Altman[33] and Intraclass Correlation Coefficient (ICC)[34] statistics to compare
381     the manual and automatic landmark indications. The Bland-Altman method is preferred over
382     correlation or regression as it is less influenced by the variance of the sample and ICC is
383     preferred because it tests both the degree of correlation and agreement between methods. We
384     additionally compared estimates of centroid size calculated using each method and performed
385     an ANOVA on the centroid size calculations, with individual, observer, method, and individual x
386     observer as predictors, to determine if variation in centroid size could be attributable to variation
387     in landmarking method.
388
389     We utilized several methods to determine if the variance structures produced by the two
390     methods were similar. Fitting a MANOVA estimates the variance explained, in correlated
391     outcome variables, by various factors included in the model. Here, we performed MANOVAs
392     separately on the GPA-aligned average manual landmark indications from each observer ($A_{ML}$
393     and $B_{ML}$) as well as on the GPA-aligned automatic landmark indications trained using the
394     average of each observer's three landmark placements ($A_{Auto}$ and $B_{Auto}$), with image and
395     observer as predictors in both tests. By comparing the results of these two tests, we can
396     determine how the explanation of shape variance changes given a different landmarking
397     method. To directly determine if any variance in shape was attributable to landmarking method,
398     we combined the average manual landmark placements of each observer with the automatic
399     placements trained using each of these averages and aligned them using GPA ($A_{ML}$, $B_{ML}$, $A_{Auto}$,
400     and $B_{Auto}$). We then tested the shape variation in this combined space as the response in a
401     MANOVA, with individual, observer, method, and individual x observer as factors.
402
403     <u>Reliability</u>
404     We calculated the manual landmarking intra-observer error, the variation between indications
405     taken at different times by the same individual, as the standard deviation between the *x*, *y*, and
406     *z* coordinates of each observer's manual landmarking indications. The inter-observer error, the
407     difference between manual landmark indications made by different individuals, was calculated
408     as the standard deviation between each observer's average *x*, *y*, and *z* coordinates ($A_{ML}$ vs.
409     $B_{ML}$). As an additional method to understand the variation present in the manual landmark
410     indications only, we performed a MANOVA after GPA-aligning the six manual landmarking
411     indications[35]. Study individual, observer, and landmarking iteration were used as factors and
412     landmark configuration as the response.
413
414     To determine if the automatic indication process was more or less variable than manual
415     landmarking, we compared the inter-observer error calculated using only the manual landmarks
416     ($A_{ML}$ vs. $B_{ML}$) to the standard deviation between one observer's manual landmarks and the
417     automatic landmarks trained using the other observer's manual placements ($A_{ML}$ vs. $B_{Auto}$ and
418     $A_{Auto}$ vs. $B_{ML}$), as if the automatic indications replaced the manual indications in a calculation of
419     inter-observer error. A paired T-test was used to determine whether the "inter-observer errors"
420     calculated using the automatic indications were significantly different than the error calculated
421     using only the manual indications. Standard deviation values calculated using both automatic
422     placements ($A_{Auto}$ vs. $B_{Auto}$) were compared to manual landmarking inter-observer error to
423     illustrate the variance of automatic landmark indications. Levene's test[36] was performed to
424     determine if the variances of the inter-observer errors calculated using the manual landmarks

425  were equal to the standard deviation between the automatic landmarks (the null hypothesis).
426  Levene's test was chosen because the distribution of standard deviation values was non-
427  normal.
428
429  All validation analyses were performed in R using the Geomorph[37], BlandAltmanLeh
430  (https://cran.r-project.org/web/packages/BlandAltmanLeh/BlandAltmanLeh.pdf), and ICC
431  (https://cran.r-project.org/web/packages/ICC/ICC.pdf) packages, as well as packages for data
432  manipulation (readxl, reshape2, plyr, car, data.table, dplyr, broom) and graphing (ggplot2,
433  GGally, GGpubr). Centroid sizes were calculated using Geomorph and MANOVAs for shape
434  variation were implemented using the ProcD.lm function from Geomorph[37,38].
435

436  **Data Availability Statement**
437
438  The informed consent with which the data were collected does not allow for dissemination of
439  identifiable data to persons not listed as researchers on the IRB protocol. Thus, the full surface
440  3D facial images used for validation cannot be made publicly available. In the interest of
441  reproducibility, we have provided the 19 manual and automatic landmarks used for validation as
442  well as the code used to analyze them. These data are available in the following GitHub
443  repository: https://github.com/juliedwhite/MeshMonkValidation/. The MeshMonk code and
444  tutorials are available at https://github.com/TheWebMonks/meshmonk.
445

446  **References**
447

448  1.  Houle, D., Govindaraju, D. R. & Omholt, S. Phenomics: The next challenge. *Nat. Rev.*
449      *Genet.* **11,** 855–866 (2010).
450  2.  Walter, T. *et al.* Visualization of image data from cells to organisms. *Nat. Methods* **7,**
451      S26–S41 (2010).
452  3.  Fagertun, J. *et al.* 3D facial landmarks: Inter-operator variability of manual annotation.
453      *BMC Med. Imaging* **14,** 35 (2014).
454  4.  Toma, A. M., Zhurov, A. I., Playle, R., Ong, E. & Richmond, S. Reproducibility of facial
455      soft tissue landmarks on 3D laser-scanned facial images. *Orthod. Craniofacial Res.* **12,**
456      33–42 (2009).
457  5.  Weinberg, S. M., Scott, N. M., Neiswanger, K., Brandon, C. A. & Marazita, M. L. Digital
458      three-dimensional photogrammetry: Evaluation of anthropometric precision and accuracy
459      using a Genex 3D camera system. *Cleft Palate-Craniofacial J.* **41,** 507–518 (2004).
460  6.  von Cramon-Taubadel, N., Frazier, B. C. & Mirazon-Lahr, M. The problem of assessing
461      landmark error in geometric morphometrics: Theory, methods, and modifications. *Am. J.*
462      *Phys. Anthropol.* **134,** 24–35 (2007).
463  7.  Wong, J. Y. *et al.* Validity and reliability of craniofacial anthropometric measurement of
464      3D digital photogrammetric images. *Cleft Palate-Craniofacial J.* **45,** 232–239 (2008).
465  8.  Snyders, J., Claes, P., Vandermeulen, D. & Suetens, P. Development and comparison of
466      non-rigid surface registraion and extensions. (KU Leuven, 2014).
467  9.  Claes, P. A robust statistical surface registration framework using implicit function
468      representations: application in craniofacial reconstruction. (KU Leuven, 2007).
469  10. Claes, P., Walters, M. & Clement, J. Improved facial outcome assessment using a 3D
470      anthropometric mask. *Int. J. Oral Maxillofac. Surg.* **41,** 324–330 (2012).
471  11. Hutton, T. J., Buxton, B. F., Hammond, P. & Potts, H. W. W. Estimating Average Growth
472      Trajectories in Shape-Space Using Kernel Smoothing. *IEEE Trans. Med. Imaging* **22,**
473      747–753 (2003).
474  12. Andresen, P. R. & Nielsen, M. Non-rigid registration by geometry-constrained diffusion.
475      *Med. Image Anal.* **5,** 81–88 (2001).

13. Claes, P., Hill, H. & Shriver, M. Towards DNA-based facial composites: preliminary results and validation. *Forensic Sci. Int.* **13,** 208–216 (2014).

14. Matthews, H. *et al.* Estimating age and synthesising growth in children and adolescents using 3D facial prototypes. *Forensic Sci. Int.* **286,** 61–69 (2018).

15. Imaizumi, K. *et al.* Three-dimensional analyses of aging-induced alterations in facial shape: a longitudinal study of 171 Japanese males. *Int. J. Legal Med.* **129,** 385–393 (2015).

16. Blanz, V. & Vetter, T. A morphable model for the synthesis of 3D faces. in *Proceedings of the 26th annual conference on computer graphics and interactive techniques* 187–194 (ACM Press/Addison-Wesley Publishing Co, 1999). doi:10.1145/311535.311556

17. Baynam, G. *et al.* Phenotyping: Targeting genotype's rich cousin for diagnosis. *J. Paediatr. Child Health* **51,** 381–386 (2015).

18. Hammond, P. *et al.* Discriminating Power of Localized Three-Dimensional Facial Morphology. *Am. J. Hum. Genet.* **77,** 999–1010 (2005).

19. Hutton, T. J., Buxton, B. F. & Hammond, P. Automated registration of 3D faces using dense surface models. in *British Machine Vision Conference* (eds. Harvey, R. & Bangham, A.) 1–10 (Citeseer, 2003). doi:10.5244/C.17.45

20. Wei, R., Claes, P., Walters, M., Wholley, C. & Clement, J. G. Augmentation of linear facial anthropometrics through modern morphometrics: A facial convexity example. *Aust. Dent. J.* **56,** 141–147 (2011).

21. Aldridge, K., Boyadjiev, S. A., Capone, G. T., DeLeon, V. B. & Richtsmeier, J. T. Precision and error of three-dimensional phenotypic measures acquired from 3dMD photogrammetric images. *Am. J. Med. Genet.* **138A,** 247–253 (2005).

22. Claes, P. *et al.* Sexual dimorphism in multiple aspects of 3D facial symmetry and asymmetry defined by spatially dense geometric morphometrics. *J. Anat.* **221,** 97–114 (2012).

23. Li, M. *et al.* Rapid automated landmarking for morphometric analysis of three-dimensional facial scans. *J. Anat.* **230,** 1–12 (2017).

24. Subburaj, K., Ravi, B. & Agarwal, M. Automated identification of anatomical landmarks on 3D bone models reconstructed from CT scan images. *Comput. Med. Imaging Graph.* **33,** 359–368 (2009).

25. De Jong, M. A. *et al.* An automatic 3D facial landmarking algorithm using 2D gabor wavelets. *IEEE Trans. Image Process.* **25,** 580–588 (2016).

26. De Jong, M. A. *et al.* Ensemble landmarking of 3D facial surface scans. *Sci. Rep.* **8,** 1–11 (2018).

27. Besl, P. J. & McKay, N. D. A method for registration of 3-D shapes. *IEEE Trans. Pattern Anal. Mach. Intell.* **14,** 239–256 (1992).

28. Redert, A., Kaptein, B., Reinders, M., van den Eelaart, I. & Hendriks, E. Extraction of semantic 3D models of human faces from stereoscopic image sequences. *Acta Stereol.* **18,** 255–264 (1999).

29. Bro-Nielsen, M. Medical image registration and surgery simulation. (Technical University of Denmark, 1996). doi:0909-3192

30. Davies, R. H., Twining, C. J., Cootes, T. F., Waterton, J. C. & Taylor, C. J. A minimum description length approach to statistical shape modeling. *IEEE Trans. Med. Imaging* **21,** 525–537 (2002).

31. Heike, C. L., Upson, K., Stuhaug, E. & Weinberg, S. M. 3D digital stereophotogrammetry: a practical guide to facial image acquisition. *Head Face Med.* **6,** 18 (2010).

32. Hille, E. *Analytic Function Theory, Volume I.* (AMS Chelea Publishing Company, 1982).

33. Altman, D. G. & Bland, J. M. Measurement in Medicine: The Analysis of Method Comparison Studies. *Stat.* **32,** 307–317 (1983).

34. Fisher, R. A. *Statistical methods for research workers*. *Biological Monographs and*

527      *Manuals* (Oliver & Boyd, 1925). doi:10.1056/NEJMc061160
528    35.    Rohlf, F. J. & Slice, D. Extensions of the procrustes method for the optimal
529          superimposition of landmarks. *Syst. Zool.* **39,** 40–59 (1990).
530    36.    Levene, H. in *Contributions to Probability and Statistics: Essays in Honor of Harold*
531          *Hotelling* (eds. Olkin, I. & Hotelling, H.) 278–292 (Stanford University Press, 1960).
532    37.    Adams, D. C. & Otárola-Castillo, E. Geomorph: An r package for the collection and
533          analysis of geometric morphometric shape data. *Methods Ecol. Evol.* **4,** 393–399 (2013).
534    38.    Collyer, M. L., Sekora, D. J. & Adams, D. C. A method for analysis of phenotypic change
535          for phenotypes described by high-dimensional data. *Heredity (Edinb).* **115,** 357–365
536          (2015).

537

546
**Author Contributions**

549 JW performed all landmark based analyses and landmarked the 3D scans used for validation
550 with AZ. PC and AO performed the parameter tuning on the facial data and provided the
551 automatic landmark indications. JW, AO, and HM wrote the first draft of the manuscript under
552 supervision of PC. HM, YF, and TP provided input and images using mandible scans. PC and
553 JW conceptualized the design of the study. OE, SV, and MS provided input throughout the
554 analyses and writing process. JS developed the MeshMonk code.

555

561

570
571
572
573
574
575
576
577

578 **Figures**



579
580 **Figure 1. Facial template registration.** The template (left), built as the average of more than
581 8000 admixed facial scans, can easily wrap onto any face (three example faces on the right),
582 accurately representing its particular traits. This allows for the explanation of any face in the
583 template's coordinates, enabling a spatially-dense analysis between any registered surfaces.
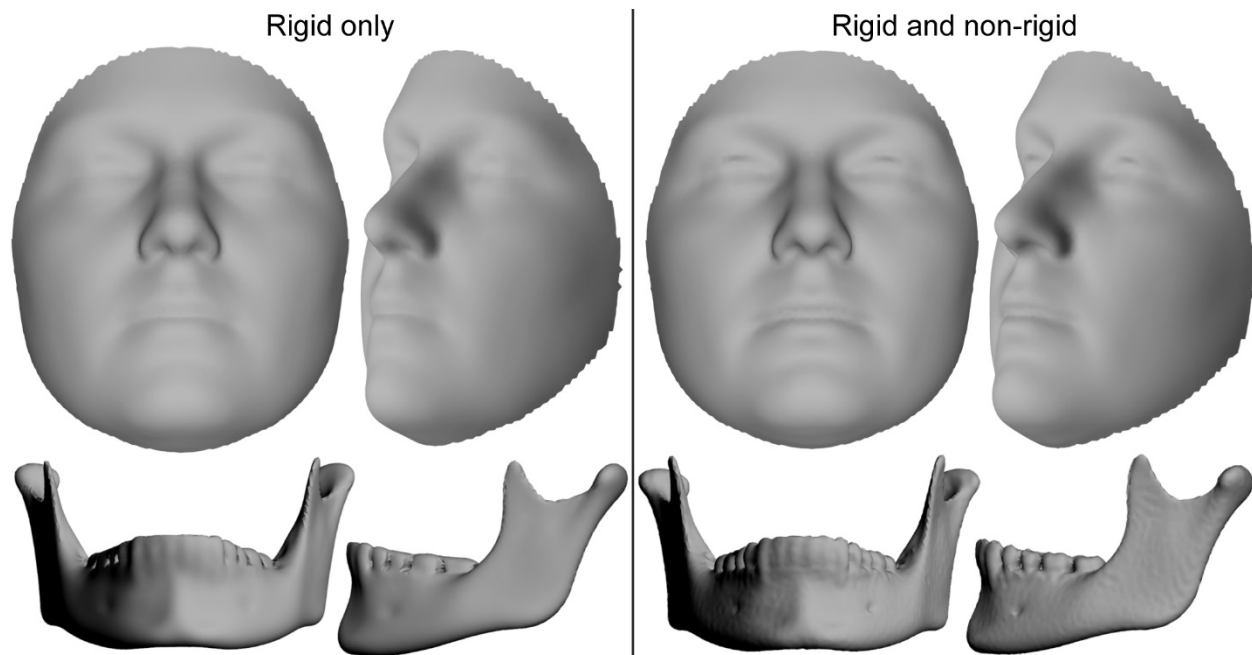
Rigid only | Rigid and non-rigid



584
585 **Figure 2. Comparison of rigid and non-rigid registration algorithms.** Sample averages
586 using the 41 validation faces and 100 mandible scans. Scans were registered using rigid
587 registration only (left) and then simply mapped exactly to their closest point on the target
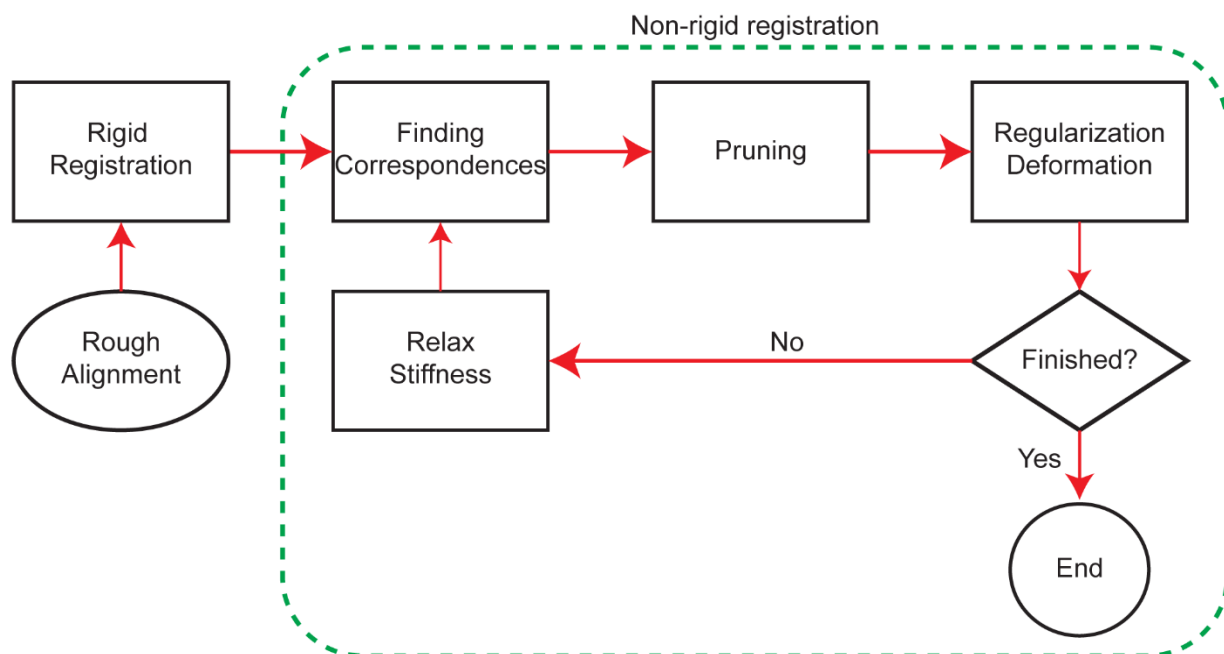588 surfaces or mapped using rigid plus non-rigid (visco-elastic) registration (right).
589



590
591 **Figure 3. Schematic of the MeshMonk's surface registration algorithm.** MeshMonk uses an
592 initial rigid registration based on the ICP algorithm. This step might require an initial rough
593 alignment to ensure similar orientation, which can be done by placing few landmarks on the
594 target surface. Then, the symmetrical weighted k-neighbor correspondences are found, and
595 outliers are detected and removed. Finally, the visco-elastic transformation is applied. This is
596 performed in an iterative manner, until either a pre-set number of iterations or a pre-set amount
597 of coverage (e.g. a pre-defined root mean squared distance of all template points to the target

14

598 surface after the transformation) has been reached. Otherwise, the correspondences are
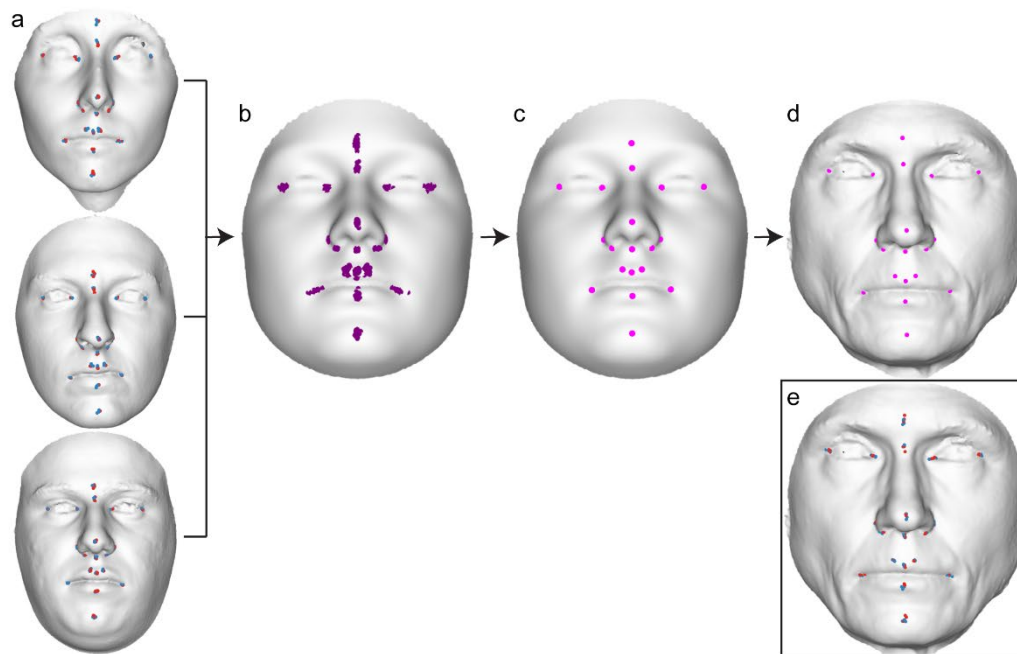599 updated and the non-rigid registration starts over.
600



601
602 **Figure 4. Depiction of automatic landmark indication. (a)** Each facial scan was manually
603 landmarked six times, three times each by two observers (red and blue points). **(b)** These
604 iterations were then averaged together and are placed on the template (purple points). **(c)** The
605 average of all but the test face (N=40) placements on the template, serving as the foundation for
606 the automatic landmark placements (magenta points). **(d)** Coordinate conversions, described in
607 more detail in the Supplemental Methods, is used to subsequently transfer the automatic
608 landmark placements from the template to the target (left-out) surface, serving as the automatic
609 landmark indication for the target surface (magenta points). **(e)** The manual landmark
610 indications from two observers (red and blue points) for the shown example face, for
611 comparison to the automatic indication in (d).

612
613 **Tables**
614
615 **Table 1. Root mean squared error between manual and automatic landmarks**. Root mean
616 squared error (mm) between the manual and automatic landmark indications. Values are
617 presented for each axis, averaged across all faces, as well as averaged across the axes
618 (mean).

| Landmark | $A_{ML}$ vs. $A_{Auto}$ | | | | $B_{ML}$ vs. $B_{Auto}$ | | | | $C_{ML}$ vs. $C_{Auto}$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | X | Y | Z | Mean | X | Y | Z | Mean | X | Y | Z | Mean |
| Alar curvature left | 0.17 | 0.54 | 0.59 | 0.44 | 0.19 | 0.65 | 0.76 | 0.53 | 0.16 | 0.52 | 0.61 | 0.43 |
| Alar curvature right | 0.18 | 0.53 | 0.67 | 0.46 | 0.18 | 0.58 | 0.61 | 0.46 | 0.17 | 0.52 | 0.57 | 0.42 |
| Chelion left | 1.23 | 0.70 | 0.64 | 0.86 | 1.26 | 0.74 | 0.66 | 0.88 | 1.11 | 0.71 | 0.61 | 0.81 |
| Chelion right | 0.93 | 0.70 | 0.53 | 0.72 | 1.15 | 0.65 | 0.62 | 0.81 | 0.98 | 0.66 | 0.55 | 0.73 |
| Crista philtri left | 0.69 | 0.85 | 0.44 | 0.66 | 0.89 | 1.01 | 0.51 | 0.80 | 0.75 | 0.89 | 0.45 | 0.70 |
| Crista philtri right | 0.66 | 0.95 | 0.50 | 0.70 | 1.00 | 1.13 | 0.47 | 0.87 | 0.76 | 1.00 | 0.44 | 0.73 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Endocanthion left* | 0.84 | 0.64 | 0.53 | 0.67 | 0.83 | 0.62 | 0.42 | 0.62 | 0.78 | 0.54 | 0.40 | 0.57 |
| *Endocanthion right* | 1.05 | 0.74 | 0.62 | 0.80 | 1.09 | 0.62 | 0.45 | 0.72 | 1.04 | 0.65 | 0.50 | 0.73 |
| *Exocanthion left* | 0.92 | 0.78 | 0.91 | 0.87 | 0.97 | 0.75 | 0.88 | 0.87 | 0.91 | 0.74 | 0.88 | 0.84 |
| *Exocanthion right* | 0.93 | 0.67 | 0.93 | 0.85 | 0.98 | 0.68 | 0.97 | 0.88 | 0.94 | 0.65 | 0.95 | 0.85 |
| *Glabella* | 0.52 | 1.43 | 0.60 | 0.85 | 0.55 | 1.46 | 0.59 | 0.87 | 0.48 | 1.31 | 0.56 | 0.78 |
| *Labiale inferius* | 0.52 | 0.75 | 0.56 | 0.61 | 0.50 | 0.71 | 0.38 | 0.53 | 0.46 | 0.72 | 0.48 | 0.55 |
| *Labiale superius* | 0.57 | 0.72 | 0.31 | 0.54 | 0.59 | 0.98 | 0.37 | 0.65 | 0.59 | 0.81 | 0.33 | 0.58 |
| *Nasion* | 0.37 | 1.10 | 0.51 | 0.66 | 0.42 | 1.04 | 0.48 | 0.65 | 0.35 | 0.97 | 0.47 | 0.60 |
| *Pogonion* | 0.48 | 1.08 | 0.45 | 0.67 | 0.54 | 1.12 | 0.42 | 0.69 | 0.43 | 1.00 | 0.38 | 0.60 |
| *Pronasale* | 0.44 | 0.71 | 0.33 | 0.49 | 0.45 | 0.57 | 0.28 | 0.44 | 0.40 | 0.56 | 0.28 | 0.41 |
| *Subalare left* | 0.78 | 0.47 | 0.54 | 0.60 | 0.79 | 0.44 | 0.64 | 0.62 | 0.73 | 0.43 | 0.56 | 0.57 |
| *Subalare right* | 0.75 | 0.46 | 0.76 | 0.66 | 0.67 | 0.50 | 0.52 | 0.56 | 0.65 | 0.43 | 0.60 | 0.56 |
| *Subnasale* | 0.33 | 0.46 | 0.33 | 0.37 | 0.35 | 0.68 | 0.33 | 0.46 | 0.32 | 0.48 | 0.26 | 0.35 |
| *Mean* | 0.65 | 0.75 | 0.57 | 0.66 | 0.71 | 0.79 | 0.55 | 0.68 | 0.63 | 0.72 | 0.52 | 0.62 |

**Table 2. ANOVA of centroid sizes.** Results from an ANOVA with centroid size as the response variable and individual, observer, method and individual x observer as predictors.

| Variable | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
|---|---|---|---|---|---|
| *Individual* | 40 | 7936 | 198.39 | 130.407 | $<2 \times 10^{-16}$ |
| *Observer* | 2 | 0 | 0.23 | 0.154 | 0.857 |
| *Method* | 1 | 0 | 0 | 0.002 | 0.962 |
| *Individual x Observer* | 80 | 12 | 0.15 | 0.101 | 1.000 |
| *Residuals* | 122 | 186 | 1.52 | | |

**Table 3. MANOVA on manual and automatic landmarks.** Results from a single MANOVA using the average manual landmark indications from each observer ($A_{ML}$ and $B_{ML}$) and the automatic landmark indications using the observer level averages ($A_{Auto}$ and $B_{Auto}$).

| Variable | DF | SS | MS | $R^2$ | F | Z | Pr(>F) |
|---|---|---|---|---|---|---|---|
| *Method* | 1 | 0.0003 | 0.0003 | 0.0004 | 0.3463 | -2.2135 | 0.987 |
| *Individual* | 40 | 0.6522 | 0.0163 | 0.8778 | 20.2019 | 23.3507 | 0.001 |
| *Observer* | 1 | 0.0167 | 0.0167 | 0.0224 | 20.6396 | 11.4067 | 0.001 |
| *Individual x Observer* | 40 | 0.0085 | 0.0002 | 0.0114 | 0.2623 | 13.7253 | 0.001 |
| *Residuals* | 81 | 0.0654 | 0.0008 | 0.0880 | | | |
| *Total* | 163 | 0.7430 | | | | | |

16